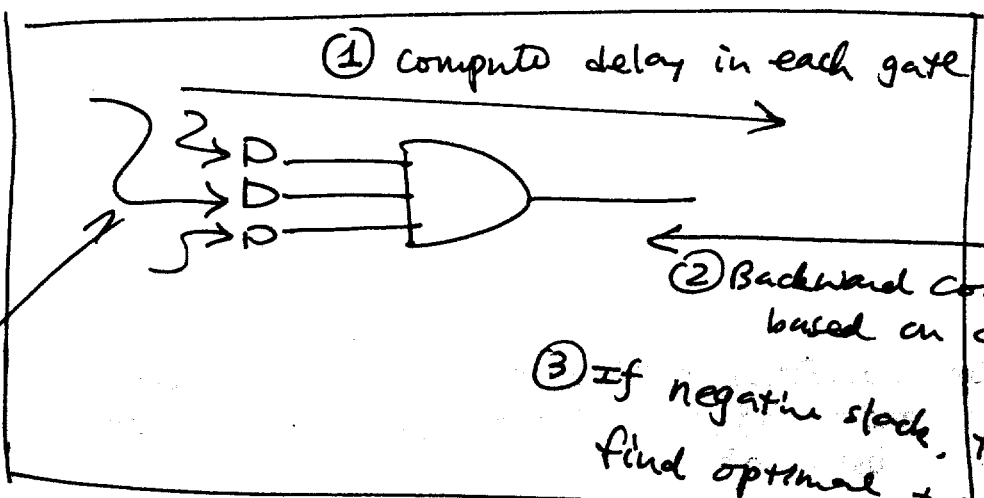


# Transistor re-ordering algorithm.

## Step I.

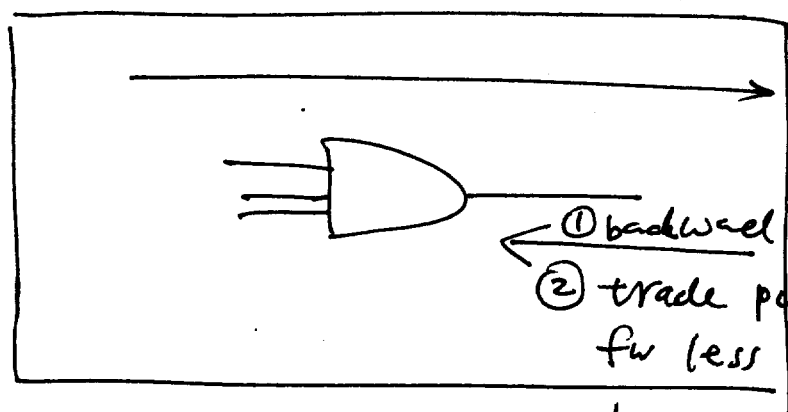


The latest arrival signal connected to the input with smallest delay.

Reorder it

- ④ Compute forward to update the delay.
- ⑤ Repeat ② → ④ until no negative slacks.

## Step 2:

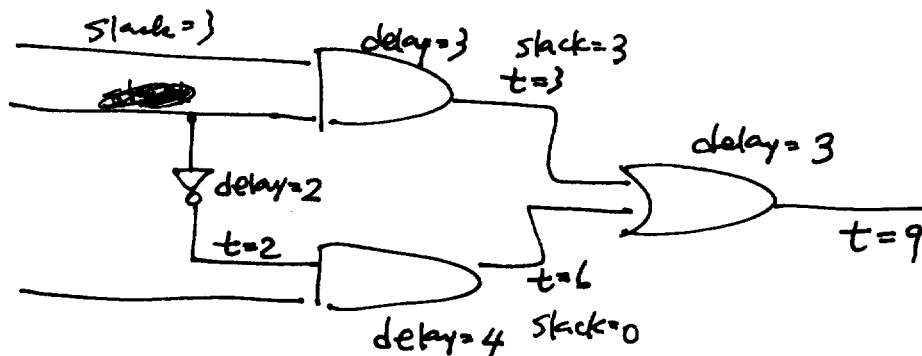


③ compute slack, make sure no neg. slack

# Experimental Results

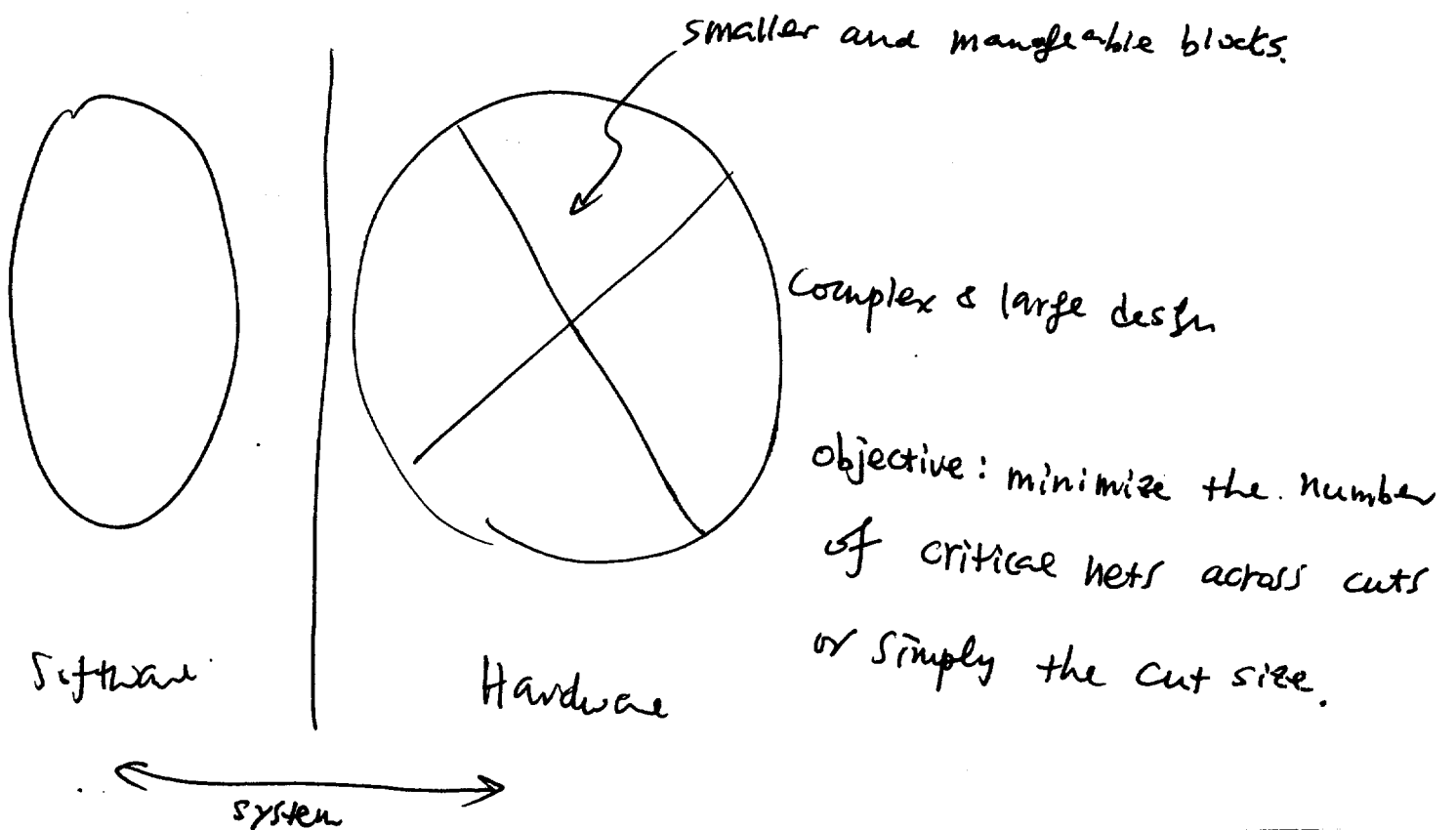
ckt	delay Reduct	power reduction
Vg2	9.5%	7.3%
rd84	6.9%	6.5%
alu4	8.2%	4.3%
	⋮	

## Stack calculation



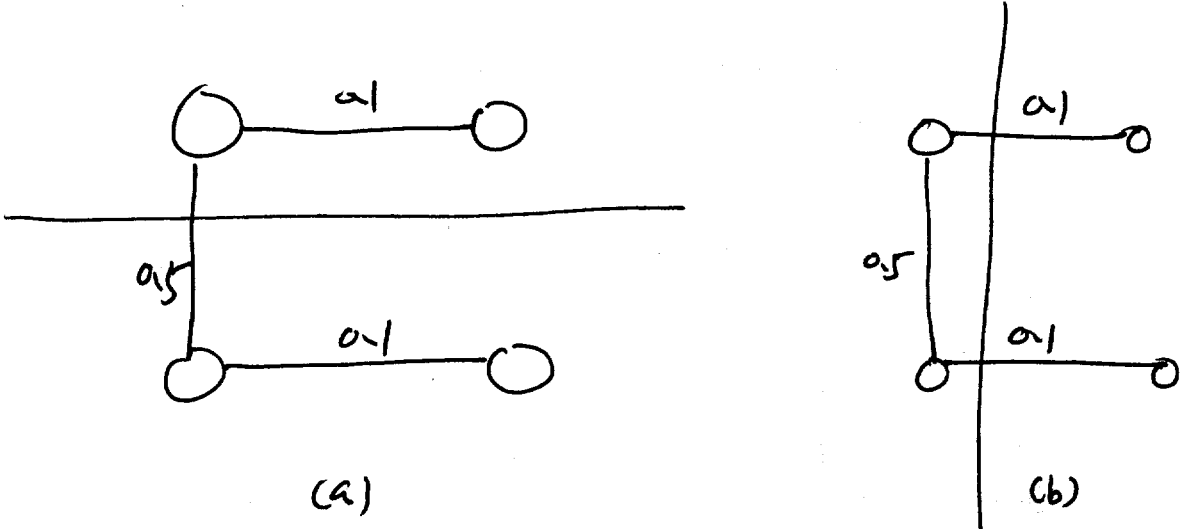
# Interconnect Layout optimization

- After device sizing, the next target is interconnect.
- In deep submicron design, the interconnect capacitance dominates the gate capacitance
  - ∴ Interconnect design is crucial.
- Discuss the interconnect issue on circuit partitioning, node clustering, placement and routing.
- Discuss wiresizing for delay and power optimization.
- ~~Circuit partitioning.~~



- For low-power design, average switching frequency of the nets are important.

Example:

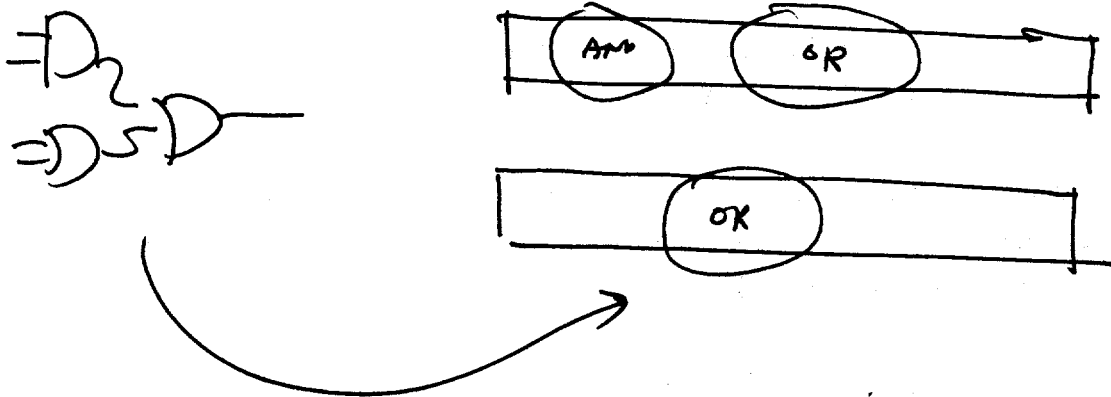


Both partitions are balanced (contains equal # of nodes).

(b) has smaller capacitive power dissipation

∴ lower switching activity across the two blocks.

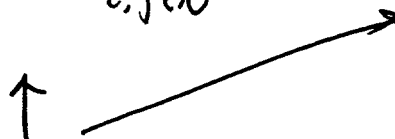
- First, ~~assess~~ estimate switching activity for each net. Then, apply a partitioning method with the objective of minimizing switching capacities.
- partitioning will be further discussed in Design automation of VLSI.



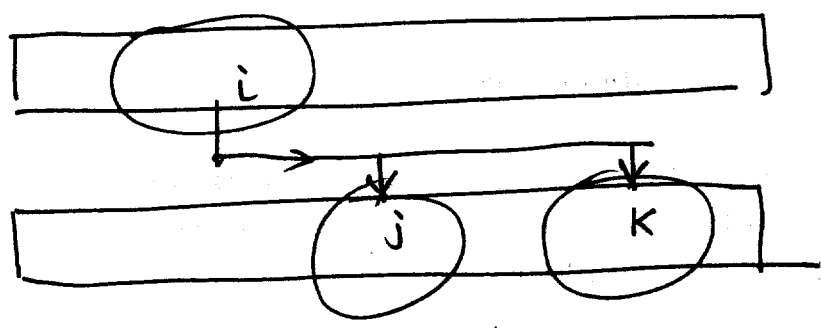
- Map the gates in a circuit onto positions on a layout surface.
- Conventional method: try to minimize the total net length of the final placement, meeting various delay constraints.
- For low-power placement: try to minimize the total <sup>effective</sup> switching capacitance, instead of total wirelength (capacitance), subject to timing constraints.

objective function:

$$L_2 = \sum_{\text{net } N} f_N^2 \sum_{\substack{\text{cells} \\ i, j \in N}} \{ (x_i - x_j)^2 + (y_i - y_j)^2 \}$$



Quadratic net length model.



$$f_N^2 [(x_i - x_j)^2 + (y_i - y_j)^2] + f_N^2 [(x_i - x_k)^2 + (y_i - y_k)^2]$$

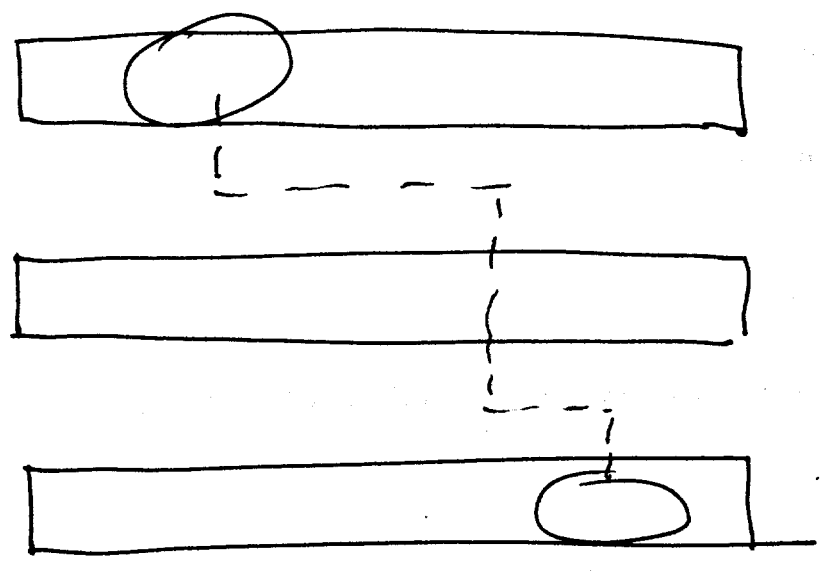
~~to find the minimum~~

- The problem can be even more complex when timing delay is added as a constraint.
- Either the lumped capacitance model, the pathlength delay model or the Elmore delay model can be used to formulate the delay model.
- An average of 10% power reduction obtained when compared with the minimum net length solution.

### Routing.

- Routing completes the electrical connection between cells
- Contains global routing and detailed routing.

• Global Routing: Assign each net to a set of routing regions



Objective: minimize the total wire length.

• Detailed Routing: Determine the actual wiring geometries and layer assignments within each region.

Objective: Minimize the channel height and total wirelength.

• Low-power routing: Minimize the total effective switching capacitance of all nets while meeting the delay constraints.

• just concentrate on global routing.

- To minimize the total effective switching capacitance, we must route high switching activity net with minimum wirelength, and allow detour for low activity nets to avoid channel congestion.
- Should avoid detour for critical nets.
- It was reported that only marginal improvement in power dissipation is observed for global routing.
- In deep submicron design, power dissipation due to ~~coupling~~ effect (cross talk) has become more significant.
- Try to find proper spacing and wire segment assignment ordering

---

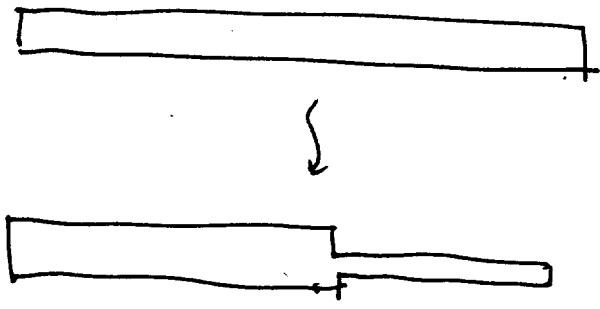
↑ Keep wider spacing for high activity nets.  
↓ and assign them to different layers.

---



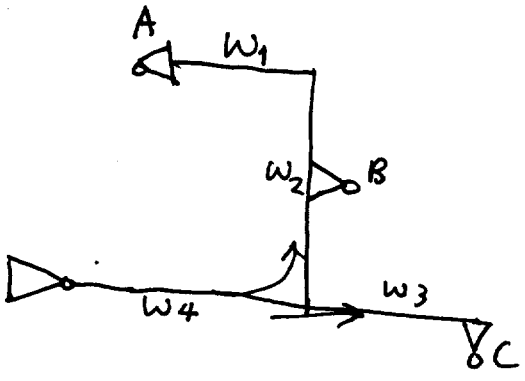
# Wiresizing optimization.

- partitioning, placement and routing mainly concentrate on interconnect length and switching activity.
- There is room for optimization for a fixed-length interconnect by sizing.



- Very important for deep submicron technology, since the delay of interconnect dominates.

## Definitions:



Routing tree:  $T$

Sink( $T$ ): The set of sinks in  $T$  (A, B, C)

$W$ : the wiresizing solution (Wire width assignment for each segment of  $T$ ).

$t_i(w)$ : the Elmore delay from the source to sink  $S_i$  under  $w$ .

~~Objective function for wiresizing:~~

$$t(w) = \sum_{S_i \in \text{Sink}(T)} \lambda_i t_i(w)$$

$\lambda_i$ : measures the criticality of sink  $S_i$ .

E.g. large  $\lambda_i$  for timing critical sinks.

~~Tradeoff between routing area (also interconnect capacitance and power dissipation) and performance.~~

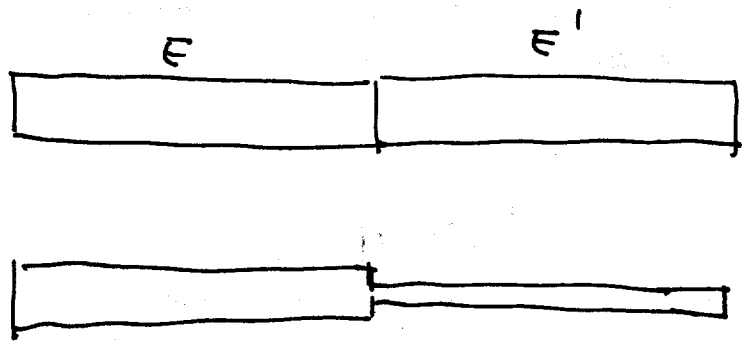
$$\alpha \text{Area}(w) + \beta \cdot t(w)$$

$\alpha, \beta$  are weights to indicate the relative importance of area and performance.

• Consider either  $t(w)$  or  $\alpha \text{Area}(w) + \beta t(w)$ , there are three important properties for wiresizing:

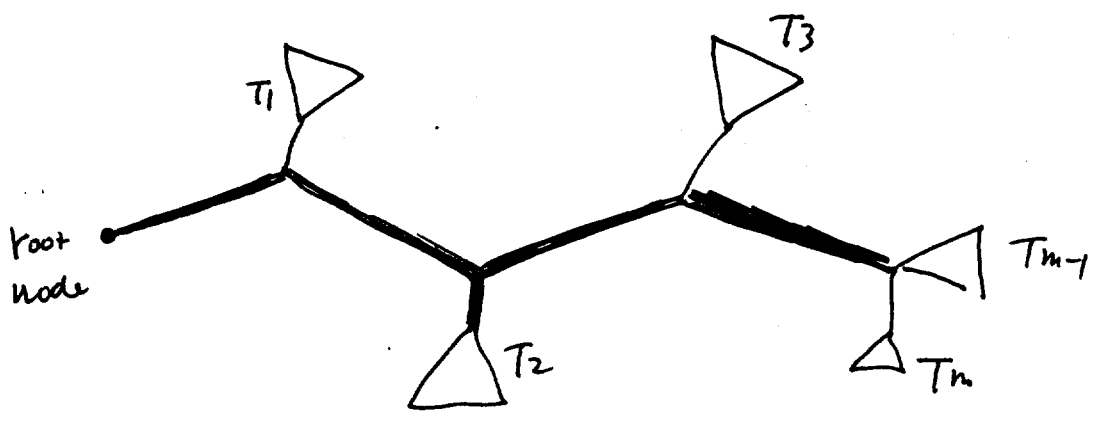
Monotone, separability, and dominance

There exists an optimal wiresizing solution which is ~~Monotonic~~. Given any pair of wire segments  $E$  and  $E'$  such that  $E'$  is the downstream of  $E$ , the width of  $E$  is no smaller than that of  $E'$ .



A very natural result since  $E$  drives more devices, so should have larger width, in general.

~~Separability~~: If the width assignment of a path  $P$  originated from the source is given, the optimal width assignment for each tree branching off  $P$  can be carried out independently.



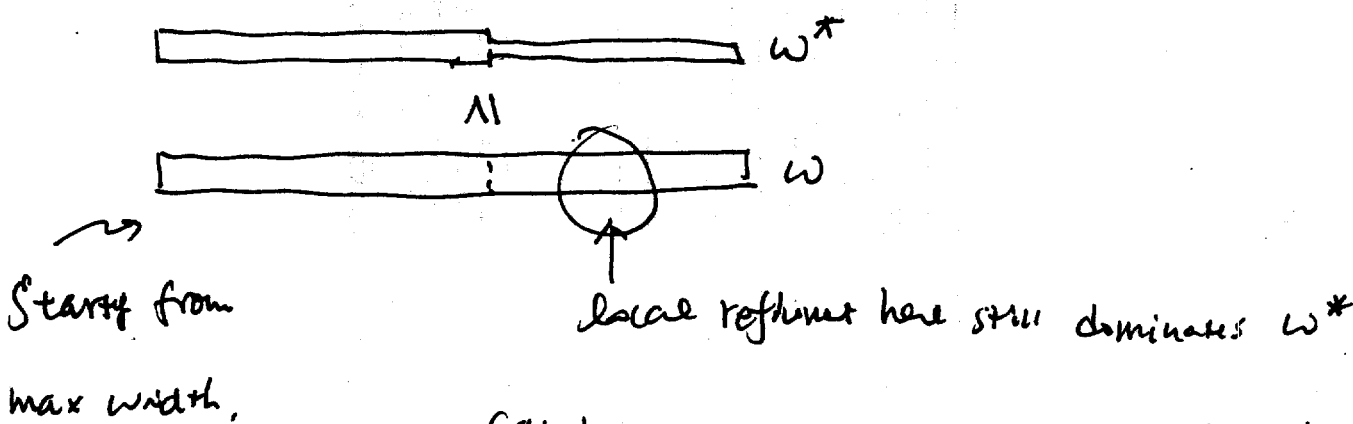
- In the above figure, based on the separability property, the optimal width assignment for  $T_1, T_2, \dots, T_m$  can be computed independently, if the width assignment of  $P$  is ~~known~~ given.
- This property is extremely important in simplifying the wire sizing problem.



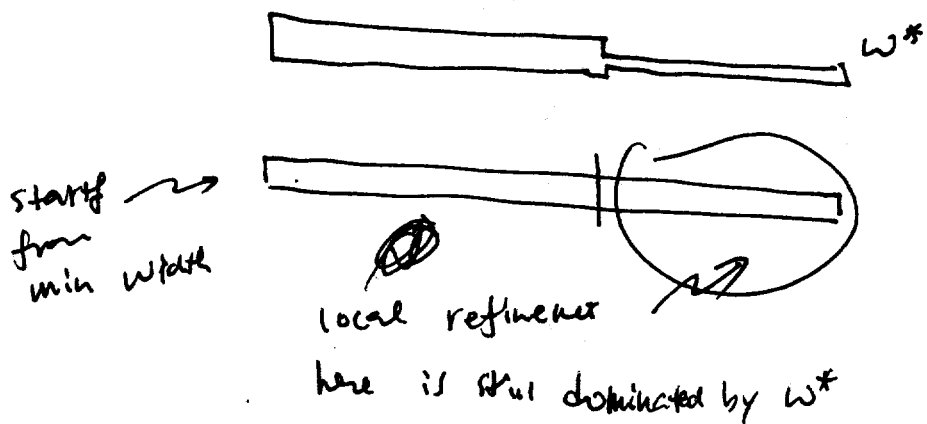
ref. Given two wire width assignments  $w$  and  $w'$ ,  $w$  dominates  $w'$  if for any segment  $E$ , the width assignment of  $E$  in  $w$  is greater than or equal to that of  $E$  in  $w'$ .

ref. Given a routing tree  $T$ , a wire width assignment  $w$  on  $T$ , and any particular segment  $E \in T$ , A local refinement of  $w$  with respect to  $E$  is the operation of optimizing the width of  $E$  while keeping the wire width assignment of other segments in  $w$  unchanged.

Thm: Let  $w^*$  be an optimal width assignment. If a width assignment  $w$  dominates  $w^*$ , then any local refinement of  $w$  still dominates  $w^*$ . Similarly, if a width assignment  $w$  is dominated by  $w^*$ , then any local refinement of  $w$  is dominated by  $w^*$ .



Can be used to find the upper bound of the optimal wire width.

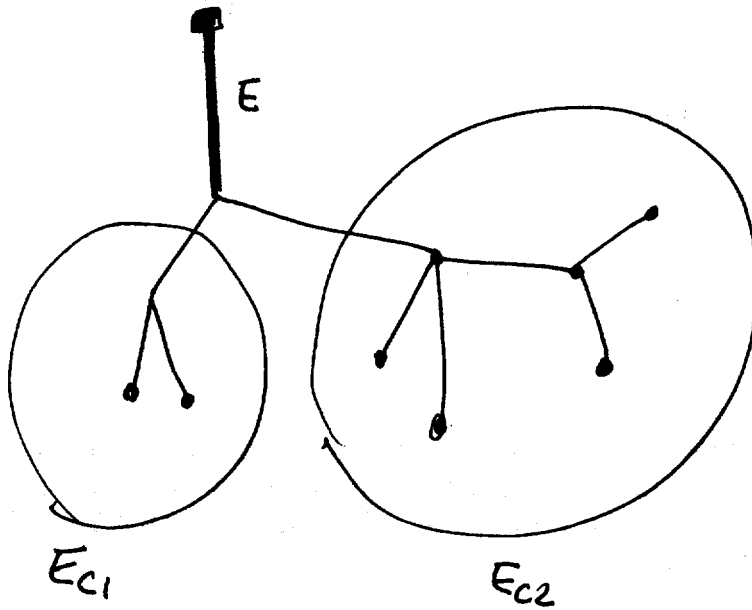


→ can be used to find the lower bound of the optimal wire width.

Based on monotone and separability properties:

(166)

optimal wiresizing Algorithm:



How to determine the width of  $E$ ?

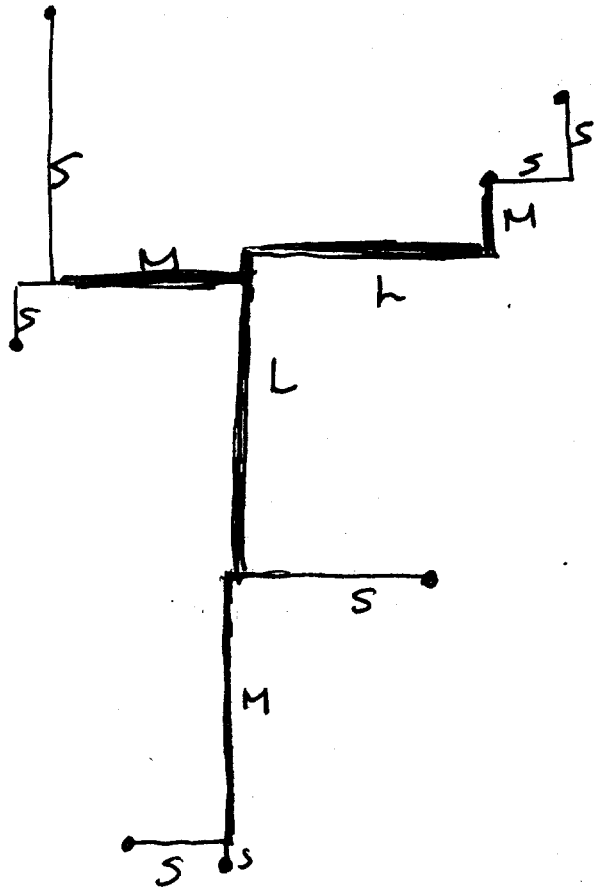
possible widths:  $w_1, w_2, \dots, w_r$ .

For each possible width assignment  $w_k$  of  $E$  ( $1 \leq k \leq r$ ), we determine the optimal assignment for each subtree  $E_{c1}, \dots, E_{c2}$ .

independently by recursively applying the same procedure to each  $E_{c_i}$  with  $\{w_1, w_2, \dots, w_k\}$  as the set of possible widths. The optimal assignment for  $E$  is the one which gives the smallest total delay.

(use  $E$  | more delay model)

A possible optimal wire sizing.



3 wire sizes

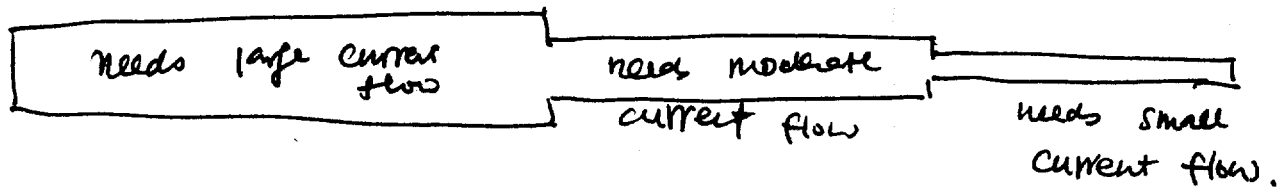
$h, M, S$

• Without this technique, designers generally use conservative <sup>(767)</sup>

design: End up ~~with~~ using oversized, uniform-width wires

i. Waste too much area & power, causes longer signal delay.

• Optimal wiresizing can be used to "taper" wires for area and power dissipation without performance degradation.



• A computer simulation result demonstrates:

Area reduction:  $10\%$  compared with manually tapered design.  
reduced by

clock power ~~is~~ reduced by  $6\%$ .

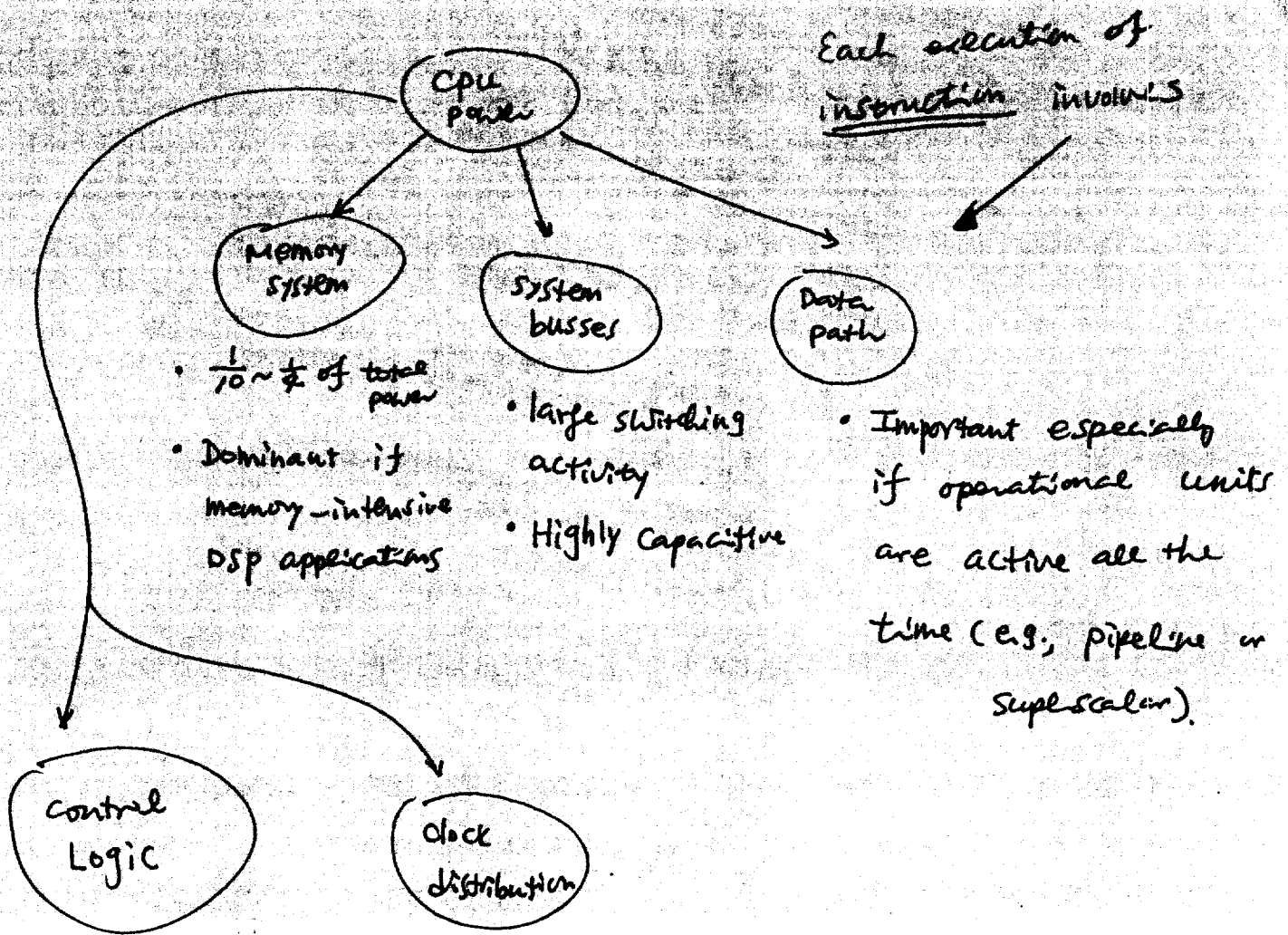
average delay is delay is reduced by  $3\%$ .

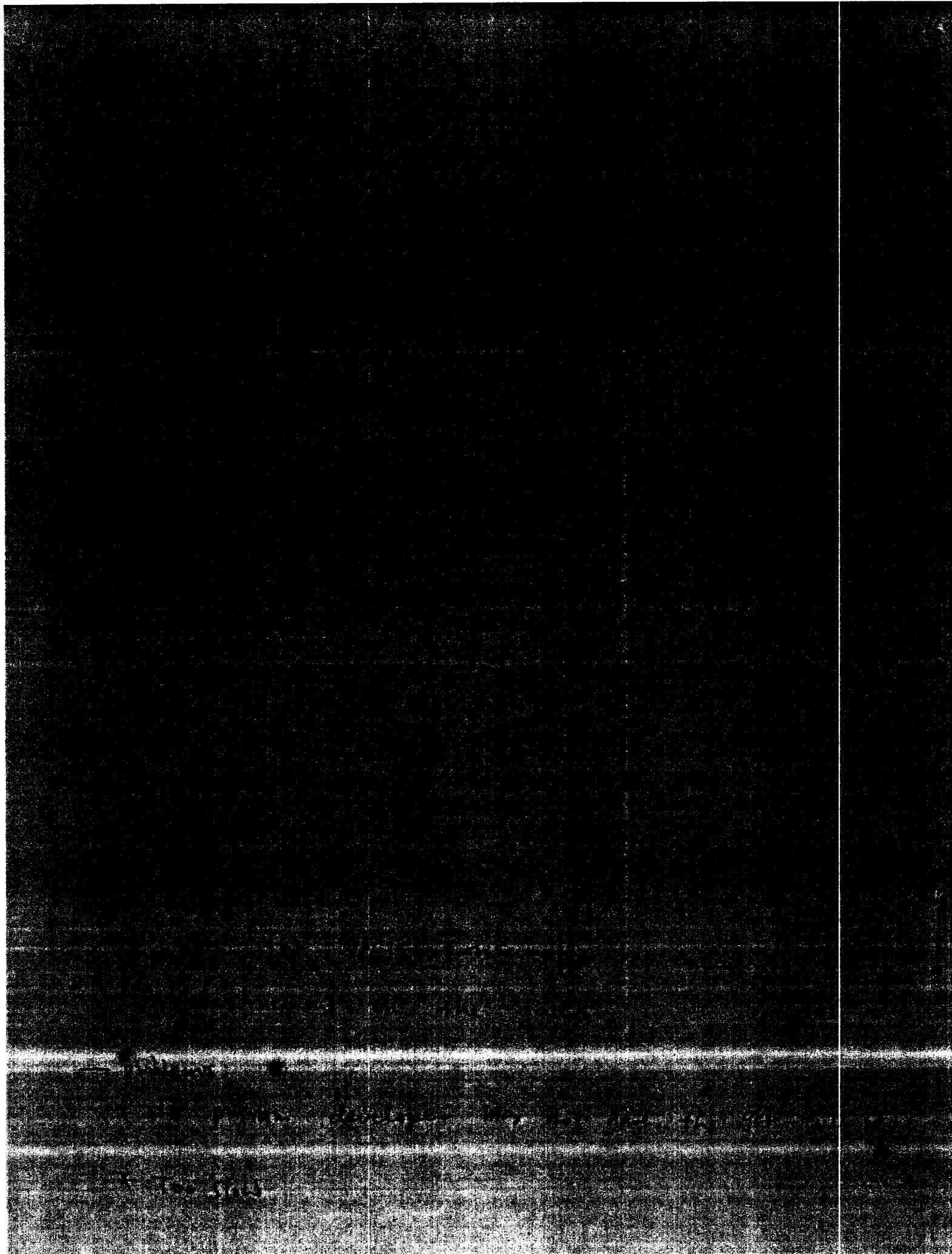


• Major emphases:

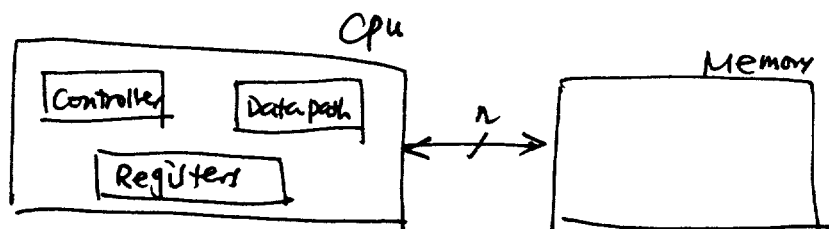
- ① Minimize memory accesses
- ② optimal selection and sequencing of machine instructions
- ③ exploit the low power features of some processors.

• Sources of Software power dissipation





• Architectural-Level power Estimation.



Model: ① power dissipation at each major component

② Which components will be activated by the program.

- Less precise, but much faster than gate-level estimation.

• Bus Switching Activity

- use bus activity to represent the overall switching activity in a processor

- Requires knowledge of:

\* Bus architecture, \* op-codes for the instruction set,

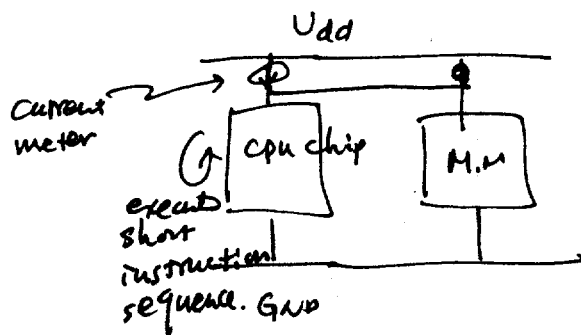
\* representative input data to a program,

\* program code and data mapping to M.M.

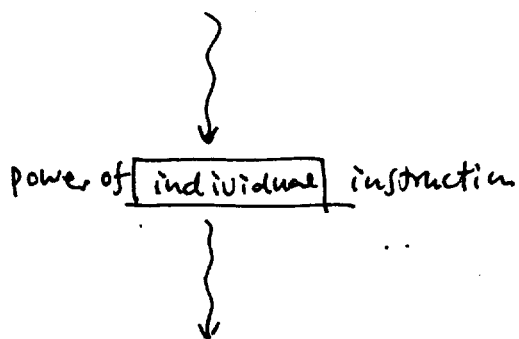
\* Simulation for switching at different buses.

• Instruction-level power analysis.

- use an empirical method for power characterization of short instruction sequences.
- Then, use the results to estimate the power dissipation of a program.
- Basic idea

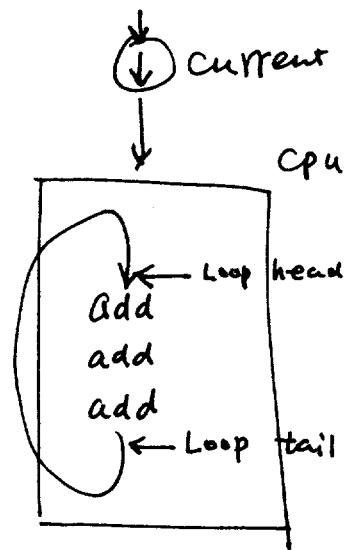


- Choice of instruction sequences for characterization is critical
- Try to estimate the Base cost of each individual instruction.
- Base cost: The portion of power that is independent of the prior state of the processor.



- Avoid estimating Pipeline Stalls, Cache misses, bus switching ... due to consecutive instructions - - - - .

• How ?



Give several instances of the instruction and execute an infinite loop.

• The loop should be as long as possible

minimize estimation error due to loop overhead.

• The loop cannot be too long.

Avoid cache misses.

• <sup>power</sup>~~Energy~~: ~~the~~ power supply voltage x average current draw.

• The effect of prior processor state should also be considered.

Example: Pipeline stall, buffer stall, cache misses.

Called Circuit state effect.

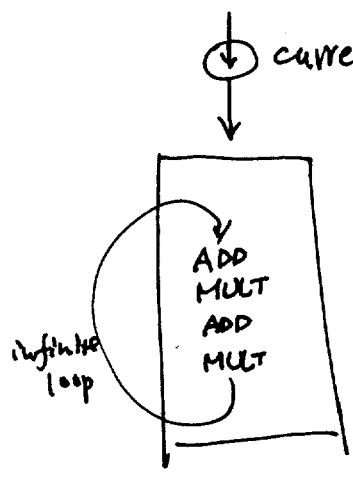
```

ADD  A ← B, C
MULT C ← A, B

```

- ① (op-code) Instruction bus ~~change~~ <sup>switch</sup>
- ② switch of control signals
- ③ mode change within ALU
- ④ Switch of data lines between ALU & registers

• How to measure circuit state effect ?



problem: We are measuring the effect of ADD to MULT. But, it is impossible to separate

ADD → MULT &

MULT → ADD.

• Instruction-level power estimate

$$E_p = \sum_i (B_i \times N_i) + \sum_{i,j} (O_{i,j} \times N_{i,j}) + \sum_k E_k$$

$\uparrow$  overall energy cost of a program  
 $\uparrow$  Base cost for type i instruction  
 $\uparrow$  # of type i instructions  
 $\uparrow$  Circuit state effect of type i instruction followed by type j.  
 $\uparrow$  # of other effects (e.g. stalls, cache misses.)  
 each occurrence of

① — ④ above

• Example:

```

DLOAD  A ← x, B ← y
LOAD   C ← z, MULT D ← A, B
ADD    A ← C, D
    
```

	Base Cost	Circuit state	effect	
		ADD	MULT	DLOAD
DLOAD	2.37	1.19	- - - -	↑ 1.06
LOAD, MULT	2.25			
ADD	0.99			0.99

Table 1 of pp. 439

∞

	Base Cost	Circuit state	
DLOAD	2.37	1.19	Each instruction requires 25ns.
LOAD-MULT	2.25	1.06	
ADD	0.99	0.99	
<hr/> total	5.61	3.24	

→ total ~~power~~ <sup>Energy</sup> = 8.83 PJ.

Ave. power = 8.83 PJ / 75ns = 118 uW.

• Automatic analysis process:

- ① Basic blocks are extracted from a program.
- ② power for each block is totalled (Base + circuit state effect)
- ③ The frequency and Cost of stalls for each block are added estimated and into the block.

④ The program profiler determines the # of time each block executes + accumulates total Cost

⑤ A cache simulator is used to estimate the frequency of cache misses, and add the power to the total program power. ①73

• Software power optimizations

• By Instruction selection & ordering

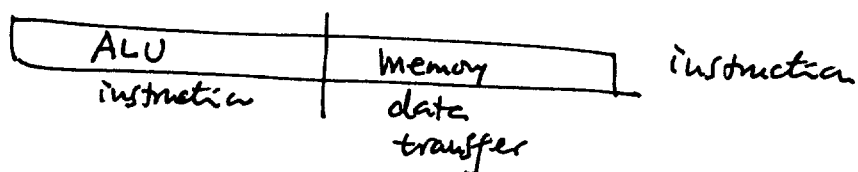
• Cache performance

Large code and small cache  $\longrightarrow$  frequent cache misses  
 $\longrightarrow$  high power penalty

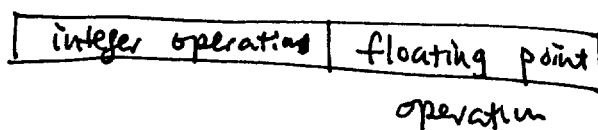
$\therefore$  should maximize code density for embedded processors

• Instruction packing

Example: Fujitsu DSP



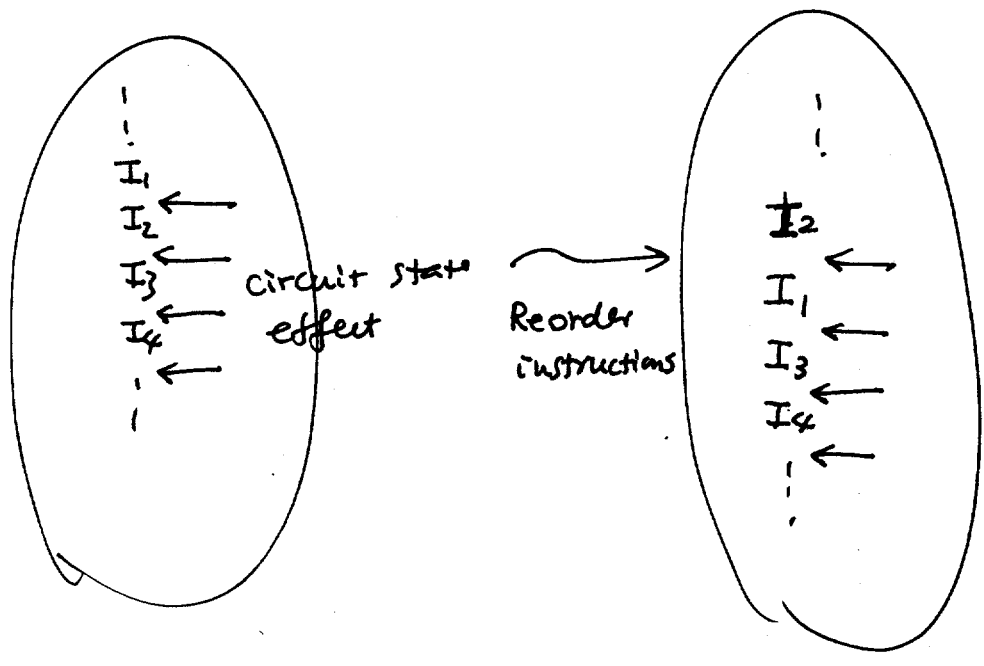
\* 50% of ~~power~~ energy reduced if <sup>^</sup> instruction fetch overhead is not duplicated each





• This method is also used in VLIW & Superscalar architectures.

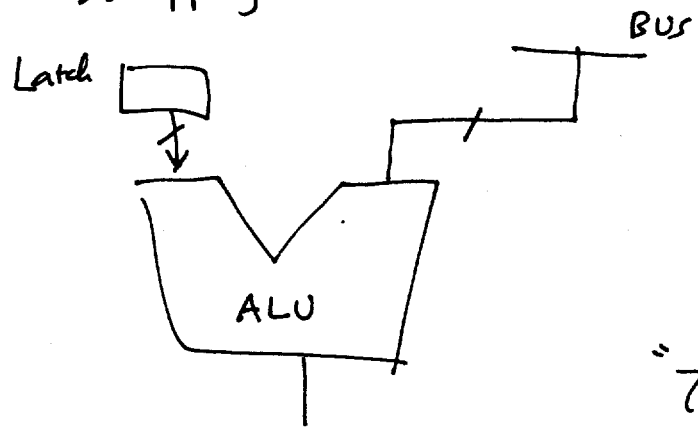
• Instruction ordering



to minimize the circuit state power

\* This technique is more significant for DSP's than general-purpose architectures.

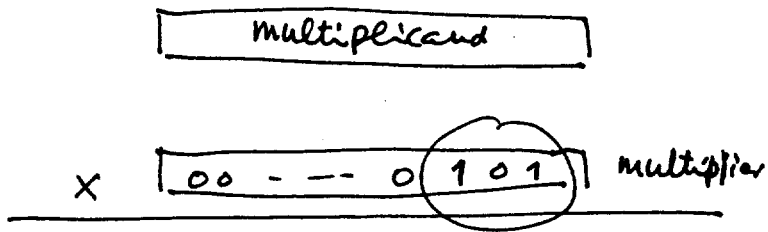
• operand swapping



$x+7$   
 $y+7$

"7" should be placed in the latch.

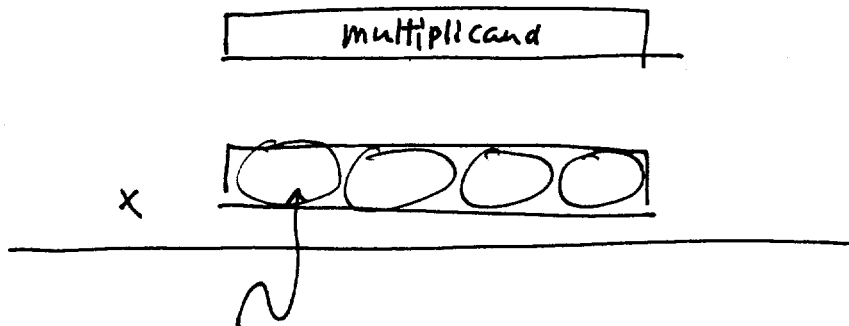
Another example: ~~Booth~~ shift multiplier



only very limited # of add & shift operations are required.

→ Reorder small weight # to the second operand (multiplier)

Booth multiplier



- \* Different bit pattern determines the # of additions & subtractions
- \* The # of additions & subtractions ~~for each bit~~ ~~is~~ is called recoding weight
- \* place the value with low recoding weight to the 2nd position.

# Minimizing Memory Access cost

- Memory is both a power and performance bottleneck.
  - It is slow & power hungry.
- Try to
  - Minimize the # of memory accesses required by an algorithm.
  - Make memory accesses as close to the process as possible.  
 register  $\rightarrow$  cache  $\rightarrow$  RAM
  - Minimize the total ~~memory~~ memory required by an algorithm
  - Make the most efficient use of available memory bandwidth.  
 e.g., try to use multiple word, parallel loads.

## Example:

```

For i=1 to N do
    B[i] = f(A[i]);
for i=1 to N do
    C[i] = g(B[i]);
    
```

$\Rightarrow$

```

For i=1 to N do
    begin
        B[i] = f(A[i]);
        C[i] = g(B[i]);
    end
    
```

↑

compiler can use a register to store B[i].

N memory transfers only.

↑ If B array is large,  
 register  $\longleftrightarrow$  memory  
 transfer  
 Waste power!!  $2N$  memory transfers

# Exploiting Low power Features of Hardware.

(P179)

- Software control over power management (instead of gated clock).

\* Software designer can use instructions to power-down some components

\* Example:

Standby mode: CPU core is stopped, but other operations are maintained.

Sleep mode: All operations except the real time clock stops.

- Advantage:

- Software designer and compiler can better determine when to remove or slow down a clock.
- pure hardware power down may make incorrect decision (based on processor activity) → hurt performance & power ∴ system restoring cost.

- Disadvantage:

needs program execution cycles,  
more costly than hardware-controlled power-down.