

Gheorghe M.Panaitescu

**SISTEME TOLERANTE LA
DEFECTE**

Note de curs

**Universitatea “Petrol-Gaze” Ploiesti
Catedra Automatică si calculatoare
2009**

CUVÂNT ÎNAINTE

Ce reprezintă această lucrarea?

Lucrarea aceasta pe care o consultați este suportul cursului *Sisteme tolerante la defecte*, un ghid al expunerilor săptămânale pe care autorul le-a ținut/le ține pentru anul IV de la specializarea Calculatoare a facultății Inginerie mecanică și electrică, de cele mai multe ori în fața unei asistențe modeste numeric. Este totodată un text de ajutor în pregătirea pentru verificarea pe care studenții o susțin la sfârșitul primului semestru al ultimului an de studii.

Desigur, nu este exclusă o lectură interesantă și interesată a unor terțe persoane, care nu sunt nici studenți și nici nu predau disciplina atât de atractivă a toleranței la defecte.

Există, desigur, o diferență între versiunea vorbită și cea scrisă a acestui curs: scrisul reclamă concizie, vorbitul permite detalierea unor aspecte care merită sau trebuie detaliate. Comunicarea verbală poate beneficia de întreruperi utile din partea celor care ascultă, ceea ce autorul și-a dorit foarte mult, dar, din păcate, s-a întâmplat foarte rar.

Despre conținutul acestei lucrări, despre surse

Nu veți găsi în paginile care urmează foarte multe lucruri originale. Pe tema originalității există în zilele noastre o (să-i spunem) dezbatere asupra dozei de noutate din cursurile scrise de universitarii români. Sunt publicate adesea opinii ale unor oameni care n-au scris vreodată un rând dar dau cu generozitate indicații altora, mai ales când vor să arate cât de incorect este profesorul X în utilizarea surselor bibliografice. Vor fi urme de adevăr în cele scrise ocazional în acele articole de campanie, altfel foarte combative. Nu sunt cel chemat a arbitra în această chestiune. Nu mă pot reține însă de la a face câteva precizări despre lucrarea de față.

Subiectele tratate aici urmează linia care se regăsește din plin în cursurile prezentate studenților unor universități din Statele Unite ale Americii. Ca să scutesc pe cititor, indiferent cine ar fi el, de orice efort detectivistic, fac trimiteri la cursurile similare ținute la University of Massachusetts (cursul ECE 655), la University of Wisconsin (cursul ECE 753), la University of Idaho (cursul CS449_549) sau la University of Illinois (cursul ECE 542); un “search” pe Internet pe codul cursurilor și pe numele universităților poate duce cititorul direct la aceste surse *principale* utilizate în versiunea curentă a acestor Note de curs. Am scris “principale” pentru că în afara slide-urilor sau paginilor postate la locurile indicate, alte sute de pagini au fost parcurse pentru lămurirea unor

detalii, deoarece, este stiut, este absurd să predai o disciplină pe care nu o înțelegi...

Am apelat la surse americane pentru motivul lesne de înțeles: tehnologia calculatoarelor și preocuparea pentru siguranța calculului nu sunt altundeva mai înalte decât în America.

Se poate adăuga ceva original la ceea ce industria și comunitatea academică americană a creat în domeniu? Poate de aici înainte, după lectura acestor Note de curs...

Domeniul este pasionant, punctele de vedere sunt multiple și variate, bibliografia este practic nepuizabilă. De aceea, această versiune a lucrării (a patra, a concea?) este radical diferită de acelea din anii trecuți și nu mă îndoiesc că este diferită și de cea (posibilă) de anul viitor...

Un *feedback* cu observații asupra acestui curs, la adresa de e-mail a autorului ar fi aproape sigur foarte util pentru realizarea versiunii viitoare.

Ce consecințe ar putea avea lectura acestei lucrări?

Cititorul va găsi aici noțiuni fundamentale despre cele patru tipuri de redundanțe care se folosesc separat sau în combinații pentru realizarea toleranței la defecte: redundanțele hardware, redundanțele software, redundanțele informaționale și redundanțele temporale.

Nu este exagerată credința că după audierea cursului de *Sisteme tolerante la defecte* și după lectura acestor Note, studentul absolvent va fi mai aproape de mecanismele utilizate pentru a prelungi la limita posibilului funcționarea corectă a unui calculator, a unui program de calcul, pentru a mări siguranța serviciului unei baze de date, a comunicării în rețelele de calculatoare.

Gheorghe M.Panaitescu
gmpanaitescu@upg-ploiesti.ro

Ploiesti, 5 octombrie 2009

C U P R I N S

CUVÂNT ÎNAINTE	3
CUPRINS	5
INTRODUCERE	7
TOLERANTA LA DEFECTE HARDWARE	17
TOLERANTA LA DEFECTE SOFTWARE	35
REDUNDANTE INFORMATIONALE	55
REDUNDANTE TEMPORALE	73
SISTEME DE DISCURI TOLERANTE LA DEFECTE	107
REPLICAREA DATELOR PENTRU TOLERANTA LA DEFECTE	123
TOLERANTA LA DEFECTE ÎN RETELE	133
ANEXA 1: ELEMENTE DE TEORIA PROBABILITĂȚILOR SI DE STATISTICĂ MATEMATICĂ	155
B I B L I O G R A F I E	165

INTRODUCERE

Definirea obiectului disciplinei

Ideal, un sistem tolerant la defecte este un sistem capabil a executa corect sarcinile lui de calcul sau de altă natură indiferent dacă apar defecte *hardware* sau *software*.

În practică, niciodată nu poate fi garantată în orice împrejurare executarea acelor sarcini fără întrerupere.

Discutiile de genul celor care urmează se limitează uzual la tipuri de defectări și de erori care sunt mai probabil să apară.

Aplicatii în care toleranța la defecte este necesară

Există în lumea reală aplicații considerate pe drept cuvânt critice sub aspect vital. Este vorba, de plidă, de aeronave, de reactoarele nucleare, de unele instalații chimice și de echipamentele medicale. O funcționare proastă a unui calculator atașat unor asemenea aplicații poate duce la catastrofe. Probabilitatea de defectare a sistemelor de calcul din această categorie trebuie să fie extrem de scăzută, sub unu la un miliard pe ora de operare.

Ambianța aspră este alt element care impune un grad de a tolera defectele. Un sistem de calcul care operează în condiții vitrege, în care sunt prezente perturbările electromagnetice, bombardamentul cu particule încărcate sau neutre electric, sau altele asemenea trebuie să fie tolerant la defecte. Un număr mare de disfuncții datorate factorilor de mediu, altminteri destul de probabile, poate face ca sistemul să nu producă rezultate utile decât dacă are încorporată încă din faza de concepție o doză de toleranță la defecte.

Din alt unghi, unele sisteme foarte complexe constau în milioane de dispozitive elementare. Fiecare dispozitiv fizic component are o anumită probabilitate de a esua. Un număr foarte mare de dispozitive implică o probabilitate de mers neconform încă mai mare. Sistemul poate manifesta defectări cu o așa frecvență încât el poate deveni inutilizabil.

Măsuri ale toleranței la defecte

Este foarte important a avea definite *măsuri* care să cuantifice nivelul toleranței la defecte. Ca în orice domeniu, o măsură este o abstracție matematică care exprimă numai un aspect al naturii obiectelor.

În domeniul siguranței în funcționare a sistemelor, sunt desigur unele măsuri considerate deja tradiționale. Se consideră mai întâi că sistemul poate fi în una

din două stări mutual exclusive: *functional* sau *disfunct*. De exemplu, un bec este fie bun, fie ars; un fir este fie continuu, fie întrerupt. În legătură strictă cu aceste două stări se utilizează două măsuri deja tradiționale: *fiabilitatea* și *disponibilitatea*.

Fiabilitatea, notată cu $R(t)$ este probabilitatea ca sistemul să fie functional pe durata intervalului $[0, t]$, el fiind functional la momentul $t = 0$.

Disponibilitatea, notată cu $A(t)$ este cota parte din timp în care sistemul este functional în intervalul $[0, t]$. Ca măsură complementară, se definește o disponibilitate punctuală $A_p(t)$: probabilitatea ca sistemul să fie functional la momentul t . O măsură înrudită este și *MTTF – Mean Time To Failure* – care este timpul mediu cât sistemul este functional înainte ca el să se defecteze pentru ca apoi să fie reparat sau înlocuit.

Sunt însă necesare și alte măsuri în completarea acestora. Presupunerea că sistemul poate fi doar în stările “functional” sau “disfunct” este foarte restrictivă. De exemplu, un procesor cu una din cele câteva sute de milioane de porti blocată în valoarea logică 0 și restul functionale poate fi conjunctural și functional, și disfunct. Poarta aceasta defectă poate afecta iesirea procesorului o dată la 25.000 de ore de utilizare. Procesorul nu este lipsit de orice defectiune dar nu poate fi calificat ca *disfunct*. Pentru a caracteriza astfel de stări, sunt necesare mijloace de apreciere cantitativă suplimentare pe lângă deja tradiționalele *fiabilitate* și *disponibilitate*. În continuare sunt discutate câteva măsuri capabile să cuprindă și alte nuanțe.

Fiabilitatea prin capacitate este tot o probabilitate, probabilitatea ca o anumită capacitate a sistemului (măsurabilă, de pildă un randament) la timpul t să depășească un prag dat. O altă extindere a gamei de măsuri ia în considerare totul din perspectiva aplicației. Aceasta duce la definirea măsurii cunoscută ca *performabilitate*.

Relativ la *capacitatea* unui sistem, fie aceasta capacitatea de calcul. Fie, de exemplu, un sistem cu N procesoare, sistem care se degradează *cu gratie*, adică nu brusc ci gradual. Sistemul se recuperează din starea de disfuncție a unora dintre procesoare și este utilizabil atât timp cât cel puțin un procesor este functional. Fie P_i probabilitatea ca i procesoare din cele N să fie functionale. Fiabilitatea sistemului este o sumă a acestor probabilități după indicele i , excluzând, desigur, valoarea $i = 0$.

$$R(t) = \sum_{i=1}^N P_i$$

Fie c capacitatea de calcul a unui procesor (de pildă numărul de task-uri de dimensiune fixă pe care le poate executa). Capacitatea de calcul a i procesoare este atunci $C_i = i \cdot c$. Capacitatea de calcul a sistemului este în consecință

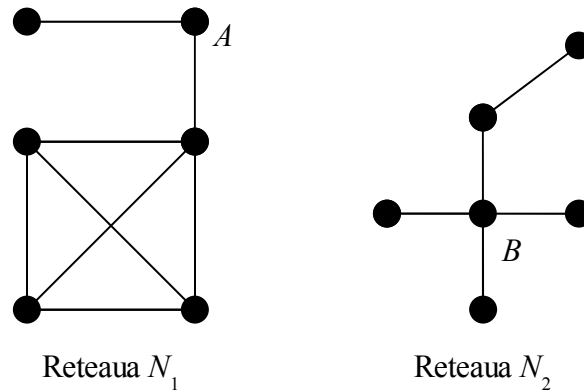
$$\sum_{i=1}^N C_i P_i .$$

Relativ la *performabilitate*, în legătură cu aplicația la care se referă se definesc niste *niveluri de îndeplinire*, L_1, L_2, \dots, L_n . Fiecare dintre acestea reprezintă un nivel de calitate a serviciului efectuat prin acea aplicație. De exemplu, L_i –

sistemul i nu “cade” pe durata T a misiunii lui. *Performabilitatea* este un vector $[P(L_1) P(L_2) \dots P(L_n)]$ în care $P(L_i)$ notează probabilitatea ca sistemul să funcționeze suficient de bine pentru a permite aplicației să atingă nivelul de îndeplinire L_i .

Măsuri ale conectivității unei rețele

O măsură a conectivității trebuie să se focalizeze pe rețeaua care leagă procesoarele. *Conectivitatea* clasică se referă la nodurile și liniile unei rețele de procesoare și reprezintă numărul minim de noduri, respectiv de linii, care trebuie să cadă înainte ca rețeaua să devină neconexă. Această măsură indică vulnerabilitatea rețelei la a deveni neconexă. O rețea care poate deveni neconexă prin disfuncția unui singur nod (poziționat critic) este potențial mult mai vulnerabilă decât o alta care poate deveni neconexă numai dacă mai multe noduri cad.



Exemple de conectivitate

Măsura rezilienței unei rețele

Conectivitatea clasică distinge între două stări ale rețelei: conexă și neconexă. Conectivitatea clasică nu spune nimic despre modul cum rețeaua se degradează pe măsură ce nodurile cad înainte ca rețeaua să devină neconexă. Sunt două măsuri posibile ale *rezilienței*:

- Distanța medie pe perechi de noduri
- Diametrul rețelei (maximul distanței între două noduri) cu probabilitățile disfuncțiilor pe noduri și/sau pe linii date

Alte măsuri pentru rețele

Ce se întâmplă când rețeaua încetează a mai fi conexă? O rețea care se divide într-o componentă (conexă) majoră și mai multe bucăți (conexe) mărunte poate încă să funcționeze. Dar care este funcționalitatea unei rețele care se divide într-

un număr mare de subrețele mici? Șansele de funcționare sunt mai reduse. De aceea, o altă măsură a rezilienței la defectare a unei rețele este probabilitatea distribuirii componentei majore (cele mai mari) la căderea unui nod, unei linii.

Redundanțe

Redundanțele sunt centrale în realizarea toleranței la defecte. Redundanta se poate defini ca rezultatul încorporării în sistem, prin proiect, a unor module suplimentare, în așa mod încât funcția sistemului să nu fie periclitată de apariția unei disfuncții. În continuare sunt studiate patru tipuri de redundanțe.

- A. *Redundante hardware*. Acest gen de redundanțe se crează atunci când se adaugă elemente hardware suplimentare pentru a depăși efectele disfuncției unei componente. Redundanta hardware poate fi *statică*, caz în care se realizează mascarea imediată a unei disfuncții. Un exemplu: se utilizează trei procesoare în loc de unul și fiecare execută aceeași funcție. Iesirea majoritară a acestor procesoare elimină iesirea greșită a unuia singur. În cazul redundanței hardware *dinamice*, componentele suplimentare sunt activate numai la apariția disfuncției unei componente curent activă, în funcțiune. Redundanta hardware *hibridă* este o combinație de metode specifice redundanțelor statice și redundanțelor dinamice.
- B. *Redundante software*. Redundanțele software se asigură prin echipe de programare multiple. Se scriu versiuni diferite de software pentru aceeași funcție, pentru aceeași aplicație. Redundanta software se bazează pe speranța că o astfel de diversitate dă siguranța că nu toate variantele de program vor esua pe același set de date de intrare.
- C. *Redundante informaționale*. Redundanțele informaționale se realizează prin adăugarea de biți la cei originali. Erorile în biți pot fi detectate și chiar corectate. S-au dezvoltat și se utilizează codurile detectoare și corectoare de erori. Redundanta informațională atrage după sine uneori redundanțe hardware menite a prelucra biții suplimentari.
- D. *Redundante temporale*. Redundanțele temporale înseamnă timp suplimentar pentru ca execuțiile esuate să poată fi repetate. Cele mai multe dintre esecuri sunt tranzitorii și efectul lor se atenuează după un timp. Dacă există timp suficient la dispoziție, modulul disfuncțiv tranzitoriu poate recupera și poate reface calculele afectate.

Clasificarea defectelor hardware

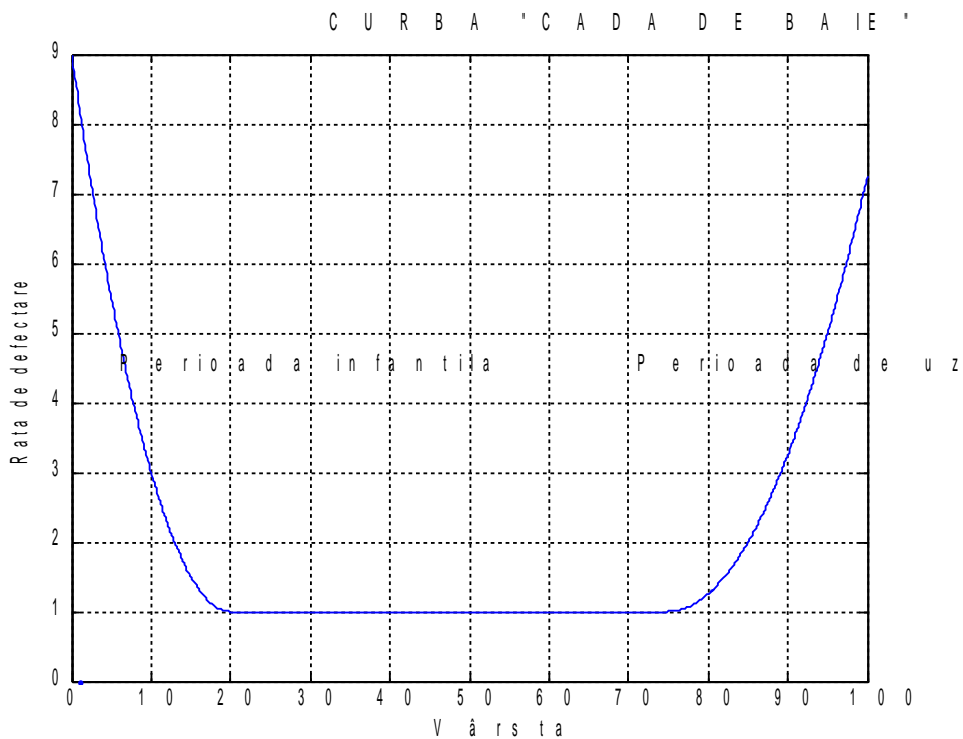
Sunt trei tipuri de defecte hardware.

Unele sunt *tranzitorii* și dispar după relativ puțin timp. Un exemplu ar putea fi o celulă de memorie al cărei conținut este schimbat datorită unei interferențe electromagnetice. Rescrierea ei cu conținutul corect face ca eroarea să dispară. Alte defecte sunt *permanente* și nu dispar niciodată, iar componenta trebuie reparată sau înlocuită.

Există și defecte *intermitente*, caz în care defectul componentei pendulează între o stare activă și o stare benignă. Se poate exemplifica cu cazul unei conexiuni slabe/slăbite.

Rata defectărilor

Rata la care o componentă manifestă disfuncții depinde de vârstă, de parametrii (micro)climatici locali, de eventualele socuri fizice cum ar fi cele date de tensiunea de alimentare, dar și de tehnologia de realizare concretă a acelei componente. Dependența de vârstă a ratei de defectare este ilustrată uzual de curba cadă-de-baie (v.figura).



La tinerete rata defectărilor este ridicată: sunt șanse bune ca unele module cu defecte de fabricație să treacă de controlul de calitate și să ajungă pe piață. Pe măsura trecerii timpului aceste unități/module sunt eliminate și pentru mare parte din viață sistemul manifestă o rată de defectare practic constantă. Pe măsură ce sistemul tinde să devină foarte vechi, efectele îmbătrânirii preiau primplanul și rata defectărilor devine din nou variabilă și crescătoare.

O formulă empirică pentru λ , rata de defectare

Rata de defectare se poate estima cu relația

$$\lambda = \pi_L \pi_Q (C_1 \pi_T \pi_V + C_2 \pi_E)$$

în care:

π_L – un factor de învățare care arată cât de matură este tehnologia

π_Q – un factor de calitate a fabricației între 0,25 și 20

π_T – un factor de temperatură (de la 0,1 la 1000), proportional cu $\exp(-E_a/kT)$ cu E_a energia de activare [eV] asociată cu tehnologia, k constanta lui Boltzmann și T temperatura în K .

π_V – un factor de stres prin tensiune pentru dispozitivele CMOS (de la 1 la 10, depinzând de tensiunea de alimentare și de temperatură); nu se aplică altor tehnologii (se pune egal cu 1)

π_E – un factor de soc ambiant: de la cca. 0,4 (mediu de aer condiționat) la 13 (ambiant dur)

C_1, C_2 – factori de complexitate, dependenti de numărul de porți pe chip și numărul de pini în pachet.

Detalii se găsesc în manualele de tip *Handbook*.

Impactul mediului ambiant

Un dispozitiv care lucrează în spațiul interaștrăl, spațiu care este plin de particule încărcate energetic și care este capabil să supună dispozitivele la variații drastice de temperatură este de așteptat să aibă o clacă mai des. Similar, calculatoarele din automobile (temperaturi ridicate și vibrații) și din aplicațiile industriale sunt susceptibile să se defecteze mai frecvent.

Defect și eroare

Un *defect* în sensul toleranței la defecte poate fi o disfuncție (mai curând locală) de hardware sau o greșală de software sau de programare. O *eroare* este o manifestare a unui defect (fault). Ca exemplu se consideră un circuit de adăuție în binar cu una din liniile de ieșire “agățată” la 1. Acesta-i un defect dar încă nu este o eroare. Defectul produce o eroare dacă circuitul este folosit și dacă rezultatul pe linia defectă ar trebui să fie 0 și nu 1.

Propagarea defectelor și erorilor

Atât defectele cât și erorile se pot difuza în sistem. Dacă de pildă un chip scurtcircuitază alimentarea cu energie la masă, asta poate determina ca chipurile apropiate să clacheze și ele. Erorile se difuzează pentru simplul motiv că o ieșire (eronată) a unui element procesor este utilizată frecvent ca intrare a altor elemente procesoare. Exemplul circuitului de adunare este ilustrativ: rezultatul eronat al circuitului poate fi utilizat în alte calcule, desigur cu propagarea erorii.

Zone de continentă

Pentru a limita situatii de genul mentionat putin mai sus, proiectantii încorporează în sistem niste asa-numite *zone de continență*. Cum? Prin bariere care reduc sansa ca un defect sau o eroare dintr-o zonă să se propage într-o altă zonă. O zonă de continentă poate fi creată uneori prin asigurarea de alimentări independente pentru fiecare zonă. Proiectantul unui sistem încearcă prin aceasta să izoleze electric o zonă de celelalte. O zonă de conținere sau de continență a erorilor poate fi creată prin utilizarea de module redundante si prin *voting* pe iesirile lor.

Timpu până la cădere – modelul analitic

Se consideră modelul următor:

Sistemul are N componente identice, toate operationale la momentul $t = 0$. Fiecare componentă rămâne operatională până când apare un defect. Toate defectele sunt permanente si apar în cele N componente independent unul de celălalt. Pentru o componentă oarecare, fie T durata de viață a acelei componente, adică timpul până la aparitia defectului considerat fatal. T este o variabilă aleatoare cu densitatea de probabilitate $f(t)$ si functia de repartitie (cumulativă) $F(t)$.

Într-o interpretare probabilistică, $F(t)$ este probabilitatea ca o componentă să esueze înainte de timpul t , $F(t) = \Pr(T \leq t)$. Functia $f(t)$ este un gen de rată momentană a căderii: $f(t)\Delta t = \Pr(t \leq T \leq t + \Delta t)$ cu Δt oricât de mic.

Ca orice functie densitate de probabilitate

$$\int_0^{\infty} f(t)dt = 1 \text{ si } f(t) \geq 0 \text{ pentru orice } t \geq 0$$

Cele două functii definite aici sunt în relatiile

$$f(t) = dF(t)/dt, F(t) = \int_0^t f(\tau) d\tau$$

Fiabilitatea si rata (hazardul) defectării

Fiabilitatea unei singure componente este

$$R(t) = \Pr(T > t) = 1 - F(t)$$

Probabilitatea (infinitesimală) de cădere a unei componente la momentul t , $p(t)$ este probabilitatea conditionată ca acea componentă să cadă la momentul t fiind functională în orice moment înainte de t .

$$\begin{aligned} p(t) &= \Pr(t \leq T \leq t + dt / T \geq t) = \\ &= \Pr(t \leq T \leq t + dt) / \Pr(T \geq t) = f(t)dt / [1 - F(t)] \end{aligned}$$

Rata defectărilor (sau rata hazardului) pentru o componentă la timpul t , este definită ca raportul $p(t)/dt$ si se notează cu

$$h(t) = f(t) / [1 - F(t)]$$

si deoarece $dR(t)/dt = -f(t)$

$$h(t) = -1/R(t) dR(t)/dt$$

Rata defectărilor constantă

Dacă rata defectărilor este constantă în timp, $h(t) = \lambda$, atunci relatiile

$$dR(t)/dt = -\lambda R(t); R(0) = 1$$

sunt o ecuatie diferentială si conditia ei initială. Solutia ecuatiei este

$$R(t) = e^{-\lambda t}$$

Ceea ce conduce la

$$f(t) = \lambda e^{-\lambda t} \text{ si } F(t) = 1 - e^{-\lambda t}$$

O rată a defectărilor constantă se obtine dacă si numai dacă T , durata de viață a componentei are o distributie exponentială.

Timpul mediu până la cădere

MTTF (Mean Time To Failure) este durata medie de viață, media variabilei aleatoare T

$$MTTF = E(T) = \int_0^{\infty} t f(t) dt$$

cu $E(\cdot)$ notatia pentru speranta matematică (espérance, expected value) a variabilei aleatoare trecută ca argument. Deoarece $dR(t)/dt = -f(t)$

$$MTTF = -\int_0^{\infty} t \frac{dR(t)}{dt} dt = [-t R(t)]_0^{\infty} + \int_0^{\infty} R(t) dt$$

Dar $-t R(t) = 0$ pentru $t = 0$. La fel si pentru $t = \infty$ deoarece $R(\infty) = 0$. Asadar

$$MTTF = \int_0^{\infty} R(t) dt$$

Dacă rata defectărilor este constantă, λ , atunci

$$R(t) = e^{-\lambda t} \text{ si } MTTF = \int_0^{\infty} e^{-\lambda t} dt = 1/\lambda$$

Distributia Weibull

În multe calcule de fiabilitate, rata defectărilor λ se presupune constantă sau, echivalent, se admite o repartitie exponentială pentru durata de viață T . Sunt însă cazuri în care această presupunere este inadecvată. Ca exemple se pot lua zonele infantilă si de senectute din curba “cadă-de-baie”. În asemenea cazuri legea de distributie *Weibull* pentru T poate fi mai potrivită.

Distributia Weibull are doi parametri, λ si β . Functia densitate de probabilitate pentru durata de viață a unei componente este

$$f(t) = \lambda \beta t^{\beta-1} e^{-\lambda t^{\beta}}$$

Rata defectărilor în cazul distribuției Weibull este

$$h(t) = \lambda \beta t^{\beta-1}$$

Rata defectărilor $h(t)$ este descrescătoare, constantă sau crescătoare după cum constanta $\beta < 1$, $\beta = 1$ sau $\beta > 1$, ceea ce o face adecvată și pentru copilăria componentei, și pentru maturitatea ei, dar și, respectiv, pentru vârste avansate.

MTTF pentru distribuția Weibull

Funcția de fiabilitate pentru distribuția Weibull este

$$R(t) = e^{-\lambda t^\beta}$$

Timul mediu până la cădere pentru aceeași distribuție este

$$MTTF = \Gamma(1/\beta) / (\beta \lambda^{1/\beta})$$

cu $\Gamma(x)$ funcția Gamma, funcția euleriană de prima specie.

Cazul particular $\beta = 1$ face din distribuția Weibull o distribuție exponențială, cu rata defectărilor constantă.

TOLERANTA LA DEFECTE HARDWARE

Structuri canonice

Din componente se pot construi diverse structuri mai mult sau mai puțin complexe. Structurile complexe pot fi construite din câteva structuri de bază. Se admite tacit sau explicit independența statistică a căderilor pe care le pot manifesta componentele individuale. Structurile de bază sunt trei: sistemele *serie*, sistemele *paralel* și sistemele *M-din-N*.

Sistemele serie

Un sistem serie constă într-un set de componente conectate astfel încât căderea unei componente provoacă căderea întregului sistem (v.figura¹).



Fiabilitatea unui sistem serie, $R_s(t)$ este produsul fiabilităților celor N module. Într-adevăr, sistemul este funcțional dacă și numai dacă toate modulele componente sunt funcționale. Dacă $R_i(t)$ este fiabilitatea componentei i a sistemului atunci, folosind formula pentru intersecția de evenimente mutual independente rezultă

$$R_s(t) = \prod_{i=1}^N R_i(t)$$

În particular, dacă fiecare modul i are o rată de defectare constantă, λ_i , atunci

$$R_i(t) = e^{-\lambda_i t}$$

și

$$R_s(t) = e^{-\lambda_s t} = e^{-\sum \lambda_i t}$$

cu $\lambda_s = \sum \lambda_i$ rata de defectare constantă a sistemului serie în ansamblu.

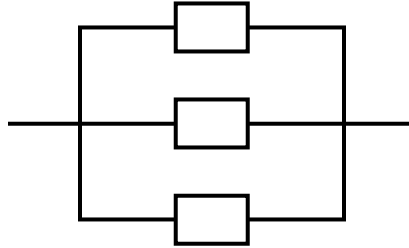
Conform relației generale de calcul, timpul mediu până la cădere pentru un sistem serie este

$$MTTF_s = 1/\lambda_s = 1/\sum \lambda_i$$

¹ Este convenabilă reprezentarea sistemelor prin așa-numitele *diagrame de semnal*. Diagramele de semnal au o singură intrare, plasată uzual în stânga și o singură ieșire, plasată uzual în dreapta diagramei. Dreptunghiurile reprezintă module funcționale. Starea sistemului este “funcțional” dacă există un drum de la intrare la ieșire. Absența oricărui drum intrare-ieșire corespunde stării “disfuncț”.

Sisteme paralele

Un sistem paralel este alcătuit dintr-un set de module conectate astfel încât pentru ca sistemul să cadă trebuie să cadă toate componentele (v.figura)



Fiabilitatea unui sistem paralel, $R_p(t)$ rezultă din relația

$$1 - R_p(t) = \prod_{i=1}^N (1 - R_i(t))$$

care exprimă probabilitatea eșecului concomitent al tuturor componentelor și deci al sistemului. Din nou se utilizează relația de calcul al probabilității unei intersecții de evenimente mutual independente, fiecare cu probabilitatea particulară egală cu $(1 - R_i(t))$. Din relația de mai sus rezultă fiabilitatea pentru structura paralelă

$$R_p(t) = 1 - \prod_{i=1}^N (1 - R_i(t))$$

Din nou, în particular, dacă fiecare modul i are o rată de defectare constantă λ_i , atunci

$$R_i(t) = e^{-\lambda_i t}$$

și

$$R_p(t) = 1 - \prod_{i=1}^N (1 - e^{-\lambda_i t})$$

Ca exemplu, pentru un sistem paralel cu două module

$$R_p(t) = e^{-\lambda_1 t} + e^{-\lambda_2 t} - e^{-(\lambda_1 + \lambda_2)t}$$

și timpul mediu până la cădere pentru un astfel de sistem este

$$MTTF_{P(2)} = \frac{1}{\lambda_1} + \frac{1}{\lambda_2} - \frac{1}{\lambda_1 + \lambda_2}$$

Timpul mediu până la cădere pentru un sistem paralel alcătuit din N module identice cu rata de defectare λ oarecare este

$$MTTF_{P(N)} = \sum_{i=1}^N \frac{1}{i\lambda}$$

relație verificată și de formula anterioară, dacă $\lambda_1 = \lambda_2 = \lambda$.

Sisteme M -din- N

Un sistem M -din- N este un sistem alcătuit din N module componente, identice în ceea ce privește funcția de îndeplinit. Un astfel de sistem rămâne funcțional atât timp cât cel puțin M componente sunt funcționale. Identitatea funcțională a celor N module se extinde de obicei și asupra fiabilității lor.

Cel mai bun exemplu și totodată un sistem foarte utilizat este *triplexul* (*TMR* – Triple Modular Redundant). Este vorba aici de trei componente identice care primesc aceeași intrare, suplimentat cu un dispozitiv de votare, *voter*, asupra ieșirilor lor. Acesta este un sistem ($M = 2$)-din- $(N = 3)$: atât timp cât cele mai multe (majoritatea) dintre componente produc rezultate corecte, sistemul este funcțional.

Fiabilitatea unui sistem M -din- N

Cele N componente sunt identice. Dacă $R(t)$ este fiabilitatea individuală a unui subsistem atunci fiabilitatea sistemului este probabilitatea ca $N - M$ sau mai puține subsisteme componente să fi căzut până la momentul t

$$R_{M\text{-din-}N}(t) = \sum_{i=0}^{N-M} C_N^i [1 - R(t)]^i [R(t)]^{N-i}$$

cu

$$C_N^i = \frac{N!}{i!(N-i)!}$$

Căderi corelate în sistemele M -din- N

Independența statistică a căderilor de componente este cheia fiabilității ridicate a sistemelor M -din- N . Căderile corelate pot diminua considerabil fiabilitatea sistemului. Dacă P_{cor} este probabilitatea ca sistemul să cadă în totalitate atunci

$$R_{M\text{-din-}N_{\text{cor}}}(t) = (1 - P_{\text{cor}}) \sum_{i=0}^{N-M} C_N^i [1 - R(t)]^i [R(t)]^{N-i}$$

Moduri de corelare în sistemele M -din- N

Dacă sistemul nu este proiectat cu grijă, factorul de cădere corelată poate domina probabilitatea generală de cădere. Există moduri diferite de corelare între căderile componentelor, care nu înseamnă în mod necesar o cădere generală antrenată de căderea unei singure componente. Ratele de defectare corelate sunt extrem de dificil de estimat. Din acest motiv, de aici înainte căderile componentelor vor fi considerate statistic independente: un modul nu alterează starea de funcționare a celorlalte în momentul căderii lui sau ulterior.

TMR (Triple Modular Redundant) – cluster redundant cu trei module

Sistemul TMR este probabil cel mai important sistem de tipul M -din- N . Este vorba, așa cum s-a mai spus, de cazul valorilor $M = 2$, $N = 3$ și sistemul este funcțional dacă cel puțin două din componente sunt funcționale. Un voter reține și eventual transmite totdeauna ieșirea majoritară. Nu trebuie omis faptul că voterul poate cădea el însuși. Fiabilitatea lui este notată $R_{\text{vot}}(t)$ și este strict subunitară. Expresia fiabilității sistemului TMR este

$$R_{\text{TMR}}(t) = R_{\text{vot}}(t) \sum_{i=0}^1 C_3^i [1 - R(t)]^i [R(t)]^{3-i} = R_{\text{vot}}(t) \{3[R(t)]^2 - 2[R(t)]^3\}$$

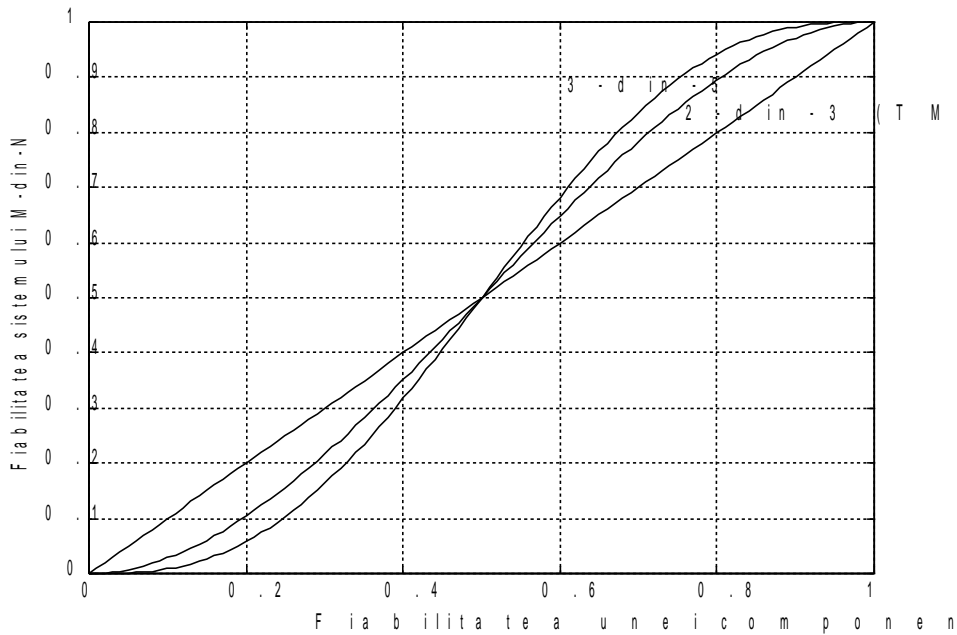
Pentru un sistem TMR cu module care au rate de defectare constante, λ , și în consecință fiabilități modulare $R(t) = e^{-\lambda t}$ identice, dacă defectarea voterului este exclusă, $R_{\text{vot}}(t) = 1$, se obține expresia particulară

$$R_{\text{TMR}}(t) = 3e^{-2\lambda t} - 2e^{-3\lambda t}$$

și timpul mediu până la cădere

$$MTTF_{\text{TMR}} = \int_0^{\infty} R_{\text{TMR}}(t) dt = 5/(6\lambda) < 1/\lambda = MTTF_{\text{simplex}}$$

Pare surprinzător: o structură mai sofisticată pare a fi mai puțin viabilă (nu fiabilă, viabilă!) decât o componentă a ei. Adevărul exprimat de aceste valori este dincolo de orice îndoială. Dacă se mai adaugă și posibilitatea reală de defectare a voterului, diferența dintre TMR și modulul unic (simplex) se mărește.



Pentru sublinierea avantajelor pe care le prezintă schemele redundante de tipul TMR sau NMR sunt necesare unele explicații în termeni de fiabilitate.

NMR (N-Modular Redundant) – cluster redundant din N module

Sistemul NMR este un cluster M -din- N cu N impar și cu $M = (N + 1)/2$. Se admite din nou că rata de defectare a voterului este neglijabilă, $R_{\text{vot}}(t) = 1$.

Graficul alăturat arată dependența fiabilității globale de fiabilitatea unui modul pentru cazurile $N = 3$ și $N = 5$. După cum se observă, pentru fiabilități ale unui singur modul $R < 0,5$, în ambele cazuri redundanța sistemului devine dezavantajoasă. Observația este generalizabilă și pentru numere N mai mari. Uzual însă $R \gg 0,5$ și atunci triplexul și sistemul 3-din-5, ca orice sistem M -din- N oferă câștiguri de fiabilitate semnificative.

TMR – compensarea defectelor

Ipoteza de lucru de până acum este “fiecare cădere a voterului produce o ieșire eronată și orice cădere a două module componente este fatală”. Sunt posibile însă și alte situații și iată un contraexemplu: un modul produce permanent un 1 logic și un al doilea modul generează permanent un 0 logic; sistemul TMR va funcționa adecvat relativ la acest bit. Este vorba aici de o compensare a defectelor. O situație similară poate apărea cu referire la anumite defecte în circuitul voterului.

Un alt exemplu de compensare a defectelor este acela al “erorilor care nu se suprapun” (non-overlapping faults): un modul poate avea un sumator cu defect și un alt modul poate avea un multiplicator cu defect. Dacă circuitele sunt disjuncte, este improbabil ca ele să genereze simultan ieșiri gresite.

Votere

Un voter primește intrările X_1, X_2, \dots, X_N de la un cluster M -din- N și generează o ieșire reprezentativă. Voterul cel mai simplu face comparația ieșirilor bit-cu-bit și ieșirea rezultă bit-cu-bit pe baza votului majoritar. Această schemă funcționează numai când toate procesoarele funcționale generează ieșiri care se potrivesc bit-cu-bit în totalitate. Dar pentru aceasta procesoarele trebuie să fie identice și să utilizeze același software. Altminteri, două ieșiri corecte pot diferi întrucâtva prin bitii cel mai puțin semnificativi.

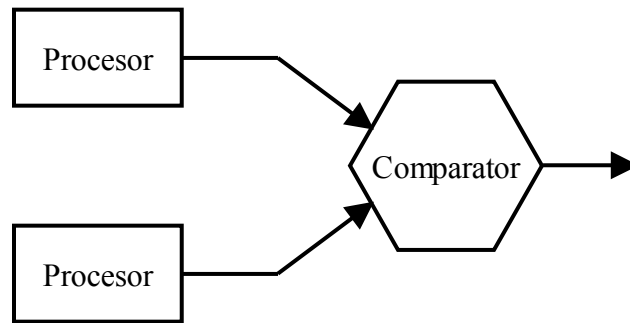
Votarea pluralității

Două ieșiri X și Y sunt declarate *practic* identice dacă $|x - y| < \delta$ pentru un δ precizat. Un voter cu k -pluralitate caută o mulțime de cel puțin k ieșiri *practic* identice și ia una din ele, oarecare (sau mediana lor), ca reprezentativă. De

exemplu dacă $\delta = 0,1$ si sunt cinci iesiri 1,10 1,11 1,32 1,49 si 3,00, submultimea {1,10 1,11} va fi retinută de un voter de pluralitate 2.

Sisteme duplex

Într-un sistem duplex (v.figura), ambele procesoare execută același task.



Dacă iesirile sunt concordante, rezultatul este considerat corect. Dacă iesirile sunt diferite, nu poate fi identificat procesorul care a gresit. Un software de nivel mai înalt trebuie să decidă cum trebuie tratată eroarea. Acest lucru se poate face utilizând una din metodele date în continuare.

Prima metodă constă în aplicarea unor așa-numite teste de acceptare. Test de acceptare înseamnă o verificare a domeniului de valori pentru fiecare din iesirile procesoarelor. De exemplu, presiunea într-un boiler trebuie să fie într-un anumit domeniu de valori cunoscut. Se folosesc asadar informații de natură semantică relativ la task-ul în derulare pentru a preciza care valori ale iesirii indică prezenta unei erori.

Cum se poate selecta gama de valori acceptabilă? Sunt două criterii concurente: *sensibilitatea* și *specificitatea*.

Domeniul de acceptare poate fi îngust ceea ce echivalează cu o probabilitate înaltă de a identifica o ieșire eronată, dar și o probabilitate ridicată de a aprecia gresit o ieșire bună ca fiind eronată (falsul pozitiv). Domeniul de acceptare poate fi larg și atunci probabilitățile sunt reduse pentru ambele situații arătate în fraza anterioară. În contextul acestor situații posibile (nuantat!) *sensibilitatea* este probabilitatea ca testul să recunoască o ieșire eronată ca o ieșire eronată, iar *specificitatea* este probabilitatea ca testul să recunoască o ieșire corectă ca o ieșire corectă. Cu aceste definiții, un domeniu de acceptare îngust se asociază cu o sensibilitate înaltă dar cu specificitate scăzută. Un domeniu de acceptare larg înseamnă o sensibilitate scăzută dar o specificitate înaltă.

A doua metodă apelează la o anumită testare la care sunt supuse ambele procesoare. Procesorul care nu trece testul este identificat ca fiind cel cu defect. Testele practicate în realitate nu sunt niciodată perfecte. De aceea se vorbește despre *acoperirea* testului care este același lucru cu sensibilitatea testului, adică cu probabilitatea ca testul să identifice un procesor defect ca defect.

Transparenta testului, o altă caracteristică a testului aplicat celor două procesoare este complementara acoperirii testului și exprimă probabilitatea ca testul să treacă un procesor defect în categoria “corect”.

Metoda a treia este așa-numita recuperare târzie (forward recovery). Această metodă utilizează un al treilea procesor pentru repetarea calculului executat de duplex. Dacă numai unul din cele trei procesoare este defect atunci cel care este în discordie cu procesorul al treilea este apreciat ca defect.

Practic este posibil să se folosească o combinație a acestor trei metode. Testul de acceptare este cel mai rapid de executat dar adesea este și cel mai puțin sensibil.

Fiabilitatea unui duplex.

Cele două procesoare active sunt presupuse a fi identice, fiecare cu fiabilitatea $R(t)$. Durata de viață a duplexului este timpul scurs până când ambele procesoare cad. Se utilizează și un factor de acoperire c , care reprezintă probabilitatea ca un procesor defect să fie diagnosticat corect, identificat și deconectat. Dacă $R_{duplex}(t)$ este fiabilitatea sistemului duplex, atunci

$$R_{duplex}(t) = R_{comp}(t) [R^2(t) + 2cR(t)(1 - R(t))]$$

S-a ținut seamă și de $R_{comp}(t)$, care este fiabilitatea comparatorului.

Pentru un duplex cu ratele de defectare ale procesoarelor componente constante (λ) și cu un comparator ideal, adică cu $R_{comp}(t) = 1$, fiabilitatea duplexului are expresia

$$R_{duplex}(t) = e^{-2\lambda t} + 2ce^{-\lambda t} (1 - e^{-\lambda t})$$

și durata lui medie de viață este

$$MTTF_{duplex} = 1/(2\lambda) + c/\lambda$$

Duplex cu redundante

Prin duplex cu redundante se înțelege un duplex cu două procesoare identice active și un număr (nelimitat) de procesoare în rezervă, înlocuitoare potențiale ale celui care se defectează. Când un procesor cade, căderea lui este detectată cu probabilitatea P_d și un procesor nou îl înlocuiește pe cel defect. Probabilitatea ca acest proces² de substituție să producă căderea întregului duplex este $1 - P_s$. Procesul indentificator al erorii este considerat instantaneu, rezervele sunt presupuse a fi totdeauna funcționale.

Imediat mai jos se prezintă un model pentru un duplex cu redundanțe. Se admite că fiecare procesor are o rată de defectare constantă λ . Durata de viață a unui procesor are, asadar, o distribuție exponențială cu parametrul λ . Timpul

² De cele mai multe ori cuvintele *procesor* și *proces* exprimă același lucru în sensul că un procesor găzduiește un proces. Alături, termenul proces semnifică altceva decât un proces pe un procesor. Cititorul va putea să distingă singur între cele două situații. Explicațiile vor fi detaliate numai acolo unde se consideră că este cazul.

între două defectări succesive ale unui procesor oarecare este și acesta distribuit exponențial cu același parametru λ . Se notează:

$M(t)$ – numărul de defectări ale unuia dintre procesoarele din duplex în intervalul $[0, t]$;

$N(t)$ – numărul de defectări ale sistemului duplex în intervalul $[0, t]$.

Relativ la distribuția lui $M(t)$ se parcurge raționamentul care urmează.

Fie Δt un interval de timp mic, atât de mic încât probabilitatea ca în scurtul răstimp Δt să se producă mai mult de o cădere să fie neglijabilă. $M(t + \Delta t) = n$ fie dacă $M(t) = n - 1$ și în intervalul Δt are loc o cădere, fie dacă $M(t) = n$ și în intervalul Δt nu se produce nici o cădere. Se scrie

$$\Pr[M(t + \Delta t) = n] \approx \Pr[M(t) = n - 1]\lambda\Delta t + \Pr[M(t) = n](1 - \lambda\Delta t)$$

având în vedere o firească proporționalitate (aproape riguroasă) a probabilității de manifestare a unui defect în acel scurt interval, cu durata Δt și cu parametrul λ . În relație se identifică cu ușurință și probabilitatea evenimentului contrar, lipsa producerii unui defect în acel interval Δt , care este $1 - \lambda\Delta t$. Această relație produce prin trecere la limită, $\Delta t \rightarrow 0$, ecuația diferențială

$$d\Pr [M(t) = n]/dt = -\lambda\Pr[M(t) = n] + \lambda\Pr[M(t) = n - 1]$$

care se poate scrie pentru fiecare n , întrucâtva diferit pentru $n = 0$.

Sistemul de ecuații diferențiale rezultat completat cu condițiile inițiale $\Pr[M(0) = n] = 0$ pentru $n \geq 1$ și $\Pr[M(0) = 0] = 1$ poate fi rezolvat, iar soluția lui este

$$\Pr[M(t) = n] = e^{-\lambda t} (\lambda t)^n / n! \text{ pentru } n = 0, 1, 2, \dots$$

Prin urmare, $M(t)$ are o distribuție Poisson de parametru λt . Se admite promptitudinea, chiar instantaneitatea în ceea ce privește restabilirea funcționalității sistemului, astfel încât după o defectare și apariția unei erori el funcționează din nou un interval de timp fără erori, fără manifestarea vreunui defect.

Ca o aplicație a celor stabilite mai sus, urmează calculul fiabilității pentru duplexul cu redundanță.

Duplexul se compune din două procesoare, ceea ce face rata de defectare dublă, 2λ . Rata de defectare a comparatorului se consideră neglijabilă. Probabilitatea a n căderi ale duplexului în intervalul $[0, t]$ este

$$\Pr[N(t) = n] = e^{-2\lambda t} (2\lambda t)^n / n! \text{ pentru } n = 0, 1, 2, \dots$$

Pentru ca duplexul să nu cadă în ansamblu său, fiecare din căderile particulare trebuie detectate și înlocuite cu succes. Aceasta se produce cu probabilitatea $c = P_d P_s$, produsul a două probabilități definite mai devreme, cea a detectării corecte a modulului defect și cea de funcționare a duplexului, probabilități a două evenimente independente. Pentru n căderi (independente) probabilitatea devine c^n .

Asadar, fiabilitatea unui duplex în condițiile menționate este

$$\begin{aligned} R_{\text{duplex}}(t) &= \sum_{n=0}^{\infty} \Pr(n \text{ căderi})c^n = \sum_{n=0}^{\infty} \exp(-2\lambda t)(2\lambda t)^n c^n / n! = \\ &= \exp(-2\lambda t) \sum_{n=0}^{\infty} (2\lambda t c)^n / n! = \exp(-2\lambda t) \exp(2\lambda t c) \end{aligned}$$

$$R_{duplex}(t) = \exp[-2\lambda(1 - c)t]$$

Iată acum o cale alternativă de stabilire a formulei pentru fiabilitatea unui duplex. Se admite că procesoarele cad la rata λ și atunci rata căderii duplexului este 2λ . Probabilitatea ca fiecare cădere să fie tratată cu succes este c și $(1 - c)$ este probabilitatea ca duplexul să cadă. Căderile care fac duplexul disfuncționar apar la o rată $2\lambda(1 - c)$. Funcția de fiabilitate a sistemului este aceeași ca mai sus, adică $\exp[-2\lambda(1 - c)t]$.

Fiabilitatea unor sisteme mai complexe

Printre sistemele mai complexe se află sistemele NMR în care procesoarele care cad sunt identificate și înlocuite dintr-un stoc indefinit de rezerve identice. Dacă rezervele sunt finite ca număr, atunci sumele din expresia fiabilității sunt limitate la acel număr de rezerve și nu infinite.

Se menționează și alte variante ale sistemelor duplex: un procesor este activ în timp ce al doilea este în *standby*. Se admite că procesoarele pot fi reparate atunci când devin nefuncționale.

În cazul sistemelor mai complexe, argumentele de natură combinatorială pot fi insuficiente pentru calculul fiabilității. Dacă ratele de defectare sunt constante, pentru evaluarea fiabilității se pot utiliza modelele Markov.

Lanturi Markov (scurtă introducere)

Modelele Markov furnizează o tratare structurată pentru deducerea funcțiilor de fiabilitate ale sistemelor complexe. Un lant Markov este un proces stohastic $X(t)$, o secvență de variabile aleatoare indexate de timpul t , cu o structură probabilistică specială. Pentru ca un proces stohastic să fie un lant Markov, comportarea lui viitoare trebuie să depindă numai de starea lui prezentă și nu de alte stări trecute. $X(t + s)$ depinde de $X(t)$, dar dându-se $X(t)$ și $X(t + s)$ acestea nu depind nemișcat de orice $X(\tau)$ pentru $\tau < t$. Dacă $X(t) = i$ se spune că la timpul t lantul este în starea i .

În continuare se tratează exclusiv lanturi Markov cu timp continuu ($0 \leq t \leq \infty$) și cu stări discrete, $X(t) = 0, 1, 2, \dots$

Lanturile Markov au o interpretare probabilistică consistentă. Relația

$$\Pr[X(t + s) = j | X(t) = i, X(\tau) = k] = \Pr[X(t + s) = j | X(t) = i] \quad (\tau < t)$$

exprimă prin probabilități condiționate dependența stării următoare exclusiv de starea precedentă. Sistemul, odată ajuns în starea i , rămâne în acea stare un timp aleator care este distribuit exponențial cu parametrul λ_i . Se poate spune că sistemul modelat de un lant Markov are o rată constantă λ_i de a părăsi starea i .

Probabilitatea ca lantul/sistemul să treacă din starea i în starea j se notează cu P_{ij} . Rata trecerii din starea i în starea j este atunci $\lambda_{ij} = P_{ij} \lambda_i$ și deoarece

$$\sum_{j \neq i} P_{ij} = 1 \quad \text{rezultă că} \quad \sum_{j \neq i} \lambda_{ij} = \lambda_i.$$

Probabilități asociate stărilor

Se notează cu $P_i(t)$ probabilitatea ca procesul/sistemul să fie la momentul t în starea i , starea lui inițială (la $t = 0$) fiind i_0 . Se scriu ecuații diferențiale pentru $P_i(t)$, ($i = 0, 1, 2, \dots$). Dacă la momentul t sistemul se află în starea j , după un interval scurt Δt sistemul poate ajunge în starea i în unul din următoarele cazuri:

- El a fost la momentul t în starea i , asadar $j \equiv i$, și nu și-a schimbat starea după Δt . Acest fapt are asociată probabilitatea $P_i(t)(1 - \lambda_i \Delta t)$.
- El a fost la momentul t în starea j , (oarecare, $j \neq i$) și a trecut în intervalul Δt din acea stare în starea i , eveniment care se petrece cu probabilitatea $P_j(t)\lambda_{ji}\Delta t$.

În acest raționament s-a admis că intervalul Δt este suficient de mic pentru ca producerea a mai mult de o tranziție pe durata lui să fie practic exclusă. Însușind acum termenii asociați tuturor tranzițiilor posibile, pentru diferiți j , se obține relația

$$P_i(t + \Delta t) \approx P_i(t)(1 - \lambda_i \Delta t) + \sum_{j \neq i} P_j(t)\lambda_{ji}\Delta t$$

După puțină prelucrare algebrică, trecerea la limită $\Delta t \rightarrow 0$ produce următoarele ecuații diferențiale în probabilitățile $P_i(t)$:

$$\frac{dP_i(t)}{dt} = -\lambda_i P_i(t) + \sum_{j \neq i} \lambda_{ji} P_j(t)$$

cu i (și j) parcurgând mulțimea de stări posibile ale lanțului. Și deoarece $\sum_{j \neq i} \lambda_{ij} = \lambda_i$ se poate rescrie

$$\frac{dP_i(t)}{dt} = -\sum_{j \neq i} \lambda_{ij} P_i(t) + \sum_{j \neq i} \lambda_{ji} P_j(t)$$

Acum aceste ecuații pot fi rezolvate pentru $i = 0, 1, 2, \dots$, cu condițiile inițiale $P_{i_0}(0) = 1$ și $P_i(0) = 0$ pentru $i \neq i_0$.

Un exemplu: Duplex cu unul din module în *standby*: un procesor este activ și unul este în *standby* (în rezervă caldă) și este conectat când cel activ cade. Rata căderii pentru procesorul activ este λ , constantă. Se ia în considerare și un factor de acoperire c , care este probabilitatea ca o cădere a procesorului activ să fie detectată corect și conectarea rezervei să fie reușită. Lanțul Markov în reprezentare grafică este dat în figură.

Ecuațiile diferențiale pentru duplexul cu un modul în *standby* sunt

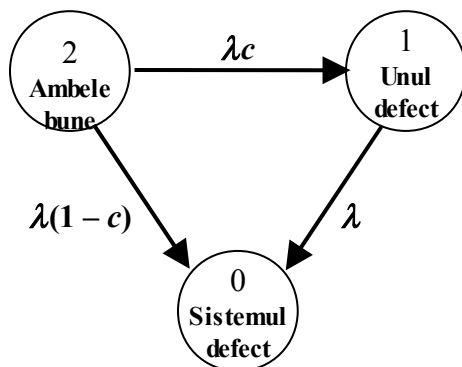
$$\frac{dP_i(t)}{dt} = -P_i(t) \sum_{j \neq i} \lambda_{ij} + \sum_{j \neq i} \lambda_{ji} P_j(t) \quad i = 0, 1, 2$$

cu detalierea

$$\begin{aligned} \frac{dP_0(t)}{dt} &= \lambda P_1(t) + \lambda(1-c)P_2(t) \\ \frac{dP_1(t)}{dt} &= -\lambda P_1(t) + \lambda c P_2(t) \end{aligned}$$

$$\frac{dP_2(t)}{dt} = -\lambda P_2(t)$$

si conditii initiale sunt: $P_2(0) = 1, P_1(0) = P_0(0) = 0$.



si conditii initiale sunt: $P_2(0) = 1, P_1(0) = P_0(0) = 0$.

Pentru a evalua fiabilitatea unui duplex care are un modul în *standby*, mai întâi solutia ecuatiilor diferentiale

$$P_0(t) = 1 - P_1(t) - P_2(t)$$

$$P_1(t) = c\lambda t \exp(-\lambda t)$$

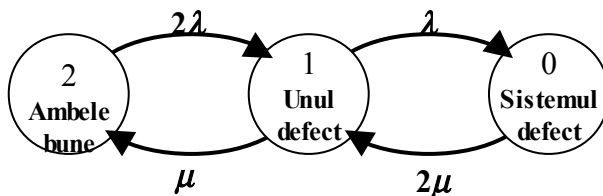
$$P_2(t) = \exp(-\lambda t)$$

apoi functia de fiabilitate

$$R_{sistem}(t) = 1 - P_0(t) = P_2(t) + P_1(t) = \exp(-\lambda t) + c\lambda t \exp(-\lambda t)$$

Propunem cititorului ca exercitiu obtinerea acestei expresii pe baza unor considerente combinatoriale.

Este adus acum în discutie un sistem duplex care beneficiază de reparatii. Este vorba de două procesoare active, fiecare cu rata de defectare λ si cu rata reparatiilor μ . Modelul Markov al unui astfel de sistem este dat în figura alăturată



Ecuatiile diferentiale în forma generală

$$\frac{dP_i(t)}{dt} = -P_i(t) \sum_{j \neq i} \lambda_{ij} + \sum_{j \neq i} \lambda_{ji} P_j(t)$$

se particularizează pentru cazul în discutie, pentru $i, j = 0, 1, 2$, sub forma

$$\frac{dP_2(t)}{dt} = -2\lambda P_2(t) + \mu P_1(t)$$

$$\frac{dP_1(t)}{dt} = 2\lambda P_2(t) + 2\mu P_0(t) - (\lambda + \mu)P_1(t)$$

$$\frac{dP_0(t)}{dt} = \lambda P_1(t) - 2\mu P_0(t)$$

Condițiile inițiale sunt $P_2(0) = 1, P_1(0) = P_0(0) = 0$

Soluția acestor ecuații diferențiale este

$$P_2(t) = \mu^2 / (\lambda + \mu)^2 + 2\lambda\mu / (\lambda + \mu)^2 e^{-(\lambda + \mu)t} + \lambda^2 / (\lambda + \mu)^2 e^{-2(\lambda + \mu)t}$$

$$P_1(t) = 2\lambda\mu / (\lambda + \mu)^2 + 2\lambda(\lambda - \mu) / (\lambda + \mu)^2 e^{-(\lambda + \mu)t} - 2\lambda^2 / (\lambda + \mu)^2 e^{-2(\lambda + \mu)t}$$

$$P_0(t) = 1 - P_2(t) - P_1(t)$$

adică probabilitățile celor trei stări posibile ca funcții de timp.

În sistemele duplex fără reparații, indicatorul cel mai important referitor la funcționarea (sau nefuncționarea) lor este fiabilitatea. Pentru sistemele duplex cu reparații, toți cei trei indicatori – fiabilitatea, disponibilitatea și disponibilitatea punctuală – sunt la fel de importanți. Evaluarea acestora se face în maniera descrisă imediat.

- Disponibilitatea punctuală – $A_p(t) = \text{Pr}(\text{sistemul la timpul } t \text{ este operational}) = 1 - P_0(t)$.
- Fiabilitatea – $R(t) = \text{Pr}(\text{sistemul este operational pe intervalul } [0, t])$ – prin eliminarea tranziției de la starea 0 la starea 1, rezolvând ecuațiile diferențiale noi care rezultă.
- Disponibilitatea – $A(t)$ – se ia media pe intervalul $[0, t]$ a proporției din timp în care sistemul este operational. Acesta este un indicator mai relevant pentru sistemele cu reparații.

Accesibilitatea staționară

Se calculează accesibilitatea staționară, $A(\infty)$ sau A , care este proporția din timp, pe durată îndelungată, cât sistemul este operational.

Se calculează mai întâi probabilitățile staționare, $P_2(\infty), P_1(\infty)$ și $P_0(\infty)$ sau P_2, P_1, P_0 . Probabilitățile staționare se pot calcula în două moduri: fie prin trecerea la limită, $t \rightarrow \infty$, în expresiile $P_i(t)$, $i = 0, 1, 2$, fie prin anularea tuturor derivatelor $dP_i(t)/dt = 0$, $i = 0, 1, 2$ și rezolvarea sistemului de ecuații algebrice liniare în P_i cu adaosul $P_2 + P_1 + P_0 = 1$. Rezultă

$$A = 1 - P_0 = P_1 + P_2$$

Starea staționară a unui duplex care admite reparații

Probabilitățile staționare sunt

$$P_2(t) = \mu^2 / (\lambda + \mu)^2$$

$$P_1(t) = 2\lambda\mu / (\lambda + \mu)^2$$

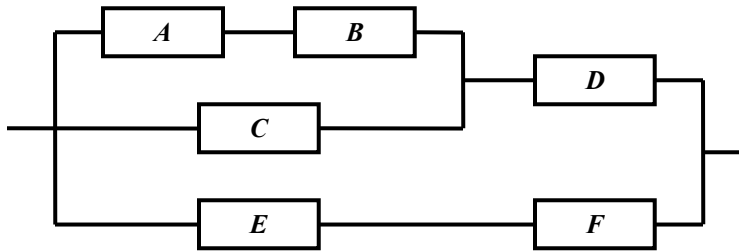
$$P_0(t) = \lambda^2 / (\lambda + \mu)^2$$

si accesibilitatea stationară este

$$A = A(\infty) = P_2 + P_1 = 1 - P_0 = (\mu^2 + 2\lambda\mu)/(\lambda + \mu)^2 = 1 - \lambda^2/(\lambda + \mu)^2$$

Structuri complexe

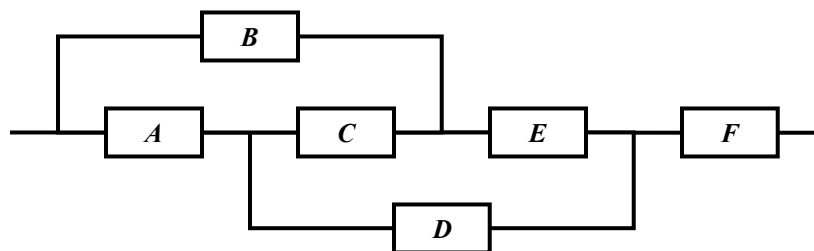
Multe din structurile mai complexe pot fi reduce în etape la structurile canonice serie sau paralel. Pentru ilustrarea acestei idei se dă sistemul din figura care urmează.



Modulele A, B sunt legate serie si reuniunea lor (A, B) este un subsistem pentru care se poate calcula fiabilitatea prin relatia consacrată. Modulul C si “modulul” (A, B) sunt asezate paralel si, din nou, există formula de calcul de la sistemele paralel. În continuare, modulul D si “modulul” (A, B, C) sunt serie, tot așa, separat, modulele E si F sunt înseriate si evaluările sunt posibile cu formulele cunoscute. În cele din urmă, “modulele” (A, B, C, D) si (E, F) sunt în paralel si din nou se poate evalua fiabilitatea ansamblului. Prin urmare, într-un număr finit de pasi se poate evalua prin calcule simple fiabilitatea sistemului alcătuit din modulele $\{A, B, C, D, E, F\}$ si analog pentru orice sistem decompozabil succesiv în subsisteme serie si/sau paralel.

Sisteme care nu sunt nici serie, nici paralel, nici succesiv reductibile la tipurile serie si/sau paralel.

Iată acum un sistem mai complicat. Se poate verifica usor că tentativa de tratare prin descompunere după reteta dată la exemplul introductiv este sortită eșecului.

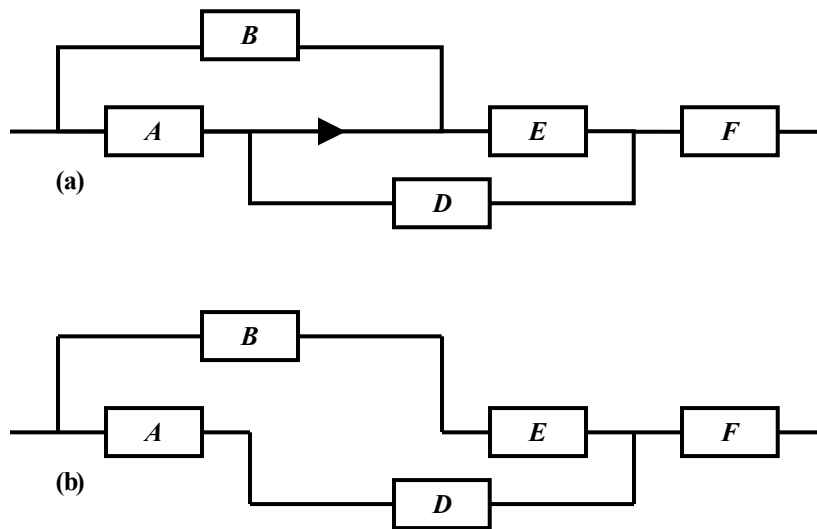


Si în acest caz, fiecare drum/cale de la stânga la dreapta reprezintă o configurație care permite sistemului să opereze corect, de pildă $A-D-F$.

Fiabilitatea sistemului poate fi calculată prin dezvoltarea în jurul unui singur modul i , $i \in \{A, B, C, D, E, F\}$ conform relației cu probabilități condiționate

$$R_{sistem} = R_i \Pr(\text{sistemul lucrează}/i \text{ este funcțional}) + (1 - R_i) \Pr(\text{sistemul lucrează}/i \text{ este nefuncțional})$$

Se figurează mai jos două grafuri derivate din cel de mai sus: în varianta (a) modulul i este funcțional, în varianta (b) modulul i este nefuncțional, în particular $i \equiv C$. Modulul C a fost selectat pentru ca cele două structuri reprezentate de cele două grafuri să fie mai apropiate de structurile simple decompozabile serie-paralel, paralel-serie etc.



Dezvoltarea în jurul unui modul poate fi repetată (etapizat) în jurul altui modul și apoi a altui modul s.a.m.d., până când grafurile (diagramele de semnal) rezultate sunt de tipurile mai simple serie-paralel, paralel-serie etc. Figura (a) de pildă necesită încă o dezvoltare în jurul modului E . Graful din figura (a) nu trebuie văzut ca fiind compus – între altele – dintr-o grupă (A, B) în paralel, conectată în serie cu grupa (D, E) la rândul ei în paralel. O asemenea tratare ar admite calea $B-C-D-F$ care nu este o cale validă decât dacă modulul C ar fi bidirecțional simetric. Acest detaliu este subliniat în figură prin săgeata atașată arcului care înlocuiește pe C când modulul este funcțional.

Dezvoltare evaluărilor în jurul modulelor C și E din exemplul dat urmează etapele de mai jos:

$$R_{sistem} = R_C \Pr(\text{sistemul lucrează}/C \text{ este operational}) + (1 - R_C)R_F[1 - (1 - R_A R_D)(1 - R_B R_E)]$$

Termenul al doilea s-a putut scrie deja după dezvoltarea din jurul modului C pentru că graful (b) este posibil a fi cuprins imediat în una din schemele simple

decompozabile serie-paralel, paralel-serie etc. Prin dezvoltare în jurul lui E rezultă și factorul necunoscut al primului termen:

$$\begin{aligned} \text{Pr(sistemul lucrează/C este operational)} &= \\ &= R_E R_F [1 - (1 - R_A)(1 - R_B)] + (1 - R_E) R_A R_D R_F \end{aligned}$$

Printr-o simplă înlocuire și puțină prelucrare algebrică se ajunge la expresia

$$\begin{aligned} R_{\text{sistem}} &= R_C [R_E R_F (R_A + R_B - R_A R_B) + (1 - R_E) R_A R_D R_F] + \\ &+ (1 - R_C) R_F (R_A R_D + R_B R_E - R_A R_D R_B R_E) \end{aligned}$$

Exemplu: $R_A = R_B = R_C = R_D = R_E = R_F = R$ conduce la

$$R_{\text{sistem}} = R^3(R^3 - 3R^2 + R + 2)$$

Dacă structura sistemului este mai complicată, calculele pot deveni la rândul lor foarte complicate. În asemenea situații se pot stabili limitele superioară și inferioară pentru fiabilitatea sistemului. O margine superioară este

$$R_{\text{sistem}} \leq 1 - \Pi(1 - R_{\text{calea}_i})$$

cu R_{calea_i} fiabilitatea subsistemului alcătuit din modulele serie de pe calea i . În produsul din formulă apar factori asociați tuturor căilor paralele. Pentru sistemul exemplificat mai sus: căile care fac sistemul funcțional sunt în număr de trei, $A-D-F$, $B-E-F$ și $A-C-E-F$ și atunci

$$R_{\text{sistem}} \leq 1 - (1 - R_A R_D R_F)(1 - R_B R_E R_F)(1 - R_A R_C R_E R_F)$$

În particular, dacă $R_A = R_B = R_C = R_D = R_E = R_F = R$ atunci

$$R_{\text{sistem}} \leq R^3(R^7 - 2R^4 - R^3 + R + 2)$$

Marginea inferioară a fiabilității se calculează pe baza așa-numitelor *multimi de tăietură minimă* a diagramei-graf a sistemului. O multime de tăietură minimă este o listă minimală de module constituită astfel încât eliminarea (datorată defectării) a tuturor modulelor din acea listă să facă sistemul disfuncțional. În cazul în discuție, multimi de tăietură minimală sunt $\{F\}$, $\{A, B\}$, $\{A, E\}$, $\{D, E\}$ și $\{B, C, D\}$.

Marginea inferioară a fiabilității sistemului din exemplul dat este partea dreaptă a inegalității

$$R_{\text{sistem}} \geq \Pi(1 - Q_{\text{tăietura}_i})$$

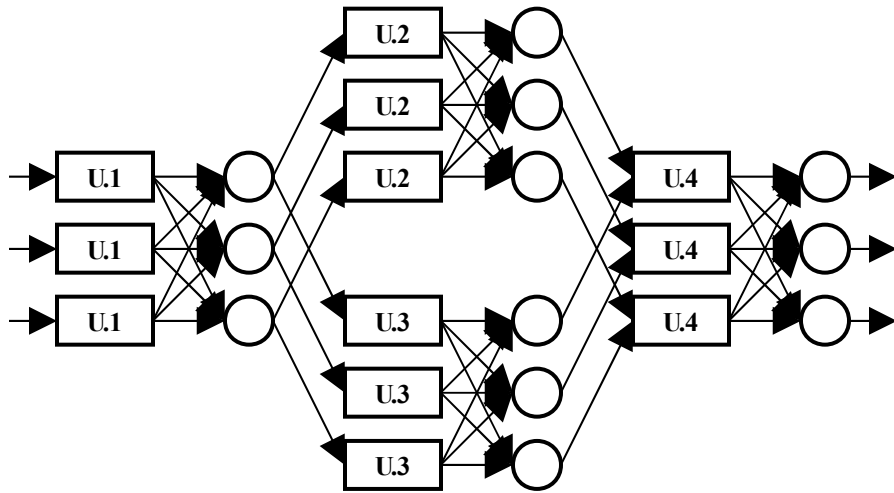
cu $Q_{\text{tăietura}_i}$ probabilitatea ca modulele din tăietura minimă i să fie toate defecte.

În particular, dacă $R_A = R_B = R_C = R_D = R_E = R_F = R$, atunci

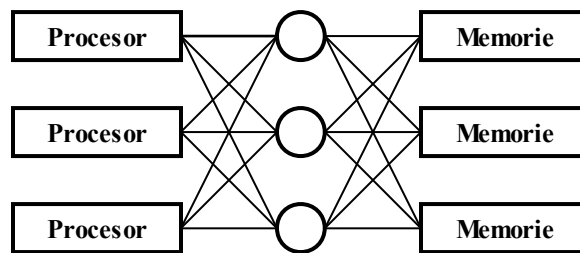
$$R_{\text{sistem}} \geq R^5(24 - 60R + 62R^2 - 33R^3 + 9R^4 - R^5)$$

Este readus acum în discuție sistemul NMR (N -Modular Redundant). Figura alăturată ilustrează cazul unor redundante modulare la nivel de unități.

Voterele (reprezentate prin cercuri) nu mai sunt elemente atât de critice ca în sistemul NMR discutat anterior; situația cu un singur voter defect nu este mai rea decât situația în care o singură unitate este defectă. Nivelul la care se aplică replicarea și votarea poate fi încă mai scăzut cu relativa risipă a unor votere aditionale care fac să crească dimensiunea sistemului și întârzierea lui.



Iată acum un sistem de procesoare si memorii triplicat.



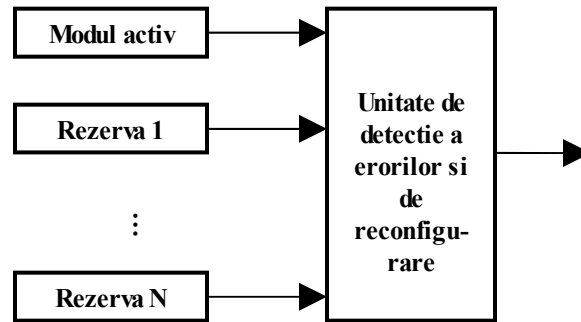
Toate comunicatiile (în ambele directii) între procesoarele triplicate si memoriile triplicate merg prin vot majoritar (voterele, de data aceasta bidirectionale sunt reprezentate si aici prin cercuri). Această organizare are desigur o fiabilitate mai înaltă decât un vot majoritar simplu în structura triplicată procesoare-memorii.

Redundante active/dinamice

Variatiile de mai sus ale sistemelor NMR (*N*-Modular Redundant) au neajunsul că angajează prea mult hardware, este adevărat, pentru a masca instantaneu erorile.

Rezultatele eronate pot fi acceptate temporar dacă sistemul poate detecta erorile si se poate reconfigura de unul singur.

Solutia este înlocuirea modului defect printr-o rezervă care funcționează corect. Iată în figura alăturată o schemă cu redundante active (sau dinamice).



Fiabilitatea sistemului cu rezerve active se evaluează observând că dacă toate modulele de rezervă sunt active (sub tensiune, alimentate) ele au aceeași rată de defectare și, întocmai ca la sistemul fundamental paralel, fiabilitatea sistemului capabil să detecteze disfuncțiile și să se reconfigureze este

$$R_{dinamic}(t) = R_{det}(t)[1 - (1 - R(t))^N]$$

În această expresie, s-a notat cu $R(t)$ fiabilitatea unui modul și cu $R_{det}(t)$ fiabilitatea unității de detectare și reconfigurare.

Fiabilitatea pentru cazul rezervelor în *standby*, adică pentru o foarte puțin probabilă cădere a modulelor de rezervă (de pildă pentru că nu sunt alimentate, nu sunt puse sub tensiune din motive de economie de energie) fiabilitatea unui sistem cu un modul activ și altul ca rezervă în *standby* este

$$R_{dinamic}(t) = R(t) + cR(t) [1 - R(t)]$$

cu c un așa-numit *factor de acoperire*, altfel spus o probabilitate ca modulul activ să fie corect diagnosticat ca defect și să fie deconectat, iar cel de rezervă, funcțional, să fie conectat cu succes.

Prin generalizare, în cazul cu N rezerve

$$R_{dinamic}(t) = R(t) \sum_{k=0}^N c^k [1 - R(t)]^k$$

Este de discutat aici dacă momentele căderii modulului activ au sau nu vreo importanță (exercitiu).

Redundanta hibridă

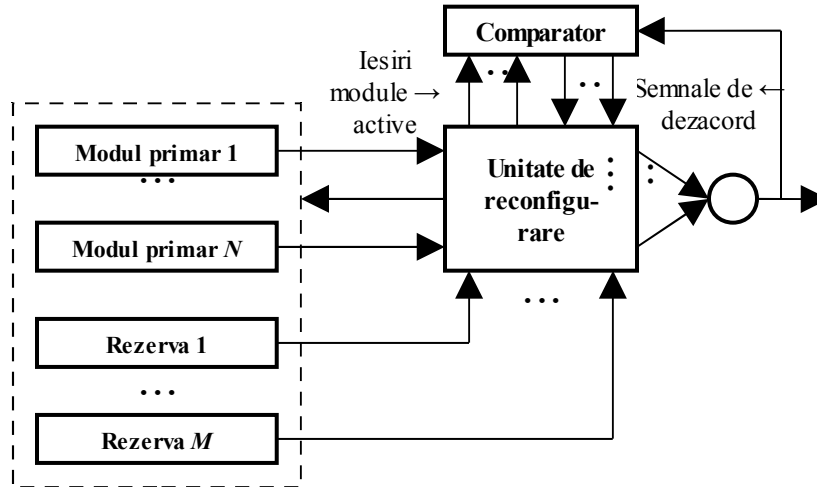
După cum s-a arătat mai devreme, un sistem de tipul NMR maschează defectele permanente și intermitente dar fiabilitatea lui pentru durate lungi de utilizare scade sub aceea a unui singur modul. Redundanta hibridă depășește acest inconvenient prin adăugarea de module de rezervă pentru înlocuirea modulelor active imediat ce ele devin disfuncționale.

Un sistem hibrid constă într-un miez de N procesoare (NMR) și M rezerve.

În cazul unei redundanțe hibride, fiabilitatea obținută printr-o grupare de tipul TMR susținută de M rezerve are expresia

$$R_{hibrid}(t) = R_{voter}(t)R_{reconf}(t) \{1 - mR(t)[1 - R(t)]^{n-1} - [1 - R(t)]^n\}$$

cu $n = M + 3$ numărul total de module, cu $R_{voter}(t)$ și $R_{reconf}(t)$ fiabilitățile voterului și circuitelor de comparare și reconfigurare. Se presupune că orice defect al voterului sau al subsistemului de comparare și reconfigurare produce căderea sistemului. În practică, nu toate defectele din aceste circuite sunt fatale: fiabilitatea va fi, de aceea, mai mare.



Dacă se urmărește o evaluare mai precisă pentru $R_{hibrid}(t)$ se cere o analiză detaliată a voterului și a circuitelor de comparare și de reconfigurare precum și a modurilor în care acestea pot cădea.

Se mai vorbește în literatură de redundanța modulară *sift-out* (cernută, cu purjare). Ca și în sistemele NMR toate cele N module sunt active. Un voter lucrează pe ieșirile tuturor modulelor (încă) operationale. În afara voterului, un circuit de comparare și de comutare compară ieșirea fiecărui modul cu ieșirile celorlalte module care sunt (încă) operationale. Un modul a cărei ieșire diferă de celelalte ieșiri este imediat deconectat. Este un sistem întrucâtva mai simplu decât cel cu redundanță hibridă. Un astfel de sistem poate să nu fie făcut prea agresiv în procesul de purjare (cernere – *sifting-out*) deoarece imensa majoritate a defectelor pot fi tranzitorii și pot dispărea într-un timp rezonabil. Eliminarea unui modul se face numai dacă el produce sustinut ieșiri eronate, incorecte pe o perioadă de timp potrivit aleasă.

Arhitecturi triplex-duplex

În această tratare a toleranței la defecte, procesoarele sunt legate împreună pentru a forma duplexuri. Apoi din aceste duplexuri ca “module” se face un triplex. Când procesoarele dintr-un duplex sunt în dezacord, ambele sunt scoase din sistem. Aranjarea triplex-duplex permite o identificare mai simplă a procesoarelor cu probleme. Mai departe, triplexul poate continua să funcționeze chiar dacă numai un duplex a rămas funcțional, deoarece aranjamentul duplex permite detectarea erorilor, a funcționării defectuoase.

TOLERANTA LA DEFECTE SOFTWARE

Cauze ale erorilor software

Proiectarea si scrierea software-ului sunt operatii nu lipsite de dificultate. Erorile software sunt datorate atât unor cauze care apartin de esenta realizării aplicatiilor cât si unor cauze accidentale.

Dificultățile de esență sunt generate de:

- Înțelegerea unei aplicatii complexe si a conditiilor de operare.
- Structura care poate cuprinde un număr de stări extrem de mare cu reguli de tranzitie de la o stare la alta foarte complexe.
- Modificările frecvente de software: sunt adăugate noi aspecte pentru a-l adapta la cerintele schimbătoare ale aplicatiilor.
- Platformele hardware si sistemul de operare care se pot schimba în timp si software-ul trebuie ajustat în consecință.
- Utilizarea ocazională a software-ului pentru a pune de acord măcar partial (paper over) incompatibilitățile între componentele în interacțiune ale unui sistem.

Dificultățile accidentale au ca sursă erorile umane.

Consideratiile privind costurile nu pot fi ocolite. De cele mai multe ori se recurge la software COTS (Commercial Off-The-Shelf) care nu este proiectat pentru aplicatii de mare fiabilitate.

Software-ul contine aproape inevitabil defecte. De aceea se face tot posibilul ca rata de defectare să scadă sau, pentru a rezolva problema defectelor de software, se folosesc tehnici de tolerare.

O metodă este verificarea formală a corectitudinii software-ului, dar ea nu este practică pentru module mari de software. Testele de acceptare sunt de asemenea utilizate în asa-numitele wrappers (ambalaje) si în blocurile de recuperare. Aceste teste reprezintă mecanisme importante în toleranta la defecte.

Exemplu: dacă un termometru indică -40°C în miezul unei zile de vară, termometrul respectiv este suspect de functionare proastă.

Verificările temporizate, un alt instrument utilizat în toleranta la defecte, constau în a pune un timer *watchdog* la timpul asteptat (expected) de executie. Dacă timerul este depășit, se presupune o posibilă disfuncție de hardware sau de software. Verificările temporizate se pot folosi în paralel cu alte genuri de teste de acceptare.

Testele de acceptare se pot realiza în moduri diferite.

O posibilitate constă în verificarea ieșirii: testul de acceptare este sugerat uneori în mod natural prin sortare, prin luarea rădăcinii pătrate, prin factorizarea unui număr mare sau prin soluția unor ecuații.

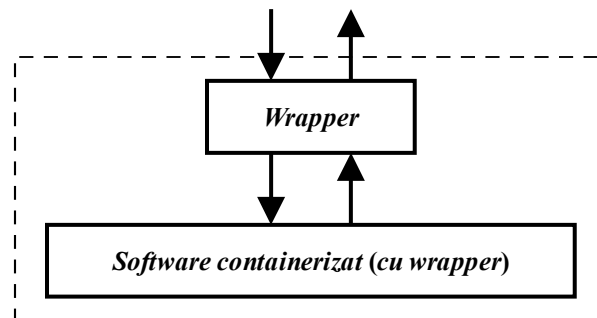
O altă cale o constituie verificările statistice: de exemplu, multiplicarea matricilor de întregi $n \times n$, $C = A \cdot B$. Tratarea naivă ar consuma un timp $O(n^3)$. Metoda lui Strassen generează la întâmplare un vector de întregi cu n elemente, R . Se calculează $M_1 = A \cdot (B \cdot R)$ și $M_2 = C \cdot R$. Dacă $M_1 \neq M_2$ atunci s-a produs o eroare; dacă $M_1 = M_2$ atunci corectitudinea este înalt probabilă. Operația se poate repeta prin generarea unui alt vector; complexitatea este $O(mn^2)$, cu m numărul de verificări.

De asemenea, o posibilitate utilizată frecvent este verificarea de domeniu (range). Se stabilesc limite acceptabile pentru ieșire. Dacă ieșirea este în afara limitelor: eroare. Limitele sunt fie presupuse, fie funcții simple de intrări (probabilitatea de testare greșită software trebuie să fie mică). Exemplu: un satelit de telemăsurare care preia imagini termice ale pământului; se fixează limite pentru domeniul de temperaturi. Mai departe, se utilizează corelațiile spațiale. Diferențele excesive între temperaturile unor zone adiacente sunt suspecte.

Este necesară echilibrarea între *sensibilitate* și *specificitate*. Sensibilitatea este probabilitatea condiționată ca testul să esueze când sunt date ieșiri eronate. Specificitatea este probabilitatea condiționată ca eroarea să fie adevărată, reală când testul de acceptare semnalează o eroare.

Limite ale domeniului mai înguste fac să crească sensibilitatea dar fac să crească și rata de alarme false ceea ce înseamnă o descreștere a specificității.

O soluție o reprezintă containerele (wrappers), care sunt interfețe menite să amelioreze robustețea modulelor de software. Aceste containere sunt utilizate și pentru kernelul sistemului de operare, și pentru ceea ce se numește *middleware*, și pentru software-ul de aplicație.



Intrările sunt interceptate de wrapper. Wrapper-ul le transmite mai departe sau semnalează o excepție. Ieșirile sunt filtrate de wrapper în mod similar.

Exemplu: utilizarea de software COTS³ pentru aplicatii de mare fiabilitate; componentele COTS sunt containerizate (wrapped) pentru a reduce rata de aparitie a disfuncțiilor software. Împachetarea aceasta previne intrări care (1) sunt în afara domeniului specificat sau (2) sunt cunoscute ca producătoare de căderi. Similar, iesirile sunt trecute printr-o testare de acceptare.

Exemplul 1: Rezolvarea depășirilor de buffer. Limbajul C nu execută verificări de domeniu pentru masive (arrays), ceea ce poate produce pagube accidentale sau malitioase. Scrierea unui secvente mari într-un buffer mic face ca bufferul să dea pe-afară, iar zone de memorie din afara bufferului să fie modificate prin suprascriere. Pagubele accidentale se materializează ca defecte de memorare. Pagubele malitioase se manifestă prin suprascrierea unor porțiuni din zonele *stack* și *heap* ale programului, o tehnică bine cunoscută de hackeri. Atacul cu distrugerea stivei (stack-smashing): un proces cu privilegii în rădăcină depune adresa lui de retur în stivă. Programul rău-voitor (malitios) suprascrie această adresă de revenire și controlul este redirectat la o locație de memorie unde hackerul a depus codul său atacator. Codul atacator are acum privilegii în rădăcină și poate distruge sistemul.

În cazul unui container (wrapper) pentru protecția la depășire de buffer, toate apelurile *malloc* ale unui program containerizat sunt interceptate de container. Containerul reține și urmărește poziția de început al memoriei alocate și dimensiunea. Scrierile sunt interceptate pentru a verifica potrivirea pe spațiul alocat, în limitele acestuia. Dacă apar nepotriviri, containerul nu permite scrierea și semnalează o eroare de depășire a spațiului alocat.

Exemplul 2: Verificarea corectitudinii gestionarului (scheduler) de taskuri. Containerul în jurul gestionarului (scheduler) de taskuri este un sistem tolerant la defecte în timp real. Astfel de gestionări nu utilizează programări *round-robin*⁴; poate utiliza EDF (Earlier Deadline First); execută taskul cel mai timpuriu între taskurile gata de rulare. Poate fi supus unor reguli de prioritate (taskuri dintr-o anumită parte a execuției nu pot fi întrerupte (*preemptable*)).

Un container poate verifica dacă algoritmul de planificare/gestionare este executat corect, adică dacă s-a reținut/selectat taskul cu cel mai apropiat termen de execuție și că orice task care sosește cu un termen mai timpuriu primează/prevalează asupra taskului în execuție (admitând că ultimul este *preemptable*).

Containerul cere informații despre taskurile care sunt gata de rulare și termenele lor de încheiere și totodată și dacă sunt curent în execuție și sunt *preemptable*.

³ **Commercial off-the-shelf (COTS)** este un termen pentru software și hardware (dar și pentru alte produse din tehnologia informației) care sunt gata făcute și sunt larg accesibile publicului prin vânzare, închiriere etc.

⁴ **Round-robin** este unul din cei mai simpli algoritmi de planificare/gestionare pentru un sistem de operare. Algoritmul atribuie ordonat fiecărui proces un interval de timp egal, fără vreo prioritate. Planificarea round-robin este simplă și ușor de implementat și evită situațiile în care un proces nu are acces niciodată la resursele necesare. Planificarea/gestionarea *round-robin* poate fi aplicată și în alte împrejurări cum ar fi gestionarea pachetelor în rețelele de calculatoare.

Exemplul 3: Utilizarea de software cu erori (bugs) cunoscute. Se cunoaste din testări sau din informatii provenite de la utilizatori că software-ul esuează pentru un anumit set de intrări, S . Containerul interceptează intrările acelui software si verifică dacă intrările sunt în S . Dacă nu sunt în S , sunt înaintate modulului de software pentru executarea calculelor. Dacă sunt în S , returnează sistemului o exceptie potrivită. Suplimentar, uneori redirectează intrarea spre un anumit cod alternativ, scris special pentru a manipula această intrare, acest gen de intrări.

Exemplul 4: Utilizarea unui container pentru a verifica iesirea corectă. Include un test de acceptare pentru filtrarea iesirii. Dacă iesirea trece testul, ea este transmisă în afară; dacă nu, este ridicată o exceptie si sistemul trebuie să trateze cum se cuvine iesirea suspectă.

Printre factorii care favorizează succesul în operatia de containerizare se numără:

- (a) Calitatea testelor de acceptare care este dependentă de aplicatie si are impact direct asupra abilității containerului de a opri iesirile eronate.
- (b) Disponibilitatea informatiilor necesare de la componenta containerizată: dacă acea componentă este gen “black box” (se observă numai răspunsul la o intrarea dată), containerul va fi întrucâtva mărginit; exemplu: un container de gestionare este imposibil de realizat fără informatii despre starea taskurilor care asteaptă să fie rulate.
- (c) Extinderea pe care o are testarea software-ului containerizat: testarea extinsă identifică intrări pentru care software-ul esuează.

Reîntinerirea software-ului

Exemplu: re-bootarea unui PC. Pe măsură ce procesul evoluează în executie el achizitionează memorie si face inchideri de fisiere (file-locks) fără a le elibera adecvat. De asemenea, spatiul de memorie tinde să devină din ce în ce mai fragmentat. Procesul poate deveni susceptibil la defecte si se poate opri din executie. Pentru a preveni aceste evenimente, se opreste procesul cu premeditare, se curăță starea sa internă de partea parazitara si apoi se reporneste.

Reîntinerirea poate fi din timp în timp (periodică) sau predictivă. Perioada de reîntinerire trebuie să echilibreze beneficiile si costurile.

Model analitic pentru reîntinerirea periodică

$\bar{N}(t)$ – numărul de disfunctii asteptat statistic (expected) în intervalul $[0, t]$.

K_f – costul căderii.

K_r – costul fiecărei reîntineriri.

P – perioada reîntineririi.

Prin adunarea costurilor de reîntinerire si costurilor disfunctiei, costul general asteptat pe o perioadă între două reîntineriri este

$$C_{re}(P) = \bar{N}(P) K_f + K_r$$

si rata costului reîntineririi este

$$C_{rata} = C_{re}(P)/P = [\bar{N}(P) K_f + K_r]/P$$

Exemple pe modelul analitic:

(1) – rata de defectare constantă în timp – $\bar{N}(P) = \lambda P$

Rata costului este $C_{rata} = \lambda K_f + K_r/P$; optimă este valoarea $P^* = \infty$ – nu este indicată reîntinerirea.

Reîntinerirea este profilaxia unei rate de defectare crescătoare cu vârsta – aici nu-i vorba de nici un fel de îmbătrânire.

(2) $\bar{N}(P) = \lambda P^2$

Rata costului este $C_{rata} = \lambda K_f P + K_r/P$; optimă este valoarea

$$P^* = \sqrt{K_r / (\lambda K_f)} .$$

(3) $\bar{N}(P) = \lambda P^n, n > 1$

Rata costului este $C_{rata} = \lambda K_f P^{n-1} + K_r/P$; valoarea P optimă este

$$P^* = [K_r / (n-1)\lambda K_f]^{1/n} .$$

Concluzii: P^* crește odată cu raportul K_r/K_f ; rata de creștere se micșorează pe măsură ce crește n ; creșterea lui λ are ca efect descreștere a lui P^* .

Perioada de reîntinerire optimă. Estimarea parametrilor K_r, K_f și $N(t)$ se poate face experimental executând simulări ale software-ului. Sistemul poate fi făcut adaptiv: valori initiale prin default și ajustarea valorilor pe măsură ce sunt recoltate date statistice noi.

Reîntinerirea predictivă se face prin monitorizarea caracteristicilor sistemului – volumul de memorie alocată, numărul de fișiere mentinute blocate (locked) etc. Acestea sunt premise pe baza cărora se poate prezice momentul când sistemul va cădea.

Exemplu: un proces consumă memorie la o rată anumită, iar sistemul estimează când el va termina memoria; reîntinerirea trebuie să aibă loc tocmai înaintea clacării prezise.

Software-ul care implementează reîntinerirea predictivă trebuie să aibă acces la suficientă informație asupra stării procesului pentru a face astfel de predicții.

Dacă software-ul predictiv este parte a sistemului de operare, astfel de informații se colectează ușor. Dacă acesta este un pachet care rulează peste sistemul de operare fără privilegii speciale, trebuie constrâns să utilizeze interfețele oferite de sistemul de operare.

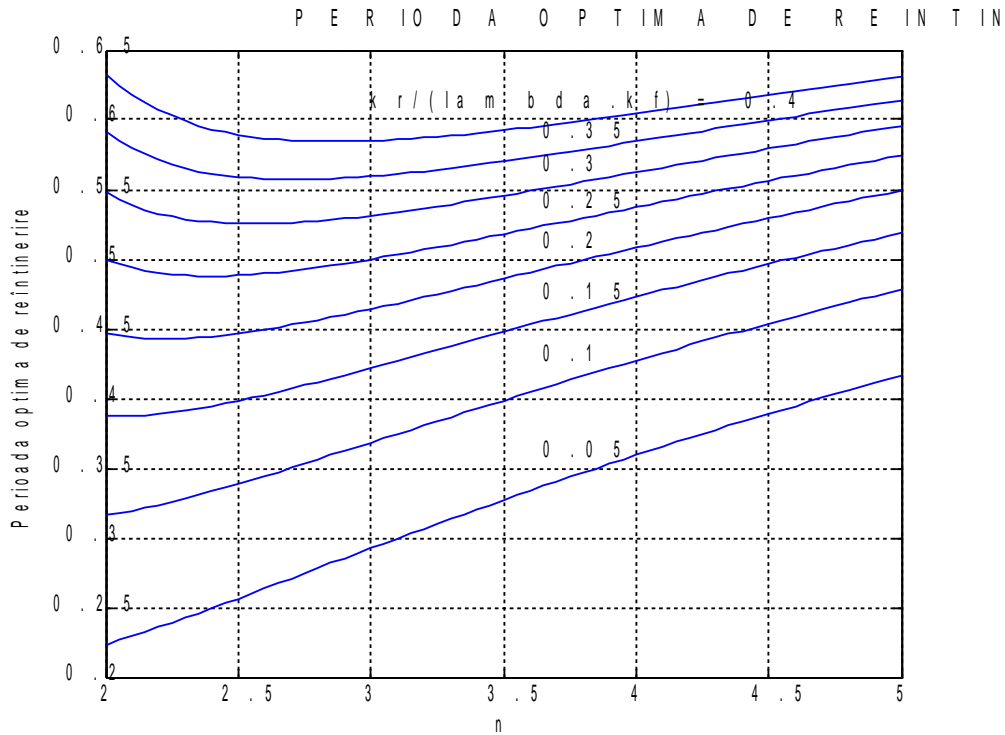
Exemplu: Sistemul de operare Unix oferă următoarele utilități de colectare a informațiilor de stare:

iostat – extrage procentajul de utilizare a CPU la nivelurile utilizator și sistem, precum și un raport asupra utilizării fiecărui dispozitiv I/O.

netstat – indică conexiunile în rețea, tabelele de rutare și un tabel al tuturor interfețelor de rețea.

nfsstat – furnizează informații despre statisticile de kernel ale serverului de fișiere din rețea.

vmstat – furnizează informații despre utilizarea procesorului, despre activitatea cu memoria și paginarea⁵, despre trape⁶ și despre acțiunile I/O.



Predicția momentului căderii prin metoda celor mai mici pătrate

Cu informația despre stare colectată, urmează a identifica tendințele și a prezice momentul căderii.

Exemplu: urmărind alocarea memoriei pentru un proces, sunt colectate valorile $\mu(t_0), \mu(t_1), \dots, \mu(t_p)$, cu $\mu(t_i)$ memoria alocată la momentul t_i .

Se poate face o regresie polinomială

⁵ Dacă un program face referire la o locație de memorie care se află într-o pagină de memorie virtuală inaccesibilă, calculatorul generează o eroare de pagină. Hardware-ul trece controlul la sistemul de operare într-un loc unde poate încărca din memoria auxiliară pagina cerută (de pildă un fișier paginat de pe disc) și pune un flag pe *on* pentru a spune că pagina este accesibilă. Hardware-ul va lua apoi locația de start a paginii, va adăuga offset-ul la bitii de ordin inferior din registrul de adrese și va accesa locația de memorie dorită.

⁶ În calcul și în sistemele de operare, o trapă este un gen de întrerupere sincronă provocată de condiții excepționale (de pildă împărțirea cu zero, accesul la memorie nepermis) într-un proces utilizator. Uzual, o trapă produce o comutare la modul kernel în care sistemul de operare execută o anumită acțiune înainte de a returna controlul procesului care l-a generat. În unele utilizări, termenul de trapă se referă specific la o întrerupere intenționată a iniția o comutare contextuală la un program monitor sau la un program depanator (debugger). Într-un SNMP (Simple Network Management Protocol), o trapă este un tip de PDU (Protocol Data Unit) utilizat pentru a raporta o alertă sau alt eveniment asincron relativ la un subsistem gestionat.

$$f(t) = m_n t^n + m_{n-1} t^{n-1} + \dots + m_1 t + c$$

cu coeficienti care să minimizeze

$$\sum_{i=0}^p [\mu(t_i) - f(t_i)]^2$$

Regresia cea mai simplă este desigur cea liniară: $f(t) = mt + c$.

Polinomul poate fi utilizat pentru extrapolări (prudente) în viitor și prognoze asupra momentului când procesul va ajunge în criză de memorie.

Cele mai mici pătrate ponderate

În regresia standard, fiecare punct $\mu(t_i)$ are pondere egală în stabilirea coeficientilor de regresie. O variantă mai bună este aceea cu ponderi diferite: se selectează w_0, w_1, \dots, w_p și apoi se calculează coeficienții din $f(t)$ pentru a minimiza

$$\sum_{i=0}^p w_i [\mu(t_i) - f(t_i)]^2$$

Ponderile permit diferentierea efectelor unor anumite puncte. De exemplu, ponderile $w_0 < w_1 < \dots < w_p$ fac ca datele mai recente să influențeze mai mult coeficienții de regresie decât datele mai vechi.

Regresia (potrivirea de curbe) este vulnerabilă la impactul câtorva, putine puncte care sunt neuzual de mari sau neuzual de mici. Prin acest efect relația poate rezulta distorsionată. Sunt accesibile tehnici de a face potrivirea mai robustă prin reducerea impactului unor astfel de puncte.

Tratamentul combinat

Reîntinerirea predictivă lucrează cu un timer resetat la momentul reîntineririi. Dacă timerul expiră, reîntinerirea se face indiferent de momentul prezis pentru producerea următoarei căderi.

Nivelul reîntineririi poate fi cel al aplicației sau cel al nodului. Depinde de locul unde s-au degradat resursele sau sunt pe cale de epuizare. Reîntinerirea la nivelul aplicației constă în suspendarea unei aplicații individuale, efectuarea unei curățiri a stării sale (prin colectarea de *garbage*, re-initializarea structurilor de date etc.) și apoi repornirea acelei aplicații. Reîntinerirea la nivel de nod se face prin re-bootarea nodului și afectează toate aplicațiile în execuție pe acel nod.

Programarea în N versiuni

N echipe de programatori independente dezvoltă software cu aceleași date initiale, cele N versiuni se execută în paralel, se execută un vot asupra ieseirilor. Dacă programele sunt dezvoltate independent este foarte puțin probabil ca ele să clacheze la aceleași intrări.

Ipoteze de lucru: căderile sunt independente statistic; probabilitatea de a claca a unei versiuni particulare este p .

Probabilitatea de a avea exact m căderi din N versiuni este

$$P(N, m, p) = C_N^m p^m (1-p)^{N-m}$$

Probabilitatea căderii sistemului (pentru N impar) este $\sum_{m=(N+1)/2}^N P(N, m, p)$

Problema *comparării consistente*. Programarea în N versiuni nu este simplu de implementat. Chiar dacă toate versiunile sunt corecte, atingerea unui consens este dificilă.

Exemplu: V_1, \dots, V_N – N versiuni scrise independent pentru calculul mărimii x și compararea ei cu c .

x_i – valoarea lui x calculată cu versiunea V_i .

B_i – variabilă booleană: $B_i \equiv (x > c)$.

Compararea cu c se spune că este *consistentă* dacă $B_i = B_j$ pentru orice pereche (i, j) , cu $i, j = 1, \dots, N$.

Cerinta de consistentă.

Exemplu: se calculează o funcție de presiune și temperatură $f(p, t)$; dacă $f(p, t) \leq c$ se întreprinde acțiunea A_1 , dacă $f(p, t) > c$ se întreprinde acțiunea A_2 . Fiecare versiune produce tipul de acțiune de întreprins; ideal, dacă toate versiunile sunt consistente, atunci ele indică aceeași acțiune.

Versiunile sunt scrise independent, se utilizează algoritmi diferiți de evaluare a funcției $f(p, t)$, valorile vor fi întrucâtva diferite.

$c = 1,0000$; $N = 3$.

Toate cele trei versiuni operează corect și valorile generate sunt 0,9999, 0,9998 și 1,0001. $B_1 = B_2 = \text{false}$ – acțiunea A_1 , $B_3 = \text{true}$ – acțiunea A_2 .

Apare o inconsistentă, cu toate că versiunile sunt toate corecte.

În problema consistenței se formulează următoarea

Teoremă: Cu excepția algoritmului trivial care aplică orice intrare pe același număr (aplicatia identică), nu există vreun algoritm care să garanteze că doi întregi de n biti care diferă cu mai puțin decât 2^k vor fi aplicați pe aceeași ieșire de m biti, cu $m + k \leq n$.

Demonstratie: Se face prin reducere la absurd. Se presupune că există un astfel de algoritm (care nu aplică orice intrare pe același număr). Se alege cazul $k = 1$. Numerele 0 și 1 diferă prin mai puțin de 2^k astfel că ambele vor fi aplicate pe același număr, să spunem L . Similar, 1 și 2 vor fi aplicate în L . La fel 3, 4, ... vor fi aplicate de acest algoritm în L . Acest fapt contrazice ipoteza inițială.

Un rezultat similar se menține și pentru numere reale care diferă puțin unul de altul (exercitiu).

Problema comparării consistente are nuanțele expuse imediat.

Dacă versiunile nu concordă, ele pot fi cu defecte sau fără. Versiunile-ecșec multiple pot produce ieșiri gresite identice datorită defectelor corelate. Votul sistemului va alege în astfel de cazuri o ieșire gresită. Problema se poate ocoli prin obligarea versiunilor să decidă asupra unei valori de consens a unei variabile. Înainte de a verifica dacă $x > c$, versiunile cad de acord asupra unei

valori x de utilizat adăugând cerința de a specifica ordinea comparațiilor pentru comparații multiple. Aceasta poate reduce diversitatea de versiuni, crescând potențialul pentru eșecuri corelate.

Această schemă poate să degradeze performanța: versiunile care termină evaluarea mai devreme trebuie să aștepte.

Semnale de încredere. Fiecare versiune calculează $|x - c|$; dacă rezultatul este sub un δ pozitiv precizat, versiunea anunță o joasă încredere în ieșirea sa. Voterul dă ponderi reduse versiunilor de încredere scăzută.

Problemă: dacă o versiune funcțională are $|x - c| < \delta$, sunt șanse mari ca același lucru să se întâmple și pentru alte versiuni, ale căror ieșiri vor fi devalorizate de voter. Frecvența apariției acestei probleme și durata ei depind de natura aplicației.

De exemplu, în aplicații unde calculul depinde numai de intrările cele mai recente și nu de valori trecute, problema consensului poate apărea mai rar și poate dispărea repede.

Dependenta versiunilor

Eșecurile corelate între versiuni pot provoca creșterea probabilității generale de cădere cu ordine de mărime.

Exemplu: $N = 3$, poate tolera până la o versiune greșită pentru orice intrare; fie $p = 0,0001$ – o ieșire incorectă la fiecare 10.000 de execuții.

Dacă versiunile sunt stochastic independente, probabilitatea de a eșua a sistemului de trei versiuni este $p^3 + 3p^2(1 - p) = 3 \cdot 10^{-8}$.

Se presupune acum o dependență stochastică și existența unui defect comun pentru două versiuni, defect care se manifestă o dată la un milion de execuții provocând căderea sistemului.

Probabilitatea eșecului pentru sistemul de trei versiuni crește la peste 10^{-6} , mai mult de 30 de ori față de situația când defectele nu sunt corelate.

Modelul cu versiuni corelate. Spațiul intrărilor este divizat în regiuni: probabilități diferite de preluare de intrări din regiunea care produce eșecul unei versiuni.

Exemplu: algoritmul poate avea instabilitate numerică într-o parte a spațiului intrărilor; rata disfuncțiilor în acea parte este mai mare decât cea medie.

Ipoteză de lucru: versiunile sunt independente stochastic în fiecare subspațiu dat S_i ; $\Pr[\text{atât } V_1 \text{ cât și } V_2 \text{ să esueze} | \text{intrare din } S_i] =$

$$= \Pr[V_1 \text{ esuează} | \text{intrare din } S_i] \cdot \Pr[V_2 \text{ esuează} | \text{intrare din } S_i]$$

Probabilitatea necondiționată ca o versiune să esueze

$$\Pr[V_1 \text{ esuează}] = \Pr[V_1 \text{ esuează} | \text{intrare din } S_i] \cdot \Pr[\text{intrare din } S_i]$$

Probabilitatea necondiționată ca ambele să esueze

$$\begin{aligned} & \Pr[\text{atât } V_1 \text{ cât și } V_2 \text{ să esueze}] = \\ & = \Pr[\text{atât } V_1 \text{ cât și } V_2 \text{ să esueze} | \text{intrare din } S_i] \cdot \Pr[\text{intrare din } S_i] \neq \\ & \neq \Pr[V_1 \text{ esuează}] \cdot \Pr[V_2 \text{ esuează}] \end{aligned}$$

Exemplul numeric 1:

Se identifică două subspații de intrare S_1 și S_2 , fiecare cu probabilitatea 0,5.
 Probabilitățile condiționate de eșec sunt:

Versiunea	S_1	S_2
V_1	0,01	0,001
V_2	0,02	0,003

Probabilitățile necondiționate de eșec sunt:

$$\Pr[V_1 \text{ esuează}] = 0,01 \cdot 0,5 + 0,001 \cdot 0,5 = 0,0055$$

$$\Pr[V_2 \text{ esuează}] = 0,02 \cdot 0,5 + 0,003 \cdot 0,5 = 0,0115$$

Dacă versiunile ar fi independente, probabilitatea ca ambele să esueze pentru aceeași intrare ar fi

$$0,0055 \cdot 0,0115 = 6,33 \cdot 10^{-5}$$

Probabilitatea reală de eșec conjugat este întrucâtva mai mare

$$0,01 \cdot 0,02 \cdot 0,5 + 0,001 \cdot 0,003 \cdot 0,5 = 1,02 \cdot 10^{-4}$$

Cele două versiuni sunt corelate pozitiv: ambele sunt mai predispușe la eșec în S_1 decât în S_2 .

Exemplul numeric 2.

Probabilitățile condiționate de eșec sunt:

Versiunea	S_1	S_2
V_1	0,01	0,001
V_2	0,003	0,02

Probabilitățile necondiționate de eșec sunt la fel ca în exemplul de mai sus.

Probabilitatea de eșec combinat este

$$0,01 \cdot 0,003 \cdot 0,5 + 0,001 \cdot 0,02 \cdot 0,5 = 2,5 \cdot 10^{-5}$$

Mai puțin decât probabilitatea combinată evaluată mai sus sau decât produsul probabilităților individuale.

Tendințele de cădere sunt corelate, în acest caz, negativ: V_1 este mai bun în S_1 decât în S_2 , invers pentru V_2 ; V_1 și V_2 își maschează deficiențele reciproc. Ideal ar fi ca versiunile multiple să fie corelate negativ. În practică corelațiile sunt pozitive deoarece versiunile rezolvă aceeași problemă.

Cauzele corelării versiunilor

Specificatii comune: comune sunt și erorile în specificatii (elementele de pornire) care se propagă la software.

Dificultatea intrinsecă a problemei: algoritmi pot fi mai dificil de implementat pentru unele intrări, producând declansarea defectelor pentru aceleși intrări.

Algoritmi comuni: algoritmul însuși poate conține instabilități în anumite regiuni ale spațiului intrărilor. De aceea, versiuni diferite pot avea instabilități în aceeași regiune.

Factori culturali: în interpretarea unor specificații ambigue, programatorii fac greșeli similare.

Platforme software și hardware comune: dacă se utilizează același hardware, același sistem de operare și același compilator, defectele lor pot declanșa un eșec corelat.

Din aceste motive și din altele, se recomandă obținerea independenței versiunilor prin diversitatea incidentală.

Se creează obligația ca aceia care dezvoltă module software diferite să lucreze independent unul de altul. Se impune interdicția ca echipele care lucrează la module diferite să comunice între ele. Întrebările privind ambiguitățile din specificații sau oricare altă problemă trebuie adresate unei autorități centrale care are obligația de a face corecturile și actualizările necesare pentru toate echipele. Inspectia software-ului trebuie să fie foarte atent coordonată astfel încât inspectorii unei versiuni să nu scape direct sau indirect informații despre o altă versiune.

Obținerea independenței versiunilor prin metode de diversitate forțată.

Specificatiile inițiale se fac deliberat diverse. Versiunile sunt și ele create deliberat cu capacități diferite. Se folosesc limbaje de programare diferite. Până și instrumentele de dezvoltare și compilatoarele se diversifică. La fel sistemele de operare și hardware se aleg diferite.

Specificatii inițiale diverse. Multe disfuncții de software sunt datorate cerințelor formulate prin specificație. Diversitatea poate începe chiar din faza de specificare: specificațiile pot fi exprimate într-un formalism diferit. Erorile de specificare nu vor coincide pe versiuni, fiecare specificare va declanșa un profil de defecte de implementare diferit.

Versiuni cu capacități diferite. Un exemplu poate fi o versiune rudimentară care furnizează ieșiri mai puțin precise dar încă acceptabile. O a doua versiune mai simplă, mai puțin predispusă la defecte și bazată pe un algoritm mai robust. Dacă cele două versiuni nu concordă, o a treia versiune poate ajuta la determinarea versiunii corecte din cele două de până atunci. Dacă a treia versiune este foarte simplă, pentru a dovedi corectitudinea pot fi folosite metodele formale.

Limbaje de programare diferite. Limbajele de programare afectează calitatea software-ului.

Exemple:

- Assemblerele sunt mai predispuse la erori decât limbajele de nivel mai înalt.
- Natura erorilor este diferită: în programele scrise în C se depășește ușor memoria alocată, este imposibil ca memoria să fie gestionată strict în vreun limbaj.
- Utilizarea fără defect a pointerilor în Fortran este imposibilă deoarece Fortranul nu are posibilitatea de a defini pointeri.
- Fată de C sau Fortran, Lisp-ul este un limbaj mai natural pentru o parte din algoritmii de inteligență artificială.

Limbajele diferite de programare au biblioteci diferite si compilatoare diferite; programele vor avea esecuri necorelate (sau, chiar mai bine, corelate negativ).

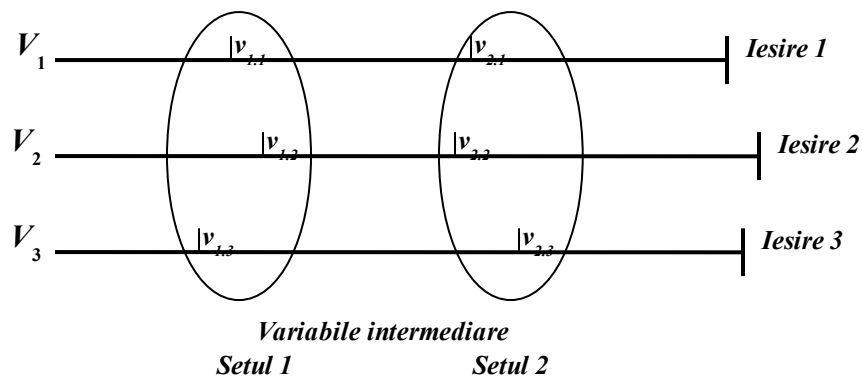
Alegerea limbajului de programare. Problemă ridicată de programatori: este necesar ca toate versiunile să folosească cel mai bun limbaj pentru rezolvarea unei probleme sau unele versiuni ar trebui scrise în alte limbaje mai puțin potrivite?

Dacă se utilizează același limbaj, rata de defectare individuală este mai redusă dar esecurile posibile sunt corelate pozitiv. Dacă limbajele utilizate sunt diferite, unele rate de defectare particulare pot fi mai mari dar rata de defectare generală a sistemului de N versiuni poate fi mai redusă dacă esecurile sunt mai slab corelate. Echilibrul (*tradeoff*) este dificil de rezolvat; nu există vreun model analitic, de aceea este necesară o experimentare extinsă.

Instrumente de dezvoltare si compilatoare diverse. Diversitatea de instrumente si de compilatoare face posibilă o “diversitate de notatii” care poate să reducă corelatiile pozitive între disfuncții. Utilizarea de instrumente de dezvoltare diverse si de compilatoare diverse (posibil purtătoare de defecte ele înseși) pentru versiunile diferite pot produce fiabilități mai mari.

Sisteme de operare si hardware diverse. Iesirile sistemului depind de interacțiunea între software-ul de aplicatie si platforma sa, sistemul de operare si procesorul. Atât procesorul cât si sistemul de operare au specificitatea lor prin defectele (*bugs*) pe care le contin. Este asadar o idee bună de a completa proiectarea software diversă cu diversitatea hardware si de sistem de operare. Urmează natural rularea fiecărei versiuni pe procesoare si sisteme de operare diferite.

Testarea spate-n-spate. Compararea valorilor unor variabile intermediare sau a iesirilor pentru aceeași intrare reprezintă o cale de a identifica erorile, erori care uzual nu coincid.



Variabilele intermediare furnizează o observabilitate crescută asupra comportării programelor. Dar definirea variabilelor intermediare de observat constrânge pe cei care dezvoltă programe să producă aceste variabile, ceea ce reduce diversitatea în programare si independenta.

Versiune unică sau versiuni multiple (N)? Este o întrebare care nu are un răspuns simplu. O presupunere de ordin practic: dezvoltând N versiuni, cheltuielile cresc de N ori față de cele pentru o singură versiune. Unele părți ale procesului de dezvoltare de software pot fi comune. De pildă, dacă toate versiunile utilizează aceeași specificatie, trebuie dezvoltat și pus la punct numai un set de specificații.

Gestiunea unui proiect pe N versiuni impune overhead aditional.

Costurile pot fi reduse prin identificarea celor mai critice porțiuni ale codului și prin dezvoltarea de versiuni diferite numai pentru acestea.

Fiind dat un timp total și un buget, sunt două alegeri posibile: (a) dezvoltarea unei singure versiuni utilizând întregul buget sau (b) dezvoltarea a N versiuni.

Nu există un model satisfăcător pentru a alege între cele două posibilități.

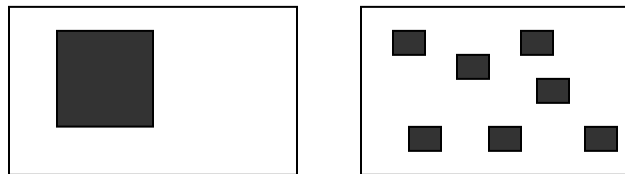
Rezultate de natură experimentală. Sunt puține studii experimentale asupra eficienței programării în versiuni multiple. Rezultatele publicate sunt generate numai de activitatea derulată în universități. Iată un studiu făcut la universități din Virginia și din California:

- 27 de studenți au scris un cod pentru aplicații anti-rachetă.
- Unii nu aveau nici o experiență în industria software, iar alții aveau peste zece ani de experiență.
- Toate versiunile au fost scrise în Pascal.
- Prin metode de verificare a ipotezelor statistice standard s-au identificat 93 de defecte corelate: dacă versiunile ar fi fost independente statistic, numărul de defecte corelate așteptat ar fi fost nu mai mare de 5.
- Nu s-a identificat nici o corelație între calitatea programelor produse și experiența programatorilor.

Diversitatea de date

Spațiul intrărilor unui program poate fi divizat în regiuni asociate cu defecte și regiuni fără vreo asociere cu defectele programului. Programul esuează dacă și numai dacă i se aplică o intrare din regiunea asociată cu defecte.

Se consideră (nu tocmai nerealist) un spațiu de intrare cu două dimensiuni (v.figura). În ambele cazuri, regiunile cu defecte ocupă o treime din spațiu.



Dacă se perturbă întrucâtva intrarea, valorile noi de intrare pot cădea în regiunea fără vreo asociere cu defectele. Această procedură este numită *diversitate de date*.

Un modul de software care utilizează teste de acceptare reface calculul cu intrări perturbate și reverifică ieseșile. Este vorba aici de o redundanță masivă: se aplică seturi ușor diferite de intrări pe diferite versiuni de software și se face o votare pe rezultate.

Perturbarea despre care se vorbește mai sus poate fi explicită și/sau implicită. Explicită este atunci când se adaugă unui subset de intrări selectat un termen-deviație redus. Implicit este când se colectează intrări pentru program astfel încât ne putem aștepta ca ele să fie ușor diferite.

Exemplul 1: Controlul prin software al proceselor industriale; intrările sunt presiunea și temperatura unui boiler. La fiecare secundă se măsoară (p_i, t_i) și se servesc ca intrări regulatorului. Măsurătorile la momentul i nu sunt mult diferite de cele de la momentul anterior $i - 1$. Perturbarea implicită poate să consistă în utilizarea perechii (p_{i-1}, t_{i-1}) în loc de (p_i, t_i) . Dacă perechea (p_i, t_i) este în regiunea cu defecte, perechea (p_{i-1}, t_{i-1}) poate că nu este de-acolo.

Perturbarea explicită, intrări reordonate.

Exemplul 2: Se adună numerele în virgulă mobilă a, b, c : mai întâi se face $a + b$, apoi se adună c . Valori: $a = 1.1E+20, b = 5, c = -1.1E+20$. În funcție de precizia utilizată, $a + b$ poate fi $1.1E+20$ și rezultatul $a + b + c = 0$.

Se schimbă ordinea intrărilor la a, c, b : atunci $a + c = 0$ și $a + c + b = 5$.

În exemplul 2 re-exprimarea este exactă, ieseșea poate fi utilizată așa cum este (dacă trece testul de acceptare sau votul).

În exemplul 1 re-exprimarea este inexactă, este probabil ca $f(p_i, t_i) \neq f(p_{i-1}, t_{i-1})$. Se utilizează ieseșea brută ca o alternativă degradată dar acceptabilă sau se încearcă corectarea înainte de utilizare, de pildă prin dezvoltarea Taylor

$$f(t) = f(t_0) + \sum_{n=1}^{\infty} \frac{(t - t_0)^n f^{(n)}(t_0)}{n!}$$

SIHFT (Software Implemented Hardware Fault Tolerance)

Variația din diversitatea de date poate fi utilizată pentru rezolvarea căderilor permanente hardware. Fiecare intrare este multiplicată cu o constantă k și programul este transformat pentru a corecta această multiplicare.

Exemplu: constructul **if (x=y) then ...**

Pentru valorile $x = 001, y = 000$, dacă egalitatea este verificată de hardware care are bitul x_0 agățat-la-0, se va calcula eronat $x = y$.

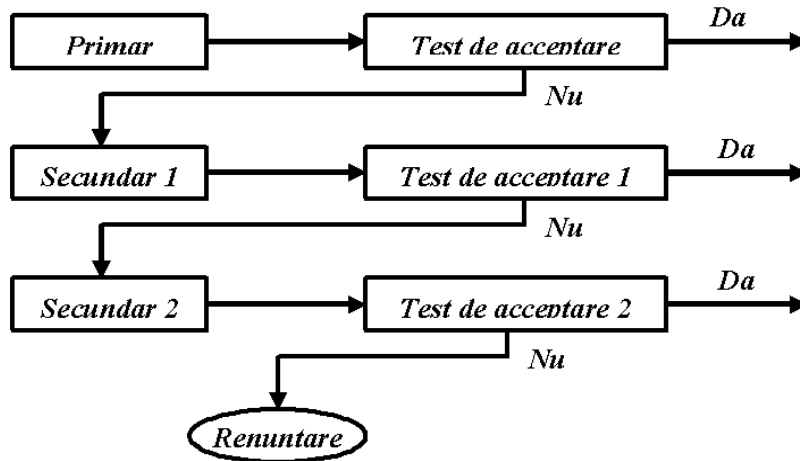
SIHFT cu $k = 2$ produce $x' = 010, y' = 000$ și defectul nu acționează, iar circuitul determină corect că $x \neq y$.

Transformarea programului în sensul de a compensa multiplicarea cu k nu este dificilă. Singura problemă este găsirea unei valori potrivite pentru k :

- (1) Ne asigurăm că este posibil a găsi tipuri de date potrivite astfel ca depășirile aritmetice în jos și în sus să nu apară.
- (2) Se alege k astfel încât să fie capabil a masca o mare cotă din defectele de hardware; se fac studii experimentale prin injectarea de defecte.

Realizarea unui bloc de recuperare

N versiuni, o executie: dacă aceasta esuează, executia e comutată la un backup.
Exemplu: o versiune primară și două versiuni secundare. Versiunea primară se execută, iesirea se dirijează la testul de acceptare. Dacă iesirea nu este acceptată, starea sistemului este restabilită (*rolled back*) și este pornită versiunea secundară 1, și tot așa în continuare. Dacă toate versiunile esuează, calculul însuși esuează. Succesul tratării prin bloc de recuperare este condiționat de independența disfuncțiilor pe versiuni diferite și de calitatea testului de acceptare.



Model analitic pentru tratarea cu bloc de recuperare. Ipoteză de lucru: versiunile diferite esuează independent.

Notatii:

E – evenimentul “iesirea unei versiuni este eronată”

T – evenimentul “testul esuează” (testul nu detectează un defect, o eroare)

f – probabilitatea eșecului unei versiuni: $f = \Pr(E)$

s – sensibilitatea testului: $s = \Pr(T|E)$

σ – specificitatea testului: $\sigma = \Pr(E|T)$

N – numărul de versiuni software.

Calculul probabilității unei disfuncții urmează schema de mai jos.

Se calculează probabilitatea complementară, aceea a succesului. Pentru ca schema să aibă succes, trebuie să aibă succes într-una din etape i , $1 \leq i \leq N$. La etapele $1, \dots, i - 1$ atât software-ul cât și testul esuează, la etpa i versiunea software este corectă și iesirea trece testarea

$$\Pr(\text{succes în etapa } i) = [\Pr(E \cap T)]^{i-1} \Pr(\bar{E} \cap \bar{T})$$

$$C = \Pr(T|E)\Pr(E) = s.f$$

$$\Pr(\bar{E} \cap \bar{T}) = \Pr(\bar{E}) - \Pr(\bar{E} \cap T)$$

$$\Pr(\bar{E}) = 1 - \Pr(E) = 1 - f$$

$$\Pr(\bar{E} \cap T) = \Pr(\bar{E} | T)\Pr(T)$$

$$\Pr(\bar{E} | T) = 1 - \sigma$$

$$\Pr(E|T) = \Pr(E \cap T) / \Pr(T)$$

în consecință

$$\Pr(T) = \frac{\Pr(F \cap T)}{\Pr(E | T)} = \frac{s \cdot f}{\sigma}$$

$$\Pr(\bar{E} \cap T) = (1 - \sigma) \frac{s \cdot f}{\sigma}$$

și

$$\Pr(\bar{E} \cap \bar{T}) = (1 - f) - (1 - \sigma) \frac{s \cdot f}{\sigma}$$

Prin substituire și prin însumare după i se obține

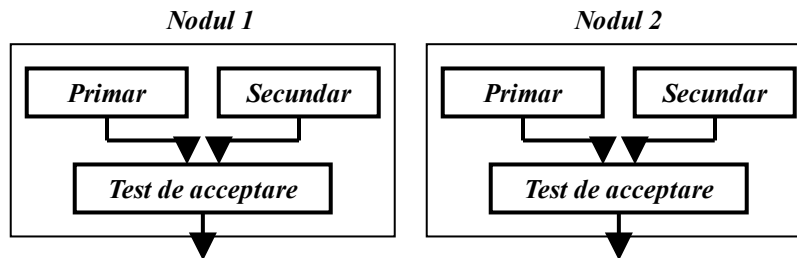
$$\Pr(\text{schema să aibă succes}) = \sum_{i=1}^n (s \cdot f)^{i-1} \left[(1 - f) - (1 - \sigma) \frac{s \cdot f}{\sigma} \right] =$$

$$= \left[(1 - f) - (1 - \sigma) \frac{s \cdot f}{\sigma} \right] \frac{1 - (s \cdot f)^n}{1 - s \cdot f}$$

$$\Pr(\text{eșec}) = 1 - \Pr(\text{schema să aibă succes})$$

Importanță critică are specificitatea σ a testului de acceptare care trebuie să fie înaltă.

Blocuri de recuperare distribuite. Două noduri poartă copii identice ale versiunilor primară și secundară.



Nodul 1 execută versiunea primară. În paralel, nodul 2 execută versiunea secundară. Dacă nodul 1 esuează la testul de acceptare, este folosită ieșirea nodului 2 (dat fiind că aceea a trecut testul). Ieșirea nodului 2 poate fi utilizată și dacă nodul 1 nu reușește să producă o ieșire într-un timp pre-specificat.

Odată ce versiunea primară esuează, rolurile versiunilor primară și secundară sunt inversate. Nodul 2 continuă să execute copia secundară care acum este tratată ca primară. Executia în nodul 1 a versiunii primare este utilizată ca un backup. Aceasta continuă până când executia din nodul 2 este semnalată ca eronată; atunci sistemul comută invers, utilizând executia din nodul 2 în rol de backup.

Revenirea (*roollback*) nu este necesară. Se economiseste timp, timp util pentru sistemele în timp real cu termene foarte ferme pentru taskuri.

Schema poate fi extinsă la N versiuni (una primară plus $N - 1$ versiuni secundare cu rulare în paralel pe N procesoare).

Tratarea exceptiilor. O excepție indică faptul că ceva s-a întâmplat în timpul execuției, ceva care merită atenție. Controlul este transferat la un *handler* de excepție, o rutină care acționează potrivit.

Exemplu: când se execută $y = a*b$, dacă apare o depășire, rezultatul este incorect; se generează un semnal de excepție.

Tratarea efectivă a excepțiilor poate aduce o creștere semnificativă a toleranței la defecte a sistemului. În multe programe, peste jumătate din liniile unui cod sunt dedicate tratării unor excepții.

Excepțiile pot fi utilizate pentru a rezolva:

- Disfuncții de domeniu sau de limite
- Evenimente extra-ordinare (nu neapărat disfuncții) care necesită o atenție aparte
- Esecuri de temporizare.

Disfuncții de domeniu și de limite. O disfuncție de domeniu se produce când sunt utilizate intrări ilegale.

Exemplu: dacă x, y sunt numere reale și se încearcă $x = \sqrt{y}$ cu $y = -1$, se produce o disfuncție de domeniu.

O disfuncție de limite apare când programul produce o ieșire sau execută o operație care este percepută ca eronată într-un anumit sens.

Exemple:

- Atingerea unui sfârșit de fișier (EOF) în timpul citirii unui fișier de date
- Producerea unui rezultat care violează un test de acceptare
- Încercarea de a tipări o linie care este prea lungă
- Generarea unei depășiri aritmetice de tip *overflow* sau *underflow*.

Evenimente extra-ordinare (neobisnuite). Excepțiile pot fi utilizate pentru a asigura tratarea specială a unor evenimente rare dar perfect normale.

Exemplu: citirea unui ultim articol al unei liste dintr-un fișier poate declanșa o excepție pentru a notifica pe cel care a invocat citirea că acela a fost ultimul articol.

Disfuncții de temporizare. În aplicațiile în timp real, taskurile au termene fixe (*deadlines*). Dacă termenele sunt depășite, se poate declanșa o excepție. Handlerul de excepție decide ce trebuie făcut în replică (de pildă comutarea la o rutină backup).

Cerinte relativ la handlerele de excepții.

- (1) Trebuie să fie ușor de programat și de utilizat. Să fie modulare și separabile de restul software-ului. Să nu fie amestecate cu alte linii de cod într-o rutină; ar putea fi greu de înțeles, de depanat, de modificat.
- (2) Tratarea excepției nu trebuie să impună un overhead substanțial asupra funcționării normale a sistemului. Excepțiile să fie invocate numai în împrejurări într-adevăr excepționale. Tratarea excepției să nu fie o povară în cazul uzual fără condiții de excepție.

- (3) Tratarea exceptiilor nu trebuie să compromită starea sistemului, să nu-l lase inconsistent.

Definiii pentru fiabilitatea software

Disfuncție: abaterea comportării software-ului de la cerințele utilizatorului.

Fiabilitate: probabilitatea de operare software fără disfuncții într-o ambianță definită, pentru o perioadă de timp specificată.

Software-ul nu se deteriorează în timp ca hardware-ul; el rămâne constant, același dacă nu sunt operate schimbări.

Dacă în perioada testării sunt detectate și eliminate defecte, fiabilitatea software crește în timp.

Notatii:

N – numărul de defecte existente la începutul testării (poate fi o variabilă aleatoare)

$M(t)$ – numărul de defecte detectate și eliminate până la momentul t

$N - M(t)$ – numărul de defecte rămase la timpul t .

Modele pentru fiabilitatea software

Se încearcă predicția ratei de defectare viitoare a software-ului ca funcție de numărul de defecte eliminate sau de numărul de defecte rămase la momentul t . Spre deosebire de modelele fiabilității hardware, acestea sunt netestate în sensul cel mai larg al notiunii.

Modelele disponibile dau adesea rezultate contradictorii.

Când testarea porneste, defectele mai “usoare” sunt capturate curând. Defectele rămase sunt din ce în ce mai dificil de capturat; fie sunt mai greu de întâlnit, fie efectul lor este mascat de calcule următoare manifestării lui.

Rata la care defectele nedescoperite încă produc disfuncții scade pe măsură ce testarea avansează. Rata de defectare este descrisă fie ca o funcție descrescătoare de $M(t)$, fie ca o funcție crescătoare de $N - M(t)$.

Modelul Jelinski-Moranda.

Rata de defectare $\lambda(t)$ este proporțională cu numărul de defecte rămase în software: $\lambda(t) = c[N - M(t)]$. Când un defect este detectat, el este imediat eliminat; timpul până la disfuncția următoare este distribuit exponențial cu parametrul $\lambda(t)$.

Problemă: nu toate defectele se manifestă egal; unele apar mai frecvent, altele mai rar și sunt dificil de (sur)prins.

Modelul Littlewood-Verall.

Timpul scurs între disfuncții este distribuit exponențial cu parametrul $\lambda(i)$ în care $i = N - M(t)$ este numărul de defecte rămase. Parametrul $\lambda(i)$ este o variabilă aleatoare distribuită cu o densitate de probabilitate de tipul *gamma*

$$f(x) = \frac{[\psi(i)]^\alpha x^{\alpha-1} e^{-\psi(i)x}}{\Gamma(\alpha)}$$

Funcția *gamma* este într-un fel o generalizare a funcției factorial

$$\Gamma(t) = \int_0^{\infty} e^{-x} x^{t-1} dx \quad \text{și} \quad \Gamma(t+1) = t\Gamma(t)$$

Funcția $\psi(i)$ și parametrul α se determină prin observarea experimentală a software-ului.

Modelul Musa-Okumoto.

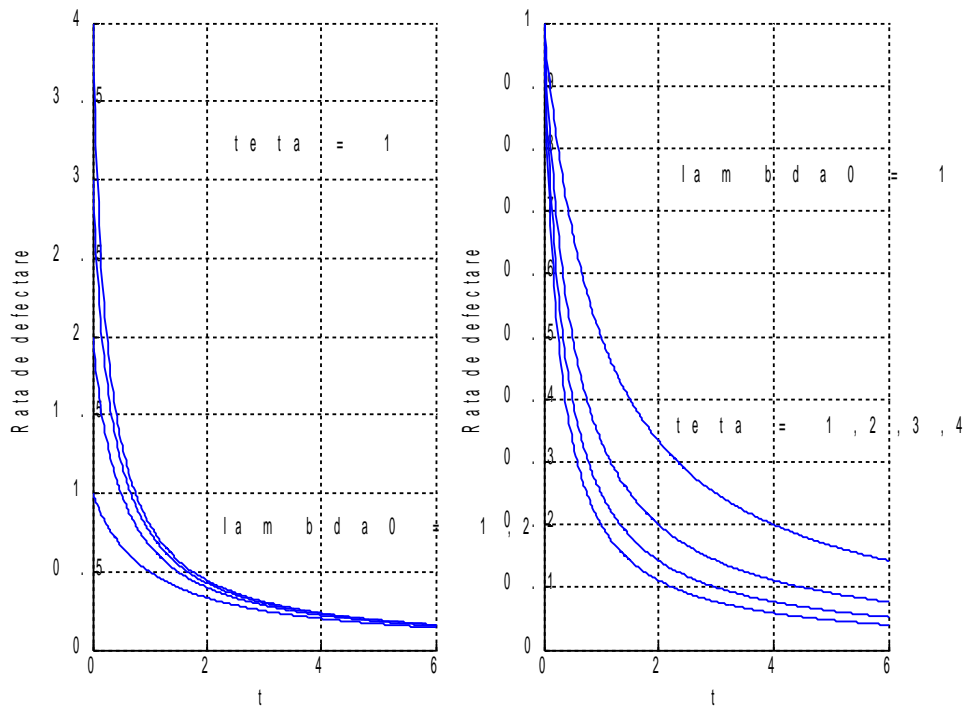
Este un model cu timpul de execuție Poisson logaritmic. Este un model mai larg utilizat pentru fiabilitatea software.

Rata de defectare după timpul t este

$$\lambda(t) = \lambda(0) \exp[-\theta \mu(t)]$$

cu θ o constantă și cu $\mu(t)$ valoarea medie a lui $M(t)$, numărul de disfuncții manifestate și înlăturate înainte de momentul t .

$$\frac{d\mu(t)}{dt} = \lambda(t)$$



$$\lambda_0 = \frac{d\mu(t)}{dt} e^{\theta \mu(t)}$$

$$\frac{de^{\theta \mu(t)}}{dt} = \theta \frac{d\mu(t)}{dt} e^{\theta \mu(t)} = \lambda_0 \theta$$

$$e^{\theta \mu(t)} = \lambda_0 \theta t + 1$$

$$\mu(t) = \theta^{-1} \ln(\lambda_0 \theta t + 1)$$

$$\lambda(t) = \frac{\lambda_0}{\lambda_0 \theta t + 1}$$

Aceasta este expresia unei foarte lente scăderi a ratei de defectare, ceea ce duce la o cantitate apreciabilă de observatii experimentale.

Selectarea modelului si estimarea de parametri. Întrebări legitime:

- (1) Care model este cel mai potrivit?
- (2) Cum se estimează parametrii modelului ales?

Datele experimentale nu sunt foarte cuprinzătoare pentru a ghida utilizatorul. Se studiază rata de defectare ca functie de testare si se ghiceste mai curând modelul pe care-l urmează; apoi se estimează parametrii lui.

Se utilizează tehnicile de estimare statistică, de pildă metoda verosimilității maxime sau metoda celor mai mici pătrate.

REDUNDANTE INFORMATIONALE

Redundanta informatională (codarea)

Un cuvânt cu d biti este codat într-un cuvânt de cod de c biti, $c > d$. Rezultă un număr de 2^c cuvinte binare, dar nu toate cele 2^c combinații de biti sunt cuvinte de cod valide.

Înainte de a putea utiliza informația este necesar să se extragă din cuvântul de cod cuvântul original. Altfel zis cei c biti trebuie decodati.

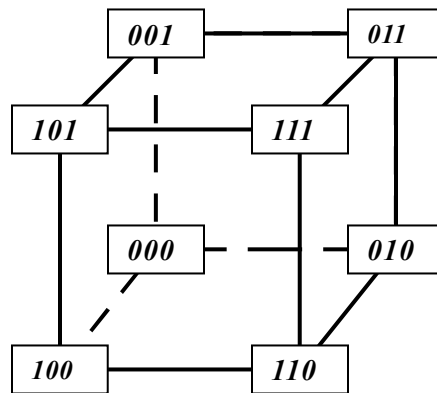
Dacă cei c biti primiti nu constituie un cuvânt de cod valid, înseamnă că s-a detectat o eroare. Pentru anumiți algoritmi de codare, unele tipuri de erori pot fi nu numai detectate ci pot fi chiar corectate la destinație, la utilizator.

Parametrii cheie în operația de codare/decodare îl reprezintă numărul de biti eronați care pot fi detectați ca eronați și numărul de biti eronați care pot fi corectati. În general, aceste două numere nu coincid.

Într-o privire foarte generală, bitii suplimentari necesari pentru a crea posibilitatea de detectare/corectare a erorilor fac să crească timpul consumat, și nu numai prin operațiile de codare și de decodare ci și prin durata transmiterii unui număr mai mare de biti.

Distanța Hamming

Distanța Hamming între două cuvinte (de cod) este dată de numărul de poziții binare prin care cele două cuvinte diferă. Evident, distanța Hamming se poate defini numai pentru cuvinte de lungimi identice.



În figura alăturată, distanța Hamming dintre două cuvinte care sunt separate (legate) de o muchie este 1. Două cuvinte separate prin două muchii (ca cel mai scurt parcurs de la unul la altul) sunt distanțate Hamming cu 2. Cea mai mare distanță Hamming, 3, este între cuvintele situate în extremitățile unei diagonale a cubului.

Odată definită distanța Hamming, se poate defini distanța unui cod. *Distanța unui cod* este distanța Hamming minimă între două cuvinte valide distincte ale codului.

De exemplu, codul de patru cuvinte {001, 010, 100, 111} are o distanță egală cu 2. Acest cod poate detecta o eroare pe un singur bit (inversarea unui bit).

Un alt exemplu: codul cu două cuvinte {000, 111} are o distanță de 3. Acest cod poate detecta orice eroare de un bit sau de doi biti (unul sau doi biti inversați). Dacă o eroare dublă în același cuvânt de cod este suficient de improbabilă, un alt mod de a spune “are o probabilitate foarte mică față de eroarea unică, singulară”, codul poate corecta un bit eronat.

Exemplele de mai sus sugerează vag, dar sugerează, condițiile necesare pentru a putea face o detecție sau o corecție a cuvântului de cod recepționat. Regulile respective sunt:

- Pentru a detecta într-un cuvânt de cod până la t biti eronați, distanța acelui cod trebuie să fie de cel puțin $t + 1$
- Pentru a corecta până la t biti eronați într-un cuvânt de cod, distanța codului trebuie să fie de cel puțin $2t + 1$.

Aceste condiții sunt demonstrate riguros în multe surse, în particular în lucrarea Gh.M.Panaitescu *Transmiterea și codarea informației. Note de curs* (2007), p.81 (<http://ac.upg-ploiesti.ro/gpanaitescu/tci.pdf>).

Codare și redundanță

Codul {000, 111} poate fi utilizat pentru a codifica date de un bit: 0 poate fi codat ca 000 și 1 ca 111. Codul acesta lucrează, într-un fel, identic cu TMR (sistemul Triple Modular Redundant) și anume prin “vot” majoritar. Sub acest aspect, multe tehnici redundante hardware pot fi considerate scheme de codare. Și duplexul poate fi considerat un cod ale cărui cuvinte constau în două cuvinte identice alăturate. Pentru datele de un bit, cuvintele de cod sunt 00 și 11.

Separabilitatea unui cod

Un cod este *separabil* dacă are câmpurile pentru bitii de date și bitii de codare în secvențe separate. La decodare bitii de date sunt recuperați prin simpla ignorare a bitilor adăugați pentru codare. Asadar, bitii de codare pot fi ei înșiși procesați separat pentru a verifica corectitudinea acelor de date.

Un cod *neseparabil* are bitii de date și bitii de codare integrați și extragerea datelor din cuvântul de cod necesită o procesare specifică.

Coduri prin paritate

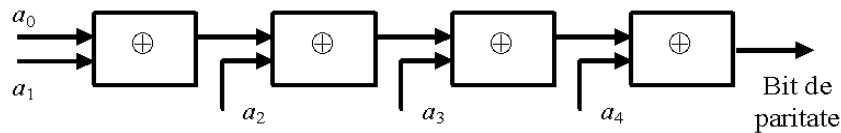
Cele mai simple coduri separabile sunt codurile cu paritate (constantă). Un cuvânt al unui cod cu paritate (constantă) conține n biti informaționali și un bit suplimentar care menține paritatea. Pentru un cod cu paritate cu sot (sau fără sot), bitul suplimentar este fixat astfel ca numărul total de unități binare din cuvântul de $(n + 1)$ biti să fie totdeauna par (sau impar). Fracția suplimentară inevitabilă a unui cod cu paritate constantă este $1/n$.

Proprietăți ale codurilor cu paritate constantă

Un cod cu paritate constantă are o distanță Hamming de 2 și este capabil să detecteze erorile de un singur bit într-un cuvânt de cod. Dacă un bit se schimbă din 0 în 1 sau invers, paritatea se schimbă și eroarea poate fi detectată prin verificarea parității. Această paritate simplă nu poate însă corecta bitul eronat.

Circuite de codare și de decodare pentru coduri cu paritate constantă

Codorul este un sumator de biti modulo-2. Acesta generează un 0 dacă numărul de unități binare din cuvântul-date este cu sot, un 1 dacă numărul de unități binare a cuvântului-date este fără sot. Iesirea sumatorului este semnalul de paritate.



Codur

Decodorul re-generează bitul de paritate din bitii care constituie partea de date în cuvântul recepționat într-o schemă identică cu cea a codorului și îl compară cu bitul de paritate primit într-un circuit nici-exclusiv. Dacă există potrivire între cei doi biti de paritate, cel evaluat și cel primit, iesirea porții nici-exclusiv este un 1, ceea ce indică absența erorii. Lipsa coincidenței produce în același punct al schemei un 0 ceea ce indică o eroare.

Erorile pe doi biti nu pot fi detectate de un cod cu paritate constantă: paritatea nu este modificată prin inversarea a doi biti.

Toate erorile de trei biti într-un cuvânt, desigur o situație mult mai rară, sunt detectate fără a putea discerne dacă este vorba de una sau trei erori în bitii cuvântului de cod.

După cum s-a amintit mai devreme, paritatea utilizată poate fi cu sot, fie fără sot. Când se alege una, când se alege cealaltă? Uzual este preferată paritatea cu sot dar decizia depinde uneori de tipul mai probabil de erori pe toți bitii cuvântului.

Pentru paritatea pară, bitul de paritate pentru toti bitii nuli este 0 si o cădere de tipul *toti-bitii-0* trece neobservată, deoarece apare deghizată ca un cuvânt de cod valid. Prin alegerea tipului de paritate impar, eroarea de tipul *toti-bitii-0* devine detectabilă.

Dacă eroarea de tipul *toti-bitii-1* este mai probabilă, atunci, după caz, dacă numărul total de biti ($n + 1$) este par trebuie selectată paritatea impară, iar dacă ($n + 1$) este un număr impar trebuie reținută si utilizată paritatea pară.

Un sistem cu paritate constantă este si acela în care bitii de paritate sunt atribuiti separat unui byte sau oricărui alt grup de biti de lungime stabilită. Proportia de biti suplimentari se modifică de la $1/n$ la m/n (m este numărul de bytes-i sau de alte grupe egal dimensionate). Cu acest sistem se pot detecta până la m erori dacă erorile apar în grupe de biti (eventual bytes-i) diferite. Dacă penele de tipurile *toti-bitii-0* si *toti-bitii-1* sunt în linii mari la fel de posibile, atunci se selectează alternativ paritatea pară pentru un byte, paritatea impară pentru alt byte.

Codul cu paritate pentru bytes-i întretesuti

Cuvântul de cod de lungime $n = 64$, de pildă, arată astfel: $a_{63}a_{62}\dots a_1a_0$. Sunt cuprinsi aici opt biti de paritate: primul bit al fiecărui octet, cel mai semnificativ, este bit de paritate. Pentru exemplul cu $n = 64$, bitii de paritate sunt $a_{63}, a_{55}, a_{47}, a_{39}, a_{31}, a_{23}, a_{15}, a_7$. Primul bit de paritate a_{63} este cel mai semnificativ bit în lantul de 8 bytes-i considerat. Ceilalti 7 biti de paritate se atribuie astfel ca grupele corespondente de biti să fie întretesute. Schema aceasta este benefică atunci când scurtcircuitarea unor biti adiacenti este un mod de defectare comun (de exemplu pe un bus). Dacă tipul de paritate (pară – impară) este alternat între grupe, erorile unidirectionale (*toti-bitii-0* sau *toti-bitii-1*) vor fi si ele detectate.

Coduri cu paritate capabile a corecta erori

În schema cea mai simplă, informatia este organizată într-o matrice bidimensională

0	0	0	1	1	1	1
1	0	1	0	1	1	0
1	1	0	0	0	0	0
0	0	0	1	1	1	1
1	1	1	1	1	1	0
1	0	0	1	0	0	0

Bitii de la finele unei linii sunt biti de paritate pe acea linie, bitii de la baza unei coloane sunt de paritate pe acea coloană. O eroare pe un bit situată undeva, oriunde în matricea utilizată, afectează corectitudinea parității atât pe o linie cât

si pe o coloană. Identificarea liniei si coloanei care contin bitul eronat nu este o problemă, corectarea lui este imediată.

Acesta este un exemplu de *paritate suprapusă*: fiecare bit este acoperit de mai mult de un bit de paritate.

Modelul general pentru paritatea suprapusă

Scopul urmărit este identificarea vreunui bit eronat singular, presupus unicul bit eronat din cuvânt. În structura cuvântului de cod intră d biti de informație si r biti de paritate, asadar cuvântul are un total de $d + r$ biti. Admitând că erorile sunt pe un bit si pe numai unul singur, sunt $d + r$ stări cu eroare si o stare fără eroare, ceea ce dă un total de $d + r + 1$ stări distincte ale cuvântului receptionat. În consecința acestor observatii, pentru a distinge între stările enumerate sunt necesare $d + r + 1$ “semnături” de paritate distincte, tot atâtea configurații de biti. Aceste r verificări de paritate generează 2^r semnături. Asadar, r este cel mai mic întreg care satisface inegalitatea $2^r \geq d + r + 1$.

Apare întrebarea firească: cum pot fi atribuiti bitii de paritate? Iată un exemplu: Fie $d = 4$ numărul bitilor de date si $r = 3$ numărul (minim) de biti de paritate. Într-adevăr, sunt $d + r + 1 = 8$ stări în care poate fi cuvântul de cod, sunt $2^3 = 8$ adaosuri distincte la secvența de biti de date.

Tabelul alăturat arată o atribuire posibilă, nu singura, de biti de paritate stărilor cuvântului receptionat. În cuvântul $(a_3a_2a_1a_0p_2p_1p_0)$ pozițiile 0, 1 si 2 sunt consacrate bitilor de verificare a parității, restul sunt biti de informație.

Starea	Biti de paritate
Fără eroare	000
Eroare pe bitul 0 (p_0)	001
Eroare pe bitul 1 (p_1)	010
Eroare pe bitul 2 (p_2)	100
Eroare pe bitul 3 (a_0)	011
Eroare pe bitul 4 (a_1)	101
Eroare pe bitul 5 (a_2)	110
Eroare pe bitul 6 (a_3)	111

Aceiasi parametri d , r îi are si un cod consacrat si foarte utilizat: codul Hamming (7,4) corector de o singură eroare (SEC – Single Error Corrector). După cum se poate observa, bitii de paritate în afară de faptul că numerează în baza 2 cele opt variante posibile ale cuvântului receptionat, 000 pentru cuvântul fără eroare, cu o eroare pe unul din cei 7 biti poziționată diferit de la caz la caz, acesti biti de paritate mai sunt si într-o anumită relație cu bitii de informație. Astfel, p_0 acoperă pozițiile 0, 3, 4, 6, adică $p_0 = a_0 \oplus a_1 \oplus a_3$, p_1 acoperă pozițiile 1, 3, 5, 6, adică $p_1 = a_0 \oplus a_2 \oplus a_3$ si p_2 acoperă pozițiile 2, 4, 5, 6, prin $p_2 = a_1 \oplus a_2 \oplus a_3$.

Un exemplu: pentru $a_3a_2a_1a_0 = 1100$ completarea cu biti de paritate este $p_2p_1p_0 = 001$. Dacă 1100001 devine $\underline{1}000001$, cu bitul eronat subliniat, se recalculează secvența de paritate din cuvântul eronat: $p_2p_1p_0 = 111$. Diferența dintre secvențele de control al parității se obține printr-o operație “sau exclusiv” bit-cu-bit, care este 110 . Diferența aceasta se numește *sindrom* și indică bitul eronat. Bitul a_2 este eronat și informația corectă este $a_3a_2a_1a_0 = 1100$. Pentru codul Hamming (7,4) sindromul poate fi calculat direct, într-un pas, din bitii recepționați $a_3, a_2, a_1, a_0, p_2, p_1, p_0$, prin operația matricială următoare, cu toate operațiile făcute *modulo 2*.

$$\begin{bmatrix} a_3 & a_2 & a_1 & a_0 & p_2 & p_1 & p_0 \\ 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} a_3 \\ a_2 \\ a_1 \\ a_0 \\ p_2 \\ p_1 \\ p_0 \end{bmatrix} = \begin{bmatrix} s_2 \\ s_1 \\ s_0 \end{bmatrix}$$

Matricea care (pre)multiplică cuvântul recepționat reprezentat aici ca un vector de valori binare este așa-numita matrice de verificare. Cuvintele de cod verifică relația de mai sus cu un sindrom cu toate componentele nule. Prezența erorii este semnalată de un vector sindrom nenul. Identificarea bitului eronat (și corectarea lui) se face prin potrivirea sindromului pe una din coloanele matricei de verificare. Dacă sindromul este $[1\ 1\ 0]^T$, cum s-a întâmplat în exemplul dat, acesta se potrivește (numai) pe coloana a_2 . Consecința: acela este bitul eronat și trebuie corectat (prin inversare).

Într-o pre-ordonare potrivită a bitilor de informație și a bitilor de paritate, bitul eronat poate fi localizat scăzând o unitate din indicele sindromului: rezultatul este indicele bitului eronat. În exemplul prezentat mai sus, ordinea aceea este $a_3a_2a_1a_0p_2p_1p_0$.

În general, dacă $2^r > d + r + 1$ este necesar a alege pentru a fi combinații sindrom $d + r + 1$ combinații binare din cele 2^r . Combinațiile cu mulți de 1 trebuie evitate deoarece mai puțini de 1 în matricea de verificare a parității înseamnă practic circuite mai simple pentru operațiile de codare și de decodare. Selectarea matricei de verificare este de oarecare importantă. De exemplu, pentru $d = 3$ și $r = 3$, numai șapte din cele opt combinații de 3 biti sunt necesare. Sunt posibile două matrice de verificare:

$$\begin{bmatrix} a_2 & a_1 & a_0 & p_2 & p_1 & p_0 \\ 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 \end{bmatrix} \text{ și } \begin{bmatrix} a_2 & a_1 & a_0 & p_2 & p_1 & p_0 \\ 0 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Prima are prea mulți de 1 (111 pe prima coloană), a doua este mai potrivită. Pentru prima circuitele de codare necesită câte o poartă XOR (sau exclusiv)

pentru p_1 și p_2 , dar două pentru bitul p_0 . Pentru a doua matrice circuitele de codare necesită câte o poartă XOR pentru fiecare bit de paritate.

Îmbunătățirea detecției este un obiectiv frecvent în toleranța la defecte. Codul descris mai devreme poate corecta un bit eronat și poate detecta erori pe doi biți. De exemplu, 1100001 poate deveni 1010001, cu bitii subliniați, a_2 și a_1 , eronați. Sindromul este în cazul respectiv 011. Acest sindrom indică gresit că necesară o corectare a bitului a_0 . O cale de a îmbunătăți capacitatea de detecție constă în a suplimenta cu un bit secvența de verificare. Acesta este un bit de verificare a parității pentru toți ceilalți biți, de informație sau de paritate. Codul acesta este cunoscut ca un cod Hamming (8,4) de corectare a unui bit eronat/de detectare a doi biți eronați (SEC/DED – Single Error Correction/Double Error Detection).

Generarea sindromului pentru codul Hamming (8,4) este similară celei pentru codul Hamming (7,4): prin multiplicarea cuvântului cu o matrice de verificare. Rezultatul este diferit: un vector sindrom cu 4 componente.

$$\begin{matrix}
 & a_3 & a_2 & a_1 & a_0 & p_3 & p_2 & p_1 & p_0 & & \\
 \begin{matrix} a_3 \\ a_2 \\ a_1 \\ a_0 \\ p_3 \\ p_2 \\ p_1 \\ p_0 \end{matrix} & \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 \end{bmatrix} & = & \begin{bmatrix} s_3 \\ s_2 \\ s_1 \\ s_0 \end{bmatrix}
 \end{matrix}$$

Bitul p_3 este, cum s-a spus, un bit de paritate pentru toți bitii de informație și de control. O eroare de un singur bit schimbă paritatea generală și produce $s_3 = 1$. Ultimii trei biți ai sindromului indică bitul eronat, care trebuie corectat ca și mai devreme dacă $s_3 = 1$. Dacă $s_3 = 0$ și oricare alt bit al sindromului este nenul, se detectează o eroare dublă sau multiplă de ordin superior lui 2. Un exemplu:

Eroare unică – 11001001 devine 10001001

Sindromul este 1110 – indică eroare pe a_2

Eroare dublă – 11001001 devine 10101001

Sindromul este 0011 și indică o eroare (multiplă) care nu poate fi corectată

Iată acum un cod Hamming (8, 4) diferit. Codul Hamming (8, 4) prezentat și comentat mai devreme necesită calculul bitului de verificare suplimentar, ceea ce aduce un plus în consumul de timp la codare și la decodare. Pentru a evita această situație, o soluție posibilă constă în atribuirea de sindromuri cu număr impar de unități binare.

$$\begin{matrix} a_3 & a_2 & a_1 & a_0 & p_3 & p_2 & p_1 & p_0 \\ \left[\begin{array}{cccccccc} 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 & 1 \end{array} \right] \end{matrix}$$

O eroare dublă produce un sindrom cu un număr par de 1, ceea ce indică o eroare (multiplă) care nu poate fi corectată. În acest cadru, se utilizează numai 2^{p-1} combinații din cele 2^p posibile. Este necesar, prin urmare, un bit de control suplimentar peste minimul necesar pentru un cod Hamming de tipul corector de o singură eroare (SEC). Numărul total de biti de control este același cu cel necesar pentru codul Hamming original, capabil să corecteze o eroare și să detecteze o eroare dublă (SEC/DED).

Ca structură, codurile cu *parități suprapuse* sunt multe și variate. O comparație între aceste coduri este interesantă. Astfel, la codurile Hamming, pe măsură ce d crește, *overhead-ul* parității r/d , raportul biti de protecție/ biti de date, scade. Dar probabilitatea de a avea mai mult de un bit eronat în $d + r$ biti crește.

Fie f probabilitatea apariției unui bit eronat. Se presupune că erorile apar independent una de alta. Probabilitatea ca într-un câmp de $d + r$ biti succesivi să apară mai mult de un bit eronat este

$$\Phi(d, r) = 1 - (1 - f)^{d+r} - (d + r)f(1 - f)^{d+r-1} \approx 0,5(d + r)(d + r - 1)f^2$$

d	r	Overhead r/d	$\Psi(D, d, r)$
2	3	1,5000	0,5120E-16
4	3	0,7500	0,5376 E-16
8	4	0,5000	0,8448 E-16
16	5	0,3125	0,1344 E-15
32	6	0,1875	0,2250 E-15
64	7	0,1094	0,3976 E-15
128	8	0,0625	0,7344 E-15
256	9	0,0352	0,1399 E-14
512	10	0,0195	0,2720 E-14
1024	11	0,0107	0,5351 E-14

adică complementul la unitate al probabilității prezentei a cel mult o eroare, cu aproximarea valabilă pentru $f \ll 1$.

Dacă sunt D biti de date în total, probabilitatea de a avea mai mult de un bit eronat poate fi redusă prin partitionarea celor D biti în D/d porțiuni egale, cu fiecare porțiune codată separat. Este vorba aici de un compromis, un echilibru între probabilitatea erorilor nedetectabile și *overhead-ul* r/d . Probabilitatea de a avea o eroare necorectată în cel puțin una din cele D/d porțiuni este

$$\Psi(D, d, r) = 1 - [1 - \Phi(d, r)]^{D/d} \approx (D/d)\Phi(d, r)$$

cu aproximarea acceptată dacă $\Phi(d, r) \ll 1$.

Tabelul de mai sus cuprinde o comparatie numerică pentru cazul $D = 1024$ si $f = 10^{-11}$.

Codurile cu *sumă de control* sunt utilizate în principal pentru a detecta erori în transmisia de date prin rețelele de comunicatie. Ideea de bază: se însumează blocul de date care se transmite si se transmite *si* această sumă. Receptorul însumează datele primite si compară suma obtinută cu suma de verificare trimisă odată cu datele. Dacă cele două sume de control nu coincid înseamnă că transmiterea s-a făcut cu eroare.

Sunt mai multe versiuni pentru *suma de control*.

Fie d (biti) lungimea cuvintelor de date. În versiunea *simplă precizie*, suma de control rezultă ca un număr modulo 2^d . În versiunea *dublă precizie*, suma de control rezultă ca un număr modulo 2^{2d} .

Deoarece în procesul de însumare sunt retinuti numai d biti, cei mai din dreapta, este de înțeles că suma de control în simplă precizie este capabilă a cuprinde mai putine erori decât suma de control în dublă precizie. În cazul dublei precizii, sunt retinuti $2d$ biti din suma de control.

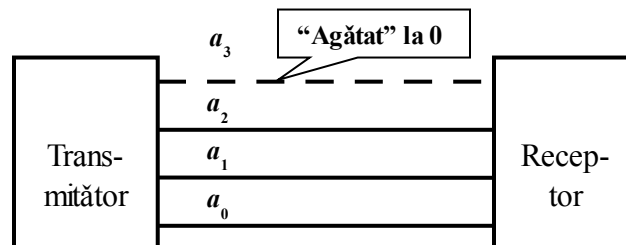
Suma de control *cu reziduu* face un transfer de la bitul al d -lea, cel mai semnificativ, la bitul cel mai puțin semnificativ, un transfer circular, ciclic, si de aceea codul este întrucâtva mai fiabil.

Suma de control *Honeywell* concatenează cuvintele în perechi pentru calculul sumei de verificare (făcută modulo 2^{2d}) si prin aceasta codul protejează la erori repetate în aceeasi pozitie a secventei.

Imediat mai jos este dată o comparatie între cele patru versiuni ale sumelor de control:

Precizie simplă	Precizie dublă	Cu reziduu	Honeywell
0000	0000	0000	
0101	0101	0101	
1111	1111	1111	00000101
0010	0010	0010	11110010
0110	00010110	0111	11110111

Un exemplu este prezentat în figura alăturată.



Precizie simplă		Honeywell	
Transmis	Receptionat	Transmis	Receptionat
1000	0000		
1011	0011	10001011	00000011
0000	0000	00000100	00000100
0100	0100		
0111	0111	10001111	00000111

În cazul sumei de control în simplă precizie, sumele pentru secvența transmisă și pentru cea calculată se potrivesc: asadar, erori nedetectate. În cazul sumei Honeywell, sumele de control pentru secvența transmisă și pentru cea calculată nu se potrivesc și erorile sunt detectate.

Toate schemele cu sumă de control au un neajuns conceptual: atunci când fac detectia de erori nu le și localizează. Dacă sunt detectate erori, întregul bloc de date trebuie retransmis.

Codurile *Berger* sunt coduri separabile. Se numără unitățile binare din cuvânt, se exprimă rezultatul în binar, se completează și se atasează rezultatul la bitii de informație.

Exemplu: Codarea secvenței 11101. Sunt patru de 1, în binar 100, 011 după complementare. Cuvântul de cod se obține prin juxtapunerea rezultatului după bitii de informație: 11101011.

d	c	<i>overhead</i>
4	3	0,7500
8	4	0,5000
16	5	0,3125
32	6	0,1875
64	7	0,1094
128	8	0,0625
256	9	0,0352
512	10	0,0195

Un cod *Berger* detectează toate erorile unidirectionale, cele care constau în faptul că una sau mai multe unități binare devin zerouri și zerourile rămân zerouri (sau invers). Totuși, dacă numărul de biți care trec din 0 în 1 este egal cu cel al celor care trec din 1 în 0, eroarea devine nedetectabilă.

Se poate evalua un overhead pentru un cod Berger. Dacă bitii de informație sunt în număr de d , sunt cel mult d de 1 și sunt necesari până la $\log_2(d + 1)$ biți pentru reprezentarea numărului de unități binare. Asadar

$$\text{overhead} = \lceil \log_2(d + 1) \rceil / d$$

Cu notația c pentru numărul de biți de verificare, se obține tabelul de mai sus.

Coduri ciclice

Codurile ciclice au această denumire derivată din proprietatea cuvintelor de cod de a rămâne în cadrul codului dacă bitii componente sunt permutati circular. Codurile ciclice sunt în același timp coduri liniare: suma a două cuvinte de cod, bit-cu-bit modulo 2, este totdeauna un cuvânt de cod.

Aceste coduri sunt uzual neseparabile, dar există și coduri ciclice separabile. Codarea constă în multiplicarea (modulo 2) a cuvântului de date cu un număr constant. Produsul este cuvântul codat. Decodarea constă în diviziunea cu aceeași constantă. Dacă restul nu este zero, de vină este o eroare.

Codurile ciclice sunt larg utilizate în stocarea și transmiterea datelor.

Teoria pe care se bazează codurile ciclice porneste de la un număr k de biti de informație care trebuie codati. Cuvântul de cod este lung de n biti obtinuti prin multiplicarea celor k biti de informație cu un număr care este ca lungime de până la $n - k + 1$ biti. Numărul multiplicator este prezentat ca un polinom, denumit și polinom generator. Unitățile și zerourile binare din multiplicatorul de $n - k + 1$ biti sunt considerate coeficienti ai unui polinom de gradul $n - k$.

Exemplu: multiplicatorul este 11001, polinomul generator este

$$G(X) = 1X^0 + 0X^1 + 0X^2 + 1X^3 + 1X^4 = 1 + X^3 + X^4$$

Un cod ciclic (n, k) este un cod ciclic care utilizează un polinom generator de gradul $n - k$ și are numărul total de biti în cuvintele de cod n . Un cod ciclic (n, k) poate detecta toate erorile singulare și toate secvențele de biti adiacenți eronați mai scurte de $n - k$ biti.

Codurile ciclice sunt utile în aplicații de genul comunicațiilor fără fir: canalele sunt uzual afectate de zgomot și de interferențe care produc erori pe secvențe de biti adiacenți.

Implementarea hardware execută multiplicarea (modulo 2) pe registre de deplasare și porți XOR (sau exclusiv). De exemplu, dacă polinomul generator este $1 + X^3 + X^4$ (care corespunde multiplicatorului 11001), circuitul de codare este alcătuit din două porți XOR și patru operații de deplasare. Registrele de deplasare sunt elemente de întârziere care mențin la ieșire intrarea lor pe durata unui ciclu.

În ceea ce privește multiplicarea modulo 2, dacă multiplicatorul este alimentat serial și se adună multiplicatul deplasat (shifted), se obține produsul

$$\begin{array}{r} 1100101 \\ 11001 \\ \hline 1100101 \\ 0000000 \\ 0000000 \\ 1100101 \\ 1100101 \\ \hline 10100011101 \end{array}$$

Al cincilea bit al produsului, de pildă, este suma modulo 2 a bitilor corespunzători ai multiplicatului (subliniați) deplasati spre stânga de 0 ori, de 3 ori și de 4 ori, deplasări de ordine egale cu ordinele bitilor nenuli ai multiplicatorului. Deplasările se realizează cu elementele de întârziere. Codul ciclic nu este separabil: bitii de informație și bitii de control nu apar separati în cuvântul de cod 10100011101.

Operațiile circuitului de codare sunt date în tabelul alăturat.

Tempo	Intrare	O_4	i_3	O_3, O_2, O_1	Iesire
1	1	0	1	000	1
2	0	1	1	100	0
3	1	0	1	110	1
4	0	1	1	111	1
5	0	0	0	111	1
6	1	0	1	011	0
7	1	1	0	101	0
8	0	1	1	010	0
9	0	0	0	101	1
10	0	0	0	010	0
11	0	0	0	001	1

Coloana i_3 conține intrarea în elementul de întârziere care produce pe O_3 . Intrarea de biti informaționali în cazul ilustrat este 1100101, iar iesirea codată este 10100011101.

Decodarea se face prin împărțirea la polinomul generator și urmează algoritmul prezentat imediat. Pentru cuvântul fără eroare și respectiv, pentru cuvântul cu eroare (eroarea este reprezentată diferit) se obține:

$$\begin{array}{r}
 10100011101 : 11001 = 1100101 \\
 \underline{11001} \\
 11010 \\
 \underline{11001} \\
 11111 \\
 \underline{11001} \\
 11001 \\
 \underline{11001} \\
 00000
 \end{array}$$

$$\begin{array}{r}
 10100\underline{1}11101 : 11001 = 1100110 \\
 \underline{11001} \\
 11011 \\
 \underline{11001} \\
 10111 \\
 \underline{11001} \\
 11100 \\
 \underline{11001} \\
 1011
 \end{array}$$

asadar, o dată rest nul, altă dată rest nenul. Restul nul arată că nu a apărut nici o eroare. La împărțirea cu 11001 practică mai sus, scăderea (aici modulo 2) cu care suntem obișnuiți la împărțirea a două numere, este identică cu adunarea.

Prezența unei/unor erori scoate de obicei cuvântul din multimea de cuvinte ale codului. Semnalarea acestei situații o face restul nenul: dacă cuvântul eronat este 10100111101, restul este, cum s-a văzut, 01011, nenul.

Să vedem ce se poate întâmpla dacă trei biti sunt eronați. Sunt date ilustrativ două cazuri. Rezultatele împărțirii cu polinomul generator sunt respectiv:

$$101\underline{11}01\underline{0}101 : 11001 = 1101101$$

$$\begin{array}{r} \underline{11001} \\ 11100 \\ \underline{11001} \\ 10110 \\ \underline{11001} \\ 11111 \\ \underline{11001} \\ 11001 \\ \underline{11001} \\ 00000 \end{array}$$

$$100\underline{11}011101 : 11001 = 1110011$$

$$\begin{array}{r} \underline{11001} \\ 10100 \\ \underline{11001} \\ 11011 \\ \underline{11001} \\ 10110 \\ \underline{11001} \\ 11111 \\ \underline{11001} \\ 00110 \end{array}$$

Situația 10111010101 în loc de 10100011101 produce un rest nul, asadar erorile rămân nedetectate. Dacă cele trei erori sunt adiacente, 10011011101, restul nenul indică prezența erorilor.

Implementarea unui circuit de împărțire ține seama de faptul că diviziunea poate fi făcută prin multiplicare în buclă feedback.

Exemplu: Fie polinomul $E(X)$ cuvântul recepționat; fie $G(X)$ polinomul generator; fie $D(X)$ cuvântul original informational. Dacă nu există eroare, se recepționează $E(X)$ și se calculează $D(X)$ prin împărțirea $E(X)/G(X)$ și restul este polinomul nul. De exemplu

$$E(X) = D(X)G(X) = D(X)(1 + X^3 + X^4) = D(X) + D(X)(X^3 + X^4)$$

și

$$D(X) = E(X) + D(X)(X^3 + X^4)$$

deoarece în aritmetica modulo 2 scăderea coincide cu adunarea.

Circuitul divizor folosește și el registre de deplasare și porți XOR. Apare în plus un circuit feedback pentru diviziunea $D(X) = E(X) + D(X)(X^3 + X^4)$. La pornire,

fiecare element de întârziere retine un 0. Circuitul produce mai întâi câtu de sapte biti (bitii de informatie) si apoi cei patru biti ai restului. Dacă restul este nenul, este prezentă cel puțin o eroare.

Operatiile divizorului sunt prezentate în tabelul care urmează; i_3 este intrarea pentru elementul de întârziere $O_3 = i_4 \oplus O_4$.

Tempo	Intrare	i_4	O_4	i_3	O_3, O_2, O_1	Iesire
1	1	1	0	1	000	1
2	0	0	1	1	100	0
3	1	1	0	1	110	1
4	1	0	1	1	111	0
5	1	0	0	0	111	0
6	0	1	0	1	011	1
7	0	1	1	0	101	1
8	0	0	1	1	010	0
9	1	0	0	0	101	0
10	0	0	0	0	010	0
11	1	0	0	0	001	0

Dacă intrarea este 10100011101, iesire este 1100101 si restul este nul. Orice eroare singulară în secventa $E(X)$ receptionată produce un rest nenul.

Într-un cod ciclic (n, k) , numărul de biti de control $n - k$ este independent de numărul k al bitilor de informatie; k poate fi crescut fără a crește si $n - k$, ceea ce lasă gradul polinomului generator neschimbat si complexitatea circuitelor de codare si de decodare rămâne la fel.

În ceea ce priveste detectarea erorilor în secventă (bursty errors), pe măsură ce k crește, abilitatea codului ciclic de a detecta erorile succesive se reduce; numai erorile în secventă de lungime până la $n - k$ sunt detectate garantat.

În multe aplicatii apare necesitatea de a face sigură detectarea tuturor erorilor succesive pe lungime de până la 16 biti. Pentru aceasta se utilizează codurile ciclice de tipul $(16 + k, k)$. Cele mai frecvent utilizate:

CRC-16 (Cyclic Redundancy Code pe 16 biti), cu polinomul generator

$$G(X) = X^{16} + X^{15} + X^2 + 1$$

si CRC-CCITT cu polinomul generator

$$G(X) = X^{16} + X^{12} + X^5 + 1$$

Coduri aritmetice

Codurile aritmetice sunt coduri care se conservă printr-un anumit set de operatii aritmetice. Această proprietate permite detectarea erorilor care pot apărea în timpul executării unei operatii aritmetice din setul de operatii definit. O astfel de detectare concurentă a erorilor poate fi totdeauna atinsă prin duplicarea procesorului aritmetic. Această din urmă metodă este totdeauna costisitoare.

Relativ la conservarea codului, se spune că un cod este conservat printr-o operație aritmetică * dacă pentru orice pereche de operanzi (X, Y) și pentru perechea de entități codate corespunzătoare (X', Y') există o operație ⊗ pentru operanzii codati astfel încât $X' \otimes Y' = (X * Y)'$. Altfel spus, operația aritmetică executată când ⊗ este aplicat operanzilor codati X' și Y' produce același rezultat ca cel obținut prin codarea rezultatului operației originare * aplicat operanzilor X și Y înainte de codare. Rezultatul operației aritmetice este codat în același mod.

Cum se face detectarea erorii? Este de așteptat ca un cod aritmetic să fie capabil a detecta toate erorile de un singur bit. O eroare pe un singur bit într-un operand sau într-un rezultat intermediar poate produce o eroare pe mai mulți biti în rezultatul final. De exemplu, la adunarea a două numere în binar, dacă secțiunea i a sumatorului este defectă, toți ceilalți n - i digiti de ordin mai înalt pot fi eronați.

Codurile aritmetice *neseperabile* sunt cele mai simple. Acestea se vor identifica în continuare sub denumirea de codurile AN. Sunt formate prin multiplicarea operanzilor cu o constantă A, $X' = A * X$ și operațiile * și ⊗ coincid: ambele sunt multiplicări.

De exemplu, dacă A = 3, fiecare operand este multiplicat cu 3. Rezultatul fiecărei operații aritmetice este verificat dacă este un întreg multiplu de 3. Toate erorile care au mărimea multiplu de A nu vor putea fi detectate.

În codarea AN, constanta A nu trebuie să fie o putere a lui 2. Un A impar este cea mai bună alegere; codul va detecta orice eroare de un bit. O astfel de eroare are o mărime de 2^i .

Dacă A = 3, se obține codul AN cel mai ieftin, care face posibilă detectarea tuturor erorilor de un bit.

Exemplu: Numărul $0110_2 = 6_{10}$. Reprezentarea în cod AN cu A = 3 este $010010_2 = 18_{10}$. O eroare în bitul cu poziția 2^3 poate da rezultatul eronat $011010_2 = 26_{10}$. Eroarea este detectabilă deoarece 26 nu este un multiplu de 3.

Există și coduri aritmetice *seperabile*. Cele mai simple sunt *codul rezidual* și *codul rezidual inversat*. În fiecare din acestea se atașează un simbol de control separat C(X) pentru fiecare operand X. Pentru codul rezidual, $C(X) = X \bmod A = |X|_A$. A este numit modulul de verificare.

Pentru ambele coduri separabile $C(X) \otimes C(Y) = C(X * Y)$ cu operațiile * sau ⊗ reprezentând fie adunare, fie multiplicare

$$|X + Y|_A = ||X|_A + |Y|_A|_A$$

$$|X * Y|_A = ||X|_A * |Y|_A|_A$$

În împărțirea $X - S = Q * D$, X este împărțitul, D este împărțitorul, Q este câtul, S este restul și verificarea este $||X|_A - |S|_A|_A = ||Q|_A * |D|_A|_A$.

Exemple: Dacă A = 3, X = 7 și Y = 5, resturile sunt $|X|_3 = 1$ și $|Y|_3 = 2$ și

$$|7 + 5|_3 = 0 = ||7|_3 + |5|_3|_3 = |1 + 2|_3 = 0$$

$$|7 * 5|_3 = 2 = ||7|_3 * |5|_3|_3 = |1 * 2|_3 = 2.$$

Dacă A = 3, X = 7 și D = 5, atunci Q = 1 și S = 2; verificarea prin rest este

$$||7|_3 - |2|_3|_3 = ||5|_3 * |1|_3|_3 = 2$$

$$\begin{array}{r}
+ \quad \underline{101} \quad z_1 \\
\quad \quad 111 \\
+ \quad \underline{011} \quad z_0 \\
1 \quad 010 \\
\hline
\quad \quad 1 \quad \text{end-around-carry} \\
+ \quad 011
\end{array}$$

Coduri aritmetice pentru operanzi cu semn

Codul trebuie să fie complementabil în raport cu R . $R = 2^n$ (complement față de 2) sau $R = 2^n - 1$ (complement față de 1), n este numărul de biti ai operandului codat.

Pentru un cod AN, $R - AX$ trebuie să se dividă cu A , adică A trebuie să fie factor al lui R . Dacă e imperativ necesar ca A să fie impar atunci $R = 2^n$ se exclude. Pentru A impar, se poate utiliza numai complementul față de 1; A trebuie să fie un factor (divizor) al lui $2^n - 1$.

Exemplu: $n = 4$, $R = 2^n - 1 = 15$ pentru complementul față de 1 și divizibil cu A pentru codul AN cu $A = 3$.

$X = 0110$ este reprezentat prin $3X = 010010$. Complementul la 1 este $101101 = 45_{10}$ care se divide cu 3.

Complementul la 2 al lui $3X$ este $101110 = 46_{10}$, care nu se divide cu 3.

Dacă $n = 5$ și se ia complementul la 1 atunci $R = 31$ care nu este divizibil cu $A = 3$. $X = 00110$ este reprezentat prin $3X = 0010010$ al cărui complement la 1 este $1101101 = 109_{10}$ – care nu este divizibil cu 3.

În cazul unui cod rezidual cu operanzi cu semn trebuie să fie satisfăcută relația $A - |X|_A = |R - X|_A$. R trebuie să fie un întreg multiplu al lui A , care, din nou, să permită numai aritmetică pe complementul unității. Prin modificarea procedurii astfel încât și complementul la 2 (cu $R = 2^n$) să poată fi utilizat, se obține

- $|2^n - X|_A = |2^n - 1 - X + 1|_A = |2^n - 1 - X|_A + |1|_A$
- Pentru formarea complementului la 2, este necesară adăugarea la codul rezidual a unui termen de corectie $|1|_A$
- A trebuie să fie un factor al lui $2^n - 1$.

Exemplu de cod rezidual: $A = 7$, $n = 6$, $R = 2^6 = 64$ pentru complementul la 2. $R - 1 = 63$ care este divizibil cu 7. $001010_2 = 10_{10}$ și are reziduul 3 modulo 7. Complementul la 2 al lui 001010 este 110110 . Complementul lui $|3|_7$ este $|4|_7$ și prin adunarea termenului de corectie $|1|_7$ se obține 5, reziduul corect modulo 7 al lui $110110 = 54_{10}$.

Corectia similară este necesară și când se adună operanzii prin complement la 2. Un transfer (carry-out) de pondere 2^n se generează și se pierde. Pentru a-l compensa, se scade $|2^n|_A$ la verificarea rezidualului. Deoarece A este un factor al lui $2^n - 1$, $|2^n|_A = |1|_A$.

Există o interdependență între unitățile principale și de verificare. În adunarea de complemente ale lui 2, se generează un transfer la care se renunță.

$$110110 = \mathbf{x} \qquad 101 = |\mathbf{x}|_7$$

+ 001101 = Y	+ 110 = Y ₇
1 000011	1 011
	1 transfer la coadă (end-around-carry)
	100
	- 1 termen de corectie
	011

Aceasta rezultă într-o interdependentă între unitățile principale și de verificare. O eroare într-o unitate principală se poate propaga la unitatea de verificare și efectul erorii este mascat. O eroare pe un singur bit este totdeauna detectabilă.

Coduri bi-reziduale

Corectarea erorii poate fi obținută prin utilizarea a două sau mai multe verificări de reziduu. Cazul cel mai simplu este cel al codului bi-rezidual.

Codarea bi-reziduală constă în două verificări de reziduu față de A_1 și față de A_2 . Cele două numere sunt $A_1 = 2^a - 1$ și $A_2 = 2^b - 1$ și sunt două verificări de reziduu de cost redus cu $n = \text{c.m.m.c.}(a, b)$ cu n numărul de biti în operanzi. Orice eroare pe un singur bit poate fi corectată.

REDUNDANȚE TEMPORALE

Conceptul de bază din redundanța temporală este repetarea calculului de două sau de mai multe ori și compararea rezultatelor pentru a observa existența unor anumite discrepanțe. Dacă se detectează o eroare, calculele pot fi executate din nou pentru a vedea dacă discordanța rămâne sau dispăre. Aceste tratări sunt bune pentru detectarea erorilor datorate unor defecte tranzitorii dar nu oferă protecție la erorile care rezultă din defecte permanente.

O altă formă de redundanță temporală, de data aceasta pentru manipularea defectelor permanente modifică maniera în care se execută calculele a doua oară. O soluție utilizează o logică alternativă pentru circuite combinatoriale auto-duale care execută o funcție pe un set de intrări la un tempo și execută aceeași funcție pe intrări complementate într-un tempo următor: ieșirea a doua ar trebui să fie complementul rezultatului primar. Dacă valoarea secundară a funcției nu este complementul primei valori înseamnă că s-a detectat o eroare.

Tratarea următoare utilizează recalculul cu operanzi deplasati (*shifted*) și este aplicabilă pe organizări hardware detaliate la biti (*bit-sliced*). În primul pas se execută calculul normal pe operanzi și rezultatul este stocat într-un registru. În pasul următor operanzii sunt deplasati la stânga cu k biti, ieșirea este deplasată la dreapta cu k biti și rezultatul se compară cu cel din calculul precedent. Orice eroare de natură aritmetică sau logică în secțiuni (*slices*) de $k - 1$ biti este astfel detectată. Necesarul de hardware aditional constă în trei unități de deplasare, registrul de stocare pentru rezultatul primului calcul și comparatorul.

O variantă a acestei metode face o recalculare cu operanzi fragmentati și cu fragmente inversate (*swapped*): operația este executată în doi pași, mai întâi în forma normală, apoi jumătățile superioare și inferioare ca magnitudine ale operanzilor sunt inversate (*swapped*) astfel că o secțiune (*slice*) de biti defectă operează pe jumătăți diferite ale operanzilor în cele două calcule. Cerințele de hardware suplimentar iau forma câtorva multiplexoare, un registru de memorare și un comparator.

Temporizatoare (timere) *watchdog*

Temporizatoarele *watchdog* au fost utilizate încă de timpuriu în sistemele digitale ca un mod ieftin de detectare a erorilor. În principiu, separat de procesul monitorizat, este implementat un temporizator. Procesul supravegheat trebuie să reseteze timerul înainte ca acesta să expire; altminteri procesul este considerat afectat de defecte.

Traditional, timerele *watchdog* sunt utilizate pentru a detecta erorile în fluxul de control, care rezultă din absența resetării la timp a timerului. Când un timer

expiră, sistemul este resetat. Alternativ, în locul resetării sistemului, se poate genera o întrerupere pentru a iniția recuperarea din eroare. Timerele watchdog pot de asemenea să fie utilizate în cam aceeași manieră cum *timeout*-urile sunt utilizate pentru monitorizarea comportamentului unui singur subsistem. Timeout-urile diferă de timerele watchdog prin aceea că ele asigură o verificare mai fină a fluxului de control.

Timerele de control pot fi implementate fie hardware (timerul este în general ceva extern care poate fi resetat cu un semnal) fie software (rulat adesea pe același procesor pe care este rulat și procesul monitorizat dar timerul este întreținut ca un proces separat).

O implementare nouă a efectului unui timer watchdog fără a utiliza un timer propriu-zis este tehnica de verificare a fluxului de control bazată pe suma de control în dublă precizie. Sumele de control în dublă precizie sunt luate dintr-un bloc de instrucțiuni fără ramificații ca suma totală a instrucțiunilor sau a unei transformate a instrucțiunilor. Înainte de fiecare bloc, valoarea sumei de control este trimisă unui buffer. Pe măsură ce instrucțiunile se execută ele sunt scăzute din buffer. Când blocul ajunge la final sau când apare o ramificație, se trimite un semnal de verificare nul. Dacă bufferul devine zero sau negativ înainte de setarea semnalului, este un semn că a apărut o eroare în fluxul de control. Dacă bufferul este pozitiv când semnalul este setat, și atunci este un semn de eroare.

Exemple de aplicații ale timerelor watchdog

Pluribus Reliable Multiprocessor. Un exemplu de sistem proiectat cu utilizarea extensivă a timerelor watchdog este multiprocesorul Pluribus. Pluribus a fost confectionat mai întâi pentru scopuri de cercetare; obiectivul principal – fiabilitate înaltă. Comportarea lui Pluribus ca un întreg nu este monitorizată dar timere hardware și software monitorizează aproape fiecare subsistem. Această tratare crește fiabilitatea generală a sistemului deoarece un subsistem care esuează din cauza unui defect intermitent sau tranzitoriu este repornit și nu este lăsat să producă căderea sistemului. Dacă Pluribus utilizează alte tehnici de detectare a erorilor, acele tehnici sunt uzual combinate cu un timer. Timerele acoperă durate de la cinci microsecunde la două minute.

Subsistemele Pluribus funcționează ciclic cu o constantă de timp caracteristică. În timpul fiecărui ciclu subsistemul execută un auto-test complet de consistență. Traversarea întregului ciclu indică faptul că subsistemul operează corect; o scurgere de timp îndelungată fără ca un timer să fie resetat indică dimpotrivă faptul că subsistemul a suferit o cădere din care nu se poate recupera prin acțiune proprie. Un exemplu de subsistem controlat astfel este lista de buffere de mesaj libere, unde sunt stocate bufferele de mesaj când nu sunt în uz. Bufferele părăsesc lista pentru cel mult două minute astfel că un timer de 2 minute este asociat cu bufferul. Dacă un timer depășește 2 minute, faptul indică o disfuncție. Bufferul nu revine în lista de buffere libere prin acțiune proprie și bufferul monitorizat este forțat să revină în listă. În acest caz, disfuncția produce o degradare a performanței de mai puțin de două minute, în timpul

căreia sistemul operează cu mai puține buffere de mesaj. Cu toate acestea, eroarea apărută nu produce vreo pierdere de date deoarece timerul facilitează o recuperare completă.

Un alt exemplu este disfuncția unor blocaje (*locks*) de excludere mutuală care sunt prevăzute pe fiecare subsistem. Un lacăt disfuncțional poate menține blocată o resursă deși nici un subsistem nu o utilizează. Un subsistem care încearcă să utilizeze acea resursă trece printr-o stare de așteptare. Deoarece lacătul este disfuncțional, resursa nu va (re)deveni niciodată liberă dar un timer de 1/15 secunde întrerupe procesul și descurcă forțat resursa. Cu toată degradarea temporală (1/15 dintr-o secundă) a performanței sistemului, el nu este afectat de eroare. Ca și în primul exemplu, nu se pierd date și este posibilă recuperarea completă.

O eroare mai serioasă, din care timerele watchdog fac o recuperare parțială, este disfuncția permanentă a unui procesor. Dacă un procesor cade, orice buffer de mesaje care nu a fost returnat listei de buffere libere va fi returnat prin operația executată de timerul care monitorizează lista, cum s-a descris mai devreme. Ca și atunci, orice resursă încuiată de procesor este descuiată, deblocată. Deși recuperarea completă nu este imediat posibilă, deoarece procesorul însuși trebuie reparat sau substituit, sistemul rămâne în funcție și eroarea poate fi limitată la acel unic procesor.

VAX-11/780. Un sistem multiprocesor proiectat pentru o varietate de aplicații comerciale, care face uz de timere watchdog este VAX-11/780. Pe acest sistem procesorul consolă rulează un proces watchdog care este resetat când este explorată (*strobed*⁷) o linie de întrerupere. Dacă această explorare făcută de un procesor nu se încheie în 200 de microsecunde, faptul indică o disfuncție și procesorul de consolă încearcă a determina motivul acestei disfuncții.

Centralele Bell System Telephone. Încă un sistem care utilizează procesoare watchdog pentru a detecta erori este sistemul de comutare programată stocat pe telefon de Bell Systems. Timere watchdog externe monitorizează operarea adecvată a programului prin amorțirea recuperării când timerele nu sunt setate periodic. Aceasta permite detectarea preventivă (înainte ca eroarea să se propage și să producă pagube severe sistemului) a problemelor cauzate de erorile software și, în consecință, cu recuperare mai ușoară. Trebuie observat că în detrimentul detecției erorilor bazate pe principiul watchdog, tehnicile de auditare software cunoscute sunt linia principală de apărare în fața erorilor.

Sojournorul martian. Un exemplu de caz în care un timer watchdog a demonstrat abilitatea sa de a detecta erori este legat de misiunea NASA *Mars Pathfinder* care utilizează vehiculul Sojourner. Sistemul de calcul care controlează vehiculul utilizează un sistem de operare în timp real pe mai multe fire (*multithreaded*) cu preîntâmpinare (*preemptive*). Task-urile sunt planificate pe bază de priorități care reflectă urgența lor relativă. Datorită unei erori de proiectare, poate apărea o condiție cunoscută ca inversare de prioritate. Pentru a

⁷ Cuvântul *strobe* (*strobed*) apare în limbajul celor care se ocupă de calculatoare în două prescurtări: RAS (**R**ow **A**ddress **S**trobe) și CAS (**C**olumn **A**ddress **S**trobe). În ambele cazuri este vorba de un semnal de ceas într-un cip de memorie, utilizat pentru a localiza cu precizie linia sau coloana unui bit particular într-o matrice organizată în linii și coloane.

ilustra inversarea de prioritate se consideră scenariul din exemplu următor: (1) un fir de prioritate scăzută obține un blocaj (*lock*) mutual exclusiv al accesului la date în indiviziune, partajate, (2) în aceste condiții un task îndelungat cu prioritate mai înaltă decât a firului de prioritate joasă este planificat datorită unei întreruperi și (3) firul cu prioritate superioară are nevoie de acces la datele blocate de task-ul de prioritate inferioară. Ca rezultat (1) task-ul de prioritate inferioară este oprit de la execuție de firul de prioritate mai înaltă și (2) task-ul de prioritate mai înaltă este și el oprit de la execuție deoarece el se blochează așteptând ca firul de prioritate inferioară să deschidă lacătul (*lock*).

Utilizarea timerului watchdog face ca scenariul de mai sus să fie detectat și sistemul să fie repornit. Cu toate acestea, restartul complet produce pierdere de date și resetările repetate limitează serios lucrul corect al sistemului de deplasare pe Marte. Problema a fost ulterior detectată și software-ul a fost remediat (*patched*) pentru a restabili comportarea potrivită.

În acest sistem metoda de recuperare s-a aplicat când timeout-ul de watchdog este un reset de sistem traditional, o măsură drastică dar robustă care reprezintă o practică inginerescă bună. Disponibilitatea sistemului este mult mai importantă decât pierderea de date datorată resetării sistemului.

Limitările timerelor watchdog

Timerele watchdog nu sunt perfecte în ceea ce privește detectarea erorilor în sistemele numerice. Motivele sunt patru:

1. Deși detectia erorilor nu este limitată la vreun model al defectelor precizat, timerele watchdog detectează numai erorile de un tip foarte special. Presupunerea este că orice eroare se va manifesta ca eroare pe fluxul de control astfel încât sistemul nu continuă a reseta timerul. Dacă apare o eroare pe fluxul de control dar programul resetează timerul la timp, eroarea va trece nedetectată.
2. Resetările timerului trebuie plasate cu grijă pentru a fi eficiente. Ele nu pot fi plasate în interiorul rutinelor de întrerupere sau al buclelor (pentru a evita buclele infinite), dar pot apărea suficient de frecvent astfel ca timerul să nu poată expira pe durata unei operații normale.
3. Pot fi verificate numai procesele cu durate de execuție relativ deterministe deoarece detectia erorilor se bazează pe de-a-ntregul pe timpul dintre resetările timerului. Dacă timpul de setare este mai scurt decât durata de execuție cea mai lungă posibilă a procesului verificat, el poate expira chiar dacă nu există vreo eroare. Pe de altă parte, dacă timpul de setare este prea lung, atunci chiar dacă apare o eroare pe fluxul de control, procesul poate avea timp suficient de a reveni la punctul la care timerul este resetat și eroarea nu va fi detectată.
4. Un timer watchdog furnizează numai o indicație de disfuncție posibilă a procesului; un proces esuat parțial poate el însuși să reseteze timerul. Acoperirea este limitată, deoarece nu sunt verificate nici datele nici rezultatele. Când este utilizat pentru a reseta sistemul, un timer watchdog

poate îmbunătăți disponibilitatea (timpul mediu de recuperare se scurtează) dar nu fiabilitatea (aparitia disfuncțiilor este la fel de probabilă și după resetare). Când disponibilitatea unui sistem digital este mai importantă decât pierderea de date în anumite condiții, utilizarea unui timer watchdog pentru resetarea sistemului la detectarea unei erori este o alegere potrivită.

Heartbeats

Heartbeats reprezintă o tratare comună a detectării disfuncțiilor din procese și noduri în ambianța calculului distribuit (rețelizat). Periodic, o entitate monitorizantă trimite un mesaj (*a heartbeat* – o bătaie de inimă) la un nod sau la un proces monitorizat și așteaptă un răspuns. Dacă nodul monitorizat nu răspunde într-un interval de timp predefinit, este declarat disfuncțiv și se inițiază acțiunea potrivită de recuperare.

Limitări ale heartbeat-ului traditional

Sunt două probleme majore asociate cu schema tradițională de heartbeat:

- Perioada de timeout este prenegociată de cele două părți sau uneori chiar *hard-coded*⁸ de programator. Valoarea de timeout predefinită nu se poate adapta la schimbările de trafic în rețea sau la variabilitatea încărcării individuale pe noduri. În cazul traficului intens, al încărcării mari a nodurilor sau al prezentei unui nod lent, valoarea de timeout poate fi prea scurtă și face ca monitorizarea nodurilor să declare un nod sănătos ca fiind defect. O astfel de alarmă falsă este de nedorit într-un mediu distribuit, în special pentru aplicații critice cum sunt cele din banking-ul comercial și sistemele de baze de date.
- Nodul monitorizat este presupus a fi sănătos dacă este capabil a răspunde la mesajul heartbeat. Acest fapt este acceptabil uzual pentru aplicațiile cu un singur fir (*single-threaded*). Cu toate acestea, într-o aplicație multifir, un fir independent al execuției este răspunzător de răspunsul la mesajul heartbeat. Operarea sănătoasă a acestui fir nu implică în mod necesar operarea sănătoasă a întregii aplicații multithread. Alte fire din proces pot fi în situație de blocaj, situație care previne progresul, avansul întregii aplicații, altfel spus, alte fire ar putea opera într-o stare coruptă care interzice procesului producerea unui serviciu corect.

Algoritmii heartbeat *adaptivi* și *inteligenti* rezolvă aceste două probleme. Un algoritm heartbeat este numit *adaptiv* dacă valoarea de timeout utilizată de monitor nu este fixată ci este negociată periodic între cele două părți pentru adaptare la schimbările traficului din rețea sau la variațiile de încărcare a nodului. Un algoritm heartbeat este denumit *inteligent* dacă entitatea care este

⁸ **hard-coded** – în jargonul programării înseamnă valori de date sau comportamente (pre)scrise direct într-un program, posibil în mai multe locuri, unde nu pot fi modificate prea ușor. Sunt mai multe posibilități, depinzând de cât de des valoarea s-ar putea schimba.

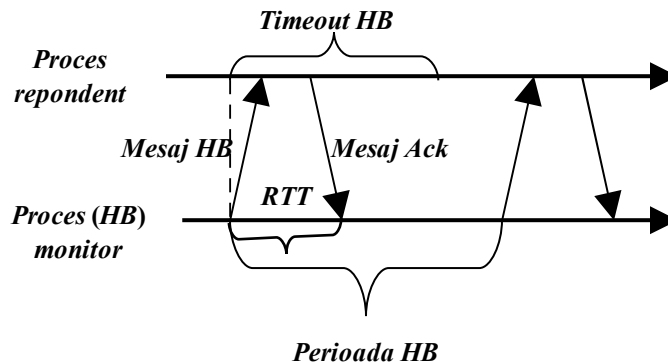
monitorizată excită un set de verificări predefinite pentru a testa robustețea întregului proces și numai apoi răspunde procesului de monitorizare.

Proiectarea heartbeat-urilor adaptive, inteligente

Pentru a ilustra conceptele de *heartbeat adaptiv* și *heartbeat inteligent*, sunt create două procese multifir independente, unul care răspunde la mesajul heartbeat, altul heartbeat-ul care monitorizează.

Heartbeat-ul monitor este entitatea care monitorizează fiind răspunzătoare de trimiterea periodică de mesaje cerere de heartbeat către nodul țintă. *Heartbeat-ul care răspunde* este entitatea monitorizată și răspunde la mesajele cerere de heartbeat trimise de monitor. Schema *adaptivă* utilizează algoritmul Jacobson care permite ajustarea valorii de timeout potrivit performanței măsurate a rețelei în termeni de timp dus-întors (RTT – round trip time) la transmiterea mesajului. Algoritmul heartbeat este făcut *inteligent* (adică capabil a verifica robustețea întregului proces) prin utilizarea unui mesaj de test *null* în interiorul procesului pentru a testa operarea sănătoasă a tuturor firelor din proces. În paragrafele următoare se prezintă implementări ale acestor două scheme.

Figura următoare reprezintă protocolul heartbeat. Periodic, monitorul heartbeat trimite un mesaj către procesul care răspunde la heartbeat, sterge contorul *ack_missed* și porneste timerul. Durata timerului este dictată de valoarea curentă a variabilei timeout asociată cu procesul care răspunde la heartbeat.



Protocolul heartbeat (HB) adaptiv

Pe de altă parte, procesul care răspunde la heartbeat răspunde cu un mesaj de confirmare a recepției heartbeat-ului. Dacă acest mesaj de confirmare (*acknowledge*) este primit de monitorul heartbeat înainte de expirarea timpului potrivit, monitorul presupune că procesul îndepărtat este viu, altminteri contorul *ack_missed* este incrementat. Dacă contorul nu a atins valoarea sa maximă, poate fi trimis încă un mesaj heartbeat de la monitor la respondent, altminteri procesul monitorizat este presupus a fi defect.

Cruciale pentru acest protocol sunt valorile de timeout si perioada de heartbeat. În general perioada poate fi fixată ca un multiplu al valorii curente de timeout. Este totusi de dorit a avea o valoare de timeout care se adaptează la timpul curent de răspuns al procesului respondent. Timpul de răspuns, cum este văzut de monitor este o functie de încărcarea curentă pe masina de departe si de timpul necesar pentru a transmite mesajul si răspunsul de confirmare a heartbeat-ului, adică timpul de răspuns este o functie de RTT.

Pentru a calcula RTT este suficient a include în mesajul heartbeat o marcă de timp a cărei valoare este timpul de trimitere. Această marcă de timp (*timestamp*) va fi trimisă înapoi la monitor de respondent; astfel, când monitorul primește o confirmare de heartbeat, poate calcula instantaneu RTT ca diferenta dintre timpul curent si acea marcă de timp. Totusi, se dovedeste că o astfel de solutie nu face față cazului unei încărcări variabile. Problema principală este cea a variabilității valorii instantanee a RTT, care poate prezenta oscilatii substantiale. Este necesară o estimare a RTT mai precisă si mai netedă.

O abordare a acestei probleme se găseste în algoritmul lui Jacobson (care s-a dovedit eficace în cazul încărcărilor variabile si este implementat curent în protocolul TCP). Pentru fiecare respondent la heartbeat, monitorul heartbeat mentine un RTT variabil care este cea mai bună estimare curentă a timpului dus-întors la destinatia în discutie. Când este trimis un mesaj heartbeat, se porneste un timer. Dacă confirmarea revine înainte de expirarea timpului, monitorul măsoară cât timp consumă sosirea confirmării; fie M această valoare. Apoi se actualizează RTT potrivit formulei

$$RTT = \alpha RTT + (1 - \alpha)M$$

cu α un factor de netezire care exprimă ce pondere se acordă valorii vechi. Tipic, α este fixat la $7/8$.

Chiar fiind dată o valoare bună pentru RTT, alegerea unei valori de timeout nu este tocmai simplă. O alegere ar putea fi $\beta \cdot RTT$, dar este greu de ales β . Mai mult, experienta a arătat că o valoare constantă pentru β este inflexibilă deoarece ea nu reuseste să acopere dispersii mari ale valorilor RTT. Jacobson propune a lua β proportional cu deviatia standard a functiei densitate de probabilitate a timpului de sosire a confirmării. În acest mod, o dispersie mare face pe β mare si invers. În particular, el sugerează utilizarea deviatiei medii care este o estimare ieftină a deviatiei standard. Algoritmul lui cere retinerea secventială a unei alte variabile netezite, D , deviatia. Ori de câte ori soseste o confirmare heartbeat, se calculează diferenta între valoarea asteptată si cea observată, $|RTT - M|$. O valoare netezită a lui D este mentinută după formula

$$D = \alpha D + (1 - \alpha)|RTT - M|$$

cu α acelasi sau diferit de cel din expresia de netezire a RTT. Cu toate că deviatia D nu este exact egală cu deviatia standard, aceasta este destul de bună si este posibil a calcula D si RTT în mod foarte eficace utilizând numai adunări, scăderi si deplasări de întregi. Valoarea initială a timeout-ului este calculată după cum urmează:

$$RTT = (7/8)RTT + (1/8)M$$

$$D = (3/4)D + (1/4)|RTT - M|$$

Utilizând aceste două numere, valoarea de timeout este calculată ca $RTT + 4D$. Mai mult, dacă timerul expiră înainte ca monitorul heartbeat să primească confirmarea de heartbeat, valoarea de timeout este dublată în maniera retractilă a pasului înapoi (*back-off*).

Implementarea si evaluarea de heartbeat inteligent

Schema heartbeat-ului inteligent este implementată în cadrul procesului heartbeat respondent. La primirea unui mesaj cerere de heartbeat de la monitorul de heartbeat, firul respondent initiază o rundă de mesaje surogat (MSG_HB_NULL_TEST) către toate firele din proces. Pentru a procesa un mesaj, un anumit fir obtine mai întâi un *lock* pentru structura de date locală. Apoi incrementează un contor în mesaj si relaxează *lock*-ul elementului structură de date. Dacă procesul operează corect, adică nu există un blocaj (*deadlock*), contorul din mesajul *null* va atinge ulterior numărul total de fire curent active în proces. Ultimul fir din lant care vede această conditie satisfăcută trimite un MSG_HB_NULL_TEST_REPLY înapoi la respondentul de heartbeat (această transmisie este executată în cadrul aceluiasi proces). La primirea acelu MSG_HB_NULL_TEST_REPLY, elementul respondent presupune că întregul proces este într-o stare de sănătate liberă de blocaje si trimite apoi către monitor un mesaj de confirmare de heartbeat. Dacă, totusi, există fire în proces în stare blocată, mesajul de test *null* este blocat si respondentul nu trimite la monitor un mesaj de confirmare. Ulterior, monitorul intră în timeout si declară disfuncț procesul monitorizat.

Checkpointarea

Disfunctii la executarea programelor

Azi calculatoarele sunt mult mai rapide, iar aplicatiile sunt din ce în ce mai complicate. Printre aplicatiile care sunt încă mari consumatoare de timp se pot mentiona:

- Actualizarea bazelor de date.
- Simulările curgerii fluidelor, necesare în particular pentru modelarea meteo si climatică.
- Optimizarea, alocarea optimă a resurselor economice (de pildă în utliarea liniilor aeriene).
- În astronomie: simularea în problema celor n corpuri si modelarea universului.
- În biochimie: studierea compusilor proteici.

Când timpul de executie a calculelor este foarte lung, atât probabilitatea de a claca în timpul executiei cât si costurile unui astfel de eveniment devin semnificative.

Un exemplu relativ la costul întreruperii executiei unui program:

Rularea completă a programului consumă T ore. Sistemul suferă disfuncții tranzitorii la o rată de λ disfuncții pe oră. Disfuncția este instantanee dar tot efortul de calcul anterior este pierdut.

E este timpul total de execuție așteptat statistic (*expected*), timp care include și orice efort de calcul pierdut din cauza disfuncțiilor apărute.

Dacă nu sunt întreruperi, cu alte cuvinte dacă nu apar disfuncții în timpul executiei (un caz cu probabilitatea $e^{-\lambda T}$), timpul total de execuție așteptat (mediu) conditionat este T .

Probabilitatea apariției unei disfuncții la momentul τ al executiei este $\lambda e^{-\lambda \tau} d\tau$.

Într-un asemenea caz, τ ore sunt risipite, programul trebuie reluat și un timp suplimentar adițional E este de așteptat pentru a încheia calculul. Timpul mediu așteptat pentru încheierea calculului este asadar $\tau + E$.

Pentru a calcula costul disfuncției se face o mediere pe toate cazurile:

$$E = T e^{-\lambda T} + \int_0^T (\tau + E) \lambda e^{-\lambda \tau} d\tau = E(1 - e^{-\lambda T}) + (1 - e^{-\lambda T})/\lambda$$

o ecuație în E cu soluția $E = (e^{\lambda T} - 1)/\lambda$.

O măsură a overhead-ului (comparat cu T) este

$$\eta = E/T - 1 = (e^{\lambda T} - 1)/(\lambda T) - 1$$

Overhead-ul relativ η depinde numai de produsul λT care este numărul mediu de disfuncții pe durata executiei programului. Acest η crește foarte rapid (exponential) cu λT .

Este de preferat, desigur, a face în așa fel ca lucrul să nu fie reluat de la început dintr-o stare cu toate disfuncțiile din nou posibile. Soluția este așa-numita *checkpointare* (utilizarea unor puncte de control).

Definiția checkpointării

Un checkpoint este un instantaneu (aproape în sensul fotografic) al întregii stări a procesului de calcul la momentul când acel checkpoint este prelevat: se retine toată informația necesară pentru repornirea procesului din acel punct.

Checkpoint-ul este salvat într-o memorie stabilă, de o fiabilitate suficientă. Cele mai utilizate ca memorie stabilă sunt de obicei discurile. Discurile pot menține date chiar dacă, de pildă, alimentarea cade (cădere care nu produce vătămări fizice ale suprafeței active). În plus, discurile pot reține cantități enorme de date la un cost foarte scăzut. Și checkpoint-urile pot fi foarte cuprinzătoare: zeci sau chiar sute de megabytes.

Adesea este utilizată ca memorie stabilă și RAM cu o baterie de backup.

Nici un mediu nu este perfect fiabil; fiabilitatea trebuie să fie suficient de înaltă pentru aplicația *at hand*.

Overhead-ul si latența unui checkpoint

Overhead-ul checkpointării este creșterea timpului de execuție al aplicației datorată preluării de checkpoint-uri.

Latenta checkpointării este timpul necesar pentru salvarea checkpoint-urilor.

Într-un sistem simplu, overhead-ul si latenta coincid, sunt identice. Dacă parte din checkpointare poate fi suprapusă cu execuția aplicației, latenta poate fi substantial mai mare decât overhead-ul. Dacă un proces face checkpointarea prin transcrierea stării sale într-un buffer intern, unitatea centrală (CPU) poate continua execuția în timp ce checkpointul este transferat din buffer pe disc.

Latenta checkpointării – un exemplu

```
for (i = 0; i < 1000000; i++)
    if (f(i) < min) {min = f(i); imin = i;}

for (i = 0; i < 100; i++) {
    for (j = 0; j < 100; j++) {
        e[i][j] += i*j/min;
    }
}
```

Prima parte a acestui scurt program face calculul celei mai mici valori a unei funcții $f(i)$ pentru $0 < i < 1000000$. Partea a doua execută o multiplicare urmată de o împărțire.

Ce se poate spune despre dimensiunea checkpointului din acest exemplu? Întrebarea este cât se poate de legitimă deoarece dimensiunea checkpointului determină latența aceluși checkpoint. Dimensiunea poate varia de la program la program și chiar în timpul execuției unui aceluși program.

Un checkpoint în prima parte a exemplului dat poate fi redus: numai contorul programului și variabilele min și $imin$; multe alte registre sunt irelevante.

Un checkpoint prelevat în partea a doua trebuie să includă ceea ce s-a calculat din masivul $e[i][j]$ până la momentul respectiv.

În general, dimensiunea checkpointului este dependentă de program și, cum s-a mai spus, poate fi de la câțiva kilobytes la mai multi gigabytes.

Sunt câteva probleme în legătură cu checkpointarea. Iată-le enumerate:

- Câte checkpointuri sunt necesare?
- În ce puncte din execuția programului se cuvine a se preleva un checkpoint?
- Cum se poate reduce overhead-ul de checkpointare?
- Cum se poate face checkpointarea în sisteme distribuite în care un control central poate exista sau poate să nu existe?
- La ce nivel (kernel/user/application) trebuie făcută checkpointarea?
- Cât de transparent față de utilizator trebuie să fie procesul de checkpointare?

Checkpointarea la nivel de kernel

Procedurile de checkpointare sunt incluse în kernel, sunt transparente pentru utilizator și nu necesită schimbări în program.

Când sistemul reporneste după cădere, kernelul este responsabil pentru administrarea operației de recuperare.

Orice sistem de operare preia checkpointuri când un proces este întrerupt de un altul; starea procesului este înregistrată astfel că execuția poate fi reluată din punctul de întrerupere fără pierdere de efort de calcul.

Multe sisteme de operare au puține checkpointuri orientate pe toleranța la defecte sau nu au deloc.

Checkpointarea la nivel de utilizator

În această variantă, biblioteca de nivel utilizator este cea pregătită să execute checkpointarea. Pentru a face checkpointarea, programele aplicației sunt conectate la aceste biblioteci.

Ca și checkpointarea la nivel de kernel, această abordare nu cere în general vreo schimbare în codul aplicației dar legarea explicită cu biblioteca la nivelul utilizatorului este cerută în mod necesar.

Biblioteca la nivel de utilizator gestionează și recuperarea din starea de disfuncție.

Checkpointarea la nivel de aplicație

Aplicația este responsabilă pentru toate funcțiile de checkpointare, codul pentru checkpointare și pentru recuperare este parte din aplicație. Acesta este cel mai important control asupra procesului de checkpointare, dar este costisitor de implementat și de depanat.

Firele (threads) sunt invizibile la nivel de kernel, dar nivelurile de utilizator și de aplicație nu au acces la informația detinută la nivel de kernel; nivelurile de utilizator și de aplicație nu pot atribui un identificator (ID) de proces particular pentru un proces de recuperare.

Nivelurilor de utilizator și de aplicație s-ar putea să nu li se permită checkpointarea unor părți din sistemul de fișiere. Ar putea în schimb să memoreze nume de fișiere și pointeri la fișiere.

Model analitic pentru latență și overhead

Overhead-ul este partea din checkpointare care nu se execută în paralel cu aplicația.

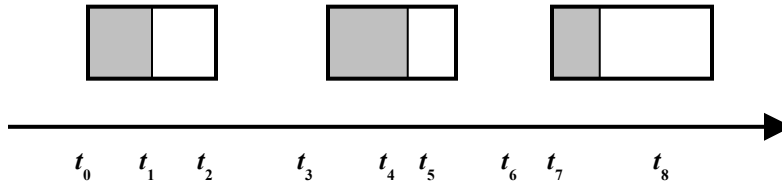
Latența este timpul între t_{start} (când operația de checkpointare începe) și t_{end} (când operația se încheie).

Overhead-ul are impact mai mare asupra performanțelor decât latența.

Checkpointul preia și reprezintă starea sistemului la t_{start} .

Overhead-ul este partea din $[t_{start}, t_{end}]$ pe durata căreia aplicația este oprită din execuție din cauza checkpointării (CPU este ocupată cu checkpointarea).

Notarea overhead-ului cu t_c face intervalul de overhead $[t_{start}, t_{start} + t_c]$.



În figură, dreptunghiurile reprezintă latența, iar partea umbrită reprezintă overhead-ul.

Dacă apare o disfuncție în intervalul $[t_{start}, t_{end}]$, checkpointul preluat este inutil și sistemul trebuie să revină la checkpointul precedent.

Exemplu: dacă disfuncția apare în intervalul $[t_3, t_5]$, se revine la checkpointul precedent, adică la starea procesului de la timpul t_0 .

Dacă t_r exprimă timpul mediu de recuperare, acesta este timpul petrecut în starea cu defect plus timpul de recuperare la o stare funcțională (adică timpul de re-bootare completă a procesorului).

Dacă la momentul τ apare o disfuncție tranzitorie, procesul redevine activ la momentul $\tau + t_r$.

Câteva notații suplimentare necesare pentru a scrie modelul analitic:

I – intervalul între checkpointuri, timpul între executarea checkpointului i și a checkpointului $i + 1$.

$E_{interval}$ – durata medie statistică (expected value) a lui I .

T – timpul cheltuit pentru executarea aplicației pe această perioadă; dacă nu apare vreo disfuncție, $I = T + t_c$.

t_l – latența, egală cu diferența $t_{end} - t_{start}$.

Dacă disfuncția apare la τ ore în intervalul I , lucrul pierdut include:

- Lucrul executat pe durata τ
- Lucrul util executat în timpul $[t_{start}, t_{end}] = t_l - t_c$.

t_r – timpul mediu de recuperare și de reluare a calculului.

Timpul suplimentar total datorat unei disfuncții la τ ore în intervalul I este

$$\tau + t_l - t_c + t_r$$

O aproximare de primul ordin pentru intervalul dintre checkpointuri

Se presupune că nu mai mult de o disfuncție poate lovi sistemul între două checkpointuri, ceea ce este o aproximare bună dacă $T + t_c$ este mic comparativ cu $1/\lambda$, timpul mediu între apariția a două disfuncții.

Timpul mediu între două checkpointuri succesive se evaluează pe baza unei priviri asupra a două cazuri:

Cazul 1: între două checkpointuri succesive nu apare nici o disfuncție. Durata scursă este $T + t_c$, iar probabilitatea cazului este $e^{-\lambda(T+t_c)}$.

Cazul 2: între două checkpointuri succesive apare o disfuncție. Probabilitatea aproximativă este $1 - e^{-\lambda(T+t_c)}$.

Timpul aditional datorat disfuncției: $\tau + t_r + t_l - t_c$.

Media timpului τ este $(T + t_c)/2$.

Timpul mediu (asteptat) aditional este $(T + t_c)/2 + t_r + t_l - t_c$.

La calculul lungimii medii a intervalului dintre checkpointuri avem: contribuția cazului 1:

$$(T + t_c)e^{-\lambda(T+t_c)}$$

și contribuția cazului 2:

$$\begin{aligned} & (1 - e^{-\lambda(T+t_c)})[(T + t_c)/2 + T + t_c + t_r + t_l - t_c] = \\ & = (1 - e^{-\lambda(T+t_c)})[3T/2 + t_c/2 + t_r + t_l] \end{aligned}$$

Prin însumarea celor două contribuții se obține:

$$E_{interval} \approx 3T/2 + t_c/2 + t_r + t_l - (T/2 - t_c/2 + t_r + t_l)e^{-\lambda(T+t_c)}$$

Sensibilitatea acestui $E_{interval}$ la t_l și t_c este dată de derivatele (partiale) în raport cu variabilele respective:

$$\begin{aligned} \partial E_{inter}/\partial t_c & \approx \frac{1}{2} + \left[\frac{1}{2} + \lambda \left(\frac{T - t_c}{2} + t_r + t_l \right) \right] e^{-\lambda(T+t_c)} \\ \partial E_{inter}/\partial t_l & \approx 1 - e^{-\lambda(T+t_c)} \end{aligned}$$

Dacă $\lambda(T + t_c) \ll 1$ atunci derivata în raport cu t_c este mult mai mare decât cea în raport cu t_l . Speranța E este mult mai sensibilă la overhead decât la latență și, de aceea, t_c trebuie menținut scăzut chiar dacă t_l crește.

Reducerea overheadului prin bufferizare

Se poate face ca procesul să scrie checkpointul în memoria principală și apoi să revină la execuția aplicației. Memoria cu acces direct (DMA) este utilizată pentru a copia checkpointul din memoria principală pe disc. DMA reclamă implicarea CPU numai la începutul și la încheierea operației.

Este posibilă o rafinare a acestei posibilități de copiere prin scriere bufferizată. Se are în vedere faptul că copierea unor porțiuni din starea procesului care nu s-au schimbat de la ultimul checkpoint este pierdere de timp. Dacă procesul nu actualizează prea frecvent unele pagini din memoria principală, mult din efortul de copiere într-o zonă buffer poate fi evitat.

Copiere cu scriere bufferizată

Multe sisteme de memorare sunt prevăzute cu biti de protecție pentru memorie (pentru o pagină de memorie fizică principală) care indică fie că pagina este read-write, fie că este read-only, fie că este inaccesibilă.

Când checkpointul este prelevat, bitii de protecție ai paginilor care aparțin procesului sunt setați ca read-only. Aplicația continuă rularea în timp ce

paginile de checkpoint sunt transferate pe disc. Dacă aplicația încearcă să actualizeze pagina, se amorsează o violare a tipului de acces. Sistemul bufferizează atunci pagina și tipul de acces se setează la read-write. Pagina bufferizată este copiată mai târziu pe disc.

Reducerea overhead-ului de checkpointare – excluderea de memorie

Două tipuri de variabile nu trebuie să fie checkpointate: cele care nu au fost actualizate și cele care sunt “moarte”. O variabilă moartă este o variabilă a cărei valoare prezentă nu va mai fi utilizată nicodată în program.

Sunt două genuri de variabile moarte: acelea care nu vor mai fi niciodată adresate de program și acelea pentru care următorul acces este pentru scriere. Provocarea este să identifice corect astfel de variabile.

Identificarea de variabile moarte

Spatiul de adresare al unui proces are patru segmente: codul, datele globale, heap-ul și stack-ul.

Identificarea în cod a variabilelor moarte este facilă: codul cu automodificare nu se mai utilizează, astfel, segmentul de cod din memorie este read-only și nu mai trebuie checkpointat.

Segmentul stack este la fel de accesibil. Conținutul de adrese reținute în locații situate sub pointerul de stack sunt evident moarte (spațiul de adresare virtual are uzual segmentul de stack la vârf, crescător în jos).

Segmentul heap: multe limbaje de programare permit programatorului să aloce și să elibereze explicit memorie (exemple: apelurile *malloc()* și *free()* utilizate în C). Conținutul listei *free* este mort prin definiție.

Unele pachete de checkpointare la nivel de utilizator (de pildă *libckpt*) furnizează programatorului apeluri la proceduri (de pildă *checkpoint_here()*) care specifică regiuni de memorie care ar trebui excluse din sau incluse în checkpointuri viitoare.

Reducerea latenței

Compresia la checkpointare are ca rezultat un volum mai mic de scris pe disc.

Cât de mult se câștigă prin compresie depinde de:

- extinderea compresiei, care este dependentă de aplicație și poate varia între 0 și 50%.
- efortul cerut de algoritmul de comprimare, care este executat de CPU și se adaugă la overhead-ul checkpointării ca și la latență.

În checkpointarea simplă secvențială unde $t_c = t_l$, compresia poate aduce avantaje.

În sistemele mai eficiente, în care $t_c \ll t_l$, utilitatea acestei soluții este discutabilă și trebuie evaluată cu grijă înainte de a fi utilizată.

Amplasarea checkpointurilor

Asezarea în timp (si în program) a checkpointurilor trebuie să aibă în vedere o echilibrare între costuri si beneficii. Se urmărește uzual minimizarea duratei de executie a unei aplicatii de durată.

Costul, considerat a fi timpul pentru stocarea unui checkpoint poate fi mare. Câteva notatii:

t_x – durata executiei fără checkpointare.

t_c – timpul mediu consumat cu prelevarea unui checkpoint.

N – o variabilă de decizie cu semnificatia numărul de checkpointuri plasate uniform în job pentru minimizarea timpului total de executie $T_{tot}(N)$.

$\tau_x = t_x/N$ – timpul (mediu) între două checkpointuri succesive.

Disfunctiile apar cu rata λ . Disfunctiile sunt tranzitorii, ele dispar după o durată de existență t_f .

La manifestarea unei disfunctii, sistemul revine la checkpointul cel mai recent.

Checkpointurile sunt presupuse a se afla într-o memorie sigură, care nu poate fi coruptă de vreo disfunctie.

Model analitic pentru amplasarea checkpointurilor

t_l – timpul total pierdut pentru fiecare disfunctie tranzitorie.

t_f – timpul cât sistemul este disfuncț.

Dacă disfunctia apare pe durata checkpointului, se asociază aparitiei o probabilitatea $P_c = t_c/(t_c + \tau_x)$. Timpul irosit este în acest caz $\tau_x + t_c/2$.

Dacă disfunctia apare pe durata executiei se lucrează cu probabilitatea $P_x = \tau_x/(t_c + \tau_x)$. Timpul irosit este acum $\tau_x/2$.

$$t_l = t_f + P_c(\tau_x + t_c/2) + P_x(\tau_x/2) = t_f + (t_c + \tau_x)/2.$$

Rezultatul obtinut este intuitiv: $(t_c + \tau_x)/2$ este jumătate din intervalul $t_c + \tau_x$.

Amplasarea optimă a checkpointurilor

Se presupune că λ este suficient de mic astfel încât probabilitatea disfunctiei în timpul revenirii (*rollback*) este neglijabilă.

Numărul de disfunctii mediu statistic (*expected*) pe durata timpului total de executie de $t_x + N t_c$ este $\lambda(t_x + N t_c)$.

Timpul total consumat:

$$T_{tot}(N) = t_x + N t_c + \lambda(t_x + N t_c)[t_f + (t_c + t_x/N)/2]$$

Selectarea lui N pentru minimizarea timpului $T_{tot}(N)$

$$dT_{tot}(N)/dN = t_c + \lambda t_c(t_c/2 + t_f) - (\lambda t_x^2)/(2N^2)$$

si prin anularea derivatei se obtine

$$N_{opt} = \frac{t_x \sqrt{\lambda}}{\sqrt{t_c(2 + \lambda t_c + 2\lambda t_f)}}$$

Numărul N_{opt} trebuie să fie un întreg, cel care prin trunchiere sau prin rotunjire în plus să minimizeze $T_{tot}(N)$. Rezultă imediat intervalul optim între două checkpointuri succesive: $\tau_{opt} = t_x/N_{opt}$.

Exercitiu: Relaxati ipoteza că probabilitatea unor disfuncții suplimentare în timpul procesului de recuperare este neglijabilă.

Amplasarea uniformă este optimă dacă checkpointarea are un cost constant pe toată durata executiei.

Dacă dimensiunea (deci durata) checkpointului variază mult de la o parte a executiei la alta, timpul optim dintre checkpointuri nu mai are un pas constant.

Plasarea optimă a checkpointurilor: un model la nivelul instructiunilor

Probabilitatea unei defectiuni pe durata executării unei instructiuni depinde de unitățile functionale utilizate și de timpul de executie al acelei instructiuni.

Se definește variabila de decizie M care reprezintă numărul de instructiuni între două checkpointuri consecutive.

Este vizată minimizarea lui W , timpul consumat cu o instructiune.

Multimea instructiunilor este partitionată în N submultimi de instructiuni similare. Pentru o instructiune de tipul i , timpul de executie este T_i , rata de defectare este λ_i , frecventa este f_i ($\sum f_i = 1$).

Se notează, de asemenea, cu s ($1 - s$) fractia de defecte permanente (tranzitorii).

Constanta μ_i reprezintă rata de “reparare” a disfuncțiilor tranzitorii pentru cazul unei instructiuni de tipul i .

În continuare sunt date alte câteva notatii utilizate în scrierea modelului la nivelul de instructiune.

Evenimentele posibile în timpul executării unei instructiuni sunt:

H^C – instructiunea este executată cu succes când este executată prima dată; probabilitatea evenimentului este p^C .

H^{RB} – instructiunea esuează, disfuncția e identificată, programul este readus (*rolled-back*) la ultimul checkpoint, instructiunea este executată; probabilitatea asociată este p^{RB} .

H^{PF} – revenirea (*roll-back*) esuează și ea, programul esuează, programul este reîncărcat și repornit; probabilitatea este p^{PF} .

$p_i^C, p_i^{RB}, p_i^{PF}$ sunt probabilități conditionate pentru instructiunile de tipul i , trecut ca indice.

Aceste probabilități conditionate se calculează și apoi se mediază

$$p^{(j)} = \sum_{i=1}^N f_i p_i^{(j)} \quad (j = C, RB, PF)$$

Sunt prezentate imediat alte notatii pentru un sistem cu rata de defectare λ și cu rata de reparare μ .

Probabilitatea lipsei defectelor în intervalul de timp $(0, t)$

$$P_0(\lambda, t) = e^{-\lambda t}$$

Probabilitatea tranzitiei de la starea lipsită de defecte la timpul 0 la starea lipsită de defecte de la momentul t

$$P_{00}(\lambda, \mu, t) = \frac{\mu}{\lambda + \mu} + \frac{\lambda}{\lambda + \mu} e^{-(\lambda + \mu)t}$$

Pentru $0 \leq t_1 \leq t$, probabilitatea tranzitiei de la starea lipsită de defecte la momentul 0 la starea lipsită de defecte la momentul t , cu cel puțin o defectare în perioada $(0, t_1)$

$$\bar{P}_{00}(\lambda, \mu, t_1, t) = P_{00}(\lambda, \mu, t) - e^{-\lambda t} P_{00}(\lambda, \mu, t - t_1)$$

Probabilitățile din modelul la nivel de instructiune

$$P_i^{(C)} = P_0(\lambda_i, T_i) = e^{-\lambda_i T_i}$$

M – numărul de instructiuni dintre checkpointuri.

m – numărul de instructiuni dintre instructiunea cu disfuncție și ultimul checkpoint = 1, ..., M , fiecare cu probabilitatea $1/M$.

$P_{i,m}^{RB}$ – probabilitatea condiționată a revenirii izbutite, fiind dat tipul i și numărul m de instructiuni executate de la ultimul checkpoint.

δ_1 – timpul de setup necesar pentru a iniția revenirea programului, inclusiv timpul necesar pentru a încărca informația salvată la ultimul checkpoint.

$$P_{i,m}^{(RB)} = \bar{P}_{00}[(1-s)\lambda_i, \mu_i, T_i, T_i] P_0(s\lambda_i, T_i) P_0(\lambda_i, \delta_1) [P^{(C)}]^{m-1} P_i^{(C)}$$

$$P_i^{(RB)} = \bar{P}_{00}[(1-s)\lambda_i, \mu_i, T_i, T_i] P_0(s\lambda_i, T_i) \frac{1 - [P^{(C)}]^M}{1 - P^{(C)}} \frac{P_i^{(C)}}{M}$$

$$P_i^{(PF)} = 1 - P_i^{(C)} - P_i^{(RB)}$$

Modelul la nivelul instructiunii servește la calculul lui W .

Cu τ timpul mediu pentru a executa cu succes o instructiune și cu T_s timpul consumat de checkpointare se poate scrie expresia timpului cheltuit pe instructiune

$$W = \tau + T_s/M.$$

Prin creșterea lui M , numărul de instructiuni între checkpointuri, se obțin creșteri ale primului termen și descreșteri ale termenului secund.

$T = \sum f_i T_i$ este durata medie de execuție a instructiunii fără defecte.

δ_2 este timpul mediu necesar diagnozei și reparației.

L notează numărul mediu de instructiuni pe program și

$$\tau = \bar{T} + P^{(RB)} \left(\delta_1 + \frac{M+1}{2} \bar{T} \right) + P^{(PF)} \left(\delta_1 + \frac{M+1}{2} \bar{T} + \delta_2 + \frac{L+1}{2} W \right)$$

Acest τ include pe W ca unul din termeni. Dacă τ este substituit în expresia lui W se obține o ecuație care se rezolvă pentru W .

Rezolvarea pentru W produce rezultatul

$$W = \frac{1}{1 - \frac{L+1}{2} P^{(PF)}} \left[\bar{T} + P^{(RB)} \left(\delta_1 + \frac{M+1}{2} \bar{T} \right) + P^{(PF)} \left(\delta_1 + \frac{M+1}{2} \bar{T} + \delta_2 \right) + \frac{T_s}{M} \right]$$

Stabilirea valorii optime pentru M , care minimizează pe W , se face iterativ. Valoarea inițială se obține prin punerea valorii 1 pentru numitor și 0 pentru δ_2 ; se ia derivata în raport cu M și se anulează pentru a obține o ecuație în M .

Valoarea inițială pentru M rezultată este aproximativ următoarea:

$$M_{opt} \approx \sqrt{\frac{2T_s}{\bar{T}(P^{(RB)} + P^{(PF)})}} = \sqrt{\frac{2T_s}{\bar{T}(1 - \sum_{i=1}^N f_i e^{-\lambda_i T_i})}}$$

Schema CARER (Cash-Aided Rollback Error Recovery)

Reducerea duratei checkpointării permite checkpointarea mai frecventă; se reduce și penalitatea la revenire după manifestarea disfuncției și recuperarea din starea disfuncțională.

Schema CARER reduce timpul cerut pentru a preleva un checkpoint prin marcarea urmei (*footprint*) procesului în memoria principală și în memoria cache ca parte a stării checkpointate. Se admite că memoria și memoria cache sunt candidate mult mai puțin la disfuncții decât procesorul.

Checkpointarea constă în memorarea registrelor procesorului în memoria principală și includerea urmei proceselor în memoria principală plus în orice linie de cache marcată ca fiind parte din checkpoint.

Biti de checkpointare pentru fiecare linie de cache

Această tehnică necesită modificări de hardware: asocierea unui bit suplimentar de checkpointare asociat cu fiecare linie de cache.

Când acest bit este 1: linia corespunzătoare nu poate fi modificată, adică linia este parte din ultimul, cel mai recent checkpoint; nu poate fi actualizată fără a obliga la prelevarea imediată a unui checkpoint.

Dacă acel bit este 0: procesorul este liber să modifice cuvântul.

Urma procesului în memoria principală și liniile marcate din cache îndeplinesc dubla sarcină de memorie și de parte a checkpointului. Schema aceasta restrânge libertatea în a decide când trebuie prelevate checkpointurile.

Checkpointarea este forțată, obligatorie atunci când:

- O linie marcată ca nemodificabilă trebuie actualizată
- Ceva în memoria principală, nu importă ce, trebuie actualizat
- Se execută o instrucțiune I/O sau apare o întrerupere externă.

Prelevarea unui checkpoint implică salvarea registrelor procesorului în memorie și setarea la 1 a bitului de checkpoint asociat cu fiecare linie de cache validă.

Revenirea la checkpointul anterior este foarte simplă: se restaurează registrele și se marchează ca invalide toate liniile din cache cu bitul de checkpoint nul.

Costul implicat este după cum urmează:

- Un bit de checkpoint pentru fiecare linie de cache.
- Fiecare rescriere a unei linii de cache în memoria principală implică prelevarea unui checkpoint.

Checkpointarea în sistemele distribuite

Sistemele distribuite sunt caracterizate de faptul că procesoarele și memoriile asociate sunt conectate printr-o rețea. Fiecare procesor poate avea discuri locale. Nu este exclusă existența unui sistem de fișiere în rețea accesibile tuturor procesoarelor.

Procesele conectate prin canale directionale: conexiuni punct-la-punct de la un proces la un altul. Canalele se presupun a fi libere de erori și programate pentru a livra mesajele în ordinea primirii.

Într-un sistem distribuit se petrec evenimente deterministe și non-deterministe.

Un eveniment non-determinist este un eveniment a cărui apariție nu poate fi prezisă pe baza stării anterioare a sistemului. Dimpotrivă, un eveniment determinist poate fi prezis.

Execuția procesului este o secvență de evenimente deterministe, întreruptă când și când de unele evenimente non-deterministe.

Exemplu: un program care controlează o supapă de presiune a unui reactor chimic este o buclă nesfârșită cu intrări furnizate cu o anumită frecvență de senzori de presiune; urmează de fiecare dată decizia de control. Valoarea unei intrări se constituie ca un eveniment non-determinist: ea nu poate fi prezisă pe baza rezultatelor observării anterioare.

O noțiune cu care se lucrează uneori este aceea de proces determinist pe porțiuni. Aceasta vrea să cuprindă o situație reală în care odată intrarea cunoscută, restul este predicibil (presupunând lipsa vreunei disfuncții).

O execuție a unui proces poate fi privită ca deterministă pe porțiuni. Se iau în considerare "bucăți" de timp. Fiecare din ele începe cu un eveniment non-determinist. Fiind date informații asupra evenimentului non-determinist și asupra stării procesului la începutul intervalului de timp, se poate prezice fiecare eveniment care se produce pe durata intervalului.

Proces/canal/starea sistemului

Starea unui proces are un sens evident: starea canalului la momentul t este descrisă prin multimea de mesaje transferate prin el până la momentul t (dar și ordinea recepționării lor).

Starea sistemului distribuit este agregatul de stări ale proceselor individuale și ale canalelor.

Starea se zice că este consistentă dacă pentru fiecare livrare de mesaj există un eveniment corespunzător de expediere a mesajului.

O stare care violează această condiție – un mesaj livrat care nu a fost încă trimis – violează cauzalitatea; un asemenea mesaj se numește *orfan*.

Inversul este consistent: starea unui sistem reflectă trimiterea unui mesaj dar nu recepționarea lui.

Despre stări consistente și stări inconsistente, mai multe în discuția următoare.

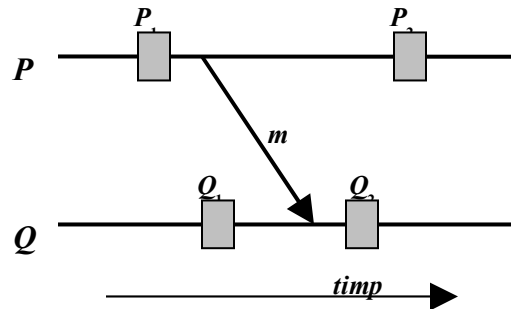
Două procese P și Q au fiecare de prelevat checkpointuri. Mesajul m este trimis de P spre Q .

Iată multimile de checkpointuri care reprezintă stări consistente ale sistemului (v.figura):

$\{P_1, Q_1\}$ – nici un checkpoint nu are stire de mesajul m .

$\{P_2, Q_2\}$ – P_2 indică faptul că m a fost expedit; Q_2 indică recepționarea lui.

$\{P_2, Q_1\}$ – P_2 indică expedierea mesajului m ; Q_1 nu are înregistrată recepția lui m .



Dimpotrivă, mulțimea $\{P_1, Q_2\}$ reprezintă o stare inconsistentă; P_1 nu are înregistrată expediția lui m , Q_2 știe de recepționarea lui m , ceea ce face din m un *mesaj orfan*.

Mulțimea de checkpointuri care reprezintă o stare consistentă a sistemului formează o linie de recuperare: sistemul poate fi făcut să revină (*roll-back*) prin ele și repornit de acolo.

$\{P_1, Q_1\}$: revenirea lui P la P_1 anulează trimiterea mesajului m și revenirea lui Q la Q_1 are semnificația neprimirii lui m de către Q .

Repornirea din aceste checkpointuri are ca efect retransmiterea mesajului m de către P pe care-l va primi Q .

$\{P_2, Q_1\}$: revenirea lui P în P_2 nu include retransmiterea lui m ; iar revenirea lui Q în Q_1 înseamnă că procesul Q nu are înregistrată deocamdată primirea lui m .

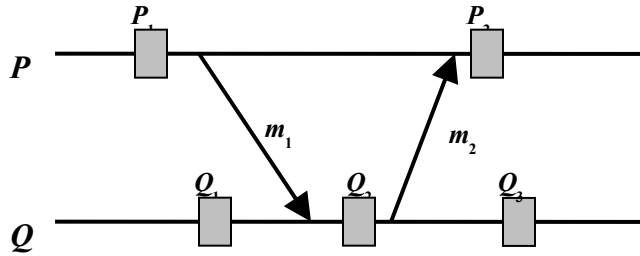
Procesul de recuperare trebuie să fie capabil să reproducă pe m pentru Q ; acest fapt poate fi făcut prin adăugarea lui m la checkpointul lui P sau prin detinerea unui registru separat al mesajelor care să înregistreze tot ce primește Q .

Uneori, checkpointurile pot fi inutile; ele nu vor face parte deloc din linia de recuperare astfel că prelevarea lor este o pierdere de timp.

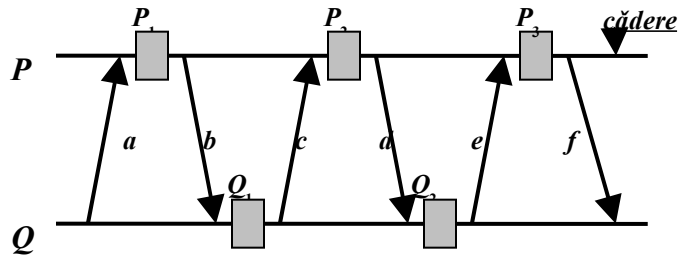
Figura următoare arată ceea ce este un checkpoint inutil.

Q_2 este un checkpoint inutil. El înregistrează recepționarea lui m_1 dar nu expediția lui m_2 .

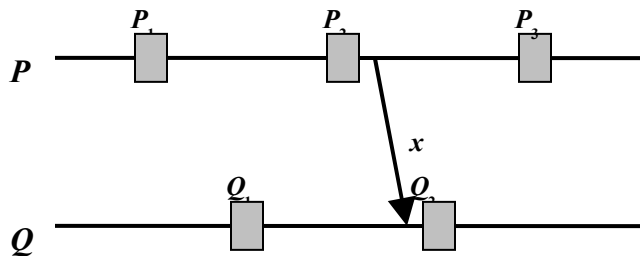
Mulțimea $\{P_1, Q_2\}$ nu poate fi consistentă (altminteri m_1 ar deveni un mesaj orfan). Similar $\{P_2, Q_2\}$ nu poate fi o mulțime de checkpointuri consistentă (deoarece altminteri m_2 ar deveni un mesaj orfan).



Dacă checkpointurile nu sunt coordonate (direct – prin transmitere de mesaje sau indirect – prin ceasuri sincronizate), o disfuncție singulară poate produce un *efect de domino*. Figura următoare ilustrează o asemenea situație.



Când P suferă o disfuncție tranzitorie, el revine la checkpointul P_3 . Deoarece mesajul f a fost trimis după preluarea checkpointului P_3 , Q trebuie să revină (*roll-back*) și el (altminteri Q ar avea un mesaj care n-a fost trimis niciodată, un mesaj orfan). P va reveni la P_2 deoarece Q a trimis un mesaj e către P . Această revenire pas cu pas continuă până când (toate) procesele au revenit la pozițiile lor de plecare. Uneori, unele mesaje se pot pierde.



Mesajele se pot pierde datorită revenirii (*roll-back*).

Se presupune, de pildă, că procesul Q revine la Q_1 după primirea mesajului x de la P . Înregistrarea primirii lui x este pierdută.

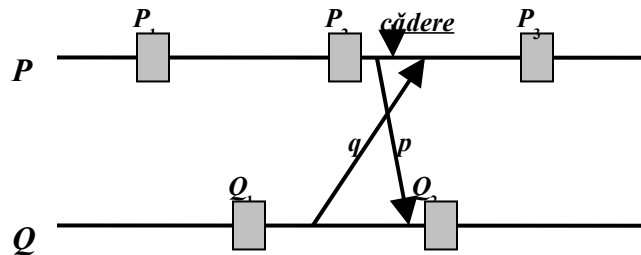
Dacă P nu revine la P_2 , este ca și când P a trimis un mesaj care n-a fost primit vreodată de Q .

Mesajele pierdute astfel nu violează cauzalitatea. Este ceva similar cu pierderea de mesaje din cauza unor probleme în rețea și se rezolvă prin retransmitere.

Totuși, dacă Q trimite un ACK (*acknowledgement*) despre x la P înainte de revenire, atunci acel ACK va fi un mesaj care devine orfan cu excepția cazului când P revine la P_2 .

Este posibil și un blocaj, blocaj care se numește *livelock*⁹. Livelock-ul este o altă problemă care poate apărea în sistemele distribuite checkpointate. Iată un exemplu:

Q trimite lui P un mesaj q , P trimite lui Q un mesaj p . P esuează înainte de primirea mesajului q . Pentru a preveni situația în care p ar deveni orfan, Q trebuie să revină în Q_1 .



Între timp, P se recuperează, revine la P_2 , trimite o copie nouă a lui p și apoi primește copia lui q care a fost trimisă înainte ca orice revenire să înceapă.

Totuși, deoarece Q a revenit, copia lui q este acum orfană și astfel P trebuie să repete revenirea sa. La rândul său, face orfană copia secundă a lui p , forțând astfel pe Q să repete și el revenirea sa.

Acest dans de reveniri poate continua indefinit.

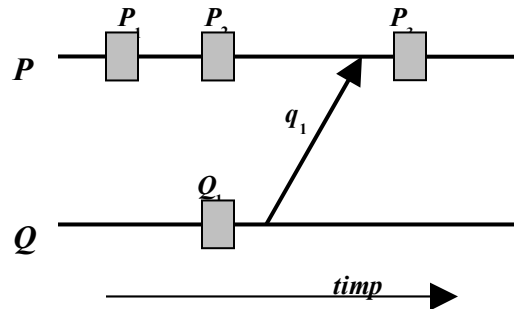
⁹ *Livelock* – o buclă nesfârșită în execuția unui program. Ea apare când un proces se repetă pe sine din cauză că el continuă să primească informații eronate. Poate apărea și atunci când un proces care apelează alt proces este el însuși apelat de acel proces și nu există o logică detectoare a acestei situații în vederea opririi ei. Un “livelock” diferă de un “deadlock” când procesul continuă să aibă loc și nu-i o așteptare într-o buclă fără efecte. Este deosebit și de *deadly embrace* care este un punct mort manifestat prin așteptarea reciprocă a răspunsului de către două elemente de program. De pildă, într-o rețea, un utilizator lucrează cu fișierul A și necesită fișierul B pentru a continua, iar un alt utilizator lucrează cu fișierul B și are nevoie de A pentru a continua; fiecare îl așteaptă pe celălalt, amândoi sunt temporar blocați.

Algoritmi de checkpointare coordonati

Cum s-a văzut, checkpointarea necoordonată poate duce la efectul de domino sau la buclare (*livelock*). Alternativa este checkpointarea coordonată.

Sunt două moduri de bază pentru coordonarea checkpointării:

- Algoritmul Koo-Toueg care are un initiator al procesului de checkpointare pe întreg sistemul.



- Un algoritm care esalonează în timp checkpointuri: esalonarea checkpointurilor poate evita încărcarea excesivă cvasisimultană a sistemului de discuri.

Comunicarea generează proceduri de checkpointare induse.

O soluție mixtă cuprinde utilizarea simultană a algoritmilor de checkpointare coordonați și necoordonați, cei din urmă pentru a se ocupa de disfuncțiile cele mai izolate.

Algoritmul Koo-Toueg

Se admite că P dorește să stabilească un checkpoint la P_3 . Acesta va înregistra q_1 primit de la Q . Pentru a preveni orfanizarea lui q_1 , Q trebuie să checkpointeze și el. Astfel, stabilirea de către P a unui checkpoint la P_3 obligă pe Q să preia un checkpoint care să înregistreze că mesajul q_1 a fost trimis.

Un algoritm pentru astfel de checkpointare coordonată cunoaște două tipuri de checkpointuri: checkpointuri tentative și checkpointuri permanente.

P înregistrează mai întâi starea lui curentă într-un checkpoint tentativă, apoi trimite un mesaj tuturor celorlalte procese de la care a primit un mesaj de când a prelevat ultimul său checkpoint. Se notează mulțimea acestor procese cu Π .

Mesajul înștiințează fiecare proces din Π (de pildă pe Q) de ultimul mesaj, m_{pq} , pe care P l-a primit de la el înainte ca checkpointul tentativă să fi fost prelevat.

Dacă m_{pq} nu a fost înregistrat într-un checkpoint de Q , pentru a preveni m_{pq} de a deveni orfan, Q este pus să preia un checkpoint tentativă pentru a înregistra expedierea lui m_{pq} .

Dacă toate procesele din mulțimea Π , care sunt necesare, confirmă preluarea unui checkpoint cum se cere, atunci toate checkpointurile tentativă pot fi convertite în permanente.

Dacă unele elemente ale lui Π sunt incapabile de un checkpoint cum se cere, P și toate elementele din Π abandonează checkpointurile tentativă și niciunul nu este transformat în permanent.

Aceasta ar putea dezamorsa (*set off*) o reacție în lanț de checkpointuri.

Fiecare element al lui Π poate depune un set de checkpointuri între procese în propria sa mulțime corespunzătoare.

Checkpointarea esalonată (*staggered*)

Algoritmul Koo-Toueg și altele asemenea pot conduce la un număr mare de procese care prelevă checkpointuri aproape simultan. Dacă ele scriu toate într-o memorie stabilă partajată, de pildă un set de discuri comune, transferul poate conduce la congestii ale accesului la disc sau ale rețelei sau concomitent ale amândurora.

Oricare din cele două tratări poate fi utilizată pentru a ne asigura că în orice moment cel mult un proces își preia checkpointul său.

(1) Se scrie checkpointul într-un buffer local, apoi se trec esalonat (*stagger*) scrierile din buffer în memoria stabilă; se presupune că bufferul are o capacitate suficientă.

(2) Se încearcă esalonarea checkpointurilor în timp.

Checkpointurile esalonate ar putea deveni inconsistente; este posibil să existe mesaje orfane în sistem.

Acest fenomen poate fi evitat printr-o fază coordonată în care fiecare proces înregistrează în memoria stabilă toate mesajele pe care le trimite urmând checkpointului precedent. Faza de înregistrare a mesajelor lansate de procese se va suprapune (*overlap*) în timp.

Dacă volumul de mesaje este mai redus decât dimensiunea checkpointurilor individuale, discurile și rețeaua vor cunoaște un flux redus.

Recuperarea din starea de disfuncție

Dacă un proces clachează, el poate fi repornit după revenire la ultimul checkpoint și toate mesajele stocate în registru sunt citite.

Această combinație de checkpoint cu registru de mesaje este denumită *checkpoint logic*.

Acest algoritim de checkpointare esalonată garantează că toate checkpointurile logice formează o linie de recuperare consistentă.

Algoritmul esalonat, *faza primă*, faza checkpointării:

pentru ($i = 0; i \leq n - 1; i++$) {
 P_i preia un checkpoint

P_i trimite un mesaj la $P_{\{(i+1) \bmod n\}}$ cu ordinul ca acesta din urmă să preia un checkpoint
}

Dacă P_0 primește un mesaj de la $P_{\{n-1\}}$ care îi ordonă să preia un checkpoint, acesta este pretextul pentru P_0 de a iniția faza a doua (de înregistrare de mesaje). El trimite un mesaj marker pe fiecare dintre canalele sale de plecare. Când procesul P_i recepționează un mesaj marker, el trece la faza a doua.

Faza a doua a algoritmului esalonat, faza înregistrării de mesaje:

if (nici un mesaj marker anterior nu a fost primit în această rundă de P_i)

then {

P_i trimite un mesaj marker pe fiecare din canalele sale de plecare

P_i înregistrează toate mesajele primite de el după checkpointul precedent

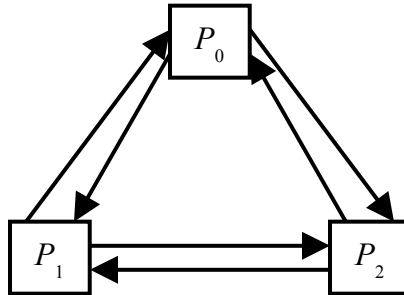
}

else

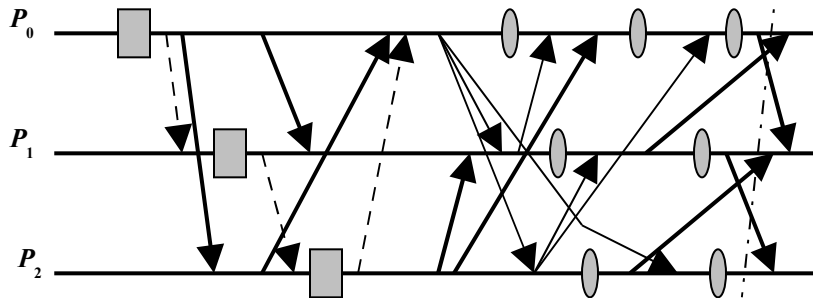
P_i actualizează registrul său de mesaje prin adăugarea tuturor mesajelor primite de el de la ultima actualizare a registrului

end if.

Exemplu de algoritm esalonat – faza 1:



Sistem cu trei procese.



Legenda figurii:

Dreptunghi: prelevare de checkpoint

Elipsă: înregistrare de mesaje

Linie continuă îngrosată: mesaje schimbate de procese

Linie continuă simplă: ordin de înregistrare de mesaje

Linie întreruptă: ordin de prelevare de checkpoint (take_checkpoint)

Linie din linii și puncte: poziție de checkpoint logic

Mesajele schimbate de procese sunt notate m_0, m_1, m_2 etc. în ordinea temporală a generării lor.

P_0 prelevă un checkpoint și trimite mesajul și ordinul take_checkpoint la P_1 .

P_1 trimite un ordin la fel către P_2 după ce prelevă propriul său checkpoint.

P_2 trimite un mesaj imperativ take_checkpoint înapoi la P_0 .

La acest punct, fiecare proces a prelevat un checkpoint și faza a doua poate începe.

Faza a doua pentru același exemplu:

P_0 trimite ordinul de înregistrare de mesaje, message_log la P_1 și P_2 care înregistrează mesajele pe care le-au primit după cel mai recent checkpoint.

P_1 și P_2 trimit mesaje similare cu ordinul message_log.

De fiecare dată când se primește un astfel de mesaj, procesul înregistrează mesajele. Dacă este pentru prima oară când procesul primește un astfel de ordin message_log, procesul trimite mesaje marker pe fiecare din canalele sale emergente.

Recuperarea

Ipoteză: fiind date checkpointul și mesajele primite, un proces se poate recupera. Putem avea mesaje orfane în raport cu checkpointurile fizice prelevate în prima fază. Nu vor exista mesaje orfane în raport cu ultimele (temporal) checkpointuri logice care sunt generate utilizând checkpointurile fizice și registrul de mesaje.

Sincronizarea temporală

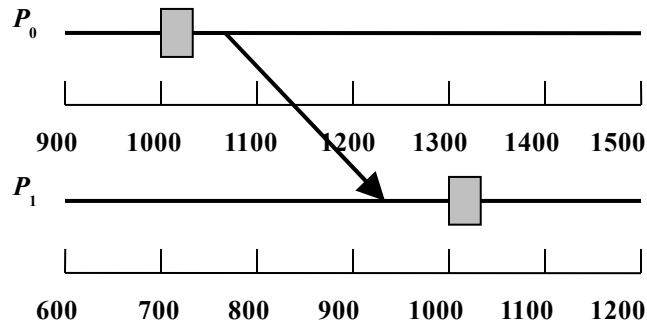
Mesajele orfane pot să nu apară dacă fiecare proces checkpointează exact în aceleși momente. Practic acest lucru este imposibil: deviatiile de ceas și timpii consumați de comunicarea mesajelor nu se pot reduce la zero.

Sincronizarea temporală poate fi încă utilizată pentru a facilita checkpointarea: trebuie însă să luăm în calcul diferențele nenule pe ceasuri.

Sincronizarea temporală (time-based): procesele sunt checkpointate la momente agreeate în prealabil.

Exemplu: se indică fiecărui proces a checkpointa când ceasul local indică un multiplu de 100 de secunde.

O astfel de procedură nu este suficientă în sine pentru evitarea apariției mesajelor orfane. Iată un exemplu de creare a unui mesaj orfan (v.figura).



Fiecare proces prelevă un checkpoint la timpul 1000 (ceasul local). Neconcordanța dintre cele două ceasuri este astfel că procesul P_0 preia checkpointul său mult mai devreme (în timp real universal) decât P_1 . Ca rezultat, P_0 trimite un mesaj către P_1 după checkpointul său, mesaj primit de P_1 înainte de a checkpointa el însuși. Acest mesaj este un potențial orfan.

Cum se poate preveni crearea unui mesaj orfan?

Să presupunem că abaterea de temporalitate între oricare două ceasuri ale sistemului este mărginită de δ și fiecare proces este programat să checkpointeze când ceasul lui local indică timpul τ .

Urmând checkpointului său, un proces P_x ar putea să nu trimită mesaje la vreun alt proces P_y până când nu este sigur că ceasul local al lui P_y arată τ . P_x poate să rămână "tăcut" pe durata $[\tau, \tau + \delta]$ cu timpii măsurați pe ceasul lui local.

Dacă timpul de livrare al mesajelor între procese are o limită inferioară β , pentru a preveni mesajele orfane, procesul P_x trebuie să rămână tăcut un interval mai scurt $[\tau, \tau + \delta - \beta]$. Dacă $\beta > \delta$, acest interval este de lungime nulă și necesitatea ca P_x să păstreze tăcerea dispăre.

Alte metode de prevenire, diferite sunt discutate imediat.

Să presupunem că mesajul m este primit de procesul P_y când ceasul său indică momentul t . Mesajul m trebuie să fi fost trimis (de P_x) nu mai târziu de β de mai sus, înainte ca ceasul lui P_y să indice $t - \beta$.

Deoarece desincronizarea ceasurilor este de cel mult δ , la acest moment, ceasul lui P_x ar fi trebuit să indice cel mult $t - \beta + \delta$.

Dacă $t - \beta + \delta < \tau$, trimiterea lui m va fi înregistrată în checkpointul lui P_x și nu poate fi orfan.

Un mesaj m recepționat de P_y când ceasul său indică cel puțin $\tau - \delta + \beta$ nu poate fi orfan.

Mesajele orfane pot fi evitate de P_y prin neutilizarea și neinclusiunea în checkpointul său la momentul τ a nici unui mesaj primit în intervalul $[\tau - \delta + \beta, \tau]$ (după ceasul P_y) până după prelevarea checkpointului său de la momentul τ .

Checkpointarea fără disc

Memoria principală este volatilă și de aceea nepotrivită pentru depozitarea unui checkpoint; totuși, cu procesoare suplimentare se poate face checkpointarea în memoria principală. De ce? Pentru că prin evitarea scrierii pe discuri, checkpointarea poate fi mai rapidă.

Checkpointarea fără discuri este cel mai bine utilizată ca un nivel într-o checkpointare pe două niveluri. În cazul utilizării tehnicilor de tipul RAID, pentru tratarea disfuncțiilor există procesoare redundante.

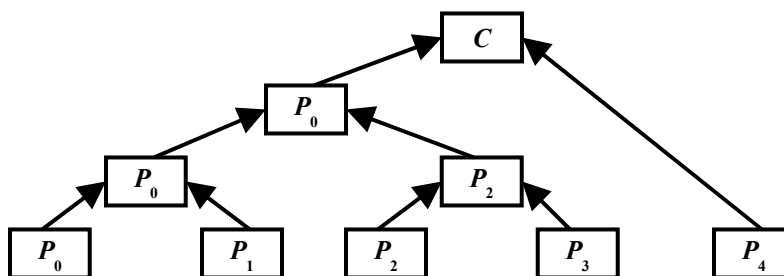
Exemplu: un sistem distribuit cu cinci procesoare care lucrează și unul suplimentar. Fiecare procesor lucrător detine checkpointul său depus în memoria proprie; procesorul suplimentar stochează informația paritară a acestor checkpointuri. Dacă un proces lucrător devine disfuncțional checkpointul său poate fi reconstituit din cele patru rămase plus checkpointul paritar.

Checkpointare fără discuri în stilul RAID

Reteaua care leagă procesoarele trebuie să aibă o bandă de trecere suficientă pentru expedierea checkpointurilor.

Exemplu: cazul cu n procesoare executante și un procesor consacrat checkpointării. Dacă toate procesoarele executante trimit checkpointurile lor la procesorul dedicat checkpointării pentru calculul parității, acesta poate deveni un *hotspot*, un loc de strângere.

Soluția constă în distribuția calculului de paritate.



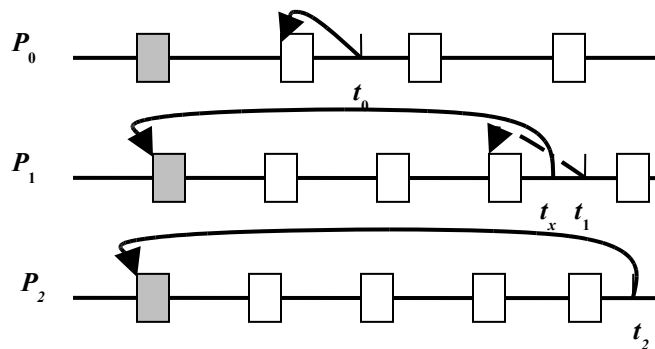
Recuperarea se poate face pe două niveluri. Coordonarea checkpointurilor previne apariția mesajelor orfane dar impune un overhead.

- Nu afectează corectitudinea dacă disfuncțiile sunt izolate, adică la orice moment cel mult un procesor este în starea de disfuncție (în recuperare).
- Vasta majoritate a disfuncțiilor sunt izolate.
- Aceasta face recuperarea din disfuncții izolate rapidă.
- Se acceptă durate mai mari de recuperare pentru disfuncții simultane.

Aceasta sugerează o schemă de recuperare pe două niveluri.

- Primul nivel: fiecare proces prelevă propriul checkpoint fără coordonare (util, potrivit numai pentru recuperarea din disfuncție izolată).
- Checkpointul nu trebuie scris pe disc, poate fi scris în memoria altui procesor.
- Nivelul al doilea: ocazional, întregul sistem face o checkpointare coordonată (cu overhead mai mare) care protejează la disfuncții care nu sunt izolate.

Un exemplu de recuperare pe două niveluri: P_0 cade la momentul t_0 ; sistemul revine la ultimul checkpoint de prim nivel; recuperarea este cu siguranță reușită.



Dreptunghiuri albe: checkpointuri de primul nivel
Dreptunghiuri umbrite: checkpointuri de nivelul al doilea

P_1 cade la momentul t_1 ; inițiază revenirea; în punctul t_x (pe durata recuperării), esuează și P_2 . Disfuncții care nu mai sunt izolate: sistemul revine în ambele procese la ultimul checkpoint de nivelul secund.

În general, cu cât sunt mai dese disfuncțiile neizolate, corelate, cu atât mai mare trebuie să fie frecvența de prelevare a checkpointurilor de nivelul doi.

Înregistrarea mesajelor

Pentru a continua calculul dincolo de ultimul checkpoint, procesul de recuperare poate cere toate mesajele primite până atunci, citite în ordinea originală.

Pentru checkpointarea coordonată, fiecare proces poate fi făcut să revină la ultimul lui checkpoint și repornit: aceste mesaje vor fi reexpediate în timpul execuției repetate.

Pentru a evita overheadul de coordonare și a lăsa checkpointurile proceselor independente, o opțiune o reprezintă înregistrarea mesajelor.

Sunt folosite două modalități de înregistrare:

- Înregistrarea pesimistă: dă siguranță că revenirea nu va difuza, adică dacă un proces esuează, nici un alt proces nu va trebui să facă revenirea pentru a asigura consistența.

- Înregistrarea optimistă: eșecul unui proces poate declanșa și revenirea altor procese.

Înregistrarea pesimistă a mesajelor reprezintă metoda cea mai simplă: primitorul unui mesaj oprește orice activitate când primește acel mesaj, înregistrează mesajul într-o memorie stabilă, apoi reia execuția.

Recuperarea procesului din starea de eșec: procesul este făcut să revină la ultimul său checkpoint și își recitește mesajele pe care le-a primit de la acel checkpoint, în ordinea corectă. Nu va exista nici un mesaj orfan, fiecare mesaj va fi ori primit înainte de ultimul checkpoint, ori salvat explicit în registrul de mesaje. Revenirea unui proces nu va declanșa revenirea altor procese.

Înregistrarea mesajelor după expeditor ține seamă de faptul că înregistrarea mesajelor într-o memorie stabilă poate aduce un overhead semnificativ.

Înregistrarea mesajelor pe baza expeditorului este potrivită vis-à-vis de căderile izolate. Fazele acestui gen de înregistrare sunt următoarele:

- Expeditorul unui mesaj îl înregistrează într-un registru și, la cerere, registrul este citit pentru a reproduce mesajul.
- Fiecare proces are contoare de trimiteri și de primiri care sunt incrementate de fiecare dată când procesul trimite sau primește un mesaj.
- Fiecare mesaj are un SSN (*Send Sequence Number*), valoarea contorului de trimiteri la momentul transmiterii.
- Unui mesaj primit i se alocă un RSN (*Receive Sequence Number*), valoarea contorului de primiri la momentul primirii.
- Primitorul trimite totodată un ACK (*acknowledge*) expeditorului, care include RSN-ul pe care l-a alocat mesajului.
- La primirea acestui ACK, expeditorul confirmă ACK-ul într-un mesaj către primitor.
- În timpul cât receptorul primește mesajul și trimite acel ACK și cât primește ACK-ul expeditorului despre ACK-ul propriu, receptorului îi este interzisă expedierea de mesaje către alte procese, ceea ce este esențial pentru funcționarea corectă la recuperare.
- Un mesaj se numește deplin-înregistrat dacă nodul expeditor știe atât SSN-ul cât și RSN-ul lui. El este partial-înregistrat dacă nodul expeditor nu știe încă RSN-ul lui.
- Când un proces revine și reporneste calculul de la ultimul checkpoint, el trimite celorlalte procese un mesaj care listează SSN-ul ultimului lui mesaj pe care l-a înregistrat în checkpointul său.
- Când acest mesaj este primit de un proces, acel proces știe care mesaje trebuie retransmise și se conformează.
- Procesul de recuperare trebuie acum să utilizeze aceste mesaje în aceeași ordine în care ele au fost utilizate înainte de cădere; este ușor de făcut aceasta pentru mesajele deplin-înregistrate deoarece RSN-urile lor sunt cunoscute și pot fi sortate după aceste numere.

Care este soarta mesajelor partial-înregistrate? Mesajele partial-înregistrate sunt cu RSN-uri necunoscute. Ele sunt trimise, dar ACK-urile lor nu vor fi primite

niciodată de expeditor. Receptorul a esuat înainte ca mesajul să-i fie livrat sau a esuat după recepția mesajului dar înainte de a trimite ACK-ul. Receptorul este oprit a trimite mesaje proprii către alte procese între primirea mesajului și trimiterea ACK-ului asociat. Ca rezultat, recepția mesajului partial-înregistrat într-o ordine diferită a doua oară nu poate afecta nici un alt proces din sistem, corectitudinea este prezervată. Este limpede, această metodă este garantată a lucra numai dacă la un moment dat, altfel oarecare, există cel mult un nod căzut.

Înregistrarea de mesaje optimistă are overhead mai redus decât înregistrarea pesimistă, dar recuperarea din starea de disfuncție este mult mai complexă.

Înregistrarea optimistă este mai curând de interes teoretic.

Când mesajele sunt primite, ele sunt scrise într-un buffer volatil care la momentul potrivit este copiat în memoria stabilă. Derularea procesului nu este întreruptă și astfel overheadul înregistrării este redus.

La cădere, conținutul bufferului poate fi pierdut conducând la necesitatea ca mai multe procese să revină (*multiple roll-back*). Este necesară o schemă specială pentru a rezolva această situație.

Checkpointarea în sistemele cu memorie partajată

Într-o variantă de CARER (*Cash-Aided Roll-back Error Recovery*) pentru sisteme multiprocesor pe bus cu memorie partajată, fiecare procesor are propriul său cache.

Este necesară o schimbare a algoritmului pentru a menține coerența de cache între memorii cash multiple. În locul unui singur bit care să marcheze o linie ca imuabilă, avem un identificator multi-bit: un identificator de checkpoint C_id pe fiecare linie, un contor de checkpointuri pe fiecare procesor, C_count , care reține numărul checkpointului curent. La fiecare prelevare de checkpoint, se incrementează contorul.

O linie modificată înainte va avea C_id -ul ei inferior contorului. Când o linie este actualizată, se pune $C_id = C_count$.

Dacă o linie a fost modificată de când a fost adusă în cache și $C_id < C_count$, linia este parte din starea checkpointului și este prin urmare inaccesibilă scrisului. Orice scrieri într-o astfel de linie trebuie să aștepte până când linia este scrisă mai întâi în memoria principală.

Dacă contorul are k biți, el revine la 0 după atingerea valorii 2^{k-1} .

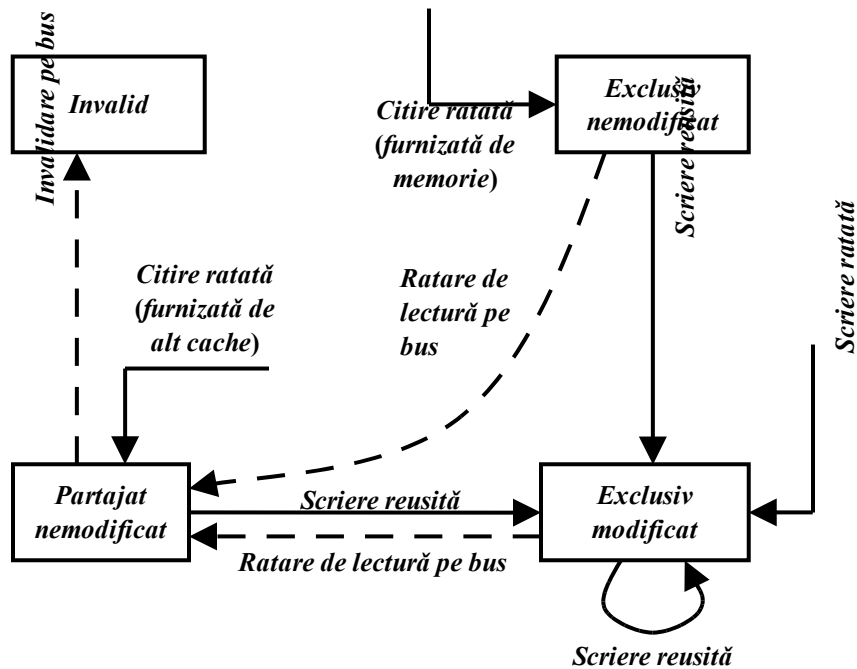
Un protocol de coerență bus-based

Modifică un algoritm de coerență a cache-ului pentru a ține cont de checkpointare. Întregul trafic între cache-uri și memorie trebuie să uzeze de bus, altfel spus, toate cache-urile pot urmări traficul pe bus.

O linie de cache poate fi în una din următoarele stări: invalidă, nemodificată și partajată, modificată și exclusivă și nemodificată și exclusivă.

Exclusivă: aceasta este singura copie validă în vreun cache.

Modificată: linia a fost modificată de când a fost adusă din memorie în cache.



Dacă un procesor are de actualizat o linie aflată în starea nemodificată partajată, ea se schimbă în starea modificată exclusivă. Alte cache-uri care detin aceeași linie trebuie să invalideze copiiile lor care nu mai sunt actuale.

Când e vorba de stări exclusiv modificate sau exclusiv nemodificate și un alt cache pune pe bus o cerere de citire, cache-ul solicitant trebuie să facă serviciul acestei cereri (numai copia curentă a acelei linii).

Produs secundar: memoria este și ea actualizată dacă este necesar. Apoi se schimbă la starea partajată nemodificată.

Ratarea de scriere a unei linii în cache conduce la starea modificată exclusivă.

Coerenta bus-based și protocolul de checkpointare nu sunt lipsite de posibile contradicții. Cum se poate modifica protocolul pentru a ține seamă de checkpointare?

Starea originală exclusiv modificată se divide acum în două:

- Exclusivă modificată
- Inaccesibilă pentru scris (*unwritable*).

Protocolul directory-based

În această tratare este menținut ca piesă centrală un director care înregistrează starea fiecărei linii. Acest director poate fi privit ca unul controlat de un anume *controller* al memoriei partajate. Controllerul manipulează toate ratările la citire și la scriere și toate celelalte operații care schimbă starea liniei.

Exemplu: dacă o linie este în starea exclusivă-nemodificată și cache-ul care detine acea linie vrea să o modifice, el notifică controllerul de intenția sa. Controllerul poate schimba starea la exclusivă-modificată. A implementa această schemă de checkpointare pe un astfel de protocol este atunci o problemă simplă.

Alte utilizări ale checkpointării

(1) Migrarea proceselor

Un checkpoint reprezintă starea procesului; migrarea procesului de la un procesor la un altul înseamnă deplasarea checkpointului și calculul poate fi reluat pe acel nou procesor; permite recuperarea din erori permanente sau intermitente.

Natura checkpointului determină dacă noul procesor trebuie să fie de același model și trebuie să ruleze același sistem de operare.

(2) Echilibrarea încărcării

O utilizare mai bună a unui sistem distribuit prin asigurarea că încărcarea cu calcule este distribuită potrivit între procesoare.

(3) Depanare (*debugging*)

Fisiere de miez (*core files*) sunt basculate (*dumped*) când un program se încheie anormal; acestea sunt în esență checkpointuri care conțin întreaga informație de stare despre procesul afectat; depanatorii (programele de depanare) pot citi fișierele *core* care pot ajuta la depanarea procesului.

(4) Instantanee

Observarea stării programului la momente discrete: înțelegere mai adâncă a comportării programului.

SISTEME DE DISCURI TOLERANTE LA DEFECTE

Pentru a spori siguranța în funcționare, sistemele de memorii pot fi structurate în așa manieră încât prin utilizarea unor redundante ele să manifeste o anumită toleranță la defecte. În alți termeni, prin măsuri prealabile, anumite subsisteme pot suplini alte subsisteme care dintr-un motiv sau altul se defectează. Iată în continuare un exemplu semnificativ în ceea ce privește toleranța la defecte a sistemelor de memorare, în particular cele pe disc.

Memorii ieftine exploatate în condiții de siguranță

Fie n dispozitive de memorare, D_1, D_2, \dots, D_n . Fiecare dintre ele conține k octeți și sunt dispozitive de stocare a datelor. Fie alte m dispozitive de memorare C_1, C_2, \dots, C_m . Și acestea conțin fiecare tot câte k octeți și sunt denumite dispozitive de verificare. Conținutul fiecărui dispozitiv de verificare se calculează din conținutul dispozitivelor de date. Problema este să calculeze conținutul dispozitivelor C_i în așa mod încât oricare m dispozitive din $D_1, D_2, \dots, D_n, C_1, C_2, \dots, C_m$ să-și defecteze, conținutul dispozitivelor defecte să poată fi reconstituit din conținutul dispozitivelor rămase în funcție.

Strategia generală

Formal, modelul defectului este acela al unei pierderi de informație prin ștergere. Dacă un dispozitiv se defectează el iese din joc și sistemul recunoaște această situație de inutilitate a acelui dispozitiv. Pierdere care diferă de apariția unor erori, caz în care defectarea se manifestă prin stocarea și restituirea unor valori incorecte care pot fi recunoscute printr-un anumit gen de codare intrinsecă.

Calculul conținutului fiecărui dispozitiv de verificare C_i necesită o funcție F_i aplicată pe conținutul tuturor dispozitivelor de date. Formulele următoare sunt un exemplu pentru $n = 8$ și $m = 2$. Conținutul dispozitivelor de verificare C_1 și C_2 se obține prin evaluarea funcțiilor F_1 , respectiv F_2 .

$$C_1 = F_1(D_1, D_2, D_3, D_4, D_5, D_6, D_7, D_8)$$

$$C_2 = F_2(D_1, D_2, D_3, D_4, D_5, D_6, D_7, D_8)$$

Metoda de codare RS-RAID (RS – de la Reed-Solomon, RAID – de la Redundant Arrays of Inexpensive/Independent Disks) divide fiecare dispozitiv de memorare în cuvinte. Fiecare cuvânt este alcătuit din w biți, număr ales de programator, dar raportat la anumite restricții. Asadar, fiecare dispozitiv conține

$$l = (k \text{ octeti}) \left(\frac{8 \text{ biti}}{\text{octet}} \right) \left(\frac{1 \text{ cuvânt}}{w \text{ biti}} \right) = \frac{8k}{w} \text{ cuvinte}$$

Funcțiile de codare F_i operează pe cuvinte cu rezultatul tot în cuvinte, ca în relațiile următoare unde $x_{i,j}$ reprezintă cuvântul j stocat pe dispozitivul de memorare X_i .

D_1	D_2	C_1	C_2
$d_{1,1}$	$d_{2,1}$	$c_{1,1} = F_1(d_{1,1}, d_{2,1})$	$c_{2,1} = F_2(d_{1,1}, d_{2,1})$
$d_{1,2}$	$d_{2,2}$	$c_{1,2} = F_1(d_{1,2}, d_{2,2})$	$c_{2,2} = F_2(d_{1,2}, d_{2,2})$
$d_{1,3}$	$d_{2,3}$	$c_{1,3} = F_1(d_{1,3}, d_{2,3})$	$c_{2,3} = F_2(d_{1,3}, d_{2,3})$
⋮	⋮	⋮	⋮
$d_{1,i}$	$d_{2,i}$	$c_{1,i} = F_1(d_{1,i}, d_{2,i})$	$c_{2,i} = F_2(d_{1,i}, d_{2,i})$

Pentru a avea notații mai simple, cu un indice mai puțin, se admite că fiecare dispozitiv reține un cuvânt și numai unul. Pe calea aceasta problema se reduce la n cuvinte-date, d_1, d_2, \dots, d_n și la m cuvinte de verificare c_1, c_2, \dots, c_m calculate din cuvintele-date în așa mod încât să fie tolerată pierderea oricăror m cuvinte din cele $n + m$. Pentru calculul unui cuvânt de verificare c_i de pe dispozitivul C_i se aplică funcția F_i cuvintelor-date

$$c_i = F_i(d_1, d_2, \dots, d_n)$$

Dacă un cuvânt-dată din dispozitivul D_j este actualizat de la d_j la d'_j atunci fiecare din cuvintele de verificare c_i trebuie recalculat prin utilizarea unei funcții $G_{i,j}$ astfel încât

$$c'_i = G_{i,j}(d_j, d'_j, c_i)$$

Când m dispozitive de memorare clachează se reconstruiește sistemul după cum urmează. Mai întâi, pentru fiecare dispozitiv defect D_j se construiește o schemă care să recupereze conținutul lui D_j din cuvintele depuse în dispozitivele functionale. Când operația aceasta este încheiată se reconstituie conținutul unor eventuale dispozitive de verificare disfuncționale C_i , cu ajutorul funcțiilor F_i .

De exemplu, presupunând că $m = 1$, paritatea $n + 1$ se poate descrie în termenii generali de mai sus. Există numai un dispozitiv de verificare C_1 și lungimea cuvântului este de 1 bit ($w = 1$). Pentru calculul cuvântului de verificare c_1 se ia paritatea prin SAU EXCLUSIV (XOR) a cuvintelor de date

$$c_1 = F_1(d_1, d_2, \dots, d_n) = d_1 \oplus d_2 \oplus \dots \oplus d_n$$

Dacă cuvântul de pe suportul D_j se schimbă din d_j în d'_j atunci c_1 se recalculează din paritatea vechiului cuvânt și din cele două cuvinte-date

$$c'_1 = G_{1,j}(d_j, d'_j, c_1) = c_1 \oplus d_j \oplus d'_j$$

Dacă un dispozitiv se defectează atunci fiecare cuvânt poate fi reconstituit prin paritatea cuvintelor de pe dispozitivele rămase în funcție

$$d_j = d_1 \oplus \dots \oplus d_{j-1} \oplus d_{j+1} \oplus \dots \oplus d_n \oplus c_1$$

Sistemul este rezistent la defectarea oricărui (unic) suport.

O reformulare a problemei sună astfel: sunt date n date d_1, d_2, \dots, d_n , toate de dimensiunea w . Se definesc functiile F si G care sunt utilizate pentru a calcula si pentru a întretine, a mentine actuale cuvintele de verificare c_1, c_2, \dots, c_m . Se face o descriere a modului cum se reconstituie cuvintele de pe orice suport esuat când numărul de dispozitive de memorare defecte nu depășeste m . Cu cuvintele-date reconstituite se recalculează cuvintele de verificare din cuvintele-date cu ajutorul functiilor F . Sistemul este refăcut în întregime.

Algoritmul RS-RAID

Trei aspecte sunt deosebite în aplicarea algoritmului. Primul constă în utilizarea matricilor Vandermonde (Alexandre-Théophile Vandermonde, 1735-1796) pentru calculul si mentinerea cuvintelor de control, al doilea este utilizarea eliminării Gauss pentru recuperarea din starea de nefunctionare si al treilea, utilizarea aritmeticii specifice câmpurilor Galois. Acestea toate sunt detaliate în continuare.

Calculul si întretinerea cuvintelor de verificare

Functiile F_i sunt prin definitie combinatii liniare ale cuvintelor-date

$$c_i = F_i(d_1, d_2, \dots, d_n) = \sum_{j=1}^n f_{i,j} d_j$$

Cu alte cuvinte, dacă se adoptă o reprezentare matricială cu D si C vectori si F_i linii într-o matrice F

$$FD = C$$

Matricea F este definită ca o matrice Vandermonde $m \times n$ cu $f_{ij} = j^{i-1}$, ceea ce face din relatia de mai sus

$$\begin{bmatrix} f_{1,1} & f_{1,2} & \dots & f_{1,n} \\ f_{2,1} & f_{2,2} & \dots & f_{2,n} \\ \vdots & \vdots & & \vdots \\ f_{m,1} & f_{m,2} & \dots & f_{m,n} \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_n \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & 2 & 3 & \dots & n \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 2^{m-1} & 3^{m-1} & \dots & n^{m-1} \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_n \end{bmatrix} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix}$$

Când unul din cuvintele-date d_j se schimbă în d_j' cuvintele de verificare trebuie schimbate în consecință. Prin scăderea porțiunii din cuvântul de verificare care corespunde lui d_j si adunarea cantității necesare pentru d_j' se obtine pentru $G_{i,j}$ definitia din relatia de mai jos

$$c'_i = G_{i,j}(d_j, d'_j, c_i) = c_i + f_{i,j}(d'_j - d_j)$$

Asadar, calcularea si întretinerea cuvintelor de verificare pot fi făcute printr-o aritmetică simplă, dar după regulile date mai departe.

Recuperarea din *crash*

Pentru a explica recuperarea aceasta se definesc matricea $A = \begin{bmatrix} I \\ F \end{bmatrix}$ si vectorul

$E = \begin{bmatrix} D \\ C \end{bmatrix}$. Apoi se scrie ecuatia

$$AD = E$$

care în formă detaliată apare ca

$$\begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 1 \\ 1 & 2 & 3 & \cdots & n \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 2^{m-1} & 3^{m-1} & \cdots & n^{m-1} \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_n \end{bmatrix} = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_n \\ c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix}$$

Se observă că fiecare suport din sistem are o linie în matricea A si în vectorul E . Dacă un dispozitiv esuează, eșecul se materializează în relația de mai sus prin ignorarea/stergera liniei care corespunde acelui dispozitiv. Rezultă o matrice A' și un vector E' cu linii mai puține, dar care verifică o ecuație asemănătoare cu cea de mai sus

$$A'D = E'$$

Dacă exact m dispozitive sunt inutilizabile atunci matricea A' este o matrice $n \times n$. Deoarece matricea F este de tipul Vandermonde orice submulțime de linii este liniar independentă. Matricea A' este asadar nesingulară și valorile care compun vectorul D pot fi calculate din ecuația matricială de mai sus prin eliminare Gauss. Prin urmare toate datele pot fi recuperate.

Cu D odată obținut, valorile oricărui suport cu date de verificare C_i esuat pot fi reconstituite. Dacă sunt mai puțin de m dispozitive cu probleme, alegerea la întâmplare a unui număr de exact n linii din A' permite eliminarea gaussiană și continuarea este evidentă. Sistemul tolerează până la m dispozitive devenite de neutilizat.

Aritmetica în câmpurile Galois

O preocupare majoră în algoritmul RS-RAID o constituie domeniul calculelor care este o mulțime de cuvinte binare de lungime fixă w . Recuperarea dintr-o eroare comisă obisnuit ar putea consta în efectuarea unor calcule modulo 2^w . Maniera aceasta însă nu funcționează deoarece împărțirea nu-i definită pentru orice pereche de elemente. De exemplu, $3:2 \text{ modulo } 4$ nu este definită. Această situație face eliminarea Gauss imposibilă în foarte multe cazuri.

Câmpurile cu 2^w elemente sunt câmpuri Galois – notate $GF(2^w)$ – un subiect fundamental în algebra abstractă. Mai jos se definesc moduri eficiente de a aduna, scădea, multiplica și împărți elemente aparținând unui câmp Galois.

Elementele unui câmp Galois $GF(2^w)$ sunt întregi de la zero la $2^w - 1$. Adunarea și scăderea sunt aplicații simple: operații XOR (SAU EXCLUSIV). De pildă, în $GF(2^4)$

$$11 + 7 = 1011 \oplus 0111 = 1100 = 12$$

$$11 - 7 = 1011 \oplus 0111 = 1100 = 12$$

Multiplicarea și diviziunea sunt mai complexe. Când w este mic, 16 sau mai mic, se utilizează tabele de logaritmi lungi de $2^w - 1$. Tabelul conține indici, o funcție *gflog* și o funcție *gfilog*. Ambele funcții sunt funcții cu valori întregi. Prima este listată pentru indici de la 1 la $2^w - 1$ și este o listă de logaritmi în câmpul Galois, a doua este definită pentru indici de la 0 la $2^w - 2$ și conține rezultatul operației inverse logaritmării. Evident, compunerea celor două funcții, în orice ordine produce funcția identitate, $gflog[gfilog(i)] = i$, $gfilog[gflog(i)] = i$. Cu aceste funcții se pot executa operațiile de multiplicare și de împărțire prin luarea logaritmilor factorilor, prin calculul sumei sau a diferenței valorilor obținute (*modulo* $2^w - 1$) și revenirea la rezultat prin operația inversă logaritmării. Iată un tabel de logaritmi pentru $w = 4$.

<i>i</i>	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
gflog(<i>i</i>)		0	1	4	2	8	5	10	3	14	9	7	6	13	11	12
gfilog(<i>i</i>)	1	2	4	8	3	6	12	11	5	10	7	14	15	13	9	

Evident, numai numerele nenule au logaritmi. Logaritmul invers al unui număr i este egal cu logaritmul invers al numărului $[i \bmod (2^w - 1)]$.

Exemple de calcul în aritmetica din $GF(2^4)$:

$$3 * 7 = gfilog[gflog(3) + gflog(7)] = gfilog[4 + 10] = gfilog[14] = 9$$

$$13 * 10 = gfilog[gflog(13) + gflog(10)] = gfilog[13 + 9] = gfilog[7] = 11$$

$$13 / 10 = gfilog[gflog(13) - gflog(10)] = gfilog[13 - 9] = gfilog[4] = 3$$

$$3 / 7 = gfilog[gflog(3) - gflog(7)] = gfilog[4 - 10] = gfilog[9] = 10$$

Asadar, o multiplicare sau o diviziune necesită trei apeluri la tabel – două pentru logaritmi și unul pentru inversul logaritmului –, o adunare sau o scădere și o operație de tip *modulo*.

Aritmetica unui câmp Galois are fundamentele date în continuare.

Un câmp $GF(n)$ este o mulțime de n elemente închisă la operațiile de adunare și multiplicare, cu un invers (opus) raportat la adunare pentru fiecare element, cu un invers raportat la operația de înmulțire pentru fiecare element nenul. Câmpul $GF(2)$ de exemplu, conține două elemente, adunarea și înmulțirea se execută *modulo* 2 (operatorii XOR și AND, respectiv). Analog, dacă n este prim atunci $GF(n)$ este mulțimea $\{0, 1, \dots, n - 1\}$ în care adunarea și înmulțirea se execută *modulo* n .

Dacă $n > 1$ nu este prim, atunci mulțimea $\{0, 1, \dots, n - 1\}$ cu adunarea și multiplicarea *modulo* n nu este câmp. De exemplu, dacă $n = 4$ atunci mulțimea

$\{0, 1, 2, 3\}$, închisă la adunare și înmulțire nu este câmp deoarece 2 nu are un invers la înmulțire, nu există a astfel încât $2a = 1 \pmod{4}$. Asadar, nu se poate face codarea cu cuvinte binare de lungime $w > 1$ cu operațiile de adunare și înmulțire modulo 2^w . În loc trebuie utilizat câmpul Galois corespunzător.

Explicațiile relative la câmpurile Galois uzează de polinoamele într-o nedeterminată cu coeficienți în $GF(2)$. Adică, dacă $r(x) = x + 1$ și $s(x) = x$ atunci $r(x) + s(x) = 1$ deoarece $x + x = (1 + 1)x = 0x = 0$. Mai mult, se pot lua polinoame de acest gen *modulo* alte polinoame conform cu identitatea: dacă $r(x) \pmod{q(x)} = s(x)$, atunci $s(x)$ este un polinom de grad inferior lui $q(x)$ și $r(x) = q(x)t(x) + s(x)$ cu $t(x)$ un polinom în x .

Dacă, de exemplu, $r(x) = x^2 + x$ și $q(x) = x^2 + 1$ atunci $r(x) \pmod{q(x)} = x + 1$.

Fie acum $q(x)$ un polinom *primitiv* de gradul w cu coeficienți în $GF(2)$. Primitiv înseamnă că polinomul nu are divizori cu coeficienți în $GF(2)$ și polinomul x este generatorul câmpului $GF(2^w)$. Cum generează x câmpul? Se consideră inițial elementele obligatorii 0, 1 și x și se continuă enumerarea elementelor obținute prin multiplicarea ultimului element cu x și reținerea rezultatului modulo $q(x)$. Enumerarea se încheie cu elementul pentru care rezultatul *modulo* $q(x)$ este egal cu 1.

Dacă, de exemplu, $w = 2$ și $q(x) = x^2 + x + 1$ atunci primele elemente sunt 0, 1 și x , iar $x^2 \pmod{q(x)} = x + 1$ și cele patru elemente ale câmpului $GF(4)$ sunt $\{0, 1, x, x + 1\}$. Un al cincilea element nu există deoarece $x(x + 1) = x^2 + x$ care luat *modulo* $q(x)$ produce 1, element deja existent.

Câmpul general $GF(2^w)$ se construiește prin găsirea unui polinom primitiv $q(x)$ de gradul w peste $GF(2)$ urmată de enumerarea elementelor generate de x . Adunarea și multiplicarea elementelor câmpului se fac după regulile adunării și multiplicării polinoamelor cu grija de a lua rezultatul totdeauna *modulo* $q(x)$. Un asemenea câmp se mai scrie ca $GF(2^w) = GF(2)[x]/q(x)$.

Acum, pentru a uza de un câmp $GF(2^w)$ în algoritmul RS-RAID este necesară definirea unei aplicații a elementelor din acest câmp pe cuvinte binare de lungime w . Un polinom $r(x)$ din $GF(2^w)$ poate fi aplicat pe un cuvânt binar b de lungime w prin punerea celui de al i -lea bit din b egal cu coeficientul puterii x^i din polinom. Pentru $GF(4) = GF(2)[x]/x^2 + x + 1$ se obține tabelul următor

Elemente generate	Elemente polinomiale	Elemente binare	Reprezentarea zecimală
0	0	00	0
x^0	1	01	1
x^1	x	10	2
x^2	$x + 1$	11	3

Adunarea elementelor se realizează prin operația SAU EXCLUSIV (XOR) bit-cu-bit. Multiplicarea este mai complicată: se iau elementele sub formă polinomială, se multiplică ca polinoame și se ia rezultatul *modulo* $q(x)$. Tabelele

de logaritmi ca acela de mai sus se bazează pe o tabelă de compunere ca aceea dată pentru cazul $GF(4)$.

Pentru alte valori w , polinoame primitive $q(x)$ se găsesc în literatură. Iată câteva:

$$\begin{aligned} w = 4: & \quad x^4 + x + 1 \\ w = 8: & \quad x^8 + x^4 + x^3 + x^2 + 1 \\ w = 16: & \quad x^{16} + x^{12} + x^3 + x + 1 \\ w = 32: & \quad x^{32} + x^{22} + x^3 + x + 1 \\ w = 64: & \quad x^{64} + x^4 + x^3 + x + 1 \end{aligned}$$

Cu elementul de pornire $x^0 = 1$, $GF(2^w)$ se completează prin enumerarea elementelor obtinute prin multiplicarea cu x a ultimului element enumerat si luarea rezultatelor *modulo* $q(x)$. Tabelul care urmează cuprinde cazul câmpului $GF(2^4)$ cu polinomul primitiv $q(x) = x^4 + x + 1$. În același tabel se observă și modul cum se generează tabelele de logaritmi și de invers-logaritmi prezentate mai devreme.

Element generat	Polinom	Exprimare binară	Exprimare zecimală
0	0	0000	0
x^0	1	0001	1
x^1	x^1	0010	2
x^2	x^2	0100	4
x^3	x^3	1000	8
x^4	$x + 1$	0011	3
x^5	$x^2 + x$	0110	6
x^6	$x^3 + x^2$	1100	12
x^7	$x^3 + x + 1$	1011	11
x^8	$x^2 + 1$	0101	5
x^9	$x^3 + x$	1010	10
x^{10}	$x^2 + x + 1$	0111	7
x^{11}	$x^3 + x^2 + x$	1110	14
x^{12}	$x^3 + x^2 + x + 1$	1111	15
x^{13}	$x^3 + x^2 + 1$	1101	13
x^{14}	$x^3 + 1$	1001	9
x^{15}	1	0001	1

Sumarul algoritmului

Fiind date n dispozitive pentru date și m dispozitive de control, algoritmul RS-RAID care le face tolerante la cel mult m defecte se aplică astfel:

1. Se alege o valoare pentru w astfel ca $2^w > m + n$. Este convenabil a se alege $w = 8$ sau $w = 16$ ceea ce conduce la cuvinte numărate în octeți (bytes). Pentru $w = 16$ suma $m + n$ poate fi până la 65535.

2. Se stabilesc tabelele cu functiile *gflog* si *gfilog* după metoda dată mai devreme
3. Se calculează matricea F care este o matrice Vandermonde $m \times n$: $f_{i,j} = j^{i-1}$ (pentru $1 \leq i \leq m$, $1 \leq j \leq n$), cu operatiile în câmpul Galois $GF(2^w)$.
4. Matricea F se foloseste la calculul si la întretinerea dispozitivelor de verificare pe baza cuvintelor depuse pe dispozitivele de date. Iarăsi, operatiile se fac în $GF(2^w)$.
5. Dacă un număr de dispozitive, mai putine de m clachează atunci ele se reconstituie în maniera care urmează. Se aleg oricare n dispozitive din cele rămase în functie si se construiesc matricea A' si vectorul E' ca mai sus. Se rezolvă apoi pentru D ecuatia $A'D = E'$. Prin aceasta datele de pe dispozitivele de stocare a datelor sunt recuperate. Acum se pot reconstitui si dispozitivele de verificare esuate, prin utilizarea matricii F .

Un exemplu. Se presupune că sunt trei suporturi de date si trei suporturi de verificare si fiecare dintre ele detine un megaoctet. Asadar, $n = 3$ si $m = 3$. Se alege $w = 4$, asa încât $2^w > m + n$. Pentru multiplicări se foloseste tabelul dat mai devreme pentru $GF(2^4)$. În aceste conditii matricea F este

$$F = \begin{bmatrix} 1^0 & 2^0 & 3^0 \\ 1^1 & 2^1 & 3^1 \\ 1^2 & 2^2 & 3^2 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \\ 1 & 4 & 5 \end{bmatrix}$$

Se pot calcula acum cuvintele de verificare prin relatia $C = FD$. Se admite că primele cuvinte stocate pe cele trei dispozitive de date sunt, respectiv, 3, 13, 9. Calculul cuvintelor de control C_1, C_2, C_3 produce valorile următoare:

$$C_1 = 1*3 + 1*13 + 1*9 = 3 + 13 + 9 = 0011 \oplus 1101 \oplus 1001 = 0111 = 7$$

$$C_2 = 1*3 + 2*13 + 3*9 = 3 + 9 + 8 = 0011 \oplus 1001 \oplus 1000 = 0010 = 2$$

$$C_3 = 1*3 + 4*13 + 5*9 = 3 + 1 + 11 = 0011 \oplus 0001 \oplus 1011 = 1001 = 9$$

Dacă, de pildă, continutul dispozitivului de date cu indicele 2 se modifică si primul număr devine 1, atunci fiecare din dispozitivele de verificare primeste valoarea $(1-13) = (0001 \oplus 1101) = 1100 = 12$, care este utilizată pentru recalcularea valorilor de verificare

$$C_1 = 7 + 1*12 = 0111 \oplus 1100 = 1011 = 11$$

$$C_2 = 2 + 2*12 = 2 + 11 = 0010 \oplus 1011 = 1001 = 9$$

$$C_3 = 9 + 4*12 = 9 + 5 = 1001 \oplus 0101 = 1100 = 12$$

Dacă D_2, D_3 si C_3 se pierd atunci, din matricea A si din vectorul E se sterg liniile care corespund dispozitivelor disfuncte pentru a obtine ecuatia $A'D = E'$

$$\begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 3 \end{bmatrix} D = \begin{bmatrix} 3 \\ 11 \\ 9 \end{bmatrix}$$

Prin eliminare Gauss se poate inversa matricea A' si se obtine

$$D = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 3 & 1 \\ 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 11 \\ 9 \end{bmatrix}$$

si valorile reconstituite sunt

$$D_2 = 2*3 + 3*11 + 1*9 = 6 + 33 + 9 = 48$$

$$D_3 = 3*3 + 2*11 + 1*9 = 9 + 22 + 9 = 40$$

Cu matricea F se poate reconstitui si

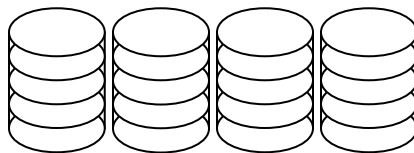
$$C_3 = 1*3 + 4*11 + 5*9 = 3 + 44 + 45 = 92$$

si sistemul este recuperat în întregime.

Există si alte sisteme RAID la care se fac scurte referiri în continuare.

RAID nivelul 0 (fără redundanță)

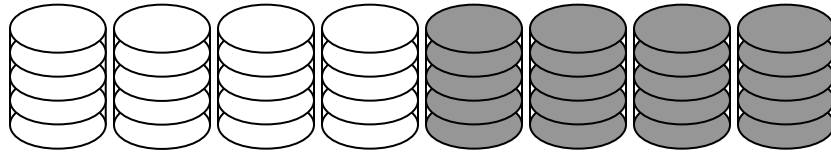
Un sistem de discuri non-redundant (sau de nivel 0) are costul cel mai scăzut din cauza lipsei oricărei redundante. Schema aceasta oferă cea mai bună performanță la scriere deoarece nu necesită vreodată actualizarea vreunei informații redundante. Surprinzător, nu are cea mai bună performanță la citire. Schemele redundante (ca aceea numită “în oglindă” sau “oglindită”, care creează duplicate ale datelor) pot executa mai bine citirile prin planificarea selectivă a cererilor pe discul cu timpul de căutare mediu cel mai scurt si întârzierea rotativă cea mai mică. Fără redundante, orice cădere a unui disc va produce pierderi iremediabile de date. Sistemele de discuri fără redundanțe sunt larg utilizate în supercalcul unde performanța si capacitatea trec ca importanță înaintea fiabilității.



Sistem de discuri fără redundanțe

RAID nivelul 1 (sisteme “în oglindă”)

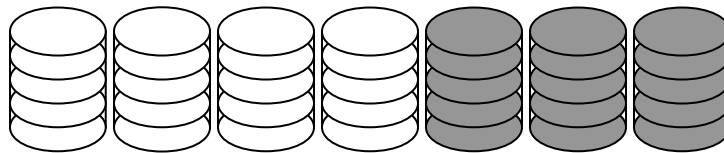
Sistemul traditional denumit “în oglindă” sau “cu umbră” utilizează de două ori mai multe discuri decât un sistem de discuri fără redundante. Ori de câte ori un articol-dată este scris pe un disc el este scris si pe un disc redundat, astfel încât există totdeauna două copii ale informației, două exemplare. Când articolul trebuie citit, el poate fi recuperat de pe disc cu întârzieri de asteptare, de căutare si rotationale mai scurte. Dacă un disc se defectează, copia este utilizată pentru serviciul cerut. Oglindirea este utilizată frecvent în aplicatii cu baze de date, când disponibilitatea sistemului si viteza tranzactiilor sunt mai importante decât eficiența stocării.



Sistem de discuri “în oglindă”

RAID nivelul 2 (memorie în stilul codurilor corectoare de erori)

Sistemele de memorii asigură recuperarea din starea de disfuncție a unor componente la costuri mai reduse decât prin oglindire, dacă folosesc codurile Hamming. Codurile Hamming fac verificări de paritate pe submultimi de componente distincte și suprapuse. În una din variantele acestei scheme, patru discuri de date necesită trei discuri redundante, unul mai puțin decât în cazul sistemului oglindă. Deoarece numărul de discuri redundante este proporțional cu logaritmul numărului total de discuri din sistem, eficiența memorării crește odată cu numărul discurilor de date.

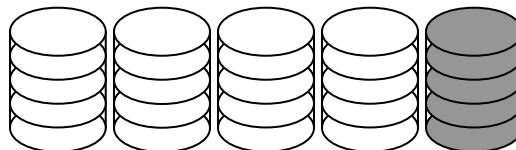


Sistem de discuri în stilul codurilor corectoare de erori

Dacă unul (și numai unul) din discuri cade, componentele mai multor discuri de paritate devin inconsistente cu datele și componenta defectă este identificată: este componenta comună tuturor subseturilor incorecte. Informația pierdută este recuperată prin regulile obișnuite ale codului Hamming utilizat.

RAID nivelul 3 (paritate cu biti intercalați)

Sistemului din paragraful anterior i se pot aduce îmbunătățiri prin observarea faptului că spre deosebire de căderea elementelor de memorie, controlerul discurilor pot identifica ușor care disc este cel cu defect. Astfel, pentru recuperarea informației pierdute se poate utiliza un singur disc de paritate și nu un set de discuri de paritate.

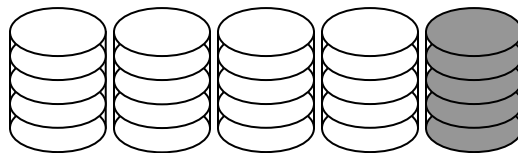


Sistem cu paritate prin biti intercalați

Într-un sistem de discuri cu paritate de biti intercalati, datele sunt conceptual intercalate bit-cu-bit pe discurile de date si se adaugă un singur disc de paritate pentru a tolera căderea unui disc (si numai a unui). Fiecare cerere de citire accesează toate discurile de date si fiecare cerere de scriere accesează toate discurile de date si discul de paritate. Astfel numai o cerere poate fi servită la un moment dat. Deoarece discul de paritate contine numai informatia de paritate si nu date, discul de paritate nu poate participa activ la citiri, ceea ce produce o usoară scădere în performanta la citire față de sistemele cu redundante care distribuie informatia de paritate si datele pe toate discurile. Sistemele de discuri cu paritate pe biti intercalati sunt utilizate frecvent în aplicatii care cer o lărgime de bandă mare dar nu viteze intrare-iesire mari. Mai sunt si simplu de implementat.

RAID nivelul 4 (paritate pe blocuri intercalate)

Există o similitudine între sistemele de discuri cu intercalare de biti si cele cu intercalare de blocuri. Deosebirea constă în obiectul operatiei de intercalare: nu biti, ci blocuri de dimensiune arbitrară. Dimensiunea acestor blocuri este denumită unitatea de realizare de benzi (*striping*). Citirile cerute, mai mici decât unitatea de *striping* accesează numai un disc de date. Cererile de scriere trebuie să actualizeze blocurile de date cerute si trebuie totodată să calculeze si să actualizeze blocul de paritate. Pentru scrieri de mare întindere care ating blocuri pe toate discurile, paritatea este calculată observînd cum diferă datele noi de cele vechi si aplicînd acele diferente pe blocul de paritate. Scrierile de mică întindere cer asadar patru operatii de intrare-iesire pe disc: una pentru a scrie articolul nou, apoi două pentru a citi vechiul articol si vechea paritate pentru a calcula noua informatie de paritate si una de scriere a noii parități. Această operatie este cunoscută ca o procedură citește-modifică-scrie. Deoarece un sistem de discuri cu paritate cu bloc intercalat are numai un disc de paritate, care trebuie actualizat la toate operatiile de scriere, discul de paritate poate deveni cu usurință un loc îngust, o strangulare. Din cauza acestei posibile limitări sunt de preferat sistemelor de discuri cu paritate pe blocuri, sistemele de discuri cu paritate pe blocuri distribuite.

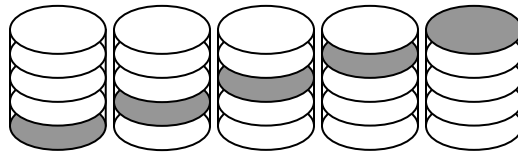


Sistem de discuri cu paritate pe blocuri intercalate

RAID nivelul 5 (paritate pe blocuri intercalate distribuite)

Sistemele de discuri cu paritate pe blocuri distribuite elimină strangularea de pe discul de paritate, care se constată la sistemele cu paritate pe blocuri intercalate,

prin distribuirea informației de paritate uniform pe toate discurile. Un avantaj suplimentar al distribuirii informației de paritate uniform pe toate discurile, frecvent trecut cu vederea, constă în distribuirea datelor pe toate discurile și nu pe toate discurile cu excepția unuia. Aceasta permite tuturor discurilor să participe la servirea operațiilor de citire spre deosebire de schemele redundante cu discuri de paritate dedicate exclusiv în care discurile de paritate nu pot participa la servirea solicitărilor de citire. Sistemele cu paritate pe blocuri intercalate distribuite au una dintre ele mai bune performanțe pentru citiri restrânse, citiri lungi și scrieri lungi, între toate sistemele de discuri cu redundanță. Cu toate acestea, cererile de citire de lungime restrânsă sunt întrucâtva ineficiente comparativ cu alte scheme redundante, cum sunt cele cu oglindire, datorită necesității de a executa operații citeste-modifică-scrie pentru actualizarea parității. Aceasta este partea slabă majoră a sistemelor RAID de nivel 5, în ceea ce privește performanțele.



Sistem de discuri cu paritate pe blocuri intercalate distribuite

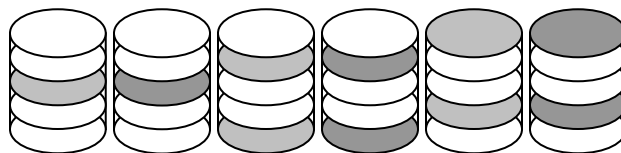
Metoda exactă utilizată pentru a distribui paritatea în sistemele cu paritate distribuită în blocuri intercalate are impact asupra performanțelor. Figura alăturată ilustrează cea mai bună distribuție a informației de paritate (discurile gri), numită distribuție de paritate simetrică la stânga. O proprietate utilă a acestui gen de distribuție constă în faptul că ori de câte ori sunt traversate secvențial unitățile de *striping*, fiecare disc este accesat în succesiune o dată înainte de a fi accesat a doua oară. Această proprietate reduce conflictele de disc atunci când sunt servite solicitările mari.

0	1	2	3	P₀
5	6	7	P₁	4
10	11	P₂	8	9
15	P₃	12	13	14
P₄	16	17	18	19

Sistem RAID de nivel 5 cu simetrie la stânga

RAID nivelul 6 (redundante P + Q)

Paritatea este verificată printr-un cod redundant capabil a corecta orice defectare singulară care se autoidentifică. Dacă sunt luate în considerare sisteme cu discuri mai numeroase, este necesar a utiliza coduri mai puternice, capabile să tolereze defecte multiple. Mai mult, când un disc într-un sistem protejat prin paritate cade, recuperarea conținutului discului defect reclamă lectura cu succes a conținutului tuturor discurilor functionale. În aceste cazuri, probabilitatea de a întâlni în cursul recuperării o eroare de citire necorectabilă poate fi semnificativă. Asadar, aplicațiile cu cerințe de fiabilitate mai severe trebuie tratate cu coduri corectoare de erori mai puternice.



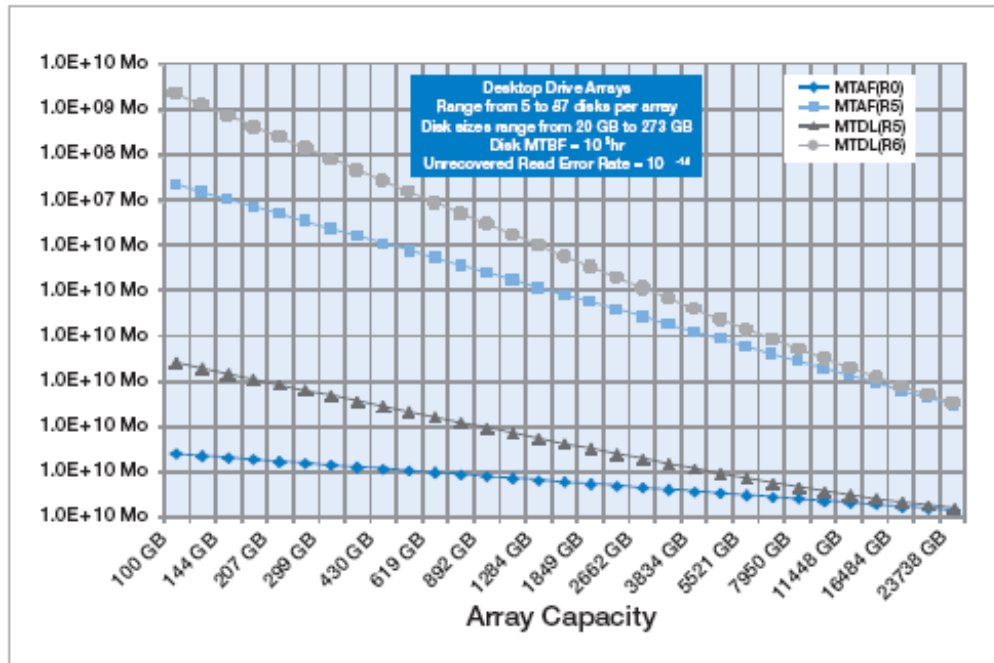
Sistem de discuri cu redundante P + Q

O astfel de schemă, denumită adesea schemă cu redundanță P + Q, folosește codurile Reed-Solomon pentru protecția față de căderea a până la două discuri utilizând cel puțin două discuri redundante. Sistemele de discuri cu redundanță P + Q sunt structural foarte asemănătoare sistemelor cu paritate distribuită bloc-intercalată și operează în mare măsură în același mod. În particular, sistemele de discuri cu redundanță P + Q execută operații scurte de scriere în maniera citește-modifică-scrie cu deosebirea că în loc de patru accesări de disc sunt aici necesare șase accesări pentru a actualiza atât informațiile P cât și cele Q. Sistemele RS-RAID prezentate cu mai multe detalii în prima parte a acestei secțiuni se încadrează în această clasă de sisteme de nivelul 6.

Alte sisteme RAID

Literatura mai menționează:

- Sistemele de **nivel 0+1** – oglindă și *striping* – cu două subsisteme 0 cu *striping* și un subsistem 1 suprapus acestora. Se utilizează și în cazul replicării datelor pentru partajarea lor.
- Sistemele de **nivel 1+0** – *striping* pe oglinzi – în care sunt create subsisteme RAID 1 și peste acestea un subsistem RAID 0 de *striping*.
- Sisteme de **nivelul 7** – un sistem brevetat de Storage Computer Corporation care adaugă secțiuni cache la sistemele de nivelul 3 și 4.
- Sisteme **RAID S** – proprietar EMC Corporation – sisteme cu paritate și *striping* utilizate în sistemele proprii de memorie Symmetrix.



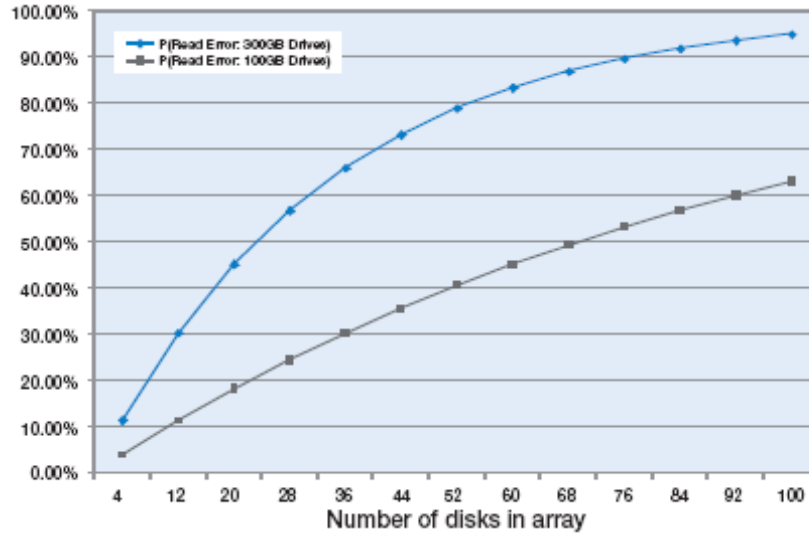
Timpul mediu până la defectare în funcție de capacitatea de memorare
 Structuri de discuri De la 5 la 87 de discuri pe structură.
 Capacitatea discurilor de la 20 la 273 GB. MTBF pentru un disc 10^6 ore.
 Rata erorilor de citire nerecuperate 10^{-14} .

Comparatii de costuri si de performanțe

Primele trei măsuri prin care se evaluează un sistem de discuri sunt fiabilitatea, performanța și costul. Sistemele RAID de la 0 la 6 acoperă o gamă largă de compromisuri între aceste trei măsuri. Este necesar a lua în considerare toate aceste trei măsuri pentru a înțelege deplin valoarea și costul fiecărei organizări a sistemelor de discuri. Despre fiabilitate – în graficul alăturat.

Linia cea mai de jos (romburi) reprezintă timpul mediu până la defectare pentru un singur disc. Într-un sistem fără redundanțe (RAID 0) defectarea produce pierderea datelor. Linia următoare (triunghiuri) arată timpul mediu până la pierderea de date, MTDL (Mean Time to Data Loss) pentru un sistem RAID 5 cu probabilitatea de a găsi un defect latent după reconstituirea informației. De observat că un sistem RAID 5 cu capacitatea totală mai mare de 5 TB ar putea pierde date de mai multe ori în decursul unui an. Pentru a ilustra impactul defectelor latente asupra calculului MTDL, sau timpul mediu până la eroarea adițională, MTAF (Mean Time to Additional Failure), curba următoare (pătrate) arată probabilitatea căderii din cauza a două erori de disc, care ignoră defectele latente, pentru un sistem RAID 5. Ignorarea impactului defectelor latente arată că MTDL pentru un sistem RAID 5 rămâne destul de bun chiar la capacități de peste 5 TB. Linia cea mai de sus (cercuri) arată MTDL pentru un sistem RAID

6 cu luarea în considerare a probabilității de a apărea și defecte latente. Linia arată MTDL pentru sistemele RAID 6, chiar ținând seamă de impactul defectelor latente, care este cu ordine de mărime mai bun decât cel pentru un sistem RAID 5 comparabil.

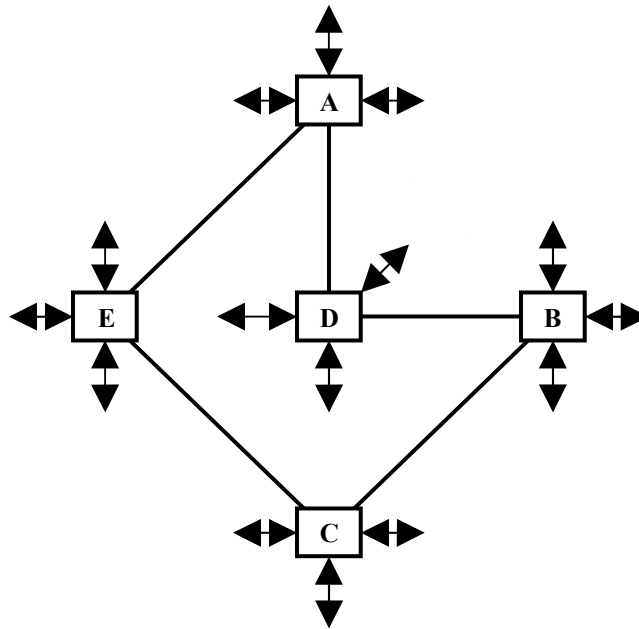


Probabilitatea erorilor de citire irecuperabile în timpul reconstrucției
 Rata erorilor nerecuperabile ale discului: 1 la 10^{14} biti cititi

Pentru înțelegerea mai exactă a modului în care defectele latente din sistemele RAID 5 afectează MTDL să examinăm probabilitatea de a întâlni un defect latent în timpul operației de reconstituire. Dacă un *controller* de RAID 5 întâlnește un defect în timpul reconstituirii, datele utilizatorului sunt pierdute deoarece discul defect și sectorul defect reprezintă două elemente lipsă, ceea ce depășește capacitatea sistemului RAID 5 de a recupera date pierdute. Figura alăturată arată probabilitatea de a găsi un defect latent în timpul reconstrucției sistemului la capacități ale discurilor variate, din ce în ce mai mari. Pentru sisteme foarte mari de discuri de mare capacitate, ar fi surprinzător a nu găsi un defect latent în timpul reconstrucției. Graficul presupune o rată a erorilor tipică pentru discuri din clasa desktop. Probabilitatea este un ordin de mărime mai mică pentru discuri din clasa *enterprise*.

REPLICAREA DATELOR PENTRU TOLERANȚA LA DEFECTE

Cópii identice ale datelor sunt detinute în mai multe noduri ale unui sistem distribuit. O asemenea replicare a datelor poate oferi performante îmbunătățite în ceea ce privește toleranța la defecte.



Replicile datelor trebuie mentinute consistente în pofida oricăror căderi din sistem.

Exemplul din figura de mai sus poate reține cinci copii în cinci noduri. Dacă nodul A este deconectat și un *write* actualizează copia din A , celelalte copii nu vor mai fi consistente cu A . Orice *read* al datelor lor va rezulta în utilizarea de date vechi, depășite. Se pune întrebarea: Câte copii ar trebui citite (scrise)? Răspunsul este dat în continuare pe baza unor scheme consacrate.

O schemă simplă cu voting neierarhizată

Se atribuie v_i voturi copiei i a datelor.

S – mulțimea tuturor nodurilor cu copii ale datelor.

V – suma tuturor voturilor cu $V = \sum_{i \in S} v_i$.

r, w – variabile astfel ca $r + w > V$; $w > V/2$.

$V(X)$ – numărul total de voturi atribuite cōpiilor din multimea X – $V(X) = \sum_{i \in X} v_i$.

Este necesară o strategie care asigură că fiecare *read* utilizează date de ultimă versiune. Pentru a face lectura completă, se citesc noduri dintr-o multime $R \subseteq S$ astfel încât $V(R) \geq r$. Pentru a face scrierea completă, se scrie în fiecare nod al unei multimi $W \subseteq S$ astfel încât $V(W) \geq w$.

Justificarea procedurii este prezentată în continuare.

O multime R pentru care $V(R) \geq r$ este numită cvorum de lectură. O multime W pentru care $V(W) \geq w$ se numeste cvorum de scriere.

Pentru oricare două multimi R și W pentru care $V(R) \geq r$ și $V(W) \geq w$, $R \cap W \neq \emptyset$ (deoarece $r + w > V$).

Orice operație de citire (*read*) este garantată a citi valoarea a cel puțin unei cōpii care a fost actualizată prin cel puțin un *write*. Mai mult, pentru oricare două multimi W_1, W_2 astfel încât $V(W_1), V(W_2) \geq w$, are loc $W_1 \cap W_2 \neq \emptyset$.

Exemplu: Dacă se atribuie un vot fiecărui nod (v. figura de mai sus), atunci suma tuturor voturilor este $V = 5$. Este necesar ca $w > 5/2$ și $r > 5 - w$. Combinațiile permise de r și w sunt (1, 5), (2, 5), (3, 5), (4, 5), (5, 5), (2, 4), (3, 4), (4, 4), (5, 4), (3, 3).

Se consideră spre exemplificare perechea $(r, w) = (1, 5)$. Atunci există un *read* posibil al oricăreia din (cele) cinci cōpii, iar un *write* trebuie să actualizeze fiecare din cele cinci cōpii.

Fiecare operație de citire obține cele mai actuale date.

Dacă $w = 5$, reținerea unui $r > 1$ încetinește operația de citire.

Dacă nodul A este (accidental) deconectat, se poate citi încă din fiecare nod dar nu se pot actualiza toate nodurile.

Se consideră acum $(r, w) = (3, 3)$; numai trei cōpii de scris dar cititul consumă timp mai îndelungat.

Dacă nodul A este (accidental) deconectat, scrierea și lectura în/din A devine imposibilă dar cele patru noduri rămase pot continua cu scrieri și citiri în modul uzual.

Performanță și disponibilitate

Alegerea celor două numere r și w afectează performanța și disponibilitatea sistemului. Dacă sunt de făcut mai multe citiri decât scrieri se alege un r mai mic pentru a accelera operațiile de citire.

Alegerea $r = 1$ impune în exemplul discutat $w = 5$; scrierea nu poate fi executată dacă fie și numai un singur nod este (accidental) deconectat.

Alegerea $r = 2$ permite $w = 4$ și scrierea poate fi încă executată dacă patru din cinci noduri sunt conectate în rețea. Este vorba aici de un *trade-off* între performanță și disponibilitate.

Analiza fiabilității și disponibilității

Disfuncțiile apar în orice nod potrivit unui proces Poisson de rată λ (se admite că legăturile nu cad). Când un nod cade, el este reparat și îi sunt încărcate date actuale. Timpul consumat cu reparatia este aleator, distribuit exponențial cu media $1/\mu$.

Exemplu: $(r, w) = (3, 3)$, ambele operații de citire și de scriere pot avea loc dacă cel puțin trei din cinci noduri sunt funcționale. Pentru a calcula fiabilitatea și disponibilitatea se folosesc pentru modelare lanțurile Markov.

Fiabilitatea pentru cazul $(r, w) = (3, 3)$ urmează etapele de mai jos.

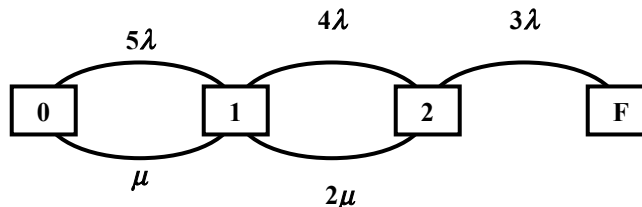
Starea: se definește ca numărul de noduri căzute, disfuncte. Se scriu următoarele ecuații diferențiale:

$$\begin{aligned}\frac{dp_0(t)}{dt} &= -5\lambda p_0(t) + \mu p_1(t) \\ \frac{dp_1(t)}{dt} &= 5\lambda p_0(t) - (4\lambda + \mu)p_1(t) + 2\mu p_2(t) \\ \frac{dp_2(t)}{dt} &= 4\lambda p_1(t) - (3\lambda + 2\mu)p_2(t) \\ \frac{dp_3(t)}{dt} &= 3\lambda p_2(t)\end{aligned}$$

Fiabilitatea la momentul t este

$$R(t) = 1 - P_F(t)$$

(indicele F se referă la disfuncția totală – Failure).

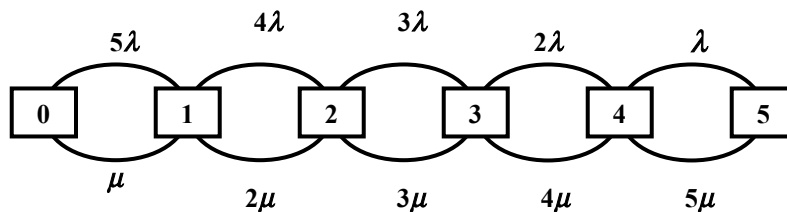


Disponibilitatea pentru același caz, $(r, w) = (3, 3)$ se calculează pe baza ecuațiilor pentru regimul staționar:

$$\begin{aligned}5\lambda \pi_0 - \mu \pi_1 &= 0 \\ -5\lambda \pi_0 + (4\lambda + \mu)\pi_1 - 2\mu \pi_2 &= 0 \\ -4\lambda \pi_1 + (3\lambda + 2\mu)\pi_2 - 3\mu \pi_3 &= 0 \\ -3\lambda \pi_2 + (2\lambda + 3\mu)\pi_3 - 4\mu \pi_4 &= 0 \\ -2\lambda \pi_3 + (\lambda + 4\mu)\pi_4 - 5\mu \pi_5 &= 0 \\ \pi_0 + \pi_1 + \pi_2 + \pi_3 + \pi_4 + \pi_5 &= 1\end{aligned}$$

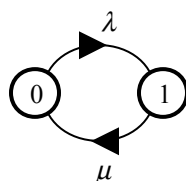
Sistemul este disponibil dacă se află în una din stările 0, 1 sau 2.

Disponibilitatea staționară este $\pi_0 + \pi_1 + \pi_2$



Disponibilitatea poate cunoaste si o a doua tratare. Fiecare nod este considerat separat. Starea unui nod este 0 sau 1, functional, respectiv căzut. Pentru un nod:

$$\pi_0 = \mu/(\lambda + \mu); \pi_1 = \lambda/(\lambda + \mu)$$



Disponibilitatea sistemului A este dată de probabilitatea ca cel puțin 3 noduri să fie functionale. Nodurile sunt presupuse a fi statistic independente si atunci

$$A = \sum_{i=3}^5 C_5^i \left(\frac{\mu}{\lambda + \mu} \right)^i \left(\frac{\lambda}{\lambda + \mu} \right)^{5-i}$$

Exercitiu: Să se calculeze disponibilitatea punctuală (momentană) $PA(t)$ în cazul aceleiasi retele.

În sistemele cu reparatii, interesul principal este legat de disponibilitate (sau de disponibilitatea punctuală) si mai puțin de fiabilitate.

Asignarea votului pentru maximizarea disponibilității

În general, nodurile au fiabilități si disponibilități diferite, variate. În plus, si legăturile pot cădea.

Disponibilitatea punctuală este probabilitatea ca *la momentul t* sistemul să fie functional, să existe cvorumuri de citire si de scriere.

O problemă care comportă discutie o constituie atribuirea de voturi nodurilor pentru a maximiza disponibilitatea punctuală.

Asignarea optimală riguroasă este dificilă. În loc este necesară o euristică.

Mai întâi câteva notatii. Pentru un anumit moment t fixat (t omis pentru simplitatea scrierii):

- Disponibilitatea punctuală pentru nodul $i - a_n(i)$
- Disponibilitatea punctuală pentru legătura $j - a_l(j)$
- $L(i)$ – multimea de conexiuni (legături) la nodul i .

Euristica asignării voturilor

Euristica 1: se pune

$$v(i) = a_n(i) \sum_{j \in L(i)} a_l(j)$$

(cu rotunjire la cel mai apropiat întreg).

Dacă suma tuturor voturilor este pară, se dă un vot în plus unuia din nodurile cu număr maxim de voturi.

Euristica 2: $k(i, j)$ – nod conectat la nodul i prin conexiunea j .

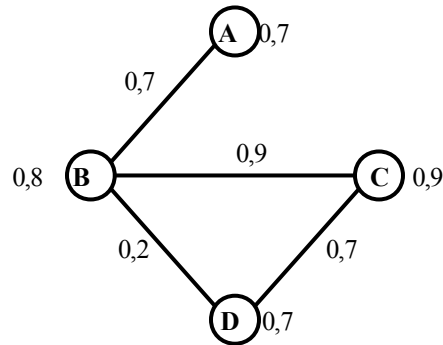
Se pune

$$v(i) = a_n(i) + \sum_{j \in L(i)} a_l(j) a_n[k(i, j)]$$

(cu rotunjire la cel mai apropiat întreg).

Din nou, dacă suma tuturor voturilor este pară, se dă un vot în plus unuia din nodurile cu număr maxim de voturi.

Exemplu pentru euristica 1



$$v(A) = \text{round}(0,7 \times 0,7) = 0$$

$$v(B) = \text{round}(0,8 \times 1,8) = 1$$

$$v(C) = \text{round}(0,9 \times 1,6) = 1$$

$$v(D) = \text{round}(0,7 \times 0,9) = 1$$

Nodul A este nesigur, nefiabil comparativ cu celelalte noduri și de aceea nu primește voturi.

Suma voturilor este 3, cvorumurile trebuie să satisfacă $r + w > 3$; $w > 3/2 \Rightarrow w = 2$ sau $w = 3$.

Dacă $w = 2$, $r = 2$ este cel mai redus cvorum. Cvorumuri posibile pentru *read* (sau *write*): BC, CD, BD.

Dacă $w = 3$, $r = 1$ este cvorumul cel mai redus. Cvorumuri posibile pentru citire: B, C, D. Un singur cvorum pentru scriere: BCD.

Exemplu pentru euristica 2

$$v(A) = \text{round}(0,7 + 0,7 \times 0,8) = 1$$

$$v(B) = \text{round}(0,8 + 0,7 \times 0,7 + 0,9 \times 0,9 + 0,2 \times 0,7) = 2$$

$$v(C) = \text{round}(0,9 + 0,9 \times 0,8 + 0,7 \times 0,7) = 2$$

$$v(D) = \text{round}(0,7 + 0,2 \times 0,8 + 0,7 \times 0,9) = 1$$

Suma voturilor este pară și de aceea B primește un vot suplimentar astfel că asignarea finală este $v(A) = 1$, $v(B) = 3$, $v(C) = 2$, $v(D) = 1$.

Suma voturilor este 7, cvorumurile pentru citire și scriere trebuie să satisfacă relațiile $r + w > 7$, $w > 7/2$, $w \in \{4, 5, 6, 7\}$.

Cvorumuri rezultate din euristica 2, cvorumuri posibile pentru $r + w = 8$.

r	w	Cvorumuri la citire	Cvorumuri la scriere
4	4	AB, BC, BD, ACD	AB, BC, BD, ACD
3	5	B, AC, CD	BC, ABD
2	6	B, C, AD	ABC, BCD
1	7	A, B, C, D	ABCD

Fiecare pereche (r, w) are asociată o disponibilitate, probabilitatea ca cel puțin un cvorum de citire și unul de scriere să existe indiferent de căderile unor noduri și/sau ale unor legături.

Listele de cvorumuri sunt identice pentru citire și scriere în cazul $(r, w) = (4, 4)$.

Alte combinații (r, w) au listele de cvorumuri *read/write* diferite.

Calculul disponibilității pentru cazul $(r, w) = (4, 4)$:

Disponibilitatea A este probabilitatea ca cel puțin unul din cvorumurile AB, BC, BD, ACD să poată fi utilizat.

Marcăm evenimentele E_1, E_2, E_3, E_4 ca fiind respectiv AB, BC, BD, ACD functionale deci utilizabile.

Evenimentele astfel definite nu sunt mutual exclusive.

$$A = \Pr(E_1 \cup E_2 \cup E_3 \cup E_4) =$$

$$= \sum_i \Pr(E_i) - \sum_{i < j} \Pr(E_i \cap E_j) + \sum_{i < j < k} \Pr(E_i \cap E_j \cap E_k) - \Pr(E_1 \cap E_2 \cap E_3 \cap E_4)$$

Calcul (parțial):

$$\Pr(E_1) = \Pr(\text{cvorumul AB utilizabil}) = 0,7 \times 0,7 \times 0,8 = 0,392$$

$$\Pr(E_2 \cap E_3) = \Pr(\text{cvorumurile BC și BD utilizabile}) = 0,8 \times 0,9 \times 0,9 \times 0,2 \times 0,7 = 0,091.$$

Exercițiu: Să se completeze calculele pentru $(r, w) = (4, 4)$.

Urmează o metodă diferită de a calcula disponibilitatea.

Sistemul are 8 module, 4 noduri și 4 conexiuni, fiecare modul poate fi functional sau disfuncțional. Rezultă un total de $2^8 = 256$ de stări mutual exclusive.

Probabilitatea fiecărei stări este un produs de 8 factori, fie $a_n(i)$, fie $[1 - a_n(i)]$, fie $a_i(i)$, fie $[1 - a_i(i)]$.

Un mod sistematic (dar cu consum de timp mare) de a calcula disponibilitatea: se listează toate stările și se adună probabilitățile acelor pentru care există un cvorum.

Pentru orice alte valori (r, w) cvorumurile pentru lectură și pentru scriere sunt diferite.

Se calculează disponibilitatea prin aflarea stărilor în care există atât cvorum pentru scriere cât și cvorum pentru citire.

Asignarea dinamică de voturi

Dacă repararea nu este suficient de rapidă, sistemul se poate degrada. Dacă sistemul se degradează suficient, nu mai există un cluster conex cu o majoritate totală de voturi. Soluția este cea a unor cvorumuri ajustabile în locul celor satice.

Se presupune că fiecare nod are atribuit exact un vot. Pentru fiecare articol-dată sunt menținute numere de versiune, incrementate cu fiecare actualizare. Această incrementare poate fi executată numai dacă se poate aduna un cvorum pentru scriere.

Notatiile pentru asignarea dinamică a voturilor este dată în continuare.

VN_i – numărul versiunii datelor în nodul i .

SC_i – cardinalitatea site-urilor de actualizare la nodul i – numărul de noduri care au participat la actualizarea a VN_i -a a acestor date. Când sistemul începe operația, SC_i este inițializat la numărul total de noduri din sistem.

S_i – mulțimea de noduri cu care nodul i poate comunica.

M – numărul maxim al versiunii din S_i .

I – mulțimea parțială de mulțimi S_i cu nodurile cu numărul de versiune M .

N – cardinalitatea maximă a site-urilor de actualizare (SC_i) de noduri din I .

Algoritmul de asignare dinamică a voturilor este următorul:

1. Dacă la nodul i sosește o cerere de actualizare, nodul i calculează următoarele cantități:
 - $M = \max\{VN_j, j \in S_i\}$ (S_i este mulțimea nodurilor cu care nodul i poate comunica, inclusiv el însuși), adică numărul de versiune maxim al datei de interes peste toate nodurile cu care nodul i poate comunica.
 - $I = \{j | VN_j = M, j \in S_i\}$, adică mulțimea tuturor nodurilor cu numărul de versiune egal cu cel maxim.
 - $N = \max\{SC_j, j \in I\}$, adică maxima cardinalitate a site-urilor de actualizare asociate cu toate nodurile din I .
2. Dacă $\|I\| > N/2$ atunci nodul i poate crește un cvorum de scriere și este admis să execute actualizarea tuturor nodurilor din I . Dacă $S_i - I \neq \emptyset$ atunci toate elementele din $S_i - I$ trebuie să aibă versiunea lor și numerele de versiune aduse la actualitate înainte ca ele să poată face această actualizare. Suplimentar, numărul de versiune al fiecărei copii a acelei date din S_i este incrementat, adică VN_j este incrementat pentru fiecare $j \in S_i$. Totodată, pentru fiecare $j \in S_i$, se pune $SC_j = \|I\|$. Acest pas în totalitate trebuie executat atomic: toate aceste operații trebuie făcute în fiecare nod din S_i sau deloc, în nici unul dintre ele.

Exemplu de atribuire dinamică a voturilor: cazul cu șapte noduri, aceleași date, cu starea la momentul t_0

	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>	<i>G</i>
<i>VN</i>	5	5	5	5	5	5	5
<i>SC</i>	7	7	7	7	7	7	7

O disfuncție face mulțimea nodurilor neconexă: cele două submulțimi disjuncte sunt $\{A, B, C, D\}$ și $\{E, F, G\}$.

E primește o solicitare de actualizare – $SC_E = 7$. *E* trebuie să găsească încă cel puțin 7/2 noduri (inclusiv el însuși) dar nu poate găsi mai mult de 3: cererea de actualizare este refuzată.

A primește o cerere de actualizare. Cererea poate fi onorată și *A*, *B*, *C*, *D* sunt actualizate.

Starea nouă rezultată este

	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>	<i>G</i>
<i>VN</i>	6	6	6	6	5	5	5
<i>SC</i>	4	4	4	4	7	7	7

Apare o altă disfuncție și submulțimile disjuncte devin $\{A, B, C\}$, $\{D\}$, $\{E, F, G\}$.

O cerere de actualizare sosește la *C*, cvorumul de scriere este 3 în nodul *C*, actualizarea se face cu succes.

Rezultă o stare nouă:

	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>	<i>G</i>
<i>VN</i>	7	7	7	6	5	5	5
<i>SC</i>	3	3	3	4	7	7	7

Votarea poate fi organizată ierarhic. Dacă V este mare, $r + w$ este de asemenea mare, operațiile făcute cu datele consumă timp îndelungat. O soluție posibilă este adoptarea unei scheme de votare ierarhice.

Pentru aceasta se construiește un arbore cu m niveluri. Toate nodurile care detin copii ale datelor sunt frunze la nivelul $m - 1$. Se adaugă noduri virtuale la nivelurile mai înalte până la nivelul 0, cel al rădăcinii. Nodurile adăugate sunt grupări virtuale de noduri reale. Fiecare nod din nivelul i va avea exact L_{i+1} descendenți.

Un exemplu de arbore pentru generare ierarhică de cvorumuri este dat în figura alăturată.

Algoritmul de generare a cvorumurilor este următorul:

Se atribuie câte un vot fiecărui nod din arbore.

Se fixează cvorumul de citire și cvorumul de scriere la nivelul i , r_i și w_i astfel încât $r_i + w_i > L_i$; $w_i > L_i/2$.

Mai departe algoritmul se utilizează recursiv.

Se marchează “de citire” rădăcina la nivelul 0.

La nivelul 1 se marchează “de citire” r_1 noduri.

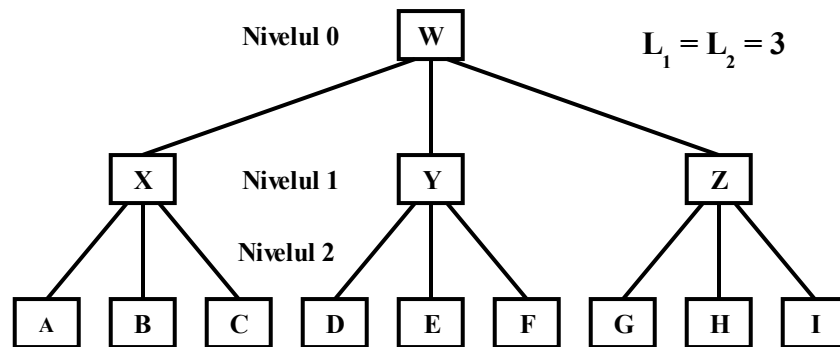
Se trece de la nivelul i la nivelul $i + 1$, se marchează “de citire” r_{i+1} descendenți ai fiecărui nod marcat “de citire” la nivelul i .

Nu se poate marca “de citire” un nod care nu are cel puțin r_{i+1} descendenți fără defecte.

Se continuă până când $i = m - 1$.

Frunzele marcate “de citire” alcătuiesc cvorumul de citire.

Formarea unui cvorum de scriere este similară.



Exemplu numeric pentru $i = 1, 2$, $w_i = 2$, $r_i = L_i - w_i + 1 = 2$.

Pornind de la rădăcină se marchează “de citire” doi din descendenții săi, să spunem X și Y .

Se marchează “de citire” doi descendenți pentru X și Y , de pildă A, B pentru X și D, E pentru Y .

Cvorumul de citire este mulțimea frunzelor marcate “de citire”: A, B, D, E .

Se presupune acum că D este defect. În această stare D nu poate fi parte a cvorumului de citire.

Trebuie ales un alt descendent al lui Y , de pildă F , pentru un nou cvorum de citire.

Dacă doi din descendenții lui Y sunt defecti, Y însuși nu poate fi marcat “de citire”; se revine și se încearcă marcarea “de citire” a lui Z .

Exercițiu: Să se listeze cvorumurile generate de $r_1 = 1$, $w_1 = 3$, $r_2 = 2$, $w_2 = 2$.

Relativ la tratarea ierarhică și non-ierarhică se pot face unele observații. Astfel: Cvorumul de citire constă în numai 4 copii. Similar, putem avea un cvorum de scriere cu 4 copii.

Pentru tratarea lipsită de ierarhizare, cu un vot pe nod, $r + w > 9$, $w > 9/2$. w este cel puțin 4, comparativ cu exact 4 în tratarea ierarhizată pe arbore.

Pentru a dovedi că tratarea ierarhică lucrează, se poate arăta că orice cvorum de citire posibil trebuie să intersecteze orice cvorum de scriere în cel puțin un nod.

Solutia cu backup primar

Un nod este desemnat ca nod primar; toate accesările se produc prin acel nod. Celelalte noduri sunt desemnate ca noduri de backup.

În operarea normală, toate scrierile pe nodul primar sunt repetate și pe backup-urile functionale. Dacă nodul primar cade, unul din nodurile de backup este ales pentru a-l înlocui.

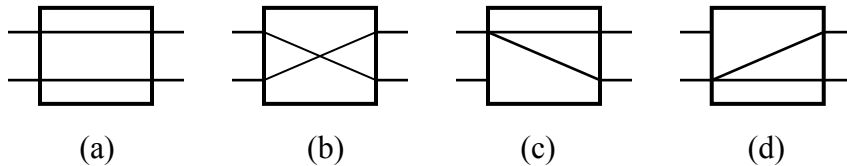
TOLERANȚA LA DEFECTE ÎN REȚELE

O particularitate a rețelilor este aceea că aproape todeauna există căi multiple care conectează sursa unui mesaj cu destinația lui. Totodată există noduri de rezervă care pot fi conectate în rețea pentru a înlocui unitățile disfuncte.

Literatura menționează câteva topologii tolerante la defecte:

- Rețelele în trepte multiple (*multi-stage*), unele dintre ele suplimentare
- Sitele (*meshes*) interstiale
- *Crossbar*-ul cu redundante
- Hipercurburile
- Rețelele punct-la-punct

Rețelele multi-stage lipsite de toleranță la defecte (rețele fluture) sunt alcătuite tipic din comutatoare 2×2 , comutatoare cu două intrări și două ieșiri.



Un comutator poate avea una din cele patru setări figurate mai sus:

- Directă – linia de intrare superioară conectată la linia superioară de ieșire, linia de intrare inferioară conectată la linia inferioară de ieșire.
- Încrucisată – linia de intrare superioară conectată la linia inferioară de ieșire, linia de intrare inferioară conectată la linia superioară de ieșire.
- Cu difuzare (broadcast) superioară – linia de intrare superioară conectată la ambele linii de ieșire.
- Cu difuzare (broadcast) inferioară – linia de intrare inferioară conectată la ambele linii de ieșire.

Rețelele fluture sunt rețele în k trepte, cu $k \geq 3$. Au 2^k intrări și 2^k ieșiri. Treptele au fiecare 2^{k-1} comutatoare. Conexiunile urmează o anumită regulă recursivă, de la intrare către ieșire.

În stratul de intrare, linia superioară a fiecărui comutator este conectată la intrările unui fluture $2^{k-1} \times 2^{k-1}$ și linia de ieșire inferioară a fiecărui comutator este conectată la intrările unui alt fluture $2^{k-1} \times 2^{k-1}$. În particular, un fluture cu două straturi, pentru stratul de intrare linia de ieșire superioară a fiecărui comutator este conectată la un comutator 2×2 și linia de ieșire inferioară este conectată la alt comutator 2×2 . Fluturele cel mai simplu, cu un strat constă dintr-un singur comutator 2×2 .

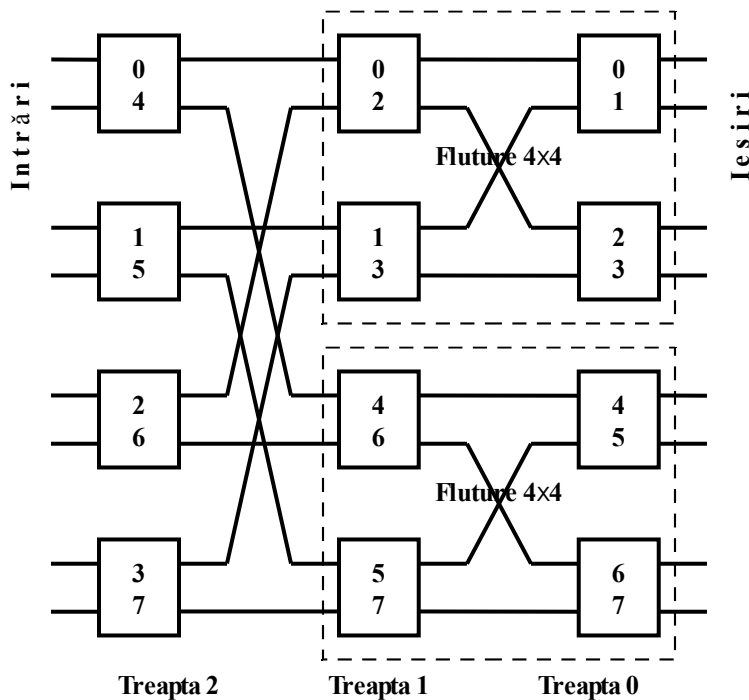


Figura alăturată reprezintă o rețea fluture mai complexă, pe care se pot observa unele detalii și se pot preciza unele notații.

Un comutator din treapta i are liniile numerotate cu numere care diferă cu 2^i . Linia de ieșire j a fiecărei trepte merge la linia de intrare j a stratului următor ($j = 0, \dots, 2^{k-1}$).

Numerele asociate oricărui comutator (*switchbox*), cu excepția celor din stratul de ieșire, sunt ambele de aceeași paritate (pare sau impare).

O rețea fluture nu este tolerantă la defecte: există o singură cale de la oricare dintre intrări la o anumită ieșire. Dacă un comutator din stratul i clachează, 2^{k-i} intrări sunt deconectate de la 2^{i+1} ieșiri.

Sistemul poate încă opera dar într-o manieră degradată.

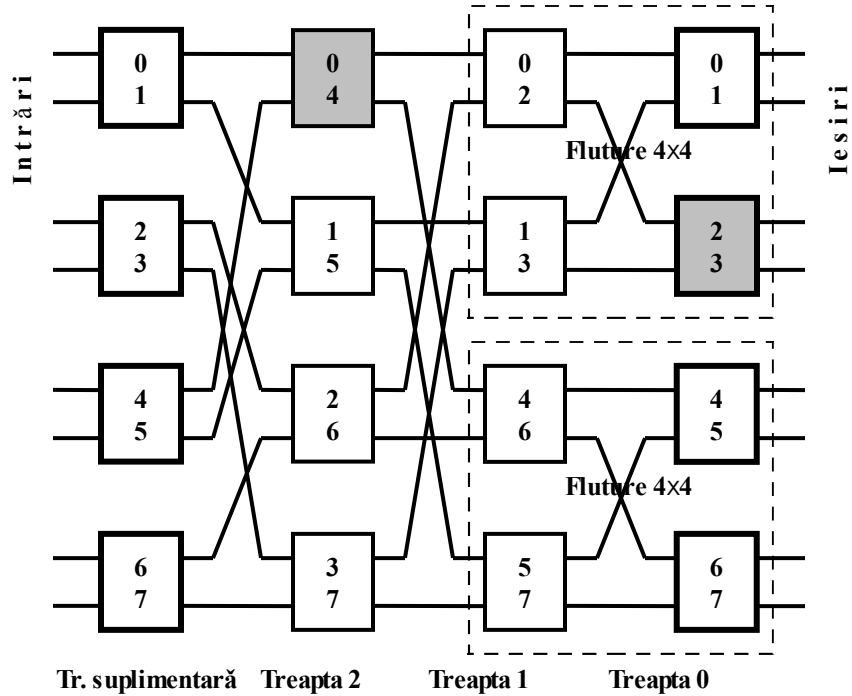
Pentru a crea o rețea tolerantă la defecte se utilizează rețele cu trepte suplimentare.

O posibilitate constă în a adăuga o extra-treaptă prin duplicarea treptei de intrare. Este necesară și o multiplexare în scopul *bypass*-ării comutatoarelor din straturile/treptele de la intrare și de la ieșire. Prin adoptarea acestei soluții, un comutator disfuncțional poate fi ocolit prin rutarea pe *bypass*.

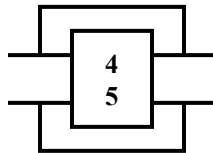
Exemple:

Comutatorul din stratul/treapta 0 cu liniile 2, 3 căzute este duplicat printr-o treaptă suplimentară. Comutatorul disfuncțional este ocolit (*bypass*-at) cu concursul unui multiplexor.

Comutatorul din stratul 2 cu liniile 0, 4 căzute: extra-stratul este setat astfel ca linia de intrare 0 să fie comutată la linia de iesire 1 si linia de intrare 4 la linia de iesire 5, prin bypassarea cutiei de comutare disfuncte.



Comutatoarele din treapta suplimentară si din treapta ultimă (0), reprezentate cu contur îngrosat, trebuie “citite” fiecare cu posibilitățile de ocolire (bypass), ca înfigura imediat următoare.



Se propune ca exercitiu demonstrarea faptului că rețeaua cu o treaptă suplimentară poate rămâne conexă în pofida disfuncției unei cutii de comutare situată oriunde în sistem.

Măsuri ale siguranței (*dependability*) unei rețele cu mai multe straturi.

Retelele cu interconectare în mai multe trepte conectează N procesoare la N unități de memorie într-o arhitectură cu memorie partajată ($N = 2^k$). În prezența elementelor cu defecte, sistemul poate opera într-un mod degradat.

Reziliența sistemului în degradare progresivă poate fi măsurată. Iată măsurile uzuale pentru reziliență:

- Lărgimea de bandă (sau banda de trecere)
- Numărul mediu de căi operationale
- Metrici ale conectivității între procesoare și memorii

Toate măsurile sunt funcții de timp și presupun că defectele apar și pot fi reparate în intervalul $[0, t]$.

Definițiile complete și detaliate ale măsurilor enumerate mai sus sunt:

Banda de trecere $BW(t)$ – numărul mediu statistic (*expected*) de procesoare la momentul t , care comunică cu (parte din) memorie.

Conectivitatea $Q(t)$ – numărul mediu statistic (*expected*) de căi procesoare-memorii operationale la momentul t ; o cale operatională include un procesor, o memorie și legăturile dintre ele, toate lipsite de defecte.

Un procesor (memorie) este accesibil(ă) (la momentul t) dacă este lipsit(ă) de defecte și este conectat(ă) la cel puțin o memorie (un procesor) lipsit(ă) de defecte.

O măsură suplimentară a conectivității este cuplul $(A_r(t), A_m(t))$ alcătuit din numărul mediu de procesoare și de memorii accesibile la momentul t .

O altă măsură este dată de perechea $(N_m(t), N_r(t))$ formată din numărul mediu de procesoare (memorii) fără defecte, la care o memorie (un procesor) accesibil este conectat(ă) la timpul t .

Sunt de observat câteva imperfecțiuni ale acestor măsuri.

Banda de trecere, $BW(t)$ nu depinde numai de condițiile rețelei ci și de volumul solicitării de memorie din partea procesoarelor.

Conectivitatea $Q(t)$ prin numărul de căi nu spune câte procesoare și câte memorii distincte sunt încă accesibile.

Perechea $(A_r(t), A_m(t))$ nu cuprinde faptul că există o rețea de interconectare total conexă și fără defecte $A_r(t) \times A_m(t)$; nu indică nici numărul de memorii fără defecte care sunt conectate în medie la un procesor accesibil.

Prin combinarea măsurilor $Q(t)$ și $(A_r(t), A_m(t))$ se obține o caracterizare mai completă a sistemului.

Dacă avem în vedere relațiile

$$N_m(t) = Q(t)/A_r(t) \text{ și } N_r(t) = Q(t)/A_m(t)$$

cuplul $(N_m(t), N_r(t))$ reprezintă o margine superioară pentru sistemul operațional mediu maximal deplin conex la momentul t .

Analiza siguranței (*dependability*)

O premisă importantă: timpul mediu între căderile (MTBF) componentelor (și posibilele reparatii) este mult mai mare decât durata medie de comunicare între un procesor și o memorie.

O altă premisă importantă: starea componentelor sistemului (cu defecte sau fără) este constantă pentru o perioadă de timp suficient de lungă, atât de lungă cât să permită analiza comportării sistemului într-o stare statistic staționară.

Sistemul este observat la un moment arbitrar t fixat pentru întreaga analiză. Toate măsurile sunt funcții de t , inclusiv probabilitățile $p_r(t)$, $p_m(t)$, $p_l(t)$ de bună funcționare pentru procesoare, memorii, legături. De regulă, pentru simplitate, timpul t este omis din notații (se scrie doar p_r , p_m , p_l).

Se notează cu p_q probabilitatea ca un procesor să aibă nevoie de o legătură cu memoria.

Analiza benzii de trecere

Banda de trecere BW este, cum s-a mai spus, numărul mediu statistic (*expected*) de procesoare în comunicare activă cu o (parte din) memorie. Se admite o ipoteză simplificatoare: destinațiile cererilor de memorie din partea procesoarelor sunt independente statistic și uniform distribuite pe cele N memorii.

În condițiile specificate, banda de trecere a rețelei este produsul numărului de memorii N cu Ψ_m , probabilitatea ca o memorie dată (de pildă memoria 0) să fie lipsită de defecte și să aibă o cerere la intrare ei.

Probabilitățile Ψ_m se calculează iterativ, urmând calea care duce la acea memorie desemnată de indicele m . Probabilitatea unei cereri pe o legătură de ieșire a unui comutator se calculează din probabilitatea ca o solicitare să fi fost acceptată la legăturile de intrare ale acestui comutator (*switch*).

Calculul benzii de trecere urmează în linii mari schema explicată în continuare.

O legătură este în starea $X = 1$ ($X = 0$) dacă are (nu are) o solicitare pentru memoria specificată; o legătură defectă este în starea $X = 0$.

Atribuirea de numere celor $k + 1$ straturi ale rețelei ($k = \log_2 N$).

- Stratul 0 este ultimul strat; legăturile de ieșire sunt conectate la memorii.
- Stratul k este primul și preia ieșirile procesoarelor.

X^i , Y^i sunt ieșirile unui comutator din stratul/treapta i .

X^{i+1} , Y^{i+1} sunt intrările unui comutator din stratul i , în același timp ieșiri a două dintre comutatoarele (diferite) din stratul $i + 1$.

Banda de trecere în cazul rețelelor fără redundanțe

Solicitările de memorie sunt distribuite uniform între memorii; o solicitare care sosește este rutată la oricare legătură de ieșire din cele două ale unui comutator, cu probabilitate egală (0,5). În calculul probabilității $\Pr(X^i = 1)$, este suficient a lua în considerare numai una din legăturile de ieșire.

O cerere către un modul de memorie poate atinge legătura de ieșire a unui comutator pe oricare dintre cele două legături de intrare

$$\begin{aligned} \Pr(X^i = 1) &= \sum_{u,v=0,1} \Pr(X^i = 1/X^{i+1} = u, Y^{i+1} = v) \Pr(X^{i+1} = u, Y^{i+1} = v) = \\ &= 0 \cdot \Pr(X^{i+1} = 0, Y^{i+1} = 0) + (1/2)p_l \Pr(X^{i+1} = 0, Y^{i+1} = 1) + \\ &+ (1/2)p_l \Pr(X^{i+1} = 1, Y^{i+1} = 0) + (p_l + (1/4)p_l^2) \Pr(X^{i+1} = 1, Y^{i+1} = 1) \end{aligned}$$

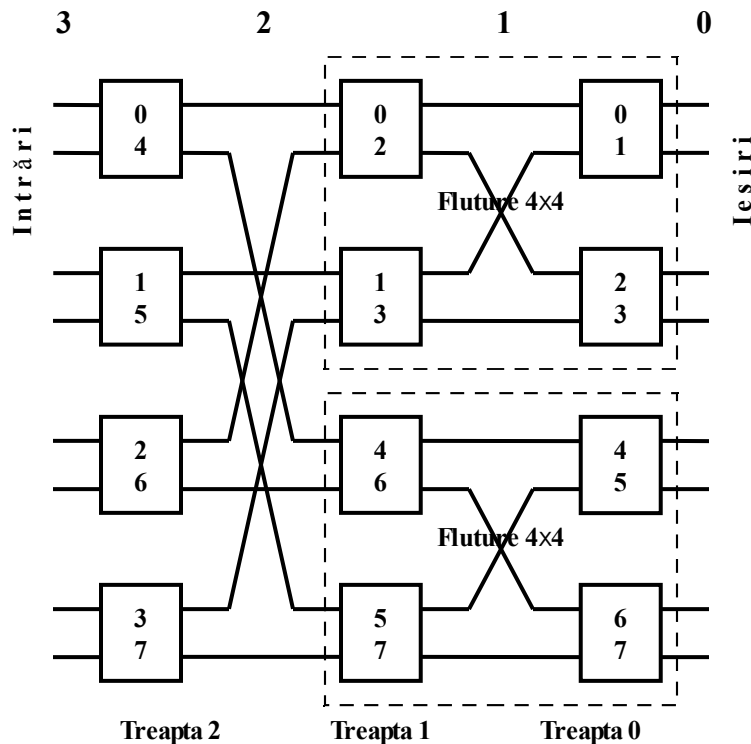
Sunt luate în considerare numai defectele legăturilor de intrare. Defectele legăturilor de iesire sunt considerate defecte ale legăturilor de intrare ale etapei următoare.

Stările intrărilor unui comutator sunt presupuse a fi independente statistic:

$$\Pr(X^i = u, Y^i = v) = \Pr(X^i = u) \Pr(Y^i = v)$$

$$\Pr(X^i = 0) = \Pr(Y^i = 0) = 1 - \Pr(X^i = 1)$$

După câteva prelucrări algebrice nu foarte complicate se poate arăta că $\Pr(X^k = 1) = p_q p_r$.



Probabilitatea $\Pr(X^0 = 1)$ se calculează recursiv.

Memoria si legătura ei de intrare pot fi fără defecte: probabilitatea unei astfel de stări este

$$\Psi_m = \Pr(X^0 = 1) p_i p_m$$

si în final

$$BW = N \Psi_m$$

Conectivitatea rețelelor de interconectare fără redundanțe

Cum s-a definit, Q este numărul mediu de căi operationale pentru perechi procesor-memorie conectate. În lipsa oricăror redundanțe există exact o cale între un procesor si o memorie.

În cuvinte, conectivitatea Q este produsul numărului de perechi procesor-memorie cu probabilitatea existenței unei căi fără defecte. Probabilitatea aceasta este $p_r p_l^{k+1} p_m$, cu $k + 1$ numărul de legături pe cale, adică numărul de trepte + 1 ($k = \log_2 N$).

Deoarece numărul de căi procesor-memorie este N^2 , rezultă

$$Q = N^2 p_r p_l^{k+1} p_m.$$

Calculul măsurilor adiționale pentru rețelele de interconectare fără redundanțe

Am notat cu A_r numărul mediu de procesoare accesibile. A_r este produsul numărului de procesoare N cu probabilitatea ϕ_r ca un procesor (de pildă procesorul 0) să fie accesibil.

O legătură este în starea $X = 0$ ($X = 1$) dacă toate (nu toate) căile de la procesor la memorie sunt defecte.

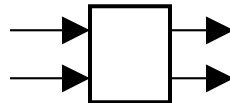
- O cale defectă este o cale cu cel puțin o legătură defectă.
- O legătură defectă este în starea $X = 0$.

Numerotarea treptelor se face și de această dată de la k la 0; X^i exprimă starea legăturii în treapta i . În particular, $\Pr(X^0 = 1) = p_m$; $\Pr(X^0 = 0) = 1 - p_m$.

O ecuație recursivă

$$\Pr(X^{i+1} = 1) = p_l [1 - \Pr(X^i = 0)^2]$$

este suportul evaluărilor curente.



În cele din urmă

$$\phi_r = p_r p_l \Pr(X^k = 1)$$

și

$$A_r = N \phi_r$$

Numărul mediu de memorii accesibile, A_m se obține pe o cale similară, prin schimbarea probabilității p_r cu probabilitatea p_m .

Urmează un exemplu numeric.

Calculul benzii de trecere, în ipoteza unor legături fără defecte, în cazul $N = 8$, $k = 3$, pentru un moment t fixat.

$p_r = 0,8$; $p_m = 0,9$; $p_l = 1$ (legături fără defecte); $p_q = 0,7$.

Calculul benzii de trecere:

$$\Pr(X^3 = 1) = p_q p_r = 0,56$$

$$\Pr(X^2 = 1) = 0,56 - 0,25 \times 0,562 = 0,536$$

$$\Pr(X^1 = 1) = 0,536 - 0,25 \times 0,5362 = 0,464$$

$$\Pr(X^0 = 1) = 0,464 - 0,25 \times 0,4642 = 0,41$$

$$BW = 0,41 N p_m = 0,41 \times 8 \times 0,9 = 2,95$$

Calculul conectivității și al măsurilor suplimentare:

$$Q = N^2 \times 0,8 \times 0,9 = 0,72N^2 = 46,08$$

$$A_r = Np_r[1 - (1 - p_m)^N] = 0,8N(1 - 0,1^N) \approx 0,64$$

$$A_m = Np_m[1 - (1 - p_r)^N] = 0,9N(1 - 0,2^N) \approx 0,72$$

$$N_r = Q/A_m = 64$$

$$N_m = Q/A_r = 72$$

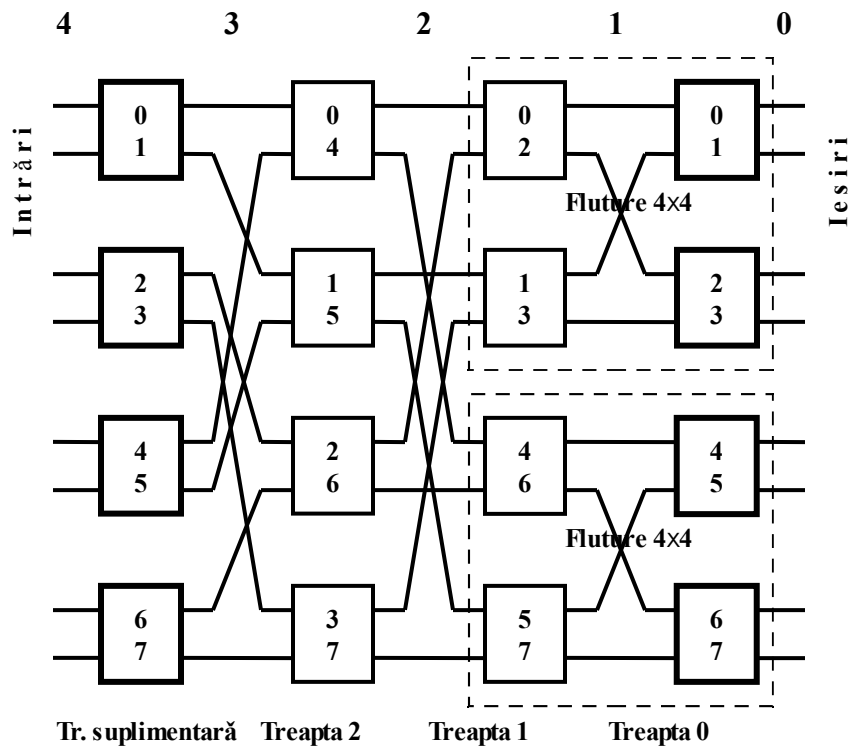
Asupra acestor ultime rezultate sugerăm cititorului o discuție.

Retele fluture și rețele fluture cu trepte suplimentare

În rețeaua fluture fără redundante, cele două intrări în orice comutator sunt considerate independente statistic. Într-o rețea cu o treaptă în plus sunt câte două căi care conectează orice pereche procesor-memorie. Legăturile sunt de data aceasta dependente și ecuațiile date mai sus nu mai sunt valide.

Așa cum s-a mai spus, primul și ultimul strat au multiplexoare și demultiplexoare pentru care analiza este diferită de aceea a etajelor interioare.

Sunt patru legături care duc la două comutatoare, două perechi sunt independente, deși legăturile din aceeași pereche sunt dependente.



Exemplu: legăturile de iesire 0 si 1 din etajul 2 sunt dependente (procesoarele 0 si 1 trimit solicitări către memoria 0 prin ambele); legăturile 2 si 3 sunt si ele dependente; perechile 0, 1 si 2, 3 sunt însă independente.

Banda de trecere pentru o rețea cu trepte suplimetare

Banda de trecere (BW), reamintim, este numărul mediu de procesoare care comunică activ cu (o parte din) memorie si este același lucru cu numărul mediu de memorii care comunică activ cu vreun procesor. Banda de trecere se obtine ca produsul numărului de memorii N cu Ψ_m , probabilitatea ca o memorie dată (de pildă memoria 0) să fie fără defect si să aibă o solicitare la intrarea ei.

Ca si altădată, Ψ_m se calculează iterativ, urmând o cale de la procesor care duce la o anumită memorie.

Legătura este în starea 1 (0) dacă ea are (nu are) o solicitare de memorie; o legătură defectă este în starea 0.

Calculul benzii de trecere pentru rețeaua din figură urmează pașii prezentați imediat.

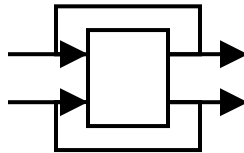
Sunt $k + 2$ trepte numerotate $k + 1, k, \dots, 0$ (cu $k = \log_2 N$). Se notează

- X^i, Y^i – starea celor două legături de iesire din etajul i
- $X^{i+1}, Y^{i+1}, Z^{i+1}, W^{i+1}$ – starea intrărilor pe legăturile din etajul i , totuna cu legăturile de iesire pentru etajul $i + 1$.

Probabilitatea ca o intrare din etajul i să aibă solicitare este calculată pe baza probabilității ca o solicitare să fie acceptată la legăturile de intrare.

Pentru primul etaj ($k + 1$ – etajul procesorului) se poate scrie

$$\Pr[X^{k+1} = 1] = p_q p_r; \Pr[X^{k+1} = 0] = 1 - \Pr[X^{k+1} = 1]$$



Pentru etajul k (procesoarele sunt independente statistic):

$$\begin{aligned} & \Pr[(X^k, Y^k) = (0, 0)] = \\ & = (\Pr[X^{k+1} = 0])^2 + q_i^4 (\Pr[X^{k+1} = 1])^2 + q_i^3 2 \Pr[X^{k+1} = 0] \Pr[X^{k+1} = 1] \\ & \Pr[(X^k, Y^k) = (0, 1)] = \\ & = (1 - q_i^3) \Pr[X^{k+1} = 0] \Pr[X^{k+1} = 1] + (1 - q_i^2) q_i^2 (\Pr[X^{k+1} = 1])^2 \\ & \Pr[(X^k, Y^k) = (1, 0)] = \Pr[(X^k, Y^k) = (0, 1)] \\ & \Pr[(X^k, Y^k) = (1, 1)] = (1 - q_i^2)^2 (\Pr[X^{k+1} = 1])^2 \end{aligned}$$

Treptele interne în rețelele cu trepte suplimentare

Expresiile de mai devreme au presupus că o solicitare este trimisă mai întâi prin conexiunea directă și se folosește conexiunea încrucișată numai dacă cea directă este disfuncțională.

- Protocol diferit: conexiunile directă și încrucișată cu probabilități egale, expresiile pentru probabilități vor fi diferite.

Pentru etajele interioare ($i = k - 1, \dots, 1$):

$$\begin{aligned} \Pr[(X^i, Y^i) = (u, v)] &= \\ &= \sum_{\substack{(s_0, s_1, s_2, s_3) = (1, 1, 1, 1) \\ (s_0, s_1, s_2, s_3) = (0, 0, 0, 0)}} \Pr[(X^{i+1}, Y^{i+1}, Z^{i+1}, W^{i+1}) = (s_0, s_1, s_2, s_3)] \\ &\cdot \Pr[(X^i, Y^i) = (u, v) | (X^{i+1}, Y^{i+1}, Z^{i+1}, W^{i+1}) = (s_0, s_1, s_2, s_3)] = \\ &= \sum_{\substack{(s_0, s_1, s_2, s_3) = (1, 1, 1, 1) \\ (s_0, s_1, s_2, s_3) = (0, 0, 0, 0)}} \Pr[(X^{i+1}, Y^{i+1}) = (s_0, s_1)] \Pr[(Z^{i+1}, W^{i+1}) = (s_2, s_3)] \\ &\cdot \Pr[X^i = u | (X^{i+1}, Z^{i+1}) = (s_0, s_2)] \Pr[Y^i = v | (Y^{i+1}, W^{i+1}) = (s_1, s_3)] \end{aligned}$$

pentru $u, v = 0, 1$.

Numai probabilitățile combinate (*joint*) ale celor două legături sunt necesare. Acestea pot fi calculate recursiv de la etajul $k + 1$ (etajul procesor) la etajul 0 (etajul de memorare).

Etajul 0 include demultiplexoare

$$\begin{aligned} \Pr[X^0 = 1 | (X^1, Y^1) = (0, 0)] &= 0 \\ \Pr[X^0 = 1 | (X^1, Y^1) = (0, 1)] &= (1/2)p_l \\ \Pr[X^0 = 1 | (X^1, Y^1) = (1, 0)] &= (1/2)(1 - q_l^2) \\ \Pr[X^0 = 1 | (X^1, Y^1) = (1, 1)] &= (1/2)(3p_l - p_l^2) - (1/4)p_l(1 - q_l^2) \\ \Psi_m &= \Pr(X^0 = 1)p_l p_m \end{aligned}$$

În final:

$$BW = N\Psi_m$$

Conectivitatea pentru o rețea cu trepte suplimentare

Q este produsul numărului de perechi procesor-memorie N^2 cu probabilitatea ca între componentele perechii să existe cel puțin o cale fără defecte.

Fiecare pereche procesor-memorie este conectată prin două căi disjuncte (se-nțelege, cu excepția celor două capete).

Probabilitatea ca cel puțin o cale să fie fără defecte este egală cu probabilitatea ca prima cale să fie fără defecte adunată cu probabilitatea ca cealaltă cale să fie fără defecte din care trebuie scăzută probabilitatea ca ambele căi să fie fără defecte.

Probabilitatea poate asuma una din cele două expresii (a se compara calea între procesorul 0 și memoria 0 cu calea între procesorul 0 și memoria 1).

Calculul conectivității urmează pașii de mai jos.

Pentru căile dintre procesorul 0 și memoria 0:

$$\begin{aligned} \Pr(\text{cel puțin o cale este fără defecte}) &= \Pr(0, 0) = \\ &= 2p_r(1 - q_i^2)p_i^{k+1}p_m - p_r p_i^{2k+2}(1 - q_i^2)^2 p_m \end{aligned}$$

Pentru căile dintre procesorul 0 și memoria 1:

$$\begin{aligned} \Pr(\text{cel puțin o cale este fără defecte}) &= \Pr(0, 1) = \\ &= p_r(1 - q_i^2)p_i^k(1 - q_i^2)p_m + p_r p_i^{k+2}p_m - p_r p_i^{2k+2}(1 - q_i^2)^2 p_m \end{aligned}$$

Jumătate din perechile procesor-memorie urmează valoarea $\Pr(0, 0)$ și cealaltă jumătate urmează valoarea $\Pr(0, 1)$.

$$Q = [\Pr(0, 0) + \Pr(0, 1)]N^2/2$$

Măsurile adiționale pentru o rețea cu trepte suplimentare

A_r și A_m sunt numărul mediu de procesoare accesibile, respectiv numărul mediu de memorii accesibile.

ϕ_r (ϕ_m) sunt probabilitățile ca un procesor (o memorie) dat(ă) să fie conectat(ă) la cel puțin o memorie (un procesor).

Pentru calcularea lui A_r se face aceeași descriere de stări: legătura este în starea $X = 0$ ($X = 1$) dacă toate (nu toate) căile de la procesor la memorii sunt defecte.

O cale defectă este o cale care conține cel puțin o legătură defectă.

O legătură defectă este în starea $X = 0$.

Numerotarea etajelor se menține, $k + 1$ (procesoarele) la 0 (memoriile).

Dacă X^i descrie starea legăturii din etajul i

$$\phi_r = p_r p_l \Pr(X^{k+1} = 1)$$

și

$$A_r = N \phi_r$$

Urmează calculul măsurii A_r .

$\Pr(X^i = 1)$ se calculează recursiv de la treapta 0 la treapta $k + 1$.

X^i , Y^i notează și acum starea celor două legături din etajul i .

Pentru etajul 0:

$$\Pr(X^0 = 1) = p_m \text{ și } \Pr(X^0 = 0) = 1 - p_m.$$

Pentru etajul 1:

$$\begin{aligned} \Pr[(X^1, Y^1) = (0, 0)] &= \{\Pr[X^0 = 0]\}^2 + 2\Pr[X^0 = 0]\Pr[X^0 = 1]q_i^3 + \\ &+ \{\Pr[X^0 = 1]\}^2 q_i^6 \end{aligned}$$

$$\begin{aligned} \Pr[(X^1, Y^1) = (0, 1)] &= \Pr[X^0 = 0]\Pr[X^0 = 1][q_i(1 - q_i^2) + q_i^2 p_l] + \\ &+ \{\Pr[X^0 = 1]\}^2 q_i^3(1 - q_i^3) \end{aligned}$$

$$\Pr[(X^1, Y^1) = (1, 0)] = \Pr[(X^1, Y^1) = (0, 1)]$$

$$\begin{aligned} \Pr[(X^1, Y^1) = (1, 1)] &= 2\Pr[X^0 = 0]\Pr[X^0 = 1]p_l(1 - q_i^2) + \\ &+ \{\Pr[X^0 = 1]\}^2(1 - q_i^3)^2 \end{aligned}$$

Pentru stările 2, ..., k , variabilele X^{i-1} , Y^{i-1} , Z^{i-1} , W^{i-1} exprimă starea celor patru legături din etajul $i - 1$.

$$\Pr[(X^i, Y^i) = (u, v)] = \sum_{s_0, \dots, s_3 = 0, 0, 0, 0}^{1, 1, 1, 1} \Pr[(X^{i-1}, Y^{i-1}) = (s_0, s_1)].$$

$$\cdot \Pr[(Z^{i-1}, W^{i-1}) = (s_2, s_3)] \cdot \Pr[X^i = u | (X^{i-1}, Z^{i-1}) = (s_0, s_2)].$$

$$\cdot \Pr[Y^i = v | (Y^{i-1}, W^{i-1}) = (s_1, s_3)]$$

Probabilitățile conditionate sunt:

$$\Pr[X^i = 0 | (X^{i-1}, Z^{i-1}) = (0,0)] = 1$$

$$\Pr[X^i = 0 | (X^{i-1}, Z^{i-1}) = (0,1)] = q_l$$

$$\Pr[X^i = 0 | (X^{i-1}, Z^{i-1}) = (1,0)] = q_l$$

$$\Pr[X^i = 0 | (X^{i-1}, Z^{i-1}) = (1,1)] = q_l^2$$

Pentru etajul suplimentar $k + 1$:

$$\Pr[X^{k+1} = 0] = \Pr[(X^k, Y^k) = (0,0)] + \Pr[(X^k, Y^k) = (0,1)](q_l + q_l^2) + \Pr[(X^k, Y^k) = (1,1)]q_l^3$$

$$\Pr[X^{k+1} = 1] = 1 - \Pr[X^{k+1} = 0]$$

$$\phi_r = p_r p_l \Pr[X^{k+1} = 1]$$

În final

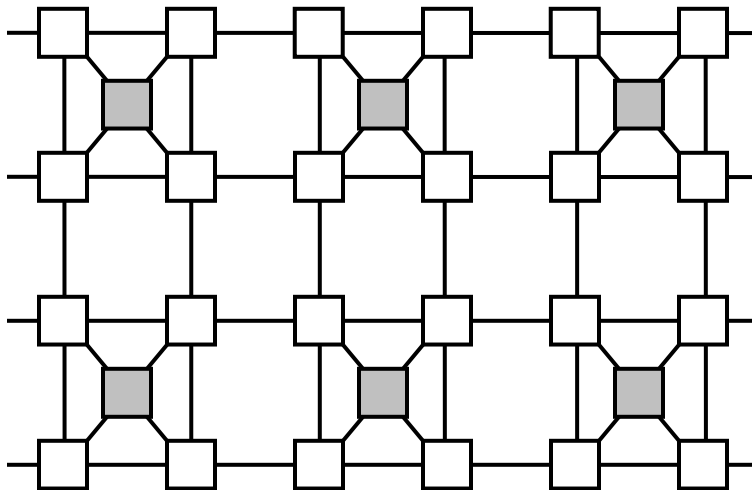
$$A_r = N \phi_r$$

Măsura A_m se calculează similar înlocuind p_r cu p_m .

Plasa (*mesh*) interstițială

O rețea convențională de genul plasă rectangulară bidimensională este incapabilă să tolereze vreun nod defect.

Redundanta interstițială (1, 4) este ilustrată de figura alăturată.



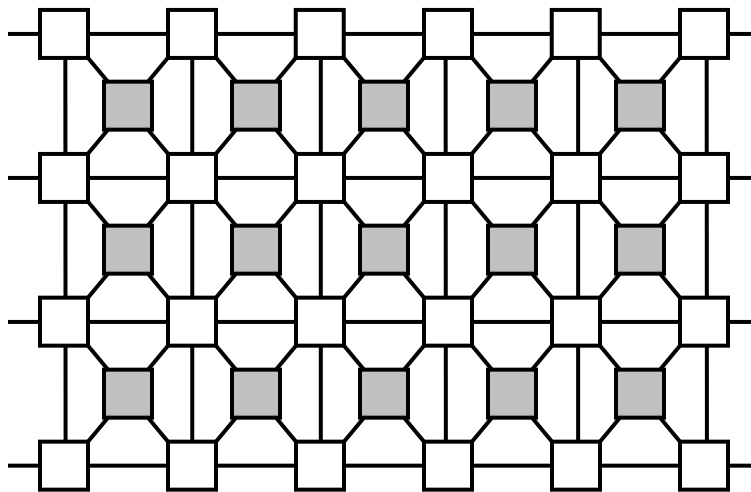
Nodurile umbrite sunt noduri de rezervă

Se observă căte un nod de rezervă adăugat pentru a înlocui oricare din cele patru noduri vecine, care a clacat. Asadar, fiecare nod primar are un singur nod de rezervă, fiecare nod suplimentar este rezervă pentru patru noduri primare. Overheadul de redundanță este 25%.

Avantajul principal rezidă în proximitatea fizică a nodului de rezervă față de nodurile primare pe care le poate înlocui. Aceasta poate reduce penalitatea de întârziere datorată utilizării unei rezerve.

Redundanta interstitială se practică uneori într-o formă diferită. Iată imediat, în figura alăturată, redundanta interstitială (4, 4).

Un nod primar are în această schemă patru noduri de rezervă și fiecare nod suplimentar este rezervă pentru patru noduri primare. O astfel de structură are un nivel de toleranță mai ridicat dar și overheadul de redundanță este mai ridicat, este de cca. 100%.



Nodurile umbrite sunt noduri de rezervă

Fiabilitatea plasei (*mesh*) cu redundanță interstitială (1, 4)

Plasa este de dimensiunile $m \times n$, cu m și n numere pare. Reteaua este alcătuită din clustere de patru noduri primare cu un nod de rezervă. Reteaua (*mesh*) are în total $mn/4$ astfel de clustere.

Fie $R(t)$ fiabilitatea unui nod primar sau de rezervă.

Fiabilitatea unui cluster este

$$R_{cluster}(t) = R^5(t) + 5R^4(t)[1 - R(t)]$$

iar fiabilitatea unei plase (*mesh*) de $m \times n$ este

$$R_{plasa}(t) = [R_{cluster}(t)]^{mn/4}$$

Dacă în cazul redundanței interstitială (1, 4), pentru funcția de fiabilitate există această expresie, pentru schemele cu redundanță interstitială (4, 4) nu există un algoritm simplu pentru calculul fiabilității.

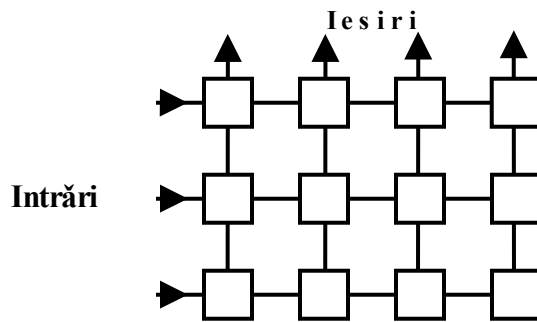
Retele *crossbar* fără redundante

Figura alăturată arată o rețea *crossbar* 3×4 (trei intrări și patru ieșiri). În general, o rețea *crossbar* $m \times n$ are n intrări, m ieșiri și mn comutatoare. Comutatoarele leagă toate perechile alcătuite dintr-o intrare și o ieșire. Rețeaua *crossbar* nu este tolerantă la defecte. Disfuncția oricărui comutator deconectează anumite perechi de noduri.

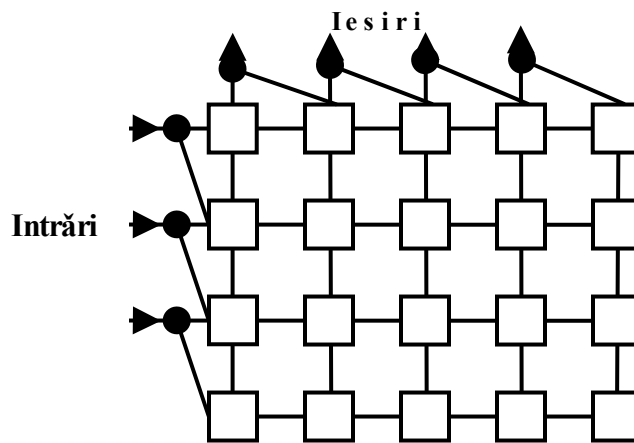
Retele *crossbar* cu redundante

Se adaugă redundante pentru a face rețeaua tolerantă la defecte. Pentru aceasta se adaugă, de pildă, o linie și o coloană de comutatoare. Conexiunile de intrare și de ieșire sunt multiplicare prin faptul că fiecare intrare poate fi trimisă pe două linii și fiecare ieșire poate fi obținută de la două coloane.

Dacă un comutator se defectează, atunci linia și coloana de care acesta aparține sunt înlocuite de linia și coloana de rezervă (v.figura).



(a)



(b)

Retele *crossbar* fără redundante (a) și cu redundante (b)

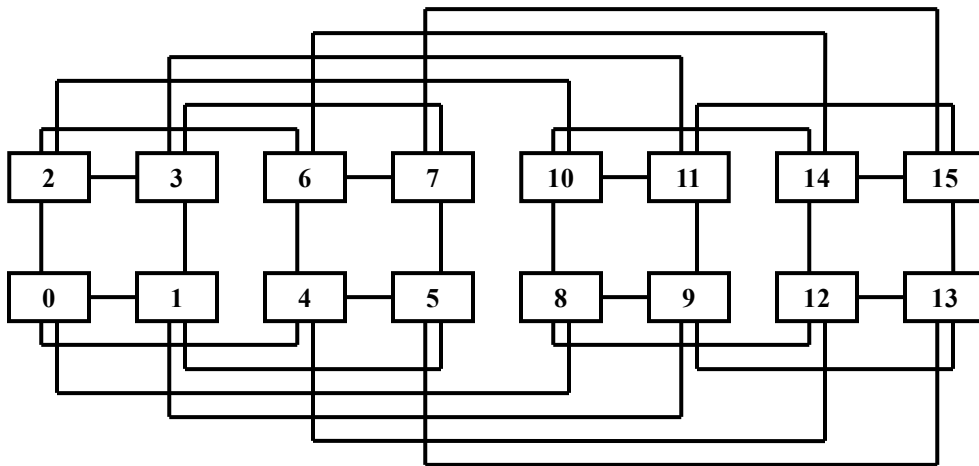
Rețele de tip hipercub

Cu H_n se notează o rețea de tip hipercub n -dimensională care are 2^n noduri. Un hipercub 0-dimensional are un singur nod. Un hipercub H_n se construiește prin conectarea nodurilor corespondente din două rețele H_{n-1} , două hipercuburi cu o dimensiune mai puțin. Muchiile adăugate pentru a conecta noduri corespondente sunt numite muchii de dimensiune $(n - 1)$.



Muchie de dimensiune 0

Exemple de hipercuburi sunt date în figura alăturată. În hipercubul H_4 din figură se disting cu ușurință hipercuburi de dimensiuni inferioare și muchii de diverse dimensiuni.



Exemple diverse rezultate din lectura figurii:

Muchie de dimensiune 0: 8-9 (diferența între numerele purtate de noduri: $1 = 2^0$).

Muchie de dimensiune 1: 4-6 (diferența între numerele purtate de noduri: $2 = 2^1$).

Muchie de dimensiune 2: 10-14 (diferența între numerele purtate de noduri: $4 = 2^2$).

Muchie de dimensiune 3: 3-11 (diferența între numerele purtate de noduri: $8 = 2^3$).

Hipercuburi H_0 : oricare nod.

Hipercuburi H_1 : oricare pereche de noduri dintre următoarele: (0, 1), (2, 3), (4, 5), (6, 7), (8, 9), (10, 11), (12, 13), (14, 15).

Hipercuburi H_2 : (0, 1, 2, 3), (4, 5, 6, 7), (8, 9, 10, 11) și (12, 13, 14, 15).

Hipercuburi H_3 : (0, 1, 2, 3, 4, 5, 6, 7) și (8, 9, 10, 11, 12, 13, 14, 15).

Rutarea în hipercuburi

Pentru a simplifica rutarea se folosește o numerotare specială. Numerele sunt exprimate în binar și dacă nodurile i și j sunt conectate de o muchie de dimensiune k , numerele pentru i și j diferă prin bitii de pe poziția k .

Exemplu: nodurile 0000 și 0010 diferă numai prin bitul de pe poziția 2^1 ; ele sunt conectate printr-o muchie de dimensiune 1.

Alt exemplu: un pachet trebuie să se deplaseze de la nodul $14 = 1110_2$ la nodul $2 = 0010_2$ într-o rețea H_4 . Rutările posibile sunt:

- $1110 \rightarrow 0110$ (dimensiune 3) $\rightarrow 0010$ (dimensiune 2)
- $1110 \rightarrow 1010$ (dimensiune 2) $\rightarrow 0010$ (dimensiune 3)

Rutarea în cazul general

Distanța între sursă și destinație este în general numărul de biți diferiți în cele două adrese (distanța Hamming). Transferul de la nodul X la nodul Y poate fi făcut prin trecerea câte o dată pe fiecare din dimensiunile prin care sursa și destinația diferă.

Dacă adresele sunt $X = x_{n-1} \dots x_0$ și $Y = y_{n-1} \dots y_0$, se definesc bitii $z_i = x_i \oplus y_i$ ($i = 0, \dots, n-1$) cu \oplus operatorul "sau-exclusiv".

Pachetul trebuie să traverseze o muchie în fiecare dimensiune pentru care $z_i = 1$.

Toleranța la defecte în rețelele hipercub

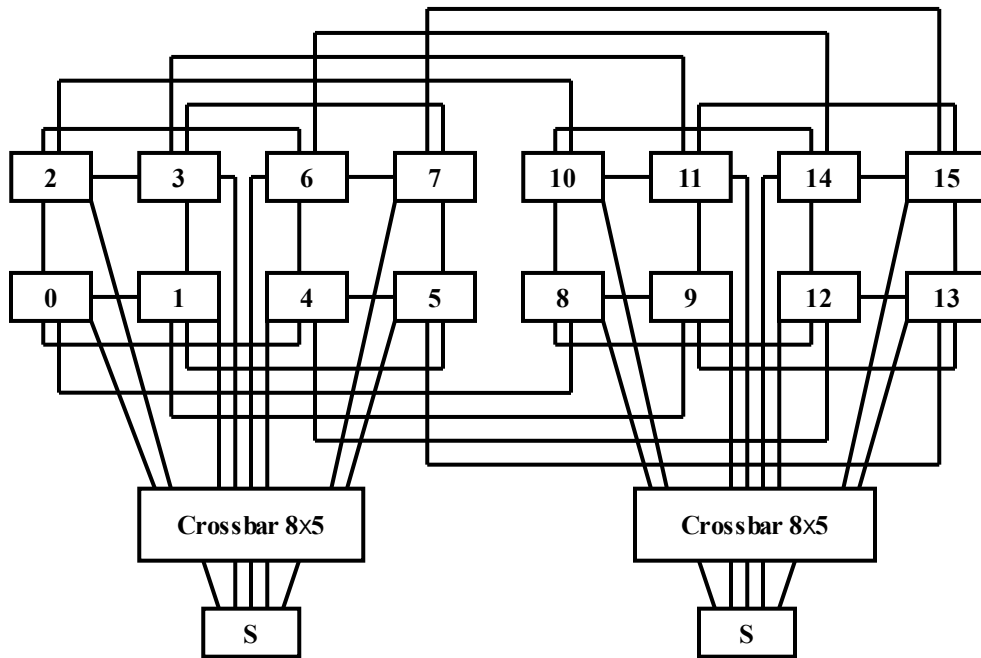
Pentru $n \geq 2$, un hipercub H_n poate tolera disfuncții ale legăturilor deoarece există căi multiple de la orice sursă la orice destinație.

Disfuncțiile nodurilor pot însă compromite operația. O modalitate de ameliorare a situației constă în creșterea numărului de porturi de comunicare ale fiecărui nod de la n la $n + 1$ și conectarea acestor porturi suplimentare prin legături adiționale la unul sau mai multe noduri de rezervă.

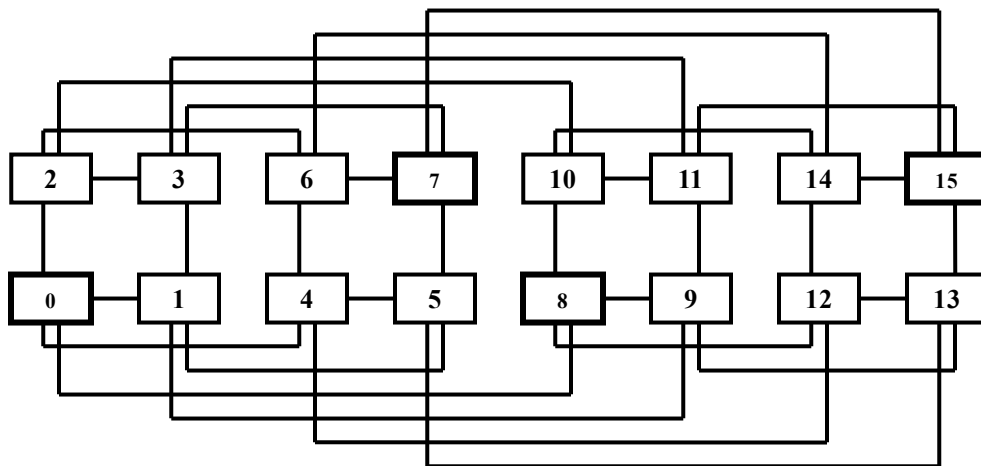
Exemplu: Se pot adăuga două noduri de rezervă, fiecare din acestea fiind o rezervă pentru 2^{n-1} noduri ale unui subcub H_{n-1} .

Nodurile de rezervă ar putea necesita 2^{n-1} porturi. Numărul de porturi poate fi redus prin utilizarea unor comutatoare *crossbar* ale căror ieșiri sunt conectate la nodul de rezervă corespunzător. Numărul de porturi ale nodului de rezervă este redus la $n + 1$, același ca pentru toate celelalte noduri.

Figura alăturată arată un hipercub H_4 cu două noduri de rezervă.



O metodă diferită de tolerare a defectelor constă în duplicarea procesoarelor din câteva (putine) noduri selectate. Fiecare procesor aditional este rezervă pentru oricare dintre procesoarele din nodurile vecine. În exemplul din figura următoare, nodurile 0, 7, 8, 15 ale unui hipercub H_4 sunt modificate prin duplicare (reprezentate îngrosat).



Fiecare nod are acum o rezervă la distanță nu mai mare de 1. Înlocuirea unui procesor defect cu unul din rezervă produce, desigur, o întârziere suplimentară în comunicare.

Rutarea în hipercuburi cu defecte

Algoritmul de rutare trebuie modificat pentru a ocoli nodurile sau legăturile defecte. Ideea de bază se poate formula astfel: se listează dimensiunile pe care un pachet trebuie să meargă și se parcurg acestea una câte una. Pe măsură ce muchiile sunt parcurse și marcate/eliminate (*crossed off*) din listă, dacă din cauza unui nod sau din cauza unei legături disfuncționale legătura dorită nu este disponibilă se alege o altă muchie din listă (dacă există una) pentru continuarea parcursului; dacă pachetul atinge un anumit nod în care găsește toate dimensiunile din lista sa căzute, el revine (*backtracks*) la nodul anterior și încercarea continuă.

Algoritmul formal de rutare utilizează următoarele notații:

TD – lista dimensiunilor pe care circulă mesajul, în ordinea parcurgerii.

TD^R – același lucru în ordine inversă (*reversed*).

$\oplus_{i=1}^k$ – operația sau-exclusiv executată de k ori, secvențial.

Exemplu: $\oplus_{i=1}^3 a_i$ înseamnă $(a_1 \oplus a_2) \oplus a_3$.

D – nod destinație, S – nod sursă, $d = D \oplus S$ (\oplus – operația sau-exclusiv se execută bit-cu-bit pe bitii corespondenți din adresele binare D și S).

$SC(A)$ – mulțimea de noduri vizitate pe un parcurs, pe fiecare din dimensiunile listate în mulțimea A .

Exemplu: la nodul 0010 – $SC(1, 3) = \{0000, 1000\}$.

e_n^i – un vector de n biți care are 1 în poziția bitului i și 0 în rest.

Exemplu: $e_3^2 = 100$.

Pachetele sunt presupuse a consta în:

- (i) d ; $d = D \oplus S$
- (ii) mesajul transmis (“încărcătura”)
- (iii) lista de dimensiuni vizitate deocamdată – TD

Θ – operația de adăugare (*append*). Scrierea $TD \Theta x$ înseamnă adăugarea la finalul listei TD a lui x .

$\text{transmit}(j)$ – rutina de trimitere a pachetului ($d \oplus e^j$, mesaj, $TD \Theta x$) pe legătura j -dimensională de la nodul curent.

Algoritm de rutare pentru hipercuburi cu defecte

```
if ( $d == 0\dots 0$ )
    destinația a fost atinsă; exit
else
    for  $j = 0$  to  $(n - 1)$  step 1 do {
        if ( $d_j == 1$ ) && (legătura în dimensiunea  $j$  din acest nod este fără
            defect) && ( $e_v^j \notin SC|TD^R$ ) {
            transmit( $j$ )
            exit
        }
    }
```

```

endif
if (există o legătură fără defect în  $SC|TD^R$ )
    fie  $h$  o astfel de legătură
else {
     $g = \max(m: \oplus_{i=1}^m e^{TD^R}(i) = 0\dots 0)$ 
    if ( $g =$  numărul de elemente din  $SC(TD)$ ) {
        nu există o cale
        exit
    }
    else
         $h =$  elementul al  $(g + 1)$ -lea din  $TD^R$ 
    endif
    transmite( $h$ )
}
end

```

Exemplu pe hipercubul H_3 :

H_3 cu nodul defect 011.

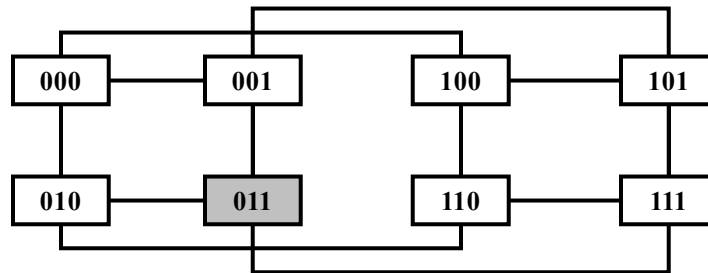
Nodul 000 trebuie să transmită un pachet la 111.

La 000, $d = 111$, trimite mesajul pe dimensiunea 0, la nodul 001.

La 001, $d = 110$ și $TD = (0)$, tentative la muchiile de dimensiune 1: imposibil.

Bitul 2 din d este tot 1. Se verifică și se stabilește că muchia de dimensiune 2 la 101 este disponibilă, mesajul este trimis la 101 și apoi la 111.

Exercițiu: Ce se întâmplă dacă sunt căzute nodurile 001 și 101?

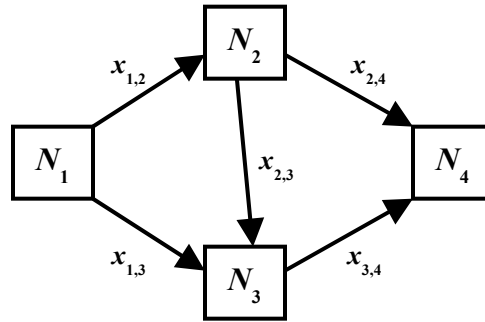


Fiabilitatea rețelelor punct-la-punct

Retelele nu sunt în mod necesar structuri regulate și de cele mai multe ori există mai multe căi între două noduri.

Se definește *fiabilitatea terminală*: probabilitatea ca să existe o cale operațională între două noduri anumite, fiind date probabilitățile disfuncțiilor pe fiecare legătură.

Exemplu: să se calculeze fiabilitatea terminală pentru perechea sursă-destinație $N_1 - N_4$ (v. figura).



Sunt trei căi de la N_1 la N_4

- $P_1 = (x_{1,2}, x_{2,4})$
- $P_2 = (x_{1,3}, x_{3,4})$
- $P_3 = (x_{1,2}, x_{2,3}, x_{3,4})$

$p_{i,j}$ ($q_{i,j}$) – probabilitatea ca legătura $x_{i,j}$ să fie bună (respectiv defectă).

Nodurile sunt presupuse a fi fără defecte. Dacă nu aceasta este situația, probabilitatea disfuncției lor este încorporată în aceea a legăturilor care pleacă din noduri.

Multimea de căi trebuie prelucrată pentru a obține o mulțime echivalentă alcătuită din evenimente mutual exclusive, altminteri unele evenimente ar putea fi luate în calcul de mai multe ori.

Evenimente mutual exclusive în cazul în discuție:

- calea P_1 funcțională;
- calea P_2 funcțională și calea P_1 disfuncțională;
- calea P_3 funcțională, căile P_1 și P_2 disfuncționale.

Dacă numim rețeaua din figură “punte” (denumire legată de forma și topologia ei), atunci fiabilitatea legăturii $N_1 - N_4$ este:

$$R_{\text{punte}} = p_{1,2}p_{2,4} + p_{1,3}p_{1,4}(1 - p_{1,2}p_{2,4}) + p_{1,2}p_{2,3}p_{3,4}(q_{1,3}q_{2,4})$$

Calculul fiabilității terminale

Pentru a calcula fiabilitatea terminală a unei rețele cu m căi P_1, \dots, P_m de la sursă la destinație se folosesc notațiile care urmează.

E_i (\bar{E}_i) – eveniment constând în operationalitatea (disfuncția) căii P_i .

$$R = \Pr(\text{existența unei căi operationale}) = \Pr\left(\bigcup_{i=1}^m E_i\right).$$

Multimea de evenimente poate fi descompusă în evenimente mutual exclusive. După descompunere, expresia evenimentului “există o cale operatională” este

$$E_1 \cup (E_2 \cap \bar{E}_1) \cup (E_3 \cap \bar{E}_1 \cap \bar{E}_2) \cup \dots \cup (E_m \cap \bar{E}_1 \cap \dots \cap \bar{E}_{m-1})$$

și

$$R = \Pr(E_1) + \Pr(E_2 \cap \bar{E}_1) + \Pr(E_3 \cap \bar{E}_1 \cap \bar{E}_2) + \dots \\ \dots + \Pr(E_m \cap \bar{E}_1 \cap \dots \cap \bar{E}_{m-1})$$

Expresia din urmă poate fi rescrisă uzând de probabilități conditionate

$$R = \Pr(E_1) + \Pr(E_2)\Pr(\bar{E}_1|E_2) + \Pr(E_3)\Pr(\bar{E}_1 \cap \bar{E}_2|E_3) + \dots \\ \dots + \Pr(E_m)\Pr(\bar{E}_1 \cap \dots \cap \bar{E}_{m-1}|E_m)$$

Problema centrală este calcularea probabilităților conditionate de forma generală $\Pr(\bar{E}_1 \cap \dots \cap \bar{E}_{i-1}|E_i)$.

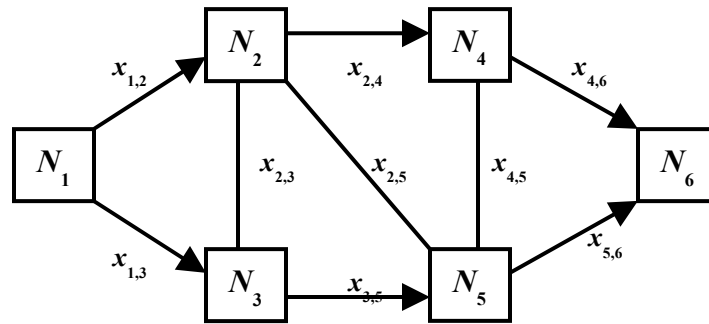
Pentru a identifica legăturile care trebuie să cadă pentru ca E_i să aibă loc dar nu E_1, \dots, E_{i-1} , se folosesc așa-numitele multimii conditionate

$$S_{j/i} = P_j - P_i = \{x|x \in P_j \text{ și } x \notin P_i\}$$

Identificarea evenimentelor disjuncte în cazul general nu este totdeauna o operație facilă.

Exemplu suplimentar pentru fiabilitatea terminală

O rețea cu șase noduri care are 9 legături unidirectionale și 3 bidirectionale.



O listă cu toate căile de la N_1 la N_6 :

- | | |
|--|--|
| $P_1 = \{x_{1,3}, x_{3,5}, x_{5,6}\}$ | $P_8 = \{x_{1,2}, x_{2,3}, x_{3,5}, x_{5,6}\}$ |
| $P_2 = \{x_{1,2}, x_{2,5}, x_{5,6}\}$ | $P_9 = \{x_{1,2}, x_{2,4}, x_{4,5}, x_{5,6}\}$ |
| $P_3 = \{x_{1,2}, x_{2,4}, x_{4,6}\}$ | $P_{10} = \{x_{1,3}, x_{2,3}, x_{2,4}, x_{4,5}, x_{5,6}\}$ |
| $P_4 = \{x_{1,3}, x_{3,5}, x_{4,5}, x_{4,6}\}$ | $P_{11} = \{x_{1,3}, x_{2,3}, x_{2,5}, x_{4,5}, x_{4,6}\}$ |
| $P_5 = \{x_{1,3}, x_{2,3}, x_{2,4}, x_{4,6}\}$ | $P_{12} = \{x_{1,3}, x_{3,5}, x_{2,5}, x_{2,4}, x_{4,6}\}$ |
| $P_6 = \{x_{1,3}, x_{2,3}, x_{2,5}, x_{5,6}\}$ | $P_{13} = \{x_{1,3}, x_{2,3}, x_{3,5}, x_{4,5}, x_{4,6}\}$ |
| $P_7 = \{x_{1,2}, x_{2,5}, x_{4,5}, x_{4,6}\}$ | |

Căile, se observă, sunt ordonate de la cea mai scurtă la cea mai lungă.

Pentru a calcula alți termeni din sumă, trebuie avută în vedere intersecția mai multor multimii conditionate.

- | |
|--|
| $P_1 = \{x_{1,3}, x_{3,5}, x_{5,6}\}$ |
| $P_2 = \{x_{1,2}, x_{2,5}, x_{5,6}\}$ |
| $P_3 = \{x_{1,2}, x_{2,4}, x_{4,6}\}$ |
| $P_4 = \{x_{1,3}, x_{3,5}, x_{4,5}, x_{4,6}\}$ |

Pentru a calcula termenul al patrulea – expresia lui P_4 – multimile conditionate sunt:

$$S_{1/4} = \{x_{5,6}\}; S_{2/4} = \{x_{1,2}, x_{2,5}, x_{5,6}\}; S_{3/4} = \{x_{1,2}, x_{2,4}\}$$

$S_{1/4}$ este inclus în $S_{2/3}$; dacă $S_{1/4}$ este cu defecte, atunci și $S_{2/4}$ este cu defecte. $S_{2/4}$ poate fi ignorat în acest caz.

Al patrulea termen din ecuația fiabilității este

$$p_{1,3}p_{3,5}p_{4,5}p_{4,6}(1 - p_{5,6})(1 - p_{1,2}p_{2,4})$$

Calculul termenului al treilea conduce la

$$S_{1/3} = \{x_{1,3}, x_{3,5}, x_{5,6}\}$$

$$S_{2/3} = \{x_{2,5}, x_{5,6}\}$$

Cele două mulțimi conditionate nu sunt disjuncte.

Evenimentul care constă în defectarea simultană a mulțimilor de arce $S_{1/3}$ și $S_{2/3}$ trebuie să fie împărțit în evenimente disjuncte:

(I) $x_{5,6}$ cu defecte

(II) $x_{5,6}$ este operational și atât $x_{1,3}$ cât și $x_{2,5}$ sunt defecte

(III) atât $x_{1,3}$ cât și $x_{5,6}$ sunt în funcțiune și atât $x_{2,5}$ cât și $x_{3,5}$ sunt defecte.

Pentru termenul al treilea rezultă expresia

$$p_{1,2}p_{2,4}p_{4,6}(q_{5,6} + p_{5,6}q_{1,3}q_{2,5} + p_{5,6}p_{1,3}q_{2,5}q_{3,5})$$

Termenii rămași se calculează similar.

Fiabilitatea terminală este suma tuturor celor 13 termeni definiți mai devreme.

ANEXA 1

ELEMENTE DE TEORIA PROBABILITĂȚILOR SI DE STATISTICĂ MATEMATICĂ

Spatiul evenimentelor

Un experiment oarecare, provocat sau spontan poate avea rezultate diverse. Aceste rezultate sunt denumite *evenimente*. Astfel, rostogolirea unui zar pe o suprafață plană orizontală (exemplu aparent banal dar des utilizat de matematicieni) poate avea ca rezultat apariția pe fața sa de deasupra a, să spunem, cinci puncte. S-a produs asadar evenimentul apariției pe fața de deasupra a cinci puncte. Tot așa, conform definiției de mai sus, extragerea valetului de cupă dintr-un pachet de cărți de joc bine amestecat este un eveniment.

Dacă E este multimea tuturor evenimentelor posibile relativ la un experiment, această multime poate fi numită, cum adesea se întâmplă, *spatiul evenimentelor*. Evenimentele unui astfel de spațiu se pot găsi în anumite relații unul cu altul și cu evenimentele acelui spațiu se pot face unele operații.

O relație importantă între evenimente este *implicatia*. Implicatia se notează $A \subset B$ și se citește *evenimentul A implică evenimentul B* , ceea ce înseamnă că producerea evenimentului A conduce automat, implicit la producerea evenimentului B ; implicatia mutuală, $A \subset B$ și $B \subset A$ este un mod de a exprima egalitatea sau echivalența a două evenimente.

Operațiile cu evenimente sunt *unare*, cu un singur eveniment ca operand, sau *binare*, cu două evenimente ca operanzi.

Operația de luare a *complementarului* sau, ceea ce este totuna, a *contrarului* unui eveniment este unară, operează cu un singur eveniment. *Reuniunea* și *intersecția* de evenimente sunt operații binare, operează pe două evenimente.

Complementarul sau contrarul unui eveniment este acel eveniment care se produce atunci când nu se produce evenimentul al cărui contrar/complementar este. Într-un exemplu foarte simplu, aruncarea unei monede cu cădere pe o suprafață plană orizontală poate avea ca rezultat afișarea deasupra fie a unei fețe, fie a celeilalte. Fiecare din cele două evenimente generate de acest experiment simplu este contrarul celuilalt. Dacă evenimentul asupra căruia se operează este A atunci evenimentul contrar lui se notează cu \bar{A} . De ce *contrar* și/sau *complementar* se va explica mai în detaliu după definirea celor două operații binare anuntate.

Reuniunea a două evenimente se notează $A \cup B$ și este evenimentul care constă în producerea a cel puțin unuia din cele două evenimente, adică sau a unui sau a celuilalt sau a ambelor evenimente.

Intersecția a două evenimente se notează $A \cap B$ și este evenimentul constând în producerea ambelor evenimente, adică atât a unui eveniment cât și a celuilalt. Există două evenimente speciale care se includ în mulțimea E . Unul este evenimentul imposibil, notat cu \emptyset , evenimentul care nu se produce niciodată. Celălalt este evenimentul sigur, notat cu E , evenimentul care se produce de fiecare dată.

O relație de forma $A \cap B = \emptyset$ exprimă incompatibilitatea mutuală a celor două evenimente A și B , în alte cuvinte producerea unuia exclude producerea celuilalt.

Acum se poate formula mai precis raportul între un eveniment și contrarul lui: $A \cap \bar{A} = \emptyset$, $A \cup \bar{A} = E$. Într-o lectură în cuvinte a acestor relații, un eveniment este incompatibil cu contrarul său, producerea unui eveniment sau a contrarului său este sigură. Este acum momentul să se aducă precizarea că contrarul contrarului unui eveniment este exact acel eveniment. Simbolic, aceasta se scrie $\bar{\bar{A}} = A$.

Mulțimea E este parțial ordonată, relația de ordine este implicația.

Mulțimea E împreună cu operațiile de luare a contrarului unui eveniment, de reuniune și de intersecție a evenimentelor se organizează ca o algebră booleană. Între evenimentele dintr-o mulțime E se disting atomi (sau evenimente elementare) și evenimente compuse. De pildă, prin aruncarea zarului se pot produce între altele evenimentele A_2 și A_5 care constau în apariția pe fața de deasupra a numărului de puncte trecut ca indice. Ambele sunt atomi sau evenimente elementare în sensul că nu sunt alte evenimente încă mai simple decât ele. Reuniunea $A_2 \cup A_5$ este însă un eveniment compus.

Fie acum Ω mulțimea tuturor evenimentelor elementare dintr-o mulțime finită E de evenimente. Evident $\Omega \neq \emptyset$. O submulțime de părți ale lui Ω , $K \subset P(\Omega)$ se organizează ca un corp dacă

$$\begin{aligned} A \in K &\Rightarrow \bar{A} \in K \\ A, B \in K &\Rightarrow A \cup B \in K \\ A, B \in K &\Rightarrow A \cap B \in K \end{aligned}$$

În aceste condiții perechea (Ω, K) este un corp de evenimente și este un σ -corp sau corp borelian de evenimente dacă orice reuniune sau intersecție finită sau infinită de evenimente din K aparține mulțimii K .

Într-un spațiu E complet și atomic, orice eveniment $A \neq \emptyset$ se poate scrie ca o reuniune de elemente din Ω

$$A = \bigcup_{\omega_i \in \Omega} \omega_i$$

Se numeste *partitie* a unui eveniment $A \in K$ o multime de evenimente $A_i \in K$, ($i = 1, 2, \dots, n$), care sunt mutual incompatibile, adică $A_i \cap A_j = \emptyset$ pentru orice pereche $A_i, A_j \in K$ cu $i \neq j$, astfel încât

$$\bigcup_{i=1}^n A_i = A$$

Dacă $A = \Omega$ atunci evenimentele $A_i \in K$, ($i = 1, 2, \dots, n$) alcătuiesc un *sistem complet de evenimente* sau o *familie exhaustivă de evenimente*.

Probabilități, probabilități conditionate

Pe multimea evenimentelor din K se definește o funcție reală P , numită *probabilitate*, care are proprietățile:

1. $P(A) \geq 0, \forall A \in K$
2. $P(\Omega) = 1$
3. $P\left(\bigcup A_i\right) = \sum P(A_i), A_i \in K, A_i \cap A_j = \emptyset, i \neq j$

Dacă ultima proprietate are loc și pentru reuniuni numerabile atunci probabilitatea P se numeste *complet aditivă* (sau σ -aditivă) pe corpul (borelian) de evenimente (Ω, K) .

Tripletul (Ω, K, P) se numeste *câmp (borelian) de probabilitate*. Dacă Ω este o multime finită atunci (Ω, K, P) este un *câmp de probabilitate discret*.

Din proprietățile de mai sus derivă alte câteva proprietăți importante ale probabilității P . Astfel

4. $P(\emptyset) = 0$
5. $P(A) = 1 - P(\bar{A})$
6. $P(A - B) = P(A) - P(A \cap B)$
7. $0 \leq P(A) \leq 1$
8. $P(A \Delta B) = P(A) + P(B) - 2P(A \cap B)$
9. $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

unde $A - B = A \cap \bar{B}$ și $A \Delta B = (A - B) \cup (B - A)$ sunt diferența, respectiv diferența simetrică a două evenimente. O extindere a relației ultime la reuniunea a n evenimente este

$$P\left(\bigcup_{i=1}^n A_i\right) = \sum_{j=1}^n (-1)^{j+1} S_j \quad \text{cu} \quad S_j = \sum_{i_1, i_2, \dots, i_j \leq n} P(A_{i_1} \cap \dots \cap A_{i_j}) \quad j \leq n.$$

Dacă $F = \{A_i\}_{i \in I}$ este o familie numerabilă de evenimente două câte două mutual incompatibile, atunci $P\left(\bigcap_{i \in I} A_i\right) = 0$.

Dacă familia $F = \{A_i\}_{i \in I}$ este și exhaustivă, adică este un sistem complet de evenimente, atunci $P\left(\bigcup_{i \in I} A_i\right) = 1$.

Evenimentele se pot afla în relație de condiționare reciprocă în sensul că un eveniment odată produs poate modifica probabilitatea de producere a altui eveniment. Probabilitatea devine astfel *condiționată*.

Relația de bază pentru calculul probabilităților condiționate este

$$P_B(A) = P(A/B) = P(A \cap B) / P(B)$$

cu evenimentul B , cel care condiționează producerea evenimentului A , trecut ca indice sau pe poziția a doua în notația $P(A/B)$, deci după caracterul “/” (sau “|”) în argumentul funcției probabilitate.

În general,

$$p(A/B) \neq P(A) \text{ și } P(B/A) \neq P(B)$$

ceea ce indică o dependență, o condiționare reală între cele două evenimente. Dacă are loc egalitatea în ambele relații atunci evenimentele nu se condiționează în nici un fel, sunt independente.

Dacă probabilitatea unei intersecții finite de evenimente este nenulă

$$P\left(\bigcap_{i=1}^n A_i\right) \neq 0$$

atunci probabilitatea respectivă se poate calcula cu formula

$$P\left(\bigcap_{i=1}^n A_i\right) = P\left(A_n / \bigcap_{i=1}^{n-1} A_i\right) \dots P(A_2 / A_1) P(A_1)$$

care se demonstrează inductiv pornind de la relația pentru două evenimente care se condiționează unul pe altul

$$P(A \cap B) = P(A/B)P(B) = P(B/A)P(A)$$

derivată simplu din formula probabilității condiționate.

Dacă $F = \{A_i\}_{i=1, \dots, n}$ este o partiție a câmpului Ω atunci probabilitatea unui eveniment oarecare se poate calcula cu relația

$$P(A) = \sum_{i=1}^n P(A_i)P(A/A_i)$$

cunoscută sub numele de *formula probabilității totale*.

Mai este de reținut *formula lui Bayes*

$$P(A_i/A) = \frac{P(A_i)P(A/A_i)}{\sum_{i=1}^n P(A_i)P(A/A_i)}$$

care în aceleași condiții, $F = \{A_i\}_{i=1, \dots, n}$ o partiție a câmpului Ω , permite calculul probabilității fiecărui eveniment al partiției condiționat de evenimentul $A \in K$, altfel oarecare.

Exemplu. În cazul zarului corect enunțat mai devreme, mulțimea Ω este alcătuită din evenimentele $A_1, A_2, A_3, A_4, A_5, A_6$. Mulțimea de părți ale lui Ω care

se constituie în corp de evenimente este mulțimea tuturor reuniunilor de 2, 3, 4, 5 sau 6 evenimente la care se adaugă evenimentele atomice, elementare deja enumerate și evenimentul imposibil \emptyset .

Prin percepție imediată se poate afirma că șansele de producere a celor șase evenimente sunt egale (șansa aceasta de producere a unui eveniment sau a altuia este măsurată de probabilitate). Se poate scrie, asadar

$$P(A_1) = P(A_2) = P(A_3) = P(A_4) = P(A_5) = P(A_6) = p$$

Evenimentul sigur Ω se poate scrie ca o reuniune

$$\Omega = A_1 \cup A_2 \cup A_3 \cup A_4 \cup A_5 \cup A_6$$

și deoarece evenimentele din reuniune sunt două câte două mutual incompatibile (nu pot apărea deasupra două fete diferite deodată), conform proprietății 3 se poate scrie

$$1 = P(\Omega) = P(A_1) + P(A_2) + P(A_3) + P(A_4) + P(A_5) + P(A_6) = 6p$$

adică $p = 1/6$. Acum se pot calcula probabilități diverse.

a) Probabilitatea apariției unui număr par de puncte este

$$P(A_2 \cup A_4 \cup A_6) = P(A_2) + P(A_4) + P(A_6) = 3(1/6) = 1/2$$

ca probabilitate a unei reuniuni de evenimente două câte două reciproc incompatibile.

b) Probabilitatea evenimentului A_4 condiționată de evenimentul reuniune de la punctul precedent, $A = A_2 \cup A_4 \cup A_6$

$$P(A_4/A) = \frac{P(A_4 \cap A)}{P(A)} = \frac{P[A_4 \cap (A_2 \cup A_4 \cup A_6)]}{P(A)} = \frac{P(A_4)}{P(A)} = \frac{1/6}{1/2} = \frac{1}{3}$$

etc.

De reținut din acest exemplu o modalitate de a evalua probabilități prin raportarea numărului de situații favorabile unui eveniment la numărul total de situații. De pildă, evenimentul A_3 din cele de mai sus se produce în proporția 1 caz favorabil din 6 posibile, adică $1/6$.

La loteria "6 din 49", se pot juca C_{49}^6 (combinări de 49 de numere luate câte 6) variante distincte. Șansa (probabilitatea) unei variante particulare de a fi câștigătoare a premiului cel mare este de $1/C_{49}^6$, o probabilitate foarte, foarte mică desigur.

Șansa de a câștiga la categoria a II-a este întrucâtva mai mare. Un bilet câștigător poate conține una din cele $C_6^5 = 6$ combinații de 5 numere din cele ieșite din urnă la extragerea duminicală, la care se adaugă unul din celelalte 43 de numere din afara extragerii. Numărul de situații favorabile câștigării unui premiu la categoria a II-a este, evident, $6 \times 43 = 258$ și probabilitatea este de $258/C_{49}^6$, încă destul de mică dar mai mare decât cea de la categoria I.

La jocul de table (foarte cunoscut în toată lumea – backgammon), evoluția disputei dintre jucători este hotărâtă pas cu pas prin aruncarea a două zaruri. Într-o anumită fază a jocului, unul dintre jucători are nevoie ca zarurile să producă o sumă a punctelor egală cu 5. Care este probabilitatea acestui eveniment? Numărul total de rezultate este 36: fiecare din cele șase fete ale

unui zar poate apărea combinată cu oricare din cele șase fete ale celuilalt zar. Suma punctelor este 5 în următoarele 4 cazuri: (1, 4), (2, 3), (3, 2) și (4, 1). Prin raportarea numărului de cazuri favorabile (4) la numărul total de cazuri (36) se obține răspunsul la întrebare: $4/36 = 1/9$.

Variabile aleatoare

Formal, o *variabilă aleatoare* este o funcție definită pe o mulțime atomică, cu valori reale, $X: \Omega \rightarrow R$, care are proprietatea

$$\{X < x\} \Rightarrow \{\omega \in \Omega / X(\omega) < x\} \in K, \forall x \in R$$

În termeni mai puțin riguroși din punct de vedere matematic, o variabilă aleatoare este o variabilă care ia valori la întâmplare dar în nici un caz haotic. Explicit sau tacit, în spatele manifestării variabilei aleatoare prin valori diverse se află un câmp de probabilitate (Ω, K, P) definit de mulțimea atomică Ω , de corpul de părți ale acesteia K și de probabilitatea P . Probabilitatea face ca unele valori pe care variabila aleatoare le poate lua să fie (eventual) mai probabile decât altele. Numărul de puncte afișate de un zar comun este o variabilă aleatoare. Cu fetele zarului, care pot fi de pildă colorate, nu neapărat “punctate”, se pot asocia și alte numere printr-o funcție X particulară. Funcția X este o altă variabilă aleatoare definită pe câmpul (Ω, K, P) .

O *variabilă aleatoare simplă* ia numai un număr finit de valori. De exemplu funcția indicator a unui eveniment $A \in K$, care se poate produce sau nu

$$\chi_A = \begin{cases} 0 & \omega \notin A \\ 1 & \omega \in A \end{cases}$$

este o variabilă aleatoare simplă care ia numai două valori, 0 și 1. Variabilele aleatoare definite în relație cu zarul sunt variabile aleatoare simple.

Dacă X este o variabilă aleatoare definită pe câmpul (Ω, K, P) atunci pentru oricare două valori $x_1, x_2 \in R, x_1 \leq x_2$ toate intervalele finite sau infinite delimitate de cele două valori corespund unor evenimente din K și, prin generalizare, pentru orice mulțime I reuniune de intervale din R , se poate calcula $P_X(I) = P[X(\omega) \in I] = P[X^{-1}(I)]$.

$P_X(I)$ reprezintă *distributia probabilistică* a variabilei aleatoare X . Se poate vorbi de P_X ca de o probabilitate definită pe câmpul (R, K_X) în care $K_X = \{I \subset R / X^{-1}(I) \in K\}$ este o mulțime de intervale ale multimii numerelor reale R , intervale care sunt imagini prin funcția X ale unor evenimente din K .

Dacă variabila aleatoare X ia valori într-o mulțime cel mult numerabilă

$$\{x_i / x_i \in R, i \in I, I \subset N^+\}$$

atunci ea se numește discretă și

$$\sum_{i \in I} P_X(x_i) = 1$$

$$P_X(J) = \sum_{x_i \in J} P_X(x_i), \forall J \in K_X$$

Dacă X poate lua toate valorile unui interval $I \in K_X$ atunci probabilitatea asociată intervalului este

$$P_X(I) = \int_I f_X(x) dx$$

și este o funcție absolut continuă. Funcția $f_X(x)$ care apare în formulă se numește *densitatea de probabilitate* sau *densitatea de repartiție* a variabilei aleatoare X , este nenegativă pentru orice x și are proprietatea

$$\int_{-\infty}^{\infty} f_X(x) dx = 1$$

Funcția

$$F_X(x) = P[X(\omega) < x] = P_X[(-\infty, x)] = \int_{-\infty}^x f_X(x) dx$$

se numește *funcție de repartiție* a variabilei aleatoare X .

Funcția de repartiție este nedescrescătoare pe întreaga axă reală

$$a < b \Rightarrow F_X(a) \leq F_X(b) \quad \forall a, b \in R$$

și este continuă la stânga în fiecare punct

$$\lim_{x \rightarrow a, x < a} F_X(x) = F_X(a) \quad \forall a \in R$$

Valoarea minimă și valoarea maximă ale unei funcții de repartiție sunt date de limitele

$$\lim_{x \rightarrow -\infty} F_X(x) = 0, \quad \lim_{x \rightarrow \infty} F_X(x) = 1$$

Eventualele discontinuități sunt de speta primă și sunt cel mult numerabile. Reciproc, orice funcție cu proprietățile de mai sus poate fi pusă în corespondență cu un câmp de probabilitate.

Pentru o variabilă aleatoare discretă

$$F_X(x) = \sum_{x_i < x} P_X(x_i)$$

iar pentru una continuă, în afară de relația scrisă deja mai sus

$$F_X(x) = \int_{-\infty}^x f_X(x) dx$$

are loc și relația

$$f_X(x) = \frac{d}{dx} F_X(x)$$

Pentru orice interval $[a, b) \subset R$ $P_X\{[a, b)\} = F_X(b) - F_X(a)$ și pentru orice valoare a , $P_X(a) = 0$.

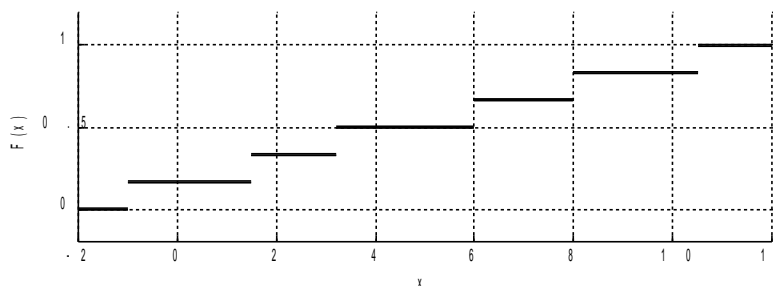
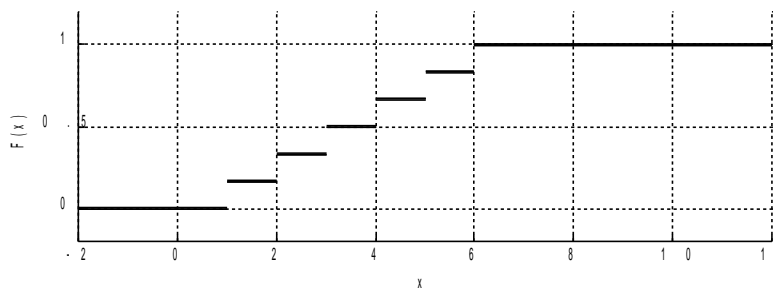
În referirea făcută puțin mai devreme la cazul zarului, s-a semnalat posibilitatea ca pe același câmp de probabilitate să se definească nu una ci mai multe variabile aleatoare. Se notează uzual cu $V(\Omega, K, P)$ mulțimea tuturor variabilelor aleatoare definite pe câmpul de probabilitate trecut între paranteze.

Dacă $X, Y \in V(\Omega, K, P)$ atunci suma, diferența, produsul celor două variabile aleatoare, modulul, puterea, în general o funcție măsurabilă Borel de oricare dintre ele sunt toate variabile aleatoare din mulțimea $V(\Omega, K, P)$.

Ori de câte ori nu este pericol de confuzie, variabila aleatoare trecută până acum ca indice al funcției de repartiție sau al funcției densitate de probabilitate/de repartiție poate lipsi din acea poziție.

Dacă se reia exemplul zarului, care la fiecare experiență este făcut să se rostogolească pe o suprafață plană, orizontală, atunci mulțimea evenimentelor elementare (atomi) Ω este alcătuită din aparițiile deasupra a celor șase fețe, marcate uzual cu unu până la șase puncte. Mulțimea de părți ale lui Ω este alcătuită din evenimentele elementare și din toate reuniunile posibile de evenimente elementare la care se adaugă evenimentul imposibil. Mulțimea K organizată ca un corp de evenimente coincide chiar cu mulțimea de părți $P(\Omega)$, iar funcția numită probabilitate ia valoarea $1/6$ pentru fiecare din evenimentele elementare deoarece fețele zarului au șanse egale de a apărea deasupra.

F u n c t i i d e r e p a r t i t i e p e n t r u d o u a v a r i a b i l e a



Funcții de repartiție pentru cazul zarului corect

Cum s-a mai spus, numărul de puncte observat pe fața de deasupra a zarului poate fi considerat o variabilă aleatoare. În acest caz funcția de repartiție se prezintă ca în graficul superior din desenul de mai sus și este, ca pentru orice variabilă aleatoare discretă, o funcție în scară.

Dar pe același câmp de probabilitate se pot defini și alte variabile aleatoare. Pe câmpul asociat zarului perfect se poate imagina, de pildă, funcția $X: \Omega \rightarrow R$ definită astfel

$$\begin{array}{cccccc} \omega_1 & \omega_2 & \omega_3 & \omega_4 & \omega_5 & \omega_6 \\ -1 & 1,5 & 3,2 & 6 & 8 & 10,5 \end{array}$$

($\omega_i \equiv A_i, i = 1, 2, 3, 4, 5, 6$) și atunci funcția de repartiție se prezintă diferit, ca în graficul inferior din aceeași figură. Asadar, mulțimea $V(\Omega, K, P)$ este extrem de bogată.

De variabilele aleatoare sunt legate câteva valori remarcabile. Una foarte importantă este *media* definită ca

$$M(x) = \int_{-\infty}^{\infty} xf(x)dx$$

care face parte din lista nesfârșită a momentelor de diferite ordine ale variabilei, acesta fiind momentul de ordinul 1.

Cu o relație similară se poate calcula media unei funcții $g(x)$ de variabila aleatoare x având în vedere caracterul aleator al valorilor funcției

$$M[g(x)] = \int_{-\infty}^{+\infty} g(x)f(x)dx$$

și dacă în particular $g(x) = x^r, r \in N$ avem tocmai *momentul de ordinul r* despre care s-a amintit.

În cazul particular $g(x) = [x - M(x)]^2$ se obține o altă valoare importantă, caracteristică variabilei aleatoare descrise de funcția de repartiție $F(x)$ sau de densitatea de repartiție $f(x)$, și anume *dispersia*. Rădăcina pătrată pozitivă a dispersiei se numește *abaterea medie pătratică* a variabilei aleatoare respective. Dispersia este momentul centrat de ordinul doi al variabilei aleatoare, unul din multiplele momente ale variabilei, centrate pe medie, de ordine diferite.

Nu numai variabilele aleatoare continue au momente, medii, dispersii, ci și cele discrete. În cazul discret, formulele de calcul contin sume în locul integralelor și valorile variabilei parcurg întreaga listă de valori posibile, iar densitatea de repartiție este înlocuită de probabilitățile asociate valorilor pe care variabila le poate lua efectiv.

Câteva legi de repartiție teoretice foarte utilizate sunt prezentate pe scurt în continuare.

Legea binomială (Bernoulli) este exprimată de relația

$$P(m) = C_n^m p^m (1-p)^{n-m}$$

cu $0 \leq m \leq n$ și p un număr în intervalul $[0, 1]$. Legea binomială este de tip discret. Variabila aleatoare este m . Are media np și dispersia $np(1-p)$. Există un model fizic legat de această lege de repartiție. Modelul îl constituie o urnă cu bile de două culori, iar evenimentele constau în rezultatele extragerii repetate a câte unei bile după care bila extrasă este reintrodusă în urnă. Variabila m reprezintă numărul bilelor de o anumită culoare din cele două, în n extrageri succesive, conform schemei cu bila returnată. Numărul p reprezintă proporția de bile de acea culoare în urnă, cu alte cuvinte probabilitatea de extragere a unei bile de culoarea respectivă.

Legea Poisson exprimată sub forma

$$P(m) = \frac{\mu^m}{m!} \exp(-\mu)$$

cu $\mu > 0$ și m natural ca variabilă aleatoare. Media variabilei m este μ , dispersia ei este de asemenea μ . Un model fizic îl reprezintă numărul dezintegrărilor radioactive, numărul de apeluri telefonice solicitate într-o centrală etc. într-un interval de timp precizat, (relativ) scurt.

Legea normală (gaussiană) care este dată de densitatea de probabilitate

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-m)^2}{2\sigma^2}}$$

în care m este media variabilei x și σ^2 este dispersia ei. Legea normală este considerată o lege limită pentru sumele de variabile aleatoare. Un fenomen afectat de foarte mulți factori aleatori care acționează prin însumare (aditiv) se prezintă de cele mai multe ori ca un fenomen aleator descris de o lege normală.

Variabilele aleatoare din expunerea teoretică sau din exemplele prezentate mai sus au fost până acum simple, adică a fost vorba în toate cazurile de o singură aplicație $X: \Omega \rightarrow R$ legată de un unic câmp de probabilitate (Ω, K, P) . Se pot imagina variabile aleatoare cu mai multe componente, variabile aleatoare sub forma unor vectori cu componente aleatoare definite relativ la un același câmp de probabilitate sau la câmpuri de probabilitate diferite. Astfel, legea următoare se referă la o variabilă aleatoare vectorială.

Legea normală multidimensională dată de densitatea de repartiție

$$f(x) = \frac{1}{(2\pi)^{\frac{n}{2}} \sqrt{\det W}} e^{-\frac{1}{2}(x-m)^T W^{-1}(x-m)}$$

cu media m un vector cu n componente și cu matricea de covarianță W o matrice $n \times n$ pozitiv definită. Pentru ca expresia dată să aibă consistența necesară trebuie definită mai exact matricea W .

Este de comentat mai întâi problema corelației a două variabile aleatoare. Două variabile aleatoare pot fi necorelate, caz în care valorile uneia nu influențează în nici un fel valorile pe care le poate lua cealaltă, dar, alternativ, pot fi mai mult sau mai puțin dependente, ceea ce înseamnă că dacă una din variabile a luat o valoare atunci legea de repartiție a celeilalte se modifică în funcție de acea valoare a primei variabile.

Fiind date două variabile aleatoare x și y de medii nule, media produsului lor $M(xy)$ se numește covarianță. Dacă covarianța este nulă se poate spune în general că cele două variabile nu sunt corelate. Dimpotrivă, dacă $M(xy) \neq 0$ variabilele sunt corelate, există o corelație între ele, există o dependență între valorile pe care ele le iau în sensul arătat puțin mai devreme. Dacă mediile sunt diferite de zero, afirmația și definiția se mențin pentru abaterile de la medie. Întrucât covarianța $M(xy)$ poate lua valori foarte diferite, pentru o apreciere cantitativă mai riguroasă a tăriei corelației se utilizează coeficientul de corelație

$$\rho = \frac{M(xy)}{\sqrt{M(x^2)M(y^2)}}$$

care ia valori în intervalul $[-1, 1]$ și în expresia căruia se disting dispersiile celor două variabile, $M(x^2)$ și $M(y^2)$. O valoare apropiată de extremele intervalului indică o corelație strânsă, o valoare apropiată de zero exprimă o corelație slabă.

Componentele unui vector aleator, privite ca variabile aleatoare simple sunt mutual mai mult sau mai puțin corelate. Se definește ca matricea a covariațiilor unui vector aleator x media produsului xx^T , adică media produsului acelui vector cu transpusul său. Se obține o matrice pătrată simetrică care are pe diagonală dispersiile componentelor pure. Aceasta este matricea W utilizată în expresia densității de repartiție a variabilei aleatoare normale multidimensionale din discuția de mai sus. Dacă matricea covariațiilor este diagonală (are toate elementele nule cu excepția celor de pe diagonală principală) atunci componentele vectorului aleator sunt mutual independente. Împărțirea fiecărui element al matricei covariațiilor cu abaterile medii pătratice ale componentelor corespunzătoare ale vectorului x produce o matrice a coeficienților de corelație, cu 1 pe diagonală, cu valori în intervalul $[-1, 1]$ în rest.

BIBLIOGRAFIE

1. T.I.Băjenescu *Fiabilitatea, disponibilitatea si mentenabilitatea sistemelor electronice complexe*, Editura de Vest 1997
2. R.K.Iyer, Cursul ECE 542, Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign, Ill., 2004
3. I.Karen, Cursul ECE 655, Department of Electrical and Computer Engineering, University of Massachusetts, Amherst Mass., 2006
4. A.Krings, Cursurile CS 449/549, Department of Computer Science, University of Idaho, Moscow, Idaho, 2005
5. Gh.M.Panaitescu *Sisteme tolerante la defecte*, Note de curs (pe suport electronic), Universitatea "Petrol-Gaze" Ploiesti, 2006
6. K.K.Saluja, Cursul ECE 753, Department of Electrical and Computer Engineering, University of Wisconsin, Madison Wisc., 2005
7. C.Stefănescu *Sisteme tolerante la defecte*, Matrix Rom, Bucuresti 1999

