- Bandwidth or throughput
  - Total work done in a given time
  - 10,000-25,000X improvement for processors
  - 300-1200X improvement for memory and disks

- Latency or response time
  - Time between start and completion of an event
  - 30-80X improvement for processors
  - 6-8X improvement for memory and disks

---

## Performance

Comparing Machines/Systems

- Response Time (latency) $= T$
  - How long does it take for my job to run?
  - How long does it take to execute a job?
  - How long must I wait for the database query?

avg/best case/worst case

- Throughput $= \mathcal{V}$
  - How many jobs can the machine run at once?
  - What is the average execution rate?
  - How many queries per minute?

avg

What do we "really" want to know?

--- Which system works best in our larger system?

--- What costs can be traded off?

Time?

- Elapsed Time
  - Counts everything *(disk and memory accesses, I/O , etc.)*  →→ "wall clock"
  - A useful number, but often not good for comparison purposes
    - E.g., OS & multiprogramming time make it difficult to compare CPUs

Depends on load, disk layout, ...

more abstract

- CPU time (CPU = Central Processing Unit = processor)
  - Doesn't count I/O or time spent running other programs
  - Can be broken up into system time and user time  → user cpu time

  OS
  I/O

- Our focus: user CPU time                    → Time CPU used for our job (+ overhead)
  - Time spent executing the lines of code that are "in" our program
  - Includes arithmetic, memory, and control instructions...

  unix  ⇒  localhost >   time my job
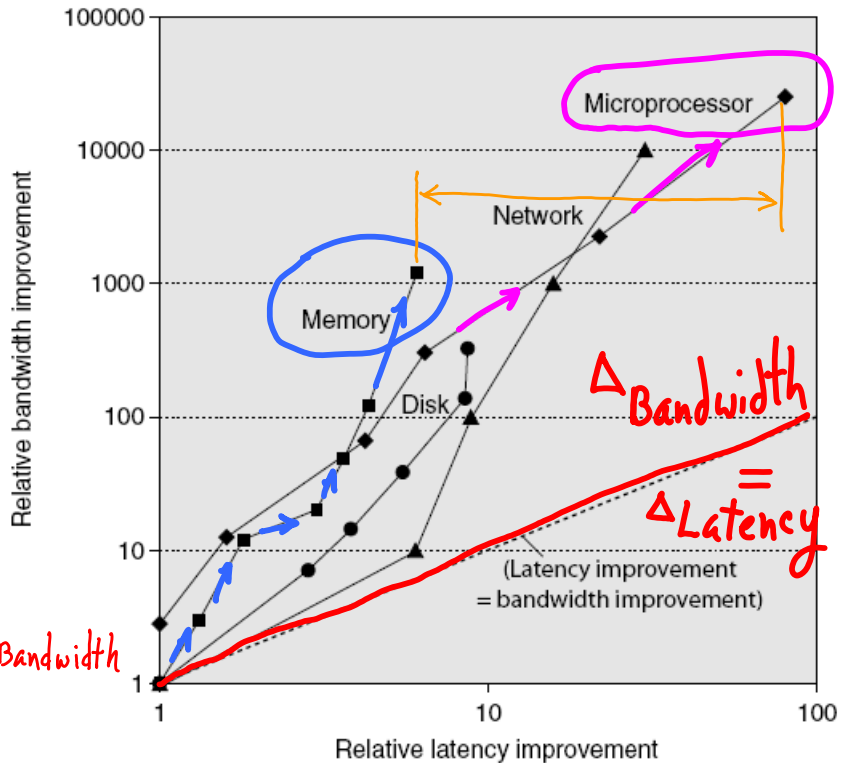
Latency vs. Bandwith
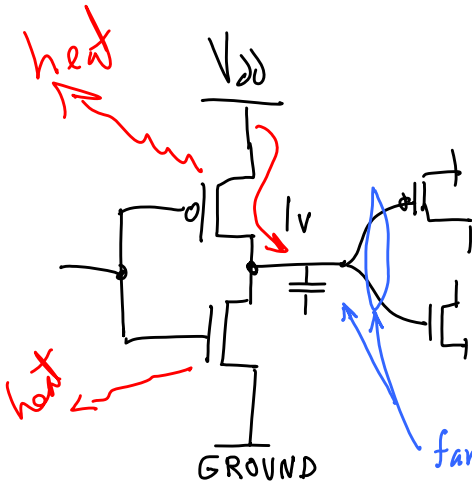
Increasing complexity/density/speed

1. V / L gets worse for each component.

2. Vproc / Lmem gets worse even faster

relative performance



Relative bandwidth improvement (y-axis: 1 to 100000)
Relative latency improvement (x-axis: 1 to 100)

Microprocessor
Network
Memory
Disk

(Latency improvement = bandwidth improvement)

$\Delta$Bandwidth = $\Delta$Latency

$$\text{Log } \frac{V_{new}}{V_{old}} = \Delta_{Bandwidth}$$

$$\text{Log } \frac{L_{old}}{L_{new}} = \Delta_{Latency}$$

heat

$V_{dd}$

Iv

hot

GROUND

fanout
= capacitive load, changing capacitance up to $V_{dd}$
or draining to GROUND.

- Dynamic energy
  - Transistor switch from 0 -> 1 or 1 -> 0
  - ½ x Capacitive load x Voltage²

- Dynamic power
  - ½ x Capacitive load x Voltage² x Frequency switched

- Reducing clock rate reduces power, not energy

$$f_{max} = \propto V$$

Rising clock edge
falling clock edge
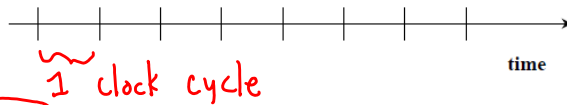
CLOCK:   0   1   0   1

one cycle, Tcycle

→ time

cpu   **Clock Cycles** ⇒ cpu time

---

- Instead of reporting execution time in seconds, we often use cycles

$$\text{CPU Time} = \left(\frac{\text{seconds}}{\text{program}}\right) = \left(\frac{\text{cycles}}{\text{program}}\right) \times \left(\frac{\text{seconds}}{\text{cycle}}\right)$$

$$T_{cycle} = \left(\frac{seconds}{Tick}\right)$$

- Clock "ticks" indicate when to start activities:



1 clock cycle                                           time

$$Freq = \left(\frac{Ticks}{sec}\right) = \frac{1}{T_{cycle}}$$

- Cycle time = time between ticks = seconds per cycle
- Clock rate (frequency) = cycles per second  (1 Hz. = 1 cycle/sec)

2 GHz clock ⇒ $Freq = \left(\frac{2 \times 10^9\ ticks}{sec}\right)$ ⇒ $T_{cycle} = \left(\frac{1\ sec}{2 \times 10^9\ Ticks}\right)$

$$= \tfrac{1}{2}\ ns$$
$$= \tfrac{1}{2}\ (1,000\ ps)$$
$$= 500\ ps$$
↖ $10^{-12}$ sec

$$CR \overset{def}{\equiv} freq$$

---

- User CPU execution time         ※ cycles × $T_{cycle}$         rewrite

$$\text{Execution Time} = \boxed{\text{Clock Cycles}}\ \text{for Pr}ogram \times \boxed{\text{Clock Cycle Time}}$$

- Since Cycle Time is 1/Clock Rate (or clock frequency)

$$\text{Execution Time} = \frac{\text{Clock Cycles for Pr}ogram}{\text{Clock Rate}} = ※cycles\left(\frac{1}{T_{cycle}}\right)^{-1}$$

$$= \frac{※cycles}{CR}$$

- The program should be something real people care about
  - Desktop: MS office, edit, compile
  - Server: web, e-commerce, database
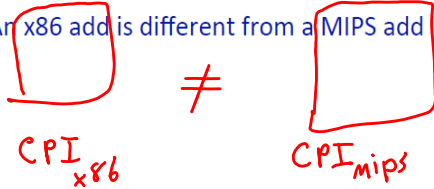  - Scientific: physics, weather forecasting

real jobs?
or
benchmarks?

# Measuring Clock Cycles

- Clock cycles/program is not an intuitive or easily determined value, so

$$\text{※ } \textit{Clock Cycles} = \textit{Instructions} \times \overbrace{\textit{Clock Cycles Per Instruction}}^{average}$$

$$= \text{※ } instr \times \left(\overline{\frac{\text{※} cycles}{instr}}\right)$$

$$= \sum_{i=1}^{\text{※} instr} \text{※} cycles_i$$

↑ each instruction is different

$$\begin{bmatrix} lc3: \\ ADD \ vs \ RTI \end{bmatrix}$$

- Cycles Per Instruction (CPI) used often

- CPI is an **average** since the number of cycles per instruction varies from instruction to instruction
    - Average depends on instruction mix, latency of each inst. type etc.

- CPIs can be used to compare two implementations of the same ISA, but is not useful alone for comparing different ISAs
    - An x86 add is different from a MIPS add

$$\Box \neq \Box$$
$$CPI_{x86} \qquad CPI_{mips}$$

$$\left(\text{※} instr \times \overline{CPI}\right) \times \left(\tfrac{1}{CR}\right)$$

Depends on
--- instruction mix,
--- system configuration,
--- data

$$Time = \text{※} cycles \times T_{cycle}$$

- Drawing on the previous equation:

$$\textit{Execution Time} = \left(\textit{Instructions} \times \overline{CPI}\right) \times \textit{Clock Cycle Time}$$

$$\textit{Execution Time} = \frac{\textit{Instructions} \times \overline{CPI}}{\textit{Clock Rate}} = \text{※} instr \left(\frac{CPI}{CR}\right)$$
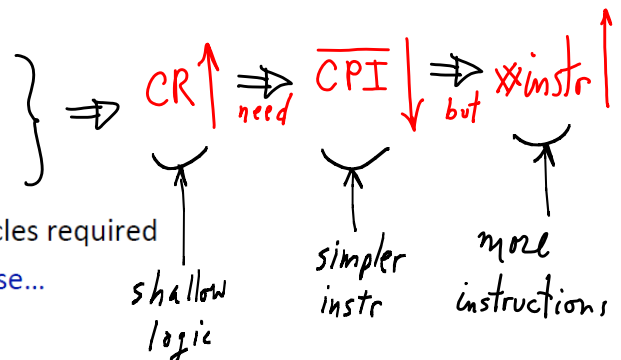
Reg → f → Reg

longest path limits CR

- To improve performance (i.e., reduce execution time)
    - Increase clock rate (decrease clock cycle time) OR
    - Decrease CPI OR
    - Reduce the number of instructions
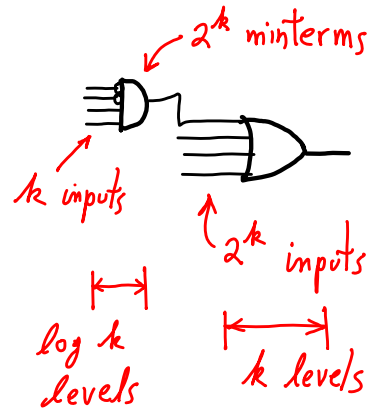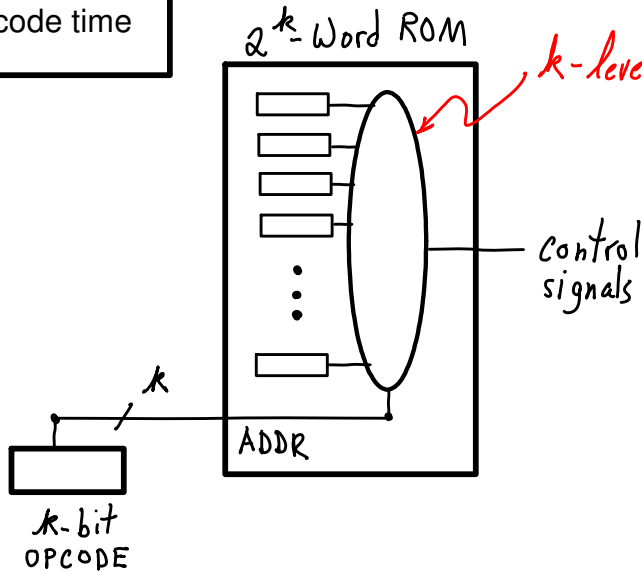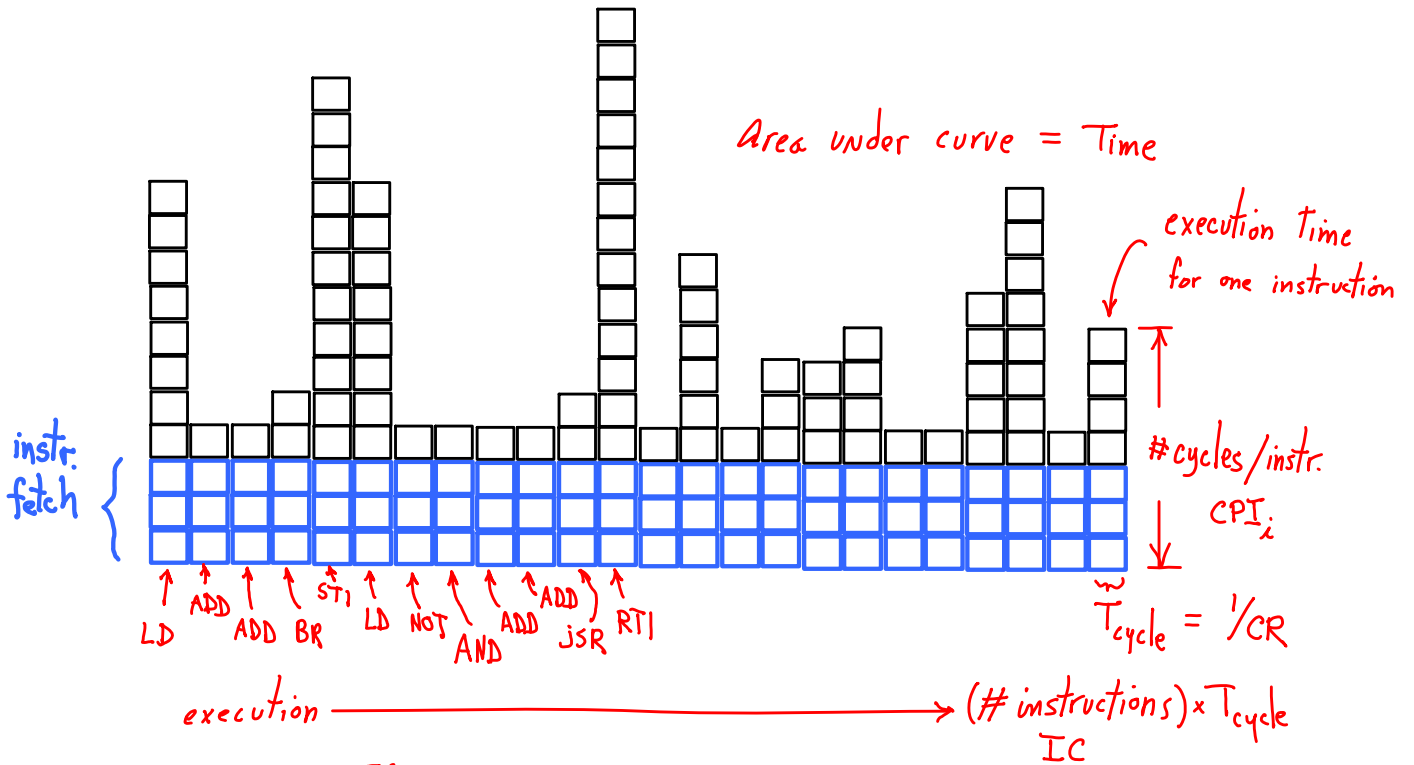- Designers balance cycle time against the number of cycles required
    - Improving one factor may make the other one worse…

reduce Time ↓

$$\left.\begin{array}{l}\end{array}\right\} \Rightarrow \quad CR\uparrow \underset{need}{\Rightarrow} \overline{CPI}\downarrow \underset{but}{\Rightarrow} \text{※} instr \uparrow$$

shallow logic        simpler instr        more instructions

[look for sweet spot]

E.G. Decode time

$2^k$-Word ROM

*k-level logic*

ADDR

control signals

*k*

*k*-bit OPCODE

*$2^k$ minterms*

*k inputs*

*$2^k$ inputs*

*log k levels*

*k levels*

*k small ⟹ few instructions*

*Area under curve = Time*

*execution time for one instruction*

instr. fetch {

*#cycles/instr.*

*$CPI_i$*

↑ ↑ ↑ ↑ $ST_i$ ↑ ↑ ↑ ↑ $ADD$ ↑ ↑
LD ADD ADD BR LD NOT AND ADD JSR RTI
      AND

*$T_{cycle} = 1/CR$*

execution ⟶ (# instructions) × $T_{cycle}$
IC

$$Time = Area = \sum_{i}^{IC} CPI_i \times T_{cycle} = IC \, \overline{CPI} \left( \frac{1}{CR} \right)$$

*shorter $T_{cycle}$, more instructions*

*and/or*

*shorter cycle, more $CPI_i$*

*Instruction mix matters!*

instr. fetch {

$\overline{(\#cycle/instr)}$

$\overline{CPI}$

$T_{cycle}$

$$Time = (\#instr) \times T_{cycle} \times (\overline{CPI})$$

$$= IC \left(\frac{1}{CR}\right)(\overline{CPI})$$

average by classes

LD/ST

BR JSR

ALU

$$time = \left[ (IC_1 \times \overline{CPI_1}) + (IC_2 \times \overline{CPI_2}) + (IC_3 \times \overline{CPI_3}) + (OTHER) \right]\left(\frac{1}{CR}\right)$$

--- Get averages by running batches of class?

--- Guessing from architecture?

# Clock Rate ≠ Performance

- Mobile Intel Pentium 4      Vs      Intel Pentium M
  - 2.4 GHz                        1.6 GHz
  - P4 is 50% faster?

- Performance on Mobilemark with same memory and disk
  - Word, excel, photoshop, powerpoint, etc.
  - *But* — Mobile Pentium 4 is only 15% faster

- What is the relative CPI?
  - ExecTime = IC • CPI/Clock rate
  - $ExecTime_M = 1.15\ ExecTime_4$
  - $IC \bullet CPI_M/1.6 = 1.15 \bullet IC \bullet CPI_4/2.4$
  - $CPI_4/CPI_M = 2.4/(1.15 \bullet 1.6) = 1.3$

$$CR_{P4} = \left(\frac{2.4}{1.6} = 1.5\right) CR_{PM}$$

$$\Rightarrow T_{PM} = (1.5)\, T_{P4}\ ?$$

$$But,\quad T_{PM} = (1.15)\, T_{P4}\ !$$

$$Is\ the\ difference$$

$$\overline{CPI}_{PM}\ versus\ \overline{CPI}_{P4}\ ?$$

$$\left\{ T_{P4} = IC_{P4}\left(\frac{CPI_{P4}}{CR_{P4}}\right) \right\}(1.15) = \left\{ T_{PM} = IC_{PM}\left(\frac{CPI_{PM}}{CR_{PM}}\right) \right\}$$

$$\frac{CPI_{P4}}{\left(CR_{P4} = (1.5)\, CR_{PM}\right)}(1.15) = \frac{CPI_{PM}}{CR_{PM}}$$

$$\boxed{\begin{array}{c} \text{Same ISA} \\ IC_{P4} = IC_{PM} \end{array}}$$

$$\frac{CPI_{P4}}{CPI_{PM}} = \frac{(1.5)}{(1.15)} = 1.304\ldots$$

$$CPI_{P4} = 1.304\ CPI_{PM}$$

$$\Rightarrow 30\%\ more\ cycles/instr\ on\ avg\ for\ P4$$

How can that be?

--- same ISA
--- same program + data
--- what is different?

Average by classes.
Average CPI?

$$\left[ \left( IC_1 \times \overline{CPI_1} \right) + \left( IC_2 \times \overline{CPI_2} \right) + \left( IC_3 \times \overline{CPI_3} \right) \right] \left( \frac{1}{CR} \right) = IC \left( \overline{CPI} \right) \left( \frac{1}{CR} \right)$$

$$\Longrightarrow \qquad \sum_{i=1}^{n} \left( \frac{IC_i}{IC} \right) \overline{CPI_i} = \sum_{i=1}^{n} \left( \%_i \right) \overline{CPI_i} = \overline{CPI}$$

$$\overline{CPI_i} \quad \%_i \implies \left\{ \overline{CPI_i} \left( \%_i \right) \right\}$$

| Instruction Type | CPI | Frequency | CPI * Frequency |
|---|---|---|---|
| ALU | 1 | 50% | 0.5 |
| Branch | 2 | 20% | 0.4 |
| Load | 2 | 20% | 0.4 |
| Store | 2 | 10% | 0.2 |

$$SUM = 1.5 = \overline{CPI}$$

- Given this machine, the CPI is the sum of CPI ✕ Frequency

- Average CPI is 0.5 + 0.4 + 0.4 + 0.2 = 1.5

- What fraction of the time for data transfer?

$$\frac{T_{LD\text{-}ST}}{T_{Total}} = \frac{Cycles_{LD\text{-}ST} * (1/f)}{Cycles_{Total} * (1/f)} = \frac{\#(cycles\ LD\text{-}ST)}{\#\ cycles} = \frac{(\#Cycles\ LD) + (\#cycles\ ST)}{(\#\ cycles) = IC \times CPI}$$

cycle Time,
$f = CR$

$$= \left[ \frac{(\#cycles\ LD)}{IC} + \frac{(\#cycles\ ST)}{IC} \right] \frac{1}{CPI}$$

$$= \left[ \frac{CPI_{LD} \cdot IC_{LD}}{IC} + \frac{CPI_{ST} \cdot IC_{ST}}{IC} \right] \frac{1}{CPI}$$

$$= \left[ CPI_{LD} \cdot \%_{LD} + CPI_{ST} \cdot \%_{ST} \right] \left( \frac{1}{1.5} \right)$$

$$= \left[ 0.4 + 0.2 \right]\left( \frac{1}{1.5} \right) = 40\%$$

## Speedup

- Speedup allows us to compare different CPUs or optimizations

$$Speedup = \frac{CPUtimeOld}{CPUtimeNew}$$

- Example
  - Original CPU takes 2sec to run a program
  - New CPU takes 1.5sec to run a program
  - Speedup = 1.333   or speedup or 33%

$T_{old} = 2\,s$

$T_{new} = 1.5\,s$

What do we mean by speedup?

$$S_{new-old} = \frac{V_{new}}{V_{old}} = \frac{\left( W_{new}/T_{new} \right)}{\left( W_{old}/T_{old} \right)} = T_{old}/T_{new} = 1.3 \quad \Rightarrow \quad V_{new} = (1.3)\ V_{old}$$

$\Rightarrow$ new is 30% faster

Assuming $W_{old} = W_{new}$

$$S = \frac{V_{old}}{V_{new}} = \frac{W/T_{old}}{W/T_{new}}$$

$$W = W_{parallel} + W_{sequential} = f \cdot W + (1-f)W$$

$$T = T_{parallel} + T_{sequential}$$

$$\left( \begin{array}{l} parallel = improvable \\ sequential = fixed \end{array} \right)$$

$$V_p = W_p / T_p \quad V_s = W_s / T_s$$

$$V_p = V_s \quad \text{\color{blue}{old}}$$

$$V_p = a V_s \quad \text{\color{blue}{new}} \quad (i.e.\ S_p = a)$$

### Amdahl's Law

---

- If an optimization improves a fraction $f$ of execution time by a factor of $a$

$$Speedup = \frac{Told}{[(1-f)+f/a]*Told} = \frac{1}{(1-f)+f/a}$$

Is f fixed? $\quad W_p = \mathcal{O}(n) \quad W_s = \mathcal{O}(1)$

- This formula is known as Amdahl's Law
- Lessons from
  - If f→100%, then speedup = a
  - If a→∞, the speedup = 1/(1-f)   ← target speed up to max %
- Summary
  - Make the (common case) fast
  - Watch out for the non-optimized component

$$T_{old} = \frac{f \cdot W}{V_s} + \frac{(1-f)W}{V_s}$$

$$= W / V_s$$

$$T_{new} = \frac{W_p / V_p} + \frac{W_s / V_s}$$

$$= \frac{f \cdot W}{a V_s} + \frac{(1-f) \cdot W}{V_s}$$

$$= \left( \frac{f}{a} + (1-f) \right) \frac{W}{V_s}$$

$$S_{new-old} = \frac{T_{old}}{T_{new}} = \frac{(W/V_s)}{(W/V_s)(f/a + (1-f))}$$

$$= \frac{1}{(f/a + (1-f))}$$

$\xrightarrow{f=1} \quad \frac{1}{(1/a + 0)} = a = S_{all-improved} = S_{max}$

$\xrightarrow{f=0} \quad \frac{1}{(0/a + 1)} = 1 = S_{none-improved} = S_{min}$

$\xrightarrow[\text{f fixed}]{a \to \infty} \quad \frac{1}{(f/\infty + (1-f))} = \frac{1}{(1-f)} = S_\infty \geq S_{actual}$

$$f = 50\% \Rightarrow S_\infty = \frac{1}{(1-\frac{1}{2})} = 2 \quad \text{max improvement for } a \to \infty \text{ is twice as fast}$$

- If a=100, what is the overall speedup as a function of f?

Speedup vs Optimized Fraction



$S \uparrow$

$f \longrightarrow$

e.g., Use 100 CPUs

for $W_{parallel}$

$cost_{new} = 100 \cdot Cost_{old}$

$S' = 17$ ?

5% sequential

$f = 95\%$, $cost = 100$, $S' = 17$   Hmm?
only very large f is worth it?

## Amdahl's Law Example

- Suppose a program runs in 100 seconds on a machine, with multiply responsible for 80 seconds of this time. How much do we have to improve the speed of multiplication if we want the program to run 4 times faster?"   $T_{old} = T_{other} + T_{mult} = 20s + 80s = 100s$

$= S' = 4_{overall}$

$S_{new-old} = 4 = T_{old}/T_{new} = \dfrac{100s}{20s + 80s/S_p} \Rightarrow 80/S_p = \dfrac{100}{4} - 20 = 5$

$S_p = \dfrac{T_{p-old}}{T_{p-new}} = \dfrac{80}{5} = 16$

- How about making it 5 times faster?

$S = 5$ ?

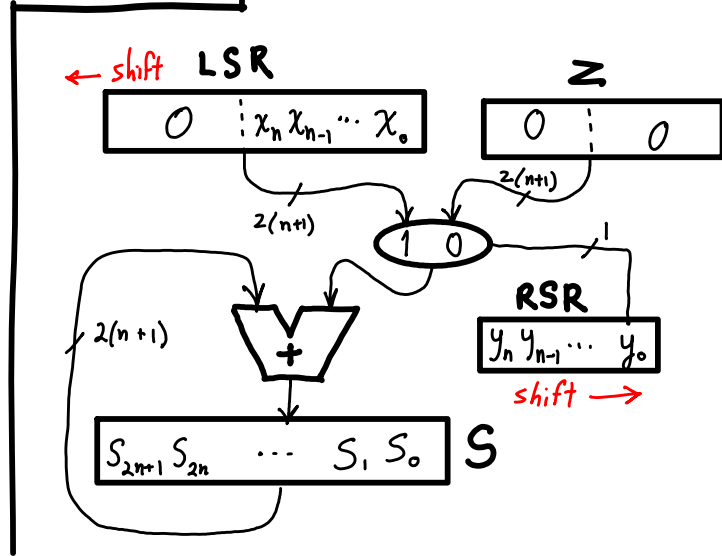$$x \cdot 7 = x(4 + 2 + 1) = 4x + 2x + x$$

↑ $y = 00\ldots0111$

2 shifts (C), 1 shift (B), 0 shifts (A)

**INT MULTIPLY**: $S = x * y$
**LSR**: partial products, initially x.
**S**: partial sum, initially 0.
**RSR**: initially y.
**Z**: all 0s

$$+ \quad \begin{array}{cccccc} 0 & 0 & \cdots & 0 & 0 \end{array} \qquad = S^0$$
$$\underline{\phantom{+} \quad x_n \; x_{n-1} \cdots x_1 \; x_0} \qquad + (A)$$
$$S'_n \; S'_{n-1} \cdots S'_1 \; S'_0 \qquad \overline{S'}$$

shift

$$+ \quad x_n \; x_{n-1} \cdots x_1 \; x_0 \; 0 \qquad + (B)$$
$$\underline{\phantom{+}} \qquad \overline{S''}$$
$$S''_{n+1} \; S''_n \cdots S''_2 \; S''_1 \; S''_0$$

shift

$$+ \quad x_n \; x_{n-1} \cdots x_1 \; x_0 \; 0 \; 0 \qquad + (c)$$
$$\underline{\phantom{+}} \qquad \overline{S'''}$$
$$S'''_{n+2} \; S'''_{n+1} \cdots S'''_3 \; S'''_2 \; S'''_1 \; S'''_0$$

What if y has a 0 bit? Then add 0 instead of shifted x: e.g., y = 0...101 add 0, not **B**.



e.g.

```
    1011
  x 0101
  ------
    1011

 + 101100
 --------
 00110111
```

(possible carry)

rewrite ⟹

```
         1011
       x 0101
       -------
```
+ (start, 0-shift) ← $Y_0 = 1$
+ 1011 (1-shift) ← $Y_1 = 0$
+ (2-shift) ← $Y_2 = 1$
+ 1011 (3-shift) ← $Y_3 = 0$
```
       -------
= 00110111
```

multiplier bit

⟹ We shift left (1011) every time, but add either the shifted (1011) or all zeroes, depending on whether $Y_i$ is 1 or 0.

$2^n$ minterms



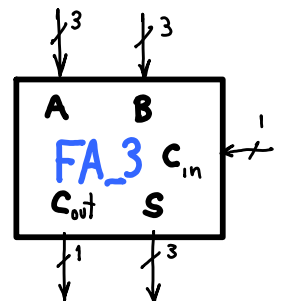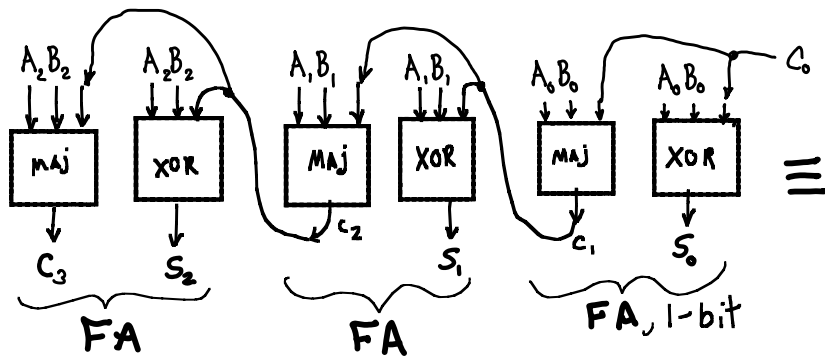## ADDER

Delay is longest path.

**2 levels per** MAJ, n bit operands.

**2n gate delays** until result is ready.



FA    FA    FA, 1-bit

**FA_3**  $A$  $B$  $C_{in}$  $C_{out}$  $S$

**3-bit Full Adder**

Shift Register Delay = 1  (all stages in Parallel)



overall delay: $(n \text{ shifts}) \times (2n \text{ delay}_{ADD}) = 2n^2$

| Improvement? | Recursive Refinement? |
| --- | --- |

(1-bit MULT)

$$\begin{array}{cccc} \times\frac{0}{0} & \times\frac{0}{1} & \times\frac{1}{0} & \times\frac{1}{1} \\ \hline 0 & 0 & 0 & 1 \end{array}$$

$\Rightarrow$  A B  AND gate  = $\boxed{\times}$  S

(2-bit MULT)

$$\begin{array}{rcccc} & & & (x_1 & x_0) \\ & & \times & (y_1 & y_0) \\ \hline & _c0 & _c0 & & x_0 y_0 \\ + & 0 & x_1 y_0 & 0 \\ + & 0 & y_1 x_0 & 0 \\ + & x_1 y_1 & 0 & 0 \\ \hline & S_3 & S_2 & S_1 & S_0 \end{array}$$

$\boxed{X-2}$ = 

(Use Recursively)

(n-bit)

$$\boxed{x_3 x_2}\ \boxed{x_1 x_0}$$
$$\times \quad \boxed{y_3 y_2}\ \boxed{y_1 y_0}$$

| | | | |
| --- | --- | --- | --- |
| 0 | 0 | $\boxed{X_L \cdot Y_L}$ | |
| 0 | $\boxed{X_H \cdot Y_L}$ | 0 | |
| 0 | $\boxed{X_L \cdot Y_H}$ | 0 | |
| $\boxed{X_H \cdot Y_H}$ | 0 | 0 | |

$+$

$$S_{HH} \quad S_{LH} \quad S_{HL} \quad S_{LL}$$



$\Rightarrow$ longest delay, $f(n)$, is,
   $4n$-bit ADD + $(n/2)$-bit MULT delay:
$$f(n) = 4(2n) + f(n/2)$$

$$f(n) = 8n + f(n/2) = 8n + 8(n/2) + f(n/4) \Rightarrow 16n\left(\frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \cdots\right) \approx 16n$$

**time improvement:** $2n^2 \Rightarrow 16n$     **Area increase:** $7n \Rightarrow ?$     $A(n) = 4A(n/2) + 4(n)$

<u>32-bit</u>  $2(2^5)^2 = 2^{11} \Rightarrow 2^4(2^5) = 2^9$  $\rightarrow$  $S'_{32\text{-bit}} = \frac{2^{11}}{2^9} = 2^2 = 4$     $S'_{64\text{-bit}} = 8$

$S'_{128\text{-bit}} = 16$

$$4 = S'_{overall} = \frac{v_{new}}{v_{old}} = \frac{W/T_{new}}{W/T_{old}} = \frac{T_{old}}{T_{new}} = \frac{W_s/V_{s\text{-old}} + W_p/V_{p\text{-old}}}{W_s/V_{s\text{-new}} + W_p/V_{p\text{-new}}}$$

<u>OUR Requirement</u>

$$= \frac{(0.2)W/V_{p\text{-old}} + (0.8)W/V_{p\text{-old}}}{(0.2)W/V_{p\text{-old}} + (0.8)W/S'_p \cdot V_{p\text{-old}}}$$

$$\left[\begin{array}{l} \underline{Assume} \\ V_{p\text{-old}} = V_{s\text{-old}} = V_{s\text{-new}} \\ V_{p\text{-new}} = S'_p \, V_{p\text{-old}} \end{array}\right]$$

$$= \frac{1}{(0.2 + 0.8/S'_p)} = 4$$

$$\Rightarrow 0.25 = 0.2 + 0.8/S'_p$$

$$\Rightarrow 0.05 = 0.8/S'_p \implies S'_p = \frac{0.8}{0.05} = \frac{80}{5} = 16 = S'_{128\text{-MULT}}$$

OK!

We can just make it if our data is 128-bit

---

$$S'_{overall} = 5? \qquad \frac{1}{(0.2 + 0.8/S'_p)} \overset{?}{=} 5 \quad \text{Can we?}$$

$$S'_{overall-\infty} = \frac{1}{(0.2 + 0.8/\infty)} = \frac{1}{0.2} = 5 \quad \text{Possible?}$$

if we eliminate MULT time?

---

# what else could we try?



shift_Right

shift_Left

$\longrightarrow$ n stages

n-stage pipeline
1 result every shift
  stage delay = $4n$   (2n-bit ADD)

(32-bit):
$$S'_{32\text{-MULT}} = \frac{2(2^5)^2}{4(2^5)} = 2^4 = 16$$

Can we keep it full?

$$\text{Latency} = 4n \times (n) = 4n^2$$

# Pipeling

a Task, consisting of things to do in order

| $S_5$ | $S_4$ | $S_3$ | $S_2$ | $S_1$ |
|---|---|---|---|---|

a sub-task

a processor: does one complete task at a Time.

| $S_5$ | $S_4$ | $S_3$ | $S_2$ | $S_1$ |
|---|---|---|---|---|
next Task

| $S_5$ | $S_4$ | $S_3$ | $S_2$ | $S_1$ |
|---|---|---|---|---|
currently processed Task

| $S_5$ | $S_4$ | $S_3$ | $S_2$ | $S_1$ |
|---|---|---|---|---|
finished task

time for task $= t$
Throughput $V = 1/t$

Break processor into pieces, one for each $S_i$

(new Tasks) → □ □ □ □ □ → (completed tasks)
$P_1$  $P_2$  $P_3$  $P_4$  $P_5$

| $S_5$ | $S_4$ | $S_3$ | $S_2$ | $S_1$ |
|---|---|---|---|---|
$T_5$

| $S_5$ | $S_4$ | $S_3$ | $S_2$ | $S_1$ |
|---|---|---|---|---|
$T_1$

| $S_5$ | $S_4$ | $S_3$ | $S_2$ | $S_1$ |
|---|---|---|---|---|
$T_2$

| $S_5$ | $S_4$ | $S_3$ | $S_2$ | $S_1$ |
|---|---|---|---|---|
$T_3$

| $S_5$ | $S_4$ | $S_3$ | $S_2$ | $S_1$ |
|---|---|---|---|---|
$T_4$

All tasks have a sub-task being processed in parallel.
Sub-task time $= t/5$.
One task completes every $t/5$.
$V = 1/(t/5) = 5/t$.
$S = 5$

## Adding Partial Products

0   X
↓32

$y_0$

0

(+) $b_0$ → $S_0$

0   X
$y_1$

32

(+) $b_0$ → $S_1$

0   X
$y_2$

32

(+)

32 → $b_0$ → $S_2$

⋮

← Partial sum, includes carry

$S = 2^{11}/2^3 = 2^8$

$\dfrac{10}{2 + 8/5} = \dfrac{10}{2 + 1/25}$

$\dfrac{\frac{10}{65}}{32} = \dfrac{320}{65} = 4.9$

## carry-save adder

$x_i$ 0

$C_{in}$

$C_{out}$

(+)

$S_i$

$X_2$ 0   $C_2$        $X_1$ 0  $C_1$       $X_0$ 0   $C_0$
$y_0$→                $y_0$→              $y_0$→

$C_2'$       $S_2$     $C_1'$      $S_1$     $C_0'$

→ $S_0$

$X_2$  0         $X_1$           $X_0$
$y_1$→      $y_1$→        $y_1$→

→ $S_1$

$X_2$ 0      $X_1$       $X_0$   $S_2$
$y_2$→     $y_2$→     $y_2$→

→ $S_2$

← 0

$S_5$      $S_4$      $S_3$

delay: ADD all partial Products = $n(2)$

$S'_{32\text{-}MULT} = 2^{11}/2^6 = 32$

$S' = \dfrac{1}{(0.2 + 0.8/32)}$

$= \dfrac{1}{(2/10) + (2^3/10)/2^5}$

$= \dfrac{10}{(2 + 1/4)} = \dfrac{40}{(8+1)} = \dfrac{40}{9} = 4 + 4/9$

## Other Ideas?

A $\xrightarrow{n}$
B $\xrightarrow{n}$

ROM

ADDR

$2^{2n}$
$2n$-bit
WORDS

↓ $2n$

A × B

How fast?

$2^{2n}$ inputs

MUX → out

$\div 2^{2n}$

delay? $2n$

area? exponential

A [ | | | | ]

⊗ ⊗ ⊗ ⊗

B [ | | | | ]

$n/4$  $n/4$  $n/4$  $n/4$

delay = $n/2$

(Combine)

Do mod arith in Parallel, combine results How?

# Evaluating Performance

- Performance best determined by running a real application *(in typical environment?)*
  - Use programs typical of expected workload
    - e.g., compilers/editors, scientific applications, graphics, etc.
- Microbenchmarks
  - Small programs, synthetic or kernels from larger applications
  - Nice for architects and designers
  - Can be misleading
- Benchmarks
  - Collection of real programs that companies have agreed on
  - Components: programs, inputs & outputs, measurements, rules, metrics
  - Can still be abused

⇒ Build compiler optimized for benchmark?
⇒ "Buggy" ⇒ skips work?

# The SPEC CPU Benchmark Suite
## (System Performance Evaluation Cooperative)

| SPEC2006 benchmark description | SPEC2006 | SPEC2000 | SPEC95 | SPEC92 | SPEC89 |
|---|---|---|---|---|---|
| GNU C compiler *(gcc →)* | | | | | gcc |
| Interpreted string processing | | perl | | | espresso |
| Combinatorial optimization | | mcf | | | li |
| Block-sorting compression *(GO →)* | | bzip2 | | compress | eqntott |
| Go game (AI) | go | vortex | go | sc | |
| Video compression | h264avc | gzip | ijpeg | | |
| Games/path finding | astar | eon | m88ksim | | |
| Search gene sequence | hmmer | twolf | | | |
| Quantum computer simulation *(QM sim →)* | libquantum | vortex | | | |
| Discrete event simulation library | omnetpp | vpr | | | |
| Chess game (AI) | sjeng | crafty | | | |
| XML parsing *(XML parse →)* | xalancbmk | parser | | | |
| CFD/blast waves | bwaves | | | | fpppp |
| Numerical relativity | cactusADM | | | | tomcatv |
| Finite element code | calculix | | | | doduc |
| Differential equation solver framework | dealII | | | | nasa7 |
| Quantum chemistry | gamess | | | | spice |
| EM solver (freq/time domain) | GemsFDTD | | | swim | matrix300 |
| Scalable molecular dynamics (~NAMD) | gromacs | | apsi | hydro2d | |
| Lattice Boltzman method (fluid/air flow) | lbm | | mgrid | su2cor | |
| Large eddie simulation/turbulent CFD | LESlie3d | wupwise | applu | wave5 | |
| Lattice quantum chromodynamics | milc | apply | turb3d | | |
| Molecular dynamics | namd | galgel | | | |
| Image ray tracing | povray | mesa | | | |
| Spare linear algebra | soplex | art | | | |
| Speech recognition *(speech recog. →)* | sphinx3 | equake | | | |
| Quantum chemistry/object oriented | tonto | facerec | | | |
| Weather research and forecasting | wrf | ammp | | | |
| Magneto hydrodynamics (astrophysics) | zeusmp | lucas | | | |
| | | fma3d | | | |
| | | sixtrack | | | |

. Kozyrakis          EE108b - Winter 2010 - Lecture 5

# Other Benchmarks

- Scientific computing: Linpack, SpecOMP, SpecHPC, …
- Embedded benchmarks: EEMBC, Dhrystone, …
- Enterprise computing
  - TCP-C, TPC-W, TPC-H
  - SpecJbb, SpecSFS, SpecMail, Streams,
- Other
  - 3Dmark, ScienceMark, Winstone, iBench, AquaMark, …

- Watch out: your results will be as good as your benchmarks
  - Make sure you know what the benchmark is designed to measure
  - Performance is not the only metric for computing systems
    - Cost, power consumption, reliability, real-time performance, …

# Summarizing Performance

- Combining results from multiple programs into 1 benchmark score
  - Sometimes misleading, always controversial…and inevitable
  - We all like quoting a single number

- 3 types of means
  - Arithmetic: for times

$$AM = \frac{1}{n}\sum_{i=1}^{n}\left(Weight_i\right)\cdot Time_i$$

  - Harmonic: for rates

$$HM = \frac{1}{\displaystyle\sum_{i=1}^{n}\frac{\left(Weight_i\right)}{Rate_i}}$$

  - Geometric: for ratios

$$GM=\left(\prod_{i=1}^{n}Ratio_i\right)^{\left(\frac{1}{n}\right)}$$

find ratio $\bar{r}$ s.t.

$r_1 \cdot r_2 \cdots r_n \, x = \left(\bar{r}\right)^n x$

$b_1$   $b_2$

$$R \Rightarrow \left( T_{R_1} , \; 10\, T_{R_1} \right)$$

**Normalize**: use reference **machine R** to get **speedups w.r.t. benchmarks (b-1, b-2)**.

**(R's time on b-2) = 10 X (R's time on b-1).**

**Combine speedups w.r.t R**:
--- Get **mean of speedups** w.r.t. R for **A**
--- Get **mean of speedups** w.r.t. R for **B**
--- Take **ratio of mean speedups.**

$$S_{A-R_1} = \frac{T_{R_1}}{100} \qquad S_{A-R_2} = \frac{10\, T_{R_1}}{4}$$

$$S_{B-R_1} = \frac{T_{R_1}}{200} \qquad S_{B-R_2} = \frac{10\, T_{R_1}}{1}$$

$$\overline{S}_{A-R} = \tfrac{1}{2}\left( \frac{T_{R_1}}{100} + \frac{10\, T_{R_1}}{4} \right)$$

$$\overline{S}_{B-R} = \tfrac{1}{2}\left( \frac{T_{R_1}}{200} + \frac{10\, T_{R_1}}{1} \right)$$

$$\Rightarrow \quad \frac{\overline{S}_{A-R}}{\overline{S}_{B-R}} = \frac{ T_{R_1}/2 \left( \frac{4 + 1000}{400} \right) }{ T_{R_1}/2 \left( \frac{1 + 2000}{200} \right) } = \left(\tfrac{1}{2}\right)\left(\frac{1004}{2001}\right) \cong \tfrac{1}{4} \; ?$$

$$= \left(\tfrac{1}{2}\right)\left(\frac{4 + 100\, r}{1 + 200\, r}\right) \quad \begin{array}{l} r \to \infty : \tfrac{1}{4} \\ r \to 0 : 2 \end{array}$$

**r** makes all the difference: **changing R or benchmarks ===> opposite conclusions?**

## Geometric Mean

$$\overline{S}_{A-R} = G(S_{A-R_1}, S_{A-R_2}) = \sqrt{\left(\frac{T_{R_1}}{100}\right)\left(\frac{10\, T_{R_1}}{4}\right)} = \sqrt{\frac{10\, T_{R_1}^2}{400}} = \frac{(\sqrt{10}\; T_{R_1})}{2 \cdot 10}$$

$$\overline{S}_{B-R} = G(S_{B-R_1}, S_{B-R_2}) = \sqrt{\left(\frac{T_{R_1}}{200}\right)\left(\frac{10\, T_{R_1}}{1}\right)} = \sqrt{\frac{10\, T_{R_1}^2}{200}} = \frac{(\sqrt{10}\; T_{R_1})}{\sqrt{2} \cdot 10}$$

$$S_{A-B} = \frac{\overline{S}_{A-R}}{\overline{S}_{B-R}} = \frac{(\sqrt{10}\; T_{R_1})}{(\sqrt{10}\; T_{R_1})} \left( \frac{T_{B_1} \cdot T_{B_2}}{T_{A_1} \cdot T_{A_2}} \right)^{1/2} = \left( \frac{T_{B_1} \cdot T_{B_2}}{T_{A_1} \cdot T_{A_2}} \right)^{1/2} \cong 0.7$$

**R** cancels. Conclusion **S$_{A-B}$ = 30 % slower?** Is this fair?
---- **on b1: S$_{A-B}$ = 200/100 = 2**
---- **on b2: S$_{A-B}$ = 1/4**

**job mix = (n1 runs of b-1) + (n2 runs of b-2)**

$$S_{A-B} = \frac{n_1 T_{B_1} + n_2 T_{B_2}}{n_1 T_{A_1} + n_2 T_{A_2}} = \frac{200 n_1 + n_2}{100 n_1 + 4 n_2} = \frac{200 + a}{100 + 4a} = \begin{cases} a \to \infty : \tfrac{1}{4} \\ a \to 0 : 2 \end{cases}$$

$$(a = n_2/n_1)$$

**Sanity check:** Given our result above, what **a** does GM **assume**?

$$\frac{200 + a}{(100 + 4a)} \approx \frac{3}{4} \implies (200 + a)4 = (100 + 4a)3 \implies 500 = 8a$$

$$\implies n_2 = 62 \, n_1 \quad \text{For every short job (b1), 62 long jobs (b2)?}$$

What if we hadn't taken the SQRT in GM?

$$\overline{S}_{A-B} = \frac{1}{2} \implies (200 + a)2 = (100 + 4a) \implies a = 50$$

---

$$\text{HM} = \frac{1}{\sum w_i/r_i} \qquad \text{where } r_1 = \frac{W_1}{T_1}, \quad r_2 = \frac{W_2}{T_2} \quad \cdots \qquad \left[ r_i = \frac{q_i}{} \right]$$

$$\left[ \begin{array}{l} \text{find } \overline{r} \text{ s.t., if we} \\ \text{did all work at same} \\ \text{rate, takes same time} \end{array} \right] \qquad \frac{(W_1 + W_2 + \cdots W_n)}{\overline{r}} = (T_1 + T_2 + \cdots T_n)$$

$$\implies \quad \overline{r} = \frac{W}{\left(\frac{W_1}{r_1}\right) + \left(\frac{W_2}{r_2}\right) \cdots \left(\frac{W_n}{r_n}\right)}$$

$$= \frac{1}{\sum \frac{(W_i/W)}{r_i}} \qquad (W_i/W) = \omega_i$$

We have reference Times:

$$W = \sum W_i$$

$$T_{R_1} = \frac{W_1}{q_{R-1}} \qquad T_{R_2} = \frac{W_2}{q_{R-2}} \qquad \left[ \text{assume } q_{R-1} = q_{R-2} = q_R \right]$$

$$W_1 = T_{R_1} q_R \qquad W_2 = T_{R-2} q_R \qquad \omega_i = \frac{W_i}{\sum W_i}$$

Use that to get weights:

$$\omega_1 = \frac{W_1}{W_1 + W_2} = \frac{T_{R_1} q_R}{(T_{R_1} q_R + T_{R_2} q_R)} = \frac{T_{R_1}}{T_{R_1} + T_{R_2}}$$

$$= \frac{1}{(1+10)} = \frac{1}{11}$$

$$\omega_2 = \frac{T_{R_2}}{T_{R_1} + T_{R_2}} = \frac{10}{1+10} = \frac{10}{11}$$

$$\overline{q}_A = \frac{1}{\sum \omega_i / q_{A-i}}$$

Given our assumption that $q_{R-1} = q_{R-2} = q_R$

$$\frac{W_1}{W_2} = \frac{T_{R-1} q_R}{T_{R-2} q_R} = \frac{T_{R-1}}{T_{R-2}}$$

$$q_{A-1} = \frac{W_1}{T_{A-1}} = \frac{W_1}{100}$$

$$q_{A-2} = \frac{W_2}{T_{A-2}} = \frac{10 W_1}{4}$$

$$= \frac{T_{R-1}}{10 T_{R-1}}$$

$$= \frac{1}{10}$$

$$W_2 = 10 W_1$$

$$\overline{q}_A = \frac{1}{\left( \frac{(1/11)}{W_1/100} + \frac{(10/11)}{10 W_1 / 4} \right)} = \frac{(11)}{\frac{100}{W_1} + \frac{10 \cdot 4}{10 W_1}} = \frac{11 W_1}{(100 + 4)}$$

↑ $T_{A-1}$   ↑ $T_{A-2}$

$$\overline{q}_B = \frac{11 W_1}{(200 + 1)}$$

$$\longrightarrow \quad S_{A-B}^{HM} = \overline{q}_A / \overline{q}_B = \frac{(200+1)}{(100+4)} \cong 2$$

---

Suppose, again, $n_1$ runs of $b-1$ and $n_2$ runs of $b_2$, w/ $a = n_2/n_1$. Comparing our "real world" performance, what $a$ does HM imply?

$$S_{A-B}^{real} = \left( \frac{200 + a}{100 + 4a} \right) = S_{A-B}^{HM} = \frac{2001}{104} \Rightarrow 104(200+a) = 201(100 + 4a)$$

$$\Rightarrow 20800 + 104a = 20100 + 804a$$

$$\Rightarrow 700 = 700a$$

$$\boxed{a = 1}$$

# Principles of Computer Design

- ## Take Advantage of Parallelism
  - e.g. multiple processors, disks, memory banks, pipelining, multiple functional units

- ## Principle of Locality
  - Reuse of data and instructions

- ## Focus on the Common Case
  - Amdahl's Law

$$\text{Execution time}_{new} = \text{Execution time}_{old} \times \left( (1 - \text{Fraction}_{enhanced}) + \frac{\text{Fraction}_{enhanced}}{\text{Speedup}_{enhanced}} \right)$$

1.02 #bytes per frame, time per file (cache, DRAM, ... )
1.03 avg CPI, CR, performance
1.04-05 CPI by class, CR, instr. mix,
1.06 compilers, avg CPI, CR, speedup, CPI by class, peak performance versus
1.07 Voltage scaling laws, C, power, GM, %change,
1.08 dynamic power, C, V
1.09 static and dynamic power, voltage dependence
1.10 multi-cores, #instructions, CPIs, execution time, power
1.11 die yield and cost
1.12 SPEC ratio from times
1.13 Faster clock, change ISA ==> fewer instructions executed, CPI vs CR
1.14 Performance measured by MFLOPS or MIPS versus overall
1.15 Amdahl's Law (improving only a fraction)
1.16 Speedup w/ communication costs