

# DISTRIBUIRANI ALGORITMI I SISTEMI

# Distribuirana deljena memorija

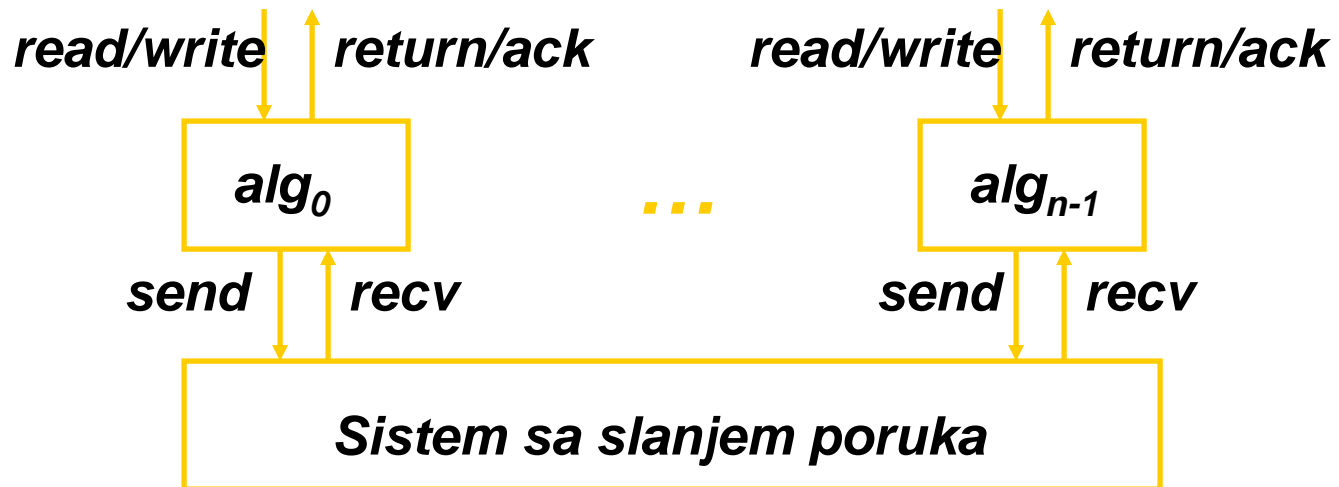
2

- Jedan model međuprocene komunikacije
- Obezbeđuje iluziju **deljenih promenljivih** na vrhu sistema za slanje poruka
- Deljena memorija se obično smatra pogodnijom platformom od sistema za slanje poruka
- Formalno, to je simulacija modela deljene memorije na vrhu modela sa slanjem poruka
- Posmatraćemo poseban slučaj:
  - ▣ nema otkaza
  - ▣ simuliraju se samo promenljive sa operacijama čitaj/piši

# Simulacija

3

*korisnici deljene memorije*



# Problemi deljene memorije

4

- Proces poziva operaciju deljene memorije (čitaj ili piši) u nekom trenutku
- Simulacioni algoritam na nekom čvoru izvršava neki kod, uz moguću razmenu poruka
- Na kraju simulacioni algoritam informiše proces o rezultatu pozvane operacije
- Znači operacije deljene memorije **nisu trenutne!**
  - ▣ Operacije (pozvane od raznih procesa) se mogu preklapati
- Koje vrednosti treba da vrate operacije koje se preklapaju sa drugim operacijama?
  - ▣ to definiše **uslov memorijske konzistentnosti**

# Sekvencijalne specifikacije

5

- Svaki deljeni objekt ima **sekvencijalnu spec.:** specificira ponašanje objekta u *odsustvu* konkurencije
- Objekt podržava **operacije:**
  - ▣ pozive
  - ▣ odgovarajuće odgovore
- Skup sekvenci operacija koje su **legalne**

# Sekvencijalna specifikacija za R/W registre

6

- Svaka operacija ima dva dela, poziv i odgovor
- Operacija **read** ima poziv  $\text{read}_i(X)$  i odgovor  $\text{return}_i(X,v)$  (indeks  $i$  odgovara procesu)
- Operacija **write** ima poziv  $\text{write}_i(X,v)$  i odgovor  $\text{ack}_i(X)$  (indeks  $i$  odgovara procesu)
- Sekvenca operacija je **legalna** ako i samo ako svako čitanje vraća vrednost poslednjeg prethodnog upisa
- Npr.:  $[\text{write}_0(X,3) \text{ack}_0(X)] [\text{read}_1(X) \text{return}_1(X,3)]$

# Uslovi memorijske konzistentnosti

7

- Uslovi konzistentnosti povezuju sekvencijalnu specifikaciju sa onim što se dešava u prisustvu konkurencije
- Proučićemo dva dobro-poznata uslova:
  - ▣ mogućnost linearizacije (linearizability)
  - ▣ sekvencijalna konzistentnost
- Razmotićemo samo registre za čitanje/upis, u slučaju kada nema otkaza

# Definicija mogućnosti linearizacije (Linearizability)

8

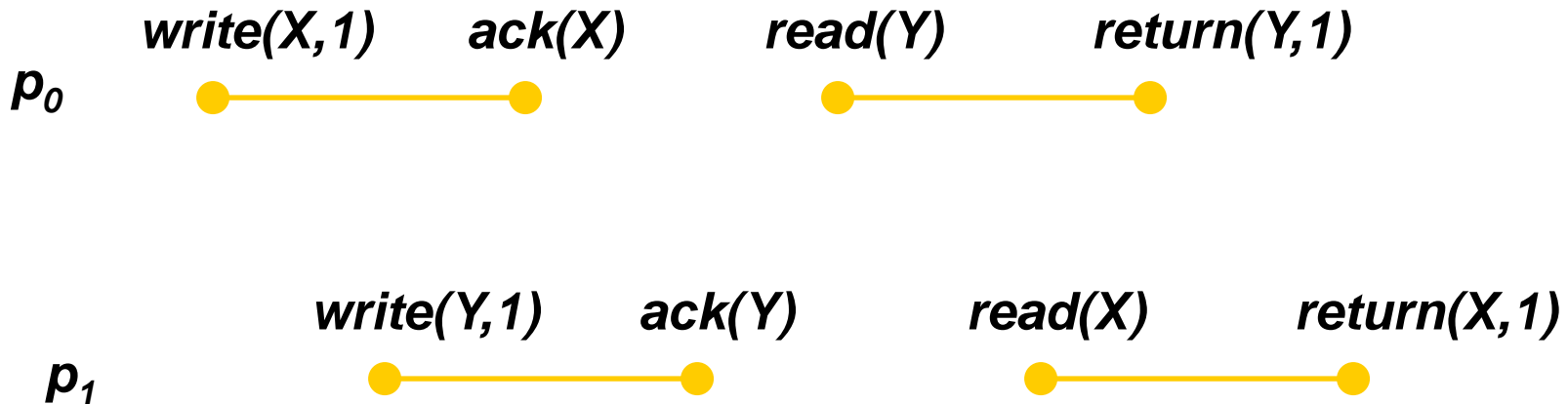
- Neka je  $\sigma$  sekvenca poziva i odgovora za skup operacija
  - ▣ poziv ne mora biti odmah praćen svojim odgovorom, može imati konkurentne, preklapajuće operacije
- $\sigma$  se **može linearizovati** ako postoji permutacija  $\pi$  svih operacija u  $\sigma$  (gde je svaki poziv odmah praćen svojim odgovorom), koja zadovoljava:
  - ▣  $\pi \mid X$  je legalna (**zadovoljava sekve. spec.**) za sve prom.  $X$ , i
  - ▣ ako se odgovor za operaciju  $O_1$  desio u  $\sigma$  pre poziva operacije  $O_2$ , onda se  $O_1$  desila u  $\pi$  pre  $O_2$  ( **$\pi$  respektuje redosled realnog vremena nepreklopljenih operacija u  $\sigma$** )



# Primeri mogućnosti linearizacije

9

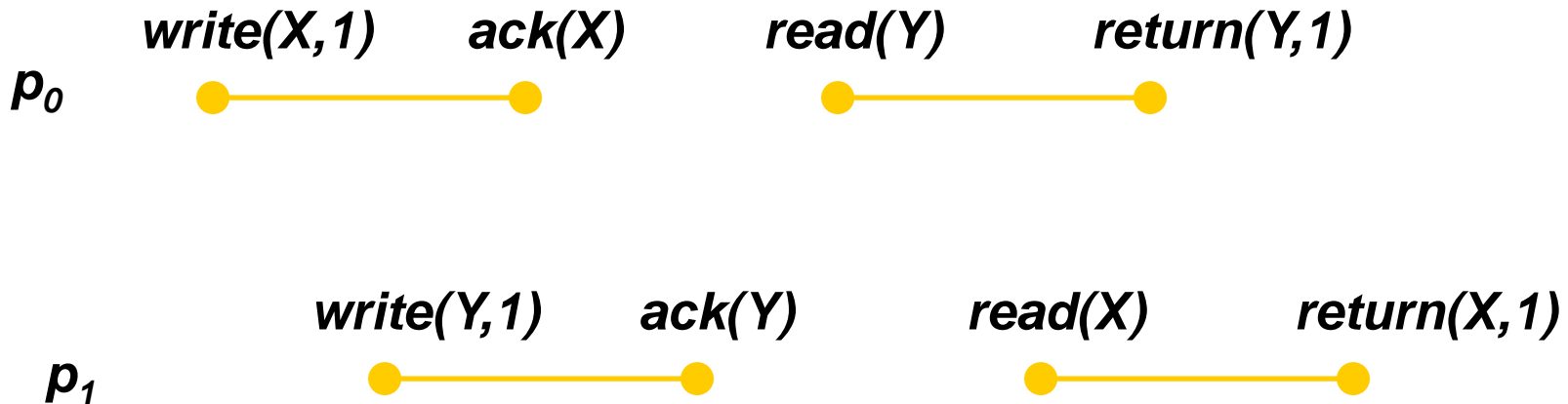
Neka postoje dve deljene promenljive,  $X$  i  $Y$ , obe inicijalno 0



# Primeri mogućnosti linearizacije

10

Neka postoje dve deljene promenljive,  $X$  i  $Y$ , obe inicijalno 0

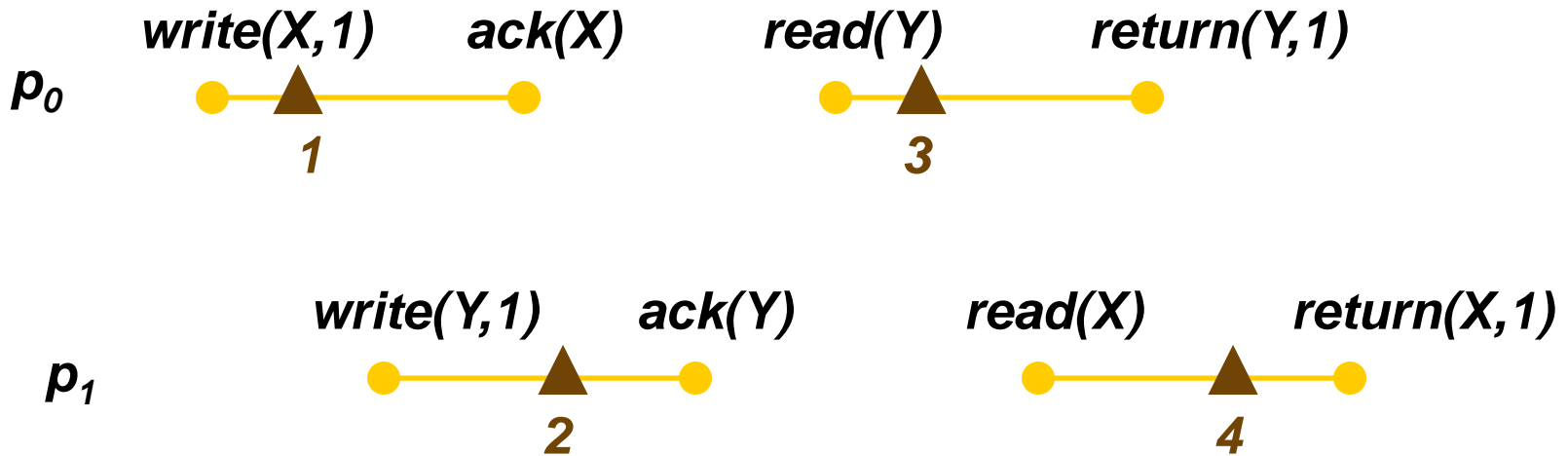


Da li se ova sekvenca  
može linearizovati?

# Primeri mogućnosti linearizacije

11

Neka postoje dve deljene promenljive,  $X$  i  $Y$ , obe inicijalno 0



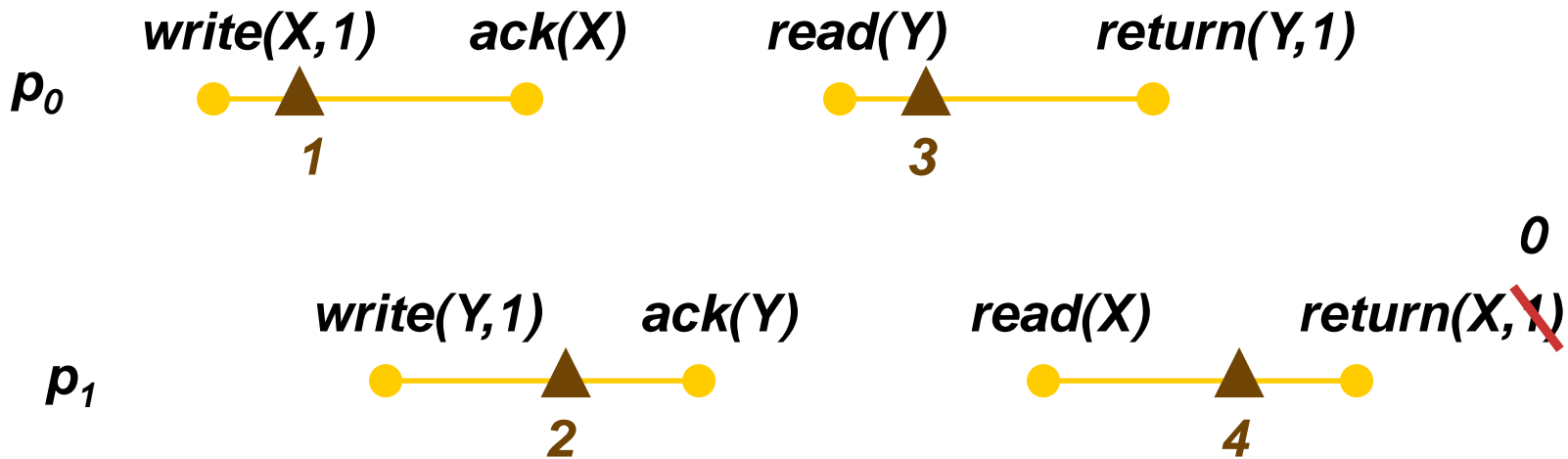
Da li se ova sekvenca može linearizovati?

Da – braon trouglovi.

# Primeri mogućnosti linearizacije

12

Neka postoje dve deljene promenljive,  $X$  i  $Y$ , obe inicijalno 0



Da li se ova sekvenca može linearizovati?

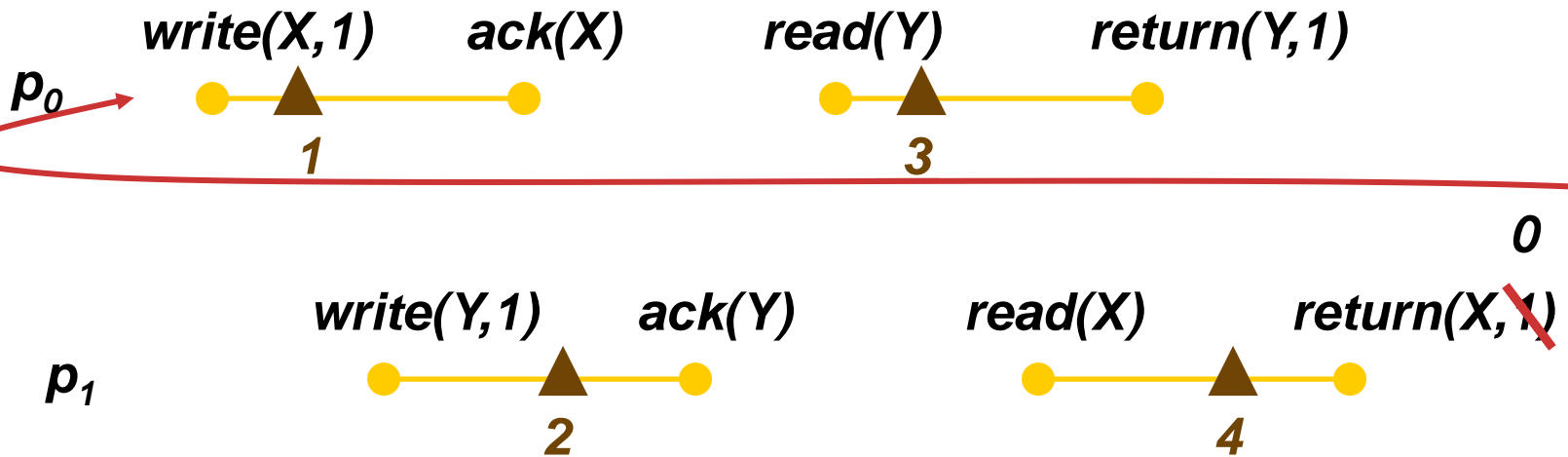
Da – braon trouglovi.

Šta ako  $read$  od  $p_1$  vrati 0?

# Primeri mogućnosti linearizacije

13

Neka postoje dve deljene promenljive,  $X$  i  $Y$ , obe inicijalno 0



Da li se ova sekvenca može linearizovati?

Da – braon trouglovi.

Šta ako  $\text{read}$  od  $p_1$  vrati 0?

Ne – vidi strelicu.

# Definicija sekve. konzistentnosti (Sequential Consistency)

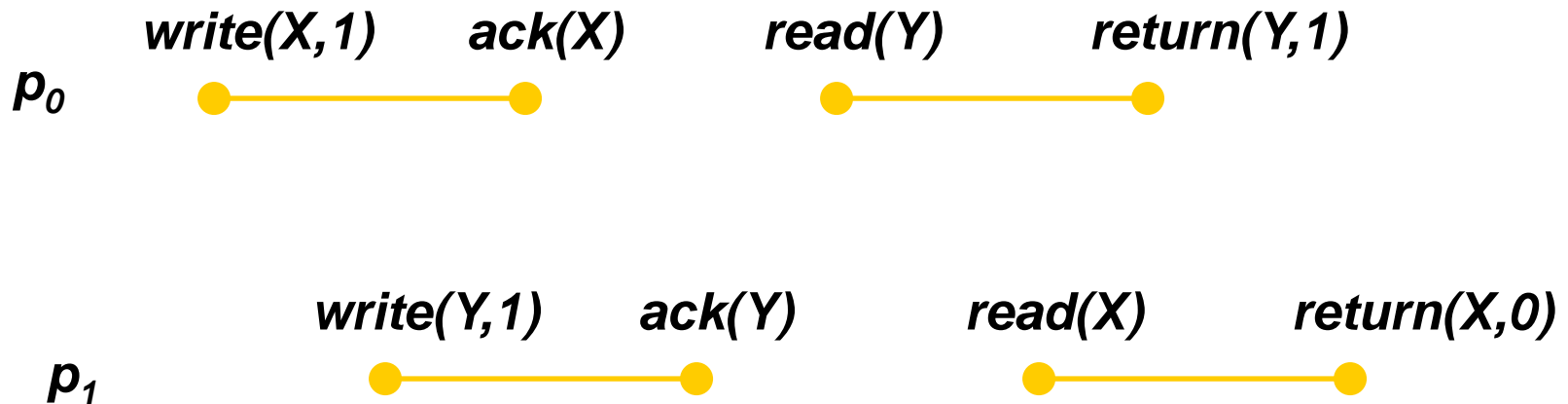
14

- Neka je  $\sigma$  sekvenca poziva i odgovora za neki skup operacija
- $\sigma$  je **sekve. konzistentna** ako postoji permutacija  $\pi$  svih operacija u  $\sigma$  takva da:
  - $\pi | X$  je legalna (zadovoljava sekve. spec.) za sve prom.  $X, i$
  - ako se odgovor za operaciju  $O_1$  desio u  $\sigma$  pre poziva operacije  $O_2$  *u istom procesu*, onda se  $O_1$  desila u  $\pi$  pre  $O_2$  ( $\pi$  respektuje redosled realnog vremena operacija *od strane istog procesa* u  $\sigma$ )

# Primeri sekvencijalne konzistentnosti

15

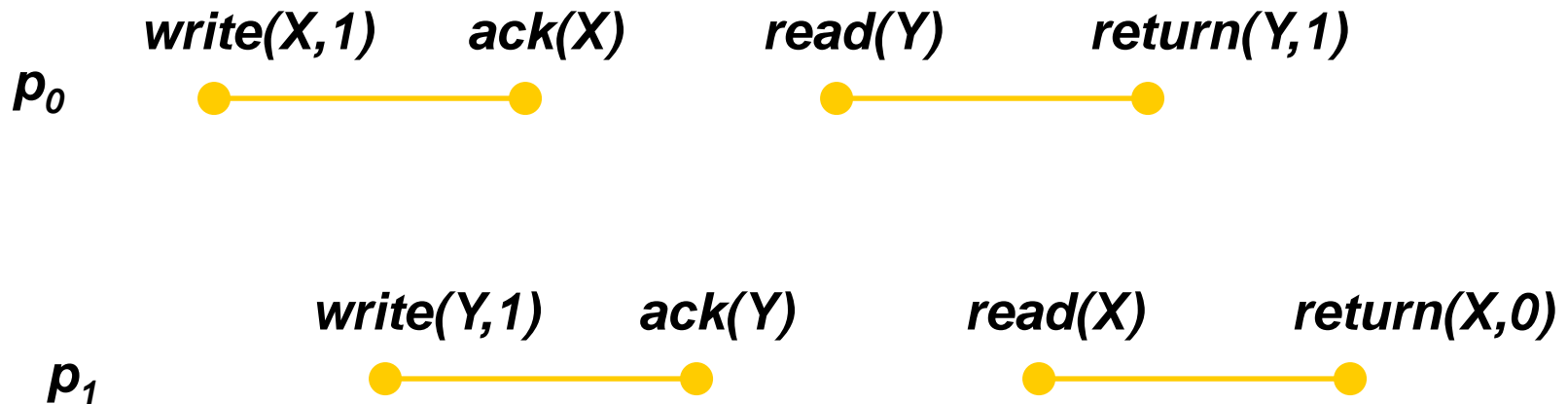
Neka postoje dve deljene promenljive,  $X$  i  $Y$ , obe inicijalno 0



# Primeri sekvencijalne konzistentnosti

16

Neka postoje dve deljene promenljive,  $X$  i  $Y$ , obe inicijalno 0



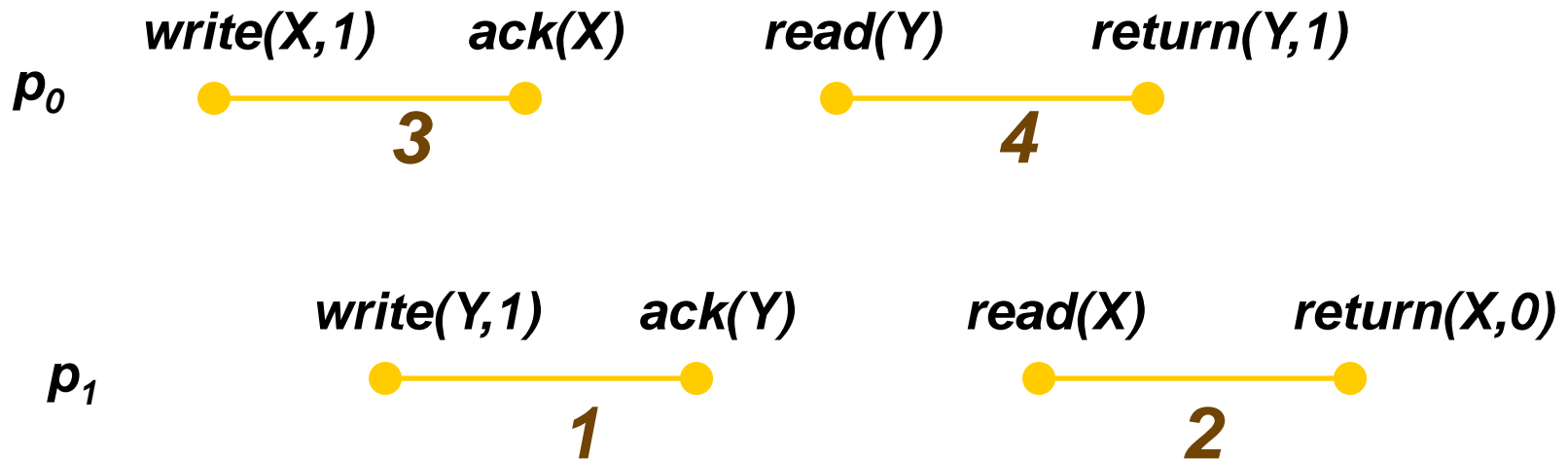
Da li je ova sekvenca sekve. konzistentna?



# Primeri sekvencijalne konzistentnosti

17

Neka postoje dve deljene promenljive,  $X$  i  $Y$ , obe inicijalno 0



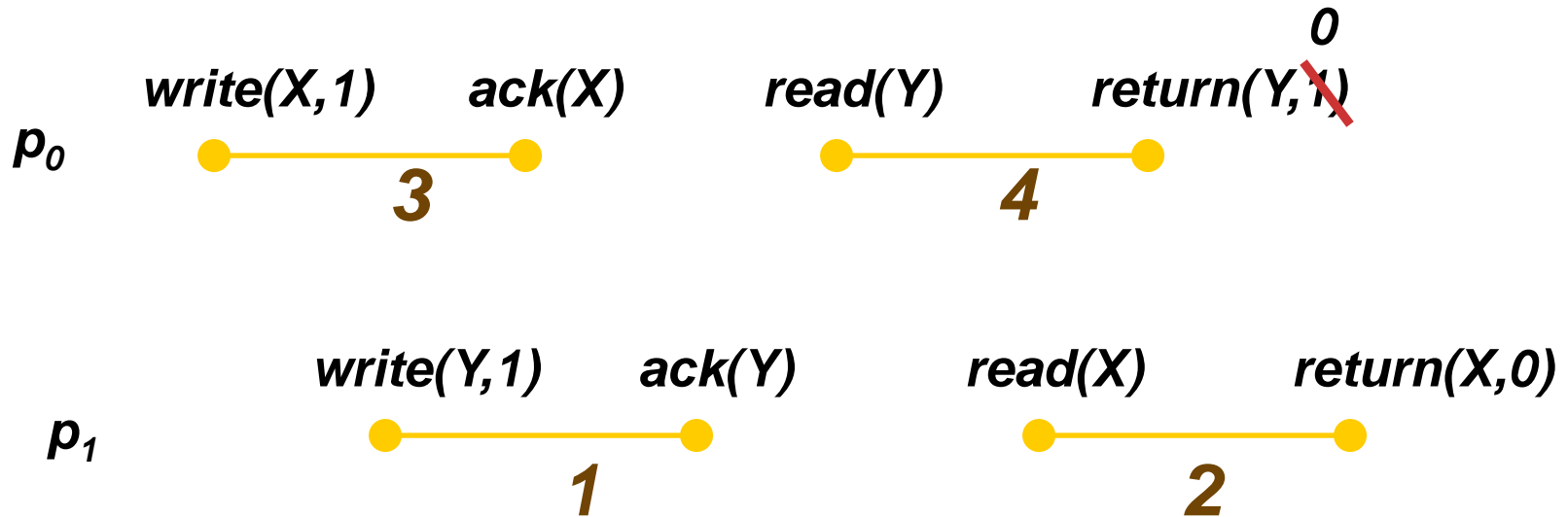
Da li je ova sekvenca sekvencijalno konzistentna?

Da – braon brojevi.

# Primeri sekvencijalne konzistentnosti

18

Neka postoje dve deljene promenljive,  $X$  i  $Y$ , obe inicijalno 0



Da li je ova sekvenca sekve. konzistentna?

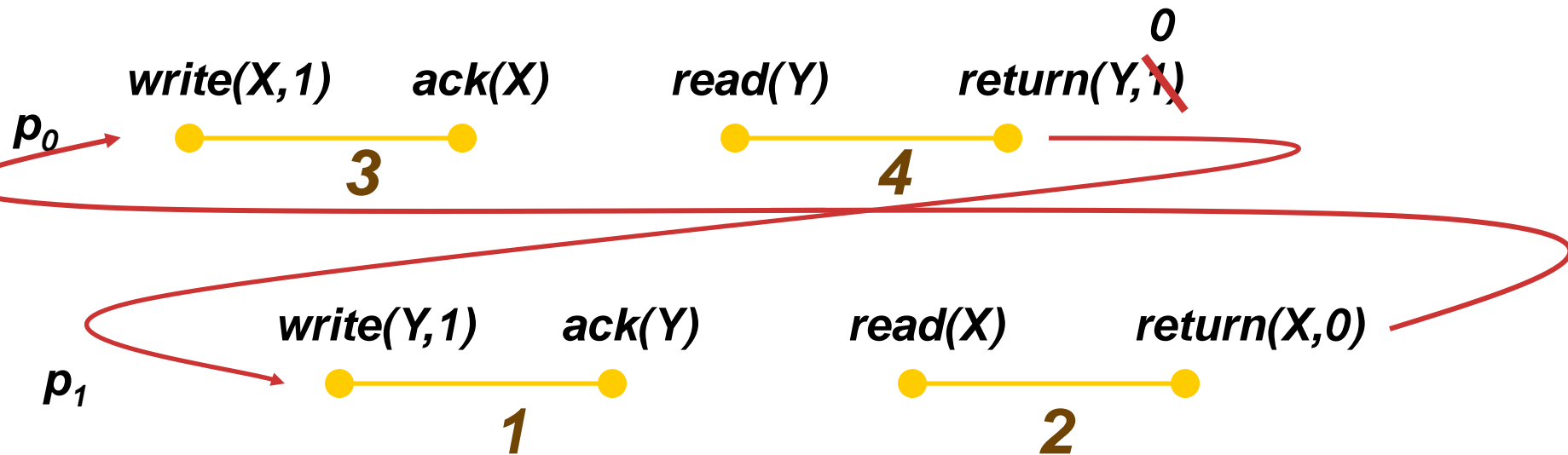
Da – braon brojevi.

Šta ako  $read$  od  $p_1$  vrati 0?

# Primeri sekvencijalne konzistentnosti

19

Neka postoje dve deljene promenljive,  $X$  i  $Y$ , obe inicijalno 0



Da li je ova sekvenca sekve. konzistentna?

Da – braon brojevi.

Šta ako  $read$  od  $p_1$  vrati 0?

Ne – vidi strelice.

# Specifikacija deljene memorije čije sekvence se mogu linearizovati

20

- Ulazi su pozivi operacija deljenih objekata
- Izlazi su odgovori iz deljenih objekata
- Sekvenca  $\sigma$  je u dopuštenom skupu ako i samo ako:
  - ▣ *Korektna interakcija*: svaki proc. ima naizmenične pozive i odgovarajuće odzive
  - ▣ *Životnost*: svaki poziv ima odgovarajući odgovor
  - ▣ *Mogućnost linearizacije*:  $\sigma$  se može linearizovati

# Specifikacija deljene memorije koja je sekvencijalno konzistentna

21

- Ulazi su pozivi operacija deljenih objekata
- Izlazi su odgovori iz deljenih objekata
- Sekvenca  $\sigma$  je u dopuštenom skupu ako i samo ako:
  - ▣ *Korektna interakcija*: svaki proc. ima naizmenične pozive i odgovarajuće odzive
  - ▣ *Životnost*: svaki poziv ima odgovarajući odgovor
  - ▣ *Sekvencijalna konzistentnost*:  $\sigma$  je sekvencijalno konzistentna

# Algoritam deljene memorije čije sekvence se mogu linearizovati

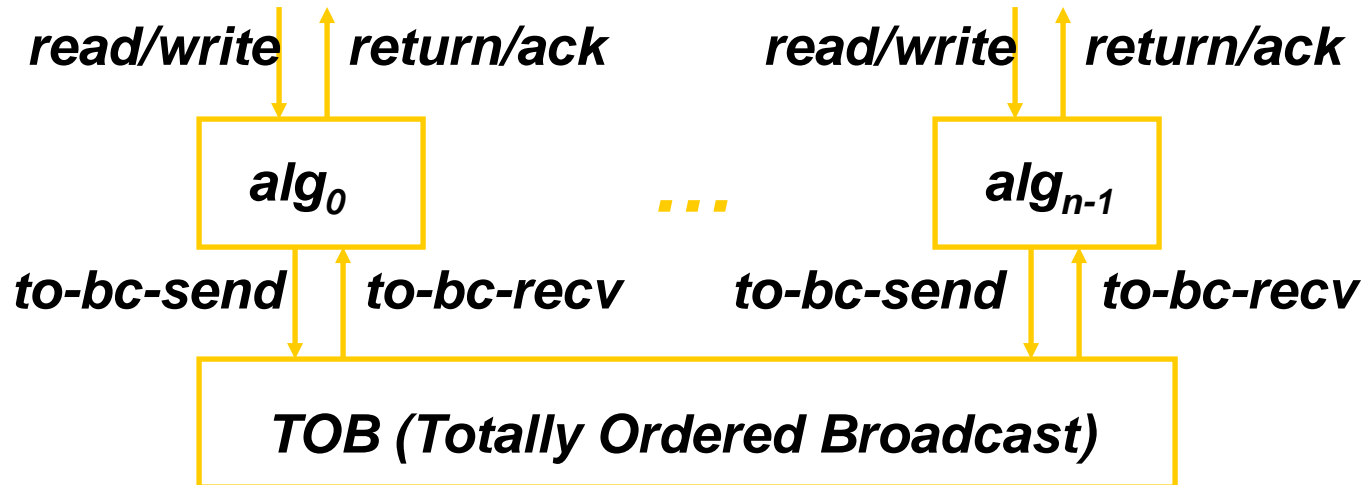
22

- Koristi slanje poruke svima (broadcast) u skladu sa totalnim redosledom (total order)
- Svaki proc održava repliku za svaku deljenu promenljivu
- Kad stigne zahtev za čitanje:
  - ▣ pošalji bcast por. koja sadrži zahtev
  - ▣ kad stigne sopstvena bcast por., vrati vrednost lokalne replike
- Kad stigne zahtev za pisanje:
  - ▣ pošalji bcast por. koja sadrži zahtev
  - ▣ nakon prijema, svaki proc ažurira svoju repliku
  - ▣ kad stigne sopstvena bcast por., odgovori sa ack

# Simulacija

23

*korisnici deljene memorije*



# Korektnost ovog algoritma

24

- Neka je  $\alpha$  bilo koje prihvatljivo izvršenje algoritma za koje važi:
  - ▣ slanje poruka svima u skladu sa totalnim redosledom radi ispravno
  - ▣ korisnička interakcija je ispravna (naizmenični pozivi i odgovori)
- Pokazati da  $\sigma$  (restrikcija  $\alpha$  na događaje na vrhu interfejsa) zadovoljava Životnost i Mogućnost linearizacije



# Korektnost ovog algoritma

25

- Životnost (svaki poziv ima svoj odgovor): po svojstvu životnosti TOB (Totally Ordered Broadcast)
- Mogućnost linearizacije: Definišimo permutaciju operacija  $\pi$  da bude u redosledu u kom su odgovarajuće bcast por. primljene:
  - ▣  $\pi$  je legalna: jer su sve operacije konzistentno poređane po TOB.
  - ▣  $\pi$  respektuje redosled operacije u realnom vremenu: ako se  $O_1$  završi pre nego  $O_2$  počne, bcast od  $O_1$  se obavlja pre bcast od  $O_2$

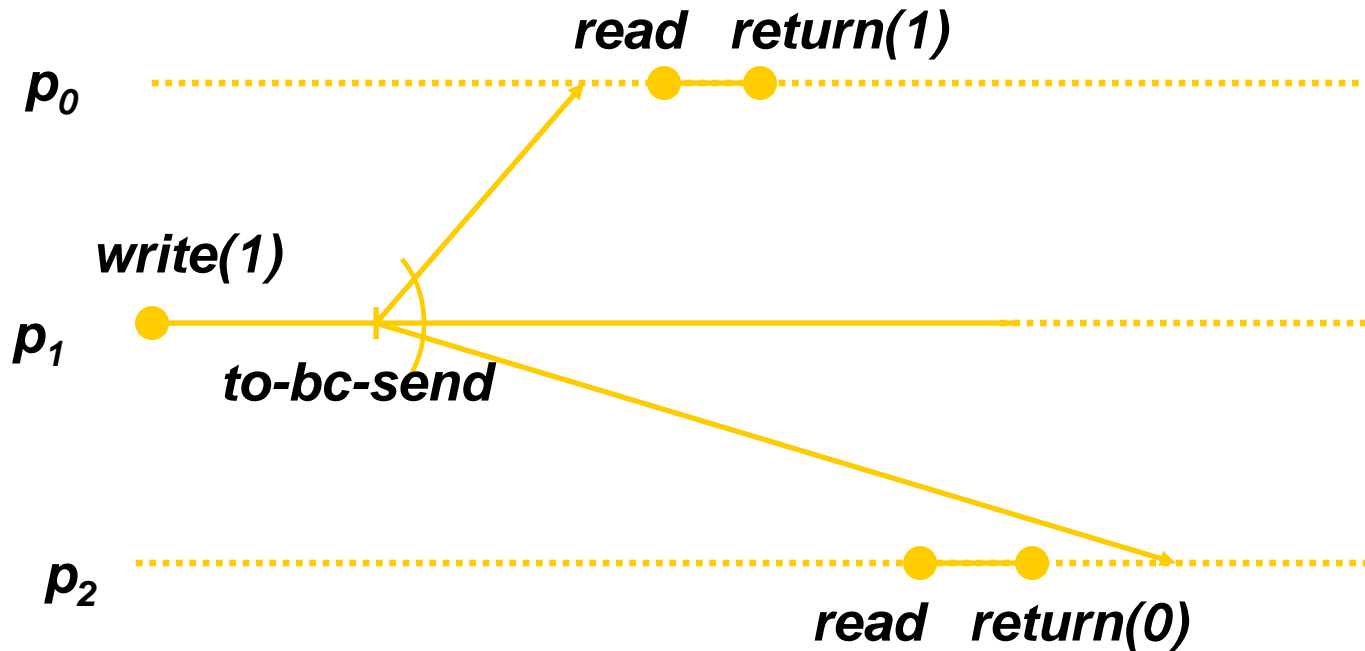
# Zašto je potrebno slanje svima (Bcast) kod operacije čitanja?

26

- Slanje svima kod read operacije ne izaziva promene replika, samo kasni odgovore na read op
- Zašto je ono ipak potrebno?
- Pogledajmo šta se dešava ako ga uklonimo

# Zašto je read bcast neophodan?

27



**Ne može se linearizovati!**

# Algoritam za sekvencijanu konzistentnost (SC algoritam)

28

- Kao predhodni algoritam, ali bez bcast za čitanja:
- Koristi TOB
- Svaki procesor održava repliku svake deljene promenljive
- Kad stigne zahtev za čitanje:
  - ▣ odmah vrati vrednost lokalne replike
- Kad stigne zahtev za pisanje:
  - ▣ pošalji bcast por. koja sadrži zahtev
  - ▣ nakon prijema, svaki proc ažurira svoju repliku
  - ▣ kad stigne sopstvena bcast por., odgovori sa ack

# Korektnost SC algoritma

29

**Lema (8.3):** Lokalne kopije svakog proc. uzimaju sve vrednosti iz write operacija, u istom redosledu, koji očuvava **redosled nepreklopljenih write operacija**

- ▣ implicira da je redosled write op po proc. očuvan

**Lema (8.4):** Ako  $p_i$  piše u  $Y$  i kasnije čita  $X$ , onda  $p_i$ -jevo ažuriranje njegove lokalne kopije  $Y$  (unutar write op) predhodi njegovom čitanju svoje lokalne kopije  $X$  (unutar read op)

# Korektnost SC algoritma

30

## (Teorema 8.5) Zašto važi SC?

- Za bilo koje prihvatljivo izvršenje  $\alpha$ , moramo doći do permutacije  $\pi$  operacija deljene mem. koja je:
  - legalna i
  - respektuje redosled operacija po proc.

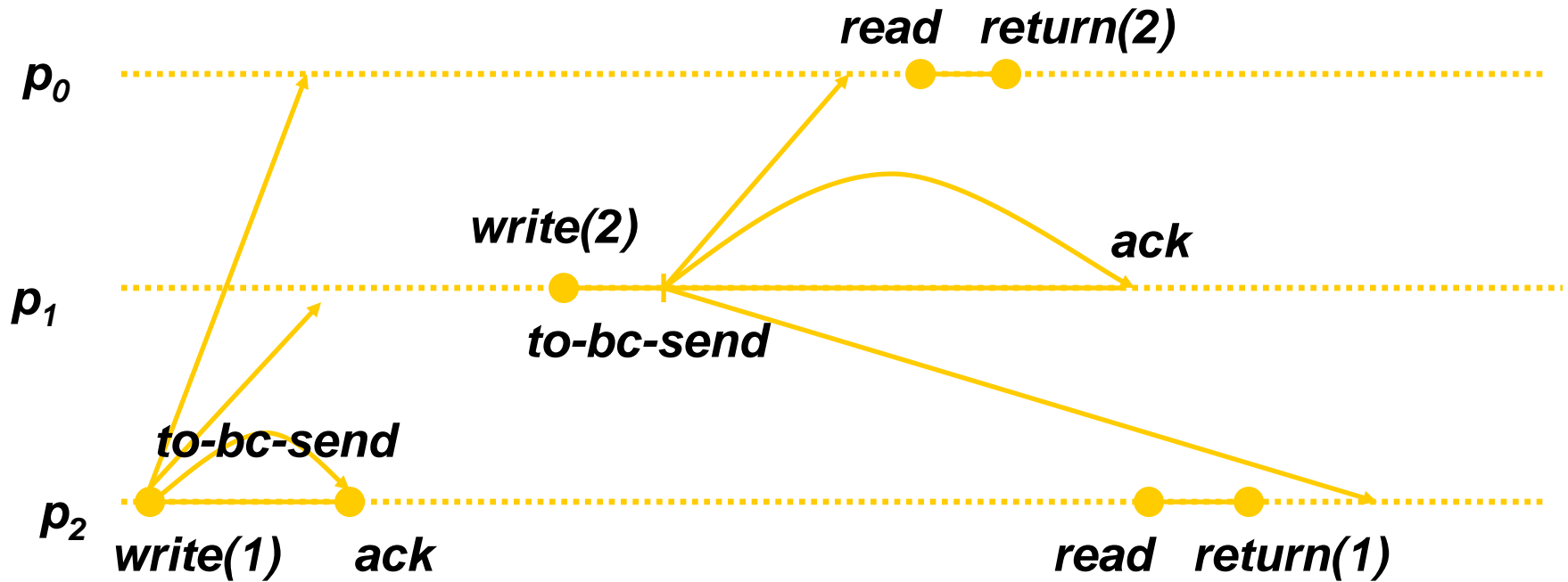
# Permutacija $\pi$

31

- Ubacite sve write op u  $\pi$  u njihovom TOB redosledu
- Razmotrimo svaki read  $R$  u  $\alpha$  u redosledu poziva:
  - ▣ neka je  $R$  čitanje  $X$  od  $p_i$
  - ▣ stavi  $R$  u  $\pi$  odmah iza poslednje od ove 2 op:
    1. operacija od  $p_i$  koja neposredno predhodi  $R$  u  $\alpha$ , i
    2. Write op „iz koje  $R$  čita“ (izaziva poslednje ažuriranje  $p_i$ -jeve lokalne kopije  $X$ , koje predhodi odgovoru za  $R$ )

# Primer permutacije

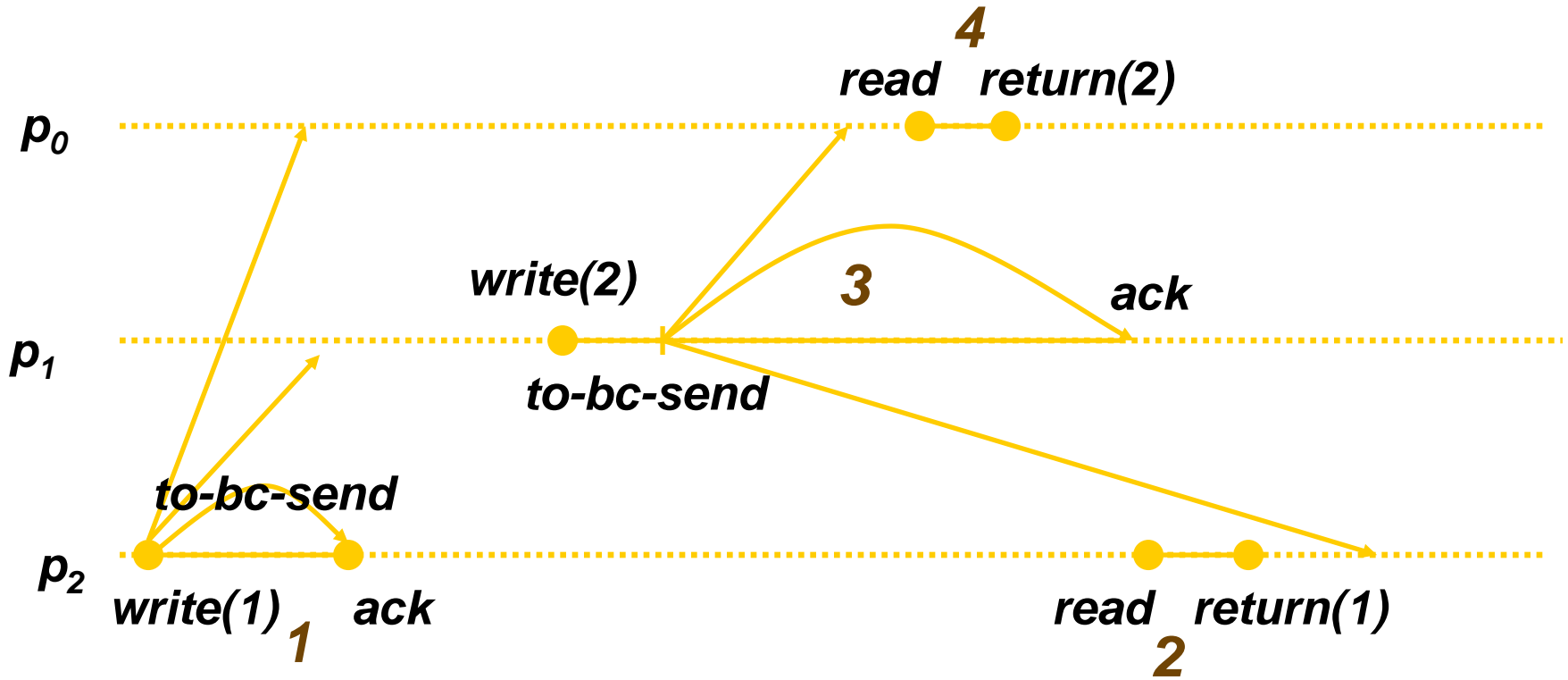
32





# Primer permutacije

33



*permutacija je data sa braon brojevima*

# Performansa SC algoritma

34

- Read operacije su implementirane „lokalno“, i ne zahtevaju međuprocesnu komunikaciju
- Zato se read op može posmatrati kao „brza“: vreme između poziva i odgovora je vreme lokane obrade
- Vreme za write op je vreme isporuke jedne to-bcast por (zavisi od implementacije to-bcast)

# Alternativni SC algoritam

35

- Postoji alternativni alg koji ima obrnutu performansu:
  - ▣ write op je lokalna/brza (iako se šalju bcast-ovi, ne čeka se da oni budu primljeni)
  - ▣ read op može zahtevati da neki bcast-ovi budu primljeni
- Kao i predhodni SC algoritam, ni ovaj ne obezbeđuje mogućnost linearizacije sekvence op deljene memorije