



**Diplomski studij**

**Informacijska i  
komunikacijska tehnologija:**

Telekomunikacije i informatika

**Računarstvo:**

Programsko inženjerstvo i  
informacijski sustavi

Računarska znanost

**Ak.g. 2008./2009.**

# Raspodijeljeni sustavi

## 10.

### Vrednovanje performansi raspodijeljenih aplikacija

Dr. Dalibor F. Vrsalović  
dalibor.f.vrsalovic@fer.hr

15.1.2009

- ◆ **Uvod**
- ◆ **Dio I: Raspodijeljeni sustavi u praksi**
- ◆ **Dio II: Vrednovanje performansi aplikacija teorijom repova**
- ◆ **Dio III: Primjer analize performansi web aplikacije**
- ◆ **Rekapitulacija**

Zašto je vrednovanje performansi sustava  
važno u praksi ?

- ◆ Organizacija prodaje putem Interneta. Aplikacija za prodaju ima sljedeće značajke:
  - ◆ Neuspješni posjeti zbog loše kvalitete usluge
    - ◆ **60 %** kupaca napušta Web stranicu aplikacije ako je odziv aplikacije **između 4 i 6 sekundi**
    - ◆ **95 %** kupaca napušta Web stranicu aplikacije ako je odziv aplikacije **veći od 6 sekundi**
  - ◆ Uspješni posjeti s ostvarenom prodajom
    - ◆ **5 %** kupaca od svih koji su posjetili Web stranicu aplikacije kupi proizvode za **prosječnu cijenu 1200 kn**

- ◆ **Projektiranje i održavanje Web aplikacije ostvaruje se u skladu s očekivanim brojem i porastom korisnika aplikacije**
- ◆ **Ako se broj posjeta Web aplikaciji poveća za 30%, 60% ili 90%:**
  - ◆ Da li će odziv aplikacije biti zadovoljavajući ?
  - ◆ U kojim uvjetima će odziv aplikacije preći u nezadovoljavajuće područje ?
  - ◆ Koliki gubitak prihoda uzrokuje gubitak kupaca zbog slabog odziva aplikacije ?
  - ◆ Koje investicije su potrebne da se uz povećanje prometa zadrži sva posao ?
  - ◆ Kada će se, uz trenutačni trend, potreba za kapacitetom udvostručiti ?

# Rezultati analize značajki Web aplikacije



## Povećanje broja korisnika

	Danas	+30 %	+60 %	+90 %
Maks. posjeta/sat	900.00	1,170.00	1,440.00	1,710.00
Vrijeme odziva (s)	2.96	3.80	5.31	8.83
Izgubljeni kupci (%)	0.00	0.00	60.00	95.00
Mogući broj prodaja / sat (kn)	45.00	58.50	72.00	85.50
Mogući prihod / sat (kn)	54,000.00	70,200.00	86,400.00	102,600.00
Stvarni prihod / sat (kn)	54,000.00	70,200.00	34,560.00	5,130.00
Izgubljeni prihod / sat (kn)	0.00	0.00	51,840.00	97,470.00

## ◆ Na ovom predavanju saznati će te kako vrednovati performanse aplikacije

- ◆ Linearna ekstrapolacija najčešće daje netočne rezultate
- ◆ Pogrešno projektirana aplikacija može imati katastrofalne posljedice na poslovanje

- ◆ **Sadržaj ovog predavanja nastao je na temelju:**
  - ◆ N.J. Gunther: "**The practical performance analyst**", *Mcgraw Hill i Authors Choice Press*, 1998 i 2000. ISBN 0-595-12674-X (poglavlje 2 i 3)
  - ◆ D.A. Menasce, V.A.F.Almeida: "**Capacity planning for web services**", *Prentice Hall*, 2002 ISBN 0-13-065903-7 (poglavlje 1 i 5)
  - ◆ D.F. Vrsalovic, et. al: "**Performance prediction and calibration for a class of multiprocessors**", *IEEE Transactions on Computers*, Volume: 37 Issue: 11 , Nov. 1988, pp. 1353 -1365

## Za one koji zele dalje istrazivati

---



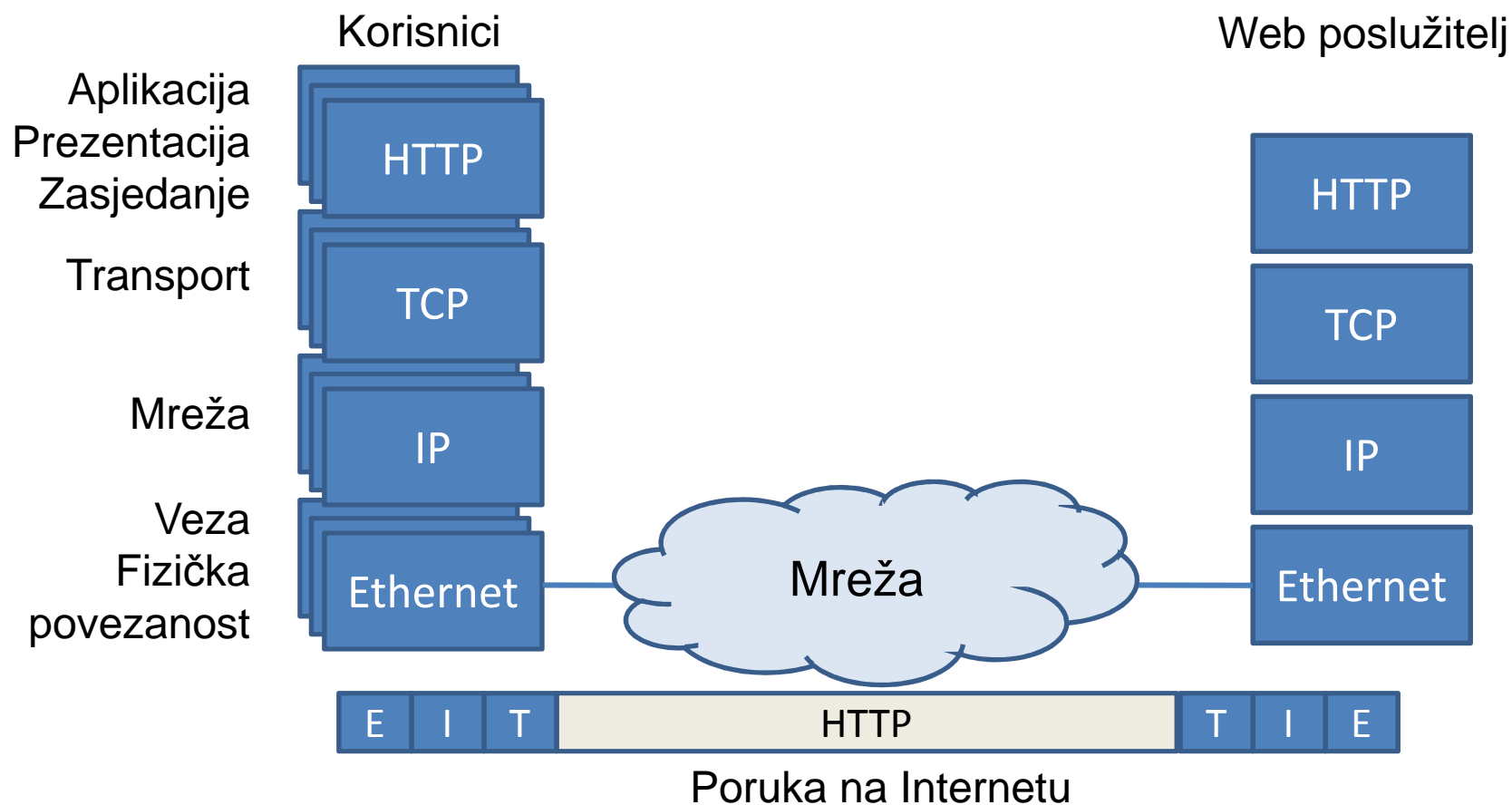
- ◆ A.O. Allen: "**Probability, Statistic, and Queueing Theory with Computer Science Applications**", Academic Press 1978.
- ◆ S. Joines, R. Willenborg, K. Hygh: "**Performance analysis for Java Web Sites**", Addison Wesley, 2003
- ◆ S. Sounders: "**High Performance Web Sites**", O'Reilly, 2007.
- ◆ T. Schlosssnagle: "**Scalable Internet Architectures**", Sams Publishing, 2007.
- ◆ D.A. Menasce, V.A.F. Almeida, L.W. Dowdy: "**Performance by Design**", Prentice Hall, 2004
- ◆ N.J. Gunther: "**Analyzing computer system performance with Perl::PDQ**", Springer, 2005. ISBN 3-540-20865-8.



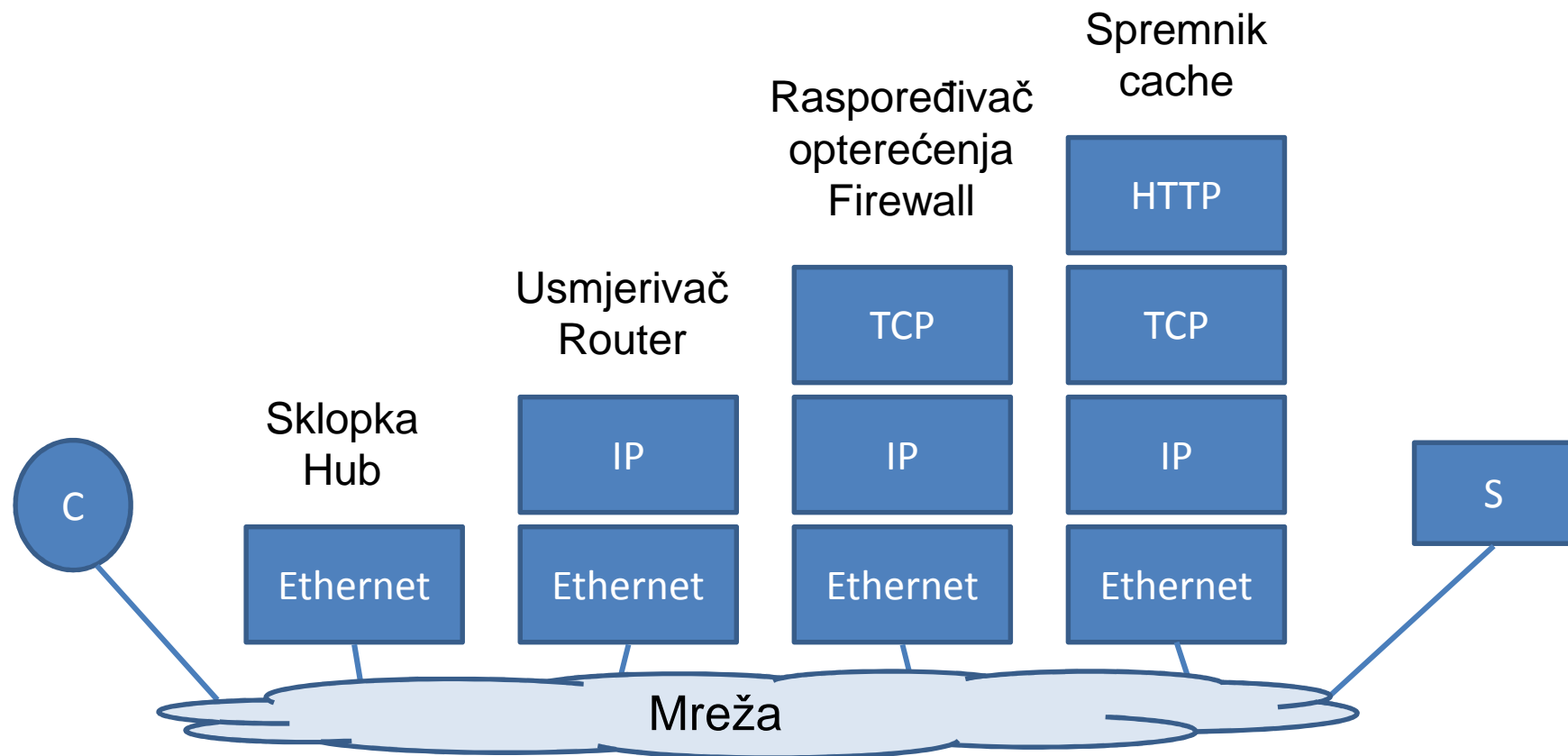
## **Raspodijeljeni sustavi u praksi**

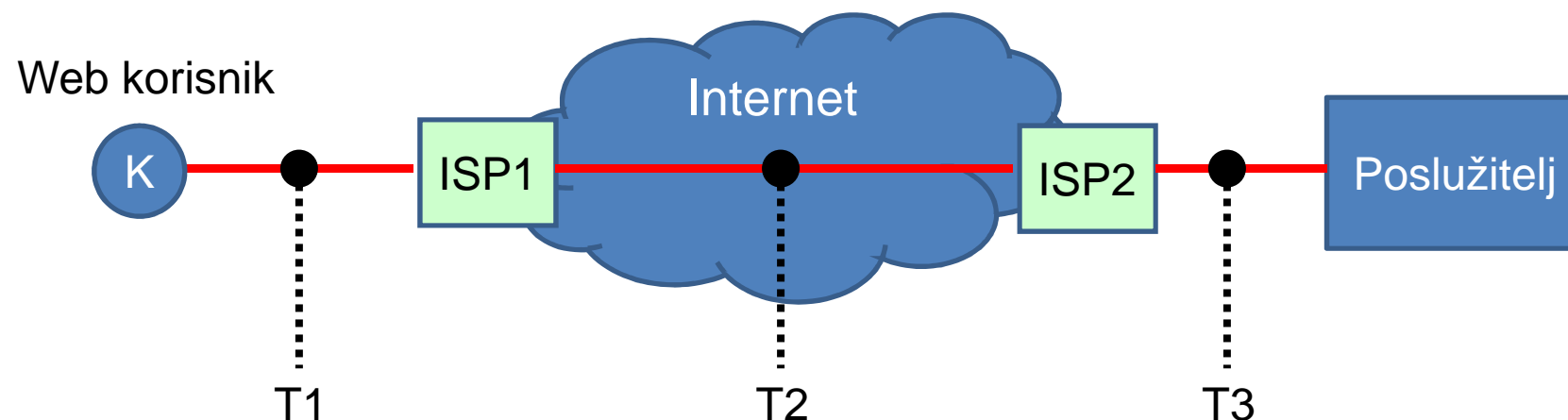
Modeliranje sustava zahtjeva razumijevanje  
sustava

- ◆ Sedam slojeva OSI skupa protokola u praksi se preslikava se na četiri razine protokola u globalnoj mreži Internet



- ◆ Komunikacijski posrednički sustavi koriste se na različitim razinama protokola



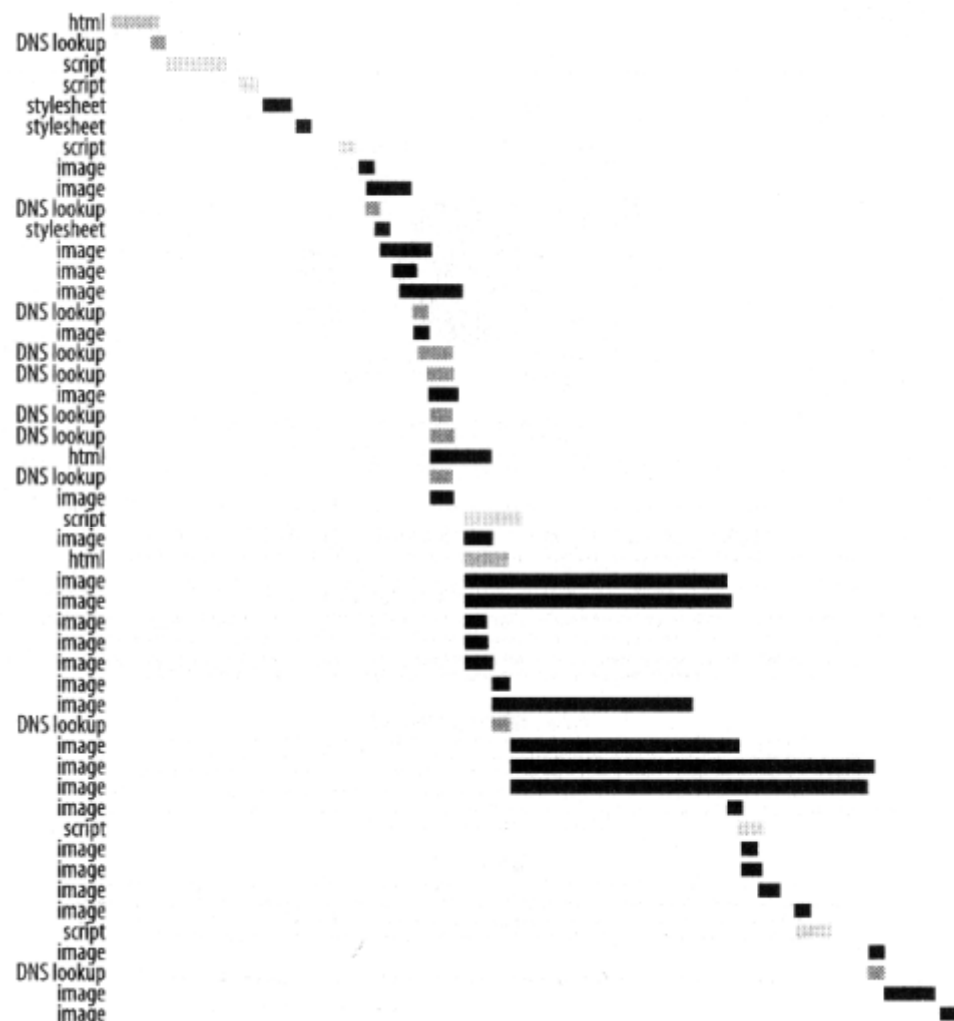


- ◆ **Pružatelj pristupa mreži Internet (ISP, Internet Service Provider)**

- ◆ **Vremena kašnjenja**

- ◆ Korisnik – ISP1 (T1)
- ◆ ISP1 – ISP2 (T2)
- ◆ ISP2 – Poslužitelj (T3)

# Primjer dobavljanja Web stranice: My Space



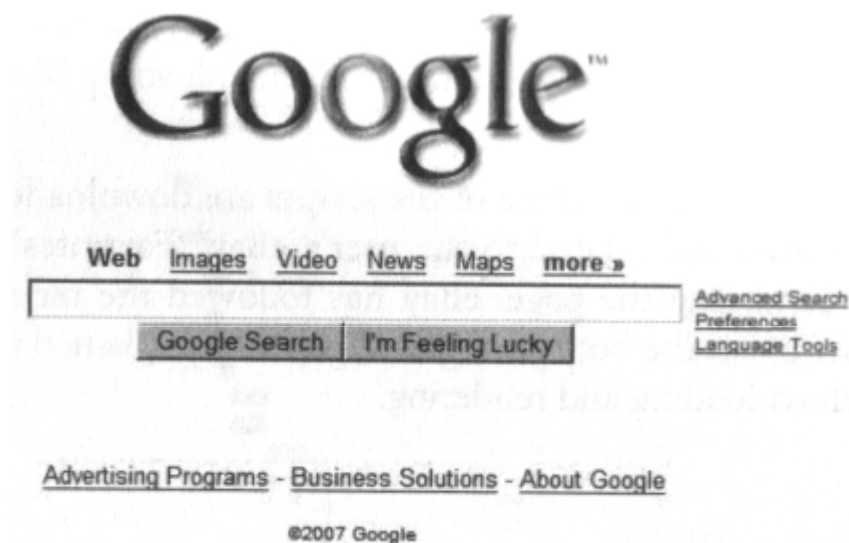
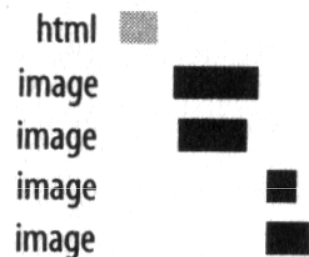
**Vremenski dijagram dohvaćanja objekata stranice**



Veličina: 205K	Vrijeme: 7.8 s
Broj zahtjeva: 39	YSlow: D

**Prikaz stranice pregledniku**

# Bolji primjer: Google



Veličina: 18K	Vrijeme: 1.7 s
Broj zahtjeva: 3	YSlow: A

**Vremenski dijagram dohvaćanja  
objekata stranice**

**Prikaz stranice pregledniku**

- ◆ **Osnovni modeli ostvarivanja razmjernog rasta kapaciteta aplikacija**
  - ◆ **Vertikalno skaliranje** podrazumijeva prijelaz na server sa većim kapacitetom
  - ◆ **Horizontalno skaliranje** podrazumijeva dodavanje paralelnih servera (obično istog kapaciteta).
  - ◆ **Skaliranje prema gore** (većem kapacitetu) je jednako važno kao i **skaliranje prema dolje** (manjem kapacitetu) zbog potrebe da se troškovi prilagode prihodima!

# Horizontalno skaliranje donosi nove izazove...

---



## ◆ Raspoređivanje opterećenja u sustavu

- ◆ Osigurava jednakomjerno opterećenje paralelnih podsustava

## ◆ Raspoređivanje podataka

- ◆ Osiguravanje koherencije (Replikacija)
- ◆ Osiguravanje podjele (Federacija)

## ◆ Protokoli za sinkronizaciju rada grupe

- ◆ Osiguravanje vremenskog slijeda
- ◆ Garantirana dostava

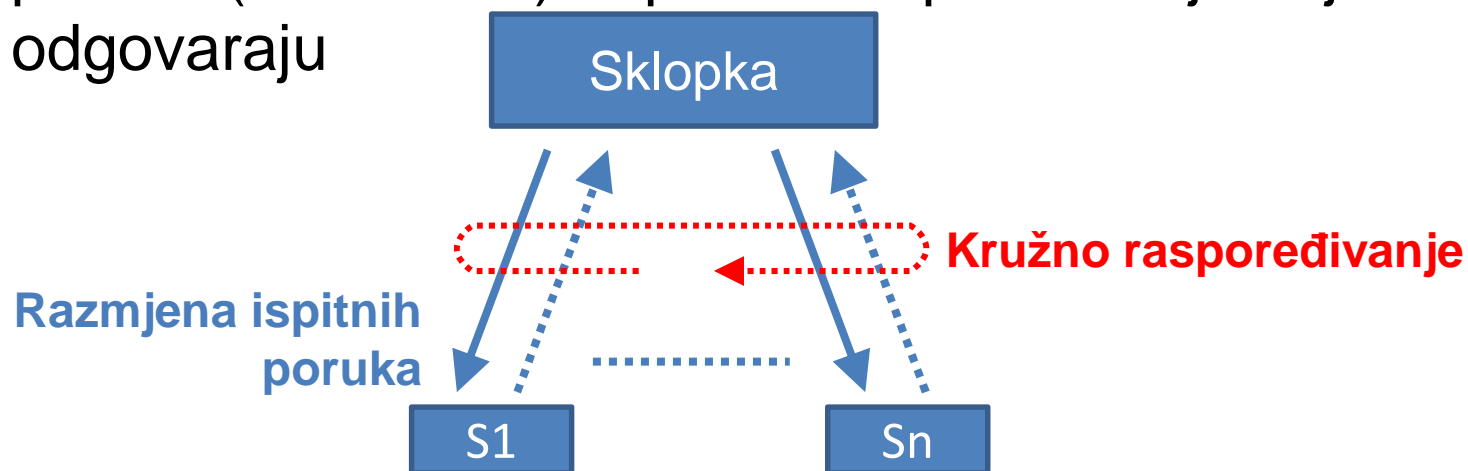


- ◆ **Visoka dostupnost sustava** (*high availability*)
  - ◆ Uobičajena implementacija kroz neosjetljivost na greške (*fault tolerance*)
    - ◆ Neosjetljivost na greške omogućava ostvarivanje dostupnosti Web aplikacija uz prisutnost grešaka u podsustavima
- ◆ Visoki stupanj eliminacije grešaka je preskup (engl. žargon – "platinum tank")
- ◆ Budući da povećanje kapaciteta obično zahtjeva ostvarenje replikacije, replikacija se može iskoristiti za toleranciju pogrešaka

- ◆ **U slučaju web aplikacija, u praksi se koriste oba termina jer se visoka dostupnost osigurava kroz neosjetljivost na greške**
  
- ◆ **U praksi se koriste dva modela organizacije sustava**
  - ◆ Aktivni → Pripravan (Active → Standby)
  - ◆ Aktivni → Aktivni (Active → Active)
  
- ◆ **Sustavi u pripravnim stanju mogu se koristiti kada nije potrebno zapamtiti stanje sustava**
  - ◆ Web serveri ne čuvaju stanje dok baze podataka moraju sačuvati sve transakcije

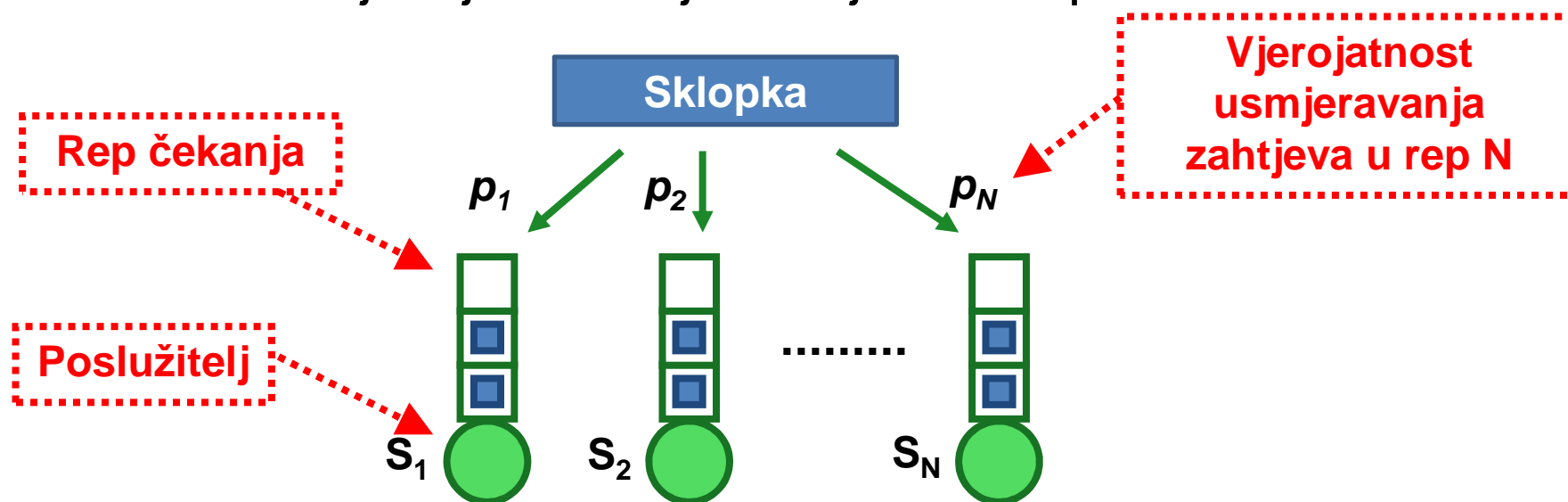
## ◆ Sklopke za raspoređivanje opterećenja

- ◆ Najčešće se ostvaruju na razini transportnog protokola (TCP/IP)
- ◆ Najčešće discipline raspoređivanja je kružno posluživanje (*Round Robin*) i najmanjem opterećenju (*Least Loaded First*)
- ◆ Sklopka provjerava stanje poslužitelja slanjem ispitnih poruka (*heart beat*) te preskače poslužitelje koji ne odgovaraju



## ◆ Elementi modela

- ◆ Sklopka koja proslijeđuje zahtjeve na skup računala
- ◆ Repovi za spremanje dolaznih zahtjeva
- ◆ Poslužitelji koji obrađuju zahtjeve u repu



- ◆ Način posluživanja određuje vjerojatnosti raspoređivanja zahtjeva na poslužitelje ( $p_1 \dots p_n$ )

## ◆ HTTP

- ◆ Obično na razini URL tako da se različite stranice dohvaćaju iz različitih spremnika

## ◆ DNS

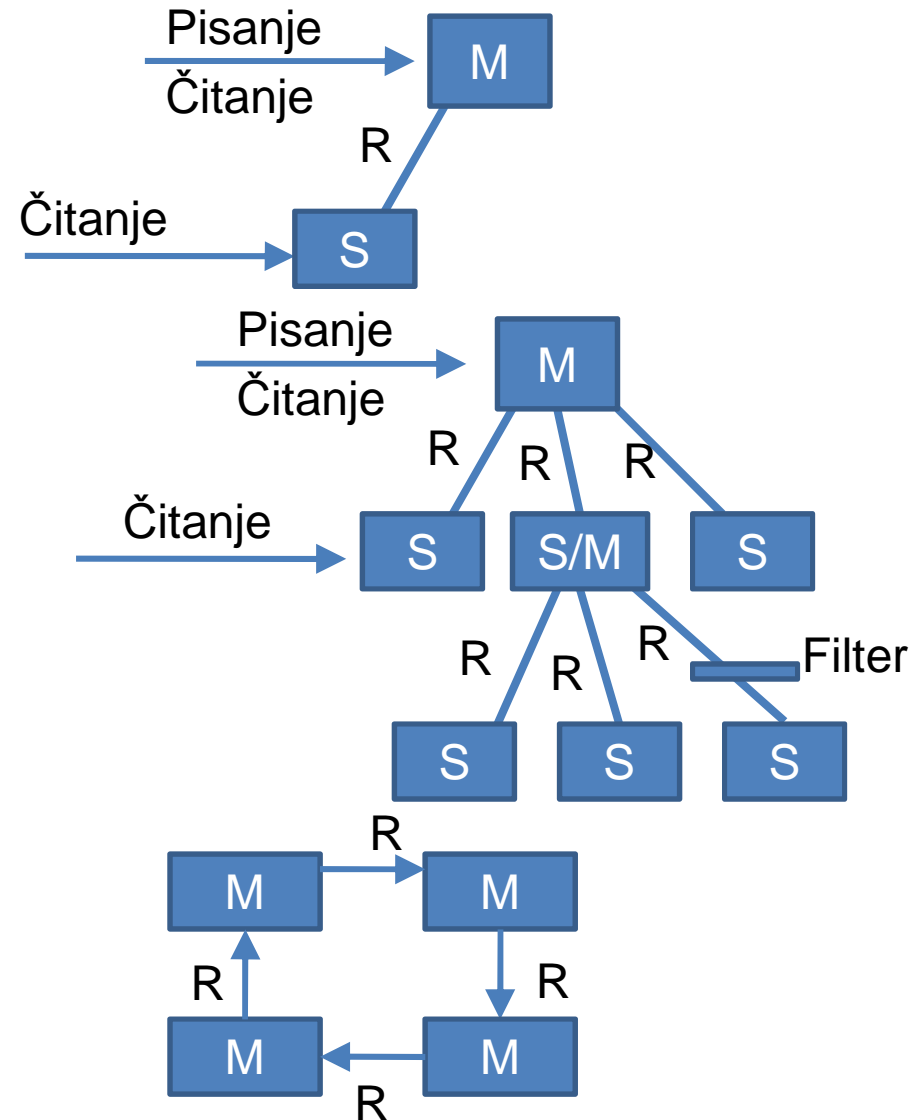
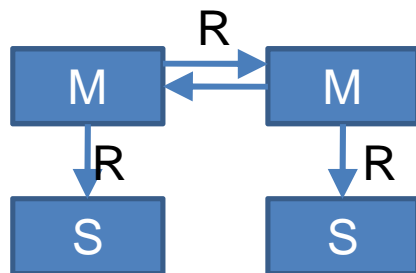
- ◆ Pruža mogućnost geografske podjele

## ◆ SMTP

- ◆ Elektronička pošta je zasnovana na ASCII protokolu te se može lako prosljeđivati prema ciljnoj adresi

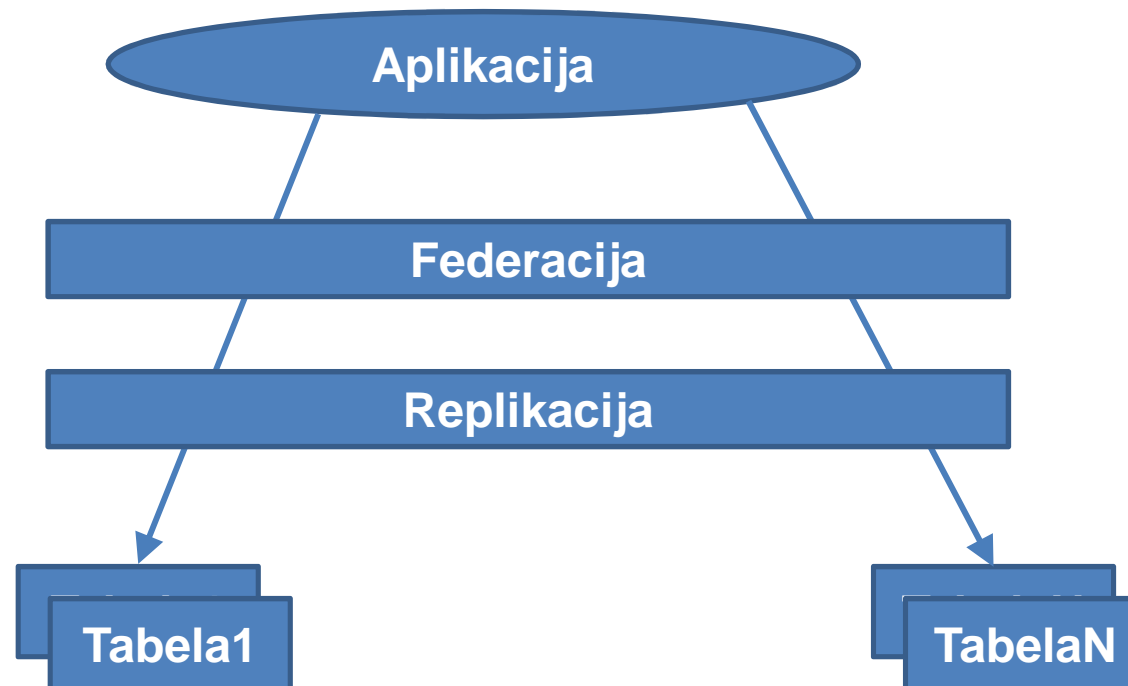
## ◆ Replikacija

- Master – slave
- Master – slaves
- Stablo
- Stablo sa filterima
- Master – master
- Master – master - slaves
- Master ring



## ◆ Particioniranje

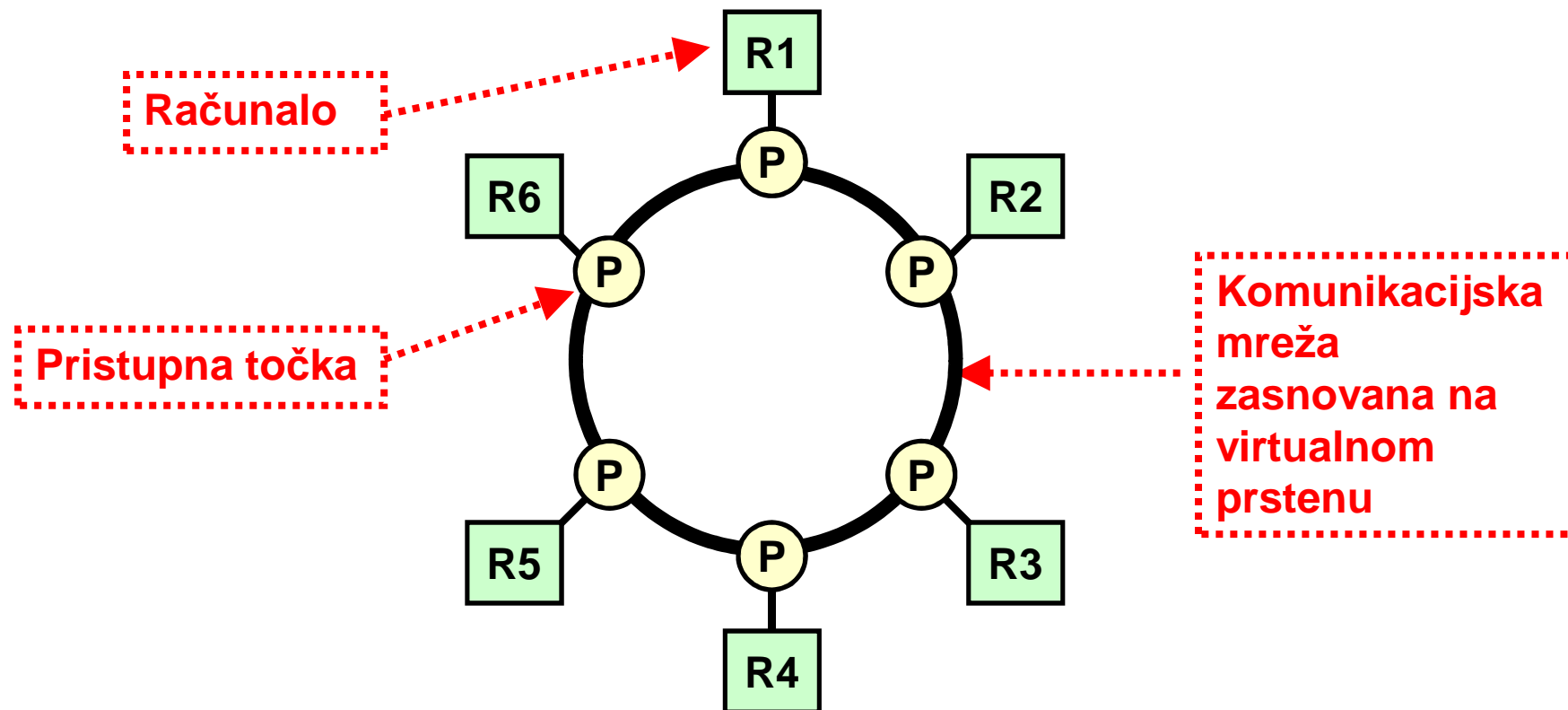
- Clustering
- Federacija
- Federacija replikacija



# Sinkronizacija i protokoli za komunikaciju u grupi

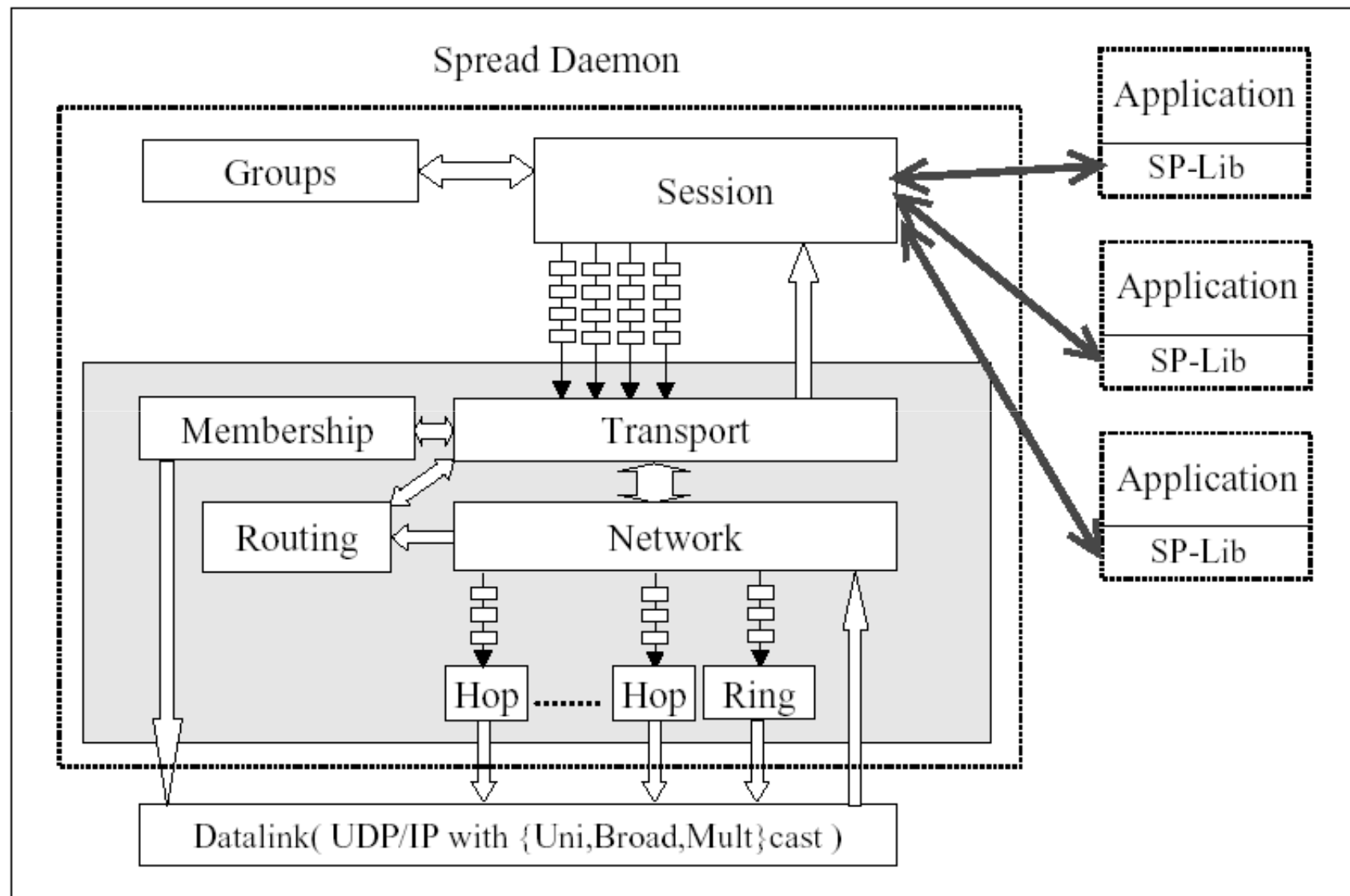


- ◆ **Primjer:** Spread - Wide Area Multicast and Group Communication Toolkit
  - ◆ Center for Networking and Distributed Systems, John Hopkins University





# Arhitektura čvora



- ◆ **Redoslijed prispijeća (*FIFO*)**

- ◆ Proces X šalje poruke A i B u slijedu
- ◆ Svi procesi dobivaju poruku A prije poruke B

- ◆ **Međuzavisnost (*Casual*)**

- ◆ Proces X šalje poruku B nakon prijema poruke A
- ◆ Svi procesi dobivaju poruku B nakon poruke A

- ◆ **Potpuna garancija (*Total*)**

- ◆ Proces X prima poruke A i B u slijedu
- ◆ Svi procesi primaju poruke A i B u istom slijedu

# Tipična arhitektura web poslužitelja



## ◆ Usmjernik (R)

*Router*

## ◆ Sklopka (S)

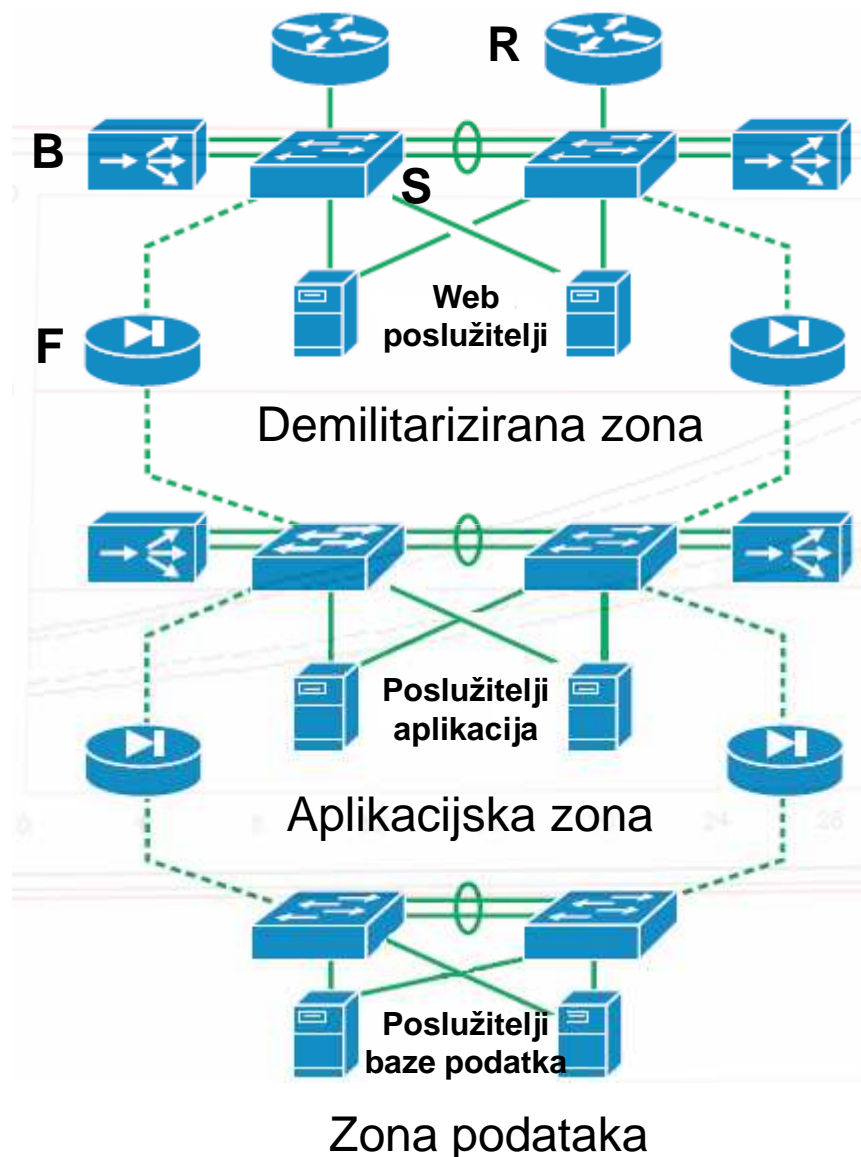
*Request Switch*

## ◆ Raspoređivač (B)

*Load balancer*

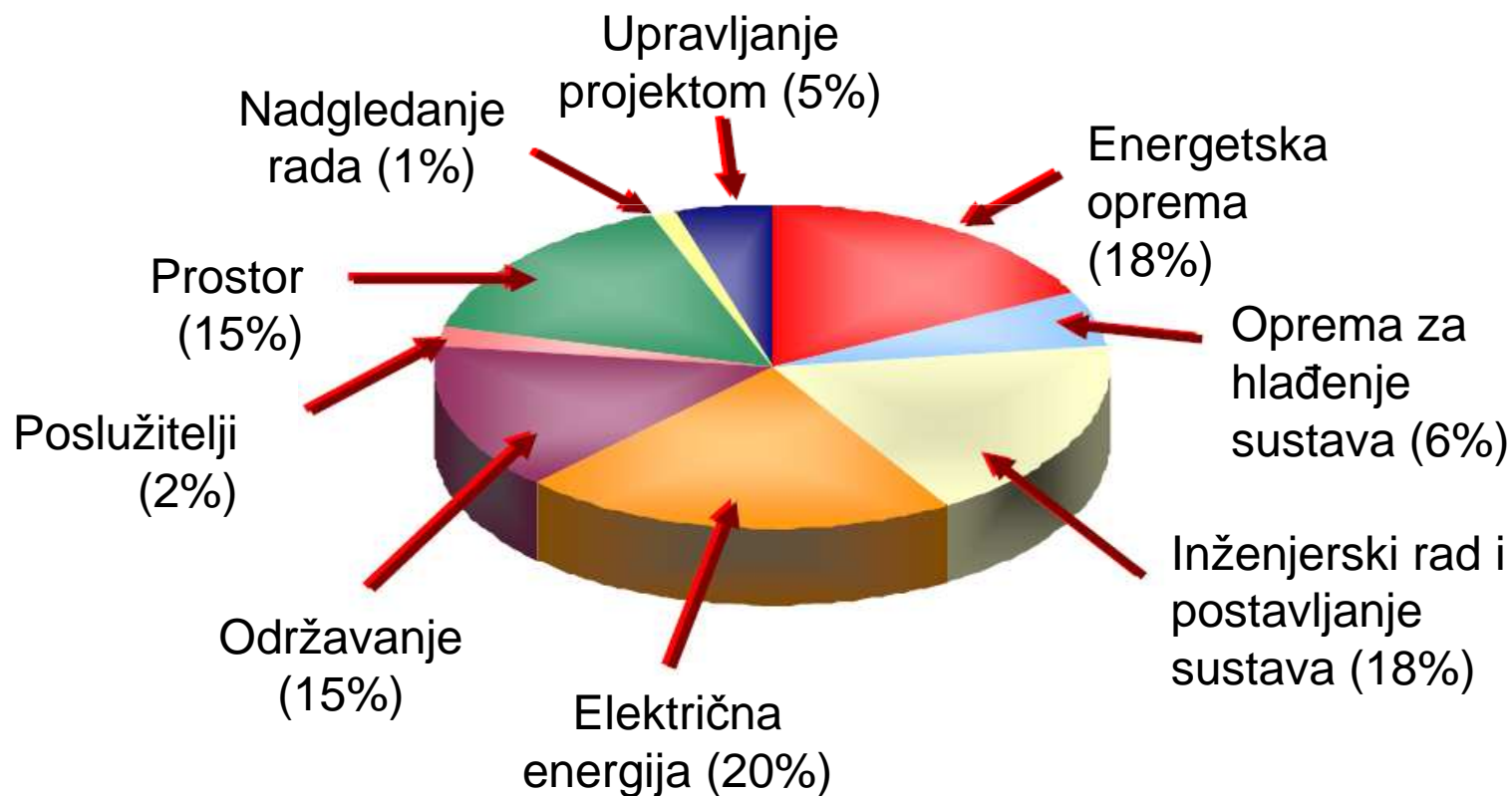
## ◆ Zaštitnik pristupa (F)

*Firewall*



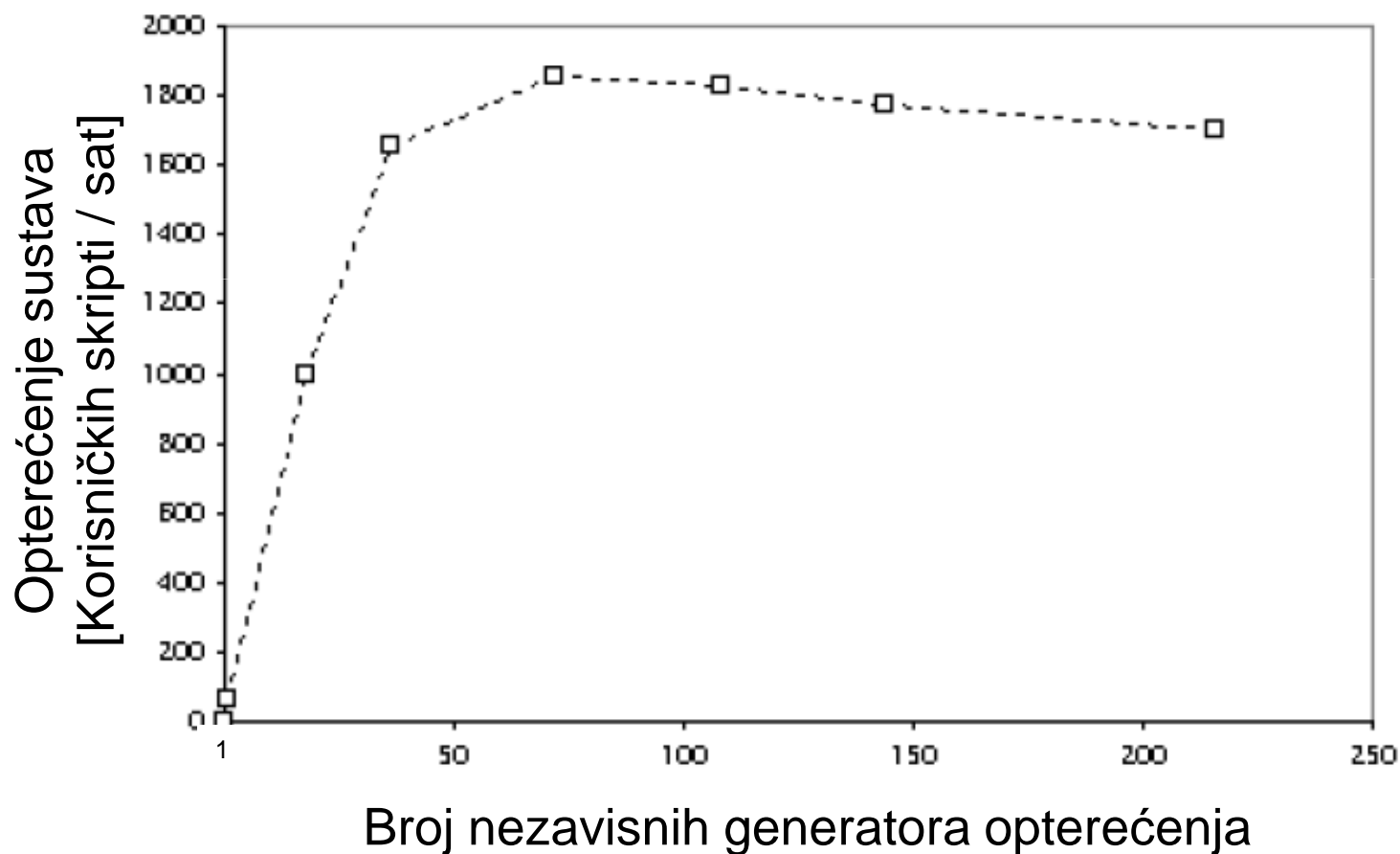
## ◆ Prosječni trošak uporabe sustava poslužitelja u nizu (rack)

◆ \$120K kroz tri godine korištenja (CAPEX/OPEX = 50/50)



- ◆ **Iznajmljivanje infrastrukture poslužitelja**
  - ◆ Udomljivanje sustava
    - ◆ Udomitelj infrastrukture  
*pruža i upravlja fizičkom infrastrukturom, kao što je zgrada, napajanje te pristup Internetu*
    - ◆ Zakupnik infrastrukture  
*postavlja i upravlja sredstvima koja se poslužuju*
  - ◆ Upravljeni sustav
    - ◆ Operacijski sustav s listom poznatih aplikacija
  - ◆ Dedicirani sustavi
  - ◆ Zajednički sustavi
- ◆ **Izgradnja vlastitog poslužitelja**

- ◆ Mjerenja se ostvaruju primjenom generatora ispitnog opterećenja (*Load generatora*)



$$S(\alpha, \beta, N)_{future} = \frac{N}{1 + \alpha[(N-1) + \beta N(N-1)]} \quad [1]$$

Parametri  $\alpha$  i  $\beta$  povezani su sa degradacijom performansi zbog čekanja na zajednička sredstva i kašnjenja zbog održavanja koherencije

Slijedeća formula omogućava modeliranje u Excel-u

$$= \frac{Nr}{1 + A1((Nr-1) + B1 * Nr * (Nr-1))}$$

$A1 = \alpha$  i  $B1 = \beta$

Nr = vrijednost N u redu r

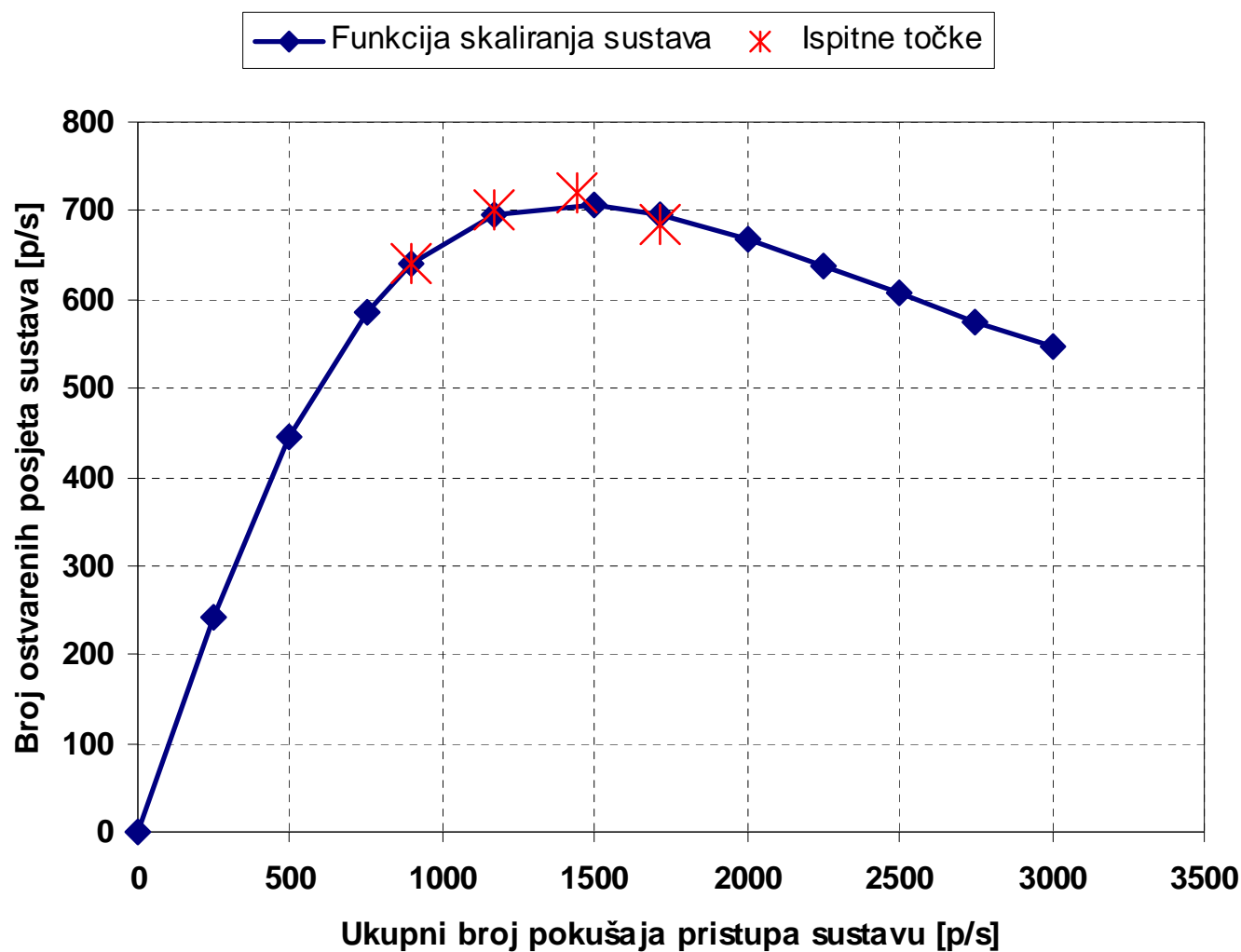
[1] Gunther, N. J., "Hit-and-Run Tactics Enable Guerrilla Capacity Planning", *IEEE IT Professional*, Jul-Aug, 2002, pp. 40-46

## ◆ Osnovni koraci:

- 1) Izmjeri performanse kao funkciju od  $N$  koristeći alate kao *WebLoad* ili *LoadRunner*
- 2) Obično je dovoljan mali uzorak (najmanje 4 točke)
- 3) Odredi  $\alpha$  i  $\beta$  primjenom regresije koristeći alat EXCEL
- 4) Koristi dobivene vrijednosti za izračunavanje skaliranja prema Excel modelu iz prethodnog teksta



# Uporaba regresije za određivanje funkcije skaliranja



Skaliranje.xls

- ◆ Ako se iskorištenje servera  $U$  mjeri u pravilnim intervalima, dugoročna potreba za kapacitetom se može ustanoviti podrazumijevajući eksponencijalni trend model

$$U_{future} = U_{now} * e^{L*W}$$

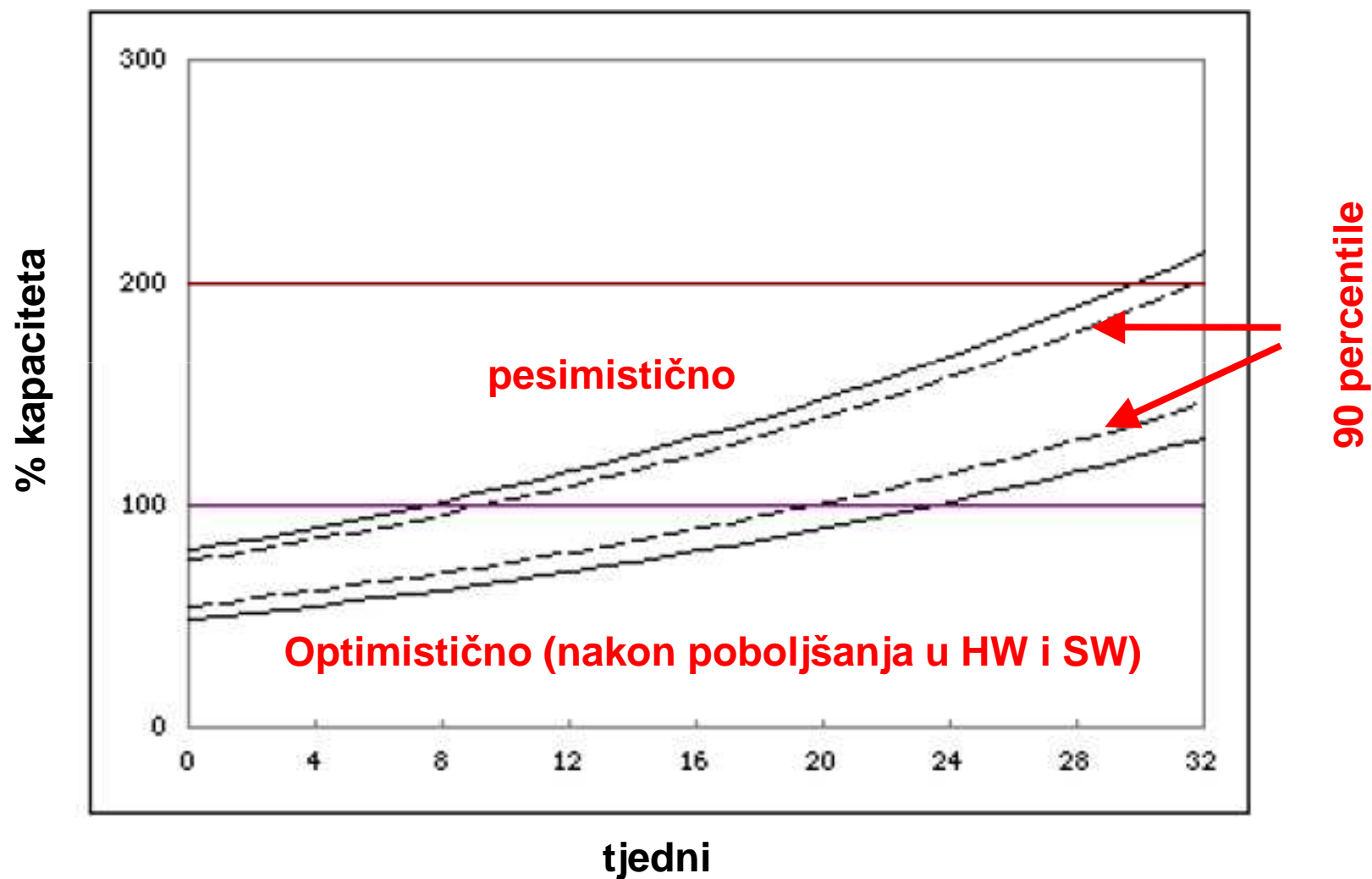
$L$  = trend rasta, određen primjenom Excel opcije "*Add Trendline*"

$W$  = broj tjedana kroz koje je trend aproksimiran

- ◆ Vrijeme do udvostručenja potrebnog kapaciteta može se izračunati primjenom:

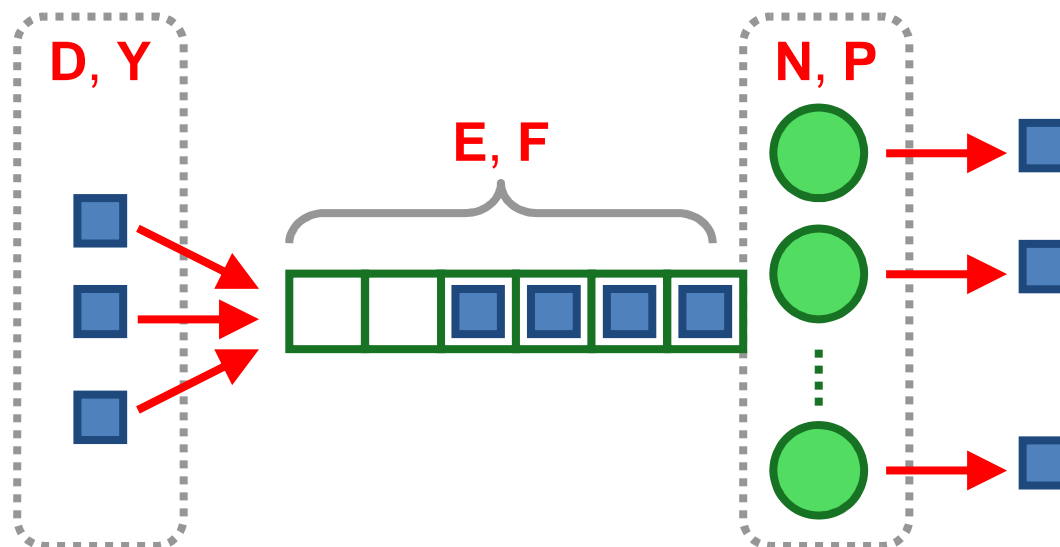
$$T_{double} = \frac{\ln(2)}{L}$$

# Prognoza kapaciteta kroz 32 tjedna



# Vrednovanje performansi sustava teorijom repova

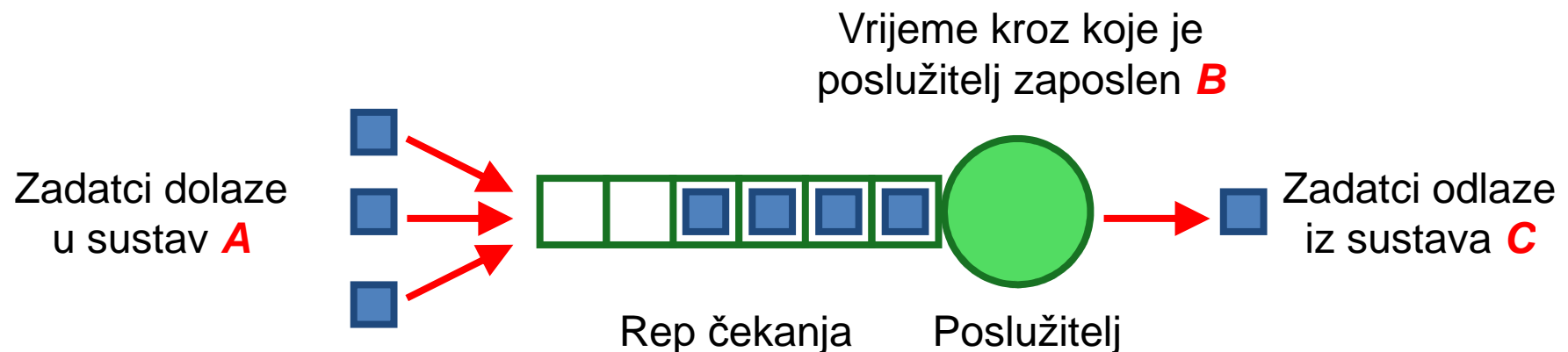
## ◆ Zapis značajki modela **D/P/N/Y/E/F** (*Kendall notacija*)



- ◆ **D** – Razdioba dolazaka zahtjeva (G=Proizvoljna, M=Exp, ...)
- ◆ **P** – Razdioba posluživanja zahtjeva (G=Proizvoljna, M =Exp, ...)
- ◆ **N** – Broj poslužitelja
- ◆ **Y** – Maksimalni broj zahtjeva (m zahtjeva ili  $\infty$  zahtjeva)
- ◆ **E** – Maksimalni kapacitet repa (m ćelija ili  $\infty$  ćelija)
- ◆ **F** – Disciplina posluživanja (FIFO, LIFO, ...)

## ◆ Model jednogprocesorskog sustava M/M/1

### ◆ Poslužitelj i rep čekanja



## ◆ Osnovne značajke modela

- ◆ **T** – Ukupno vrijeme promatranja rada sustava
- ◆ **A** – Broj dolazaka zadataka u vremenu  $T$
- ◆ **B** – Vrijeme kroz koje je poslužitelj zaposlen u vremenu  $T$
- ◆ **C** – Broj odlazaka u vremenu  $T$

◆ **Učestalost dolazaka zadatka (  $L$  [zad/s] )**

◆  $L = A / T$

◆ **Propusnost sustava (  $X$  [zad / s] )**

◆  $X = C / T$

◆ **Srednje vrijeme posluživanja (  $S$  [s / zad] )**

◆  $S = B / C$

◆ **Srednja zaposlenost poslužitelja (  $U$  [] )**

◆  $U = B / T$

◆  $U = (B / T) * (C / C) = (B / C) * (C / T) = S * X$

# Primjer 1: Posluživanje zahtjeva na disku



◆ Disk za trajno spremanje podataka ispunjava **50 zahtjeva u sekundi**. Srednje vrijeme obrade zahtjeva operacija pisanja i čitanja je **10 ms**.

◆ Kolika je prosječna zaposlenost diska?

## ◆ Rješenje

◆ Propusnost sustava  $X = 50 \text{ z/s}$

◆ Srednje vrijeme obrade zahtjeva  $S = 10 \text{ ms/z}$

◆ Prosječna zaposlenost diska  $U$

$$U = X * S = 50 \text{ z/s} * 0.01 \text{ s/z} = 0.5 ( 50 \% )$$



**pr1.c**



- ◆ **Broj zahtjeva u repu jednak je umnošku učestalosti dolazaka zahtjeva u rep i prosječnog vremena zadržavanja zahtjeva u sustavu**

- ◆  $L [z/s]$  – Učestalost dolazaka zahtjeva u rep zahtjeva
- ◆  $R [s]$  – Prosječno vrijeme zadržavanja zahtjeva u sustavu
- ◆  $Q [z]$  – Broj zahtjeva u repu

$$Q = L * R$$

- ◆ **Ako je sustav stabilan**

- ◆ Broj prispjelih zahtjeva u vremenu jednak je broju zahtjeva koji napuštaju sustav (  $L = X$  )
- ◆ Napomena: sve veličine su srednje vrijednosti!

$$Q = X * R$$

## Primjer 2: Čekanje na posluživanje zahtjeva s diska



- ◆ Disk iz prethodnog slučaja ima prosječno **1 zahtjev u repu**

- ◆ Koliko je prosječno vrijeme čekanja na obradu zahtjeva ?

### ◆ Rješenje

- ◆ Ulazni ritam zahtjeva  $L = 50 \text{ z/s}$

- ◆ Broj zahtjeva u repu  $Q = 1 \text{ z}$

- 
- ◆ Vrijeme zadržavanja zahtjeva u sustavu  $R$

$$R = Q/L = ( 1 \text{ z} ) / ( 50 \text{ z/s} ) = 20 \text{ ms}$$

- ◆ Vrijeme zadržavanja uključuje vrijeme čekanja u repu ( $W$ ) i vrijeme obrade zahtjeva ( $S$ ):  $R = W + S$

- ◆ Vrijeme čekanja na obradu  $W$

$$W = R - S = 20 \text{ ms} - 10 \text{ ms} = 10 \text{ ms}$$



pr2.c

- ◆ **Ukupno vrijeme zadržavanja zahtjeva u sustavu (R)**
  - ◆ Vrijeme obrade svih Q zahtjeva u repu ispred novog zahtjeva uvećano za vrijeme obrade novog zahtjeva
  - ◆  $R = S + W = S + S \cdot Q$
- ◆ **U stabilnom stanju, primjenom Littleovog  $Q = X \cdot R$** 
  - ◆  $R = S + S \cdot X \cdot R$ ,  $R(1 - X \cdot S) = S$ ,  $R = S / (1 - X \cdot S)$
- ◆ **Primjenom supstitucije  $X \cdot S = U$** 
  - ◆  $R = S / (1 - U)$
- ◆ **Množenje obje strane sa X ->  $[X \cdot R = (X \cdot S) / (1 - U)]$** 
  - ◆  $Q = U / (1 - U)$
- ◆ **Množenje obje strane sa S ->  $[S \cdot Q = (S \cdot U) / (1 - U)]$** 
  - ◆  $W = (S \cdot U) / (1 - U)$

## Primjer 3: Komunikacijski kanal



- ◆ Mjerenjem na pristupnoj točki mreže dobivamo **srednji protok od 125 paketa u sekundi** i **srednje vrijeme posluživanja 0.002 sekunde**.

- ◆ Što je sve moguće zaključiti o promatranom kanalu ?

- ◆ **Rješenje**

- ◆ Srednji protok  $X = 125 \text{ p/s}$

- ◆ Srednje vrijeme posluživanja  $S = 0.002 \text{ s/p}$

.....

- ◆ Prosječna zaposlenost komunikacijskog sustava  $U$

$$U = X * S = ( 125 \text{ p/s} ) * ( 0.002 \text{ s/p} ) = 0.25 (25 \%)$$

- ◆ Srednje vrijeme zadržavanja paketa u sustavu ( $R$ )

$$R = S / (1 - U) = (0.002 \text{ s/p}) / (1 - 0.25) = 0.0026666 \text{ s}$$

- ◆ Srednji broj paketa u repu ( $Q$ )

$$Q = X * R = (125 \text{ p/s}) * (0.0026 \text{ s}) = 0.333 \text{ p}$$

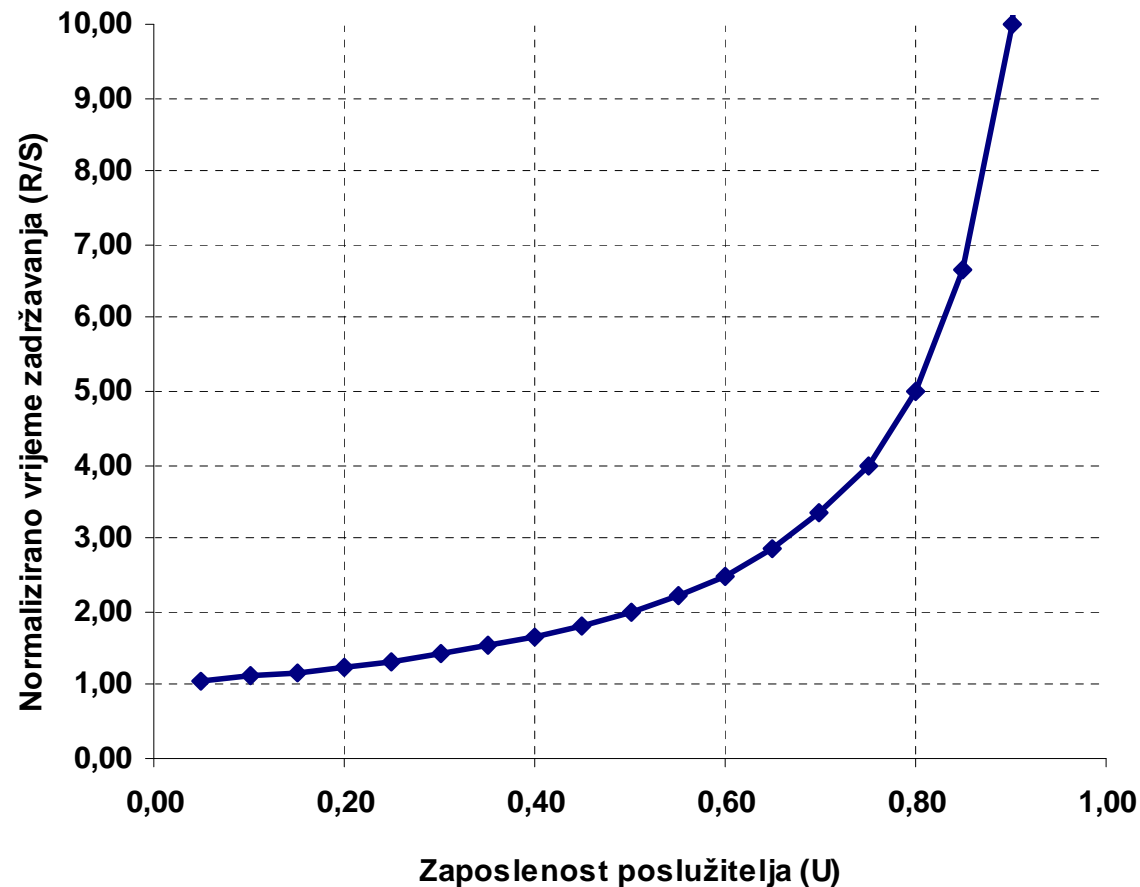


**pr3.c**

# Odziv sustava s repovima je izrazito nelinearan



- ◆ Graf normilzirano cekanje (R/S) kao funkcija opterećenja poslužitelja U



$$R/S = 1/(1-U)$$



**Odziv.xls**

## Primjer 4: Vrijeme čekanja i broj zahtjeva



- ◆ Sustav ima **prosječno vrijeme posluživanja 1 sekunda** i **učestalost dolazaka zahtjeva je 0.5 zadatka u sekundi**.

- ◆ Kolika je srednja vrijednost ukupnog vremena čekanja (**R**) i srednja vrijednost broja zahtjeva u repu (**Q**)?

### ◆ Rješenje

- ◆ Prosječno vrijeme posluživanja  **$S = 1 \text{ s/z}$**
- ◆ Učestalost pristiglih zahtjeva  **$L = 0.5 \text{ z/s}$**

- 
- ◆ Prosječna zaposlenost sustava  $U$

$$U = S * L = ( 1 \text{ s/z} ) * ( 0.5 \text{ z/s} ) = 0.5 (50 \%)$$

- ◆ Srednje vrijeme zadržavanja paketa u sustavu ( $R$ )

$$R = S / (1 - U) = ( 1 ) / ( 1 - 0.5 ) = 2 \text{ s}$$

- ◆ Srednja vrijednost broja zahtjeva u sustavu ( $Q$ )

$$Q = U / (1 - U) = 0.5 / (1 - 0.5) = 1 \text{ z}$$



**pr4.c**

## ◆ Little-ov zakon

◆  $Q = L * (R1 + R2 + R3)$

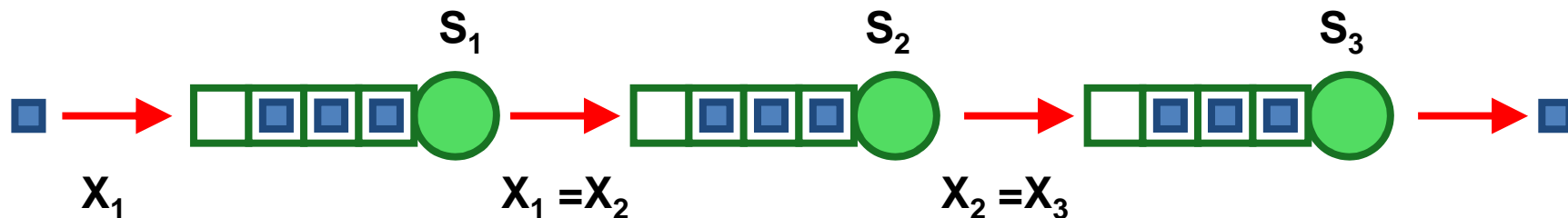
## ◆ U stabilnom stanju sustava ( $L = X$ i $X1 = X2 = X3 = X$ )

◆  $Q = X * R = X * (R1 + R2 + R3)$

## ◆ Uz $R_N = S_N / (1 - X * S_N)$

◆  $Q = X * (S_1 / (1 - X * S_1)) + (S_2 / (1 - X * S_2)) + (S_3 / (1 - X * S_3))$

◆  $R = R1 + R2 + R3$



## Primjer 5: Posluživanje u seriji



◆ Sustav sadrži 3 serijske procesne jedinice s prosječnim vremenima posluživanja 1 s, 2 s i 3 s.

◆ Koliko će biti vrijeme zadržavanja u sustavu uz ulazni ritam zahtjeva od 0.1 z/s ?

◆ Koliki će biti prosječni broj zahtjeva u sustavu ?

### ◆ Rješenje

◆ Prosječna vremena posluživanja  $S_1 = 1 \text{ s/z}$ ,  $S_2 = 2 \text{ s/z}$ ,  $S_3 = 3 \text{ s/z}$

◆ Propusnost sustava  $X = 0.1 \text{ z/s}$

◆ Vremena zadržavanja  $R_N = S_N / (1 - X \cdot S_N)$

$$R_1 = 1.11\text{s}, R_2 = 2.5\text{s}, R_3 = 4.29\text{s}$$

◆ Prosječni broj zahtjeva u repu Q

$$Q = X \cdot (R_1 + R_2 + R_3) = 0.1 \text{ z/s} \cdot (1.11 + 2.5 + 4.29) = 0.79 \text{ z}$$

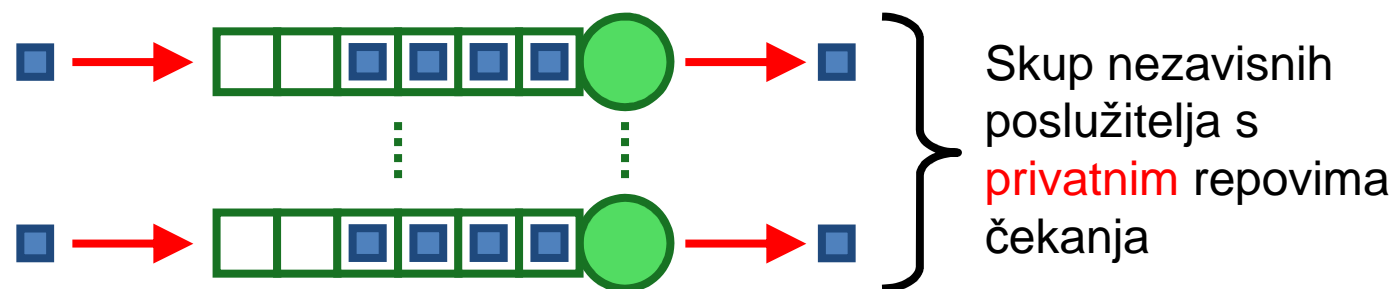


pr5.c



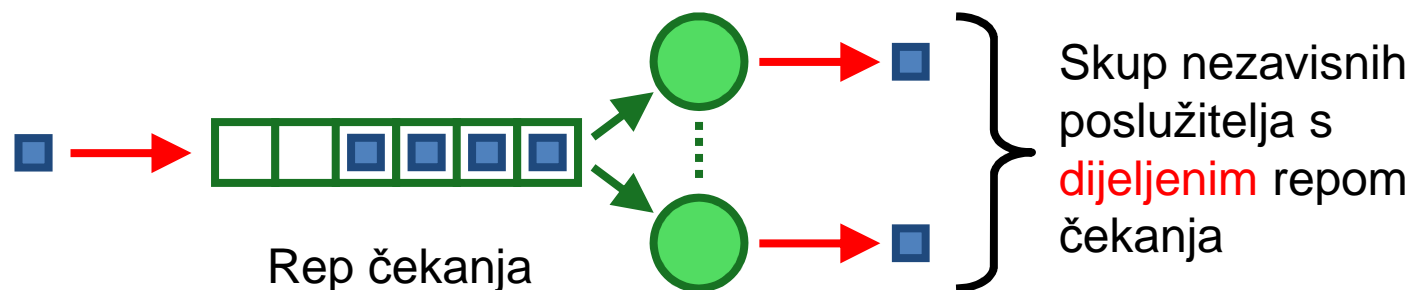
## ◆ Multi-računalo

- ◆ Model koji se primjenjuje u supermarketima



## ◆ Multi-procesor

- ◆ Model koji se primjenjuje u bankama



# Sustav više paralelnih repova čekanja

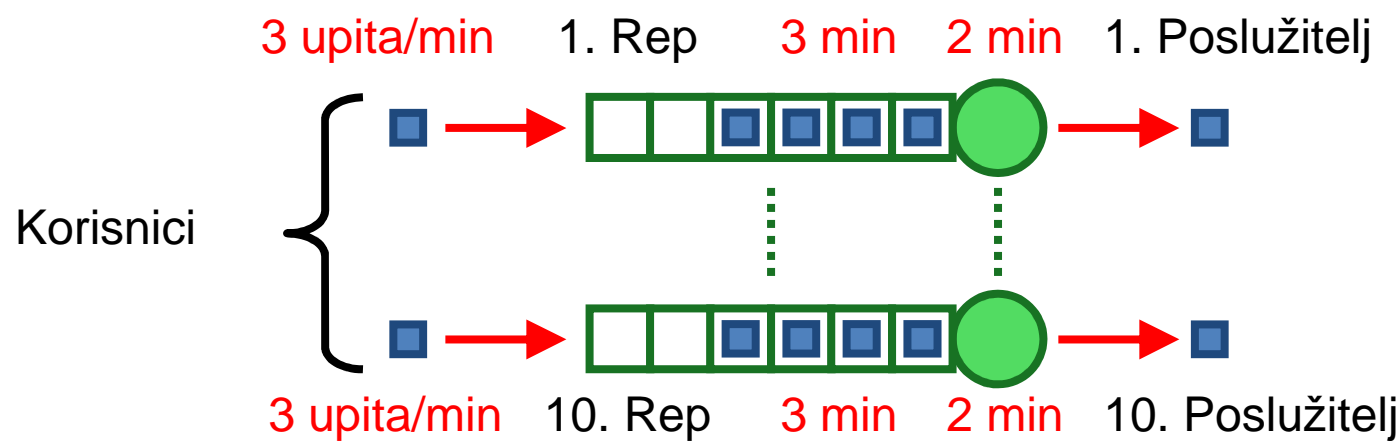


- ◆ **Ukupno vrijeme čekanja ( $R$ ) za dva repa**
  - ◆  $R = S + (S * (0.5 * Q))$
- ◆ **Primjenom Little-ovog zakona  $Q=X*R$** 
  - ◆  $R = S / (1 - (0.5 * X * S))$
- ◆ **Primjenom supstitucije  $X*S = U$** 
  - ◆  $R = S / (1 - 0.5 * U)$
- ◆ **Ukupna zaposlenost  $U$  sustava podijeljena s brojem repova  $N$  je faktor iskorištenja  $ro$  koji predstavlja vjerojatnost da je poslužitelj zaposlen**
  - ◆  $ro = U/N, R = S / (1 - ro)$
  - ◆ Za beskonačno mnogo repova  
 $N \rightarrow \infty ; ro \rightarrow 0 ; R \rightarrow S$  (Nema čekanja na posluživanje)

## Primjer 6: Aplikacija korisničke podrške



- ◆ Web aplikacija uključuje podršku korisnicima putem *chat* usluge. Kupci sami odabiru jedan od **10 repova čekanja**. Mjerenja pokazuju da zahtjevi **prosječno dolaze 3 upita u minuti** te da svaki kupac **prosječno čeka 3 minute u repu** i **prosječno provodi 2 minute u konverzaciji**.
- ◆ Koliko bi dodatnih tehničara trebalo zaposliti da se prosječno vrijeme čekanja svede na 1 minutu ?



pr6.c

## Primjer 6: Aplikacija korisničke podrške



### ◆ Rješenje

- ◆ Prosječno vrijeme posluživanja  $S = 2 \text{ min/z}$
- ◆ Broj pristiglih zahtjeva u jednom repu  $L = 3 \text{ z/min}$

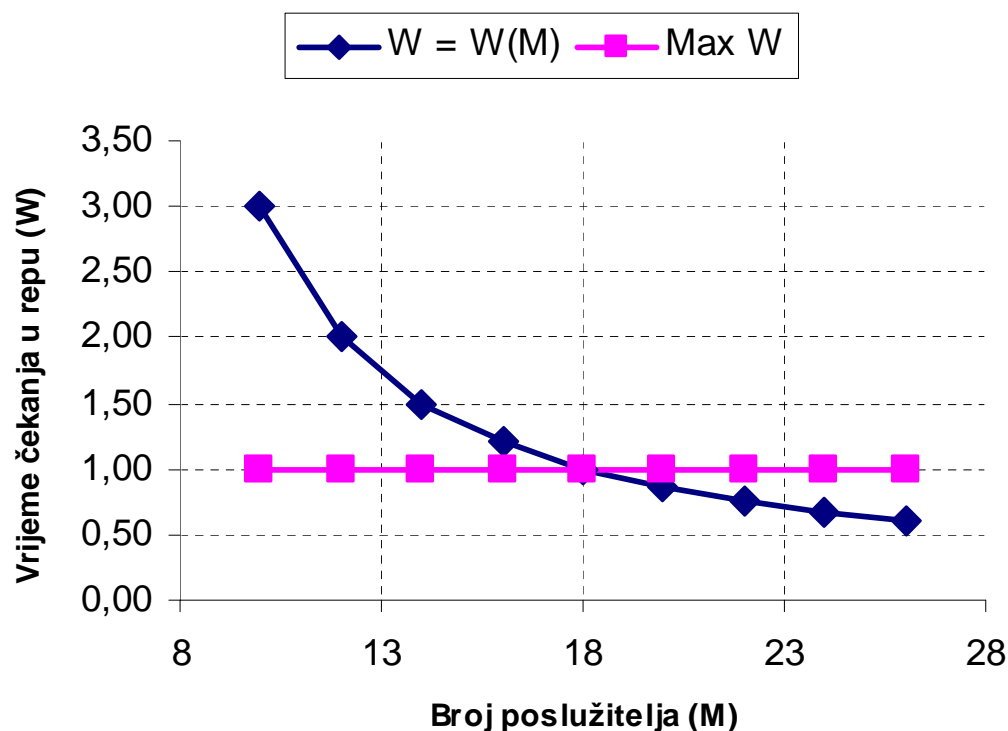
- 
- ◆ Prosječna zaposlenost sustava (U)  
$$U = S L = ( 2 \text{ min/z} ) ( 3 \text{ z/min} ) = 6$$
  - ◆ Faktor iskorištenja ( $\rho$ )  
$$\rho = U/N = 6/10 = 0.6$$
  - ◆ Srednje vrijeme zadržavanja korisnika u sustavu (R)  
$$R = S / (1 - \rho) = 2 / (1 - 0.6) = 5 \text{ min}$$
  - ◆ Srednje vrijeme čekanja u repu (W)  
$$W = R - S = 5 \text{ min} - 2 \text{ min} = 3 \text{ min}$$

## Primjer 6: Aplikacija korisničke podrške



### ◆ Rješenje za broj tehničara

- ◆ Za zadani sustav ne postoji analitičko rješenje. Rješenje se određuje primjenom numeričkih metoda ili primjenom metode pokušaja i promašaja.



- ◆ Kao rješenje dobije se da je potrebno 18 tehničara

[pr5.xls](#)

## Primjer 6: Aplikacija korisničke podrške



### ◆ Rješenje

- ◆ Broj poslužitelja (tehničara)  $N = 18$
- ◆ Prosječno vrijeme posluživanja  $S = 2 \text{ min/z}$
- ◆ Propusnost sustava  $X = 3 \text{ z/min}$

- 
- ◆ Prosječna zaposlenost sustava  $U$

$$U = X * S = ( 3 \text{ z/min} ) * ( 2 \text{ min/z} ) = 6$$

- ◆ Faktor iskorištenja  $ro$

$$ro = U/N = 6/18 = 0.33$$

- ◆ Srednje vrijeme zadržavanja korisnika u sustavu ( $R$ )

$$R = S / (1 - ro) = 2 / (1 - 0.33) = 2.985 \text{ min}$$

- ◆ Srednje vrijeme čekanja u repu ( $W$ )

$$W = 2.985 - 2 = 0.985 \text{ min}$$

## Primjer 6: Vrijeme čekanja i broj zahtjeva

---



### ◆ Zadatci za vježbu

- ◆ **Zadatak 1:** Kakvi će biti odzivi sa 10 i 18 tehničara ako publiciranje Web stranice sa odgovorima na najčešća pitanja smanji broj upita na 2 u minuti?
- ◆ **Zadatak 2:** Kakve će rezultate dati smanjenje razgovora na 1.5 minutu?

## ◆ Efektivno vrijeme posluživanja ovisi o dva čimbenika

- ◆ Dodatni poslužitelj smanjuje vrijeme posluživanja za faktor 0.5
- ◆ Vrijeme posluživanja se množi sa vjerojatnošću da je poslužitelj zaposlen  $ro = U/2$

## ◆ Zbog navedenog vrijedi

- ◆  $S(ro) = (0.5 * S) * ro$
- ◆  $R = S + Q * S(ro) = S + 0.5 * S * ro * Q$

## ◆ Primjenom supstitucije $Q = X * R$

- ◆  $R = S + (0.5 * S * ro * X * R)$

## ◆ Primjenom supstitucije $0.5 * S * X = 0.5 * U = ro$

- ◆  $R = S + R * ro$
- ◆  $R = S / (1 - ro)$  -> množenjem sa X i supstitucijom  $S * X = 2 * ro$
- ◆  $Q = 2 * ro / (1 - ro)$



# Poopćenje na sustav s N paralelnih poslužitelja



## ◆ Za sustav s N paralelnih poslužitelja vrijedi

- ◆  $ro = U/N$

- ◆  $R = S + Q \cdot (S/N) \cdot (ro^{(N-1)})$  (aproksimacija)

## ◆ Primjenom supstitucija $Q=X \cdot R$ i $S=U/X$

- ◆  $R = S + (X \cdot R) \cdot (U/(X \cdot N)) \cdot (ro^{(N-1)})$

- ◆  $R = S + R \cdot (U/N) \cdot (ro^{(N-1)})$  uz  $ro = U/N$

- ◆  $R = S + R \cdot ro^N$

- ◆  $R \cdot (1 - ro^N) = S$

- ◆  $R = S / (1 - ro^N)$

- ◆  $X \cdot R = X \cdot S / (1 - ro^N)$

- ◆  $Q = U / (1 - ro^N)$

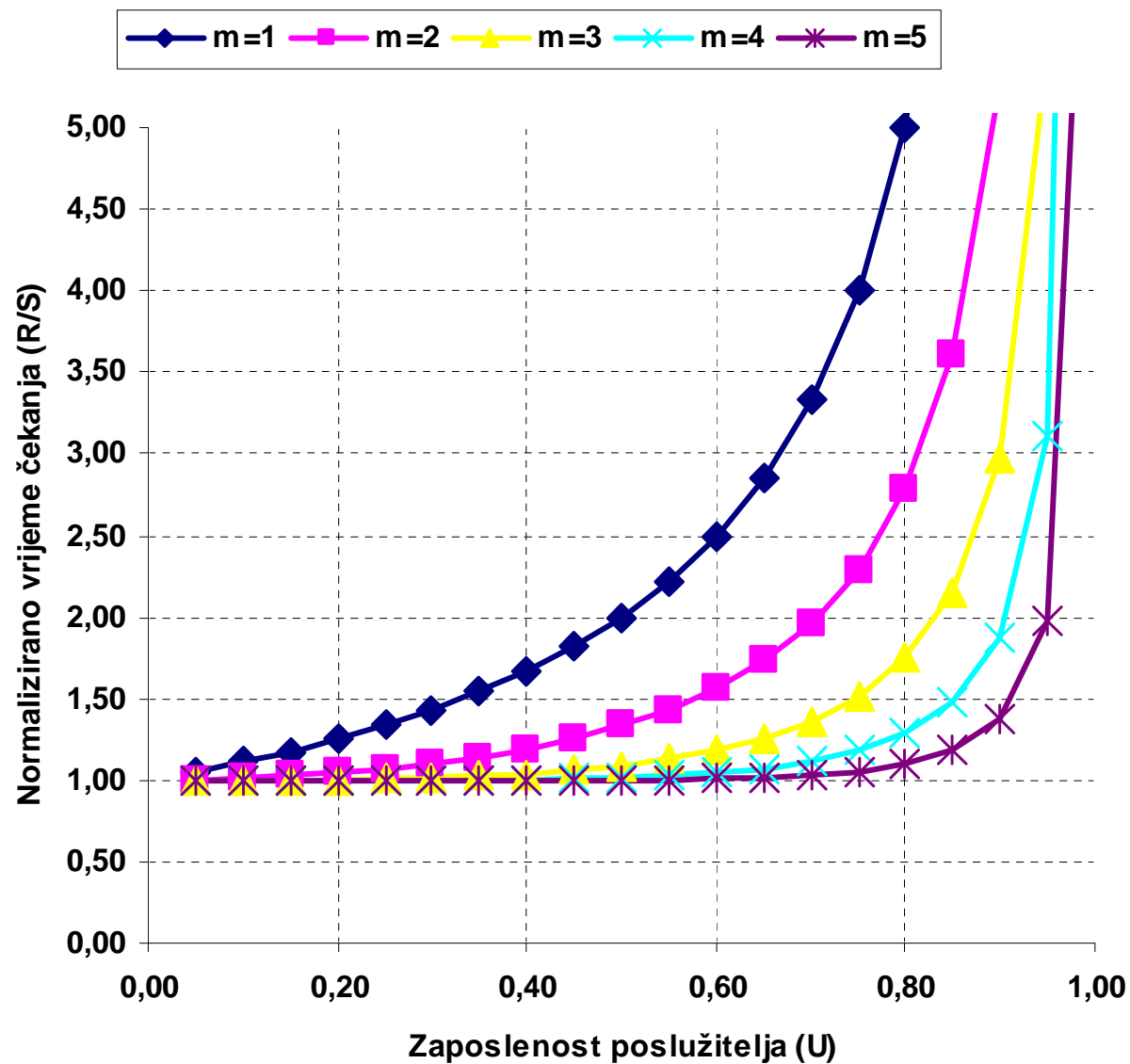
- ◆  $Q = (N \cdot ro) / (1 - ro^N)$

množenjem sa X

uz  $Q = X \cdot R$  i  $U = X \cdot S$

uz  $U = N \cdot ro$

# Usporedba sustava M/M/m



Erlang.xls

## ◆ Erlangova formula

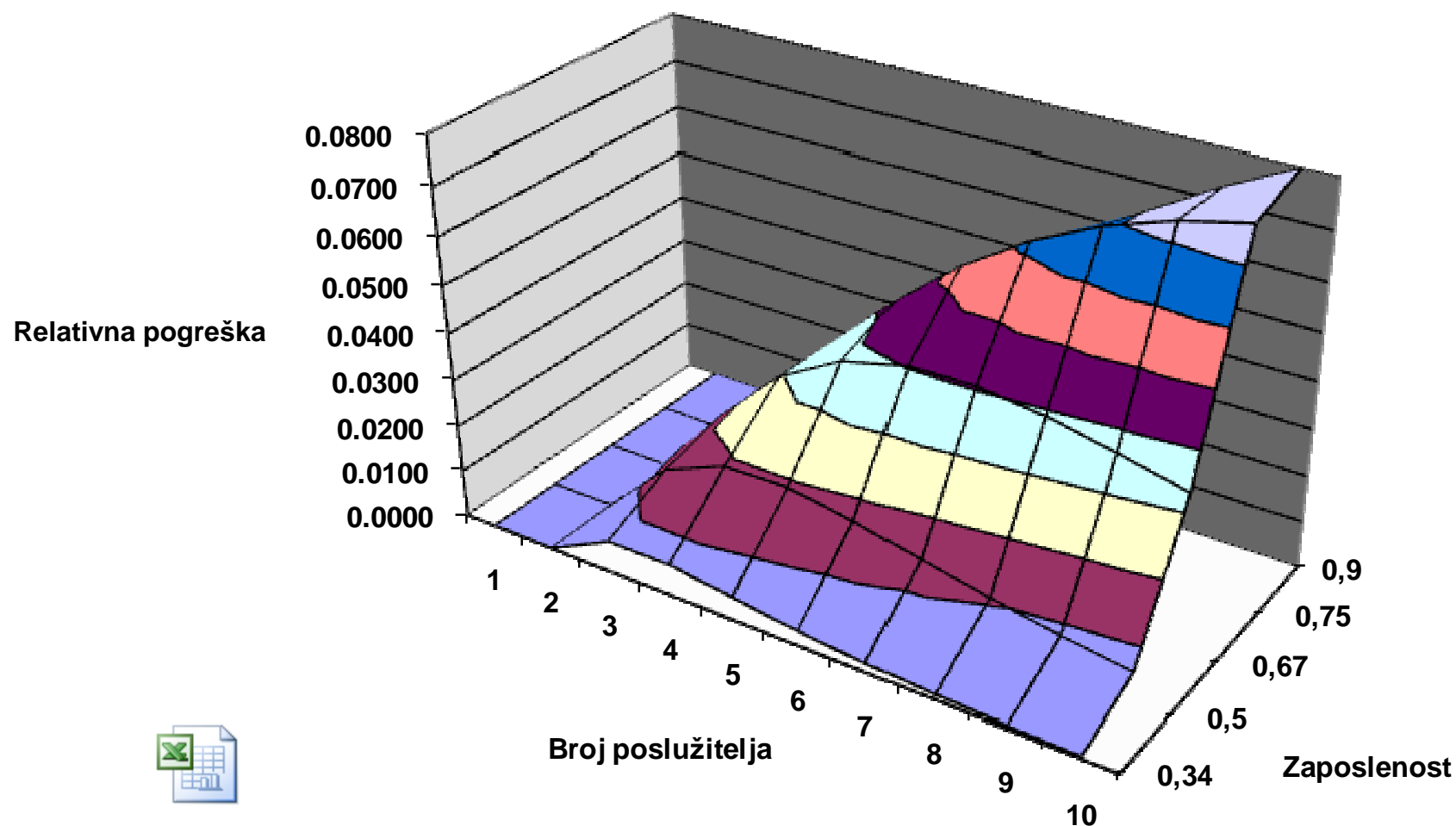
- ◆ Analitičko rješenje za vrijeme zadržavanja  $R$  u sustavu s  $N$  paralelnih poslužitelja

$$R = S * \left[ 1 + \frac{C(N, ro)}{N * (1 - ro)} \right]$$

- ◆ Koeficijent  $C(N, ro)$

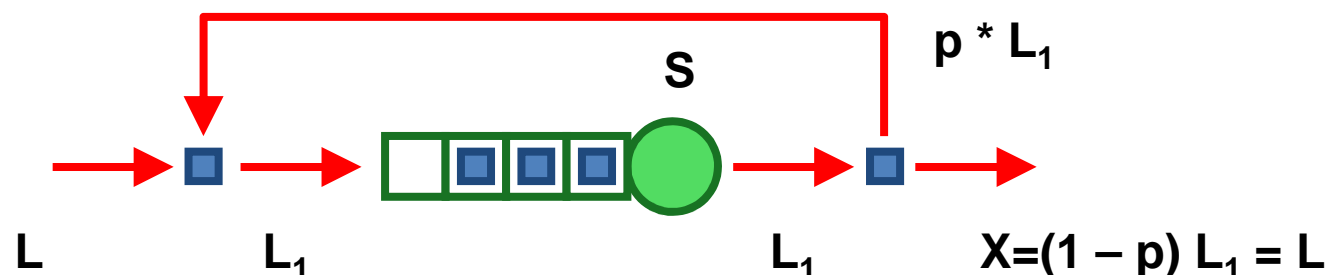
$$C(N, ro) = \frac{\frac{(N * ro)^N}{N!}}{(1 - ro) * \sum_{k=0}^{N-1} \frac{(N * ro)^k}{k!} + \frac{(N * ro)^N}{N!}}$$

## ◆ Pogreška aproksimacije



**Erlang.xls**

- ◆ Dio dolaznih zahtjeva nakon posluživanja ponovno se vraća u rep za čekanje

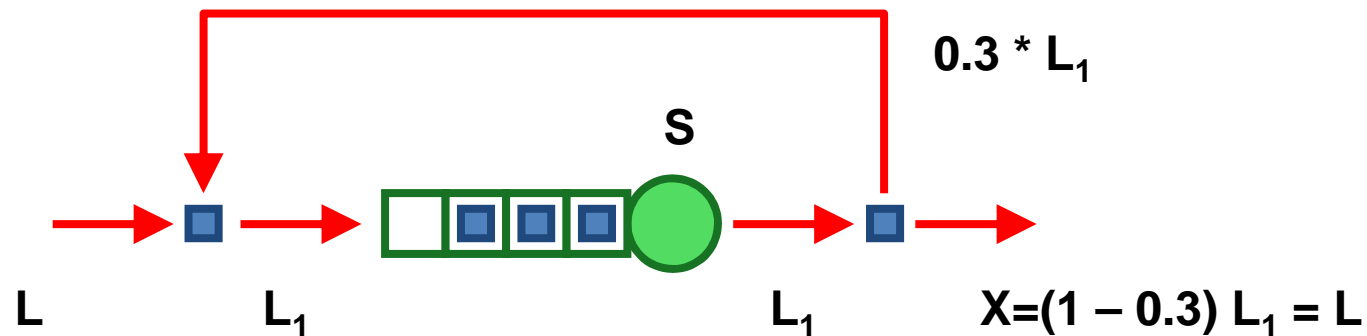


- ◆ Dva Poissonova procesa  $L_1$  i  $L_2$  rezultiraju u sa novim Poissonovim procesom  $L = L_1 + L_2$
- ◆  $L_1 = L + L_1 p = L / (1 - p)$
- ◆  $U = L_1 * S = L * S / (1 - p)$
- ◆  $R_1 = S / (1 - U)$  (zadržavanje za jedan prolaz)
- ◆ Vrijeme zadržavanja u sustavu s povratnom vezom  
 $R = R_1 * (1 + (p / (1 - p))) = R_1 / (1 - p)$

## Primjer 7: Komunikacijski kanal s pogreškom



- ◆ Paketi dolaze u komunikacijski kanal s učestalošću 0.5 paketa u sekundi i zahtijevaju 0.75 sekundi za obradu. Za 30 % paketa dogodi se pogreška pri prijenosu i takvi paketi se umeću u rep za ponovno slanje.
- ◆ Koliko vremena paket prosječno provede u kanalu ?



pr7.c

## Primjer 7: Komunikacijski kanal s pogreškom



### ◆ Rješenje

- ◆ Broj pristiglih paketa u sekundi  $L = 0.5 \text{ p/s}$
- ◆ Prosječno vrijeme obrade paketa  $S = 0.75 \text{ s/p}$
- ◆ Vjerojatnost pogreške paketa pri prijenosu  $p = 0.3$

---

◆  $L_1 = L / (1 - p) = 0.5 / 0.7 = 0.714 \text{ p/s}$

- ◆ Prosječna zaposlenost kanala  $U$

$$U = L_1 * S = 0.714 \text{ p/s} * 0.75 \text{ s/p} = 0.536 \text{ ( 53.6 \% )}$$

- ◆ Srednje vrijeme čekanja u repu  $W$

$$W = S * U / ( 1 - U ) = 0.866 \text{ s/p}$$

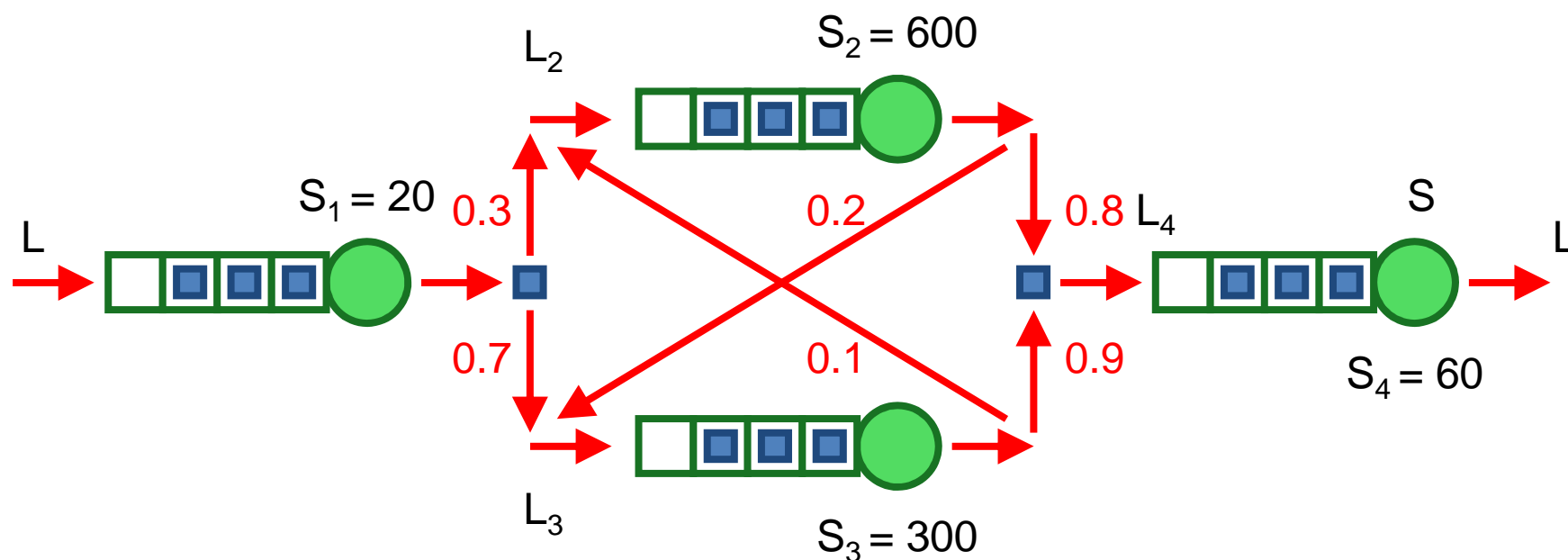
- ◆ Srednje vrijeme zadržavanja paketa u kanalu ( $R_1$ )

$$R_1 = W + S = 0.866 \text{ s/p} + 0.75 \text{ s/p} = 1.616 \text{ s/p}$$

Prosječno vrijeme u kanalu:  $R = R_1 / (1 - p) = 2.31$

## ◆ Mreža repova poruka

- ◆ Mrežna struktura proizvoljne složenosti s povratnim granama





## ◆ U stabilnom stanju sustava

◆  $U_N = X_N * S_N$  ( $X_N = L_N$  za stabilni slučaj)

◆  $U_1 = L * S_1 = 20 * L$

◆  $U_2 = 600 * L_2 = 600 (0.3 * L + 0.1 * L_3)$

◆  $U_3 = 300 * L_3 = 300 (0.7 * L + 0.2 * L_2)$

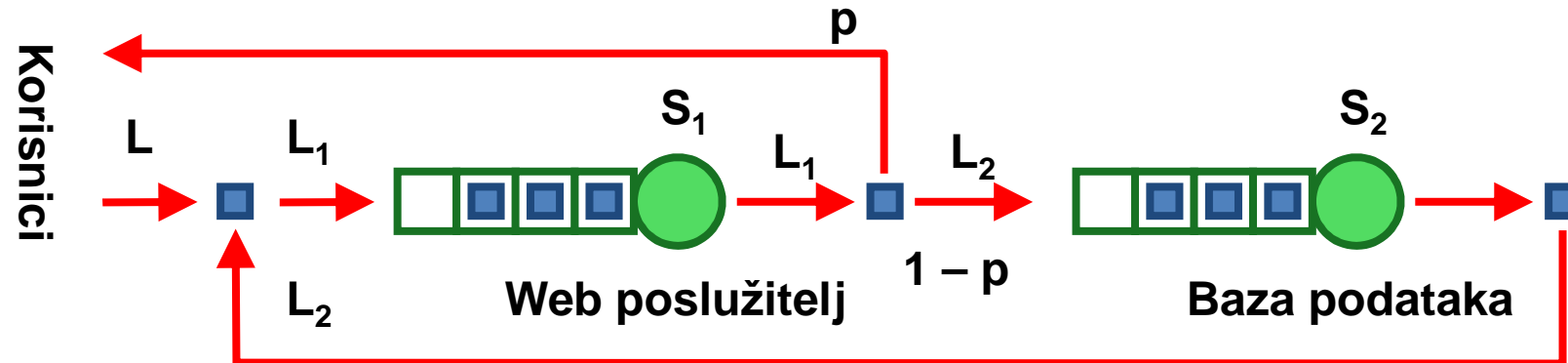
◆  $U_4 = 60 * L$

## ◆ Nakon rješenja za $L_2$ i $L_3$ i izračunavanja $U_1$ do $U_4$ , izračunavamo $Q_1$ do $Q_4$ iz:

◆  $Q_N = U_N / (1 - U_N)$

## ◆ Vrijeme zadržavanja u sustavu R

◆  $R = Q/L = (Q_1 + Q_2 + Q_3 + Q_4)/L$



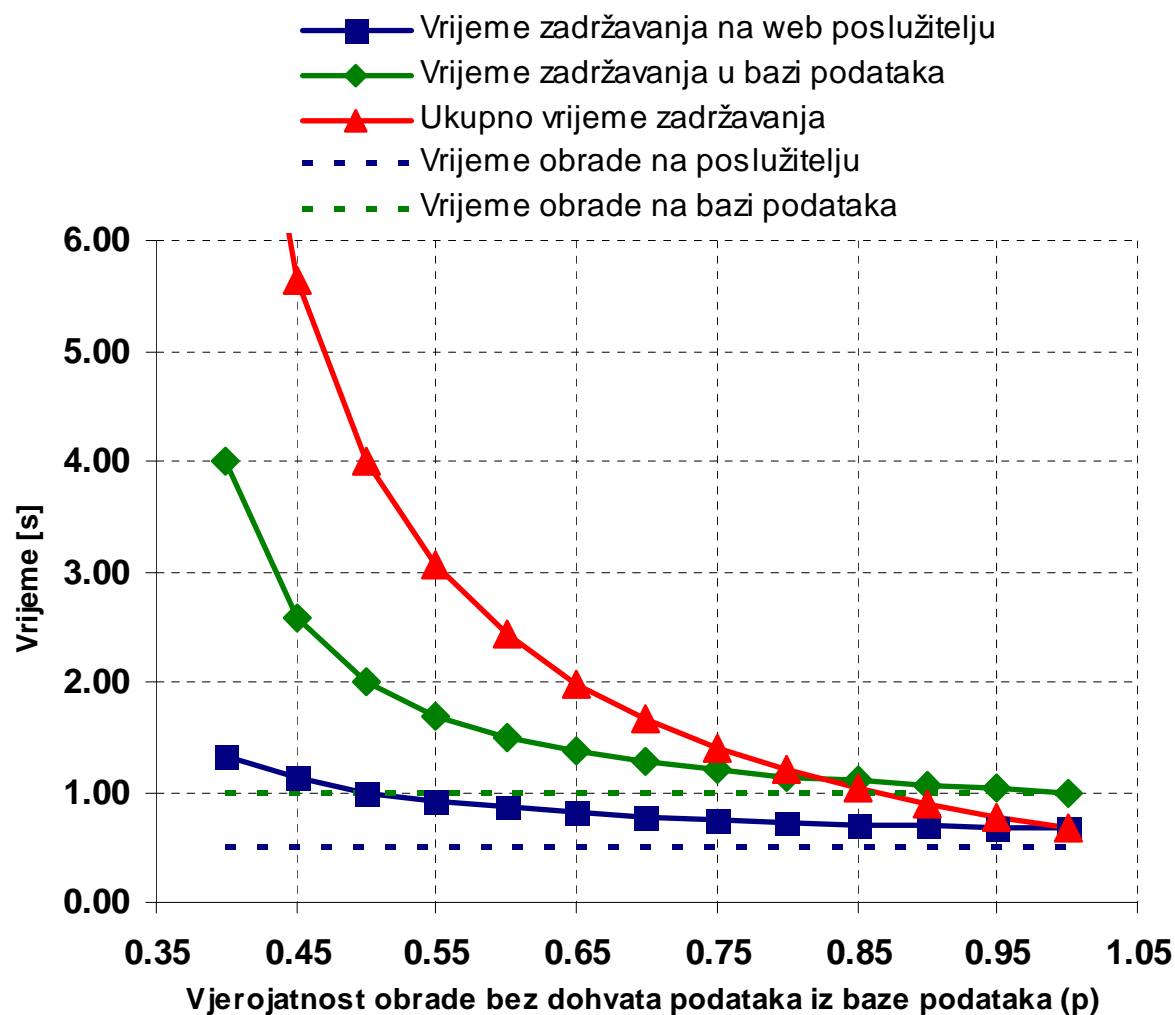
## ◆ Učestalost dolazaka zahtjeva ( $L_1$ , $L_2$ )

- ◆  $L_1 = L + L_2 = L + (1-p)L_1 = L/p$
- ◆  $L_2 = (1 - p)L_1 = ((1 - p)/p) * L$

## ◆ Vrijeme zadržavanja zahtjeva u sustavu ( $R$ )

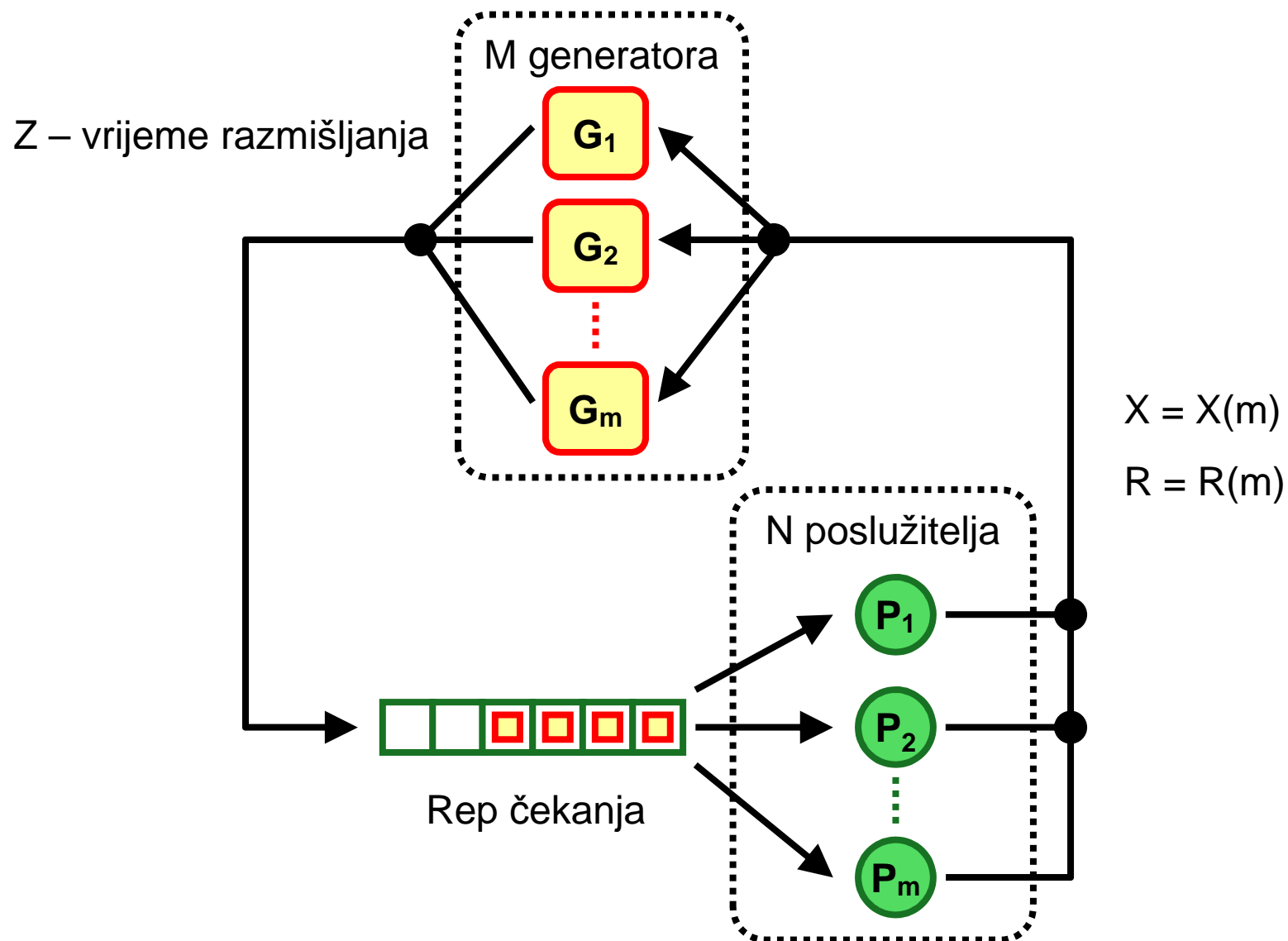
- ◆  $X_1 = L_1$ ;  $X_2 = L_2$
- ◆  $U_1 = X_1 * S_1 = L * S_1 / p$ ;  $U_2 = ((1 - p)/p) * L * S_2$
- ◆  $R_1 = S_1 / (1 - U_1)$  ;  $R_2 = S_2 / (1 - U_2)$
- ◆  $R = R_1 * (1 + (1 - p)/p) + R_2 * (1 - p)/p$

# Utjecaj parametara na ponašanje Web poslužitelja



Posluzitelj.xls

# Zatvoreni centar za posluživanje



## ◆ Propusnost sustava

- ◆ Jednaka je ritmu razmišljanja pomnoženom s brojem slobodnih mislioca (ukupan broj  $m$  minus broj u repu)
- ◆  $X(m) = (m - Q)/Z$

## ◆ Množenjem s $Z$ te primjenom supstitucije $Q = X \cdot R$

- ◆  $Z \cdot X(m) = m - X(m) \cdot R$
- ◆  $X(m) = m / (R + Z)$

## ◆ Prosječno vrijeme zadržavanja zahtjeva $R$

- ◆  $R = m / X(m) - Z$

## Primjer 8: Poslužitelj aplikacija



- ◆ Poslužitelj aplikacija omogućava skupini inženjera razvoj programa u dijeljenom vremenu. Mjerenjem su utvrđene sljedeće značajke sustava:
  - ◆ Srednji broj aktivnih razvojnih inženjera  $m = 230$
  - ◆ Srednje vrijeme između kompilacija je  $Z = 300$  s
  - ◆ Srednje iskorištenje poslužitelja je  $U = 0.48$
  - ◆ Srednje vrijeme kompilacije  $S = 0.63$  s
- ◆ Upravitelj sustava želi odrediti:
  - ◆ Propusnost sustava ( $X$ )?
  - ◆ Koliko je srednje vrijeme kompilacije ( $R$ )?



pr8.c

## Primjer 8: Poslužitelj aplikacija



### ◆ Rješenje

- ◆ Broj generatora zahtjeva  **$m=230$**
- ◆ Srednje vrijeme između kompilacija je  **$Z = 300$  s**
- ◆ Srednje iskorištenje poslužitelja je  **$U = 0.48$**
- ◆ Srednje vrijeme kompilacije  **$S = 0.63$  s/kom**

---

### ◆ Propusnost sustava (X)

$$X = U/S = 0.48 / 0.63 \text{ s} = \mathbf{0.7619 \text{ kom/s}}$$

### ◆ Srednje vrijeme zadržavanja u sustavu (R)

$$R = m / X(m) - Z$$

$$R = (230 \text{ kom} / 0.7636 \text{ kom/s}) - 300 \text{ s} = \mathbf{1.21 \text{ s}}$$

- ◆ Poznato u operacijskim istraživanjima kao "machine repair center" problem
- ◆ Kendall-ova notacija
  - ◆ M/M/N/m/m
- ◆ Egzaktno numeričko rješenje dano je izvornim kodom u jeziku C
  - ◆ repair.c



[repair.c](#)



## Primjer 8: Zatvoreni sustav s paralelnim poslužiteljima



### ◆ Proširivanje sustava iz prethodnog primjera

- ◆ Što će se dogoditi sa sustavom ako poduzeće zaposli novih 200 programera?

Odgovor:  $U = 0.88$ ;  $R = 5.009$  s

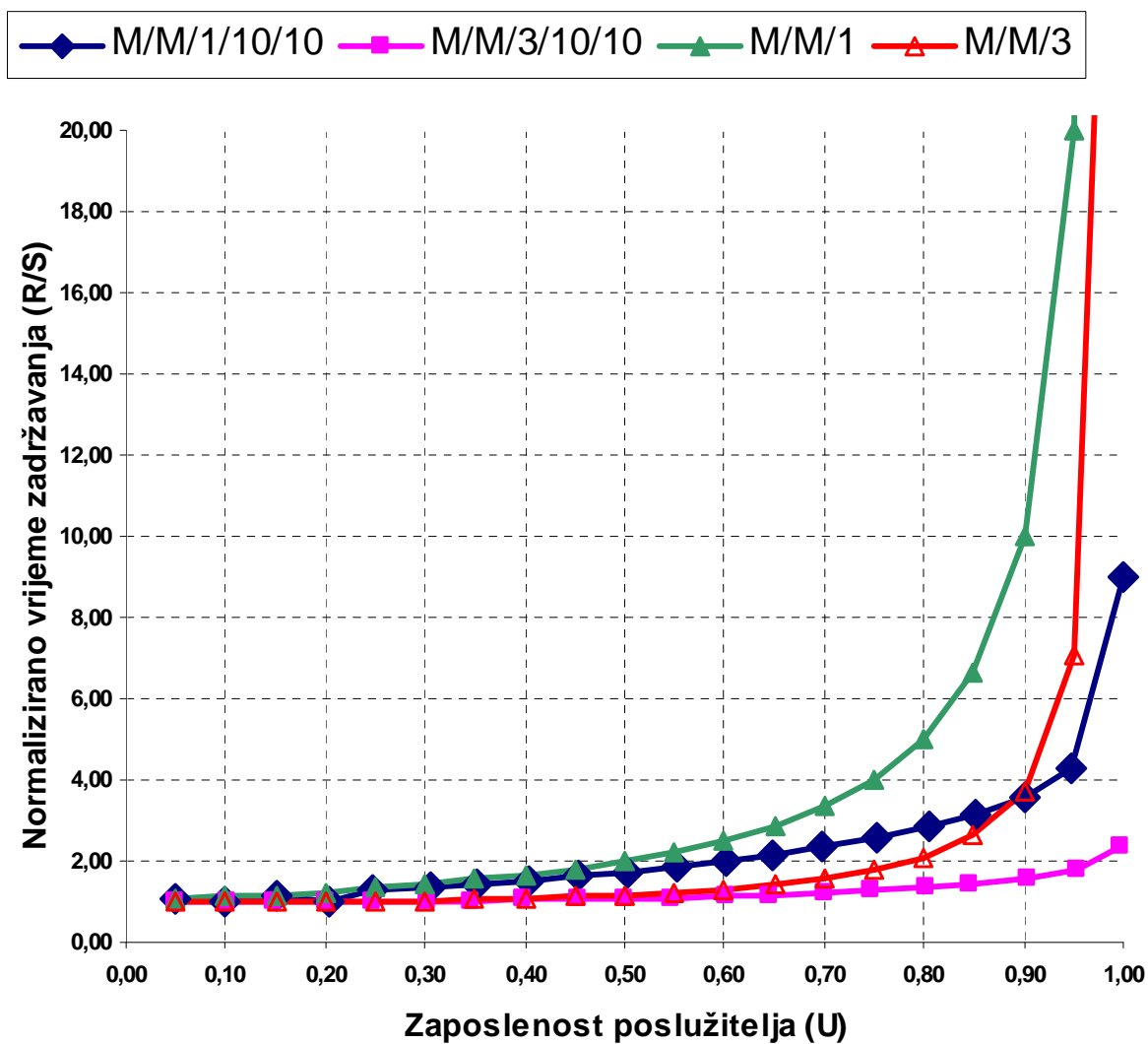
### ◆ Koliko će se situacija popraviti ako se poslužitelju doda drugi procesor sa istim značajkama?

Odgovor korištenjem *repair.c*:

- ◆ Praksa je pokazala da se u multi-procesorskim sistemima postoji dodatni teret zbog sinkronizacije procesora. Uobičajeni faktor je 3 - 5% tj. u našem slučaju uzmimo da se  $S$  rate se povećava na  $\sim 0.66s$
- ◆ Program daje slijedeće rezultate:
  - ◆  $U = 0.47$ ;  $R = 0.8465s$

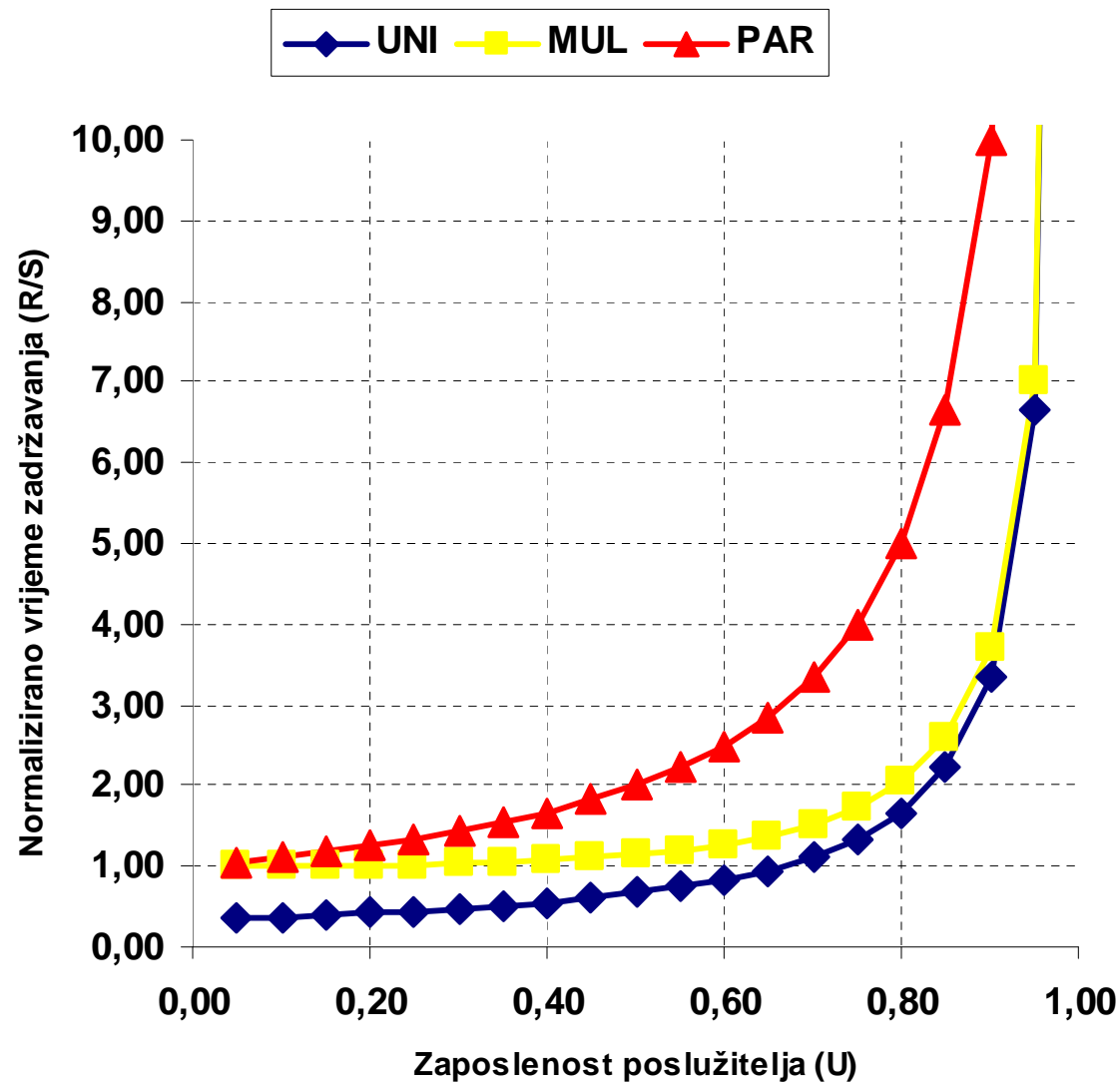
- ◆ U zatvorenom sustavu rep ne može narasti preko ukupno  $m - 1$  zahtjeva (plus jedan u izvršavanju)
  - ◆  $R = m/X - Z$
  - ◆  $R/S = m/XS - Z/S$
  - ◆  $R/S = m/U - Z/S$
  - ◆  $U = r_o * N$
  - ◆  $R/S = m/(r_o * N) - Z/S$
  
- ◆ U slučaju kada  $r_o \rightarrow 1$  i  $Z \rightarrow 0$ 
  - ◆  $R/S \rightarrow m/N$
  - ◆  $R/S$  ne raste u beskonačno jer je broj zahtjeva u repu ograničen (t.j. zatvoreni sustav)

# Usporedba otvorenih i zatvorenih sustava



Usporedba.xls

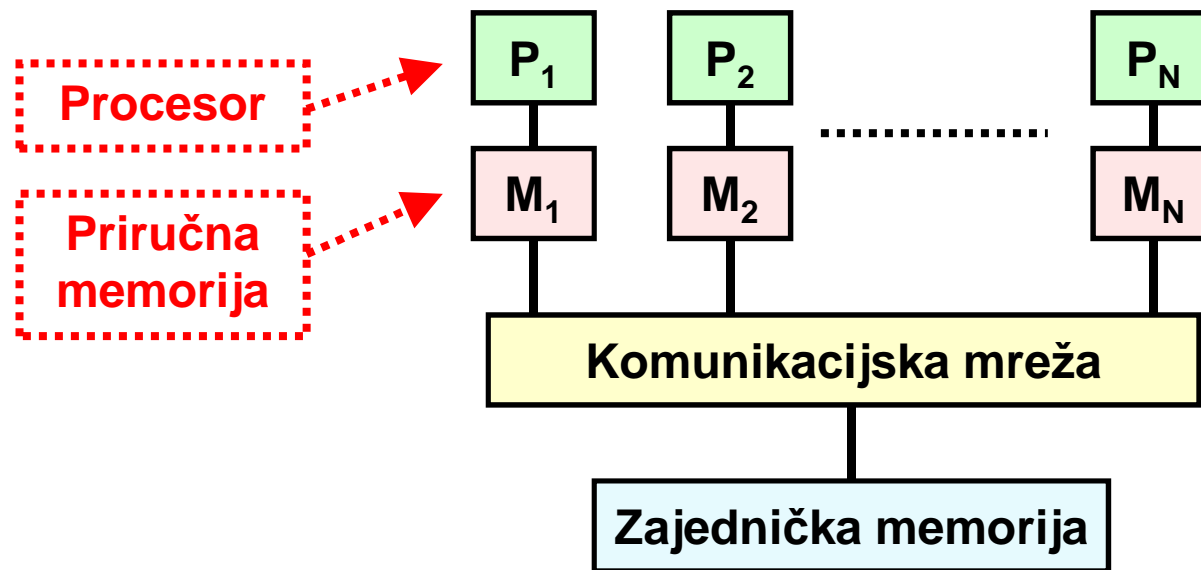
# Koji je model bolji: Multi-računalo ili Multi-procesor ?



Uniprocessor ima  
3x procesnu moć



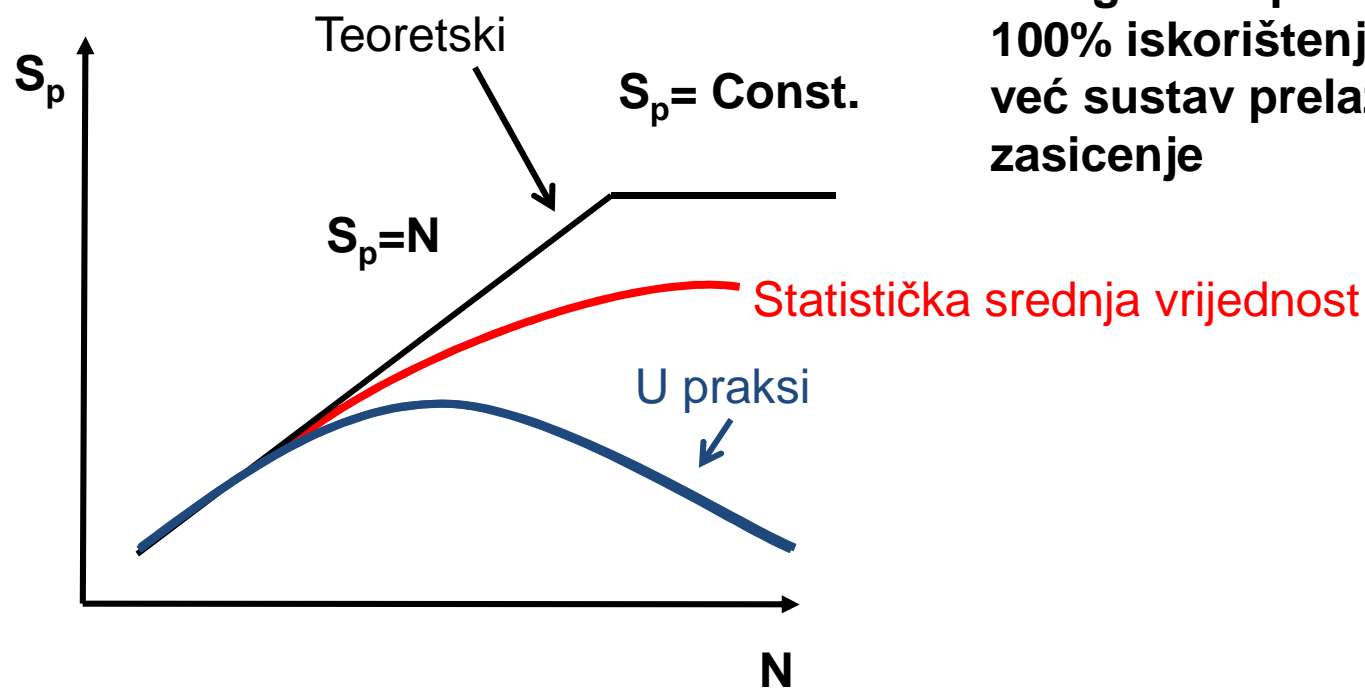
**Sustavi.xls**



## ◆ Komunikacijska mreža

- ◆ sabirnica (bus), zbir (crossbar), prsten (ring), stablo (tree), matrica (mesh), hiper-kocka (hiper-cube)

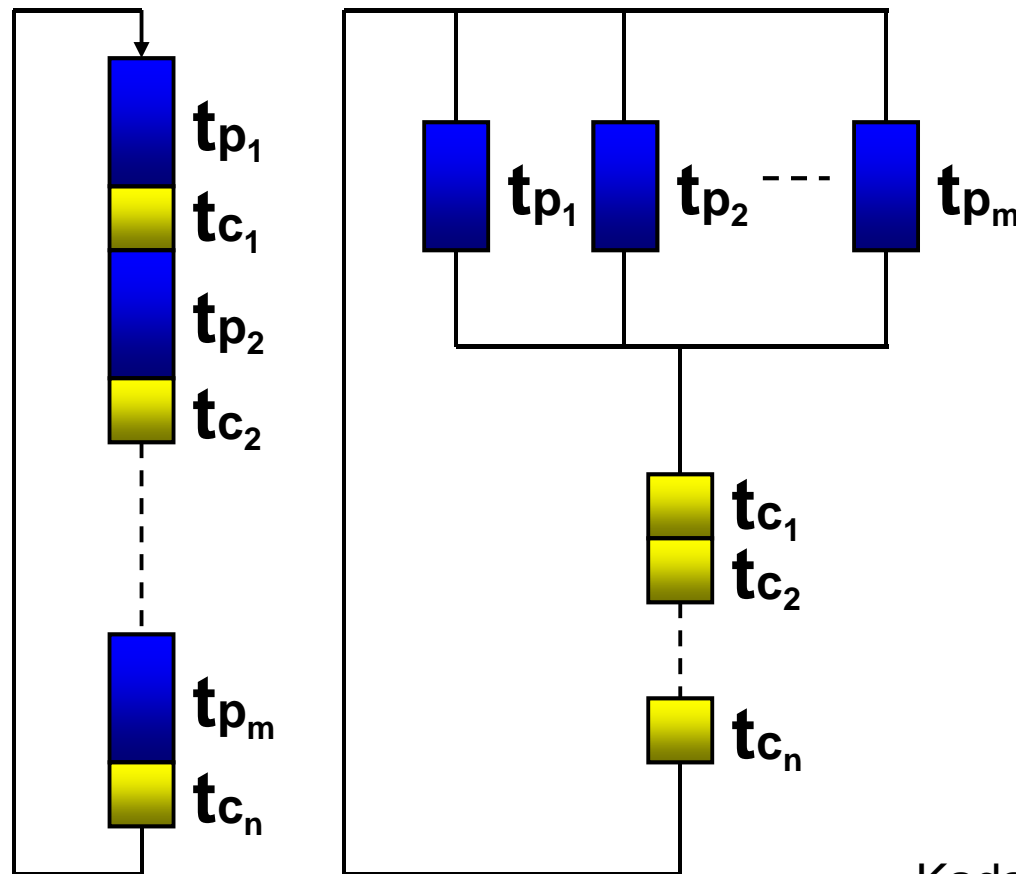
# Ubrzanje (Speedup) je nelinearna funkcija od N



Sabirnica sustava (kao i izlazna jedinica) ne omogućava punu brzinu kod 100% iskorištenja sustava, već sustav prelazi u potpuno zasícenje

- ◆ **Uzrok degradacija performansi u paralelnim sustavima**
  - ◆ Ostvarivanje sekvencijalnog pristupa zajedničkim sredstvima u aplikaciji i sklopovima
  - ◆ Kašnjenja zbog održavanja koherencije i sinkronizacije
  
- ◆ **Modeliranje paralelnih aplikacija**
  - ◆ Deterministički modeli
  - ◆ Statistički modeli

# Deterministički model



$$Tp = \sum_i tp_i$$

$$Tc = \sum_j tc_j$$

$$x = Tp / Tc$$

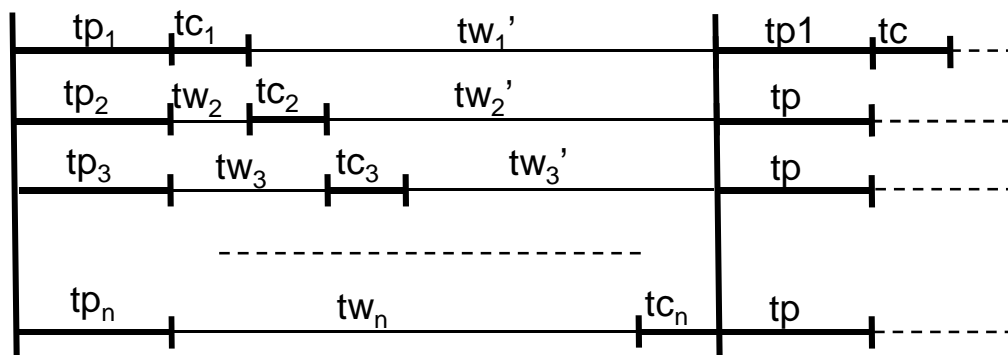
Amdahl's Law

$$S_{\max} = \frac{Tp + Tc}{Tc} = x + 1$$

Kada je  $tp = 0$ , tj. paralelno izvođenje u 0 vremena kada  $N$  teži beskonačno

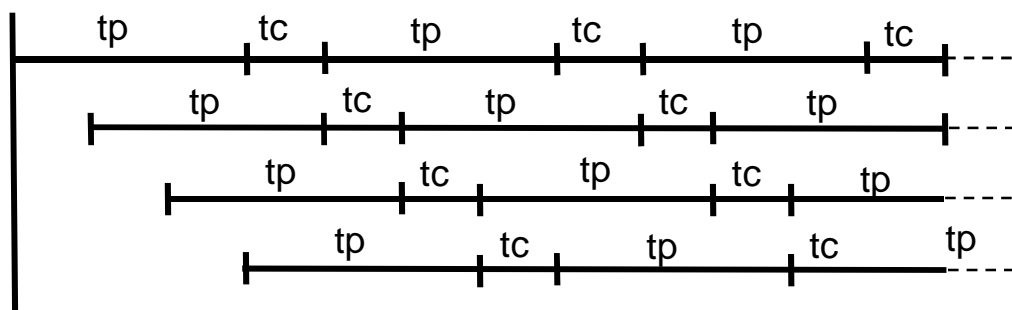


# Dva granična slučaja



**Sinkroni (Najgori slučaj)**

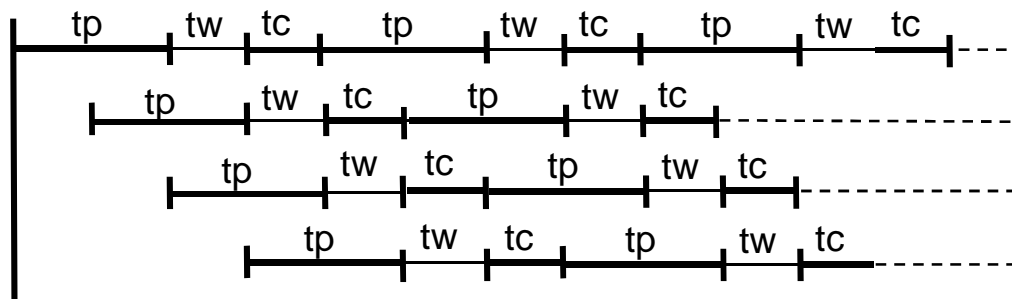
$$tw = (N - 1) * tc$$



**Asinkroni (Najbolji slučaj)**

$$tw = 0 \text{ ako } tp \geq (N - 1) * tc$$

$$tw = (N - 1) * tc - tp \text{ ako } tp < (N - 1) * tc$$



## ◆ Osnovne veličine

$$E = \frac{tp + tc}{tp + tc + tw}, \quad S = E \frac{Tp + Tc}{tp + tc}$$

$$fp(N) = \frac{Tp}{tp_N}, \quad fc(N) = \frac{Tc}{tc_N}$$

## ◆ Asinkroni slučaj

$$S = \min \left[ \frac{fpfc(x+1)}{fp + xfc}, \frac{fc(x+1)}{N} \right]$$

## ◆ Sinkroni slučaj

$$S = \frac{fpfc(x+1)}{Nfp + xfc}$$

# Rezultati za prikazane funkcije



Dekompozicija		Sinkroni slučaj			Asinkroni slučaj		
tp	tc	SP	SP <sub>max</sub>	N <sub>max</sub>	SP	SP <sub>max</sub>	N <sub>max</sub>
$N$	$N$	$\frac{(1+x)N}{N+x}$	$1+x$	$\infty$	$\min[N, 1+x]$	$1+x$	$1+x$
$N$	$\sqrt{N}$	$\frac{(1+x)N}{N^{3/2}+x}$	$\frac{2^{2/3}(1+x)}{3x^{1/3}}$	$2x^{2/3}$	$\min\left[\frac{(1+x)N}{\sqrt{N}+x}, \frac{(1+x)}{\sqrt{N}}\right]$	$\frac{(1+x)x^{1/3}}{1+x^{2/3}}$	$x^{2/3}$
$N$	$1$	$\frac{(1+x)N}{N^2+x}$	$\frac{1+x}{2\sqrt{x}}$	$\sqrt{x}$	$\min\left[\frac{N(1+x)}{N+x}, \frac{1+x}{N}\right]$	$\frac{2(1+x)}{1+\sqrt{1+4x}}$	$\frac{1+\sqrt{1+4x}}{2}$
$\log N$	$\log N$	$(1+x)\frac{\log N}{N+x}$	- *)	$N(\log N - 1) = x$ *)	$\min\left[\log N, \frac{(1+x)\log N}{N}\right]$	$\log(1+x)$	$1+x$

\*) Samo numeričko rješenje

Preuzeto iz: **Performance prediction and calibration for a class of multiprocessors**

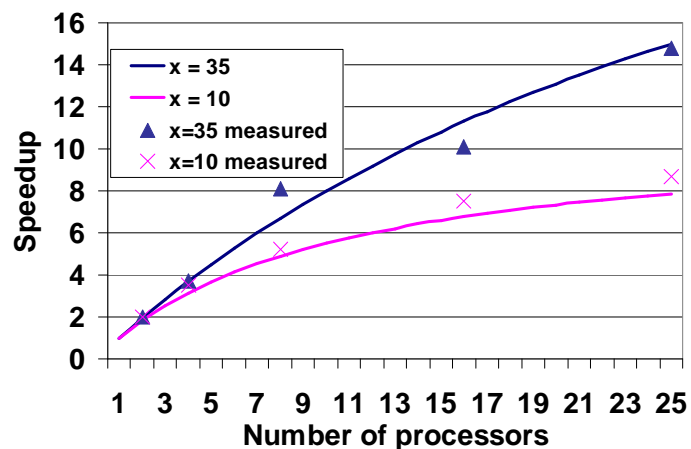
Vrsalovic, D.F.; Siewiorek, D.P.; Segall, Z.Z.; Gehring, E.F.;

IEEE Transactions on Computers, Volume: 37 Issue: 11, Nov. 1988. Page(s): 1353 -1365

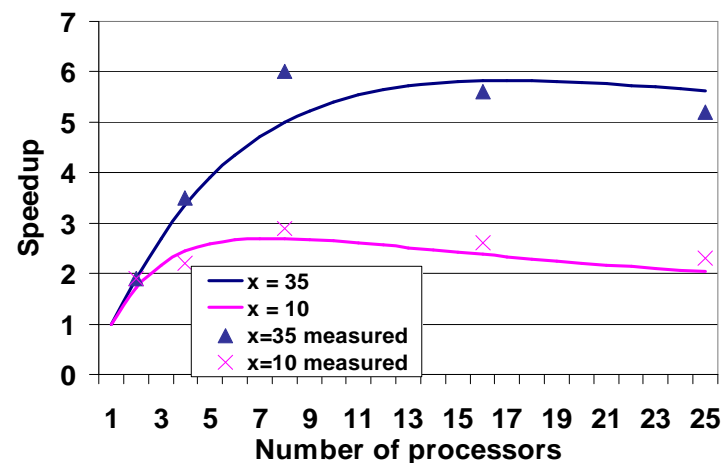
# Funkcije raspodjele opterećenja u praksi



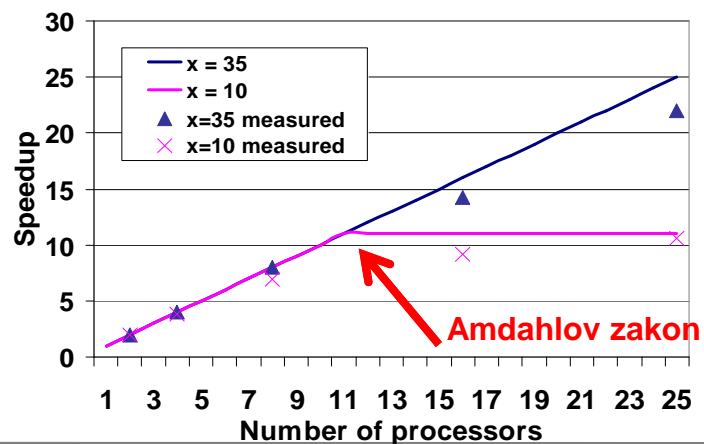
Speedup of a synchronous algorithm with an  $(N; N)$  decomposition



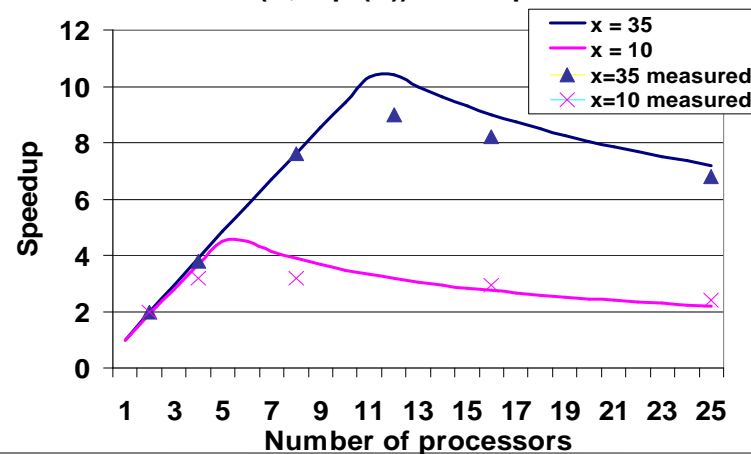
Speedup of a synchronous algorithm with an  $(N; \sqrt{n})$  decomposition



Speedup of an asynchronous algorithm with an  $(N; N)$  decomposition



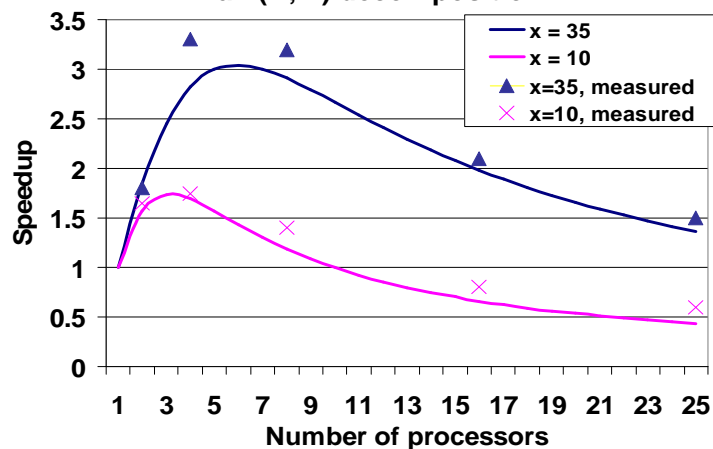
Speedup of an asynchronous algorithm with an  $(N; \sqrt{N})$  decomposition



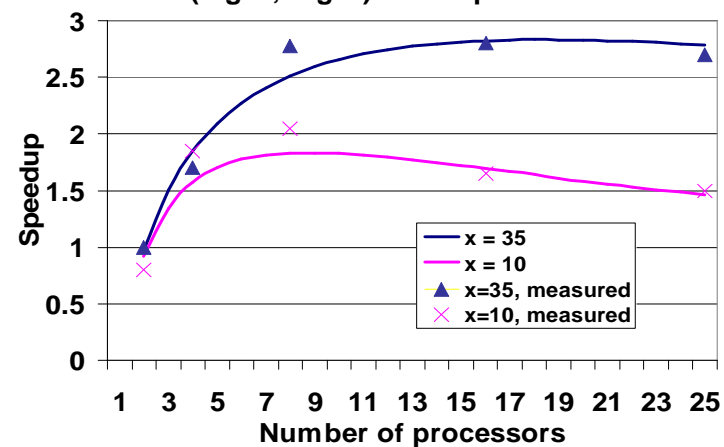
# Funkcije raspodjele opterećenja u praksi



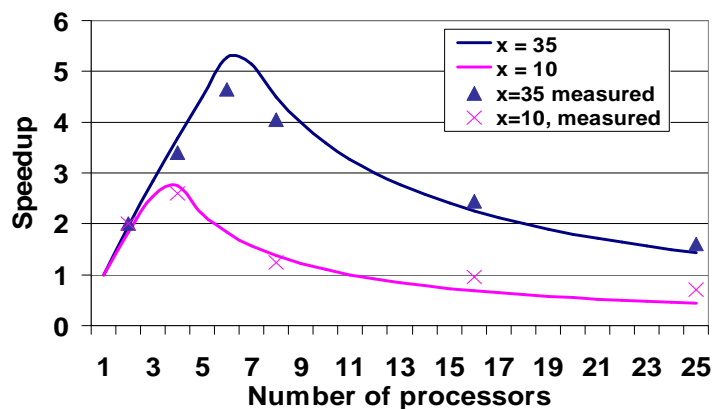
Speedup of a synchronous algorithm with an  $(N; 1)$  decomposition



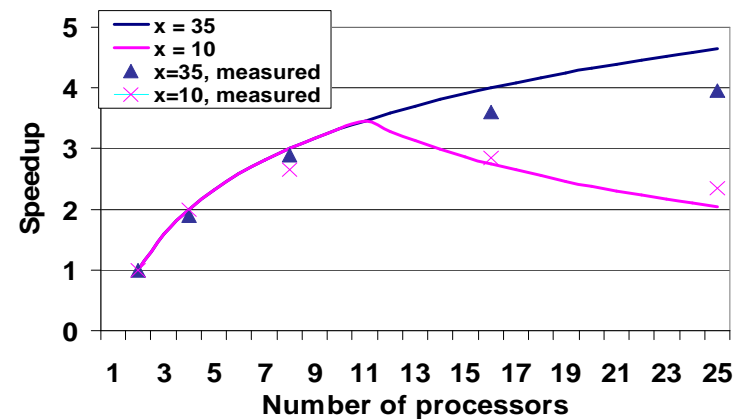
Speedup of a synchronous algorithm with a  $(\log N; \log N)$  decomposition



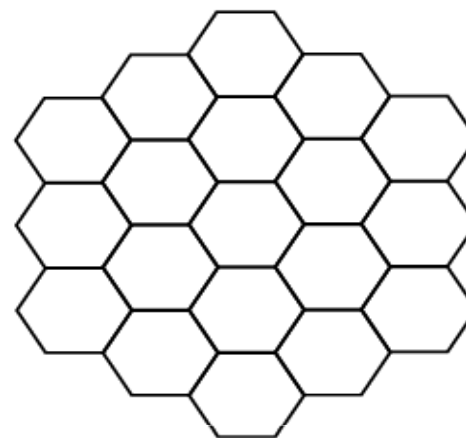
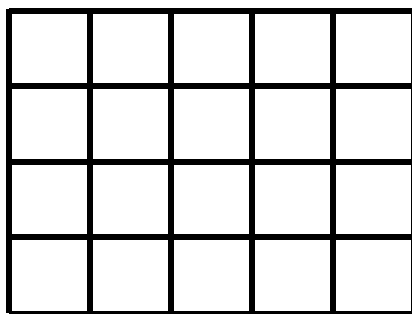
Speedup of an asynchronous algorithm with an  $(N; 1)$  decomposition



Speedup of an asynchronous algorithm with a  $(\log N; \log N)$  decomposition



## ◆ Primjer: Algoritmi najbližih susjeda



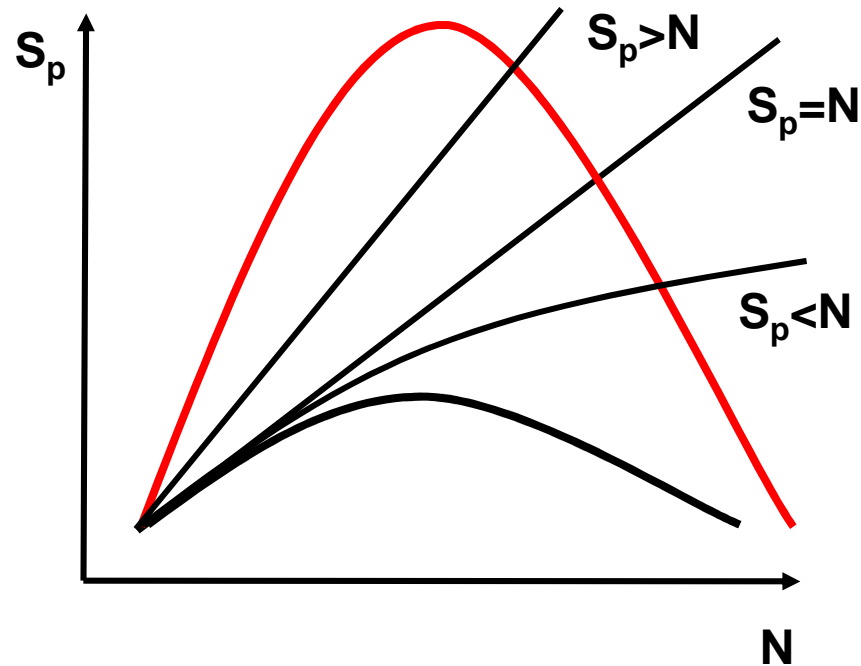
## ◆ Trend analiza ili težinska suma najbližih susjeda

$$\begin{aligned} \text{◆ } C(i,j) = & K1 \times C(i-1,j) + K2 \times C(i+1,j) + K3 \times C(i, j-1) + \\ & K4 \times C(i,j+1) - 4K5 \times C(i,j) \end{aligned}$$

◆ Za ovu klasu algoritama ubrzanje je proporcionalno omjeru površine (tp) i opsega (tc) osnovne ćelije podataka

◆ Šesterokut ima bolji omjer površine i opsega od kvadrata

# Da li je superlinearno ubrzanje moguće?



## Utjecaj primjene međuspremnika (cache)

$$tc = f(\text{veličina spremnika}) \times fc$$

## Primjena genetskih algoritama

$$P(T_i \leq T) = p$$

$$t_{i[N=2]} \leq T = p + p + p^2 > 2p$$

$$Sp = f(t_i)$$

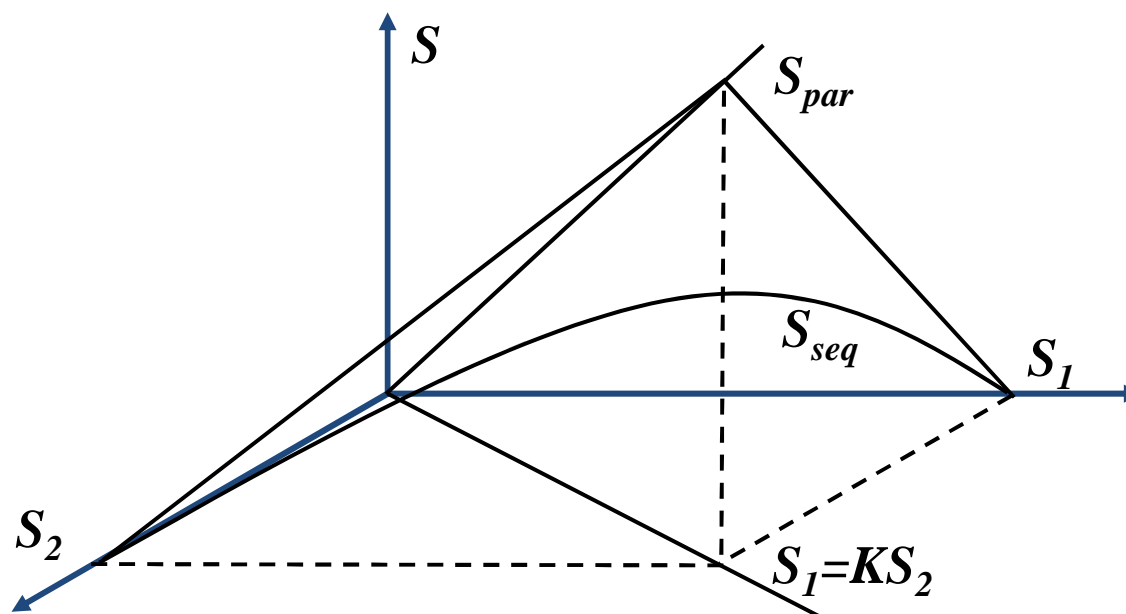
# Kombiniranje dviju paralelnih aplikacija



$$S_{seq} = \frac{T_1 + T_2}{t_1 + t_2} = \frac{S_1 S_2 (1 + K)}{S_1 + K S_2}, \quad K = \frac{T_1}{T_2},$$

$$T_i = T_p + T_c, \quad t_j = t_p + t_c + t_w$$

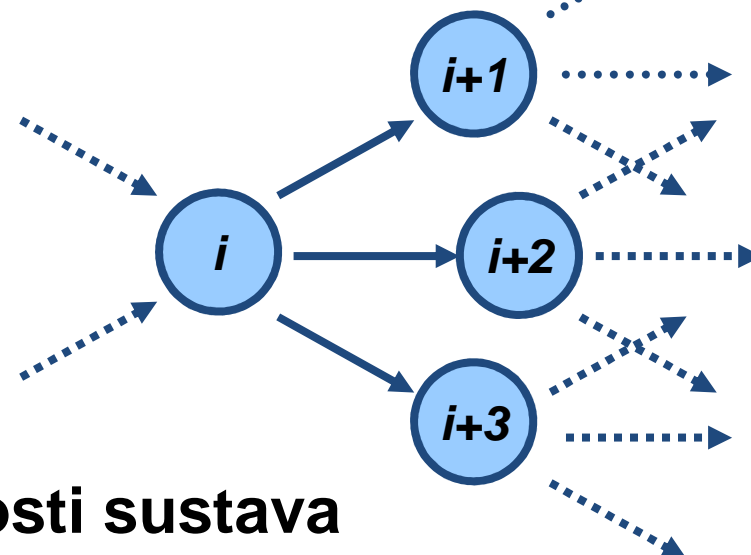
$$S_{par} = \min \left[ \frac{T_1}{t_1}, \frac{T_2}{t_2} \right] = \min[S_1, S_2]$$





## ◆ Dijagram stanja sustava

- ◆ Opisuje izmjenu stanja sustava tijekom rada



## ◆ Matrica prijelaznih vjerojatnosti sustava

- ◆ Element matrice u retku  $i$  te stupcu  $j$  određuje vjerojatnosti da će sustav iz stanja  $i$  preći u stanje  $j$

$$\begin{bmatrix} p_{0,0} & p_{0,1} & \dots & & \\ p_{1,0} & \dots & & & \\ \dots & \dots & p_{i,j} & \dots & \dots \\ & & & \dots & p_{N-2,N-1} \\ & & \dots & p_{N-1,N-2} & p_{N-1,N-1} \end{bmatrix}$$

## ◆ Faktor iskorištenja sustava ( $\rho_0$ )

- ◆ Sustav se koristi ako se nalazi u bilo kojem stanju osim početnom
- ◆ Vjerojatnost da se sustav nalazi u stanju koje nije početno:  $\rho_0 = 1 - p_0$

## ◆ Vrijeme čekanja u sustavu ( $W$ )

- ◆  $\rho_0 = U = L \cdot S$ ;  $L = \rho_0 / S$
- ◆  $L = K / (Z + W + S) = \rho_0 / S$
- ◆  $W = (K \cdot S / \rho_0) - Z - S$

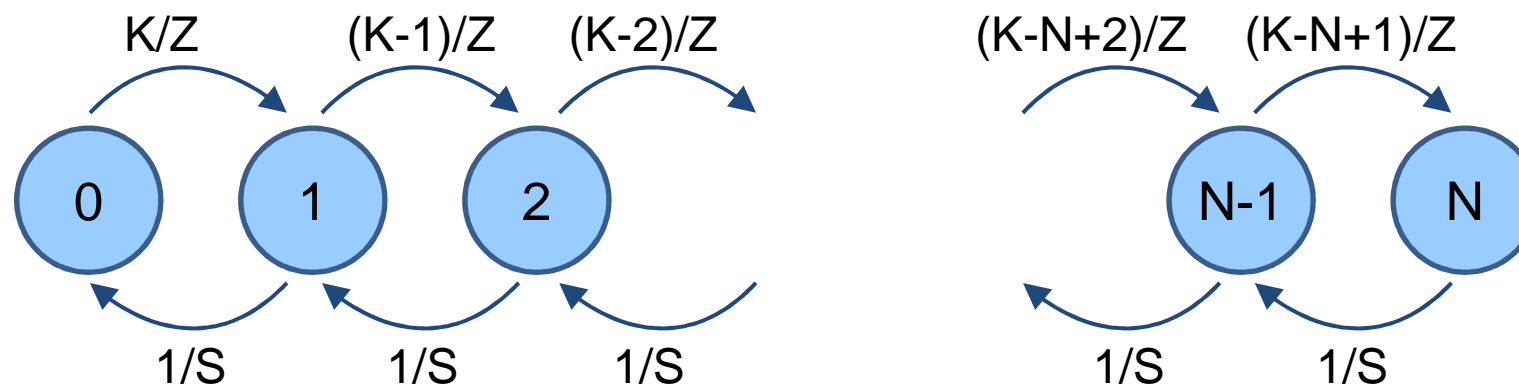
## ◆ Vrijeme zadržavanja u sustavu ( $R$ )

- ◆  $R = S + W = (K \cdot S / \rho_0) - Z$ ;  $Q = L \cdot R$

# Stohastički model paralelnog sustava (M/M/1/K/K)



Markovljev Lanac



- ◆ Vjerojatnost da se Markovljev lanac nalazi u *n*-tom ( $p_n$ ) i početnom ( $p_0$ ) stanju

$$p_n = \frac{K!}{(K-N)!} \frac{1}{(S/Z)^n} p_0$$

$$p_0 = \frac{1}{\sum_{k=0}^K \frac{(K-N)!}{(S/Z)^k}}$$

## ◆ Faktor iskorištenja sustava ( $\rho_0$ )

- ◆  $\rho_0 = 1 - p_0 = 1 - B(K, Z/S)$   
( $B(K, Z/S)$  je Erlangova formula)

## ◆ Vrijeme čekanja u sustavu ( $W$ )

- ◆  $\rho_0 = U = LS$ ;  $L = \rho_0/S$
- ◆  $L = K/(Z + W + S) = \rho_0/S$
- ◆  $W = (KS/\rho_0) - Z - S$

## ◆ Vrijeme zadržavanja u sustavu ( $R$ )

- ◆  $R = S + W = (KS/\rho_0) - Z$ ;
- ◆  $Q = L R$

## ◆ Efikasnost procesora (E)

$$◆ E = (Z + S)/(Z + S + W)$$

$$◆ E = (Z + S)/(Z + R) = (Z + S)/(Z + K \cdot S / r_o - Z)$$

$$◆ E = (Z + S) \cdot r_o / K \cdot S$$

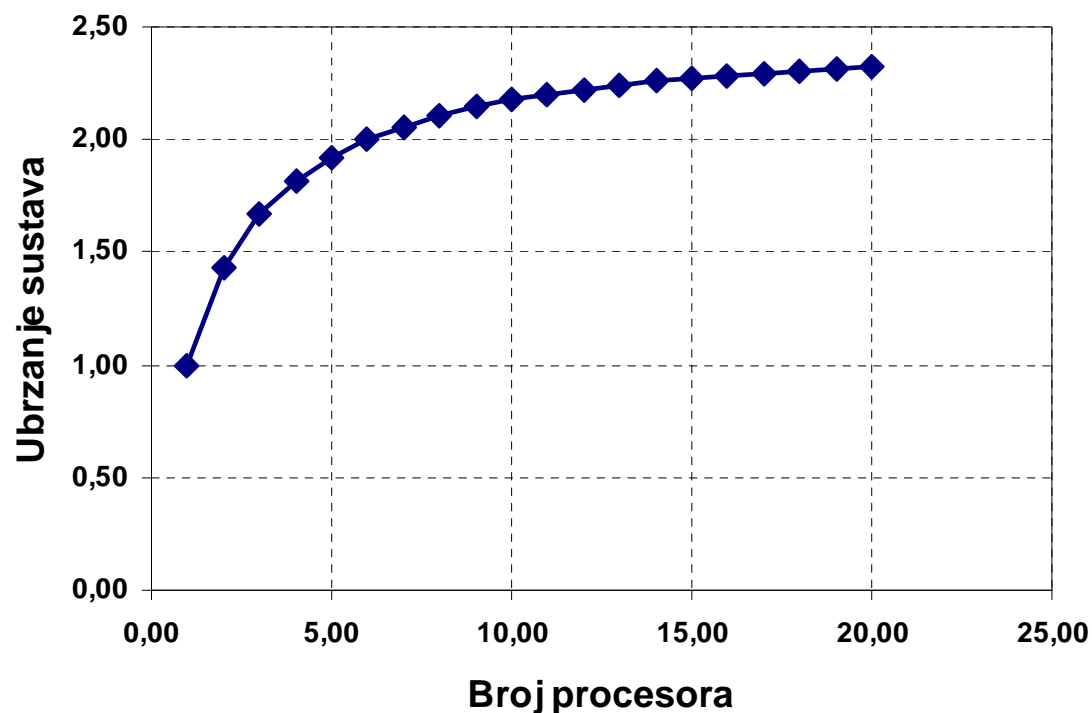
## ◆ Ubrzanje ( $S_p$ )

$$◆ \text{Ubrzanje} = K \cdot E$$

Podrazumijeva  
idealnu dekompoziciju  
tereta na K procesora



**Ubrzanje.xls**



# Primjer analize performansi web aplikacije

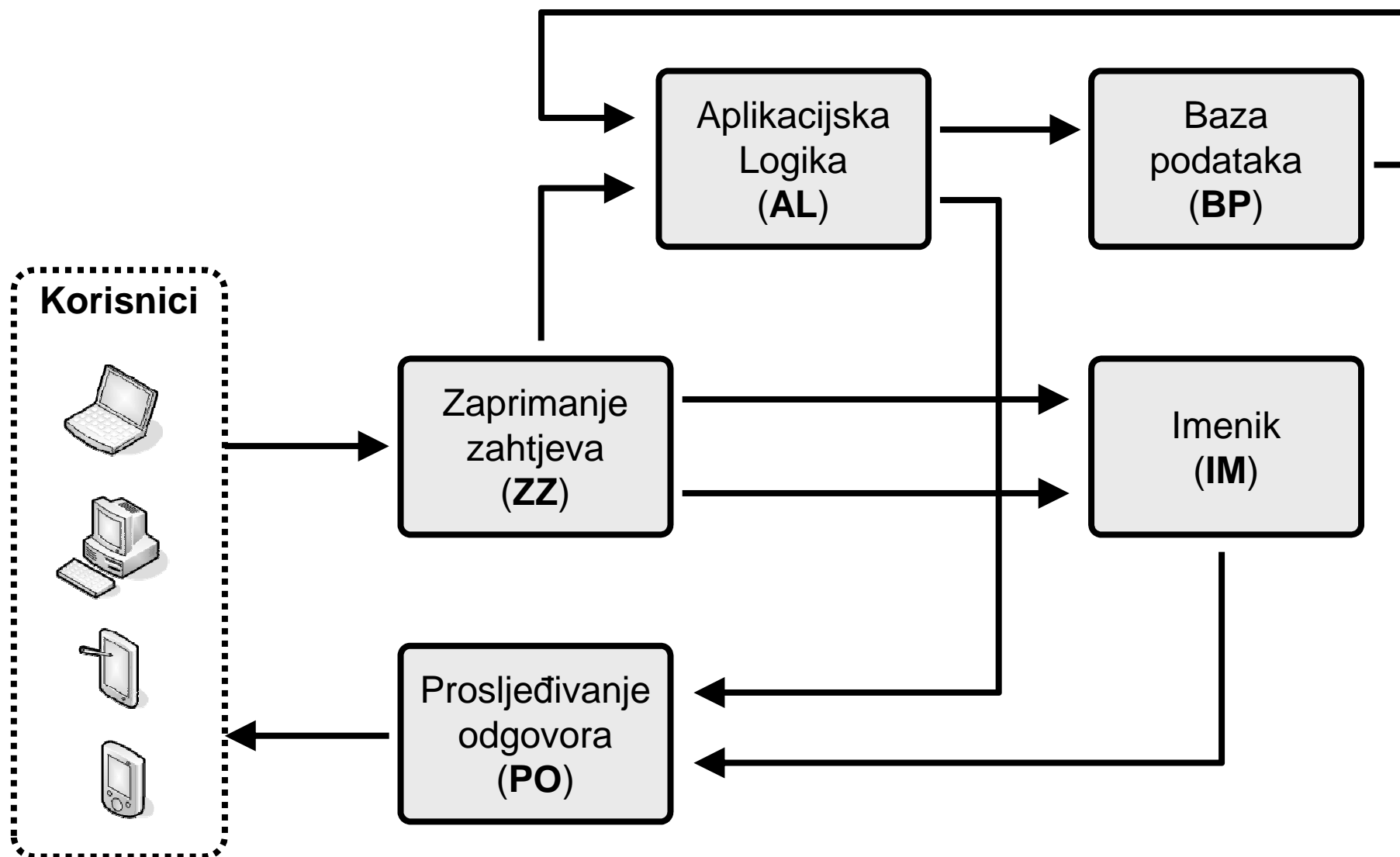
# Primjer analize raspodijeljene aplikacije

---



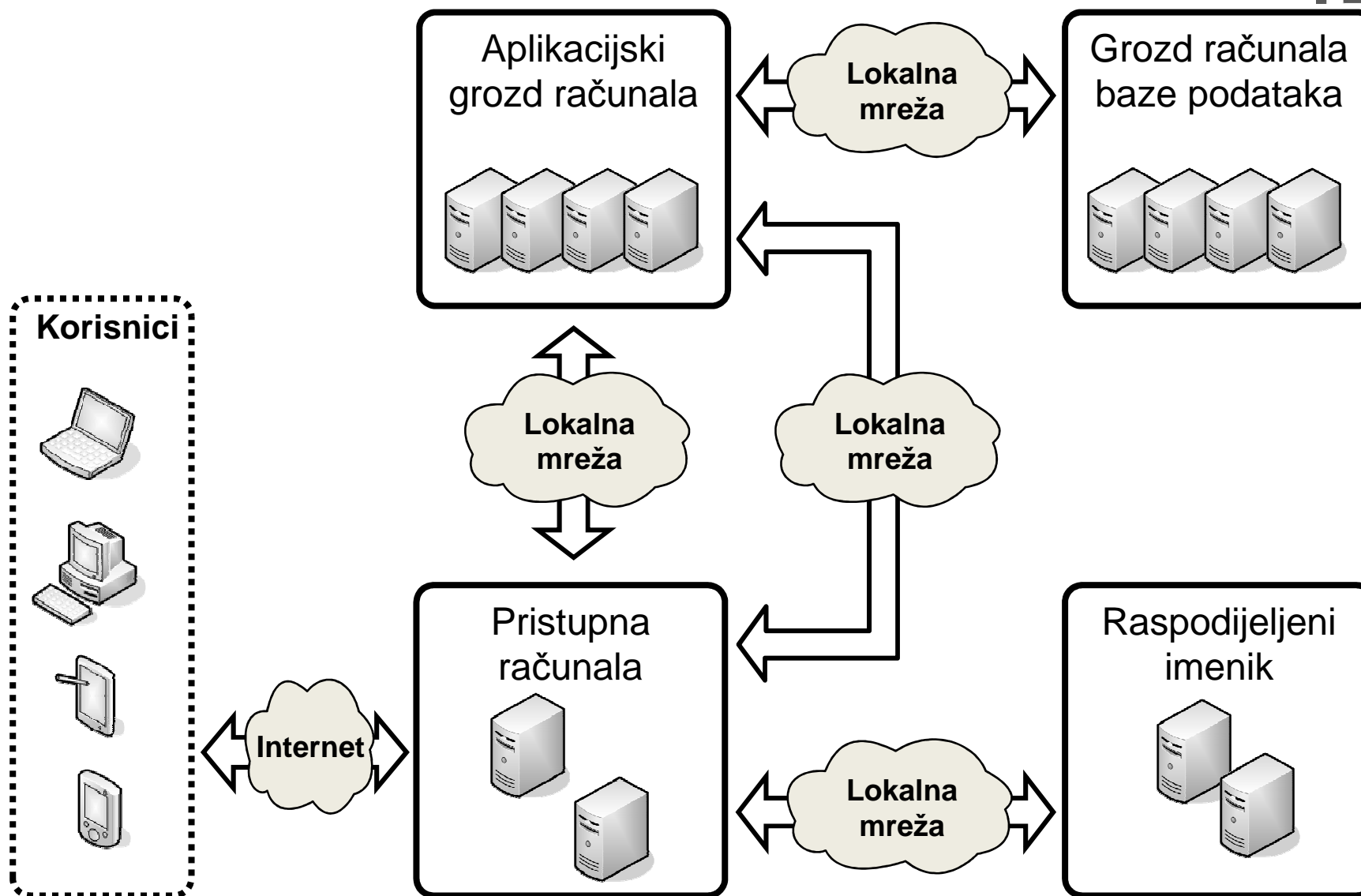
- ◆ Logička arhitektura raspodijeljene aplikacije
- ◆ Fizička arhitektura raspodijeljene aplikacije
- ◆ Model raspodijeljene aplikacije
- ◆ Vrednovanje značajki performansi aplikacije

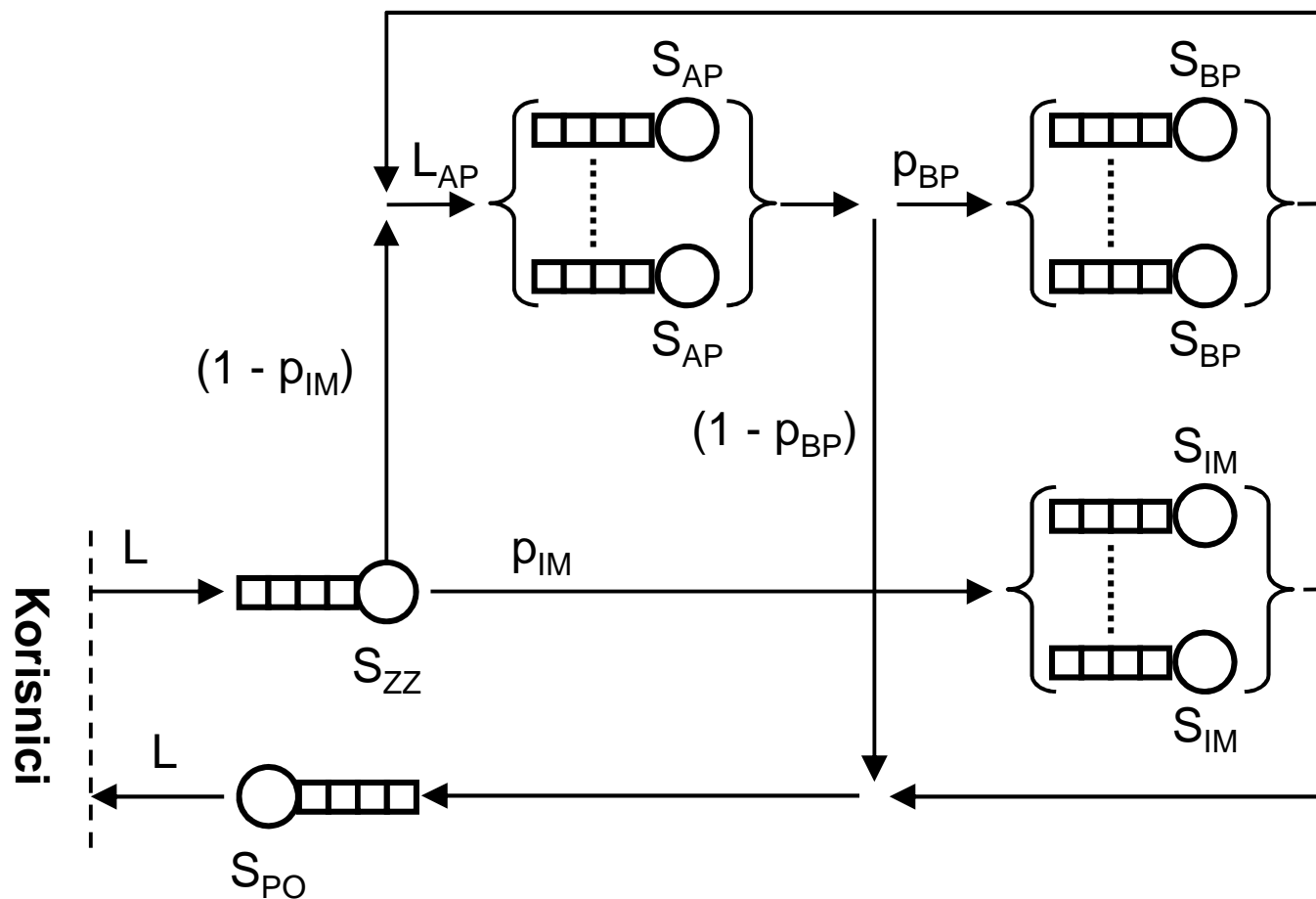
# Logička arhitektura aplikacije





# Fizička arhitektura aplikacije





# Primjer analize raspodijeljene aplikacije



- ◆ Svaki od podsustava aplikacije na jednom računalu
- ◆ Učestalost dolazaka zahtjeva na podsustave
  - ◆  $L_{IM} = p_{IM}L$ ,  $v_{IM} = p_{IM}$
  - ◆  $L_{AP} = p_{BP}L_{AP} + (1 - p_{IM})L \Rightarrow$   
 $L_{AP} = [(1 - p_{IM})/(1 - p_{BP})]L$   
 $v_{AP} = [(1 - p_{IM})/(1 - p_{BP})]$
  - ◆  $L_{BP} = p_{BP} L_{AP} \Rightarrow$   
 $L_{BP} = p_{BP} [(1 - p_{IM})/(1 - p_{BP})] L$   
 $v_{BP} = p_{BP} [(1 - p_{IM})/(1 - p_{BP})]$
  - ◆  $L_{ZZ} = L$ ,  $v_{ZZ} = 1$
  - ◆  $L_{PO} = L$ ,  $v_{PO} = 1$

## ◆ Skalirana vremena posluživanja

$$◆ D_{IM} = v_{IM} S_{IM} = p_{IM} S_{IM}$$

$$◆ D_{AP} = v_{AP} S_{AP} = [(1 - p_{IM})/(1 - p_{BP})] S_{AP}$$

$$◆ D_{BP} = v_{BP} S_{BP} = p_{BP} [(1 - p_{IM})/(1 - p_{BP})] S_{BP}$$

$$◆ D_{ZZ} = v_{ZZ} S_{ZZ} = 1 S_{BP}$$

$$◆ D_{PO} = v_{PO} S_{PO} = 1 S_{PO}$$

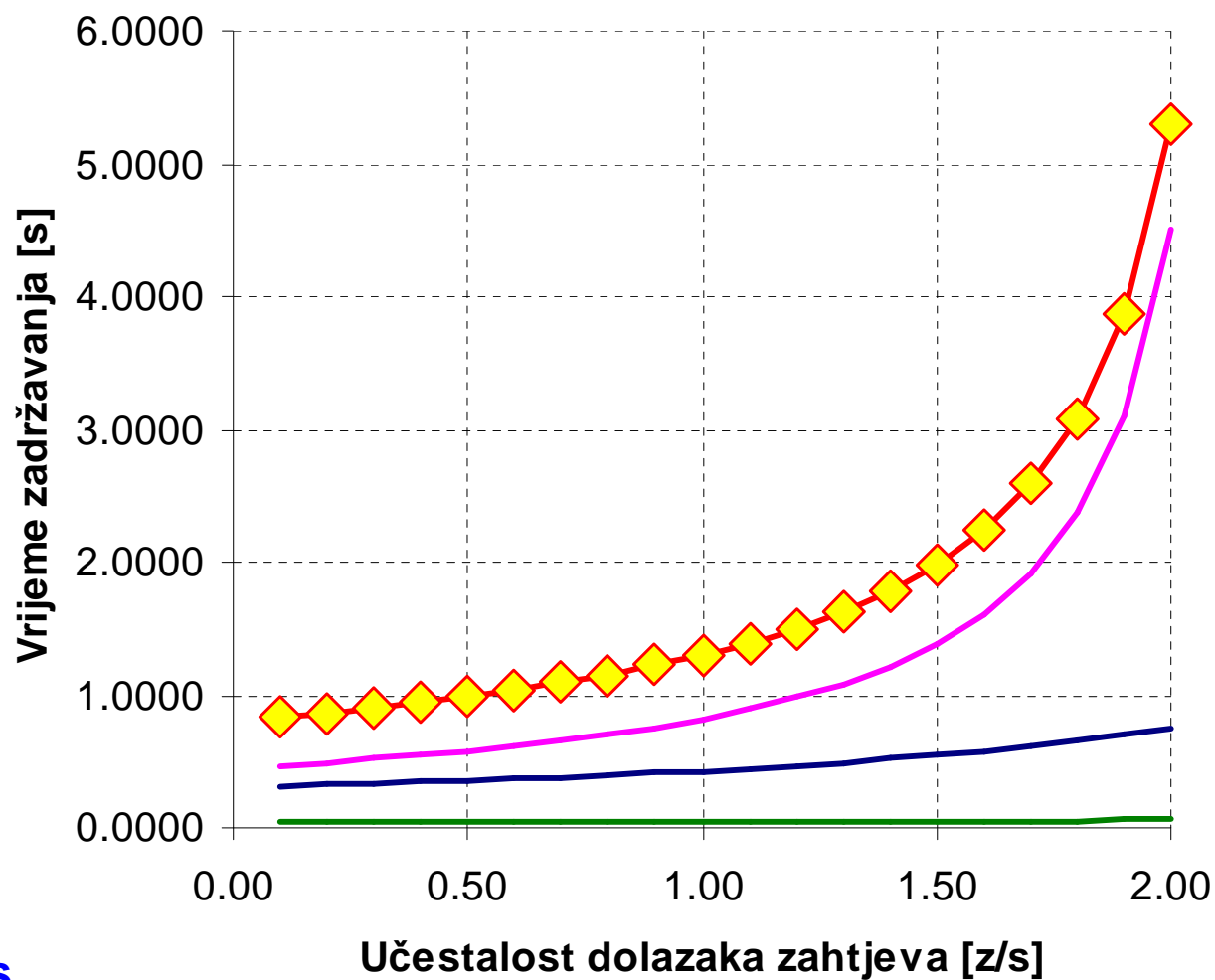
## ◆ Vrijeme zadržavanja zahtjeva

$$◆ R = D_{IM}/(1 - L D_{IM}) + D_{AP}/(1 - L D_{AP}) + D_{BP}/(1 - L D_{BP}) + \\ D_{ZZ}/(1 - L D_{ZZ}) + D_{PO}/(1 - L D_{PO})$$

# Vrijeme zadržavanja zahtjeva



— R\_AP — R\_BP — R\_IM — R — ♦ R\_PDQ



$$N_{AP} = 1$$

$$N_{BP} = 1$$

$$N_{IM} = 1$$

$$S_{AP} = 0.3$$

$$S_{BP} = 4.5$$

$$S_{IM} = 0.5$$

$$S_{ZZ} = 0.001$$

$$S_{PO} = 0.001$$

$$p_{IM} = 0.1$$

$$P_{BP} = 0.1$$



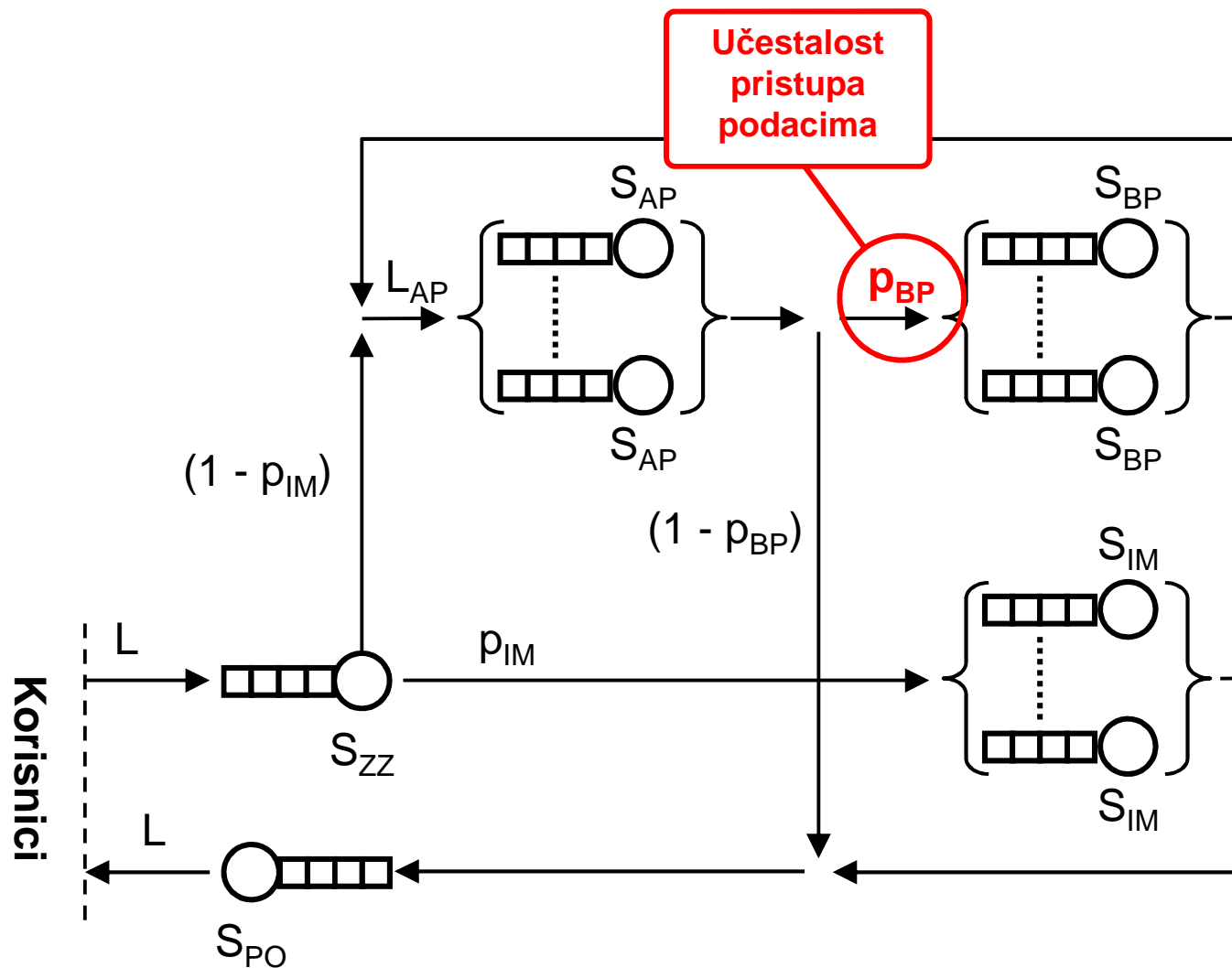
ap-Zad.c



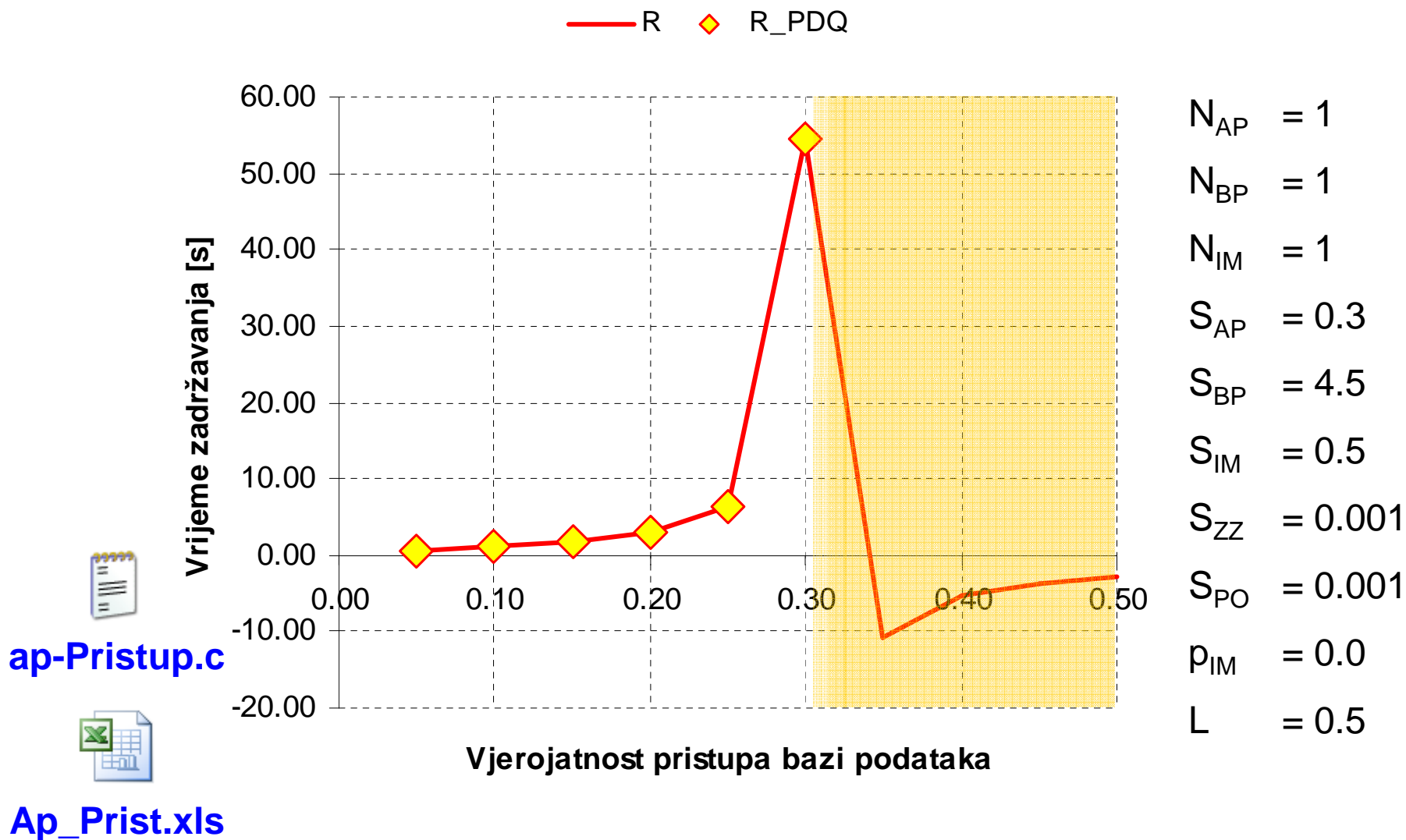
Ap\_Zad.xls

- ◆ Učestalosti pristupa podacima
- ◆ Veličina grozda baze podataka
- ◆ Promjena organizacije podataka
- ◆ Promjena stupnja sigurnosti

# Učestalost pristupa podacima

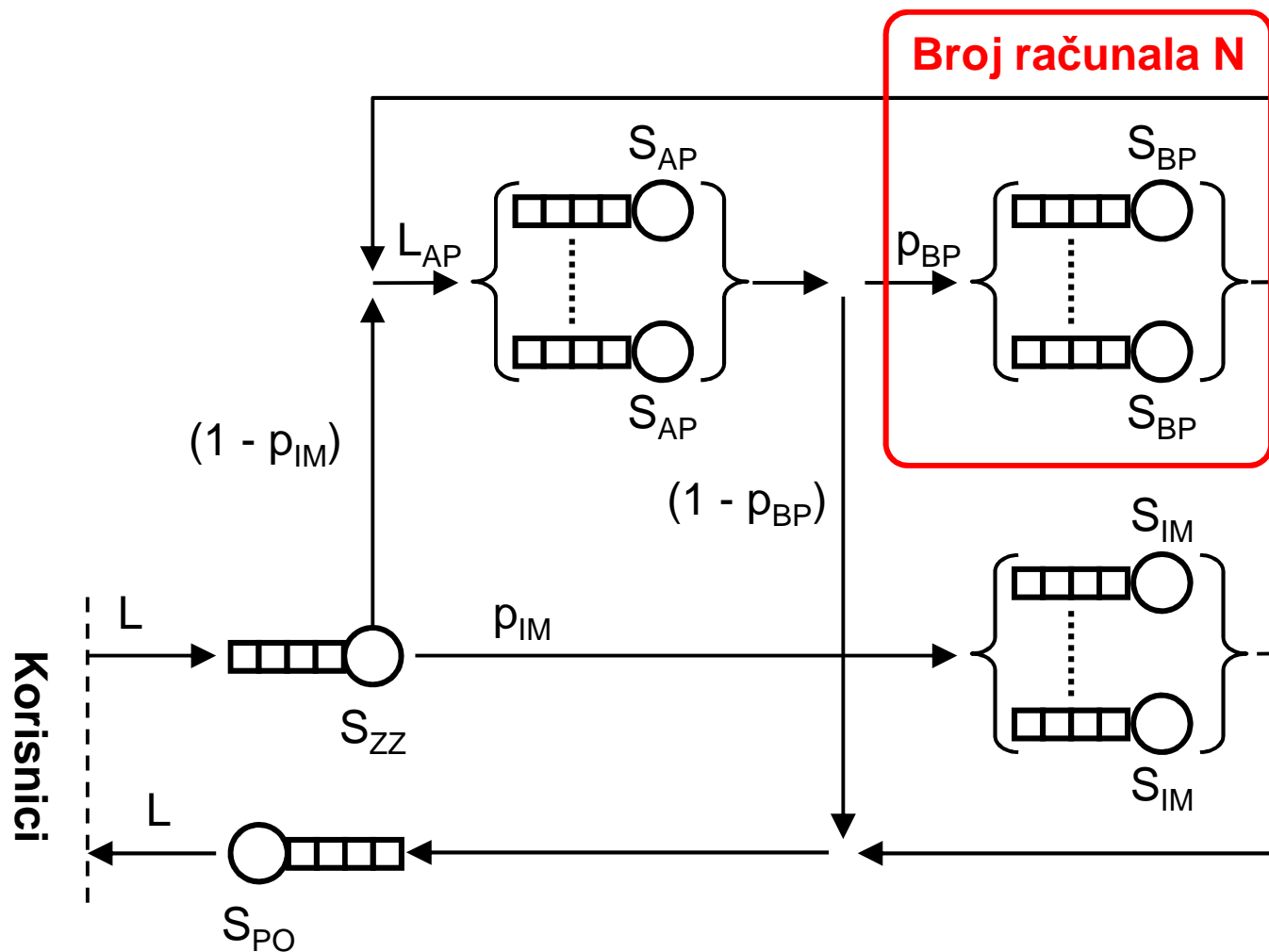


# Vrijeme zadržavanja zahtjeva





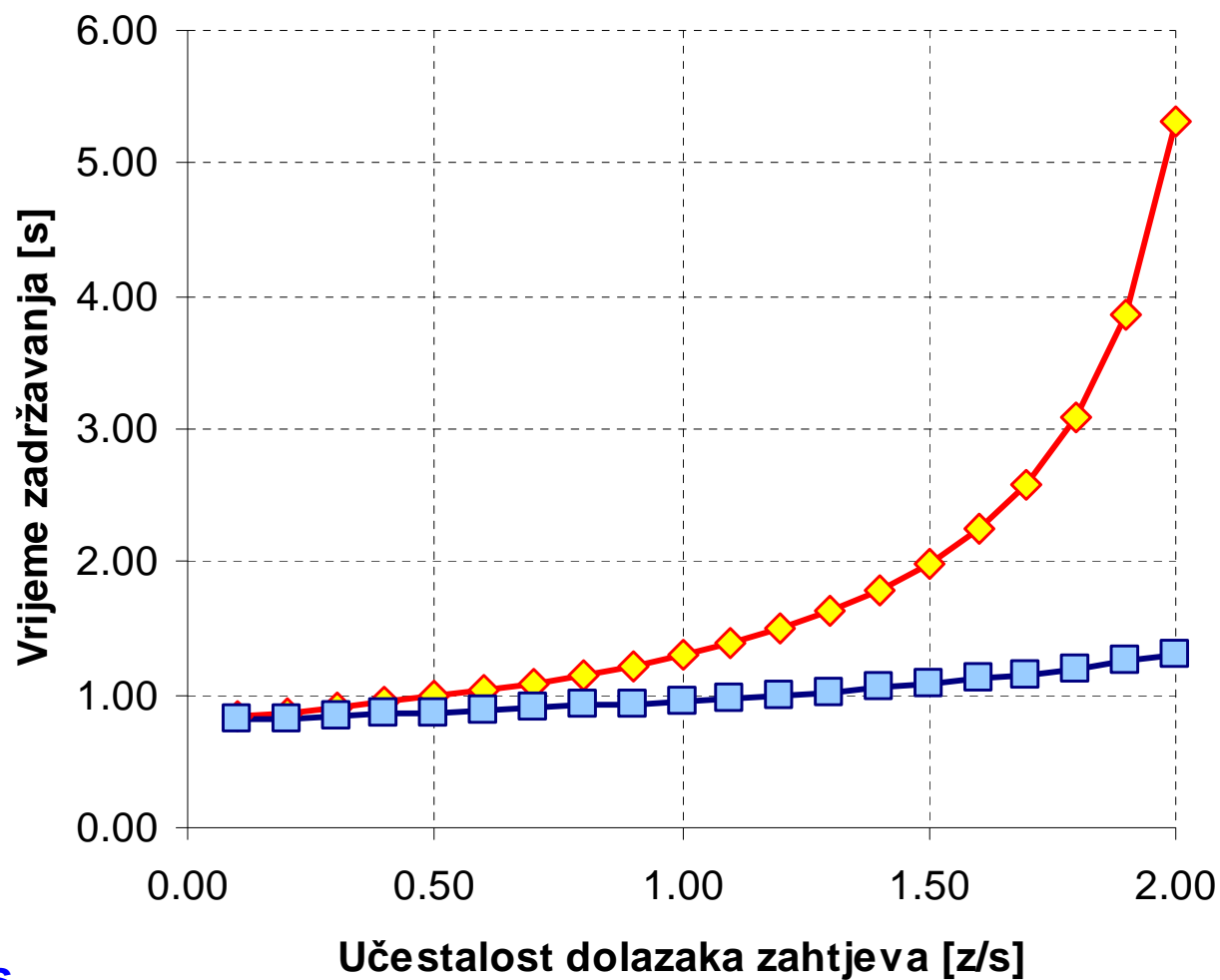
# Veličina grozda baze podataka



# Vrijeme zadržavanja zahtjeva



—◆—  $N = 1$  —■—  $N = 10$



$$S_{AP} = 0.3$$

$$S_{BP} = 4.5$$

$$S_{IM} = 0.5$$

$$S_{ZZ} = 0.001$$

$$S_{PO} = 0.001$$

$$\rho_{IM} = 0.1$$

$$P_{BP} = 0.1$$



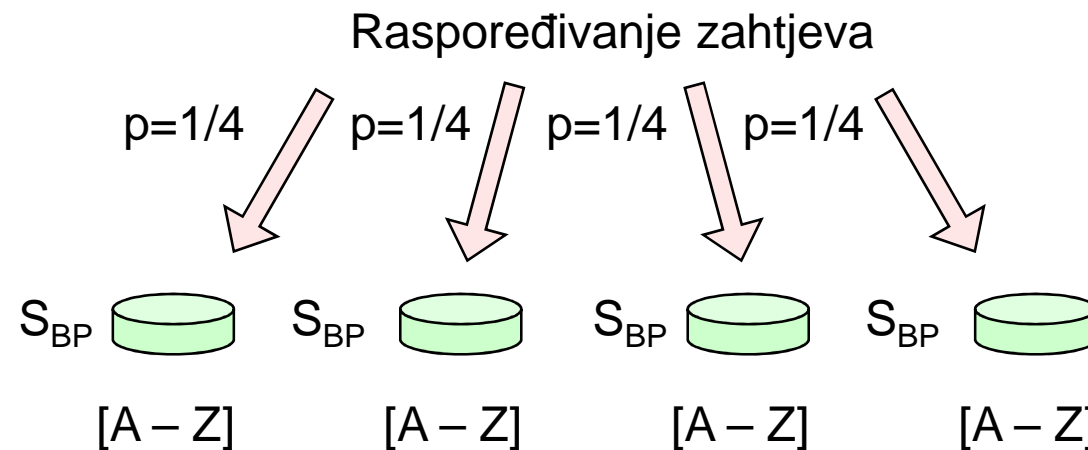
ap-Zad.c



Ap\_Rac.xls

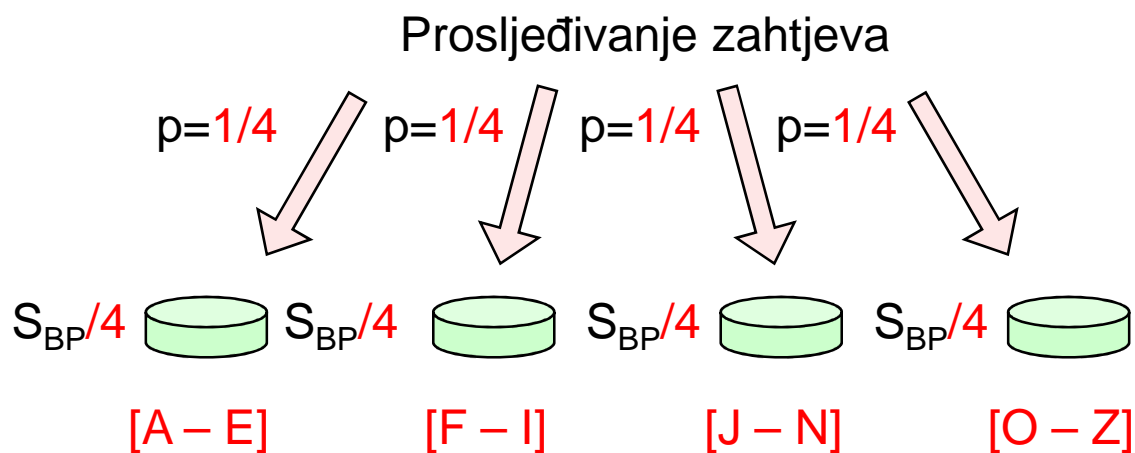
## ◆ Replikacija podataka

- ◆ Skup računala od kojih svako u spremniku sadrži kopiju cijele baze podataka
- ◆ Zahtjevi se **raspoređuju** na računala s ciljem raspoređivanja opterećenja



## ◆ Segmentacija podataka

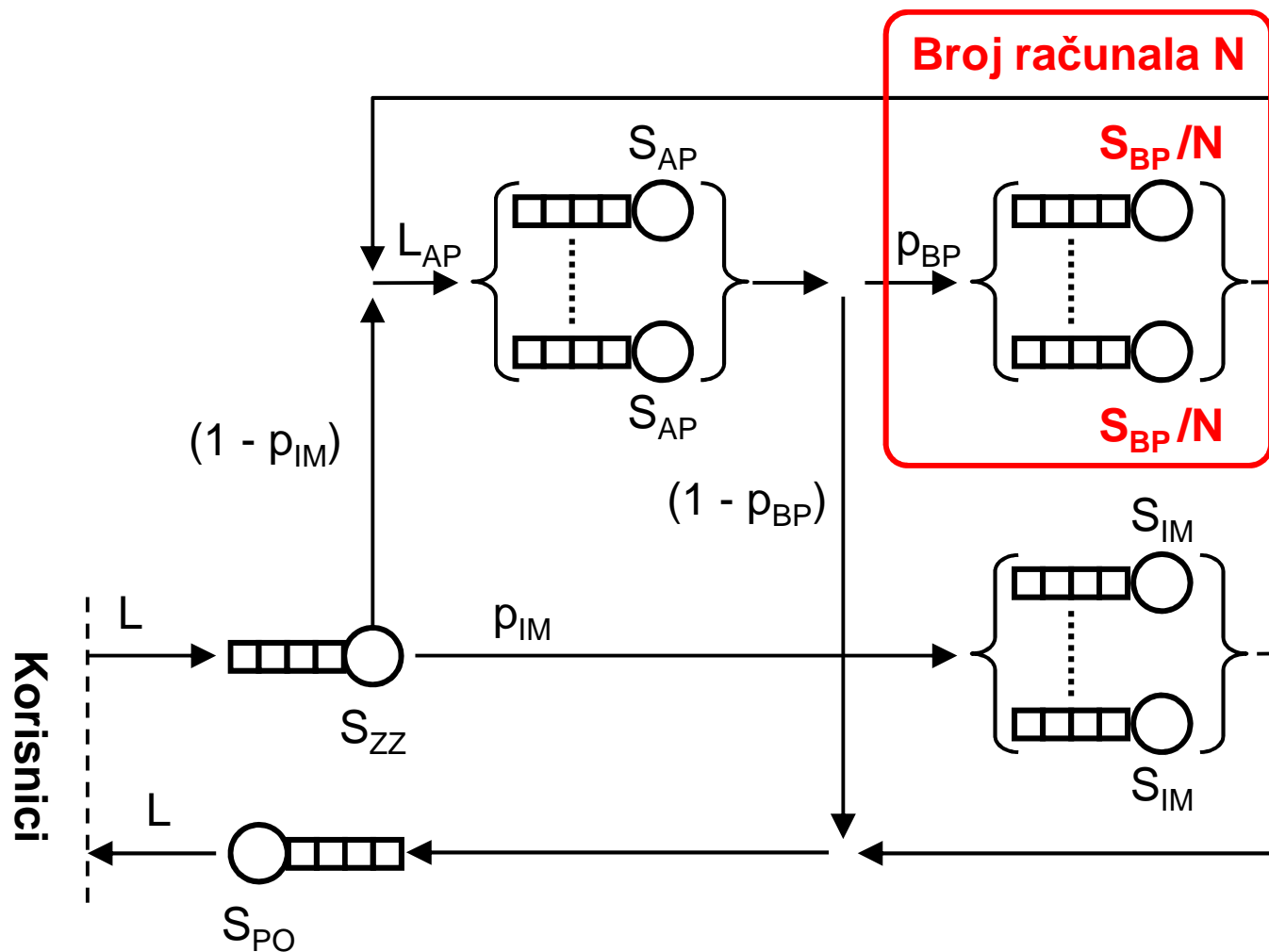
- ◆ Skup računala od kojih svako u spremniku sadrži dio cijele baze podataka
- ◆ Zahtjevi se **prosljeđuju** prema računalu s traženim zapisima



## ◆ Pretpostavke

- ◆ Uniformna raspodjela zahtjeva na zapise
- ◆ Linearna složenost obrade zahtjeva o količini zapisa

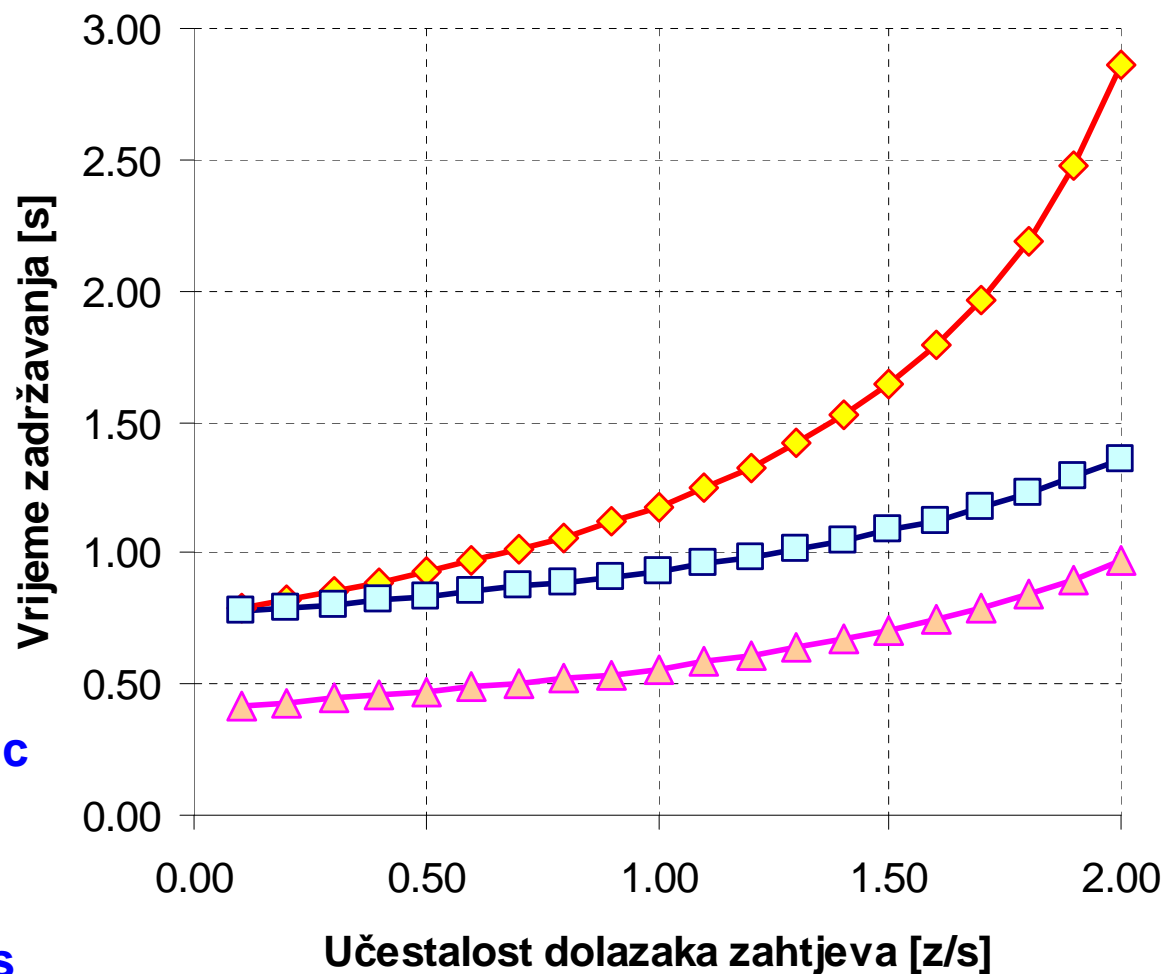
# Promjena organizacije podataka



# Vrijeme zadržavanja zahtjeva



—◆—  $N = 1$  (REP) —□—  $N = 10$  (REP) —△—  $N = 10$  (SEG)



$N_{AP} = 1$

$N_{BP} = 1$

$N_{IM} = 1$

$S_{AP} = 0.3$

$S_{BP} = 2.5$

$S_{IM} = 0.5$

$S_{ZZ} = 0.001$

$S_{PO} = 0.001$

$p_{IM} = 0.1$

$p_{BP} = 0.15$



[ap-pod.org.c](http://ap-pod.org.c)



[Ap\\_Pod.xls](#)

- ◆ **Imenik aplikacije sadrži informacije o korisnicima**

- ◆ Korisnički identiteti
- ◆ Korisnička prava pristupa

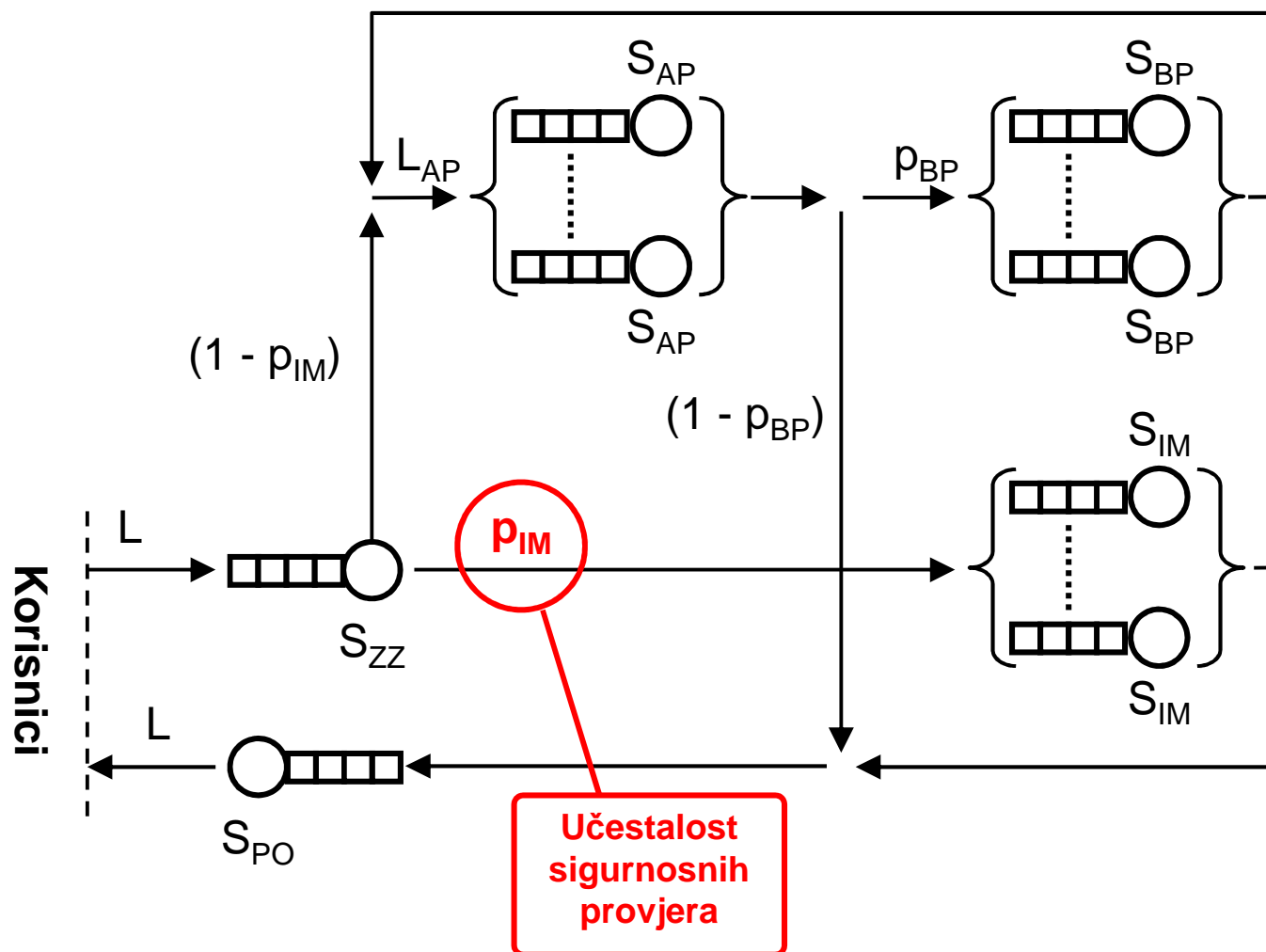
- ◆ **Sigurnosna značka**

- ◆ Određuje sigurnosne postavke korisnika aplikacije
- ◆ Značka se dohvaća iz imenika

- ◆ **Životni vijek sigurnosne značke**

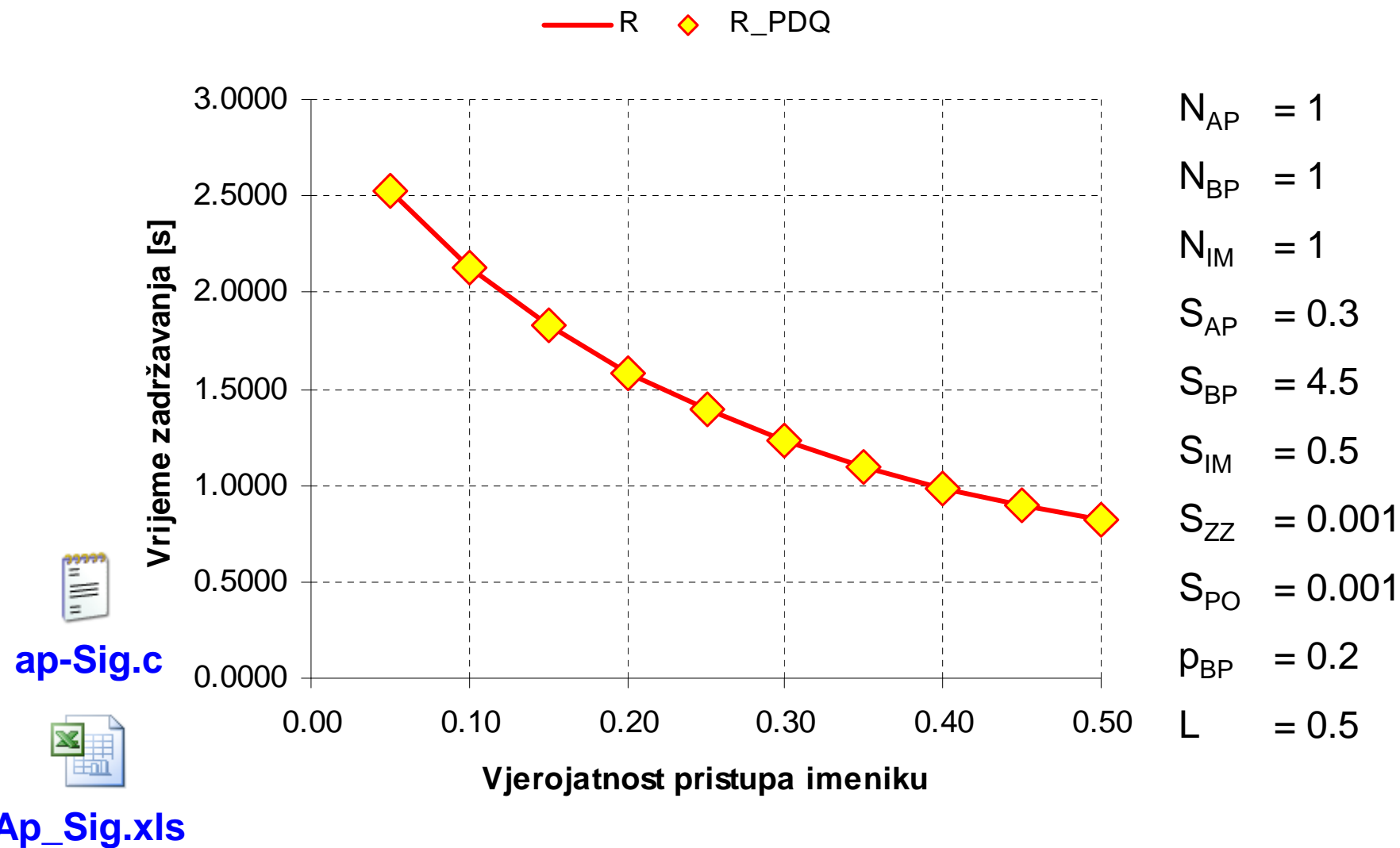
- ◆ Ograničeni broj pristupa
- ◆ Zadano vrijeme korištenja
- ◆ Ostali sigurnosni modeli

# Promjena stupnja sigurnosti



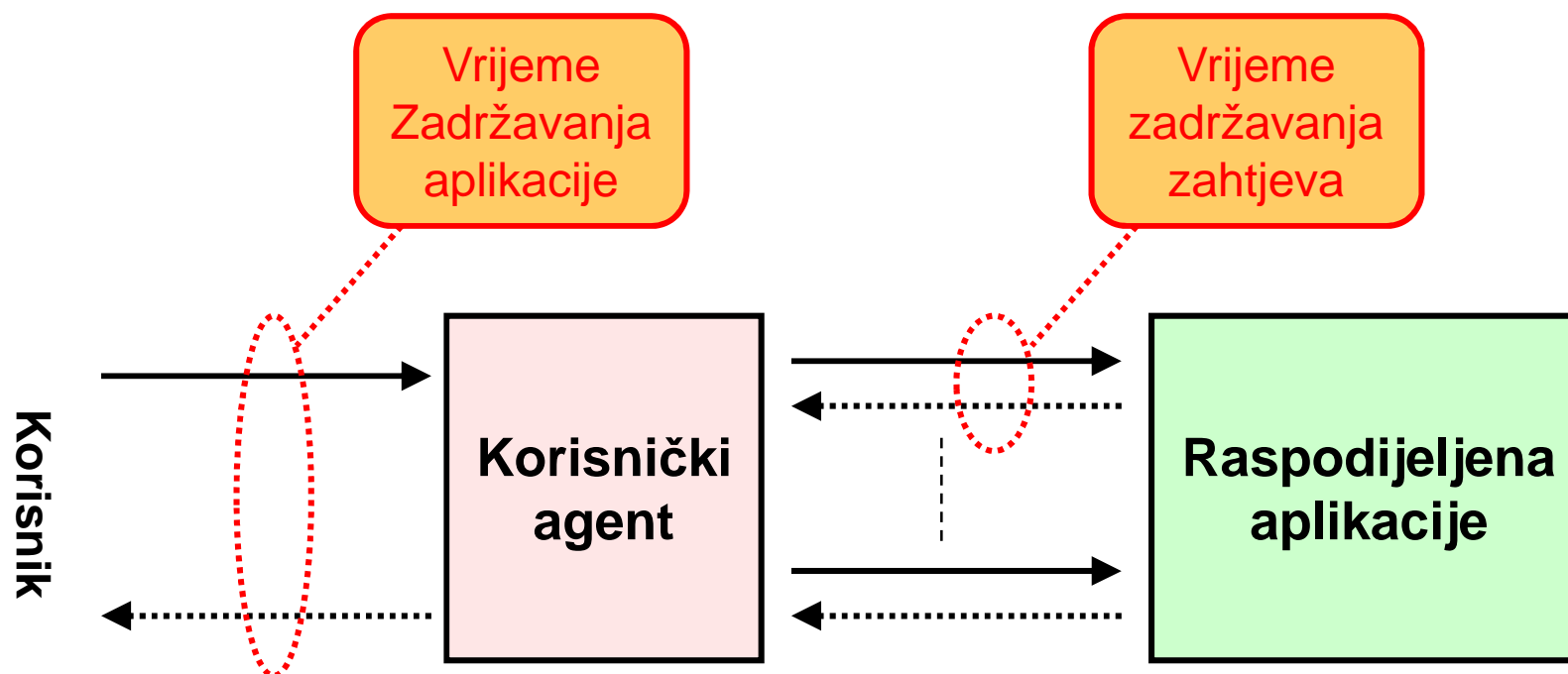


# Vrijeme zadržavanja zahtjeva



## ◆ Zašto vrijeme zadržavanja opada?

- ◆ Modelirano je vrijeme zahtjeva zadržavanja ali ne i ukupno vrijeme zadržavanja aplikacije koje doživljava korisnik



◆ **Zadatak 1:** Web aplikacija uključuje podršku korisnicima putem chat usluge. Kupci sami odabiru jedan od 10 repova čekanja u kojima upite poslužuje po jedan tehničar. Mjerenja pokazuju da zahtjevi prosječno dolaze 3 upita u minuti te da svaki kupac prosječno čeka 3 minute u repu i prosječno provodi 2 minute u razgovoru.  
*(nadogradnja primjera 6)*

**1) Kakvi će biti odzivi sa 10 i 18 tehničara ako publiciranje Web stranice sa odgovorima na najčešća pitanja smanji broj upita na 2 u minuti?**

**2) Kakve će rezultate dati smanjenje razgovora na 1.5 minutu?**

◆ **Zadatak 2:** Oblikovati proizvoljnu raspodijeljenu aplikaciju i ostvariti analizu performansi ostvarene aplikacije

**1) Definirati logičku i fizičku arhitekturu aplikacije**

**2) Izgraditi model aplikacije primjenom teorije repova**

- Odrediti analitičko rješenje funkcije zadržavanja zahtjeva u aplikaciji  $R=f(L)$

**3) Izgraditi model aplikacije za alat PDQ**

- Primjenom izgrađenog modela odrediti vrijednosti funkcije zadržavanja zahtjeva  $R=f(L)$  u nekoliko točaka

**4) Usporediti i obrazložiti dobivene rezultate**

# Rekapitulacija

- ◆ **Zašto je vrednovanje performansi važno ?**
  - ◆ Zbog potrebe za planiranjem kapaciteta koji je potreban za uspješno poslovanje
- ◆ **Koliko razina protokola uključuju Web aplikacije ?**
  - ◆ Četiri: Ethernet, IP, TCP/IP i HTTP
- ◆ **Koji su osnovni dijelovi odziva Web aplikacije ?**
  - ◆ Klijent-ISP, ISP-ISP, ISP-servis
- ◆ **Kako se ostvaruje razmjernan rast aplikacije ?**
  - ◆ Vertikalno ili horizontalno
- ◆ **Koje probleme donosi horizontalan rast ?**
  - ◆ Sinkronizacija, razdioba tereta i razdioba podataka
- ◆ **Koju mogućnost donosi horizontalan rast ?**
  - ◆ Manja osjetljivost na greške
- ◆ **Koje konfiguracije podržavaju neosjetljivost na greške ?**
  - ◆ Aktivan-pripravan i Aktivan-aktivan

- ◆ **Kako se ostvaruje razdioba tereta ?**
  - ◆ Sklopkama koje rade na IP ili aplikacijskoj razini
- ◆ **Što je replikacija ?**
  - ◆ Replikacija osigurava razdiobu istih podataka na više sustava
- ◆ **Kakve vrste replikacije postoje ?**
  - ◆ Master-slave, master-master-slave, stablo, stablo s filterima, master-master i master ring
- ◆ **Što je federacija ?**
  - ◆ Federacija dijeli podatke u različite grupe koje raspoređuje na različite sustave
- ◆ **Zašto se koriste grupni protokoli ?**
  - ◆ Za pouzdanu dostavu svim članovima grupe i garantirani slijed poruka
- ◆ **Koje su tipične garancije za slijed poruka?**
  - ◆ FIFO, posljedična i totalna

- ◆ **Koje su tipične zone u sustavu Web servisa ?**
  - ◆ Demilitarizirana, aplikacijska i zona podataka
- ◆ **Koje su dvije osnovne grupe troškova za gradnju i pogon Web sustava i kako su obično raspodijeljeni?**
  - ◆ Kapitalni i operacioni troškovi grubo raspodijeljeni 50/50 kroz tri godine
- ◆ **Koje su mogućnosti za udomljivanje Web servisa?**
  - ◆ Udomitelj infrastrukture ili vlastito udomljavanje
- ◆ **Kako se mjeri i modelira kapacitet servisnog sustava?**
  - ◆ Mjerenjem pomoću generatora tereta, te aproksimacijom pomoću polinoma
- ◆ **Kako se određuje vrijeme potrebno da se dostigne potreba za dvostrukim kapacitetom?**
  - ◆ Podrazumijevajući eksponencijalni trend



- ◆ **Koje su osnovne veličine u modelu repa čekanja ?**
  - ◆ Vrijeme obzervacije (T), broj dolazaka (A), broj odlazaka (C) i vrijeme zaposlenosti poslužitelja (B)
- ◆ **Koje su izvedene veličine ?**
  - ◆ Ulazni ritam ( $L=A/T$ ), izlazni ritam ( $X=C/T$ ), srednje vrijeme posluživanja ( $S=B/C$ ) i zaposlenost poslužitelja ( $U=B/T$ )
- ◆ **Kako se definira stacionarno stanje ?**
  - ◆  $X = L$
- ◆ **Kako glasi Little-ov zakon ?**
  - ◆ Broj zahtjeva u repu proporcionalan je ritmu dolaska zahtjeva i vremenu provedenom u sustavu ( $Q = L \cdot R$ )
- ◆ **Kako je definirano vrijeme čekanja u repu ?**
  - ◆  $W = Q \cdot S$
- ◆ **Kako je definirano ukupno vrijeme provedeno u sustavu?**
  - ◆  $R = S + W = Q/L$

- ◆ **Kako se izračunava vrijeme odziva za serijske repove?**
  - ◆  $R = (Q1 + Q2 + \dots + QN)/L$
- ◆ **Kako je definirana iskoristivost procesora  $ro$ ?**
  - ◆  $ro = U/N$
- ◆ **Kako se izračunava vrijeme odziva za paralelne repove?**
  - ◆  $R = S/(1 - ro)$
- ◆ **Sto utječe na  $R$  kod paralelnih poslužitelja?**
  - ◆  $R = S/(1 - ro^N)$  ;  $ro = U/N$  (aproksimacija)
  - ◆  $R = S(1 + (C(N,ro) / (N*(1 - ro))))$   
(egzaktno rješenje uz Erlangovu formulu)
- ◆ **Kolika je najveća greška aproksimacije?**
  - ◆ Manja od 8% u praktičnim situacijama
- ◆ **Kako se računa  $R$  za sustav sa vjerojatnošću povratne veze  $p$ ?**
  - ◆  $L1 = L/(1 - p)$ ;  $U = L1 * S$ ;  $R1 = S*(1 + U/(1-U))$ ;  $R = R1/(1 - p)$

- ◆ **Kako je definiran zatvoreni sustav ?**
  - ◆ Broj zahtjeva je ograničen
- ◆ **Kako se izračunava R za zatvoreni sustav sa m izvora zahtjeva ?**
  - ◆  $X(m) = (m - Q)/Z$  ;  $Q = X(m) * R$  ;  $R = m / X(m) - Z$  ;  $X(m) = U/S$
- ◆ **Što je efikasnost procesora a što ubrzanje u paralelnim sustavima ?**
  - ◆  $E = (t_p + t_c)/(t_p + t_c + t_w)$  ;  $S = ((T_p + T_c)/t_p + t_c) * E$
- ◆ **Kako je definirana funkcija dekompozicije ?**
  - ◆ Funkcija od N koja definira kako se paralelna aplikacija dijeli na N procesora
- ◆ **Koji je nagori i najbolji slučaj rada paralelne aplikacije ?**
  - ◆ Sinkrono i asinkrono izvršavanje
- ◆ **Koji stohastički proces modelira ponašanje paralelnog sustava ?**
  - ◆ Markovljev lanac