

# Parallel Programming | ParProg

## Zusammenfassung

### 1. MULTI-THREADING

**Parallelism** (Subprograms run simultaneously for faster programs) vs. **Concurrency** (interleaved execution of programs for simpler programs)

**Process: Program under Execution**, in own address space.

(heavyweight; Proc. Process isolation and responsiveness, Context: Interprocess communication overhead, expensive in creation, slow context switching and process termination)

**Thread: Parallel sequence** within a process. Sharing the same address space, but separate stack and registers (lightweight because most of the overhead already happened in the process creation).

**Multi-threads:** Changes made by one thread to shared resources will be seen by other threads.

**Context switch:** Required when changing threads. **Synchronous** (Thread waiting for condition) or **Asynchronous** (Thread gets released after defined time)

**Multitasking: Cooperative** (threads must explicitly initiate context switches, OS scheduler can't interrupt) or **preemptive** (scheduler can asynchronously interrupt thread via time interrupt)

**JVM Thread Model:** Java Virtual Machine is a **process** in the OS. It runs as long as threads are running (Exception: threads marked with `setDaemon(true)` will not be waited upon). Threads are realized by the **Thread class** and the **Interface Runnable**. Code to be run in a Thread is within an overridden `run()`. Need to be started manually with `start()`. Threads run in a thread-deterministic order by default.

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

**Runnable** (Thread gets released after defined time)

```
// Java
class BankAccount {
    private int balance = 0;
    // Entry in the monitor
    public synchronized void withdraw
    (int amount) {
        while (amount > balance) { // not if
            wait(); // wait on while-condition
        }
        balance -= amount;
        // release / leave monitor
        public synchronized void deposit
        (int amount) {
            while (value <= 0) { value--; }
            public synchronized void release() { value++; notify(); } }
}
```

**3. SPECIFIC SYNCHRONIZATION PRIMITIVES**

**3.1. SEMAPHORE**

Allocation of a limited number of free resources. Is in essence a **counter**. If a resource is **acquired**, counter—, if a resource is **released**, counter+. Can wait until resource becomes available. Can also acquire/release multiple permits at once atomically.

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**Semaphore** (Thread gets released after defined time)

**4. CONCURRENCY HAZARDS**

**4.1. RACE CONDITIONS**

**Insufficiently synchronized access to shared resources.** The order of events affects the correctness of the program. Leads to **non-deterministic behavior**. Can occur without data race, but data race is often the cause.

**Race Condition without data race:** The critical section is not protected. Data race is eliminated using synchronization, but there is no synchronization over larger blocks, so race conditions are still possible (e.g. non-atomic incrementing).

**4.2. DATA RACE**

Two threads in a single process access the **same variable** concurrently without synchronization, at least one of them is a **write access**. The reading threads may read the old value. **Synchronize everything!** Must help and is expensive. So no.

**4.3. THREAD SAFETY**

Dispensable cases in synchronization: **Immutable Classes** (Declaring all fields private and final and don't provide setters), **Read-only Objects** (Read-only accesses are thread-safe)

**Thread Safety:** Avoidance of Data Races. When no sharing is intended, give each thread a private copy of the data. When sharing is important, provide explicit synchronization.

**4.4. DEADLOCKS**

Happens when threads lock each other out, prohibiting both from running. Programs with potential deadlock are not considered correct. Threads can suddenly block each other.

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**Deadlock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock. Should be avoided!

**5.1.4. Pairwise sum (recursive)**

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**PairwiseSum** extends `RecursiveAction`

**GUI Premise:** No long operations in UI events, or else blocks UI. No access to UI-elements by other threads, or else incorrect (Exception in .NET & Android, Race Condition in Java/Swing).

**6.1.1. Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

**Non-Blocking UI Implementation**

## 8. GPU (GRAPHICS PROCESSING UNIT)

**End of Moores Law:** We can no longer gain performance by "growing" sequential processors. Instead, we *improve performance* by running codes in *parallel* on *multi-core CPUs* (low latency) and *many cores* for throughput. Programming is more difficult.

**GPU's** are specialized electronic circuits designed to accelerate the computation of *computer graphics*. They are faster than CPUs for suitable algorithms on large datasets. *Useful* for calculations which consist of *multiple independent sub-calculations*, not very useful for calculations where the results rely on the previous results (like Fibonacci).

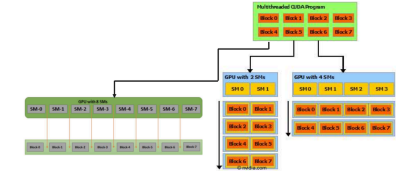
**High Parallelization:** A CPU offers few cores (4, 8, 16, 64) and is very fast. Programming is easier. A GPU offers a very large number of cores (312, 1024, 3584, 5760) and has very specific slower processors. It is optimized for throughput. Programming is more difficult.

**GPU Structure:** A GPU consists of multiple *Streaming Multiprocessors (SM)* which in turn consist of multiple *Streaming Processors (SP)* (e.g. 30 SM @ 152 SPs per SM).

**SIMD:** Single Instruction Multiple Data. The *same instruction* is executed simultaneously on *multiple cores* working on *different data elements* (vector parallelism). Saves fetch & decode instructions.

**SISD:** Single Instruction Single Data. Purely *sequential* calculations.

**SIMT:** Single Instruction Multiple Threads. The same instruction is executed in different threads over different data.



## 8.1. LATENCY VS. THROUGHPUT

**Latency:** *Elapsed time* of an event (waiting from point A to B takes one minute, the latency is one minute).  
**Throughput:** *The number of events* that can be executed per unit of time (*Bandwidth*).  
There is a *tradeoff* between latency and throughput. Increasing throughput by pipelined processing, latency must also increase. All pipeline stages must operate in *lockstep*. The *rate of processing* is determined by the *slowest step*.

**Pipelining:** Run processes in an overlapping manner.  
**Example:** A program consists of two operations: Transfer data from CPU memory to GPU memory (T<sub>1</sub> units = 20ms). Execute computation on the device (T<sub>2</sub> units = 60ms). What is the *latency* (non-pipelined)? 20 + 60 = 80ms.  
What is the *throughput* (pipelined)? Every 60ms an operation is finished. Throughput = 1/60 operations/ms.

## 8.2. CPU VS GPU

CPUs	GPUs
- Low latency	- Can execute <i>highly parallel</i> data operations
- Few but a lot of <i>optimized cores</i>	- Simple but a lot of <i>cores</i> with cache per core
- General purpose	- Very useful for problems which consist of a lot of <i>independent data elements</i>
- Architecture and Compiler help to run any code fast	- Efficiency must be achieved by <i>optimizing the program</i>

Aim: low latency per thread

Aim: high throughput

## 8.3. NUMA MODEL

**NUMA stands for Non-Uniform Memory Access.** CPUs on both shared and GPU devices each have local memories. There is *no common main memory* between the two, so *explicit transfer* between CPU and GPU is needed. There is also *no garbage collector* on the GPU.

## 8.4. CUDA

Computer Unified Device Architecture. Is a *parallel computing platform* and an API for Nvidia GPUs that allows the host program to use GPUs for general purpose processing.

### CUDA Execution steps

1. *cudaMalloc()*: GPU memory allocation
2. *cudaMemcpy()*: Data transfer to GPU (host-to-device)
3. *Kernel* <C1, N>: Kernel execution
4. *cudaFree()*: Deallocate GPU memory

### Example: Array addition

```
for (int i = 0; i < N; i++) { C[i] = A[i] + B[i]; } // sequential
(i = 0, ..., N-1): C[i] = A[i] + B[i]; // parallel using n threads
```

### CUDA Kernel

A kernel is a function that is executed *n* times in parallel by *m* different CUDA threads.

```
__global__
void VectorAddKernel(float *A, float *B, float *C) {
    int i = threadIdx.x; C[i] = A[i] + B[i];
}
```

// kernel invocation on CPU  
VectorAddKernel<<1, N>>(A, B, C); // N is amount of threads

Only the GPU knows when the task is finished.

### Boilerplate Orchestration Code

```
void CudaVectorAdd(float *A, float *B, float *C, int N) {
    size_t size = N * sizeof(float);
    float *d_A, *d_B, *d_C; // data on GPU
    cudaMalloc(&d_A, size); cudaMalloc(&d_B, size); cudaMalloc(&d_C, size); // 1.
    cudaMemcpy(d_A, A, size, cudaMemcpyHostToDevice); // 2.
    cudaMemcpy(d_B, B, size, cudaMemcpyHostToDevice); // 3.
    VectorAddKernel<<1, N>>(d_A, d_B, d_C, N); // 4.
    cudaMemcpy(d_C, C, size, cudaMemcpyDeviceToHost); // 5.
    cudaFree(d_A); cudaFree(d_B); cudaFree(d_C);
}
```

## 8.5. PERFORMANCE METRICS

The performance is either limited by *memory bandwidth* or *computation*. **Compute Bound:** Throughput is limited by calculation

(Cores are at the limit, but the memory bus could transfer more data).

**Memory Bound:** Throughput is limited by data transfer (Memory bus bandwidth is at the limit, but cores could process more data).

**Arithmetic intensity (AI):** Defined as *FLOPS* (Floating Point Operations per Second) per Byte. The higher, the better.

Number of operations = FLOPS  
Number of transferred bytes = Bytes

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW_{memory}}$

Memory Bound:  $\frac{FLOPS}{Bytes} < \frac{BW_{compute}}{BW_{memory}}$

Example: for(i=0; i<N; i++) { x[i] = x[i] + y[i] + x[i]; }  
Reads 4 bytes from memory and writes 4 bytes to memory. That makes 2 reads and 1 write (a is used twice, but read only once). In case x, y and z are ints, we have 12 (4+4) bytes transferred and 2 arithmetic ops (+, =).

The arithmetic intensity is therefore 2/12 = 1/6.

Compute Bound:  $\frac{FLOPS}{Bytes} > \frac{BW_{compute}}{BW$