

Algoritmi Paralleli e Distribuiti

Massimo Perego

Contents

Introduzione

All'interno del corso sarà presente solo la progettazione di algoritmi, non si parla di architetture.

Contesto: Generalmente parlando di algoritmi si considera un singolo esecutore, mentre con “paralleli e distribuiti” si aggiunge l'idea di avere una **pool di esecutori**, velocizzando possibilmente la risoluzione del problema, ma introducendo nuove problematiche.

Per un algoritmo non parallelo, quindi sequenziale, le **fasi** sono:

- **Progettazione:** scrittura dell'algoritmo, secondo diverse tecniche come *divide et impera* (dividere il problema in sottoproblemi e combinare le sotto-soluzioni per la soluzione globale, es mergesort), *programmazione dinamica* (utilizzare la memoria per memorizzare soluzioni parziali che saranno utilizzate per istanze più grandi, es calcolo Fibonacci), *greedy* (per problemi di ottimizzazione, c'è una funzione da massimizzare/minimizzare, il problema generale viene risolto scegliendo a ogni passo la soluzione localmente ottima, migliore per quel singolo step, es algoritmo di Kruskal), ...
- **Valutazione delle prestazioni:** in termini di tempo, spazio di memoria e in generale valutare l'algoritmo nella sua complessità
- **Codificazione:** fase di implementazione con un opportuno linguaggio di programmazione

Per gli algoritmi paralleli le problematiche sono simili, ma l'unica tecnica che può essere riciclata è quella di divide et impera.

Per gli algoritmi paralleli questi vengono valutati sempre in termini di tempo ma, al posto dello spazio, si conta il numero di processori.

Ma la pool di esecutori può avere diverse caratteristiche, ci sono due casi:

- **primo caso:** “una squadra in cui batte un solo cuore”, più processori sincroni
- **secondo caso:** “ogni membro un mondo a parte”, computer interfacciati ma ognuno ha il suo clock ed il suo tempo per risolvere le operazioni (es. diversi pc su internet), membri del pool lavorano in modo asincrono

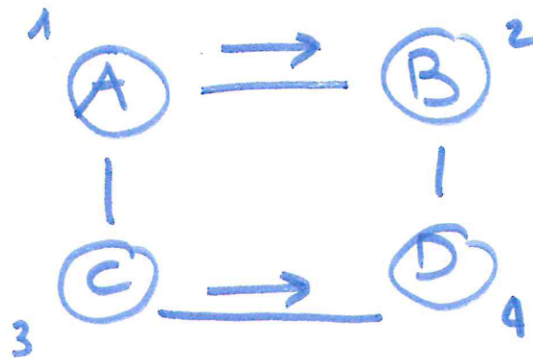
Primo caso: CPU Sincrone

I processori sono sincroni, quindi vanno all'unisono, si ha un **clock centrale** che scandisce il tempo per l'intero insieme di processori quindi le **istruzioni vengono eseguite tutte assieme** in un tempo determinato per tutti.

Può essere che condividano risorse, come ad esempio una memoria comune, oppure i processori possono essere solo collegati tra loro .

Con la **memoria condivisa** il sistema si chiama **PRAM**, l'altro si chiama **sistema parallelo a memoria distribuita** (manca quella condivisa).

Esempio:



Memoria distribuita, ogni dato è in un processore diverso, i dati sono distribuiti. Generalmente i collegamenti sono full duplex (avanti e indietro, ambo le direzioni).

Operazioni:

send(1,2)	send(3,4)
A+B	C+D
send(2,4)	
A+B+C+D	

Riga per riga:

- Il processore 1 invia A al 2, il 3 invia C al 4
- Nel processore 2 si fa la somma $A+B$, nel 4 si fa $C+D$
- Il processore 2 invia la somma al 4
- Il 4 farà il totale.

Ogni riga permette un'operazione, si tratta di un ciclo di clock, quindi per fare questa somma in totale sono 4 cicli di clock.

In genere un solo processore contiene la risposta, solitamente quello con l'indice maggiore (ultimo, nel caso dei processori 1, 2, 3 e 4 il 4 contiene la risposta).

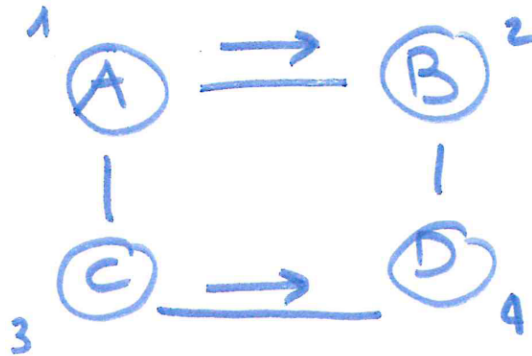
Con questo esempio non si capisce bene il guadagno dato dal parallelismo rispetto all'algoritmo sequenziale, ma si può generalizzare l'idea su un'istanza di lunghezza n , con un numero di processori adeguato (ma anche in caso di un numero di processori ridotti l'algoritmo si può adeguare alla situazione). In questo modo si può notare che il numero di passi scala in modo logaritmico con la dimensione n (prima si sommano 2 a 2, poi 4 a 4, poi 8 a 8, ...). Si passa da n passaggi con il sequenziale a $\log n$ con il parallelo.

Secondo caso: CPU Asincrone

In questa casistica **ogni processore ha i propri tempi** nonostante siano **collegati tra loro**, e ovviamente non si può assumere memoria condivisa. Non avere un clock centrale porta ad avere algoritmi distribuiti e NON paralleli come nel caso precedente.

Non si può assumere **nessun tipo di sincronizzazione gratuita**, devono **scambiarsi dei messaggi** (es: attraverso internet, se non mi viene detto nulla non so niente degli altri processori).

Esempio:



In questo caso potrebbero esserci **tempi** più lunghi e che **variano da collegamento a collegamento**, anche tra scambi diversi sullo stesso collegamento (come può accadere su internet).

Avendo velocità diverse, **un processore potrebbe dover aspettare** che gli altri terminino di elaborare/scambiare i messaggi. Serve scambiare messaggi per il coordinamento.

Teoria e realtà: motivazioni

Si studiano gli algoritmi a livello teorico, basandosi sulla realtà, esistono architetture distribuite per l'implementazione di tali algoritmi. La nascita di architetture parallele e distribuite comincia dalla seconda metà degli anni '60, di conseguenza si è sviluppata la teoria.

Esempi di architetture parallele:

1. **Supercomputer:** cluster di processori con incredibili prestazioni di calcolo (CDC, CRAY), solitamente a scopo di simulare sistemi complessi (fisici, militari, ...; tutto ciò che richiede un grande numero di operazioni aritmetiche). Di solito sono a memoria distribuita in quanto sarebbe difficile realizzarli a memoria condivisa dato che tutti i processori dovrebbero essere in grado di accedervi
2. **GPU:** più processori collegati tra loro, in genere dedicate a problemi di grafica ma permettono di risolvere bene problemi di algebra lineare quindi viene usata in senso più generale
3. **Multicore processor:** più unità di calcolo all'interno dello stesso processore, per migliorare le prestazioni ottimizzando l'assorbimento di energia (non bisogna aumentare il clock che porterebbe a un aumento dell'assorbimento) e relativi problemi di raffreddamento (meno energia meno calore); sistemi PRAM
4. **Circuiti integrati:** sistemi di calcolo formati da gate opportunamente connessi; sostanzialmente una struttura a livelli che lavorano in parallelo prendendo input dal livello precedente

Esempi di architetture distribuite:

1. **Reti di calcolatori:** Internet (dagli anni '60), connettono dispositivi in molti diversi, quindi servono protocolli di comunicazione come TCP/IP
2. **Reti mobili:** la topologia di connessione dei dispositivi cambia continuamente (gli smartphone si muovono)
3. **Reti di sensori:** dispositivi con capacità limitate, per la maggior parte del tempo in "stand by", eseguono solo quello per cui sono costruite; solitamente hanno scopo di monitoraggio degli ambienti; necessitano di meccanismi di "wake up" (a un certo punto devo svegliarli), "acknowledge" (devono dire quando finiscono) e "recovery" (memorizzare le proprie informazioni/quelle ricevute)

Negli **algoritmi paralleli** il fattore rilevante è il **tempo**. Nel sommare 4 numeri sono 4 passi, per sommarne 1000 ne servono 10. Per vedere se un algoritmo parallelo è efficiente si valuta il tempo.

Per gli **algoritmi distribuiti** quello che conta è il **coordinamento**, ovvero il **numero di messaggi da scambiare** (meno messaggi = più veloce). Il numero di messaggi definisce una sorta di tempo e permette anche di avere un'idea della congestione della rete.

Problematiche parallele e distribuite

La **progettazione** richiede **nuove idee**, mentre la il processo di **valutazione** delle prestazioni richiede **nuove misure**.

Perché nuove idee? Un algoritmo sequenziale efficiente non necessariamente porta a un algoritmo parallelo e/o distribuito efficiente, di conseguenza servono tecniche ad hoc per entrambi le tipologie.

Inoltre buoni algoritmi paralleli non sempre portano a buoni algoritmi distribuiti e viceversa. Ad esempio, nel **distribuito serve gestire problemi di ritardo e comunicazione**, che nel parallelo non si presentano.

Le misure che verranno definite non faranno riferimento a una architettura in particolare ma saranno definite solo in maniera teorica.

Per ogni paradigma si dovrà definire:

1. **Modello teorico**
2. Come **valutare le performance** degli algoritmi
3. **Semplici problemi** per apprendere le tecniche

Nel caso di **architetture parallele**:

A. Primo caso:

1. **PRAM** (comunicazione immediata): memoria condivisa
2. Risorse di calcolo: tempo, hardware (inteso come numero di processori)
3. Problemi: sommatoria, somme prefisse, ordinamento

B. Secondo caso

1. Modello a **memoria distribuita**:

- array lineari: n processori collegati all'interno di un array lineare, ovvero P_i collegato a P_{i+1} , il quale è collegato a P_{i+2} e così via
 - mesh: versione bidimensionale dell'array lineare
 - albero: albero binario completo di n foglie
 - ipercubo: per $n = 2^d$, un ipercubo (o d -cubo) è un grafo i cui vertici sono elementi $c_1, \dots, c_d \in \{0, 1\}^d$ e due vertici x e y sono estremi di un lato se le parole x e y differiscono in una sola posizione, quindi P_x e P_y sono collegati se x e y differiscono di una sola posizione. Sostanzialmente, un cubo ha 3 dimensioni ed ogni vertice è collegato ad altri 3, il d -cubo è la stessa cosa generalizzata in d dimensioni, con un processore su ogni vertice. Inoltre è una struttura ricorsiva, un $d + 1$ -cubo è formato da 2 d -cubi collegati tra loro, dove uno rappresenta la parte del $d + 1$ -cubo con bit più significativo a 0, mentre l'altro la parte con bit più significativo ad 1
2. Risorse di calcolo: tempo, hardware (processori non connessi direttamente comunicano più lentamente)
 3. Problemi: shuffle, max, ordinamento

Nel caso di **architetture parallele**:

1. Definizione del modello astratto
2. Risorse di calcolo: tempo, numero di messaggi (troppi = congestione)
3. Problemi: broadcast, wake up, traversal, spanning tree, election, routing

Il Tempo

Gli algoritmi vanno **valutati** (anche) **sul tempo**, in quanto sia nel caso di algoritmi paralleli che di algoritmi distribuiti, la risorsa tempo è cruciale.

Come nel caso sequenziale, la **definizione formale** è:

- $T(x)$ = numero di operazioni elementari su input x (istanza).
- $t(n) = \max\{T(x) | x \in \Sigma^n\}$, il massimo valore di T tra tutti gli input di dimensione n , ovvero il caso peggiore tra tutte le istanze di dimensione n . Si tratta quindi di una funzione in n , dove n è la dimensione dell'input.

Spesso non saremo interessati a una valutazione precisa ma a un **tasso di crescita della funzione**: quindi si usano le funzioni asintotiche O, Ω, Θ .

Definizioni (non formali, già viste troppe volte): siano $f, g : \mathbb{N} \rightarrow \mathbb{N}$ due funzioni definite sui numeri naturali

- $f(n) = O(g(n))$, la funzione è limitata superiormente (da un certo punto in poi, non per forza da subito) dalla funzione $c \cdot g(n)$, per un valore di c
- $f(n) = \Omega(n)$, la funzione è limitata inferiormente (da un certo punto in poi) da $c \cdot g(n)$, per un valore di c
- $f(n) = \Theta(n)$, unione delle due sopra, limite stretto, ovvero la funzione è compresa tra $c_1 \cdot g(n) \leq f(n) \leq c_2 \cdot g(n)$

Tutte valgono asintoticamente, quindi da un certo punto in poi.

Da qui in poi sarà dato per scontato il concetto di MdT (chiamata Deterministic Turing Machine DTM), non lo riporterò, è già su almeno altri 2 blocchi di appunti e mi sono rotto il cazzo di scriverlo, comunque testina, quintupla, bla bla bla.

Valutazione di $t(n)$:

1. Viene **valutata su** un particolare **modello di calcolo**, quando calcolo la funzione vado a contare il numero di operazioni
2. Va scelto il **criterio di costo: uniforme o logaritmico**

Le **operazioni da contare** sono quelle **primitive messe a disposizione dal modello di calcolo**.

Esempio: palindromi

- Input: $x \in \{0, 1\}^*$
- Output: x è palindroma

Soluzione:

```
for(i=1, j=|x|; i<j && x_i = x_j; i++, j--);  
return i >= j;
```

Two pointers, controlla se ogni lettera è uguale a quella “dall’altro lato”. Se si blocca (quindi non è palindroma) gli indici rimarranno tali che $i < j$ (e di conseguenza restituirà falso alla fine), altrimenti i supererà j , dato che se arrivano a metà il controllo sarà valido fino alla fine portando gli indici a essere invertiti rispetto alla partenza.

Tale programma:

- su RAM (memoria ad accesso casuale) ha

$$t(n) = O(n)$$

(basta leggere una volta ogni lettera, posso accedere dove voglio nella memoria)

- Su DTM (una sola testina di lettura) ha

$$t(n) = O(n^2)$$

(devo fare avanti e indietro a ogni lettera che controllo)

Inoltre bisogna **fare attenzione alla dimensione dei dati** in gioco. Esempio: Fattoriale

- Input: $n \in \mathbb{N}$
- Output $n!$

Soluzione:

```
k = 1;
for(i = 1; i <= x; k = k*i, i++);
return k;
```

Effettua semplicemente il calcolo del fattoriale.

Tale programma:

- Contando la moltiplicazione come operazione elementare di una RAM

$$t(\log n) = n$$

($\log n$ dato che n è un numero naturale che si scrive in $\log n$ bit)

- su DTM, dovendo scrivere il binario il risultato $n! \simeq n^n$ abbiamo

$$t(\log n) \geq \log n^n = n \log n$$

Criteri di costo:

- **Uniforme:** le operazioni elementari richiedono una unità di tempo
- **Logaritmico:** ogni operazione elementare ha un costo che dipende dal numero di bit degli operandi

Teoria della Complessità

Possibili tassi di crescita di $t(n)$. Dobbiamo **stabilire quando qualcosa è efficiente per algoritmi paralleli e quando per i sequenziali**.

La funzione $t(n)$ si dirà:

- logaritmica quando è $O(\log n)$
- polilogaritmica quando è $O(\log^k n)$ per una costante $k > 0$
- lineare quando è $O(n)$
- polinomiale quando è $O(n^k)$ per una costante $k > 0$
- esponenziale quando NON è $O(n^k)$ per ogni costante $k > 0$

Efficienza nel caso sequenziale

Definizione formale: un problema è efficiente se ammette un **algoritmo efficiente**, ovvero viene **risolto** da una DTM in **tempo polinomiale**.

Si dicono:

- \mathcal{P} = classe dei problemi di decisione risolti in efficiente tempo
- \mathcal{FP} = classe dei problemi generali risolti in efficiente tempo

Oppure

- \mathcal{P} = classe dei problemi di decisione risolti in tempo polinomiale
- \mathcal{FP} = classe dei problemi generali risolti in tempo polinomiale

La differenza tra \mathcal{P} e \mathcal{FP} è che i problemi di decisione richiedono solo *s* o *no*, che sostanzialmente è il calcolo di un solo bit, mentre gli altri richiedono di calcolare più bit.

Perché consideriamo il tempo polinomiale come efficiente? Esempio: considerando un processore da 4 GHz (\neq da una DTM) su un'istanza di 4000 caratteri (mezza pagina circa)

complessità $t(n)$	tempo di attesa
n	$1\mu s$
n^2	$4ms$
2^n	> 1 secolo
3^n	~ 2 secoli e mezzo

Come si può vedere i tempi esponenziali esplodono anche su istanze piccole.

Quindi: \mathcal{P} e \mathcal{FP} hanno una definizione robusta che non cambia con il modello di calcolo (non importa se sia su RAM o DTM).

Inoltre \mathcal{P} e \mathcal{FP} contengono molti problemi fondamentali, come ordinamento, prodotto di matrici, alcuni problemi di ottimizzazione, ...

Purtroppo alcuni problemi interessanti non hanno ancora soluzioni efficienti e stanno in

- \mathcal{NP} = problemi di decisione risolti in tempi polinomiale su una DTM non deterministica

Si congettura $\mathcal{P} \neq \mathcal{NP}$, ma non è ancora noto (non siamo certi che i problemi attualmente non risolvibili in tempo polinomiale non abbiano una soluzione efficiente).

I problemi che possono essere parallelizzati non possono appartenere ad \mathcal{NP} , altrimenti sarebbe come provare che $\mathcal{P} = \mathcal{NP}$.

All'atto pratico siamo interessati a una **sottoclasse** di \mathcal{P}

1. La **valutazione asintotica** dei tempi può **nascondere costanti o termini** che all'atto pratico fanno la differenza. Esempio: a volte si preferisce il QuickSort ($O(n^2)$) rispetto al MergeSort ($O(n \log n)$).
2. Da un punto di vista pratico il **grado del polinomio deve essere basso**. Es: n^{1000} anche su piccole istanze risulta peggiore di un esponenziale. Esempio pratico: in $O(n^6)$ è il tempo per un test di primalità, ma il grado è troppo alto, quindi si preferisce un test probabilistico con possibilità di errore.

1 Algoritmi Paralleli

Possono essere affrontate 3 diverse problematiche:

- **Sintesi:** progettazione degli algoritmi
- **Valutazione:** delle prestazioni degli algoritmi
- **Universalità:** quali problemi sono efficientemente risolubili e quali no

Sintesi: Come costruire algoritmi paralleli? Ci si può ispirare ad algoritmi sequenziali, ma quasi mai si può riciclare l'intero algoritmo.

Esempio: sommare due numeri binari. L'algoritmo sequenziale ottimo somma i bit per colonna.

Porta a un buon algoritmo parallelo? Ogni processore si può occupare di una colonna.

Sorge un problema con il riporto, il quale sostanzialmente fa tornare l'algoritmo a uno sequenziale, con in più uno spreco ingiustificato di hardware (ogni processore deve aspettare quello prima...).

Ma esistono algoritmi paralleli efficienti per sommare numeri binari.

Esempio 2: ordinamento. L'algoritmo sequenziale ottimo è il MergeSort.

Ma questo non porta a un buon algoritmo parallelo, mentre esistono algoritmi paralleli efficienti per l'ordinamento.

Servono nuove tecniche e punti di vista totalmente diversi da quelli usati per gli algoritmi sequenziali.

Valutazione delle prestazioni: Come si misurano le prestazioni degli algoritmi paralleli? Vanno definiti:

- **Tempo impiegato**
- **Hardware occupato** = numero di processori

Cos'è efficiente nel caso di **algoritmi paralleli**?

Si ha un **parametro** $E = \text{efficienza}$. Tiene conto sia del tempo, sia del numero di processori, permettendoci di dire se l'uso aggiuntivo di hardware è giustificato o meno.

Universalità: Si tratta di una problematica teorica, riesco a caratterizzare la classe dei problemi che ammettono algoritmi paralleli efficienti?

Nel caso sequenziale la classe di problemi si chiama \mathcal{FP} . Nel caso parallelo si **definisce la classe** \mathcal{NC} .

\mathcal{NC} = coincide con problemi risolti da algoritmi paralleli efficienti, ovvero con:

- **Tempo polilogaritmico** (quindi meglio che polinomiale, altrimenti userei un algoritmo sequenziale efficiente, dato anch'esso sarebbe polinomiale).
- **Hardware polinomiale** (richiede un numero di processori che scala polinomialmente con la lunghezza dell'input).

Esempi di problemi in \mathcal{NC} :

- somma di due numeri
- ordinamento

In realtà è ragionevole una sottoclasse di \mathcal{NC} , come per \mathcal{P} , se il grado del polinomio o le costanti nascoste sono alte, il tempo è sì polinomiale ma rimane poco pratico.

Si può **dimostrare che** $\mathcal{NC} \subseteq \mathcal{FP}$.

I problemi che possiamo parallelizzare efficientemente rientrano nella classe di problemi risolvibili con algoritmi di calcolo sequenziali efficienti.

Per dimostrarlo: per ottenere un algoritmo sequenziale a partire da uno parallelo basta eseguire le istruzioni, che nel parallelo sono divise su vari processori, in fila su un processore singolo, mantenendo il tempo polinomiale (tempo polilogaritmico moltiplicato per il numero polinomiale dei processori, quindi dominato ancora da un polinomiale) e ottenendo un algoritmo sequenziale.

Ci si può chiedere se $\mathcal{NC} = \mathcal{FP}$, ma è un problema aperto.

Si pensa di no, la questione riprende il problema “ogni algoritmo sequenziale efficiente è parallelizzabile?” Esiste un metodo per trasformare ogni algoritmo sequenziale in un algoritmo parallelo?

Se $\mathcal{NC} \neq \mathcal{FP} \implies \mathcal{FP} \setminus \mathcal{NC} =$ problemi \mathcal{P} -completi; esistono problemi in \mathcal{P} ma non in \mathcal{NC} , chiamati problemi \mathcal{P} -completi, problemi tali che la loro appartenenza a \mathcal{NC} farebbe collassare \mathcal{NC} e \mathcal{FP} (ovvero $\mathcal{NC} = \mathcal{FP}$; un problema \mathcal{P} -completo permetterebbe di risolvere tutti gli altri problemi di \mathcal{FP} , facendo coincidere le due classi).

Pertanto non sarà possibile parallelizzare i problemi \mathcal{P} -completi (duh).

L'universalità non sarà tratta ulteriormente in questo corso.

Nel mondo parallelo bisogna:

- fissare un modello di calcolo teorico da mappare su architetture reali
- fissare le risorse computazionali, ovvero tempo e numero di processori
- Fissare il parametro efficienza, per dire se una soluzione parallela è preferibile ad una efficiente soluzione sequenziale

1.1 Tipi di architetture parallele

Ci sono due tipologie:

- a **memoria condivisa**
- a **memoria distribuita**

Memoria condivisa: Ci sono dei processori che si affacciano su una **memoria centrale condivisa tra tutti**. Si ha un **clock centrale**.

La **comunicazione avviene in tempo costante**: avendo due processori P_i e P_j , P_j vuole comunicare un dato x a P_i

- P_j scrive x in memoria centrale
- P_i legge x dalla memoria centrale

In questo modo in 2 cicli di clock avviene la comunicazione. Si dice “comunicazione costante” perché non importa la dimensione del dato da trasferire. Questa è l’attuale architettura di CPU multicore.

Questa architettura permette una **forte parallelizzazione**, dato che i tempi di comunicazione sono trascurabili, rimuovendo quindi il costo dovuto alla comunicazione (o comunque se considerato, è minimo).

Memoria distribuita: Da non confondere con “architetture distribuite”, questa è un’architettura parallela a memoria distribuita.

Si ha sempre un unico **clock centrale**, processori sincroni, ma **ognuno di questi ha una memoria privata** e sono **collegati tra loro** tramite fili (con qualche topologia, solitamente non fully connected), si ha una **rete di interconnessione tra i processori**.

La comunicazione dipende dalla “distanza” tra i processori: se P_j vuole comunicare x a P_i dobbiamo chiederci quanti processori collegano P_j a P_i . Se sono collegati direttamente allora il tempo è costante.

La rete di interconnessione è fissa, può essere:

- **array lineare:** processori collegati su una linea
- **mesh:** collegati su una matrice, sia per riga che per colonna
- **ipercubo:** come prima, ma aumentando le dimensioni

1.1.1 PRAM

Parallel RAM, sta ad indicare che si hanno **RAM (Random Access Machine) che lavorano in parallelo**, si tratta di un modello teorico ma in pratica lo si ritrova nelle CPU multicore.

Si ha un **determinato numero di processori** ed una **memoria centrale** M , la quale i – lesima cella viene indicata con $M[i]$. All’inizio contiene l’input, alla fine può contenere l’output (che può anche essere in uno dei processori).

Ogni processore è una RAM (= Random Access Machine) sequenziale e mettendole assieme queste possono lavorare in parallelo.

Ogni processore è costituito da una unità di calcolo e dei registri, i quali costituiscono la memoria privata del processore, chiamati $R[0], R[1], \dots$

Tipi di istruzioni dei processori P_i :

- Operazioni aritmetico/logiche.
- Istruzioni da/per la memoria centrale
 - **store** $R[k] \ M[h]$: mette il dato nella memoria centrale
 - **load** $R[k'] \ M[h']$: carico nel registro dati dalla memoria centrale alla memoria privata

Un processore non può usare i dati x e y se questi sono in memoria centrale, vanno prima caricati nella memoria privata (la ALU è collegata solo ai registri). **Ogni processore può effettuare operazioni solo sulla propria memoria privata.**

Esempio P_j vuole comunicare a P_i il contenuto di $R[t]$. Istruzioni:

$$\begin{array}{ll} P_j : & \text{store} \quad R[t] \ M[x] \\ P_i : & \text{load} \quad R[s] \ M[x] \end{array}$$

Ovviamente sia **store** che **load** devono lavorare sulla stessa cella $M[x]$ della memoria centrale.

La comunicazione avviene in tempo $O(1)$.

Come si programma con una PRAM? Si ha un **unico programma** per tutti i processori ed una **memoria condivisa**:

1. Il tempo per ogni P_i è scandito dal clock centrale
2. Ogni P_i esegue la “stessa istruzione”

Cosa si intende con “la stessa istruzione”? Forma dell’istruzione:

```
for all  $i \in I$  pardo  
    istruzione $_i$ 
```

Con I definito come l’insieme di indici dei processori che devono eseguire quella istruzione. Si ha un **for** parallelo, quindi tempo costante e non n passi diversi, ogni istruzione è indicizzata con l’indice del processore.

I processori con indice in I eseguono **istruzione $_i$** . I processori con indice NON in I eseguono l’istruzione nulla.

Cosa si intende con “istruzione $_i$ ”? Dipende dal tipo di architettura:

- **SIMD**: Single Instruction Multiple Data
- **MIMD**: Multiple Instruction Multiple Data

SIMD: L’istruzione eseguita è la **stessa per ogni processore** $\in I$, ma ognuno di questi la effettua **su dati diversi**.

In funzione della **capacità di accedere alla memoria** M si possono avere più architetture:

1. **EREW**: Exclusive Read Exclusive Write, scrittura e lettura sono esclusivi, quindi una stessa cella di memoria può essere letta o scritta da un solo processore per volta
2. **CREW**: Concurrent Read Exclusive Write, la lettura può essere simultanea, la scrittura è esclusiva
3. **CRCW**: Concurrent Read Concurrent Write, lettura e scrittura possono essere simultanee

Per la **scrittura simultanea servono politiche di accesso** (se più processori vogliono scrivere sulla stessa cella, chi vince?):

- **common:** i processori possono scrivere solo lo stesso dato, pena l'arresto del sistema
- **random:** si sceglie un P_i a caso
- **max/min:** vince il P_i con il dato max/min
- **priority:** vince il P_i con priorità massima (più alta, va definita una priorità)

Indicando con i numeri 1,2,3 le architetture come sopra

$$Alg(1) \implies Alg(2) \implies Alg(3)$$

Ovvero un algoritmo EREW funziona su architetture CREW, ed un CREW funziona su architetture di tipo CRCW. Questo è abbastanza ovvio, ma vale il contrario? Di base no, ma si possono trasformare CRCW in CREW ed eventualmente anche EREW.

L'architettura più ragionevole e più semplice è EREW.

Esempio di algoritmo su PRAM di tipo CREW:

- numero di processori = n
- dimensione input = n (lunghezza dell'input pari al numero di processori)
- assumiamo input array A con valori distinti

Volendo cercare in A un valore particolare:

```
Cerca (A, n, x) {  
    indice=-1  
    for i=0 to n-1 pardo  
         $P_i$ :  
            if A[i]=x then  
                indice=i  
    return indice  
}
```

Ogni processore di carica nella memoria locale il un dato di A e lo confronta con il valore x cercato (il processore 1 controlla la cella 1, il 2 la cella 2, ...). La variabile **indice** parte $=-1$, ma viene settato quando trovato il valore corretto, restituito alla fine.

Si ha tempo parallelo e costante. Quanti clock sono? Vanno effettuati: `load x`, `load A[i]`, il confronto ed eventualmente la scrittura della variabile `indice`, quindi 4 cicli di clock.

Questo algoritmo è concurrent read perché la x viene letta assieme da tutti, mentre su A vengono lette sempre celle diverse. Exclusive write perché la scrittura invece viene effettuata da un solo processore, solo quando viene scritto l'indice, con il presupposto che A contenga solo valori distinti.

Se A può contenere elementi ripetuti allora questo algoritmo diventa CRCW, dato che posso trovare più valori pari a x tutti nello stesso momento.

Risorse di calcolo:

- nel **sequenziale**: $t(n)$ (tempo), $s(n)$ (spazio)
- nel **parallelo**: $p(n)$ (numero di processori), $T(n, p(n))$ (tempo, in funzione di input e numero di processori)

Definizioni informali:

- $p(n)$: numero di processori richiesti su input di lunghezza n (nel caso peggiore)
- $T(n, p(n))$: tempo richiesto da un input di lunghezza n su $p(n)$ processori. Di conseguenza $T(n, 1)$ è il tempo sequenziale

Valutazione precisa del tempo $T(n, p(n))$: Possiamo pensare all'esecuzione del programma con una sorta di matrice con:

- colonne che rappresentano ognuna uno dei $p(n)$ processori
- righe che rappresentano ognuna un passo parallelo

Quindi nella prima riga ho il primo passo parallelo e sotto la prima colonna avrò l'istruzione eseguita dal primo processore.

Il numero dei passi generalmente dipende da n , per cui diventa una funzione $k(n)$.

Il costo, in termini di tempo, per il primo processore al primo step del programma viene indicato con $t_1^{(1)}(n)$, dove la t è piccola perché è un tempo sequenziale, all'apice si trova il numero del processore e al pedice il numero del passaggio.

Vogliamo **definire $T(n, p(n))$ in funzione** di questi **tempi richiesti dalle singole RAM**, ovvero la **somma dei tempi richiesti da ognuno dei $k(n)$ passi**. Supponendo un'architettura MIMD, istruzioni diverse potrebbero richiedere tempo diverso, quindi a **ogni passo dobbiamo considerare il tempo massimo** impiegato da un singolo processore; il massimo per ogni riga, per ogni t_1 .

Indichiamo il tempo impiegato al passo i -esimo con

$$t_i(n) = \max \left\{ t_i^{(j)}(n) \mid 1 \leq j \leq p(n) \right\}$$

Di conseguenza:

$$T(n, p(n)) = \sum_{i=1}^{k(n)} t_i(n)$$

Il **tempo parallelo** è la somma del tempo a ogni passaggio:

- **T dipende da $k(n)$** (devo sommare $k(n)$ passi)
- **T dipende anche dalla dimensione dell'input** (costo logaritmico/uniforme, il tempo sequenziale a ogni passo dipende dalla dimensione dell'input in quanto elaborato da macchine sequenziali)
- **T dipende da $p(n)$** (le operazioni da eseguire sono sempre quelle, se i processori sono troppo pochi devo assegnare le operazioni rimanenti in coda ad altre su altri processori; eventualmente alcuni processori in un solo step dovranno svolgere più operazioni)

1.1.2 Broadcast P-RAM

Al posto di processori consideriamo entità, che collaborano tra loro ed una informazione i , detenuta da una sola entità, deve essere conosciuta da tutte le altre entità che compongono il sistema. Broadcast o replica, si ha un array in memoria centrale/condivisa in cui **una cella contiene il valore x** e vogliamo che **nelle celle seguenti compaia il valore x per un certo numero di elementi n** , in modo che gli altri processori possano conoscere l'informazione nella loro cella associata. Le celle della memoria condivisa vengono indicate come A .

```

Broadcast ( $x$ )
{
 $P_0$  :       $A[0] = x$ 
          for  $i = 0$  to  $\log n - 1$  do
              for  $j = 2^i$  to  $2^{i+1} - 1$  par do
 $P_j$  :       $A[j] = A[j - 2^i]$ 
          }

```

Per la condivisione ci sono $\log n$ passi paralleli, un certo numero di processori prendono il valore x e lo copiano in una delle celle successive (ogni processore lo passa ad un'altra cella, raddoppiando il numero di celle contenenti x ogni volta e ripetendo), rendendo il numero di passaggi richiesti logaritmico rispetto al numero n di copie dei dati da avere.

Questo processo potrebbe essere eseguito perché in un programma EREW ogni processore deve accedere ad una sua copia del dato non avendo letture concorrenti.

Il broadcast può essere usato per risolvere **la ricerca**:

```

Cerca (a, n, x)
{
    Broadcast(x)
    Indice=-1
    for i = 0 to n-1 par do
Pi:           if A[i]=x[i] then indice=i
    return indice
}
```

Faccio il broadcast del dato x e lo faccio cercare da tutti i processori, restituendo l'indice i del dato, se trovato. Di conseguenza impiega tempo $\log n$ (solamente 1 passo per la ricerca, $\log n$ per il broadcast).

Se tutti gli indici sono distinti l'algoritmo è EREW, altrimenti ERCW (ma solitamente si usa solo CRCW, non ha senso avere scritture concorrenti e letture esclusive). ERCW può essere trasformato in EREW (con un aumento della funzione tempo).

1.2 Efficienza

Per progettare algoritmi paralleli è necessario stabilire cosa significhi in questo contesto “buone prestazioni”, in modo da essere sicuri che questi migliorino effettivamente gli algoritmi sequenziali.

Bisogna valutare la funzione tempo dell'algoritmo parallelo $T(n, p(n))$ in confronto al tempo sequenziale $T(n, 1)$: se il tempo parallelo ha un tasso di crescita migliore di quello sequenziale allora l'algoritmo parallelo ha migliorato la soluzione.

Per il confronto tra i due tempo ci sono due possibilità:

- Il tempo **parallelo** è **circa uguale** al tempo **sequenziale**

$$T(n, p(n)) = \Theta(T(n, 1))$$

Il tempo parallelo è limitato sia superiormente che inferiormente dal tempo sequenziale, **stesso tasso di crescita**. Questo è il caso **da evitare**. Il confronto dei tempi tramite parametro efficienza esclude la possibilità di ricadere in questo caso.

- Il tempo **parallelo** **cresce più lentamente** del **tempo sequenziale**

$$T(n, p(n)) = o(T(n, 1))$$

Viene **limitato superiormente** (in qualche modo) dal **tempo sequenziale**, la funzione del tempo parallelo è sempre più piccola di quello sequenziale. Questo è il caso **da preferire**, il parametro efficienza vuole questa situazione.

1.2.1 Speed-up

Parametro che serve a definire l'efficienza, si tratta del **rapporto tra il tempo sequenziale per risolvere il problema e il tempo richiesto dall'algoritmo parallelo**:

$$S(n, p(n)) = \frac{T(n, 1)}{T(n, p(n))}$$

Esempio: se $S = 4$ l'algoritmo parallelo è 4 volte più veloce del sequenziale.

Ma in realtà se desse 4 sarebbe un caso in cui il parallelo è un Θ del sequenziale. Ponendoci nel **caso in cui** $T(n, p(n)) = o(T(n, 1))$ lo speed-up **tende a infinito**

$$S(n, p(n)) \rightarrow +\infty$$

Con un rapporto dei tempi (questo parametro), viene considerato il numero di processori? In realtà no, **non considera il numero di processori**, potremmo averne utilizzata una quantità troppo elevata/non vantaggiosa, ma questo parametro non fornisce modo di saperlo, sappiamo il rapporto tra i tempi ma non riusciamo a considerare se il numero di processori è adeguato.

Per questa mancanza, il parametro di speed-up viene poi **usato solamente per definire il parametro efficienza**.

Esempio di un problema: Soddisfacimento di formule (dove la lunghezza della formula è legata linearmente al numero di variabili coinvolte); si tratta di un problema $\in \mathcal{NP}$:

$$T(n, 1) = 2^n$$

Utilizzando un processore per ogni possibile assegnamento (sostanzialmente il funzionamento di una macchina non deterministica):

$$T(n, p(n)) = O(n)$$

$$S(n, p(n)) = \frac{2^n}{n} \rightarrow +\infty \text{ ma } p(n) = 2^n$$

Tecnicamente lo speed-up tende a infinito, ma il numero di processori è decisamente non polinomiale, costo hardware esponenziale.

1.2.2 Definizione di Efficienza

Il parametro efficienza è definito come il **rapporto tra lo speed-up e il numero di processori** $p(n)$

$$E(n, p(n)) = \frac{S(n, p(n))}{p(n)} = \frac{T(n, 1)^*}{p(n) \cdot T(n, p(n))}$$

Viene usato $T(n, 1)^*$, ovvero il miglior tempo sequenziale conosciuto per risolvere il problema o un lower bound di quest'ultimo, giusto per essere più cattivi con il nostro algoritmo parallelo.

Il parametro efficienza **può assumere i valori**

$$0 \leq E(n, p(n)) \leq 1$$

Che sia sempre ≥ 0 lo si vede dalla formula, essendo un rapporto tra valori sempre positivi.

Esempio per il problema di soddisfacimento di formule:

$$E(n, p(n)) = \frac{2^n}{2^n \cdot n} = \frac{1}{n} \rightarrow 0$$

Principio di Wyllie: Quando l'efficienza tende a 0 $E \rightarrow 0$, l'**utilizzo di processori è troppo elevato**, che magari rimarranno utilizzati la maggior parte del tempo.

Che l'efficienza sia sempre ≤ 1 invece lo si può intuire dal fatto che è un rapporto tra tempo sequenziale e tempo parallelo moltiplicato per numero di processori. Portato all'estremo, ipotizzando di poter usare n processori per risolvere in 1 passo un problema con input n il tempo sequenziale dello stesso algoritmo diventerebbe n (mettere tutti i passi uno dopo l'altro), il numero di processori è n mentre il tempo parallelo sarà 1, quindi diventa n/n (inventato da me, sotto spiegazione un po' più formale).

Per dimostrarlo, chiamo $\tilde{T}(n, 1)$ la versione sequenziale dell'algoritmo parallelo (tutti i passaggi svolti in fila), di conseguenza

$$T(n, 1) \leq \tilde{T}(n, 1) \leq p(n)t_1(n) + p(n)t_2(n) + \dots + p(n)t_{k(n)}(n)$$

l'algoritmo sequenziale “normale” è sicuramente \leq del nostro \tilde{T} , che impiega al massimo la somma del tempo di ogni passaggio moltiplicato per il numero di processori, ma questa è la definizione del tempo parallelo.

$$p(n)t_1(n) + p(n)t_2(n) + \dots + p(n)t_{k(n)}(n) = p(n) \sum_{i=1}^{k(n)} t_i(n) = p(n)T(n, p(n))$$

Il miglior tempo sequenziale è limitato superiormente dal miglior tempo parallelo moltiplicato per il numero di processori

$$T(n, 1) \leq p(n) \cdot T(n, p(n))$$

E da qua possiamo ricavare

$$\frac{T(n, 1)}{p(n)} \leq T(n, p(n))$$

Il tempo parallelo $T(n, p(n))$ ha un lower bound, ovvero il rapporto tra il miglior tempo sequenziale e il numero di processori. Il meglio che posso fare con un algoritmo parallelo è distribuire equamente il lavoro tra tutti i processori.

Spostando anche il tempo parallelo nella disequazione si ottiene

$$\frac{T(n, 1)}{p(n)T(n, p(n))} \leq 1 \implies E(n, p(n)) \leq 1$$

l'efficienza è ≤ 1 .

Di conseguenza **l'efficienza è un parametro che varia tra 0 e 1:**

- Se $E \rightarrow 0$ non va bene: implica che $p(n)$ cresce troppo velocemente, dato che $T(n, p(n)) = o(T(n, 1))$
- Il meglio da avere è $E \rightarrow k \leq 1$ dove k è una costante

Altro modo per dimostrare che $E \leq 1$, basandosi sull'idea che se $E \rightarrow 0$ allora per migliorare l'algoritmo provo a ridurre $p(n)$ senza degradare il tempo. Cambio il numero di processori da p a p/k .

Ricetta presa dalla PhD thesis di J. Wyllie.

Su un algoritmo generico, ridurre i processori a p/k significa che a ogni passo ogni processore deve svolgere in sequenza le operazioni di k processori, sostanzialmente vengono raggruppate le operazioni k alla volta. Il tempo di ogni passo sarà limitato superiormente da $k \cdot t_i(n)$ (k il tempo di ogni passo, sicuramente sarà maggiore in quanto il tempo di ogni passo è definito come il max di tutti).

Il tempo parallelo richiesto dall'algoritmo con p/k processori è limitato superiormente dalla somma dei tempi all' i -esimo passo:

$$T(n, p/k) \leq \sum_{i=1}^{k(n)} k \cdot t_i(n) = k \cdot \sum_{i=1}^{k(n)} t_i(n) = k \cdot T(n, p)$$

Quindi si ha che

$$T(n, p/k) \leq k \cdot T(n, p)$$

Di conseguenza il tempo con p/k processori è limitato superiormente dal tempo che usa tutti i p processori, moltiplicato per k .

Partendo da questa disuguaglianza si può far vedere che **l'efficienza E cresce al diminuire dei processori.**

Quanto vale $E(n, p/k)$?

$$E(n, p/k) = \frac{T(n, 1)}{\frac{p}{k} \cdot T(n, p/k)} \geq \frac{T(n, 1)}{\frac{p}{k} \cdot k \cdot T(n, p)} = E(n, p)$$

Sostituisco $T(n, p/k)$ con $k \cdot T(n, p)$ (ricordando che $T(n, p/k) \leq k \cdot T(n, p)$). Abbiamo sostituito un valore più grande, dividendo per un valore più grande ottengo un valore più piccolo, quindi la disuguaglianza diventa \geq .

$$E(N, p/k) \geq E(n, p)$$

La formula mostra che **diminuendo i processori migliora il parametro efficienza** (aumenta, rendendo vera la ricetta di Wyllie). Se il miglioramento è significativo è un altro discorso.

Considerando $k \rightarrow p$

$$1 = E(n, 1) = E(n, p/p) \geq E(n, p/k) \geq E(n, p)$$

Quindi si ottiene l'efficienza di un algoritmo sequenziale. Attenzione a mantenere $T(n, p/k) = o(T(n, 1))$ (perché $E(n, 1) = 1$, ma $T(n, p = 1) = T(n, 1)$ cioè sequenziale).

1.3 Problema Sommatoria

Perché prendiamo esempi di problemi? Motivazioni:

- **Tecnica:** scomposizione del problema in sottoproblemi e fusione dei risultati
- **Schema:** “imparando” la sommatoria posso adattare la soluzione per altre operazioni associative
- **Modulo:** sotto-problema di altri problemi più grandi

Problema della sommatoria:

- **Input:** $M[1], M[2], \dots, M[n]$, n elementi in n celle della memoria
- **Output:** $M[n] = \sum_{i=1}^n M[i]$, nell’ultima cella verrà scritta la somma di tutti i valori (presupponendo non sia importante mantenere le informazioni preesistenti nelle celle)

Algoritmo sequenziale:

```
for i = 1 to n - 1 do:  
    M[n] = M[n] + M[i]
```

Accumulo tutto in $M[n]$ sequenzialmente. Il tempo diventa $T(n, 1) = n - 1$ passaggi.

L’idea per parallelizzare è “*una somma a processore*”: somme a due a due, ogni processore somma una coppia di valori e ripeto il processo fino ad ottenere la somma finale. Servono quindi $\log n$ passi per sommare n elementi.

Questo **funziona** solo **perché** il $+$ è **associativo**:

$$((a + b) + c) + d = (a + b) + (c + d)$$

Quindi questo **schema può essere applicato ad ogni altra operazione associativa**.

Al primo passo sommo elementi a distanza 1, al secondo passo distanza 2, al terzo passo distanza 4, ..., all'ultimo passo rimarranno 2 valori da sommare a distanza $n/2$; la distanza raddoppia ad ogni passo. Ogni risultato parziale viene memorizzato nella cella di indice più alto.

Istruzioni (con k = numero del processore):

$$\begin{aligned} 1^\circ \text{ passo: } & M[2k] = M[2k] + M[2k-1] & 1 \leq k \leq \frac{n}{2} \\ 2^\circ \text{ passo: } & M[4k] = M[4k] + M[4k-2] & 1 \leq k \leq \frac{n}{4} \\ & \dots \\ \log n^\circ \text{ passo: } & M[n] = M[n] + M[n/2] & 1 \end{aligned}$$

Pseudo-codice:

```

for j = 1 to log n
    for k = 1 to n/2j par do
        M[2jk] = M[2jk] + M[2jk - 2j-1]
return M[n]

```

Questo algoritmo è EREW? Dobbiamo valutare se l'uso delle celle è **esclusivo**, ovvero se ci sono letture/scritture simultanee.

Considerando due processori a, b con $a \neq b$, a questo scopo dobbiamo verificare che $2^j a$, $2^j a - 2^{j-1}$, $2^j b$ e $2^j b - 2^{j-1}$ siano tutti diversi tra loro.

I confronti possibili quindi sono:

- $2^j a \neq 2^j b$, che vale per ogni $a \neq b$
- $2^j a \neq 2^j b - 2^{j-1}$ ovvero $2a = 2b - 1$ e non esistono valori $\in \mathbb{N}$ che soddisfano l'equazione

Quindi abbiamo dimostrato che è **un algoritmo EREW**.

Dimostrare la correttezza: Per la correttezza dobbiamo **dimostrare la proprietà**:

$$M[2^j k] = M[2^j k] + \dots + M[2^j(k-1) + 2] + M[2^j(k-1) + 1]$$

Nelle celle multiple di 2^j (con $1 \leq j \leq \log n$) ci sono 2^j valori sommati, ovvero quelli appartenenti a tutte le 2^j celle precedenti.

Per $j = \log n$ e di conseguenza $k = 1$ la proprietà sopra risulta:

$$M[n] = M[n] + \dots + M[1]$$

In $M[n]$ ci va la somma degli n precedenti valori.

Si può dimostrare **per induzione**:

- **Caso base:** $j = 1$ e $1 \leq k \leq \frac{n}{2}$. L'istruzione dell'algoritmo è

$$M[2k] = M[2k] + M[2k - 1]$$

Quindi per $j = 1$ è facile verificare che la proprietà è vera.

- **Induzione:** si presuppone che la proprietà sia vera per $j - 1$ e si dimostra che vale per j . Consideriamo l'istruzione del programma al passo j

$$M[2^j k] = M[2^j k] + M[2^j k - 2^{j-1}]$$

e dobbiamo dimostrare che questa è in realtà un'istruzione che mette nella cella di memoria $2^j k$ la somma dei suoi precedenti 2^j valori.

Raccolgo 2^{j-1} ottenendo $M[2^{j-1} 2k]$, quindi per la proprietà e per ipotesi induttiva posso dire che:

$$M[2^{j-1} 2k] = M[2^{j-1} 2k] + \dots + M[2^{j-1}(2k-1) + 1]$$

Quindi dentro $M[2^{j-1} 2k]$ sono presenti tutti i valori precedenti con un fattore moltiplicativo $(2k-1)$; è uguale alla proprietà iniziale, cambia solo che al posto di k è presente $2k$. Nella cella di indice $2k$ abbiamo la somma di tutti i precedenti 2^{j-1} valori.

Raccogliendo 2^{j-1} anche da $M[2^j k - 2^{j-1}]$ risulta $M[2^{j-1}(2k - 1)]$, che come prima può essere visto come un multiplo di 2^{j-1} , quindi applichiamo ancora la proprietà per ipotesi di induzione:

$$M[2^{j-1}(2k - 1)] = M[2^{j-1}(2k - 1)] + \dots + M[2^{j-1}(2k - 2) + 1]$$

Allora abbiamo che la cella $M[2^{j-1}(2k - 1)]$ contiene tutte le precedenti 2^{j-1} celle.

Quindi abbiamo che in $M[2^j k]$ sono contenuti 2^{j-1} valori, così come in $M[2^j k - 2^{j-1}]$, ovvero le due celle considerate all'inizio. Il totale degli elementi sommati quindi è $2^{j-1} + 2^{j-1} = 2^j$. Di conseguenza in $M[2^j k]$ vengono sommati 2^j valori.

Possiamo verificare che siano tutti valori differenti notando che sono tutti valori in sequenze decrescenti, ma l'ultimo della prima uguaglianza e il primo della seconda sono rispettivamente $M[2^{j-1}(2k - 1) + 1]$ e $M[2^{j-1}(2k - 1)]$ ovvero il secondo valore è il successivo del primo.

$$M[2^{j-1}(2k - 1) + 1] \xleftrightarrow[\text{precedente}]{\text{successivo}} M[2^{j-1}(2k - 1)]$$

Essendo sequenze decrescenti abbiamo verificato che siano tutti i valori diversi, quindi la somma delle due celle racchiude tutti i 2^j valori.

Di conseguenza

$$M[2^j k] = M[2^j k] + \dots + M[2^j(k - 1) + 1]$$

è vera.

1.3.1 Valutazione dell'algoritmo

Numero di processori: Il numero di processori impiegati è dato dal livello più “costoso”, ovvero il primo, nel quale ne vengono usati $n/2$ (anche se a ogni passo ne verranno usati la metà di quello precedente)

$$p(n) = \frac{n}{2}$$

Tempo: Il tempo impiegato deve considerare, oltre al numero di passaggi, le microistruzioni necessarie:

- LD per caricare il primo numero
- LD per caricare il secondo numero
- ADD per sommare i numeri
- ST per rimettere il numero in memoria centrale

ovvero 4 operazioni per $\log n$ passaggi, quindi

$$T(n, n/2) = 4 \log n$$

Se n non è potenza di 2 bisogna allungare l'input con degli 0 fino ad arrivare al multiplo di 2 successivo ad n . Non può peggiorare drasticamente le prestazioni in quanto in binario, chiamando t il numero di bit necessari per rappresentare n , la potenza di 2 più vicina ad n diventa t zeri con in testa un 1 e quindi sarà sempre compresa tra n e $2n$. Nel peggiore dei casi l'input si raddoppia.

Quindi le valutazioni viste per n potenza di 2 valgono anche nelle altre casistiche, però:

$$p(n) = \frac{2n}{2} = n$$

$$T(n, n) = 4 \log 2n \leq 5 \log n$$

$$\implies p(n) = O(n) \quad \text{e} \quad T(n, n) = O(\log n)$$

Efficienza: Possiamo calcolare l'efficienza

$$E(n, n) = \frac{n-1}{n \cdot 5 \log n} \sim \frac{1}{5 \log n} \rightarrow 0 \quad (\text{lentamente})$$

I processori sono un po' “sprecati”, in quanto vengono utilizzati tutti solamente al primo passo.

Proviamo a vedere se, applicando Wyllie, si può **diminuire il numero di processori** fino a sub-lineare, quindi con $p(n) = o(n)$, **mantenendo un tempo logaritmico**, anche con diverso coefficiente. Questo per avere una efficienza che tende ad una costante diversa da zero, $E \rightarrow c \neq 0$.

Quindi partendo da n numeri, al posto di avere $n/2$ processori ne abbiamo p (per ora indefinito). Ognuno dei p processori dovrà **farsi carico di** una certa quantità di **numeri da sommare** $\Delta = n/p$. Ogni processore effettuerà la somma di Δ elementi.

I **risultati** verranno messi nella **cella coinvolta dalla somma** con **indice più alto**, quindi $M[\Delta], M[2\Delta], \dots, M[p\Delta]$.

Con k indice del processore, 1° passo parallelo:

$$M[k\Delta] = M[k\Delta] + \dots + M[(k-1)\Delta + 1] \quad 1 \leq k \leq p$$

Quindi nella cella $M[k\Delta]$ verrà messa la somma dei Δ numeri affidati al processore k .

Passi **paralleli successivi**: i risultati parziali sono nelle celle

$$M[\Delta], M[2\Delta], \dots, M[p\Delta]$$

A queste celle applichiamo l'algoritmo sommatoria visto prima, in modo tale che $M[p\Delta] = M[n] = \sum_i M[i]$, ovvero le somme di tutti i gruppetti parziali. Questo **dimostra la correttezza** dell'algoritmo.

Di conseguenza applico l'algoritmo precedente per la sommatoria alle p somme parziali ottenute. Tempo n/p per il primo passo e $5 \log p$ per i restanti.

Valutazione: Per il nuovo algoritmo:

- I **processori** sono

$$p(n) = p$$

- Il **tempo** diventa

$$T(n, p) = T(1^\circ \text{ passo}) + T(\text{passi succ.}) = \frac{n}{p} + 5 \log p$$

Di conseguenza l'**efficienza**

$$E(n, p) = \frac{n - 1}{p \left(\frac{n}{p} + 5 \log p \right)} = \frac{n - 1}{n + 5p \log p} \rightarrow c \neq 0$$

Non abbiamo dato un valore a p , ma sarebbe carino se $5p \log p = n$ in modo tale che $n + n$ al denominatore abbatta il tempo sequenziale sopra, **tendendo a una costante** $c \neq 0$. Presupponendo $5p \log p = n$:

$$E(n, p) \sim \frac{n}{2n} \rightarrow \frac{1}{2}$$

Per fissare p :

$$p \log p = \frac{n}{5} \implies p = \frac{n}{5 \log n}$$

e si può facilmente verificare che il valore è corretto.

Ricapitolando: i valori diventano:

$$p(n) = \frac{n}{5 \log n}$$

$$T(n, p(n)) = 5 \log n + 5 \log n - \dots \leq 10 \log n$$

Quindi il numero di processori è sub-lineare, mantenendo il tempo logaritmico, seppur con una diversa costante.

Si può fare di meglio? Si può migliorare l'algoritmo per la sommatoria? Si può dimostrare un lower bound sfruttando il fatto che un algoritmo per la sommatoria non è altro che un albero binario e può essere sempre rappresentato come tale. Il livello con più nodi determina il numero di processori, mentre l'altezza dell'albero determina il tempo.

Di conseguenza **meglio di $\log n$ non si può fare.**

1.3.2 Sommatoria come schema per altri problemi

Una operazione iterata, con la proprietà di essere associativa:

- Input $M[1], \dots, M[n]$
- Output $Op_i M[i] \rightarrow M[n]$, ovvero effettuo l'operazione su tutti gli elementi e metto il risultato in $M[n]$

Esempi di operazioni associative: $+, *, \wedge, \vee, \oplus$.

Questo permette di avere una soluzione efficiente parallela:

$$p = O\left(\frac{n}{\log n}\right) \quad T = O(\log n)$$

Con modelli di PRAM più potenti si possono ottenere tempi migliori, fino a costante.

Problema dell' \wedge -iterato:

$$M[n] = \bigwedge_i M[i]$$

Programma

```
for  $i \leq k \leq n$  par do
    if  $M[k] = 0$  then
         $M[n] = 0$ 
```

Sostanzialmente, mi basta uno 0 ed il risultato diventa 0.

Per questo programma però serve un'architettura **CRCW** (lettura esclusiva in realtà), in quanto potenzialmente **più processori** dovranno **accedere** alla **cella** $M[n]$. Si può usare la politica common.

Le **prestazioni** che ne risultano:

$$p(n) = n$$

$$T(n, n) = 3$$

$$E(n, n) = \frac{n-1}{n \cdot 3} \rightarrow \frac{1}{3}$$

Si può usare lo stesso programma per l'**OR**, basta cambiare 0 con 1 e diventa “basta un 1 per far diventare 1 il risultato”.

La sommatoria può essere un sottoproblema di altri:

- prodotto interno di vettori: $\langle \dots, \dots \rangle$
- prodotto matrice-vettore
- prodotto matrice-matrice
- potenza di una matrice

Prodotto interno di vettori:

- Input: $x, y \in \mathbb{N}^n$, due vettori di n elementi
- Output: il prodotto interno dei due vettori, $\langle x, y \rangle = \sum_{i=1}^n x_i * y_i$

Sono n prodotti e sommarne i risultati, quindi $n - 1$ somme, il tempo sequenziale risulta $2n - 1$.

La soluzione EREW è applicare il **modulo sommatoria** con in **input** ad **ogni processore** un **gruppo di** $\Delta = n/p$ coppie di $x_i * y_i$ (sommatoria parallela con la riduzione dei processori descritta prima).

- Fase I: Δ prodotti in sequenza per processore e somma sequenziale
- Fase II: somma di p numeri in parallelo

Per la sommatoria (Fase II) serve un numero di processori:

$$p = c_1 \left(\frac{n}{\log n} \right) \quad t_{II} = c_2 \log n$$

Mentre per la Fase I:

$$p \sim \frac{n}{\log n} \implies \Delta = \frac{n}{p} = \log n \implies t_I = c_3 \log n$$

Riassumendo: $\langle x, y \rangle$ **costa**: $p(n) \sim \frac{n}{\log n}$, $t = t_I + t_{II} \sim \log n$, quindi:

$$E \sim \frac{2n - 1}{\frac{n}{\log n} \cdot \log n} \rightarrow c \neq 0$$

Ci sono costanti non riportate, viene un numero minore di 1 ovviamente.

Prodotto matrice vettore:

- Input: una matrice $A \in \mathbb{N}^{n \times n}$ ed un vettore $x \in \mathbb{N}^n$
- Output: il prodotto tra i due Ax

Ognuna delle n **righe** della matrice A viene **moltiplicata per il vettore**, quindi sono necessarie $n(2n - 1) = 2n^2 - n$ operazioni (ognuna delle n righe della matrice è un vettore moltiplicato per un altro, quindi n volte il prodotto vettoriale).

Sapendo già come fare il prodotto interno, si può usare il modulo di quest'ultimo in parallelo n -volte per risolvere il prodotto matrice-vettore.

Il vettore x è acceduto simultaneamente dai moduli per il prodotto interno \implies algoritmo CREW.

Le **prestazioni** che ne risultano:

$$p(n) = n \cdot \frac{n}{\log n} \quad T(n, p(n)) \sim \log n$$

$$E(n, p(n)) \sim \frac{n^2}{\frac{n^2}{\log n} \cdot \log n} \rightarrow c \neq 0$$

Prodotto matrice matrice:

- Input: due matrici $A, B \in \mathbb{N}^{n \times n}$
- Output: il prodotto $A \times B$

Si può risolvere usando n^2 **prodotti interni in parallelo**.

Tempo sequenziale: $n^{2.8}$ (Strassen), oppure $n^{2.37}$ (Le Gall 2014)

Ogni riga di A e colonna di B vengono accedute simultaneamente \implies algoritmo CREW. La lettura concorrente si potrebbe eliminare usando un broadcast dei dati.

Le **prestazioni** che ne risultano:

$$p(n) \sim n^2 \cdot \frac{n}{\log n} \quad T(n, p(n)) \sim \log n$$

$$E(n, p(n)) \sim \frac{n^{2.8}}{\frac{n^3}{\log n}} = \frac{n^{2.8}}{n^3} \rightarrow 0 \quad (\text{lentamente})$$

Anche se lentamente, tende a zero, quindi in questo modo non è considerato efficiente.

Potenza di matrice:

- Input: una matrice $A \in \mathbb{N}^{n \times n}$
- Output: A^n con $n = 2^k$

Si tratta di un **prodotto iterato** della **stessa matrice**, quindi posso moltiplicare $A \times A$, poi $A^2 \times A^2$, ecc.; così richiedendo un numero logaritmico di passaggi.

Sequenzialmente:

```
for i = 1 to log n do  
    A = A × A
```

Tempo sequenziale: $n^{2.8} \cdot \log n$.

Per il parallelo: sfrutta la stessa idea del sequenziale, quindi bisogna effettuare $\log n$ volte il prodotto $A \times A$, sostituendo il prodotto sequenziale con quello parallelo. Ovviamente algoritmo CREW.

Le **prestazioni** che ne risultano:

$$p(n) = \frac{n^3}{\log n} \quad T(n, p(n)) = \log n \log n = \log^2 n$$

$$E(n, p(n)) \sim \frac{n^{2.8} \cdot \log n}{\frac{n^3}{\log n} \cdot \log^2 n} = \frac{n^{2.8}}{n^3} \rightarrow 0$$

Tende a zero, anche se lentamente, quindi non è efficiente, secondo questo parametro.

1.4 Problema delle somme prefisse

Definizione del problema:

- **Input:** n celle $M[1], \dots, M[n]$
- **Output:** $\sum_{i=1}^k M[i] \rightarrow M[k]$, per ogni $1 \leq k \leq n$

Ogni cella deve contenere i valori di tutte le celle precedenti

Per semplicità verrà assunto n potenza di 2.

Algoritmo sequenziale:

```
for  $k = 2$  to  $n$  do  
     $M[k] = M[k] + M[k-1]$ 
```

Quindi tempo $n - 1$.

Prima proposta parallela: risolvo con un **modulo sommatoria** tutti i **possibili prefissi**.

Problemi:

1. **Non EREW** in quanto ogni cella è acceduta da tutti i moduli sommatoria, ma risolvibile.
2. Un algoritmo CREW su PRAM con:

$$p(n) \leq (n-1) \cdot \frac{n}{\log n} \sim \frac{n^2}{\log n}$$

$n - 1$ moduli sommatoria, ognuno composto di $n / \log n$ processori Di conseguenza

$$T(n, p(n)) \sim \log n$$

$$E \sim \frac{n-1}{\frac{n^2}{\log n} \cdot \log n} \rightarrow 0$$

L'efficienza tende a 0, quindi non eccezionale.

1.4.1 Pointer Doubling

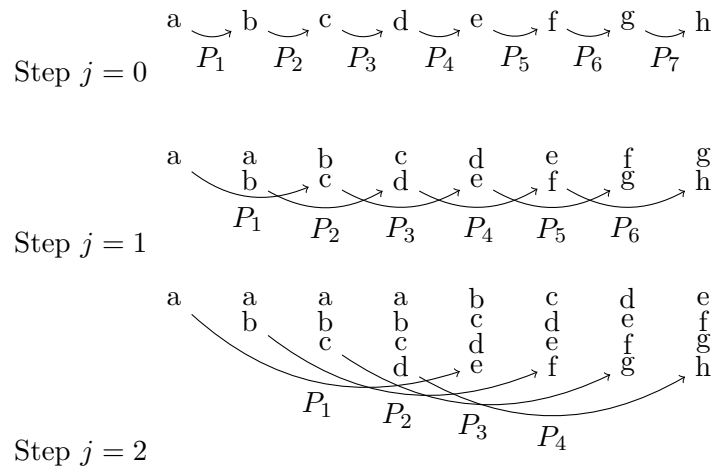
Nel '73 Kogge-Stone ha fornito un algoritmo chiamato pointer doubling.

L'idea: stabilire dei **legami** tra i numeri ed **ogni processore** si occupa di **un legame** e ne fa la somma.

Se i valori m e k sono nelle celle i e j e queste hanno un legame allora il processore i farà la **somma** $m+k$ e la **metterà in posizione** j (la più alta).

Come si continua? Prima si “**legano**” i numeri a distanza 1, poi a distanza 2, poi 4, a **potenze di 2**, fino al termine. Al passo j si legano i valori a distanza 2^j .

Esempio di passaggi:



Facendo tutte le somme in quest'ordine, ogni cella **conterrà la somma delle precedenti**, l'algoritmo **termina** una volta arrivati al punto in cui **nessun elemento ha un successore** in range, ovvero il “salto” di celle porta tutti gli elementi fuori range.

Quanti passi sono necessari? Si completa quando nessun elemento ha un successore, considerando che al j -esimo passo sono presenti 2^j elementi senza successori, l'ultimo passo dell'algoritmo è quando $2^j = n \implies j = \log n$, **servono** $\log n$ **passi**.

Ad ogni passaggio vengono **usati** $n - 2^j$ **processori**, quindi vengono usati tutti i processori con indice $1 \leq k \leq n - 2^{j-1}$.

Poniamo nel vettore $S[k]$ la posizione del **successore** di $M[k]$ (la cella legata), come viene **inizializzato** prima dell'esecuzione dell'algoritmo vero e proprio?

$$S[k] = k + 1 \text{ e } S[n] = 0$$

Con n indice dell'ultima cella (dimensione del problema).

Dato il processore P_k , quale istruzione su M deve eseguire?

$$M[S[k]] \leftarrow M[k] + M[S[k]]$$

Poi bisogna aggiornare S :

$$S[k] \leftarrow (S[k] == 0 ? 0 : S[S[k]])$$

Se $S[k] = 0$ rimane 0, altrimenti diventa $S[S[k]]$.

Codice dell'algoritmo parallelo (con M ed S già inizializzati):

```

for  $j = 1$  to  $\log n$  do
    for  $1 \leq k \leq n - 2^{j-1}$  par do
         $M[S[k]] = M[k] + M[S[k]]$ 
         $S[k] = (S[k] == 0) ? 0 : S[S[k]]$ 

```

All'inizio l'unico elemento che non ha successore, ovvero con $S[k] = 0$ è l'ultimo, ad ogni passo parallelo se ne aggiungono 2^{j-1} (si arriva a 2^j).

Questo algoritmo è EREW? La scrittura è esclusiva ma ogni processore deve **leggere** anche da una **cella che verrà usata anche dai processori adiacenti**, però **ognuno** di questi lo farà in **tempi diversi**: k accederà prima alla cella $M[k]$ poi a $M[S[k]]$, mentre $k + 1$ accederà ad $M[S[k]]$ e poi ad $M[S[S[k]]]$. Questi sono passi paralleli quindi **la stessa cella non verrà usata nello stesso momento**.
L'algoritmo è **EREW**.

Correttezza:

1. è **EREW PRAM**: P_k lavora su $M[k]$ e $M[S[k]]$: se $i \neq j \implies S[i] \neq S[j]$, quindi hanno successori diversi (se non sono 0)
2. **dimostro** $M[k] = \sum_{i=1}^k M[i]$, $1 \leq k \leq n$. Si lavora sulla proprietà del j -esimo passo:

$$M[k] = \begin{cases} M[k] + \dots + M[1] & k \leq 2^j \\ M[k] + \dots + M[t - 2^j + 1] & k > 2^j \end{cases}$$

C'è sempre la somma dei 2^j elementi precedenti, sia che si arrivi a 1 che a $t - 2^j + 1$.

Se questa proprietà è vera si ha per $j = \log n$ (ultimo passo):

$$M[k] = \begin{cases} M[k] + \dots + M[1] & k \leq 2^j = 2^{\log n} = n \\ \dots & k > 2^j = n \end{cases}$$

Si dimostra per induzione su j :

- **Base:** $j = 1$, per $k \leq 2$

$$\begin{array}{ll} k = 1 & M[1] \leftarrow M[1] \\ k = 2 & M[2] \leftarrow M[1] + M[2] \end{array}$$

per $k > 2$:

$$M[t + 1] \leftarrow M[t] + M[t + 1] = M[k - 1] + M[k] \rightarrow M[k]$$

Quindi il caso base è vero.

- **Passo induttivo:** si presuppone vera la proprietà per il passo $j - 1$ e bisogna dimostrare che vale per j .

Prima di iniziare il j -esimo passo quanto vale S ?

$$S[k] = \begin{cases} k + 2^{j-1} & k \leq n - 2^{j-1} \\ 0 & k > n - 2^{j-1} \end{cases}$$

I legami al j -esimo passo legano numeri a distanza 2^{j-1} .

Le celle con indice $\leq 2^{j-1}$ sono già a posto, in quanto la proprietà è vera per $j - 1$.

Le celle con indice compreso tra 2^{j-1} e 2^j , hanno indice t indicato come:

$$2^{j-1} \leq t \leq 2^j \implies t = 2^{j-1} + a$$

Quindi sappiamo che

$$M[a + 2^{j-1}] \leftarrow M[a] + M[a + 2^{j-1}]$$

Ma a è per forza $\leq 2^{j-1}$ e t invece è $t > 2^{j-1}$, quindi, per ipotesi induttiva, i valori di $M[a]$ e $M[a + 2^{j-1}]$ corrispondono a:

$$M[1] + \dots + M[a] + M[a + 1] + \dots + M[a + 2^{j-1}]$$

Invece, per le celle con indice $t > 2^j$, quindi $t = 2^j + a$. Gli elementi sommati in $M[a + 2^j]$ sono:

$$M[a + 2^j] \leftarrow M[a + 2^{j-1}] + M[a + 2^j]$$

Ma, considerando che l'indice è $> 2^{j-1}$, per ipotesi induttiva, corrispondono a:

$$M[a + 1] + \dots + M[a + 2^{j-1}] + M[a + 2^{j-1} + 1] + \dots + M[a + 2^j]$$

E di conseguenza la proprietà continua a essere vera.

Valutazione dell'algoritmo: Il numero di processori massimo è 1 in meno rispetto alla dimensione dell'input, usati nel primo passo:

$$p(n) = n - 1$$

Per quanto riguarda il **tempo**, si ha un ciclo **for** con $\log n$ passi, con all'interno un **par do** che fa eseguire in parallelo ai processori 2 istruzioni, somma e aggiornamento del vettore S .

La prima istruzione

$$M[S[k]] = M[k] + M[S[k]]$$

è composta da 5 microistruzioni

```
LOAD    M[k]
LOAD    S[k]
LOAD    M[S[k]]
ADD
STORE   M[S[k]]
```

Mentre l'istruzione di aggiornamento

$$S[k] = (S[k] == 0 ? 0 : S[S[k]])$$

è composta da 4 microistruzioni

```
LOAD    S[k]
JZERO
LOAD    S[S[k]]
STORE   S[k]
```

Quindi in totale sono 9 operazioni svolte $\log n$ volte, ottenendo un **tempo**:

$$T(n, n - 1) \sim 9 \log n$$

Di conseguenza l'**efficienza** è:

$$E(n, p(n)) = \frac{n - 1}{(n - 1)9 \log n} = \frac{1}{9 \log n} \rightarrow 0 \text{ (lentamente)}$$

Sfruttando **Wyllie**, come per la sommatoria, in modo da far **tendere l'efficienza** a **una costante** (far sparire la funzione $\log n$ da E).

Vogliamo migliorare l'uso dei processori, raggruppiamo di $\log n$ in $\log n$ i valori da sommare, con un processore che esegue somme in sequenza per ognuno. Di conseguenza:

$$p(n) = o\left(\frac{n}{\log n}\right), \quad T(n, p(n)) = O(\log n) \quad E \rightarrow C \neq 0$$

Il primo passo è sequenziale su $\log n$ numeri, ogni processore effettua la somma sequenziale di $\log n$ numeri, di conseguenza il numero di processori passa da lineare a $n/\log n$. Il tempo rimane logaritmico per la fase parallela, ma si aggiunge tempo logaritmico per la fase sequenziale $\sim O(\log n)$. Nel calcolo dell'efficienza $\log n$ adesso si semplifica e rimane che tende a una costante $\neq 0$.

Come per la sommatoria, anche l'**algoritmo** per le somme prefisse può essere **usato per un problema generale** “op-prefissa”, con operazione associativa:

- Input: $M[1], \dots, M[n]$, n celle
- Output: $M[k] = OP_{i=1}^k M[i]$, $1 \leq k \leq n$, in ogni cella ci sia il risultato dell'operazione sulle k celle precedenti

L'operazione deve essere associativa come ad esempio: $+, *, \wedge, \vee, \min, \max, \dots$

1.5 Valutazione di polinomi

Si tratta di un problema che prende in input un polinomio $p(x)$ di grado n e un valore α , restituendo il polinomio valutato sul valore α , quindi $p(\alpha)$.

Definizione:

- **Input:** $p(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$, α
- **Output:** $p(\alpha)$

Dati in memoria:

- il **valore** α
- $a_0, a_1, \dots, a_n \rightarrow A[0], A[1], \dots, A[n]$, n celle nella memoria condivisa, da $A[0]$ ad $A[n]$, per **contenere i coefficienti**

Algoritmo sequenziale tradizionale:

- prodotti: $\sum_{i=0}^n i \sim n^2$
- somme: n

Di conseguenza, in totale $\sim n^2$.

Ma si può fare di meglio, **miglioramento di Ruffini-Horner**: l'idea è raccogliere x in maniera iterata:

$$\begin{aligned} p(n) &= a_0 + a_1x + a_2x^2 + a_3x^3 + a_4x^4 \\ &= a_0 + x(a_1 + a_2x + a_3x^2 + a_4x^3) \\ &= a_0 + x(a_1 + x(a_2 + a_3x + a_4x^2)) \\ &= a_0 + x(a_1 + x(a_2 + x(a_3 + a_4x))) \end{aligned}$$

Generalizzando:

$$p(x) = a_0x(a_1 + \dots a_{n-2} + x(a_{n-1} + a_nx) \dots)$$

Questa forma del polinomio suggerisce un algoritmo, **sostituendo** α a **partire** dall'**ultimo coefficiente** del polinomio, partendo dalla parentesi più interna (e chiamo questo valore ottenuto p). Quindi, sostituisco α , sommo un coefficiente e ripeto fino ad arrivare alla fine

$$p = a_n \cdot \alpha + a_{n-1} \rightarrow p = a_{n-2} + p \cdot \alpha \rightarrow \dots$$

Quindi:

$$p = a_j + p \cdot \alpha$$

Codice per l'algoritmo sequenziale Ruffini-Horner:

```
Input  ( $\alpha$ )
p= $a_n$ 
for  $i = 1$  to  $n$ 
    p =  $a_{n-i} + p \cdot \alpha$ 
Output(p)
```

Prestazioni sequenziali: 2 operazioni per il numero di iterazioni

$$T(n, 1) = 2n$$

Lineare, meglio che n^2 di prima.

Possibile **algoritmo parallelo**:

- **Costruisco** il **vettore** delle **potenze di α** : Q

$$Q[k] = \alpha^k \quad 0 \leq k \leq n$$

Memorizzo le potenze di α nel vettore Q .

- Eseguo il **prodotto interno** $\langle A, Q \rangle$

$$\langle A, Q \rangle = \sum_{k=0}^n A[k] \cdot Q[k]$$

dove A è il vettore dei coefficienti. Questo effettivamente è valutare il polinomio su α .

- **Restituisco** $\langle A, Q \rangle$.

Il prodotto interno parallelo lo abbiamo già visto ed è efficiente, rimane da capire se è parallelizzabile efficientemente la creazione del vettore Q .

Per **creare** Q :

- **Pongo** α in **tutti gli elementi di** Q da 1 a n

$$Q[1] = \alpha, Q[2] = \alpha, \dots, Q[n] = \alpha$$

Non si considera la cella $Q[0]$, che deve contenere 1. Questo è un problema di replica.

- Applico il **prodotto-prefisso** su Q :

$$Q[1] = \alpha, Q[2] = \alpha^2, \dots, Q[n] = \alpha^n$$