

The background of the slide is a dense field of 3D-rendered numbers in various shades of blue. The numbers are of different sizes and are scattered across the frame, creating a sense of depth and data. Some numbers are in the foreground, appearing larger and more prominent, while others are in the background, appearing smaller and more faded. The overall effect is a vibrant, digital landscape of numbers.

Big Data Course

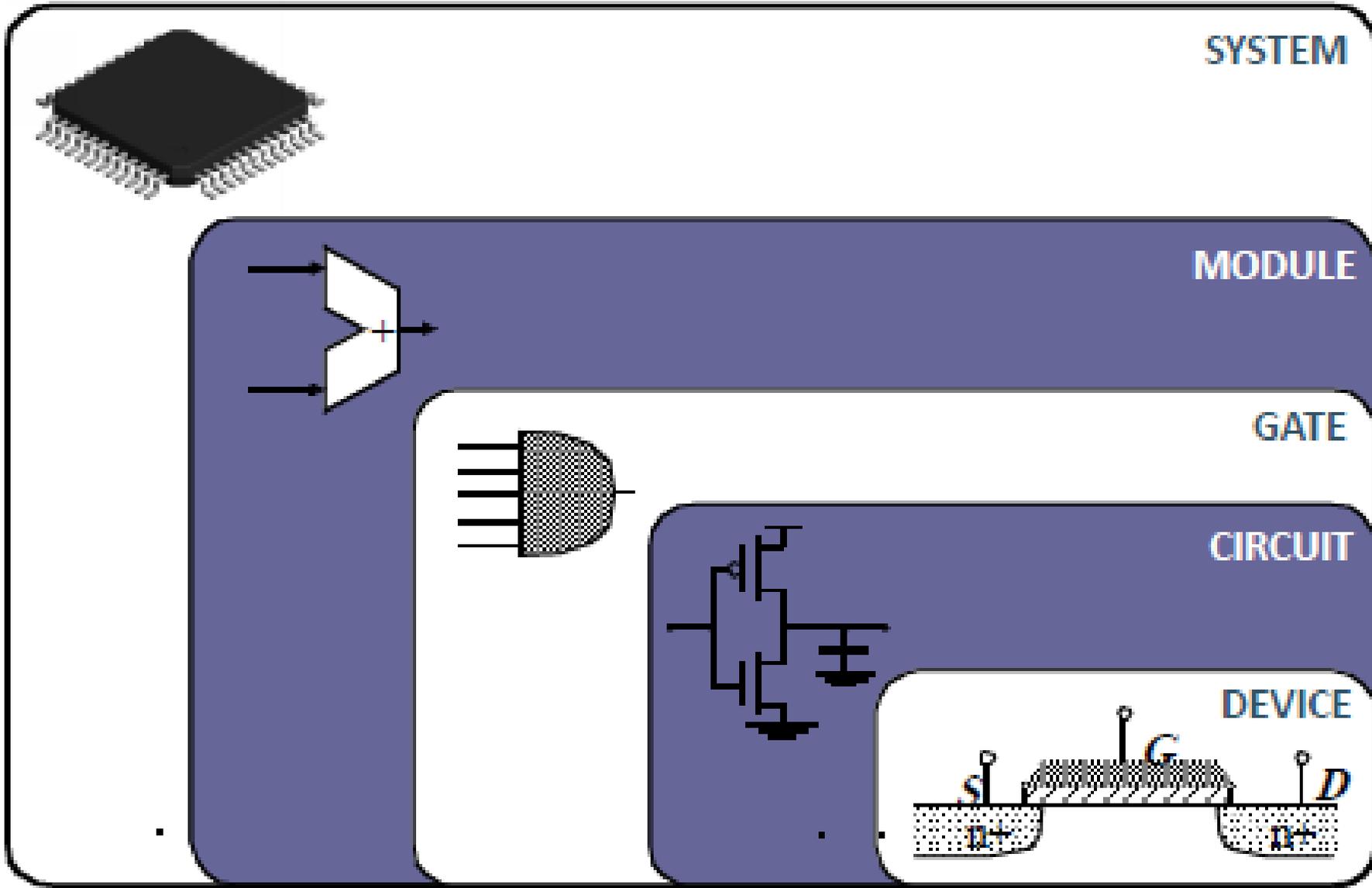
Chap 5- Cache and GPU

Electrical Engineering department of
Amirkabir University of technology
Dr. Mohammadreza Pourfard
August 2025

- **Architecture:** programmer's view of computer
 - Defined by instructions & operand locations
- **Microarchitecture:** how to implement an architecture in hardware

Application Software	
Operating Systems	
Architecture	
Micro-architecture	
Logic	
Digital Circuits	
Analog Circuits	
Devices	
Physics	

Application Software	programs
Operating Systems	device drivers
Architecture	instructions registers
Micro-architecture	datapaths controllers
Logic	adders memories
Digital Circuits	AND gates NOT gates
Analog Circuits	amplifiers filters
Devices	transistors diodes
Physics	electrons



جایگاه درس مدار منطقی و معماری کامپیوتر

This Course

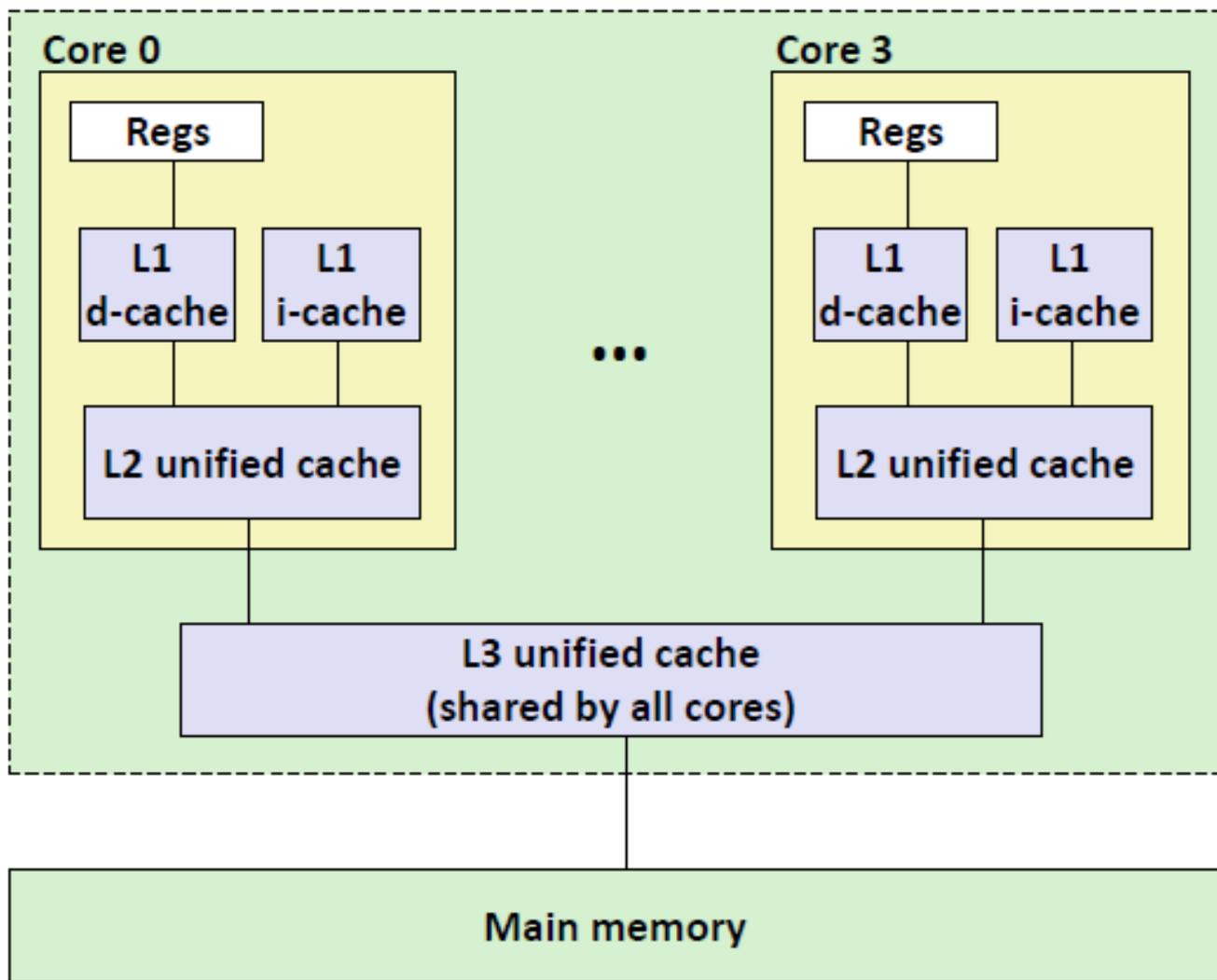
VLSI Course

مفهوم

Cache

Intel Core i7 Cache Hierarchy

Processor package



L1 i-cache and d-cache:

32 KB, 8-way,
Access: 4 cycles

L2 unified cache:

256 KB, 8-way,
Access: 10 cycles

L3 unified cache:

8 MB, 16-way,
Access: 40-75 cycles

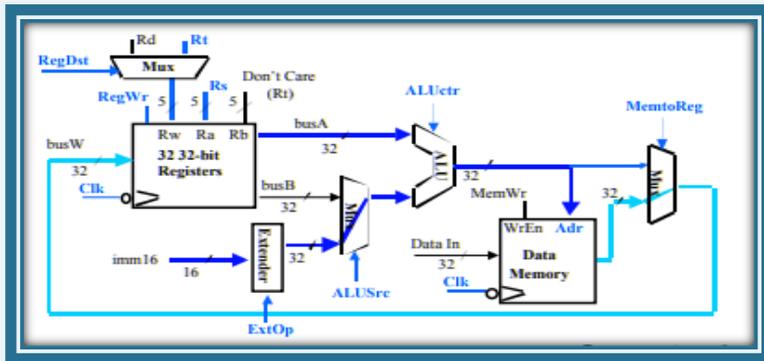
Block size:

64 bytes for all caches

```
temp = v[k];
v[k] = v[k+1];
v[k+1] = temp;
```

```
lw $15, 0($2)
lw $16, 4($2)
sw $16, 0($2)
sw $15, 4($2)
```

```
0000 1001 1100 0110 1010 1111 0101 1000
1010 1111 0101 1000 0000 1001 1100 0110
1100 0110 1010 1111 0101 1000 0000 1001
0101 1000 0000 1001 1100 0110 1010 1111
```



زبان برنامه نویسی سطح بالا

کامپایلر

زبان برنامه نویسی اسمبلی

اسمبلر

زبان ماشین

مفسر ماشین

مشخصه سیگنال کنترل



مقدمه

الگوهای برنامه نویسی

انواع دسته بندی زبان های برنامه نویسی

انواع حوزه های کاری برنامه نویسی

مقایسه ی زبان های برنامه نویسی مختلف

مفهوم Cache

- در کلان داده (Big Data) مفهوم کش کردن (Caching) به معنای ذخیره‌سازی موقت داده‌های پرکاربرد یا پرتکرار در حافظه سریع‌تر مثل RAM تا در درخواست‌های بعدی، سیستم مجبور نباشد دوباره آن داده‌ها را از منبع اصلی (مثل دیسک یا پایگاه داده توزیع شده) بخواند. این کار سرعت پردازش، پاسخ‌گویی، و کارایی کل سامانه را به‌طور چشمگیری افزایش می‌دهد.

مثال کاربردی در کلان داده

- وقتی در Spark داده‌ای را می‌خوانی و چند مرحله محاسبه روی آن انجام می‌دهی، اگر همان داده را چند بار در مراحل بعد نیاز داشته باشی، می‌توانی آن را **cache** کنی تا Spark مجبور نباشد دوباره از فایل اصلی در HDFS یا S3 بخواند.
- مثلاً در PySpark:

Copy code 

python

```
df = spark.read.parquet("hdfs://data/sales/")  
df.cache() # داده در حافظه نگه داشته می‌شود  
df.count() # می‌شود cache اولین بار داده بارگذاری و
```

مزایای Caching در Big Data

- افزایش سرعت پردازش (کاهش زمان اجرای job ها) ⚡
- کاهش بار روی سیستم ذخیره‌سازی (مثل HDFS یا S3) 💾
- تسریع در محاسبات تکراری یا **iterative algorithms** (مثلاً در یادگیری ماشین) ↻
- استفاده بهینه از حافظه توزیع شده بین نودها 🧠

◆ انواع Cache در معماری کلان داده:

نوع Cache	محل قرارگیری	مثال
In-memory cache	حافظه RAM گره‌ها	Spark Cache, Redis
Distributed cache	بین چند سرور	Hadoop DistributedCache, Ignite
Query cache	در سطح پایگاه داده	Hive, Presto, BigQuery
Result cache	نتایج پرس‌وجوی پرتکرار	Druid, ElasticSearch
Application cache	در لایه برنامه یا API	Redis, Memcached

♦ تفاوت Cache با Storage:

ویژگی	Cache	Storage
هدف	سرعت	ماندگاری
دوام داده	موقتی	دائم
محل	RAM یا SSD	دیسک یا HDFS
سرعت دسترسی	بسیار بالا	پایین تر

Caching .۱



Apache Hadoop

۱. Caching در Apache Hadoop

- نقش:
- Hadoop به صورت سنتی روی دیسک (HDFS) کار می کند، نه حافظه. بنابراین caching در Hadoop به معنی ذخیره ی موقت داده های پر کاربرد در حافظه ی نودها (DataNodes) است تا از خواندن مکرر از دیسک جلوگیری شود.
- ابزار و روش ها:
- از نسخه های جدید Hadoop، قابلیتی به نام **DistributedCache** یا **Centralized Cache Management** اضافه شده است.
- داده هایی که در چند Job استفاده می شوند (مثلاً lookup tables) را می توان در حافظه نودها cache کرد.
- هدف:
- کاهش I/O دیسک در MapReduce و افزایش سرعت Job های تکراری.

Caching .۲



Apache Spark

۲. Caching در Apache Spark

- نقش:
- Spark برای محاسبات در حافظه (In-Memory Computing) طراحی شده است، پس caching یکی از اجزای اصلی آن است.
- **نحوه استفاده:**
- در Spark، هر RDD، DataFrame یا Dataset را می‌توان با cache یا persist کرد:

Copy code

python

```
rdd.cache() # در حافظه ذخیره می‌شود  
rdd.persist(StorageLevel.MEMORY_AND_DISK) # در حافظه یا در صورت پر بودن، روی دیسک
```

- **کاربردها:**
- در الگوریتم‌های یادگیری ماشین که داده‌ها بارها تکرار می‌شوند (iterative).
- در تحلیل‌های تعاملی (Interactive Analytics).
- در pipeline‌های ETL با چند مرحله وابسته به یک منبع داده واحد.
- **مزایا:**
- چندین برابر سریع‌تر از Hadoop MapReduce.
- Spark خود به صورت خودکار حافظه نودها را بین jobها تقسیم می‌کند.

Caching .۳



Apache Kafka

۳. Caching در Apache Kafka

- نقش:
- Kafka برخلاف Hadoop و Spark، سیستم پردازش یا تحلیل مستقیم داده نیست بلکه پیام‌رسان توزیع شده (streaming platform) است.
- محل استفاده از cache:
- Kafka از caching در چند لایه استفاده می‌کند:
- **Page Cache** سیستم عامل:
داده‌ها قبل از نوشتن روی دیسک یا هنگام خواندن از آن در حافظه kernel نگه داشته می‌شوند.
- **Broker Cache**
برای کاهش latency در ارسال پیام‌ها به consumerها.
- **Streams API Cache**
در Kafka Streams، نتایج stateful processing (مثل join یا aggregation) در حافظه cache می‌شوند تا نیاز به نوشتن مداوم در state store نباشد.
- هدف:
- کاهش تأخیر در پردازش جریان داده‌ها
- جلوگیری از overload روی دیسک
- پایداری و کارایی بالا در throughput

◆ مقایسه‌ی نهایی

ویژگی	Hadoop	Spark	Kafka
نوع سیستم	Batch Processing	In-Memory & Batch/Streaming	Streaming & Messaging
هدف از Cache	کاهش I/O دیسک	تسریع محاسبات تکراری	کاهش latency پیام‌ها
محل Cache	DataNode Memory	Executor Memory	Stream Cache و OS Page Cache
کنترل توسط کاربر	محدود	کامل (cache/persist)	تا حدی (configurable)
دوام داده	تا پایان Job	تا زمانی که job زنده است	تا زمان retention message

تعريف

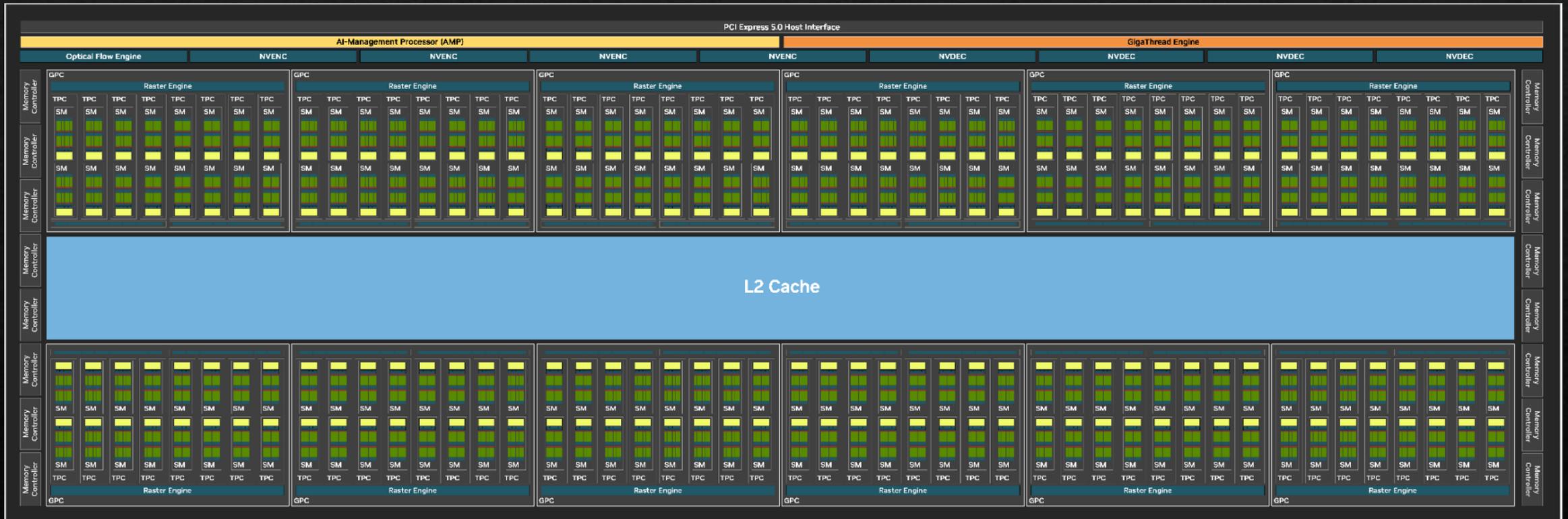
»

SM

»

GPU

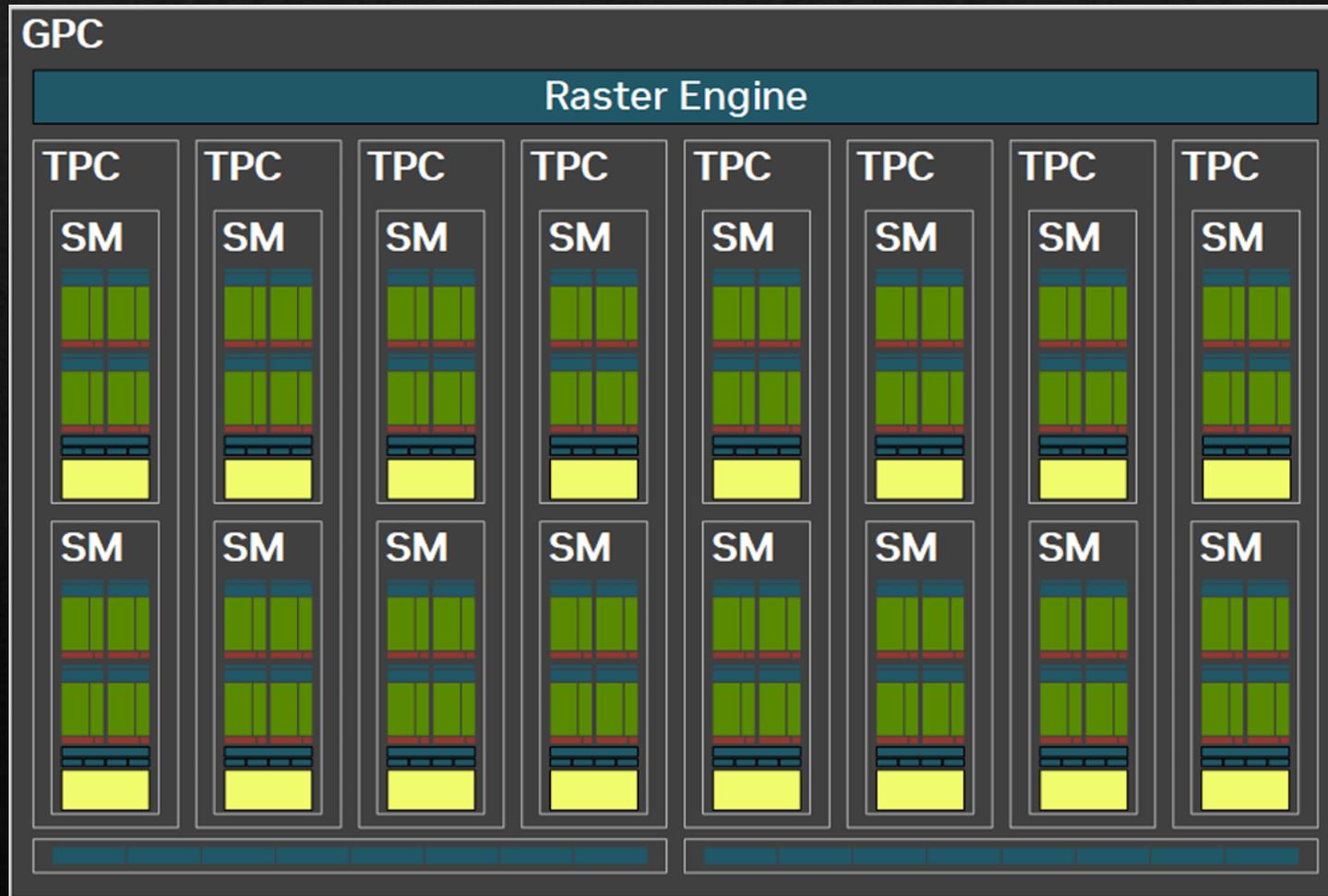
GB202 Die Shot



GB202 Components

- ◇ GPCs
- ◇ Memory Controllers.
- ◇ Cache.
- ◇ AMP & GigaThread Engine.
- ◇ NVENC / NVDEC.
- ◇ Optical Flow Engine.
- ◇ PCI Express 5.0 Host Interface.

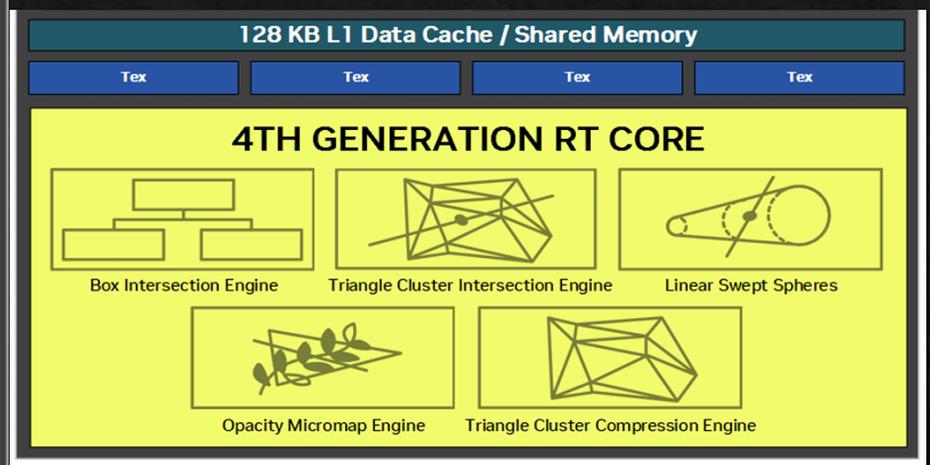
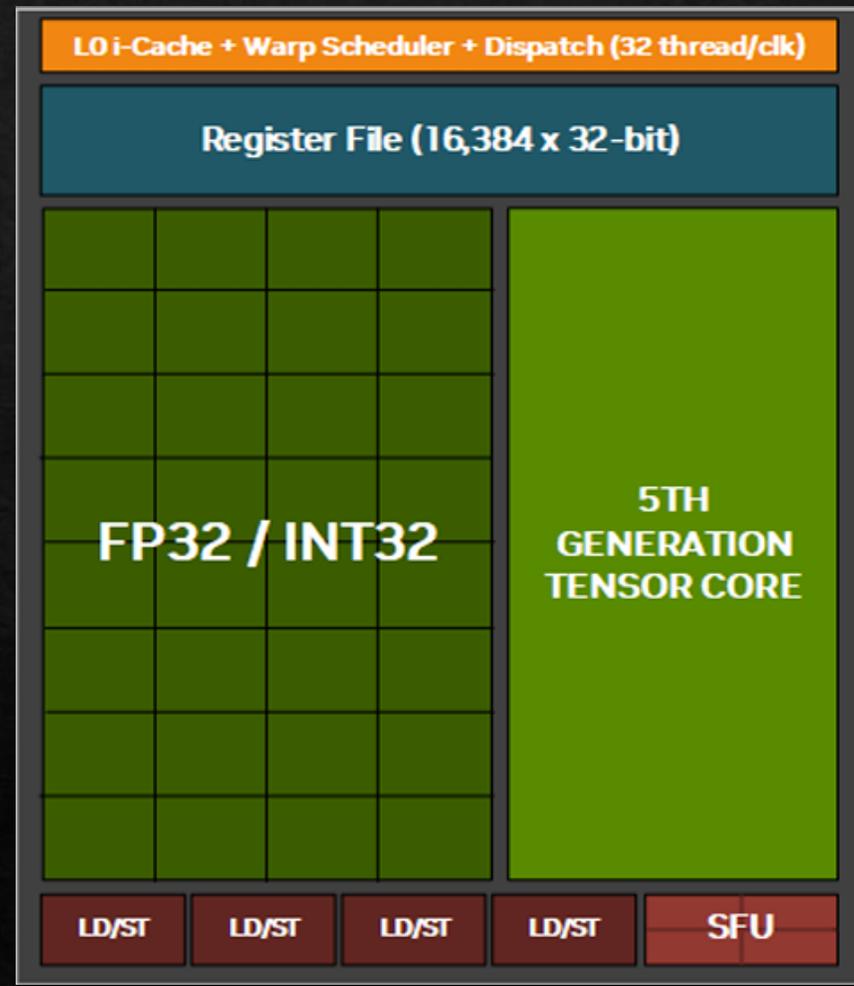
GPC



GPC cont.

- ◇ Fundamental building block of the GPU's architecture.
- ◇ Stands for **G**raphics **P**rocessing **C**lusters.
- ◇ The GB202 has 12 GPCs.
- ◇ Contains:
 - ◇ Raster Engines.
 - ◇ TCPs.
 - ◇ ROPs.
- ◇ Provides scalable graphics and compute performance.

SM



تعریف SM

- **SM (Streaming Multiprocessor)** در واقع واحد پردازش اصلی GPU است.
- هر GPU از دهها تا صدها SM تشکیل شده، و هر SM خودش شامل دهها هسته (CUDA Cores)، حافظهها، و واحدهای محاسباتی اختصاصی است.
- به زبان ساده:
- هر SM را می توان مثل "یک مینی پردازندهی چند هسته ای" درون GPU تصور کرد.

ساختار یک SM (در معماری‌های مدرن مثل H200)

هر SM شامل اجزای زیر است:

مؤلفه	توضیح
CUDA Cores	هسته‌های محاسباتی برای عملیات عددی (مثل جمع، ضرب، ماتریس).
Tensor Cores	واحدهای مخصوص برای شتاب‌دهی به عملیات ماتریسی (در یادگیری عمیق).
Shared Memory / L1 Cache	حافظه‌ی سریع و محلی که بین threadهای همان SM مشترک است.
Register File	سریع‌ترین حافظه، مخصوص هر thread.
Warp Scheduler	زمان‌بند اجرای threadها (واحد کنترل).
Load/Store Units	واحدهای دسترسی به حافظه (برای خواندن و نوشتن داده).
Special Function Units (SFU)	برای محاسبات خاص (مثل سینوس، ریشه دوم، لگاریتم).

SM

- The fundamental execution unit of NVIDIA GPUs.
- Includes:
 - CUDA Cores(FP32/INT32 ALUs): handle general math & logic.
 - Tensor Cores: accelerate AI/ML workloads
 - RT Core: for ray/triangle intersection.
 - Warp Schedulers & Dispatch Units: decide which warps (groups of 32 threads) run each cycle.
 - Texture Units: perform texture fetches and filtering.
 - Load/Store Units (LD/ST): handle memory access (global/local memory reads/writes).
 - Register File: private storage per thread.
 - Shared Memory / L1 Cache: memory accessible by all threads in the SM.

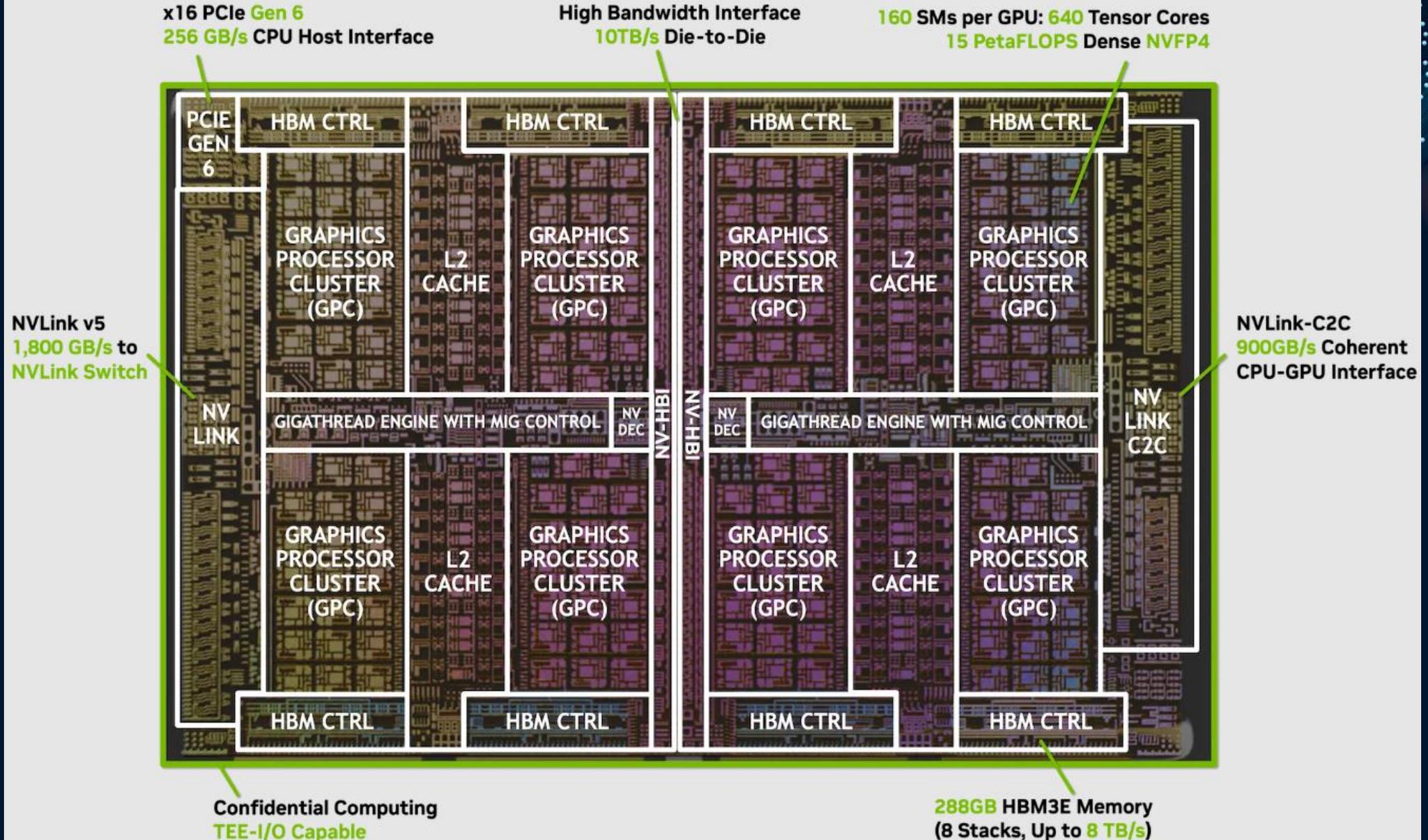
SM Cont.

- Example of SMIT(Single Instruction Multiple Threads)

```
for (i=0; i<32; i++)  
    C[i] = A[i] + B[i];
```

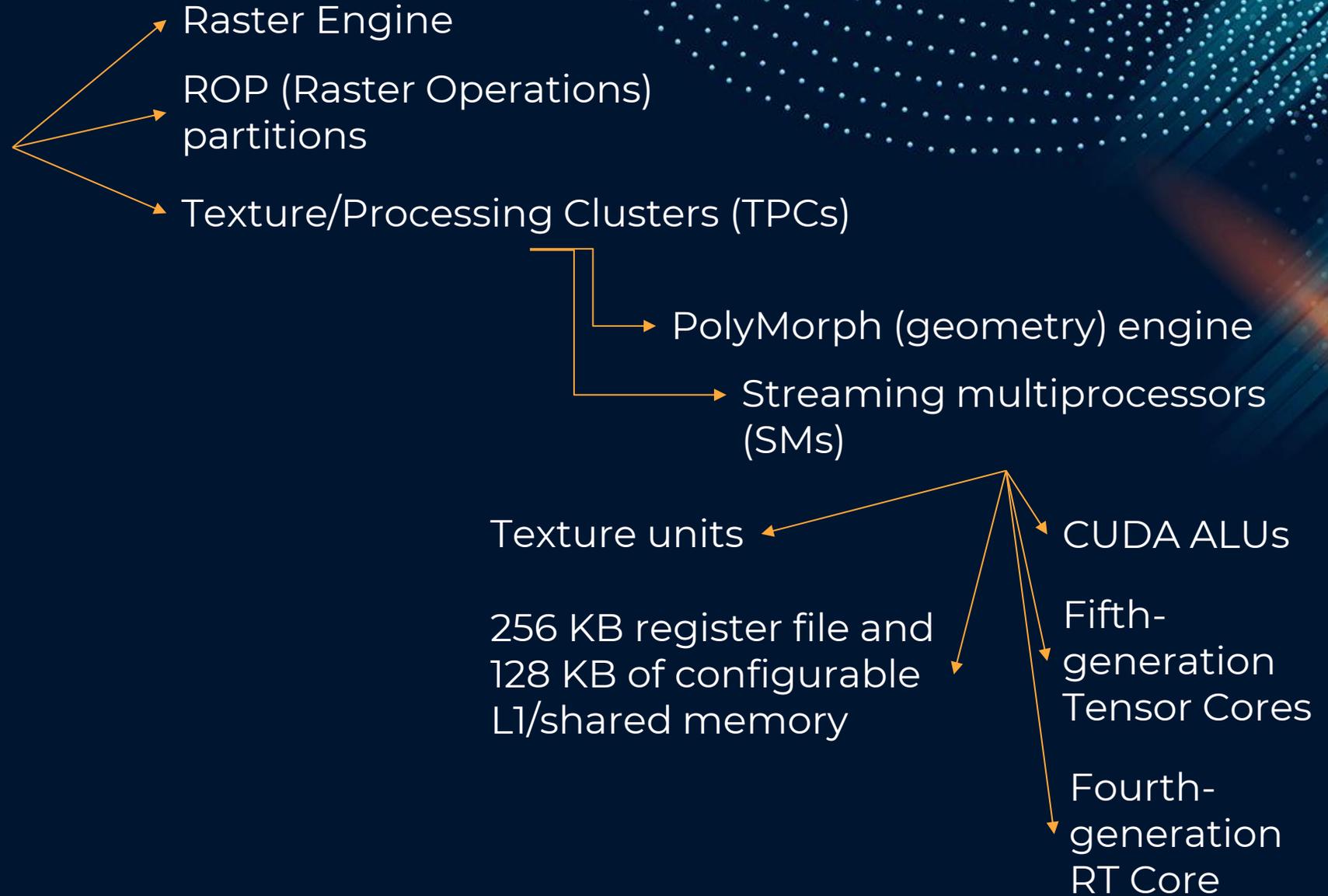
- One ADD instruction is issued.
- All 32 threads perform the addition simultaneously, each with different A[i], B[i].
- 64 warps per each SM(4 can be active in each cylce).
- 4 warp schedulers in each SM.

NVIDIA Blackwell Ultra GPU



بلاک های معماری بلک ول

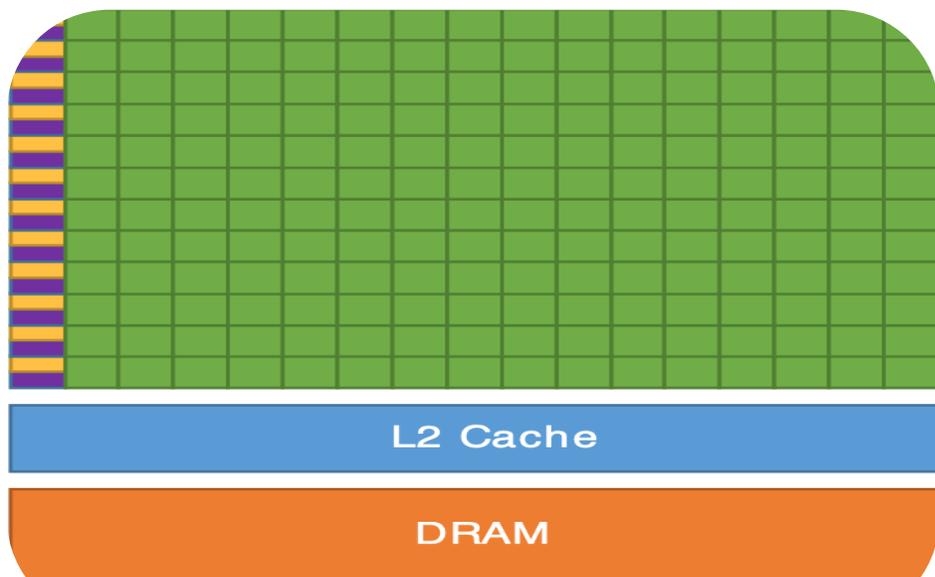
GPC
(Graphics Processing Clusters)



انواع حافظه در

GPU

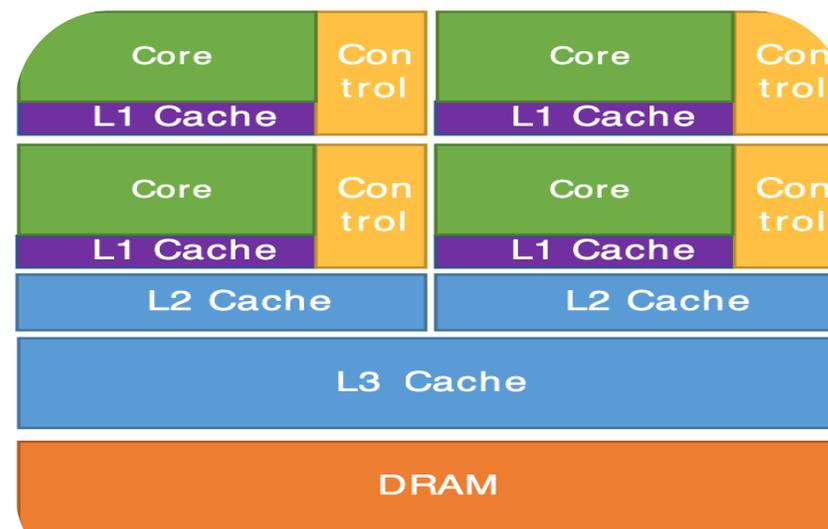
معماری پردازنده‌ها (GPU و CPU)



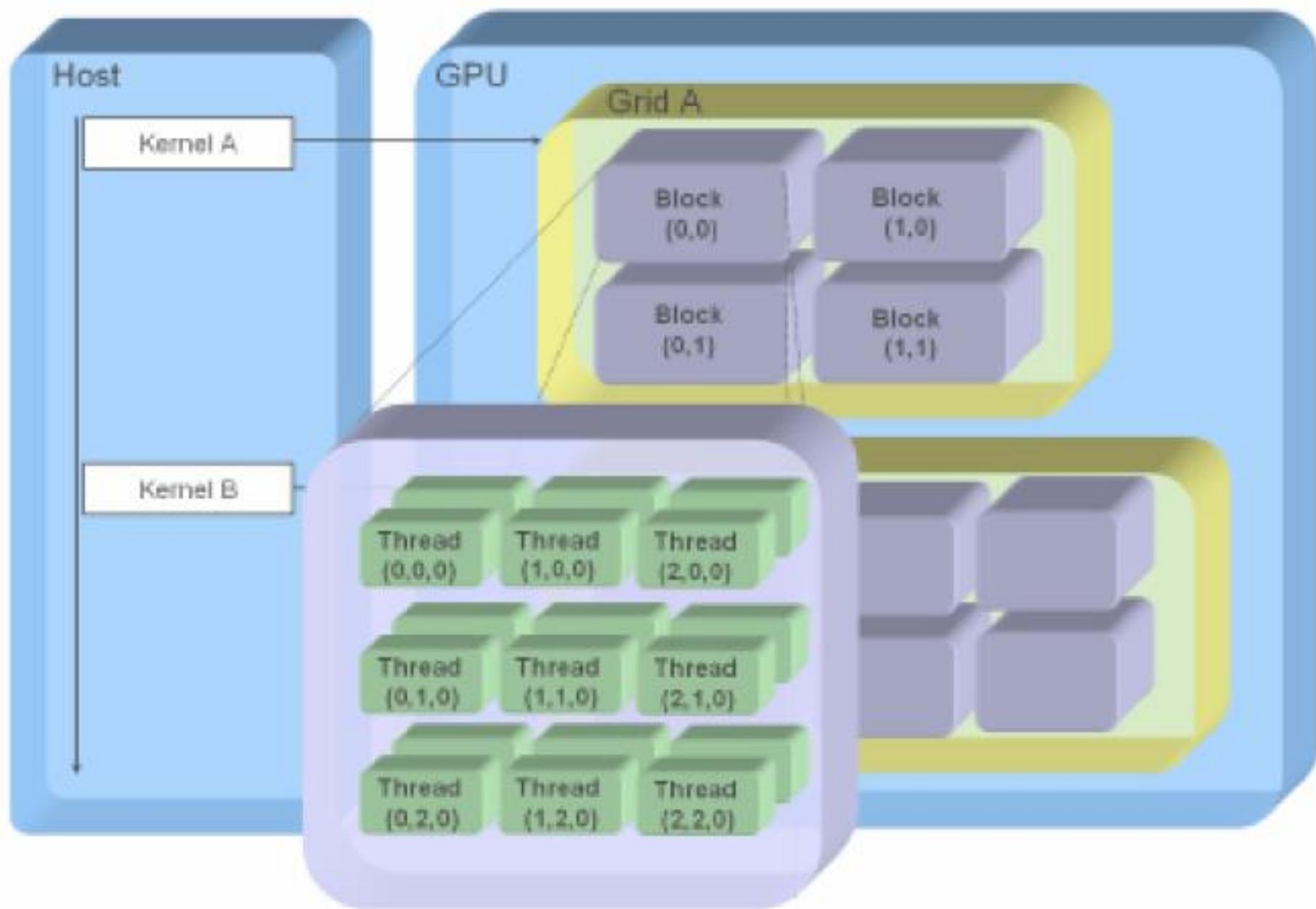
GPU

- دارای هزاران هسته
- طراحی شده برای پردازش موازی
- مناسب برای workloadهایی با داده‌های مستقل (مثل خوانش‌های NGS)
- بهره‌برداری از حافظه shared برای کاهش latency

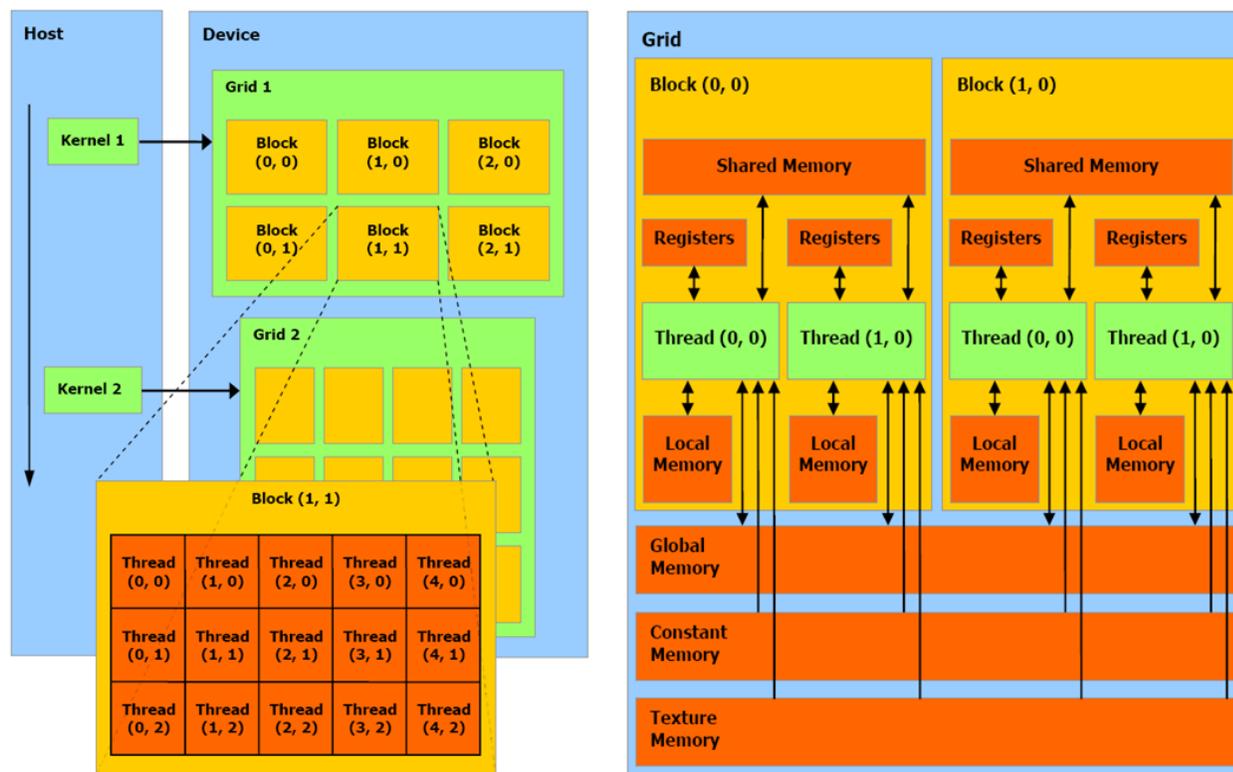
- مناسب برای پردازش‌های ترتیبی
- دارای تعداد محدودی هسته (4-64 هسته)
- کش‌های بزرگ و بهینه‌سازی شده برای latency پایین



CPU



اجزای اصلی GPU



- Grid ۱
- Block ۲
- Thread ۳
- Memory ۴

۱. Global Memory (حافظه‌ی سراسری)

- **تعریف:**
- Global Memory در واقع **حافظه‌ی اصلی GPU** است (اغلب از نوع HBM3 یا HBM3e در H200).
- همه‌ی **Threads** (رشته‌ها) در همه‌ی **Blocks** به آن دسترسی دارند.
- **ویژگی‌ها:**
- **ظرفیت زیاد:** چند ده گیگابایت (مثلاً ۱۴۱ GB در H200).
- **سرعت نسبتاً پایین‌تر:** دسترسی به آن صدها چرخه‌ی کلاک طول می‌کشد.
- **دسترسی جهانی:** همه‌ی threadها از همه‌ی blocks می‌توانند به آن بنویسند یا بخوانند.
- **پایدار:** داده‌ها تا پایان اجرای کرنل در دسترس می‌مانند.
- **معایب:**
- تاخیر بالا (High Latency)
- پهنای باند محدود نسبت به حافظه‌های نزدیک‌تر (مثل L1 cache یا shared memory)

۲. Shared Memory (حافظه‌ی اشتراکی)

- تعریف:
- Shared Memory نوعی حافظه‌ی کوچک‌تر اما بسیار سریع‌تر است که فقط میان **thread**های یک **block** مشخص مشترک است.
در واقع می‌توان آن را حافظه‌ی سطح **L1** برای **GPU threads** دانست.
- ♦ ویژگی‌ها:
- ظرفیت محدود: معمولاً ۱۰۰ تا ۲۵۶ کیلوبایت در هر SM (Streaming Multiprocessor).
- سرعت بسیار بالا: تأخیر چند چرخه‌ی کلاک (خیلی سریع‌تر از global memory).
- اشتراک در سطح **block**: فقط **thread**های همان **block** می‌توانند به آن دسترسی داشته باشند.
- برنامه‌پذیر: توسط برنامه‌نویس **CUDA** به صورت صریح تعریف و استفاده می‌شود.
- ♦ مثال در **CUDA**:

Copy code

cpp

```
__global__ void example(int *input, int *output) {
    __shared__ int cache[256]; // block تعریف حافظه اشتراکی در هر

    int tid = threadIdx.x;
    cache[tid] = input[tid]; // داده را از حافظه سراسری می‌خوانیم
    __syncthreads(); // threadها همگامسازی

    output[tid] = cache[tid] * 2; // Local محاسبه با داده‌ی
}
```

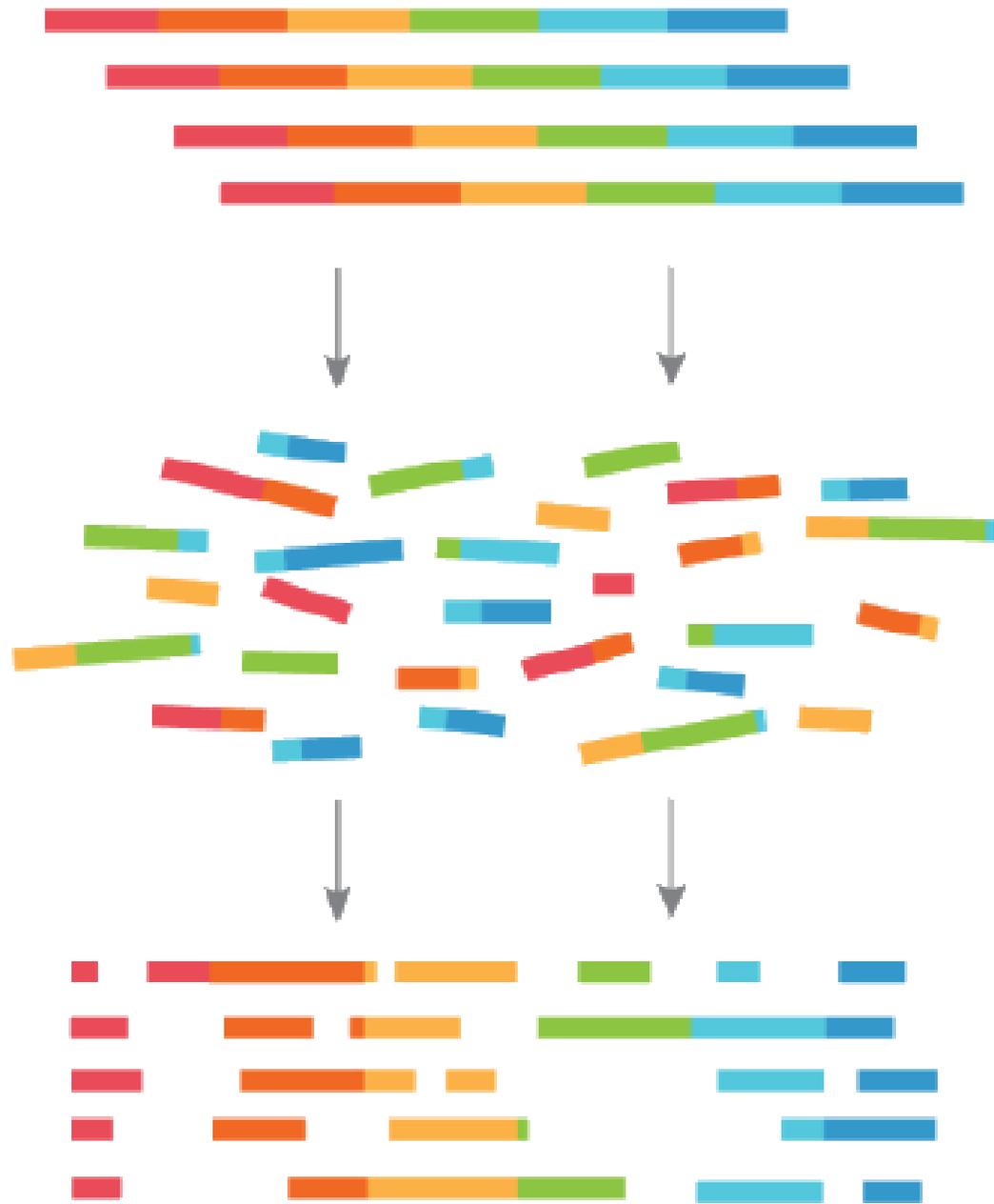
⚡ تفاوت‌های کلیدی

ویژگی	Global Memory	Shared Memory
محدوده‌ی دسترسی	همه‌ی threadها در همه‌ی blocks	فقط threadهای همان block
ظرفیت	بسیار زیاد (GBها)	بسیار کم (KBها)
سرعت	کند (صدها cycle)	سریع (چند cycle)
کنترل توسط برنامه‌نویس	غیرمستقیم	مستقیم (با <code>__shared__</code>)
محل فیزیکی	در خارج از SM (در HBM)	درون SM
دوام داده	در طول اجرای کرنل	تا پایان block
کاربرد	ذخیره‌ی داده‌های اصلی یا ورودی‌های بزرگ	بافر موقت برای محاسبات محلی و تکراری

نتیجه نهایی

- محاسبات خیلی زمانبر با تعداد تکرار زیاد اگر در **Shared memory** بجای **Global memory** قابل اجراست، بهتر است با تغییر کد در آنجا انجام شود تا سرعت به طرز قابل توجهی افزایش پیدا کند.
- حافظه **Shared** به نوعی نقش **Cache memory** را در **GPU** ایفا می کند و بین **Thread** های یک بلوک مشترک است.
- پایپ لاین ها هم باید به نحوی طراحی شوند که هر گروه از داده فقط داخل یک **Shared memory** هدایت شود.

مشالی از کار سرد کلان داده
در سردانش داده رشیک



Template DNA

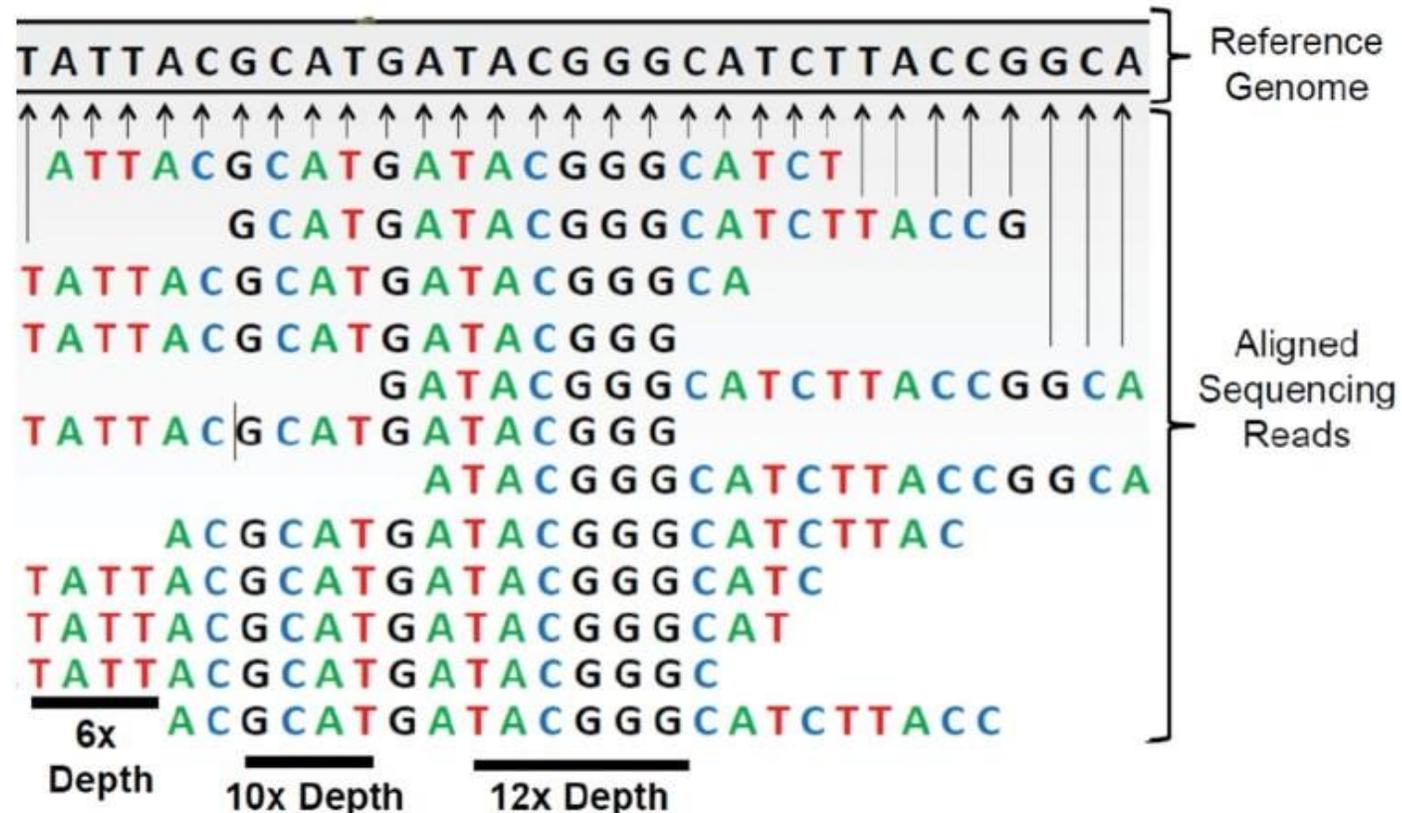
Short reads from
Sequencer
(FASTQ)

Mapped reads
(SAM/BAM-file)

ATGTTCCGATTAGGAAACCATCTGTAACGTGTTTCATTCAGTAAAAGGGAGGAAA

Sequence Alignment

- All reads are aligned to a reference genome like GRCh38 and hg19 in this step.
- Two alignment files **SAM** ([Sequence Alignment Map](#)) and **BAM** ([Binary Alignment Map](#)) are produced.



SV Types

 Destructive (non-balanced)

 Non-destructive (balanced)

