

Sintesi ad alto livello dei sistemi digitali

Linguaggi di descrizione dell'hardware

M. Favalli



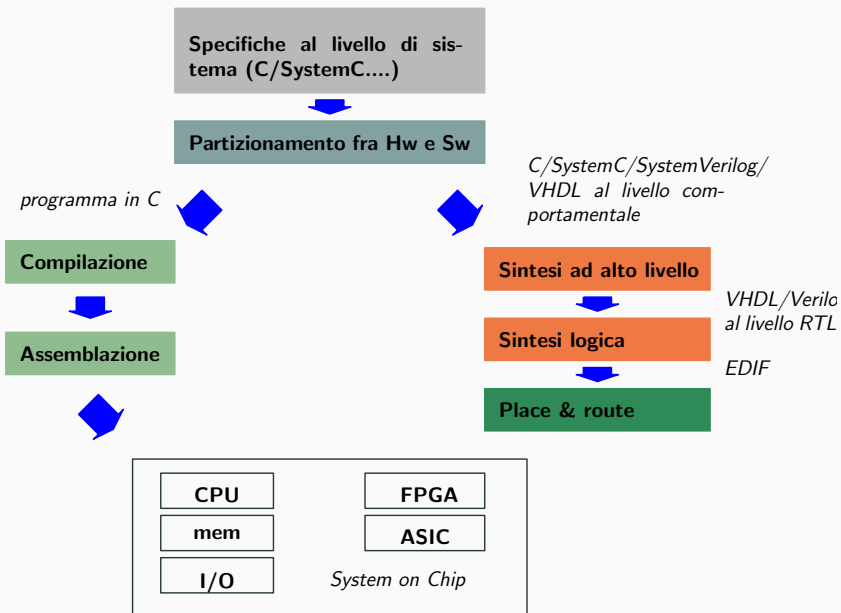
DE Department of
Engineering
Ferrara

- Presentare i concetti base della sintesi ad alto livello
- Introdurre le strutture dati utilizzate nella sintesi ad alto livello
- Introdurre i problemi relativi ad allocazione, scheduling e binding
- Descrivere alcuni algoritmi di sintesi ad alto livello
- Fornire alcuni esempi

Sommario

- 1 Introduzione
- 2 Strutture dati per HLS
- 3 Temporizzazioni
- 4 Ottimizzazione
- 5 Algoritmi

Introduzione



Realizzazione hardware di un algoritmo

- motivazioni e applicazioni
- parallelo con il software
- obiettivi
- Ottimizzazione prima della sintesi
- Strutture dati per rappresentare un algoritmo
 - Control Data Flow Graph (CDFG)
- Applicazione alla valutazione di un'espressione aritmetica
 - Data Flow Graph (DFG)
 - costo - risorse
 - allocazione e scheduling
 - prestazioni
 - controllo
- Istruzioni di selezione e di ciclo
- Flusso di progetto a basso livello

Motivazioni per l'HLS

- Algoritmo
 - indipendente dalla realizzazione hardware
 - non specifica le risorse da utilizzare
 - specifica solo in parte l'ordine delle operazioni
- CPU (semplice architettura di Von Neumann)
 - versatilità
 - sequenzialità
- Applicazioni (elaborazione dei segnali, AI)
 - specializzazione
 - specifiche:
 - prestazioni elevate
 - costi ridotti
- Obiettivi non raggiungibili con una CPU che:
 - "spreca" risorse per garantire la versatilità
 - una (o più) unità multifunzionali (ALU)
 - complessa circuiteria di controllo

Motivazioni

- L'architettura di Von Neumann per l'HLS è un caso particolare di sistema digitale
- Cambiamento di prospettive nella realizzazione di un algoritmo

Software		Hardware
Un numero limitato di unità di elaborazione polifunzionali	↔	Eventualmente numerose unità funzionali specializzate
Sequenzialità	↔	Parallelismo

- Realizzazioni hardware specializzate richiedono la disponibilità di:
 - tecnologia
 - linguaggi di descrizione degli algoritmi (HDL al livello behavioral)
 - strumenti EDA

- **Tecnologie microelettroniche per ASICs (Application Specific Integrated Circuits)**
 - Gate-Array
 - FPGA
 - Standard-Cells
- **Linguaggi di descrizione hardware consentono di descrivere sia:**
 - algoritmi (livello behavioral)
 - architetture (livello Register Transfer Level - RTL)
- **Strumenti di EDA**
 - sintesi
 - analisi (verifica formale, simulazione)

- L'algoritmo é una sequenza di passi (istruzioni) che porta al calcolo di un risultato
- Supponiamo che ciascuna istruzione possa:
 - calcolare un risultato parziale
 - decidere quale istruzione verrà eseguita di seguito
- Sia nel caso dell'hardware che in quello del software si ha un **data-path** e un **controllo**

Parallelo fra hardware e software

Costi

	Software	Hardware
Controllo	Istruzioni binarie in memoria caricate in sequenza da una FSM con tre stati: fetch, decode, execute	Macchina a stati finiti
	Istruzioni condizionali, salti (cicli)	Mappati direttamente sulla FSM
Data-path	Valutazione sequenziali delle operazioni con una ALU	Operazioni mappate su adder, moltiplicatori, divisori con la possibilità di parallelismo
	Variabili contenute in un register file con l'indirizzo di risultati e operandi contenuto nelle istruzioni	Registri non necessariamente organizzati in un register file e accessibili tramite una logica dedicata di multiplexing

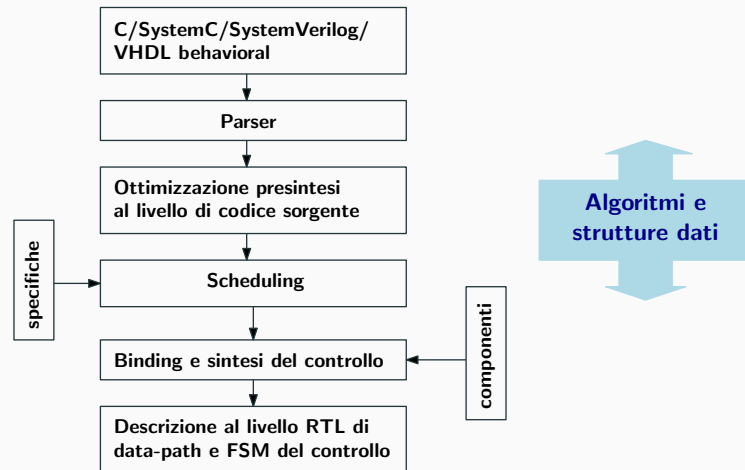
- L'ottimizzazione di costi e prestazioni é alla base della realizzazione hardware di algoritmi e in generale della sintesi ad alto livello
- A costi e prestazioni contribuiscono sia il controllo che il data-path
 - considereremo prevalentemente esempi in cui il data-path domina entrambi gli aspetti
- Area su silicio che non é nota prima della sintesi fisica
 - per esplorare diverse soluzioni, la sintesi ad alto livello necessita di stime approssimate
 - costo approssimato per i diversi componenti (c) blocchi funzionali, registri e multiplexer
 - costo stimato: $\sum_{\forall c} A(c)$

- Anche queste sono valutabili accuratamente solo dopo la sintesi fisica e quindi la sintesi ad alto livello usa modelli approssimati
 - **prestazioni dinamiche (throughput e latenza)**
 - consumo di potenza
 - collaudabilità, affidabilità

- Le specifiche di un ASIC possono riguardare
 - la banda nel caso di sistemi che operano su flussi di dati continui
 - la latenza nel caso di sistemi reattivi
 - l'area del chip o di un modulo (se il sistema è parte di un SOC)
- Considerazioni di sistema dettano a volte la frequenza di clock per cui latenza e banda possono essere date in maniera relativa
- Le specifiche possono essere date come
 - vincoli su prestazioni e costi (es. D_{max} e A_{max})
 - un vincolo su una quantità con la richiesta di minimizzare l'altra
 - minimo ritardo per una data area massima
 - minima area per un certo ritardo massimo

- Blocchi funzionali ciascuno caratterizzato da
 - tipo di dato (interi e float di diverse dimensioni)
 - operazione svolta (+, -, *, /, >, = ...)
 - stima dell'area occupata
 - stima del ritardo massimo e critical path
 - eventualmente alcuni blocchi possono realizzare più funzioni (ALU)
- Registri, tipicamente con reset e write enable
- Steering logic (multiplexer, in alcuni casi realizzati con bus)
- Logica combinatoria e FF per realizzare la FSM di controllo

- L'insieme di realizzazioni funzionalmente corrette della descrizione comportamentale del componente costituisce lo spazio di progetto
 - ogni punto è caratterizzato da prestazioni e costo
 - i tool di EDA devono essere in grado di esplorarlo in maniera efficiente
- Tipicamente non esiste un ottimo assoluto che minimizzi area e latenza
- I punti di Pareto (soluzioni non dominate) danno un'idea dello spazio di progetto
- Tassonomia di riferimento
 - rete sincrona
 - inizialmente senza pipelining



Strutture dati per HLS

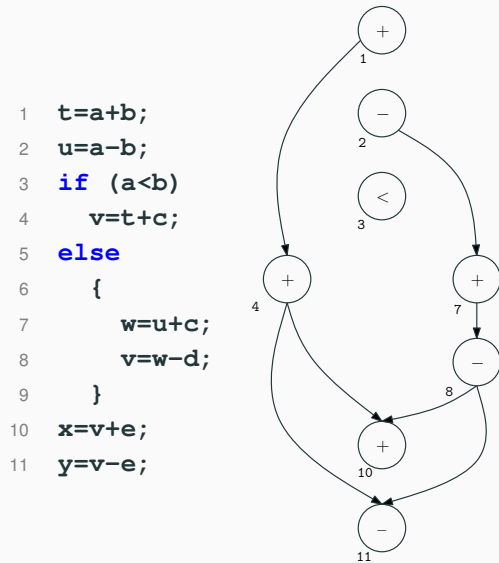
Strutture dati per rappresentare un algoritmo

- Semplice algoritmo descritto mediante un comune linguaggio C/JAVA
- Interpretazione restrittiva di tipo imperativo e sequenziale riferita ad un'architettura di tipo Von Neumann
- In realtà il codice può avere un'interpretazione molto più generale (parallelismo) che è sfruttata sia dalla sintesi dell'hardware che dai compilatori per le CPU attuali
- Sia i compilatori (software) che gli strumenti di sintesi (hardware) necessitano di una struttura dati astratta per rappresentare l'algoritmo

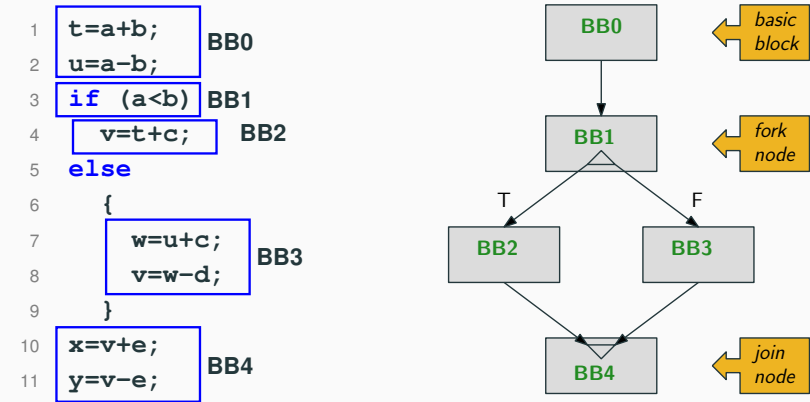
Data-Flow graph (DFG)

```
1 t=a+b;
2 u=a-b;
3 if (a<b)
4   v=t+c;
5 else
6   {
7     w=u+c;
8     v=w-d;
9   }
10 x=v+e;
11 y=v-e;
```

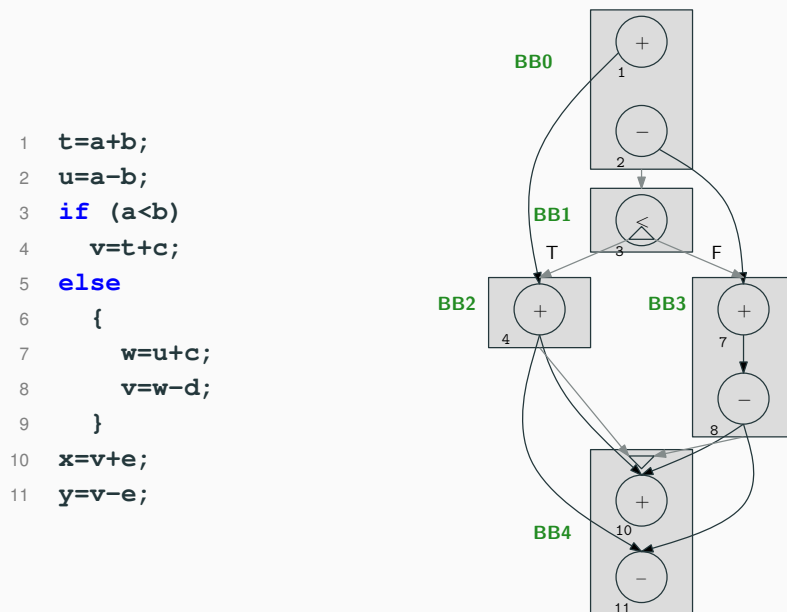
- Il significato di questo codice dal punto di vista di una CPU scalare è evidente
- In realtà l'algoritmo contiene delle informazioni più generali
- Per esplicitarle è necessario ricorrere a un modello più generale
- Si osserva che esistono istruzioni che si limitano ad alterare il flusso dei dati e altre che alterano il controllo cambiando l'ordine di esecuzione delle operazioni
 - data-flow
 - control-flow



Struttura dati astratta che rappresenta le dipendenze che esistono fra i dati elaborati dalle varie operazioni supponendo che queste siano binarie



Struttura dati astratta che rappresenta le possibili alternative nel flusso di esecuzione dell'algoritmo



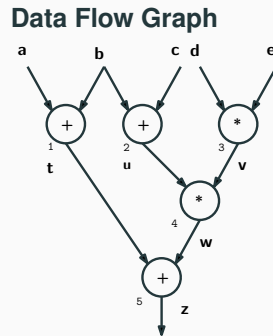
- Le strutture dati viste sono molto generali e possono essere utilizzate sia per HLS sia per la compilazione in CPU con parallelismo
- Si parte con il DFG avendo come obiettivo architetture come i DSP il cui costo é dominato dal data-path
- Ipotesi
 - rete sincrona
 - blocchi funzionali di tipo combinatorio
- Si vedranno per prima cosa alcuni concetti base utilizzati per esplorare lo spazio di progetto cercando di ottimizzare costi e prestazioni (ritardo)
- In seguito vedremo anche alcuni euristici

- Si consideri il seguente algoritmo corrispondente alla valutazione dell'espressione aritmetica

$$z = a + b + ((b + c) * (d * e))$$

```

1  t=a+b;
2  u=b+c;
3  v=d*e;
4  w=u*v;
5  z=t+w;
    
```



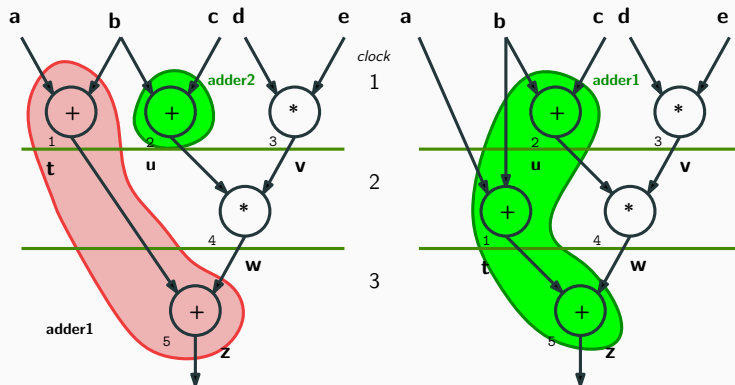
- Si noti l'ipotesi che tutti gli ingressi siano disponibili contemporaneamente

- Il DFG ci dice:
 - l'ordine relativo di priorità fra le operazioni
 - quali tipi di risorse sono necessari
- Il DFG non dice:
 - quali e quante risorse utilizzare (**allocazione**)
 - quale risorsa utilizzare per una data operazione (**binding**)
 - quando effettuare le operazioni (**scheduling**)
- Allocazione e scheduling sono legate fra loro e determinano costi e prestazioni
- Ottimizzazione

Data Flow Graph - II

Data Flow Graph - III

- Esempio di correlazione fra allocazione, scheduling e binding
- L'operazione 1 può essere eseguita sia nel primo che nel secondo ciclo di clock
- Nel primo caso sono necessari due addizionatori
- Nel secondo è sufficiente un singolo addizionatore: **riduzione di area a parità di prestazioni**



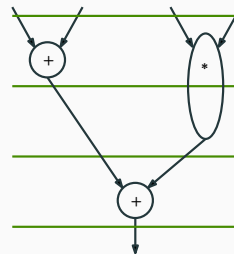
- Nell'esempio precedente si utilizza l'ipotesi di *ciclo singolo*: ciascun blocco funzionale esegue l'operazione associata in un ciclo di clock
- Sotto questa ipotesi, il periodo di clock è principalmente determinato dal blocco più lento (il moltiplicatore in questo caso)
- Questo fa sì che alcuni cicli di clock siano sovradimensionati
- Per risolvere questo problema si possono utilizzare due tecniche dette **multicycling** e **chaining**

Temporizzazioni

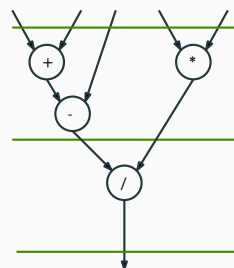
- Sotto l'ipotesi di ciclo singolo, il **cammino critico** di un DFG é il percorso che va da un ingresso a un uscita con il maggior numero di nodi
- Si noti che é un concetto diverso da quello di cammino critico in una rete combinatoria
 - il cammino critico in una rete combinatoria determina il periodo di clock minimo
 - il cammino critico di un DFG determina la latenza relativa minima del DFG

Multipling e chaining

multicycling: se il ritardo in un componente é maggiore del periodo di clock si può fare un'operazione multiciclo agendo sul comando di WE del registro ove viene scritta la sua uscita



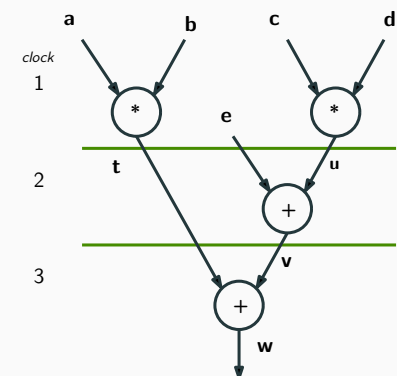
chaining: se la somma dei ritardi di cammino critico di due componenti é minore del periodo di clock, le due operazioni corrispondenti possono essere concatenate



Esempi su chaining e multicycling

- **Algoritmo:**
 $u = a * b; v = c * d; w = e + v; z = w + u;$
- **Ritardi**
 - adder 3.4ns
 - moltiplicatore 11ns
 - FF e MPX 0.5ns
- **Specifiche:** latenza minima non considerando il costo
- **Realizzazione combinatoria**
 $D = 17.8ns$

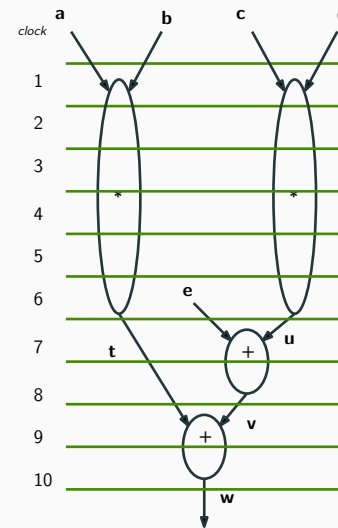
- **Realizzazione a ciclo singolo:**
 $T > 11.5ns$
- **Latenza:** $D = 3T$



- Si vuole esplorare lo spazio dei possibili periodi di clock dal punto di vista della latenza
- Chiaramente c'è un impatto sul consumo di potenza che qui non viene considerato

T (ns)	latenza (ns)	technique
2	20 (10T)	multicycling
3	21 (7T)	multicycling & chaining
4	20 (5T)	multicycling
5	25 (5T)	multicycling
6	24 (4T)	multicycling
....		
12	24 (2T)	chaining
12	36 (3T)	single cycle

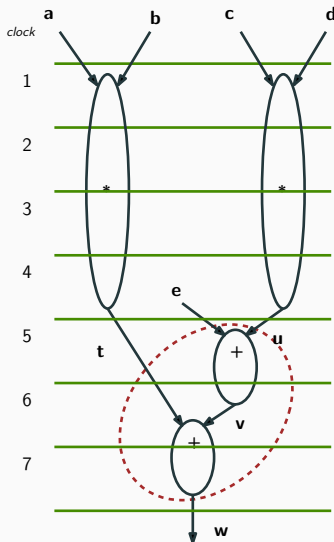
- Si osserva una dipendenza non lineare e non monotona
- Tutte le realizzazioni richiedono 2 moltiplicatori e un adder, a parte l'ultima che richiede due adder



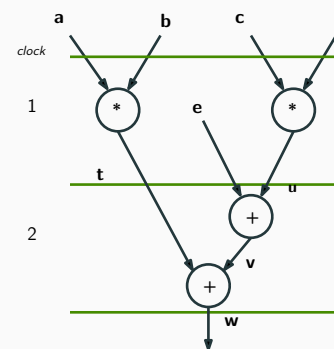
periodo di clock $T = 2ns$
 latenza $D = 10T$
 risorse 2 moltiplicatori, 1 adder

Esempio di multicycling e chaining

Esempio di chaining



- Periodo di clock $T = 3ns$
- Per le moltiplicazioni sono richiesti 4 periodi di clock
- Le somme richiedono 3.9ns e quindi sarebbe necessario un multicycling di due periodi di clock per ciascuna di esse
- Concatenando due sommatore come un sommatore a 3 operandi (tratteggiato), si può applicare il multicycling e risparmiare un periodo di clock
 - latenza $D = 7T$
 - risorse 2 moltiplicatori, 2 adder



- Periodo di clock $T = 12ns$
- Latenza $D = 2T$
- Risorse 2 moltiplicatori e 2 adder

Ottimizzazione

- Si é vista la dipendenza di costi e prestazioni da allocazione e scheduling
- Si ha uno spazio di possibili soluzioni che va esplorato alla ricerca di una soluzione che soddisfi le specifiche di progetto
- Si utilizzano le seguenti ipotesi
 - solo operazioni di tipo data-path
 - ciclo singolo
- E le seguenti specifiche
 - area minima per una data latenza
 - latenza minima per una data area
- Il problema ha un costo esponenziale nel numero di operazioni

Spazio delle possibili soluzioni

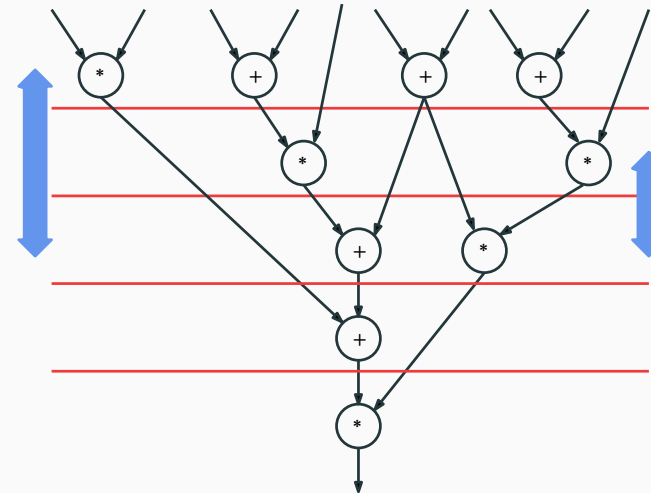
Allocazione

- Esistono dei punti particolarmente interessanti
 - soluzione con area minima (determinata dal minimo numero di blocchi funzionali che possono svolgere le operazioni del DFG)
 - soluzione con latenza minima (determinata dal cammino critico del DFG)
- L'esplorazione dello spazio di possibili soluzioni può essere fatta in maniera esatta o mediante metodi euristici se il numero di operazioni é molto grande
- Si noti che per semplicitá stiamo trascurando costi e ritardi legati a multiplexer e flip-flop: anche questi possono essere soggetto di ottimizzazioni

- Il numero di risorse da considerare per ogni tipo di blocco funzionale può
 - essere fissato come specifica iniziale (si noti che questa é una semplificazione perché potrebbe essere specificato il costo senza distinguere fra i diversi tipi di risorse)
 - essere calcolato per una certa allocazione
- **Dato uno scheduling, il numero di risorse di un certo tipo (+, *, ...) può essere calcolato come il massimo numero di quel tipo di risorse per ogni ciclo di clock**

- É l'operazione principale nell'ottimizzazione del DFG di un algoritmo
- Per prima cosa individuamo i gradi di libertà che abbiamo a disposizione con l'ipotesi di fare uno scheduling a latenza minima ottimizzando le risorse
- A ogni operazione del DFG può essere associata una **mobilità** che descrive i suoi possibili scheduling
 - le operazioni sul critical path non hanno mobilità
 - il concetto é relativo, la mobilità di un'operazione dipende da quella delle altre
- Introduciamo due tipi di scheduling che ci aiutano a valutare quantitativamente la mobilità

Ogni operazione viene programmata non appena sono disponibili i suoi operandi

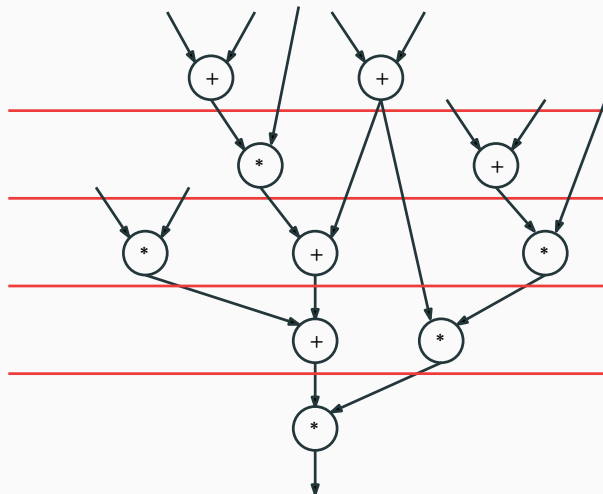


latenza = 5T

allocazione
2 adder
2 moltiplicatori

Scheduling ALAP (As Late As Possible)

La programmazione di ogni operazione viene dilazionata fino a quando un ulteriore ritardo risulterebbe in un aumento della latenza

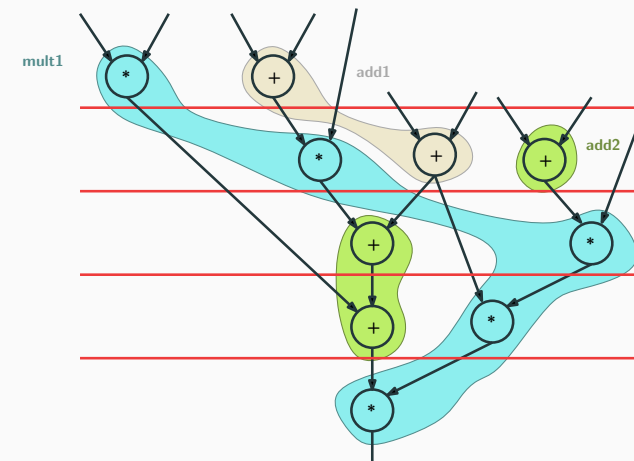


latenza = 5T

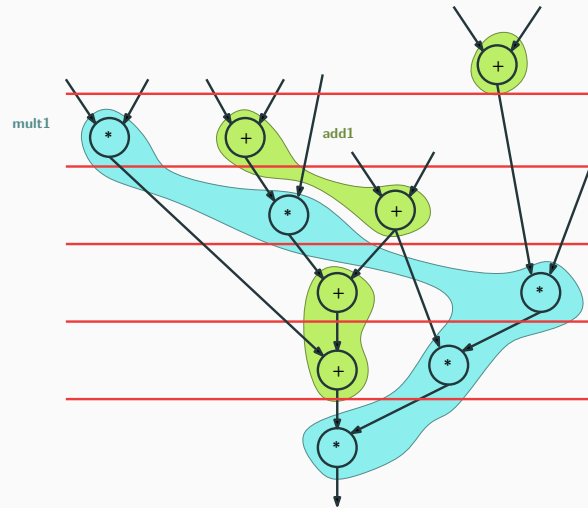
allocazione
2 adder
2 moltiplicatori

Scheduling ottimizzato

- Lo scheduling ASAP e ALAP tipicamente non presentano l'ottimo nell'allocazione delle risorse
- Infatti, questo scheduling consente di utilizzare un singolo moltiplicatore e due adder



- Allocazione di 1 adder e 1 moltiplicatore
- Si osserva che é necessario aggiungere un periodo di clock



- Abbiamo visto come scheduling e allocazione siano fortemente correlati e come il binding sia una conseguenza di questi passi di progetto
- Per arrivare a capire il ruolo della HLS ci manca ancora
 - la descrizione di approcci sistematici allo scheduling
 - una strategia per definire i dettagli implementativi corrispondenti a un certo scheduling e alla corrispondente allocazione
 - registri
 - multiplexer
 - input/output
 - controllo

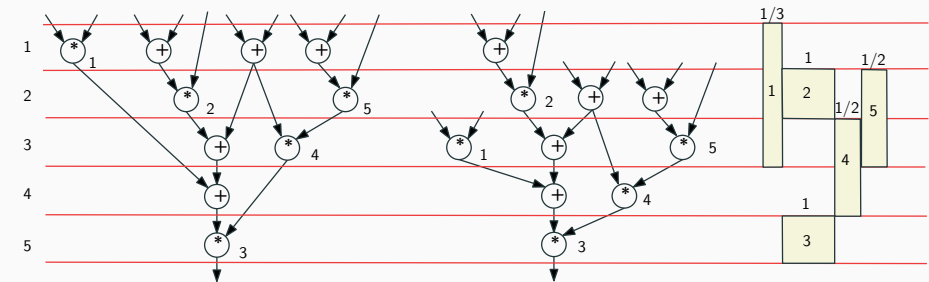
Algoritmi di scheduling

Algoritmi

- Iniziamo con il vedere un paio di tecniche euristiche per l'ottimizzazione dello scheduling di un DFG
 - **Force Directed Scheduling** per la minimizzazione delle risorse per una data latenza
 - **List Based Scheduling** per la minimizzazione della latenza per una data allocazione
- Queste tecniche hanno una complessità computazionale inferiore a quelle esatte
- Chiaramente possono produrre risultati corrispondenti a minimi locali di costo e prestazioni

- Euristico per minimizzare l'area in presenza di vincoli sulla latenza (tipicamente dati dal cammino critico)
- In pratica, la densità operazioni dello stesso tipo in un certo ciclo di controllo viene vista come una forza elastica in grado di attrarre/respingere un'operazione di quel tipo
- Tale forza viene minimizzata per ottenere una distribuzione bilanciata delle risorse
- L'algoritmo serve sia per unità funzionali che per registri e logica di steering

- Vengono fatti prima uno scheduling ASAP e uno ALAP
- Time frame di un'operazione: intervallo fra lo scheduling ASAP e quello ALAP
- Si suppone che la probabilità di programmare un'operazione in tale intervallo sia distribuita uniformemente
- Considereremo le moltiplicazioni perché dominanti dal punto di vista del costo

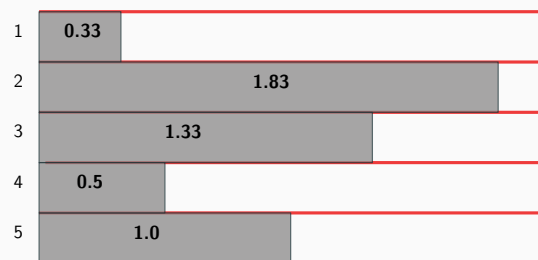


Grafo di distribuzione delle operazioni

- Concorrenza in ogni ciclo di controllo (i) delle operazioni dello stesso tipo (t)

$$DG(i) = \sum_{op=t} Prob(opn, i)$$

- Maggiore è $DG(i)$, maggiore sarà la probabilità che il ciclo i sia quello che definisce il numero di risorse per operazioni di tipo t



Calcolo della forza

- Per ogni operazione si definisce una forza per ogni ciclo di controllo (j) corrispondente a uno scheduling fattibile per tale operazione

$$Force(j) = DG(j) - \sum_{i=t}^b \frac{DG(i)}{b - t + 1}$$

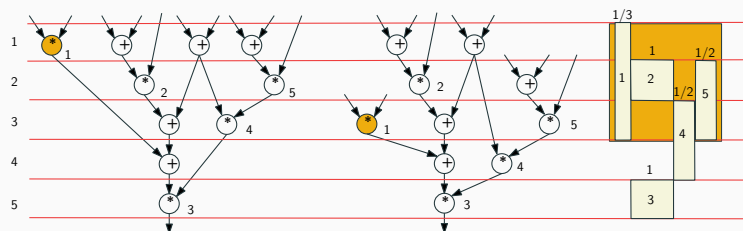
- Dove
 - b è il ciclo di controllo dell'operazione nello scheduling ALAP
 - t è il ciclo di controllo dell'operazione nello scheduling ASAP
 - DG è il grafo di distribuzione per il tipo dell'operazione considerata
- Una forza negativa indica che l'operazione è "attratta" dal ciclo j , mentre una forza positiva indica che è "respinta" dal ciclo j

- Si consideri la moltiplicazione 1 nell'esempio di DFG
- Lo scheduling di tale operazione può essere da $t = 1$ a $b = 3$

$$Force(1) = 0.33 - 1.16 = -0.83$$

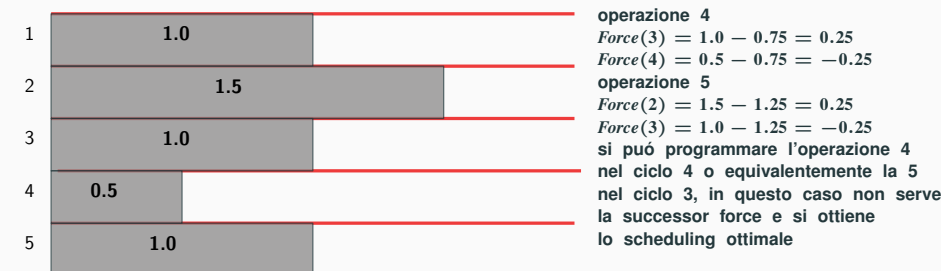
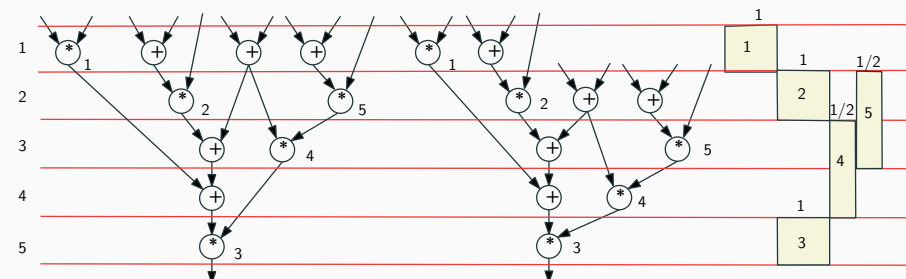
$$Force(2) = 1.83 - 1.16 = 0.33$$

$$Force(3) = 1.33 - 1.16 = 0.17$$
- Da cui si vede che risulta conveniente utilizzare il ciclo $j = 1$



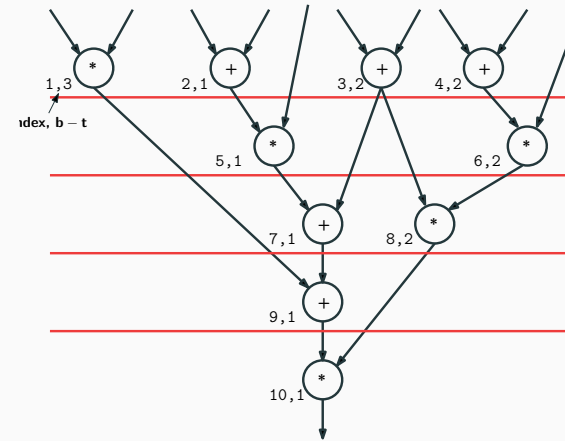
- A ogni passo l'algoritmo sceglie l'operazione con la minor forza e la programma
- Dopo questa operazione, viene aggiornato DG e vengono ricalcolate le forze e se lo scheduling non è completo, si torna al passo precedente
- In realtà il problema è più complesso perché ogni scheduling cambia la mobilità di alcune operazioni
- Bisogna aggiungere una self-force e una predecessor-successor force
 - ad esempio se si programma l'operazione 5 nel ciclo di controllo 3, l'operazione 4 è fissata nel ciclo 4 e questo viene tenuto in conto da una successor force

- Per le operazioni 4 e 5 si hanno i seguenti risultati
 - operazione 4: $Force(3) = 1.83 - 0.91 = 0.92$,
 $Force(4) = 0.5 - 0.91 = -0.41$
 - operazione 5: $Force(2) = 1.83 - 1.58 = 0.25$,
 $Force(3) = 1.33 - 1.58 = -0.25$
- La moltiplicazione 1 ha il minor valore di forza ($Force(1) = -0.83$) e quindi viene programmata nel primo ciclo di controllo
- A questo punto, tenendo fissa l'operazione 1, si ricalcola il grafo di distribuzione e le forze



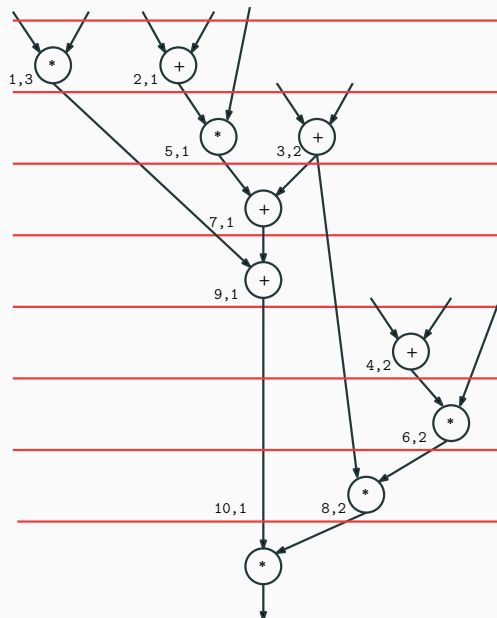
- Diversamente dal force directed scheduling, il list based scheduling cerca di minimizzare la latenza per un ben definito ammontare delle risorse
- In pratica, é una generalizzazione del metodo ASAP
- Si mantiene una lista di operazioni pronte (i cui predecessori sono già stati programmati)
- A ciascuna operazione pronta si associa una funzione di priorità
 - ad esempio l'inverso della mobilità ($b - t$)
- A ciascuna iterazione, per un dato tipo di operazione, vengono programmati i nodi con la priorità piú alta dilazionando gli altri
- Gli operatori a piú alta mobilità sono quelli programmati dopo (greedy)
- Sono state sviluppate diverse funzioni di costo

- Allocazione: 1 adder e 1 moltiplicatore
- Si calcolano lo scheduling ASAP e quello ASAP annotando ciascun nodo con $b - t$ (inverso della mobilità)



passo	lista di nodi pronti	scheduling
1	1, 2, 3, 4	1(*), 2(+)
2	3,4,5	5(*),3(+)
3	7,4	7(+)
4	9,4	9(+)
5	4	4(+)
6	6	6(*)
7	8	8(*)
8	10	10(*)

Esempio - II

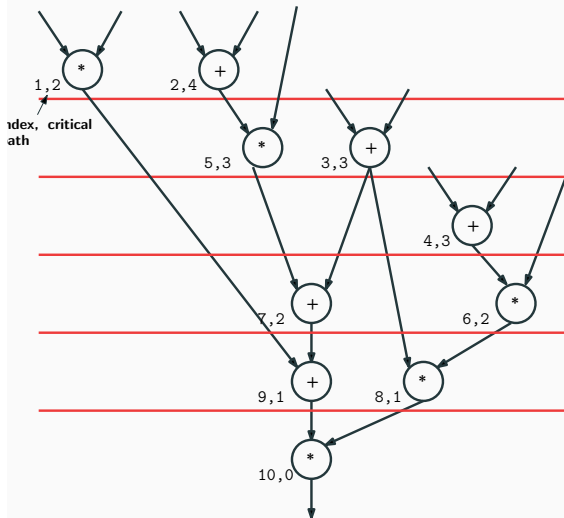


- latenza uguale a 8 cicli di controllo
- la scelta di programmare 3 invece di 4 al passo 2 é arbitraria
- una scelta diversa sarebbe risultata in una minore latenza
- trattandosi di un euristico é normale che l'algoritmo sia arrivato a un minimo locale

List based scheduling basato sul critical path

- Il critical path di un nodo é il numero massimo di nodi fra tale nodo e l'uscita
- La funzione di priorità é la lunghezza del cammino critico
- Ogni nodo é annotato con il suo indice e il suo cammino critico
- Anche in questo caso si tratta di un euristico
- Questa tecnica viene utilizzata anche dai compilatori nel caso di CPU superscalari

- Allocazione: 1 adder e 1 moltiplicatore



passo	lista di nodi pronti	scheduling
1	1, 2, 3, 4	1(*), 2(+)
2	3,4,5	5(*),3(+)
3	7,4	4(+)
4	7,6	6(*), 7(+)
5	8,9	8(*),9(+)
6	10	10(*)

- L'ottimizzazione dello scheduling può essere formulata come un problema di Integer Linear Programming
- Si tratta di un insieme di disequazioni lineari con una funzione lineare di costo (area o latenza) da minimizzare
- Una volta descritto il problema in questo modo, si può utilizzare un solver/optimizer ILP per avere la soluzione