

Parallel: mehrere Apps gleichzeitig ausgeführt. (2 cores)  
Concurrent: mehrere Apps machen «gleichzeitig» Fortschritt (1 core), werden in Wahrheit aufgeteilt und sequenziell ausgeführt (Context Switch).  
Thread States: running, waiting, ready

## Multi-Thread Programmierung (Java)

```
var myThread = new Thread(() -> { ... });
myThread.start(); // startet Thread & JVM called run()
class SimpleLogic implements Runnable { public void run () { ... } }
var myThread = newThread(new SimpleLogic());
Thread Start und Ende: Richtiger Thread wird erst bei start() erzeugt. Führt run() des Runnable Interface auf. Thread endet beim Verlassen der run()-Methode.
```

Sub-Klasse von Thread:  
class SimpleThread extends Thread { public void run() { ... } }
var myThread = new SimpleThread(); // Konstruktor mit Args möglich  
myThread.start(); myThread.join(); // Warten auf Thread-Ende

## Datenstrukturen

Thread safe (keine Data Races, Lock-Free, nur atomic Operations):  
- Old Java 1.0 Collections (Vector, Stack, HashTable)  
- java.util.concurrent (ConcurrentHashMap, ConcurrentLinkedQueue, CopyOnWriteArrayList, Updates während Iteration sind möglich/weitewise nicht sichtbar).  
Nicht Thread safe:  
- HashSet, TreeSet, ArrayList, LinkedList, HashMap, TreeMap  
ConcurrentLinkedQueue<Integer> q=new ConcurrentLinkedQueue<>();  
q.add(1); int element1 = q.poll();

## Thread Passivierung: Static Methoden Thread-Classes

Thread.sleep(ms): Laufender Thread geht in Wartezustand. Nach Zeitablauf wird er wieder ready  
Thread.yield(): Laufender Thread gibt Prozessor frei, wird aber direkt wieder ready. Nicht mehr nötig, da moderne Hardware preemptive ist.  
Preemptive Scheduling: Thread gibt Ressourcen frei, wenn fertig  
Weitere Thread Methoden

```
static Thread currentThread(): gibt aktiven Thread zurück
void setDaemon(boolean on): Thread als Daemon markieren (default ist false) → in Java JVM wird nicht gewartet bis Thread Terminiert ist (fire & forget)
myThread.interrupt(): Request zum Thread von aussen stoppen, Verhalten kann vom Coder gewählt werden, verwenden wenn Cancel-Policy bekannt
getId(), getName(), isAlive(), getState()
```

## Threads in Java

Heap enthält von mehreren Threads gesehene Shared Ressourcen & Code  
Memory Cost: Jeder Thread hat seinen eigenen Stack und seine eigenen Register. Full register-backup at preemption per thread

Thread States: blocked, new, runnable, terminated, timed, waiting, waiting

## Threads in DotNet

Exceptions in Threads führen zu Abbruch von App.  
Keine Fairness-Flags, keine Lock & Conditions

ReaderWriterLockSlim für Upgradable Read/Write

## Synchronization

Immutability: (Unveränderlichkeit): Objekte mit nur lesendem Zugriff  
Confinement: (Einsperrung): Objekt gehört nur einem Thread zur Zeit  
synchronized-Methode: Objekt erhältet Mutual Exclusion (mutex) Lock static synchronized: Ganze Klasse wird gelockt  
Im Hintergrund wird bei synchronized ein Monitor-lock verwendet.

happens-before relationship: Am Ende von einem synchronized Block ist garantiert, dass alle Änderungen auf das Objekt für alle Threads sichtbar sind.

## Kein synchronization notwendig

Object Confinement: Objekt ist nur «eingehüllt» in synchronized Struktur und braucht somit keine innere synchronization Handhabung. **Read-only Objects**

## Thread Safety Java und DotNet

DotNet: Collection nicht ThreadSafe, außer unter namespace

System.Collections.Concurrent: BlockingCollection<>, ConcurrentDictionary<TKey,TValue>, ConcurrentQueue<>, ConcurrentStack<>

## Monitor

Jedes Object hat ein Lock, nur 1 Thread kann es acquiren, Methoden oder Objects (mit synchronized(object)...)) können synchronisiert werden sleep() & yield() geben Link nicht frei, stattdessen wait() verwenden

Fairness: Nicht zwingend fair, Überholproblem und keine garantierter Reihenfolge trotz synchronized (kein FIFO)

Uniform Waiters: Alle Threads warten auf die gleiche Bedingung.

One-In-One-Out: Bedingung gilt jeweils nur für einen. Nur ein einziger warternder Thread kann weitermachen

Negativ: Effizienz (vielen Bedingungen bei notifyAll()), keine shared Locks, unfair

Signal & Continue: Signalisierender Thread behindert Monitor:

nach Notify() NotifyAll() läuft er im Monitor weiter. Erst nach Ende von synchronized Block können andere Threads den Monitor lock kriegen. Aufgeweckter Thread kommt in äussernen Warteraum, muss neu um Monitor-Eintritt kämpfen.

notify() gebrauchen wenn: Nur eine semantische Bedingung existiert (Uniform Waiters) oder nur einer warternder Thread kann fahren (One-In/One-Out), notify() weckt Thread aus inner waiting room auf.

notifyAll() z.B. bei One-In, Multiple-Out

Spurious Wakeup: Fälschliches Aufwecken eines Threads in vereinzelten Betriebssystemen (z.B. in POSIX Thread API spezifiziert). Schlechtes Design. OS Implementierung vereinfacht, darf Benutzung kompliziert

Waiting Rooms: Jeder Monitor hat einen eigenen einzigen outer waiting room. Pro Bedingung gibt es je einen inneren waiting room. Immer waiting room befinden sich die Threads, die bereit zur Ausführung sind, im inner waiting room befinden sich die, die auf eine Bedingung warten.

Monitor Beispiel:

```
class BoundBuffer<T> {
    private Queue<T> queue = new LinkedList<>();
    private int limit = 1;
    public synchronized void put(T item) throws InterruptedException {
        while(queue.size() == limit) { wait(); }
        queue.add(item);
        notifyAll();
    }
    public synchronized T get() throws InterruptedException {
        while(queue.size() == 0) { wait(); }
        var item = queue.remove();
        notifyAll();
        return item;
    }
}
```

In Dot Net Monitor: FIFO Queue, kein spurious wakeup, Signal & Continue Problem ebenfalls, pulse() / pulseAll() analog Java

private decimal balance;

```
private object syncObj = new();
public void Withdraw(decimal amount) {
    lock(syncObj) {
        while (amount > balance) { Monitor.Wait(syncObj); }
        balance -= amount;
    }
}
public void Deposit(decimal amount) {
    lock(syncObj) {
        balance += amount; Monitor.PulseAll(syncObj); //notifyAll
    }
}
```

## Semaphore

Allgemeine Semaphore (Zähler zwischen 0 bis N)  
new Semaphore(N)  
• Bis zu N Threads können gleichzeitig akquiriert haben  
• Für limitierte Ressourcen, Quotas, Service Throttling etc.

Binäre Semaphore (Zähler nur 0 oder 1)

new Semaphore(1)  
Für mutual exclusion (1 offen, 0 geschlossen)

```
acquire(): Bezieht freie Ressource, wartet, wenn keine verfügbar (Zähler <= 0), sonst Zähler aktualisieren, Passieren, am besten vor try { ... }
acquire(int permits): (Zähler = permits) Blockiert bis Zähler == permits ist
release(): Gibt Ressource frei und benachrichtigt Wartende, vrijgate release(): Gibt Pernits-Zähler + 1 zurück
release(int permits): Zähler += permits, am besten in finally
class BoundBuffer<T> {
    /* lowerLimit # in Queue, upperLimit=verfügbare Kapazität Queue
    upperLimit.counter + lowerLimit.counter = Kapazität Buffer */
    private Queue<T> queue = new LinkedList<>();
    private Semaphore upperlimit = new Semaphore(CAPACITY, true);
    private Semaphore lowerlimit = new Semaphore(0, true);
    private Semaphore mutex = new Semaphore (1, true);
    public void put(T item) throws InterruptedException {
        mutex.acquire();
        queue.add(item); mutex.release();
        or monitor: synchronized(queue) { queue.add(item); }
        public T get() throws InterruptedException {
            upperlimit.acquire();
            mutex.acquire(); queue.add(item); mutex.release();
            or monitor: synchronized(queue) { queue.remove(); }
            upperlimit.release();
            return item;
        }
    }
}
```

Faire Semaphore mit neu Semaphore(N, true);  
• Benutzt FIFO-Warteschlange für Fairness  
• Langsamer als unfair Variante

## Lock & Conditions (Java)

Monitors mit mehreren Wartezeiten für verschiedene Bedingungen.

Lock-Objekt: Sperrre für Eintritt in Monitor, äußerer Warteraum  
Condition-Obj: Wait&Signal für bestehende Bedingung. Thread bleibt im inneren Warteraum

Fairness: new ReentrantLock(true) → fair, default unfair

```
private Queue<T> queue = new LinkedList<>();
private Lock monitor = new ReentrantLock(true);
private Condition nonFull = monitor.newCondition();
private Condition nonEmpty = monitor.newCondition();
public void put(T item) throws InterruptedException {
    monitor.lock();
    try {
        while (queue.size() == Capacity) { nonFull.await(); }
        queue.add(item); nonEmpty.signal();
    } finally { monitor.unlock(); }
}
```

## Read Write Lock

Ermöglicht gleichzeitige Lesezugriffe, da diese kein Lock brauchen, Read/Write, Write/Write, Write/Read sind nicht erlaubt

```
private Collection<String> names = new HashSet<>();
private ReadWriteLock rlock = new ReentrantReadWriteLock(true);
// true -> fairer Lock
public boolean exists(String pat) {
    rlock.readLock().lock(); //shared lock (andere reads erlaubt)
    try {
        return names.stream().anyMatch(n->n.matches(pat));
    } finally { rlock.readLock().unlock(); }
}
public void insert(String name) {
    rlock.writeLock().lock(); // exclusive lock
    try { names.add(name); }
    finally { rlock.writeLock().unlock(); }
}
```

## CountDownLatch

Synchronisationsprimitiv mit Count Down Zähler

• Threads können warten, bis Zähler <= 0 wird

• Threads können runterzählen

Latches sind nur einmalig verwendbar, **kein countUp()** wegen Race Cond.: man würde nicht so schon oben await()'s durchgekommen sind beim Aufruf await(): warten bis CountDown = 0 ist

CountDownLatch: Zähler um 1 dekrementieren

```
var ready = new CountDownLatch(N); // Warte auf N cars
var start = new CountDownLatch(1); // Einer gibt Signal
N Cars: ready.countDown(); ready.await(); start.countDown();
Race Control: ready.await(); start.await();
Java Future
```

Cyclic Barrier - Treffpunkt für fixe Anzahl an Threads

Anzahl Threads muss im Konstruktur angegeben werden. Keine RaceControl benötigt. Reset auf ursprüngliche Anzahl Threads, nachdem Barriere hochging.

cyclicBarrier = new CyclicBarrier(2); // 2 Biker

class Biker extends Thread {
 private final int nr; public Biker(int nr) { this.nr = nr; }
 public void run() {
 System.out.println(nr + " started");
 System.out.println(nr + " reachedCommonPointWithBikersForOthers");
 barrier.await(); System.out.println(nr + " continueJourney");
 }
}

biker1 = new Biker(1); var biker2 = new Biker(2);

biker1.start(); biker2.start(); biker1.join(); biker2.join();

In Dot Net Monitor: FIFO Queue, kein spurious wakeup, Signal & Continue Problem ebenfalls, pulse() / pulseAll() analog Java

private decimal balance;

## Exchanger

2 Threads warten aufeinander (wie new CyclicBarrier(2)) & tauschen Objekt aus  
var exchanger = new Exchanger<Integer>();
for (int i = 0; i < 2; i++) { new Thread(() -> {
 balance -= amount;
}) }
for (int i = 0; i < 5; i++) {
 try { int out = exchanger.exchange(in);
 System.out.println("I got " + out);
 } catch (InterruptedException e) { e.printStackTrace(); }
}

Semaphore - Lock & Condition: L&C hat mehrere waiting rooms

Semaphore - CountDownLatch semaphore: raufl und runter, blockiert bei 0; CountdownLatch: Einweg zu 0, blockiert bei >0 und lässt los bei

CountDownLatch - CyclicBarrier: CyclicBarrier is reusable, CD not

## Concurrency Hazards

### Data Race (spezifische Race Condition)

Nebenläufiger unsynchronisierter Zugriff mit mindestens einem schreibenden Zugriff (Read-Write, Write-Read, Write-Write). Variable wird gelesen, ist aber nicht korrekt (anderer Thread hat sein eigenes Resultat noch nicht in shared memory geladen). Gelöst durch volatile → Java garantiert, dass beim Lesen von volatile Variablen der aktuellste Wert vom shared memory gelesen wird.

### Race Condition (Grund für Lost Update)

Timing oder Keine Regeln bestimmen Korrektheit vom Code. Passiert, wenn Operation nicht atomic ist, z.B. count--; anderer Thread könnte zwischen Auslesen und Schreiben von count diesen schon verändert haben. Konflikt passiert, da die Funktion (nicht wie bei Data Race) mit local & shared memory einen Wert ausliest und zwischenspeichert. Lösung: atomic

### Race Condition ohne Data Race: z.B. bei verschachtelten synchronized Funktionen. Oder: Thread A soll Thread B signalisieren, dass er starten soll. Wenn A failed zu signalisieren, warte B für immer.

### Lost Updates

Wert wird bei nicht synchronisierten Schreibzugriffen überschrieben.

### Deadlock

Gegenseitiges aussperren von Threads.

Zyklus im Graph:

```
public synchronized void transfer(BankAccount to, int amount) {
    balance -= amount;
    to.deposit(amount); // implizit geschachtelter lock
    public synchronized void deposit(int amount) {
        balance += amount; // Thread 1 a.transfer, 2 b.transfer
        a.transfer(b,2); // same synchronized(a){synchronized(b){...}}
        b.transfer(a,5); // same synchronized(b){synchronized(a){}}
```

### Livelock

Threads haben sich gegenseitig permanent blockiert. Führen aber noch Warteinstruktionen aus. Verbrauchen CPU während Deadlock. (Schleifer)

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

Lösung Lock Hierarchy: Locks bekommen ein Level, man kann einen Lock nur acquiren, wenn der Lock vom Level darüber acquired ist

Lösung wenn Hierarchy keinen Sinn macht: weniger fein locken, z.B. die ganze Bank anstatt einzelner Konten locken.

Fairness-Problem: 100% Fairness: alle Threads kommen gleichmäßig vorwärts. 0% Fairness: nur 1 Thread läuft und andere (fast) nie. Starvation:

### Starvation

Thread erhält nie die Chance, um auf eine Ressource zuzugreifen, obwohl Ressource kein Deadlock unterliegt. Andere Threads können ihn dauernd überholen und ihm die Ressource weg schaffen.

Verhindern: faire Synchronisationskonstrukte, faire Semaphores, Lock & Condition, Read-Write Lock. Threads können Prioritäten zugewiesen werden, aber Scheduling ist vom OS abhängig.

### Code-Beispiel mit Starvation:

```
class Switch {
    private AtomicBoolean on = new AtomicBoolean(true);
    private void toggle() {
        on = on.get();
        while (!on.compareAndSet(state, !state)) { }
    }
}
```

neu SignumTask<array> array = signum(array[1]); }

// Aufruf für so viele Threads wie array lang ist (für T=1):

new SignumTask<array> (0, array.length).invoke();

### Keine Überparallelisierung

If (upper - lower > THRESHOLD) { sequential code }

Anzahl Tasks -> Keine Überparallelisierung

Java: für THRESHOLD=20 und array läng: 100: 100/(50,25,25,25) => 12,(12,12,12,12,12,12) => 15 Tasks

⇒ Abhängig von Array-Länge und THRESHOLD

.Net: Anzahl (freier) logischer Prozessoren (Range Partitioning)

Loop partitioning abhängig von Anzahl freier Worker Threads

### C# Net Task Parallel Library (TPL)

```
// Nicht parallelisierbar, wenn auf gleiches Element in Loop ==> ausgeführt wird → Data Race
Parallel.For(0, array.length, (row) => {
    Parallel.For(0, array[0].length, (col) => {
        array[row][col] = (row + 1) * (col + 1); });
}); // Inkrementation um 2:
Parallel.For(0, array.length / 2, idx => {
    int i = 2 * idx;
}); // oder:
Parallel.For(0, array.length, i => SomeTask(i));
Parallel.ForEach(list, file => Convert(file));
```

ThreadPool.SetMaxThreads()

// Task starten / abwarten

TaskInt<task> task = Task.Run(() -> { ...return total; });
task.Wait(); // blockiert (falls keine Rückgabe)

task.Result; // blockiert um liefert Resultat

// kann parallel ausgeführt werden, implizite Barrier am Ende

Parallel.Invoke(() -> res1 = MergeSort1(m, ...));
 () -> res2 = MergeSort2(m, ...));
 // PLINQ: Reihenfolge von Aserufen wird nicht wegen .AsParallel() beibehalten, deshalb .AsOrdered(), um Reihenfolge zu behalten

inPutList.AsParallel().AsOrdered().Select(x =>

IsPrime(x)).ToList();

### Asynchronität

#### Java - CompletableFuture

Paralelle Action muss ein CompletableFuture eingesetzt werden. Main Thread wartet nicht auf CompletableFutures die noch am Laufen sind. Async-Threads laufen im ForkJoinPool.commonPool()

Führt asynchrone Action aus, liefert Rückgabe: CompletableFuture<String>

= CompletableFuture.supplyAsync(() -> "Hello");

f.runAsync(() -> doIt()); führt asynchrone Action aus, ohne Rückgabe.

## Daemon Workers

Normalerweise wartet JVM auf Threads bevor sie terminiert. Wenn Thread ein Daemon ist, bricht JVM diese Threads unkontrolliert ab, sobald der main Thread endet. Daemon Effekt kaput mit .join() bei jedem Thread.

Thread daemonThread = new Thread();
daemonThread.setDaemon(true);

## Java Thread Pool manuelle Implementation

public class ThreadPool {
 private BlockingQueue<Runnable> queue; //custom implementation

private int size; private Thread thread;

public ThreadPool(int size) {
 if (size < 1) {
 throw new IllegalArgumentException("Must be > 0");
 }
 this.size = size;
 queue = new BlockingQueue<>(size);
 for (int i = 0; i < size; i++) {
 Thread thread = new Thread(() -> {
 while (true) {
 Runnable task = queue.remove();
 if (task != null) {
 task.run();
 } else {
 break;
 }
 }
 });
 threads.add(thread);
 }
}

return CompletableFuture.allOf(futures.toArray(new ...

CompletableFuture[]));

// real array now } } Aufruf (blocking): analyzeAsync(..).get();

## DotNet

Caller-centerd (Caller wartet auf Task-Ende): int result = task.Result;

Callee-centric: Task (Über gibt Resultaten der Nachfolger, Task Continuation): Task.Run(t1, WhenAny(t1, t2).ContinueWith(t3).Wait()

Aufrufer von asynchrone Methode ist nicht zwingend während gesamter Ausführung blockiert. Synchron Ausführung bis zu await von aufrufendem Thread.

dannach wird async Methode parallel ausgeführt, die aufrufende Methode wird suspended bei der awaiten Task completed ist, Code nach await = Continuation Rückgabewert: void Task (erlaubt Warten auf Ende), Task<T> Continuation: Compiler: await in nur in asynchrone (Error), asynchrone muss await enthalten (Warnung).

## GUI Thread Model

### Non-Blocking UI (NET)

Caller ist normaler Thread: TPL-Worker-Thread führt Continuation aus

Caller ist UI-Thread: UI-Thread führt Continuation als Event aus

var url = "textfield.text";
Task.Run(url, () => {
 var text = await DoItAsync(url);
 Dispatcher.InvokeAsync(() => { label.Content = text; });
});

// oder mit await/await (Namenskonvention in Violet):

async Task<Void> DoItAsync() {
 return await DownloadAsync(url);
}

### Non-Blocking UI (Java)

button.addActionListener(event -> log(list));

void log(List<String> pending) {
 if (pending.size() == 0) {
 statusLabel.setText("done");
 } else {
 String message = pending.remove(0);
 ForkJoinPool.commonPool().submit(() -> {
 var text = download(url);
 SwingUtilities.invokeLater(() -> {
 var text = download(url);
 SwingUtilities.invokeLater(() -> {
 textArea.setText(text);
 });
 });
 });
 }
}

Rekursive Aufruf (um Reihenfolge zu erhalten):

button.addActionListener(event -> log(list));

void log(List<String> pending) {
 if (pending.size() == 0) {
 statusLabel.setText("done");
 } else {
 String message = pending.remove(0);
 ForkJoinPool.commonPool().submit(() -> {
 var text = download(url);
 SwingUtilities.invokeLater(() -> {
 var text = download(url);
 SwingUtilities.invokeLater(() -> {
 textArea.setText(message + " logged");
 log(pending);
 });
 });
 });
 }
}

## UI Thread: blau, Worker Thread: rot

Wartender UI Aufruf: SwingUtilities.invokeLater(() -> ...);

invokeAndWait(): der aufrufende Thread wird blockiert

invokeLater(): der aufrufende Thread läuft parallel weiter

## Memory Model

### Lock-Free Programmierung

Korrekte nebenläufige Interaktionen ohne Locks

Garanten des Speichermodells nutzen: Ziel: Effiziente Synchronisation

### Probleme (Java Weak Memory Model)

Weak Consistency: Instruktionen werden in verschiedenen Reihenfolgen von

verschiedenen Threads ausgeführt. Ausnahme: Synchronisationen/Speicherbarrieren (Memory Fences)

Grund: Optimierungen von Compiler, Laufzeitssystem und CPUs. Instruktionen werden umgeordnet/wegeoptimiert

### Java Memory Model: Minimale spezifizierte Garantien:

#### Atomicity (Unfallbarkeit)

Separates Lesen / Schreiben ist atomar für primitive Datentypen bis 32 Bit, Objekt-Referenzen, long und bei double mit volatile Keyword atomar

#### Synchronität (Sichtbarkeit)

Locks, Release & Acquire: Änderungen vor Release sind bei acquire sichtbar

volatile: Variable wird in main memory anstatt local beim Thread gespeichert. Beim volatile wird read ist garant, dass alle vorherigen Leseen sichtbar ist. Beim volatile read ist garant, dass alle vorherigen Writes sichtbar sind. volatile Lese- und Schreibzugriffe werden bei Java nicht umgeordnet (starke Ordnung). Race Conditions können immer noch entstehen wenn nicht atomar. Keine Data Races mehr.

atomic: Operationen werden nicht von anderen Threads unterbrochen

final Variablen: Nach Ende von Konstruktor sichtbar

Thread Star/Ende: Start happen-before join(): Return in irgendeinem anderen Thread

#### AtomicInteger

private AtomicInteger count = new AtomicInteger();

addAndGet(int); // Addiert n und gibt neuen Wert zurück

## Volatile in Java

Liveck Beispiel fix: Keine Umrundungen in Java, wegen volatile:

```
volatile boolean a = false, b = false;
Thread 1
a = true;
while ((b) {
    Thread 2
    b = true;
    while ((a) {
```

## Dot Net Memory Model

**Unterschied zu Java:** kein Atomicity für long/double durch volatile (Atomare Instruktionen durch Interlocked class (sowie double/long)).  
Nicht definiert, aber implizit durch Ordering. Ordering: nur Half und Full Fences.

**Halt Fence:** Zugriffe vor volatile write bleiben davor.

Zugriffe nach volatile read bleiben danach.

**Gar keine Umrundung über Memory (Full)**  
**Fence:** Thread.MemoryBarrier(); in Kombination mit volatile Variablen! (.NET) → bei write dann und bei read davorsetzen.

## Unterschied volatile Java VS Dotnet

Atomarität: In Java sind volatile double/long atomar, in .NET nicht.  
Sichtbarkeit: In .NET nicht definiert. **Ordnung:** volatile/volatile sind in Java stark geordnet, in .NET sind in einer Richtung unordnerbar. In .NET gibt es eine definierte partielle Ordnung bei volatile Zugriffen in jedem Fall; bei Java ist dies nur bei einer Schreiben-Lesen Beziehung zwischen Threads.

## Performance Scaling

GPU: High Throughput, CPU: Low Latency

### Performance

**Moore's Law:** alle zwei Jahre die doppelte Menge an Transistoren. Atommengrenzen fordern neue Wege -> Parallelisierung mit mehr Cores

**Performance** ist definiert durch Speicherbandbreite, Rechenbandbreite und Latenzzeit-latency (Gesamteinheit = Speicherzugriff + Rechenzeit)

**Compute bound:** nicht genügend Prozessorleistung. Prozessor ist bottleneck. Memory ist kein bottleneck vom Setting her.

**Memory bound:** mehr Zeit für den Speicherzugriff im Vergleich zu Berechnungen.

**Conclusion:** Software, die compute-bound ist, wird bevorzugt, da es möglich ist, diese Zeit zu verkürzen, indem die Berechnungen optimiert werden. Dies funktioniert nicht für memory-bound Software.

### Arithmetic Intensity

**Arithmetische Intensität (FLOPs/Byte):** z.B. 32Bit INTs  
Code:  $z[i] = y[i] + y[i]$  → 4Byte + 4Byte + 4Byte und ein "\*" = 1 FLOP = 1/12 FLOPs/Byte. Höher ist besser für eine effiziente parallele Nutzung.  
**Throughput (operations/sec)** beschreibt, wie viele Instruktionen / Operationen in bestimmter Zeit gemacht werden können (Frequenz).

**Latency** Zeit, welche für eine einzelne Instruktion / Operation von Anfang bis Ende benötigt wird (Einheit: sec).

T Zeit für eine Instruktion ist:  $\text{Memory Access} + t_{\text{computation}}$

**Compute Bound** bedeutet, dass  $t_{\text{Memory Access}} > t_{\text{computation}}$  gilt.

**Memory Bound** bedeutet, dass  $t_{\text{Memory Access}} < t_{\text{computation}}$  gilt.

Wir können dies auch mit der arithmetischen Intensität ausdrücken:

Al: operations = Bandwidth<sub>compute</sub>, compute bound: #ops / bytes = BW<sub>compute</sub> bytes = Bandwidth<sub>memory</sub>  
Je höher arithmetischen Intensität ist, desto **besser** können wir moderne Parallelisierung durch CPU / GPU ausnutzen.

### Roofline Model

Zieht locality, bandwidth, parallelization paradigms, multicore, etc. in Betracht und zieht memory oder compute-bound hinzu.  
Nach dem Ridge Point: Leistung ist compute-bound (minimales Arithmetic Intensity (AI)), um beste Leistung zu erreichen!

FLOPs=Floating Point Operations per Second  
 $\beta$ =Bandwidth,  $\alpha$ =Arithmet. Intensity,  $\pi$ =Peak performance  
⇒ Ridge point ist operational, i.e. Arithmet. Intensity:  $I = \frac{\pi}{\beta}$

⇒ Attainable Performance P mit given  $\beta$  and  $I$ :  $P = \min(\pi, \beta * I)$  (flops/s)

⇒ Higher Performance:  $P_{\text{faster}} = \text{AI}$  values how long the faster processor stays faster (intersection point)

**GPUs are useful for:** Matrix & Vektor, Rendering, Hash-Cracking

### Weak Scaling (Gustafson's Law)

Die Anzahl Prozessoren und die Problemgröße werden gleichzeitig skaliert → perfekt scaling: gleiche workload pro Prozessor.

**Speedup:**  $s = p = N = (1 - s) * N$

Example: We have 64 Processors, 5% of the program is serial.

scaled weak speedup =  $0.05 + (1 - 0.05) * 64 = 60.85$

### Strong Scaling (Amdahl's Law)

Die Anzahl Prozessoren wird erhöht, während das Problem gleich bleibt → weniger Arbeit pro Prozessor.

$T = \text{total time} = (1 - p)T + Tp$

p = part of the program can be parallelized, s = part serial, N = # of Cores

$$T_p = \frac{T}{N} + (1 - p) * T$$

This law ignores the parallel overhead such as task startup time, interprocess interaction, idling due load imbalance / synch, etc.

$$\text{Speedup: } \frac{T_p}{N} + (1-p)*T = \frac{1-p}{N} + 1 - p = \frac{1-p}{N} + \frac{N-1}{N} = \frac{N+1-p}{N}, \quad s = (1-p)$$

⇒ Max speedup:  $N = \infty$

$$T =$$

$$N * T_N$$

Efficiency:  $90\%$  can be parallel using 8 processors:

$$p = 0.9, s = 0.1, N = 8 \Rightarrow \text{Speedup} = \frac{1}{0.1 * 8} = 4.7$$

### Berechnungen & Formeln

Die Maximale Thread Anzahl lässt sich aus den Informationen vom Device Query (cudaGetDeviceProperties()) berechnen:

Device 0: "NVIDIA RTX A4000"

Multiprocessor count 48  
Maximum threads per multiprocessor 1536  
Maximum threads per block 1024  
Woraus wir nun ableiten können: Limit<sub>total</sub> = min(I<sub>max</sub>, I<sub>0</sub>) = XXX (siehe unten)  
⇒ 200x90pixels → 1 Thread pro Pixel: 180'000 threads  
⇒ **Block Dimension (quadratisch):**  $32 \times 32 \rightarrow (1 \text{ Block} = 1024 = 32 \times 32)$   
⇒ **Grid Dimension:**  $(\frac{200}{32} \times \frac{90}{32}) = 7 \times 29$

⇒ **idle threads (zuviel gestartet):**  $7 \times 29 - 1024 = 200 + 900 \text{ pixels} = 27'872$

Max Threads=I<sub>A</sub> = Num<sub>Threads</sub> / MaxThreadsPerBlock =  $48 \times 1024 = 49'152$

### GPU Structure

SM: Streaming Multiprocessors beinhalten Streaming Processors (=1 GPU core)

Thread wird auf Core, Thread Block auf 1 SM & Grid auf GPU ausgeführt

**NUMA (Non-Uniform Memory Access)**

Daten müssen von Host (CPU) nach Device (GPU) übertragen werden und wieder zurück nach Berechnung.

**Flynn's Classical Taxonomy (SISD, SIMD, MISD, MIMD)**

SISD: Single Instruction Single Data, Sequenzielle Ausführung von Instruktionen

SIMD: Single Instruction Multiple Data (auf SM ausgeführt)

MISD: Multiple Instruction Single Data, Mehrere Nodes führen verschiedene Operationen auf den gleichen Daten auf. **MIMD:** Multiple Instruction Multiple Data, tasks, die von verschiedenen Prozessoren ausgeführt werden, können zu unterschiedlichen Zeiten beginnen oder enden.

**CUDA (Computer Unified Device Architecture) (GPU)**

Device Query: Infos über GPU: Max N of Threads per Block, Max Block Dim, Max Grid Size, Bus Width, Warp Size

**Whole block runs on 1 SM.** All threads in a block must complete.

The more sm in a GPU, the more blocks can be executed in parallel.

The SM also perform Context Switching.

Jeder Thread hat ein private memory, jeder block hat ein shared memory und ein grid hat ein global memory.

Für Variablen- und Methodendeklarationen:

• **\_global:** Läuft auf Device wird vom Host aufgerufen

• **\_device:** Läuft auf Device und wird vom Host aufgerufen (global)

• **\_device\_** shared : device aber nur im Block shared

• **\_host:** Läuft auf Host und wird auch vom Host aufgerufen

• **\_shared:** Memory zwischen Threads in Block

**Memory Access** in den Device Global Memory dauert viel länger als in den Shared Memory ( \_ shared ) eines SM (ca. Faktor 125). Aus diesem Grund ist das Shared Memory oftmals vorzuhaben.

**Memory Coalescing** Eine Coalesced Memory Transaction findet statt, wenn alle Threads in einem Warp nachfolgende Speicherstellen aus dem global memory laden. Denn dann fasst die Hardware die langsamsten access des global memory zusammen (burst) und ist somit schneller. Sieht oftmals so aus:

data[blockIdx.x \* blockDim.x + threadIdx.x]

Synchronisierung in CUDA kann mit **\_syncthreads()** erreicht werden, was alle Threads in einem Block zum Synchronisieren zwingt. Kann nur in einem if-else Block verwendet werden, wenn alle Threads die gleiche Abweigung nehmen (ansonsten Fehlerhälften).

**Wars:** Hardware gruppirt Threads, die dieselbe Anweisung ausführen, in Wards. Ein Ward sind 32 Threads innerhalb eines Thread-Blocks (SIMD).

**Data Partitioning**

threadIdx.x/y/z: Thread index innerhalb Block

blockIdx.x/y/z: Block index innerhalb Grid

gridDim.x/y/z: Grid size

Für Matrizen braucht man dim3, Block size für 1024 = dim3(32, 32)

Indexierung in memory:  $C[i + k] = sum; // equivalent to C[i, j]$

**Boundary Checks**

Wenn wir den boilerplate Code brauchen, um die aufergerundete Gridsize zu bestimmen, dann erstellen wir: gridSize.x \* gridSize.y \* gridSize.z % N mehr Threads, als wir eigentlich wollen. Somit könnten andere Threads beeinflusst werden. Um das zu beheben, muss die Grösse des Problems (B.Z. Arraygröße) bekannt sein und nach einem Overhead geprüft werden:

if ( $i < N$ ) {  $C[i] = A[i] + B[i];$  }

**Thread Divergence**

Happens when threads of the same warp take different branches in the code. E.g., for the signum task we have 3 branches ( $x > 0, x = 0, x < 0$ ) - if elements of a warp have different signs, different branches are being taken. Then, threads of different branches have to wait because alternating the instructions of different branches are being executed. For 3 branches the worst case is 3 times slower.

Number of Threads

VectorAddKernel<dim3(8, 4, 2), dim3(16, 16)>>>(d\_A, d\_B, d\_C);

Start 8 \* 4 \* 2 \* 16 = 16'384 Threads

VectorAddKernel<dim3(16, 1024)>>>(A, B, C, 4096);

Start Grid mit 4 Blocks, welche je 1024 Threads haben, 4 \* 1024 = 4096 Threads. 1024 = 16 \* dim3(4, 1, 1), 2D: dim3(4, 2) = dim3(4, 2, 1), 3D: dim3(4, 2, 2)

**Error Handling**

void handleCudaError(cudaError error) {

if (error == cudaSuccess) {

printf(stderr, "CUDA: %s\n", cudaGetStringError(error));

exit(EXIT\_FAILURE); }

handleCudaError(cudaMemcpy(...)); // Bei normalen CUDA Befehlen

SomeKernel<...>(>>...);

cudaDeviceSynchronize();

handleCudaError(cudaGetLastError()); // Nach Kernel Operation

**Clusters**

A High-Performance Cluster (HPC) uses a hybrid memory model and can run a program on multiple nodes. A node is a standalone computer. Nodes are networked together to form a supercomputer.

- Shared Memory (= NUMA) between Nodes

- Shared Memory for Cores inside a Node (Symmetric Multiprocessing)

**Unified Memory**

Automatischer Speicher Transfer vom Hauptspeicher zum GPU, neue Regeln:

- cuMemAllocManaged(&A, size);

- SomeKernel<...>(>>A);

- cudaDeviceSynchronize();

- cudaFree(A);

**Code**

**PairwiseSum**

// gridSize = 1024; blockSize=(len/2+blockSize-1)/blockSize;

// Uses strideless coalescing

\_global void pairwiseSum(int\* array, int length){

int i = 2 \* (blockIdx.x \* blockDim.x + threadIdx.x);

if (i < length - 1) {

array[i] += array[i + 1];

array[i + 1] = 0; }

}

**Message Passing Interface (MPI)**

Based on Actor/CSP principle. Difference: MPI starts all processes on start.

Similarities: isolated memory (NUMA) and async. inter-node communication.

When more than one processor wants to exchange data, then they need to send messages and communicate with each other. All variables inside a process are private. All processes start and terminate synchronously. Process works in their (virtually) isolated address space. No shared memory between nodes (NUMA).

But there is shared memory for all cores inside a node. Synchronization with barriers possible.

**Communication Modes**

Point to point: One sender, one receiver. Relies on matching send/receive.

Collective communications:

int idx = blockIdx.x \* blockDim.x + threadIdx.x;  
// Make sure we do not go out of bounds (boundary check):  
if (idx < n) { c[idx] = a[idx] + b[idx]; }

int main(int argc, char \*argv[]) { // CPU (Host), N=100000000

int blockSize = 1024; // max

int gridSize = (N + blockSize - 1) / blockSize;

size\_t size = N \* sizeof(double);

\_d\_a, \_d\_b, \*d\_c; // h\_b, \*d\_b, ...

cuMemAlloc(&d\_a, size); // same for d\_b and d\_c

cuMemCopy(d\_a, d\_a, size, cudaMemcpyHostToDevice);

cuMemCopy(d\_b, d\_b, size, cudaMemcpyHostToDevice);

vecAdd(&gridSize, blockSize); // d\_a, d\_b, d\_c, n;

cuMemFree(d\_a); // same for d\_b and d\_c

return 0; }

**Multi Kernel Grid – Matrix Multiplication**

global

void matrixMultKernel(float \*a, float \*b, float \*c) {

\_shared float Asub[TILE\_SIZE][TILE\_SIZE]; // shared memory,

\_shared float Bsub[TILE\_SIZE][TILE\_SIZE]; // i.e. "cache"

int tx = threadIdx.x % tileSize;

int col = blockIdx.x \* tileSize + tx;

int row = blockIdx.y \* tileSize + ty;

int noffiles = (A\_COLS + tileSize - 1) / tileSize;

float sum = 0;

for (int t = 0; t < noffiles; t++) {

if (row < C\_ROWS & col < C\_COLS) {

for (int k = 0; k < C\_SIZE; k++) {

if (t \* TILE\_SIZE + k < A\_COLS) {

sum += Asub[ty][k] \* Bsub[k][tx]; } } }

\_syncthreads(); //matrices are done before calculating sum

if (row < C\_ROWS & col < C\_COLS) {

C[row \* C\_COLS + col] = sum; // 2D -> 1D }

// Aufruf -> receiverRank (durch **printf** in der obigen Funktion)

return A, B, C;

dim3 threadsPerBlock(TILE\_SIZE, TILE\_SIZE);

dim3 blocksPerGrid((C\_ROWS + tileSize - 1) / tileSize, (C\_COLS + tileSize - 1) / tileSize);

dim3 blockDim(TILE\_SIZE, TILE\_SIZE, 1);

dim3 gridDim((C\_ROWS + tileSize - 1) / tileSize, (C\_COLS + tileSize - 1) / tileSize, 1);

MPICH2 MPI API

int MPI\_Init(&argc, &argv);

int MPI\_Finalize();

int MPI\_Comm\_rank(MPI\_COMM\_WORLD, &rank);

int MPI\_Comm\_size(MPI\_COMM\_WORLD, &size);

int MPI\_Send(void \*send\_data, int count, MPI\_Datatype datatype, MPI\_Comm\_rank(MPI\_COMM\_WORLD, &rank), MPI\_Comm\_rank(MPI\_COMM\_WORLD, &dest), int tag);

int MPI\_Recv(void \*recv\_data, int count, MPI\_Datatype datatype, MPI\_Comm\_rank(MPI\_COMM\_WORLD, &src), MPI\_Comm\_rank(MPI\_COMM\_WORLD, &dest), int tag, MPI\_Status \*status);

int MPI\_BARRIER(MPI\_COMM\_WORLD);

int MPI\_WTIME();

int MPI\_WTIMEF();

int MPI\_WCSIZE(MPI\_COMM\_WORLD);

int MPI\_WCSIZEF(MPI\_COMM\_WORLD);

int MPI\_WCSIZE(MPI\_COMM\_WORLD, &size);

int MPI\_WCSIZEF(MPI\_COMM\_WORLD, &size);

int MPI\_WCSIZE(MPI\_COMM\_WORLD, &size);

int MPI\_WCSIZEF(MPI\_COMM\_WORLD, &size);