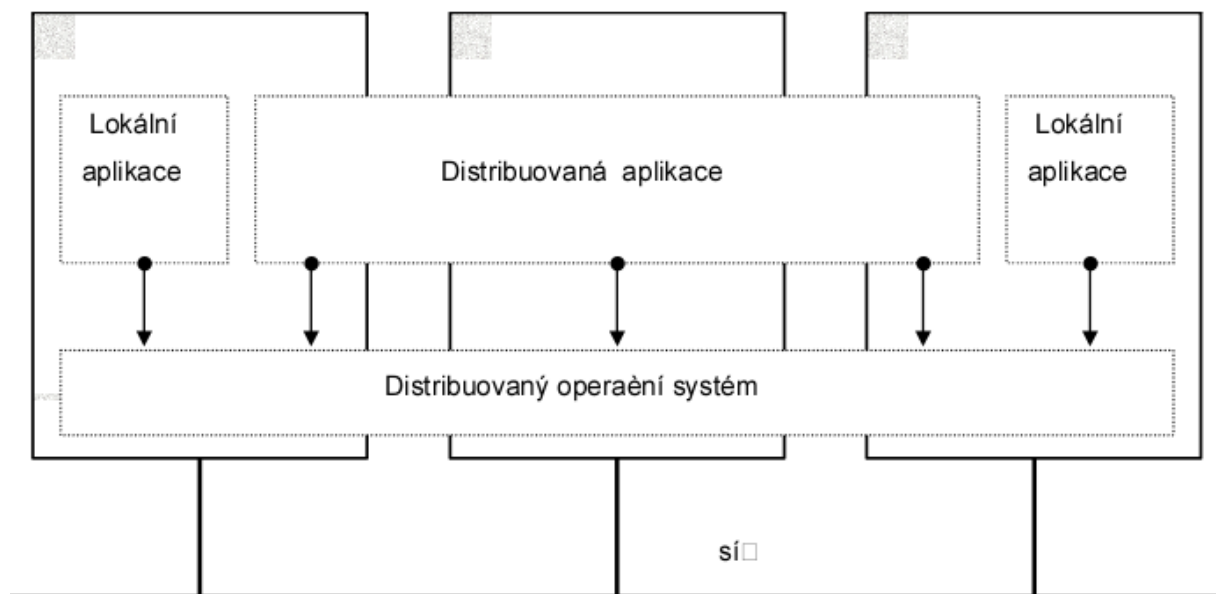


DS (distribuovaný systém)

- je systém propojení množiny nezávislých uzlů, který poskytuje uživateli dojem jednotného systému
- hardwarová nezávislost
 - uzly jsou tvořeny nezávislými “počítači” s vlastním procesorem a pamětí
 - jedinou možností komunikace přes komunikační (síťové) rozhraní
- dojem jednotného systému
 - uživatel (člověk i software) komunikuje se systémem jako s celkem
 - nemusí se starat o počet uzlů, topologii, komunikaci apod.
- celá škála různých řešení
 - multiprocesory na jedné základové desce
 - celosvětový systém pro miliony uživatelů
 - moderní technologie: cluster, cloud, grid, ...

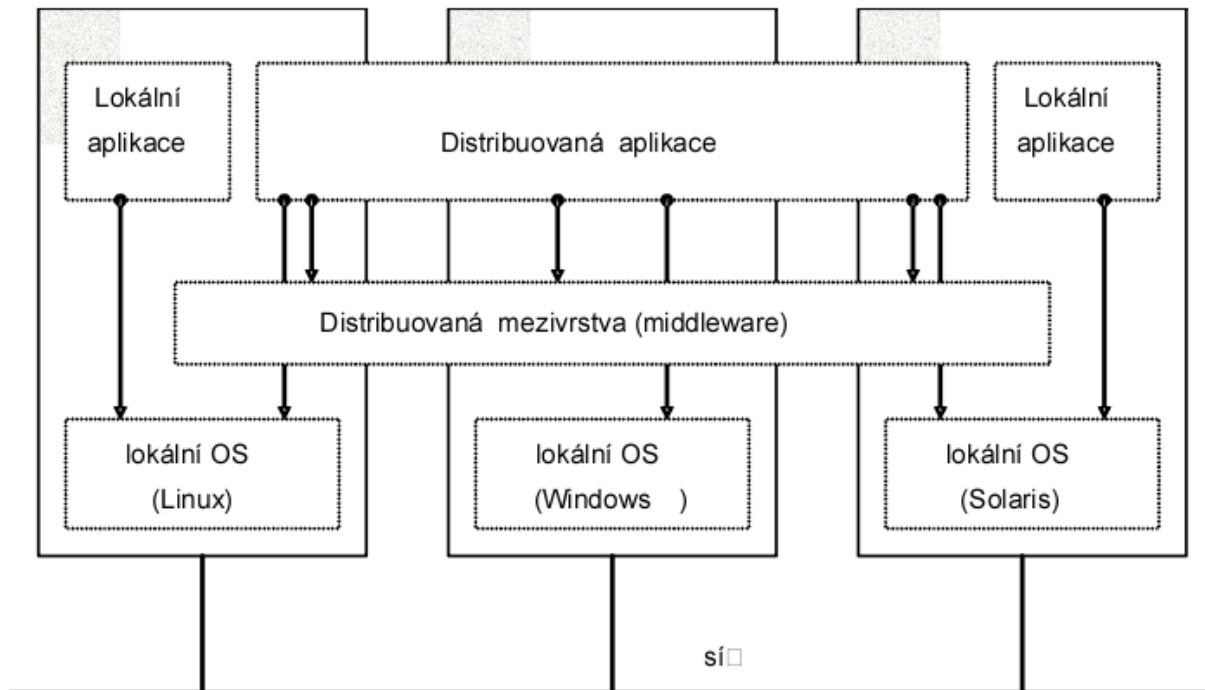
DS v 199x

- operační systém vyvinutý speciálně pro potřeby DS
- “zlatý věk”: 90. léta - Amoeba, Sprite, V, Chorus, Mosix ... T4
- součástí jádra OS podpora distribuovanosti, komunikace, synchronizace
- aplikace: transparentní využití distribuovaných služeb
- neexistuje rozdíl mezi lokální a distribuovanou aplikací
- nesplněná očekávání, žádný systém nedotažen do masově použitelného stavu



DS v 200x

- softwarový průmysl - jiný směr
- dominance MS Windows, UNIX / Linux
- použití některých principů a mechanismů distribuovaných systémů
- libovolný lokální OS, rozšiřující vrstva - distribuované prostředí + další služby
- prostředí pro distribuované aplikace – middleware (OSF DCE, CORBA, DCOM, Globe, ...)



DS v 201x

- service-oriented computing
- SaaS, PaaS, IaaS
 - software as a service, platform as a service, infrastructure as a service
- nutná podpora virtualizace
- kompletní infrastruktura připravená k použití
- Windows Azure, Amazon Elastic Compute Cloud, Google App Engine, ...

Distributed computing

- obecné principy distribuovaných systémů aplikovatelné na různá prostředí
 - cluster / grid (LHC-CERN) / cloud / sky computing
 - distributed data processing
 - distribuované databáze, NoSQL databáze, Big Data, Hadoop, data mining
 - mobilní systémy
 - bezdrátové sítě, routing, lokalizace, replikace
 - ubiquitní a pervasivní systémy
 - mikrouzly, každodenní sci-fi, samoorganizace, kolektivní znalost a rozhodování
 - distribuovaní agenti, inteligentní brouci
 - mobilita, migrace, task-oriented entities, kooperace, distribuovaná inteligence
 - peer-to-peer sítě
 - sdílení MM obsahu, P2P výpočty, distribuované ukládání, replikace, důvěryhodnost
 - content management
 - publish-subscribe, media distribution, streaming
 - ...

Motivace a cíle návrhu DS

- motivace
 - ekonomika
 - síť běžných PC × supercomputer
 - rozšiřitelnost
 - přidání uzlu × upgrade centrálního serveru
 - spolehlivost
 - výpadek jednoho uzlu × výpadek CPU
 - výkon
 - technologické limity
 - distribuovanost
 - 'inherentní distribuovanost
- cíle návrhu
 - transparentnost – přístupová, lokační, migrační, replikační, ...
 - přizpůsobivost – autonomie, decentralizované rozhodování, ...
 - spolehlivost
 - výkonnost
 - rozšiřitelnost – rozdílné metody pro 100 a 100.000.000 uzlů

Transparentnost

- na úrovni komunikace s uživatelem, komunikace procesů s jádrem
- stupeň transparentnosti x výkonnost systému
- přístupová
 - proces má stejný přístup k lokálním i vzdáleným prostředkům
 - jednotné služby, reprezentace dat, ...
- lokační
 - uživatel (proces) nemůže říci, kde jsou prostředky umístěny
 - nevyhovuje: počítač:cesta/soubor
- exekuční / migrační
 - procesy mohou běžet na libovolném uzlu / přemísťovat se mezi uzly
- replikační
 - uživatel se nemusí starat počet a aktualizaci kopií objektů
- perzistentní
 - uživatel se nemusí starat o stav objektů, ke kterým přistupuje
- konkurenční
 - prostředky mohou být automaticky využívány více uživateli
- paralelismová
 - různé akce mohou být prováděny paralelně bez vědomí uživatele

Přizpůsobivost

- autonomie
 - každý uzel je schopný samostatné funkčnosti
- decentralizované rozhodování
 - každý uzel vykonává rozhodnutí nezávisle na ostatních
- otevřenost
 - lze zapojit různé komponenty vyhovující rozhraní
 - oddělení interface x implementace
- migrace procesů a prostředků
 - procesy i prostředky mohou být přemístěny na jiný uzel

Rozšiřitelnost

- škálovatelnost, scalability
- rozdílné metody pro 100 uzlů × 100 milionů uzlů
- princip: vyhnout se čemukoliv centralizovanému
 - žádný uzel nemá úplnou informaci o celkovém stavu systému
 - uzly se rozhodují pouze na základě lokálně dostupných informací
 - výpadek jednoho uzlu nesmí způsobit nefunkčnost algoritmu
 - nelze spoléhat na existenci přesných globálních hodin
 - 'každou celou hodinu si všechny uzly zaznamenají poslední zprávu'

Další problémy

- koordinace a synchronizace
- geografická odlehlost, administrace

Techniky škálovatelnosti

- asynchronní komunikace
- distribuce
- caching a replikace

Výkonnost

- teorie: více uzlů vyšší výkon
- praxe: výrazně nižší než lineární nárůst výkonu
- příčiny: komunikace, synchronizace, komplikovaný software
- granularita výpočetní jednotky

Chyby návrhu distribuovaných systémů

- síť je spolehlivá
- síť je zabezpečená
- síť je homogenní
- topologie se nemění
- nulová latence
- neomezená kapacita sítě
- jeden administrátor

Paralelní architektury

- multiprocesory
 - bus
 - sdílená sběrnice mezi procesory a pamětí
 - levná, běžně dostupná
 - max. desítky uzlů
 - cache – synchronizace
 - switched
 - možné zapojení i většího počtu uzlů
 - mřížka - crosspoint switch
 - sběrnicová mřížka - procesory x paměťové bloky
 - nákladné – kvadratický počet přepínačů
 - omega network
 - postupně přepínaná cesta mezi procesorem a pamětí
 - při větším počtu úrovní pomalé
 - počet přepínačů $n * \log n$

- multicomputery

- bus

- vlastní paměť každého uzlu → výrazně nižší komunikace
 - možnost teoreticky neomezeného počtu uzlů
 - levná, běžně dostupná
 - 'normální' počítače propojené 'normální' sítí
 - lze i jiné provedení: processor poo
 - rack full of cpus in the machine room, which can be dynamically allocated to users on demand

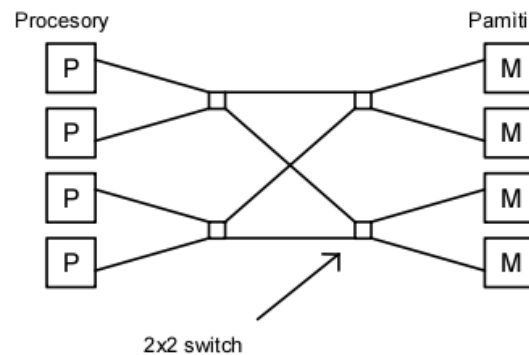
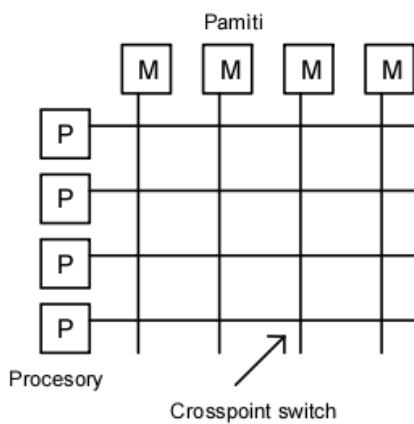
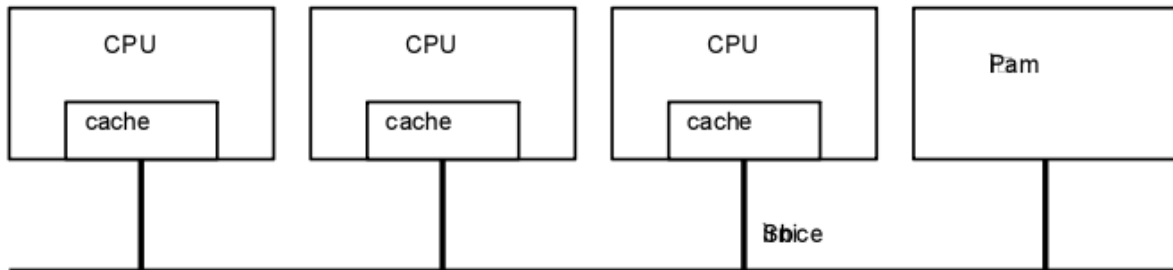
- switched

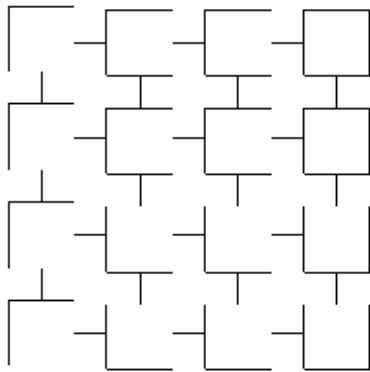
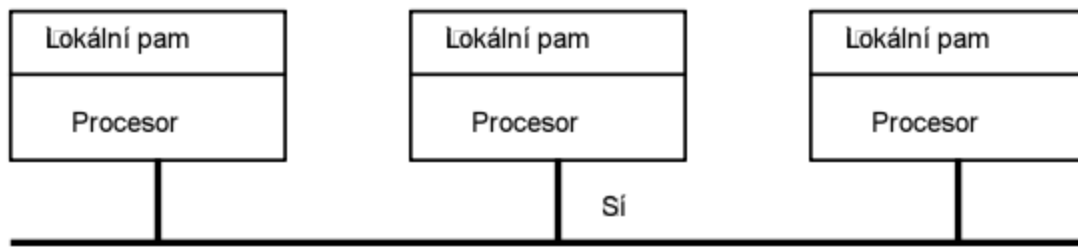
- mřížka

- relativně jednoduchá (dvourozměrná) implementace
 - vhodné pro řešení 'dvourozměrných' problémů – grafy, analýza obrazu, ...

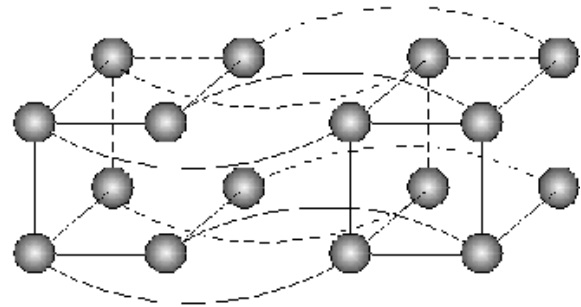
- hyperkrychle

- n-rozměrná krychle, častý rozměr 4
 - každý uzel přímo propojen se sousedy ve všech rozměrech
 - supercomputers - tisíce až miliony uzlů
 - 2012 IBM Sequoia Blue Gene - 1.5 mil cores, 1.5 PB RAM, 8 MW Linux



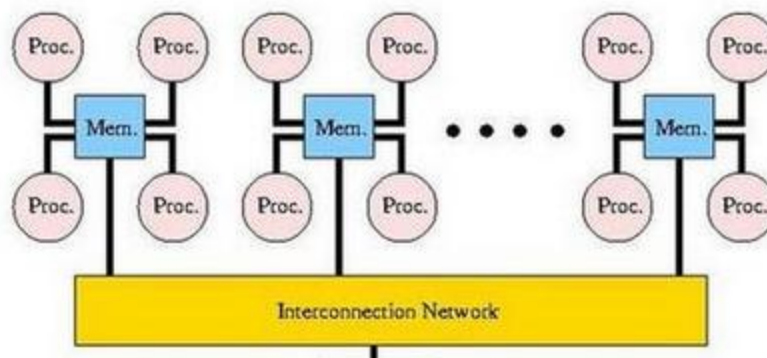


(a)



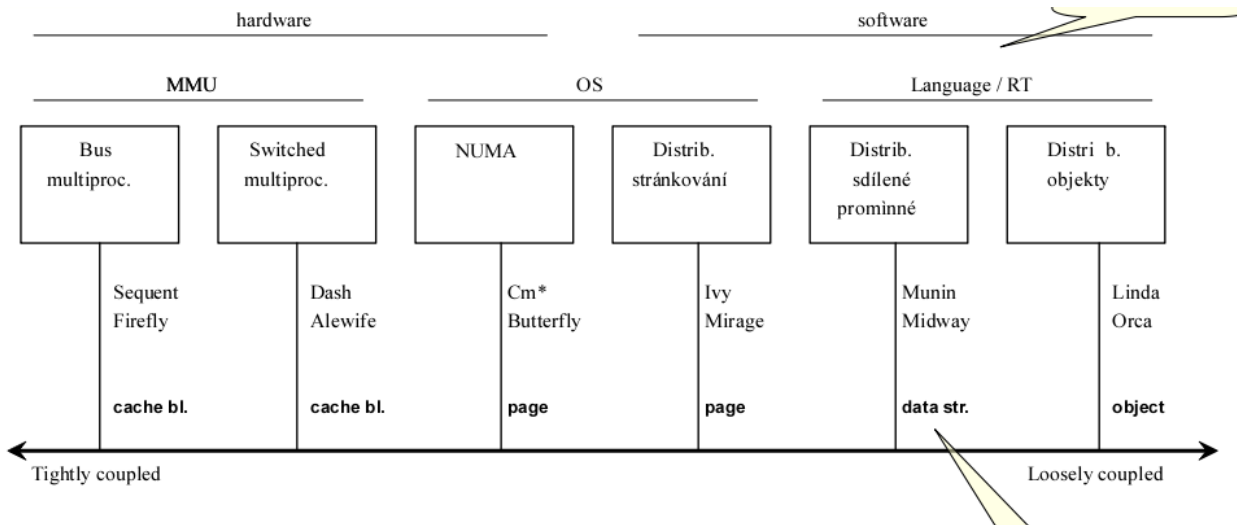
(b)

- NUMA
 - non uniform memory access
 - S-COMA - Simple Cache-Only Memory Architecture
 - požadavek na nelokální data je propagován vyšší vrstvě
 - ccNUMA - Cache Coherent NUMA
 - globální adresace, konzistence cache
 - NUMA-faktor - rozdíl latence lokálních a vzdálených paměťových přístupů



- GPGPU

Mechanismy sdílení paměti

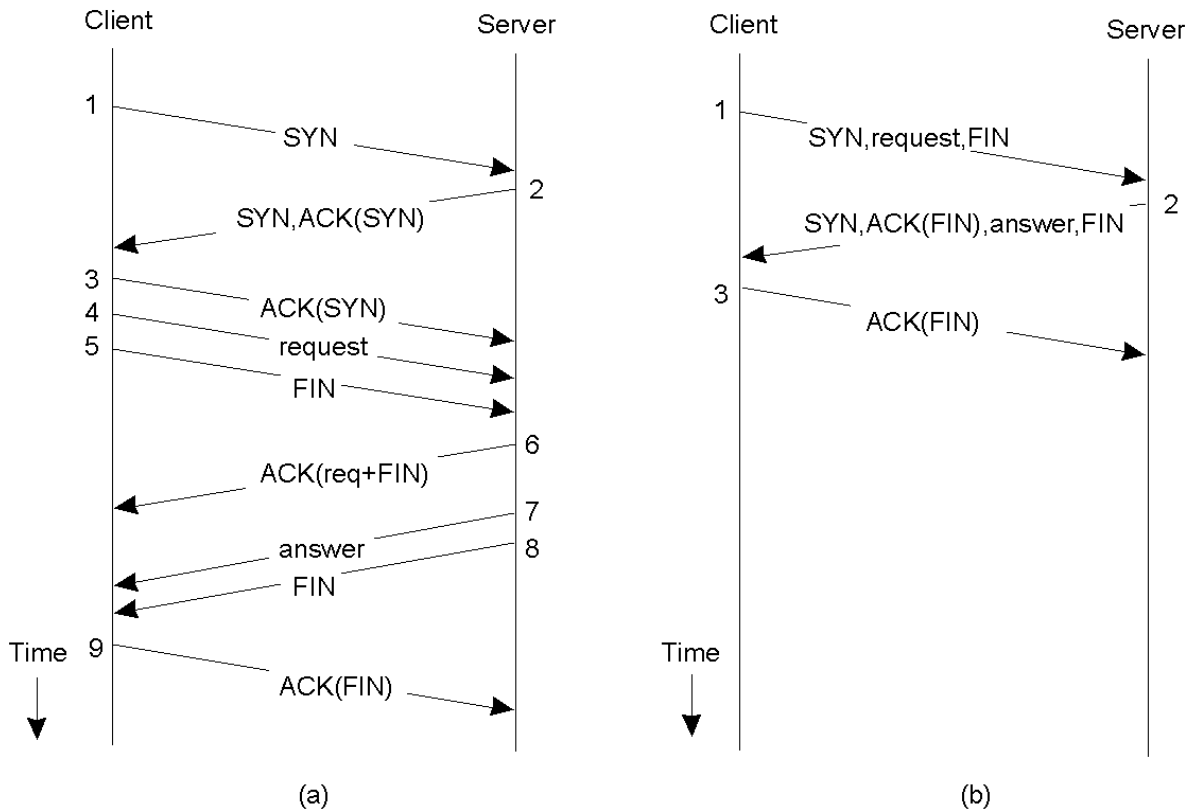


Komunikace

- absence sdílené paměti → zasílání zpráv
 - RPC, doors, RMI, M-queues, ...
- jednotný komunikační mechanismus
 - klient-server model
 - efektivita - lokální / vzdálený přístup

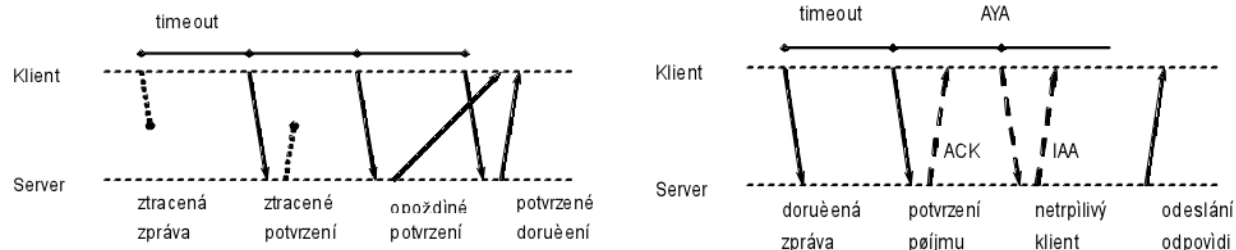
Síťová komunikace

- TCP - standard pro spolehlivý vzdálený přenos
 - ale: většina komunikace je v rychlé a spolehlivé LAN
 - nevýhoda - složitost protokolu, menší výtěžnost sítě
- specializované protokoly pro (spolehlivé) lokální síť
 - optimistické
 - vysoký výkon za cenu složitějšího zotavení z chyb
 - T/TCP, FLIP, ...



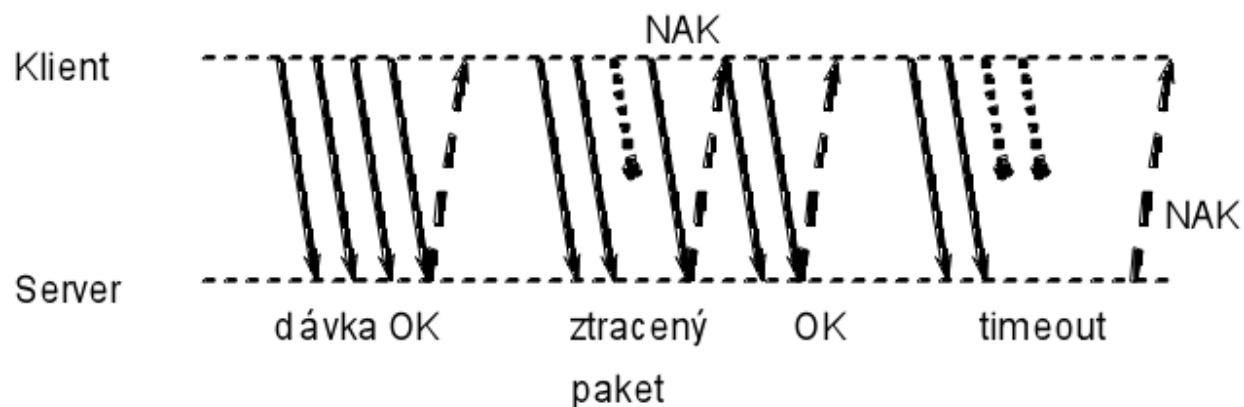
Spolehlivost síťové komunikace

- nespolehlivý unicast – to co nám dává hardware, počítáme s best effort
 - v IPv4 se pakety mohou na routerech fragmentovat po cestě, v IPv6 vyrazí s nějakou délkou a s tou i dorazí (pokud po cestě nějaká část sítě neumí danou délku přenést, tak je zahodí)
- ochrana proti poškození
 - duplikace, checksums, parita, křížová parita, CRC
- ochrana proti ztrátě – potvrzování
- ochrana proti duplikaci – unikátní ID paketů
 - při TCP handshake se dohodne náhodné počáteční a pak roste
- duplicita zpráv - číslování, zahození - poslední zpráva / okno
- předbíhání zpráv - číslování, zpráva o nedoručení - režie, okno
- ztráta zprávy - klient vyslal žádost a neví
 - ztráta žádosti / ztráta odpovědi / příliš pomalá komunikace
 - typické řešení: potvrzování, timeout, opakování
- potvrzování → zvýšení režie → odpověď = potvrzení příjmu
- při delším zpracování - timeout serveru nebo klienta
- služební zprávy - ACK, NAK, AYA, IAA, CON, ...



Přenos dlouhých zpráv

- jednotka přenosu = zpráva
 - příliš dlouhá → rozdělení na pakety
- co potvrzovat
 - každý paket → režie při správném přenosu
 - celá zpráva → režie při špatném přenosu
- možné řešení - dávky (buřty, blast, burst)
 - potvrzování dávky (definovaný počet paketů)
 - v pořádku celá dávka → ACK
 - předbíhání, výpadek → NAK
 - přenos celé dávky / od místa výpadku
 - ztráta posledního (-ích) paketů timeout
 - pevné dávky - jednodušší
 - při vyšší zátěži sítě velká chybovost
 - při volné síti nízká výtěžnost
 - dynamické dávky - úprava velikosti podle aktuální situace
 - n-krát správný přenos dávky → zvětšení
 - špatný přenos → zmenšení
 - vysoká výtěžnost aniž by jeden přenos omezoval ostatní
 - konsensus o velikosti dávky



Nespolehlivost serveru

- klientovi nepřichází odpověď
 - ztratila se zpráva s žádostí/odpovědí
 - server je zahlcen/je pomalý
 - server spadnul/nefunguje
- obecně nelze zjistit co se stalo
 - i když lze, pro některé služby je to náročné - např. transakce

Sémantika zpracování

- idempotentní služby
 - nevadí opakované provedení služby
 - sečti 1 + 1
- exactly once sémantika
 - pro některé služby nelze (tisk na tiskárně)
- at-least-once sémantika
 - služba se určitě alespoň jednou provede
 - nelze zajistit, aby se neprovedla vícekrát
- at-most-once sémantika
 - služba se určitě neprovede vícekrát
 - nelze zajistit, že se provede

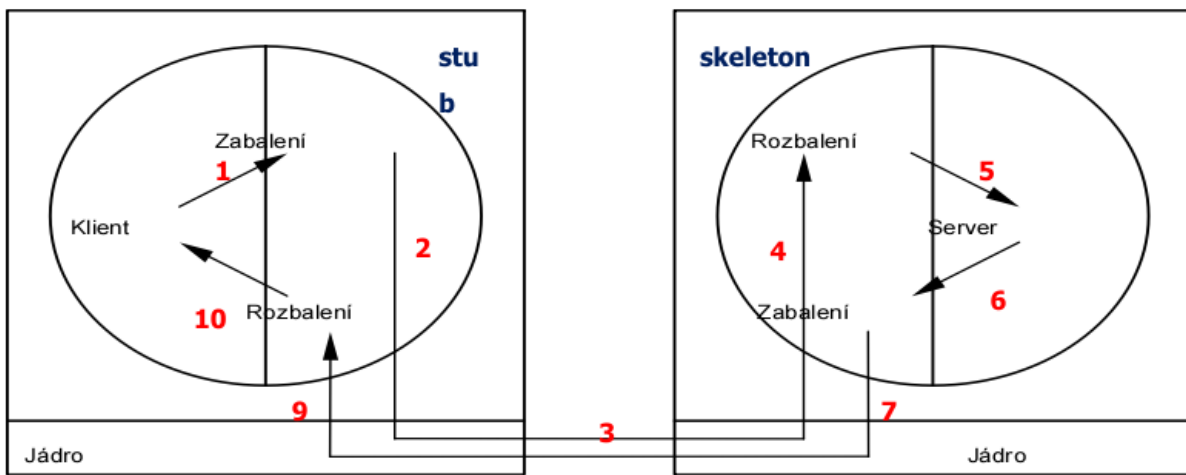
Havárie klienta

- zůstanou sirotci
 - procesy na serveru, které čekají na odpověď klienta, který ale spadnul
- typicky se to neřeší
 - specializované dlouhotrvající výpočty
- možnosti řešení:
 - exterminace
 - předtím, než pošle klient stub RPC zprávu, zalogue co bude dělat
 - log se udržuje na disku nebo na nějakém médiu, které přežije pad
 - po restartu se podívá do logu a sirotek se explicitně zabije
 - reinkarnace
 - čas se rozdelí do posloupnosti epoch
 - když se klient restartuje, započne novou epochu - broadcastuje to zprávou všem strojům
 - když taková zpráva přijde, všechny vzdalene výpočty daného klienta jsou zabity
 - mohlo se stát, že došlo rozdělení site (partition) a nejaci sirotci, přežili, pokud ale budou odpovídat, bude v odpovědi špatná epocha a tím se prokázou

- expirace
 - každý proces má přiděleno kvantum času K
 - pokud nestíhá, musí explicitně požadovat o další - což je nepříjemnost
 - po restartu však může klient počkat K času a pak si může být jist, že všichni sirotci budou zabiti
 - vhodné K

Vzdálené volání pomocí zpráv (RPC)

- myšlenka je v tom, že na klientu se zavolá funkce, která se provede na serveru
- meziprocetová komunikace většinou příliš nízkourovňová
- idea: přiblížit zažitým mechanismům volání procedur
- RPC - remote procedure call
 1. klientova funkce zavola beznym způsobem klientsky stub
 2. stub vytvori zpravu, zavola jadro (marshalling)
 3. jadro posle zpravu uzlu se serverem
 4. vzdalene jadro preda zpravu skeletonu - server-side stubu
 5. skeleton rozbali parametry a zavola vykonnou funkci (unmarshalling)
 6. server zpracuje pozadavky a beznym způsobem se vrati do skeletonu
 7. skeleton zabali vystupni parametry do zpravy a zavola jadro
 8. jadro posle zpravu zpet klientovi
 9. klientovo jadro preda zpravu stubu
 10. stub rozbali vystupni parametry a vrati se ke klientovi



- spojení
 - nejprve se vygeneruje UUID rozhraní
 - programátor dodá definici rozhraní (interface definition file)
 - IDL kompilátor udělá header, co se nainkluduje do klienta i serveru
 - IDL kompilátor taky udělá klientskou a serverovou část, která se stará o komunikaci (client stub, server skeleton)
 - programátor dodá implementaci serverové části
 - celé se to zkompiluje a slinkuje
 - server si zaregistruje u nějakého directory serveru to, že dělá tuhle službu
 - klient se rozběhne, když si chce zavolat RPC, tak:
 - normálně zavolá fci
 - vygenerovaný stub si vytvoří zprávu (marshalling), předá jádru jádro, look-upne službu na directory serveru, předá jí jádru serveru, to ji dá skeletonu, ten si ji rozbálí (unmarshalling), a putuje to podobně zpět
 - rozhraní může být i ze signatury v "normálním" jazyce, jen s nějakými doplňujícími informacemi
 - IDL – popisuje rozhraní funkcí, z něj se pak generují skeletony a stuby; pak mapování do jazyků (je to prog. jazyk)
 - varianta s objekty – máme proxy objekty u klienta a servanty na serveru
 - implementace tohoto objektového cirkusu je třeba RMI
- problémy
 - reprezentace dat
 - endiany
 - nejde to jen tak, že reverzujeme poradi bytu, protože normalne nevime o datech nic a mohli bychom takto reverzovat string, což nechceme → nutné informace navíc
 - floaty
 - různá kodování řetězců
 - přístup ke globálním proměnným
 - neexistence sdílené paměti
 - možné řešení pomocí DSM (distr. sdílení paměti), ale je to těžkopadné
 - předávání ukazatelů, dynamické struktury, předávání polí
 - jde to například tak, že pokud víme že pointer ukazuje do pole a víme jeho velikost tak uděláme jeho kopii a dáme ji do zprávy
 - problém, že obecně nevím kam pointery ukazují nebo jak jsou ty objekty velké
 - může se pak napr. zpětně vyžádat po klientovi kopie daného objektu
 - ztrata transparency, režie navíc
 - skupinová komunikace
 - komunikační chyby
 - bezpečnost

- RPC je užitečné a prakticky použitelné (a používáno), ale není zcela transparentní, takže nutno dávat pozor na omezení
- další mechanismy: MPI, message queues

Skupinová komunikace

- jeden odesílatel – více příjemců
- atomicita
 - doručení všem členům nebo nikomu
 - evidence a změna členství
- synchronizace
 - doručování od různých odesílatelů
 - příjem vs. doručení zprávy
 - doručovací protokol - sémantika
- uzavřená skupina
 - pouze členové skupiny
 - vhodné pro kooperativní algoritmy
- otevřená skupina
 - zasílat zprávy může kdokoli
 - distribuované služby, replikované servery
- problém překrývajících se skupin

Synchronizace

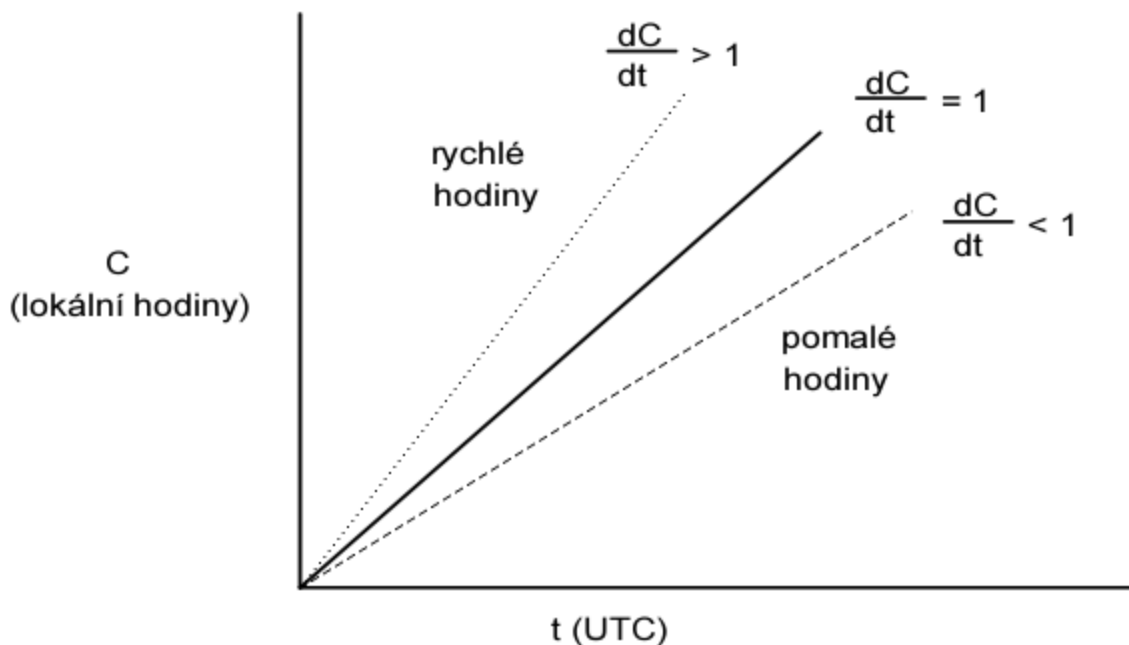
- logické a fyzické hodiny a jejich synchronizace
- distribuované vyloučení procesů
- volba koordinátora
- doručovací protokoly, virtuální synchronie, změny členství ve skupinách
- detekce globálního stavu
- distribuovaný konsensus
- neexistence sdílené paměti → zprávy
- rozprostřená informace mezi uzly
- rozhodování na základě lokálních informací
- vyloučení havarijních komponent
- neexistence společných hodin

Fyzické hodiny

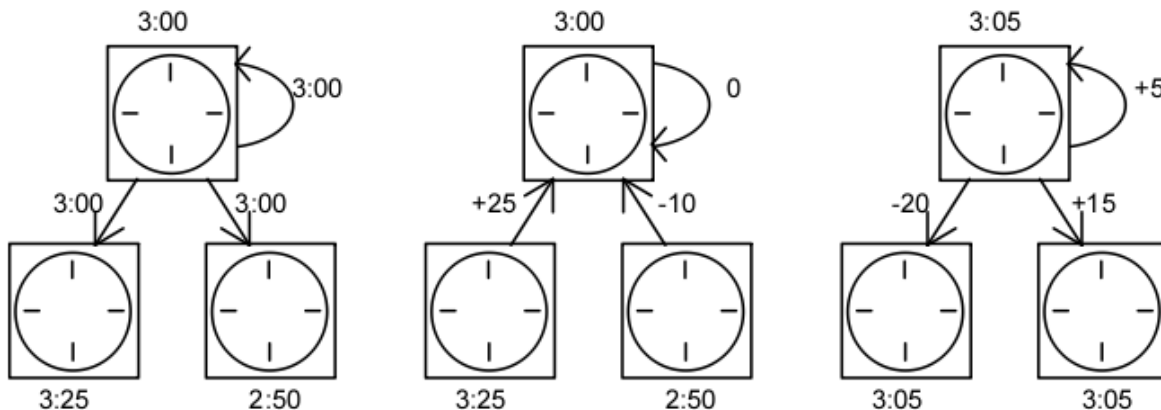
- jak sesynchronizovat lokální hodiny jednotlivých komponent systému
- jak sesynchronizovat jednotlivé hodiny mezi sebou
- astronomické měření → solární sec. = 1/86400 solárního dne
 - avšak slapové tření, atmosferické jevy → nepřesné
- 1948 - atomové hodiny - 1s ≈ 9 mld přechodů atomu cesia 133
 - nezávislé na okolních přírodních jevech, mnohem přesnější
- 1950 - TAI - Bureau International de l'Heure – průměr asi 50 laboratoří
 - international atomic time
 - 1 TAI den je však o 3 ms kratší než solární den
 - řešení: rozdíl >800 ms – vložená sekunda – Universal Coordinated Time UTC
- vysílání UTC (krátké vlny, satelit) – nepřesnost cca 0.1 - 10 ms

Synchronizace fyzických hodin

- jde určit po jaké době je nutné synchronizovat k udržení nějaké maximální odchylky od správného času
 - míra přesnosti ρ : $1-\rho \leq dC/dt \leq 1+\rho$, přesné hodiny $dC/dt = 1$



- Cristianův algoritmus
 - jeden pasivní time server (UTC)
 - periodické dotazy v intervalech $< \delta/2p$ (kde delta je požadovaná snesitelná odchylka)
 - $T = T_UTC + (T1 - T0 - I)/2$
 - $T0$ čas zadosti, $T1$ čas odpovědi
 - pokud známe dobu zpracování (interruptu) I , můžeme ji zohlednit a zpřesnit odhad
 - nikdy nepřehřizovat najednou !
 - (dis)kontinuita
 - postupná změna - serie měření
 - zrychlování nebo zpomalování - dynamicky
- Berkeley algoritmus
 - máme aktivní time server (time daemon) - třeba na jednom ze strojů
 - demon se periodicky ptá každého stroje, jaký má čas a na základě všech odpovědí vypočítá průměrný čas
 - každému stroji pak pošle, o kolik si mají přenastavit čas



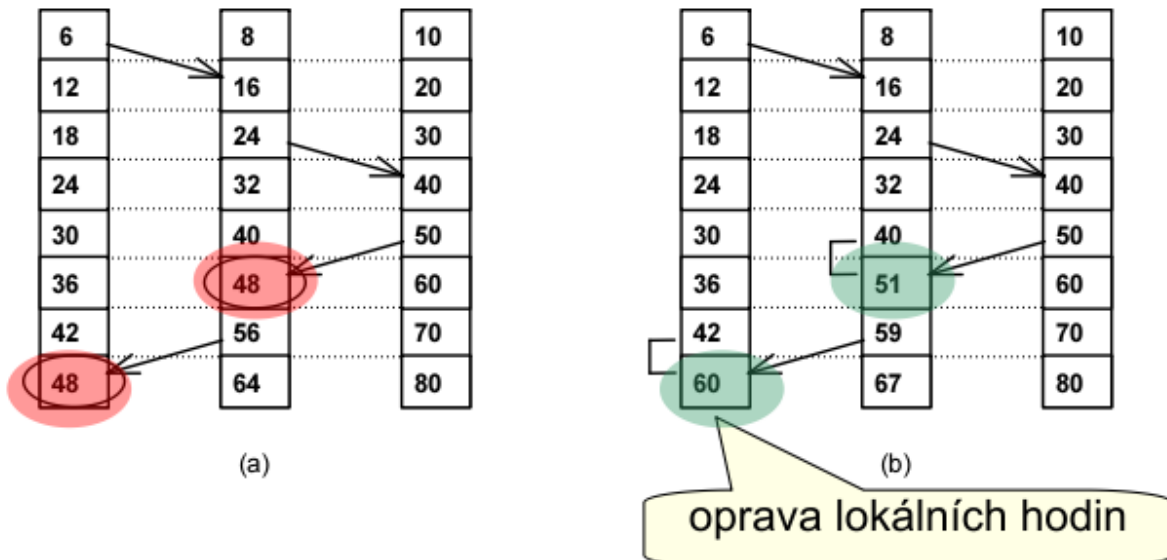
- distribuovaný algoritmus (averaging)
 - předtím metody byly centralizované (se známými negativy)
 - zde žádný server ani koordinátor
 - rozdělení času do resynchronizačních intervalů pevné délky
 - i -ty začíná na $T0 + iR$ a končí $T0 + (i+1)R$
 - $T0$ domluvený moment v minulosti
 - R systémový parametr
 - na začátku každého intervalu každý stroj broadcastuje svůj čas
 - každý stroj po odeslání času zapne svůj lokální timer a v nějakém intervalu S čeká na odpovědi od ostatních
 - když vše dorazí, spočítá se průměr (příp. se zahodí extremy)
- intervalový čas
 - čas není okamžik, ale interval (některé údaje jsou neporovnatelné)
 - v DCE (Distributed Computing Environment - DS z 1990s)

Logické hodiny

- problém: fyzické hodiny nelze dostatečně přesně sesynchronizovat
- Lamport: důležité pořadí, nikoliv přesný čas
 - nekomunikující procesy nemusí být sesynchronizovány
- kauzalni zavislost
 - jestliže $p: e1 \rightarrow p \ e2$ potom $e1 \rightarrow e2$
 - $\rightarrow p$ v rámci procesu p
 - forall $m: \text{send}(m) \rightarrow \text{rcv}(m)$
 - jestliže $(e1 \rightarrow e2 \ \& \ e2 \rightarrow e3) \Rightarrow e1 \rightarrow e3$
 - happens-before relace $a \rightarrow b$
 - všechny procesy se shodnou, že a proběhl před b
 - konkurentni udalosti: $e1 \not\rightarrow e2 \ \& \ e2 \not\rightarrow e1$
 - nic o nich nemůže být receno
- udalost a , čas $C(a)$: pokud $a \rightarrow b \Rightarrow C(a) < C(b)$

Synchronizace logických hodin

- podle přijímání zpráv
 - čas. značka zprávy m : T_m
 - proces i vysílá v čase $C_i(a)$ zprávu $m \Rightarrow T_m = C_i(a)$
 - proces j přijme zprávu m v čase $C_j(b) \Rightarrow C_j = \max(C_j(b), T_m+1)$
 - zuplneni
 - událost a v procesu i , událost b v procesu j
 - $C(a)=C(b) \ \& \ P_i < P_j \Rightarrow C'(a) < C'(b)$
 - 'byrokratické' uspořádání



- platí, že jestliže $a \rightarrow b \Rightarrow C(a) < C(b)$
- neplatí ale opačná implikace, což by se však hodilo
 - jak tedy zjistit kauzalitu událostí \Rightarrow vektorové hodiny

Vzájemné vyloučení procesů

- na jednoprosesorovém systému jsou kritické sekce chráněny semaforemi, monitory a podobnými konstrukcemi
- v DS pouze zprávy
- nutno zajistit korektnost
- řešení:
 - centralizované
 - permission-based
 - časové značky
 - Lamport - $3n$
 - Ricart-Agrawala - $2n$
 - volby
 - Maekawa - \sqrt{n}
 - Agrawala-El Abadi - $\log n$
 - token-based
 - token-passing - Suzuki-Kasami
 - strom - Raymond
 - token ring

Centralizovaný algoritmus

- jeden proces vybrán jako koordinátor
- chce-li někdo do kritické sekce → pošle koordinátorovi zprávu, pokud tam aktuálně nikdo není, pošle mu potvrzení
 - pokud není, tak jsou možné implementace:
 - nepošle mu nic - blokuje ho
 - pošle mu deníček a po opuštění aktuálního mu pošle OK
 - každopádně si udržuje frontu, kam si vkládá procesy, kteří chtějí do dané kritické sekce
- centralizovaná komponenta - ideově nevhodné
 - server může představovat performance bottleneck
 - výpadek serveru - ztráta informací, nutnost volby nového koordinátora
 - výpadek klienta - problém vyhledání (starvace)

Lamport algoritmus

- idea: proces vyšle žádost a čeká až:
 - dorazí odpovědi od všech procesů a zároveň
 - všechny žádosti v jeho frontě mají větší časovou značku (čili on má nejstarší)
- komunikační složitost $3n$

Událost	Akce procesu p	
Žádost Mp	send Mp={req,p,Tp}	
Přijetí žádosti Mi	add Mi; send {ack,p,Tp}	uložení žádosti a odpověď s vlastním časem
Přijetí potvrzení Ai	add Ai	
Podmínka vstupu Mp	$\forall i \neq p \exists \{ack,i,T_i\} : T_p < T_i$ & $\forall \{req,i \neq p, T_i\} : T_p < T_i$	od všech přišlo novější potvrzení neviduju starší žádost
Uvolnění Rp	send {rel,p,Tp}	
Přijetí uvolnění Ri	del {req,i,Tk} : Tk < Ti	po přijetí uvolnění smažu starší žádosti tohoto procesu

Ricart & Agrawala

- vylepsení Lamporta
- proces chce vstoupit do kritické sekce:
 - zašle žádost s TS a čeká na odpovědi s potvrzením
- proces přijme zprávu se žádostí:
 - jestliže není v kritické sekci a ani do ní nechce vstoupit, pošle potvrzení
 - jestliže je v kritické sekci, neodpovídá, požadavek si zařadí do fronty
 - jestliže není v kritické sekci, avšak chce do ní vstoupit, porovná TS žádosti s vlastní žádostí:
 - vlastní žádost má nižší časovou značku (= starší), neodpovídá a zařadí žádost odesílatele do fronty
 - v opačném případě (žádost odesílatele je starší) pošle zpět potvrzení
- po opuštění kritické sekce pošle proces potvrzení všem procesům, které má ve frontě
- komunikační složitost $2n$

Princip voleb

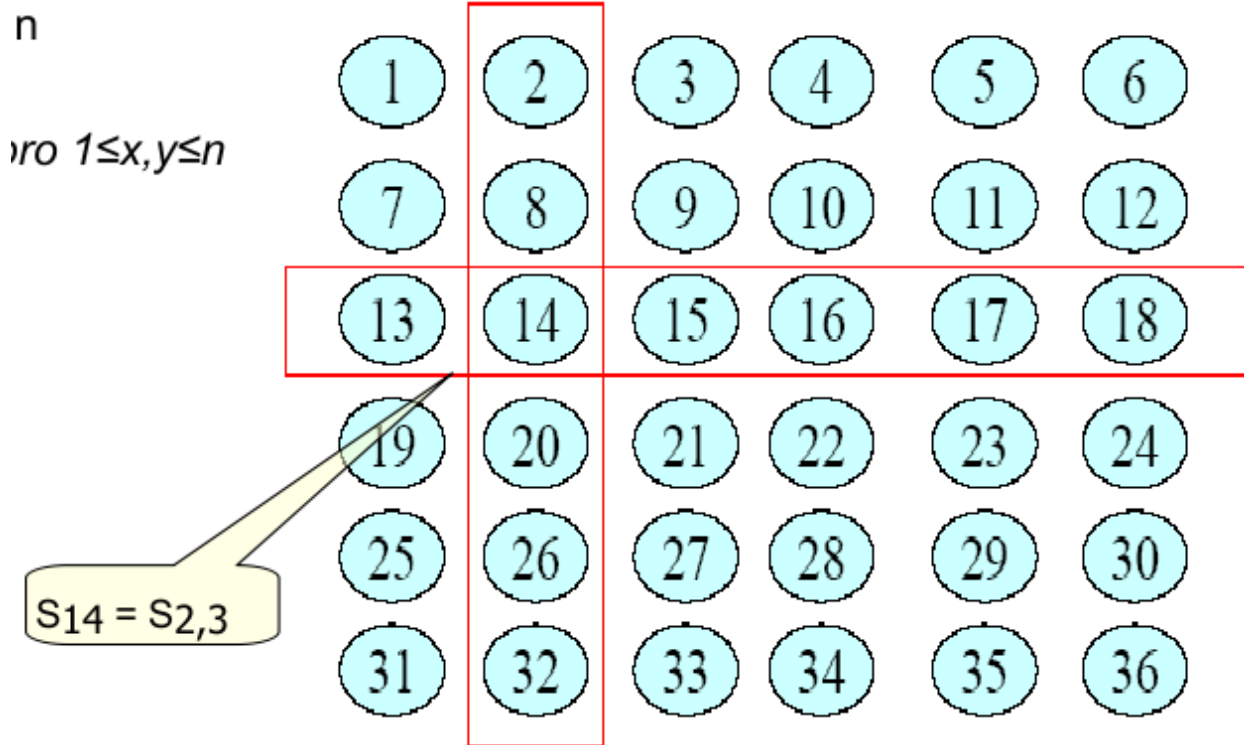
- základní princip:
 - procesy se snaží získat hlasy ostatních procesů
 - každý proces má v jeden okamžik právě jeden hlas
 - může ho dát sobě nebo jinému procesu
 - před vstupem do kritické sekce žádost o hlasy
 - pokud proces dostane víc hlasů než jakýkoliv jiný, může vstoupit
 - pokud dostane méně, čeká dokud nedostane dostatečný počet
 - problém: jak hlasovat a počítat výsledky, kdy už proces ví že vyhrál
- naivní volby
 - každý proces zná všechny ostatní
 - při žádosti jiného procesu dá proces hlas (pokud ještě nehlasoval)
 - relativně odolný proti výpadkům
 - vydrží výpadek až poloviny procesů

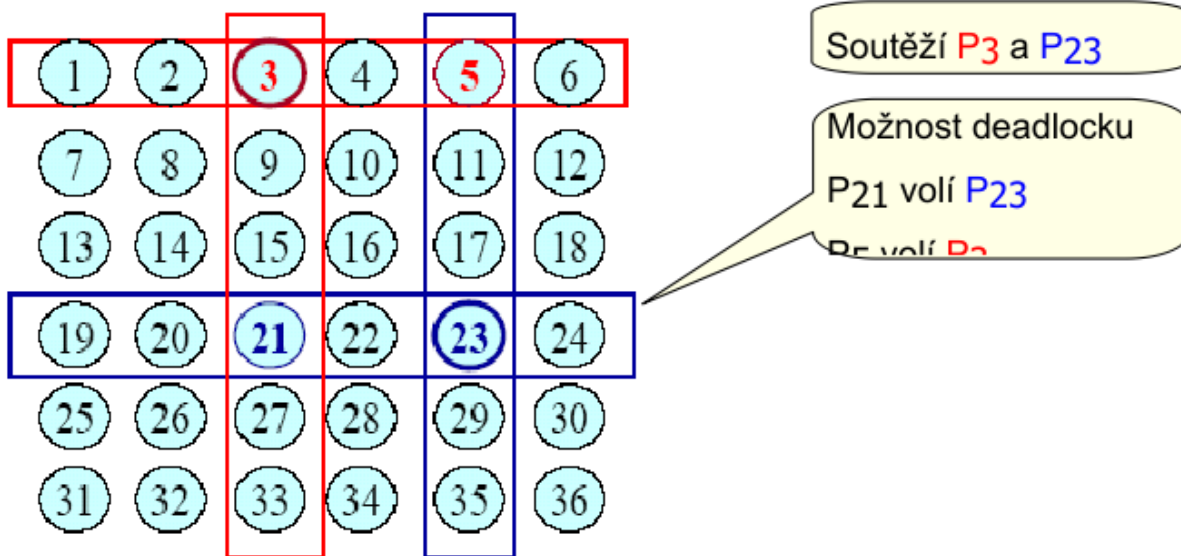
- složitost $O(n)$
 - není lepší než TS
- možnost deadlocku!
 - více procesů může dostat stejný počet hlasů

Maekawa

- organizace procesů pro optimalizaci komunikační složitosti
- každému procesu p je přiřazen volební okrsek S_p (voting district)
- pro vstup do kritické sekce je nutné získat všechny hlasy z vlastního okrsku
- podmínky pro okrsky:
 - forall $p, q: S_p \cap S_q \neq \emptyset$
 - každá dvojice kandidátů má alespoň jednoho společného voliče
 - zaručí korektnost
 - každý proces může volit pouze jednou, nemohou být současně zvoleny dva procesy
 - forall $p, q: |S_p| = |S_q| = K$
 - velikost volebních okrsků je konstantní, procesy potřebují stejný počet hlasů
 - zaručí spravedlnost
 - forall $p, q: |S_i: p \in S_i| = |S_j: q \in S_j| = D$
 - p je obsažen ve stejném počtu D volebních okrsků
 - každý proces má stejnou zodpovědnost
- důsledek: komunikační složitost $O(|S_p|)$
 - chceme tedy minimalizovat velikost okrsků
- jak velké mají být okrsky, v kolika okrscích má být každý proces ?
 - K velikost okrsku, N # okrsku, M # procesu, D (min) # okrsku, ve kterých je každý proces
 - důsledek podmínek: $N \cdot K = D \cdot M$ - TODO: Proc?
 - důležité je uvědomit si, že:
 - proces \approx okrsek - protože proces musí získat všechny hlasy svého okrsku
 - tedy $\sum \text{procesu} \approx \sum \text{okrsku} \rightarrow M = N \rightarrow K = D$
 - protože každý q z okrsku S_p je obsažen v dalších $D-1$ okrscích, tak při zachování podmínky průniku dostáváme max počet okrsku $N = (D-1) \cdot K + 1$
 - a tedy $M = (K-1) \cdot K + 1 \rightarrow K = O(\sqrt{M})$
 - jak rozdělit procesy do okrsku
 - optimalně je to složité - vyžaduje restrukturalizaci při změně členství
 - suboptimálně (při zachování $O(\sqrt{M})$ jednoduché)
 - $M = n^2, P_{i,j}, S_{i,j}, 1 \leq i, j \leq n$
 - $S_{i,j} = (U_{P_i, x}) \cup (U_{P_y, j})$, pro $1 \leq x, y \leq n$
 - může dojít k deadlocku

- řešení je použít Lamportovy hodiny
 - při příjmu žádosti procesu r s TS_r procesem p
 - p má volný hlas: potvrzení ACKr
 - p dal již hlas jinému procesu q s $TS_q < TS_r$: zařadit do fronty
 - p dal již hlas jinému procesu q s $TS_q > TS_r$: pošle zprávu REJECT procesu q
 - race condition?
 - když q ještě není v kritické sekci - vrati hlas
 - přidělení procesu s nejnižší TS
 - pokud už je, vrati ho při opuštění CS





Agrawal & El Abbadi

- volby v $O(\log n)$
- kvorum je cesta od kořene k listu
- fault tolerance - nahrazení havarovaného uzlu jeho podstromy

Raymond

- token-based, $O(\log n)$
- kostra, nezakorenený strom
- každý uzel si udržuje referenci na souseda směrem k token holderu
- při přenosu tokenu se přesměrují reference - rekonfigurace

Token ring

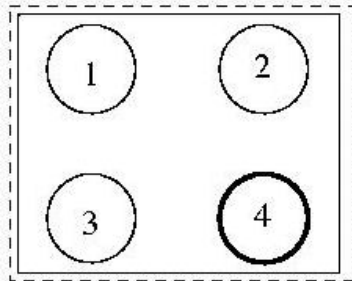
- token-based, komunikační složitost 1 až ∞
- logický kruh procesu, nějak očíslovane, každý ví kdo následuje
- přenos zprávy obsahující povolení ke vstupu do kritické sekce
 - dostane token
 - chce vstoupit - provede CS a posle token dál
 - nechce - posle dál
- problémy:
 - dlouhý čekací čas na token
 - token se ztratí - nutno regenerovat
 - špatně se detekuje i díky minulemu problému
 - výpadek procesu
 - detekce - soused mu chce dát token a nejde to
 - vymaze se z kruhu, rekonfigurace

Bully algoritmus

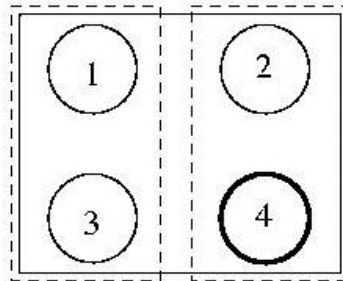
- volba koordinatora
- predpokladem je synchronni spolehlivy system
- $T_m \sim$ message propagation time, $T_p \sim$ message handling time
- omezená doba přenosu zprávy T_m
- omezená doba zpracování zprávy a odpovědi T_p
- detektor havárie v čase: $2T_m + T_p$ (velmi silný požadavek)
- když se proces rozhodne volit, zašle zprávu všem procesům s vyšší identifikací (číslo procesu)
 - když přijde odpověď, proces končí
 - když nepřijde nic, proces vyhrál, je novým koordinátorem, pošle zprávu všem ostatním
- když proces přijme zprávu o volbě, vrátí zpět odpověď a vyšle žádosti všem vyšším procesům
- volba ve dvou kolech
 - přijde víc odpovědi od procesů s vyšším PID a ty pak mají mezi sebou vyřešit kdo je koordinátor ve druhém kole
- při překročení časových limitů možnost více koordinátorů !

Invitation algoritmus

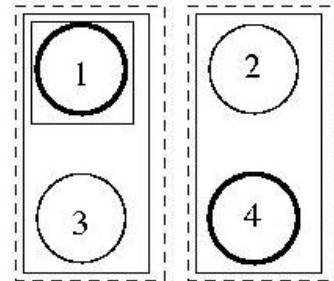
- volba koordinatora
- asynchronni system - nelze predpokladat nic o delce udalosti
- pouziti v "reálnem" prostredí
 - procesor: může havarovat (fail-stop), neomezená doba reakce
 - komunikace: rozpojení skupiny na segmenty, ztráta zpráv
 - nelze spolehlivě detekovat havárii: absence odpovědi neznamena havárii
- koordinátor je vázán na skupinu, skupiny lze štěpit
 - všichni členové stejné verze skupiny (pohled, view) vidí stejného koordinátora
- koordinátor: pravidelná výzva AreYouCoordinator
 - pokud dojde k příjmu AYC koordinátorem \Rightarrow sjednocení do skupiny koordinátora s vyšším ID
- pokud člen skupiny neobdrží AYC svého koordinátora (dostatečně dlouho)
 - prohlásí se za koordinátora vlastní nové skupiny
 - posílá pozvání ke sjednocení ostatním novým skupinám
- konzistence relativní vzhledem ke skupině
 - procesy se shodují na členství ve skupině a na nějaké hodnotě
 - koordinátor vyzve ostatní k připojení
 - pokud se nepřipojí, jsou konzistentní samy se sebou



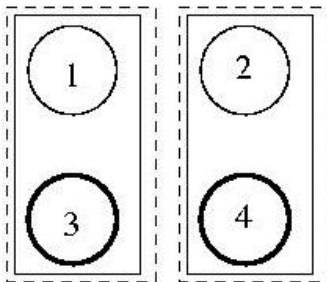
One partition. One group
Coordinator = Processor 4



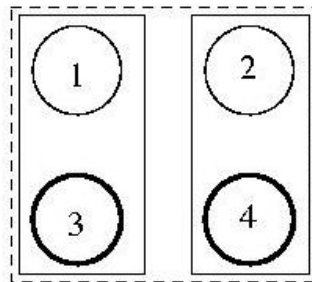
The network is partitioned, but
no one notices



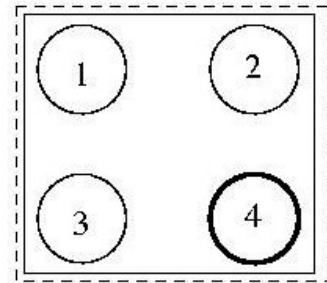
Processor 1 notices the partitioning,
and declares itself a coordinator.
It calls out and reaches processor 3.



Processor 3 realizes the partitioning
and becomes a coordinator.
Processor 1 retires from this role.



The network partition is repaired,,
but none of the processors notice



Either processor 3 or processor 4,
discover that the network has been
repaired. Processor 4 responds by
announcing that it is coordinator
and restoring global consistency.

Kruhový algoritmus

- volba koordinatora
- proces se rozhodne volit (koordinator je nedostupný), pošle zprávu následníkovi
- zpráva obsahuje čísla procesů - kandidáti na koordinátora
 - každý při obdržení zprávy přidá sebe na seznam
- po návratu k iniciátorovi obsahuje zpráva nového koordinátora
- následuje fáze oznámení nového koordinátora
 - změní se typ zprávy
 - koordinátor je ten s nejvyšším id
 - zároveň zpráva obsahuje seznam žijících - tak jak se zapsali na seznam
- stačí znát následníky a mít možnost zjistit následníka nedostupného uzlu
- složitost $O(n^2)$
 - optimalizace - koordinator jen sveho okoli $2^k \rightarrow O(n \log n)$
 - implementacne narocnejši
 - vyhodne pro velke skupiny a vysoke frekvenci konkurencnich voleb - caste vypadky

Doručovací protokoly skupinové komunikace

- globální uspořádání
 - zprávy jsou doručovány v pořadí odeslání
 - nelze – neexistence globálních hodin
- sekvenční uspořádání
 - všechny uzly doručí zprávu ve stejném pořadí
 - doručení nezávisí nutně na času odeslání
 - sekvencer, dvoufázový distribuovaný alg.
- atomický multicast
 - spolehlivé sekvenční doručení
- kauzální uspořádání
 - kauzálně vázané zprávy ve správném pořadí
 - konkurenční zprávy v libovolném pořadí
- kauzálně vázané zprávy
 - obsah vázané zprávy může být ovlivněn
 - není podstatné, jestli obsah byl ovlivněn
- konkurenční zprávy
 - ty, které nejsou kauzálně vázané

Kauzální doručování

- kauzální závislost - viz logické hodiny (zpet)
- kauzální uspořádání doručovaných zpráv
 - $\text{dest}(m)$ množina procesů, kterým je zaslána zpráva m
 - $\text{deliver}_p(m)$ je událost doručení zprávy m procesu p
 - $m_1 \rightarrow m_2 \Rightarrow \text{forall } p \text{ in } (\text{dest}(m_1) \cap \text{dest}(m_2)) : \text{deliver}_p(m_1) \rightarrow_p \text{deliver}_p(m_2)$
 - jestliže je m_2 kauzálně závislá na m_1 , potom na všech procesech, kterým se doručí obe tyto zprávy, se zachová pořadí tohoto doručení

Vektorové hodiny

- máme vektor délky n , kde n je # procesu ve skupině
- uvažujeme časovou značku procesu $VT(p)$ a značku zprávy $VT(m)$
- pravidla aktualizací časových značek
 - $\text{start} \rightarrow VT(p_i) = 0^n$
 - dojde-li k nějaké interní události v procesu p_i , inkrementuje si vlastní logické hodiny - tedy: $++VT(p_i)[i]$
 - vždy, když proces p_i připravuje posílání zprávy, tak si inkrementuje vlastní hodiny a celý svůj vektor zabalí do zprávy: $\text{send } p_i(m) \sim VT(m) = ++VT(p_i)[i]$
 - p_j si při doručení zprávy m upraví $VT(p_j)$:
 - po složkách maxima z přijatého a svého vektoru
 - $\text{forall } k: VT(p_j)[k] = \max(VT(p_j)[k], VT(m)[k])$
 - obecná pravidla nerikají KDY k doručení dojde !
 - libovolné dvě zprávy nemají stejné vektory:
 - $\text{forall } m_i, m_j (i \neq j): VT(m_i) \neq VT(m_j)$

- porovnávání časových značek
 - ne každá dvojice VT musí být porovnatelná - tj. mohou být konkurentní
 - $VT1 \leq VT2 \Leftrightarrow \text{forall } i: VT1[i] \leq VT2[i]$
 - $VT1 < VT2 \Leftrightarrow VT1 \leq VT2 \ \&\& \ \text{exist } i: VT1[i] < VT2[i]$

Doručovací protokol pro jednu skupinu

- využití VT
- odeslání zprávy m procesem p_i
 - klasicky
 - $VT(p_i)[i]++$
 - $VT(m) = VT(p_i)$
- přijetí zprávy m od procesu p_i procesem p_j
 - proces $p_j \neq p_i$ pozdrží doručení m a posle jej dále, pokud dojde ke splnění následujících podmínek:
 - $VT(m)[i] = VT(p_j)[i] + 1$ a
 - $VT(m)[k] \leq VT(p_j)[k]$ pro $k \neq i$
 - tj. toto jsou podmínky doručení
 - např. proces p_v poslal zprávu s článkem $\sim a$ a proces p_j posílá reakci na tuto zprávu $\sim m$
 - první podmínka $\sim m$ je následující zpráva, kterou proces p_j čeká od p_i
 - druhá podmínka $\sim p_j$ neviděl žádné zprávy, které neviděl p_i ve chvíli posílání zprávy m, tzn. p_j již viděl zprávu a (viděl již vše, co viděl p_i při odesílání zprávy m)
- doručení zprávy
 - po doručení si proces p_j upraví VT
 - klasicky
 - forall $k: VT(p_j)[k] = \max(VT(p_j)[k], VT(m)[k])$

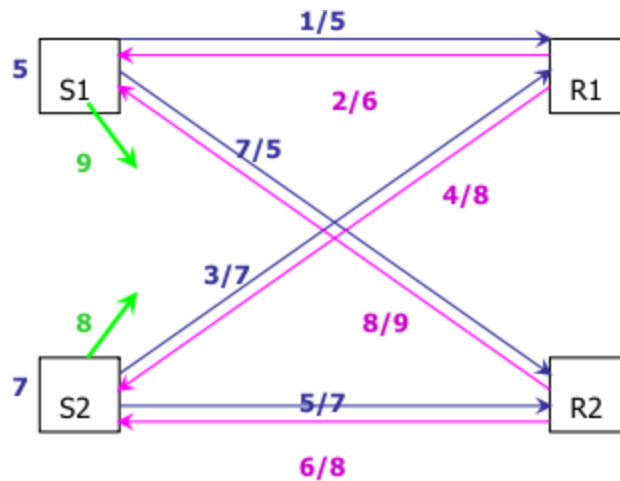
Doručovací protokol pro překrývající se skupiny

- maticový protokol
 - vektory vektorových hodin
- VTa - časová značka skupiny g_a
- $VTa[i]$ - počet multicastů procesu p_i do g_a (zpráv, které vyšle do skupiny)
- s odesílanou zprávou posílá proces VT všech skupin, kterých je členem
- odeslání zprávy m procesem p_i do skupiny g_a
 - $VTa(p_i)[i]++$
 - $VT(m) = \cup VTb(p_i)$, ze p_i in g_b (všechny skupiny, do kterých p_i patří)
- přijetí zprávy m od procesu p_i procesem p_j ze skupiny g_a
 - proces $p_j \neq p_i$ pozdrží doručení m a posle jej dále, pokud dojde ke splnění následujících podmínek:

- $VTa(m)[i] = VTa(pj)[i] + 1$
 - forall k (pk in ga & $k \neq i$): $VTa(m)[k] \leq VTa(pj)[k]$
 - kauzalita vzhledem ke skupina ga , kam byla zprava odeslana
 - forall b (pj in gb): $VTb(m) \leq VTb(pj)$
 - kauzalita vzhledem k ostatnim skupinam prijemce (pj)
- doruceni zpravy
 - po doruceni si proces pj upravi VT
 - klasicky

Distribuvany total-order protokol

- totalni nebo sekvenčni dorucovani
 - sekvenčni dorucovani formou centralizovane komponenty - sekvencer
 - efektivnejsi
 - komunikacne jednoduchsi
 - nutnost vyberu koordinatora
 - serializace/usporadani dle prijmu sekvencera
 - moznost prioritniho dorucovani
- bez ohledu na typu dorucovaciho protokolu pouziteho pro zpravu je navic nutne, aby byly zpravy pri doruceni doruceny vsem clenum skupiny ve stejnem poradí
- stačí skalární hodiny (časové značky) - logické hodiny
- při příjmu zprávy potvrzení odesílateli $TS(R_i)$
- odesílatel po příjmu všech potvrzení odešle finalizační zprávu s $TSF = \max(TS(R_i))$
- po příjmu finalizační zprávy doručí příjemce zprávy podle TSF



Spolehlivé doručování

- naivní definice - všem členům skupiny bude doručena zpráva
- naivní implementace - poslu spolehlivým kanálem všem členům skupiny zprávu
 - problémy: výpadek příjemce/odesílatele
- transakce
 - doručení pouze po commitu
 - funguje, ale příliš striktní
 - nepřijatelné důsledky:
 - odesílatel neobdrží potvrzení - abort celého procesu
 - i když by stačilo něco méně drastického, např. redukce skupiny
 - havárie odesílatele po odeslání všech zpráv

Flooding algoritmus

- při příjmu každé doposud nepřijaté zprávy ji každý uzel přepošle **všem** ostatním
- nakonec všichni zbývajcí členové zprávu přijmou
- spolehlivý, neefektivní

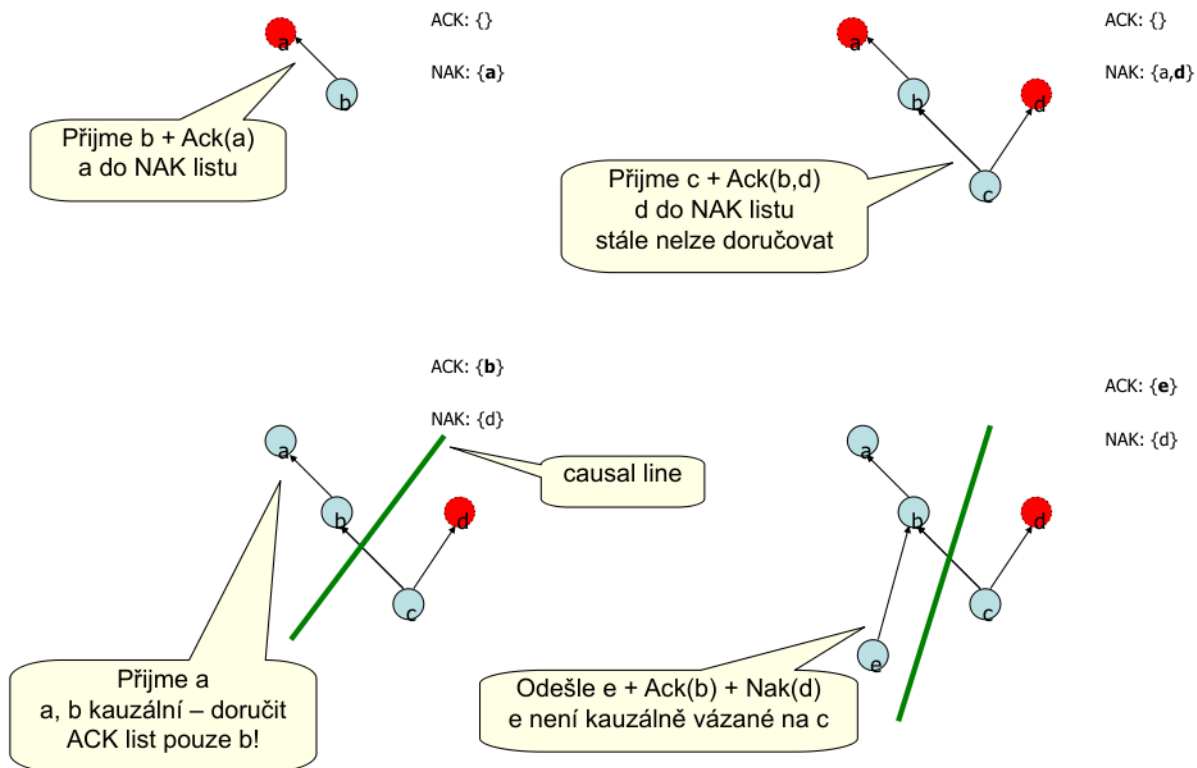
Algoritmus s potvrzováním

- jak to udelat:
 - p_i odesílatel zprávy, p_j in L příjemce zprávy, p_x havarovaný uzel
 - p_i odešle zprávu for all p in L , zprávu si uchová až do obdržení $Ack(p)$ for all p in L nebo do zjištění že p_x havaroval
 - p_j po příjmu zprávy odešle $Ack(p_j)$ uzlu p_i , zprávu si uchová až do zjištění že zprávu přijaly all p in L
 - jestliže p_j zjistí, že p_i havaroval, odešle zprávu for all p in L o kterých neví, že zprávu přijaly
- problém s tím, jak zjistit, které uzly zprávu přijaly a které ne

Trans algoritmus

- protokol pro spolehlivé kauzální doručování
- pro potvrzování se používá graf závislosti - DAG
- používá se invariant:
 - p rozesílá $Ack(m)$ když přijal m a všechny kauzálně předcházející zprávy
 - transitive acknowledgements
 - není nutné potvrzovat zprávy kauzálně předcházející m
- stabilní zpráva je ta, kterou přijali všichni ze skupiny
- DAG G bude graf kauzality skupiny
- G_p - všechny zprávy které p přijal a ještě nejsou stabilní (ještě je nemají všichni)
- jestliže m potvrzuje zprávu m' pak (m, m') in G
 - hrany orientovány ve směru potvrzení (ne ve směru kauzality)
 - vypočítáváme si graf na základě doslych potvrzení
- zdrojem informací pro negativní potvrzování je G (nepřijatá zpráva)
- jak p_j zjistí které uzly zprávu přijaly

- možnosti:
 - rozsilani Ack(p) nejen odesilatel, ale vsem - neefektivni, mnoho zbytecných potvrzení
 - vyuziti kauzality zprav, piggybacking potvrzení
 - D muze po prijmu a, b+Ack(a), c+Ack(b) odvodit, ze B prijal a, C prijal a i b
 - D muze po prijmu b+Ack(a), c+Ack(b) odvodit, ze B prijal a, C prijal a i b, A odeslal a -- vyzadame si zaslani
 - pri korektnim kauzalnim dorucovani jsou potvrzeni zprav tranzitivni (viz invariant)
- implementace:
 - ack_list, nak_list - seznam zprav pro doruceni a seznam neprijatých zprav
 - undelivered_list - seznam prijatých, ale jeste nedorucených zprav
 - zpravu preposilame vcetne ack_listu, nak_listu

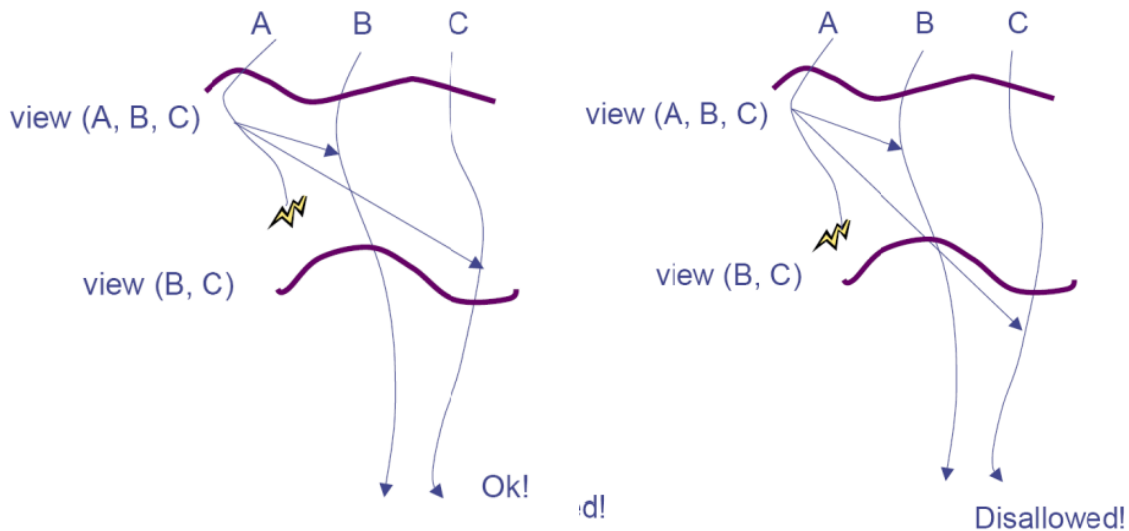


- pohled na protokol - posouvame stable line a causal line grafem
 - nove prichazi zprava vne causal line
 - pokud vsechny kauzalne predchazejici zpravy jsou uvnitr, pak i tuto zpravu lze dat dovnitr

- pozorování:
 - je-li zpráva doručena spolehlivému procesu, pak někdy každý nehavaruující proces zprávu dostane
 - zprávy jsou doručovány v kauzálním uspořádání
 - G reprezentuje graf závislostí v kauzálním uspořádání
 - všechny procesy vytvářejí stejný graf závislostí
 - pořadí přidávání uzlů však může být různé
 - jestliže procesor havaruje, paměťová náročnost je neomezená!
 - slabina - pro praktické použití musí být Trans protokol doplněn nějakým protokolem pro změnu členství ve skupinách

Virtuální synchronie

- group view / delivery list / ... - množina uzlů ve skupině
 - množina procesu naležící skupině, kterou měl odesílatel při posílání zprávy m
- znacení: L , L_i , L^x , L_i^x - globalní pohled, lokální pohled procesu i, verze pohledu x, verze lokálního pohledu x pro proces i
- máme p, q in L_x , pohled se mení na L_{x+1}
- $\text{install } L_{px} < \text{deliverp}(m) < \text{install } L_{px+1} \Rightarrow \text{install } L_{qx} < \text{deliverq}(m) < \text{install } L_{qx+1}$
 - pokud je zpráva m odeslána skupině s pohledem L_x před změnou na L_{x+1}
 - všechny uzly z L_x doručí m před provedením změny na L_{x+1}
 - nebo žádný uzel z L_x , který provede změnu na L_{x+1} zprávu m nedoručí
- podmínka vzájemné konzistence
 - $p \text{ in } L_q \Rightarrow q \text{ in } L_p$
 - všechny uzly ve skupině udržují stejný pohled L
 - instalují pohledy ve stejném pořadí



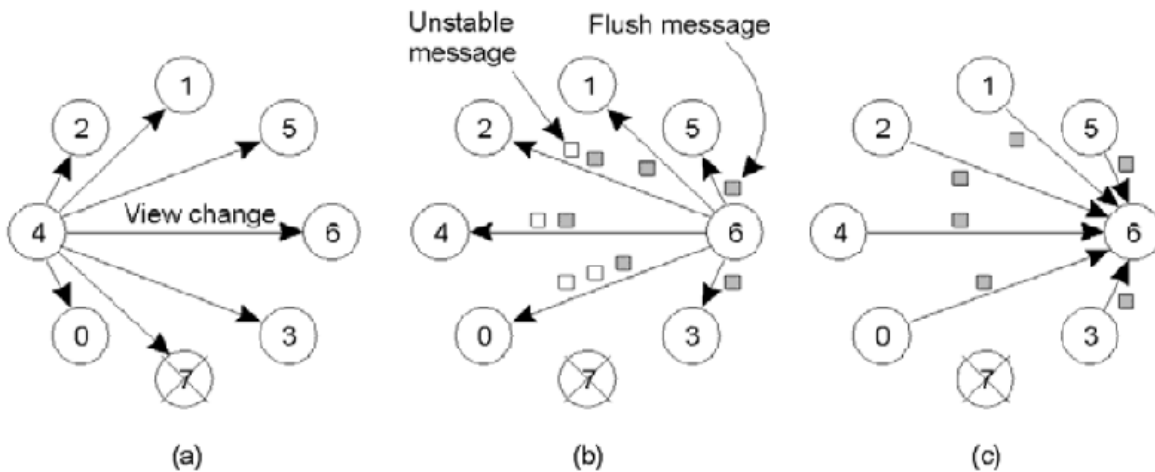
Transis algoritmus

- spolehlivý kauzalni multicast + členství ve skupinách
- protokol je monotonní
 - nemožnost dosažení distribuovaného konsensu (později)
 - paranoia - jedno podezření stačí k vyloučení
 - k vyloučení stačí větší zpoždění - nerozeznatelné od havárie
 - jednosměrnost - vyloučené procesy se již do algoritmu nevrací
 - lze se ale explicitně znovu připojit
- idea - konzistentní změny pohledů, doručování v rámci pohledů
 - při detekci havárie zpráva FAULT(q), při přijetí její přeposlání
 - kauzální hranice pohledu
 - vynutí doručení předcházejících zpráv
 - pozdržení zpráv kauzálně následujících
 - konkurentní detekce různých havárií
 - stanou-li se dvě havárie zároveň → společná hranice
 - doručení obou zpráv
 - doručení pozdržených zpráv
 - kauzální doručení zpráv havarovaného procesu vzhledem ke změně pohledu
 - zprávy odeslané kauzálně před zjištěním havárie se doručí v L_x
 - zprávy odeslané konkurentně se zjištěním havárie se zahodí
 - zprávy odeslané kauzálně po zjištění havárie se zahodí

Protokol ISIS

- spolehlivé kauzalni doručování
- základem doručovacích protokolů je způsob přenosu informace o doručení
 - trans alg. ~ kauzalni potvrzení
 - ISIS ~ maticové hodiny
- idea - každý proces zná $VT[]$ všech procesů - udržuje matici $MT_{pi[j]}[k]$
 - co proces pi ví o doručení zpráv procesu pj od pk
 - tj proces si ukládá:
 - vektor s odeslanými zprávami
 - při přijetí nějaké další zprávy od příjemce čas odeslání poslední zprávy od něj, kterou příjemce dostal
 - a zároveň vektory co mu došly od všech ostatních
- pj si při příjmu zprávy od pi aktualizuje $MT_{pj[j]}$ (co vi příjemce o sobě) a $MT_{pj[i]}$ (co vi příjemce o odesílateli)
 - $MT_{pj[j]}[i] = MT_m[i]$
 - $MT_{pj[i]}[*] = MT_m[*]$
- stabilní zprávy (přijaté všemi členy skupiny) - byly doručeny MT od všech členů skupiny
- srovnání s trans
 - trans potřebuje potvrdit max. M jiných zpráv - lze nahradit MT
 - trans lze považovat za formu komprimace MT

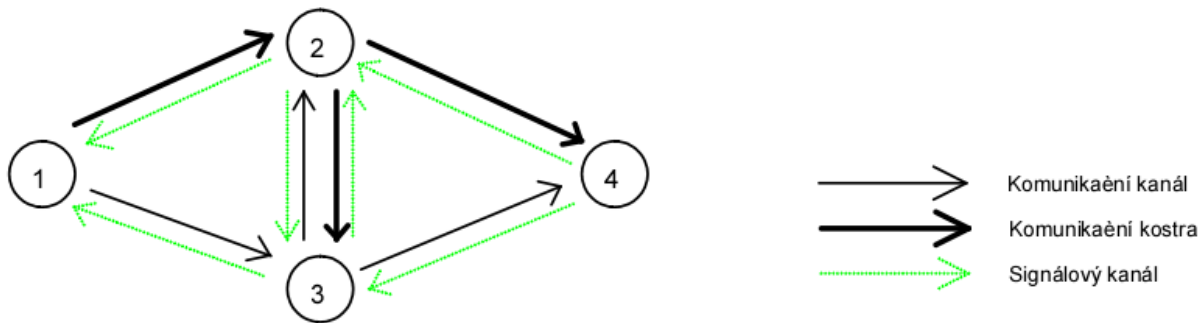
- optimistický algoritmus
 - relativně nízká režie při správném přenosu, vysoká režie při chybách
- cílem je, že všechny zprávy zaslané do pohledu L jsou doručeny všem žijícím procesům v L před instalací nového pohledu
- každý proces v L udržuje zprávu dokud není stabilní - detekce pomocí MT
- po příjmu zprávy o instalaci nového pohledu proces
 - přepośle všechny nestabilní zprávy
 - pošle flush message - potvrzení instalace
 - zatím pohled neinstaluje
- po příjmu flush od každého procesu může instalovat nový pohled
- zprávy od havarovaných procesů
 - každý proces udržuje seznam 'havarovaných' procesů aktuálního pohledu
 - s každou zprávou se rozesílá seznam, při příjmu se sjednotí s vlastním
 - zprávy od havarovaných procesů se zahazují
- reálné nasazení
 - New York Stock Exchange, Swiss Stock Exchange
 - French Air Traffic Control System
 - Reuters, Bloomberg



Ukončení distribuovaných systémů

- orientovaný graf uzlů se vstupními / výstupními kanály + signální kanály pro signalizaci ukončení
- iniciační proces začíná výpočet, vyšle podél svých výstupních hran zprávu
- proces se při příjmu zprávy dostane do stavu výpočet
 - může dále posílat zprávy podél výstupních hran
 - může přijímat zprávy ze vstupních hran
 - když proces skončí výpočet, dostane se do stavu ukončen, zprávy neposílá, může přijímat
 - při přijetí zprávy se proces zrestartuje do stavu výpočet

- algoritmus ukončení
 - všechny procesy jsou ve stavu ukončen, pak se v konečném čase toto všechny procesy dozví
 - stačí, když se to dozví iniciační uzel (může ostatním uzlům poslat zprávu)



Dijkstra-Scholten algoritmus

- pro strom
 - každý listový proces při přechodu do stavu ukončen pošle signál svému otci
 - proces dostane signály od všech svých synů, pošle signál svému otci
 - iniciační proces dostane všechny signály - konec výpočtu
- pro DAG
 - s každou hranou je asociován čítač - tzv. deficit
 - rozdíl mezi počtem zpráv došlých tímto datovým kanálem a počtem signálů poslaných signálním kanálem zpět
 - končící proces vyšle každým signálním kanálem tolik signálů, aby deficit = 0
- obecný graf
 - problém: neexistují listy, které by mohly rozhodnout o ukončení výpočtu
 - řešení: výpočet vytvoří dynamicky kostru grafu
 - otec = uzel, od kterého přišla první zpráva
 - algoritmus ukončení pro jeden proces:
 - 1. poslat signál podél všech vstupních hran kromě hrany k otci
 - 2. čekat na signál od všech výstupních hran
 - 3. poslat signál otci
 - iniciační proces dostane všechny signály - konec výpočtu

Značkový algoritmus

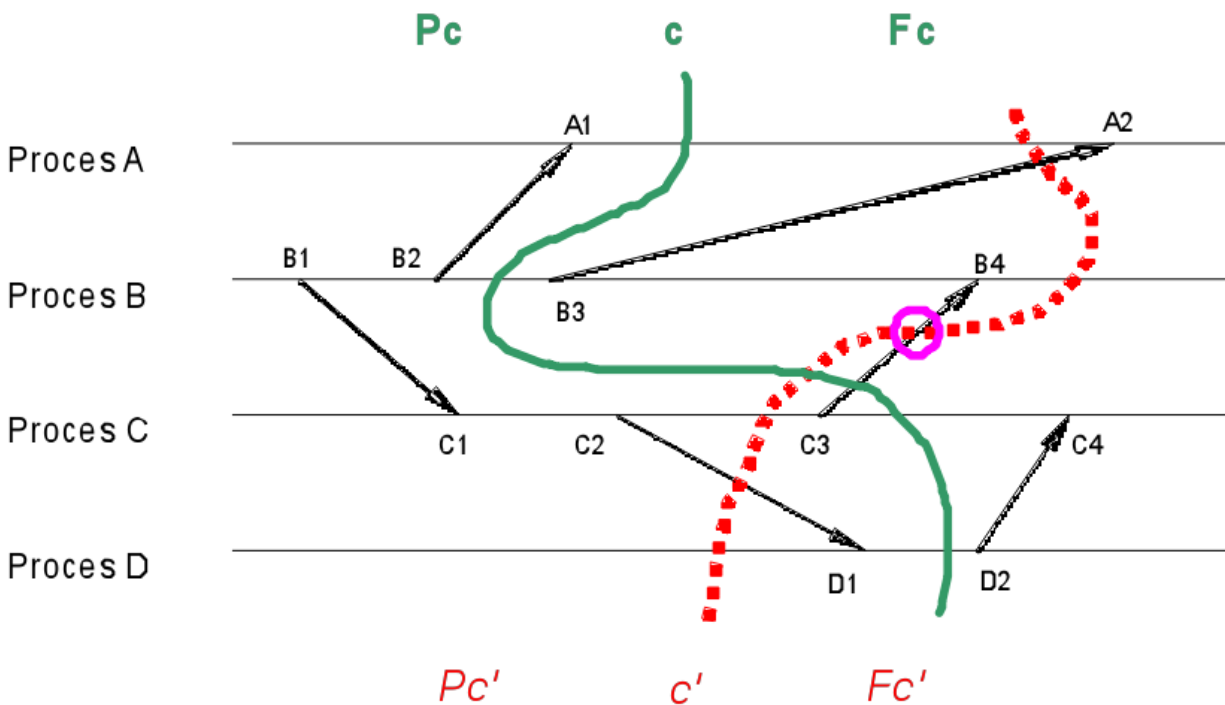
- značka je speciální zpráva odlišitelná od běžných zpráv
- iniciační uzel vyšle žádost o ukončení (značku, že je kanál prázdný) všem výstupním hranám
- uzel přijme první značku:
 - připraven - propagace značky výstupním hranám
 - nepřipraven - vysílá negativní signál (zamítnutí)

- příjem další značky: signál / zamítnutí
- příjem zamítnutí: zamítnutí první žádosti
- příjem všech potvrzení: signál první žádosti
- speciální případ obecnějšího algoritmu - detekce globálního stavu

Konzistentní stav

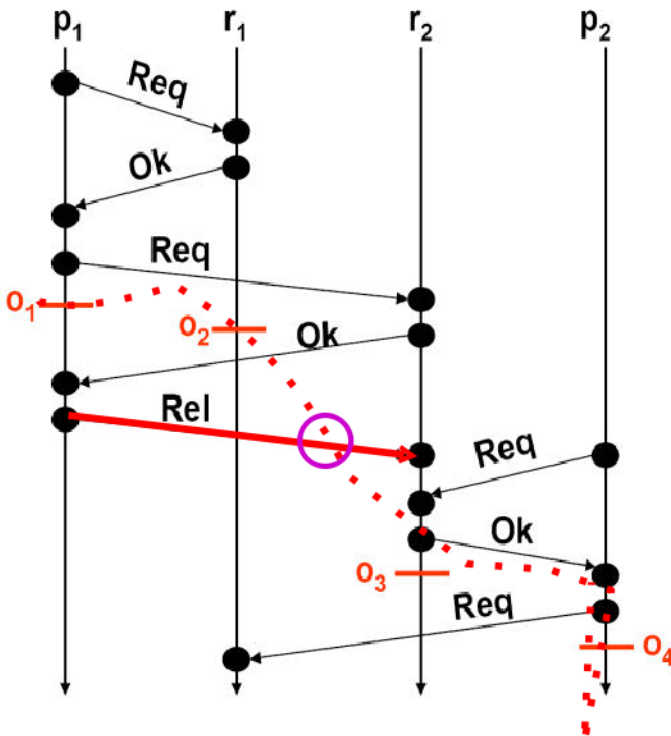
- máme množinu událostí v systému E
- řez c je disjunkttní rozdělení E na P_c a F_c (past/future): $P_c \cup F_c = E$ & $P_c \cap F_c = \emptyset$
- konzistentní řez c: $(a \rightarrow b \text{ \& } a \text{ in } F_c) \Rightarrow b \text{ in } F_c$
 - minulost nelze ovlivnit událostmi z budoucnosti
- stav dist. procesu je množina událostí, které se v procesu udály
- konzistentní stav $S = P_c$, kde c je konzistentní řez
- máme-li S konzistentní stav a hranu $e \Rightarrow S' = S \cup e$ je konzistentní stav (je dosazitelný z S pomocí e)
- posloupnost událostí $s = (e_1 \dots e_n)$ se nazývá rozvrh, pokud je S_n dosazitelné z S přes události (hrany) $e_1 \dots e_n$

c □ konzistentní řez



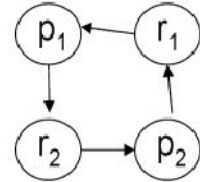
$c' \square$ nekonzistentní řez

Detekce falešného deadlocku



Inference from

- o_1 : p_1 waits for r_2
- o_2 : r_1 waits for p_1
- o_3 : r_2 waits for p_2
- o_4 : p_2 waits for r_1



From $O=\{o_1, o_2, o_3, o_4\}$
the deadlock detector
concludes there is a
deadlock!

Detekce globálního stavu

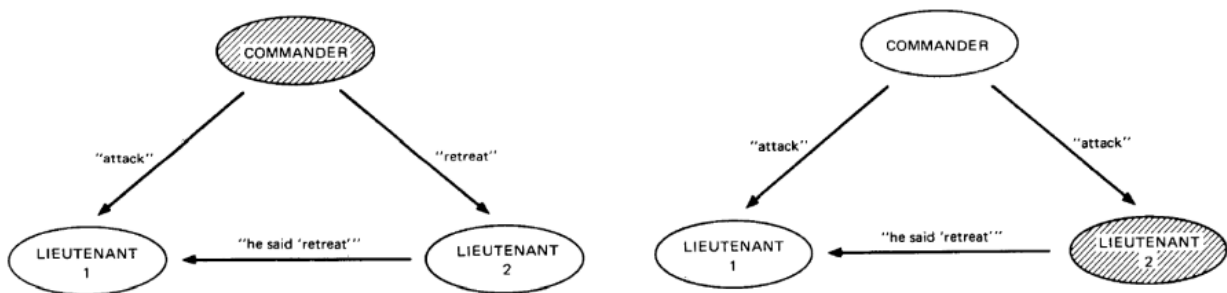
- udržujeme globalni stav DS
 - zahrnuje lokální stav všech proc. + zpravy, které jsou aktuálně v oběhu (nedoručene)
- užitečné napr. pro detekci ukončení běhu/deadlocku, garbage collection, obecně detekce glob. vlastností
- jak? - udržovat distribuovaný snapshot
 - stav, ve kterém systém mohl být
 - odrazí konzistentní glob. stav
 - pokud jsme zaregistrovali, že P obdržel zprávu od Q, měli bychom také zaznamenat, že Q zprávu poslal
 - je ale povoleno, že Q zprávu poslal a obdržení od P jsme ještě nezaregistrovali
 - konzistentní/nekonzistentní rezy
 - poslední události, které byly u procesů zaznamenány
 - předpoklad: DS je kolekce procesu navzájem propojených
 - takže napr. proces musí nejprve iniciovat TCP spojení, pokud chce komunikovat
 - jakýkoliv proces může iniciovat algoritmus - necht je to P
 - P začne zaznamenáním svého lok. stavu

- následně posílá zprávu každým svým výstupním kanálem, která zná, že příjemce se má podílet na zaznamenávání glob. stavu
- když proces Q obdrží takovou zprávu, tak pokračuje na základě toho, zda má zaznamenány svůj lok. stav či ne
 - nemá - pošle zprávu do výstupních kanálů a zaznamená svůj lok. stav
 - má - mezi první zprávou co dostal a tou co teď došla si zaznamenával přicházející zprávy a teď když mu došla druhá zpráva, tak si svoji frontu společně se svým lokálním stavem zaznamenal, jako akt. stav kanálu
- proces skončil svou úlohu, pokud obdržel zprávu všemi přicházejícími kanály a každou zpracoval
- v tuto chvíli může celý snapshot sjednotit a poslat iniciátorovi jako odpověď

Distribuovaný konsensus

- organizace replikovaných procesů do skupin pomáhá zlepšovat odolnost proti chybám
- klasifikace problému
 - byzantský konsensus ($1 \rightarrow 1$)
 - iniciátor vysílá jedinečnou hodnotu, všechny nechybové procesy se na ni musí shodnout
 - konsensus ($n \rightarrow 1$)
 - iniciujících uzlů je několik a mají svoje iniciační hodnoty, všechny nechybové procesy se musí shodnout na jedné společné hodnotě
 - interaktivní konzistence ($n \rightarrow n$)
 - iniciujících uzlů je několik a mají svoje iniciační hodnoty, všechny nechybové procesy se musí shodnout na množině/vektoru společných hodnot
- problém dvou armád
 - početnější armáda je rozdělena
 - úspěch pouze při synchronizovaném útoku, jinak drtivá porážka
 - obě části musí mít jistotu, že druhá část začne útok také
 - komunikace pouze nespolehlivým kurýrem - možnost zajetí
 - striktní řešení neexistuje
 - $A1 \rightarrow A2$: attack -- A1 neví, jestli A2 dostal zprávu
 - $A2 \rightarrow A1$: A2 neví, jestli A1 dostal potvrzení
 - $A1 \rightarrow A2$: A1 neví, jestli A2 dostal potvrzení
 - ...
 - praktická řešení
 - agresivní
 - první generál vyšle větší množství zpráv oznamujících čas útoku, zaútočí tak jako tak
 - předpokládá se, že při vysokém počtu poslaných zpráv alespoň jedna projde
 - druhý generál zde ani nemusí odpovídat

- pravdepodobnostni
 - prvni generál vyšle větší množství (N0) zpráv
 - kromě času útoku přidá i to, kolik jich poslal
 - druhý generál odpoví na každou, kterou obdržel (N1)
 - počet odpovědí vrácených prvnímu generálovi je N2
 - oba generálové znají přibližnou míru úspěšnosti, že posel zprávu pronese
- problem byzantskych generalu
 - předpoklady:
 - uzly mohou libovolně havarovat
 - havarovaný uzel se může chovat zákeřně!
 - spolehlivá komunikace (časově neohraničená)
 - úkol:
 - někteří generálové jsou zrádci
 - všichni loajální generálové se musejí rozhodnout shodně
 - každý generál se rozhoduje na základě informací obdržených od ostatních generálů
 - budeme uvažovat 1 generál, ostatní důstojníci (búno)
 - jak generál tak důstojníci mohou být zrádci
 - generál vydá rozkaz, důstojníci ho předají ostatním
 - rozkaz bude vydán na základě většiny
 - cíl:
 - C1: všichni loajální důstojníci vydají stejný rozkaz
 - C2: je-li generál loajální, pak každý loajální důstojník vydá rozkaz generála
 - 3 uzly, 1 zradce
 - řešení neexistuje
 - obecně je řešení jen pro: $n > 3m$ uzlu, kde m je # zradcu



- 4 uzly, 1 zradce
 - zradce general
 - vyda 3 stejne rozkazy - nelze, to by nebyl zradce
 - vyda 2 stejne rozkazy
 - dustojnici si vzajemne preposlou vetsinovy rozkaz
 - dosahneme C2

- vyda 3 ruzne rozkazy
 - dustojnici se shodnou na tom, ze je zradce
 - dosahneme C1
- zradce dustojnik
 - nejhorsí pripad: preda falesny rozkaz vsem ostatnim
 - loajalni dustojnici dostanou vetsinu spravnych rozkazu (C2)

Distribuovaná sdílená paměť

- paralelní výpočty
 - multiprocesory
 - malý počet procesorů
 - příliš komplikované a drahé
 - multicomputery
 - snadno dostupné
 - programování a synchronizace není prakticky (inženýrsky) dobře zvládnutá
 - RPC – problémy
 - pokus o řešení – distribuovaná sdílená paměť (DSM)
 - 1986 Li & Hudak – DSM - množina uzlů sdílející adresový prostor

Konzistenční modely DSM bez SP

- specifikace co implementace musí splňovat vzhledem k operacím čtení a zápisu
- klasifikace od nejvíc pedantského po nejvolnější
- striktní konzistence
 - jakékoliv čtení z paměti z adresy x vrátí hodnotu uloženou při posledním zápisu na adresu x
 - zajišťuje absolutní časové uspořádání
 - nejsilnější, na jednoprocessorových systémech je tradičně zajištěna
 - všechny zápisy okamžitě všude viditelné
 - podmínka: musí existovat přesný globální čas
 - ideální pro programování, v DS nedosažitelné
- sekvenční konzistence
 - výsledek výpočtu je stejný, jako kdyby všechny operace všech procesorů byly vykonávány v nějakém sekvenčním uspořádání a operace každého procesoru jsou vykonávány v pořadí specifikovaném programem
 - snadno implementovatelná
 - příjemná pro programování
 - povoleno libovolné prokládání instrukcí na různých procesorech
 - všechny procesy vidí stejné pořadí změn paměti
 - změny nejsou propagovány okamžitě, není zaručena velikost zpoždění
 - výkonnost není příliš velká
 - platí že čas čtení a zápisu dohromady musí trvat alespoň tak jako přenos jednoho paketu → pomalé

- kauzální konzistence
 - kauzálně vázané zápisy musí být viděny všemi procesy ve stejném pořadí, konkurenční zápisy mohou být viděny v různém pořadí
 - rozlišuje události, které jsou potenciálně závislé a které ne
 - kauzálně závislé zápisy
 - pokud $W(x)a$ a (druhý proces) $R(x) W(y)b$, pak $W(y)b$ je kauzálně závislý na $W(x)a$
 - analogie k zasílání zpráv: zápis ~ odeslání, čtení ~ přijetí
 - složitější implementace
 - vyžaduje udržování grafu závislostí zápisů a čtení
- PRAM konzistence (FIFO konzistence)
 - pipelined RAM
 - zápisy jednoho procesu jsou pipelinovány, nemusí se čekat vždy na dokončení každého než se začne nový
 - zápisy prováděné jedním procesem jsou viděny ostatními procesy v tom pořadí, ve kterém byly prováděny, zápisy různých procesů mohou být viděny různými procesy různě
 - srovnání se sekvenční konzistencí:
 - v PRAM neexistuje jednotný pohled na rozvrh
 - snadná na implementaci
 - nezáleží na pořadí, v němž různé procesy vidí přístupy k paměti
 - je nutné dodržet pořadí zápisů z jednoho zdroje
- slow memory
 - zápisy jedním procesem do jednoho místa musí být viděny ve stejném pořadí
 - nejslabší model - lokální zápis, pomalá nesynchronizovaná propagace
 - neposkytuje žádnou synchronizaci

Konzistenční modely DSM s SP

- doposud uvedené konzistenční modely velmi restriktivní
 - vyžadují propagaci všech zápisů všem procesům
 - málo efektivní
 - ne všechny aplikace vyžadují sledování všech zápisů, natož pak jejich pořadí
 - typická situace:
 - proces v kritické sekci ve smyčce čte a zapisuje data
 - ostatní procesy nemusí jednotlivé zápisy vidět, není nutné, aby byly propagovány
- paměť ale neví, že proces je v kritické sekci, musí propagovat všechny zápisy
 - řešení: nechat proces ukončit kritickou sekci a poté rozeslat změny ostatním
 - používá se synchronizační proměnná
- slabá konzistence
 - přístup k synchronizačním proměnným je sekvenčně konzistentní
 - všechny procesy vidí všechny přístupy k SP ve stejném pořadí

- přístup k SP není povolen, dokud neskončí všechny předchozí zápisy
 - před přístupem k SP budou dokončeny všechny předchozí zápisy
- přístup k datům není povolen, dokud nebyly dokončeny všechny předchozí přístupy k SP
 - při přístupu k obyčejným proměnným jsou dokončeny všechny předchozí přístupy k SP
- provedením synchronizace před čtením se zajistí aktuální verze dat
- odpadá nutnost propagace všech zápisů, ale paměť při přístupu k SP nerozezná vstup/výstup z CS → musí se vždy vykonat akce pro oba případy
- výstupní konzistence
 - před přístupem ke sdílené proměnné musí být úspěšně ukončeny předchozí Acq() procesu.
 - před provedením Rel() musí být ukončeny všechny předchozí zápisy i čtení prováděné procesem.
 - Acq() a Rel() musí být PRAM konzistentní.
 - po Acq() jsou všechny lokální kopie aktuální
 - po Rel() jsou propagovány změny ostatním procesům
 - při správném párování Acq() a Rel() je výsledek jakéhokoliv výpočtu ekvivalentní sekvenčně konzistentní paměti
 - implementace:
 - eager release consistency (horlivá)
 - po Rel() se propaguje všem procesům
 - optimalizace přístupové doby
 - lazy release consistency (líná)
 - po Rel() se nic nepropaguje, až po Acq() jiného procesu
 - optimalizace síťového přenosu
- vstupní konzistence
 - Acq k SP není povolen, dokud nebyly provedeny všechny aktualizace chráněných sdílených dat procesu
 - exkluzivní přístup procesu k SP (= zápis) je povolen pouze v případě, že žádný jiný proces nepřistupuje k SP, a to ani neexkluzivně (= čtení)
 - po exkluzivním přístupu k SP si příští neexkluzivní přístup libovolného procesu k SP musí vyžádat aktuální kopii dat od vlastníka SP
 - sdílená data jsou vázána na SP, při přístupu se synchronizují pouze tato data
 - přístup k datům a SP může být exkluzivní (RW) nebo neexkluzivní (RO)
 - každá SP má vlastníka - proces, který k ní naposledy přistupoval
 - vlastník může opakovaně vstupovat a vystupovat z k.sekce bez nutnosti komunikace
 - proces, který není vlastníkem, musí požádat o vlastnictví

Distribuované stránkování

- obdoba virtuální paměti - přístup k nenamapované stránce → přerušení, obsluha, načtení
- adresový prostor rozdělen na části (chunks - pages), které jsou rozprostřeny mezi všechny procesory v systému
- když procesor referencuje adresu, která není lokální, nastane výjimka, DSM software načte příslušnou část a pokračuje se dál
- problémy
 - replikace → konzistence
 - jak udržovat data konzistentní
 - namapování read-only, při zápisu synchronizační akce
 - replikovat všechny části
 - pokud dochází pouze ke čtení je to OK a mohou se replikovat i RW kusy
 - pokud je ale replikována kopie náhle změněna, dochází k nekonzistenci
 - řeší se aplikací některého konzistenčního protokolu
 - invalidace vs. aktualizace
 - nalezení stránky
 - jak nalézt distribuovaná data, vlastníka stránky
 - centralizovaný manager
 - replikovaný manager (indexace spodními n bity nebo hash)
 - broadcast
 - správa kopií
 - co dělat s kopiemi při čtení/zápisu
 - jak je najít, když musí být invalidovány
 - broadcast page number a říct všem procesorům, které ji mají, aby ji invalidovali
 - musí být spolehlivé
 - copyset
 - vlastník stránky udržuje množinu lokací kopií
 - pokud musí být stránka invalidována, posílá zprávu procesorům, které ji mají a čeká na ACK
 - uvolňování stránek
 - kterou stránku uvolnit
 - nevlastněná read-only kopie
 - vlastněná replikovaná kopie - přenos vlastnictví
 - lokální heuristika - LRU / ...
 - falešné sdílení
 - nezávislá data na jedné stránce

Sekvenčně-konzistentní distribuované stránkování

- využívá protokolu sekvenční konzistence pro zachování konzistentních dat
- stránky mají vlastníka, který na ně může zapisovat, ostatní mají kopie ke čtení

Kauzálně-konzistentní distribuované stránkování

- základní idea – graf závislostí
- vektorové hodiny:
 - VTS ~ jednotka granularity (stránky)
 - na kterých stránkách závisí obsah
 - VTP ~ procesy
 - z kterých stránek znám data
- $VT[i] \approx$ stránka i
- výpadek stránky – přenos dat, $VTP = \max(VTS, VTP)$
- zápis do stránky – $VTS = \text{inc}(VTP)$
 - jak OS pozná zápis? - změnou mapování
- aktualizace VTP – zneplatnění stránek i : $VTS_i[i] < VTP[i]$
- porovnání kauzality - obdobné jako dříve - srovnání vektorů, pokud jsou porovnatelné
- problémy:
 - velká prostorová režie (1 MB = 256 stránek = 512 B / stránku)
 - propagace konkurentních zápisů
 - zámky, bariéry, timeouty

Distribuované sdílené proměnné

- implementace na úrovni knihoven
 - potenciálně replikovaná distribuovaná databáze
 - typicky konzistenční model se synchronizačními proměnnými
- výhody
 - potenciálně lepší výkonnost
 - eliminace false sharing
- nevýhody
 - nepodporováno přímo operačním systémem
 - nutnost implementace pro různé jazyky
 - nutnost rekompile
- příklad implementace: Munin 1990 Benett & spol.
 - použití ordinary / shared / synchronization promenných
 - shared promenne: read-only, migratory, write-shared, conventional
 - read-only
 - při výpadku je v adresáři nalezen vlastník, žádost o read-only kopii dat
 - migratory
 - acquire/release protokol implementující eager release consistency
 - při opuštění kritické sekce se propagují změny
 - data chráněná SP migrují na uzel v kritické sekci

- write-shared
 - stránky iniciálně r/o, při zápisu do lokální kopie s původním obsahem jsou r/w, označí se dirty
 - po release se porovná stránka s původní, změny se propagují
 - propagace na ne-dirty stránku akceptace, jinak word-po-wordu porovnání
 - sjednocení dat / konflikt - runtime-error
 - conventional
 - jako distribuované stránkování - single writer/many readers
 - sekvenční konzistence
- distribuované objekty
 - díky zapouzdření flexibilnější - komunikace a synchronizace v metodách
 - distribuovaná data potomkem základní distribuované třídy
 - Class Definition Language - automatické generování hlaviček a kódu
 - základní "distribuovaná" třída, dědění vlastností, CORBA, Java RMI, ...

Identifikace a správa prostředků

- kteřý objekt má být použit (identifikace)
- kde je tento objekt umístěn (adresa)
- jak se k němu můžeme dostat (cesta)
- trvanlivost identifikaci
 - dynamická jména - docasna, transient
 - typicky vazana na zivotnost procesu
 - deskriptory otevrenych objektu, porty, adresy DSM
 - statická jména - trvala, persistent
 - nezavisla na procesech
 - jména souboru, kapability
- prostor jmen
 - separátní prostory - file system, registry, URL, id procesu, ...
 - jednotný - distribuovaný name server
 - lokální a replikační transparentnost
- nestrukturovaná jména (UUID) versus strukturovaná (ms.mff.cuni.cz)
- systémová jména
 - interní identifikace objektů, optimalizace pro strojové zpracování
 - typicky binární řetězce pevné délky i struktury
 - chráněné objekty jádra - porty
 - access control matrix - user x resource → effective right
 - tak jednoduché to není, taková matice by mohla být obrovská
 - capabilities - uživatel drží oprávnění k prostředkům
 - access control list - prostředek (objekt) má seznam uprávněných uživatelů

Kapability

- datová struktura umožňující jednoznačnou identifikaci objektu
- obsahuje navíc přístupová práva pro držitele capability
- ochrana objektů, šifrování, redundance
- k jednomu objektu typicky několik různých kapabilit
 - vlastník, písař, čtenář, ... další druhy služeb
- uživatelským procesům znemožněno vlastní generování kapabilit i změny práv
- výhody
 - snadný test oprávnění přístupu
 - flexibilita - každý správce prostředků může nadefinovat vlastní druhy práv
- nevýhody
 - problematická kontrola propagace
 - možnost neautorizovaného přístupu
 - copy bit, čítač - chráněný způsob přenosu
 - review - seznam oprávněných uživatelů
 - revocation - odejmutí práva
 - zrušení objektu a vytvoření nového, notifikace ostatních uživatelů
 - garbage collection
- varianty:
 - s podpisem
 - hlavní část - identifikace objektu a přístupových práv
 - platná capability je rozšířena o podpis, který je vypočítán z obsahu hlavní části
 - celá capability zašifrována tajným klíčem
 - ochrana proti odvození podpisové funkce uživatelskými procesy
 - řidkost: ze všech binárních čísel velikosti capability jsou platné jen ty, jejichž podpis odpovídá zbytku capability
 - s redundantní kontrolou
 - port serveru a identifikátor objektu jsou volně přístupné, ale část s přístupovými právy zašifrována
 - přístupová práva a náhodně vygenerované binární číslo pevné délky - zakódováno
 - číslo má u sebe i server, který capability vytvořil → verifikace
 - ochrana serveru
 - port, dostatečně velké náhodně generované číslo
 - ochrana přístupových práv
 - zašifrování pole s přístupovými právy spolu s redundantní kontrolou
 - uživatelskému procesu znemožněna změna přístupových práv
 - proces nemůže svá práva změnit, nezná dešifrovací funkci, ani ji nemůže odvodit
 - náhodně generovaná redundance
 - služba serveru - vygenerování capability s menšími právy

- klient kapabilitu pošle serveru, pomocí čísla check, které si uložil při generování kapability a masky práv, kterou mu klient předloží vygeneruje nový check a vytvoří restringovanou kapabilitu

Distribuovaná správa jmen

- name servers, directory servers, lookup servers
- adresare jsou množina položek <jmeno,hodnota>, kde hodnota může být:
 - primitivní (čísl, řetězce, binární data)
 - perzistentní reference (trvalé odkazy na objekty, kapability)
 - tranzientní reference (na živé objekty, porty, kanály)
 - odkazy na jiné adresáře
 - operace jako SET, LOOKUP(jmeno) - změna kontextu, LOOKUP(složené jméno - hierarchie v URL napr.)
 - vyhledávání typicky přes server, ten publikuje jména a vrací reference na objekty (to může být třeba nějaká kapabilita nebo něco jako handle do FS)
- klientsky proces má několik základních "adresaru"
 - file system, objekty, služby a servery, registry
- server spravuje objekty - jejich identifikace a operace nad nimi
 - každý objekt má lokální id spravované serverem
 - na něm definované operace / služby
 - server publikuje řídicí port a porty pro otevřené objekty
 - proces:
 - server publikuje kapability na name server
 - zaregistruje své služby u registration serveru a vytvoří kanál
 - klient požádá NS o kapabilitu
 - klient požádá RS o open(cap)
 - požadavek se preposle serveru
 - mezi klientem a serverem se vytvoří komunikační kanál
 - poskytované služby
 - porty ~ objekty, služby ~ metody
 - vytvoření objektu, transient reference - service_port->create_obj(...)
 - operace nad objektem - obj_port->op()
 - freeze - vytvoření kapability (perzistentního odkazu)
 - cap = obj_port->freeze()
 - zrušení transient reference (ne objektu) - obj_port->drop()
 - uložení kapability - ns->add("path/to/cap.obj", cap)
 - vytvoření restring. kapability z dosavadní kap.
 - cap = service_port->cap_restrict(cap,0101)

Distribuovaná správa prostředků

- centralizované řešení - resource manager
 - pokusy o distribuovanou správu
 - prakticky nepoužitelné - distribuované vyloučení procesů
- migrace - identifikace, komunikace
- replikované servery – aktualizací protokol
- problém deadlocku
 - časté řešení - pštroší algoritmus
 - detekce - větší problém než u centralizovaných systémů
 - WFG – Wait-For-Graph
 - orientovaný graf (procesů, transakcí)
 - $P1 \rightarrow P2$ - proces P1 je blokován procesem P2
 - orientovaná kružnice – uvážnutí

Algoritmy detekce deadlocků

- korektnost
 - každý existující deadlock je v konečném čase detekován
 - detekovaný deadlock musí existovat
 - nutné dávat pozor na vnořené deadlocky - phantom deadlock
- modely - single, and, or m-out-of-n, and-or
 - and - proces muze naraz požadat o více prostředku, je zasekly, dokud nema vsechny
 - or - jak vyse, ale staci mu jeden, aby se rozjel
 - single unit - muze zadat pouze o jeden prostredek, ve WFG vystupni hrany stupne 1, cyklus znaci deadlock
 - and-or - požadavky o prostředky reprezentovany predikatem, kde jsou atomy prostředky, když je formule splnitelna, proces se rozjede
 - m-out-of-n - muze zadat o n, staci když bude mít m

Centralizovaný algoritmus na DD

- jeden uzel vybrán jako ten, co detekuje, vsechny zavislosti, které v systemu vyvstanou, se posilaji jemu
- na centralnim uzlu se analyzuje a spravuje WFG
- prenasime informace bud po kazde zmene, po danem intervalu nebo na expl. vyzadani
- pro prevenci falesneho uvaznuti kvuli zpozdeni zprav se používaji logicke hodiny + kauzalni dorucovani
- rozsireni na hierarchicky algoritmus: koordinatori a podrizeni - koord. resi deadlocky podrizenych, uzel resi deadlocky lokalne

Path-pushing algoritmus na DD

- základní idea je, že každý uzel periodicky shromažďuje informace o lokálních závislostech a staví si svůj lokální WFG
- sousedním uzlům posílají externí žádosti, z kusu WFG, které jim nepatří
- tyto výsledky pak zahrnou do své části
- toto se opakuje dokud nějaký uzel nemá dostatečné info o tom prohlásit, zda doslo k deadlocku či nikoliv

Edge-chasing algoritmus na DD

- zaslání spec. zpráv (probes) podél WFG
 - proces odesle zprávu všem, na které čeká
 - pokud se zpráva vrátí, je to deadlock (existuje orient. kružnice)
 - mezitím se to však mohlo odblokovat
 - řešením je aging - postupné zvyšování priority těch, kteří již dlouho čekají

Diffusing computation na DD

- těm, na které čekám, se posílají pingy, oni je vrací, pokud jsou taky zablokováni, pokud dostanu všechny své pingy zpět, mám deadlock
- vhodné u složitějších modelů (OR, m-out-of-n)

Detekce globálního stavu pro DD

- jestliže WFG' následuje WFG a proces je ve WFG v deadlocku, bude i ve WFG' (existuje-li deadlock, je i v konzistentním rezu)
- při příjmu znacky - tj. v okamžiku rezu - uzel zaznamená lokální WFG
- dochází k lokální kontrakci WFG, externí závislosti se posílají iniciatorovi

Distribuované procesy

- správa procesů v distribuovaných systémech
- sdílet výpočetní sílu systému
- rozdělovat zátěž na jednotlivé procesory
- provádět operace na vzdálených procesech
- synchronizovat procesy a vést evidenci stavu

Vzdálené spouštění procesů

- nalezení volného uzlu
 - alokační algoritmy
- znalost vlastní zátěže
 - procento využití procesoru
 - periodické měření - vyhlazení (nějaký klouzavý průměr)
- spuštění vzdáleného procesu
 - transparentnost
 - přenesení kódu a dat – nezajímavé
 - heterogenní prostředí
 - vytvoření prostředí odpovídající domovskému uzlu
 - kontexty (fs, naming, environment)
 - systémová volání – vzdálená / přesměrování
- problém - hostitelský uzel přestane být volný (uživatel se vrátil z oběda, potřeba restartu, ...)
 - doběhnutí / zabití / čas na uložení / uzavření / migrace

Klasifikace alokačních algoritmů

- zda jsou předem známy údaje o procesech, podle kterých se algoritmus rozhoduje
 - deterministické / heuristické
- do jaké míry se provádí optimalizace
 - optimální / suboptimální
- co se snaží optimalizovat
 - využití CPU / čas odpovědi / přidělování prostředků / komunikační složitost, ...
- zda umožňují přemístění již rozběhnutého procesu
 - migrační / nemigrační
- rozdíly mezi HW uzlů, speciální požadavky procesů na hw vlastnosti
 - homogenní / heterogenní
- podle charakteru algoritmu
 - centralizované / distribuované
- podle jakých informací se provádí rozhodnutí o vzdáleném spuštění procesu
 - lokální / globální
- kdo iniciuje vzdálené spuštění procesu
 - odesílatel / příjemce / symetrický
- *podtržené - co se používá v praxi*

Alokace procesorů

- hierarchický algoritmus
 - manažeři skupin, při neúspěchu žádost vyšším místům
- distribuovaný heuristický algoritmus
 - k náhodných výběrů cíle
 - server / receiver initiated
- deterministický grafový
 - minimalizace komunikace
 - nutnost znalosti komunikační složitosti
 - optimální deterministický algoritmus – tok v sítích
- bidding (obchodní, nabídkový) algoritmus
 - procesy kupují výpočetní sílu, procesory ji nabízejí
- up-down algoritmus
 - optimalizace na stejnoměrné sdílení výkonu
 - koordinátor má tabulku se záznamem pro každý uzel obsahující “trestné body”
 - při každé významné akci (vytvoření procesu, ukončení procesu, tik hodin):
 - zpráva koordinátoru, který provede změny:
 - každý proces běžící na jiném uzlu - plus trestné body
 - každý neuspokojený požadavek - mínus trestné body
 - jestliže nic z tohoto - směrem k nule
 - v okamžiku uvolnění uzlu se vezme ten proces z fronty (neuspokojených požadavků), jehož vysílající uzel má nejméně trestných bodů
 - vlastnosti:
 - homogenní, nemigrační, spravedlivý, heuristický, centralizovaný, suboptimální, globální, symetrický

Migrace procesů

- jde nám o korektní a transparentní přenesení procesu během výpočtu
- motivace
 - vyvažování zátěže; optimalizace - I/O, komunikační; přemístění; shutdown
- korektnost
 - ostatní procesy nejsou migrací podstatně ovlivněny
 - po ukončení stav odpovídá stavu bez migrace
- transparentnost
 - proces o migraci neví a nemusí spolupracovat
 - zůstanou zachovány vazby na komunikující procesy
 - není narušena komunikace
 - ... v konečném důsledku

- problémy k řešení
 - přenesení rozpracovaného stavu ...
 - přenesení adresového prostoru ...
 - reinicializace vazeb na ostatní procesy
 - komunikace s ostatními procesy
 - reziduální dependence - žádná, domácí, průběžná, dočasná
 - vícenásobná migrace - viskozita

Implementace migrace - přenos procesu

- vyjmutí ze stavu, zmražení, speciální stav
- oznámení příjemci o migraci, alokace procesu
- přenos stavu - registry, zásobník, stav procesu
- přenos kódu / adresového prostoru
- přesměrování / doručení zpráv
- dealokace procesu, vyčištění
- vazby na nové jádro, nastartování procesu
 - přesunutí části stavu spolu s procesem (obsah VM)
 - forwardování (některých) požadavků (komunikace s konzolí)
 - použití odpovídajícího prostředku na cílové stanici (fyzická paměť)
- dokončení přenosu vazeb, dočištění
 - dočasná reziduální dependence, přesměrování zpráv
- co s virtuální pamětí
 - přenesení celé VM při migraci
 - výhody - eliminace reziduálních dependencí
 - nevýhody - prodloužení doby zmražení procesu, mnohdy zbytečné přesouvat celý obsah virtuálního adresového prostoru
 - pre-copying
 - výhoda - proces je zmražen pouze po dobu přesunu malého množství dat
 - nevýhoda - některé části se kopírují vícekrát - prodloužení celkové doby migrace
 - copy-on-reference
 - nejprve se přenesou stav procesu potřebný pro běh (registry, kanály, ...)
 - přesun adresového prostoru odložen
 - stránky na cílové stanici označeny jako neprezentní (jako by byly odloženy na disk)
 - při přístupu na stránku se vyvolá výjimka, obsluha obsah přenesou
 - nutnost evidence různých zdrojů dat pro každou stránku
 - swap, vícenásobná migrace, DSM, ...
- co se zprávami
 - přesměrování dočasné/trvalé
 - migrace kanálů, spojení front
 - vytvoření nových kanálů / přepojení / zrušení starých

- implementace v praxi
 - DEMOS/MP
 - IPC pomocí linků spojených s procesy, migrace v jádře
 - max transparence, jednotné komunikační rozhraní bez ohledu na umístění procesu
 - Charlotte
 - IPC pomocí linků nezávislých na umístění procesů, možnost přesouvat linky
 - migrační politika v uživ. procesu, migrace v jádru
 - snaha o max. fault-tolerance
 - sběr statistiky o procesu, počítání
 - V
 - transparence, minimalizace doby zmrazení (precopying)
 - MOSIX
 - jeden z mála reálně používaných
 - rozsahle vyvažování zátěže
 - rozšíření struktury procesu o různé statistiky
 - čas od poslední migrace, čas na procesoru, ...
 - UNIX pro distrib. prostředí (používaný 1988 až do 2004!)
 - Sprite
 - rychlost a transparence
 - předpoklad - hodně nevyužitých uzlů během nějaké pauzy, nutno pak zase odmigrovat po návratu vlastníka
 - VM: kombinace přesunu všeho + copy-on-reference
 - T4
 - byl by efektivní pro použití jako konvenční OS
 - silné prostředí pro použití jako DS
 - meziproc. komunikace, přenosové protokoly, vzdalena komunikace, name services, podpora pro high-level distrib. služby
 - splnil účel pro výukové/výzkumné účely, nesplnil účel pro každodenní využití (chyběly aplikace)

Load balancing

- nutno porovnávat zatížení procesorů (nějak konzistentně), pak vybrat co bude migrovat a kam
- rozhodnutí o okamžiku migrace
 - jak porovnávat zatížení uzlů
 - udržování konzistence údajů
 - volba migrujícího procesu
 - volba příjemce
 - přenos a běh procesu
- algoritmy

- párový algoritmus
 - vytvářejí se páry, které se vzájemně vyvažují
 - A pošle B žádost o vytvoření páru se seznamem procesů
 - B odmítne / vytvoří pár / migruje
 - zatíženější uzel vybere proces podle míry vylepšení
 - $k_i = A_i / (B_i + AB_i)$
 - významné zlepšení stavu – migrace a další proces, jinak konec
- vektorový algoritmus
 - používá se pevný vektor zátěže, první složka $L(0)$ vlastní zátěž, na ostatních složkách ostatní uzly
 - periodicky každý uzel zjistí vlastní zátěž a polovinu vektoru pošle náhodnému uzlu
 - při příjmu proloží hodnoty s vlastním vektorem + přičte se komunikační zátěži
 - použito v MOSIXu
- bidding algoritmus
 - “inteligentní” alg.
 - používá se vyhodnocovací bunka - excitatory, inhibitory, výstup
 - procesy se pravidelně vyhodnocují, při výstupu pod určitý prah se migruje:
 - broadcastuje žádost o nabídku (RFB) do nějaké max vzdálenosti d
 - odpovědi s nabídkou se zkontrolují o zátěži
 - bereme nejlepší nabídku, pokud žádná nedorazí, zvětšíme d
 - obtížná kvantifikace vlastností procesu
- UI/math-based
 - zpětné učení, bayesian decisions -- neprosadily se - moc komplikované
- centralizovaný / hierarchický algoritmus
 - centralizovaný manažer, zná zátěž všech uzlů, velí
 - hierarchicky organizované skupiny, nadřazení manažerů
- lokální algoritmus
 - zná jen lokální zátěž
 - nastaven prah - pokud ji proces překročí, pošle n uzlům žádost o přijetí, vybere si nejlepší
 - velmi jednoduché, velmi suboptimální, velmi naprd

Replikace

- replikace souborů - udržování více kopií na více fileserverech
 - spolehlivost (reliability) - při havárii serveru nejsou ztracena data
 - dostupnost (availability) - k souborům lze přistupovat i při výpadku serveru
 - výkon (performance) - lze přistupovat k nejbližším datům, rozdělení výkonu
- typy:
 - explicitní replikace
 - 'uživatel' se sám stará o udržování konzistence
 - odložená replikace
 - zápis do primární repliky, aktualizace sekundárních
 - skupinová komunikace
 - zápisy simultánně zasílány všem dostupným replikám
- aktualizací protokoly
 - problém aktualizací kopií
 - uvažujeme, že soubor je replikován na N serverech
 - možnosti:
 - primární kopie vítězí
 - většinové/vážené hlasování
 - pro aktualizaci souboru musí klient požadovat většinu serveru (přes polovinu) a získat od nich schválení
 - pokud chce klient číst soubor musí shromáždit read quorum N_R , podobně při zápisu N_W , kde platí:
 - $N_R + N_W > N$ &
 - zabranuje read-write konfliktum (unrepeatable reads)
 - $N_W > N/2$
 - zabranuje write-write konfliktum (overwriting uncommitted data)
 - toto je vážené hlasování, většinové je když $N_R = N_W$
 - např. ROWA schema (Read One, Write All) je korektní
 - hlasování s duchy
 - bezdatový (dummy, ghost) server, obsahuje pouze verze, žádná data
 - neúčastní se hlasování o čtení, pouze o zápisu
 - dynamická kvóra

Klientocentrické konzistenční modely

- u konzistenčních modelů DSM byl pohled na sdílená data různými procesy
- replikovaná databáze
 - zápisy málo časté, masivně paralelní čtení
 - není potřeba vzájemná synchronizace klientů
 - WWW + cache, mobile computing, News, ...
- u klientocentrického modelu je pohled na replikovaná data jedním procesem
- eventuelní konzistence
 - po ukončení všech zápisů budou všechny repliky v konečném čase aktualizovány
 - tolerance poměrně vysoké úrovně nekonzistence
 - problém: může se stát, že při připojení k jedné replice uživatel vidí zprávy, po připojení k jiné replice některé zprávy (které již viděl) 'ještě' nevidí
- monotónní čtení
 - po přečtení hodnoty x všechna další čtení vrátí stejnou nebo novější hodnotu
 - příklad: při připojení k jiné replice uživatel vidí všechny dosud přečtené zprávy
- monotónní zápis
 - zápis proměnné je proveden před jakýmkoliv následným zápisem této proměnné
 - někdy je důležité propagovat všechny write operace ve správném pořadí na všechny kopie
 - příklad: CVS commit na různých replikách
- čtení vlastních zápisů
 - zápis proměnné je proveden před jakýmkoliv následným čtením této proměnné
 - příklad: aktualizování webových stránek, nechceme si prohlížet kopie z cache
- zápisy následují čtení
 - zápis proměnné po předchozím čtení této proměnné je proveden na stejné nebo novější hodnotě
 - příklad: zápis odpovědi do news groups se provede tam, kde je i přečtená hodnota
- naivní implementace
 - každému zápisu je přiřazen globálně jednoznačný identifikátor WID
 - přiřazuje replika kde byl zápis proveden klientem: $WID = repl_id + loc_id$
 - pro každého klienta udržujeme dvě množiny identifikátorů WID:
 - read-set
 - WID relevantní pro read operace vykonané klientem
 - write-set
 - WID relevantní write operacím vykonaným klientem
 - monotónní čtení
 - při čtení replika serveru ověří podle read-set aktuálnost svých zápisů
 - při chybějících zápisech provede synchronizaci nebo forwarduje čtení
 - po čtení si klient aktualizuje read-set podle repliky ze které četl

- monotónní zápis
 - při zápisu replika ověří podle write-set aktuálnost svých zápisů
 - chybějí-li nějaké, zapíše je
 - po zápisu si klient aktualizuje write-set
 - problem monotónního čtení a zápisu - neomezený růst read/write setu
- čtení vlastních zápisů
 - při čtení replika ověří podle write-set aktuálnost svých zápisů
 - jiné možné řešení - forward čtení na aktuální repliku
- zápisy následují čtení
 - aktualizace repliky podle read-set
 - aktualizace read-set i write-set klienta
- efektivější implementace
 - seskupení do relací / sessions
 - typicky vázané na aplikaci nebo modul
 - při ukončení / restartu smazání množin
 - neřeší problém trvale běžících aplikací
 - nebo reprezentace množin pomocí vektorových hodin
 - umožňují efektivní operace nad množinami
 - klient posílá své aktuální VT při read/write požadavku
 - replika serveru při zápisu přiřadí WID a lokální TS(WID)
 - každá replika si udržuje $RCV(i)[j]$ - TS poslední operace zápisu přijatá replikou S_i od S_j
 - při přijetí žádosti o čtení nebo zápis replika vrátí aktuální $RCV(i)$
 - read-set i write-set jsou reprezentovány vektorovými hodinami
 - obecná pravidla:
 - $VT(A)[i] = \text{maximální TS operací z } A \text{ iniciovaných na replice } S_i$
 - sjednocení: $VT(A+B)[i] = \max(VT(A)[i], VT(B)[i])$
 - test podmnožiny: $A \text{ in } B \Leftrightarrow VT(A) \leq VT(B)$
 - po přijetí TS si klient aktualizuje read-set nebo write-set
 - čtení: $VT(RS)[j] = \max(VT(RS)[j], RCV(i)[j]) \text{ for all } j$
 - zápis: $VT(WS)[j] = \max(VT(WS)[j], RCV(i)[j]) \text{ for all } j$
 - read/write-set reprezentuje poslední operace zápisu, které klient viděl/zapisoval

Epidemické protokoly

- implementace eventuální konzistence
- optimalizace komunikace ve VELMI rozsahlych systemech, neresi konflikty
- teorie epidemie - infekcni nakaza
 - rozsireni infekce co nejrychleji, na co nejvetsi pocet uzlu
- servery jsou: infected (rozšiřují "epidemie"), removed (data mají ale nerozšiřují), susceptible (data nemají)
- entropie
 - server P vybere nahodne server Q k vymene dat
 - mozne vymeny:
 - push
 - server P pushne pouze svoje zmeny na Q
 - pri velke infekci mala pst, ze se zmeny rozsiri
 - pull
 - server P pullne zmeny z Q
 - zpusob se zmeny rozsiruji pomalu, nakonec cela mnozina
 - push/pull
 - vyhody obou postupu
- gossiping
 - pokud byl P aktualizovan, kontaktuje nějaký další server Q, aby šířil update; jestliže Q už update má, P se s pravděpodobností $1/k$ nastaví jako removed
 - nezaručuje rozsireni vsem
 - kombinuje se s periodickym pullem
- problem mazani dat
 - naivne: smazani dat → smazani vseh info o datech
 - dusledek: entropie jinym kanalem data zase obnovi
 - reseni: certifikaty smrti
 - zaznam o smazani
 - problem: neomezeny narust certifikatu
 - reseni: v konecnem case certifikat zanikne, podminkou je konecny cas rozsireni infekce
 - asi neco jako kdyz si vymazu bookmark v chrome a na jinem PC, kde jeste je se sesynchronizuje a je zpatky
- aplikace - agregace dat
 - podproblém: spočítejte průměr
 - uzel i: x_i = iniciální (libovolná) hodnota
 - epidemicky: $(x_i, x_k) = (x_i + x_k) / 2$
 - $x_i \rightarrow \text{avg for all } n (x_n)$
 - iniciátor $x_i = 1$, ostatní $x_i = 0$
 - $x_i \rightarrow 1 / N$