

Distribuované operační systémy

Filip Zavoral

Pomocné ucební texty k stejnojmenné přednášce

v. 12.3.2001

Upozornění: tyto ucební texty nejsou přesným obrazem na přednášce probírané látky – některé pasáže tu jsou popsány nebo zobrazeny velmi stručně či heslovitě, naopak jiné pasáže jsou rozebrány podrobněji než bývá obvyklé na přednášce. Proto je nelze považovat za plnohodnotný a úplný zdroj informací pro účely složení zkoušky.

Obsah

1. ÚVOD.....	6
1.1 POTENCIÁLNÍ VÝHODY DISTRIBUOVANÝCH SYSTÉMU.....	6
1.2 NEVÝHODY DISTRIBUOVANÝCH SYSTÉMU.....	7
1.3 DEFINICE A FUNKCE DISTRIBUOVANÝCH OPERAČNÍCH SYSTÉMU.....	8
1.3.1 <i>Transparentnost</i>	8
1.3.2 <i>Prizpusobivost</i>	9
1.3.3 <i>Spolehlivost</i>	10
1.3.4 <i>Výkonnost</i>	10
1.3.5 <i>Rozšířitelnost</i>	11
1.4 ARCHITEKTURY A MODELY DISTRIBUOVANÝCH SYSTÉMU.....	11
1.5 HARDWARE.....	12
2. MEZIPROCESOVÁ KOMUNIKACE.....	15
2.1 KLIENT/SERVER MODEL.....	15
2.1.1 <i>Request/reply protokol</i>	16
2.1.2 <i>Struktura zpráv</i>	17
2.1.3 <i>Synchronní a asynchronní komunikace</i>	17
2.1.4 <i>Prímé zasilání zpráv vs. zasilání zpráv pres schránku</i>	19
2.1.5 <i>Prímé a nepřímé adresování - linky a porty</i>	20
2.2 SPOLEHLIVOST KOMUNIKACE.....	20
2.2.1 <i>Spolehlivost serveru</i>	21
2.2.2 <i>Spolehlivost klienta</i>	22
2.2.3 <i>Vztah komunikace a virtuální paměti</i>	22
2.3 REMOTE PROCEDURE CALL.....	23
3. SKUPINOVÁ KOMUNIKACE.....	25
3.1 UNICAST, BROADCAST & MULTICAST.....	25
3.2 PRACOVNÍ SKUPINY.....	25
3.3 ATOMICITA.....	26
3.4 USPOŘÁDÁNÍ ZPRÁV.....	27
3.5 DORUCOVACÍ PROTOKOLY.....	29
3.5.1 <i>Kauzální závislost</i>	29
3.5.2 <i>Kauzální uspořádání dorucovaných zpráv</i>	29
3.5.3 <i>Vektorové hodiny</i>	29
3.5.4 <i>Doruovací protokol pro jednu skupinu</i>	30
3.5.5 <i>Doruovací protokol pro překrývající se skupiny</i>	31
4. SYNCHRONIZACE.....	32
4.1 USPOŘÁDÁNÍ UDÁLOSTÍ A LOGICKÉ HODINY.....	32
4.2 SYNCHRONIZACE FYZICKÝCH HODIN.....	34
4.2.1 <i>Cristianuv algoritmus</i>	35
4.2.2 <i>Berkeley algoritmus</i>	35
4.2.3 <i>Distribuovaný algoritmus</i>	36
4.3 VZÁJEMNÉ VYLOUCENÍ PROCESU.....	36
4.3.1 <i>Centralizovaný algoritmus</i>	36
4.3.2 <i>Lamport</i>	37
4.3.3 <i>Ricart & Agrawala</i>	38
4.3.4 <i>Suzuki & Kasami</i>	38
4.3.5 <i>Token ring</i>	38
4.3.6 <i>Porovnání jednotlivých algoritmu</i>	39
4.3.7 <i>Vyloučení procesu s prioritami</i>	39
4.4 VOLBA KOORDINÁTORA.....	39
4.4.1 <i>Bully algoritmus</i>	39
4.4.2 <i>Kruhový algoritmus</i>	40
4.4.3 <i>Kruhový algoritmus - $O(n \log n)$</i>	40
4.5 SYNCHRONIZACE UKONČENÍ A DETEKCE GLOBÁLNÍHO STAVU.....	41

4.5.1	<i>Dijkstra-Scholten (DS) algoritmus</i>	41
4.5.2	<i>Znackový (TM) algoritmus</i>	42
4.5.3	<i>Detekce globálního stavu</i>	42
5.	TRANSAKČNÍ ZPRACOVÁNÍ	44
5.1	VLASTNOSTI TRANSAKČÍ.....	44
5.1.1	<i>Vnorené transakce</i>	45
5.1.2	<i>Zotavení z chyb</i>	45
5.2	IMPLEMENTACE TRANSAKČÍ.....	46
5.2.1	<i>Stabilní pamet</i>	46
5.2.2	<i>Lokální pracovní prostor</i>	46
5.2.3	<i>Intenční seznam</i>	47
5.3	TRANSAKČNÍ KOMUNIKAČNÍ PRIMITIVA.....	47
5.3.1	<i>Dvoufázový commit (potvrzování)</i>	48
5.4	KONTROLA KONKURENCE.....	49
5.4.1	<i>Zámky</i>	49
5.4.2	<i>Dvoufázové uzamykání</i>	49
5.4.3	<i>Optimistická kontrola konkurence</i>	50
5.4.4	<i>Casové značky</i>	50
6.	DISTRIBUOVANÁ SDÍLENÁ PAMET (DSM)	52
6.1	KONZISTENČNÍ MODELY.....	52
6.1.1	<i>Striktní konzistence (strict consistency)</i>	52
6.1.2	<i>Sekvenční konzistence (sequential consistency)</i>	52
6.1.3	<i>Kauzální konzistence (causal consistency)</i>	54
6.1.4	<i>PRAM konzistence (PRAM consistency)</i>	54
6.1.5	<i>Slabá konzistence (weak consistency)</i>	55
6.1.6	<i>Uvolňovací konzistence (release consistency)</i>	56
6.1.7	<i>Prístupová konzistence (entry consistency)</i>	56
6.1.8	<i>Shrnutí konzistenčních modelů</i>	57
6.2	DISTRIBUOVANÉ STRÁNKOVÁNÍ.....	57
6.2.1	<i>Nalezení vlastníka stránky</i>	58
6.2.2	<i>Nalezení kopií</i>	58
6.2.3	<i>Alokace stránek</i>	58
6.3	DISTRIBUOVANÉ SDÍLENÉ PROMENNÉ.....	59
6.3.1	<i>Munin</i>	59
6.4	DISTRIBUOVANÉ OBJEKTY.....	59
7.	IDENTIFIKACE OBJEKTU	60
7.1	IDENTIFIKAČNÍ SYSTÉM.....	60
7.1.1	<i>Jména</i>	61
7.1.2	<i>Struktura jmen</i>	61
7.1.3	<i>Cesty</i>	62
7.1.4	<i>Adresy a mapování</i>	63
7.2	SYSTÉMOVÁ JMÉNA.....	63
7.2.1	<i>Interní identifikátory</i>	63
7.2.2	<i>Lokálně implementační identifikátory</i>	64
7.2.3	<i>Mapování systémových jmen</i>	64
7.3	KAPABILITY.....	64
7.3.1	<i>Kapability s podpisem</i>	65
7.3.2	<i>Kapability s redundancí</i>	65
7.4	UŽIVATELSKÁ JMÉNA.....	66
7.4.1	<i>Servry jmen</i>	66
7.4.2	<i>Správa prostoru jmen</i>	66
7.4.3	<i>Rozklad jmen</i>	66
7.4.4	<i>Agenti</i>	67
7.4.5	<i>Vyhledávání jmen</i>	68
8.	PROCESY	69
8.1	VLÁKNA, NITE A THREADY.....	69
8.1.1	<i>Použití vláken</i>	70

Distribuované operační systémy	4
8.1.2 Implementace vláken	70
8.1.3 Vlákná a zprávy	71
8.2 SYSTÉMOVÉ MODELY	71
8.2.1 Workstation model	71
8.2.2 Procesor pool model	72
8.3 VZDÁLENÉ SPOUŠTENÍ PROCESU	72
8.4 ALOKACE PROCESORU	73
8.4.1 Klasifikace alokacních algoritmu	73
8.4.2 Implementace alokacních algoritmu	74
8.4.3 Deterministický grafový algoritmus	74
8.4.4 Up-down algoritmus (Mutka-Livny)	74
8.4.5 Hierarchický algoritmus	75
8.4.6 Distribuovaný heuristický algoritmus	76
8.4.7 Bidding (obchodní) algoritmus	76
9. MIGRACE PROCESU	77
9.1 MECHANISMUS MIGRACE	77
9.1.1 Doručování zpráv	77
9.1.2 Stav procesu	78
9.1.3 Komunikace s okolím během migrace a po migraci	78
9.1.4 Implementace migrační strategie a vlastní migrace	79
9.1.5 Vícenásobná migrace	79
9.1.6 Transparence	79
9.1.7 Reziduální dependence	80
9.1.8 Virtuální pamet	80
9.1.9 Cíle návrhu migrace	81
9.1.10 Kdy migrovat	81
9.1.11 Výber cíle migrace - migrační strategie	81
9.2 PREHLED NEKTERÝCH SYSTÉMU	82
9.2.1 DEMOS/MP	83
9.2.2 Charlotte	84
9.2.3 V	86
9.2.4 MOSIX	87
9.2.5 Sprite	89
9.2.6 Migrace objektu - systém Emerald	92
9.3 KOMUNIKACE PRI MIGRACI	93
9.4 VYVAŽOVÁNÍ ZÁTEŽE	95
9.4.1 Párový algoritmus (Bryant & Finkel)	95
9.4.2 Vektorový algoritmus (Mosix)	96
9.4.3 Bidding algoritmus (Stankovic & Sidhu)	96
9.4.4 SLA algoritmus (Stankovic)	97
9.4.5 BDT algoritmus (Stankovic)	97
9.4.6 Centralizovaný algoritmus	97
9.4.7 Lokální algoritmus	97
9.4.8 Porovnání distribuovaných vyvažovacích algoritmu	97
10. SPRÁVA PROSTREDKU	98
10.1 SPRÁVCI PROSTREDKU	98
10.1.1 Centralizovaná správa prostredku	98
10.1.2 Distribuovaná správa prostredku	99
10.1.3 Správa prostredku pomocí agentu	99
10.2 ZABLOKOVÁNÍ (DEADLOCK)	100
10.2.1 Modely deadlocku	100
10.2.2 Algoritmy detekce deadlocku	101
10.2.3 Prevence deadlocku	102
10.3 OCHRANA PROSTREDKU	103
10.3.1 Klasifikace bezpecnosti DoD - the Orange Book	103
11. SPRÁVA SOUBORU	104
11.1 DISKOVÉ SLUŽBY	104
11.1.1 Blokovaná struktura souboru	104

11.1.2	<i>Systémová identifikace</i>	104
11.1.3	<i>Vyrovnávací pamet (cache)</i>	104
11.2	SOUBOROVÉ SLUŽBY.....	105
11.2.1	<i>Sémantika sdílení souboru</i>	106
11.3	ADRESÁROVÉ SLUŽBY	107
11.3.1	<i>Prostor jmen adresáru</i>	108
11.4	REPLIKACE.....	108
11.4.1	<i>Aktualizační protokoly</i>	108

1. Úvod

Až do poloviny 80. let byly počítače velké a drahé, dokonce i minipočítače stály desítky tisíc dolarů. Většina organizací měla jeden centrální počítač (mainframe), na který bylo připojeno, podle výkonu, několik až několik desítek terminálů. Od začátku 80. let se ale situace začala výrazně měnit. V té době se začaly hromadně vyrábět výkonné mikroprocesory. Bežně dostupnými se staly minipočítače založené nejprve na 8-bitových, později na 16- a 32-bitových a v současné době už i na 64-bitových mikroprocesorech. Mnoho z nich mělo výkonnost srovnatelnou se středně velkým střediskovým počítačem, avšak za zlomek jeho ceny. Druhou změnou, která ovlivnila vývoj počítačového průmyslu, byl vývoj rychlých lokálních sítí (LAN). Tyto sítě dokáží přenášet informace mezi počítači řádově za milisekundy, větší objemy dat jsou přenášeny rychlostí desítek milionů bitů za sekundu.

Důsledkem vývoje těchto technologií byl od poloviny 80. let nárůst sítí osobních počítačů, na kterých jsou provozovány **sítové a distribuované operační systémy**. V síťových operačních systémech má každý počítač svůj operační systém. Každý z těchto operačních systémů má nadstavbu umožňující komunikaci s ostatními systémy v síti. Uživatelé mají na zřeteli existenci více počítačů - mohou se explicitně přihlásit ke vzdálenému serveru, kopírovat soubory z jednoho počítače na druhý apod.

Distribuovaný systém poskytuje vyšší stupeň transparentnosti a sdílení prostředků nežli síťové systémy. Distribuované systémy je možno rozdělit na **distribuované systémy na uživatelské úrovni** (distribuované aplikace, distribuované výpočetní prostředí) a **distribuované operační systémy**. V distribuovaných systémech na uživatelské úrovni je podpora distribuce umístěna v softwarové vrstvě nad (nedistribuovaným) operačním systémem. Příkladem takového systému je OSF Distributed Computing Environment (DCE). Distribuované operační systémy implementují podporu distribuovaného zpracování přímo v jádru operačního systému. Celý systém se tváří jako tradiční jednoprocessorový systém, protože jeho jednotlivé komponenty jsou fyzicky rozmístěny na jednotlivých počítačích. V pravém distribuovaném operačním systému se uživatel nemusí (často ani nemusí) zabývat určováním na kterém počítači či procesoru jeho programy běží nebo kde jsou fyzicky uloženy jeho soubory. O to se automaticky a (v ideálním případě) efektivně stará sám operační systém. Příkladem takového systému je např. Amoeba, Chorus, Spring apod.

Architektura většiny distribuovaných systémů je založena na modelu **klient/server**. V asymetrickém klient/server modelu jsou vybrané počítače určeny pro jednotlivé servery (např. file server, autentifikační server apod.), zatímco na ostatních počítačích běží uživatelské programy. V symetrickém klient/server modelu může být každý počítač jak serverem pro jiné počítače, tak i klientem jiných serverů. Tento model je pro použití v distribuovaných operačních systémech vhodnější pro vyšší stupeň decentralizace, lepší využití dostupných prostředků a lepší rozširitelnost.

1.1 Potenciální výhody distribuovaných systémů

Asi před čtvrt stoletím platilo zhruba pravidlo, že výpočetní síla CPU je přímo úměrná druhé mocniny ceny. Zaplacením dvojnásobné ceny se výkon počítače zvýšil čtyřikrát. To vedlo k velkému rozšíření sálových počítačů, protože většina společností kupovala nejvýkonnější dostupné počítače. Avšak s rozvojem mikroprocesorové technologie toto pravidlo přestalo platit - za několik stovek dolarů je možno získat procesor s větším výkonem než měly největší sálové počítače začátku 80. let. Zaplacením dvojnásobku ceny se výkon počítače zvětší jen nepatrně. Proto se začaly používat architektury založené na propojení většího množství levných procesorů do jednoho systému. Takto vytvořené systémy mají mnohem lepší poměr ceny a výkonu nežli centralizované systémy.

Propojením velkého množství procesorů lze dokonce dosáhnout takového celkového výkonu systému, jaký nelze dosáhnout použitím jednoho sálového počítače za jakoukoliv cenu. S použitím

současné technologie je možné vytvořit systém složený z 1000 CPU o výkonu 20 MIPS (milionu instrukcí za sekundu). Celý systém pak může mít výkon až 20000 MIPS. Aby tohoto výkonu dosáhl jediný procesor, musel by vykonat 1 instrukci za 0.05 ns, což je rychlost rádově 100krát větší nežli u nejrychlejších současných procesorů.

Dalším důvodem pro vytváření distribuovaných systémů je to, že některé aplikace jsou svým charakterem samy distribuované. Například u továrního systému řídicího roboty a stroje podél celé výrobní linky je rozumné, aby každý robot nebo skupina robotů byly řízeny vlastním počítačem. Hardwarovým a softwarovým propojením těchto počítačů dostaneme průmyslový distribuovaný systém. Podobně v bankovním propojením počítačů ve všech pobočkách banky spolu s peněžními automaty dostaneme bankovní distribuovaný systém.

Další potenciální výhodou distribuovaných systémů je spolehlivost. Při rozložení výpočetní zátěže na několik procesorů způsobí výpadek jednoho procesoru přerušení práce pouze jednoho počítače, zatímco ostatní pracují dál. V ideálním případě 5% vypnutých počítačů způsobí 5% snížení výkonu celého systému. U centralizovaných sálových počítačů většinou výpadek některého důležitého obvodu způsobí výpadek celého systému. U některých kritických aplikací, např. řízení jaderného reaktoru nebo řízení letového provozu, může být spolehlivost nejdůležitějším důvodem pro výběr distribuovaného řídicího systému.

Rozšiřitelnost systému je další výhodou oproti centralizovaným systémům. V případě, že nedostává kapacita a výkon sálového počítače, jsou možná pouze dvě řešení - vyměnit celý počítač za nový nebo přidat druhý sálový počítač. Oba dvě řešení mají své velké nevýhody. Oproti tomu u distribuovaných systémů stačí přidat několik dalších komponent do systému (např. procesory, disky, terminály apod.), aby bylo dosaženo požadovaného výkonu.

Konečně distribuované systémy umožňují vzdálené provádění úloh a rozdělení zátěže pokud možno rovnoměrně mezi propojenými počítači. Navíc, pokud jsou v systému procesory různých typů nebo výkonu, lze každé úloze přiřadit nejvhodnější procesor.

Ekonomika	Mikroprocesory poskytují lepší poměr ceny a výkonu
Výkon	Distribuovaný systém může dosáhnout lepšího celkového výkonu než potenciálně jakýkoliv sálový počítač
Inherentní distribuovanost	Některé aplikace jsou svým charakterem samy distribuované
Spolehlivost	Při selhání části systému může zbytek systému pokračovat v činnosti
Rozšiřitelnost	Kapacita a výkon systému může být po částech zvyšován
Sdílení	Sdílení společných dat a periférií
Vzdálený běh procesu	Sdílení a vyvažování výpočetové zátěže

Tab. 1 - Výhody distribuovaných systémů

1.2 Nevýhody distribuovaných systémů

Přestože distribuované systémy mají mnoho výhod, mají i některé slabiny. Jednou z největších slabín je software. V současné době s tvorbou a použitím distribuovaného softwaru není příliš zkušeností. Stále existují otevřené otázky o vhodném operačním systému, programovacích jazycích, vhodných aplikacích a nástrojích pro jejich tvorbu. Není jasno o tom, co vše mají uživatelé vedet o distribuci, která rozhodování má provádět systém sám a která rozhodnutí ponechat na uživateli apod.

Software pro distribuované systémy, obzvláště distribuované operační systémy, musí řešit mnohem více problémů a provádět mnohem více rozhodování než jejich centralizované protějšky a ekvivalenty. Proto je pro jejich běh zapotřebí dostatečně kvalitní hardware - výkonné počítače s dostatkem operační paměti i diskového prostoru.

Další problémy způsobuje vlastní propojení jednotlivých počítačů. Komunikace může být nespolehlivá, mohou se ztrácet jednotlivé zprávy, což vyžaduje další software, který tyto problémy

bude řešit. Navíc při intenzivní komunikaci se může síť zahltnit. Pokud síť nebude dostatečně pružná, znegují se téměř všechny výhody, které byly použitím distribuovaného systému dosaženy.

Snadné sdílení dat vyvolává další problém - bezpečnost. Je zapotřebí zajistit, aby utajovaná data nebyla použita neoprávněným uživatelem, aby data jednoho uživatele nebyla změněna nebo dokonce smazána jiným uživatelem apod. Nedožité následky nedostatečně vyřešené bezpečnosti jsou zřejmě na první pohled např. v bankovníctví, v oblasti výzkumu zvláštní důležitosti, ve vojenství apod.

Pres výše uvedené problémy prevažuje názor, že výhody distribuovaných systémů prevažují jejich nevýhodami, a že distribuované systémy budou hrát v nejbližších letech stále důležitější roli.

Software	Málo existujícího software, velké požadavky na hardware
Propojení do sítě	Síť může být zahlcena
Bezpečnost	Snadný přístup k datům vyvolává nutnost utajení

Tab. 2 - Nevýhody distribuovaných systémů

1.3 Definice a funkce distribuovaných operačních systémů

Distribuovaný operační systém je systém, který běží na množině procesorů, které nesdílejí společnou paměť, přičemž poskytuje uživateli dojem jednoho počítače.

V literatuře bývá tato vlastnost někdy označována jako **jednotný obraz systému (single system image)** nebo také **virtuální uniprocessor**. Kromě této základní vlastnosti by distribuovaný operační systém měl:

- efektivně řídit přidělování prostředků dostupných na síti
- skrýt fyzické rozmístění prostředků
- poskytovat mechanismus ochrany přístupu systémových prostředků před neautorizovanými uživateli
- umožňovat spolehlivou a bezpečnou komunikaci

Základem každého distribuovaného systému je jednotný, globální mechanismus meziprocetové komunikace (IPC) takový, aby každý proces mohl jednotně komunikovat s jakýmkoliv jiným procesem. Tato komunikace by měla být zcela transparentní bez ohledu na fyzické umístění jednotlivých procesů. Komunikace procesů musí být v některých případech synchronizována. Protože však v distribuovaných systémech neexistuje společná paměť ani společné hodiny, navíc při komunikaci vznikají prodlevy, nelze použít synchronizační metody určené pro centralizované operační systémy. Z tohoto důvodu by každý distribuovaný systém měl poskytovat jednotný synchronizační mechanismus.

Vzhledem k tomu, že úlohy jednoho uživatele mohou běžet na několika procesorech a zároveň úlohy několika uživatelů mohou běžet na jednom procesoru, musí distribuovaný operační systém zabezpečit jednotný systém ochrany a autentifikace. Ze stejného důvodu je nezbytné, aby pro všechny propojené počítače existoval jednotný identifikační systém, který by zaručoval jednotný přístup ke všem prostředkům systému (např. souborů) z libovolného počítače.

Pro dosažení všech výše uvedených vlastností je zapotřebí při tvorbě distribuovaného systému zabezpečit následující podmínky.

1.3.1 Transparentnost

Pravděpodobně nejdůležitější podmínkou dosažení iluze jednotného systému je transparentnost. Té může být dosaženo ve dvou úrovních - na úrovni komunikace s uživatelem nebo na úrovni komunikace procesů s jádrem systému. Dosažení transparentnosti na systémové úrovni klade mnohem

větší nároky na jádro operačního systému, avšak výsledný výkon a schopnosti celého systému bývají výrazně lepší.

Idea transparentnosti může být vztažena na několik oblastí distribuovaných systému.

Prístupová transparentnosť	Proces má stejný přístup k jak k lokálním tak ke vzdáleným prostředkům
Lokální transparentnosť	Uživatel (proces) nemůže říci kde jsou jednotlivé prostředky umístěny
Migrační transparentnosť	Prostředky se mohou přemísťovat mezi jednotlivými počítači
Exekuční transparentnosť	Procesy mohou běžet na libovolném procesoru zapojeném do systému
Replikací transparentnosť	Uživatel nemůže říci kolik kopií daného objektu existuje
Konkurenční transparentnosť	Prostředky mohou být automaticky využívány zároveň několika uživateli
Paralelismová transparentnosť	Různé činnosti mohou být prováděny paralelně bez vědomí uživatele

Tab. 3 - Druhy transparentnosti v distribuovaném systému

Prístupová transparentnosť znamená, že proces má stejnou metodu přístupu jak k lokálním tak ke vzdáleným prostředkům. Např. přístup k nějakému souboru je stále stejný i když je soubor fyzicky umístěn na jiném počítači. **Lokální transparentnosť** zabezpečuje, že uživatelský proces nemusí znát aktuální umístění jednotlivých prostředků. Ve jménech prostředků nesmí být zakódováno jejich umístění, např. systém je jméno souboru tvaru *pocítac1:cesta/jméno_souboru* nesplňuje podmínku lokální transparentnosti.

Migrační transparentnosť znamená, že prostředky musí být schopné přemístění na libovolné jiné místo bez změny jména a bez vlivu na procesy, které právě daný prostředek využívají nebo používají. Operační systém musí zabezpečit, aby po dobu přemísťování prostředků (např. souboru na jiný server) byly požadavky pozdrženy a poté přesmerovány na nové místo. Podobným způsobem funguje i **exekuční transparentnosť** - procesy musejí být schopny běžet na libovolném procesoru zapojeném do systému. Tato vlastnost umožňuje přemísťování běhu programu a sdílení a vyvažování výpočtové zátěže.

Jestliže distribuovaný systém splňuje podmínku **replikací transparentnosti**, pak umožňuje vytváření kopií souboru nebo jiných prostředků bez vědomí uživatele. Replikované prostředky jsou velmi důležité zejména u rozsáhlých systému, kde přístup např. k souboru fyzicky umístěnému na jiném kontinentě by trval nepřiměřeně dlouhou dobu. Replikace prostředků je také významná z důvodu spolehlivosti - v případě, že jedna z kopií selže, ostatní kopie mohou být nadále využívány.

Distribuované systémy většinou umožňují přístup více uživatelům najednou. Systém s **konkurenční transparentností** by měl umožňovat přístup uživatele k libovolným prostředkům bez ohledu na to, kolik uživatelů tyto prostředky právě využívá.

Nejzajímavější, ale zároveň nejméně probádaná vlastnost je **paralelismová transparentnosť**. Pokud je v systému velké množství volných procesorů, uživatel by měl mít možnost je využívat co možná nejvíce podle jeho potřeb. Operační systém by měl teoreticky sám využít všech dostupných procesorů a provádět danou úlohu na tolika procesorech, aby celkový čas provádění byl co nejmenší. Bohužel, současné vědomosti jsou tomuto ideálu na hony vzdáleny - v naprosté většině současných systému úloha, která by měla využívat zároveň několik procesorů, musí být speciálně naprogramována s explicitními požadavky na procesory.

1.3.2 Prizpůsobivost

Další klíčovou podmínkou dobře fungujícího distribuovaného systému je prizpůsobivost (flexibilita). Vzhledem k tomu, že distribuovaný systém běží na mnohem složitější struktuře hardwaru

než klasické operační systémy, měl by být daleko méně citlivý na změnu prostředí. Prizpusobivost může být opět vztažena na několik oblastí distribuovaných systému.

Autonomie	Každý počítač je potenciálně schopný samostatné funkčnosti
Decentralizované řízení a rozhodování	Každý procesor vykonává rozhodnutí nezávisle na ostatních
Migrace procesu a prostředku	Procesy i prostředky mohou být přemístěny na jiný počítač
Vyvažování výpočtové zátěže	Úlohy mohou být přemístěny na méně vytížený procesor

Tab. 4 - Různé vlastnosti prizpusobivých distribuovaných systému

Jednotlivé počítače zapojené do systému by měly být **autonomní**, tj. měly by být schopny samostatné funkčnosti. To je ve většině moderních distribuovaných systému zajištěno architekturou **microkernel**. Na každém počítači je implementováno pouze nejnútnejší jádro systému obsahující meziprocesovou komunikaci, správu vnitřní paměti, nízkoúrovňovou správu procesu a nízkoúrovňové služby pro styk s periferiemi. Ostatní služby operačního systému, jako např. filesystem, jsou implementovány jako servery na uživatelské úrovni. Tímto způsobem je např. možné vytvořit distribuovaný operační systém s několika fileservery - jeden podporující souborové služby např. MS-DOSu, druhý umožňující přístup k souborům na UNIXové stanici apod.

Důležitou vlastností prizpusobivých distribuovaných systému je decentralizované řízení další činnosti a decentralizované rozhodování. V případě, že každý procesor vykonává svá rozhodnutí nezávisle na ostatních procesorech, je celý systém mnohem méně citlivý na různé poruchy sítě, přidávání a odstraňování pracovních stanic, výpadky serveru a pod.

Plnohodnotný distribuovaný systém by měl umožňovat **migraci procesu a prostředku** mezi jednotlivými počítači. V případě, že nějaký server spravuje větší množství prostředku, než na co stací jeho kapacita, lze část prostředku přemístit na jiný počítač. Stejně tak procesoru, na kterém běží velké množství aktivních procesu, lze odlehčit přemístěním určitého množství procesu na méně vytížený procesor (sdílení zátěže, load sharing). Migrace také umožňuje přemístění prostředku z počítače, který má být ze systému dočasne či trvale vyrazen, a tím zachování plné funkčnosti celého systému.

Transparentní migrace procesu je podmínkou toho, aby distribuovaný systém umožňoval **vyvažování výpočtové zátěže** (load balancing), některými autory nazývané též distribuované plánování (distributed scheduling). Vyvažování zátěže spočívá v automatické distribuci procesu v celém systému tak, aby výsledná zátěž na jednotlivých procesorech byla zhruba stejná.

1.3.3 Spolehlivost

Jednou z původních motivací pro tvorbu distribuovaných systému bylo dosažení vyšší spolehlivosti nežli u klasických centralizovaných systému. Základní ideou zvýšení spolehlivosti je to, že v případě výpadku některých komponent systému si ostatní rozdělí jejich práci. Jestliže např. nespolehlivost jednoho serveru je 1%, pak by teoreticky nespolehlivost čtyř sprážených serveru měla být $0.01^4 = 0.00000001$, tedy jedna milióntina procenta. Současné distribuované systémy však mají služby jednotlivých serveru natolik vzájemne propojeny, že zvyšování spolehlivosti probíhá mnohem pomaleji, pokud vůbec. V jedné často citované poznámce Leslie Lamport, klasika v oblasti distribuovaných systému, definoval distribuovaný systém jako: *"Systém, který pro mě odmítá vykonávat jakoukoliv činnost, protože počítač, o kterém jsem nikdy neslyšel a ani nepotřebuji jeho služby, byl náhodne vypnut"*. Přestože tato poznámka byla silne nadnesena, je pravdou, že v oblasti spolehlivosti distribuovaných systému je ještě mnoho oblastí zasluhujících si dukladný výzkum.

1.3.4 Výkonnost

K tomu, aby byl systém prakticky použitelný, nestací jen dosažení flexibility, transparentnosti a spolehlivosti. Když systém bude mít všechny tyto vlastnosti, ale bude nepoužitelne pomalý, bude

k nicemu. Konkrétně beh určité aplikace v distribuovaném systému by nemel být výrazne pomalejší než v klasickém jednoprocessorovém systému stejné třídy. Bohužel, současné distribuované systémy toto pravidlo příliš nespĺnují - pro dosažení odpovídající rychlosti behu aplikace je zapotřebí výkonnejší hardware. Zpomalení je způsobeno jednak pomalejším přenosem zpráv po síti, jednak složitějším softwarem.

Pri návrhu systému je zapotřebí dát velký pozor na granularitu výpocetní jednotky. V prípade, že se budou vzdáleně provádět příliš malé kousky programu, jako napr. scítání dvou čísel, neúmerne vzroste potreba komunikace a nekolikanásobne se zvýší režie potřebná pro inicializaci, rozbehnutí a ukončení jednotlivých částí programu. Na druhé strane přenášení příliš velkých kusu programu na vzdálený počítač může trvat nepřiměřene dlouho.

1.3.5 Rozširitelnost

Všechny současné distribuované systémy jsou konstruovány pro práci nekolika desítek až stovek procesoru. Potenciálně je však možné, že budoucí komplexní distribuované systémy budou o nekolik rádu větší a řešení, která fungovala dobře pro 200 procesoru budou naprosto nevhodná pro 200?000?000 procesoru. Prestože současné znalosti o takto obrovských systémech jsou minimální (a zkušenosti s nimi prakticky žádné), jeden základní princip je znám: vyhnout se cemukoliv centralizovanému - serverum, tabulkám, algoritmus apod. Základními charakteristikami distribuovaných algoritmu jsou:

1. Žádný počítač nemá úplnou informaci o celkovém stavu systému
2. Procesy se rozhodují pouze na základe lokálně dostupných informací
3. Výpadek jednoho počítače nesmí způsobit nefunkčnost algoritmu
4. Nelze spoléhat na existenci přesných globálních hodin

Význam prvních třech charakteristik je zřejmý. Čtvrtá je možná méně zřejmá, avšak stejně důležitá. Jakýkoliv algoritmus založený na myšlence "Přesně ve 12:00 si všechny počítače zaznamenají velikost výstupní fronty" musí nutně selhat, protože nelze všechny hodiny přesně sesynchronizovat. Na lokální síti je sice možné sesynchronizovat hodiny na rozdíl nekolika milisekund, avšak v celostátním či celosvetovém měřítku je to nemožné.

koncept	Příklad
Centralizované komponenty	Jednotný mail server pro všechny uživatele
Centralizované tabulky	Jednotný telefonní seznam
Centralizované algoritmy	Smerovací algoritmus na základe úplné informace

Tab. 5 - Potenciální úzká místa velkých distribuovaných systému

1.4 Architektury a modely distribuovaných systému

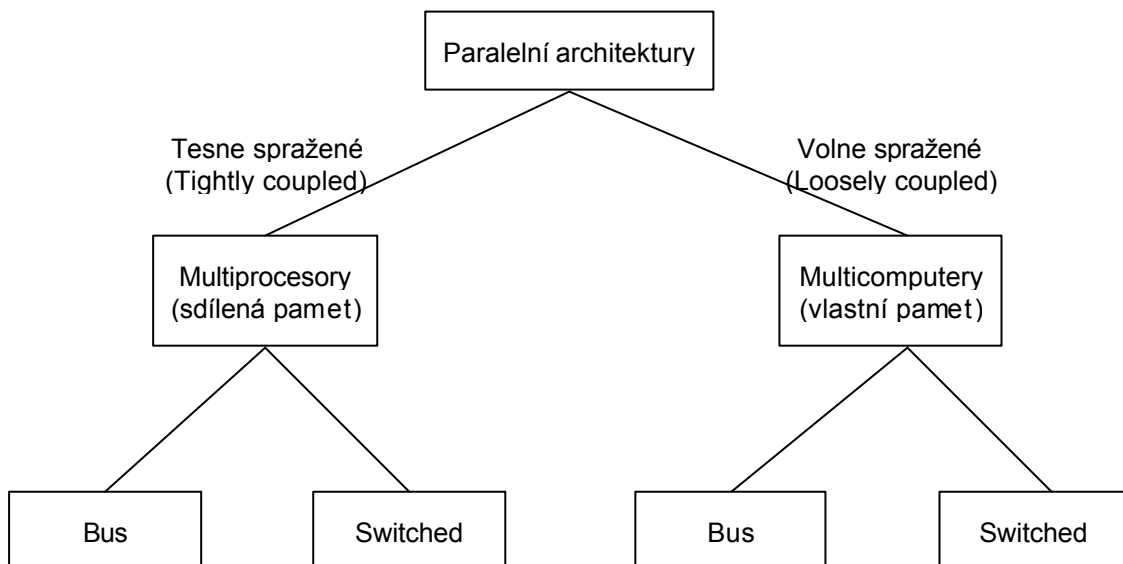
workstation / processor pool / hybrid architecture
 process / object / UNIX-like model

1.5 Hardware

Přestože jsou všechny distribuované systémy tvořeny určitým množstvím procesorů, je několik rozdílných možností jejich organizace a propojení. Podle přístupu k paměti je můžeme rozdělit na **multiprocesory**, které sdílejí společnou paměť, a **multicomputery**, které společnou paměť nesdílejí.

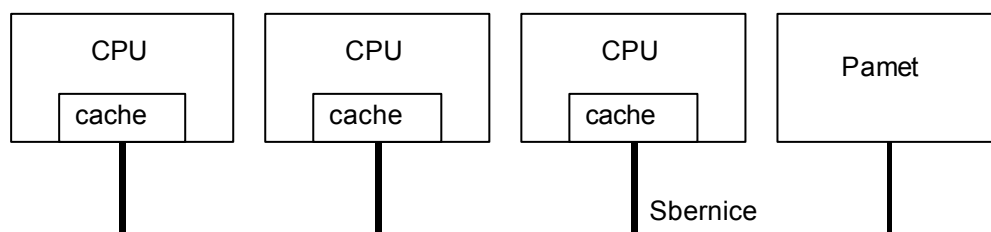
U multiprocesoru existuje jeden společný adresový prostor pro všechny propojené procesory. To znamená, že když jeden procesor zapíše na nějakou adresu určitou hodnotu, druhý procesor z této adresy přečte totéž číslo. Naproti tomu každý procesor multicomputeru má vlastní adresový prostor. Změna na jakékoli adrese nijak neovlivní adresový prostor ostatních procesorů.

Každá z těchto dvou kategorií může být dále rozdělena podle způsobu propojení jednotlivých uzlů. **Sbennicová architektura** (bus architecture) využívá jednotné médium, např. kabel, sbennici, či jiné médium, které propojuje všechny uzly. Na podobném principu funguje např. kabelová televize. **Prepínacová architektura** (switched architecture) naopak používá přímé propojení mezi jednotlivými uzly. Zprávy se zasílají mezi přímo propojenými uzly, přičemž v každém meziuzlu se provádí rozhodnutí o dalším postupném cíli zprávy. Na tomto principu funguje veřejná telefonní síť.



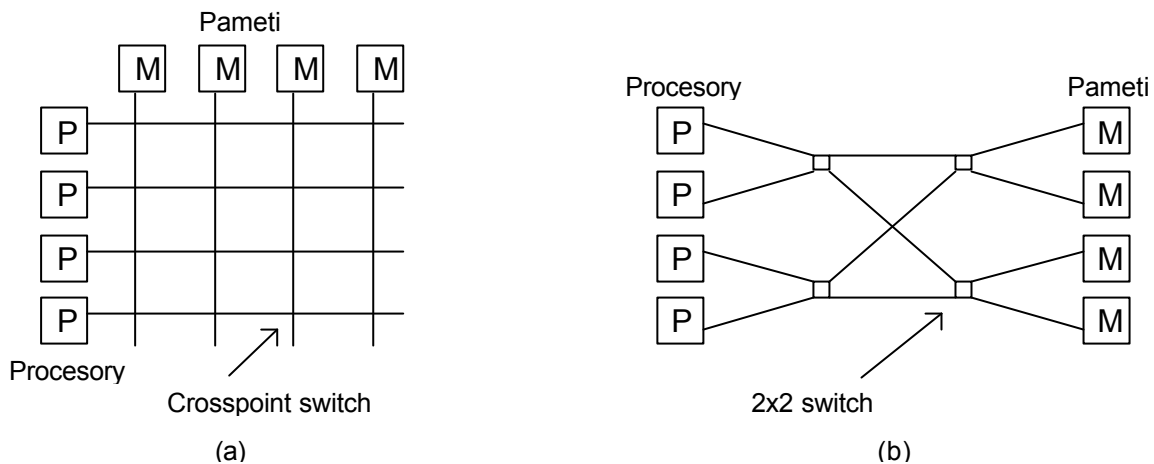
Obr. 1 - Taxonomie paralelních hardwarových architektur

Multiprocesory se sbennicovou architekturou jsou tvořeny několika procesory propojenými spolu s paměťovým modulem společnou sbennicí. Tímto způsobem je možné efektivně zapojit 32 až 64 procesorů.



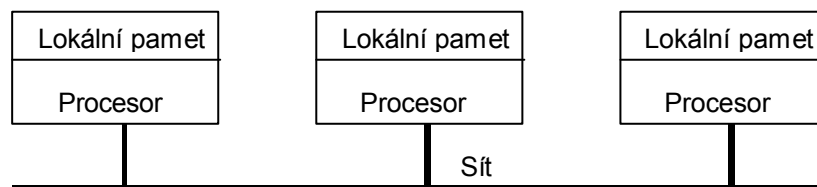
Obr. 2 - Multiprocesor se sbennicovou architekturou

Pro zapojení více než 64 procesoru je zapotřebí použít jinou architekturu. Jedno možné řešení je rozdělit paměť na několik pametových modulu a propojit je s procesory sbernicovou mřížkou. V průsečících mřížky jsou prepínace (**crosspoint switch**) určující další cestu. Nevýhodou této architektury je kvadratický počet prepínaců, což při velkém počtu procesoru a pametových modulu může být příliš nákladné. Druhou možností je použít síť typu omega (**omega network**), kde počet prepínaců roste rychlostí $n \cdot \log n$. Nevýhodou tohoto uspořádání je prodleva při větší hloubce prepínaců.



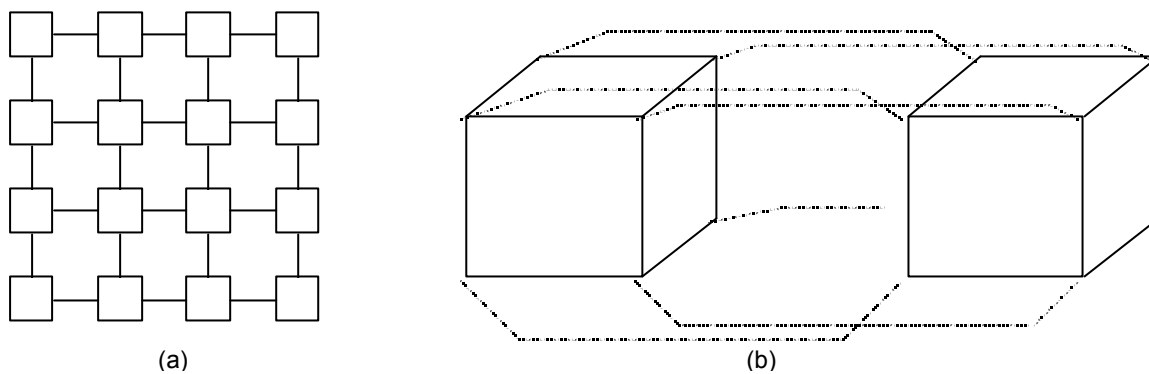
Obr. 3 - (a) Crossbar switch (b) Omega switch

Na rozdíl od rozsáhlých multiprocessoru je vytvoření multicomputeru mnohem méně nákladné. Každý procesor má přímý přístup do vlastní paměti. Vzhledem k tomu, že komunikují pouze procesory mezi sebou, je komunikace rádove nižší než u multiprocessoru. Na sbernici je proto možné umístit větší množství procesoru (pracovních stanic).

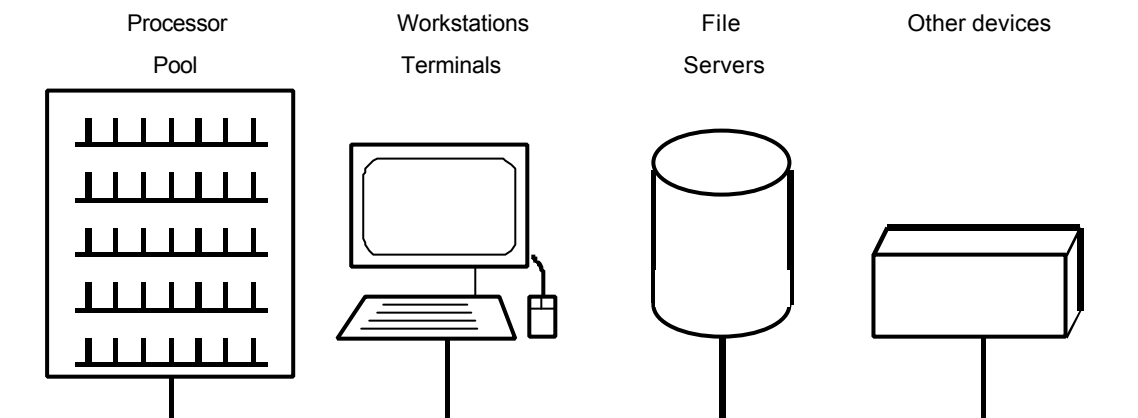


Obr. 4 - Multicomputer sestavený z pracovních stanic propojených sítí

Poslední možností jsou multicomputery s prepínacovou architekturou. Nejvíce používané topologie jsou **mřížka** a **hyperkrychle**. U mřížky je každý uzel propojen se čtyřmi svými sousedy. Jejich implementace je vzhledem k přirozené dvourozměrnosti relativně jednoduchá. Multicomputery s mřížkovou topologií jsou vhodné zejména pro úlohy z oblasti teorie grafu, analýzu obrazu apod. Hyperkrychle je n -rozměrná krychle (podle této definice je mřížka 2-rozměrná hyperkrychle). Nejčastěji používané hyperkrychle jsou 4-rozměrné, kde každý uzel obsahující procesor je propojen se dvěma uzly v každém rozměru, tj. celkem s osmi sousedními procesory. Přidáním každého dalšího rozměru se počet přímo propojených sousedů zvětší o 2. Bežně dostupné jsou hyperkrychle s 1024 procesory, v současné době se začínají používat i hyperkrychle obsahující 16386 procesoru.



Obr. 5 - (a) Mřížka (b) Hyperkrychle



Obr. 6 - Komponenty distribuovaného systému

2. Meziprocesová komunikace

Asi nejvýraznějším rozdílem mezi jednoprocessorovým a distribuovanými systémy je meziprocesová komunikace (IPC). V klasickém jednoprocessorovém systému naprostá většina implementací IPC implicitně předpokládá existenci společné paměti. Typickým příkladem je problém producent - konzument, kde jeden proces zapisuje do společného bufferu a druhý proces z něj čte. Dokonce i nezákladnější forma synchronizace, semafor, vyžaduje pro svoji funkci jedno slovo společné paměti. Oproti tomu v distribuovaných systémech nelze předpokládat existenci společné paměti, tudíž veškeré formy meziprocesové komunikace a synchronizace musejí být založeny na jiném principu - na principu **zasílání zpráv**. Jestliže proces A chce komunikovat s procesem B, nejdříve vytvoří zprávu ve svém adresovém prostoru. Poté zavolá systémovou službu, která převezme zprávu a pošle ji po síti počítači, na kterém běží proces B. Jakkoli je tato idea jednoduchá, je zapotřebí vyřešit mnoho problémů, aby tato vzdálená komunikace fungovala a komunikující počítače si vzájemně rozumely.

Pro usnadnění vzájemné komunikace mezi počítači různých výrobců vytvořila Mezinárodní standardní organizace (ISO) sedmivrstevný referenční model, označovaný jako ISO/OSI (Open Systems Interconnection). Úlohou tohoto referenčního modelu je vymezení jednotlivých vrstev a specifikace úkolů, které by měly tyto vrstvy řešit. Každá vrstva se zabývá jednou úrovní komunikace a poskytuje služby bezprostředně vyšší vrstvě.

Aplikační	Společné části síťových aplikací
Prezentací	Konverze dat, komprese, šifrování
Relační	Navazování, udržování a rušení relací
Transportní	Sestavování paketu, kvalitnější služby
Síťová	Iluze všeobecné propojenosti
Linková	Přenos bloku dat mezi uzly
Fyzická	Přenos jednotlivých bitů

Tab. 6 - Referenční model ISO/OSI

Pro použití v distribuovaných systémech je však plná implementace protokolu na všech vrstvách příliš neefektivní. Vzhledem k tomu, že každá vrstva si přidává ke zprávě svoji vlastní hlavičku, je časová i prostorová režie pro každou zprávu příliš velká. Při použití v rychlých lokálních sítích může tato režie výrazně zpomalovat chod celého systému. Proto naprostá většina distribuovaných systémů založených na lokálních sítích buď vůbec nepoužívá takto detailně vrstvené protokoly, anebo implementuje pouze omezenou množinu protokolu.

2.1 Klient/server model

Komunikační model v distribuovaných systémech by měl kromě vlastního přenosu dat mezi jednotlivými počítači specifikovat i to, jakým způsobem je celý systém strukturován. Toto řeší v naprosté většině případů tzv. klient/server model, který je postaven na základní myšlence vytvořit celý operační systém jako množinu spolupracujících procesů, nazývaných **servry**, které poskytují služby svým uživatelům, **klientům**. Počítače, na kterých běží servrovské i klientské procesy, mají většinou stejné jádro systému (microkernel), přičemž jak klienti tak servery běží jako uživatelské procesy. V případě, že je systém postaven na symetrickém klient/server modelu, mohou na libovolném počítači běžet zároveň servrovské procesy i klienti.

2.1.1 Request/reply protokol

Aby se zabránilo příliš velké režii mnohovrstevných protokolů typu ISO/OSI, je většinou klient/server model založen na jednoduchém **request/reply protokolu** (žádost/odpověď). Klient pošle serveru zprávu se žádostí o nějakou službu (request message). Server provede požadovanou činnost a pošle zpět zprávu (reply message), která obsahuje případná požadovaná data a návratový kód, který specifikuje, zda požadovaná operace byla provedena a v případě, že ne, i důvod neprovedení.

Hlavní výhodou tohoto protokolu je jeho jednoduchost. Klient pošle zprávu, dostane odpověď a to je vše. Není zapotřebí navazovat a rušit spojení, není třeba vykonávat jiné vedlejší akce, nejsou zapotřebí vyšší vrstvy protokolu. Další výhodou je efektivita - soustava protokolu je kratší, typicky jsou implementovány 3 až 4 vrstvy. Fyzická a linková vrstva starající se o přenos paketu od klienta k serveru a zpět je většinou implementována v síťové, např. ethernetové kartě. Na úrovni transportní vrstvy bývá request/reply protokol, který definuje množinu platných požadavků a množinu možných odpovědí. Další vrstvy většinou nejsou zapotřebí.

Díky této jednoduché struktuře mohou být komunikační služby poskytované jádrem systému redukovány na dvě systémová volání - jedno pro posílání zpráv, druhé pro přijímání. Tyto systémové služby mohou být vyvolány prostřednictvím knihovnických procedur `send` a `receive`. První procedura vezme připravenou zprávu a pošle ji určenému procesu, druhá procedura způsobí zablokování procesu do té doby, než dorazí nějaká zpráva.

```
main()
{
    struct message m1, m2;

    for (;;)
    {
        receive (MY_SERVER_NUMBER, &m1);
        switch (m1.opcode)
        {
            case SERVICE1:do_something (&m1, &m2);      break;
            case SERVICE2:do_something_else (&m1, &m2);break;
            case SERVICE3:do_another_service (&m1, &m2);break;
            default:      m2.retval = ERR_BAD_OPCODE;
        }
        send (m1.client, &m2);
    }
}
```

Zjednodušený příklad serveru


```

any_function()
{
    struct message m;

    ...
    m.opcode = SOME_OPCODE;
    m.arg1 = some_argument;
    m.arg2 = another_argument;
    send (MY_SERVER, &m);
    receive (MY_CLIENT_NUMBER, &m);
    if (m.retval == OK)
        read_data (&m.data);
    else
        error_handler (m.retval);
    ...
}

```

Zjednodušený příklad klienta

2.1.2 Struktura zpráv

Hlavicka, strukturovaná/nestrukturovaná data

Adresát
Odesílatel
Identifikace zprávy
Identifikace klienta
Císlo služby
Císlo objektu
Typ dat
Pocet elementu
Data
Typ dat
Pocet elementu
Data
...

Adresát
Odesílatel
Identifikace zprávy
Identifikace klienta
Císlo služby
Císlo objektu
Velikost dat
Data

Tab. 7 - Příklad formátu strukturované a nestrukturované zprávy

2.1.3 Synchronní a asynchronní komunikace

Výše popsaná komunikační primitiva nazýváme **synchronní primitiva** (blocking primitives, synchronous primitives). Když proces volá proceduru `send`, specifikuje cílový proces a buffer, který má být poslán. Po dobu, co se přenáší zpráva, je proces blokován (suspendován). Instrukce následující po volání procedury nejsou vykonávány až do té doby, dokud zpráva nedorazí k příjemci. Podobně `receive` zablokuje proces do té doby, než dorazí očekávaná zpráva.

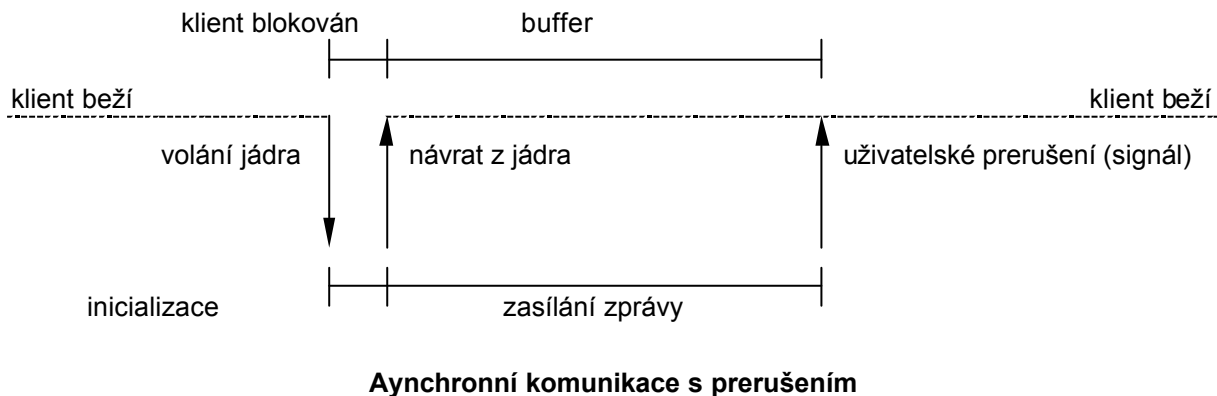
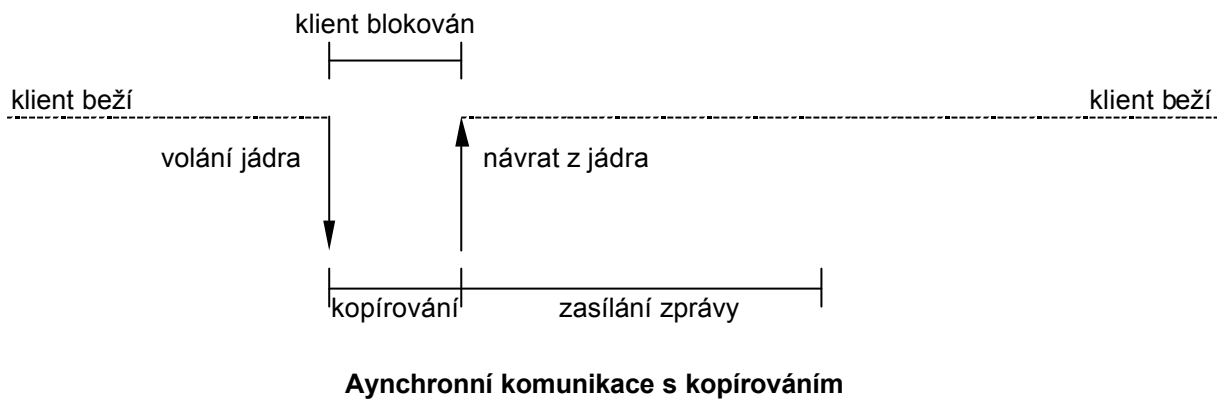
Alternativní formou komunikace jsou **asynchronní primitiva** (nonblocking primitives, asynchronous primitives). V tomto případě `send` vrátí kontrolu procesu ihned bez čekání na odeslání zprávy. Výhodou tohoto schématu je vyšší stupeň paralelismu - proces může pokračovat ve svých

výpočtech po dobu přenosu zprávy. Avšak tato výhoda je vykoupená jednou vážnou nevýhodou - proces nemůže změnit buffer se zprávou dokud zpráva není odeslána. Přepsání části zprávy během přenosu by mělo nepředvídatelné následky. Navíc proces nemá ani ponětí o tom, kdy byl přenos dokončen, tudíž nikdy neví, kdy lze bezpečně použít buffer pro jiné účely.

Existují dva způsoby řešení. První způsob spočívá v tom, že si jádro okopíruje zprávu do svých vnitřních bufferů, čímž umožní procesu další bezpečný běh. Nevýhodou tohoto řešení je kopírování každé zprávy do vnitřních bufferů jádra, což při velkém počtu zpráv může způsobit výrazné snížení výkonu celého systému.

Druhé řešení je založeno na myšlence přerušení procesu poté, co byla zpráva úspěšně odeslána. Toto řešení sice nevyžaduje žádné kopírování, nicméně programy založené na uživatelských přerušeních jsou poněkud nepřehledné, obtížné na odladení a plné kritických sekcí, čímž se stávají velice těžko udržovatelné a rozšiřitelné.

Stejně jako `send` může být i `receive` synchronní nebo asynchronní. Asynchronní `receive` pouze říká jádru, kam má umístit došlou zprávu. Zde nastává stejný problém - jak se proces dozví, že nějaká zpráva došla. Možností je několik. První možností je zavedení explicitního primitiva `wait`, které pozastaví proces do té doby, než zpráva dorazí. Další možností je zavedení primitiva `test`, které otestuje, zda je nějaká zpráva připravena, nebo `conditional_retrieve`, které buď vybere došlou zprávu nebo se vrátí s informací, že žádná zpráva není k dispozici. Poslední možností je opět zavedení uživatelského přerušení, které bude informovat proces o došlé zprávě.



Obr. 7 - Synchronní a asynchronní komunikační primitiva

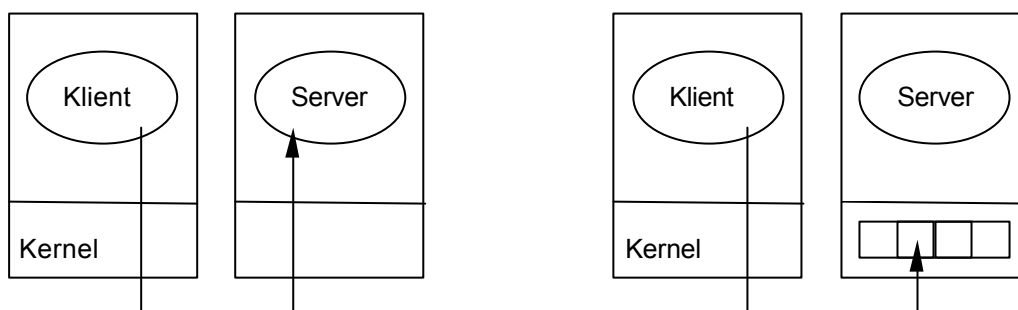
2.1.4 Prímé zasilání zpráv vs. zasilání zpráv pres schránku

Kromě volby mezi synchronní a asynchronní komunikací mají tvůrci distribuovaných systémů možnost zvolit i způsob, jak bude zpráva doručena příjemci. Výše popsaná primitiva jsou příkladem přímého zasilání zpráv (unbuffered primitives). Při tomto způsobu komunikace přijímající proces sdělí systému, kam do adresového prostoru příjemce má došlou zprávu umístit. Tímto způsobem funguje komunikace spolehlivě tak dlouho, dokud přijímající proces stihá zavolat `receive` před tím,

než libovolný proces zavolá `send`. Problém nastane, když nějaký proces posílá zprávu, ale příjemce ještě nestihl zavolat `receive`. Jádro v tomto případě neví, kam má došlou zprávu uložit.

Jedno možné řešení je došlou zprávu zahodit, informovat odesílatele o tom, že server ještě není připraven a doufat, že před tím, než dojde další zpráva, bude příjemce připraven. Toto řešení je snadné na implementaci, avšak klient, resp. vysílající jádro, může zprávu posílat nekolikrát, než je přijata. Dále je možné, že jádro po nekolika neúspěšných pokusech o spojení oznací server chybne za nefunkční nebo nedostupný.

Druhý přístup k tomuto problému je vytvoření bufferu v jádru systému - **schránky** (mailbox) - kam budou umísťovány neočekávané zprávy. V případě, že příjemce v dostatečně krátkém intervalu zavolá `receive`, obdrží zprávu ze schránky. Jestliže se příjemce neozve příliš dlouhou dobu, jádro zprávu zničí. Toto na první pohled elegantní řešení však problém neresí, pouze jej oddaluje. Schránky jsou totiž konečné a může dojít k jejich preplnění. V tom případě se jádro musí rozhodnout, zda novou zprávu přijmout a zahodit nějakou starší, anebo nechat schránku nedotčenou a novou zprávu odmítnout. V některých systémech je možné v případě preplnění příjemcovy schránky zablokovat vysílající proces do té doby, než se schránka uvolní.



Obr. 8 - (a) Přímé zasílání zpráv (b) Zasílání zpráv přes schránku

2.1.5 Přímé a nepřímé adresování - linky a porty

Meziprocesovou komunikaci je možné též delit podle toho, kam je zpráva zasílána. V případě, že odesílatel specifikuje jako příjemce přímo identifikaci přijímajícího procesu, jde o komunikaci přímou. Druhá technika - nepřímé adresování - využívá nějaké přesně určené místo, **port**, odkud příjemce zprávu dostane.

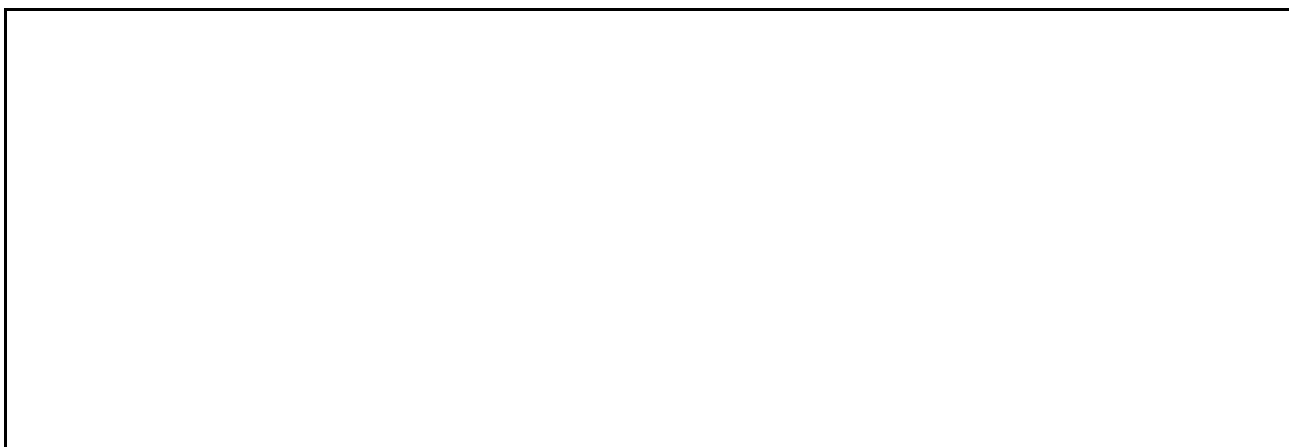
V případě přímého adresování odesílající proces specifikuje jako příjemce proces, který zprávu dostane. Mezi procesy se vytvoří virtuální spojení, tzv. **link**. Spojení je automaticky vytvořeno mezi každou dvojicí komunikujících procesů, každé spojení je asociováno právě se dvěma procesy.

Technika nepřímého adresování využívá pro doručení zprávy tzv. **porty**. Port je jednoznačně identifikovatelný objekt jádra, který umožňuje ukládání a vybírání zpráv. Vysílající proces zašle zprávu nějakému portu, přijímající proces zprávu z portu dostane. K jednomu portu může být přihlášeno více příjemců, mezi dvěma komunikujícími procesy může být navázáno spojení přes několik portů. Z implementačního hlediska je port fronta zpráv konečné délky, avšak některé "naléhavé" zprávy mohou frontu preskocit.

2.2 Spolehlivost komunikace

Doposud jsme předpokládali, že vlastní komunikace je spolehlivá, tj. když klient vyšle zprávu, tak ji server dostane, a to v nezmenené podobě. Bohužel tomu tak není - zprávy se mohou po cestě

ztrácet, může se změnit jejich pořadí apod. Proto je nutné pro zabezpečení spolehlivé komunikace zajistit potvrzování příjmu zpráv nebo paketu.



Obr. 9 – Potvrzování zpráv

reply jako ack
 potvrzování každého paketu / zprávy / dávky
 plovoucí dávky
 are you alive
 vícenásobné pakety
 timeouts

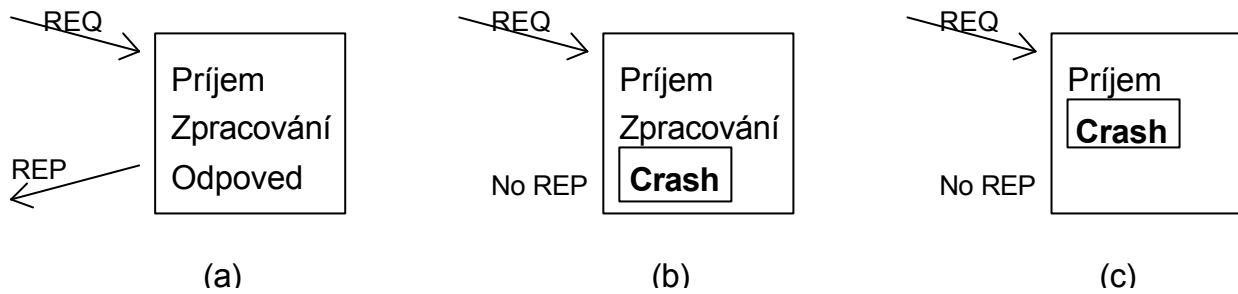
Kód	Typ zprávy	Odesílatel	Příjemce	Význam
REQ	Request	Klient	Server	Klient požaduje službu
REP	Reply	Server	Klient	Odpověď serveru na žádost klienta
ACK	Ack	kdokoliv		Potvrzení příjmu paketu
NAK	Negative ack	kdokoliv		Oznámení o nepřijmutí paketu
AYA	Are you alive?	Klient	Server	Žádost o potvrzení serveru, zda je funkční
IAA	I am alive	Server	Klient	Potvrzení funkčnosti serveru
TRA	Try again	Server	Klient	Server nemá místo pro přijmutí zprávy

Tab. 8 - Základní typy paketu používané v klient/server modelu

2.2.1 Spolehlivost serveru

Další problém komunikace v distribuovaném prostředí je spolehlivost, resp. nespolehlivost samotných počítačů, a to jak na straně serveru, tak na straně klienta. Představme si situaci, kdy klient odešle žádost, ale odpověď delší dobu nepřichází. Důvodů může být několik:

- ? ztratila se zpráva se žádostí
- ? ztratila se zpráva s odpovědí
- ? server je příliš pomalý nebo zahlcen větším počtem žádostí natolik, že ještě zpracovává žádost klienta
- ? server nefunguje



Obr. 10 – Havárie serveru

Některé operace mohou být bezpečně nekolikrát zopakovány, aniž by došlo k nějaké škodě. Např. čtení prvních n bajtů daného souboru může být provedeno nekolikrát, aniž by došlo k nějakým postranním efektům. Žádosti, resp. služby, splňující tuto vlastnost se nazývají **idempotentní**. Pokud se zasílá nějaká idempotentní žádost a nedošla odpověď, lze žádost opakovat. Ne všechny žádosti však mohou být idempotentní - např. žádost o vybrání milionu dolarů z banky rozhodně idempotentní není. Proto je nutné komunikaci zabezpečit tak, aby i tyto žádosti byly vyřizovány korektně.

První metodou je přiřazení jednotlivým zprávám jednoznačné identifikační číslo. Server tím pádem může rozeznat originální žádost od její kopie a v případě přijetí kopie zpracovávané žádosti tuto kopii ignorovat. Podobným způsobem lze řešit i ztracené zprávy s odpovědí. Dokud nedojde od klienta potvrzení o přijetí odpovědi, server si udržuje seznam vyřízených žádostí a odeslaných odpovědí a na případnou urgenci může reagovat opětovným odesláním odpovědi.

Nejsložitější situace nastane v případě, že přestane fungovat server. V případě, že nefungoval ani předtím, než zpráva došla, pak ho po několika pokusech o navázání komunikace prohlásí klient za nedostupný. Horší situace je, když server přestane fungovat poté, co zprávu přijal. Problém spočívá v tom, že obecně nelze spolehlivě určit, zda byl požadavek zpracován či ne.

Existují principiálně dva způsoby řešení této situace. První technika, tzv. **at-least-once** sémantika ("alespon jednou"), zaručuje, že požadavek (pokud to bude možné) bude zpracován, tj. po znovunastartování serveru bude požadavek opakován anebo bude předán jinému serveru zajišťujícímu stejné služby. Druhá technika, tzv. **at-most-once** sémantika ("nejvýše jednou"), zaručuje, že požadavek nebude zpracován více než jednou.

Žádná z těchto možností není příliš přitažlivá, ideální by byla sémantika typu **exactly once** ("právě jednou"), to však obecně nelze zaručit. Představme si tiskový server, který právě zpracovává požadavek na vtištění nějakého textu a přitom zhavaruje. Havárie mohla nastat mikrosekundu před vlastním tiskem i mikrosekundu po tisku, nikdo však nedokáže rozhodnout, který z těchto případů nastal.

2.2.2 Spolehlivost klienta

Havárie klienta - orphani (sirotci)

- exterminace - klient si pamatuje co rozbehl a zabije to
- reinkarnace - klient po bootu nastartuje novou epochu
- expirace - ex. kvantum, server si sám říká o další kvantum

2.2.3 Vztah komunikace a virtuální paměti

Lokální zprávy:

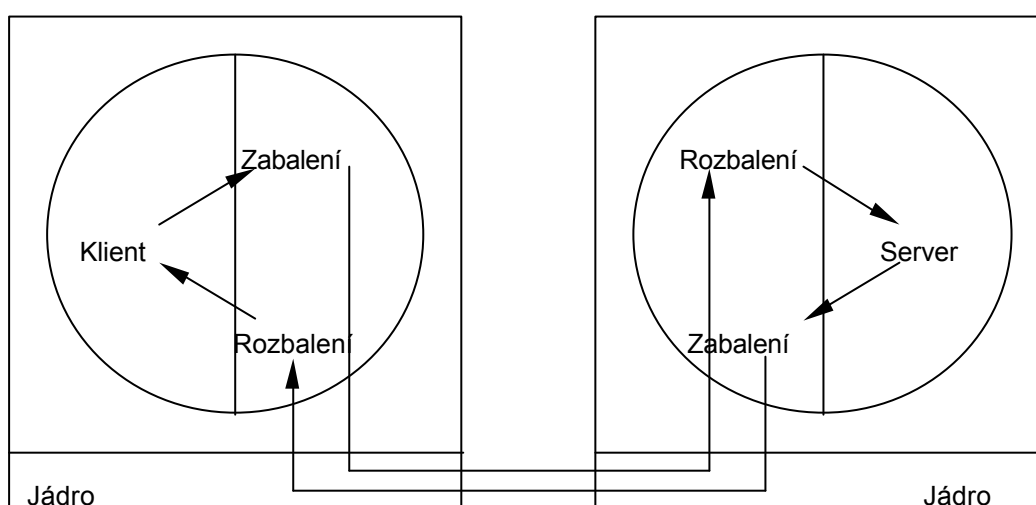
- krátké (<16B) - v registrech
- střední (<4KB) - komunikační stránka
- dlouhé - mapování

Vzdálené zprávy:

- krátké a střední zprávy - OK
- dlouhé: problém efektivity: kopírování "prázdných" stránek, vzdálené sdílení, ...

2.3 Remote procedure call

Vzdálené volání procedur, **RPC** (Remote Procedure Call), je přístup k meziprocesové komunikaci založený na základní myšlence lokálních programu - volání procedur. Idea RPC je založena na pozorování, že model "klient zašle zprávu servru, poté se zablokuje a čeká na odpověď od servru" velmi přesně odpovídá dobře známému mechanismu volání procedur. Proto cílem RPC je umožnit, aby distribuované programy mohly být psány zcela stejným stylem jako konvenční programy pro centralizované operační systémy.



Obr. 11 - Remote procedure call

1. Klientova funkce normálním způsobem zavolá klientský stub
2. Stub vytvoří zprávu a zavolá jádro
3. Jádro pošle zprávu jádru počítače, kde běží server
4. Vzdálené jádro předá zprávu servrovu stubu
5. Stub rozbálí parametry a zavolá server
6. Server zpracuje požadavky a normálním způsobem se vrátí do stubu
7. Servruv stub zabalí výstupní parametry do zprávy a zavolá jádro
8. Jádro pošle zprávu klientovu jádru
9. Klientovo jádro předá zprávu stubu
10. Stub rozbálí výstupní parametry a vrátí se ke klientovi

Problémy RPC:

- ? reprezentace dat - indiáni, floaty, ...

- ? přístup ke globálním proměnným
- ? předávání pointeru, dynamické struktury
- ? copy/restore - velikosti polí
- ? variabilní parametry - printf
- ? group communication
- ? komunikační chyby - odlišná sémantika

3. Skupinová komunikace

3.1 Unicast, broadcast & multicast

Doposud jsme předpokládali, že jeden proces posílá zprávu jinému procesu - tzv. **unicast**. V distribuovaných systémech však často nastanou situace, kdy toto nejjednodušší schéma komunikace nepostacuje. Často je zapotřebí zaslat nějakou zprávu potenciálně všem propojeným počítačům - **broadcast**. V tomto případě vysílající proces nemusí uvádět adresu příjemce - zpráva dojde všem uzlům zapojeným do systému. Způsob komunikace pomocí broadcastu je používán např. při hledání nějakého serveru nebo pro informování všech uzlů o nějaké výjimečné události.

V některých aplikacích je zapotřebí tzv. **multicastu**, který umožňuje příjem pouze vybrané podmnožiny propojených uzlů. Na rozdíl od broadcastu zde vysílající proces je schopen nějakým způsobem identifikovat množinu příjemců. Existují principiálně tři možnosti očekávání odpovědi:

1. **Bez odpovědi** - Vysílající proces neočekává žádnou odpověď, přijímající proces nemusí odpovídat.

2. **Alespon jedna odpověď** - Vysílající proces očekává alespon jednu odpověď. Pro úspěšné zakončení komunikace stačí jedna došlá a potvrzená zpráva.

3. **Všechny odpovědi** - Vysílající proces očekává odpovědi od všech příjemců. Tento způsob komunikace je vhodný pro implementaci spolehlivé skupinové komunikace.

3.2 Pracovní skupiny

Pracovní skupina je množina procesu, které v nějakém systému spolu spolupracují. Klíčová vlastnost skupiny je komunikace - jestliže je nějaká zpráva zaslána skupině jako celku, pak ji dostane každý člen této skupiny. Na rozdíl od jednoduché komunikace, kde spolu komunikují pouze dva procesy, skupinová komunikace je typu one-to-many, tj. jeden proces je vysílající a jeden nebo více procesů jsou příjemci.

Pracovní skupiny jsou dynamické struktury. Mohou být vytvářeny nové skupiny, staré skupiny mohou zanikat, proces může být začleněn do již existující skupiny anebo ze skupiny odstraněn. Proces může být zároveň i ve více skupinách. Důvodem pro vytváření skupin je potřeba zacházet s množinou procesu jako s nedílitelným celkem. Např. libovolný proces může komunikovat se skupinou serveru aniž by znal jejich aktuální počet nebo umístění.

Implementace skupinové komunikace závisí na použitém síťovém hardware. Nejjednodušší způsob je použít multicasting. Bohužel ne všechny síťové hardware tuto službu poskytují. Druhou možností je použít broadcast. Ty uzly, které do dané skupiny nepatří, budou zprávu ignorovat. Nevýhodou řešení pomocí broadcastu je zbytečné zatežování uzlů, které nejsou ve skupině. Nejméně efektivní, ale v případě absence multicastu či broadcastu jedinou možnou, technikou je použití unicastu.

Skupiny mohou být buď uzavřené nebo otevřené. V systémech s uzavřenými skupinami mohou zasílat správy skupině jen ty procesy, které jsou jejími členy. Procesy, které členy skupiny nejsou, jí zasílat zprávy nemohou. V systémech s otevřenými skupinami mohou zprávy skupině zasílat jakékoliv procesy. Uzavřené skupiny bývají většinou používány pro paralelní výpočty. Např. množina procesu řešící šachovou úlohu může být zorganizována jako uzavřená skupina. Otevřené skupiny bývají většinou používány pro replikované nebo jinak spolupracující servery.

Dalším hlediskem, podle kterého mohou být pracovní skupiny rozděleny, je vnitřní organizace skupiny. V některých skupinách si jsou všechny procesy rovnocenné. Každý proces může zasílat zprávy ostatním, rozhodnutí jsou přijímána kolektivně. Jiný typ organizace skupiny je hierarchická organizace. V tomto modelu jsou všechny žádosti směřovány na koordinátora, který rozhodne, který proces se žádosti ujme. Tímto způsobem lze vytvářet hierarchické stromy serveru.

Aby zprávy mohly být skupině skutečně zaslány, je zapotřebí, aby každá skupina byla nějakým způsobem adresována. Jednou metodou adresování je přiřadit každé skupině jednoznačnou adresu. Jestliže použita síť podporuje multicast, pak skupinová adresa může být asociována s multicastovou adresou a každá zpráva skupině je pak tvořena jedním multicastem. V případě implementace pomocí broadcastu každé jádro obdrží zprávu, zjistí z ní adresu skupiny a v případě, že zpráva je určena nějakému lokálnímu procesu, ji přijme. Konečně když neexistuje ani multicast ani broadcast, pak nezbyvá nic jiného, než zaslat unicasty všem členům skupiny. V tomto případě buď musí vysílající jádro znát všechny členy skupiny, anebo být schopno si tyto členy zjistit.

Druhou možností adresování je uvádět přímo v hlavičce zprávy explicitní seznam příjemce. Tento přístup má tu nevýhodu, že vysílající proces musí v každém případě přesně znát všechny členy skupiny a jejich adresy. Dále, pokud se skupina nějakým způsobem změní, pak každý proces zasílající zprávu skupině, musí změnit svůj seznam členů.

Třetí možnou metodou adresování je tzv. **predikátová adresace**. Při použití tohoto způsobu adresace je zpráva adresovaná libovolným výše uvedeným způsobem doplněna o predikát (logický výraz). Jádro po přijetí zprávy vyhodnotí tento výraz a zprávu přijme jen v případě, že výraz má pravdivou hodnotu. Tímto způsobem lze např. zaslat zprávu jen těm počítačům, které mají více než 16MB paměti apod.

3.3 Atomicita

Charakteristikou skupinové komunikace je vlastnost, která by se dala nazvat "**všechno nebo nic**". Většina systému, které podporují skupinovou komunikaci, je implementována tak, že když je skupině zaslána zpráva, pak buď všechny procesy korektně zprávu dostanou anebo zprávu nedostane žádný proces skupiny. Situace, kdy některé procesy zprávu dostanou a jiné ne, není přípustná. Tato vlastnost se nazývá **atomicita**, nebo také **atomický broadcast**.

Atomicita je velmi výhodná pro psaní distribuovaných programů. Když jakýkoliv proces pošle zprávu nějaké skupině, nemusí se již starat o to, kdo zprávu dostal a kdo ne. Například v replikovaném databázovém systému by fakt, že některé procesy neobdržely zprávu o změně stavu dat, mohl způsobit nekonzistenci celé distribuované databáze.

Aby bylo zabezpečeno, že vysílaná zpráva dorazí ke všem příjemcům a ti ji akceptují, je zapotřebí zasílání potvrzení - jádro každého uzlu, který zprávu přijme, potvrdí vysílajícímu jádru příjem. V případě přetížení některých uzlů nebo lokálního výpadku nebo přetížení sítě vysílající jádro zopakuje zprávu. Tato metoda funguje ovšem jen za předpokladu, že žádný z uzlů, ať už vysílající nebo přijímající, nezahavuje.

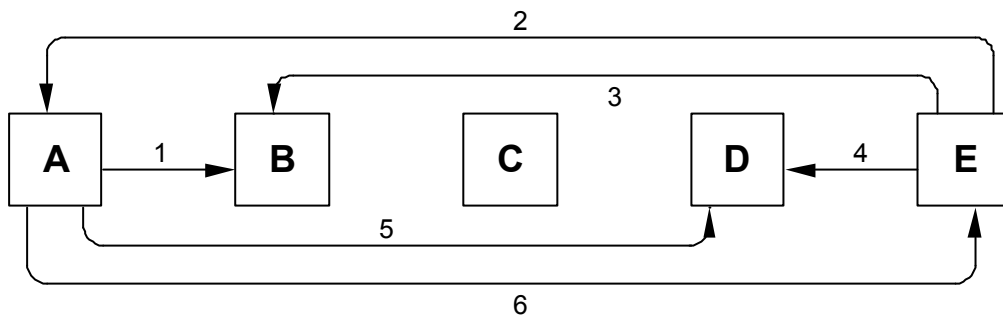
Ve většině distribuovaných systémů je však zapotřebí, aby dodržely podmínku odolnosti proti poruchám (**fault tolerance**), tj. aby atomicita byla zajištěna i v případě výpadku některých uzlů. Výše uvedený mechanismus však tuto podmínku nezaručuje - v případě odmítnutí zprávy příjemcem z důvodu přetížení a následné havárie odesílatele příjemce zprávu nedostane, což je přesně nechtěná situace - někteří členové skupiny zprávu dostali, zatímco jiní nikoliv.

Naštěstí existuje algoritmus, který demonstruje přinejmenším možnost atomického broadcastu. Odesílatel pošle zprávu všem členům skupiny, nastaví si časovač a v případě nutnosti některé zprávy zopakuje. Když příjemce dostane zprávu, kterou ještě nedostal, opět ji rozešle všem členům skupiny (s případným zopakováním některých zpráv). V případě příjmu již přijaté zprávy příjemce nedělá nic a zpráva je ignorována. Tímto způsobem se zabezpečí, že všichni členové skupiny dostanou zprávu

bez ohledu na to, kolik uzlu mezitím zhavaruje a kolik zpráv se po síti ztratí. Existují efektivnější algoritmy pro zabezpečení atomicity, ty budou uvedeny později.

3.4 Uspřádání zpráv

Druhou podmínkou správného fungování skupinové komunikace je, aby zprávy docházely všem členům skupiny ve stejném pořadí. Na obrázku je znázorněna situace, ke které by nemelo dojít. Procesy A a E se zhruba ve stejný okamžik rozhodly poslat zprávu skupině obsahující procesy A, B, D a E. Čísla udávají pořadí, ve kterém jsou jednotlivé zprávy doručovány. Procesu B došla jako první zpráva od procesu A (c. 1) a jako druhá zpráva od procesu E (c. 3). Naopak procesu D došla jako první zpráva od procesu E (c. 4) a jako druhá zpráva od procesu A (c. 5).



Obr. 12 - Konkurenční zasilání zpráv s rozdílným usporádáním příjmu

Nejlepším způsobem je zabezpečit, aby zprávy byly doručeny přesně v tom pořadí, v jakém byly vyslány. Jestliže se tedy proces rozhodne poslat zprávu nějaké skupině, pak zprávy od všech procesů, které se rozhodnou později, budou doručeny až po doručení první zprávy. V tomto případě by všichni příjemci obdrželi zprávy ve stejném pořadí, říkáme, že zprávy splňují **globální časové usporádání**.

Jelikož toto absolutní usporádání není vždy snadné nebo dokonce možné implementovat, většina systémů nabízí slabší variantu - **konzistentní časové usporádání**. Jestliže dva procesy se v přibližně stejném čase rozhodnou zaslat zprávu, systém vybere jednoho z nich jako prvního a doručí jeho zprávu jako první. V tomto případě se může stát, že jako první bude zvolena zpráva, která ve skutečnosti nebyla první, avšak jelikož to stejně nikdo nemůže zjistit, chování systému by na tom nemelo záviset. Je tedy zaručeno, že zprávy budou doručeny všem členům skupiny ve stejném pořadí, avšak toto pořadí nemusí být v některých případech shodné s reálným pořadím, jak byly odeslány.

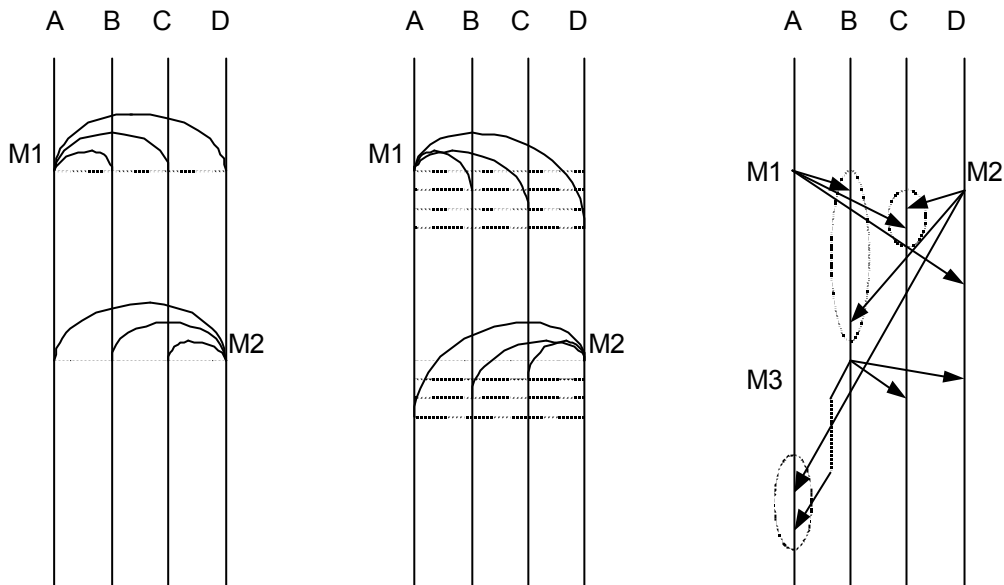
Skupinová komunikace může být založena na pojmech kauzalita a synchronie. **Synchronní systém** je systém, ve kterém se události dějí striktně sekvencně, tj. jakákoliv událost (např. zaslání zprávy) potřebuje ke svému dokončení nulový čas. Zpráva tedy dorazí ke všem příjemcům ve stejný čas, kdy byla odeslána. Z hlediska vnějšího pozorovatele je systém jeví jako množina diskrétních událostí, které se v čase nepřekrývají.

Z technických důvodů je nemožné vytvořit plně synchronní systém, proto existují volnější formy synchronie. **Slabe synchronní systém** je takový systém, ve kterém události trvají konečný předem neznámý čas, avšak všechny události se jeví všem procesům ve stejném pořadí. Slabe synchronní systémy splňují podmínku konzistentního časového usporádání.

V některých systémech může být z důvodu výkonnosti systému vhodnější ještě volnější typy synchronie. **Virtuálně synchronní systém** je systém, v kterém některé zprávy za určitých přesně definovaných podmínek nemusí dojít různým příjemcům ve stejném pořadí. Říkáme, že dvě události jsou **kauzálně vázané**, když chování druhé události může být jakýmkoliv způsobem ovlivněno první událostí. Když proces A zašle zprávu procesu B, ten zprávu přijme a pak pošle zprávu procesu C, je druhá zpráva kauzálně vázaná s první zprávu, protože její obsah může být ovlivněn obsahem první

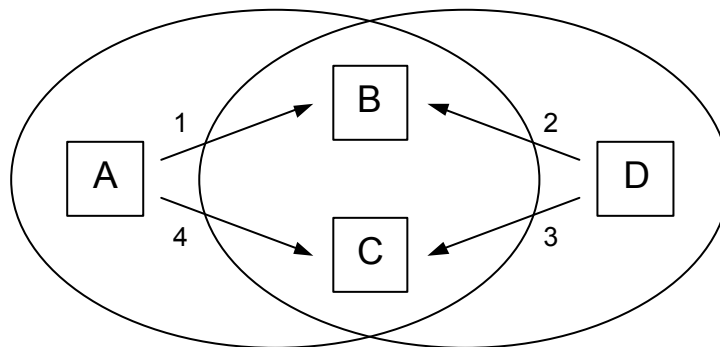
zprávy. Jestli je skutečne druhá zpráva ovlivněna nebo ne je irelevantní, podstatné je, že ovlivněna být mohla. Dve události, které nejsou vázané, se nazývají **konkurentní**. Když proces A zašle zprávu procesu B a zhruba v téže okamžiku proces C zašle zprávu procesu D, pak tyto zprávy jsou konkurentní, neboť se v žádném případě nemohou ovlivnit.

Virtuální synchronie spočívá v tom, že kauzálně vázané zprávy jsou doručeny ve stejném (správném) pořadí. Jestliže jsou zprávy konkurentní, pak jejich pořadí doručení může být pro různé procesy různé.



Obr. 13 - (a) Synchronní systém (b) Slabá synchronie (c) Virtuální synchronie

Jak již bylo zmíněno dříve, jednotlivé procesy mohou být členy více skupin. To může vést k dalšímu druhu nekonzistence. Na obrázku je znázorněna situace, kdy v jedné skupině jsou procesy A, B a C, v druhé skupině jsou procesy B, C a D. Jestliže zprávy dorazí v tom pořadí, jak je znázorněno, pak přestože v rámci jednotlivých skupin jsou zprávy doručeny ve stejném pořadí, procesy B a C dostanou zprávy určené pro různé skupiny v odlišném pořadí.



Obr. 14 - Rozdílné usporádání zpráv v prolínajících se skupinách

Globální i konzistentní časové usporádání doručování zpráv v rámci jedné skupiny nezaručuje žádnou koordinaci mezi jednotlivými skupinami. Některé systémy umožňují dobře definované časové usporádání doručování zpráv mezi prolínajícími se skupinami, zatímco jiné nikoliv.

3.5 Doručovací protokoly

3.5.1 Kauzální závislost

Necht $e_1 \stackrel{p}{\rightarrow} e_2$ je usporádání událostí v rámci procesu p , $\text{send}(m)$ resp. $\text{rcv}(m)$ je odeslání resp. příjem zprávy m . *Kauzální závislost* značená \rightarrow je relace definovaná takto:

1. jestliže $\rightarrow p: e_1 \stackrel{p}{\rightarrow} e_2$, potom $e_1 \rightarrow e_2$
2. $\rightarrow m: \text{send}(m) \rightarrow \text{rcv}(m)$
3. jestliže $e_1 \rightarrow e_2$ & $e_2 \rightarrow e_3$, potom $e_1 \rightarrow e_3$

Definice říká, že událost e_3 kauzálně závisí na e_1 ($e_1 \rightarrow e_3$), jestliže existuje orientovaná cesta (složená ze šipek \rightarrow nebo \rightarrow) od události e_1 k události e_3 .

3.5.2 Kauzální usporádání doručovaných zpráv

Kauzální usporádání doručovaných zpráv definujeme takto:

Necht $\text{dest}(m)$ je množina procesu ve skupině, do které je zaslána zpráva m , $\text{deliver}_p(m)$ je událost doručení zprávy m procesu p . Pak

$$m_1 \rightarrow m_2 \stackrel{p}{\rightarrow} (\text{dest}(m_1) \rightarrow \text{dest}(m_2)) : \text{deliver}_p(m_1) \rightarrow \text{deliver}_p(m_2)$$

Význam definice: Jestliže m_2 je zpráva kauzálně závislá na m_1 , potom ve všech procesech, kterým budou doručeny obe tyto zprávy bude zachováno pořadí tohoto doručení.

3.5.3 Vektorové hodiny

Casová značka (vektorové hodiny) procesu p_i značená $\text{VT}(p_i)$ je vektor délky n , kde n je počet procesu ve skupině. Casová značka zprávy m je značena $\text{VT}(m)$. Pro práci s casovými značkami platí následující pravidla:

1. při startu p_i je $\text{VT}(p_i)$ nulový
2. $\rightarrow \text{send}_{p_i}(m): \text{VT}(p_i)[i]++, \text{VT}(m) = \text{VT}(p_i)$
3. proces p_j , který doručuje zprávu m obsahující $\text{VT}(m)$, upraví $\text{VT}(p_j)$ takto:
 $\rightarrow k \in 1..n: \text{VT}(p_j)[k] = \max(\text{VT}(p_j)[k], \text{VT}(m)[k])$
4. $\rightarrow m_i, m_j (i \neq j): \text{VT}(m_i) \rightarrow \text{VT}(m_j)$

Porovnávání casových značek:

1. $\text{VT}_1 \leq \text{VT}_2 \stackrel{?}{\iff} \forall i: \text{VT}_1[i] \leq \text{VT}_2[i]$
2. $\text{VT}_1 < \text{VT}_2 \stackrel{?}{\iff} \text{VT}_1 \leq \text{VT}_2 \ \& \ \exists i: \text{VT}_1[i] < \text{VT}_2[i]$

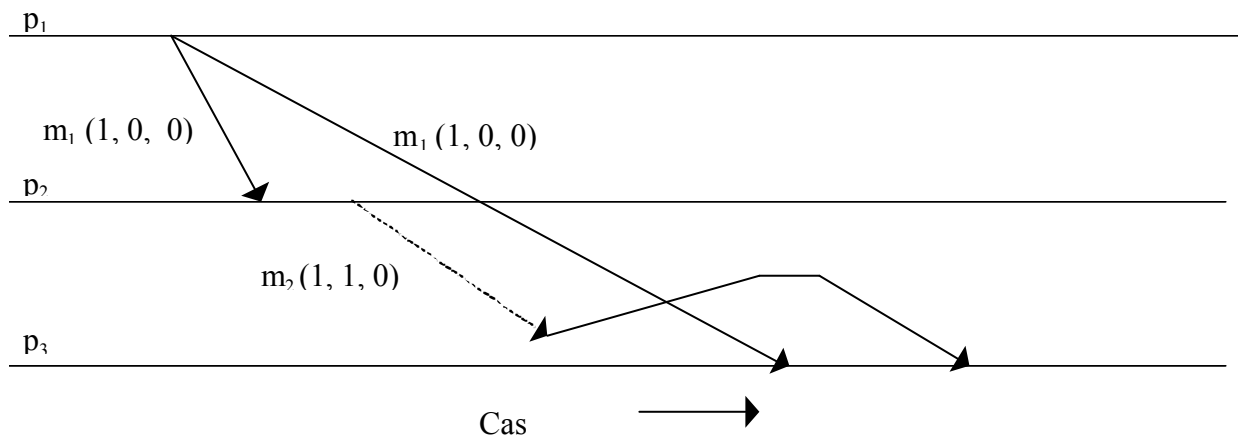
3.5.4 Doručovací protokol pro jednu skupinu

Protokol zaručuje kauzální usporádání doručovaných zpráv. Ke každé vyslané zprávě je přidán aktualizovaný vektor $VT(m)$, který říká, kolik multicastu jednotlivých procesů kauzálně předchází zprávě m .

Protokol:

1. událost odeslání v procesu p_i
 $VT(p_i)[i]++$, označit tímto novým VT odeslanou zprávu
2. přijetí procesem p_j zprávy m vyslané procesem p_i
 proces p_j různý od p_i pozdrží doručení m dokud neplatí:
 - a) $VT(m)[k] = VT(p_i)[k] + 1$ pro $k = i$
 - b) $VT(m)[k] \leq VT(p_j)[k]$ jinak
3. doručení
 po doručení si proces p_j upraví VT podle pravidla 3 pro práci s VT

Bodem 1 protokol zaručuje, že každá odeslaná zpráva bude mít unikátní hodnotu VT. Bodem 2, který je zde klíčový, je zaručeno, že libovolná zpráva kauzálně předcházející zprávě m , bude procesem p_j doručena před zprávou. Bodem 3 je zaručena aktualizace VT.



Obr. 15 - Kauzálně usporádané doručování

Na počátku mají všechny procesy skupiny (p_1, p_2, p_3) nastaven svoji lokální VT na hodnotu $(0, 0, 0)$. Proces p_1 vysílá zprávu m_1 , upraví tedy $VT(p_1)$ na $(1, 0, 0)$ podle bodu 1 protokolu, dále připojí tuto časovou značku k vyslané zprávě a multicastuje zprávu ostatním procesům skupiny.

Proces p_2 obdrží zprávu m_1 a po overení podmínek z bodu 2 protokolu zprávu doručí, přičemž upraví hodnotu lokálního VT(p_2) na $(1, 0, 0)$ (podle pravidla 3 pro práci s VT).

Nyní se proces p_2 rozhodne vyslat skupině zprávu m_2 , což učiní podle bodu 1 protokolu.

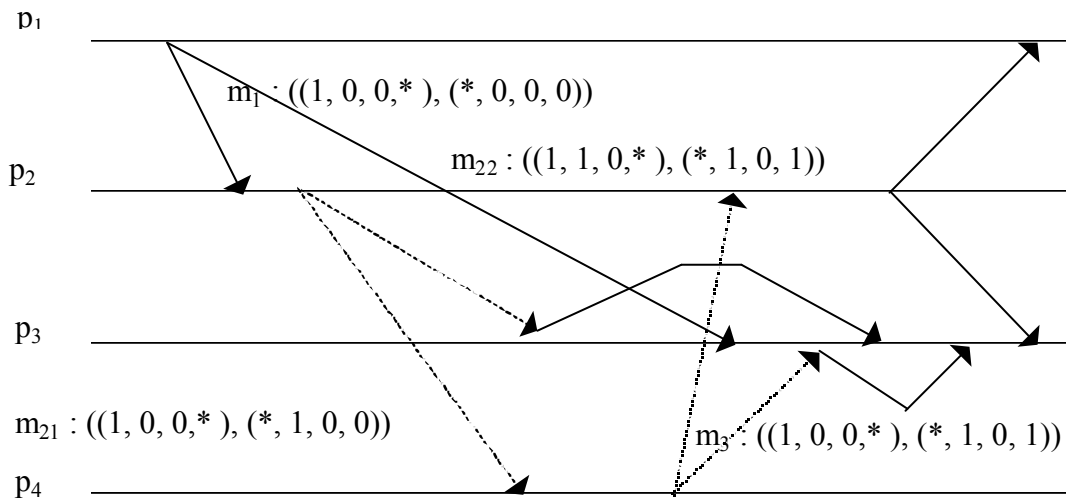
Proces p_3 však nejprve přijme zprávu m_2 , podle bodu 2b) protokolu. Tedy zprávu pozdrží. Mezitím přijme zprávu m_1 , kterou po overení platnosti podmínek bodu 2 doručí a upraví si svoji

lokální VT(p_3) podle pravidla 3 pro práci s VT. Nyní proces p_3 může doručit i zprávu m_2 , protože jsou splněny podmínky bodu 2 protokolu 1.

3.5.5 Doručovací protokol pro překrývající se skupiny

Základem rozšíření doručovacího protokolu pro překrývající se skupiny je rozšíření zasílaných časových značek. Casovou značku spojenou se skupinou g_a značíme VTa . $VTa[i]$ tedy vyjadřuje počet multicastu procesu p_i do skupiny g_a . Se zprávou, kterou proces odesílá, posílá VT všech skupin, kterých je členem.

1. odeslání procesem p_i do skupiny g_a
 $VTa(p_i)[i]++$, označit odesílanou zprávu touto casovou značkou a casovými značkami ostatních skupin a odeslat ji.
2. Přijetí zprávy m procesem p_j (od p_i , VT(m), do g_a)
 Proces p_j ? p_i pozdrží m dokud neplatí:
 - a) $VTa(m)[i] = VTa(p_i)[i] + 1$
 - b) ? k (p_k ? g_a & k ? i): $VTa(m)[k] \leq VT(p_i)[k]$
 - c) ? b (p_j ? g_b): $VTb(m) \leq VTb(p_j)$
3. Doručení
 Po doručení si proces p_j upraví VT podle pravidla 3 pro práci s VT.



Obr. 16 - Kauzálně usporádané doručování v překrývajících se skupinách

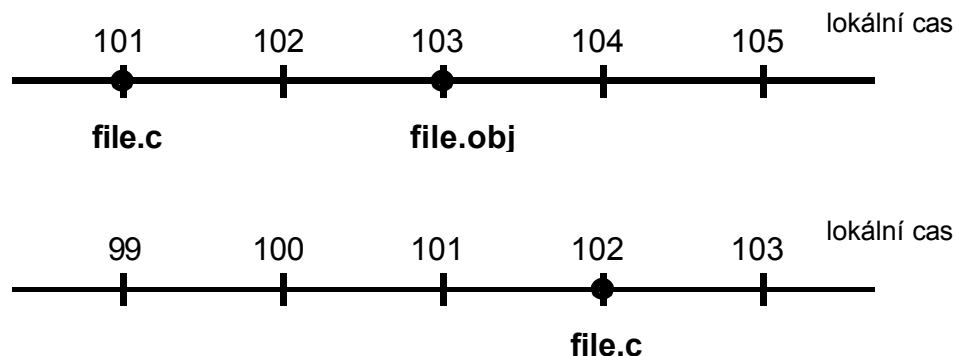
Skupinové komunikace se účastní procesy p_1, p_2, p_3, p_4 , přičemž $p_1..p_3$ jsou členy skupiny g_1 a $p_2..p_4$ jsou členy skupiny g_2 . Proces p_1 pošle zprávu m_1 do skupiny g_1 doplněnou o casovou značku. Proces p_2 zprávu m_1 přijme a následně pošle svůj vlastní multicast m_{21} do skupiny g_2 , ke kterému připojí aktualizovanou casovou značku. Proces p_3 je proto nucen pozdržet doručení m_{21} do doby, než doručí m_1 i přesto, že každý z multicastu směřoval do jiné skupiny. Stejně je pozdržena i zpráva m_3 procesu p_4 , která směřovala do skupiny g_2 . Zpráva m_{22} procesu p_2 je pak jen obyčejná zpráva doplněná o aktualizovanou casovou značku.

4. Synchronizace

V centralizovaných systémech jsou všechny synchronizační problémy, jako kritické oblasti, vzájemné vyloučení procesu apod., obecně řešeny metodami využívajícími semaforey a monitory. Tyto metody nejsou pro použití v distribuovaných systémech příliš vhodné, neboť implicitně předpokládají nějakou formu sdílené paměti. Např. dva procesy interagující pomocí semaforu musí mít k tomuto semaforu přístup. V distribuovaných systémech to z důvodu absence sdílené paměti není možné. Dalšími důvody pro použití zcela nových metod a synchronizačních algoritmů jsou obecné vlastnosti distribuovaných algoritmů:

- informace potřebná pro rozhodování je rozprostřena mezi několika uzly
- procesy se rozhodují na základě lokálních informací
- měly by být vyloučeny takové komponenty, jejichž selhání způsobí havárii systému
- neexistují společné hodiny ani přesné měření globálního času

4.1 Uspřádání událostí a logické hodiny



Obr. 17 - Casový konflikt nesynchronizovaných hodin

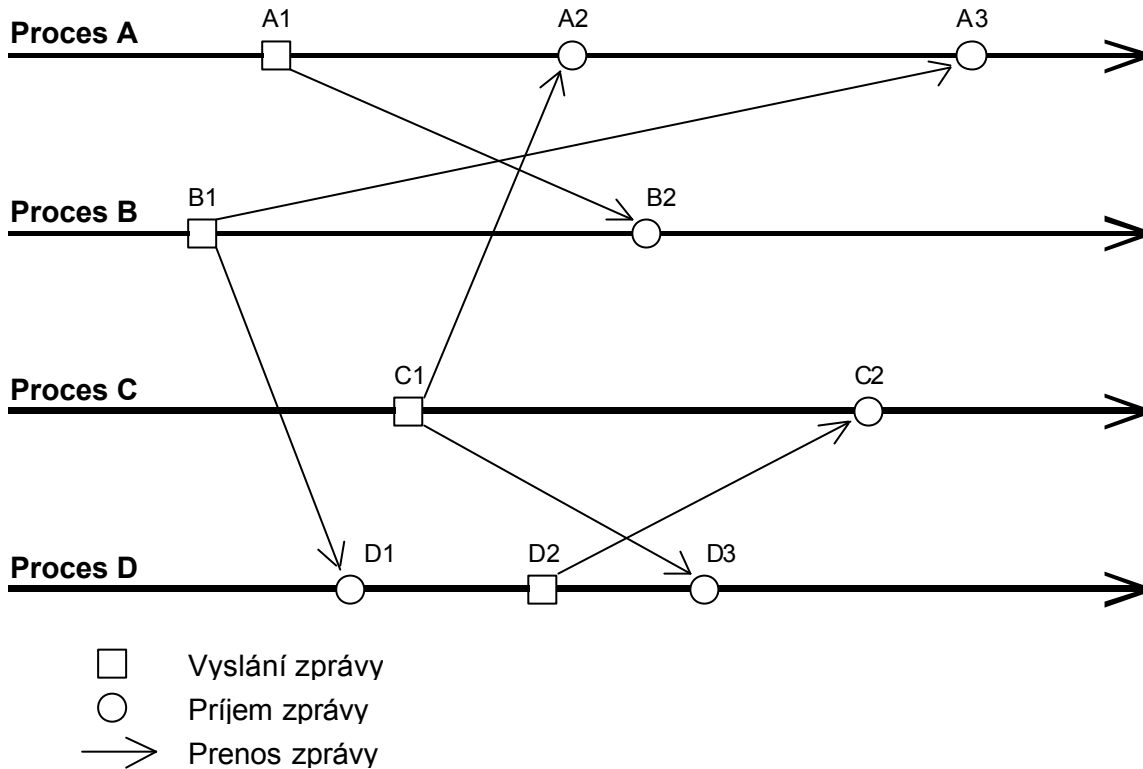
L. Lamport ve svých pracech ukázal, že lze sesynchronizovat všechny hodiny tak, aby byl zaručen jednoznačný časový standard. Své řešení založil na myšlence, že synchronizace nemusí být absolutní. Jestliže dva procesy spolu nijak nekomunikují, pak není nutné, aby jejich hodiny byly synchronizovány, protože případné rozdíly stejně nejsou pozorovatelné a nezpůsobují žádné problémy. Další jeho podstatnou ideou byl fakt, že není důležité, aby se procesy shodly na přesném case, ale aby se shodly na pořadí, v jakém se staly jednotlivé události.

Měření času s těmito vlastnostmi se nazývá logické hodiny. Logické hodiny zajišťují konzistentní měření času v celém systému, nezajišťují však žádným způsobem měření reálného času.

K zajištění synchronizace je definována relace "**predchází**", značená $a \prec b$. Tato relace znamená, že všechny procesy se shodují na tom, že událost a se udála před událostí b . Relace $a \prec b$ je definována těmito axiomy:

- ? 1. Jestliže a a b jsou události v jednom procesu a a se udála před b , pak $a \prec b$
- ? 2. Jestliže a je událost vyslání zprávy a b je událost přijetí této zprávy, pak $a \prec b$
- ? 3. Jestliže $a \prec b$ a $b \prec c$, pak $a \prec c$

Relace "predchází" je ostré částečné usporádání. Jestliže se události x a y staly ve dvou procesech, které spolu nekomunikují, pak neplatí ani $x \prec y$ ani $y \prec x$; říkáme, že události jsou konkurentní, tj. nedá nic říci o tom, která se událost se stala první.

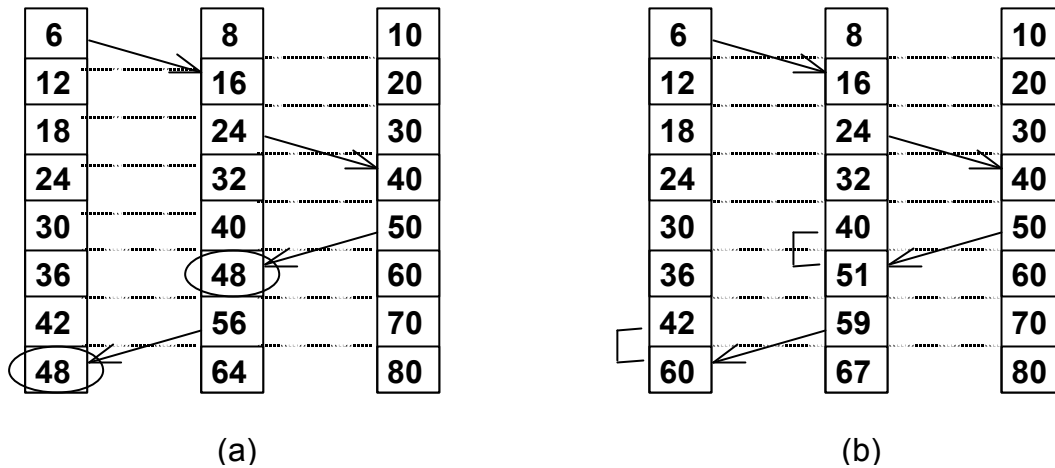


Obr. 18 - Události ve čtyřech konkurentních distribuovaných procesech

Necht čas, ve kterém se stala událost a a na kterém se shodly všechny procesy, je označen $C(a)$, lokální čas procesu i $C_i(a)$. Potom jestliže $a \prec b$, pak $C(a) < C(b)$.

Synchronizace podle přijímání zpráv - časová značka zprávy m T_m :

- Proces i vysílá v case $C_i(a)$ zprávu m ; $T_m = C_i(a)$
- Proces j přijme zprávu m v case $C_j(b)$. Pak $C_j = \max(C_j(b), T_m+1)$



Obr. 19 - (a) Procesy s lokálními hodinami (b) Korekce hodin dle Lamportova algoritmu

V některých situacích bývá nutné, aby se žádné dvě události nestaly ve stejný okamžik. Proto se zavádí relace \prec , která je definována temito axiomy:

Nechť událost a se stala v procesu i , událost b v procesu j a P_n je libovolná jednoznačná ohodnocovací funkce procesu n (napr. interní číslo procesu). Potom

1. Jestliže $C(a) < C(b)$, pak $a \prec b$
2. Jestliže $C(a) = C(b)$ a $P_i < P_j$, pak $a \prec b$

Je zřejmé, že jestli $a \prec b$, pak $a \prec b$. Tedy relace \prec je zúplněním relace \prec .

4.2 Synchronizace fyzických hodin

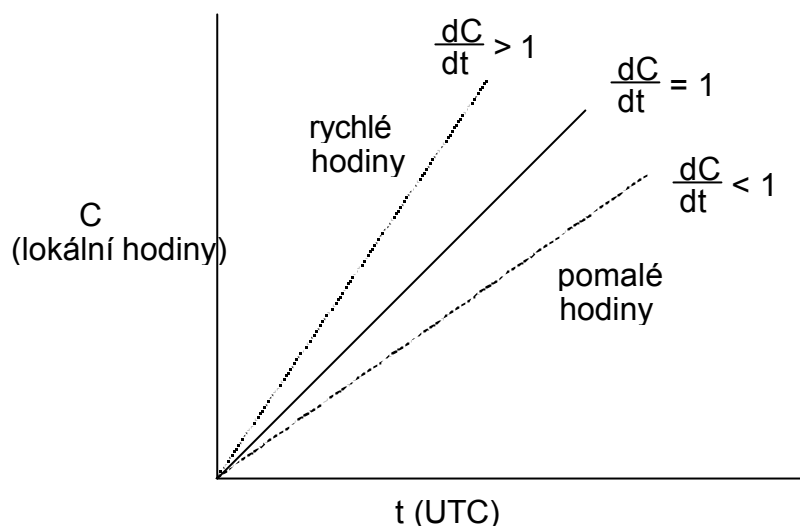
Přestože logické hodiny poskytují konzistentní usporádání událostí, časové hodnoty přiřazené jednotlivým událostem nikterak neodpovídají reálnému času. V systémech, kde je nutné, aby hodiny byly nejen vzájemně synchronizovány, ale aby se nelišily od reálného času, je nutné zavést externí fyzické zdroje času. Z důvodu výkonnosti a spolehlivosti je obvykle použito více časových generátorů, což způsobuje další synchronizační problémy:

1. Jak sesynchronizovat lokální hodiny jednotlivých komponent systému
2. Jak sesynchronizovat jednotlivé hodiny mezi sebou

Na každém počítači běží vnitřní hodiny, jejich hodnotu označme C . Jestliže skutečný čas je t , pak hodnota hodin na počítači p je $C_p(t)$. U zcela přesných hodin by platilo $C_p(t) = t$ pro všechna t , tedy dC/dt by ideálně mělo být 1. Avšak vnitřní hodiny neměří čas zcela přesně; každé hodiny mají od výrobce zaručenou míru přesnosti ϵ , pro kterou platí:

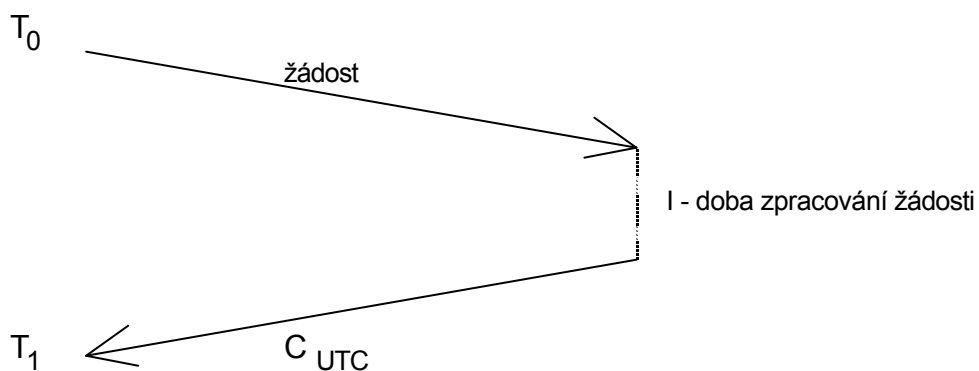
$$1 - \epsilon \leq \frac{dC}{dt} \leq 1 + \epsilon$$

Dvoje hodiny stejného typu se tedy mohou odlišovat od skutečného času maximálně o hodnotu $2\epsilon t$, kde ϵt je doba od poslední synchronizace. Jestliže je tedy zapotřebí, aby se hodiny vzájemně nerozcházely o více než ϵ , pak musí být synchronizovány v intervalech ne větších než $\epsilon/2\epsilon$.



Obr. 20 - Rozdílné rychlosti lokálních hodin

4.2.1 Cristianuv algoritmus



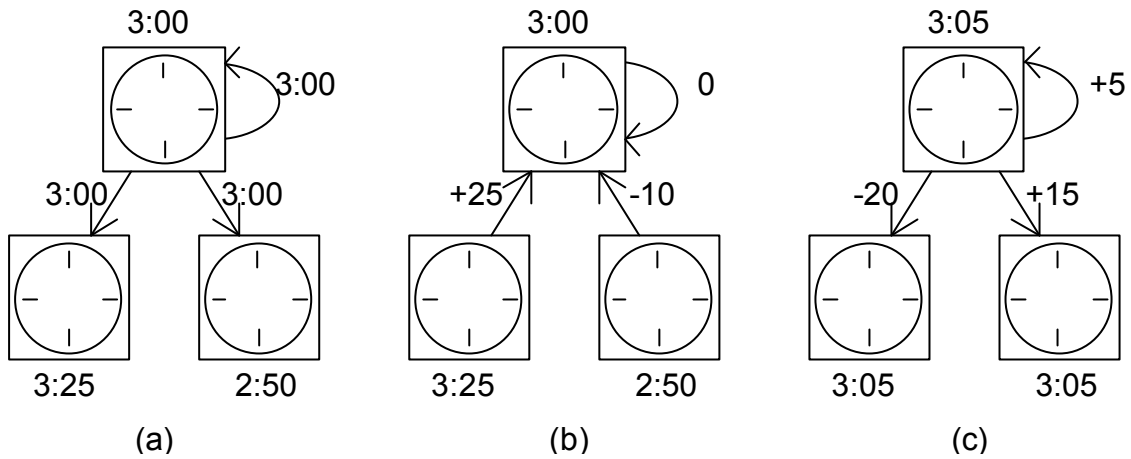
Obr. 21 - Žádost o přesný čas

- jeden pasivní time server
- postupná zmena; zrychlování nebo zpomalování
- $T = T_{UTC} + (T_1 - T_0 - I)/2$

4.2.2 Berkeley algoritmus

(Berkeley UNIX)

- aktivní time server
- periodicky se ptá ostatních na rozdíl času, počítá prumer, vrátí rozdíl



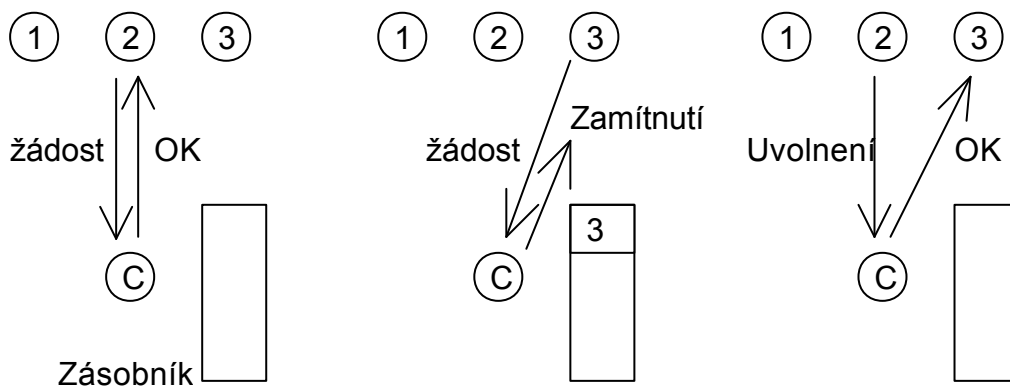
Obr. 22 - (a) Casový server zahájí synchronizaci hodin (b) Odpovědi obsahující časový rozdíl (c) Zprávy s určením opravy lokálních hodin

4.2.3 Distribuovaný algoritmus

- resynchronizací intervaly pevné délky
- i-tý interval: $T_0+iR \dots T_0+(i+1)R$
- broadcast (ve fyzicky nestejný čas), spočítání průměru, případné zahození extrému

4.3 Vzájemné vyloučení procesu

4.3.1 Centralizovaný algoritmus



Obr. 23 - (a) Žádost procesu o vstup do kritické sekce (b) Zamítnutí konkurentního procesu (c) Povolení vstupu do uvolněné kritické sekce

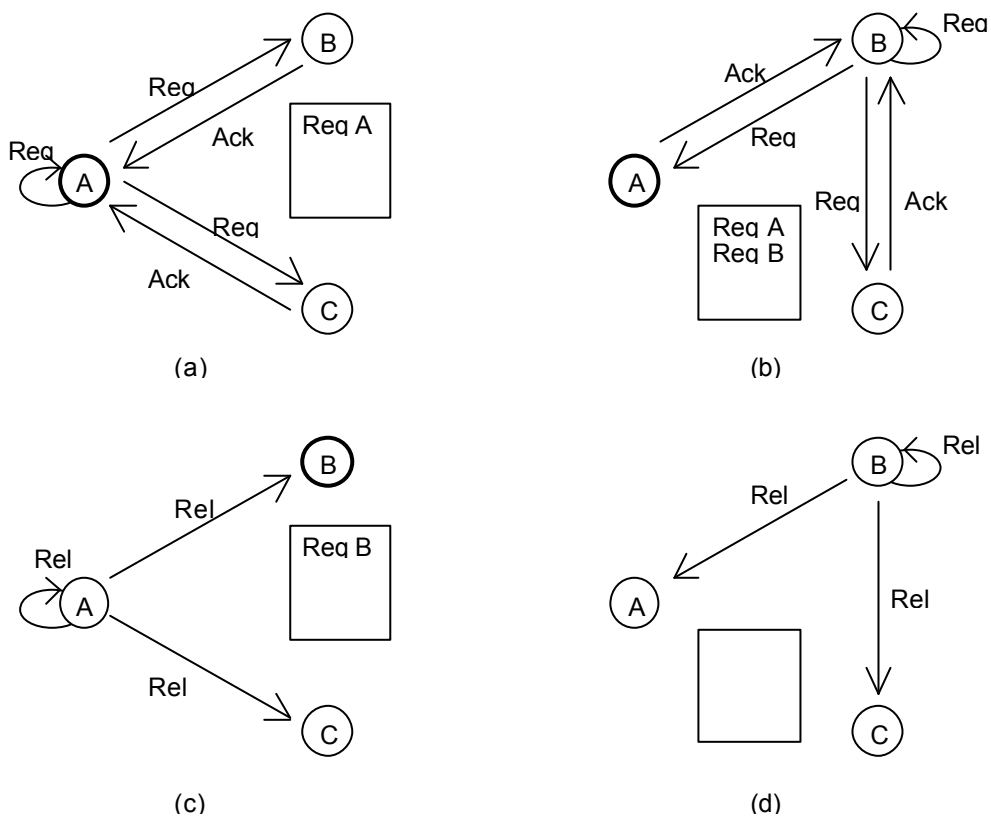
4.3.2 Lamport

Proces p , časová značka odeslání jeho zprávy T_p zprávy typu req, ack, rel, akce se zprávami send, add, del.

Proces vyšle žádost a čeká až dorazí odpovědi od všech procesů a všechny žádosti v jeho frontě mají větší časovou značku.

Událost	Akce
Žádost M_p	send $M_p = \{req, p, T_p\}$
Přijetí žádosti M_i	add M_i ; send $\{ack, p, T_p\}$
Přijetí potvrzení A_i	add A_i
Podmínka vstupu M_p	$? i? p ? \{ack, i, T_i\} : T_p < T_i \ \& \ ? \{req, i? p, T_i\} : T_p < T_i$
Uvolnění R_p	send $\{rel, p, T_p\}$
Přijetí uvolnění R_i	del $\{req, i, T_k\} : T_k < T_i$

Tab. 9 - Reakce na události v Lamportově algoritmu



Obr. 24(a) Vstup procesu A do kritické sekce (b) Žádost procesu B o vstup do kritické sekce (c) Uvolnění kritické sekce procesem A a vstup procesu B do kritické sekce (d) Uvolnění kritické sekce

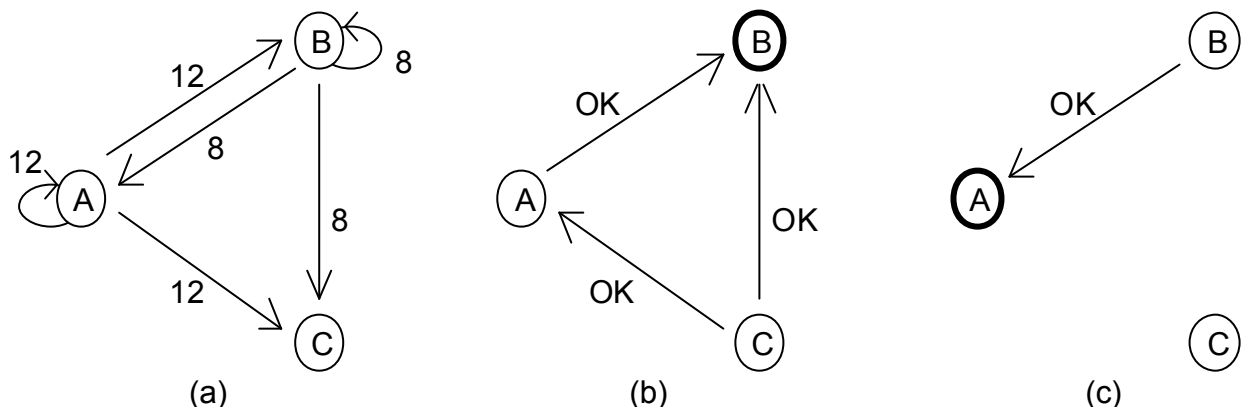
4.3.3 Ricart & Agrawala

Když chce proces vstoupit do kritické sekce, zašle žádost s časovou značkou ostatním procesům a čeká na všechny došlé odpovědi s potvrzením možnosti vstupu.

Když proces přijme zprávu se žádostí, pak:

1. Jestliže příjemce není v kritické sekci a ani do ní nechce vstoupit, pak pošle zpět zprávu s potvrzením.
2. Jestliže je příjemce v kritické sekci, pak neodpovídá, požadavek si zaradí do fronty.
3. Jestliže příjemce ještě není v kritické sekci, avšak chce do ní vstoupit, porovná časovou značku přijaté žádosti s časovou značkou vlastní žádosti. Pokud vlastní žádost má nižší časovou značku, tj. byla vyslána dříve, pak neodpovídá a zaradí žádost odesílatele do fronty. V opačném případě, kdy byla žádost odesílatele odeslána dříve, pošle příjemce zpět zprávu s potvrzením.

Po opuštění kritické sekce proces pošle zprávu s potvrzením všem procesům, které má ve frontě.



Obr. 25 - (a) Procesy A a B ve stejný okamžik žádají o vstup do kritické sekce
 (b) Proces B díky neménší časové značce vstupuje do kritické sekce, proces A čeká
 (c) Po opuštění procesem B může do kritické sekce vstoupit proces A

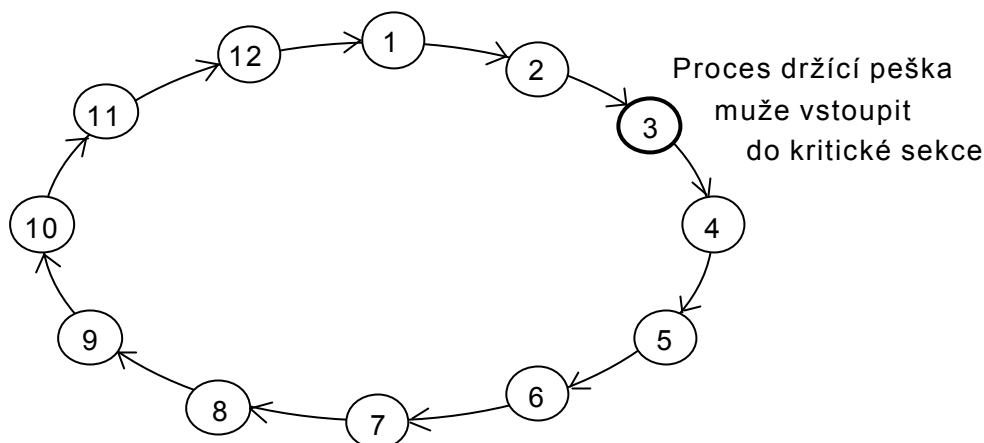
Problémy:

- n points of failure -> ack
- either group communication or selfmaintenance
- n points of overloading

4.3.4 Suzuki & Kasami

Prenos zprávy obsahující povolení ke vstupu do kritické sekce. Zpráva obsahuje frontu požadavků, která se updatuje lokálními došlými požadavky při výstupu z kritické sekce. Komunikační složitost: $2*(N-1)$

4.3.5 Token ring



Obr. 26 - Softwarově konstruovaný token ring

4.3.6 Porovnání jednotlivých algoritmu

Algoritmus	Pocet zpráv	Prodleva	Problémy
Centralizovaný	3	2	Havárie koordinátora
Lamport	$3(n-1)$	$3(n-1)$	Výpadek libovolného procesu
Ricart & Agrawala	$2(n-1)$	$2(n-1)$	Výpadek libovolného procesu
Token ring	1 až ?	0 až $n-1$	Ztráta peška, výpadek procesu

Tab. 10 - Porovnání algoritmu pro vzájemné vyloučení procesu

4.3.7 Vyloučení procesu s prioritami

- Prostředí vyžadující priority
- Real-time systems w/ deadlines

Uspořádaná fronta - podle priorit / podle času, který zbývá na dokončení akce

Předávání peška s frontou na žádost (není cirkulující)

Žádost - broadcast všem

Přijetí žádosti - zarazení do fronty

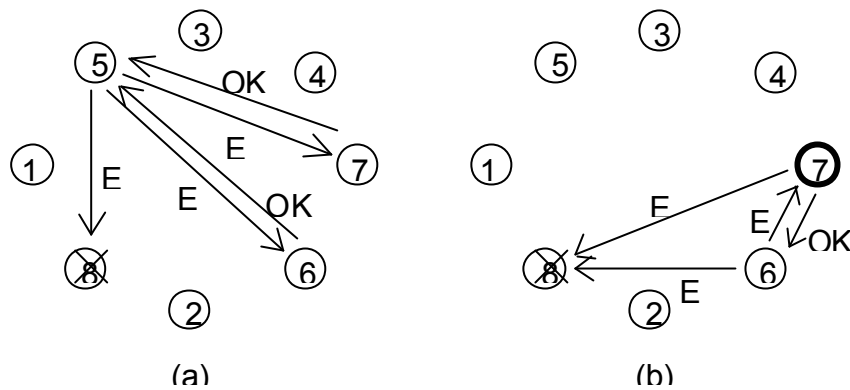
Přijetí peška - vstup do kritické sekce, merge fronty, posláni prvním

4.4 Volba koordinátora

4.4.1 Bully algoritmus

Když se proces rozhodne volit, zašle zprávu všem procesům s vyšší identifikací (číslo procesu apod.). Když přijde odpověď, proces končí, když neprijde nic, proces vyhrál, je novým koordinátorem

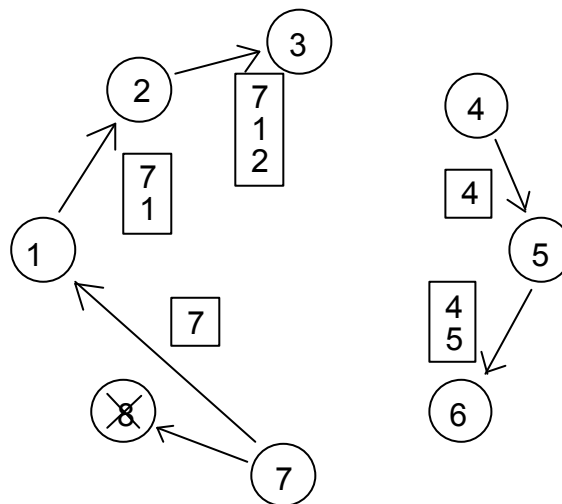
a pošle o tom zprávu všem ostatním. Když proces přijme zprávu o volbě, vrátí zpět odpověď a sám vyšle žádosti všem vyšším procesum. Volba se zabezpečí ve dvou kolech.



Obr. 27 - Bully algoritmus

4.4.2 Kruhový algoritmus

Když se proces rozhodne volit, pošle zprávu následníkovi, zpráva obsahuje čísla procesu (staci odesílatel a nejvyšší živý). Až se zpráva dostane zpět k odesílateli, vyšle dokola zprávu s oznámením koordinátora. Stačí znát následníky a zjistit následníka spadlého uzlu.



Obr. 28 - Kruhový algoritmus

4.4.3 Kruhový algoritmus - $O(n \log n)$

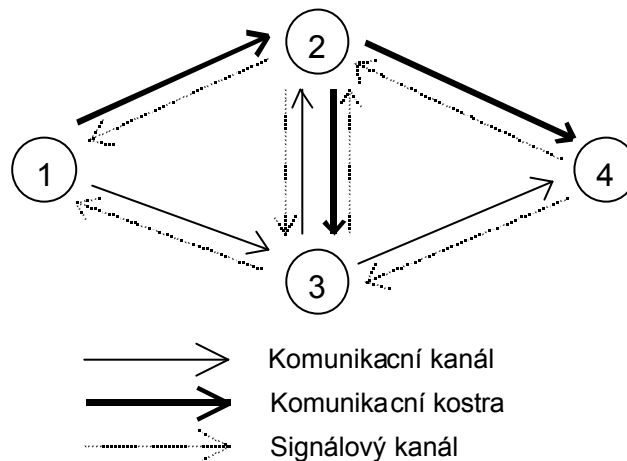
Klasický kruhový algoritmus má v nejhorším případě složitost $O(n^2)$. Chytřejší algoritmus má složitost $O(n \log n)$. Využívá pojem okolí - k -okolí procesoru p je množina procesoru ve vzdálenosti max. k . Velikost této množiny je tedy $2k+1$.

Algoritmus pracuje ve fázích, v n -té fázi se procesor pokouší být koordinátorem svého 2^n -okolí. Pouze procesory, které se stanou lokálními koordinátory, postupují do další fáze.

4.5 Synchronizace ukončení a detekce globálního stavu

Sekvenční programy obvykle končí provedením poslední instrukce. Konkurenční a distribuované procesy jsou však často konstruovány metodou nekonečné smyčky - proces čeká na zprávu, kterou zpracuje. V centralizovaných systémech je jednoduché určit, že procesy již nejsou nikým využívány - fronta připravených úkolů je prázdná. V distribuovaném systému však toto není možné, neboť neexistuje způsob, jak zjistit v jeden časový moment celkový stav systému. Existuje však několik algoritmu, které problém ukončení řeší.

Předpokládejme množinu procesu propojených jednosměrnými komunikačními kanály (obousměrná komunikace může být implementována dvojicí jednosměrných kanálů). Kanály jsou spolehlivé a zachovávají pořadí zpráv. Každý uzel tedy má množinu vstupních a množinu výstupních kanálů (hran). Ke každému komunikačnímu kanálu existuje signální kanál, který vede v opačném směru. Signální kanál neslouží k přenosu dat, ale k zaslání signálu potřebných pro daný algoritmus.



Obr. 29 - Komunikační a signálové kanály

Předpokládejme existenci iniciačního procesu, ze kterého je přes komunikační kanály přístupný každý proces. Iniciační proces tedy nemá žádné vstupní hrany. Výpočet začíná okamžikem, kdy iniciační proces vyšle podél svých výstupních hran zprávu. Každý proces se při příjmu zprávy dostane do stavu *výpočet*. V tomto stavu může dále posílat zprávy podél výstupních hran a může přijímat zprávy ze vstupních uzlů. Jakmile proces výpočet skončí, dostane se do stavu *ukončen*, žádné další zprávy už neposílá, může však zprávy přijímat. Při přijetí zprávy se proces zrestartuje a dostane se opět do stavu *výpočet*.

Algoritmus synchronizace ukončení má za úkol zabezpečit následující: V případě, že všechny procesy jsou ve stavu *ukončen*, pak se v konečném čase tuto skutečnost všechny procesy dozví. Z praktického hlediska postací, když se to dozví iniciační uzel, neboť ten již může tuto zprávu ostatním uzlům poslat.

4.5.1 Dijkstra-Scholten (DS) algoritmus

Jestliže graf procesu je strom, pak každý listový proces při přechodu do stavu *ukončen* pošle signál svému otci. Jakmile proces dostane signály od všech svých synů, pošle signál svému otci. Jestliže všechny signály dostane iniciační proces, je distribuovaný výpočet ukončen.

Tento algoritmus lze rozšířit na acyklický orientovaný graf. Ke každé hraně asociován cíťac - tzv. *deficit*, který udává rozdíl mezi počtem zpráv došlých tímto datovým kanálem a počtem signálu poslaných signálním kanálem v opačném smeru. Koncíí proces vyše každým signálním kanálem tolik signálu, aby deficit byl nulový.

V prípade, že grafem procesu je obecný (orientovaný) graf, pak tento algoritmus nefunguje, nebot neexistují listy, které by mohly samy rozhodnout o ukončení výpoctu. Tento problém se dá vyřešit vytvořením kostry grafu. Behem výpoctu se vytvoří kostra grafu tak, že uzel, od kterého přišla první zpráva, se označí za otce. Algoritmus ukončení pro jeden proces se pak uskuteční ve třech krocích:

1. Poslat signál podél všech vstupních hran krome hrany k otci.
2. Cekat na signál od všech výstupních hran.
3. Poslat signál otci.

Až dojde dostatečný signál korenu celého grafu, lze distribuovaný výpočet ukončit - napr. rozesláním zprávy o ukončení celému grafu.

4.5.2 Znackový (TM) algoritmus

Iniciací uzel vyše žádost o ukončení (znacku, že je kanál prázdný) všem výstupním hranám.

Prijmutí první znacky:

- připraven - propagace znacky výstupním hranám
- nepřípraven - negativní signál (zamítnutí)

Prijmutí další znacky: signál / zamítnutí

Prijmutí zamítnutí: zamítnutí první žádosti

Prijmutí všech potvrzení: signál první žádosti

4.5.3 Detekce globálního stavu

Znackový algoritmus je speciálním případem mnohem obecnějšího algoritmu - zjišťování globálního stavu systému (Chandy & Lamport, Distributed snapshot).

Množina událostí v systému $E = \{e\}$

Rez c je rozdělení E na P_c a F_c : $P_c \cup F_c = E$ & $P_c \cap F_c = \emptyset$

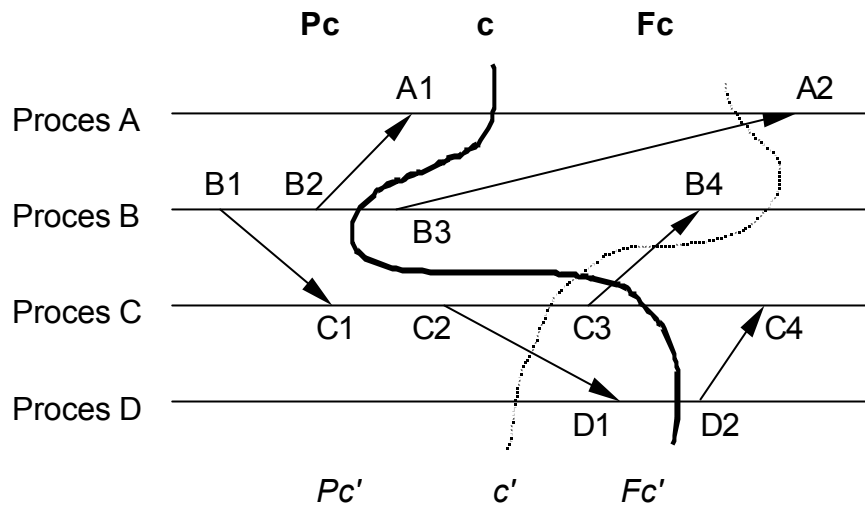
Konzistentní rez c : $a \in P_c \Rightarrow b \in F_c$ & $a \in F_c \Rightarrow b \in P_c$

Stav procesu (systému) je množina událostí, které se v procesu (systému) udály

Konzistentní stav systému $S = P_c$, kde c je konzistentní rez

Necht S je konzistentní stav a e taková událost, že $S' = S \cup e$ je konzistentní stav. Ríkáme, že S' je **dosážitelný** z S , značíme $S \xrightarrow{e} S'$.

Posloupnost událostí $s = (e_1, e_2, \dots, e_n)$ se nazývá **rozvrh (schedule)**, jestliže $S \xrightarrow{e_1} S_1 \xrightarrow{e_2} S_2 \dots S_{n-1} \xrightarrow{e_n} S_n$. Značíme $S \xrightarrow{s} S_n$. Zřejmě platí $S \xrightarrow{s} S_n \Rightarrow S \xrightarrow{s} S_n$.



Obr. 30 - Rez distribuovaného systému - konzistentní (c) a nekonzistentní (c')

Algoritmus

Stav uzlu je množina přijatých a odeslaných zpráv, stav kanálu je množina zpráv, které byly tímto kanálem odeslány, ale ještě nebyly doručeny.

Iniciátor vyšle všem výstupním uzlům značku.

Príjem první značky:

?? Uzel si zapamatuje poslední přijaté a odeslané zprávy = stav uzlu

?? Stav všech příchozích kanálu označí za prázdný

?? Vyšle všem výstupním uzlům značku

Príjem zpráv od uzlu, od kterých ještě nepřišla značka:

?? Zapamatuje si čísla zpráv

Príjem znacek od dalšího uzlu:

?? Stav kanálu = zaznamenané zprávy došlé od uzlu mezi přijutím první značky a značky od tohoto uzlu

Algoritmus končí po přijutí všech znacek. Zaznamenané stavy uzlu a kanálu definují konzistentní stav systému.

5. Transakční zpracování

Všechny doposud zminované synchronizační techniky byly nízkourovňové, vyžadující znalost problematiky semaforu, vzájemného vyloučení, kritických sekcí apod. Pro vývoj distribuovaných aplikací by však byla vhodnější mnohem vyšší úroveň abstrakce, tak, aby se programátor nemusel starat o technické detaily synchronizace. Takovýto mechanismus v distribuovaných systémech existuje, nazývá se **atomické transakce**, nebo také jen krátce transakce.

Pojem transakce pochází ze světa obchodu. Dvě strany nejprve navazují kontakt, vyjednávají, sepisují smlouvy, mění podmínky, a teprve až po oboustranném podepsání je transakce uzavřena. Pokud kdykoliv do tohoto okamžiku libovolná strana z jednání odstoupí, transakce se neuskuteční a vše se vrátí do stavu, jaký byl před započítím transakce. Jakmile je však smlouva podepsána, nelze se vrátit zpět a transakce musí být provedena.

Počítacový model je obdobný. Jeden proces začne transakci, během jejího zpracování mohou být prováděny různé operace. Až se proces rozhodne transakci uzavřít, oznámí to spolupracujícím procesům. Jestliže všechny procesy souhlasí s uzavřením, výsledky transakce se projeví v celém systému, stanou se permanentními. Jestliže alespoň jeden z procesů nesouhlasí s ukončením transakce (nebo zhavaruje před jejím dokončením), stav systému se vrátí přesně do té podoby, jakou měl před započítím transakce a všechny provedené změny budou zrušeny. Tato vlastnost "všechno nebo nic" (all-or-nothing property) významně zjednodušuje vývoj distribuovaných aplikací.

Metodou transakčního zpracování fungovalo počítačové zpracování dat v 60. letech. Na jedné pásece byl zaznamenán nějaký stav, na jinou pásku (nebo pásky) se postupně zaznamenávaly požadované změny. Ve vhodných intervalech (např. jednou denně, každou hodinu apod.) se daly obe pásy zpracovat a počítač na základě stavu jedné pásky a požadovaných změn na druhé pásece vyrobil třetí pásku s novým stavem. Výhodou tohoto způsobu zpracování byla možnost opakování celé akce v případě havárie libovolné části systému.

5.1 Vlastnosti transakcí

Transakce je sekvence operací na jednom nebo více objektu, která transformuje jeden konzistentní stav na jiný konzistentní stav. Transakce je uzavřená mezi značkami **BeginTransaction** a **EndTransaction** a má následující vlastnosti:

- ?? Atomicita - pro vnější svět je transakce nedelitelná
- ?? Konzistence - transakce neporušuje vnitřní konzistenci
- ?? Izolovanost (sekvencnost) - konkurenční transakce se navzájem neovlivňují
- ?? Trvanlivost - po úspěšném ukončení transakce jsou změny trvalé

Tyto vlastnosti se často označují zkratkou ACID (Atomicity, Consistency, Isolation, Durability).

První vlastnost, **atomicita**, zajišťuje, že každá transakce se buď provede kompletně celá, anebo se neprovede vůbec. Dokud se transakce zpracovává, je pro ostatní svět neviditelná. Například pokud by se prováděla transakce přidání deseti bajtů k nějakému souboru, tak po celou dobu zpracování této transakce ostatní procesy vidí původní stav souboru a teprve po úspěšném ukončení transakce vidí soubor i s přidávanými deseti bajty. Žádné mezistavy nejsou pro ostatní procesy viditelné bez ohledu na to, jak dlouho transakce trvala.

Konzistence zajišťuje, že transakce zachovává systém konzistentní. Během zpracování transakce může být systém dočasně nekonzistentní, avšak po ukončení transakce musí být převeden do konzistentního stavu.

Další vlastnost, **izolovanost**, nebo také **sekvencnost**, zabezpečuje, že v případě konkurentního zpracování dvou nebo více transakcí bude konečný výsledek takový, jako kdyby transakce byly prováděny sekvencne v nejakém (systémově závislém) pořadí. Mejdme tři konkurentne vyvolané transakce. Kdyby tyto transakce byly provedeny sekvencne, byl by výsledek 1, 2 nebo 3 podle toho, která transakce by se vykonala jako poslední. V tabulce jsou znázorneny tři možné rozvrhy. Rozvrh 1 je sekvencní. Rozvrh 2 sice není sekvencní, avšak splňuje podmínku sekvencnosti, neboť výsledek (v tomto případě 3) je stejný, jako kdyby byly transakce prováděny sekvencne. Rozvrh 3 nespĺňuje podmínku sekvencnosti, neboť výsledek (v tomto případě 5) neodpovídá žádnému možnému sekvencnímu pořadí transakcí.

<pre>BeginTransaction x = 0; x = x + 1; EndTransaction</pre>	<pre>BeginTransaction x = 0; x = x + 2; EndTransaction</pre>	<pre>BeginTransaction x = 0; x = x + 3; EndTransaction</pre>
--	--	--

Rozvrh 1	x=0	x=x+1	x=0	x=x+2	x=0	x=x+3	OK
Rozvrh 2	x=0	x=0	x=x+1	x=x+2	x=0	x=x+3	OK
Rozvrh 3	x=0	x=0	x=x+1	x=0	x=x+2	x=x+3	ilegální

Tab. 11 - Možné rozvrhy tří transakcí

Trvanlivost znamená, že poté, co byla transakce úspěšně ukončena, jsou její výsledky trvalé. Libovolné pozdější výpadky či jiné havárie nemohou dokončenou transakci zrušit.

5.1.1 Vnorené transakce

Transakce mohou obsahovat **vnorené transakce**. Pokud se vnorená transakce neprovede, neprovede se ani nadřazená transakce. Muže však nastat i opacný případ - vnorená transakce se úspěšně ukončí, avšak nadřazená transakce je později zrušena. V tomto případě musí být vnorená transakce zrušena a to i přes to, že je již úspěšně ukončena. Pravidlo permanence zde neplatí, permanence se vztahuje jen na transakce nejvyšší úrovně.

Vzhledem k tomu, že transakce mohou být vnorovány do libovolné úrovně, je zapotřebí jistá režie. Presto je však sémantika jednoznačná - při startu libovolné (pod)transakce se vytvoří (logická) kopie všech objektu, ke kterým má transakce přístup. Behem zpracování transakce mohou být tyto lokální objekty libovolne meneny. V případě zrušení transakce se zruší i tyto kopie. V případě úspěšného ukončení transakce se nové verze lokálních objektu prenesou do prostoru svého otce. Tím se výsledky transakce stanou viditelné.

5.1.2 Zotavení z chyb

Transakční systém by mel být odolný vůči různým chybám a poruchám. V případě havárie libovolné části systému v okamžiku, kdy transakce způsobila nekonzistenci stavu, je nutné uvést systém do konzistentního stavu. Existují dve možnosti:

1. Návrat do posledního konzistentního stavu, tj. všechny zmeny vyvolané transakcí se zruší. Tento způsob se nazývá **zpetné zotavení** (backward recovery).
2. Nastavení výsledného konzistentního stavu, tj. v případě dostatečného množství informací následně dokončení transakce. Toto se nazývá **dopředné zotavení** (forward recovery). Dopředné zotavení je možné provést jen v případě, že výpadek nastal po okamžiku dokončení transakce, avšak před tím, než mohly být zmeny provedeny.

5.2 Implementace transakcí

Transakční systém musí poskytovat spolehlivé služby v systému složeném z libovolného počtu nezávislých procesů, z nichž libovolný počet může v náhodných okamžicích zhavarovat. Komunikace může být obecně nespolehlivá, zprávy se mohou ztrácet, avšak předpokládá se spolehlivý protokol zabezpečený nižšími vrstvami, který zajišťuje timeouty, potvrzování a případné opakování ztracených paketů.

5.2.1 Stabilní paměť

Paměť můžeme rozdělit podle její spolehlivosti do tří kategorií. První kategorií je obyčejná paměť RAM, která je při výpadku proudu nebo vypnutí počítače vymazána. Druhou kategorií je disková paměť, která není citlivá na výpadky proudu či procesoru, avšak data se mohou ztrácet vinou špatných sektorů, poruch čtecí a záznamové hlavy apod.

Třetí kategorií je stabilní paměť, která je konstruována tak, aby přežila cokoliv kromě kalamit typu zemetření či povodeň. Stabilní paměť může být implementována pomocí dvou sprážených disků. Každý blok na druhém disku je přesnou kopií bloku prvního disku. Když je nějaký blok změněn, provede se změna nejprve na prvním disku a teprve po dokončení a zkontrolování se změna provede na druhém disku.

V případě výpadku proudu jsou oba disky porovnány blok po bloku. Jestliže se bloky vzájemně liší, dá se předpokládat, že k výpadku došlo mezi zápisem na první a druhý disk. V tomto případě se obsah bloku prvního disku překopíruje na druhý disk.

Další potenciální problém je ztráta dat způsobená špatným sektorem nebo poruchou hlavy. Data v bloku, kde nesouhlasí kontrolní součet, mohou být přepsána daty z odpovídajícího bloku druhého disku.

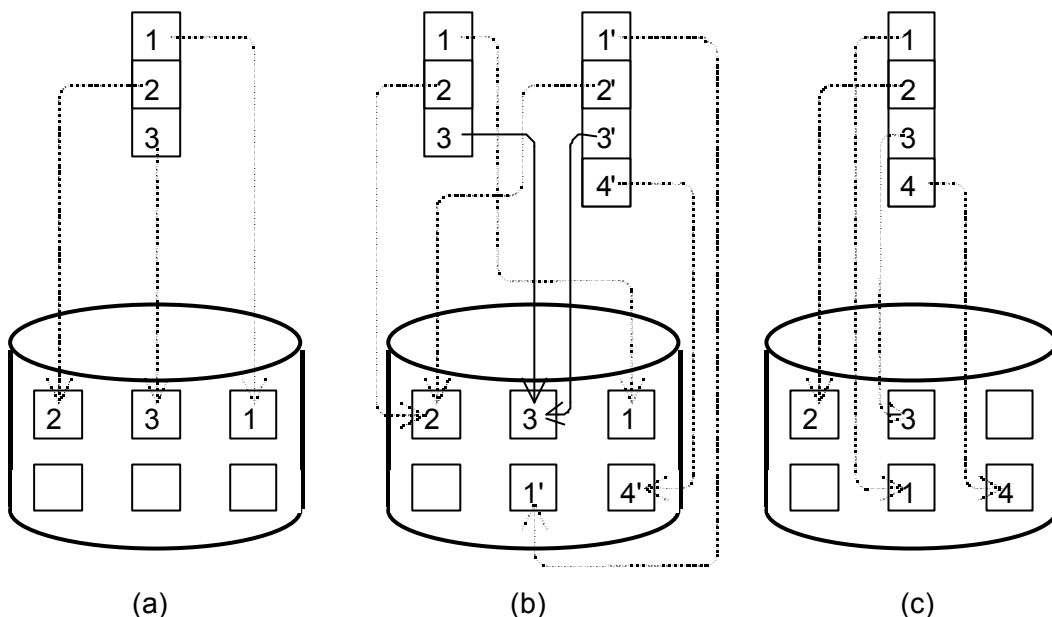
Stabilní paměť je vhodná pro aplikace vyžadující vysoký stupeň spolehlivosti. Pokud jsou již data zapsána, pravděpodobnost, že budou ztracena, je extrémně nízká.

5.2.2 Lokální pracovní prostor

Jednou z metod implementace transakcí je tzv. lokální pracovní prostor. Když proces začne provádět transakci, vytvoří si vlastní pracovní prostor obsahující všechny objekty, kterých se transakce bude týkat. Po dobu trvání transakce se změny provádějí pouze na těchto lokálních objektech. Přímá implementace této metody je však příliš neefektivní. Transakce se často provádějí na velkých databázích. Kopírovat při každé transakci celou databázi je zbytečně náročné na čas i diskový prostor.

První optimalizace je založena na myšlence nekopírovat objekty, které nejsou transakcí meneny. Dokud není takový objekt menen, používá transakce reálnou verzi objektu. Při startu transakce obdrží proces ukazatele na rodičovské objekty, v případě transakce nejvyšší úroveň ukazatele na reálná data. V případě kopie do jiného pracovního prostoru se všechny podržené ukazatele přesměrují na novou kopii.

Druhá optimalizace spočívá v kopírování jen těch částí objektu, které jsou meneny (bloky souboru, databázové záznamy apod.). Při startu transakce se vytvoří pro každý objekt blok indexu (ukazatele), které určují, kde je daná část objektu umístěna. Stejným způsobem je řešeno i přidávání bloku či záznamu k objektu.



Obr. 31 - (a) Index a diskové bloky tříblokového souboru
 (b) Transakce měnící obsah bloku 1 a přidávající blok 4 (c) Situace po ukončení transakce

5.2.3 Intencní seznam

Další používanou metodou implementace je intencní seznam. Server během provádění transakce zaznamenává požadované změny - vytváří si tzv. intencní seznam. Do tohoto seznamu jsou zaznamenávány všechny změny, které požaduje klient provést na objektech, ke kterým má server přístup. Obsahuje identifikaci transakce, která změnu požaduje, identifikaci objektu, který má být změněn a jeho starou a novou hodnotu. Objekt může být změněn teprve až je záznam úspěšně proveden.

V případě, že je transakce úspěšně ukončena, je tato skutečnost zaznamenána do intencního seznamu. Vlastní data zůstávají stejná, neboť již byla během transakce změněna. Jestliže je transakce zrušena, je seznam použit pro obnovení původního stavu. Všechny zaznamenané akce se provedou v obráceném pořadí, čímž se objekt dostane do stejného stavu, v jakém byl před započatím transakce. Tato akce je nazývána **rollback**.

Intencní seznam může být také použit pro obnovu stavu po nějakém výpadku. Jestliže server zhavaroval po zapsání do intencního seznamu, avšak před změnou vlastních dat, může server po nastartování akci dokončit a požadované změny provést.

5.3 Transakční komunikační primitiva

Komunikační schéma používané v transakčním zpracování se někdy nazývá **Commit, Concurrency and Recovery** (potvrzování, konkurence a zotavení, CCR). CCR protokol je typu master/slave, kde master iniciuje transakce a slave je zodpovědný za provedení požadovaných změn. Jelikož slave může být master vzhledem k jiným procesům, může se vytvořit exekucní strom.

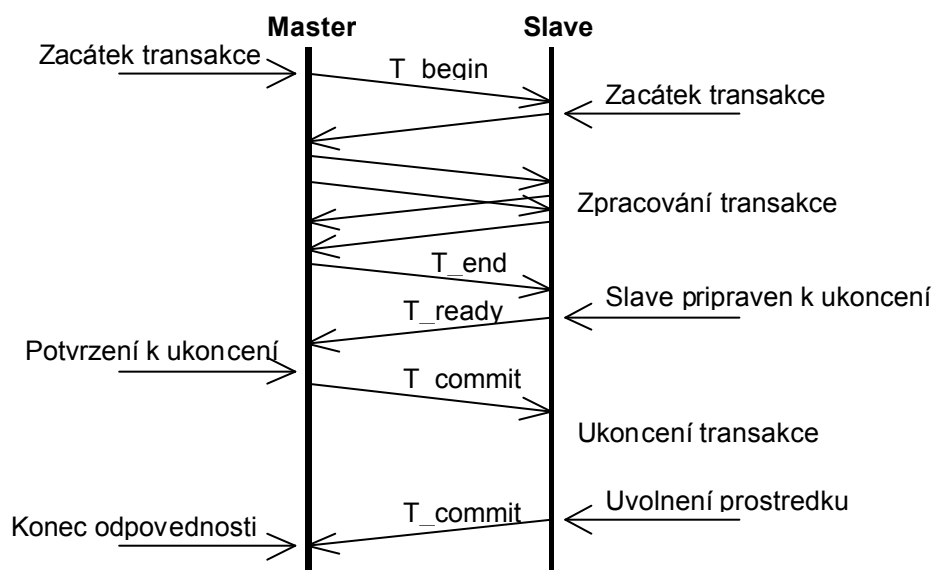
Události, které se vyskytují při komunikaci mezi masterem a slavem:

1. Zacátek transakce - master vyšle příkaz **T_begin**.
2. Master a slave si vzájemně předávají transakčně specifické zprávy. Pokud se kdokoliv rozhodne transakci zrušit, pošle příkaz **T_abort**.

3. Konec transakce (pokud není zřejmý ze zasílaných zpráv) oznámí master zasláním příkazu **T_end**.
4. Slave, v případě, že je připraven transakci dokončit, pošle zprávu **T_ready**.
5. Master pošle příkaz **T_commit**.
6. Slave provede požadované akce a potvrdí ukončení transakce zprávou **T_commit**.

T_begin	M? S	Zacátek transakce
T_end	M? S	Explicitní konec transakce
T_ready	M? S	Připravenost na commit
T_commit	M? S	Commit - potvrzení transakce
T_abort	M? S	Zrušení transakce

Tab. 12 - Základní transakční komunikační primitiva



Obr. 32 - Základní komunikační schéma CCR protokolu (dvoufázový commit)

5.3.1 Dvoufázový commit (potvrzování)

Pro zabezpečení atomicity uzavření transakcí byl vyvinut protokol nazývaný **dvoufázový commit** (dvoufázové potvrzování).

Necht se transakce vyvolané jedním procesem (master) účastní n dalších procesů (slaves). V první fázi master posílá požadavky ostatním procesům, ty podle nich vykonávají požadované akce a případně vracejí odpovědi. Až se master rozhodne transakci ukončit, pošle všem procesům zprávu o tomto záměru a vyčká na odpovědi. Každý slave po obdržení této žádosti zjistí, zda je schopen transakci úspěšně uzavřít. V případě, že ano, uzamkne přístup k modifikovaným objektům a vrátí potvrzení. V případě, že není schopen transakci uzavřít, pošle zpět zprávu oznamující zrušení transakce. Jestliže master obdržel alespoň jednu zápornou odpověď, transakci zruší.

Jestliže master obdržel všechny potvrzující odpovědi, začíná druhá fáze. Vyšle všem slavům commit, ti provedou požadované změny a uvolní objekty. Po ukončení všech potřebných akcí vyšlou potvrzení o ukončení transakce. Poté, co master obdrží odpovědi od všech slavů, je transakce ukončena.

5.4 Kontrola konkurence

Jestliže je více transakcí prováděno současně, je třeba zabezpečit korektní chování transakčního systému. Takovýto algoritmus se nazývá kontrola (řízení) konkurence.

5.4.1 Zámky

Nejstarším a doposud nejrozšířenějším algoritmem je uzamykání (locking). Ve své nejjednodušší podobě, když proces chce číst nebo zapisovat nějaká data, nejprve objekt (soubor) uzamkne, čímž k němu znemožní ostatním procesům přístup. Zámky mohou být implementovány buď pomocí centralizovaného správce zámku anebo distribuovane, kdy každý server přistupující k transakčním objektům, zpracovává své lokální objekty. V obou případech si správce zámku udržuje seznam uzamčených objektů a nepovolí k nim přístup žádnému jinému procesu. Zámky jsou normálně používány transakčním systémem, vzhledem k uživatelům jsou transparentní.

Toto základní schéma je však příliš omezující, existuje několik vylepšení, která zvyšují průchodnost transakčního systému. Zámky mohou být rozděleny na zámky pro čtení a zámky pro zápis. Zámky pro čtení zabezpečují, že po dobu, co je objekt uzamčen, nedojde k jeho změně. Jeden objekt může být uzamčen zároveň několika zámky pro čtení, protože po celou dobu uzamčení nebude změněn, nemůže však být zároveň uzamčen pro zápis. Objekt uzamčený pro zápis nemůže být zároveň uzamčen jakýmkoliv jiným zámkem. Toto schéma je často označováno termínem "many readers / single writer" (více čtenářů / jeden zapisovatel).

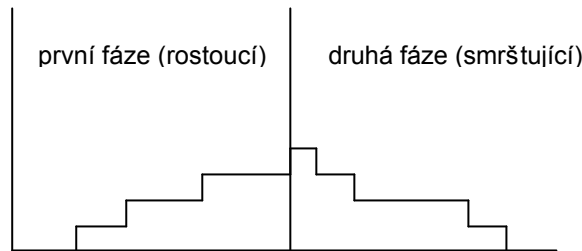
	Read	Write
Read	OK	-
Write	-	-

Tab. 13 - Kompatibilita zámku

Dalším hlediskem, které ovlivňuje celkový výkon systému je **granularita zámku**. Zámky mohou být vztaženy buď na celý soubor nebo databázi (hrubá granulace) anebo na jeden záznam či blok souboru (jemná granulace). Čím je granulace jemnější, tím je zámeček přesnější a může být dosaženo většího stupně paralelismu. Na druhé straně jemně granulované zámky vyžadují větší režii a snadněji dochází k deadlockům.

5.4.2 Dvoufázové uzamykání

Uzamykání a odemykání objektu v okamžiku potřeby může vést k nekonzistencím a deadlockům. Proto je většina transakčních systémů implementována pomocí dvoufázového uzamykání (two-phase locking, 2PL). V dvoufázovém uzamykání jsou v první fázi všechny požadované objekty uzamčeny, poté jsou na nich provedeny požadované operace, a nakonec, v druhé fázi, jsou všechny zámky odstraněny. Druhá fáze se kvůli konzistenci dat často provádí až při ukončení transakce, ať již úspěšně nebo neúspěšně. Lze dokázat, že dvoufázové uzamykané transakce vždy splňují podmínku sekvencnosti.



Obr. 33 - Dvoufázové uzamykání

I dvoufázové uzamykání může vést k deadlockum - když se dvě transakce pokoušejí uzamknout stejný objekt, avšak v opačném pořadí. Obvyklá technika v tomto případě je uzamykání objektu v nějakém jednotném uspořádání. Dalšími možnými metodami odstranění deadlocku jsou explicitní udržování grafu uzamykání a hledání cyklu nebo pomocí timeoutu.

5.4.3 Optimistická kontrola konkurence

Zcela jiný přístup ke kontrole konkurence je tzv. optimistická kontrola konkurence (Kung & Robinson). Idea této techniky je prostá - každá transakce provádí cokoliv bez ohledu na jakékoliv jiné transakce. Pouze v případě, že dojde ke konfliktu, je transakce abortována.

Implementace této metody je nejuvhodnější pomocí lokálního pracovního prostoru. Každá transakce ve svém lokálním prostoru provádí všechny požadované změny. V okamžiku commitu se zkontroluje, zda k některým objektům nepřistupovala zároveň jiná transakce. V případě, že došlo ke konfliktu, se lokální pracovní prostor bez náhrady zruší a transakce se provede znovu.

Velkou výhodou optimistické kontroly konkurence je to, že nedochází k deadlockum. Zároveň v případě řídkých konfliktů umožňuje maximální stupeň paralelismu. Nevýhodou této metody je skutečnost, že v případě většího zatížení a tím pádem větší pravděpodobnosti konfliktu se příliš mnoho transakcí provádí nekolikrát, dokonce při velkém zatížení může dojít k zahlcení systému.

5.4.4 Casové značky

Dalším možným řešením problému konkurenčního přístupu transakcí k společným datům je použití časových značek. Každá transakce při startu dostane svoji časovou značku T . Každému souboru jsou přiřazeny dvě časové značky TR a TW , které určují transakci (resp. její časovou značku), která ze souboru naposledy četla, resp. do něj zapisovala.

Před přístupem k souboru transakce porovná svoji časovou značku s časovými značkami souboru. V případě, že časová značka transakce je větší (tj. transakce je mladší), je vše v pořádku. V opačném případě mohlo dojít ke konfliktu.

	$T < TR$	$T < TW$
ctení	OK, novější transakce četla správná data	abort, data jsou prepsána novější transakcí
zápis	abort, novější transakce četla stará data	OK, ale data nezapisovat, novější transakce už zapsala nová data

Tab. 14 - Kombinace možných konfliktů časových značek pro čtení a zápis

Kombinace možných konfliktů jsou uvedeny v tabulce. Nejzajímavější je situace, kdy starší transakce chce zapisovat do souboru, který je již prepsán novější transakcí. V tomto případě se transakce abortovat nemusí a zápis do souboru se ignoruje, neboť soubor by byl stejně prepsán novější transakcí.

Casové znacky mají oproti uzamykání jednu podstatnou výhodu - nedochází zde k deadlockum. Na druhou stranu mají ponekud odlišné chování - tam, kde při zamykání stací počkat, dochází u časových znacek k abortum.

6. Distribuovaná sdílená paměť (DSM)

Page based / Shared variables / Object based DSM

6.1 Konzistenční modely

Znacení:

W(x)a - zápis hodnoty a do proměnné x

R(x)a - při čtení z proměnné x je vrácena hodnota a

S - synchronizace

Acq - vstup do kritické sekce

Rel - výstup z kritické sekce

P_i - proces i

6.1.1 Striktní konzistence (strict consistency)

Jakékoli čtení z paměti z adresy x vrátí hodnotu uloženou při posledním zápisu na adresu x.

Nejsilnější, na jednoprocessorových systémech je tradičně zajištěna:

Pr. 1:

```
a=1; a=2; print(a); // vytiskne vždy 2
```

Všechny zápisy jsou **okamžitě** všude viditelné, musí existovat přesný globální čas. Ideální pro programování, v distribuovaném systému však nedosažitelné.

Pr. 2:

P1: W(x)1

P1: W(x)1

P2: R(x)1

P2: R(x)0 R(x)1

Striktní konzistence

Paměť, která není striktně konzistentní

6.1.2 Sekvenční konzistence (sequential consistency)

Výsledek jakéhokoliv výpočtu je stejný jako kdyby všechny operace všech procesorů byly vykonávány v nějakém sekvencním uspořádání a operace každého jednotlivého procesoru jsou vykonávány v pořadí specifikovaném programem.

Sekvenční konzistence je o něco slabší model, avšak je snadno implementovatelná a příjemná pro programování. Jestliže procesy běží na různých procesorech, je povoleno libovolné prokládání jejich instrukcí, avšak s podmínkou, že všechny procesy vidí stejné pořadí změn paměti. Změny nejsou propagovány okamžitě, nic se nemluví o velikosti zpoždění (sec, min), je pouze zaručena kauzalita.

Pr. 3:

P1: W(x)1

P1: W(x)1

P2: R(x)1 R(x)1

P2: R(x)0 R(x)1

Dvakrát spuštěný tentýž program

Pr. 4:**P1:**a=1;
print(b,c);**P2:**b=1;
print(a,c);**P3:**c=1;
print(a,b);šest instrukcí - $6! = 720$ možných sekvencí, pouze $3 \cdot (5!/4) = 90$ nenarušuje konzistencia=1;
print(b,c);
b=1;
print(a,c);
c=1;
print(a,b);a=1;
b=1;
print(a,c);
print(b,c);
c=1;
print(a,b);b=1;
c=1;
print(a,b);
print(a,c);
a=1;
print(b,c);b=1;
a=1;
c=1;
print(a,c);
print(b,c);
print(a,b);

Výstup:

001011

101011

010111

111111

Signatura:

001011

101011

110101

111111

Jestliže spojíme výstupy procesu P1, P2 a P3 (v tomto pořadí), dostaneme **signaturu** - řetězec charakterizující konkrétní proložení instrukcí jednotlivých procesů, a tím i pořadí pametových referencí.

Celkem může existovat $2^6 = 64$ možných signatur, avšak ne všechny odpovídají sekvencní konzistenci - např. 000000, 001001 neodpovídají žádnému sekvencnímu proložení instrukcí

Pr. 5:

P1: W(a)1 R(b)0 R(c)0

P2: W(b)1 R(a)1 R(c)0

P3: W(c)1 R(a)0 R(b)1

Signatura 001001 je nemožná v sekvencně konzistentním modelu

Sekvencní konzistence je příjemný model pro programování, avšak jeho výkonnost není příliš velká. Bylo dokázáno (Lipton & Sandberg, 1988) že je-li čas čtení r , čas zápisu w a čas přenosu zprávy (paketu) t , pak vždy platí $r+w \geq t$. Jinak receno - optimalizace protokolu pro čtení znamená delší čas pro zápis a naopak.

6.1.3 Kauzální konzistence (causal consistency)

Zápisy, které jsou potenciálně kauzálně vázané, musí být viděny všemi procesy ve stejném pořadí. Konkurenční zápisy mohou být viděny v různém pořadí.

Rozlišuje události, které jsou potenciálně závislé a které nikoliv. Pokud jeden proces $W(x)$ a a druhý proces $R(x)$ $W(y)$ b pak $W(y)$ b je kauzálně závislý na $W(x)$ a. Operace, které nejsou kauzálně závislé jsou konkurenční.

Pr. 6:

P1:	$W(x)1$		$W(x)3$	
P2:		$R(x)1$	$W(x)2$	
P3:		$R(x)1$		$R(x)3$ $R(x)2$
P4:		$R(x)1$		$R(x)2$ $R(x)3$

Kauzální konzistence

Příklad 6 zobrazuje situaci, kdy $W(x)2$ a $W(x)3$ jsou konkurenční (tj. nejsou kauzálně závislé). Procesy P3 a P4 mohou potom změnu x číst v opačném pořadí. Tato sekvence splňuje podmínku kauzální konzistence, nespĺňuje sekvenci a striktní.

Pr. 7:

P1:	$W(x)1$		P1:	$W(x)1$
P2:		$R(x)1$ $W(x)2$	P2:	$W(x)2$
P3:		$R(x)2$ $R(x)1$	P3:	$R(x)2$ $R(x)1$
P4:		$R(x)1$ $R(x)2$	P4:	$R(x)1$ $R(x)2$

PRAM konzistence

Kauzálně konzistentní sekvence

Na druhém podobrázku příkladu 7 není $W(x)2$ kauzálně závislý na $W(x)1$ a proto obrácené čtení $R(x)1$ a $R(x)2$ je vzhledem ke kauzální konzistenci korektní. Na prvním podobrázku je $W(x)2$ kauzálně závislý na $W(x)1$ (pres $R(x)1$) a proto obrácené pořadí čtení porušuje kauzální konzistenci.

Implementace kauzálně konzistentní paměti vyžaduje udržování grafu závislostí zápisu na čtení.

6.1.4 PRAM konzistence (PRAM consistency)

Zápisy prováděné jedním procesem jsou viděny ostatními procesy v tom pořadí, ve kterém byly prováděny, avšak zápisy různých procesů mohou být viděny různými procesy různě.

PRAM (Pipelined RAM) konzistence je snadná na implementaci - nezáleží na pořadí, v němž různé procesy vidí přístupy k paměti. Jediné, co je třeba dodržet, je pořadí zápisu z jednoho zdroje. Jinak receno - všechny zápisy generované různými procesory jsou konkurenční.

Pr. 8:

P1: a=1; if(b==0) kill(P2);	P2: b=1; if(a==0) kill(P1);
--	--

V případě PRAM konzistence je možný výsledek, že oba procesy budou zabity, a to v případě, že P1 přečte b dříve, než uvidí jeho zápis a P2 přečte a dříve, než uvidí jeho zápis. Tento výsledek není možný při libovolném uspořádání instrukcí, neodpovídá tedy sekvencí konzistenci, avšak vzhledem k PRAM konzistenci je korektní.

6.1.5 Slabá konzistence (weak consistency)

Výše uvedené konzistenční modely jsou stále pro mnoho aplikací velmi restriktivní (a tedy málo efektivní), neboť vyžadují propagaci všech zápisu všem procesum. Ne všechny aplikace vyžadují sledování všech zápisu, natož pak nějaké jejich poradí. Např. proces v kritické sekci může v rychlé smyčce číst a zapisovat hodnoty nějakých promenných. Ostatní procesy nemají důvod (nebo ani možnost) jednotlivé zápisy vidět, proto není nutné, aby byly všechny propagovány. Paměť však nemá možnost vedet, že nějaký proces je v kritické sekci a tudíž musí propagovat všechny zápisy.

Řešením je nechat proces ukončit kritickou sekci a poté zajistit rozeslání změn všem ostatním procesum. Toho může být dosaženo použitím speciálního druhu promenné - **synchronizační promenné**, která bude použita pro synchronizační účely.

1. **Přístup k synchronizačním promenným je sekvencně konzistentní.**
2. **Přístup k synchronizační promenné není povolen dokud neskončí všechny předchozí zápisy.**
3. **Přístup k datům není povolen dokud nebyly dokončeny všechny předchozí přístupy k synchronizačním promenným.**

Bod 1 říká, že všechny procesy vidí všechny přístupy k synchronizační promenné ve stejném pořadí. Bod 2 zaručuje, že před přístupem k synchronizační promenné budou dokončeny všechny předchozí zápisy. Bod 3 říká, že při přístupu k obyčejným promenným jsou dokončeny všechny předchozí přístupy k synchronizačním promenným. Provedením synchronizace před čtením dat proto zajistí jejich aktuální verzi.

Pr. 9:

P1: W(x)1 W(x)2 S	P1: W(x)1 W(x)2 S
P2: R(x)2 R(x)1 S	P2: S R(x)1
P3: R(x)1 R(x)2 S	
Slabá konzistence	Porušení slabé konzistence

Na prvním podobrázku příkladu 9 jsou události odpovídající slabé konzistenci - procesy P2 a P3 mohou před synchronizací vidět paměť v libovolném pořadí. Druhý podobrázek ukazuje porušení slabé konzistence - proces P2 po synchronizaci vidí neaktuální hodnotu promenné x.

6.1.6 Uvolňovací konzistence (release consistency)

Slabá konzistence má tu nevýhodu, že paměť při přístupu k synchronizační proměnné nepoznává, zda se jedná o vstup nebo o výstup z kritické sekce. Pokaždé proto musí vykonat akce potřebné pro oba případy. K rozlišení těchto případů používá uvolňovací konzistence dvě operace - acquire (požadavek) a release (uvolnění).

1. **Před běžným přístupem ke sdílené proměnné musí být úspěšně ukončeny předchozí požadavky procesu.**
2. **Před provedením uvolnění musí být ukončeny všechny předchozí zápisy i čtení prováděné procesem.**
3. **Požadavky a uvolnění musí být PRAM konzistentní.**

Jestliže proces vydá požadavek, pak je zaručeno, že poté jsou všechny lokální kopie aktuální. Po uvolnění jsou propagovány všechny změny ostatním procesům. Jestliže jsou všechny podmínky splněny a požadavky a uvolnění se používají správe spárované, výsledek jakéhokoliv výpočtu je ekvivalentní sekvencí konzistentní paměti.

Pr. 10:

P1: Acq W(x)1 W(x)2 Rel

P2: Acq R(x)2 Rel

P3: R(x)1

Uvolňovací konzistence

Eager release consistency - po release se propagují změny všem procesům.

Lazy release consistency - po release se nic nepropaguje, až při acquire jiného procesu - časově znacky.

6.1.7 Přístupová konzistence (entry consistency)

1. **Požadavkový přístup procesu k synchronizační proměnné není povolen dokud nebyly provedeny všechny aktualizace chráněných sdílených dat procesu.**
2. **Exkluzivní přístup procesu k synchronizační proměnné je povolen pouze v případě, že žádný jiný proces nepřistupuje k synchronizační proměnné, a to ani neexkluzivně.**
3. **Po exkluzivním přístupu k synchronizační proměnné si příští neexkluzivní přístup libovolného procesu k synchronizační proměnné musí vyžádat aktuální kopii dat od vlastníka synchronizační proměnné.**

Všechna sdílená data jsou vázána na synchronizační proměnnou. Při přístupu k synchronizační proměnné jsou synchronizována pouze ta data, která jsou vázána na synchronizační proměnnou.

Každá synchronizační proměnná má v každém okamžiku vlastníka - proces, který k ní naposledy přistupoval. Vlastník může opakovaně vstupovat do a vystupovat z kritické sekce bez nutnosti komunikace. Proces, který není vlastníkem, musí požádat o vlastnictví. Při přenosu vlastnictví se aktualizují hodnoty dat vázaných na synchronizační proměnnou. Přístup k datům a synchronizačním proměnným může být exkluzivní (read & write) nebo neexkluzivní (read only).

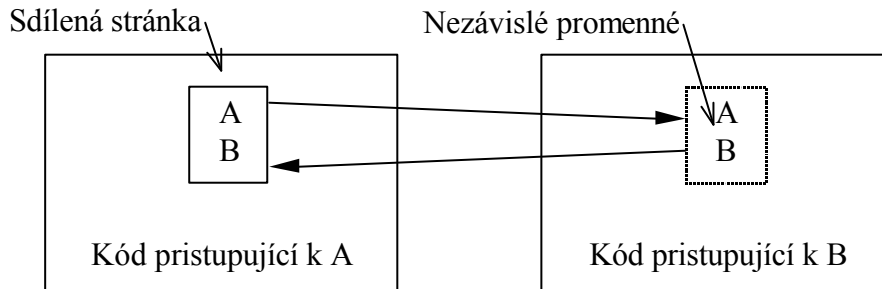
6.1.8 Shrnutí konzistencních modelu

	Konzistence	Vlastnosti
a	Striktní	Absolutní časové usporádání
	Sekvenční	Všechny události jsou vidět ve stejném pořadí
	Kauzální	Kauzálně vázané události jsou vidět ve stejném pořadí
	PRAM	Události jednoho procesu jsou vidět ve stejném pořadí
b	Slabá	Sdílená data jsou konzistentní po synchronizaci
	Uvolnovací	Sdílená data jsou konzistentní po opuštění kritické sekce
	Prístupová	Sdílená data vázaná na kritickou sekci jsou konzistentní při vstupu do kritické sekce

Tab. 15 - Přehled základních konzistencních modelu
(a) bez synchronizačních operací (b) se synchronizačními operacemi

6.2 Distribuované stránkování

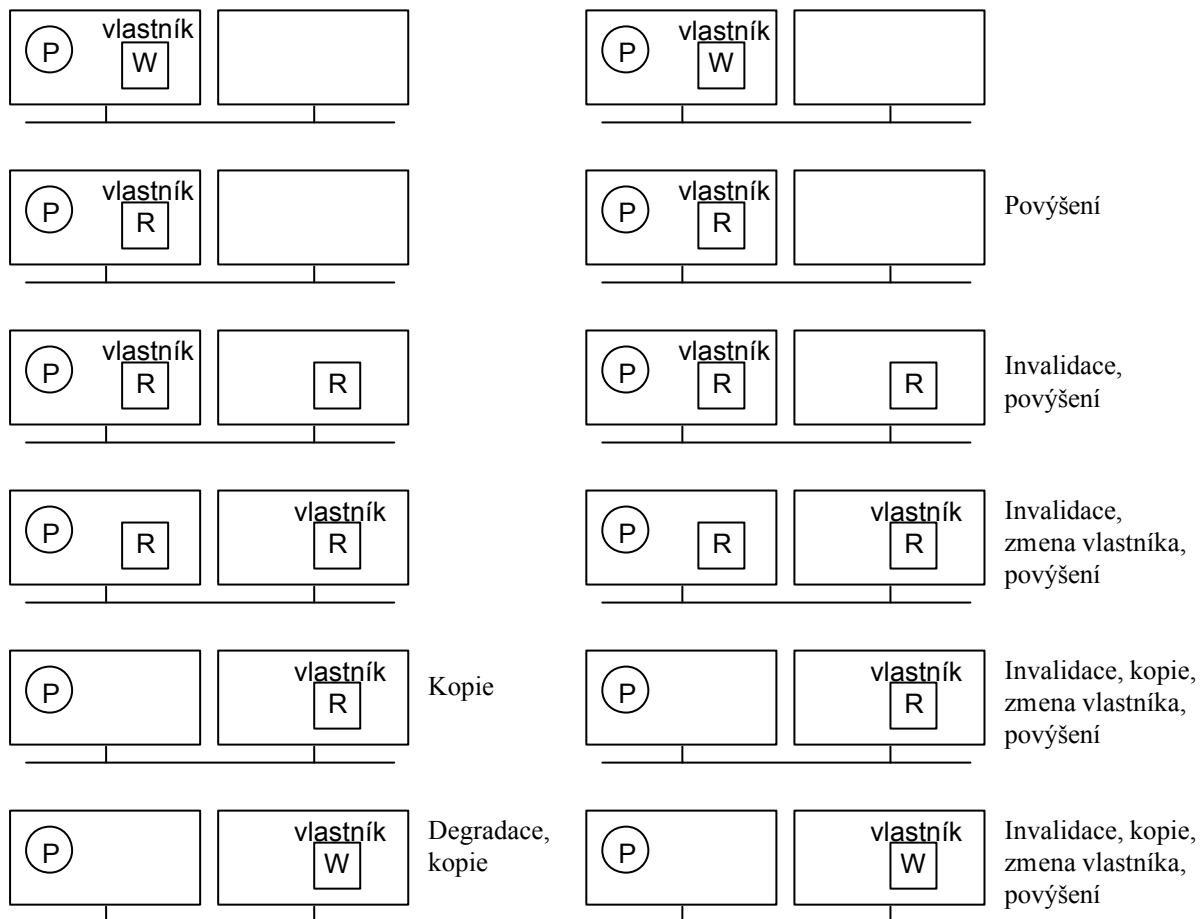
Základní schéma distribuovaného stránkování je obdobou virtuální paměti - je-li referencována stránka, která není na lokálním stroji, je vyvoláno přerušení a stránka musí být nactena ze stroje, kde se právě nalézá. Problémy: replikace -> konzistence, nalezení stránky, správa kopií, falešné sdílení.



Obr. 34 - Falešně sdílená stránka obsahující nezávislé promenné A a B

Jestliže jsou stránky (predevším z důvodu výkonnosti) replikovány, je třeba udržovat konzistenci dat. Typická implementace je namapování všech stránek read-only a v případě zápisu provádět potřebné synchronizační akce. Teoreticky jsou možné dva způsoby - invalidace nebo aktualizace. Distribuované stránkování v naprosté většině používá invalidaci - vzhledem k velikosti přenášených dat a době přenosu.

Na obrázku jsou rozebrány možnosti umístění a vlastnictví stránky při čtení a zápisu a potřebné synchronizační akce.



Obr. 35 - (a) čtení ze stránky (b) zápis do stránky

6.2.1 Nalezení vlastníka stránky:

- broadcast
- centralizovaný manager
- replikovaný manager - indexovaný spodními n bity adresy stránky / hash

6.2.2 Nalezení kopií:

- broadcast
- vlastní stránky udržuje copyset - množinu lokací kopií

6.2.3 Alokační politiky stránek:

1. not owned read copy
2. owned replicated copy - transfer of ownership
3. local policy (LRU)

6.3 Distribuované sdílené promenné

Implementace na úrovni knihoven - potenciálně replikovaná distribuovaná databáze.

Výhody: potenciálně lepší výkonnost, eliminace falešného sdílení.

Nevýhody: Nepodporováno přímo operačním systémem, nutnost implementace pro různé jazyky, nutnost rekompile.

6.3.1 Munin

ordinary / shared / synchronisation variables

Eager release consistency - při opuštění kritické sekce se propagují změny.

sdílené promenné: read-only, migratory, write-shared, conventional

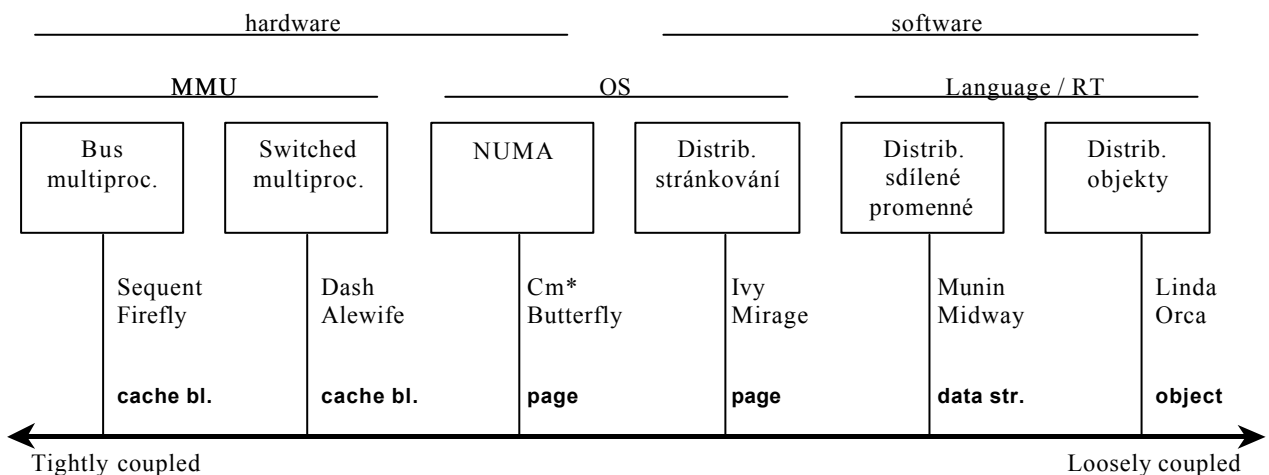
- read-only: když dojde k výpadku, v adresáři promenných je nalezen vlastník, který je požádán o read-only kopii dat.
- migratory: acquire/release protokol implementující release consistency. Data chráněná synchronizačními přístupy migrují na stroj, který je zrovna v kritické sekci.
- write-shared: všechny stránky jsou iniciálně read-only, při zápisu stránka vytvoří kopii s původním obsahem a změní se na rw a označí se dirty. Po release se porovná stránka s původní a změny se propagují. Když přijde propagace změny na ne-dirty stránku, změny se akceptují, v opačném případě se word-po-wordu zkontroluje, zda nedošlo ke konfliktu. Když ne, data se sjednotí, když ano, dojde k runtime-erroru.
- konvenční sdílená data (nepatřící do žádné z výše uvedených kategorií): chovají se jako distribuované stránkování - single writer / many readers a sekvencí konzistence.

6.4 Distribuované objekty

Díky enkapsulaci flexibilnější - komunikace a synchronizace v metodách

Class Definition Language - automatické generování hlaviček a kódu

Základní "distribuovaná" třída, dedení vlastností, CORBA



Obr. 36 – Mechanismy sdílení paměti

7. Identifikace objektu

Komunikace mezi procesy a využívání vzdálených objektu vyžaduje, aby jádro systému zabezpečilo jednotné pojmenování - identifikaci objektu. Objektem je zde míněna libovolná entita, ke které má být zajištěn přístup. Objekty mohou být aktivní (napr. procesy) nebo pasivní (napr. tiskárny). Pasivní objekt bývá spravován nějakým aktivním objektem.

K tomu, aby nějaký objekt mohl být přístupný, je zapotřebí vedet kteřý objekt má být použit (identifikace), kde je tento objekt umístěn (adresa) a jak je možné se k nemu dostat (cesta). Identifikační systém musí být navržen efektivně, aby prodleva a počet zpráv mezi započítáním komunikace a vlastním přenosem zpráv byly co nejmenší. Vlastní komunikace může začít teprve poté, co je z jména nějakým způsobem odvozena cesta k danému objektu. Adresa je často chápána jako mezistav potřebný k určení požadované cesty.

7.1 Identifikační systém

Identifikační systém má za úkol přidělovat jména, mapovat jména na adresy a starat se o nalezení nejefektivnější cesty. Měl by podporovat alespoň dvě úrovně jmen - systémové a uživatelské identifikátory. Systémové identifikátory objektu typu 0C03A489F20 jsou obvykle pro běžné uživatele necitelné a naopak, uživatelský identifikátor typu 'Toto je dopis mému šéfovi' je pro běžnou práci systému příliš nepohodlný. Systém tedy musí sám vhodným způsobem převádět uživatelská jména na systémová a naopak.

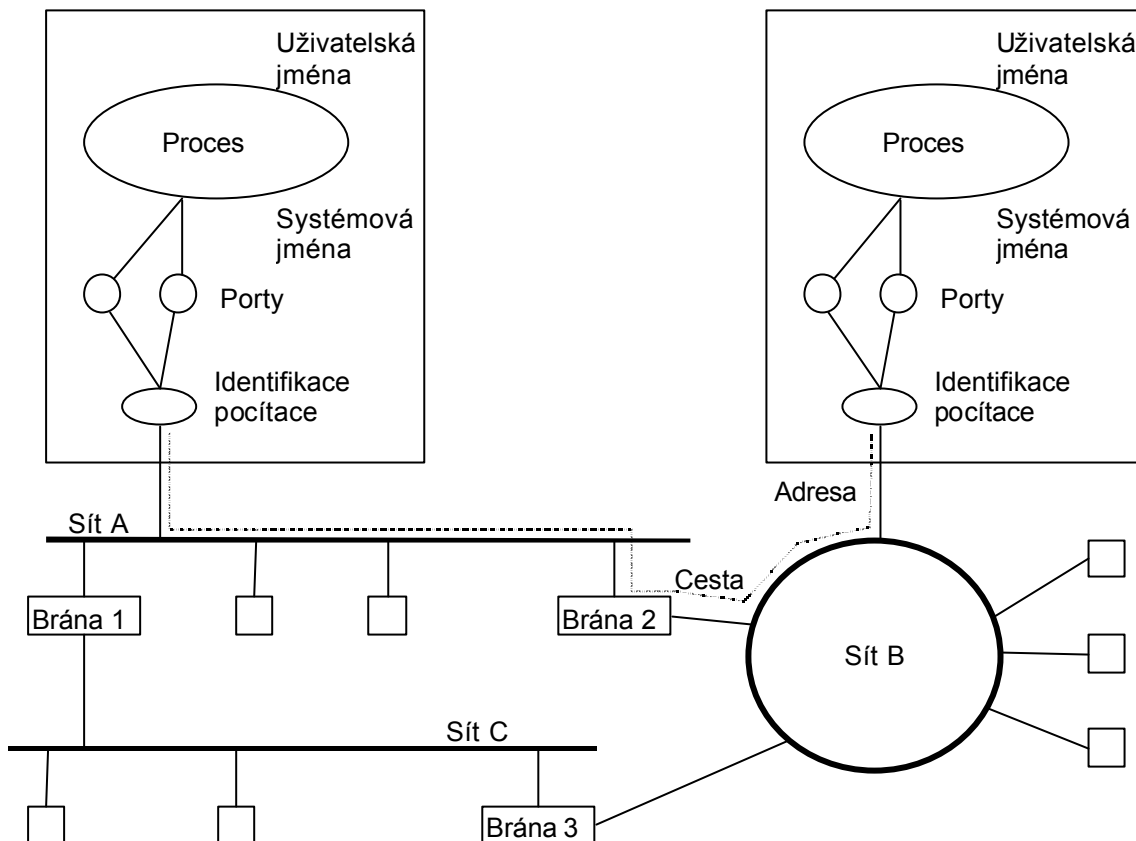
Prostor jmen musí být udržován decentralizovaně, tj. identifikační systém každého počítače musí jména přidělovat jen na základě lokálních informací. Takto přidělovaná jména musí být jednoznačná, tj. žádné dva objekty nemohou mít stejné jméno. Proto jsou jména jednotlivými jádry přidělována buď podle nějakého pevně daného algoritmu, nebo jsou jednotlivým počítačům přidělovány oblasti prostoru jmen, které může využívat.

Prostor jmen může být buď jednotný pro všechny objekty systému, nebo mohou být pro různé druhy objektu separátní prostory jmen. Napr. v klasických centralizovaných operačních systémech bývá prostor jmen souboru oddělený od identifikace ostatních objektu. Naopak v objektově orientovaných systémech bývá častěji jednotný prostor jmen.

Některé distribuované systémy podporují migraci objektu. U takových systému musí být jména objektu nezávislá na aktuálním umístění, resp. vazba mezi jménem a adresou musí být dynamická. V případě přesunu objektu na jiný počítač musí být transparentně zjištěna nová adresa a nalezena přístupová cesta k objektu.

Identifikační systém by měl podporovat vícenásobné kopie (repliky) daného objektu. Vazba mezi uživatelským jménem a adresou tedy není 1:1, ale 1:N. V případě nedostupnosti jedné kopie by měl systém transparentně namapovat jinou dostupnou kopii téhož objektu. Systém by také měl umožňovat různá lokální uživatelsky definovaná pojmenování pro stejné objekty a zároveň několika různým objektům (skupině objektu) umožnit sdílet jedno uživatelské jméno. Celkově je tedy mapování mezi uživatelskými jmény a adresami objektu typu M:N.

Návrh a konkrétní implementace identifikačního systému je vždy kompromisem mezi výkonností, pružností, snadností mobility objektu, pametovou náročností a snadností uchování a údržby konzistence.



Obr. 37 - Hierarchie identifikací v distribuovaném systému

7.1.1 Jména

Jméno jednoznačně identifikuje objekt nebo skupinu objektů. Může to být buď textový řetězec nebo nějaká binární hodnota. Binární jména jsou vhodnější pro použití jako systémové identifikátory, řetězová jména jsou naopak vhodnější pro uživatelskou identifikaci.

Jména mohou být rozdělena podle několika kritérií:

- ?? podle struktury - nestrukturovaná, strukturovaná, popisná
- ?? podle trvanlivosti - statická (trvalá), dynamická (dočasná)
- ?? podle počtu referencovaných objektů - individuální, skupinová
- ?? podle rozsahu platnosti - globální, lokální pro uzel, lokální pro proces

Statická jména trvale označují nějaký objekt. Dynamická jména jsou naopak přiřazována objektům jen (relativně k životnosti objektu) po krátký čas - např. po dobu komunikace s jiným objektem. Dynamická jména jsou v naprosté většině systémová - systém je sám vytváří a po ukončení doby platnosti je ruší.

Individuální jména jsou jména, která jednoznačně identifikují jeden objekt; všechny objekty mají navzájem různá jména. Pro účely skupinové komunikace, spolupráce replikovaných objektů apod. jsou často využívána skupinová jména, která identifikují nějakou přesně vymezenou skupinu.

7.1.2 Struktura jmen

Jména mohou být podle své struktury rozdělena na nestrukturovaná (flat), strukturovaná (partitioned) a popisná (descriptive).

Nestrukturovaná jména nemají, jak již název napovídá, žádnou vnitřní strukturu, avšak v rámci svého prostoru jmen jednoznačně identifikují objekty. Nestrukturovaná jména buď přidělují centralizovaný server jmen, nebo se vypocítávají podle nějakého algoritmu zajišťujícího jednoznačnost.

Další, v poslední době často využívanou, možností je přidělovat jména náhodně. V případě, že šířka jména, tj. počet bajtu, které jméno tvoří, je dostatečně velká a generátor náhodných čísel spolehlivý, je pravděpodobnost, že dva objekty dostanou přiděleno stejné jméno, velice nízká. Při binárních jménech o šířce 64 bitu (8 bajtu) se stejné jméno jedním uzlem vygeneruje jednou za 10^{19} cyklu. I při rychlosti generování milion adres ze vteřiny by to trvalo 10^{13} vteřin, což je asi půl milionu let. I kdyby adresy generovalo neustále paralelně tisíc počítačů, trvalo by zhruba 500 let, než by vygenerovaly stejnou adresu.

Strukturovaná jména jsou obvykle složena z posloupnosti jednoduchých jmen oddělených separátory. Jednotlivé části jména pak udávají oblast, podoblast, ... a tak dále až konečné (buď první nebo poslední) jméno udává identifikaci objektu v jeho lokálním prostoru. Separátory bývají speciální znaky jako např. '@', '.', '/', '%' apod., tyto znaky pak nemohou být součástí částečného jména. Příkladem strukturovaného jména je internetová adresa, která může být tvaru `uživatel@oddelení.organizace.oblast.stát`, kde část jména před znakem @ je identifikace koncového uživatele, část jména za znakem @ hierarchicky určuje adresu uživatele. Jiným příkladem strukturovaného jména je identifikace souboru v UNIX-like systémech. Jméno tvaru `/adresář1/adresář2/adresář3/soubor` určuje cestu z korenového adresáře k souboru a identifikaci souboru ve svém adresáři.

Strukturovaná jména mohou být absolutní a relativní. Absolutní jméno specifikuje celou cestu od nějakého korene až ke koncovému objektu. Absolutní jména by měla být stejná z libovolného uzlu systému. Relativní jména mohou nejvyšší části cesty vynechat, ty mohou být za určitých okolností z kontextu zřejmé. Tímto způsobem např. funguje systém aktuálních adresářů - ty soubory, které leží v adresářovém stromě pod aktuálním adresářem, mohou být adresovány relativně z tohoto adresáře; cesta od korene k aktuálnímu adresáři může být vypuštěna. Relativní jména je však třeba používat opatrně, neboť změnou kontextu se jména mohou stát neplatnými anebo ukazovat na jiný objekt. Absolutní jméno je v podstatě relativní jméno od korenového kontextu. V některých systémech však korene nemusí být jednoznačný, resp. korenu může být několik. V tom případě prostor jmen netvoří strom, ale les, v obecnějším případě dokonce libovolný orientovaný graf.

Popisná jména jsou tvořena množinou atributu, která je pro každý objekt jednoznačná. Podle této definice jsou strukturovaná jména zvláštní formou popisných jmen s pevnou strukturou atributu.

Atributy objektu mají typ a hodnotu, kde typ označuje formát a význam atributu. Příkladem atributového jména může být

Jméno: Novák; Oddelení: KSI; Organizace: MFF UK; Místo: Praha

Sjednocením několika atributů dostaneme jednoznačnou identifikaci objektu přesto, že jednotlivé atributy mohou být nejednoznačné. Popisná jména mohou být snadno použita i jako skupinová jména.

7.1.3 Cesty

Distribovaný systém bývá většinou složen z několika propojených sítí. K tomu, aby jednotlivé objekty spolu mohly komunikovat, je třeba znát cestu, kudy zpráva musí projít. Cestou rozumíme posloupnost jmen (adres), které určují jednotlivé úseky mezi počátečním a koncovým uzlem. Algoritmu, který mapuje jména na cesty říkáme smerování (routing). Smerovací algoritmy můžeme rozdělit podle několika hledisek:

- ?? podle místa, kde se provádí rozhodování o dalším úseku cesty
- ?? podle frekvence změn dat potřebných pro smerování
- ?? podle způsobu spolupráce jednotlivých uzlů

Jsou tři možné způsoby, kde se může provádět rozhodování o dalším úseku cesty:

- ?? zdrojové smerování
- ?? postupné smerování
- ?? záplavové smerování

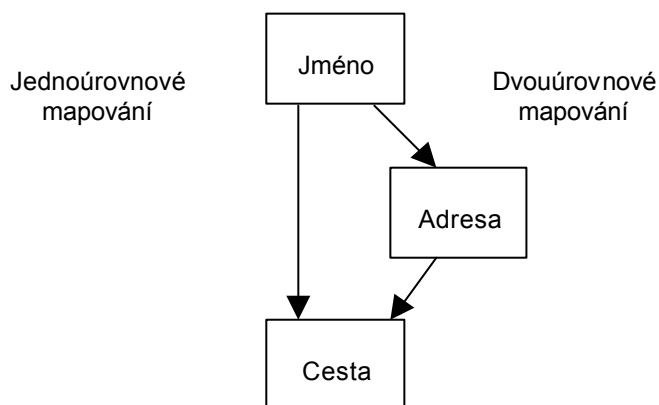
Podle frekvence zmen potřebných dat lze smerovací algoritmy rozdelit na statické (pevné) a dynamické (adaptivní).

Centralizované / lokální / distribuované smerování

7.1.4 Adresy a mapování

Adresa může být považována za jakýsi mezistav mezi jménem a cestou. Na rozdíl od jména je adresa pevně svázána s umístěním objektu. Jméno jednoznačně identifikuje objekt, kdežto adresa jednoznačně identifikuje nějaké umístění. Podobně jako jména, i adresy mohou být strukturované nebo nestrukturované - napr. trojice < síť, uzel, soket > jednoznačně určuje adresu.

Mapování jmen na cesty může být buď přímé (jednoúrovňové), kdy je cesta vygenerována bezprostředně ze jména, nebo zprostředkované (dvouúrovňové), kdy je nejprve jméno namapováno na adresu a poté adresa namapována na cestu.



Obr. 38 - Jednoúrovňové a dvouúrovňové mapování jmen na adresy a cesty

7.2 Systémová jména

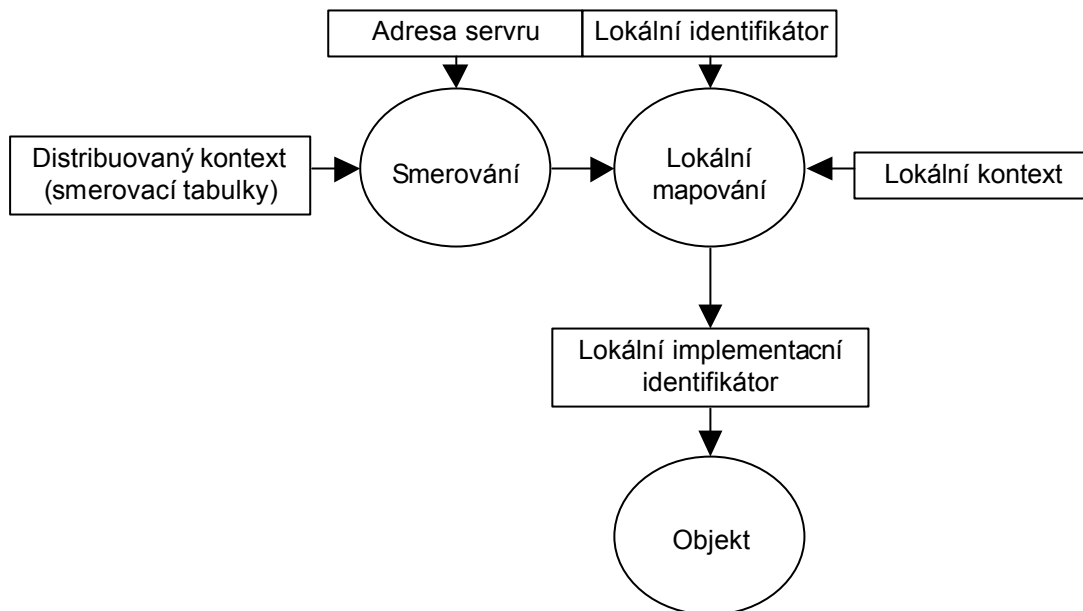
Systémová jména jsou generována a využívána systémem k vnitřní identifikaci objektu. Většinou to jsou binární řetězce pevné délky, které jsou pro uživatele necitelné. Systémová jména jsou založena buď na jednoznačných interních identifikátorech nebo na speciálních druzích jmen, které poskytují podporu ochrany objektu, tzv. kapabilitách.

7.2.1 Interní identifikátory

Jednoznačné interní identifikátory jsou v distribuovaných systémech často využívány. Typicky jsou tyto identifikátory nestrukturované, mohou však být i strukturované.

Nezávažnějším problémem nestrukturovaných interních identifikátorů je správa prostoru jmen - vytvoření globálního identifikátoru a zaručení jeho jednoznačnosti.

7.2.2 Lokálně implementační identifikátory



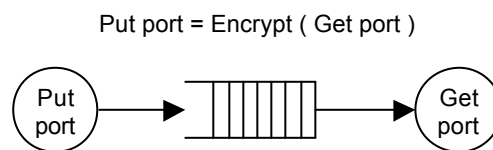
Obr. 39 - Mapování jednoznačného globálního identifikátoru

7.2.3 Mapování systémových jmen

Identifikace portu

?? globální - nutnost ochrany

?? lokální pro proces - ochrana většinou jádrem



Obr. 40 - Vstupní a výstupní port

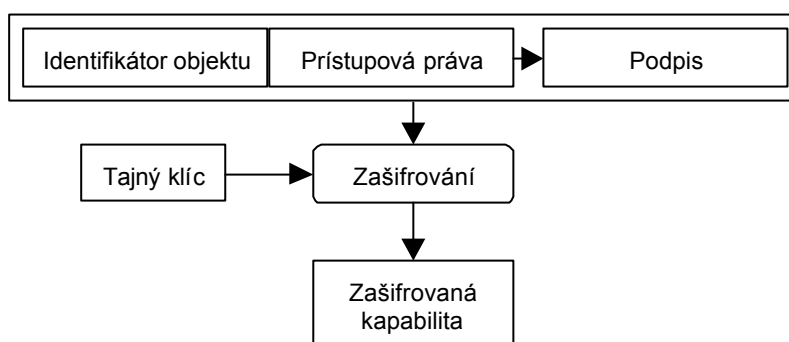
7.3 Kapability

Speciálním druhem systémových jmen jsou tzv. kapability (capabilities). Kapabilita je zvláštní "lístek" (bumážka) umožňující jednoznačnou identifikaci objektu a zároveň obsahující množinu přístupových práv, která umožňují držitelům kapability přístup k objektu a operace nad ním.

Z důvodu bezpečnosti musí být uživatelským procesem znemožněno vlastní generování kapabilit jakož i změny přístupových práv. Nejčastěji používaný způsob implementace je založen na řídkosti kapabilit a kódování.

7.3.1 Kapability s podpisem

Platná kapabilita je sestavena z dvou částí - identifikace objektu a přístupových práv. Platná kapabilita je rozšířena o podpis, který je vypočítán z obsahu hlavní části. Trojice <objekt, práva, podpis> je platná jen tehdy, pokud podpis odpovídá hlavní části kapability. Aby se zabránilo uživatelským procesum odvodit podpisovou funkci, je celá kapabilita zašifrována tajným klíčem. Řídkost je zde chápána v tom smyslu, že ze všech binárních čísel délky velikosti kapability jsou platné jen ty, jejichž podpis odpovídá zbytku kapability.

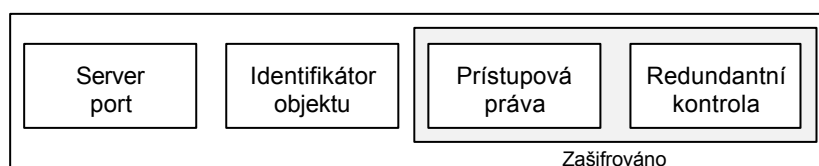


Obr. 41 – Kapabilita

7.3.2 Kapability s redundancí

Jiným způsobem ochrany kapabilit je redundantní kontrola. Platná kapabilita je složena ze čtyř částí. První dvě části, serverový port a identifikátor objektu, jsou volně přístupné. Další pole obsahující přístupová práva, je spolu s náhodně vygenerovaným binárním číslem pevné délky zakódováno.

Ochrana je zde složena ze dvou částí - ochrany serveru a ochrany přístupových práv. Ochrana serveru je dána portem, což je dostatečně velké náhodně generované číslo. Proces, který toto číslo nezná, nemůže službu serveru využít. Ochrana přístupových práv je zajištěna zašifrováním pole s přístupovými právy spolu s redundantní kontrolou. Tím je uživatelskému procesum znemožněna změna přístupových práv - proces nemůže svá práva změnit, neboť nezná dešifrovací funkci, ani ji nemůže (kvůli náhodně generované redundanci) z ničeho odvodit.



Obr. 42 - Kapabilita v Amoebě

7.4 Uživatelská jména

7.4.1 Servry jmen

Mapování uživatelských jmen na systémová je typicky prováděno servry jmen (name servers). Lokální name server udržuje veškeré informace o dostupných objektech na každém uzlu systému. Komunikační složitost je minimální, avšak údržba prostoru jmen a udržování konzistence je značně náročné. Proto jsou lokální name servery vhodné pouze pro malé a relativně statické prostory jmen.

Zcela opačným extrémem je centralizovaný name server. V tomto případě je celý prostor jmen spravován jedním name serverem. Údržba prostoru jmen a udržování konzistence je triviální, ale v případě výpadku name serveru je celý systém nefunkční. Zároveň pro větší systémy může být name server úzkým místem omezujícím výkon systému - v případě většího počtu požadavků je centrální server nemusí být schopen dostatečně rychle zpracovat. Proto se centrální name servery dají rozumně použít pouze u malých systému.

Vylepšením centralizovaného name serveru je hierarchická struktura name serveru. V případě, že požadavek nemůže být nějakým name serverem vyřešen, je předán name serveru vyšší úrovně. Ten většinou došlé požadavky sám nereší, ale zasílá je k vyřešení name serveru nižší úrovně, který je tento požadavek schopen vyřešit.

Distribuované name servery žádnou pevnou organizaci většinou neudržují. Každý lokální name server většinou udržuje informace o lokálních objektech a pro vzdálené objekty, resp. pro jejich podstromy, udržuje odkaz na vzdálený name server.

Pro minimalizaci komunikace bývají nejméně často užívaná jména lokálními servery cachována. Jestliže přijde požadavek na rezoluci nějakého jména, pak lokální server prohledá cache, a v případě nalezení jména ihned odpoví. K předání požadavku vzdálenému name serveru pak dochází jen u jmen, které se do cache již nevešly, nebo které ještě nebyly resolvovány. Jiným způsobem minimalizace komunikace je replikace údaje o důležitých nebo často používaných jménech. Takové jméno může být resolvováno libovolným name serverem, který o něm udržuje údaje. V tomto případě je však správa prostoru jmen podstatně náročnější.

7.4.2 Správa prostoru jmen

name distribution - přidělování autority spravovat část prostoru jmen

name resolution - výběr atributu objektu zadaného jménem

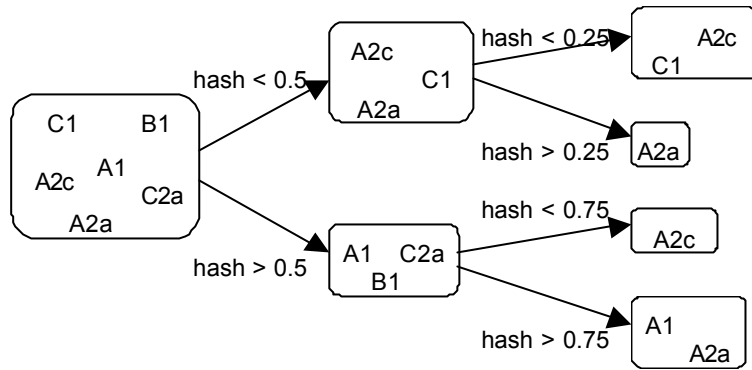
jména lokální / doménová / globální

jména absolutní / relativní / kontextová

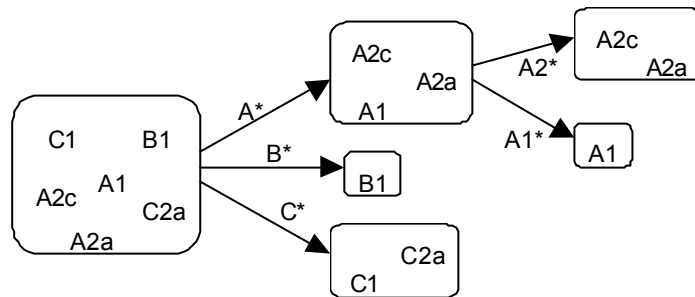
7.4.3 Rozklad jmen

Způsoby rozkladu:

- | | |
|-----------------|------------------------------------|
| ?? hierarchický | podle posloupnosti částecných jmen |
| ?? aritmetický | podle výsledku ohodnocovací funkce |
| ?? syntaktický | podle syntaktického rozkladu jména |
| ?? atributový | pro popisná (atributová) jména |



Obr. 43 - Aritmetický rozklad

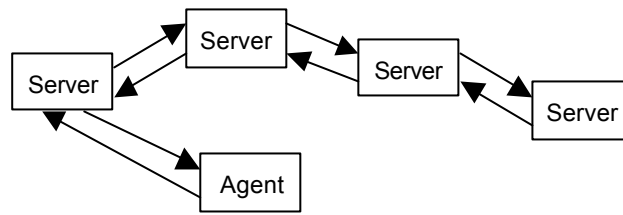


Obr. 44 - Syntaktický rozklad

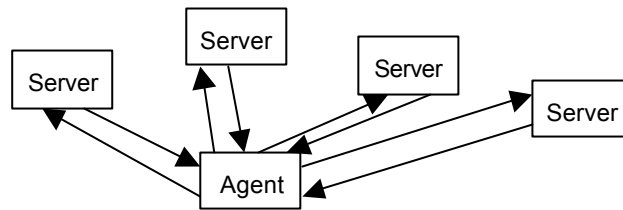
7.4.4 Agenti

lokální name server / exkluzivní (prilinkovaná knihovna) / sdílená knihovna

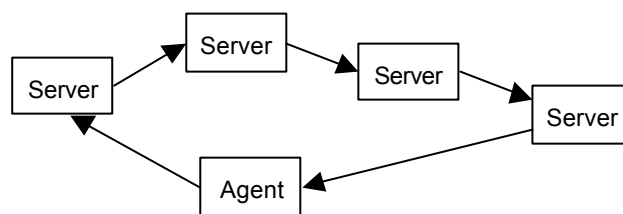
7.4.5 Vyhledávání jmen



(a)



(b)



(c)

Obr. 45 - Styly vyhledávání jmen

8. Procesy

Cílem správy procesu v distribuovaných systémech je, kromě úkolu běžných z jednoprocessorových systému, sdílet výpočetní sílu systému, rozdelovat zátěž na jednotlivé procesory, provádět operace na vzdálených procesech, synchronizovat procesy a vést evidenci jejich stavu.

8.1 Vlákna, nite a thready

Per proces	Per thread
Adresový prostor	Cítac instrukcí
globální promenné	Zásobník
deskriptory (soubory, ...)	Registry
signály, semaforey, ...	Stav
práva, úcty, ...	(globální promenné threadu)

Tab. 16 - Entity spravované procesem a threadem

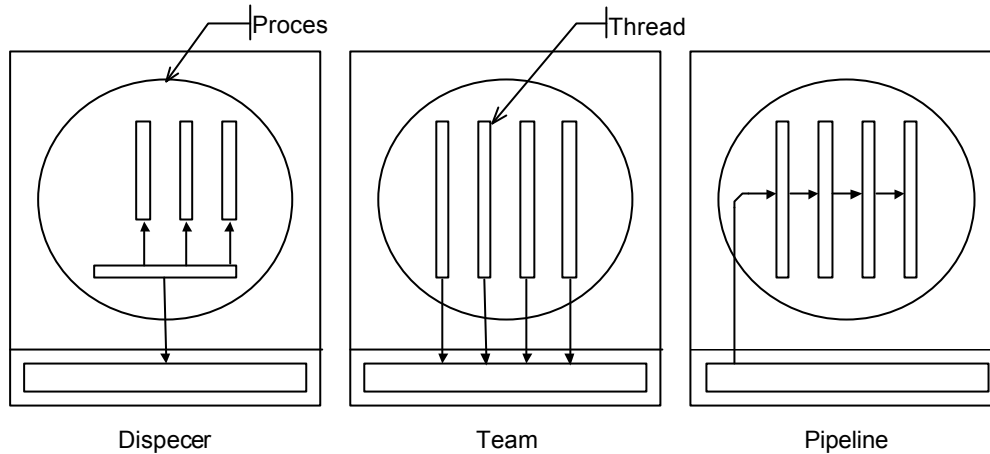
globální data threadu	globální data threadu	globální data threadu
Zásobník	zásobník	zásobník
sdílená dynamická data		
sdílená statická data		
kód		

Tab. 17 - Organizace pameti vícethreadového procesu

Vzhledem k tomu, že všechny thready jednoho procesu sdílejí adresový prostor, mají přístup ke všem globálním promenným. Často by však bylo třeba, aby každý thread měl své vlastní globální promenné, avšak lokální v daném threadu (tj. neprístupné z ostatních threadu). Příkladem možného využití může být chybová promenná `errno`. To je globální promenná, do které je při jakékoliv chybě uložen kód této chyby. Ve vícethreadovém prostředí však může dojít ke kolizi - thread A vyvolá chybu 1, do `errno` se uloží 1. Bezprostředně poté thread B vyvolá chybu 2, `errno` se prepíše hodnotou 2. Thread A pak při přístupu k `errno` přečte nesprávnou hodnotu.

Vzhledem k tomu, že snad žádný programovací jazyk nemá konstrukci umožňující práci s globálními daty threadu, musí být přístup k takovým datům implementován na uživatelské (knihovní) úrovni - napr. funkcemi `create_global`, `write_global` a `read_global`.

8.1.1 Použití vláken



Obr. 46 - Organizace vláken v procesu

8.1.2 Implementace vláken

Existují v zásadě dvě možnosti, jak vlákna implementovat, a to buď přímo v jádře, anebo jako knihovni nadstavbu na uživatelské úrovni.

Knihovna na uživatelské úrovni

Na uživatelské úrovni se systém threadu implementuje jako kolekce knihovních funkcí. Při pokusu threadu o zavolání systému se zavolá knihovni funkce. Ta kromě toho, že zavolá již přímo příslušnou systémovou službu, navíc provede případnou změnu kontextu (tj. v podstatě pouze registru) threadu.

Výhody:

- ?? nadstavba nad již existujícím softwarem (UNIX)
- ?? rychlé prepínání kontextu
- ?? možnost vlastního plánování

Nevýhody:

- ?? nutnost obalení
 - ?? blokováná (synchronní) systémová volání - ostatní thready jsou blokovány
 - ?? neblokovaná (asynchronní) systémová volání - nutnost zmen v jádru
- ?? výpadky stránek - ostatní thready jsou blokovány
- ?? není preemptivní
- ?? globální promenné
- ?? reakce na signály
- ?? nereentrance standardních knihoven

Implementace vláken v jádru

Jaderná implementace threadu nevyžaduje žádnou behovou podporu - veškerá rozhodování, plánování, prepínání kontextu, atd. se dějí přímo v jádru.

Výhody:

- ?? preemptivní plánování
- ?? jednodušší interakce s jádrem

Nevýhody:

- ?? větší jádro
- ?? pomalejší zpracování - vytváření/rušení threadu, ...

8.1.3 Vlákna a zprávy

V některých distribuovaných systémech je implementována velmi těsná vazba mezi příjmem zpráv a thready - tzv. implicitní receive, nebo také pop-up thread. Myšlenka spočívá v tom, že místo toho, aby zpráva čekala ve frontě, než si ji někdo vyzvedne, tak se automaticky vytvoří nový (resp. vezme se nějaký volný) thread, z dat zprávy se mu vytvoří zásobník a automaticky se spustí. Po skončení zpracování zprávy a odeslání odpovědi thread opět automaticky zanikne.

8.2 Systémové modely

Procesy (resp. jejich thready) běží na procesorech. V tradičních jednoprocessorových systémech není problém určit procesor, na kterém daný proces pobeží. Avšak v multiprocessorových a distribuovaných systémech je metoda alokace procesoru jednou ze základních otázek designu celého systému. Existují dva zásadně odlišné modely alokace procesoru - workstation model a procesor pool model. Kromě těchto dvou základních forem existují ještě různé hybridní formy jako mezistavy slucující výhody (a často i nevýhody) jednotlivých systémů.

8.2.1 Workstation model

Workstation model je založen na myšlence, že každý uživatel má "svůj" počítač, který je mu plně k dispozici. Kromě uživatelských uzlů bývají v takovýchto systémech ještě uzly dedikované pro nějaké servery - napr. fileservr.

Vzhledem k tomu, že naprostá většina distribuovaných systémů využívá jakýsi fileservr (at už centralizovaný, distribuovaný, apod.), je otázkou, co s lokálními disky jednotlivých uživatelských stanic (pokud ovšem nejsou bezdiskové). Možností využití lokálního disku je několik, každá má své výhody i nevýhody.

	Výhody	Nevýhody
Bezdiskové stanice	Levné, snadná údržba, flexibilita	Vytížení sítě, zátěž fileserveru
Stránkování, dočasné soubory	Snižuje vytížení sítě	Větší náklady (na disky)
Stránkování, dočasné soubory, systémové programy	Ještě více snižuje vytížení sítě	Větší náklady, údržba systémových programů
Stránkování, dočasné soubory, systémové programy, file caching	Redukce zátěže fileserveru	Udržování konzistence cache
Součást distribuovaného filesystemu	Výrazná redukce zátěže ostatních fileserveru	Udržování konzistence celého filesystemu
Kompletní lokální filesystem	Eliminace potřeby fileserveru	Ztráta transparentnosti

8.2.2 Procesor pool model

Druhou možností je abstrahovat od pojmu “uživateluv počítač/procesor” a přistupovat ke všem procesorům dostupným systému jednotně. Hardwarově je tento přístup realizován většinou polem procesorů, ke kterým jsou lokální sítě připojeny uživatelské terminály (většinou grafické X-terminály). Tímto způsobem se snáze a rovnoměrněji rozloží výpočtová zátěž mezi dostupné procesory, avšak je zapotřebí speciální, ne vždy dostupný, hardware.

8.3 Vzdálené spuštění procesu

Aby mohla být efektivně využita veškerá výpočetní síla celého systému, je občas třeba spustit nový proces na jiném procesoru.

Problémy:

1. Jak nalézt volný počítač
2. Jak proces vzdáleně spustit (tak aby to proces sám nepoznal)
3. Co se stane, když počítač přestane být ‘volný’

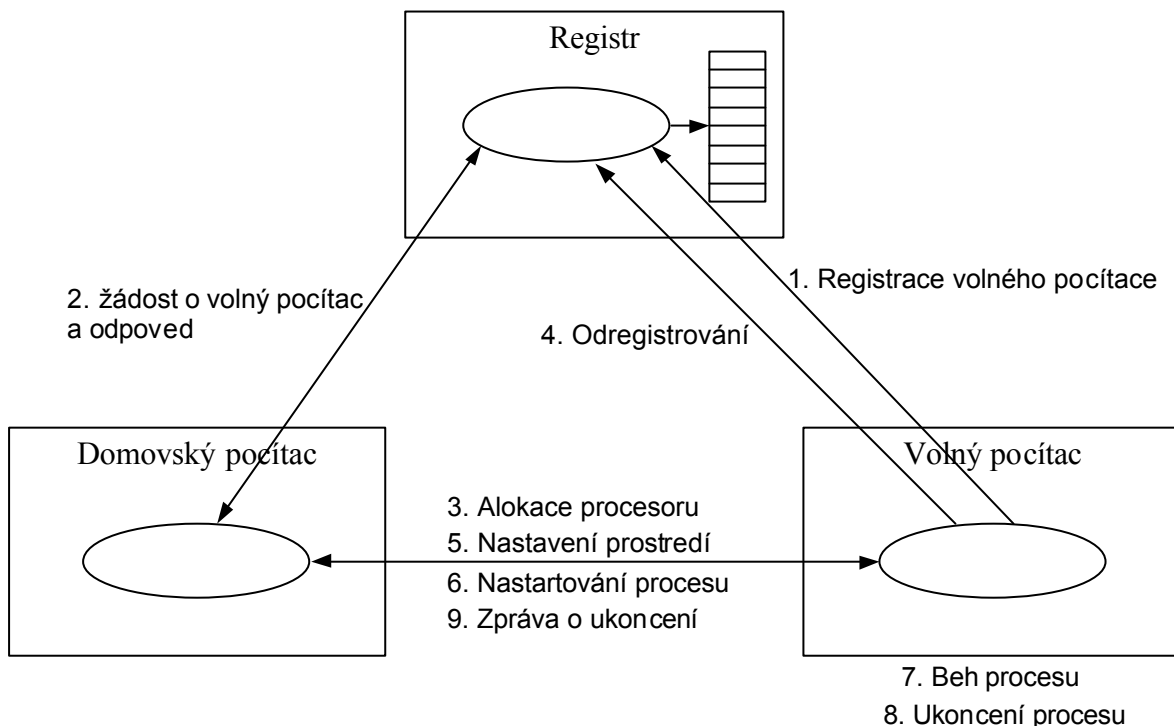
První zajímavou otázkou je, co to je vlastně volný procesor. Pro první přiblížení to může být procesor v počítači, na kterém není nikdo přihlášen. To, že na daném počítači není nikdo přihlášen, však nemusí nutně znamenat, že procesor je volný a naopak, může být volný procesor v počítači s přihlášeným uživatelem. Prvnímu případu odpovídají různé démony (mail daemon, clock daemon, atd.), které dostatečně vytežují procesor, naopak druhému případu odpovídá (v některých institucích zcela běžná) situace, že uživatel přijde k počítači, přihlásí se, přečte si noviny, dojde si na oběd apod. a procesor nemá dlouhou dobu co na práci.

Jiným možným kritériem “volnosti” procesoru je stav, kdy nikdo dostatečně dlouho nevytvárel uživatelské signály, tj. nepsal nic na klávesnici ani nepohyboval myši, a zároveň aktivně nebeží žádný uživatelský proces.

Další krok, který se musí provést, je spuštění procesu. Přenesení kódu a statických dat procesu je poměrně jednoduché (snad s výjimkou heterogenních procesorů, kde je třeba jistých konverzí). Zajímavější je nastavení prostředí na hostitelském počítači tak, aby výpočet probíhal stejně, jako kdyby proces byl spuštěn na domovském počítači. Je třeba přenést množinu kontextu (tj. aktuální adresáře filesystemu a nejen jeho), různé proměnné domovského prostředí (environmental variables), pokud takové existují, apod. Dalším problémem jsou systémová volání - některá systémová volání by se měla provádět na hostitelském počítači (napr. `sbrk`, `nice`, `profil` apod.), zatímco jiná by měla být přesmerována na domovský počítač (napr. čtení z klávesnice, zápis na terminál apod.). Jiná volání mohou být podle kontextu volána lokálně anebo přesmerována na domovský počítač (přístup

k souborům apod.). Proto je zapotřebí vytvořit kanály zpět na domovský počítač a některá systémová volání přesměrovat na tyto kanály.

Zajímavá je otázka, co se má stát, když hostitelský počítač přestane být volný (např. uživatel se vrátil z oběda). První možností je nechat všechny cizí procesy dobehnout s tím, že uživatel má prostě smulu a výkonnost “jeho” počítače se patrně dočasně o něco sníží. To však v některých případech může být nepříjemné. Druhým extrémem je nekompromisní ukončení procesu; to však má dvě nevýhody: jednak může být ztraceno značné množství práce, jednak může být systém v nekonzistentním stavu. Proto jako jedna z možností je nechat procesu jistý čas na dokončení potřebných akcí a dovedení systému do konzistentního stavu. To však vyžaduje od takto vzdáleného procesu jistou spolupráci. Poslední možností je migrace procesu v rozbehnutém stavu na jiný počítač, to je však o několik rádu složitější akce než vzdálené spuštění.



Obr. 47 - Nalezení a využití volného počítače

8.4 Alokace procesoru

Každý distribuovaný systém je složen z množiny procesoru, at již jsou organizovány jako pole procesoru, uživatelské stanice nebo nějak jinak. V každém případě je zapotřebí v okamžiku vzniku procesu určit, na kterém procesoru pobeží. Algoritmy řešící tento problém se nazývají algoritmy alokace procesoru.

8.4.1 Klasifikace alokacních algoritmu

Algoritmy alokace procesoru lze hodnotit podle mnoha různých hledisek (podtržené varianty jsou běžnější nebo považovány za prakticky použitelnější):

- ?? zda berou v úvahu rozdíly mezi jednotlivými procesory a speciální požadavky procesu na vlastnosti procesoru - homogenní / heterogenní
- ?? zda umožňují přemístění již rozbehnutého procesu - migrační / nemigrační
- ?? co se snaží optimalizovat - využití CPU / čas odpovědi / spravedlivé přidělování prostředku / komunikační složitost, ...
- ?? zda jsou předem známy údaje o procesech, podle kterých se algoritmus rozhoduje - deterministické / heuristické
- ?? podle charakteru algoritmu - centralizované / distribučované
- ?? podle toho, do jaké míry se provádí optimalizace - optimální / suboptimální
- ?? podle jakých informací se provádí rozhodnutí o vzdáleném spuštění procesu - lokální / globální
- ?? podle toho, kdo iniciuje vzdálené spuštění procesu - odesílatel / příjemce / symetrický

8.4.2 Implementace alokacních algoritmu

Základem pro jakýkoliv alokacní algoritmus je znalost vlastní zátěže. To však není tak jednoduché, jak se na první pohled zdá. Při prvním přiblížení by to mohl být počet spuštěných procesů. To však o vytíženosti procesoru příliš nevyovídá - jeden výpočetně náročný proces může vytížit procesor daleko více nežli třeba deset procesů, které však většinu času čekají na nějaký signál.

Druhou možností je počet aktivně běžících procesů. Tento počet se však velmi rychle mění - různé demony jsou periodicky ožívovány a po vykonání několika instrukcí jsou opět suspendovány, procesy interaktivně komunikující s uživatelem typicky čekají na nějaký signál. Přesnější výpočet zátěže může být založen na periodickém měření míry vytíženosti procesoru. To může být implementováno periodicky vyvolávanou prázdnou smyčkou, kdy čas provedení určitého počtu cyklu je přímo úměrný zátěži procesoru ostatními procesy. Problém tohoto řešení je v tom, že když jádro provádí nějakou důležitější akci, je zakázáno přerušení, tudíž i uživatelské časové signály. Proto někdy měření zátěže nedojde a vytížení procesoru může být podhodnoceno.

Dalším problémem alokacních algoritmu je režie vlastního algoritmu. V případě, že by nějaký alokacní algoritmus měl výrazně snižovat výkon procesoru (řádově o několik až desítky procent), pak je často lepší zvolit jednodušší algoritmus, který sice nedosahuje teoreticky optimálního výsledku, avšak v důsledku podstatně menší režie je celkově efektivnější.

8.4.3 Deterministický grafový algoritmus

Příkladem centralizovaného deterministického optimálního alokacního algoritmu je algoritmus minimalizující vzdálenou komunikaci. Každou dvojici vzájemně komunikujících procesů meje ohodnocenu velikostí vzájemné komunikace. Pak celý systém může být reprezentován grafem s ohodnocenými hranami. Úkolem algoritmu je nalézt takové vzájemné přiřazení procesů a procesoru, aby komunikace mezi jednotlivými procesory byla minimální. Úloha se pak redukuje na rozdělení grafu na podgrafy takové, aby součet ohodnocení hran mezi podgrafy byl minimální. Nevýhoda tohoto algoritmu je, že potřebuje znát předem pro každou dvojici procesů jejich komunikační složitost, proto je (kromě teoretického zkoumání) použitelný jen pro velmi speciální případy.

8.4.4 Up-down algoritmus (Mutka-Livny)

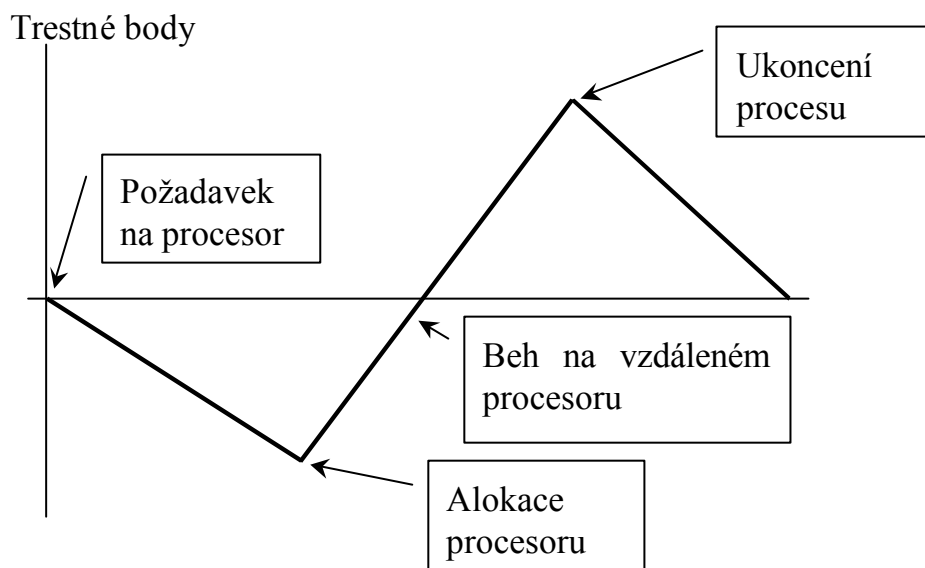
Jiný typ alokacního algoritmu, nazývaného up-down algoritmus (česky snad schodovitý), publikovali v roce 1987 Mutka & Livny. Pro každou doménu existuje jeden koordinátor, který udržuje tabulku se záznamem pro každý uživatelský počítač. Při každé významné akci, tj. vytvoření procesu, ukončení procesu a tik hodin (presněji každý n -tý tik), se pošle zpráva koordinátoru, který provede změny v příslušném záznamu.

Algoritmus je optimalizován na stejnomerné sdílení výkonu všemi uživateli. K tomuto účelu tabulka pro každý procesor obsahuje "trestné body". Kladné trestné body znamenají, že daný procesor vyslal nějaký proces na jiný procesor, zatímco záporné trestné body znamenají, že existují neuspokojené požadavky na běh vzdáleného procesu.

Při každé významné události se tabulka mění následujícím způsobem:

- ?? za každý proces běžící na jiném počítači - plus trestné body
- ?? za každý zatím neuspokojený požadavek na vzdálený procesor - mínus trestné body
- ?? jestliže nic z tohoto - směrem k nule

Při pokusu o běh procesu na vzdáleném počítači se pošle zpráva koordinátoru. V případě, že je některý procesor volný, alokuje se pro tento proces. V opačném případě se požadavek pozdrží a uloží do fronty. V okamžiku uvolnění některého procesoru se vezme ten proces z fronty neuspokojených požadavků, jehož vysílající procesor má nejméně trestných bodů.



Obr. 48 - Nalezení a využití volného počítače

Podle výše uvedených klasifikací lze tento algoritmus charakterizovat jako homogenní, nemigranční, spravedlivý, heuristický, centralizovaný, suboptimální, globální a symetrický.

8.4.5 Hierarchický algoritmus

Centralizované algoritmy jsou většinou špatně použitelné v rozsáhlých systémech, kde jsou výhodnější hierarchické nebo distribuované algoritmy. Příkladem hierarchického algoritmu je alokační algoritmus použitý v systému MICROS. Je založen na logické hierarchii procesoru nezávislé na fyzickém rozmístění jednotlivých počítačů. Procesory na nejnižší úrovni vykonávají vlastní práci, procesory na vyšších úrovních jsou manažeri. V případě, že manažer dostane požadavek na přidělení procesoru, zjistí, zda-li některý z jemu podřízených procesorů je volný. V případě, že ano, pak ho alokuje, v případě, že ne, pošle požadavek vyššímu manažerovi.

Zajímavý je problém, co se stane, když nějaký manažer vypadne. Musí se nahradit někým z jeho podřízených. Který z nich to bude si mohou určit podřízení sami (např. nějakým elekčním algoritmem), nebo ho může určit manažer vypadlého procesoru. Pro zabezpečení větší stability mohou být manažeri nejvyšší úrovně replikováni, tj. manažer nejvyšší úrovně nemusí existovat; je několik vysokých manažerů se stejnými pravomocemi.

8.4.6 Distribuovaný heuristický algoritmus

k náhodných výberu cíle

server initiated / receiver initiated / kombinovaný

8.4.7 Bidding (obchodní) algoritmus

Procesy kupují výpočetní sílu, procesory ji nabízejí

9. Migrace procesu

Migrací procesu rozumíme korektní a transparentní přenesení již spuštěného procesu kdykoliv během jeho výpočtu ze zdrojového na cílový počítač. "Korektnost" znamená, že ostatní procesy nejsou touto migrací nijak ovlivněny (kromě nepřímých vlivů jako např. změna dostupného výkonu procesoru apod.) a že po ukončení migrace stav systému odpovídá stavu, jako kdyby byl migrovaný proces spuštěn na cílovém počítači. Vzhledem k migrovanému procesu "transparentnost" znamená, že proces o migraci ani o tom, že běží na novém počítači, nemusí vůbec vědět a nemusí ani při migraci spolupracovat (tj. nemusí obsahovat kód podporující nebo umožňující migraci). Vzhledem k procesům komunikujícím s migrovaným procesem to znamená, že zůstanou zachovány všechny vazby na migrovaný proces a že budou doručeny všechny zprávy, a to i ty, které byly vyslány během vlastní migrace.

Pro zavedení mechanismu migrace procesu existuje několik důvodů:

- ?? vyvažování zátěže
- ?? optimalizace - I/O, komunikační, ...
- ?? přemístění serveru
- ?? shutdown

Před tím, než dojde k vlastní migraci, je třeba určit který proces má být migrován a kam má být přenesen. To může být provedeno buď ručně, např. přemístění nějakého serveru před vypnutím počítače, nebo automaticky, např. nějakým algoritmem vyvažování zátěže.

9.1 Mechanismus migrace

Vlastní mechanismus migrace je prováděn v několika krocích. Jelikož v různých systémech mohou být migrační mechanismy implementovány značně rozdílně, je třeba se na tyto kroky dívat spíše jako na logické události, než na výkonné akce.

1. Zmražení procesu - všechny thready migrovaného procesu musí být vyjmuty z behových front. Často bývá takto vyjmutý proces označen v systémových tabulkách "v migraci".
2. Inicializace cílového počítače - cílový počítač je informován, že mu přibude proces. Jádro připraví pro nový proces potřebné struktury, tj. vytvoří prázdný proces.
3. Přenesení stavu procesu - přenesení kódu, zásobníku threadu, datových segmentu, vazeb na jádro a ostatní procesy apod. na cílový počítač.
4. Reinicializace komunikačních kanálů - přesměrování komunikačních kanálů na cílový počítač a doručení zpráv odeslaných během migrace.
5. Vycištění zdrojového počítače - jsou zrušeny všechny jaderné struktury použité pro migrovaný proces.
6. Spuštění procesu na cílovém počítači - od tohoto okamžiku proces zcela normálně běží dál, v ideálním případě ani nezaregistruje, že byl přenesen.

9.1.1 Doručování zpráv

Je třeba implementovat mechanismus, který během migrace a krátce po ní doručí ve správném pořadí všechny zprávy určené pro migrovaný proces.

První možností je zaslat všem potenciálním odesílatelům zprávu, že proces je migrován. Procesy, které by chtěly zasílat zprávu migrovanému procesu, budou suspendovány až do ukončení migrace a navázání nových spojení. Problém tohoto řešení je v tom, že potenciální odesílatelé nemusí být zdrojovému jádru známi. Navíc takovéhoto odesílatelu může být velké množství, často v podstatě jakýkoliv proces celého systému. Pak by taková metoda byla komunikacně velmi náročná, a to i v případě, že během migrace s migrovaným procesem nikdo nekomunikuje.

Druhou možností je nechat ve zdrojovém jádře odkaz na cílový počítač s tím, že všechny zprávy budou přesmerovány. To však vytváří reziduální dependence se všemi nepříjemnými důsledky.

Další možností je všechny zprávy došlé během migrace a po ní vracet zpět odesílateli s poznámkou "adresát odmigrován" a s adresou cílového počítače. Tento způsob však také vytváří reziduální dependence, proto je možné implementovat slabší variantu, že adresa odmigrovaného procesu je známa pouze určitý časový interval a po uplynutí tohoto intervalu je případným odesílatelům zpráva vrácena s poznámkou "adresát neznámý". Je pak na odesílateli, aby si zjistil novou adresu příjemce.

copy on reference - imaginary segment

prískoky - viskozita

9.1.2 Stav procesu

Každý proces je charakterizován jednak svým vnitřním stavem (např. obsah registru), ale i stavem spojeným s ostatními procesy nebo částmi systému (komunikační kanály, ...). Problémy, které se musí řešit při návrhu migrace jsou zejména jak

1. zjistit stav procesu - nejedná se jenom o stav procesu samotného (např. obsah registru), ale i stav týkající se komunikace s ostatními procesy, eventuální otevřené soubory, ...
2. vyjmout proces ze zdrojového počítače
3. přenést proces na cílový počítač
4. vložit proces na cílový počítač

Stav procesu jako takový se skládá z velkého množství součástí. Při migrování procesu z jedné stanice na druhou je potřeba u každé části stavu rozhodnout, jakým způsobem zajistit, aby daná část stavu byla procesu k dispozici i na cílové stanici (z tohoto hlediska by bylo přesnější v bode 3 výše uvedeného výčtu uvést: „přenést proces a zajistit, aby měl na cílové stanici všechny prostředky, které měl na stanici zdrojové.“

Možnosti jak toho dosáhnout jsou tři:

1. Přesunout danou část stavu spolu s procesem.
Neodiskutovatelným zástupcem této skupiny je obsah virtuální paměti.
2. Připravit část stavu pro forwardování požadavku, které se jí týká
Zde je možné si vymyslet například komunikaci s konzolí systému. Stejně tak se forwardovací přístup volí tam, kde je nemožné stav bezpečně vyextrahovat.
3. Použít odpovídající prostředek na cílové stanici.
Například fyzická paměť.

9.1.3 Komunikace s okolím během migrace a po migraci

Protože okolní procesy netuší, že se s migrovaným procesem „něco děje“ (a po pravdě receno ani není vhodné, aby se to mohly dozvedet), nemohou po dobu migrace daného procesu přestat s posíláním zpráv pro migrovaný proces. Migrovaný proces je ovšem z poloviny na starém počítači (kam jsou mu

posílány zprávy) a z poloviny již na novém (kam by mu zprávy mely být posílány po migraci). Je několik možností, jak se s tímto problémem vyrovnat, budeme je podrobněji rozebírat.

9.1.4 Implementace migrační strategie a vlastní migrace

Součástí migrační strategie (tj. procesu, během něhož se rozhoduje co a kam odmigrovat) je zjišťovat stav zátěže jednotlivých počítačů, sledovat, které počítače nejsou jejich uživateli využívány a na základě tohoto pozorování rozhodnout o migraci.

Druhá věc je vlastní realizace migrace, což je proces, během kterého dochází k přesunu stavu procesu na cílový počítač. Je rozumné tyto dvě části migračního mechanismu od sebe oddělit jednak z důvodu modularity, jednak z důvodu bezpečnosti.

Migrační strategie se typicky ponechává na některém specializovaném procesu, který běží na uživatelské prioritě, zatímco migrace samotná je často (alespon z části) implementovaná v jádře systému - odtud bezpečnostní důvody pro rozdělení.

Co se realizace migrace týče, je několik možností, jak ji implementovat.

1. V jádře systému
2. Vne jádra, tj. jako uživatelský proces - přináší komplikace, protože ne vše si může uživatelský proces dovolit
3. Jako součást procesu, v run-time knihovně - jde přímo proti požadavkům na transparentci

9.1.5 Vícenásobná migrace

Bežný pohled na migraci říká: „Vezmi proces, odmigruj jej.“ Problém se týká skeptického pohledu na část výroku „vezmi proces“. Proc jeden, proc ne třeba dva či několik najednou?

Vícenásobná migrace se snaží řešit problémy spojené s migrováním několika procesů zároveň. Tyto problémy se týkají zejména meziprocesové komunikace, protože je potřeba speciálně ošetřit, zda náhodou nemigrují dva procesy, které spolu komunikují. Jinak se totiž řeší, když migruje jeden z komunikující dvojice a druhý zůstává na zdrojovém počítači.

9.1.6 Transparentce

Transparentce migrace znamená, že „nikdo o migraci neví a nemůže zjistit, že k ní dochází“. Význam slova „vedet o migraci“ se liší s tím, koho se migrace týká.

Migrovaný proces - neobsahuje kód, který se týká migrace, nemusí ošetřovat situace, zda došlo k jeho migraci, či nikoliv, ...

Ostatní procesy - stejně jako migrovaný proces neobsahují kód řešící otázky migrace

Uživatel - Program nesmí kvůli migraci změnit chování a musí vydat takové výsledky, jako kdyby celou dobu běžel na zdrojovém počítači.

Často zmiňovaným tématem souvisejícím s otázkami transparentce je dilema „systémová volání“ a „posílání zpráv“ jako základní kameny operačního systému. Panuje běžný názor, že při použití meziprocesové komunikace založené na posílání zpráv se zjednoduší otázky transparentce. Do jaké míry jsou tyto názory pravdivé, je diskutabilní.

Bežný pohled totiž říká, že u zpráv stačí pouze přesmerovat komunikační kanál a je vše v pořádku. Narozdíl od toho se příliš nedá mluvit o přesmerování systémových volání.

Přesmerování komunikačního kanálu nelze ovšem v žádném případě provést slepe, aniž vím, co se nachází na druhém konci. Příkladem může být správce virtuální paměti. Po přesunu procesu na novou

stanici je na místě očekávat, že o pamet bude proces žádat toho správce, který běží na nové stanici a nikoliv toho starého. Slepým přesmerováním komunikačního kanálu by požadavek došel na původní stanici a to evidentně nebylo cílem.

Naproti tomu není zase tak veliký problém implementovat „přesmerování“ volání jádra. Znamená to v extrémním případě všechna volání jádra forwardovat na domácí stanici. Dokonce jeden z prvních distribuovaných operačních systémů Remote UNIX tuto variantu používal.

9.1.7 Reziiduální dependence

Jestliže proces, který odmigroval, stále ke své činnosti potřebuje spolupráci zdrojového počítače, jedná se o reziiduální dependenci. Problém, který přináší, je právě závislost procesu na jiném počítači, než na kterém proces momentálně běží. V případě, že se přeruší spojení obou zainteresovaných počítačů, může to mít pro proces fatální důsledky. To komplikuje i situaci, kdy se snaží systém odsunout procesy z počítače, který bude v blízké budoucnosti od sítě odpojen, popř. vypnut. Další vada reziiduální dependence je ta, že se musí jiný počítač zabývat existencí procesu na jiném stroji. V těchto okamžicích je snaha reziiduálních dependencí se vyvarovat.

V jiné situaci se reziiduální dependence hodí, zejména tehdy, když uživatel zjišťuje, které procesy spustil. Tehdy by mu systém měl ukázat i procesy, které odmigrovaly, tj. systém musí o nich mít informace, které jsou ovšem reziiduálními dependencemi.

9.1.8 Virtuální pamet

Virtuální pamet tvoří často největší část stavu procesu. Na tom, jakým způsobem se bude provádět přesun virtuální pameti, velmi závisí celková doba migrace procesu. Tvůrci systému postupně vymysleli několik možností, jak přesunout obsah pameti pokud možno rychle a efektivně.

- ?? Prvotní přímocárý přístup k přesunu virtuální pameti byl přesunout veškerou pamet najednou jako součást presunu stavu procesu. Výhody - po provedení presunu stavu může zdrojová stanice zapomenout, že na ní kdy proces bežel, tedy eliminace reziiduálních dependencí. Nevýhody - prodloužení doby, po kterou je proces zmražen (behem presunu stavu proces typicky nebeží). Navíc je mnohdy zbytečné přesouvat veškerý obsah virtuálního adresového prostoru (viz například veliký proces, který okamžitě po odmigrování volá `exit()`).
- ?? Pre-copying. Tento postup spocívá v tom, že proces běží behem doby, po kterou se přesouvá obsah jeho virtuální pameti. Tesne pred koncem kopírování se proces zmrazí, provede se dokopírování zbytku adresového prostoru (ty části, které proces ještě stacil zmenit) a presun stavu (obsah registru, ...). Výhoda - proces je zmražen pouze po dobu presunu malého množství informací. Nevýhoda - některé části adresového prostoru se kopírují vícekrát, což prodlužuje celkovou dobu migrace.
- ?? Copy-on-reference. Copy on reference spocívá v tom, že se nejprve přeneše ten stav procesu, který je potřebný pro beh procesu - obsah registru, komunikační kanály, ... Presun adresového prostoru se odloží na později. Stránky adresového prostoru na cílové stanici jsou oznaceny jako neprezentní (stejne jako by byly odswapovány na disk), ale je u nich poznámka, že se mají přenést ze zdrojové stanice. Při přístupu na takovou stránku (on reference) se její obsah přeneše (copy) a na zdrojové stanici se smaže.
- ?? Další možností je kombinace výše uvedených variant. Jedna z nich bude popsána v odstavci venovanému operačnímu systému Sprite.

9.1.9 Cíle návrhu migrace

1. Oddělení migrační strategie od realizace migrace
2. Co možná největší nezávislost implementace migrace na ostatních částech systému
3. Transparence
4. Autonomie - systém sám rozhoduje o tom, co a kam migrovat
5. Možnost evikce procesu z počítače, jehož majitel se vrátil
6. Vyvarovat se reziduálních dependencí, kde je to jenom možné
7. Spolehlivost
8. Efektivita
9. Jednoduchost
10. Fault-tolerance

9.1.10 Kdy migrovat

Možnosti jsou vesměs dvě. Při spuštění procesu nebo při jeho běhu.

První varianta je méně náročná na čas, protože například nemusí docházet k přesunu virtuální paměti mezi počítači, nicméně není to postacující varianta pro implementaci migrace z důvodu občasné potřeby provést migraci na již běžícím procesu (pr. evikce).

Druhá možnost je obecnější, ale vyžaduje složitější mechanismus.

9.1.11 Výber cíle migrace - migrační strategie

Pro rozhodování co a kam se bude migrovat je potřeba provádět sber velkého množství nejručnějších dat. Tato data se porizují různými způsoby podle toho, jakou mají povahu:

- ?? Při nějaké události - event sampling - například sber informací o posílání zpráv, využívání výstupu na konzoli, ...
- ?? V intervalech - zátěž CPU, zátěž síte, ...
- ?? Periodická statistika - slouží k vypocítávání dlouhodobějších prumeru

Je pravda, že na rozhodnutí jaký proces kam odmigrovat má vliv celá rada faktorů. Kromě těch, které se týkají samotného procesu a jeho cinnosti (využití procesoru, síte, cetnost vytváření potomku) jsou to také faktory týkající se hardwarové konfigurace, stavu diskových keší, doby, po kterou je stanice uživatelem nepoužita, ...

Tyto informace je potřeba sbírat, sledovat a vycházet z nich při rozhodování o migraci.

Návrh funkce migračního politika, tj. té části migračního mechanismu, který toto rozhodování provádí, musí zohlednit následující požadavky:

- ?? Výkonnost - nelze dlouho přemýšlet o volbě procesu a jeho cíle pro migraci, neboť potom může být pozde.
- ?? Rozširitelnost - mechanismus by měl počítat s velkým počtem stanic, které bude obsluhovat. Bežným požadavkem je obslužení rádove stovek stanic pripojených do síte.
- ?? Fault-tolerance - výpadek jedné komponenty by nemel ohrozit rozhodování ve zbylé funkční části síte.

?? Spravedlivost - rozhodovací systém by nemel dovolit, aby procesy, jejichž zdrojem je jedna "ciperná stanice" ovládly síť do té míry, že by jiné stanice nemohly rozumne provádět svoji činnost.

?? Jednoduchost

Je potřeba vzít v úvahu i možnost, ve které jsou data potřebná pro rozhodování uložena na jiném místě, než rozhodovací „mozek“.

Pro návrh implementace rozhodovacího mechanismu je tedy několik možností.

Sdílený soubor

V souboru, který je uložený na nejakém centrálním fileserveru, se ukládají záznamy obsahující identifikátor stanice, její průmernou zátěž v posledních nekolika (5-15) minutách, doba od posledního interaktivního použití stanice (signalizace idle stanic), příznak, zda stanice přijímá cizí procesy, počet cizích procesu na stanici, ...

Každá stanice aktualizuje v souboru ty záznamy, které se jí týkají. Při rozhodování o cíli migrace se soubor použije jako zdroj informací.

Toto je příklad situace, kdy jsou data relevantní pro migraci uložena na jiném místě než rozhodovací mechanismus.

+ používají se již existující vyvinuté mechanismy (synchronizace akcí nad souborem)

- pád fileserveru, pomalost přístupu k datům

Centrální server

Modifikací předchozího modelu je varianta centrálního serveru. Ta přichází s tou změnou, že informace potřebné pro rozhodování jsou uloženy v paměti rozhodovacího serveru a server sám o sobě vykazuje jistou inteligenci. Mechanismus rozhodování tedy může být implementován jako součást serveru.

+ Komunikace je orientována spojově, inteligence, vyšší rychlost.

Distribuované servery

Každý server shromažďuje část informace a tyto servery mezi sebou komunikují.

- komunikacne náročné řešení, nebo na úkor komunikace zaostává informovanost jednotlivých serveru, náročnější udržování konzistence, výpočetní overhead

+ vyšší míra spolehlivosti

Je možné použít pravdepodobnostní model, který je implementován v systému MOSIX. Ten využívá stárnutí dat.

9.2 Přehled některých systému

V následující části probereme některé zajímavé aspekty implementace migrace procesu v operačních systémech DEMOS/MP, Charlotte, V, MOSIX a Sprite. Zmíníme se také o objektovém systému Emerald, který implementuje migraci objektu a některé zajímavé operace nad nimi.

9.2.1 DEMOS/MP

V systému DEMOS/MP, vyvinutém na univerzitě v Berkeley v r. 1983, je implementována meziprocesová komunikace pomocí linky spojených s procesy. Tyto linky spravuje jádro systému. Systém nabízí plnou transparentnost, k čemuž mu pomáhá uniformní komunikační interface, který je nezávislý na poloze komunikujících procesů. Jádro samo může posílat právy jako každý proces, což zjednodušuje přístup ke komunikaci ze strany jádra.

Přístup k předávání zpráv určených pro migrovaný proces během migrace, který byl autory zvolen, je forwarding.

Vlastní migrační síla je implementována v jádru systému.

Vzhled identifikace procesu:

Poslední známý počítač, kde proces běžel	Globální unikátní PID	Vytvářející počítač	Lokální unikátní PID
--	-----------------------	---------------------	----------------------

Fáze migrace

1) Negotiation - domluva

Zdrojová stanice se dotazuje cílové, zda přistupuje na migraci za daných podmínek uvedených ve zprávě. Odpoví-li cíl ano, je vše v pořádku, jinak se vybírá jiná cílová stanice.

2) Presun

U presunu je zajímavé sledovat, která stanice je v dané fázi presunu aktivní. Tyto informace budou uvedeny vždy v závorce.

?? Vyjmutí procesu (*Zdroj*)

Proces je označen jako „v migraci“ a vyjmut z run-queue.

?? Podrobnější informace pro cíl (*Zdroj*)

Posílá se například velikost procesu a další informace o alokovaných prostředcích.

?? Alokace stavu na cíli (*Cíl*)

Vytvoří se nové datové struktury pro proces na cílové stanici a alokují se potřebné prostředky. Nový proces vytvořený na cílové stanici má PID migrovaného procesu.

?? Presun stavu (*Cíl*)

Tohle je zajímavé, protože v tomto případě si k sobě cílová stanice proces „tahá“, protože je běžnější postup, kdy je cílová stanice pasivní a zdrojová naopak data „tlací“.

?? Presun adresového prostoru a paměti (*Cíl*)

?? Forward čekajících zpráv na cílovou stanici (*Zdroj*)

Posílají se jednak ty zprávy, které procesu přišly před započítím migrace a které už si nestacil vyzvednout a dále ty zprávy, které procesu přišly během migrace.

Před posláním těchto zpráv změní jádro adresu procesu ve struktuře identifikující proces na novou adresu.

(Prepíše se adresa posledního známého působíště procesu)

?? Vycištění stavu na zdroji (*Zdroj*)

Zdrojová stanice smaže veškeré struktury týkající se odmigrovaného procesu kromě informace o jeho novém působíšti. Tato je později využívána pro forwarding zpráv.

?? Restart procesu (*Cíl*)

Forwardování zpráv

Z hlediska forwardování jsou 3 druhy zpráv podle toho, jak (od-) migrovanému adresátovi přicházely.

1. Odeslané před migrací, ale nepřijaté před migrací.
- Přeneseny v rámci bodu 6 přesunu při migraci
2. Odeslané po migraci za použití staré adresy
- Jádro „starého“ počítače, tj. nekdejšího zdrojového počítače, tyto zprávy průběžně forwarduje na pozici, kterou má poznamenanou ve své informaci o odmigrovaném procesu (viz bod 6 migrace). Přestože se nabízejí jiné možnosti řešení (vracet tyto zprávy zpět - „adresát nenalezen“, nebo obnovit adresy u všech procesu komunikujících s procesem migrovaným), zvolil DEMOS/MP forwardovací strategii.
3. Odeslané po migraci s využitím nové adresy
- Tyto zprávy neciní problémy.

9.2.2 Charlotte

Systém Charlotte teží ze svého návrhu meziprocesorové komunikace - zasílání zpráv. Komunikace probíhá pomocí linku a je naprosto nezávislá na umístění procesu. Dalšími vlastnostmi IPC jsou možnost komunikovat přes více kanálu zároveň, možnost přerušit zprávu (message cancelation) a možnost přesouvat linky.

Této vlastnosti Charlotte využije tím, že po odmigrování procesu nezustanou na zdrojové stanici žádné reziduální dependence (v systému DEMOS/MP zbyla nová adresa procesu kvůli forwardování zpráv).

Návrh migrace v Charlotte

Návrháři systému se snažili o dosažení následujících cílů:

1. Migrační politiku implementovat jako uživatelský proces
migrační sílu zahrnout do jádra systému
2. Dosáhnout maximálního stupně transparence
3. Umožnit preempci procesu, který byl odmigrován na stanici, jejíž uživatel se poté vrátil
4. Snaha o maximální fault-toleranci.

Migrační strategie - sber statistiky

Uživatelský proces poverený vykonáváním migrační strategie za tímto účelem sbírá jisté informace

O počítači:

- ?? Počet procesu
- ?? Počet komunikacních linku na tomto počítači
- ?? Množina nejpoužívanějších komunikacních linku
- ?? Celkové využití CPU
- ?? Celkové využití síťové komunikace

O každém procesu:

- ?? Doba, po kterou proces již bežel
- ?? Stav procesu

?? Využití CPU

?? Využití síťové komunikace

Perioda sberu informací je 50-80 ms, informace se okolním počítačům oznamují jednou za 100 intervalu sberu (5-8 s).

Podpora migrace

V jádre systému jsou vycelenena 3 vlákna, která se o migraci starají.

1. Vlákno obhospodaruující „migrační“ volání jádra
2. Vlákno vykonávající jednotlivé fáze migrace
3. Vlákno podporující sber statistiky.

Interface jádra týkající se migrace obsahuje čtyři služby:

1. Zaciná/prestává se sbírat statistika
2. Migrate Out
3. Migrate In
4. Zruš migraci, je-li to možné

Prubeh migrace

Vlastní migrace probíhá ve třech fázích.

Proces je po celou dobu migrace suspendován, zamražen (frozen), interakce procesu s externím prostředím jsou odloženy.

1) Negotiation - domluva

Tato fáze se skládá z výměny informací mezi politickými procesy na zainteresovaných stanicích. Když se tyto dohodnou, politický proces na zdrojové stanici zavolá jádro „Migrate Out“.

?? Jaderné vlákno pošle na cílovou stanici podrobnější informace o procesu (rozvržení pametových segmentu, tabulka komunikacních linku, množina aktivních linku, využití CPU a síte).

?? Jsou-li na cílové stanici prostředky pro nový proces, pošle jádro na cílové stanici podrobné informace migračnímu politikovi cílové stanice.

?? Politik cílové stanice zavolá „Migrate In“

?? Probehne alokace prostredku

?? Kernel cílové stanice provede fork interního procesu určeného pro příjem migrovaných procesu

?? Poté se na zdrojovou stanici posílá odpověď. Bud commit (migrace akceptována) nebo refused.

Před odesláním commitu je možné migraci prerusit.

2) Vlastní presun procesu

Odehrává se ve třech krocích:

?? Presun obrazu procesu

obraz se na cílové stanici uskladní ve strukturách pripravených v bode 1) migrace

?? Nastavení komunikačních linku

Kernely na druhých koncích všech linku jsou obeznámeny s novou pozicí migrovaného konce spoje. Dojdou-li potvrzení na všechny zprávy, může se stehovat.

Všechny zprávy, které až doposud procesu přišly na zdrojovou stanici, se tam ukládaly do bufferu. S procesem se na cílovou stanici přesunou pouze hlavičky těchto zpráv bez vlastního datového obsahu. Až bude chtít v budoucnu tyto zprávy proces přijmout, jádro nasimuluje buffer-cache-miss a obsah zprávy se přenesou.

Od této chvíle může nová stanice přijímat zprávy pro migrovaný proces, který ovšem ještě nebyl přenesen. Tyto zprávy jsou na cílové stanici bufferovány.

Behem nastavování linku se objevuje jejich nekonzistence, protože některé linky jsou již nastaveny a jiné ještě ne. V důsledku toho chodí zprávy určené procesu jak na zdrojovou stanici tak i na stanici cílovou.

?? Presun stavu

Stehování provádí zdrojová stanice (narodil od DEMOS/MP)

Presunují se deskriptory procesu, komunikačních linku, událostí a zpráv. Tyto se odesílají v paketech.

3) Vycištění zdrojové stanice - Clean-up

?? Odesílatel smaže všechny datové struktury týkající se odmigrovaného procesu (až na několik zpráv, které se smažou až po jejich úplném přenesení na cílovou stanici).

?? Příjemce rozebalí přijaté informace a zaradí proces do run-queue.

9.2.3 V

V je systém, který byl vyvinut na univerzitě ve Stanfordu.

Základní vlastnost systému, důležitá z hlediska migrace, je síťově transparentní exekucní prostředí a meziprocesová komunikace (mírné komplikace působí hardwarové device drivers).

Návrh migrace

V systému V není předmětem migrace proces, ale „logical host“, který proces obsahuje. Logical host je adresový prostor, kde může být několik V-procesu.

V systému V se návrháři snažili minimalizovat dobu, po kterou je migrovaný proces zamražen. Tuto dobu je výhodné minimalizovat kvůli možným timeoutům procesu čekajících na odpověď od migrovaného procesu.

Dobu migrace se podařilo minimalizovat na úkor celkové doby migrace. Mechanismus, který byl použit, se nazývá „precopying“.

Při migraci dochází k atomickému přesunu procesu. Význam je ten, že proces není nikdy viditelný v obou instancích (nové i staré)

Meziprocesová komunikace má schopnost zotavení ze ztráty zprávy. Této vlastnosti se při migraci využívá.

Mechanika migrace

Systém V implementuje migraci ve čtyřech krocích.

1) Inicializace cíle

?? V cílové stanici se vytvoří nový logical host (s jiným identifikátorem, což umožňuje paralelní existenci starého i nového logického hostu. (na konci fáze 3 se nové ID přepíše na původní).

2) Precopying stavu

?? V této fázi proces na zdrojové stanici, který běží s vysokou prioritou (ve smyslu častějšího přístupu k CPU), provádí kopírování stavu migrovaného procesu, zejména jeho adresového prostoru. Migrovaný proces stále běží, takže se může stát, že po zkopírování nějaké oblasti paměti dojde k její změně. Tato oblast se musí přenést znovu. Proces probíhá do té doby, než zbyde malé množství paměti, kterou je třeba přesunout.

3) Dokončení kopírování

?? Migrovaný logický host je zamražen

?? Dokončí se kopírování posledních zmodifikovaných oblastí paměti

?? Ošetří se problémy se zprávami od ostatních procesů. Zde se rozlišuje podle toho, zda byla zpráva žádostí či odpovědí.

?? Žádosti - schovávají se u adresáta (migrovaného procesu) ve frontě. Ta se ovšem po migraci maže, tedy odesílatel musí poslat svoji žádost znovu na novou stanici (po migraci je nová adresa broadcastem oznámena všem).

?? odpovědi - jsou okamžitě zahozeny

?? Zahazování zpráv je umožněno díky vlastnostem IPC.

?? Kopírování stavu z jádra zdrojové stanice a z manažera programu na zdrojové stanici

?? Změna identifikátoru logického hostu, tedy vzniknou dvě identické kopie log. hostu.

4) Odmražení, přesmerování odkazu

?? Smaže se stará kopie log. hostu, nová se odmrazí.

?? PID se naváže na novou kopii log. hostu.

?? Provede se navázání na nový počítač (zejména cache)

?? Broadcast nové adresy log. hostu.

Spice a Accent

Tyto systémy jsou zajímavé tím, že (podobně jako systém V) přesouvají adresový prostor v situaci, kdy migrovaný proces běží. Systém V používá precopying, systémy Spice a Accent si vybraly opačný přístup - copy on reference.

Tento přístup podle měření autora systému Spice a Accent ušetří přesun 21% až 96% adresového prostoru, ale za cenu velkého množství reziduálních dependencí, mnohdy rozprostřených po více stanicích (například když proces „obehne“ celou síť, na každé stanici zanechá část svého virtuálního adresového prostoru).

9.2.4 MOSIX

Operační systém MOSIX je systém UNIX adaptovaný na distribuované prostředí. Migrace procesu je v systému rozsáhle využívána, protože je to základní prostředek pro vyvažování zátěže.

Protože není v systému UNIX migrace procesu standardně implementována, museli tvůrci systému MOSIX provést některé změny oproti standardní implementaci UNIXu. Zajímavé jsou zejména ty změny, které se týkají procesu samotného, tedy změny ve struktuře, v níž si systém uchovává informace o procesu.

Ve struktuře přibýly následující položky:

?? Čas, který uplynul od posledního pokusu odmigrovat daný proces

?? Cas strávený na současném procesoru.

?? Příčina migrace

- ?? explicitní žádost voláním `migrate()` ;
- ?? vyvažování zátěže - load balancing
- ?? intenzivní vzdálené I/O operace
- ?? rozsáhlé vytváření nových procesů voláním `fork()` ;
- ?? shutdown stanice, na které proces běží

?? Nové flagy procesu

- ?? SMOVE - v migraci
- ?? SHERE - nesmí odmigrovat

?? Počet vytvořených procesů - dětí a informace o posledních operacích `fork()`

?? Statistika meziprocesové komunikace

Návrh migrace

Každý proces má svůj domácí procesor, tj. procesor, na kterém byl proces vytvořen. Není možnost vzdáleného `forku`. Při vytváření procesu se nový vytvoří vždy na procesoru, na kterém zrovna běží otec.

Migrace procesu je implementována voláním jádra `migrate(destination, lock)`.

`destination` - cílový procesor

`lock` - provést uzamčení procesu na cílovém procesoru?

Toto volání má dva speciální případy. Prvním z nich je změna stavu uzamčení na současném procesoru. Tehdy nedojde k migraci, proces se pouze uzamkne/odemkne. Druhý speciální případ je, když o migraci rozhodne samotné jádro systému.

Stejný postup je používán jak pro migraci mezi stanicemi v rámci sítě, tak i pro migraci procesu mezi procesory jedné stanice.

Vlastní migrace

Zdrojová stanice volá rutinu `passto()`, která přenesení proces na cíl. Tato rutina provede následující kroky:

(Rutiny, jejichž názvy začínají S znamenají, že jde o komunikaci RPC mezi stanicemi)

- 1) Testuje se, je-li možné proces odmigrovat
- 2) Volání rutiny `Smakeproc()` ze zdrojové stanice na cílovou stanicí.
 - ?? Cílová stanice zjistí, jestli si může dovolit příjem procesu
 - ?? Zajistí opravu komunikačních kanálů
 - ?? Kontrola, zda na cílové stanici neběží proces, který má stejné PID jako proces migrovaný (to by znemožnilo migraci, ale není to tak častý případ, může nastat v důsledku pádu nějaké stanice)
 - ?? Nastaví se pametové oblasti
 - ?? Čeká na doručení procesu. Až bude proces doručen, rutina jej probudí a vrátí se.

- 3) Jestliže se migrace provádí v důsledku vyvažování zátěže, bylo vyvažování po dobu provádění prvních dvou kroků vypnuto. Teď se znovu zapíná.
- 4) Přesun dat procesu - přesun adresového prostoru.
- 5) Volání `Spasproc()` na cílovou stanici. Toto volání nastaví dosud nenastavené položky ve struktuře procesu.
- 6) Zdrojová stanice volá na cílovou stanici `Sactivate()`. Toto volání rozbehne odmigrovaný proces.
- 7) Na zdrojové stanici se smažou lokální data týkající se procesu a uvolní se paměť, kterou proces zabíral.

9.2.5 Sprite

Operační systém Sprite je vybudován na jádře podobném 4.3 BSD UNIXu. Jádra systému na jednotlivých stanicích spolu komunikují pomocí RPC.

Předpoklady, ze kterých vyšli tvůrci systému Sprite, když začínali navrhovat migraci procesu, byly zejména tyto:

- ?? V síti, na které pracovali, bylo po většinu času mnoho „idle“ stanic, jejichž výpočetní výkon ležel ladem.
- ?? Přestože je hezké využít výkon „idle“ stanic, jakmile se vrátí vlastníci stanic, je potřeba uvolnit výpočetní sílu stanic pro něj, tj. odmigrovat všechny cizí procesy. Stejně tak jako cíl migrace je možné použít pouze stanici, o které se ví, že je „idle“.
- ?? Velká většina programu běží velmi krátce, tedy migrace procesu musí být provedena velmi rychle.
- ?? Sprite je systém založený na UNIXu a na systémových voláních (narozdíl od systému Charlotte, který je založen na komunikaci pomocí posílání zpráv).
- ?? Protože již existovala jistá podpora využití zdroje systému po síti, snažili se tvůrci migrace maximálně tuto podporu využít (přístup k souborům, zařízením, „svetové“ identifikace procesu, ...)

Cíle, které si tvůrci vytyčili, byly hlavně rychlost (už byla zmínovaná výše) a dosažení maximální míry transparency.

Pro každý proces je definována domácí stanice, což je ta stanice, na které by proces bežel, kdyby migrace v systému neexistovala. Je to stanice, která o procesu stále ví, tj. pořád má informaci o momentální pozici procesu.

Návrh migrace

Stejně jako u předchozích systémů, i v systému Sprite se návrháři rozhodli implementovat maximální množství funkcí mimo jádro systému. Z toho vyplynulo již standardní rozložení funkcí.

Vlastní migraci provádí jádro systému, sber statistiky, sledování činnosti uživatele na stanici a rozhodování provádí uživatelská aplikace - load average daemon.

Migrace procesu je část systému, která zasahuje svým vlivem téměř do všech částí systému. S tímto problémem se setkali i návrháři systému Sprite. Zjistili, že změny, ke kterým docházelo v průběhu vývoje migračního mechanismu, způsobují nekonzistenci jednotlivých verzí jader (už jsme mluvili o tom, že část implementace migrace je obsažena v jádře). V důsledku toho se objevily problémy s komunikací odlišných verzí migračních mechanismů. Návrháři zavedli číslování verzí jader podle

zmen v migračním mechanizmu. Pouze jádra se stejným číslem verze mohou mezi sebou vymenovat procesy.

Původní idea návrháru byla taková, že se budou všechna volání jádra, která proces provede na nové stanici, forwardovat na domácí stanici, kde budou vykonána. To zajišťuje vysokou míru transparency, ale za cenu reziduálních dependencí. Později od této ideje Sprite ustoupil a zvolil kombinaci všech tří možností pro přístup k prostředkům procesu po migraci (přenos, forward, podobný prostředek na cílové stanici).

Presun virtuální paměti

Některé možnosti presunu virtuální paměti byly rozebrány dříve. Operační systém Sprite zvolil variaci několika výše uvedených způsobů. Uvedená kombinace je spojení presunu veškeré paměti najednou a „copy-on-reference“.

V první fázi presunu virtuální paměti se všechny modifikované stránky uloží na fileserver. Potom se presune stav procesu na cílovou stanici. Paměť se potom způsobem copy-on-reference přesouvá na cílovou stanici z fileserveru tak, jak ji proces potřebuje.

Nevýhodou tohoto postupu je, že se paměť přesouvá po síti dvakrát. Podle autora systému Sprite však k migraci dochází většinou při volání funkce exec, takže efekt dvojího presunu paměti je v této souvislosti irelevantní.

Problémy v systému Sprite nastávají tehdy, když procesy sdílejí úsek paměti pro zápis. Procesy, které sdílejí paměť pro zápis, by se mohly přesouvat všechny zároveň, ale vzhledem k tomu, že mohou tvořit rozsáhlou a složitou skupinu, není tato možnost využita a platí zákaz migrace procesu sdílejících pamětí pro zápis. (Tuto situaci by mohlo vyřešit použití distribuované sdílené paměti, ale toho Sprite nevyužívá.)

Reziduální dependence a domácí stanice

Otázku reziduálních dependencí vyřešili návrháři systému Sprite tak, že pro každý proces je definována jeho domácí stanice, která má neustále o procesu přehled. Po odmigrování procesu z jeho domácí stanice na této zůstane jistá informace, ale po odmigrování procesu ze stanice jiné (např. v důsledku evikce) po sobě proces žádné reziduální dependence nezanechá.

Součástí stavu procesu je v systému Sprite i Process Control Block - PCB. Otázkou je, jak zaručit, aby proces měl svůj PCB na cílové stanici, protože na té stanici běží. Ale i na domácí stanici je potřeba udržovat informace o momentální poloze procesu. Za tímto účelem je PCB ve dvou kopiích. Jedna kopie je na stanici, kde proces momentálně běží, druhá je na jeho domácí stanici, přičemž domácí stanice má většinu položek v PCB nevyužitou.

Migrační rutiny

Stav procesu se skládá z mnoha různorodých součástí, z nichž každá vyžaduje specifický přístup. Například presun obsahu registru vyžaduje jiné zacházení než presun informace o otevřeném komunikačním kanálu.

Tvůrci migrace pro operační systém Sprite přemýšleli nad tím, jak v této různorodosti najít společné rysy a tím i společný přístup k zacházení s částmi stavu migrovaného procesu. Výsledkem jejich snažení se stal koncept **migračních rutin**.

Každá součást přenášeného stavu má definovány právě 4 migrační rutiny.

1. Predmigrační rutinu - **pre-migration routine**

Její úkolem je připravit danou část stavu na migraci. Návrátovou hodnotou je velikost dat, která se budou přesouvat mezi stanicemi. (Pro virtuální paměť je to například příprava modifikovaných stránek procesu k zápisu na disk.)

2. Zabalovací rutinu - **encapsulation routine**
Tato rutina zabalí svojí část stavu procesu do připraveného bufferu. Jejimi vedlejšími efekty jsou například komunikace s file-serverem v případě zabalování informací o otevřených souborech nebo zápis modifikovaných stránek adresového prostoru migrovaného procesu v případě zabalování virtuální paměti.
3. Rozbalovací rutinu - **de-encapsulation routine**
Je volána na cílové stanici, slouží k „nainstalování“ stavu primigrovaného procesu.
4. Pomigrační rutinu - **post-migration routine**
Zdrojová stanice ji volá po ukončení přenosu odmigrovaného procesu. Tato rutina vycistí danou část stavu procesu, uvolní alokovanou paměť, ...

Vlastní průběh migrace

1. Procesu, o kterém se rozhodlo, že se bude migrovat, je poslán **signál**. V důsledku toho proces musí provést trap do jádra a navíc je v tomto případě v dobře definovaném stavu (není uprostřed nějakého choulostivého jaderného volání).
2. Jedná-li se o migraci z domácí stanice ven na jinou stanici, dotazuje se domácí stanice cílové, zda tato žije a zda je možné na ni poslat migrovaný proces. Při této příležitosti si také vyžádá verzi implementace migračního mechanismu. Dopadnou-li odpovědi na **otázky** dobře (je možné migrovat), alokuje se na cílové stanici struktura PCB (Process Control Block) a na zdrojovou stanici se vrátí identifikátor nově vytvořeného procesu. Migruje-li se z některé stanice sítě na domácí stanici (například v důsledku evikce), odpadá většina otázek a také alokace PCB. Jako identifikátor procesu na cílové stanici (tedy na stanici domácí - proces migruje domů) se použije jeho normální PID.
3. Pro každou část přesunovaného stavu procesu se zavolá **pre-migrační rutina**. Tato rutina inicializuje jednotlivé části stavu pro budoucí přenos. Výsledkem jejího volání je velikost dat, která se budou přenášet.
4. Podle výsledku volání pre-migrační rutiny se na zdrojové stanici naalokuje **buffer**, do kterého se bude proces zabalovat. Velikost bufferu závisí zejména na velikosti paměti používané procesem - na velikosti jeho tabulek stránek.
5. Pro každou část stavu (modul) je volána příslušná **encapsulační rutina**. Tato rutina do bufferu alokovaného v kroku 4 zabalí příslušný modul stavu procesu.
6. Přesun stavu je realizován jedním RPC. V tomto okamžiku se pro každý modul na cílové stanici volá jeho **deencapsulační rutina**.
7. Na zdrojové stanici se opět pro každý modul volá **post-migrační rutina** (je-li pro modul definována). Tato rutina vycistí danou část stavu.
8. Na zdrojové stanici je provedena **dealokace** bufferu a cílová stanice se informuje o tom, že proces může být spuštěn.

Transparence

Možností, jak dosáhnout transparence, využili návrháři Spritu několik.

1. Kde to bylo možné, tam zmenili volání jádra takovým způsobem, aby byla nezávislá na tom, zda je proces prováděn na domácí stanici či někde jinde. K tomu pomohl jednotný prostor jmen, jednotný filesystém, „svetové“ identifikátory procesu, ... Odpovídá to vlastně použití podobného prostředku na cílové stanici.
2. Kde nebylo možné použít první možnost, snažili se přesunout danou část stavu na cílovou stanici.
3. V případě, že selhaly snahy o obe předchozí varianty, nezbylo, než provádět forwardování akcí na domácí stanici - slepé forwardování.

Tímto způsobem je kuriózně ošetřeno například volání `gettimeofday()`, neboť v systému Sprite neexistuje rozumná synchronizace hodin jednotlivých stanic.

Forwardování se neprovádí pouze ze stanice, kde proces běží, na stanici domácí, jak je tomu v případě `gettimeofday()`, ale i ze stanice domácí na stanici, kde se proces nachází (například signály, ...).

4. Jistá skupina volání musela být ošetřena vzájemnou kooperací obou stanic. Jak té domácí, tak té, na které proces běží. Příkladem budiž volání `fork()`, `exit()`.

Ze 106 volání jádra systému Sprite je 91 volání ošetřeno způsobem 1) nebo 2), 11 volání je slepe forwardováno a zbylá čtyři se řeší komplikovanějším způsobem kooperace obou zainteresovaných stanic.

9.2.6 Migrace objektu - systém Emerald

Systém Emerald byl vyvíjen zároveň s programovacím jazykem Emerald. Jednotkou manipulace je objekt. Objekt se skládá ze čtyř částí:

1. „Svetove“ (tj. v rámci sítě) unikátní jméno.
2. Reprezentace - data. Data jsou buď primární (skutečný datový obsah objektu), nebo reference na jiné objekty.
3. Operace, které je možné na objektu provádět.
4. Proces, který je k objektu přiřazen. Je-li přiřazen, jedná se o aktivní objekt, v opačném případě je objekt pasivní.

Koncept mobility je obsažen v programovacím jazyku Emerald a v behové podpoře. Jazyk dává k dipozici následující primitiva:

?? LOCATE objekt nalezne momentální polohu objektu

?? MOVE objekt TO uzel presun objektu

?? FIX objekt AT uzel upevní objekt na uzlu

?? UNFIX objekt uvolní objekt

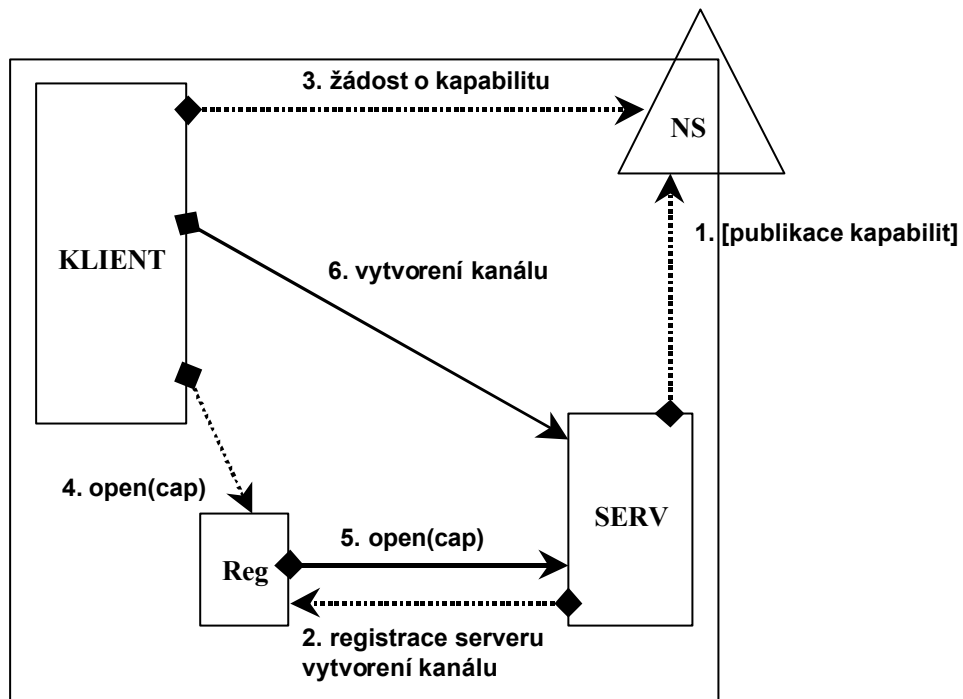
?? REFIX objekt AT uzel atomické provedení
UNFIX objekt;
MOVE objekt TO uzel;
FIX objekt AT uzel;

Lze explicitně říci, co všechno se má presunout. Stací společně svázat objekty, které se potom presunují najednou.

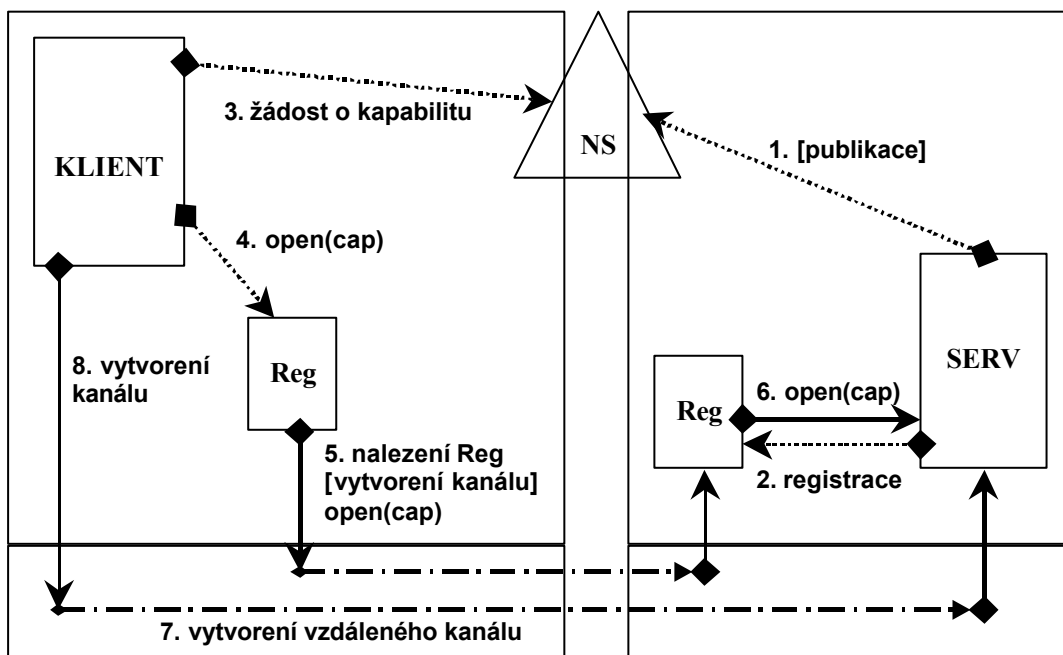
Presun objektu

1. Zdroj vyrobí zprávu, která obsahuje objekt a tabulky pro překlad adres závislých na poloze objektu.
2. Zpráva se odešle na cíl.
3. Cíl alokuje paměť, zkopíruje data a vyrobí vlastní translacní tabulku podle té, co obdržel.

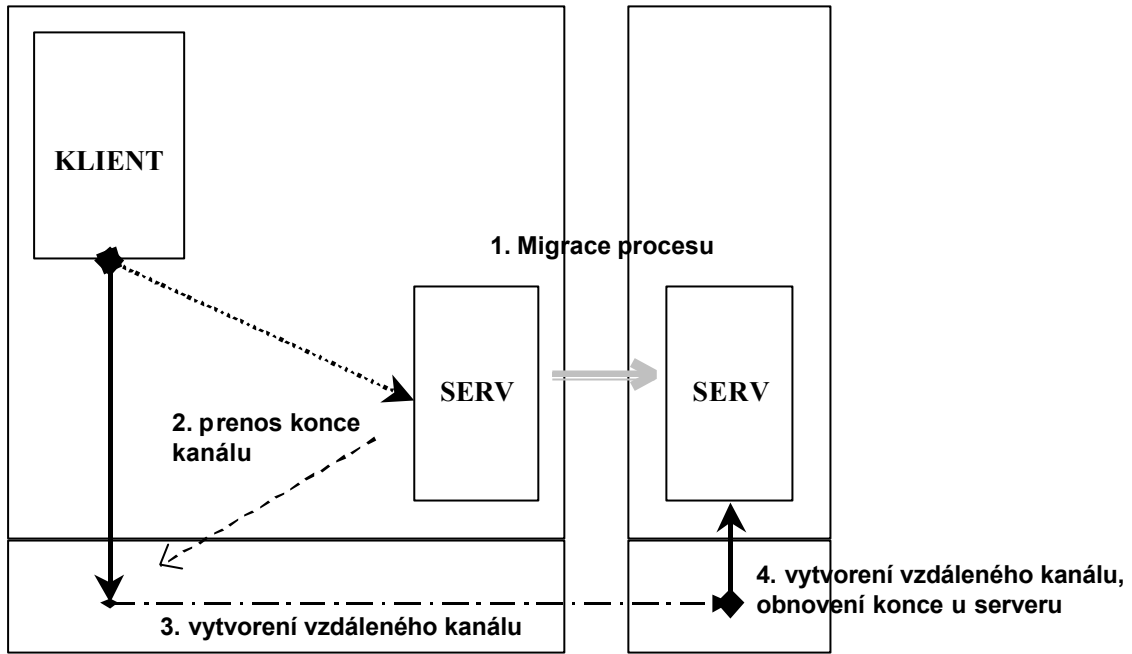
9.3 Komunikace při migraci



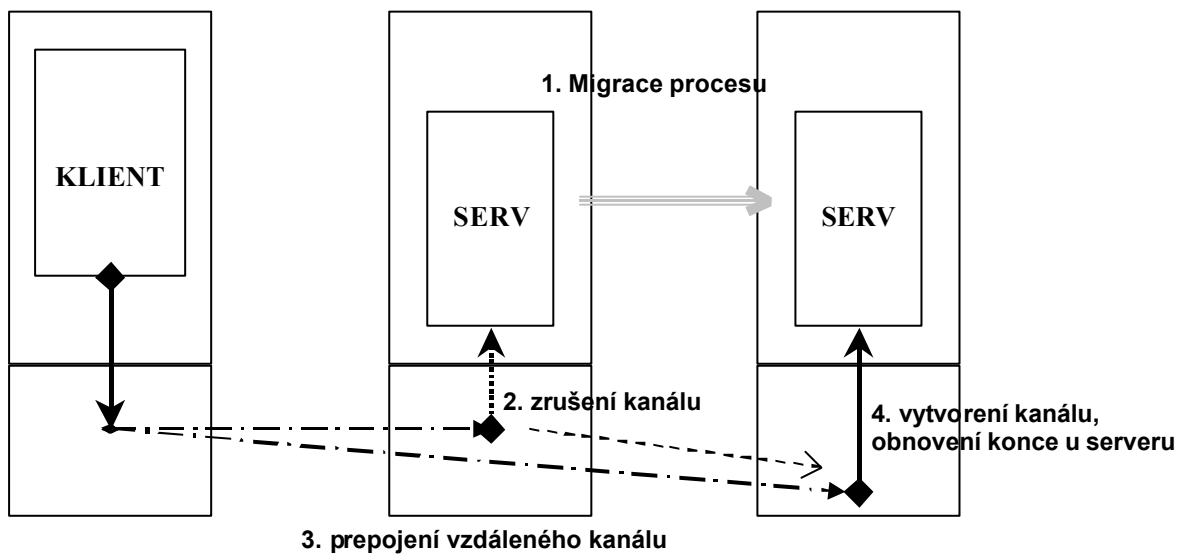
Obr. 49 - Lokální komunikace klient - server



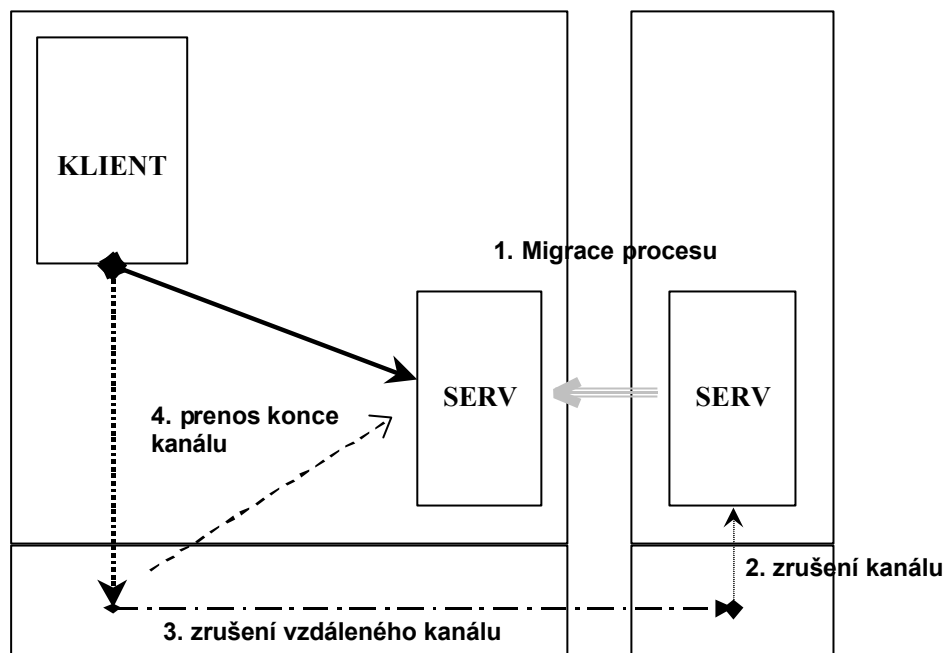
Obr. 50 - Vzdálená komunikace klient - server



Obr. 51 - Migrace serveru vzhledem k lokálnímu klientu



Obr. 52 - Migrace serveru vzhledem k vzdálenému klientu



Obr. 53 - Migrace serveru ke klientu

9.4 Vyvažování zátěže

Kontinuální prerozdělování práce mezi procesory, dlouhodobé distribuované plánování, load balancing.

Otázky k vyřešení:

1. Rozhodnutí o okamžiku migrace
2. Jak porovnávat zatížení procesoru
3. Udržování konzistence údajů
4. Volba migrujícího procesu
5. Volba příjemce
6. Přenos procesu
7. Vzdálený běh procesu

9.4.1 Párový algoritmus (Bryant & Finkel)

Dynamicky se vytvářejí páry, které se vzájemně vyvažují

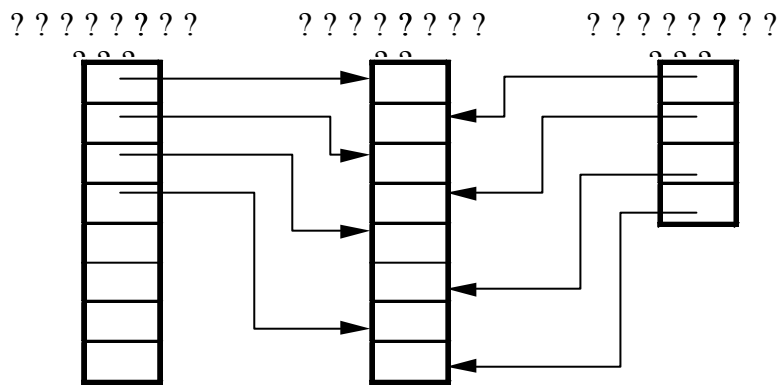
1. Počítac A pošle nějakému sousedovi B žádost o vytvoření páru. Žádost obsahuje výpis běžících procesů a jejich základní charakteristiky.
2. B po přijetí zprávy:
 - (a) odmítne žádost; A musí zkoušet jinde
 - (b) vytvoří s A pár; dokud nebude tento pár zrušen, A i B odmítají další žádosti
 - (c) oznámí počítači A, že právě migruje; A čeká, až se B rozhodne

- Po vytvoření páru zatíženější procesor (dejme tomu A) vybere migrující proces podle míry vylepšení $k_i = T_{Ai} / (T_{Bi} + T_{ABi})$, kde
 T_{Ai} je očekávaný čas odpovědi procesu i na počítači A
 T_{Bi} je očekávaný čas odpovědi procesu i na počítači B
 T_{ABi} je čas přenosu z A na B
- V případě, že migraci nelze dosáhnout (významného) zlepšení, pár je zrušen. Jinak je vybraný proces přenesen a pokračuje se od bodu 3.

9.4.2 Vektorový algoritmus (Mosix)

Každý počítač má pevný malý vektor zátěže L ostatních počítačů, každý počítač sám periodicky měří vlastní zátěž ($L_0 =$ vlastní zátěž). Každý časový interval každý počítač provede:

- Zjistí vlastní zátěž
- Zvolí náhodně jiný počítač a pošle mu první polovinu svého vektoru
- Při příjmu části vektoru zátěže (L') každý počítač provede $L_{2i} = L_i$, $L_{2i+1} = L'_i$
- Při požadavku na migraci se k zátěži jednotlivých procesorů přičte komunikační režie a vybere se nejvýhodnější procesor



Obr. 54 - Výpočet nového vektoru zátěže

9.4.3 Bidding algoritmus (Stankovic & Sidhu)

Využívá McCulloch-Pittsovu vyhodnocovací proceduru. Tato metoda je založena na tzv. vyhodnocovací bunce, která má excitátory, inhibitory a jednohodnotový výstup. Výstup je součet hodnot excitátoru anebo nula v případě nastavení libovolného inhibitoru.

- Všechny procesy jsou pravidelně vyhodnocovány, vstupy jsou vlastnosti procesu a sítě.
- Jestliže výstup je vyšší než prahová hodnota, proces je v pořádku. Jestliže je výstup nižší, ale nenulový, měl by být přenesen. Jestliže je výstup nula, je některý z inhibitorů nastaven a proces nemůže být přenesen.
- Jestliže má být přenesen alespoň jeden proces, počítač vyšle až do vzdálenosti d broadcast typu 'žádost o nabídku' (RFB) obsahující charakteristiky nabízeného procesu. Každý počítač, který obdrží RFB, tento proces ohodnotí a v případě nadprahové hodnoty vrátí odeslateli odpověď nabídku.

4. Po určité době t se všechny odpovědi zkorigují o cenu přenosu a nejlepší nabídka je považována za potenciálního příjemce procesu. Jestliže po dobu t žádná nabídka nedošla, zvětší se vzdálenost d a opakuje se bod 3.

Problém: obtížná kvantifikace vlastností procesu

9.4.4 SLA algoritmus (Stankovic)

Stochastic learning automata - zpetné ucení

9.4.5 BDT algoritmus (Stankovic)

Bayesian Decision Theory - posílání globálních stavu

9.4.6 Centralizovaný algoritmus

Jeden centrální server, zná zátěž všech počítačů, velí

9.4.7 Lokální algoritmus

Každý počítač zjišťuje pouze lokální zátěž. Jestliže ta překročí prahovou hodnotu, náhodně se vybere n počítačů, tem se zašle žádost o migraci a první / nejlepší odpověď se akceptuje.

9.4.8 Porovnání distribuovaných vyvažovacích algoritmu

Z hlediska efektivity a režie vlastního algoritmu je pro reálné distribuované systémy nejzajímavější vektorový algoritmus, který umožňuje získat globálnější představu o zátěži systému s dostatečně malou režii. Ostatní algoritmy se ukazují být pro reálné použití příliš náročné.

Důležitým hlediskem je i komunikační složitost jednotlivých algoritmu. U párového a bidding algoritmu počítače posílají detailní popisy jednotlivých procesu. To zvyšuje zátěž sítě a při větším množství počítačů síť zahlcuje. Naopak ostatní algoritmy (vektorový, SLA, BDT) přenášejí pouze globální stavovou informaci, čímž jsou z hlediska komunikační složitosti efektivnější.

Dalším problémem je aktualizace údajů. Pouze párový algoritmus zajišťuje, že rozhodnutí jsou prováděna na základě aktuálních informací - počítače komunikují po párech. Zajišťuje to ovšem pouze lokální optimalizaci.

10. Správa prostředku

Prostředkem (resource) rozumíme jakoukoliv relativně stabilní hardwarovou nebo softwarovou část systému využitelnou uživateli, resp. jejich procesy. Prostředky můžeme rozdělit na:

1. **fyzické prostředky**: permanentní fyzické části systému - procesor, paměť, disky, I/O zařízení, interní zařízení (hodiny, časovace), apod.
2. **logické prostředky**: softwarové entity - procesy, soubory, sdílená data, apod.

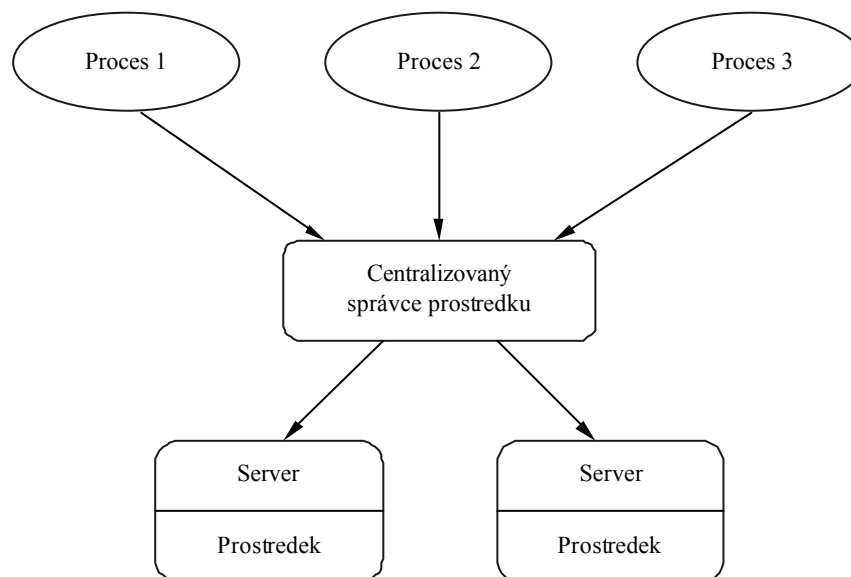
V distribuovaných systémech jsou prostředky fyzicky decentralizované, stavová informace je rozptýlena mezi jednotlivými uzly. Některé prostředky mohou být také z důvodu vyššího výkonu systému replikovány na více uzlech. To vše má za následek vyšší režii správy prostředku v distribuovaných systémech oproti centralizovaným systémům.

10.1 Správci prostředku

Prostředky jsou většinou pasivní, jsou spravovány servery - správci prostředku (resource managers). Správce prostředku má za úkol:

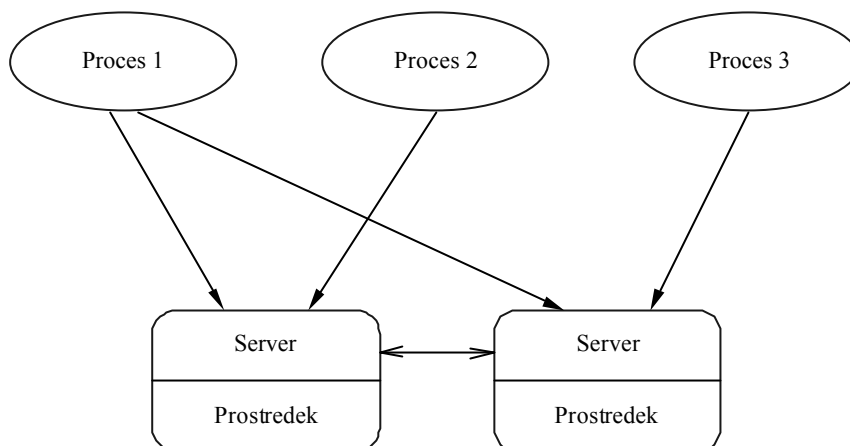
1. Udržovat spravované prostředky - lokaci, konzistenci, . . .
2. Pridělovat prostředky uživatelským procesům na jejich žádost
3. Dbát o ochranu prostředku před nepovolaným zásahem

10.1.1 Centralizovaná správa prostředku



Obr. 55 - Centralizovaná správa prostředku

10.1.2 Distribuovaná správa prostředku



Obr. 56 -Distribuovaná správa prostředku

10.1.3 Správa prostředku pomocí agentu

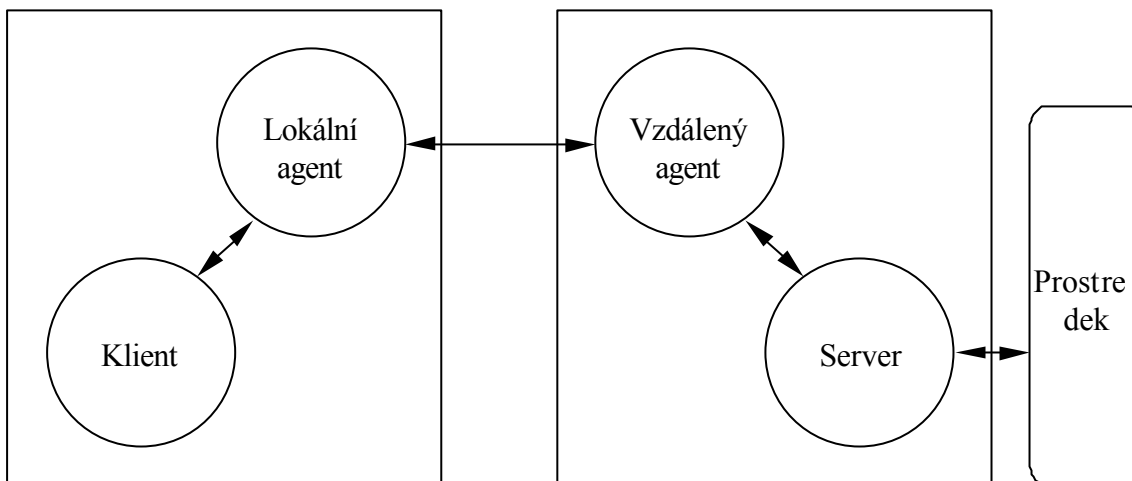
Uživatelské procesy komunikují se servery prostřednictvím agentu. Lokální agent komunikuje přímo s uživatelským procesem (klientem), vzdálený (remote) agent komunikuje s vlastním serverem poskytujícím služby nějakého prostředku.

Úkoly lokálního agenta:

1. Přijímat od klienta požadavky na prostředky
2. Lokalizovat vhodného vzdáleného agenta
3. Vyjednávat s vzdáleným agentem o zapůjčení prostředku
4. Vyvolat služby požadované klientem
5. Spravovat výjimečné stavy - odmítnutí požadavku, migraci prostředku, havárie apod.

Úkoly vzdáleného agenta:

1. Přijímat žádosti klientských agentu
2. Rozhodovat o vyhovění žádosti
3. Vytvorit případné exekucní prostředí pro server - nastavení práv, navázání spojení s klientem apod.
4. Spravovat výjimečné stavy



Obr. 57 - Správa prostředku pomocí agentu

10.2 Zablokování (deadlock)

Zablokování (deadlock) v distribuovaných systémech je podobný problém jako v centralizovaných systémech; ba ještě větší. Hůře se jim předchází, hůře se odstraňují, hůře se dokonce i detekují, protože potřebná informace je rozprostřena mezi několika počítači.

V literatuře jsou někdy rozlišeny dva druhy deadlocku - komunikační deadlock a deadlock prostředku. Komunikační deadlock nastává např. v případě, že proces A vyšle zprávu procesu B a čeká na odpověď, B pošle zprávu procesu C a C pošle zprávu zpět procesu A. Protože však proces A čeká na odpověď od procesu B, zprávu od procesu C nepřijímá, tudíž všechny procesy jsou zablokovány. Vzhledem k tomu, že oba druhy deadlocku mají stejné chování a navíc procesy lze považovat za speciální druh prostředku, nebudeme dále tyto dva druhy rozlišovat.

Podobně jako v klasických nedistribučovaných systémech je i v distribuovaných systémech nejpopulárnější a nejpoužívanější tzv. pštosí algoritmus, tj. problém ignorovat a případné řešení deadlocku nechat na uživateli, resp. na uživatelské aplikaci (např. databáze), která si vše hlídá sama.

Dalšími způsoby vypořádání se s deadlocky jsou detekce a zotavení (vznik deadlocku je připuštěn, po jeho zjištění je řešen násilím) nebo prevence (vznik deadlocku je vyloučen).

10.2.1 Modely deadlocku

Graf prostředku (WFG, wait-for-graph)

- ?? Model s jedním prostředkem - proces čeká pouze na jeden prostředek
- ?? Konjunkční (and) model - musí být uspokojeny všechny požadavky
- ?? Alternacní (or) model - musí být uspokojen alespoň jeden požadavek
- ?? Model vícenásobných prostředků (multiple resources) - musí být uspokojeno n požadavků z k dostupných prostředků

10.2.2 Algoritmy detekce deadlocku

Bylo vyvinuto mnoho algoritmu, avšak mnohé z nich byly špatné, a to i přes to, že byla autory “dokázána” jejich správnost. Chyby v důkazech byly v naprosté většině způsobeny absencí společné paměti a komunikačními prodlevami.

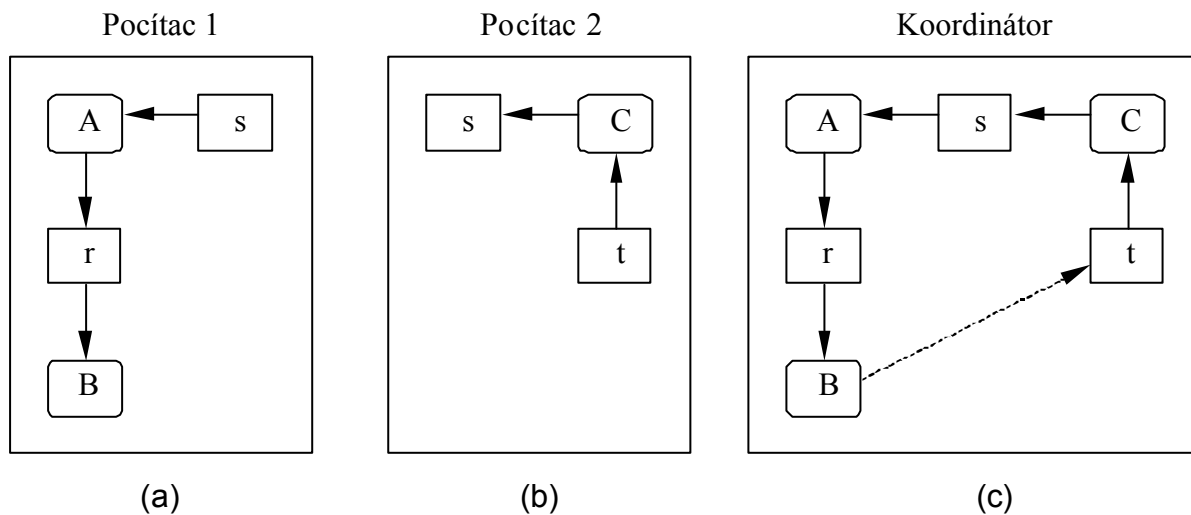
Centralizovaný algoritmus

Existuje centralizovaný server, který udržuje graf prostředku. Jednotlivé počítače informují server:

- ?? Po každé změně lokálního stavu
- ?? V pravidelných intervalech
- ?? Na požádání serveru

Problém - mohou vznikat tzv. falešné deadlocky kvůli zpoždění zpráv. Příkladem může být situace na obrázku. Na počítači 1 proces B drží prostředek r, proces A drží prostředek s a čeká na r. Na počítači 2 proces C drží prostředek t a čeká na s. Proces B může provést uvolnění r a požadavek na t. Zpráva od prostředku t (na počítači 2) však může koordinátoru dojít dříve než zpráva o uvolnění prostředku r (na počítači 1) a koordinátor může zjistit deadlock, i když ve skutečnosti k deadlocku nedošlo.

Řešení - globální timeordering, server se při deadlocku zeptá, jestli někdo nevyslal zprávu, která by předcházela zprávě o požadavku, který zavinil deadlock.



Obr. 58 - Centralizovaný algoritmus

(a) Lokální stav na počítači 1 (b) Lokální stav na počítači 2

(c) Možný stav koordinátora po akcích B.release(r), B.alloc(t) - falešný deadlock

Hierarchický algoritmus

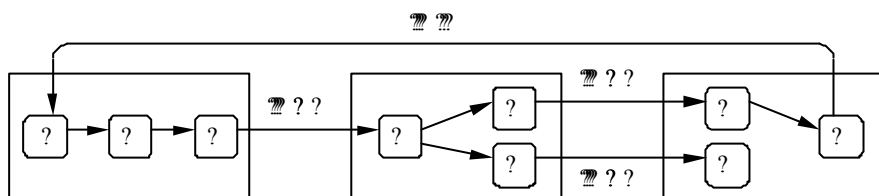
Centralizované řešení s jedním serverem není příliš vhodné (z mnohokrát opakovaných důvodů) pro rozsáhlejší systémy. Existuje proto třída algoritmu, které detekují deadlocky pomocí hierarchické struktury koordinátoru. Každý uzel řeší deadlocky lokálně, koordinátor nejnižší úrovně řeší deadlocky

jemu podřízených uzlu, a teprve když deadlock míří do jemu nepodřízených uzlu, předá řešení koordinátoru vyšší úrovně. Příkladem takového algoritmu je Menasce & Muntz '79.

Distribučovaný algoritmus - CMH

Příkladem distribučovaného algoritmu detekce deadlocku je algoritmus Chandy-Misra-Haas (CMH). Iniciátor pošle všem procesum, kterými je blokován, zprávu, která obsahuje identifikaci iniciátora, identifikaci odesílajícího procesu a identifikaci příjemce. Každý proces, který takovou zprávu obdrží, ji vyšle opět všem procesum, kterými je blokován, anebo zprávu zruší jestliže takové nejsou. Jestliže zpráva dojde zpět k iniciátoru, pak graf obsahuje kružnice a procesy jsou v deadlocku a třeba cyklus násilně přerušit.

Jestliže se několik procesů spontánně rozhodne být zároveň iniciátorem, algoritmus bude stále fungovat, avšak kdyby měl být při zjištění deadlocku zabít proces, který algoritmus inicioval, docházelo by zbytečně k zabití více procesů - overkill. Řešením tohoto problému mohou být jednoznačné časové značky přidělené procesum. Zabít potom nebude iniciátor, ale proces s největší časovou značkou, kterému iniciátor zašle zprávu o deadlocku.



Obr. 59 - Distribučovaný detekční algoritmus

Detekce globálního stavu

Detekce deadlocku může být založena na detekci globálního stavu systému. Algoritmus (Chandy & Lamport) nezjišťuje aktuální stav systému, ale pouze konzistentní stav. Graf prostředku (WFG) má však jednu podstatnou vlastnost: pokud vyloučíme násilné odebrání prostředku či ukončení procesu, pak

Jestliže WFG G a proces p je v deadlocku v WFG, pak je v deadlocku i v WFG'

Proto jestliže je zjištěna kružnice v WFG', který byl zjištěn algoritmem detekce globálního stavu, pak došlo k deadlocku a je třeba některý proces zabít.

10.2.3 Prevence deadlocku

Prevence deadlocku v distribučovaném systému je založena na stejném principu jako u centralizovaných systému - vytvořit takový systém alokace a držení prostředku, aby k deadlockum vůbec nemohlo dojít.

První použitelnou (alespon v některých systémech) metodou je jednoznačné uspořádání všech prostředku. Každý proces potom může alokovat prostředky pouze v rostoucí posloupnosti, tj. mohou být alokovány pouze ty prostředky, které jsou v daném uspořádání ohodnoceny výše, nežli všechny dosud naalokované prostředky daného procesu. Tímto způsobem nemohou vznikat cykly, tedy ani deadlocky.

Další dva způsoby prevence deadlocku jsou založeny na tom, že každé transakci je při startu přidělena jednoznačná časová značka. V okamžiku alokace prostředku držení jinou transakcí jsou porovnány časové značky obou transakcí. Zablokování je povoleno pouze starší transakci (tj. s menší

casovou značkou), mladší transakce je zrušena. Potom budou zretezeny transakce s rostoucí casovou značkou, nemuže tudíž dojít k deadlocku. Tento algoritmus bývá často nazýván wait-die (starší čeká, mladší umre). Druhým způsobem řešení konfliktu je ukončení transakce, která drží konfliktní prostředek. V případě, že prostředek drží starší transakce, mladší transakce čeká, v případě že prostředek drží mladší transakce, je tato ukončena a prostředek je přidelen starší transakci. Tento algoritmus je nazýván wound-wait.

	starší ? mladší	mladší ? starší
wait-die	starší čeká	mladší umre
wound-wait	mladší umre	mladší čeká

Tab. 18 - Prevence deadlocku

10.3 Ochrana prostredku

Prístupová matice - kapability / Access Control List (prístupový seznam)

Zabezpečení - kryptografie - verejné tajné klíče, ...

10.3.1 Klasifikace bezpecnosti DoD - the Orange Book

- ? **Trída D** - Nezabezpečené
- ? **Trída C** - Volitelná ochrana
 - ? **C1** - Zabezpečení objektu pred nepovolaným prístupem
 - ? **C2** - Login, izolace bezpecnostních mechanismu
- ? **Trída B** - Povinná ochrana - všechny objekty musí být nutne zabezpeceny
 - ? **B1** - Úrovnová ochrana - každý objekt má prirazenu klasifikaci úrovně, uživatel má prístup k objektum klasifikovaným menší úrovni ochrany, než jakou má uživatel
 - ? **B2** - Celý systém musí být napsán s využitím formálního bezpecnostního modelu, který musí být zdokumentován. Musí být identifikován každý kanál umožňující potenciálně únik informací a tento musí být zabezpecen
 - ? **B3** - Musí být implementován referenční monitor a bezpecnostní domény - seznamy uživatele a skupin s prístupovými právy k objektu a bez jakýchkoliv práv prístupu
- ? **Trída A** - Dokazatelná ochrana - vyžaduje kompletní formální dukaz bezpecnosti

11. Správa souboru

Správa souboru v distribuovaných systémech je implementována množinou file serveru, což jsou typicky běžné uživatelské procesy. V některých systémech jsou file servery nebo některé jejich části implementovány přímo v jádru systému. Jednotlivé file servery mohou poskytovat různé služby - v jednom systému mohou být vedle sebe jak server poskytující běžné souborové služby, tak i server poskytující bezpečné transakční služby. Podobně může být v jednom systému server umožňující přístup k UNIXovým souborům a zároveň server poskytující služby filesystému MS-DOSu.

Identifikační a adresárové služby		
Adresárový modul		Kontrola přístupu
Umístovací modul	Lokální modul	Relokační modul
Replikační služby		
Konzistence dat		Multiple copy update
Transakční služby		
Transakční operace	Modul zotavení	Kontrola konkurence
Souborové služby		
Přístup k souborům		Přístup k atributům
Blokové a diskové služby		
Modul správy bloku	Ovládací zařízení	Caching

Tab. 19 - Možná struktura distribuovaného filesystemu

servery dedikované / uživatelské

jeden fileserver / oddělené souborové a adresárové služby

servery stavové / bezstavové (stateles)

více serverů - soubory disjunktní / replikace

11.1 Diskové služby

11.1.1 Blokovaná struktura souboru

11.1.2 Systémová identifikace

UFID - identifikace serveru ? interní číslo souboru (? zabezpečení)

UFID ? číslo bloku ? lokální / vzdálený blok

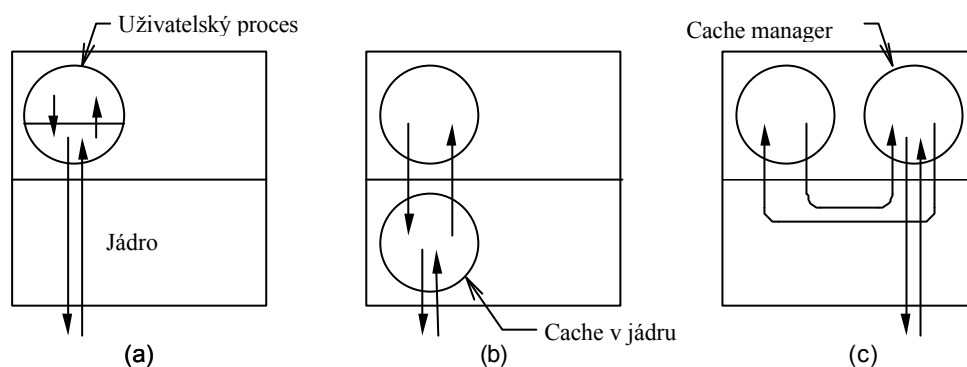
11.1.3 Vyrovnávací paměť (cache)

Cache může být v zásadě umístěna na několika různých místech: v paměti serveru, v paměti klienta nebo na disku klienta.

Cache v paměti serveru snižuje počet přístupů na disk. Vzhledem k tomu, že paměť fileserveru bývá až řádově menší nežli kapacita disku, je třeba rozhodnout, co bude umístěno v cache. Důležitým rozhodnutím musí být, jak velké kusy dat bude cache udržovat. Jednou možností je udržovat celé soubory. V tomto případě může být přístup na disk prováděn najednou, anebo po velmi velkých kusech, a tudíž (relativně) velmi rychle. Druhou možností je cachovat soubory po blocích nebo stránkách. Toto řešení má naopak výhodu v tom, že paměť dostupná pro cache je daleko lépe využita. Cache v paměti serveru je pro klienty zcela transparentní - server si sám serializuje požadavky a provádí přístup na disk, klient se o kopii souboru v paměti vůbec nedozví.

Aby se snížila zátěž sítě, musí být cache na klientské straně. Zde jsou principálně dvě možnosti, kam cache umístit - do paměti anebo na disk. Cachování souboru na lokálním disku je výrazně pomalejší a je proto efektivní pouze u přístupu po pomalých sítích (např. WAN). Na druhou stranu má tu výhodu, že umožňuje cachovat relativně velký počet velkých souborů.

Pro umístění cache v paměti klienta jsou v zásadě tři možnosti: v adresovém prostoru procesu, v jádru, nebo jako specializovaný cache manager.



Obr. 60 - Caching v paměti klienta

(a) Cache v adresovém prostoru procesu (b) Cache v jádru (c) Cache manager

Protože cache způsobují dočasnou nekonzistenci mezi daty udržovanými lokálně a daty na fileserveru, případně daty u jiných klientů, je třeba nějakým způsobem konzistenci zajistit. Jedním řešením je použít write-through cache - v případě zápisu do cachovaného souboru se zápis pošle zároveň i fileserveru. To však neresí problém více klientů přistupujících ke stejnému souboru. Navíc tato metoda nijak nešetří provoz sítě při zápisu do souboru.

Jinou metodou je posílat zápisy serveru ne při každé změně dat, ale až po určité době - tzv. delayed-write. Tento způsob mimo jiné efektivně řeší problém dočasných souborů. Soubor, který je třeba jen krátkou dobu na uchování mezivýsledku, tak nemusí být vůbec na fileserver přenášen. Odkládání zápisu má však jednu nevýhodu - zvláštní sémantiku. To, co jiné procesy vidí, silně závisí na case. Dalším řešením konzistence cachí je write-on-close, tj. zápis až při zavření souboru.

Zcela jiným způsobem řešení konzistence je centralizovaná kontrola. Fileservr má odkazy na všechny cachované soubory a konzistenci řeší sám.

11.2 Souborové služby

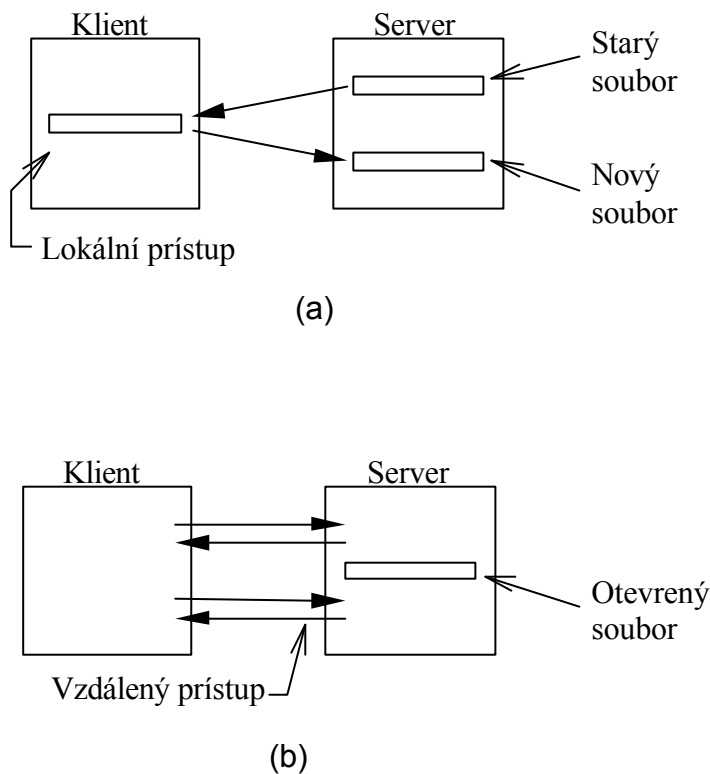
struktura souboru

?? posloupnost bajtu

?? záznamy

- ?? typovaný
- ?? virtuální adresový prostor

- atributy - per soubor, ne per adresář entry
- mutabilní / imutabilní
- kapability / přístupové seznamy (Access Control List)
- upload model / remote access model



Obr. 61 - (a) Upload/download model (b) Vzdálený přístup

11.2.1 Sémantika sdílení souboru

Jestliže více procesů současně přistupuje k některým souborům, je třeba přesně definovat vzájemné chování čtení a zápisu. V centralizovaných systémech bývá sémantika sdílení jednoznačná - každé čtení vrátí hodnotu uloženou při posledním zápisu. Dosažení této sémantiky však v distribuovaných systémech není jednoduché. Prvním problémem je neexistence přesného globálního usporádání událostí. Proto čtení na jednom procesoru, které se událo několik mikrosekund po zápisu na druhém procesoru, nemusí vrátit naposledy zapsanou hodnotu. Druhým problémem je cachování. Každý z konkurenčních procesů si sám cachuje soubor a změny provádí lokálně. Možné řešení by bylo implementovat write-through cache, tj. všechny změny by se ihned promítly do souboru spravovaném fileserverem. To však ještě nezaručuje plné sdílení - procesy, které již mají soubor načachovaný, by se stejně změnu dat nedozvedely.

Jiným řešením tohoto problému je uvolnění sémantiky. Výše popsané chování zůstává, avšak je prohlášeno za korektní. Presněji receno veškeré změny otevřeného souboru jsou viditelné pouze lokálně, zviditelní se až při zavření souboru. Takováto sémantika bývá nazývána relacní (session semantics).

Zcela odlišným způsobem je řešeno sdílení imutabilních souborů. Zde již neexistují akce “otevření souboru na zápis” apod., ale pouze CREATE a READ. Jestliže nějaký proces vytvoří nový soubor a uloží ho pod již existujícím jménem, pak stará kopie již není dostupná (alespon ne pod starým jménem). Otázkou je, co se stane, jestliže jiný proces čte soubor, který byl přepsán nově vytvořeným. Bežným řešením je nechat procesu otevřený starý soubor, a to i když na něj byl odstraněn odkaz ze všech adresářů. Jiným řešením je po zjištění, že soubor byl změněn, další pokus o čtení vrátit s chybou.

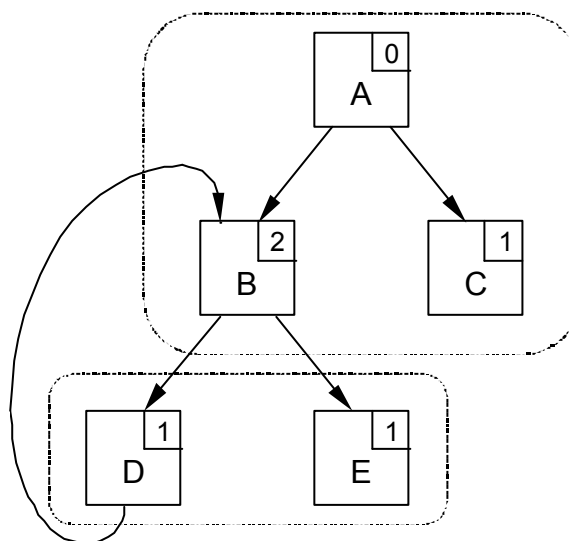
Pokud systém podporuje transakce, je možno sémantiku sdílení souborů přenechat transakčnímu manažeru, který sám implementuje potřebné mechanismy (zamykání, časové značky apod.).

Centralizovaná sémantika	Každá operace je ihned viditelná
Relacní sémantika	Změny jsou viditelné až po zavření souboru
Imutabilní soubory	Soubory nelze měnit
Transakce	Zabudované transakční mechanismy

Tab. 20 - Způsoby sdílení souborů

11.3 Adresářové služby

- mapování uživatelských jmen na systémová
- hierarchický systém souborů
- linky, graf adresáře - ! mazání linku



Obr. 62 - Graf adresáru

11.3.1 Prostor jmen adresáru

- jednoduchá jména - napr. uzel : adresár/jméno
- montování do lokálního prostoru
- jednotný celosvětový prostor jmen / kontexty

11.4 Replikace

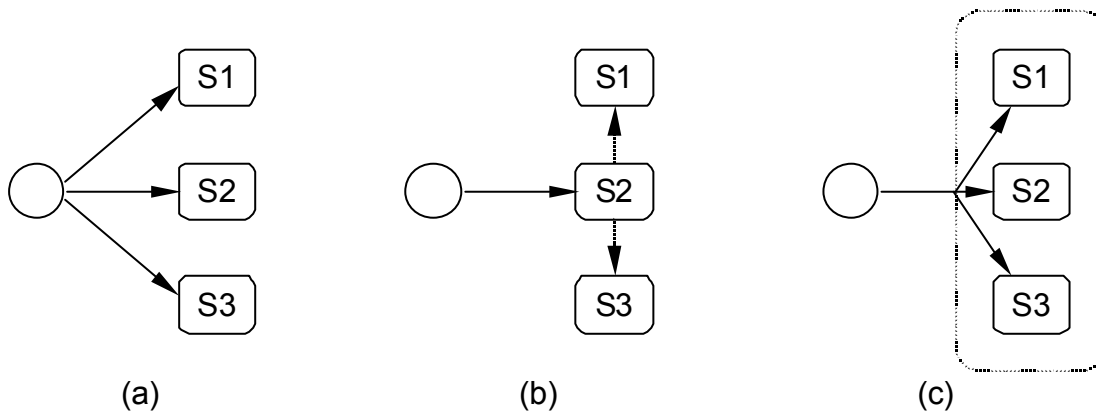
Distribuované filesystemy často umožňují replikaci souboru, tj. udržování více kopií všech nebo jen některých souborů na více filesereverech. To má několik dobrých důvodů:

1. Spolehlivost (reliability) - při výpadku jednoho serveru nejsou ztracena data
2. Dostupnost (availability) - k souborům lze přistupovat i v případě výpadku jednoho nebo několika serverů
3. Výkon (performance) - lze přistupovat k nejbližším datům, výkon fileserveru je rozdělen mezi několik serverů

?? **explicitní replikace** - uživatel se sám stará o udržování konzistence

?? **odložená replikace** - zápis do primární repliky, aktualizace sekundárních replik později

?? **skupinová komunikace** - zápisy jsou simultánně zasílány všem dostupným replikám



Obr. 63 - (a) Explicitní replikace (b) Odložená replikace (c) Skupinová komunikace

11.4.1 Aktualizační protokoly

Při zápisu replikovaného objektu je třeba zajištění konzistence replik. Pokud by všechny servery spravující repliky byly vždy dostupné, pak by dosažení konzistence nebyl žádný problém. To je však prakticky nesplnitelný požadavek - vždy je možnost, že některý server je nefunkční nebo nedostupný. Proto je nutné použít protokolů umožňující zápisy i v neideálních podmínkách, kdy je dostupná pouze část potřebných serverů.

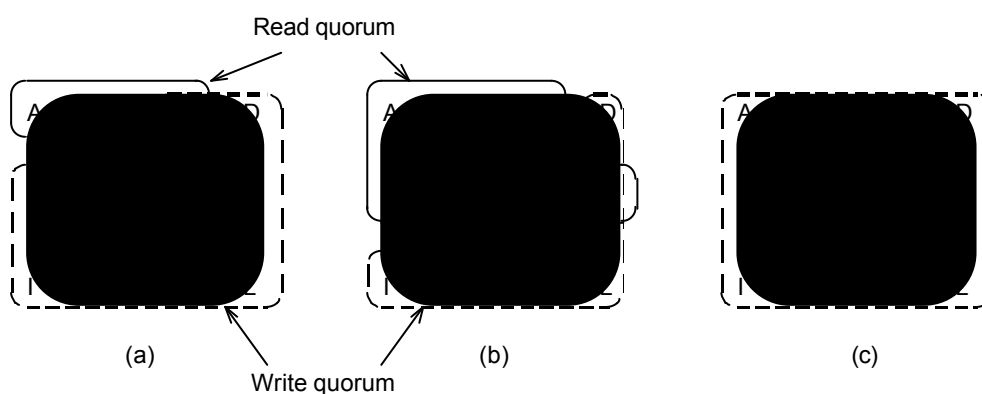
Nejjednodušší metodou je tzv. **primární replika**. Všechny změny se provádějí na jedné vybrané replice a tento server se sám stará o aktualizaci ostatních. Toto řešení má jednu vážnou nevýhodu - v případě nedostupnosti primární repliky nelze s objektem pracovat.

Řešení nabízí tzv. hlasování (voting). Základní myšlenkou tohoto přístupu je získání dostatečného počtu hlasu pro povolení přístupu k objektu. Klient před přístupem k objektu pošle žádosti všem dostupným serverům a čeká na povolení. Jakmile obdrží dostatečný počet hlasu (potvrzení přístupu), může provést požadovanou operaci. Jestliže dostatečný počet hlasu nezíská, nelze operaci provést. Podle toho, co přesně znamená "dostatečný počet hlasu", rozeznáváme několik druhů protokolu.

Nejjednodušším aktualizacním protokolem založeným na hlasování je tzv. **většinové hlasování** (majority voting). Je založeno na tom, že pro jakoukoliv operaci (ctení nebo zápis) je třeba získat většinu (nadpolovicní množství) z celkového počtu serverů. Při zápisu se všem zapsaným replikám přiřadí nové číslo verze (všem replikám stejné). Při čtení klient zažádá o hlasy, odpovědi obsahují čísla verze příslušné repliky. Jelikož jak pro čtení, tak pro zápis bylo třeba nadpolovicního počtu hlasu, je zaručeno, že v odpovědi je alespoň jedna replika s aktuálním číslem verze. Klient si pak vybere ze serveru, jejichž odpověď obsahovala aktuální číslo verze.

Zobecněním většinového hlasování je **vážené hlasování** (weighted voting). Pro zápis je třeba tzv. **write quorum** (počet hlasu), pro čtení **read quorum**. Necht celkový počet replik (serverů) je N , read quorum N_r , write quorum N_w . Pak musí platit $N_w + N_r > N$. (Z tohoto pohledu je většinové hlasování pouze speciálním případem váženého, kdy $N_w = N_r$). To, že čtecí a zápisové quorum mohou být různé, se dá využít pro optimalizaci přístupu. Většinou totiž platí, že čtení bývá mnohem častější než zápis. Proto když čtecí quorum bude menší nežli zápisové, bude režie pro čtení, jakožto častější operaci, menší.

Singulárním případem je $\text{read quorum} = 1$. Pak pro čtení stačí přístup k jednomu libovolnému serveru, zápis však potřebuje přístup ke všem serverům. V případě, že by některý server nebyl dostupný, pak by zápis nešel provést. Tento problém řeší tzv. **hlasování s duchy** (voting with ghosts). Kromě normálních (datových) serverů je vytvořen jeden bezdatový (dummy) server, který slouží k náhradě nedostupných datových serverů. Tento server (duch) se neúčastní hlasování o čtení (protože neobsahuje žádná data), ale může se účastnit hlasování o zápisu, resp. jeho hlas je počítán do write quora. Klient pak při pokusu o zápis uspeje, jestliže je alespoň jeden server reálný.



Obr. 64 - Hlasovací algoritmus pro 12 serverů
(a) $N_r=3, N_w=10$ (b) $N_r=7, N_w=6$ (c) $N_r=1, N_w=12$