# Performance

## Metrics:

① **Latency:** $t_{stop} - t_{start}$ (seconds). A task on a cycle has latency.
- speed/performance = $\frac{1}{latency}$, speedup = $\frac{latency_{before}}{latency_{after}}$ { >1: faster / <1: slower.
- importance: waiting / human vision / due soon

② **Bandwidth:** work per time, bytes/second, instruction/seconds, etc.
- latency = $\frac{work}{bandwidth}$. e.g. write 10MB at 10MB/s : 1s.

## Key Equations

① **Amdahl's Law:** optimizations × (generally) uniformly affect the entire program
- Speedup$_{total}$ = $\frac{1}{\frac{X_1}{S_1} + \frac{X_2}{S_2} + \cdots + 1 - X_1 - X_2 - \cdots - X_n}$    $X_i$: fraction of applications    $S_i$: speedup.
- corollaries: common $\xrightarrow{optimize}$ uncommon $\xrightarrow{final}$ new common ⇒ { global optimize / divide up dif
  * make common case fast until there is no common

② **CPU Performance Equation**
- Execution time = IC × CPI × CT = $\frac{cycles}{instr}$ × instructions × $\frac{seconds}{cycles}$
- CPI: average cycles, CT = $\frac{1}{frequency}$ (G: $10^9$, M: $10^6$), IC: dynamic instructions

A) **Reduce Instruction Count:** Algorithm/Compiler    { common sub elimination / constant propagation / function in-line
  { 00: no optimization, variable on stacks, lots of loads/stores
  { 01: × loads/stores, fewer movs, nothing on stack
- inline → constant → loop unroll → partial evaluation (substitute computation)
- Impacts: program input, instruction set

B) **Improve CPI:** $\frac{total\ cycles}{total\ instructions}$ = average latency of an instruction
  { Program: float/memory, larger CPI
  { Inputs: execute different parts of the program
  { Compile: which instruction. e.g. 00 Vs 01

C) **Affect CT:** processor design (cycle time$_{min}$ = delay max along critical path)
  Manufacturing variation (transistor chars). software policy

③ **Power and Energy Consumption**
A) Power = $\frac{Energy}{Time}$ = $\frac{J}{s}$ = Watt ; $P = V^2 FaC + Pidle$    (voltage, activity factor, frequency, capacitance (area), idle power consumption)
- P increases with area (c) and switching factor (a)
- Pidle proportional to the area of the chip
- V and F are linearly related. $P = F^3 aC + Pidle$ : F↑, Power↑

B) Energy = Power × Time = $(V^2 FaC + Pidle) × (CPI × IC × CT)$
- Clock freq ↑ : E↑, P↑, ET↓

## Benchmark   well-defined computation used to quantify system characteristics
  { Micro: measure specific aspect
  { Macro: stored app using standard inputs · characteristics { Realistic, portable, easy use / own datasets, useful measure / capture compiler/architecture
  { PoC: run a user-specific program
- Spec2017: industry consortium (vejo), open-source, all in C/C++/Fortran

# x86 Assembly

## ISA
Hardware – ISA { instruction format / instruction operation → software / interaction with memory

| | 16 | 32 | 64 | |
|---|---|---|---|---|
| ax | eax | rax | accumulator |
| bx | ebx | rbx | base |
| cx | ecx | rcx | counter |
| | 8 | 16 | 16 | 32 | 64 bits |

## x86
- instruction src1 src2 dst ;
- %<reg> : register, $imm : immediate, label: label.
- %eax : reg value ; (%eax): Mem(R[eax]) ; n(%eax): Mem(R+n);
- compare : cmpl %eax %ebx : %eax – %ebx ; e.g. cmp %eax, %ebx → %ebx >= %eax ; jge label
  · ZF: 01 if s1 and s2 are equal ; SF: 1 if s2 < s1.
- Function call : ret, return address ; leave: restore curr stack pointer
- Types { CISC: complex & compiler, human readable, conforms to add features [x86, VAX] / RISC: compiler operable. [Alpha, MIPS, SPARC]. fixed: simpl, high-level / interface × implementation ✓
- x86 → RISC-like instructions → execution

# Efficient Execution

**Instruction Execute:** fetch → Decode → RF read → { Arith / mem / branch } → next PC { single instr → 1 cycle / CPI↑, energy↑

**Pipeline:** use all stages at a time    cycles for an instruction remain constant
- limitation: IC [ISA, compiler], CPI=1 [Best case], CT [critical logic path in slowest stage]

① **Control hazard:** fetch stage × know next    $\frac{CPI_n × IC_n + CPI_n × IC_n}{IC_n + IC_n}$
- Problem: Perform CPI × b by scaling, Energy ~ a is small
  ⇒ **Speculative execution:** predict on prior, detect misprediction
  E.g. 1-bit local predictor T[A%$2^n$] == 1 predict T else NT ; T[A/$2^n$]: taken? 1: 0.
- Misspeculation: CPI↑, pipeline flushing
  · IC determined by non-speculative, compiler/ISA; CPI: cycles wasted Pidle
  · Energy: a=1 since all are busy. CPI↑ wasted energy; Non-speculative worse

② **Data hazard:** Forward, corrections by pass the reg file

③ **x86 terrible:** micro-op, decouple instructions for flexibility. average 1.3 uops/x86. Decode decompose instruction into uops → decode queue

## More performance
① **Deeper pipeline:** Fetch, Decode, then CT↓, CPI↑ { Data × resolve by forward / more than double    ① × poorer    pipeline adds overhead, ↑ gives a energy
  · Impact: ET↓ (CT ↓50%), Energy (aCT × extra over, switching). Power↑ (F doubles)

② **Wider pipeline:** two wide (rstr. data depends)
  · Impact: ET↓ (CPI ↓50%). Energy (aCT × bigger reg file) ↑ a bit
  · wider: more hazards, idle writes, overhead: long-latency complications

③ **Out-of-Order Execution:**
A) **Data Dependency** – instruction reads correct values/order
  { RAW / WAW / WAR } false dependency: no data flows but A must execute before B · RAW, WAW, and/or WAR-dependent instr
  · critical path: longest sequency of RAW, WAW, WAR
  · formulas: $CPI = \frac{CP}{\#\ instructions}$, average $ILP = \frac{1}{CPI}$, $IC = ILP × CP = \frac{1}{CPI} × CP$
    $ET = IC × CPI × CT = IC × \frac{CP}{IC} × CT = CP × CT$
  · effective latency: its contribution to CP    only RAW constr.
B) **Register renaming:** RAT ~ Architecture → Physical registers.    parallelism
  · ILP: examine longer region of the instruction stream    (no: log(pr)) × ar
C) **Out of order issue** [Tomasulo's Algo]
  · Inputs { physical reg / broadcast from ALU after the instruction finishes → scheduler ⇒ ALU



  · Scheduler: issue instructions to ports (pipeline). { reg value available / free, compatible port / priorities: old instrs

## Modern Processor
  · F [x86] D [Decode] → uops → [RoB] uops → rename → scheduler → ALU
    { front / Speculation / Retirement control    RAT
  · smaller uops retired each cycle → branch misprediction = [...]
  · Branch prediction ↑, RoB utilization ↑

# Cache
formulas { # block offset bits = $\log_2$ (block-size) bytes
  { # entries = $\frac{cache\ size\ (KB)}{block\ size × n}$ bytes ,  tag = address – # index – # tag
  { index = $\log_2$ (# entries)   * block size: data transferred on miss
  tag, index, offset    block size ↑ ⇒ miss ↓, # blocks ↓    MB = 1024 × 1024 = $2^{20}$

**Localities:** temporal ~ near in time, access same very soon ; spatial ~ near in space, next access is close to last access.

## cache misses
  { compulsory: first access to the data    { cache isn't big enough
  { capacity: address requested AND fully-associative able of line time exp a miss
  { conflict: indexing caused collisions. add ADD full-ass cache of m exp a hit

**Set Associative**   index → set [block1 block2]

## store policies
  · miss: write allocate (bring cache into cache), write no-allocate (notify lower cache)    notify lower cache check when evicted
  · hit: write through (tell lower level cache when not changed), write back. mark as dirty

## Prefetching
  H: Stream buffer { stride: difference / confidence . impact: { idle bandwidth delay after request
  S: PREFETCH: insert your own preferences.    { energy → lower ↑ energy

## Locality
High locality: miss (write allocate), hit (write back)
Low locality: miss (write no-allocate), hit (write through)

## Levels
  L1: write-through, write allocate ; L2 write back and allocate ; L3/Memory ~ write through    limited on more hit

## Virtual Memory

virtual address | page num | page offset |

### Address translation
- map from virtual to physical
- cache virtual in memory physical address

index → valid | physical page num |
→ physical p num | offset |

### Translation look-aside
- cache for address translation
- Access . TLB miss, TLB exception [OS], TLB fill . page fault [OS replace page table]
  TLB fill (valid) , L1 M, L2 M, L3 M, L3 fill , L2 fill , L1 fill .

## Memory Level Parallelism

critical path : long latency instructions influence critical page → interdependency?

### modern superscalar :
- non-blocking cache
- Memory Status History Register : track outstanding misses

limit { finite load
        finite MSHR
        ROB . physical register

### prefetcher . predictable access
{ finite resources ⇒ only a few accesses ahead
  idle bandwidth .

## Multi-core

### Unicore  MV decreasing ; performance $O(f)$ , Power $O(n)$ ; lack of ILP.

### parallel architecture
{ multiprocess : multiple CPUs tightly coupled enough to cooperate on single problem
  multithread : single CPU core that can execute multiple threads simultaneously
  multicore : CPU cores coexist on a single processor chip.

{ Multi-socket : two or more share same system ( RAM , Disk , Ethernet )

### Policy
- Inclusive : everything in L1 must also be in lower level of caches
- LRU : different for different levels of caches ; L3 has few information
- LFU : # times sth is used ; LFUDA lowers counters periodically so old .

### Cache Coherency
- shared memory : synchronisation + mutex
- MESI



coherency miss: miss on sharing data

## Floating Point

- Format
  32-bit.    | sign | mantissa | exponent |
             32 bit   24 bit     8 bit
  64 bit       1 bit   53 bit     11 bit

$$(-1)^s \times (\tfrac{m}{2^{24}} + 1) \times 2^{(e-127)}$$

- FP in Hardware is much faster.

- Dedicated FP physical register file . large ALU , [ Silicon Area vs Instruction ]
  fully piped : CPI = 1 / # ports

- Methods : Unrolling ⇒ ILP ; tiling ⇒ MLP

## Vector

- Single Instruction Multiple Data
- Impact : static instructions ↑ dynamic ↓ . wops ↓ . Energy ↓ , Power ↓
- Compiler / language / Intrinsics

## SMT  fetch instructions from different thread context [PC, reg file]

### Mechanism
- scheduler / ALU process Uops ,
  two mandatory : fetch , ROB . each thread RAT , load/store queue track
- small overhead

## Security

- Meltdown : reveal any memory mapped to their address space
  - rely : aggressive speculation that causes exception
    aggressive memory permission check in the TLB .
    memory mapping optimisations . high resolution timer .

- Myspectre : reveal any memory in the kernel
  aggressive speculation past branches .
  careful preparation of cache
  one of order execution
  high resolution timer .