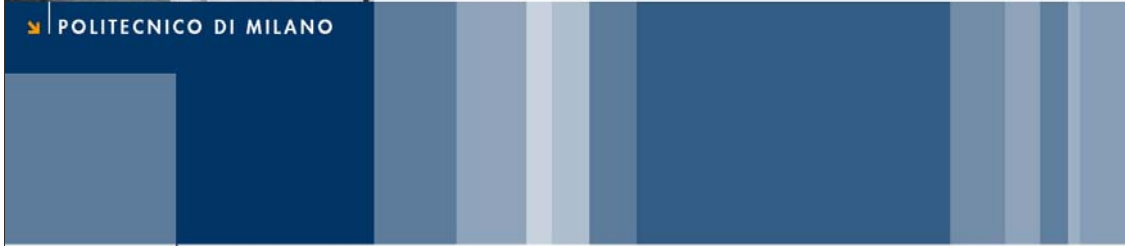




 POLITECNICO DI MILANO

*Paolo Cremonesi*

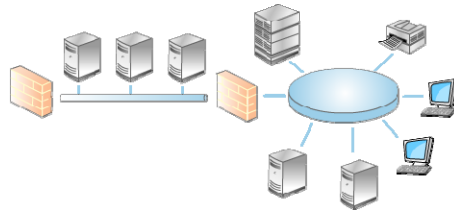
# Impianti Informatici



Modelli a rete di code:  
modelli aperti e modelli chiusi



- Modelli aperti, classi aperte
  - Numero potenzialmente illimitato di richieste nel sistema
- Admission control



$$N = 1024$$

- Modelli chiusi, classi chiuse
- Modelli misti

La lezione che andremo a presentare ci consentirà di apprendere alcune semplici tecniche che, assieme alle leggi dell'analisi operativa viste nella lezione precedente, consentono la risoluzione completa di un modello a reti di code.

- Prima di procedere, dobbiamo però operare una ulteriore distinzione tra le reti in forma prodotta che consideriamo

- Le reti che abbiamo considerato finora sono dette reti aperte, in quanto il nostro sistema comunica con il mondo esterno e da questi può ricevere

- un numero a priori illimitato di richieste da processare

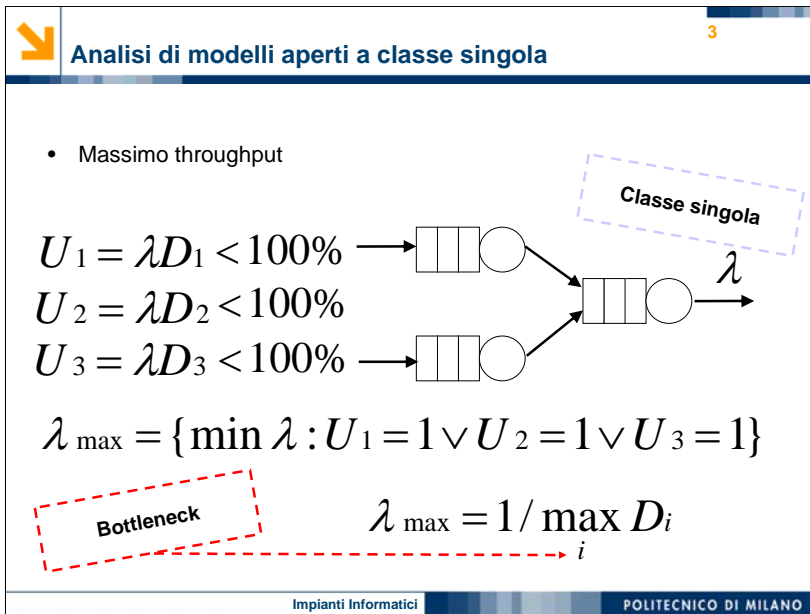
Tuttavia, un'ampia famiglia di componenti e sistemi reali ha delle limitazioni circa il massimo numero di richieste presenti nel sistema. Tali vincoli, sono specificati tramite politiche di admission control che non permettono l'accesso a nuove richieste se il sistema è pieno.

- ad esempio un router con un buffer di 1024 richieste, non può accettarne di nuove fino a quando non ha smaltito almeno una richiesta di quelle in attesa al suo interno. Inoltre, poiché al liberarsi di un slot del buffer è probabile che il router riceva subito una nuova richiesta da servire,

- il comportamento del router è meglio descritto da una coda che processa una popolazione costante  $N=1024$

- I modelli a reti di code che descrivono sistemi con numero di richieste al proprio interno costante sono detti modelli chiusi

- Infine chiamiamo modelli misti quei modelli con più classi di richieste, alcune aperte, ovvero senza limiti al numero di richieste in attesa nel sistema, in parte chiuse, ovvero con popolazione costante,



• Cominciamo allora a vedere le tecniche risolutive esistenti per i modelli a reti di code aperti.

• Per non complicarci la vita, supporremo inizialmente che tutte le richieste appartengono ad un'unica classe, ovvero che il nostro modello sia a classe singola. Nei modelli a classe singola è solitamente più comodo omettere gli indici  $c$  che denotano la classe delle richieste

• Ad esempio,  $\lambda$  indicherà il tasso di arrivo dell'unica classe di servizio nella rete

• Ricordiamo, inoltre, che per la legge del bilanciamento di flusso  $\lambda$  equivale anche al throughput della rete e, dunque, essendo  $\lambda$  un parametro noto richiesto per la specifica del modello, possiamo sempre assumere di conoscere il throughput di un modello aperto

• La prima informazione di performance che vogliamo ottenere è il massimo throughput che la rete può sostenere, o equivalentemente, il massimo tasso di arrivo medio che può ricevere

In generale, prescindendo dal tempo di risposta sperimentato dalle richieste, l'unico vincolo esistente che limita la capacità della rete deriva dalla legge dell'utilizzo

• Per le tre code del sistema rappresentato, la legge esprime per ciascuna il legame tra il throughput e la domanda globale di servizio alla stazione

• Ma poiché l'utilizzo è la percentuale di tempo speso per servire le richieste, esso non può mai superare il 100%

• Avremo cioè che il fattore limitante del throughput è l'utilizzo delle stazioni, e il throughput massimo sarà avrà per il primo valore di  $\lambda$  che porta ad utilizzo 1, ovvero ad utilizzo del 100%, uno qualunque tra i server della rete. In questo caso si parla di saturazione di un server della rete.

• Con semplici passaggi algebrici si vede che il minimo dei  $\lambda$  è dato dall'inverso del massimo valore delle  $D_i$ , tale valore è detto massimo throughput o tasso di arrivo di saturazione

• la stazione  $i$  per la quale l'utilizzo sale al 100% è detta la stazione bottleneck della rete, cioè il collo di bottiglia delle prestazioni

• Nel caso in cui esistano più stazioni bottleneck con lo stesso valore delle  $D_i$  con  $i$  si parla di reti con bottleneck multipli

Calcolo del tempo di risposta di una rete aperta

4

---

- Misura additiva
 
$$R = \sum_i R_i$$
- Pesi dei tempi di risposta
 
$$R = p_1 R_1 + p_2 R_2$$
- Probabilità catturate dalle visite
 
$$R_1 = f(D_1)$$

$$R_2 = f(D_2)$$

Impianti Informatici
POLITECNICO DI MILANO

Vediamo ora come affrontare il problema della determinazione dei tempi di risposta delle stazioni e del tempo di risposta medio della rete.

- Si noti che esiste un legame fra queste quantità, in quanto è sufficiente andare a calcolare la somma dei tempi di attraversamento dei singoli server, cioè dei loro tempi di risposta  $R_i$ , per ottenere il tempo di risposta della rete
- In altri termini, il tempo di risposta della rete è una misura additiva
- Si noti anche che questa espressione non è intuitiva, in quanto se ad esempio la rete è costituita da due server in parallelo, il tempo di risposta del job dipende da quale dei due server viene selezionato
- Ovvero, in media, si avrà  $R = p_1 R_1 + p_2 R_2$ , dove  $R_1$  ed  $R_2$  sono i tempi di risposta dei singoli server
- Ma come abbiamo detto nella scorsa lezione,  $p_1$  e  $p_2$  sono catturate implicitamente nelle visite e dunque sarà sufficiente esprimere  $R_1$  ed  $R_2$  in funzione, rispettivamente, di  $D_1$  e  $D_2$  invece che di  $S_1$  e  $S_2$  per tenere conto anche della struttura di rete

Calcolo del tempo di risposta dei server

5

$$R_i = D_i + W_i$$


$$W_i = D_i Q_i = D_i X R_i = D_i \lambda R_i$$

$$R_i = D_i + D_i \lambda R_i \quad R_i = \frac{D_i}{1 - \lambda D_i} = \frac{D_i}{1 - U_i}$$

Impianti Informatici
POLITECNICO DI MILANO

Andiamo quindi a dettagliare come si possa effettuare il calcolo dei tempi di risposta di ciascun server.

- Consideriamo quindi un server  $i$  con una domanda di servizio  $D$  con  $I$
- a cui giunga un flusso di richieste con tasso di arrivo  $\lambda$
- Il tempo di risposta sarà in generale dato dal tempo necessario alla richiesta per attraversare il server  $D_i$ , più il tempo speso ad aspettare che le richieste già in coda siano tutte servite, tempo che indichiamo come  $W$  con  $I$ . Dobbiamo quindi cercare di esprimere  $W_i$  in funzione dei soli  $D$  con  $i$  e  $\lambda$
- Come intuitivo, se in media una richiesta trova al suo arrivo  $Q_i$  richieste, ciascuna delle quali impiega  $D_i$  per essere servita, allora il tempo  $W_i$  sarà dato da  $D$  con  $i$  per  $Q$  con  $I$
- Ma dalla legge di Little abbiamo  $Q_i = X \cdot R_i$  e imponendo  $X$  e  $\lambda$ , otteniamo che  $W_i$  è uguale a  $D$  con  $I$  per  $\lambda R_i$  e reinserendo nella formula originale otteniamo  $R_i = D_i + D_i \lambda R_i$  e portando  $R_i$  a primo membro concludiamo che  $R_i = D_i / (1 - \lambda D_i)$  che per la legge dell'utilizzo può essere anche scritto come  $D_i / (1 - U_i)$ . Questa formula è sufficiente per calcolare il tempo di risposta del modello aperto dalla conoscenza delle  $D_i$  di tutti i server e del tasso di arrivo  $\lambda$  alla rete



Calcolo della lunghezza di coda

6

---

$$R_i = \frac{D_i}{1 - U_i} \quad Q_i = \lambda R_i = \frac{\lambda D_i}{1 - U_i} = \frac{U_i}{1 - U_i}$$

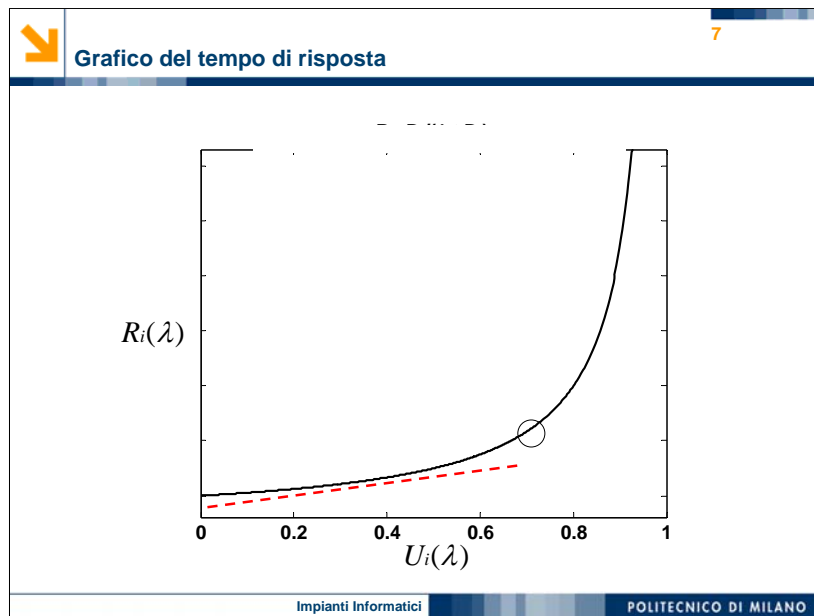
$$R(\lambda) = \sum_i R_i(\lambda) \quad R_i(\lambda) = \frac{D_i}{1 - U_i(\lambda)}$$

$$Q_i(\lambda) = \frac{U_i(\lambda)}{1 - U_i(\lambda)}$$

Impianti Informatici
POLITECNICO DI MILANO

La lunghezza di coda può immediatamente calcolata banalmente dalla legge di Little partendo

- dalla formula precedente.
- Poiché dalla legge di Little abbiamo  $Q_i = \lambda R_i$ , ricaviamo immediatamente  $Q_i = \lambda D_i / (1 - U_i)$  e raccogliendo anche in questo caso con la legge dell'utilizzo abbiamo
- $Q_i = U_i / (1 - U_i)$
- Ricapitolando, e esplicitando la dipendenza dei diversi termini dal tasso di arrivo  $\lambda$  della rete, le formule che consentono la risoluzione di un modello aperto sono
- Tempo di risposta di rete  $R(\lambda)$  uguale alla somma degli  $R_i(\lambda)$  di tutte le stazioni
- $R_i(\lambda) = D_i / (1 - U_i(\lambda))$  e
- $Q_i(\lambda) = U_i / (1 - U_i(\lambda))$



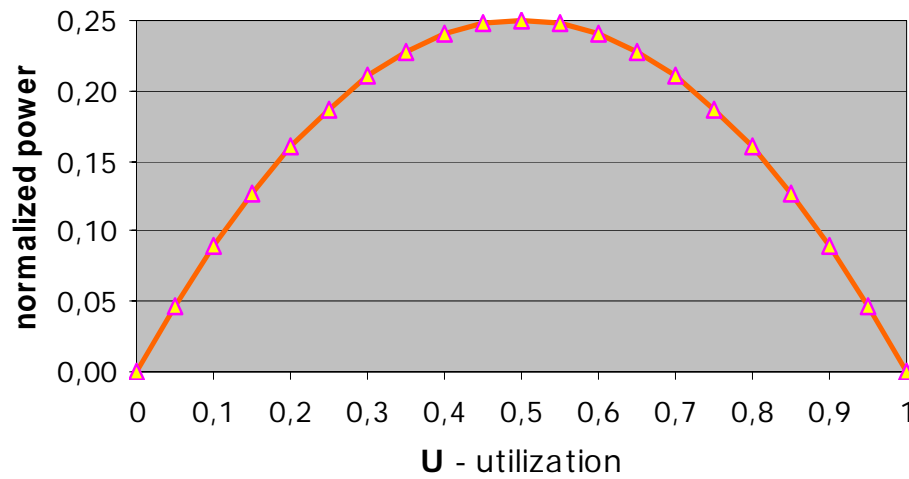
Grazie alle formule trovate possiamo disegnare l'andamento del tempo di risposta di una singolo server in funzione dell'utilizzo della stazione  $U_i(\lambda)$ .

- Quello che possiamo notare e' che inizialmente, all'aumentare dell'utilizzo, la curva subisce una crescita pressoché lineare. Questo vuol dire che al raddoppiare del traffico in ingresso, si ha inizialmente un aumento dei tempi di risposta proporzionale al nuovo carico arrivato
- Tuttavia, quello che possiamo vedere dalla parte destra del grafico e' che successivamente la stazione non riesce più a smaltire le richieste in coda e il tempo di risposta esplode in modo incontrollato, tendendo verso un asintoto a + infinito quando il denominatore  $1 - U_i$  del tempo di risposta va a zero. Questo avviene quando l'utilizzo della stazione  $U_i$  si avvicina ad 1, cioè al 100%, e questo effetto, come già visto nelle slide precedenti, e' detto saturazione del server i.
- Il punto in cui convenzionalmente si assume lo switching tra i due comportamenti e' quello in cui il server ha utilizzo pari al 75%. Dunque, componenti che abbiano utilizzi medi superiori a tale soglia sono da considerarsi un ostacolo alle prestazioni della rete.

Si tenga presente che il comportamento della lunghezza di coda e' analogo a quello del



$$\Phi \equiv \frac{\lambda}{R} = \frac{\lambda(1 - \lambda D)}{D} \quad \phi' = 0 - U = 50\%$$





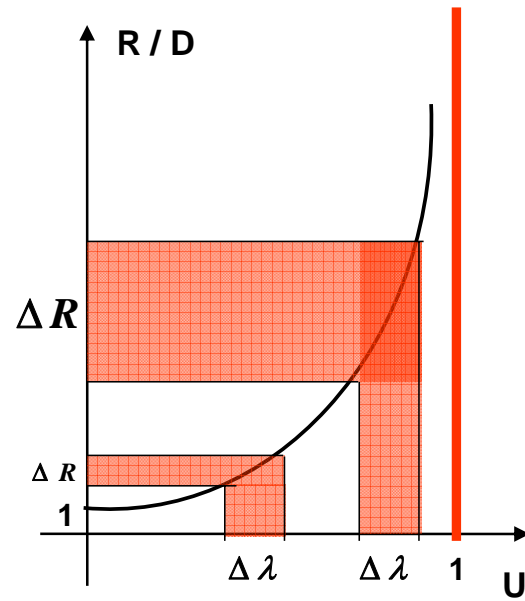


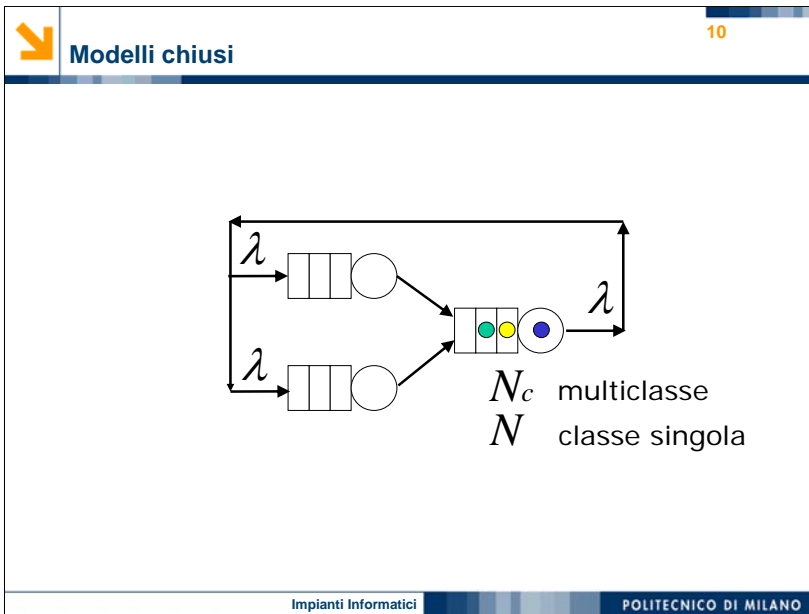
$$\Delta\lambda \Rightarrow \Delta R = \frac{D^2}{(1-U)^2} \Delta\lambda$$

$$U=0.2 \Rightarrow \Delta R = 1.56 D^2 \Delta\lambda$$

$$U=0.5 \Rightarrow \Delta R = 4 D^2 \Delta\lambda$$

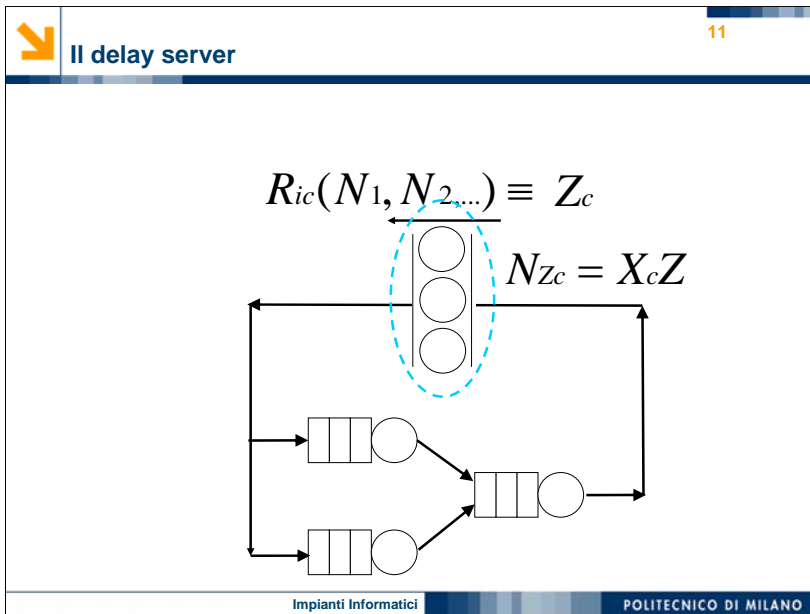
$$U=0.9 \Rightarrow \Delta R = 100 D^2 \Delta\lambda$$





In questa lezione consideriamo i modelli a reti di code chiusi.

- Come visto, la principale differenza rispetto ai modelli aperti è che nei modelli chiusi le richieste sono in numero costante e continuano a ciclare nella rete
- 
- Il numero di richieste nella rete è detto popolazione del modello
- e per una classe di richieste  $c$  viene indicato con la lettera  $N_c$
- o con la lettera  $N$  per i modelli a classe singola



- In alcuni casi le richieste impiegano alcuni secondi prima di compiere un nuovo ciclo nel modello. Per modellare questo effetto si inserisce un delay nella rete.
- un delay è uno speciale server dove il tempo di risposta ed il tempo di servizio coincidono sempre ovvero, dove nessuna richiesta si accoda in quanto viene sempre servita all'istante di arrivo.
- La domanda di servizio del delay è indicata con la lettera  $Z_c$  per una classe  $c$  ed è spesso detta think time delle richieste di classe  $c$
- Infine si noti che per la legge di Little il numero di richieste di classe  $c$  in attesa presso il delay è  $N_{zc} = X_c(N) * Z$



- Si applicano a **modelli chiusi a classe singola**
- $D_{max}$  = tempo di servizio della risorsa più lenta
- $D_{tot}$  = somma dei tempi di servizio di tutte le risorse
- $N$  = numero di utenti presenti nel sistema

