

# 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, own address space (heavy-weight: Proc. Process isolation and responsiveness, Cons: 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 can 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, scheduler can't interrupt) or *pre-emptive* (scheduler can asynchronously interrupt thread via time interrupt).

**JVM Thread Model:** JVM is a *process* in the OS. It runs as long as threads are running (Exception: threads marked as done with *setDaemon(true)* will not be waited upon). Threads are realized by the *class* and the *interface* *Runnable*. Code to be run in a Thread is within the overridden run()

#### 1.1. STARTING A THREAD

As a **anonymous function (Lambda)**: var myThread = new Thread() -> { /\* thread behaviour \*/; };

As a **named function**: var myThread = new Thread() -> AssynFunction(); myThread.start();

With **explicit runnable implementation**:

```
class MyThread implements Runnable {
    @Override
    public void run() { /* thread behavior */ }
    var myThread = new Thread(new MyThread());
    myThread.start();
}
```

In **C#**: var myThread = new Thread(() => { ... }); myThread.Start(); ... myThread.Join();

**Multi-Thread Example (no synchronization)**

```
public class MultiThreadTest {
    public static void main(String[] args) {
        var a = new Thread(() -> multiPrint("A"));
        var b = new Thread(() -> multiPrint("B"));
        a.start(); b.start(); System.out.println("main finished");
    }
    static void multiPrint(String label) {
        for (int i = 0; i < 10; i++) { System.out.println(label + " " + i); }
    }
}
```

The **printout** of this function varies. It can be all possible combinations of A's and B's due to the non-deterministic scheduler. **Thread Join:** Waiting for a thread to finish (t2.join() blocks as long as t2 is running).

var a = new Thread() -> multiPrint("A"); var b = new Thread() -> multiPrint("B"); System.out.println("threads start"); a.start(); b.start(); ... a.join(); b.join(); System.out.println("threads joined");

**Thread States:** *Blocked* (Thread is blocked and waiting for a monitor lock), *New* (Thread has not yet started), *Runnable* (Thread is runnable (ready to run or running)), *Terminated* (Thread is terminated), *Timed Waiting* (Thread is waiting with a specified waiting time), *sleeps* (Thread.sleep(ms)), *Waiting* (Thread is waiting).

**Yield:** Thread is done processing for the moment and hints to the scheduler to release the processor. The scheduler can ignore this. Thread enters into ready-state. (Thread.yield())

**Interrupts:** Threads can also be interrupted from the outside (myThread.interrupt(), Thread can decide what to do upon receiving an interrupt). If the thread is in the sleep(), wait() or join() methods, a *InterruptedException* is thrown. Otherwise a flag is set that can be checked with *isInterrupted()* / *isInterruptedOrThrow()*.

**Exceptions:** Exceptions thrown in run() can't be propagated to the Main thread. The exception needs to be handled within the code executed on the thread.

**Thread Methods:** *currentThread()*: (Reference to current thread), *setDaemon(true)*: (Mark a demon thread), *getId()/getName()*: (Get thread ID/name), *isAlive()*: (Tests if thread is alive), *getState()*: (Get thread state)

### 2. MONITOR SYNCHRONIZATION

Threads run arbitrarily. **Restriction of concurrency** for deterministic behavior.

**Communication between threads:** Sharing access to fields and the objects they refer to. Efficient, but poses problems: *Thread interference* and *memory consistency errors*.

**Critical Section:** Part of the code which must be executed by only 1 thread at a time for the values to stay consistent. Implementation with *mutual exclusion*.

```
class BankAccount {
    private int balance = 0;
    public void deposit(int amount) {
        // enter critical section
        synchronized(this) {
            this.balance += amount;
        } // exit critical section
    }
}
```

**Consistent:** Implementation with *mutual exclusion*.

**synchronized:** Body of method with the *synchronized* keyword is a critical section. Guarantees *memory consistency* and a *happens-before relationship*. Important for two invocations of a synchronized method on the same object to interleave. Other threads are *blocked* until the current thread is done with the object. Every object has a *Lock* (Monitor-Lock). Maximum 1 thread can acquire the lock. Entry of a synchronized method acquires the lock of the object, the exit releases it.

public synchronized void deposit(int amount) { this.balance += amount; }

Can also be used within a method, the object that *should be locked* must be specified. *synchronized(this)* { this.balance += amount; }

**volatile:** *synchronized* block: End of the block, return, unhandled exceptions

**2.1.1. Monitor Lock**

A monitor is used for a *small mutual exclusion*. Only one thread operates at a time in Monitor. All non-private methods are synchronized. Threads can *wait* in Monitor for condition to be fulfilled. Can be *inefficient* with different waiting conditions, has *fairness problems* and *no shared locks*.

**Recursive Lock:** A thread can acquire the same lock through recursive calls. Lock will be free by the last release.

**Busy Wait:** Running yield or sleep in a loop doesn't release the lock and is inefficient. Use wait. wait(): Waits on a condition. Temporarily releases Monitor-Lock so that other threads can run. Needs to be wrapped into a while loop to check if the while condition has been met.

**WakeUp signal:** Signalling a condition in thread in Monitor. notify() signals any waiting thread (infinite if all threads wait for the same thing, so it does not matter which one comes next - uniform waits or if only one single thread can continue like in a turnstile), notifyAll() wakes up all threads (i.e. one deposit can satisfy multiple withdrawals, does not guarantee fairness). If a thread is finished up, it goes from the *inner waiting room* (waiting on a condition) into the *outer waiting room* (Thread has not started yet) where it waits for entry to the Monitor. There is no shortcut.

if (this.monitorStateException is thrown if notify, notifyAll or wait is used outside synchronized.

```
// Java
class BankAccount {
    private int balance = 0;
    // Entry in the monitor:
    public synchronized void withdraw(int amount) {
        threads interruptedException {
            while (amount > balance) { // not if
                wait(); // wait on a condition
            }
            balance -= amount;
        } // release / leave monitor
        public synchronized void deposit(int amount) {
            balance += amount;
            notifyAll(); // Wakes up all waiting
        }
        threads in monitor inner waiting area
    }
}
```

### 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*, count--, if a resource is *released*, count++. Can wait until resource becomes available. Can also acquire/release multiple permits at once atomically.

```
public class Semaphore {
    private int value; public Semaphore(int initial) { value = initial; }
    public synchronized void acquire() throws InterruptedException {
        while (value <= 0) { wait(); } value--;
    }
    public synchronized void release() { value++; notify(); }
}
```

**General Semaphore:** new Semaphore(N) (Counts from 0 to N for limited permits to access a resource)

**Binary Semaphore:** new Semaphore(1) (Counter 0 or 1 for mutual exclusion, open/locked)

**Fair Semaphore:** new Semaphore(N, true) (Uses FIFO Queue upon fair. Slower than non-fair version)

Semaphores are *powerful*, any synchronization can be implemented. But relatively low-level.

```
class BoundedBuffer<T> {
    private Queue<T> queue = new LinkedList<>();
    private Semaphore upperLimit = new Semaphore(capacity, true); //How many free?
    private Semaphore lowerLimit = new Semaphore(0, true); // how many full?
    public void put(T item) throws InterruptedException {
        upperLimit.acquire(); // No. of free places - 1
        synchronized (queue) { queue.add(item); }
        lowerLimit.release(); // // No. of full places + 1
    }
    public T get() throws InterruptedException {
        T item;
        lowerLimit.acquire(); // No. of full places - 1
        synchronized (queue) { item = queue.remove(); }
        upperLimit.release(); // // No. of free places + 1
        return item;
    }
}
```

#### 3.2. LOCK & CONDITION

Monitor with *multiple waiting lists* and conditions. Independent from Monitor locks.

**Lock-Object:** Lock for entry in the monitor (owner waiting room)

**Condition-Object:** Wait & Signal for a specific condition (owner waiting room)

**ReentrantLock:** Class in Java, *alternative to synchronized*. Allows multiple locking operations by the same thread and supports nested locking (Thread is able to re-enter the same lock).

**Condition:** Factors out the Object monitor methods (wait, notify and notifyAll) into distinct objects to give the effect of having multiple wait-sets per object, by combining them with the use of arbitrary lock implementations. A *Condition* replaces the use of the *Object monitor methods*.

**condition.await():** Throws an *InterruptedException* if the current thread has its interrupted status set on entry to this method or is interrupted while waiting (final try frees the lock in case of interrupt).

**Buffer with Lock & Condition**

```
class BoundedBuffer<T> {
    private Queue<T> queue = new LinkedList<>();
    private Lock<T> monitor = new ReentrantLock(true); // fair queue
    private Condition nonFull = monitor.newCondition();
    private Condition nonEmpty = monitor.newCondition();
    ...
    public void put(T item) throws InterruptedException {
        monitor.lock(); // Lock queue
        try { while (Queue.size() > Capacity) { nonFull.await(); }
            queue.add(item); nonEmpty.signal(); } finally { monitor.unlock(); } }
    } // signalAll() if uniform waiters
    public T get() throws InterruptedException {
        monitor.lock(); // wait for queue to be filled & signal to other queue
        try { while (Queue.size() == 0) { nonEmpty.await(); }
            T item = queue.remove(); nonFull.signal(); return item; }
        finally { monitor.unlock(); } // // always release Lock, even after Exception
    }
}
```

#### 3.3. READ-WRITE LOCK

Mutual exclusion is *unnecessary for read-only* threads. So one can *select* allow parallel reading access, but implement mutual exclusion for write access.

	Parallel	Read	Write
Read		Yes	No
Write		No	No

ReadWriteLock rwl = new ReentrantReadWriteLock(true); // true for fairness

rwl.readLock().lock(); // shared Lock  
// read-only accesses  
rwl.readLock().unlock();  
rwl.writeLock().lock(); // exclusive Lock  
// write (and read) access  
rwl.writeLock().unlock();

#### 3.4. COUNT DOWN LATCH

Synchronization with a counter that can only *count down*. Threads can wait until counter <= 0, or they can count down. The Latches can only be used once.

var ready = new CountDownLatch(N); var start = new CountDownLatch(1);

ready.countDown(); // wait for N cars var ready.await(); // wait for all cars ready

start.await(); // await race start start.countDown(); // start the race

#### 3.5. CYCLIC BARRIER

Meeting point for fixed number of threads. Threads wait *until everyone arrives*. is *reusable*, threads can synchronize in multiple rounds at the same barrier (simplifies example above).

var start = new CyclicBarrier(N); start.await(); // all cars race as they're here

#### 3.6. EXCHANGER

Rendez-Vous: Barrier with *information exchange* for 2 parties. Without exchange: new CyclicBarrier(2), with exchange: Exchanger, exchange something. The Exchanger blocks until another thread also calls exchange(), returns argument x of the other thread.

```
var exchanger = new Exchanger<Integer>();
for (int i = 0; i < 2; i++) { count++ } // odd n. of exch.: last one blocks
new Thread(() -> {
    for (int i = 0; i < 2; i++) {
        try {
            int out = exchanger.exchange(in);
            System.out.println(Thread.currentThread().getName() + " got " + out);
        } catch (InterruptedException e) { }
    }
}).start(); }
```

### 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 (i.e. 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*.

**Synchronize Everything?** May not help and is expensive. So no.

#### 4.3. THREAD SAFETY

**Disposable classes in synchronization:** *Immutable Classes* (Declaring all fields private and final and don't provide setters. Read-only Objects (Read-only accessers are thread safe).

**Confinement:** Object belongs to only one thread at a time. *Thread Confinement* (Object belongs to only one thread), *Object Confinement* (Object is encapsulated in other synchronized objects)

**Thread safe:** A data type or method that behaves correctly when used from multiple threads as if it was running in a single thread without any additional coordination (Java collection collections).

**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.

```
// Thread 1 // Thread 2
synchronized(ListA) { synchronized(ListB) {
    synchronized(ListB) { synchronized(ListA) {
        ListA.addAll(ListA); ListA.addAll(ListB);
    } }
}
```

Both threads in this scenario have *locked each other out*, the program cannot continue.

**LiveLock:** Threads have blocked each other permanently, but still execute wait instructions and therefore consume CPU during deadlock.

```
// Thread 1 // Thread 2
b = false; while(a) { ... b = true; a = false; while(b) { ... a = true; }
```

#### 4.4.1. Resource Graph

Thread T waits for Lock T requires Lock of Resource R



Deadlocks can be identified by *cycles in the resource graph*.

**Deadlock Avoidance:** Introduce *linear blocking order*, lock nested only in ascending order. Or use *coarse granular locks*.

(When ordering does not make sense, e.g. block the whole Bank to block all accounts)

#### 4.5. STARVATION

A thread may not get chance to access a resource. **Avoidance:** Use fair synchronization constructs. (Aging, Release policies in previous synchronization constructs. Monitor and Thread priorities have a fairness problem.)

#### 4.6. PARALLELISM CORRECTNESS CRITERIA

**Safety:** No race conditions and no deadlocks, **Liveness:** No starvation

#### 4.7. NET SYNCHRONIZATION PRIMITIVES

Monitor with sync object: private synchronized sync = new(); Lock(sync){ ... }. Uses Monitor.wait(sync), Monitor.putAll(sync). Uses fair FIFO-Queue. Lacks: No fairness flag, no Lock & Condition. **Additional:** ReadWriteLocks for Upgradable Read/Write, Semaphores can also be used at OS level, Mutex. Collections are *not* Thread-safe.

### 5. THREAD POOLS

Threads do have a cost. Many threads *slow down* the system. There is also a **Memory Cost**, because there is a stack for each thread. **Recycle** threads for multiple tasks to avoid this.

**Tasks:** Define *potentially parallel* work packages. Passive objects describing the functionality.

**Thread Pool:** Tasks are queued. A much smaller number of *working threads* grab tasks from the queue and execute them. A task must run to completion before a thread can grab a new one.

**Scalable Performance:** Programs with tasks run faster on parallel machines. This allows the exploitation of parallelism without thread costs. The number of threads can be *adapted* to the system. (Rule of thumb: # of Worker Threads = # processors + 1 (Pending IO Calls))

Any task must *complete execution* before its worker thread is free to grab another task. Exception: nested tasks.

**Advantages:** Limited number of threads (Too many threads slow down the system or exceed available memory), **Thread recycling** (avoid thread creation and release), **Higher level of abstraction** (Disconnect task description from execution), **Number of threads configurable** on a per-system basis.

**Limitations:** Task must not wait for each other (nested sub-tasks), results in deadlocks (if one task T1 is waiting for something the T2, behind him in the Queue should wait, but T2 waits for T1, so dead task occurs)

#### 5.1. JAVA THREAD POOL

```
Task launch
var threadPool = new ForkJoinPool(poolSize);
Future<Integer> future = threadPool.submit() -> { // submit task into pool
    int value = ...; // long calculation + return value;
};
```

#### 5.1.1. Future<T>

Represents a *future result* that is to be computed (*asynchronous*). Acts as proxy for the result that may be not available yet because the task has not finished. Usage via .submit() (submit task into pool), .get() (waits if necessary for computation to complete and then returns its result) and .cancel() (Attempts to cancel execution of this task, removes it from queue). Task that has an unhandled exception occurs. It is included in the ExecutionException thrown by get().

#### 5.1.2. Fire and Forget

Task are started *without retrieving results* later (submit() without get()). Task is run, but Exceptions do not get caught.

#### 5.1.3. Count Prime Numbers

```
// Sequential
int counter = 0; for (int n = 2; n <= N; n++) { if (isPrime(n)) { counter++; } }
// Parallel and Recursive
class CountTask extends RecursiveTask<Integer> { // RecursiveAction: no return value
    private final int lower, upper;
    public CountTask(int l, int u) { this.lower = l; this.upper = u; }
    protected Integer compute() {
        if (lower == upper) { return 0; }
        if (lower + 1 == upper) { return 1; }
        int middle = (lower + upper) / 2;
        var left = new CountTask(lower, middle);
        var right = new CountTask(middle, upper);
        left.fork(); right.fork();
        return right.join() + left.join();
    }
}
```

int result = new CountTask(2, N).invoke(); // invokeAll() to start multiple tasks

#### 5.1.4. Pairwise sum (recursive)

```
class PairwiseSum extends RecursiveAction {
    private final int[] array;
    private final int lower, upper;
    private static final int THRESHOLD = 1; // configurable
    public PairwiseSum(int[] array, int lower, int upper) {
        this.array = array; this.lower = lower; this.upper = upper;
    }
    protected void compute() {
        if (upper - lower > THRESHOLD) {
            int middle = (lower + upper) / 2;
            invokeAll(
                new PairwiseSum(array, lower, middle),
                new PairwiseSum(array, middle, upper));
        } else {
            for (int i = lower; i < upper; i++) {
                array[2*i] += array[2*i+1]; array[2*i+1] = 0;
            }
        }
    }
}
```

#### 5.2. Java ForkJoin Pool

**Jobs:** Get submitted into the *global queue*, which distributes the jobs to the *local queues* of each worker thread. If one thread has no work left, it can "*steal*" work from another thread's local queue instead of the global queue. This *distributes* the scheduling work over idle processors.

#### 5.3. JAVA FORK JOIN POOL

**Special Features:** Fire-and-forget tasks may not finish, worker threads run as daemon threads. Automatic degree of parallelism (Default: As much worker threads as Processors).

#### 5.3.1. NET TASK PARALLEL LIBRARY (TPL)

Preferred way to write multi-threaded and parallel code. Provides public types and APIs in System.Threading and System.Threading.Tasks namespaces. *Efficient default thread pool* (tasks are queued to the ThreadPoo, supports algorithms to provide load balancing, tasks are lightweight), has *multiple abstraction layers* (Task Parallelization: use task explicitly, Data Parallelization: use parallel statements and queries using task implicitly), Asynchronous Programming and PLINQ.

```
// Task with return value in C#
Task<int> task = Task.Run(() -> {
    int total = ...; // some calculation
    return total;
});
Console.WriteLine(task.Result); //Blocks until task is done and returns the result
// Waited Task
var task = Task.Run(() => {
    var left = Task.Run(() -> Count(leftPart));
    var right = Task.Run(() -> Count(rightPart));
    return left.Result + right.Result;
});
static Task<int> Count(..., int) { ... }
```

#### 5.3.1. Parallel Statements in C#

Execute *independent statements potentially in parallel* (Once a task for each item, implicit barrier at the end).

```
Parallel.Invoke(
    () => MergeSort(l,m),
    () => MergeSort(m,r)
);
```

Execute *loop-bodies potentially in parallel* (Once a task for each item, implicit barrier at the end).

```
Parallel.ForEach(task,
    file => Convert(file));
Parallel.For(0, array.Length,
    i => DoCompute(array[i]));
```

**Parallel Loop Partitioning:** Loop with lots of quickly executing bodies. It would be *inefficient* to give each iteration a parallel thread. Instead, TPL *automatically groups multiple bodies* into a single task. There are *4 kinds of partitioning schemes* (Range: equally sized, Chunk: partitions start small and grow bigger, Stripe: n-size portions, where n = Small number, sometimes 17 and Hash: partitioning (default: if input is divisible then range partitioning, else chunk partitioning)).

**PLINQ:** Set of technologies based on the integration of SQL-like query capabilities directly into C#. PLINQ is a parallel implementation of LINQ. Benefits from simplicity and reusability of LINQ with the power of parallel programming by creating segments from its data. Analog to Java Stream API.

from book in bookCollection.AsParallel() where book.Title.Contains("Concurrency") select book ISBN // Random order

firstNumber in InputList.AsParallel().AsOrdered() select InputPrime(number)

// Maintains order but is slower

#### 5.3.2. PLINQ

**LINQ:** Set of technologies based on the integration of SQL-like query capabilities directly into C#. PLINQ is a parallel implementation of LINQ. Benefits from simplicity and reusability of LINQ with the power of parallel programming by creating segments from its data. Analog to Java Stream API.

from book in bookCollection.AsParallel() where book.Title.Contains("Concurrency") select book ISBN // Random order

firstNumber in InputList.AsParallel().AsOrdered() select InputPrime(number)

// Maintains order but is slower

#### 5.3.3. Thread Injection

TPL adds new worker threads *at runtime* every time a work item completes or every 500ms.

**Hit Climbing Algorithm:** Maximize throughput while using as few threads as possible. Measures throughput by scaling number of worker threads. Avoids deadlock with task-dependencies (but forcefully ends non-desired threads. Deadlocks with ThreadLocal, Setters/Threads() are still possible). We should keep *parallel tasks short* to better profit from this automatic performance improvement.

### 6. ASYNCHRONOUS PROGRAMMING

**Unnecessary Synchrony:** Blocking method calls are often used without need (long running calculations, I/O calls, database or file accesses). With an *asynchronous call*, other work can continue while waiting on the result of the long operation.

var task = Task.Run(LongOperation); // other work + int result = task.Result;

**Kinds of Asynchronism:** *Call-centric* ("pass" caller waits for the task end and gets the result, blocking call), *Caller-Centric* ("pass" task hands over the result directly to successor / follow-up task).

**Task Continuations:** Define task whose task is linked to the end of the predecessor task.

```
// C# .NET // Java (there can be multiple Apply/AcceptAsync calls)
Task.ContinueWith( // CompletableFuture
    .Run(task2) .supPLYAsync() -> LongOP) // runAsync for return void
    .ContinueWith(task2) .thenApplyAsync(v -> v2) // returns value
    .ContinueWith(task2) .thenAcceptAsync(v -> ... .println(v)); // returns void
```

**Multi-Continuation:** Continue when *all tasks* are finished: Task.AllOf(task1, task2).ContinueWith(continuation);

Continue when *any* of the tasks is finished (other threads get lost after first thread is done): Task.AnyOf(task1, task2).ContinueWith(continuation);

(Exceptions in fire & forget task get ignored, i.e. Task.Run() { ...; throw e; })

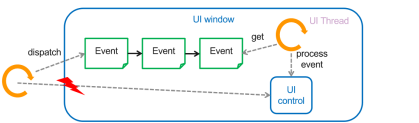
**Exception Handling:** Synchronously Wait() for the whole task-chain at the end. Register for un-observed exceptions with TaskScheduler.UnobservedTaskException (Receives unhandled exceptions from fire & forget tasks). This should be executed as soon as the task object is dead (Garbage Collection).

**Java CompletableFuture:** Modern asynchronous programming in Java. Also has *Multi-Continuation* with CompletableFuture.allOf(future1, future2) and CompletableFuture.anyOf(future1, future2).

#### 6.1. NON-BLOCKING GUI'S

Use case: If a UI is doing a long task, it should not freeze.

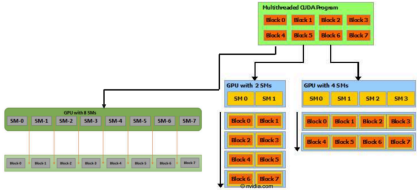
**GUI Thread Model:** Only *single-threaded* (Only a special UI-thread is allowed to access UI-components). The UI thread loops to process the event queue.



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

## 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 code in *parallel* on *multi-core CPUs* (low latency) and *many-core or massively parallel co-processors (GPUs)* (high throughput).  
**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 (e.g. 4, 8, 16, 64) and is very fast. Programming is easier. A GPU offers a very large number of cores (512, 1024, 3840, 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 SM4, 8-192 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.  
**SIMT:** Single Instruction Single Data. Purely *sequential* calculations.  
**SIMD:** 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, starting from point A to B takes one minute, the latency is one minute.  
**Throughput:** The number of events that can be executed at the same time (bandwidth). There is a *tradeoff* between latency and throughput. Increasing throughput by pipelined processing, latency must often 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 (7 units = 20ms). Execute computation on the device (7 units = 60ms).  
What is the latency (non-pipelined)? 20 + 60 = 80ms.  
What is the latency (pipelined)? Every 60ms an operation is finished.  
Throughput =  $\frac{1}{60}$  operations/ms.

## 8.2. CPUs VS GPU

CPUs	GPUs
<ul style="list-style-type: none"><li>- Low latency</li><li>- Few but <i>optimized</i> cores</li><li>- <i>General purpose</i></li><li>- Architecture and Compiler help to run any code fast</li></ul>	<ul style="list-style-type: none"><li>- Can execute <i>highly parallel</i> data operations</li><li>- Simple but a lot of cores with cache per core</li><li>- very useful for problems which consist of a <i>lot of independent data elements</i></li><li>- Efficiency must be achieved by <i>optimizing the program</i></li></ul>
Aim: low latency per thread	Aim: high throughput

## 8.3. NUMA MODEL

NUMA stands for *Non-Uniform Memory Access*. CPUs on host and GPU devices each have local memory. 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 GPU that allows the host program to use GPUs for general purpose processing.

## CUDA Execution

1. *cudaMalloc*: GPU memory allocate
2. *cudaMemcpy*: Transfer results from GPU to CPU
3. *cudaMemcpy*: Data transfer to GPU (HostToDevice) CPU (DevicetoHost)
4. *Kernel <<<1, 1>>>*: Kernel execution
5. *cudaFree*: Deallocate GPU memory

## Example: Array addition

```
(N = N; B = B; A = A; i++) { C[i] = A[i] + B[i]; } // sequential
(i = 0 ; N < i; 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.

```
// kernel definition on GPU
__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, 1, 1>>>, A, B, C); // N is amount of threads
```

Only the GPU knows when the task is finished.

## Boilerplate Orchestration Code

```
void CudaVectorAdd(float* h_A, float* h_B, float* h_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>>>);
    cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_B, h_B, size, cudaMemcpyHostToDevice);
    VectorAddKernel(<<<1, 1, 1>>>(<<<d_A, d_B, d_C, N>>>);
    cudaMemcpy(h_C, d_C, size, cudaMemcpyDeviceToHost);
    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). This is *better and reached if AI Kernel* - AI GPU.  
**Memory Bound:** Throughput is limited by data transfer (Memory bus bandwidth is at the limit, but cores could process more data).  
**Arithmetic intensity:** Defined as FLOPS (Floating Point Operations per second) per Byte. The higher, the better.  
 $\frac{\text{Number of operations}}{\text{Number of transferred bytes}} = \text{FLOPS/Bytes}$

Read `for(<<4, 1, 1>>> i++) { x[i] = x[i] + y[i] + x[i]; }`

Example `for(<<4, y, 2>>> m) { z[m] = x[m] + y[m] + x[m]; }`

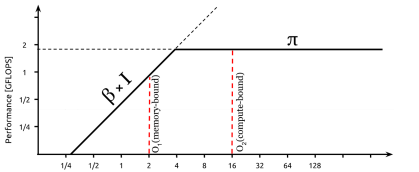
Read `x` and `y` from memory, write `z` to memory. That's 2 reads and 1 write (a 4 byte twice but read only once). In case `x`, `y` and `z` are ints, we have 12 (3 a 4 bytes transferred and 2 arithmetic ops (+, +)).

The arithmetic intensity is therefore  $\frac{2}{12} = \frac{1}{6}$ .

$\frac{\text{FLOPS}}{\text{Bytes}} = \frac{\text{FLOPS}}{\text{Bytes}} = \frac{\text{FLOPS}}{\text{Bytes}}$

## 8.5.1. Roofline model

Provides performance estimates of a kernel running on differently sized architectures. Has three parameters: Peak performance, peak bandwidth vs. arithmetic intensity.  
**Peak performance** is derived from benchmarking FLOPS or GLOPS (ops/10<sup>9</sup> FLOPS). The *peak bandwidth* from manuals of the memory subsystem. The *ridge point* above where the horizontal and diagonal lines meet = minimum AI required to achieve the peak performance.



## 9. GPU ARCHITECTURE

Because there are *so many cores* on GPUs, it is possible to run many threads in parallel *without context switches*. This allows better parallelism without a performance penalty.

## 9.1. COMPIRATION

**Just-in-time Compilation:** The *NVCC compiler* compiles the non-CUDA code with the host C compiler and translates code written in CUDA into *PTX instructions* (assembly language represented as ASCII text). The graphics driver compiles the PTX into *executable binary code*. The assembly of PTX code is *postponed until application runtime*, at which time the target GPU is known. The *disadvantage* of this is the *increased application startup delay*. However, thanks to cache this only happens once (warmup).

**Programming Interface: Runtime** (The *cuda*rt library provides functions that execute on the host to de-allocate device memory, transfer data, etc.) **or Driver API** (The *cuda* driver API is implemented in the device. dL1 or cuda, so level is copied on the system during installation of the driver. This provides an additional level of control by exposing low-level concepts such as CUDA contexts, Open overalls).

**Asynchronous Execution:** The command pipeline in CUDA works asynchronous, commands and data can be transferred from/to the GPU at the same time.

## 9.2. CUDA SIMT EXECUTION MODEL

Single instruction, multiple Threads. The kernel is executed *N* times in parallel by *N* different CUDA threads.

**Blocks:** Threads are *grouped* in blocks. The host can define how many threads each block has (up to 1024). Threads in one block can *interact* with each other but not with threads in other blocks.

**Execution Model:** *One thread runs on one virtual scalar processor* (one GPU core). *One block runs on one virtual multi-processor* (one GPU Streaming Multiprocessor). Blocks must be *independent*.

**Thread Pool Abstraction:** The compiled CUDA program has e.g. 8 CUDA blocks. The *runtime* can choose how to allocate these blocks to multiprocessors. For a larger GPU with 8 SMs, each SM gets one CUDA block. This enables performance scalability without code changes.

**Guarantees:** CUDA guarantees that *all threads in a block run on the same SM at the same time* and that the blocks in a kernel *finish* before any block from a new, *dependent kernel is started*.

**Mapping:** One SM can run several *concurrent* blocks depending on the resources needed. Each *kernel* is executed on *one device*. CUDA supports running *multiple kernels on a device* at one time.

## 9.3. CUDA KERNEL SPECIFICATION

**Specifying Kernel:** `VectorAddKernel<<<gridDim_x, blockDim_x, blockDim_y>>>(<<<A, B, C>>>);`

Dimensions can be 1D, 2D or 3D and specified via `dim3` which is a structure designed for storing block and grid dimensions: `dim3 gridDim(1, 1, 1);` (Unassigned components are set to 1)

`VectorAddKernel<<<dim3(gridDim), dim3(blockDim)>>>(<<<A, B, C>>>);`

**Number of blocks in a grid:** `dim3 n = dim3(x, y, z);`

**Number of threads in a block:** `dim3 blockDim(x, y, z);`

**1D Grid:** We can simply use `VectorAddKernel<<<1, 1, 1>>>` creates 1 block with N threads.

**2D Grid:** `dim3 gridDim(3, 3);` `dim3 blockDim(3, 3);` `VectorAddKernel<<<gridDim, blockDim>>>`

**3D Grid:** `dim3 gridDim(3, 2, 1);` `dim3 blockDim(4, 3, 1);` `VectorAddKernel<<<gridDim, blockDim>>>`

**Device Limits:** Max threads per block: 1024, Max thread dimensions per block: (1024, 1024, 64)

**Max grid size:** (2<sup>14</sup> 783 567, 65 535, 65 535)

**Calculation Examples**

`VectorAddKernel<<<dim3(8, 4, 2), dim3(16, 16)>>>(<<<A, B, d_C>>>);`

**Amount of blocks:** 8 \* 4 \* 2 = 64

**Amount of threads per block:** 16 \* 16 = 256

**Threads in total:** 64 \* 256 = 16 384

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

**Need to round up because for 1025 threads we need 2 blocks**

## 10. GPU PERFORMANCE OPTIMIZATIONS

**Hardware:** A scalable array of multithreaded *Streaming Multiprocessors* (SMs), the *threads* of a thread block execute *concurrently* on one multiprocessor, *multiple thread blocks* can execute *concurrently* on one multiprocessor. When thread blocks *terminate*, new blocks are launched on the free multiprocessors.

## 10.1. MATRIX ADDITION

```
__global__
void MatAddKernel(float* A, float* B, float* C) {
    int column = blockIdx.x * blockDim.x + threadIdx.x;
    int row = blockIdx.y * blockDim.y + threadIdx.y;
    if (row < B_ROWS && col < A_COLS) { // boundary checking
        C[row * A_COLS + col] = A[row * A_COLS + col] + B[row * A_COLS + col];
    }
}
```

```
const int A_COLS, B_COLS, C_COLS = 4;
const int A_ROWS, B_ROWS, C_ROWS = 4;
dim3 block = {2, 2}; dim3 grid = {4, 2};
MatAddKernel<<<grid, blockDim>>>(<<<A, B, C>>>);
```

## 10.2. MATRIX MULTIPLICATION

**Parallelization:** Every Thread computes one element of the result matrix *C*. Can be parallelized because results do not depend on each other.

```
__global__
void multiply(float* A, float* B, float* C) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j = blockIdx.y * blockDim.y + threadIdx.y;
    if (i < N && j < M) { // boundary checking
        float sum = 0;
        for (int k = 0; k < K; k++) {
            sum += A[i * K + k] * B[k * M + j];
        }
        C[i * M + j] = sum;
    }
}
```

## 10.3. MAPPING THREADS / BLOCKS TO GPU WARPS

**Warps:** Blocks are split into warps (Warp = 32 threads) and all threads within execute the same code. If there aren't enough threads to fill a warp, "empty" threads are added. A number of warps constitutes a *thread block*. A number of thread blocks are assigned to a *Streaming Multiprocessor*. The whole GPU consists of several SM.

Thread blocks are scheduled in *parallel* or *sequentially*. Once a thread block is *launched* on a SM, *all of its warps are resident* until their execution finishes. Therefore, a *new block on a SM is not launched until there is a sufficient number* of free registers and shared memory for *all warps* of the new block.

**Warp Execution:** All threads in a warp execute the same instruction *simultaneously*. A SM can accommodate all warps of a block, but only a *subset is running in parallel* at the same time (1 to 24).

**Divergence:** Not all threads of a warp may branch the same way. The branches do *not run simultaneously*, so the other threads need to wait until one branch is finished. So branches within one warp should be *avoided* because of *performance problems* (Branches are born at if / switch / ...)

*// bad case, divergence in same warp // good case, all in warp in same branch*  
*if (threadIdx.x > 1) { else { } if (threadIdx.x / 32 > 1) { else { }*

**DRAM (Dynamic Random Access Memory):** *Global memory of a CUDA device* is implemented with DRAMs. If a GPU kernel accesses data from *consecutive locations*, the DRAMs can supply the data at a *much higher rate* than if a random sequence of locations were accessed.

**Memory Coalescing:** Thread access patterns are critical for performance. If the threads in a warp *simultaneously access consecutive memory locations*, their reads can be combined into a single access (*burst*). Else there are *expensive individual accesses*.

**Coalesced Accesses:** Read/Write the burst in one transaction per warp burst section, swapped read/write within the same burst, only individual elements in the warp burst accessed.

**Coalescing in Use:** `data[Expression without threadIdx.x + threadIdx.x]`

**Coalescing with Matrices:** Matrices get linearized to a 1D array. The row of the matrix should be the longer side so that there are as many coalescing accesses as possible.

## 10.4. MEMORY MODEL

All threads have the access to the same *global memory*. Each thread block has *shared memory* visible to all threads of the block and with the same lifetime as the block (plus higher bandwidth and lower latency than local or global memory but longer latency and lower bandwidth than registers which are private to a thread). Each thread has *private local memory* (in device memory, high latency and low bandwidth, same as global). Constant, texture and surface memory also reside in device memory.

**Memory Hierarchy:** *Shared Memory* (per SM, fast, shared between threads in 1 block, 48 KB, *shared*), *Local Memory* (per thread, fast, private, 48 KB, *shared*), *Global Memory* (in GPU, slow, accessible to all threads, in 64, *cacheable*)

**Registers** (private to a thread, fastest but very small memory).

**Constant memory:** Constant variables are stored in the *global memory* but are *cached*.

**Shared Memory Declaration:** With keyword *\_\_shared\_\_*. A *static array size* is necessary. Limited memory, around 48KB. Multidimensionality is allowed.

**Fast Matrix Multiplication:** By *reducing global memory traffic*. Partition data into subtasks called tiles which fit into shared memory (the row & column that should be multiplied and the result cell). The kernel computation on these tiles must be able to run independently of each other. Because the shared memory is *limited*, load the tiles in several steps and calculate the *intermediate result* from this.

*// Loop over d\_M and d\_N tiles required to compute d\_P element*  
*for (int m = 0; m < M; m++) {*  
*for (int n = 0; n < N; n++) {*  
*for (int k = 0; k < K; k++) {*  
*for (int l = 0; l < L; l++) {*  
*for (int i = 0; i < I; i++) {*  
*for (int j = 0; j < J; j++) {*  
*for (int o = 0; o < O; o++) {*  
*for (int p = 0; p < P; p++) {*  
*for (int q = 0; q < Q; q++) {*  
*for (int r = 0; r < R; r++) {*  
*for (int s = 0; s < S; s++) {*  
*for (int t = 0; t < T; t++) {*  
*for (int u = 0; u < U; u++) {*  
*for (int v = 0; v < V; v++) {*  
*for (int w = 0; w < W; w++) {*  
*for (int x = 0; x < X; x++) {*  
*for (int y = 0; y < Y; y++) {*  
*for (int z = 0; z < Z; z++) {*  
*for (int a = 0; a < A; a++) {*  
*for (int b = 0; b < B; b++) {*  
*for (int c = 0; c < C; c++) {*  
*for (int d = 0; d < D; d++) {*  
*for (int e = 0; e < E; e++) {*  
*for (int f = 0; f < F; f++) {*  
*for (int g = 0; g < G; g++) {*  
*for (int h = 0; h < H; h++) {*  
*for (int i = 0; i < I; i++) {*  
*for (int j = 0; j < J; j++) {*  
*for (int k = 0; k < K; k++) {*  
*for (int l = 0; l < L; l++) {*  
*for (int m = 0; m < M; m++) {*  
*for (int n = 0; n < N; n++) {*  
*for (int o = 0; o < O; o++) {*  
*for (int p = 0; p < P; p++) {*  
*for (int q = 0; q < Q; q++) {*  
*for (int r = 0; r < R; r++) {*  
*for (int s = 0; s < S; s++) {*  
*for (int t = 0; t < T; t++) {*  
*for (int u = 0; u < U; u++) {*  
*for (int v = 0; v < V; v++) {*  
*for (int w = 0; w < W; w++) {*  
*for (int x = 0; x < X; x++) {*  
*for (int y = 0; y < Y; y++) {*  
*for (int z = 0; z < Z; z++) {*  
*for (int a = 0; a < A; a++) {*  
*for (int b = 0; b < B; b++) {*  
*for (int c = 0; c < C; c++) {*  
*for (int d = 0; d < D; d++) {*  
*for (int e = 0; e < E; e++) {*  
*for (int f = 0; f < F; f++) {*  
*for (int g = 0; g < G; g++) {*  
*for (int h = 0; h < H; h++) {*  
*for (int i = 0; i < I; i++) {*  
*for (int j = 0; j < J; j++) {*  
*for (int k = 0; k < K; k++) {*  
*for (int l = 0; l < L; l++) {*  
*for (int m = 0; m < M; m++) {*  
*for (int n = 0; n < N; n++) {*  
*for (int o = 0; o < O; o++) {*  
*for (int p = 0; p < P; p++) {*  
*for (int q = 0; q < Q; q++) {*  
*for (int r = 0; r < R; r++) {*  
*for (int s = 0; s < S; s++) {*  
*for (int t = 0; t < T; t++) {*  
*for (int u = 0; u < U; u++) {*  
*for (int v = 0; v < V; v++) {*  
*for (int w = 0; w < W; w++) {*  
*for (int x = 0; x < X; x++) {*  
*for (int y = 0; y < Y; y++) {*  
*for (int z = 0; z < Z; z++) {*  
*for (int a = 0; a < A; a++) {*  
*for (int b = 0; b < B; b++) {*  
*for (int c = 0; c < C; c++) {*  
*for (int d = 0; d < D; d++) {*  
*for (int e = 0; e < E; e++) {*  
*for (int f = 0; f < F; f++) {*  
*for (int g = 0; g < G; g++) {*  
*for (int h = 0; h < H; h++) {*  
*for (int i = 0; i < I; i++) {*  
*for (int j = 0; j < J; j++) {*  
*for (int k = 0; k < K; k++) {*  
*for (int l = 0; l < L; l++) {*  
*for (int m = 0; m < M; m++) {*  
*for (int n = 0; n < N; n++) {*  
*for (int o = 0; o < O; o++) {*  
*for (int p = 0; p < P; p++) {*  
*for (int q = 0; q &lt*