
Skupinové zasílání zpráv, Multicasting

PA 150 ◊ Principy operačních systémů

Jan Staudek

<http://www.fi.muni.cz/usr/staudek/vyuka/>



Verze : podzim 2020

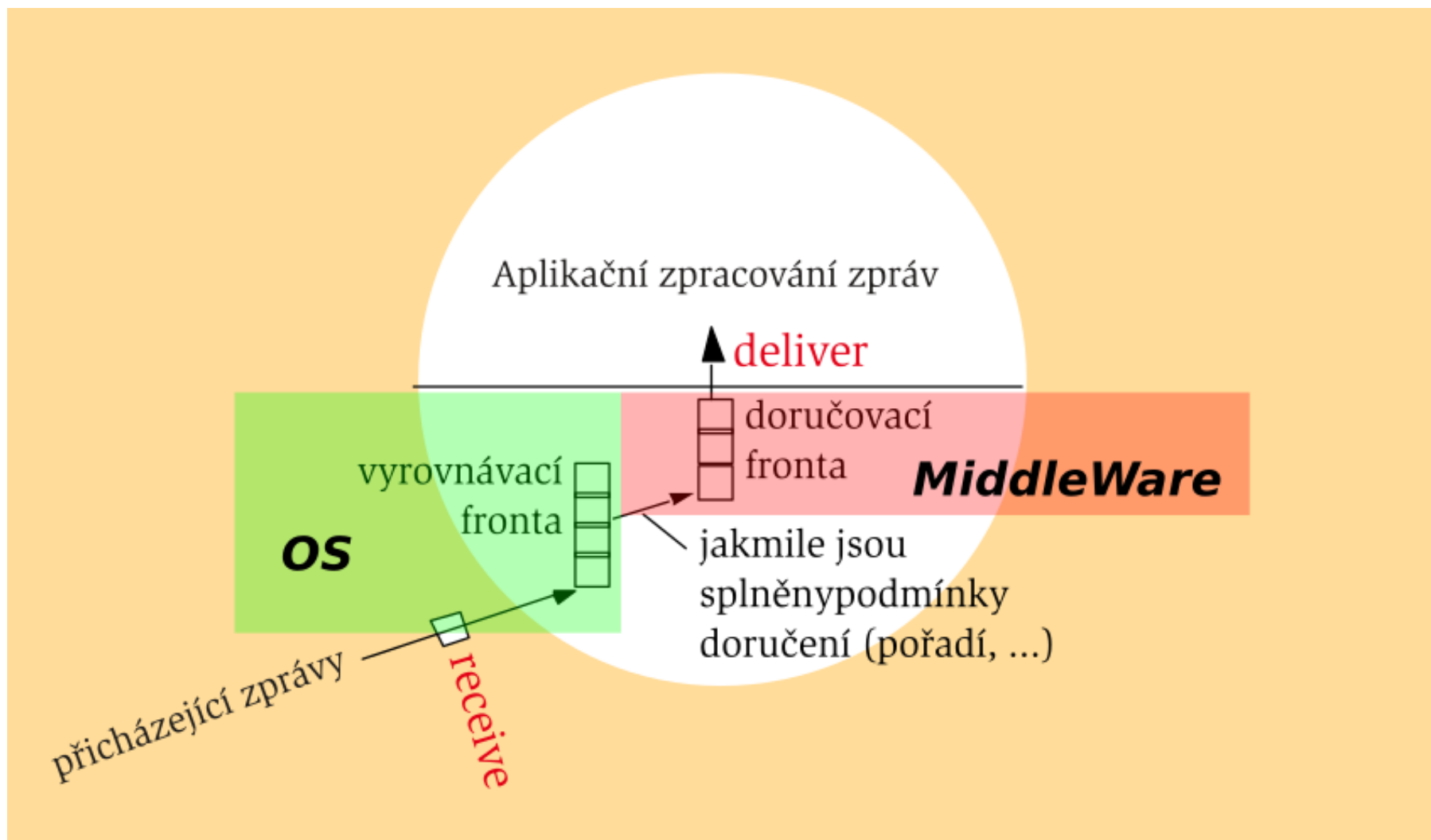
Model skupinové komunikace (Multicast Communication)

- Řešení problému sdělení zprávy **všem** procesům náležejících do **skupiny procesů**, složení skupiny je známé
- Procesy mezi sebou komunikují spolehlivými 1:1 kanály
- Procesy mohou krachovat (krach = výpadek = *fail-stop*, proces buď běží správně nebo je nečinný)
- Vysílající proces **vysílá zprávu m skupině procesů g** jedinou operací, *multicast*(g, m)
 - ✓ ta je implementovaná v middleware násobnými operacemi na úrovni služeb OS, *send*(p, m), kde p jsou id procesů ze skupiny g , $p \in g$
- Zpráva m je unikátně identifikovatelná, nese sebou jedinečný id odesílatele a jednoznačný id skupiny cílových procesů, $g = \text{group}(m)$

Problém skupinové komunikace (Multicast Communication)

- Uzel procesu v adresované skupině **zprávu přijme** (operací *receive*) na úrovni služeb OS, middleware ji doručí (operací *deliver*) aplikačnímu procesu
- Procesu ve skupině g přijatou **zprávu m doručí** operace *deliver(m)*, komplementární operace k operaci *multicast*, s vlastnostmi:
 - ✓ Zpráva m musí být doručena všem procesům aktivním ve skupině g
 - ✓ Pokud zprávu m získá jeden ze skupiny procesů g , musí ji získat všechny procesy aktivní ve skupině g

Problém skupinové komunikace (Multicast Communication)



Bázové řešení skupinové komunikace

□ *B-multicast*(g, m)

Basic-multicast – se (na úrovni middleware) implementuje násobným provedením spolehlivých operací $send(p, m)$ z úrovně OS ke všem procesům p ze skupiny g

- ✓ Provedení spolehlivých operací $send(p, m)$ lze zparalelnit pomocí vláken odpovídajících ve vysílajícím procesu procesům ve skupině g
- ✓ Pak ale potvrzení přijmutí zprávy mají tendenci se kumulovat, vysílač zprávy je nemusí stačit zpracovávat, opakuje proto vysílání, zpráv, což způsobuje příchod dalších potvrzení – *ack-implosion*

□ *B-delivered*(m), komplementární operace k *B-multicast*, zprávu m pro proces p přijímá v uzlu procesu p OS operací $receive(m)$ a middleware ji doručuje procesu p

Bázové řešení skupinové komunikace

- **Bázové řešení skupinové komunikace** nezaručuje splnění požadavku doručení zprávy všem procesům ve skupině
 - ✓ pokud vysílající uzel vypadne během vysílání operacemi $send(p, m)$ implementujícími B -multicast(g, m), některé procesy zprávu dostanou, jiné ne a procesy ve skupině g nemají šanci tuto skutečnost (kdo zprávu dostal a kdo ne) zjistit
- **Spolehlivé řešení skupinové komunikace** požaduje:
 - zpráva je doručena všem procesům ve skupině,
 - pokud ji získá jeden ze skupiny procesů,
 - musí ji získat všechny procesy ve skupině

Spolehlivé (*Reliable*) řešení skupinové komunikace

- Požadované vlastnosti operací *R-multicast*(g, m) a *R-deliver*(m)
 - ✓ **integrita** – spolehlivá 1:1 komunikace procesů
korektně se chovajícímu procesu se doručuje zpráva m **nejvýše jednou**
(zprávy lze jednoznačně identifikovat:
např. id zdroje + pořadí zprávy ve zdroji)
 - ✓ **validita** –
jestliže korektní proces skupinově rozešle zprávu m ,
zpráva m nakonec bude doručena
 - ✓ **dosažení shody** –
jestliže je zpráva m doručena jednomu korektnímu procesu skupiny,
je doručena všem korektním procesům ve skupině
společně s validitou se jedná o záruku živosti algoritmu –
pokud *send* odvozený z *multicast* vyšle alespoň jednou zprávu
k procesu skupiny, nesmí výpadek originálního vysílače ohrozit
doručení zprávy ke všem procesům

Spolehlivé (*Reliable*) řešení skupinové komunikace

- Spolehlivost se dosahuje nadstavbou základního řešení

On initialization

Received := {}; // stavová proměnná v každém procesu



For process p to R-multicast message m to group g

B-multicast(g, m); // p ∈ g is included as a destination



On B-deliver(m) at process q with g = group(m)

if (m ∉ Received)

then

Received := Received ∪ {m};

if (q ≠ p) then B-multicast(g, m); end if

R-deliver m;

end if

Spolehlivé (*Reliable*) řešení skupinové komunikace

- Je zajištěna **validita** a **shoda**
zprávu by některý z procesů ve skupině neobdržel pouze
kdyby žádný proces skupiny neudělal *B-multicast*,
tj. kdyby žádný nefungoval
- **integritu** podporuje podpůrný komunikační systém řešící
B-multicast
- Spolehlivé rozesílání je zaručeno i v asynchronním systému,
o čase se nic nepředpokládá
- Každá zpráva je každému procesu skupiny g zaslána $|g|$ -krát,
:-((

Skupinová komunikace se zachováním pořadí zpráv

- Jestliže je důležité aby se doručování zpráv se odehrávalo v definovaném pořadí, existují 3 typy řazení multicastu:
 - ✓ **FIFO řazení** –
provede-li korektní proces $multicast(g, a)$ před $multicast(g, b)$, pak každý korektní proces v g , který získává zprávu b , získá zprávu a před získáním zprávy b
 - ✓ **kauzální řazení** –
jestliže $multicast(a)$ předchází $multicast(b)$ v g z hlediska zasílání zpráv mezi procesy ve skupině g , pak každý korektní proces v g , který získává zprávu b , získá zprávu a před získáním zprávy b
 - ✓ **totální řazení** –
jestliže některý korektní proces v g získává zprávu a před získáním zprávy b , pak každý korektní proces v g , který získává zprávu b , získá zprávu a před získáním zprávy b

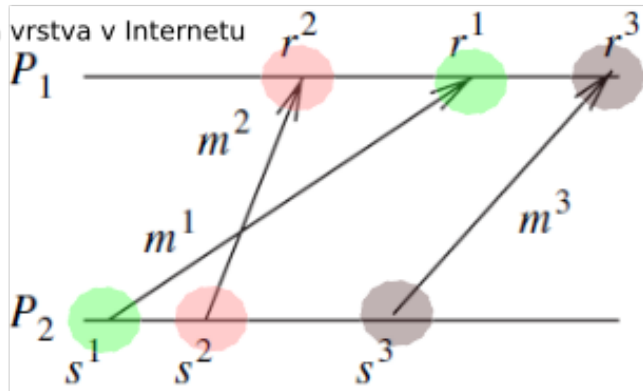
Skupinová komunikace se zachováním pořadí zpráv

- Kauzální řazení přirozeně implikuje FIFO řazení
- Kauzální řazení a FIFO řazení jsou částečná uspořádání
- Definice řazení nepředpokládá ani neimplikuje spolehlivost doručení
 - ✓ Např. při totálním řazení jestliže je korektnímu procesu p doručena zpráva a po ní zpráva b , procesu q může být doručena a bez doručení b nebo jiných zpráv řazených za a

Paradigmata doručování zpráv v asychr. DS

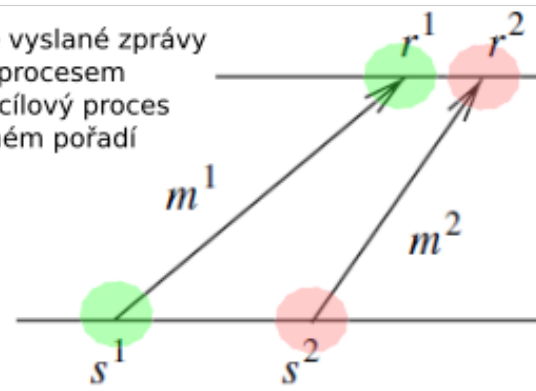
non-FIFO komunikace

např. síťová vrstva v Internetu



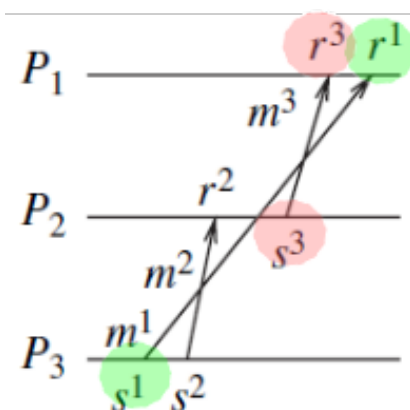
FIFO komunikace

Po sobě vyslané zprávy jedním procesem přijímá cílový proces ve stejném pořadí



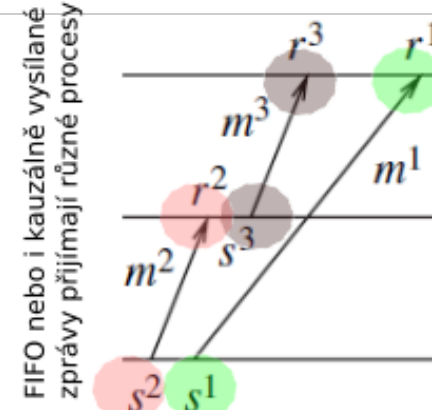
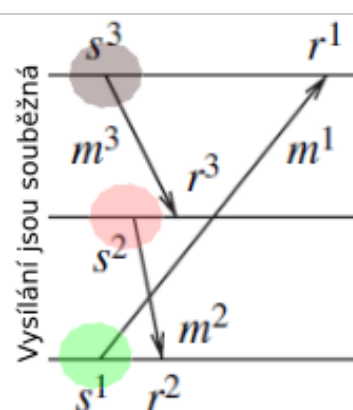
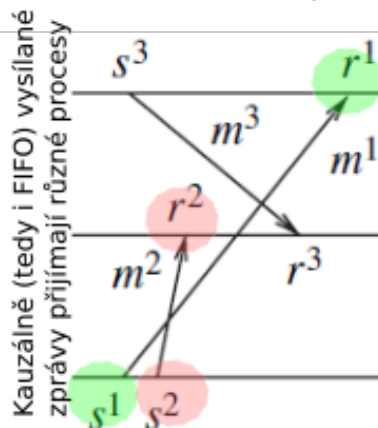
Komunikace nesplňující kauzalitu:

$s1 \rightarrow s3, r3 \rightarrow r1$

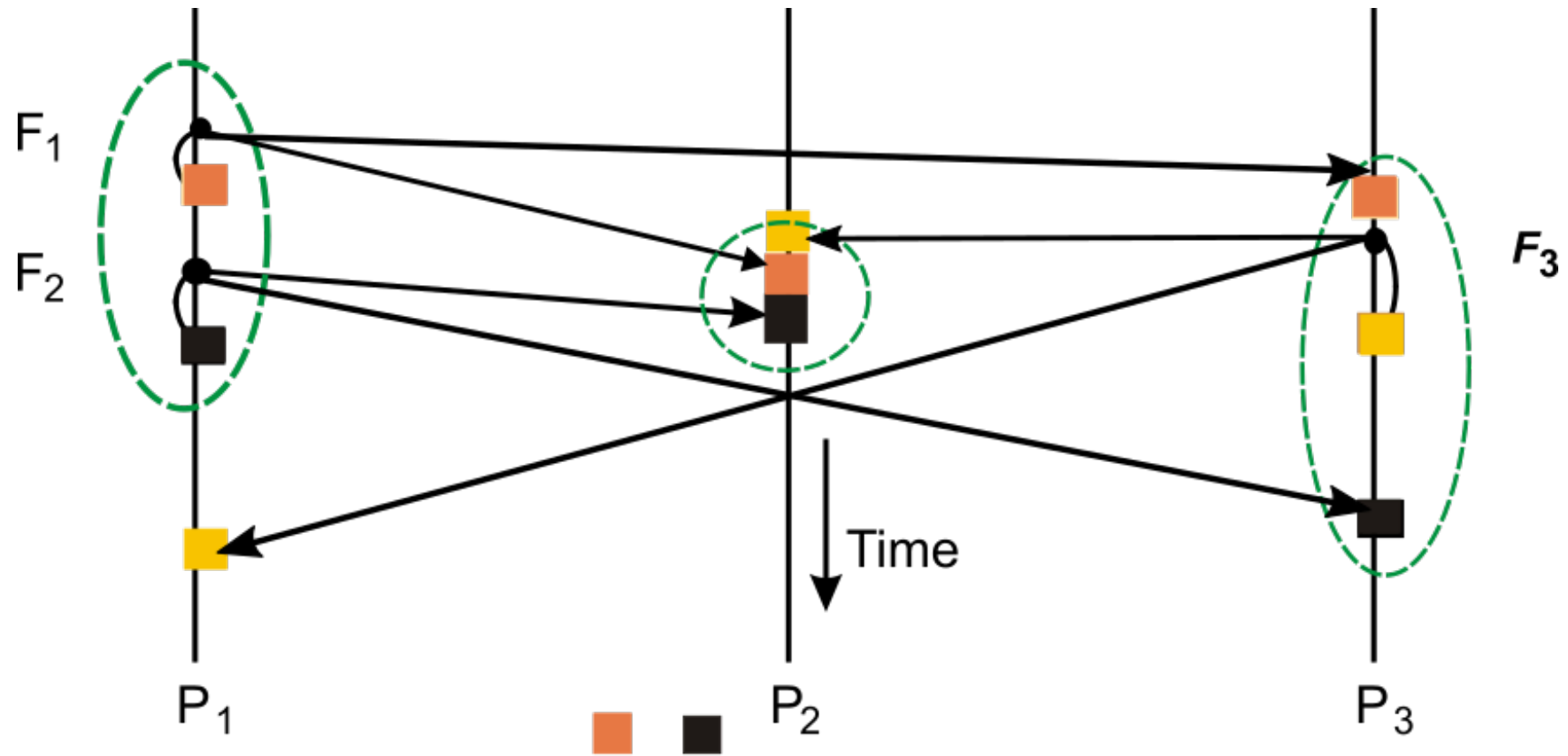


Komunikace splňující kauzalitu

Pokud po sobě vyslané zprávy přijme jeden cílový proces, pak tyto zprávy přijímá ve stejném pořadí



Skupinová komunikace se zachováním FIFO pořadí zpráv



FIFO uspořádané zprávy F1 a F2

F3 je na F1 a F2 nezávislá



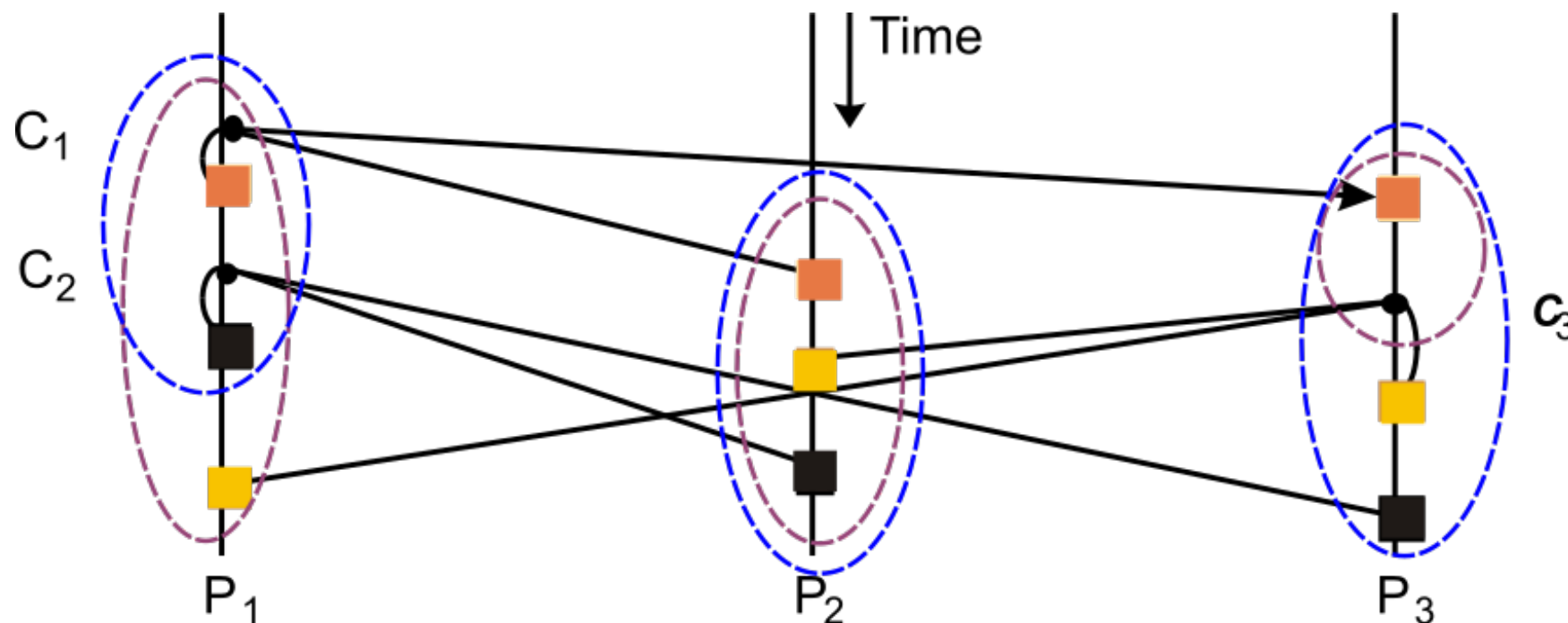
Zpráva F2 se rozesílá později než zpráva F1.
Všechny procesy ve skupině P1, P2 a P3
přijímají zprávy v pořadí F1, F2



Implementace FIFO řazení

- Pro každý proces p ze skupiny g middleware (MW) udržuje proměnné
 - ✓ S_g^p – čítač zpráv vyslaných procesem p do skupiny g
 - ✓ R_g^q – pro každý proces q ze skupiny g číslo poslední zprávy ze zpráv zaslaných do g procesem q a doručených procesu p
- Při vysílání operací *multicast* přidává MW ke zprávě $S = S_g^p$ a poté S_g^p inkrementuje o 1, zprávy z p pořadově čísluje
- Při přijetí zprávy od q pro p MW kontroluje R_g^q a S ze zprávy
 - ✓ jestliže $S = R_g^q + 1$, MW zprávu doručí p a nastaví $R_g^q = S$
 - ✓ jestliže $S > R_g^q + 1$, MW zprávu zapamatuje ve vyrovnávací frontě a doručí ji procesu p až bude platit $R_g^q + 1 = S$
- Jestliže se použije spolehlivé rozesílání, *R-multicast*, získáme spolehlivé FIFO rozesílání

Skupinová komunikace se zachováním kauzálního pořadí



Kauzálně uspořádané zprávy C1 a C3

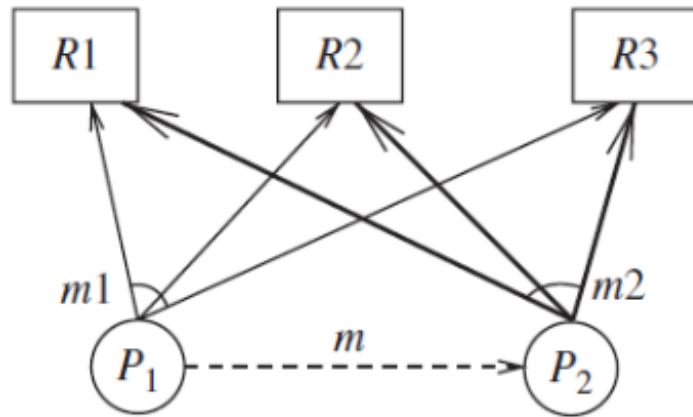
multicast C3 následuje **po** multicast C1

Kauzálně uspořádané zprávy C1 a C2

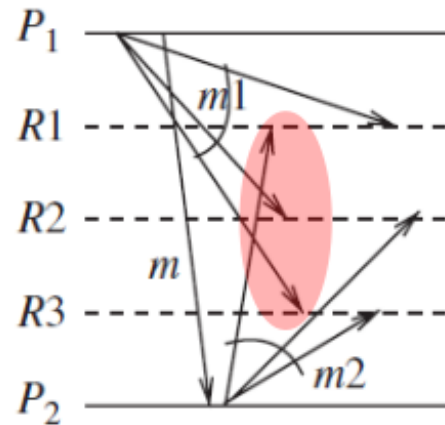
multicast C2 následuje **po** multicast C1
(kauzální řazení implikuje i FIFO řazení)

Mezi zprávami C3 a C2 kauzální vztah neexistuje

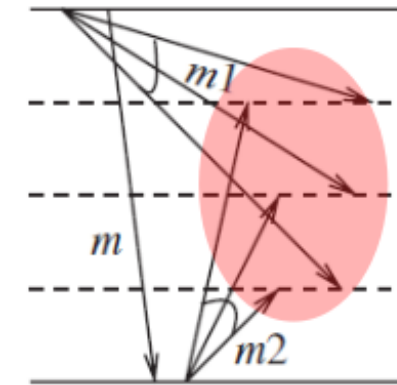
Skupinová komunikace se zachováním kauzálního pořadí



(a)



(b)



(c)

P_1 a P_2 jsou procesy běžící ve dvou různých uzlech DS a korigují 3 replikace jednoho zdroje ve 3 různých uzlech DS

Korekce prováděné P_2 musí proběhnout po korekcích prováděných P_1 , protože vysílání P_2 kauzálně následuje po vysílání P_1 díky vyslání zprávy m z P_1 do P_2 po vyslání zpráv m_1

Porušení kauzality

Implementace kauzálního řazení

□ Bázová idea

- ✓ Respektování relace *stalo-se-před* na bázi logického času hnaného rozesíláním a doručováním (*multicast*) zpráv
- ✓ Místo prostého časového razítka TS , ve zprávách použijeme **vektorové časové razítko V**
- ✓ Vektorové časové razítko V má tolik prvků kolik je procesů ve skupině, iniciálně jsou všechny prvky V nulové
- ✓ Každý proces p_i udržuje své vlastní vektorové časové razítko V_i
- ✓ Když proces rozesílá zprávu procesům ve skupině, nejprve inkrementuje svůj čas ve svém V a poté zprávu s V rozešle
- ✓ Když se přijme zpráva m pro proces q od procesu p , může mu být doručena až když se mu doručí zprávy, které ji kauzálně předcházejí
tj. až se q doručí všechny předchozí zprávy vyslané procesem p a až se q doručí všechny zprávy, které byly doručeny procesu p do doby, kdy p rozeslal m

Implementace kauzálního řazení vektorovými čas. razítky

Algorithm for group member p_i ($i = 1, 2, \dots, N$)

On initialization

$V_i^g[j] := 0$ ($j = 1, 2, \dots, N$);

To CO-multicast message m to group g

$V_i^g[i] := V_i^g[i] + 1$;

$B\text{-multicast}(g, \langle V_i^g, m \rangle)$;

p_j získal všechny zprávy, které získal p_j do doby rozeslání této zprávy

On $B\text{-deliver}(\langle V_j^g, m \rangle)$ from p_j , with $g = \text{group}(m)$

place $\langle V_j^g, m \rangle$ in hold-back queue;

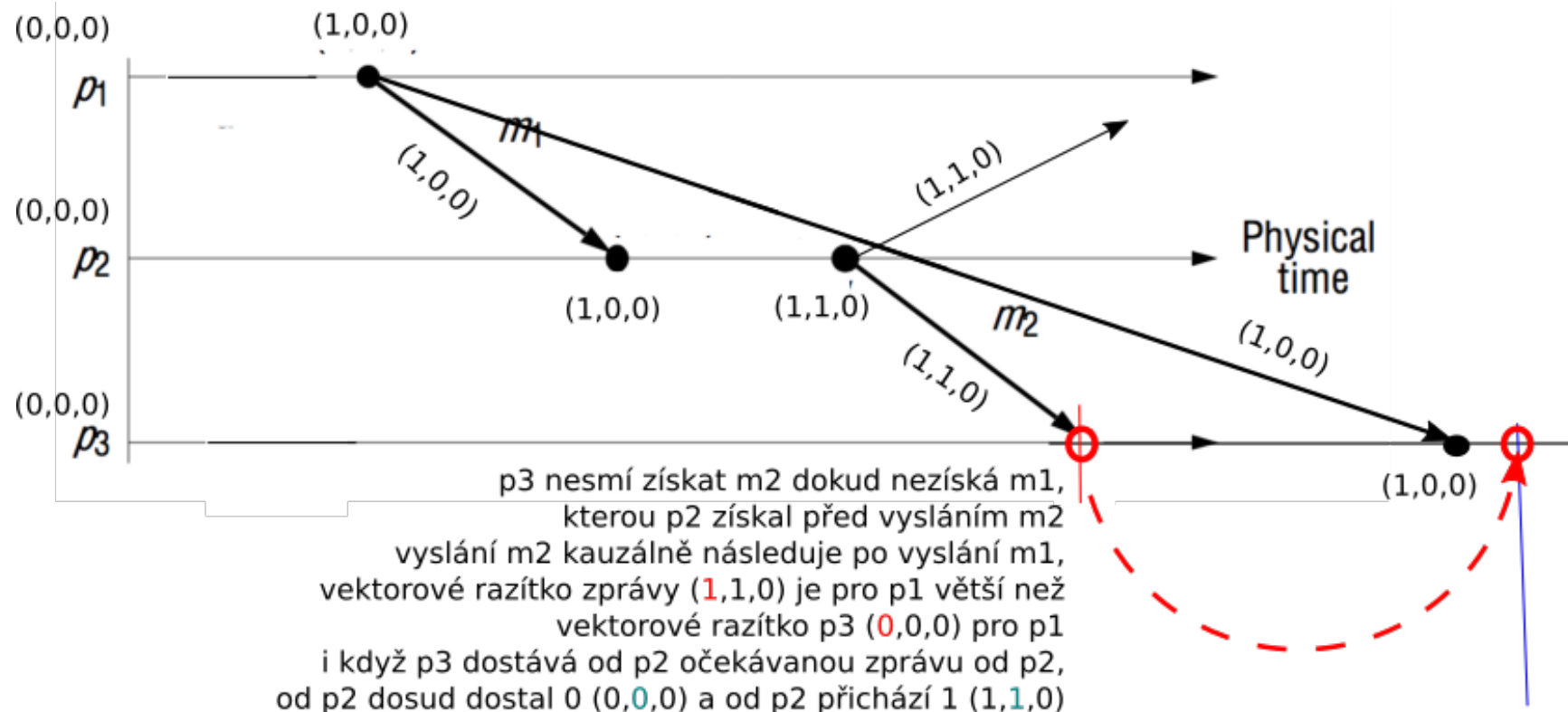
wait until $V_j^g[j] = V_i^g[j] + 1$ and $V_j^g[k] \leq V_i^g[k]$ ($k \neq j$);

$CO\text{-deliver } m$; // after removing it from the hold-back queue

$V_i^g[j] := V_i^g[j] + 1$; p_i získal od p_j všechny předcházející zprávy vyslané p_j

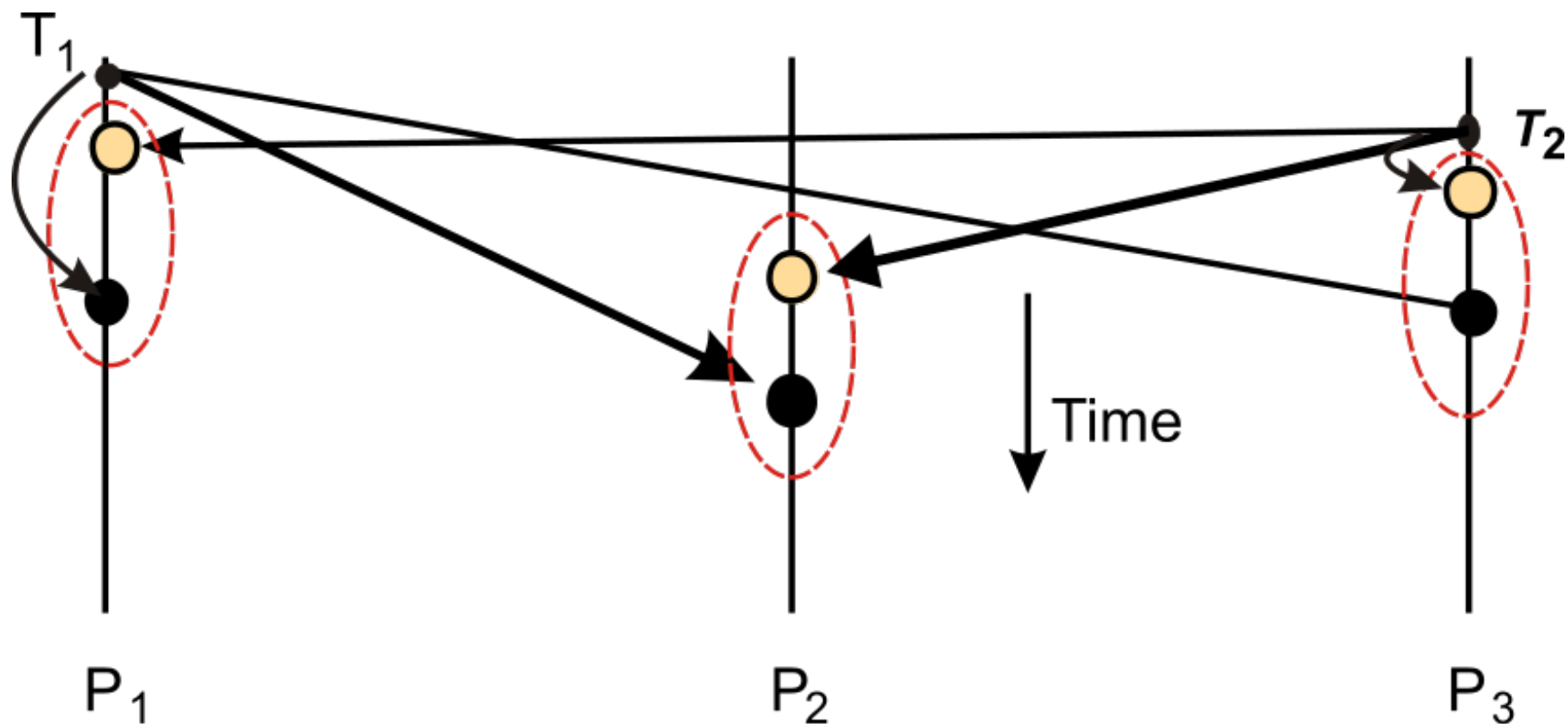
- ✓ Jestliže se použije spolehlivé rozesílání, $R\text{-multicast}$, získáme spolehlivé kauzální rozesílání

Kauzálního řazení vektorovými čas. razítky



p_3 už smí získat m_2 protože už získal m_1 ,
 kterou p_2 získal před vysláním m_2
 vyslání m_2 kauzálně následuje po vyslání m_1
 vektorové razítko zprávy $(1,0,0)$ není pro p_1 větší než
 vektorové razítko zprávy $(1,0,0)$ pro p_1

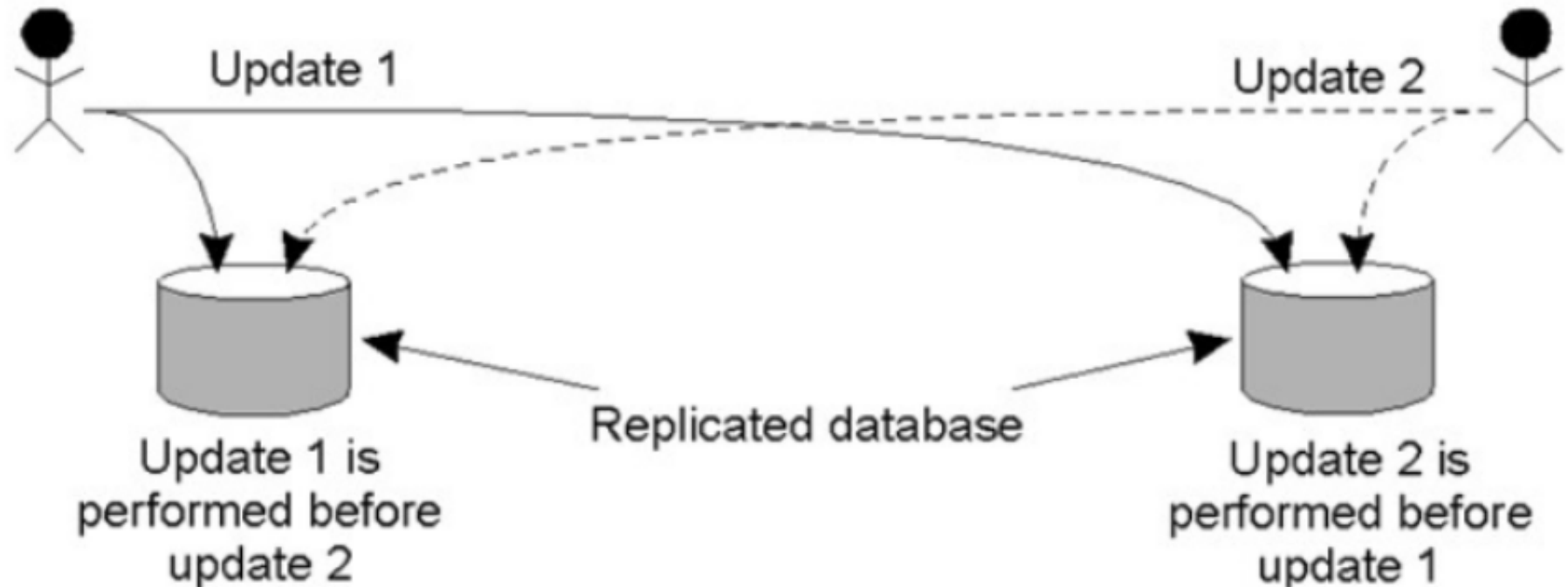
Skupinová komunikace se zachováním totálního pořadí



Totálně řazené zprávy T₁ a T₂

Poněvadž P₁ získal zprávu T₂ před zprávou T₁,
získává zprávu T₂ před zprávou T₁
každý korektní proces ve skupině P₁, P₂, P₃

Ilustrace dopadu nedodržení totálního pořadí



- Obě replikace databáze se ocitnou v nekonzistentním stavu

Implementace totálního řazení,

- Analogické řešení jako FIFO řazení, čítač rozesílaných zpráv ale **musí být ve skupině procesů jedinečný**
- centralizované řešení
 - ✓ DS podporuje pouze FIFO kanály
 - ✓ Zprávy rozesílá všem procesům ve skupině jeden centrální proces *CP*
 - ✓ Proces *P* rozesílající zprávu *M* pošle tuto zprávu *CP*
 - ✓ *CP* rozešle všem procesům ve skupině zprávu $\{P, M\}$

Implementace totálního řazení, *sequencerem*

- Řešení totálního řazení pomocí *sequenceru*
 - ✓ Čítač rozesílaných zpráv udržuje jeden proces ve skupině - *sequencer*
 - ✓ Každý proces, který chce rozesílat zprávu ji s jedinečným indexem rozešle procesům ve skupině a *sequenceru*
 - ✓ jedinečný index – např. id procesu + pořadové číslo zprávy v procesu
 - ✓ přijatou zprávu middleware umístí do vyrovnávací fronty
 - ✓ *sequencer* určí jemu doručené zprávě její pořadové číslo a toto číslo s indexem zprávy rozešle procesům ve skupině
 - ✓ Middleware pro každý proces ve skupině doručuje přijaté zprávy z vyrovnávací fronty v pořadí pořadových čísel
 - ✓ *Sequencer* z hlediska komunikační zátěže a spolehlivosti je úzkým profilem řešení

Implementace totálního řazení, *sequencerem*

1. Algorithm for group member p

On initialization: $r_g := 0$;

To TO-multicast message m to group g

B-multicast($g \cup \{\text{sequencer}(g)\}$, $\langle m, i \rangle$);

On B-deliver($\langle m, i \rangle$) with $g = \text{group}(m)$

Place $\langle m, i \rangle$ in hold-back queue;

On B-deliver($m_{\text{order}} = \langle \text{"order"}, i, S \rangle$) with $g = \text{group}(m_{\text{order}})$

wait until $\langle m, i \rangle$ in hold-back queue and $S = r_g$;

TO-deliver m ; // (after deleting it from the hold-back queue)

$r_g = S + 1$;

2. Algorithm for sequencer of g

On initialization: $s_g := 0$;

On B-deliver($\langle m, i \rangle$) with $g = \text{group}(m)$

B-multicast(g , $\langle \text{"order"}, i, s_g \rangle$);

$s_g := s_g + 1$;

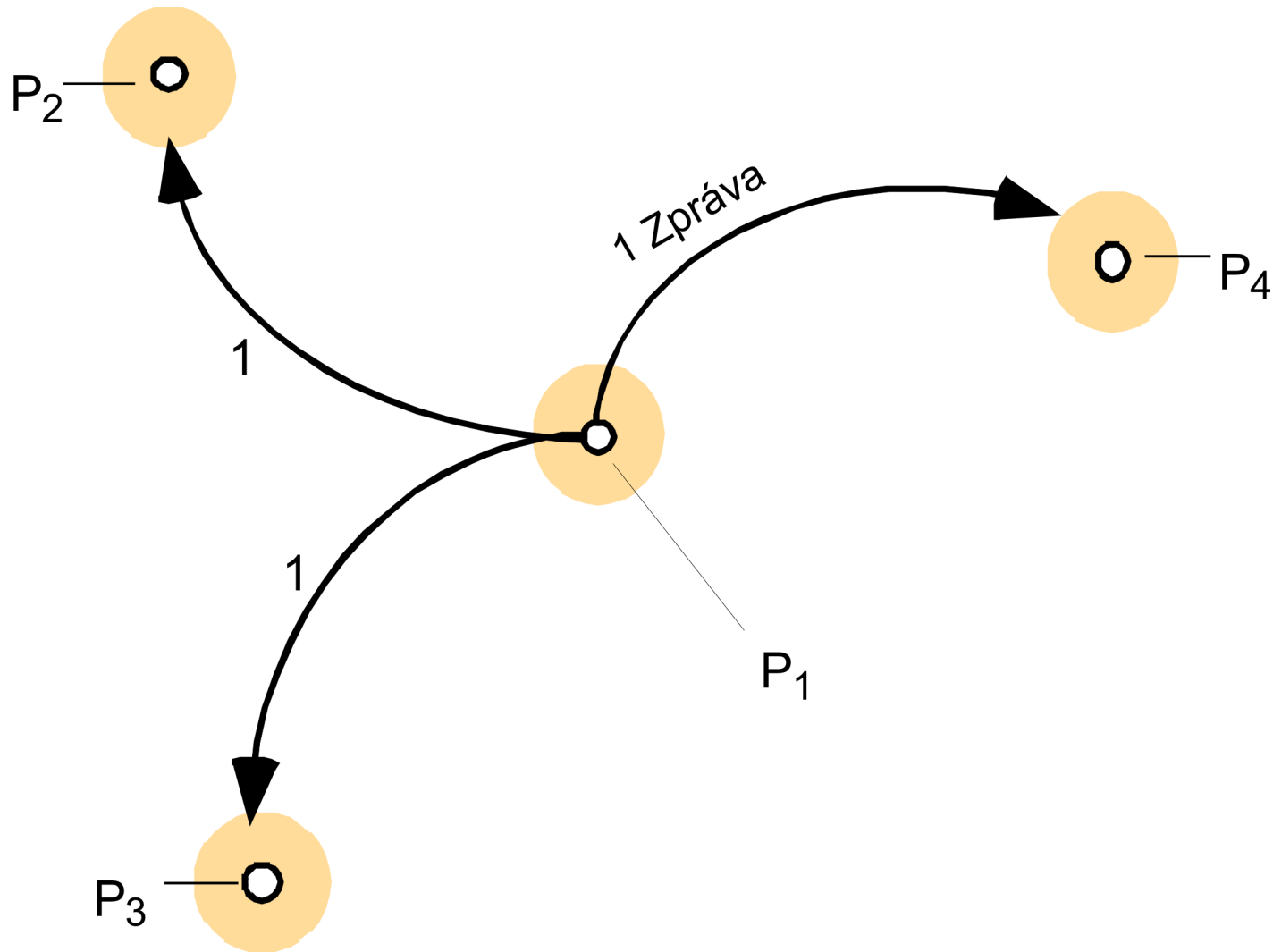
Implementace totálního řazení, *dohodou*

- Třífázový distribuovaný algoritmus
- Princip algoritmu:
 - ✓ Každý proces q si pamatuje v proměnné A_g^q nejvyšší jemu známé dohodnuté pořadové číslo zprávy některým procesem rozeslané do g
 - ✓ každý proces q vytváří a udržuje v proměnné P_g^q jím navrhované pořadové číslo zprávy
 - ✓ proces p rozešle základním rozesláním (*B-multicast*) všem procesům ve skupině jedinečně indexovanou (i) rozesílanou zprávu $\langle m, i \rangle$, procesy si ji zapamatují na úrovni *receive*

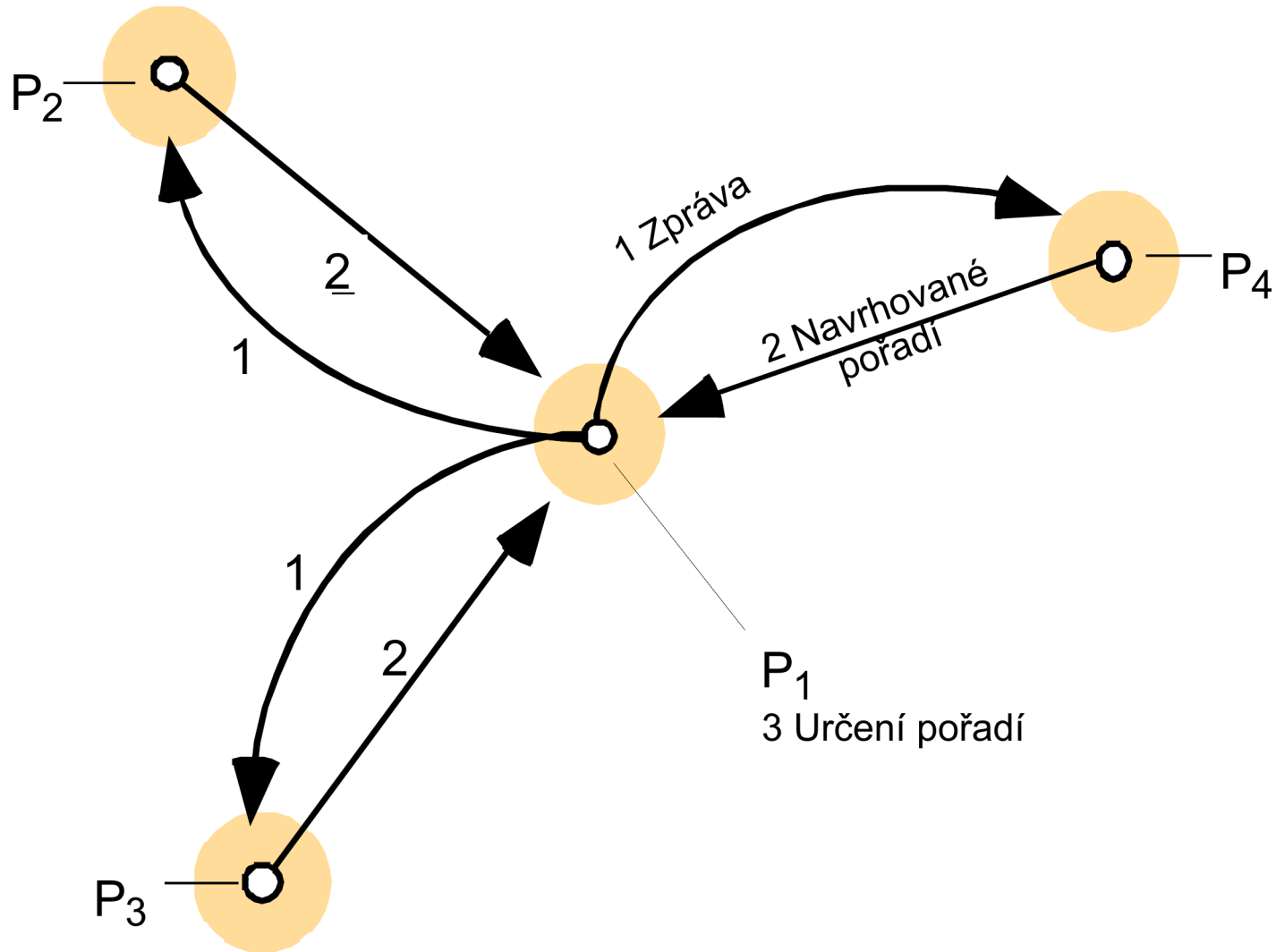
Implementace totálního řazení, *dohodou*

- ✓ proces q , který přijal zprávu $\langle m, i \rangle$ vrátí procesu p **jím navrhované pořadové číslo** zprávy $\langle m, i \rangle$ jako hodnotu $P_g^q = \max(A_g^q, P_g^q) + 1$
- ✓ proces p vybere největší navrhované číslo zprávy $\langle m, i \rangle$, označme je a , a toto základním rozesláním rozešle všem procesům skupiny jako **dohodnuté pořadové číslo** zprávy $\langle m, i \rangle$
- ✓ každý proces q ve skupině si nastaví $A_g^q = \max(A_g^q, a)$
- ✓ zpráva $\langle m, i \rangle$ je pak doručena každému procesu ve skupině v takto dohodnutém pořadí

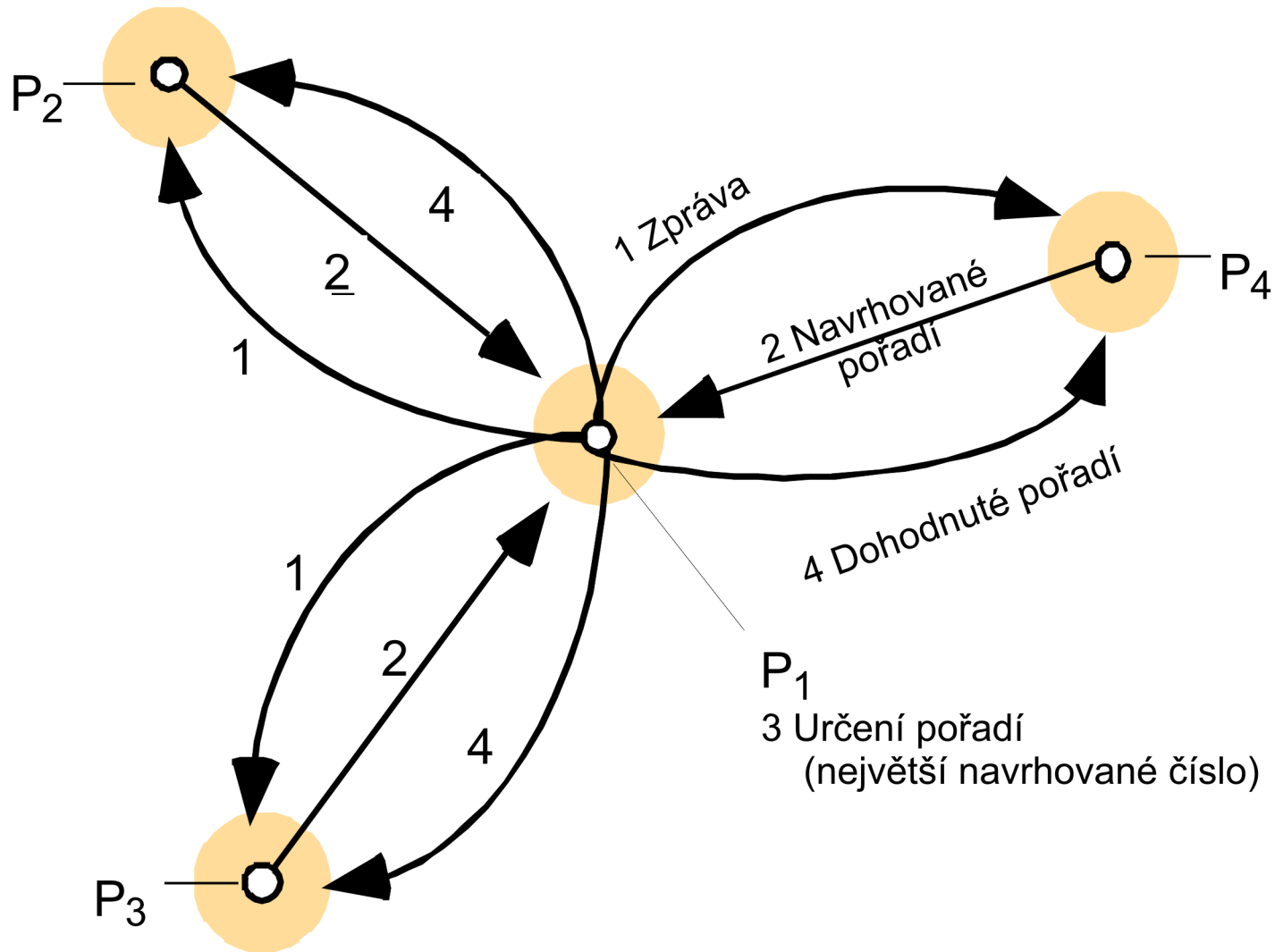
Implementace totálního řazení, *dohodou*



Implementace totálního řazení, *dohodou*



Implementace totálního řazení, *dohodou*



Implementace totálního řazení, *dohodou*

