# Introduction to SIMD

(Theory Part)

# Demo

# Overview

- What is SIMD?

- How does SIMD work on a hardware level?

- How does SIMD make programs run faster?

- How to write SIMD?

# Why should I learn about SIMD?

- Many areas of game programming are performance sensitive

- Tool in your toolbox

- Internship / job interview

# Single instruction, multiple data (SIMD)

A way to utilize the CPU hardware to process multiple values simultaneously for better performance.

# Areas of usage

Performance sensitive applications that process data in batches in a way that can be parallelized.

- Multimedia
- 3D Graphics
- Math / Physics

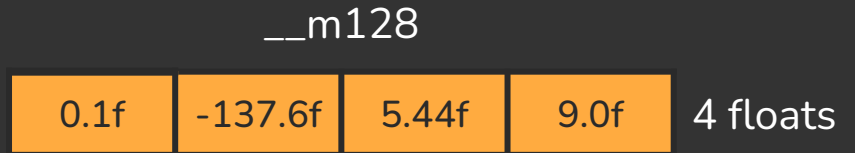## simdjson speed (C++)

twitter.json: **2.4 GB/s** on 3.4 GHz Skylake

|  | speed |
|---|---|
| simdjson (C++) | 2.4 GB/s |
| RapidJSON (C++) | 0.65 GB/s |
| Jackson (Java) | 0.35 GB/s |
| readLines C++ | 1.5 GB/s |
| disk | 2.2 GB/s |

# Vector/vectorization in the context of SIMD

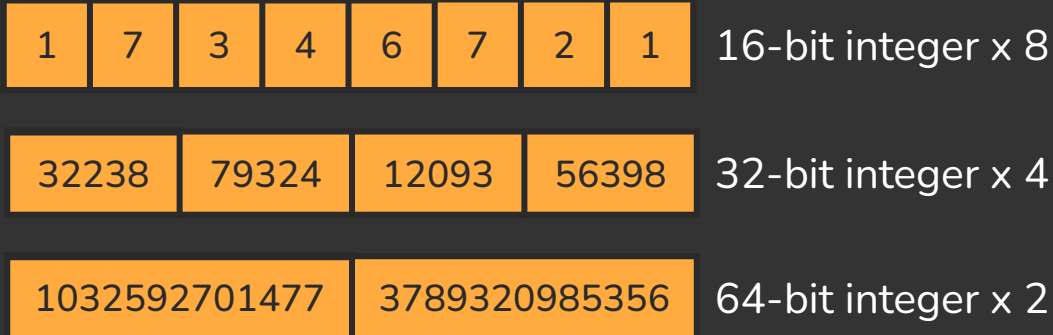Think of a SIMD vector as a fixed-size array.

No relation to:
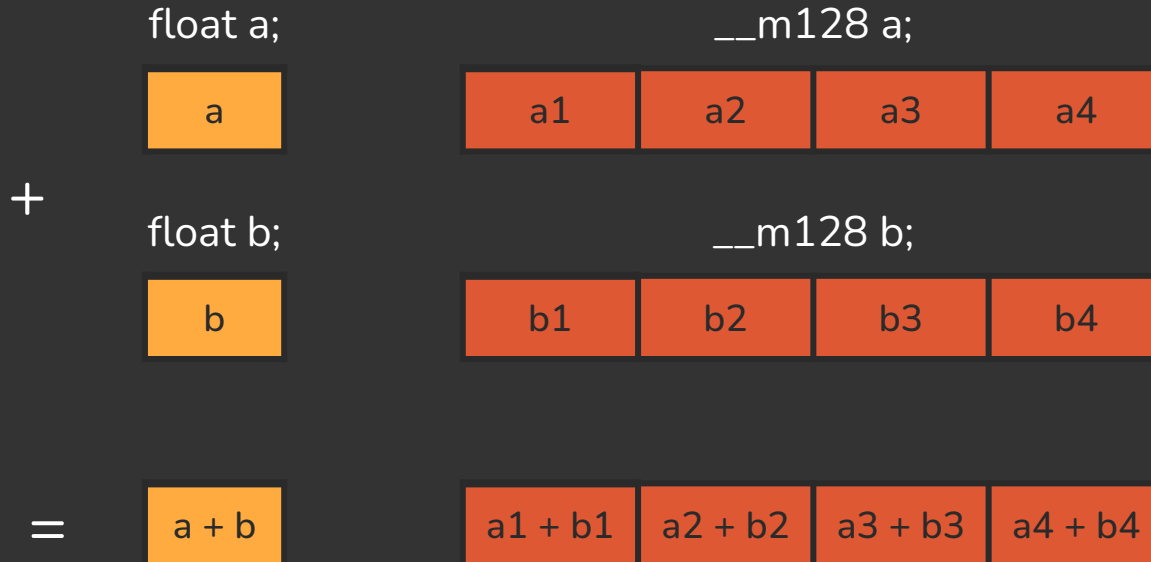
- std::vector

- Mathematical vector

- Vector2, Vector3

__m128

| 0.1f | -137.6f | 5.44f | 9.0f | 4 floats |

# Lanes

The number of elements in these vectors is not fixed. The proper term to use is lanes.

__m128i

| 1 | 7 | 3 | 4 | 6 | 7 | 2 | 1 |

16-bit integer x 8

| 32238 | 79324 | 12093 | 56398 |

32-bit integer x 4

| 1032592701477 | 3789320985356 |

64-bit integer x 2

# Scalar vs. Vector

float a;

| a |
|---|

__m128 a;

| a1 | a2 | a3 | a4 |
|----|----|----|----|

+

float b;

| b |
|---|

__m128 b;

| b1 | b2 | b3 | b4 |
|----|----|----|----|

=

| a + b |
|-------|

| a1 + b1 | a2 + b2 | a3 + b3 | a4 + b4 |
|---------|---------|---------|---------|

# CPU Architecture Overview

Instructions

CPU

# What is an instruction?

```
Source Code:
a = a * b;
```

↓ Compiler

```
Assembly Code:
mul eax, ebx
```

↓ Assembler

```
Machine Code:
11110111 11100011
```

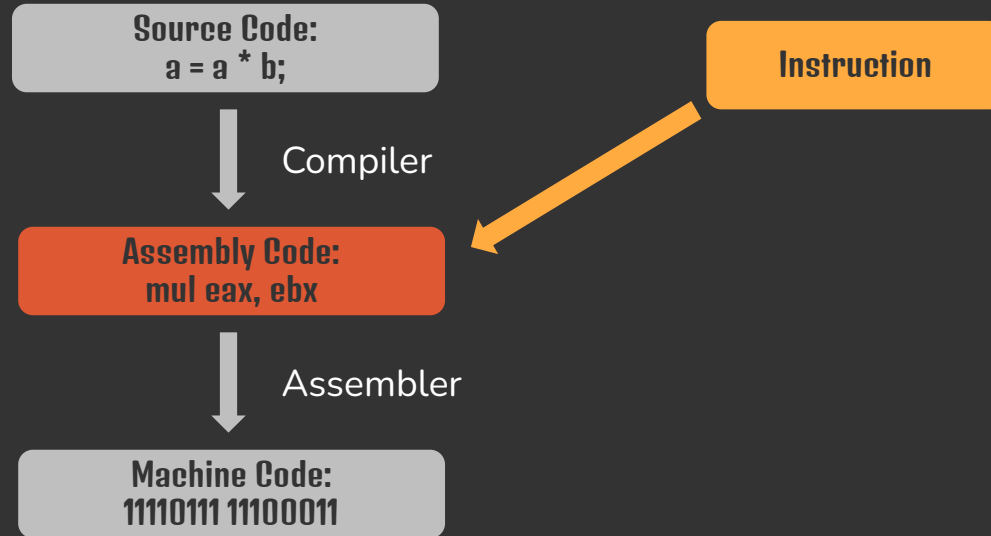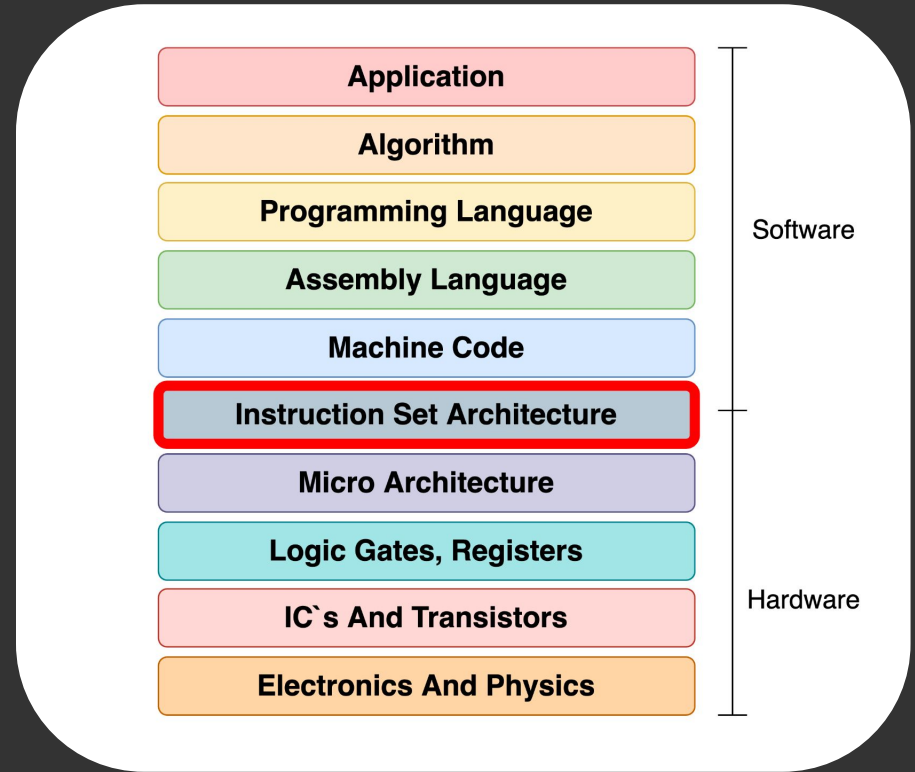**Instruction**

the NOTLANDERS

# Instruction Set Architecture (ISA)

- x86 - Desktop, PS5, Xbox
  - Intel
  - AMD
- ARM - Mobile, Switch
  - Apple Silicon
  - Qualcomm Snapdragon

| Application |
| --- |
| Algorithm |
| Programming Language |
| Assembly Language |
| Machine Code |
| **Instruction Set Architecture** |
| Micro Architecture |
| Logic Gates, Registers |
| IC`s And Transistors |
| Electronics And Physics |

Software

Hardware

the NOTLANDERS

# Instruction set extensions

SIMD instructions aren't generally part of the base instruction set. Instead, they come in the form of extensions.

- x86
  - SSE
  - AVX
- ARM
  - NEON

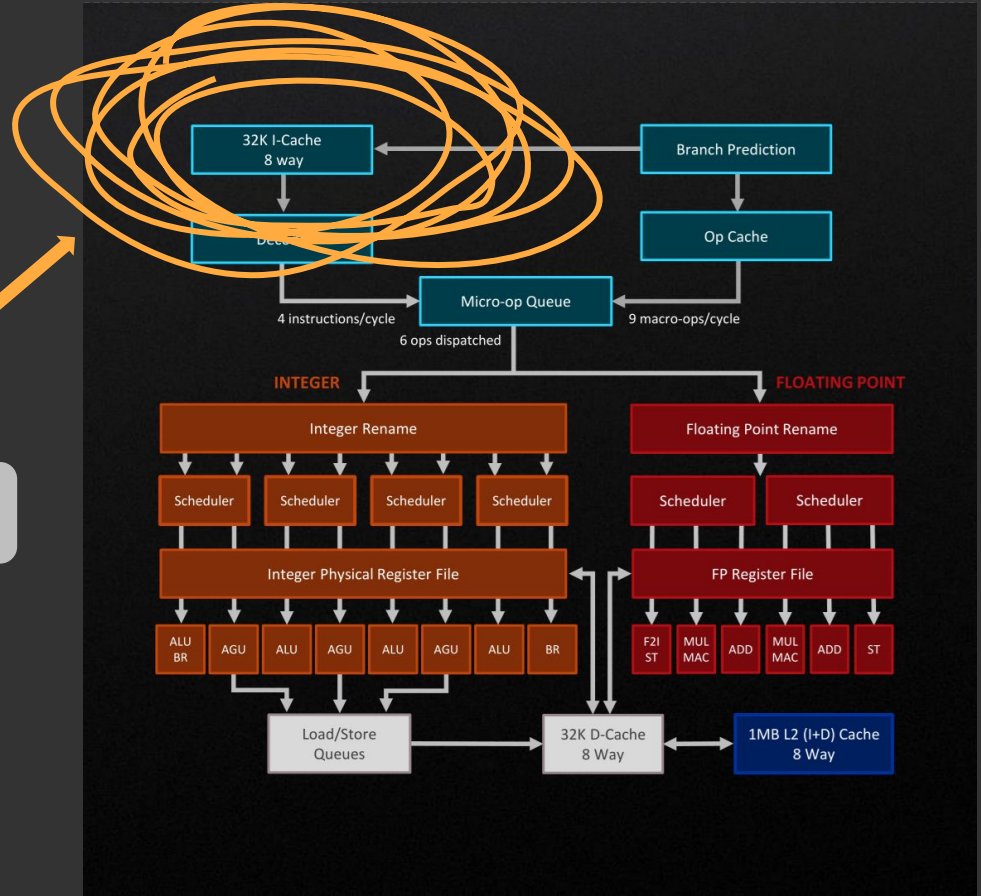# Anatomy of an x86 instruction

## add eax, ebx

| Mnemonic | Destination | Source |

# Instruction cycle

Fetch → Decode → Execute

# Fetch

Retrieve the next instruction(s) from memory.

*Instruction Cache*

```
Fetch  →  Decode  →  Execute
```

# Decode

Interpret the instruction, breaking it down into specific operations.

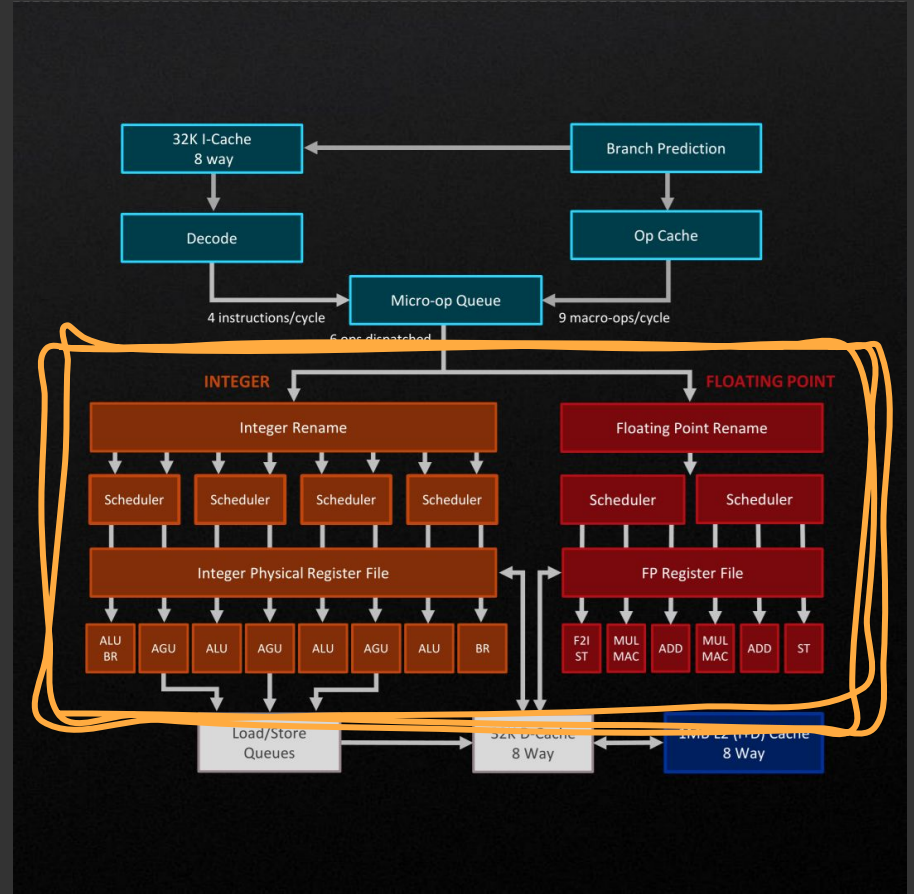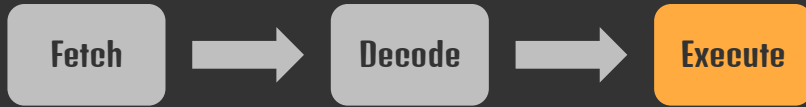*Decoders*

Fetch → Decode → Execute

# Micro-operations (micro-ops / μops)

https://uops.info/html-instr/IDIV_R32.html

# Execute

Carry out the operations specified by the instruction.

Execution Engine



Fetch → Decode → Execute

# Execute

- Registers - circuits that temporarily store inputs and outputs

- Execution units - circuits that perform operations on data in registers

Fetch → Decode → Execute

# Registers

A register is a super fast storage location inside the CPU. It's like a workspace for the CPU that temporarily hold data for processing.
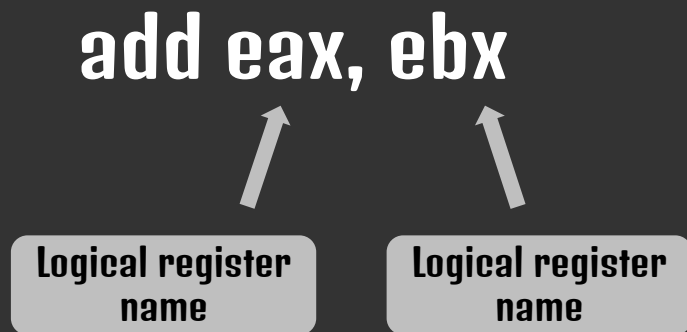
## add eax, ebx

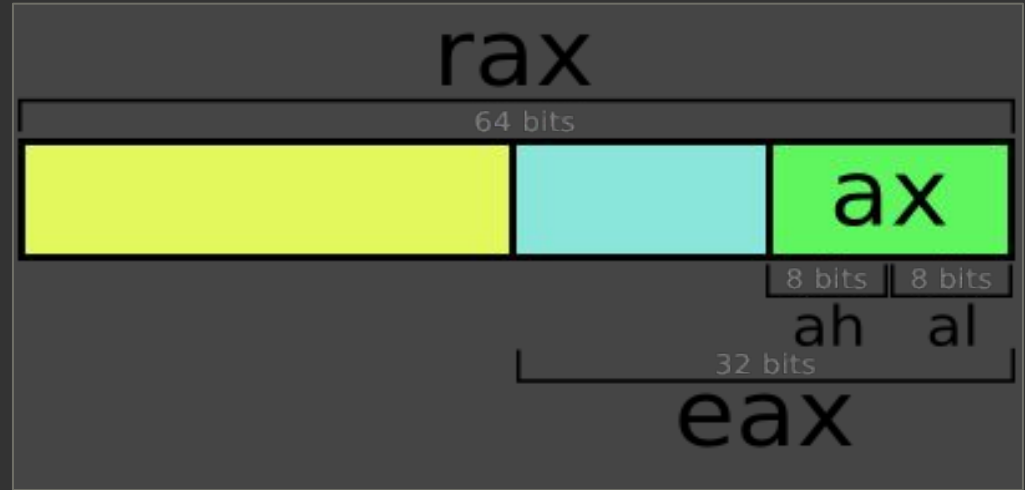| Register | Register |
| --- | --- |



the NOTLANDERS

# Register renaming

On Intel 8086 logical registers mapped directly to physical registers.

Modern CPUs have hundreds of physical registers and use register renaming.

add eax, ebx

Logical register name

Logical register name

# Register names

Different logical register names can be used to access different parts of a register.
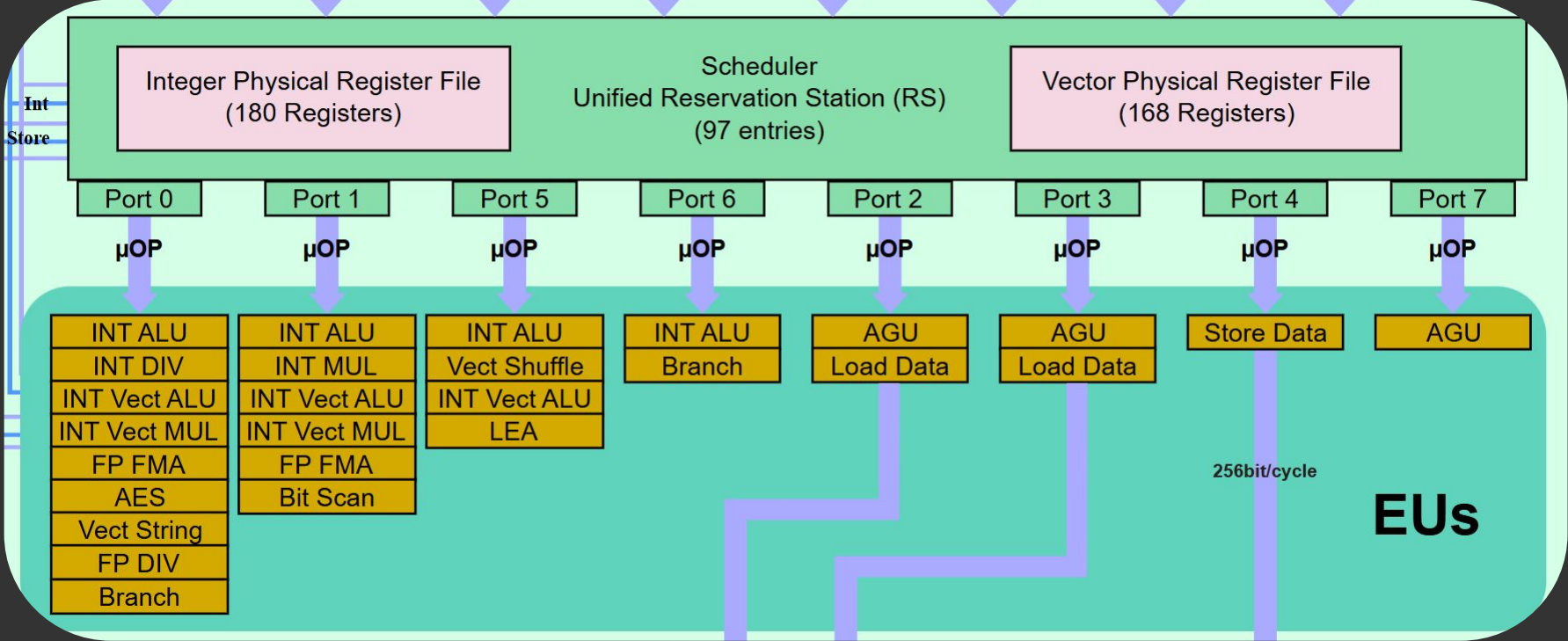
# Execution units

Specialized components for executing different types of operations.

- Arithmetic logic unit (ALU)
- Address generation unit (AGU)
- Load-store unit (LSU)
- Branch execution unit (BEU)

# Execution units



Scheduler
Unified Reservation Station (RS)
(97 entries)

Integer Physical Register File
(180 Registers)

Vector Physical Register File
(168 Registers)

Int
Store

| Port 0 | Port 1 | Port 5 | Port 6 | Port 2 | Port 3 | Port 4 | Port 7 |
|--------|--------|--------|--------|--------|--------|--------|--------|
| µOP | µOP | µOP | µOP | µOP | µOP | µOP | µOP |

**Port 0:** INT ALU, INT DIV, INT Vect ALU, INT Vect MUL, FP FMA, AES, Vect String, FP DIV, Branch

**Port 1:** INT ALU, INT MUL, INT Vect ALU, INT Vect MUL, FP FMA, Bit Scan

**Port 5:** INT ALU, Vect Shuffle, INT Vect ALU, LEA

**Port 6:** INT ALU, Branch

**Port 2:** AGU, Load Data

**Port 3:** AGU, Load Data

**Port 4:** Store Data

**Port 7:** AGU

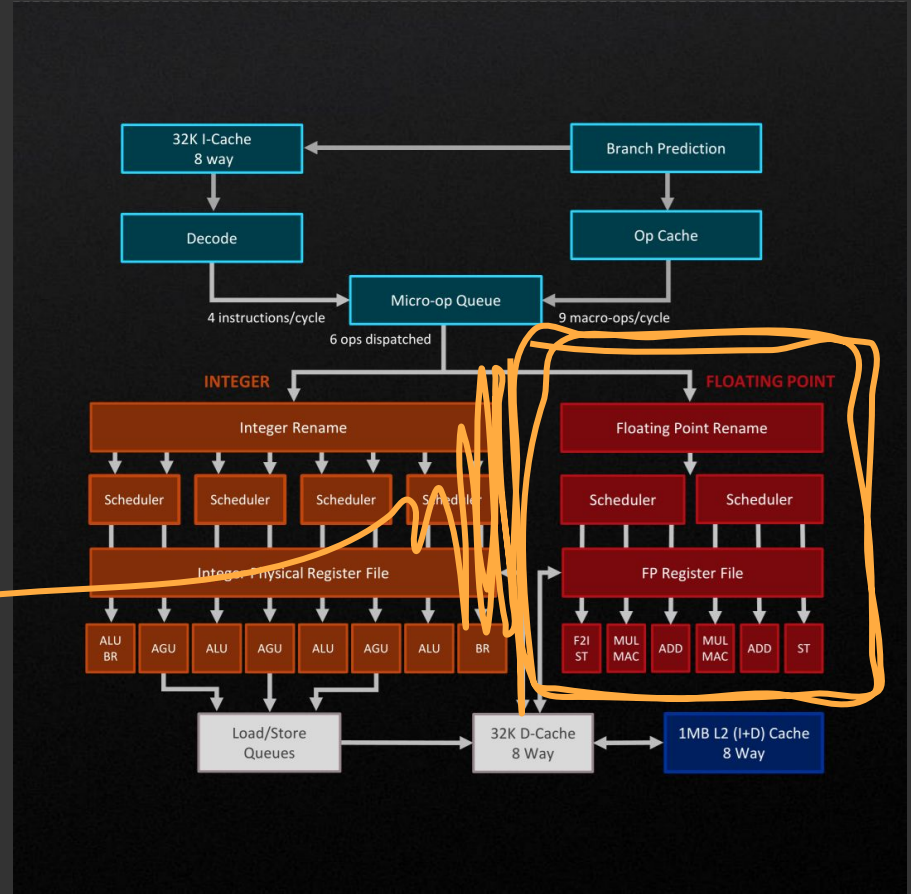256bit/cycle

EUs

the NOTLANDERS

# Where is the SIMD hardware?

# Floating point path = SIMD path

Old Intel CPUs used to have a floating point coprocessor for which they designed an extension called x87.

Modern compilers will generate SSE instructions for scalar floating point operations*, though x87 is still supported.
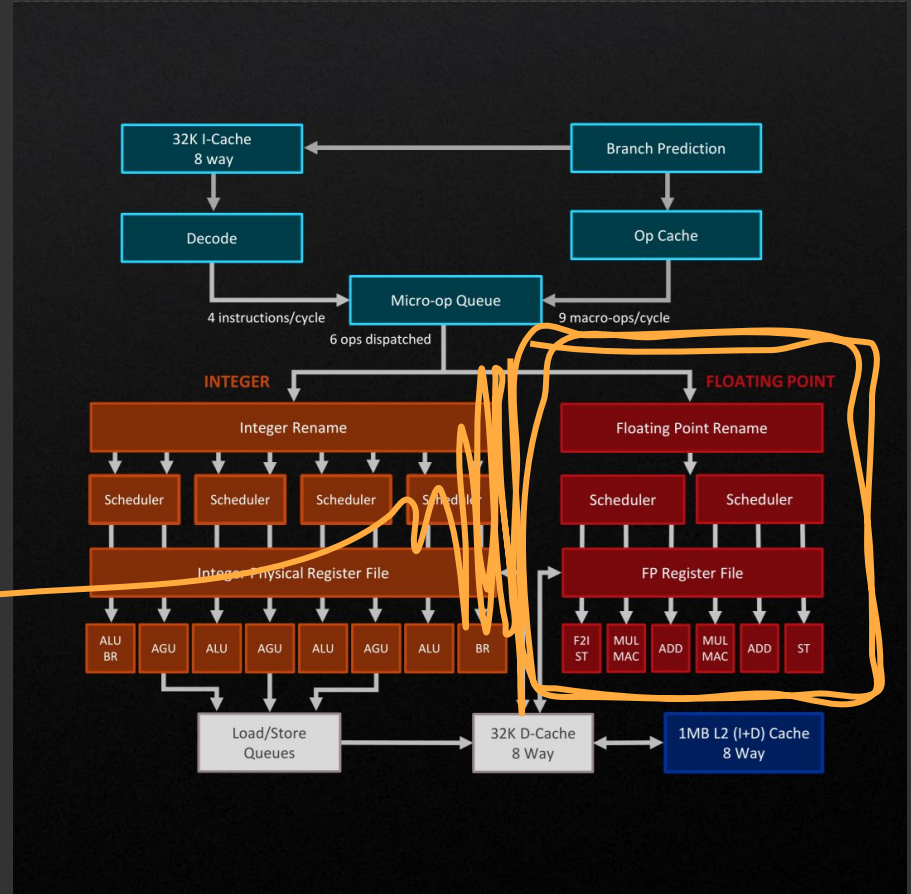
*SIMD hardware*



*I'm sure there are exceptions

# Wide hardware

While components that handle scalar integer instructions are typically 64-bit wide, hardware that handles SIMD instructions is 256-bit (or even 512-bit) wide on modern CPUs.
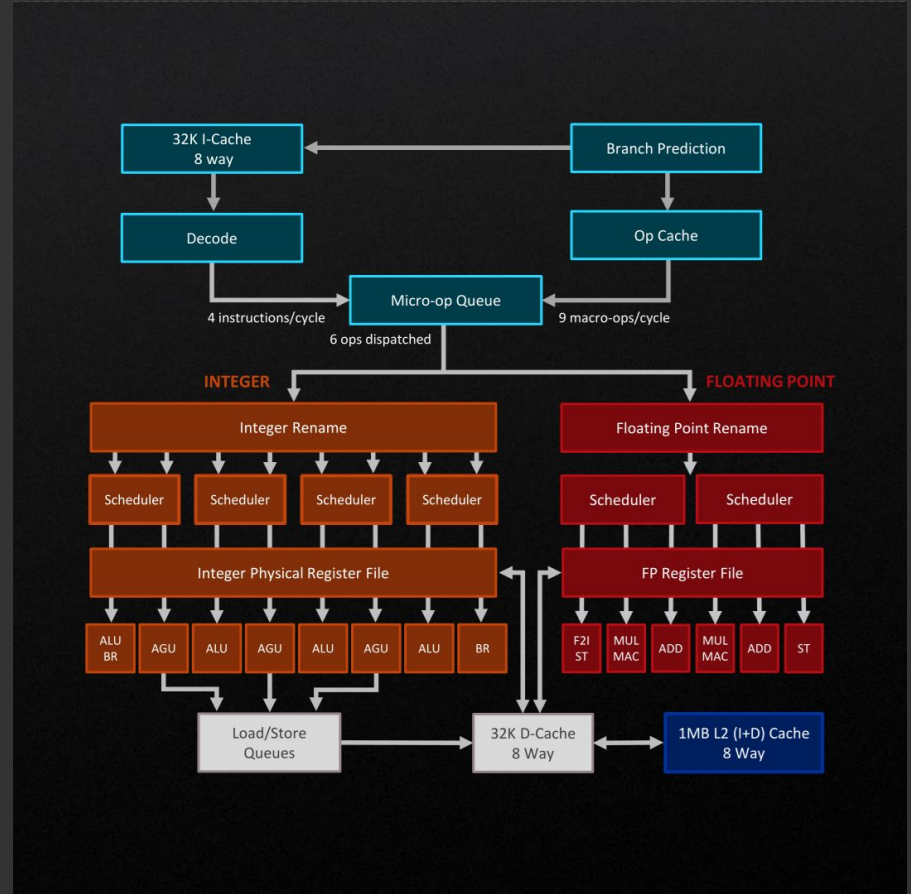
SIMD hardware

# CPU Architecture Overview
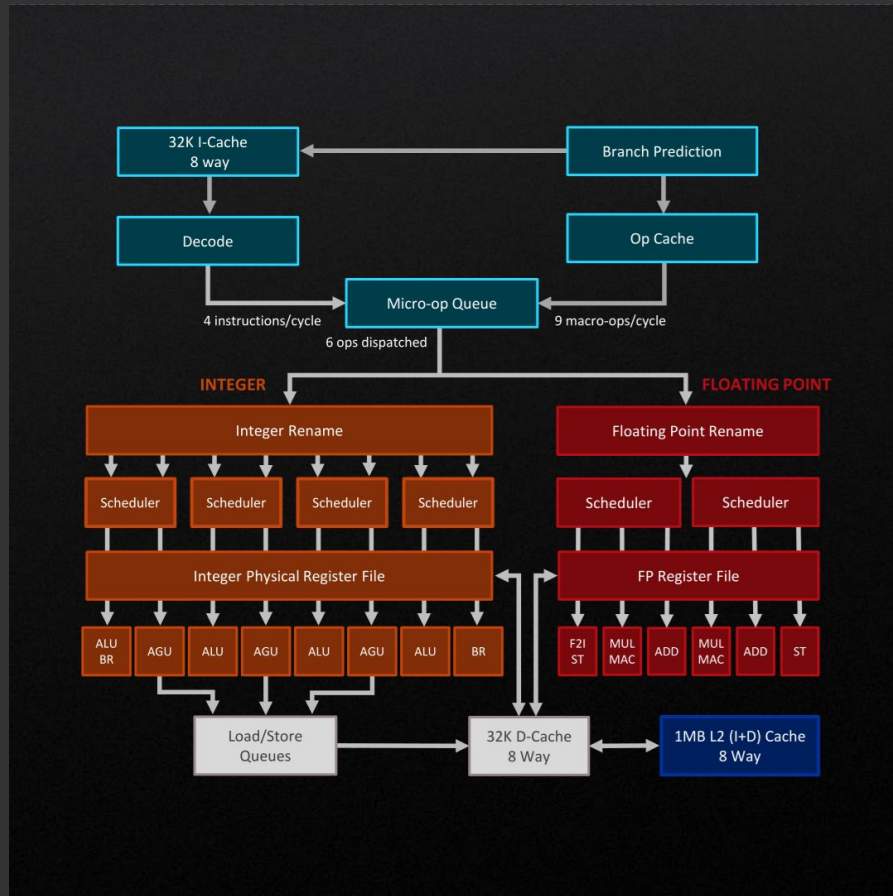
Instructions →

**CPU**

# Performance

1. Number of instructions
   a. Algorithmic complexity
   b. Waste

2. Speed at which instructions go through the CPU
   a. Data locality
   b. Multi-threading
   c. SIMD

Instructions

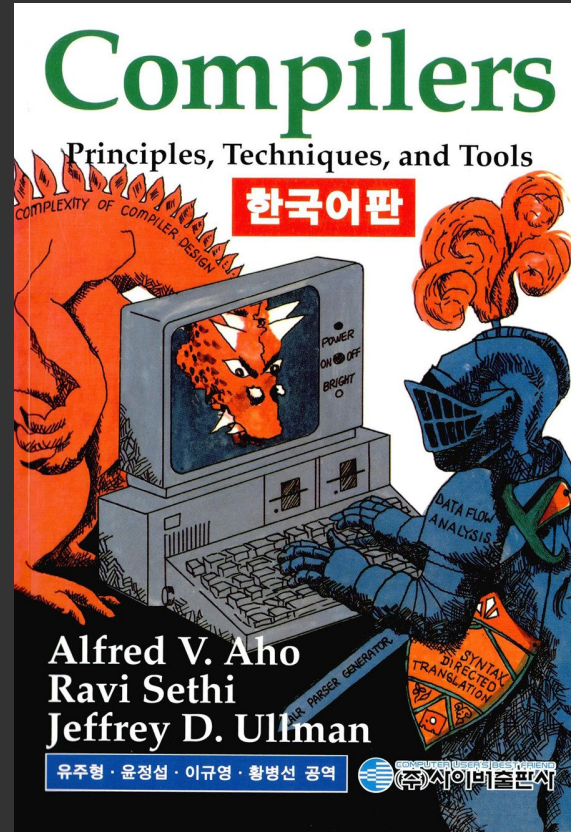**CPU**

# How does SIMD make my code faster?

- Faster execution - do more useful work at the same time by utilizing vector hardware.

- Reduced cost of decoding - only one instruction to decode instead of multiple.

- Maximizes cache bandwidth - load more data from memory per cycle.

# Auto-vectorization

If certain criteria are met, a compiler may be able to vectorize your code.

Free optimization? Yes, but compilers usually need some help to be able to auto-vectorize.

# Auto-vectorization

1.  Loop vectorizer unrolls loops and writes them as SIMD
2.  Block vectorizer (a.k.a. SLP) merges multiple scalars into a vector in a block of code

https://www.intel.com/content/dam/develop/external/us/en/documents/31848-compilerautovectorizationguide.pdf

Normal loop:

```
for (int i = 0; i < 1024; i++)
{
        a[i] = i;
}
```

Unrolled loop:

```
for (int i = 0; i < 256; i += 4)
{
        a[i] = i;
        a[i + 1] = i + 1;
        a[i + 2] = i + 2;
        a[i + 3] = i + 3;
}
```

NOTLANDERS

# Back to demo

# Learn More

Articles:

http://const.me/articles/simd/simd.pdf

https://mcyoung.xyz/2023/11/27/simd-base64/

Talks:

https://gdcvault.com/play/1022248/SIMD-at-Insomniac-Games-How

Courses:

https://www.computerenhance.com/

Books:

https://www.bokus.com/bok/9780128203316/computer-organization-and-design-risc-v-edition/

# Glossary

ARM

Auto-vectorization

AVX

CPU backend

CPU frontend

Execution engine

Execution unit

Fetch-Decode-Execute

Instruction set architecture

Microarchitecture

Register

Register renaming

SIMD

SSE

Vector

x86

# Thanks for listening!

Please fill out the feedback form!