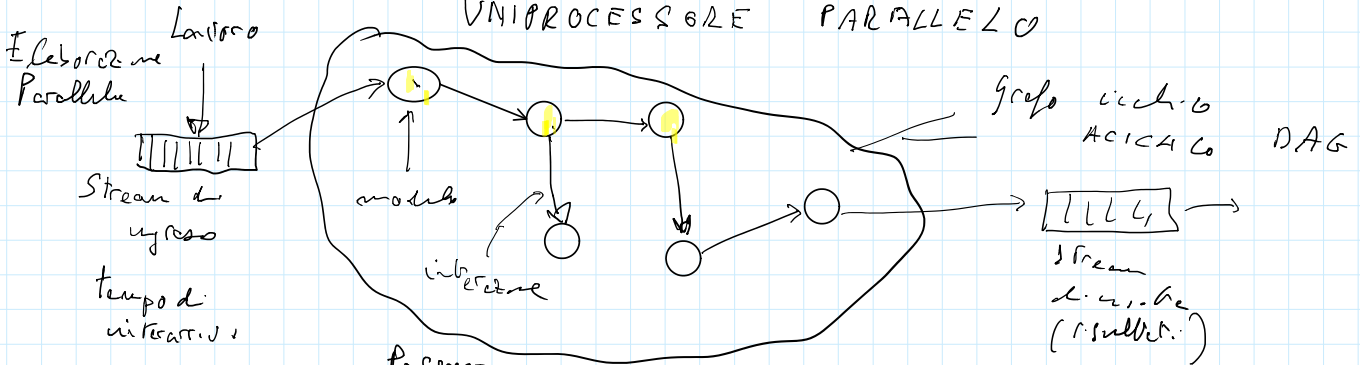


# ARCHITETTURE PARALLELE

## UNIPROCESSORE PARALLELO



possono operare in parallelo se i dati  $\neq$  della stream o uscite  
ideale = n° di moduli

Grado di parallelismo  
reale  $\equiv$  n° di moduli che operano contemp.  
po reale < po ideale  $\Rightarrow$  inefficiente

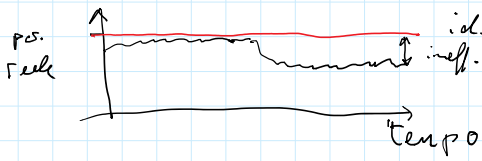
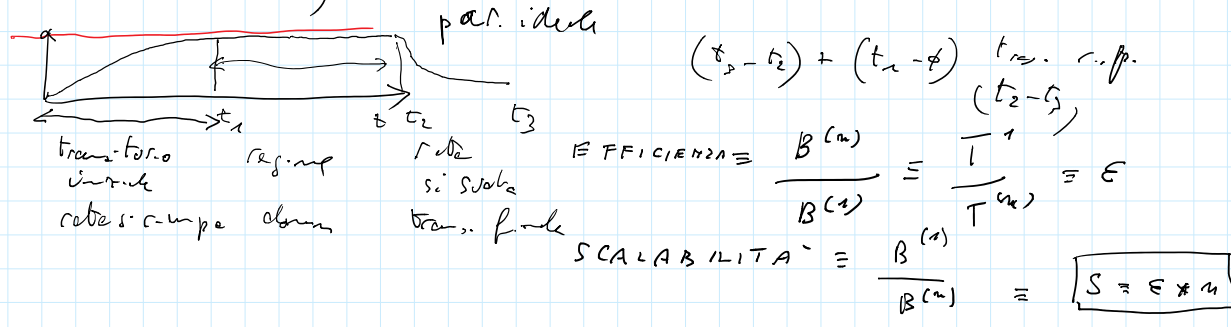
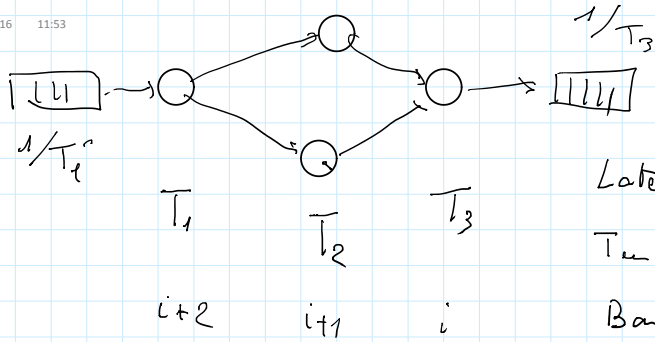


Fig. front-end  $\equiv$  Ret. diverse procedure lo stesso risultato (stesso output dello stesso input)

- Latenza  $\equiv$  tempo per elab. una richiesta nello stream di ingressi  $\begin{cases} \rightarrow \text{min} \\ \rightarrow \text{max} \\ \rightarrow \text{medio} \end{cases} \quad (n)$
- Tempo di servizio  $\equiv$  tempo tra due pacchetti nello stream di ingressi  $= T^{(n)}$   $n = \text{ps. id.}$
- Banda di elab.  $= 1/T^{(n)}$

- Tempo di compl. per  $m$  richieste,  $m \gg n$   $T_c^{(m)} = m * T^{(n)}$   
 (possiamo avere dei transistor.)





$$\text{Latency} = T_1 + T_2 + T_3 \quad (T_1 \geq T_2 \geq T_3)$$

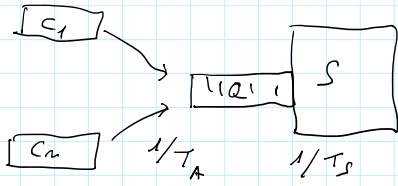
$$\text{Tempo di servizio} = T_1 \quad T^{(4)} = T_1$$

$$\text{Bandwidth} = \frac{1}{T_1}$$

$$B^{(4)} = T_1 + 2T_2 + T_3 \quad (\text{per hyp})$$

$$E = \frac{1/T_1}{1/(T_1 + T_2 + 2T_3)}$$

## Teoria delle Code



$\rho \equiv$  grado di congestione

Tempo di uscita

$$\equiv \frac{T_s}{T_A} \approx \frac{\lambda}{\mu}$$

$T_A$  se  $\rho < 1$

$T_s$  se  $\rho \geq 1$

Banda richiesta  $\approx \frac{1}{T_A}$

Banda fornita  $= \frac{1}{T_s}$

$c_1, \dots, c_n =$  clienti.

FIFO  $\equiv$   $\lambda$  serve i clienti nell'ordine di arrivo

tempo di servizio  $\equiv T_s \equiv$  tempo medio  
 $\sigma_s \equiv$  varianza

$$\mu = 1/T_s = \text{frequenza del servizio}$$

tempo di interarrivo  $\equiv T_A$

$$\left(1/T_A\right) = \text{frequenza di arrivo} = \lambda$$

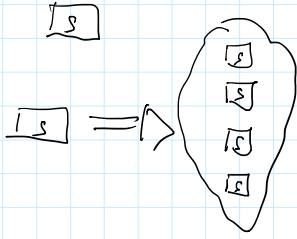
Banda fornita  $\geq$  Banda richiesta  $\Rightarrow$

$$\rho \leq 1$$

lunghezza coda media cresce con  $\rho$

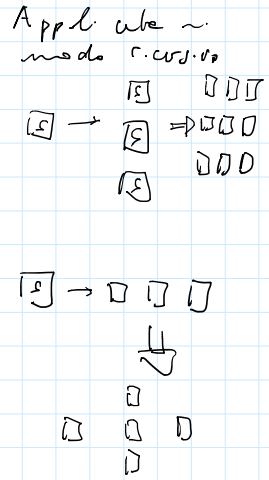
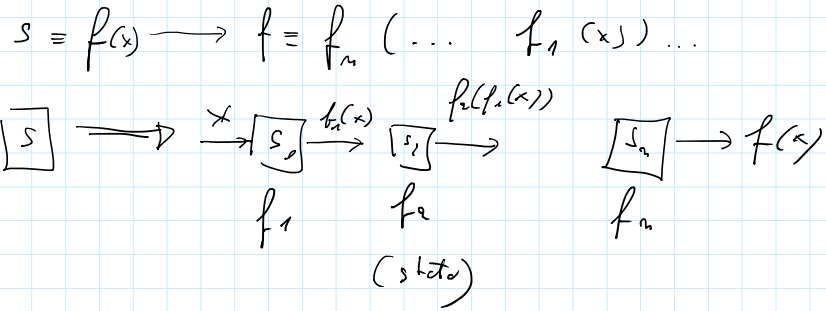
Forme di pipeline

1) replicazione



n copie dello stesso servite (stato??)

2) pipeline



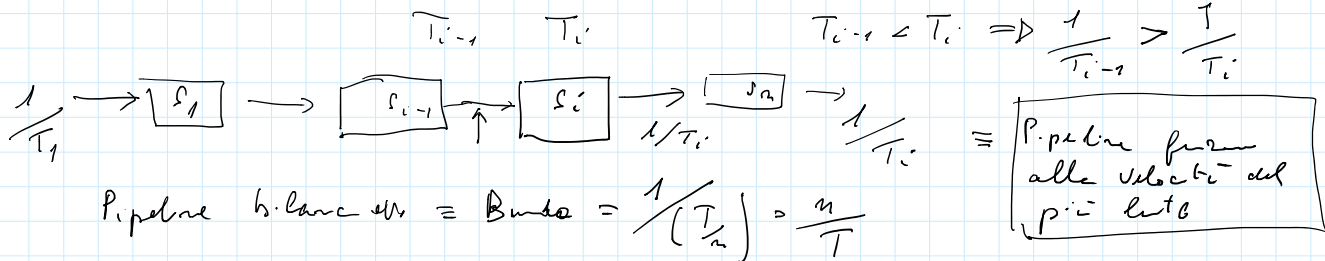
## Pipeline

Bilanciamento =  $f_1 \rightarrow T_1, f_2 \rightarrow T_2$ 

Nessuno è il più lento degli altri:

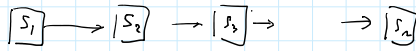
Bilanciato se  $T_1 = T_2 = \dots = T_n = T/n$   $T \equiv$  tempo di  $f$ 

$$T_i \gg T_1 = T_2 = T_{i-1} = T_{i+1} = \dots = T_n$$

Latenza pipeline =  $T$

Non bilanciato  
 $T_p = \max(T_1, \dots, T_n)$

Stadi del pipeline



Bilanciato:  $T/n = T_p$

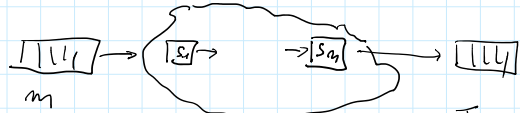
Tempo di riempimento del pipeline

$(n-1) \left( \frac{T}{n} \right)$

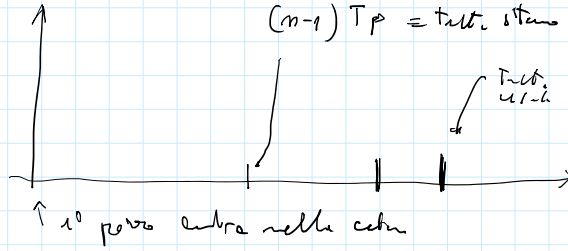
Dopo Transitorio produce valori in uscita

frequ. =  $\frac{n}{T}$ , Tempo  $\frac{T}{n}$  per elemento

m elementi nella coda in ingresso



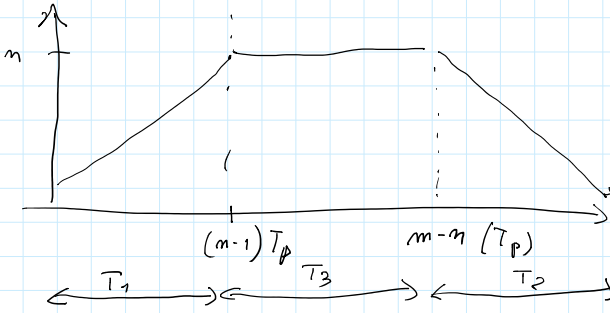
Tempo di completamento per m elem.



$T_p(n-1) \quad m \times T_p$   
 $m \times T_p + (n-1) T_p \approx n T_p$   
Regime + start      Trans. end

$m \gg n-1$

n° moduli  
alt. h



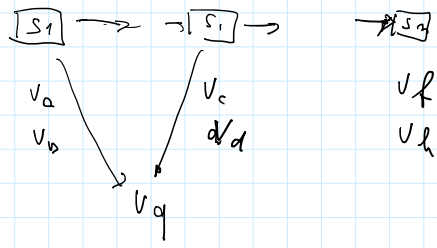
$T_1 + T_2 + T_3 \approx T_3$   
 $\Rightarrow \frac{T_1 + T_2}{T_3} \approx \phi$

$S \leftarrow \begin{matrix} v_1 \\ \vdots \\ v_n \end{matrix}$

$V_q$  condivisa =  $v_i - v_j$   
 sulle operazioni.  
 Si può lavorare su  $V_q$   
 solo se una cache rilascia il  
 blocco trattato da  $S_{i+k}$

Dipendenza

Si può operare  
 su una cache  
 solo se una cache  
 precedente è stata  
 elaborata da un'unità  
 successiva



Dipendenza tra stadi del pipeline  
 $\Rightarrow$  due stadi devono leggere/scrivere  
 lo stesso bit che  
 stato condiviso tra gli stadi.

- w-w write after write
- r-w read after write
- w-r write after read