

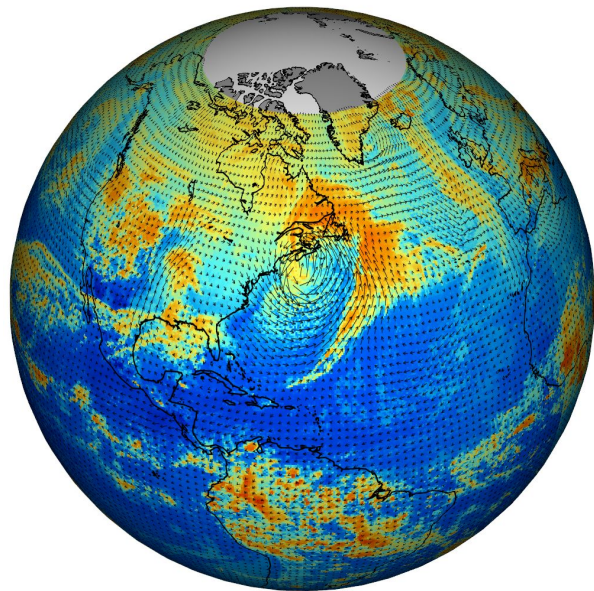
# Introduction to Parallel Programming

---

By Natalia Frumkin & Vijay Thakkar

**Para... what?**

# Climate Simulation Model



- Modeling the exchange of matter + energy over time
  - Wind and water currents
  - Atmospheric pressure
  - Weather tracking
- How is it done?
  - Discretize globe into fine grid
  - Equations characterize energy movement
  - *Solving equations repeatedly*

→ Highly parallel task

**Problem:** “Too Much Computation”

Running a simulation in series from start to finish would take forever!

**Notice:** These individual elements can be done concurrently! The operation is the same!

**Idea:** Parallelize operations for all elements

# Big Idea: Parallelization

## **1. Decompose this problem into smaller parts**

Discretize the Grid

## **2. Do those on multiple concurrent processing streams**

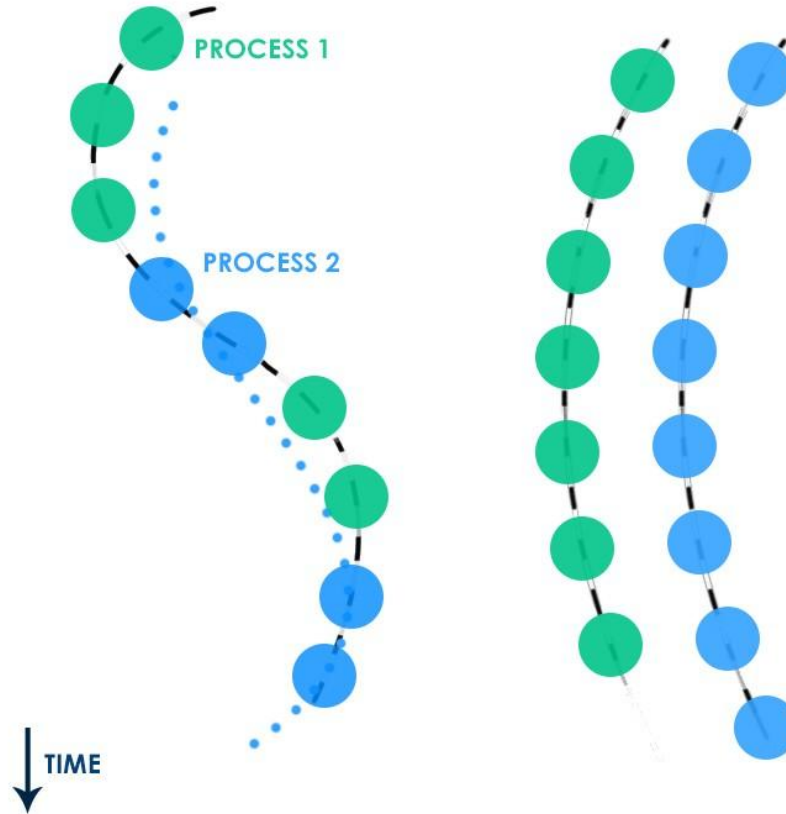
Separate system model into local and global phenomena

## **3. Run these concurrent streams in parallel**

Spread Compute across multiple processors

# Concurrent vs. Parallel

**Interrupts**



**Simultaneous**

## Consider Single Precision $ax + y$ :

```
void saxpy(int n, float a, float * restrict x, float * restrict y)
{
    for (int i = 0; i < n; ++i)
        y[i] = a*x[i] + y[i];
}

// Perform SAXPY on 1M elements
saxpy(1<<20, 2.0, x, y);
```

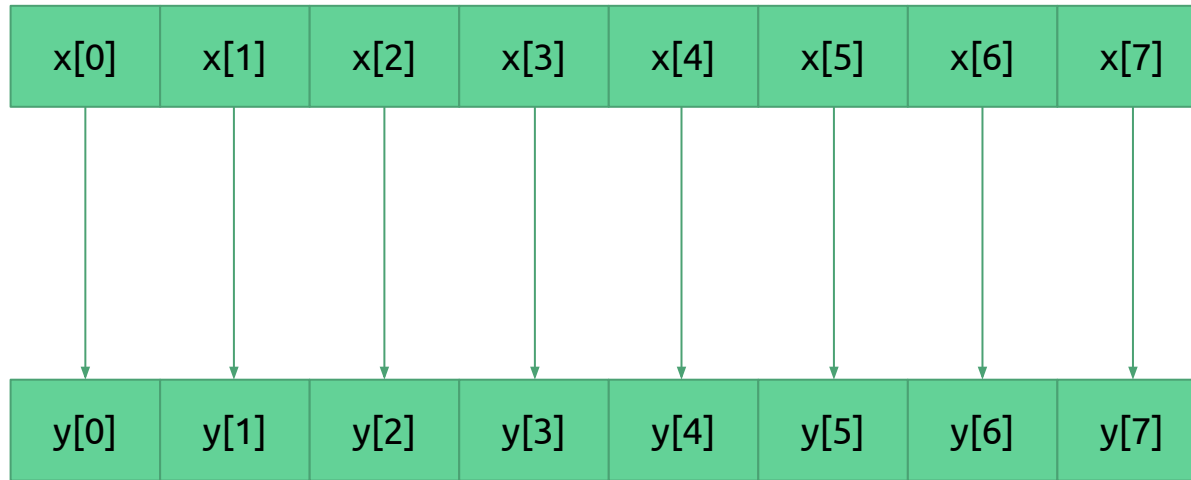
# Serially...





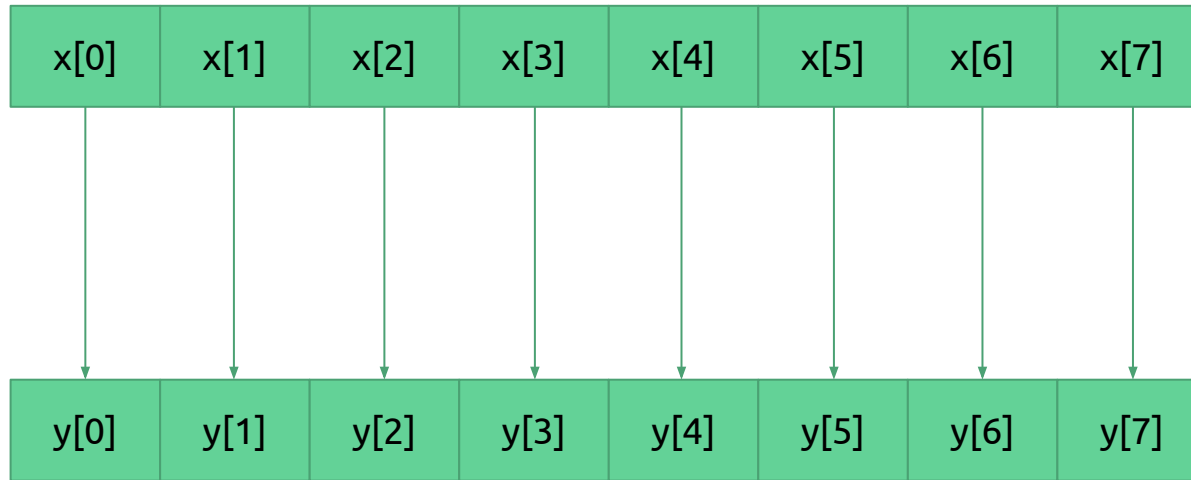
# Concurrently...

Op1a Op1b Op1c Op1d Op2a Op2b Op2c Op2d

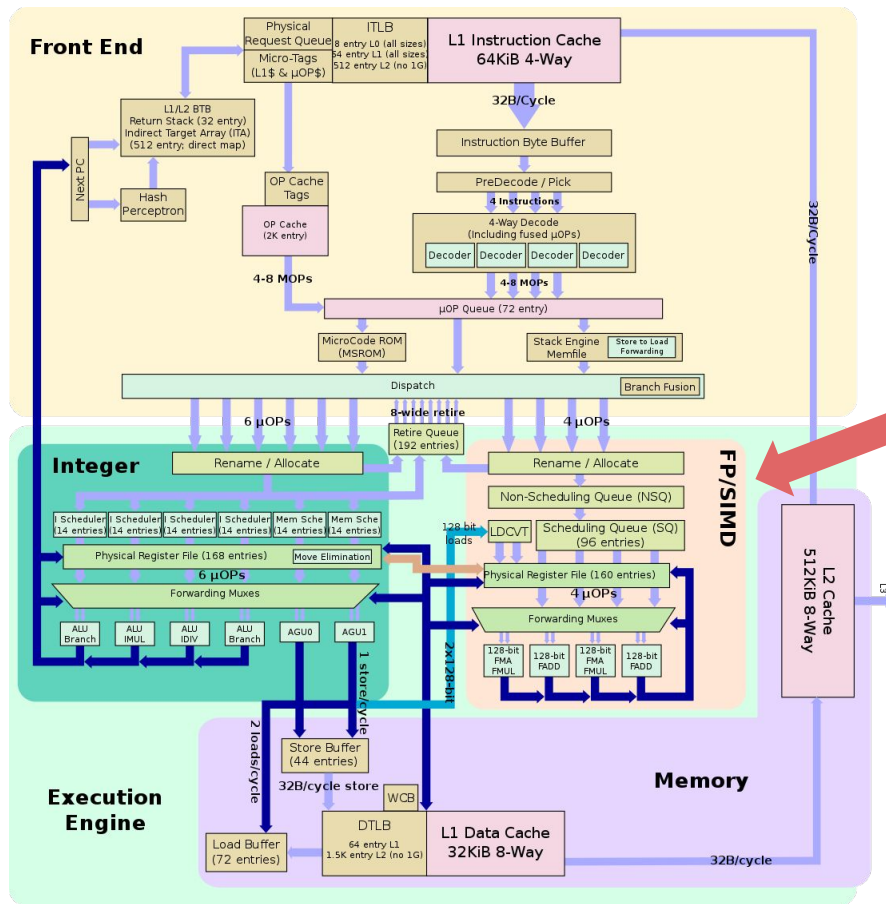


# In Parallel...

Each operation assigned to a different execution unit



# Hardware Approach: Vectorization



# Hardware Approach: Vectorization

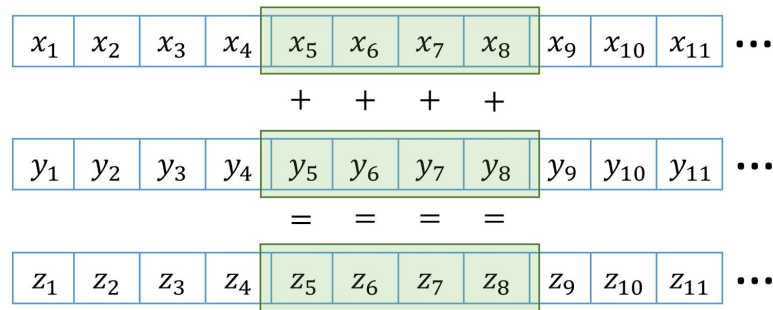
## Fine Grained Parallelism:

Do the same math operation on multiple elements in the same instruction

## SIMD:

Single Instruction - Multiple Data

aka. Vectorization



## Case Study: AVX Extensions

```
1 void saxpy_avx256(size_t n, float a, float *x, float *y, float *z) {  
2     __m256 a_vec = _mm256_set1_ps(a);  
3  
4     for (size_t i = 0; i < n; i += sizeof(__m256) / sizeof(float)) {  
5         __m256 x_vec = _mm256_load_ps(&x[i]);  
6         __m256 y_vec = _mm256_load_ps(&y[i]);  
7         __m256 res_vec = _mm256_add_ps(_mm256_mul_ps(a_vec, x_vec), y_vec);  
8         _mm256_store_ps(&z[i], res_vec);  
9     }  
10 }
```

## Case study: OpenMP `#pragma omp simd`

```
14
15 void saxpy_omp(size_t n, float a, float *__restrict__ x, float *__restrict__ y,
16 | | | | | | | float *__restrict__ z) {
17 #pragma omp simd aligned(x, y, z : 64)
18 | for (size_t i = 0; i < n; i++) {
19 | | z[i] = (a * x[i]) + y[i];
20 | }
21 }
22
```

## Case study: GPU parallel code

```
15  __global__ void saxpy_cuda(size_t n, float a, float *x, float *y, float *z) {
16      int idx = (blockIdx.x * blockDim.x) + threadIdx.x;
17      if (idx < n)
18      |   z[idx] = (a * x[idx]) + y[idx];
19  }
```

```
48      cudaMalloc(&d_x, n * sizeof(float));
49      cudaMalloc(&d_y, n * sizeof(float));
50      cudaMalloc(&d_z, n * sizeof(float));
51      cudaMemcpy(d_x, x, n * sizeof(float), cudaMemcpyHostToDevice);
52      cudaMemcpy(d_y, y, n * sizeof(float), cudaMemcpyHostToDevice);
53      cudaMemcpy(d_z, z, n * sizeof(float), cudaMemcpyHostToDevice);
54      cudaDeviceSynchronize();
55
56      auto start = high_resolution_clock::now();
57      for (auto i = 0; i < itrs; i++) {
58      |   saxpy_cuda<<<(n + 255) / 256, 256>>>(n, a, d_x, d_y, d_z);
59      |   cudaDeviceSynchronize();
60      }
61      auto end = high_resolution_clock::now();
```

# Reading List

- <https://colfaxresearch.com/skl-avx512>
- <http://svmoore.pbworks.com/w/file/fetch/70583970/VectorOps.pdf>
- <https://software.intel.com/sites/landingpage/IntrinsicsGuide/>
- <https://www.youtube.com/playlist?list=PLAwTw4SYaPnFKojVQrmyOGFCqHTxdfv2>
- <https://medium.com/@smallfishbigsea/basic-concepts-in-gpu-computing-3388710e9239>
- <https://devblogs.nvidia.com/easy-introduction-cuda-c-and-c/>
- <https://devblogs.nvidia.com/six-ways-saxpy/>