

# SIMD Processing

---

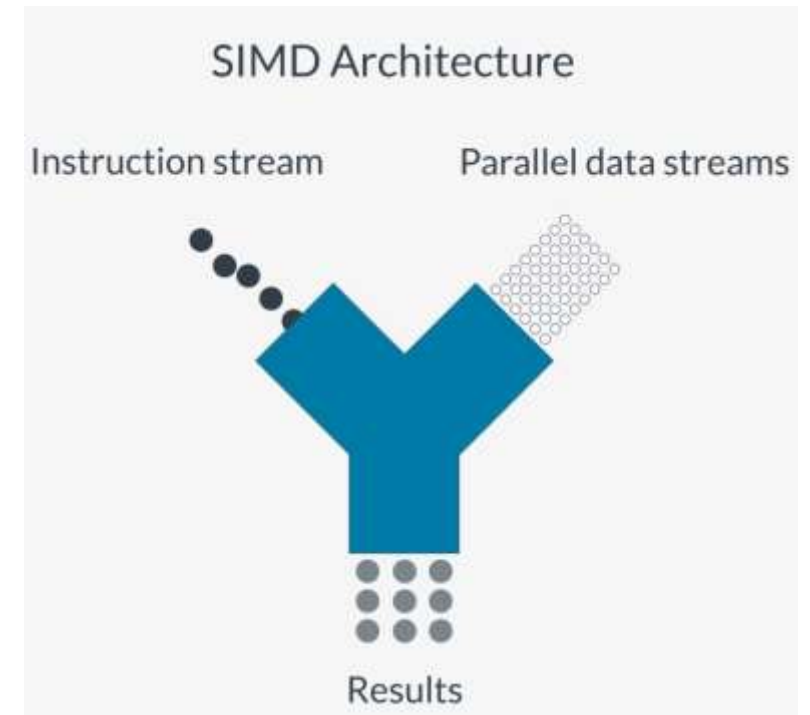
BY: TUSHAR BENDARKAR

(MASTERS IN DATA SCIENCE & ANALYTICS)

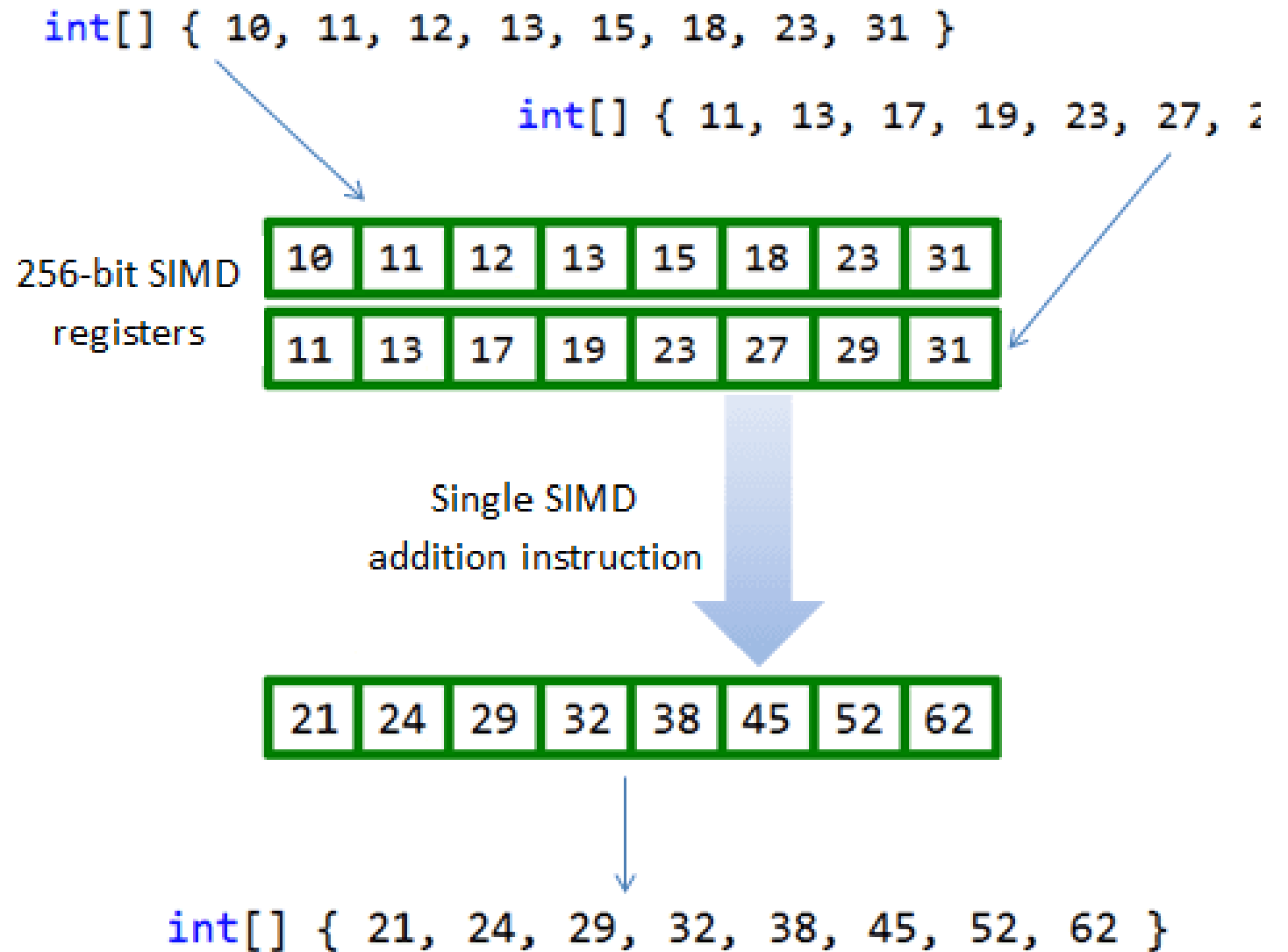
# What is SIMD ?

---

**Single Instruction, Multiple Data (SIMD)** is a parallel computing paradigm and class of architecture that enables the execution of the **same instruction on multiple data points simultaneously**. This approach is especially beneficial for tasks that involve repetitive mathematical operations on large datasets, such as **image processing and scientific simulations**. By allowing multiple cores to perform vectorized operations, SIMD helps modern CPUs operate more efficiently and maximize their potential.



For example, when adding two arrays element by element, a sequential execution would process each element one at a time. In contrast, SIMD processes multiple elements in parallel, significantly accelerating the computation.



# Key Features of SIMD:

- **Parallel Execution:** SIMD performs a single operation on multiple data points simultaneously.
- **Vectorization:** Data is processed in "vectors" (groups of elements), allowing the processor to handle many data points in fewer instructions.
- **Efficiency in Repetitive Tasks:** Ideal for operations that apply the same operation across a dataset, such as matrix multiplication, dot products, and filtering.

Scalar Operation

$$\begin{array}{l} A_1 \times B_1 = C_1 \\ A_2 \times B_2 = C_2 \\ A_3 \times B_3 = C_3 \\ A_4 \times B_4 = C_4 \end{array}$$

SIMD Operation

$$\begin{array}{l} A_1 \\ A_2 \\ A_3 \\ A_4 \end{array} \times \begin{array}{l} B_1 \\ B_2 \\ B_3 \\ B_4 \end{array} = \begin{array}{l} C_1 \\ C_2 \\ C_3 \\ C_4 \end{array}$$

# How SIMD Works ?

---

Modern CPUs come with specialized hardware for SIMD operations, often called **vector processing units (VPUs)**. These VPUs contain registers that store multiple data points, and the CPU can execute instructions on them simultaneously. For instance, a 256-bit SIMD register can hold eight 32-bit floating-point numbers, and operations can be performed on all of them at once.

This technique reduces the number of instructions the CPU needs to issue, which minimizes memory access overhead and improves throughput.

SIMD can be divided into two main types:

- **Short vector SIMD:** Used in common consumer processors like Intel's AVX (Advanced Vector Extensions) and SSE (Streaming SIMD Extensions).
- **Wide vector SIMD:** Found in specialized hardware like GPUs, used in deep learning and gaming.

# Hardware behind SIMD

---

The hardware behind SIMD (Single Instruction, Multiple Data) is fundamentally designed to accelerate parallel processing by allowing a single instruction to be executed on multiple pieces of data simultaneously. SIMD hardware is found in both CPUs and GPUs, and its architecture includes specialized components like **vector processing units (VPUs)** and **vector registers** to handle this parallel data.

## 1. Vector Processing Units (VPUs)

The **VPU is a dedicated unit within the CPU that handles vector operations**, which are central to SIMD. These units execute a single instruction on multiple data elements stored in **vector registers**. When a program uses SIMD operations, the VPU will apply the instruction (e.g., addition, multiplication) across all elements of a vector simultaneously.

### Example VPUs:

- **Intel AVX (Advanced Vector Extensions) and SSE (Streaming SIMD Extensions):** These are instruction sets in Intel and AMD CPUs that support SIMD, each with an associated VPU capable of handling vectorized instructions on 128-bit, 256-bit, or even 512-bit data.
- **ARM NEON:** ARM processors, common in mobile devices, use NEON instructions to perform SIMD operations, allowing faster processing of multimedia and gaming applications.

## 2. Vector Registers

SIMD hardware uses special registers, called **vector registers**, that **hold multiple data elements**. These registers are wider than typical scalar registers, so they can store several data points at once.

For example:

- **128-bit registers** (e.g., in SSE): Can hold four 32-bit floats or two 64-bit doubles.
- **256-bit registers** (e.g., in AVX): Can hold eight 32-bit floats or four 64-bit doubles.
- **512-bit registers** (e.g., in AVX-512): Can hold sixteen 32-bit floats or eight 64-bit doubles.

When a SIMD operation is executed, **data in these registers is processed in parallel, allowing a single instruction to act on all the elements at once**. This dramatically reduces the number of instructions and memory accesses required for large datasets.

## 3. Instruction Set Extensions (e.g., AVX, SSE, NEON)

SIMD-enabled CPUs support specific instruction set extensions, which define the SIMD operations that can be executed by the processor. These extensions specify how many data points can be processed in parallel and provide instructions for common tasks, like addition, subtraction, multiplication, and logical operations.

### Popular SIMD Instruction Sets:

- **Intel SSE (Streaming SIMD Extensions)**: The first widely adopted SIMD instructions, introduced in the 1990s. SSE uses 128-bit registers.
- **Intel AVX (Advanced Vector Extensions)**: An extension of SSE that introduced 256-bit registers, allowing more data to be processed per instruction. AVX-512 further expands this with 512-bit registers.
- **ARM NEON**: Used in ARM processors, NEON provides 128-bit SIMD instructions, enabling efficient parallel processing in mobile devices.

# Python, NumPy, and SIMD

---

Python, being an interpreted language, is not inherently optimized for SIMD. However, libraries like **NumPy** leverage lower-level languages like C, which have access to SIMD instructions through compiler optimizations and CPU-specific extensions.

## How NumPy Uses SIMD ?

NumPy is designed to work efficiently with arrays and matrices, and it implements many functions using highly optimized C code that benefits from SIMD. Operations like array addition, multiplication, and mathematical transformations (like sine, cosine, etc.) are automatically vectorized.



# SIMD Performance in NumPy

---

```
import numpy as np
import time

# Create two large arrays
N = 10**6
a = np.random.rand(N)
b = np.random.rand(N)

# Pure Python Loop (no SIMD)
start = time.time()
c = [a[i] + b[i] for i in range(N)]
end = time.time()
print(f"Pure Python loop time: {end - start:.6f} seconds")

# NumPy array addition (with SIMD)
start = time.time()
c = a + b
end = time.time()
print(f"NumPy (SIMD) time: {end - start:.6f} seconds")
```

Pure Python loop time: 0.720108 seconds  
NumPy (SIMD) time: 0.027219 seconds

As you can see, NumPy's SIMD-enabled addition is much faster than the pure Python loop. This speedup can be attributed to the use of vectorized operations in NumPy, which minimizes loop overhead and takes advantage of SIMD instructions.

## Trigonometric Function (e.g., Sine)

```
#Trigonometric Function (e.g., Sine)

# Define a large array
angles = np.random.rand(10_000_000) * 2 * np.pi

# Without SIMD (using Python Loops)
start = time.time()
sines = np.zeros_like(angles)
for i in range(len(angles)):
    sines[i] = np.sin(angles[i])
end = time.time()
print(f"Without SIMD: {end - start:.5f} seconds")

# With SIMD (using NumPy)
start = time.time()
sines = np.sin(angles)
end = time.time()
print(f"With SIMD: {end - start:.5f} seconds")

Without SIMD: 19.74245 seconds
With SIMD: 0.24002 seconds
```

## Matrix multiplication

```
# Generate two large random matrices
size = 500
matrix_a = np.random.rand(size, size)
matrix_b = np.random.rand(size, size)

# Without SIMD (manual matrix multiplication)
def manual_matrix_multiply(a, b):
    result = np.zeros((a.shape[0], b.shape[1]))
    for i in range(a.shape[0]):
        for j in range(b.shape[1]):
            for k in range(a.shape[1]):
                result[i, j] += a[i, k] * b[k, j]
    return result

start = time.time()
manual_result = manual_matrix_multiply(matrix_a, matrix_b)
end = time.time()
print(f"Without SIMD: {end - start:.2f} seconds")

# With SIMD (NumPy's dot product)
start = time.time()
numpy_result = np.dot(matrix_a, matrix_b)
end = time.time()
print(f"With SIMD: {end - start:.2f} seconds")

Without SIMD: 136.50 seconds
With SIMD: 0.02 seconds
```

# SIMD Extensions and NumPy Optimization

---

Modern CPUs from Intel, AMD, and ARM have SIMD instruction sets like AVX, AVX2, and NEON that are utilized by NumPy to improve performance. Compilers like **GCC** and **Clang** automatically vectorize code when using these instruction sets. NumPy's core is compiled with flags that enable SIMD optimizations, which means it can take advantage of AVX/AVX2 instructions on CPUs that support them.

## Key SIMD Extensions:

- **SSE (Streaming SIMD Extensions):** An older SIMD standard available on most Intel and AMD CPUs.
- **AVX (Advanced Vector Extensions):** A newer SIMD extension offering wider registers (256-bit) for even more parallelism.
- **NEON:** SIMD extension for ARM-based processors, used in mobile devices.

# Maximizing SIMD Gains in Python/NumPy

---

While NumPy automatically uses SIMD for many operations, here are a few tips to further maximize performance:

**Use NumPy arrays over Python lists:** NumPy arrays are stored contiguously in memory and optimized for SIMD, while Python lists are not.

**Use built-in NumPy functions:** Whenever possible, rely on NumPy's vectorized operations rather than custom loops. Functions like `np.add`, `np.multiply`, and `np.dot` are optimized for SIMD.

**Use contiguous arrays:** Ensure that your arrays are contiguous in memory (`args.flags['C_CONTIGUOUS']`). Non-contiguous arrays may not fully benefit from SIMD.

**Use larger arrays:** SIMD speedups are more noticeable with larger datasets where the overhead of parallel execution is minimized.

# Applications

---

## 1. Video and Image Processing

- SIMD is widely used for **accelerating video encoding/decoding, real-time image filters, and photo editing** applications.
- Operations like **resizing, color transformations, or convolution for edge detection** are vectorized for faster processing.

## 2. Machine Learning Inference

- SIMD accelerates matrix multiplications and vector operations, which are the backbone of deep learning inference.
- Libraries like **TensorFlow Lite, ONNX, and PyTorch** use SIMD to optimize model deployment on devices.

## 3. Scientific Simulations

- SIMD accelerates simulations in physics, climate modeling, and fluid dynamics by applying the same operations to vast arrays of data points (e.g., particles, temperature grids).

## 4. Gaming

- SIMD accelerates physics calculations, rendering, and AI in modern game engines.
- Used for **processing large arrays of vertices or pixels** in real-time for 3D graphics.

# Summary

---

SIMD processing is an essential technology for improving computational performance, especially when working with large datasets that involve repetitive operations. Python, through libraries like NumPy, leverages SIMD by executing vectorized operations on large arrays efficiently. This leads to significant performance improvements compared to traditional for-loops in pure Python.

To get the most out of SIMD:

- **Stick to NumPy's vectorized operations for heavy computations.**
- **Use modern hardware with AVX/AVX2 or NEON SIMD support.**

By following these principles, you can harness the power of SIMD and drastically reduce computation times in your Python code.

*Thank you* 😊

---