

Merricks.

- latency =  $\frac{\text{work}}{\text{bandwidth}}$ . e.g. write 10MB at 10MB/s : 1s

- corollaries: common  $\xrightarrow{\text{optimize}}$  uncommon  $\xrightarrow{\text{find}}$  new common  $\rightarrow$  global optim.  
 & make common case  $\xrightarrow{\text{there}}$  until  $\text{there}$  is no common  $\rightarrow$  divide up d

- inline  $\rightarrow$  constant  $\rightarrow$  loop unroll  $\rightarrow$  partial evaluation (substitute computation)
- Impacts: program input, instruction set

c) Affect CT: processor design (cycle time<sub>min</sub> = delay max along critical path)  
Manufacturing variation (transistor chans) - software policy

- clock freq  $\uparrow$  :  $E \uparrow$ ,  $P \uparrow$ ,  $ET \downarrow$

- Spec 2017: industry standard (vector), open-source, all in C/C++/Perl

x66

-Function call: ret: return address; leave: restore call stack pointer variable

RISC: compiler optimizable. [Alpha, MIPS, SPARC]. fixed: m8k, k11, d8c

- x86 → RISC-like instructions → execution

### Efficient Execution

Instruction Execution: fetch → Decode → RF read → {<sup>raise</sup> mem <sub>branch</sub> → next PC} ↓  
CPU, energy?

Pipeline use all stages at a time  
cycles for an instruction remain constant

- limitation: IC [ISA, compiler], CPI=1 [Base case], CT [critical logic path in chosen steps]
- ↳ Control hazard: fetch stage - branch mispredict


- Problem: Perform a CP1-6 by setting, Energy  $\sim \alpha$  is small  
in conventional reactions

E.g. 1-bit local predictor  $T[A \% 2^N] = 1$  predict 1 else NT;  $T[A \% 2^N] = \text{taken? } 1 : 0$

- Mis-speculation: CPI ↑, pipeline flushing
- IC determined by non-speculative, compiler/ISA; CPI: cycles wasted

② ~~Proto handler~~: Forward connections by pass the req file

③ terrible: micro-op. decode instructions for handling, average 1.3 uops/x86  
Decode & decompose instructions into uops  $\rightarrow$  decode queue

① Deeper pipeline:  , then  $CPI \downarrow$ ,  $CPI \uparrow$  (ca: of misprediction)  
pipeline stalls, overhead, ↑ mispredicts, ↓ correct by forward

- Impact: ET  $\downarrow$  (CT  $\downarrow$  50%), Energy (ACJ extra area, switching), Power (F devices)

- Impact:  $ET \downarrow$  (CPI  $\downarrow$  50%), Energy (ac + higher reg file)  $\uparrow$  a bit

⑤ Out-of-Order Execution : ...

A) Data Dependency - instruction read, correct values / order  
(RAW)

- $\left\{ \begin{matrix} \text{RAW} \\ \text{WAR} \end{matrix} \right\}$  false dependency: no data flows but A must execute before B

- critical path: longest sequence of activities
- formulas:  $CPI = \frac{CP}{\text{# innovations}}$ , average  $ILP = \frac{1}{CPI}$ ,  $IC = ILP * CP = \frac{1}{CPI} * CP$

$$E_T = I_C \times \frac{C_P}{I_C} \times C_T = C_P \times C_T$$

• attractive interest: its contribution to  $C_P$

B) Register renaming: RAT ~ Architecture  $\rightarrow$  Physical registers, parallelism








c) Out of order issue [Tommaso's Algo]

- Inputs: physical req broadcast from ALU after the insns finishes

[illegible]

all prior branches re-order before  
 } reg value available  
 } v-free, compatible proc

\* Scheduler : use instructions to parts (pipelined).  
Modern Processors : use instructions to parts (pipelined).  
 Pipeline: 1. Fetch 2. Decode 3. Execute 4. Write Back 5. Branch

F → D →  →  →  →  → ALU  
 ↑   

- smaller uops retired each cycle  $\rightarrow$  branch misprediction. = [uops]
- Branch resolution  $\uparrow$  . ROB utilization  $\uparrow$

**Cache**

- # block offset bits =  $\log_2(\text{block-size})$  bytes
- # entries = cache size (KB) / 1024

index =  $\text{base} + (\# \text{ entries}) \times \text{block size}$

localities. temporal ~ near in time, access same very soon; spatial ~ near in

space, next access is close to last access.

- completing: first access to the data
- capacity: address requested AND fully-associative cache isn't big enough

conflict: indexing covered collections. addrs ADD full-ass cache of in exp on bit

store policies

- miss: write allocate (bring cache into cache), write an allocate (bring into or to)
- hit: write through (tell know local cache not was changed), write back: means as dirty

Prefetching

H: Stream buffer { stride: difference  
                        { confidence

import: { idle hardware delay other request

5: PREFETCH: insert your own preferences. change:  $\rightarrow$  low + energy