

Bryan Melanson

How to Not Fail
Computer Architecture

While never going to class

Contents

1	Introduction	1
1.1	RISC	1
1.2	Dennard Scaling	1
1.3	Multiple Processors	2
1.4	Classes of Parallelism	2
1.5	Application Parallelism	2
1.6	Flynn's Taxonomy	2
1.7	Warehouse Scale Computers	3
1.8	Instruction Set Architecture	3
1.9	RISC-V Instruction Formats	3
1.10	Latency and Bandwidth	4
1.11	Dynamic State Energy and Power	4
1.12	Die Yield	4
1.13	Wafer Yield	5
1.14	Dependability	5
1.15	Performance	5
1.15.1	Execution Time	5
1.15.2	Processor Performance Equation	6
2	Instruction Set Architecture	6

1 Introduction

1.1 RISC

RISC (Reduced Instruction Set Computer) is an architecture which allowed for the exploitation of instruction level parallelism and the use of caches. This increased performance couldn't be matched by other architectures, and Intel 80x86 was forced to translate its instructions into RISC-like form internally to adopt RISC innovations. In low-end applications, this has led to RISC architecture like ARM becoming dominant.

1.2 Dennard Scaling

Power density was a constant for a given area of silicon, even as the number of transistors increased, due to their small dimension. This held up until a certain level in 2004, when the current and voltage levels could no longer ensure the dependability of integrated circuits.

1.3 Multiple Processors

The end of Dennard Scaling meant that improve performance more processors would be implemented, since single core chips would be inefficient with the increasing number of transistors. This meant that designs would no longer rely on **Instruction Level Parallelism**.

1.4 Classes of Parallelism

- **Data Level Parallelism** - data items operated on at the same time
- **Task Level Parallelism** - tasks performed independently in parallel

1.5 Application Parallelism

- **Instruction Level Parallelism** - exploits data-level using pipelining and speculative execution
- **Vector Architectures** - exploits data-level parallelism using single instructions on a collection of data in parallel
- **Thread Level Parallelism** - data or task-level parallelism allowing for interaction between parallel threads
- **Request Level Parallelism** - parallelism between decoupled tasks

1.6 Flynn's Taxonomy

- **Single Instruction Stream, Single Data Stream** - a single processor operates on sequential computation but can exploit ILP
- **Single Instruction Stream, Multiple Data Stream** - the same instruction executed by multiple processors, uses data level parallelism
- **Multiple Instruction Streams, Single Data Stream** - no implementation to date
- **Multiple Instruction Streams, Multiple Data Stream** - processors have their own instructions and own data operated on, using task level parallelism. Used in warehouse level computing, exploiting request level parallelism where many independent tasks can proceed in parallel

1.7 Warehouse Scale Computers

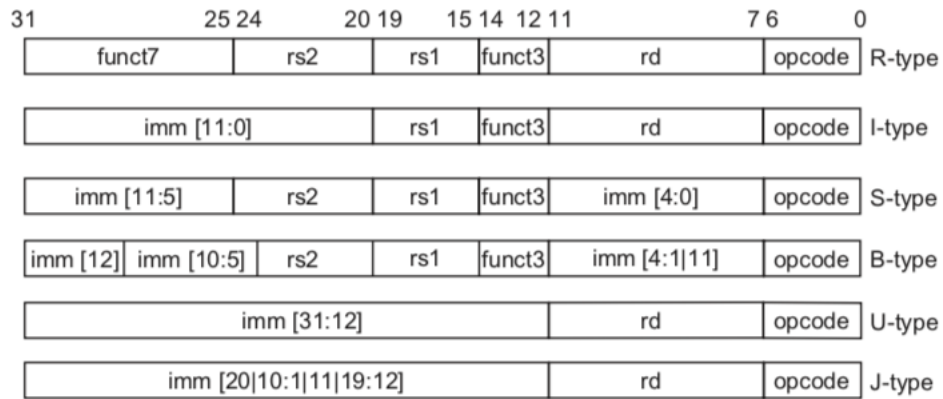
Warehouse Scale Computers **WSCs** are collections of servers or **clusters** which Software as a Service (SaaS) uses as a single larger computer, capable of considerable processing power. Price-performance and power are critical.

1.8 Instruction Set Architecture

The old view of computer architecture focused on **Instruction Set Architecture** (ISA) design involves decisions regarding registers, memory addressing, addressing modes, instruction operands, available operations, control instructions and instruction encoding.

Computer architecture now involves specific requirements of the target machine, designs meant to maximize performance within constraints, ISA, microarchitecture, and requires consideration of how compilers work.

1.9 RISC-V Instruction Formats



The above represents the format of different binary instruction types. The values *rs2*, *rs1*, *rd*, etc. represent the binary values which will specify which registers are targetted by the instruction. Values *funct* and *opcode* will describe which operation to execute, and *imm* fields will hold immediate values as arguments.

1.10 Latency and Bandwidth

Latency and Bandwidth are valuable performance metrics. Latency is the responsiveness with which commands can be executed; the time between start and completion of an event. Bandwidth is the total work done in a given time. Bandwidth and capacity have been more valued in the past, and in general show a greater growth rate - the growth of bandwidth is generally the square of the improvement in latency.

1.11 Dynamic State Energy and Power

Dynamic Energy is the switching of transistors to conserve energy. This value is therefore proportional to the product of the capacitive load and the square of the voltage:

$$\text{Energy}_{dynamic} \propto \text{Capacitive Load} \cdot \text{Voltage}^2$$

This represents the energy of a pulse from $0 \rightarrow 1 \rightarrow 0$, or vice versa, so the value of each single transition is:

$$\text{Energy}_{dynamic} \propto \frac{1}{2} \text{Capacitive Load} \cdot \text{Voltage}^2$$

And the Dynamic Power is:

$$\text{Power}_{dynamic} \propto \frac{1}{2} \text{Capacitive Load} \cdot \text{Voltage}^2 \cdot \text{Frequency Switched}$$

1.12 Die Yield

$$\text{Die Yield} = \left(1 + \left(\frac{\text{Defects/Area} \cdot \text{Area of Die}}{N}\right)\right)^{-N} \quad (1)$$

Die Yield is a value which defines how likely it is that a die will be non-defective during the fabrication process. This value is a probability < 1 , and the probability that a die will be defective is $1 - \text{Die Yield}$.

When converting between units of area (mm^2 , cm^2), be sure to repeat conversion between units twice.

$$1 \text{ mm}^2 = \frac{1}{10\text{cm}} \cdot \frac{1}{10\text{cm}} = \frac{1}{100} \text{ cm}^2$$

1.13 Wafer Yield

Wafer Yield is a value which describes how many dies (chips) can be created from a single wafer. This calculation is done using the size of a chip, and the diameter of each wafer. It takes into effect the square or rectangle shape of a chip and losses due to rounded sides of a wafer.

$$\text{Dies per wafer} = \frac{\pi \cdot (\text{Wafer Diameter}/2)^2}{\text{Die area}} - \frac{\pi \cdot \text{Wafer diameter}}{\sqrt{2} \cdot \text{Die area}} \quad (2)$$

These values can be used to calculate the cost of each die:

$$\text{Cost of die} = \frac{\text{Cost of wafer}}{\text{Dies per wafer} \cdot \text{Die yield}}$$

1.14 Dependability

- **Mean Time to Failure (MTTF)** - How long is continuous service delivered, the reciprocal is the rate of failure. One failure for every length of time.
- **Mean Time to Repair (MTTR)** - How long is service interruption
- **Availability** = $\text{MTTF} / (\text{MTTF} + \text{MTTR})$ - Service is continuous for x amount of time for every y length, $\text{MTTF} + \text{MTTR}$.

1.15 Performance

Performance metrics can differ between computers and their purposes - WSCs will want to increase bandwidth, and mobile devices will want to increase the speed with which commands are executed.

$$\frac{\text{Execution time}_Y}{\text{Execution time}_X} = n = \frac{\text{Performance}_X}{\text{Performance}_Y}$$

1.15.1 Execution Time

Amdahl's Law states that the performance improvement to be gained from using some fasted mode of execution is limited by the fraction of time the faster mode can be used.

$$\text{Speedup} = \frac{\text{Performance for entire task using enhancement when possible}}{\text{Performance for entire task without using the enhancement}}$$

The speedup depends on two factors:

1. The fraction of computation time in the original computer that can be converted to take advantage of the enhancement - $\text{Fraction}_{\text{enhanced}}$
2. How much faster the task would run if the entire program were enhanced - $\text{Speedup}_{\text{enhanced}}$

$$\text{Exec Time}_{\text{new}} = \text{Exec Time}_{\text{old}} \cdot ((1 - \text{Fraction}_{\text{enhanced}}) + \frac{\text{Fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}})$$

The above equation cancels in several points to be equal to

$$\text{Execution Time}_{\text{old}} \cdot \text{Old Code} + \text{Execution Time}_{\text{new}} \cdot \text{New Code}$$

$$\text{Speedup}_{\text{new}} = \frac{\text{Execution Time}_{\text{old}}}{\text{Execution Time}_{\text{new}}}$$

$$\text{Speedup}_{\text{new}} = \frac{1}{(1 - \text{Fraction}_{\text{enhanced}}) + \frac{\text{Fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}}}$$

1.15.2 Processor Performance Equation

$$\text{CPU Time} = \text{CPU Clock Cycles for a Program} \cdot \text{Clock Cycle Time}$$

Clock cycles can also be defined as **Instruction Count** · **Clock Cycle Time**

$$\text{CPU Time} = \text{Instruction Count} \cdot \text{Cycles per Instruction} \cdot \text{Clock Cycle Time}$$

2 Instruction Set Architecture

Memory is **byte addressable** in major architecture types - all bytes have a unique address, which allows for units as small as an ASCII character to be accessed without using larger accesses and masking.

ISA's should allow access commands for loading and writing byte, half words, and words so that units such as integers, and longs can be accessed and modified with ease.

Byte	8 bits
Half Word	16 bits
Word	32 bits
Double Word	64 bits

Little Endian orders bytes in a larger object such that the byte ending 000 will be found at the LSB position

7 6 5 4 3 2 1 0

Big Endian orders bytes so that the byte ending 000 will be located in the MSB position

0 1 2 3 4 5 6 7