

Design of Parallel & High Performance Computing
Reasoning about Performance II

Roofline Model - Balance Principles

5. Roofline Model (Williams et al. 2008)

resources in a microarchitecture
that bound performance:

- peak performance π [flops/cycle]
- memory bandwidth β [bytes/cycle]

associated program
features

- work W [flops]
- data transferred
cache \leftrightarrow mem Q [bytes]

for a given problem:

$$\min W = \text{work/time complexity}$$
$$\min Q = \text{I/O complexity}$$

Operational intensity I : Given a program, assume
empty cache:

$$I(u) = W(u)/Q(u)$$

Intuition: high $I \leftrightarrow$ compute bound
low $I \leftrightarrow$ memory bound

Example:
asymptotic bounds
on I

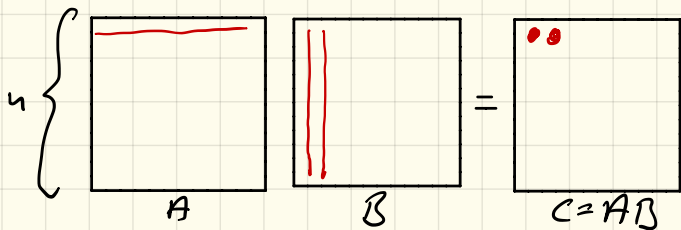
vector sum	$y = x + y$	$O(1)$
matrix-vector product	$y = Ax$	$O(1)$
fast Fourier transform		$O(\log u)$
matrix-matrix product	$C = AB$	$O(u)$

Operational Intensity: Example Matrix Multiplication

Assumptions: cache size $\ll n$, cache block = 8 doubles, 1 cache

We want to estimate Q . $W = 2n^3$ flops

1.) Triple loop:



1. entry of C : $n + 8n$ doubles

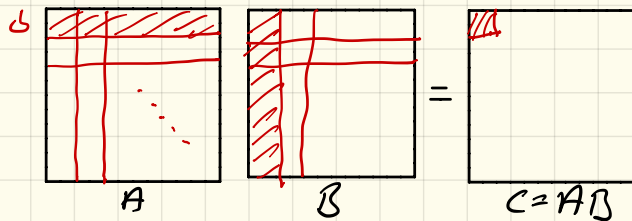
2. entry of C : same

...

total: $9n^3$ doubles

$$I(n) = O(1)$$

2.) Blocked: $8/b$, $3b^2 \leq \gamma$



1. block of C : $2nb$ doubles

2. block of C : same

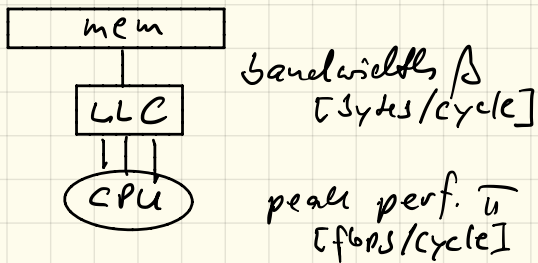
...

total: $2nb \cdot \left(\frac{n}{b}\right)^2 = \frac{2n^3}{b}$ doubles

$$b = \sqrt{\frac{\gamma}{3}} \Rightarrow I(n) = O(\sqrt{\gamma})$$

Roofline Model (Williams et al. 2008)

Computer:



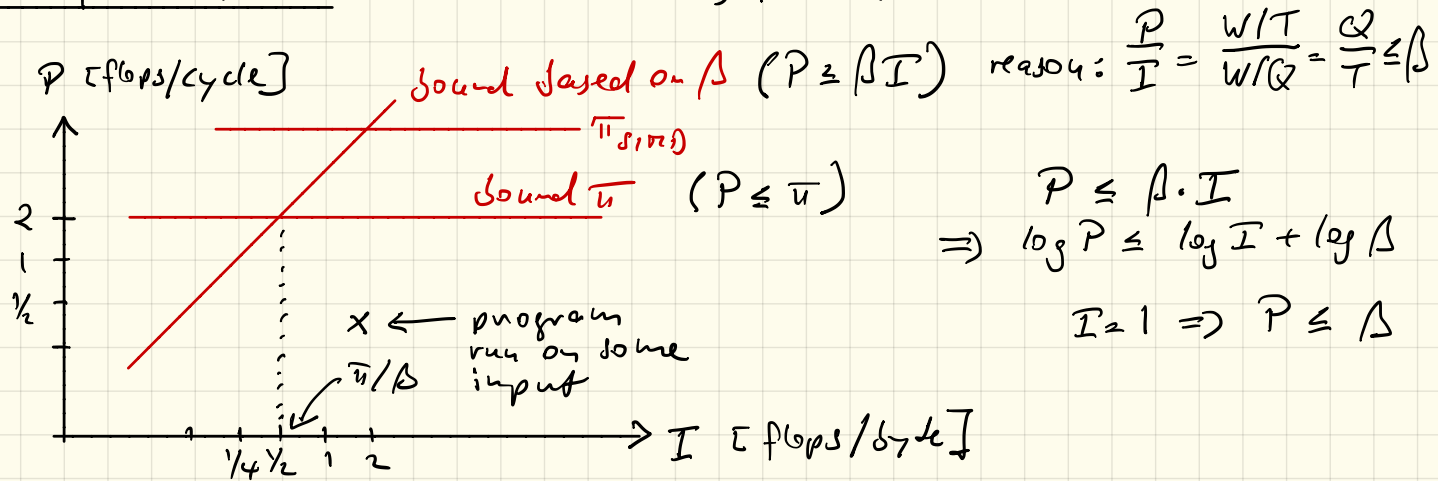
Program:

$$I = W/Q \quad [\text{flops/byte}]$$

$$T = \text{runtime} \quad [\text{cycles}]$$

$$P = W/T \quad (\text{performance}) \quad [\text{flops/cycle}]$$

Roofline plot: (example $\bar{u} = 2$, $\beta = 4$)



$$P \leq \beta \cdot I$$

$$\Rightarrow \log P \leq \log I + \log \beta$$

$$I \geq 1 \Rightarrow P \leq \beta$$

Operational Intensity: Upper bounds

$$I(u) = W(u) / Q(u)$$

#flops / #bytes
mem \leftrightarrow cache

x, y : length n , $A, B, C = n \times n$, α : scalar

		asymptotic	best	exact
$daxpy$	$y = \alpha x + y$	$O(1)$	$O(1)$	$\leq 1/12$
$dgemv$	$y = Ax + y$	$O(1)$	$O(1)$	$\leq 1/4$
fft		$O(\log n)$	$\Theta(\log n)$	—
$dgemm$	$C = AB + C$	$O(n)$	$\Theta(\sqrt{n})$	$\leq \frac{n+1}{16}$

$n =$ cache size

should be sharp if $3n^2 \cdot 8 \leq n$

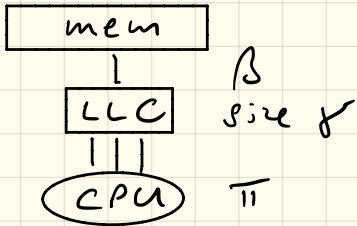
Example: $n = 12773$

$\Rightarrow n \leq 700$

back to slides

B. Balance Principles I (Kung 86)

Computer:



Program:

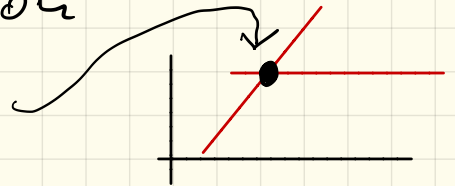
Q_γ : data transfer mem \leftrightarrow LLC

W : work = # ops

The computer is called **balanced** if
compute time = data transfer time

i.e., assuming perfect utilization

$$\frac{W}{\pi} = \frac{Q_\gamma}{\beta} \iff \frac{\pi}{\beta} = \frac{W}{Q_\gamma}$$



assume $\pi \rightarrow a\pi$ increases; how to rebalance?

a.) $\beta \rightarrow a \cdot \beta$ ✓

b.) $\gamma \rightarrow x \cdot \gamma$

but β grows slower than π !
what is x ?

Example 1: Matrix multiplication $C = AB + C$
algorithm with optimal $I = \frac{W}{Q_F} = \Theta(\sqrt{r})$

$$\frac{\overline{W}}{B} = \frac{W}{Q_F} = \Theta(\sqrt{r}) \quad \overline{w} \rightarrow a\overline{w}, \quad r \rightarrow a^2 \cdot r$$

Example 2: FFT / Sorting
algorithm with optimal $I = \Theta(\log r)$

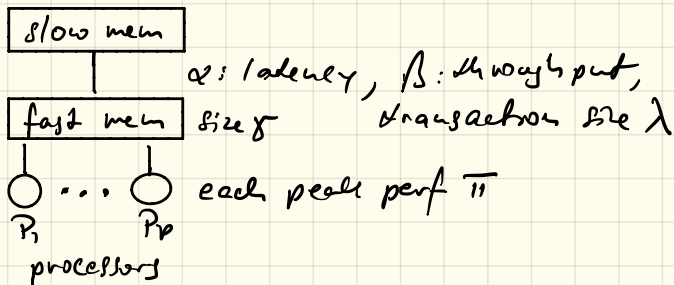
$$\frac{\overline{W}}{B} = \frac{W}{Q_F} = \Theta(\log r) \quad \overline{w} \rightarrow a\overline{w}, \quad r \rightarrow r^a$$

Both are unrealistic.

7. Balance principles II (Czechowski et al. 2011)

Goal: more detailed principles for multicores,
algorithm/architecture co-design, assessment of HW trends

Computer:



Algorithm:

PRAM: W : work, D : depth

$Q_{p,\lambda}$: mem. transfers
of size λ

assumptions: optimal W , W/Q

How to get $Q_{p,\lambda}$ from Q_{λ} ?

- there are general bounds
and some direct results

Processor is balanced if

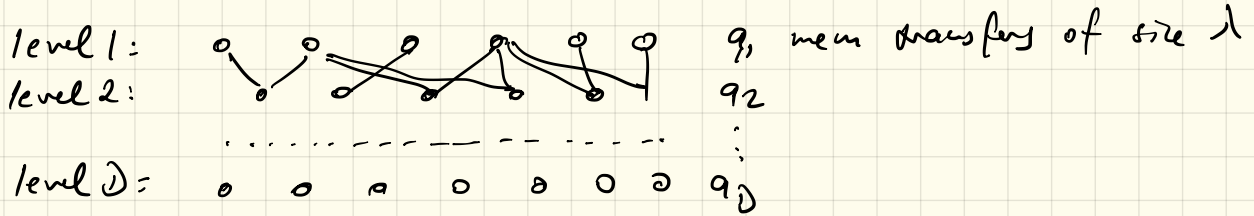
$T_{mem} \leq T_{comp}$ ("compute-bound")

Using:

$$\frac{Q \cdot \lambda}{\beta} \leq \frac{W}{p \cdot \pi} \Leftrightarrow \frac{p \pi}{\beta} \leq \frac{W}{Q \lambda}$$

Derivation principles:

1.) Estimate T_{mem} : Idea divide DAG in levels



$$\Sigma = Q_{\lambda, \lambda}$$

In each level α - β model: $\alpha + \frac{q_{i-1} \cdot \lambda}{\beta}$

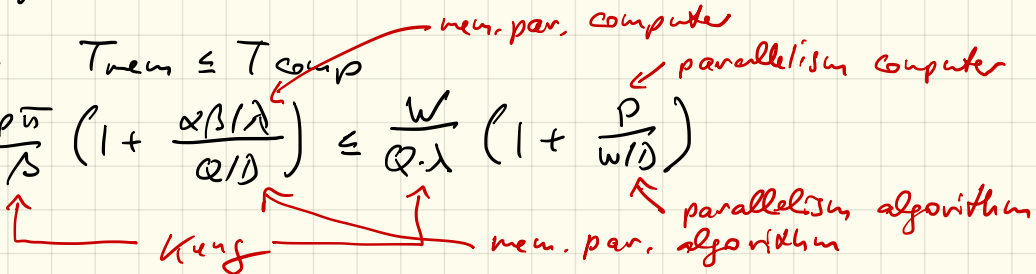
$$\Rightarrow T_{mem} \approx \sum_{i=1}^D \left(\alpha + \frac{q_{i-1} \cdot \lambda}{\beta} \right) = \alpha \cdot D + \frac{Q}{\beta} \quad Q = Q_{\lambda, \lambda}$$

2.) Estimating T_{comp} : Brent's theorem

$$T_{comp} \approx \left(D + \frac{W}{P} \right) \cdot \frac{1}{\pi}$$

Balance: $T_{mem} \leq T_{comp}$

$$\Leftrightarrow \frac{p \bar{u}}{\beta} \left(1 + \frac{\alpha \beta \lambda}{Q/D} \right) \leq \frac{W}{Q \cdot \lambda} \left(1 + \frac{P}{W/D} \right)$$



Example 1: matrix multiplication

$$\text{use } Q \geq \frac{W}{\sqrt{L} \cdot \lambda \sqrt{r p}} \quad (\text{Irony et al. 2004}) \quad \Rightarrow I(u) = O\left(\sqrt{\frac{r}{p}}\right)$$

$$\Rightarrow \text{balance principle } \frac{p u}{\beta} \in O\left(\sqrt{\frac{r}{p}}\right)$$

$$\begin{array}{l} \pi \rightarrow a \cdot \pi \\ p \rightarrow a \cdot p \end{array} \quad \text{rebalance: } \begin{array}{l} \beta \rightarrow a \beta \text{ on } \gamma \rightarrow a^2 \gamma \\ \text{": } (\beta \rightarrow a \beta \text{ and } \gamma \rightarrow a \gamma) \text{ on } (\gamma \rightarrow a^2 \gamma) \end{array}$$

Example 2: bounding FFT

$$\Rightarrow \frac{p u}{\beta} \leq O\left(\log \frac{r}{p}\right)$$

Back to slides