

Little's Law (1961)

Example 1: In a Starbucks, on average

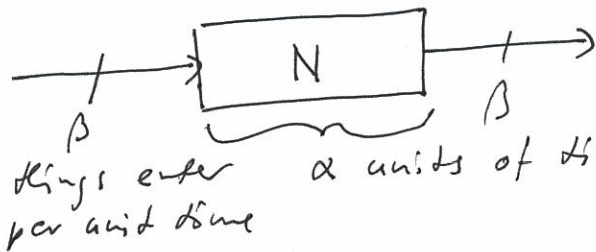
- every minute 2 customers enter and leave
 - every customer spends 8 minutes in the store
- How many people are inside? $16 = 2 \times 8$

Example 2: In your wine cellar, on average

- there are 600 bottles
 - you drink and buy 50/year
- How long is every bottle in the cellar? $12 = 600/50$

Little's Law: Given a stable (arrival rate = departure rate) system

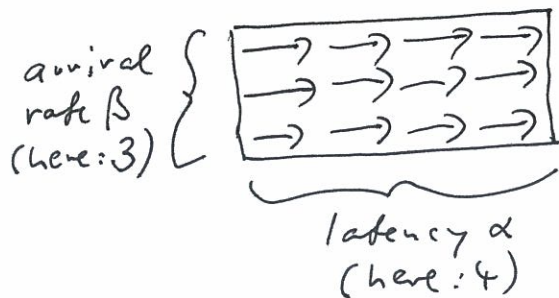
$N = \# \text{ things in system}$



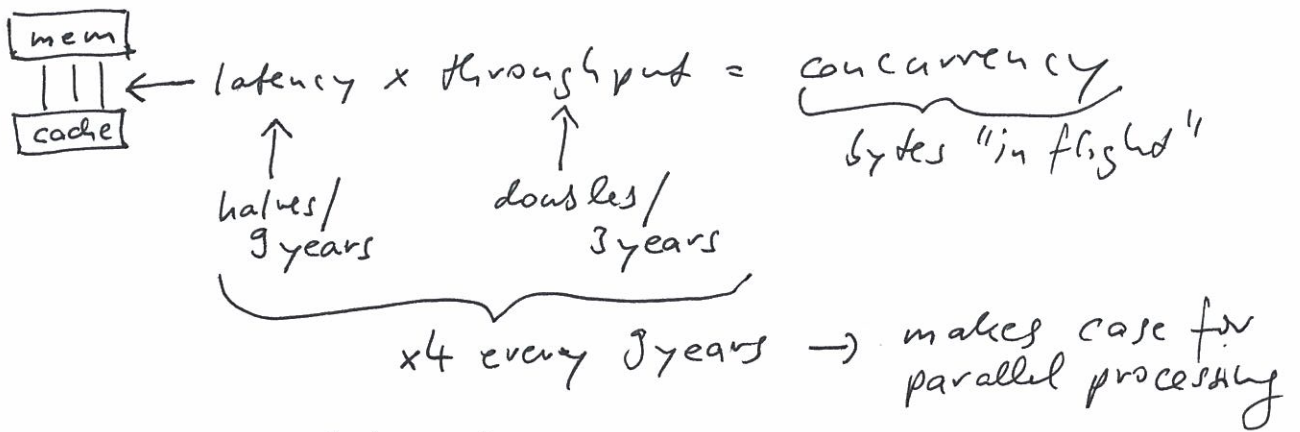
Then: $N = \alpha \beta$

Seems trivial, but the interesting aspect is the independence of input/output distribution.

Useful visualization:



Example: Memory system



Core 2: $\beta = 2 \text{ B/cycle}$
 $\alpha \approx 100 \text{ B/cycle}$
 $\Rightarrow N = 200 \text{ B} = 50 \text{ floats or ints}$

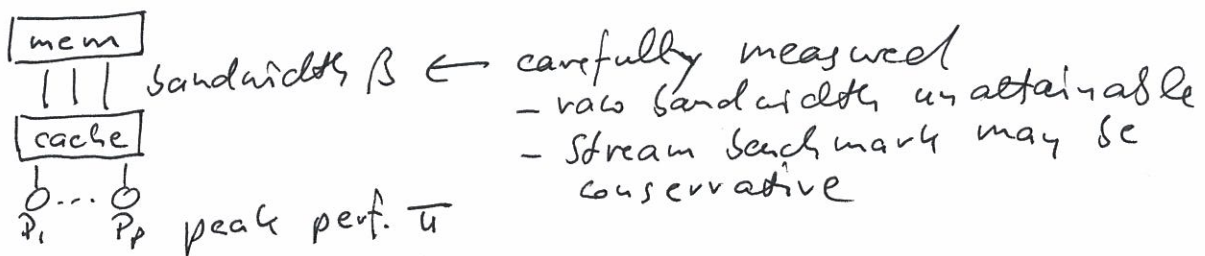
Performance Bounds: Roof line model

(Williams et al. 2008)

Which resources in a microarchitecture bound performance?

- peak performance [ops/cycle]
- memory bandwidth [bytes/cycle]
- <other>

Platform model:



Algorithm/program model:

$$\text{operational intensity } \rho = \frac{W}{Q} = \frac{\# \text{ ops}}{\# \text{ bytes transferred mem} \leftrightarrow \text{ cache}} \left[\frac{\text{ops}}{\text{byte}} \right]$$

"ops" could be

- floating point ops (adds/mults) for numerical code
- comparisons (for sorting)
- <other>

W: derive on paper or instrument your code or perf. counters

Q: perf. counters

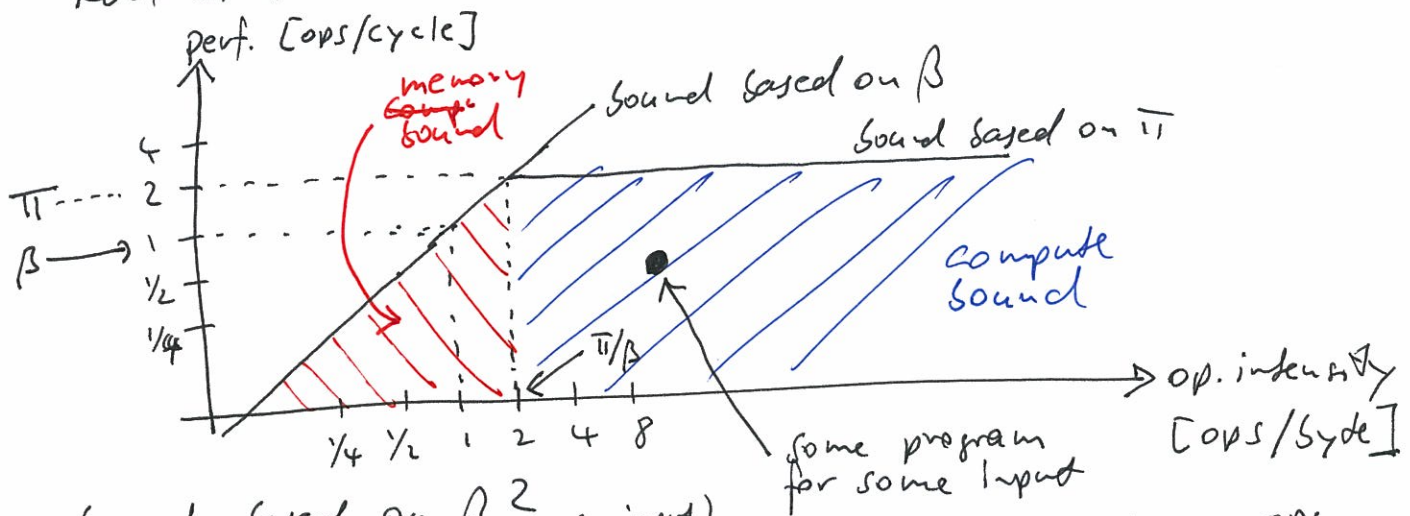
[roofline model assumes cold cache]

For some functions, asymptotic results on optimal op. intensity are available.

Assume code size of m bytes. Then the optimal $W/Q = W/Q_m$ is

FFT/sorting: $\Theta(\log m)$ (m-way mergesort)
 Matrix mult.: $\Theta(\sqrt{m})$ (suitably blocked)

Roofline model: (example $\pi = 2, \beta = 1$)



Sound based on β ? (on input)

- assume program has op. intensity of $x \frac{\text{ops}}{\text{byte}}$
- the program gets at most $\beta \frac{\text{bytes}}{\text{cycle}}$

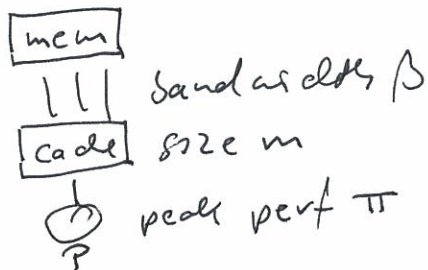
$$\Rightarrow \text{perf} = y \leq \beta \cdot x$$

$$\Leftrightarrow \log_2(y) \leq \log_2(x) + \log_2(\beta) \quad (\text{slope is } = 1)$$

$$x=1 \Rightarrow y \leq \beta$$

Balance principles I (Kung 86)

platform model:



algorithm model:

Q_m : data transfer
mem \leftrightarrow code

W : work = #ops

The processor is called balanced for an algorithm (on an input) if compute time = data transfer time, i.e., assuming perfect utilization:

$$\frac{W}{\pi} = \frac{Q_m}{\beta} \iff \frac{\pi}{\beta} = \frac{W}{Q_m} \quad (\text{compare to roofline model!})$$

If π is increased by $a \times$ ^{"times"}, β can also be increased by $a \times$ to rebalance. But π grows faster than β !

Hence this question: If $\pi \rightarrow a\pi$, how much to increase m to rebalance?

Consider algorithms with optimal $\frac{W}{Q_m}$

Case 1: Matrix mult.

$$\frac{W}{Q_m} = \Theta(\sqrt{m})$$

$\Rightarrow \frac{\pi}{\beta}$ increased $\times a \rightarrow m$ needs increase $\times a^2$

Case 2: FFT/sorting

$$\frac{W}{Q_m} = \Theta(\log m)$$

$\Rightarrow \frac{\pi}{\beta}$ increased $\times a \rightarrow m$ needs increase $\times a$

Both are unrealistic.

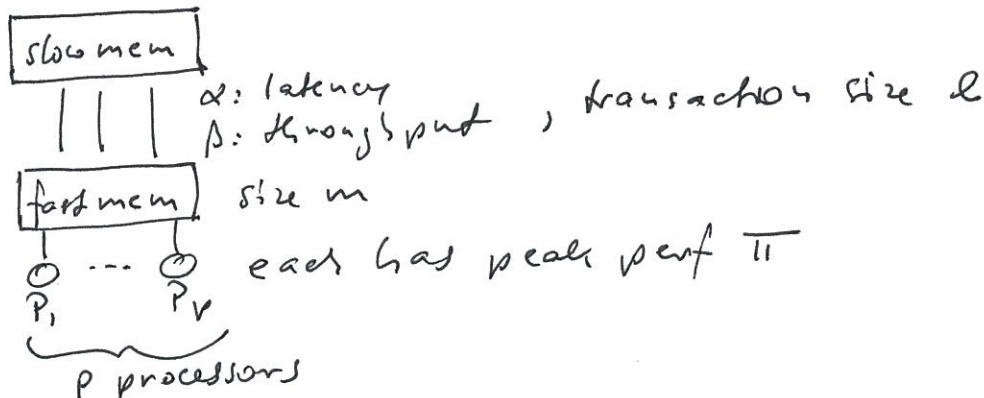
Balance principles II (Czeczowski et al. 2011)

Goal: More detailed balance principles for multi cores

Applications:

- relationship between hardware and algorithm parameters
- reasoning about architecture needs

Platform model:



Algorithm model:

PRAT: $W(n)$, $D(n)$

Memory transfers of size l : $Q_{m,e}(n)$

Assumptions: - optimal $W(n)$
- optimal $W(n)/Q_{m,e}(n)$

Question: How to get $Q_{p,m,e}(n)$ (i.e., for p processors)

Answer: - there are general bounds based on $Q_{m,e}(n)$ and a given scheduler
- some direct results are available

~~Def~~: Processor is balanced for algorithm if

$$T_{mem} \leq T_{comp} \quad (\text{Kung used "="})$$

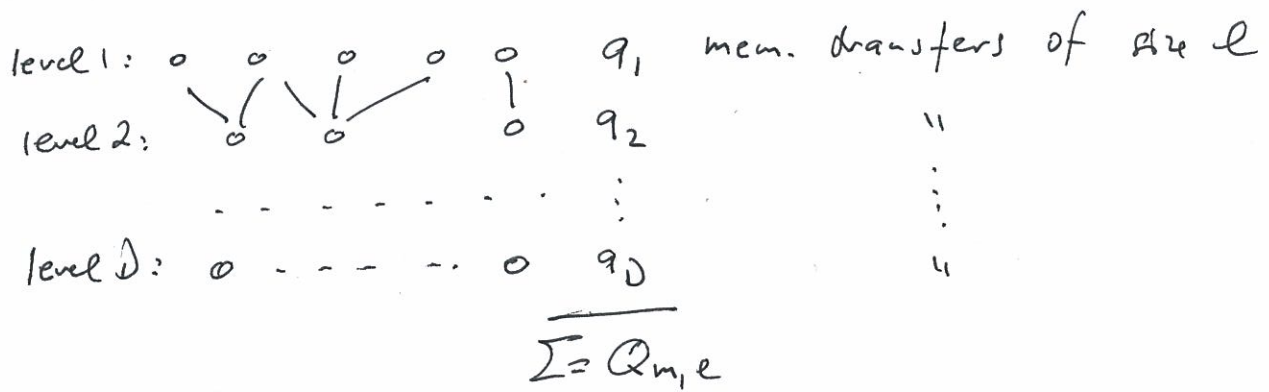
This is equivalent to "compute-bound".

Kung would get:

$$\frac{Q \cdot l}{\beta} \leq \frac{W}{p \pi} \iff \frac{p \pi}{\beta} \leq \frac{W}{Q \cdot l}$$

Estimating T_{mem} :

Idea: divide DAG into levels



In each level use α - β model: $\alpha + \frac{q_i \cdot l}{\beta}$

$$\Rightarrow T_{mem} \approx \sum_{i=1}^D \left(\alpha + \frac{q_i \cdot l}{\beta} \right) = \alpha D(n) + \frac{Q_{p,m,e}(n)}{\beta}$$

Estimating T_{comp} : Brent's theorem

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

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

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

Annotations for the balance principle equation:
 - $\frac{\alpha \beta l}{Q \cdot D}$: mem. parallelism (Kung)
 - $\frac{p}{W \cdot D}$: comp. parallelism of algorithm

Example: Matrix mult.

$$Q \geq \frac{W}{\sqrt{2} \cdot l \sqrt{m/p}}$$

(Irony et al. 2004)

$$\Rightarrow \text{Balance principle } \frac{p \bar{u}}{\beta} \leq \sqrt{\frac{m}{p}}$$

Doubling p requires doubling both β and m
 or $m \rightarrow 8m$ to rebalance

Note: doubling p is very different than doubling π
 (even though both double peak perf.)

Other example: Comparison-based sorting (+ work stealing)

$$\frac{p \pi}{\beta} \leq O\left(\log \frac{m}{p}\right)$$