

Retour GTC 2022 : GPU Technical Conference

Pierre Aubert

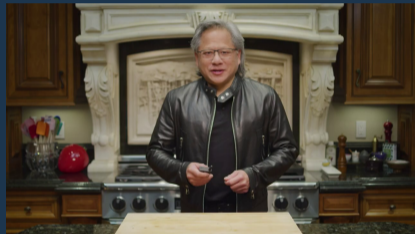


GTC 2022 <https://gtc22.event.nvidia.com/>

This is not possible to summarize **987** talks !

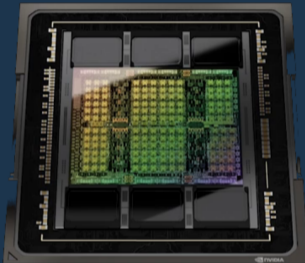
Only focus on :

- ▶ **Hardwares** evolution
- ▶ **Compilers** evolution
- ▶ **Programming Languages** evolution
- ▶ **Job Communication**
- ▶ **Data Compression/Decompression**
- ▶ **Machine Learning / Deep Learning** evolution and use
- ▶ How all (will) work together (with a bit of **Forum ORAP**)

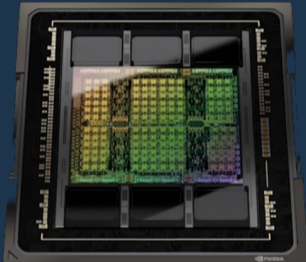


Keynotes video link

H100 (Hopper)



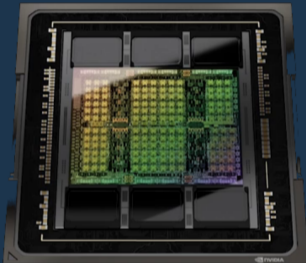
H100 (Hopper)



H100 SXM5 and PCIe

4 nm (A100 7 nm)

H100 (Hopper)



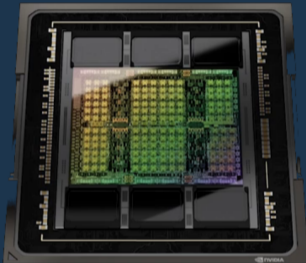
H100 SXM5 and PCIe

Hardware evolution

4 nm (A100 7 nm)

Compute Capabilities 9.0

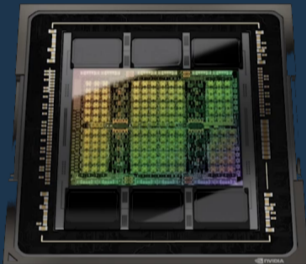
H100 (Hopper)



H100 SXM5 and PCIe

Hardware evolution

H100 (Hopper)



H100 SXM5 and PCIe

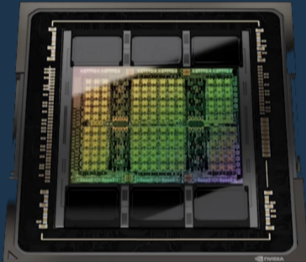
4 nm (A100 7 nm)

Compute Capabilities 9.0

132 SMs (tensor core Gen4) 8448 Cores

Hardware evolution

H100 (Hopper)



H100 SXM5 and PCIe

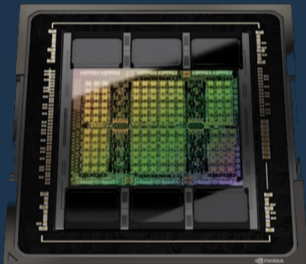
4 nm (A100 7 nm)

Compute Capabilities 9.0

132 SMs (tensor core Gen4) 8448 Cores

80 GB DRAM

H100 (Hopper)



H100 SXM5 and PCIe

4 nm (A100 7 nm)

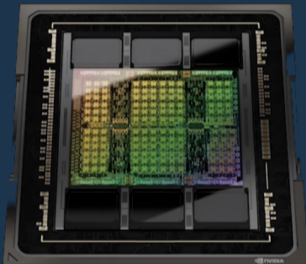
Compute Capabilities 9.0

132 SMs (tensor core Gen4) 8448 Cores

80 GB DRAM

3 TB/s Bandwidth (HBM3)

H100 (Hopper)



H100 SXM5 and PCIe

4 nm (A100 7 nm)

Compute Capabilities 9.0

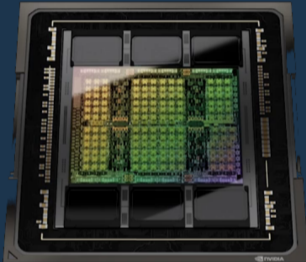
132 SMs (tensor core Gen4) 8448 Cores

80 GB DRAM

3 TB/s Bandwidth (HBM3)

NVLink Gen 4 (900 GB/s)

H100 (Hopper)



H100 SXM5 and PCIe

4 nm (A100 7 nm)

Compute Capabilities 9.0

132 SMs (tensor core Gen4) 8448 Cores

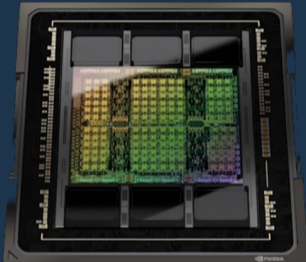
80 GB DRAM

3 TB/s Bandwidth (HBM3)

NVLink Gen 4 (900 GB/s)

MIG Gen 2 with confidential computing
(all instances have Image/Video decoding)

H100 (Hopper)



H100 SXM5 and PCIe

4 nm (A100 7 nm)

Compute Capabilities 9.0

132 SMs (tensor core Gen4) 8448 Cores

80 GB DRAM

3 TB/s Bandwidth (HBM3)

NVLink Gen 4 (900 GB/s)

MIG Gen 2 with confidential computing
(all instances have Image/Video decoding)

A lot a computing precisions

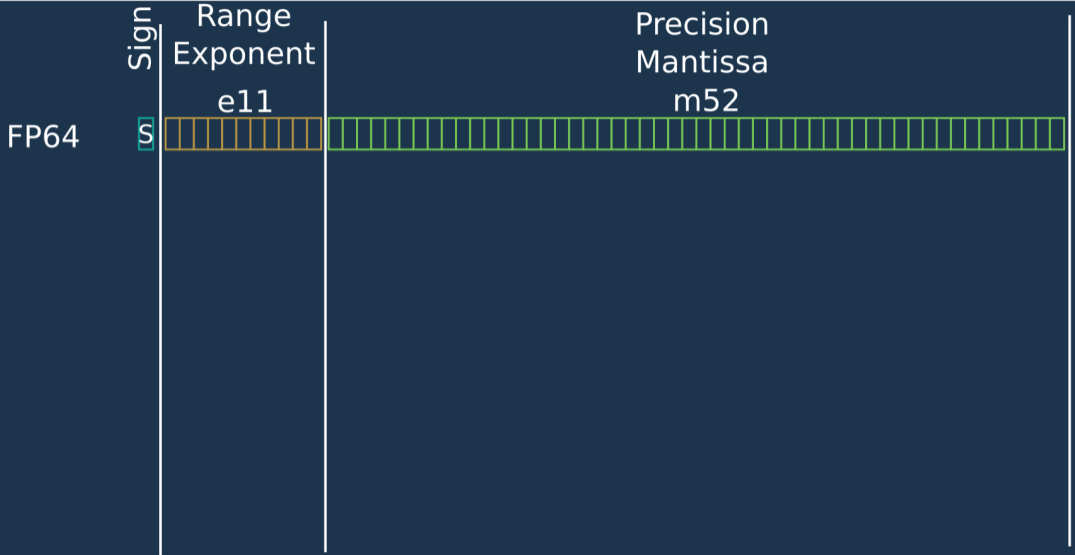
Computing Precision

Sign

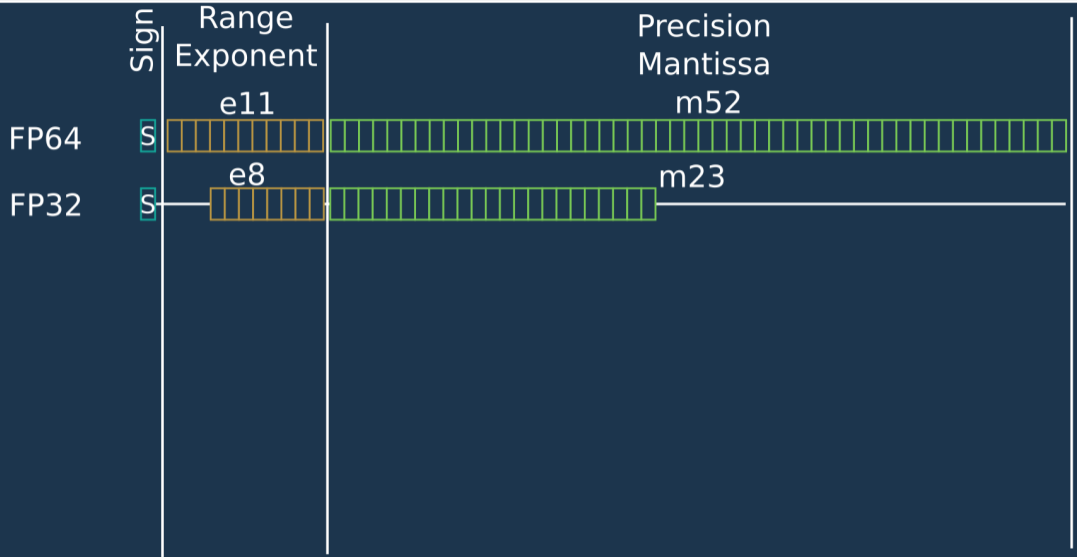
Range
Exponent

Precision
Mantissa

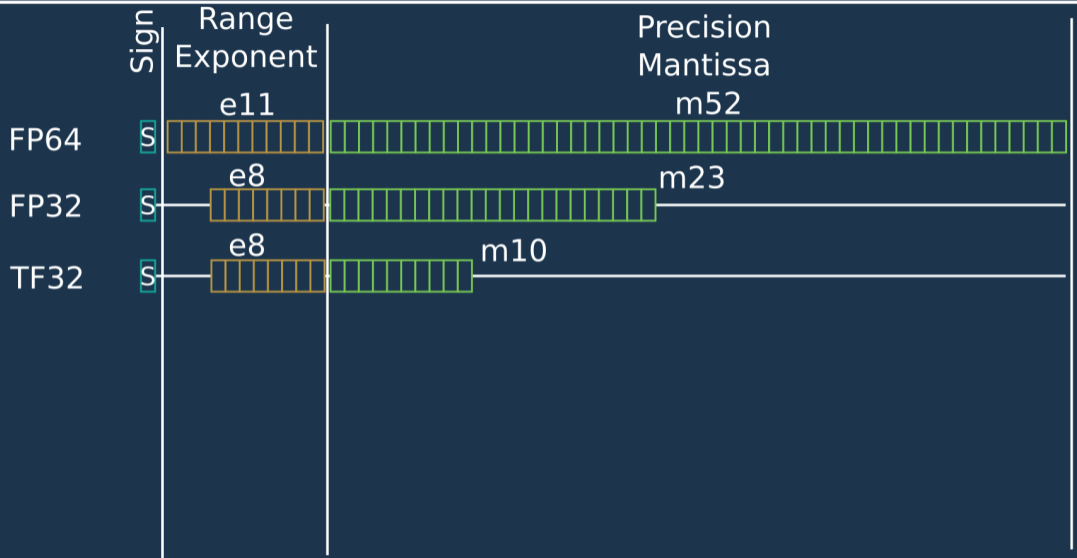
Computing Precision



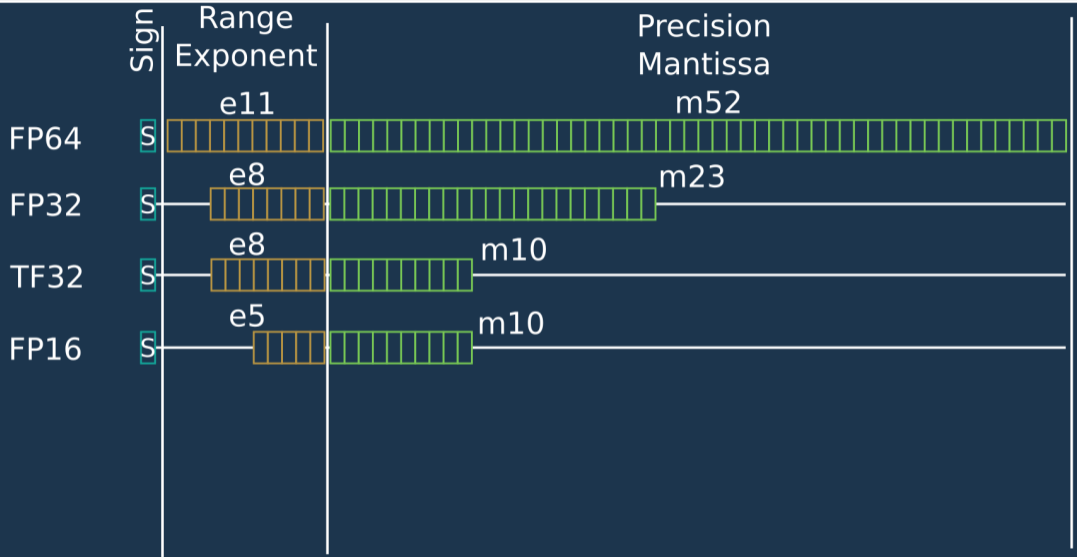
Computing Precision



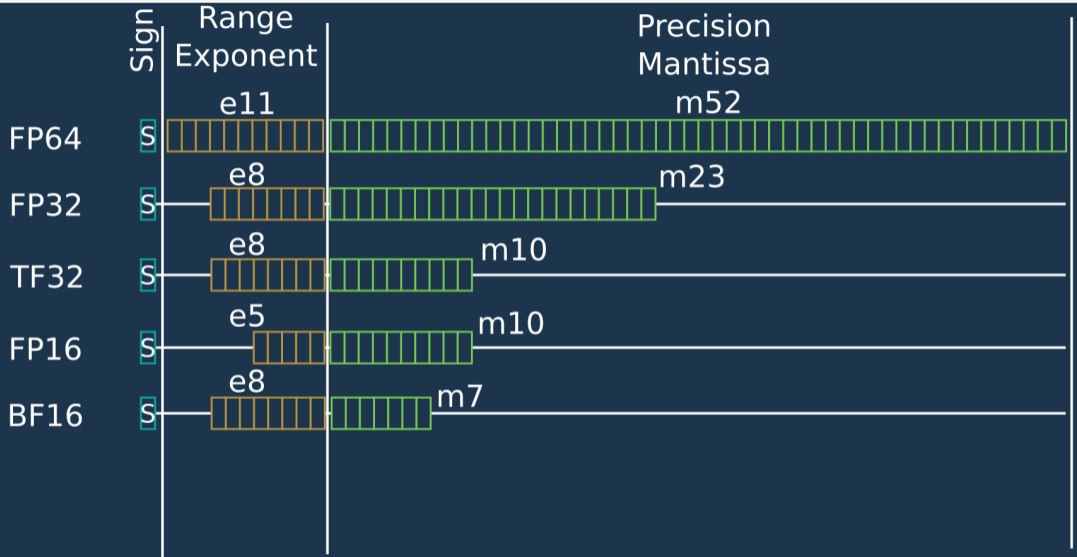
Computing Precision



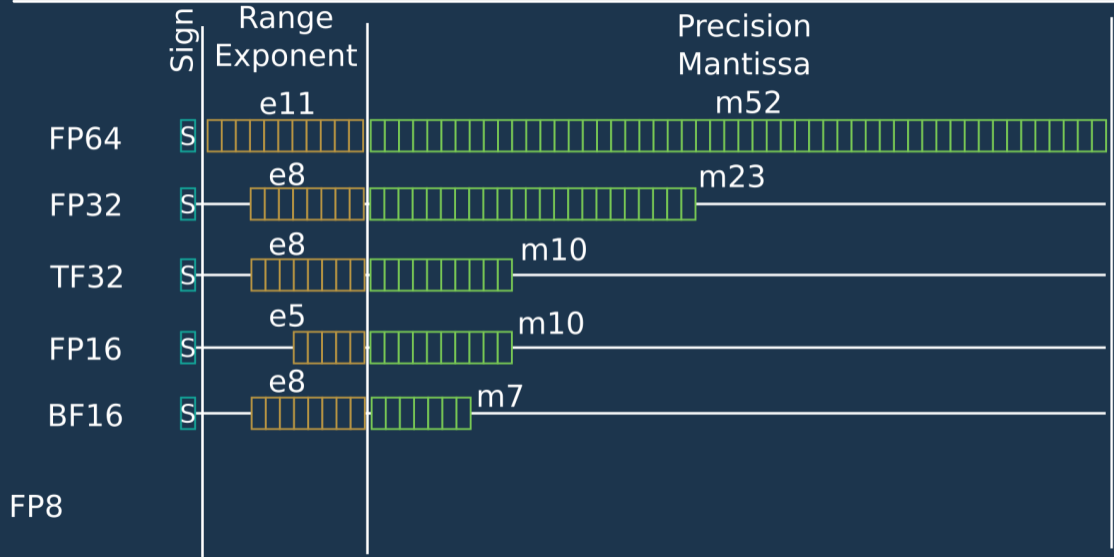
Computing Precision



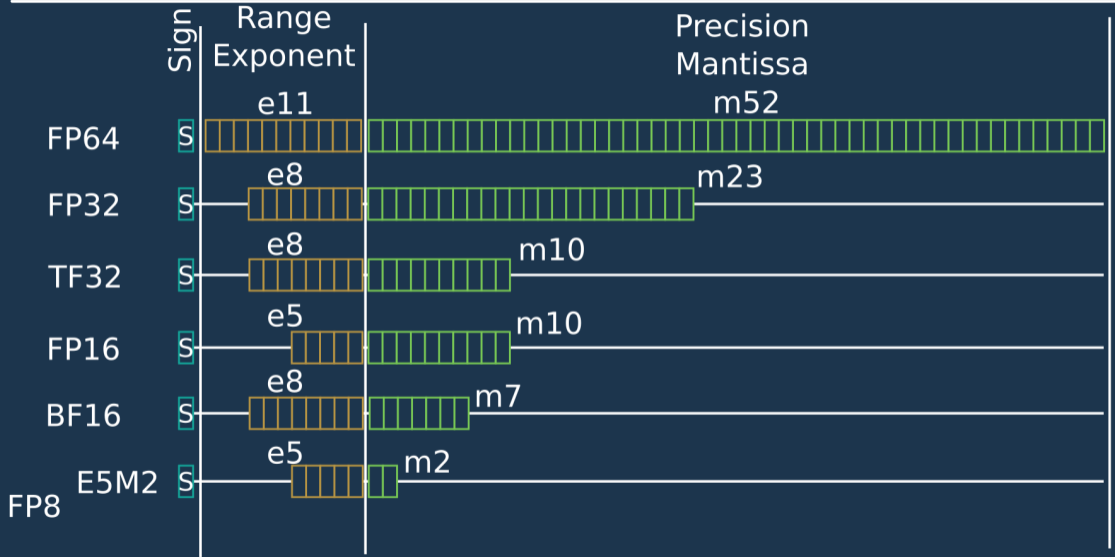
Computing Precision



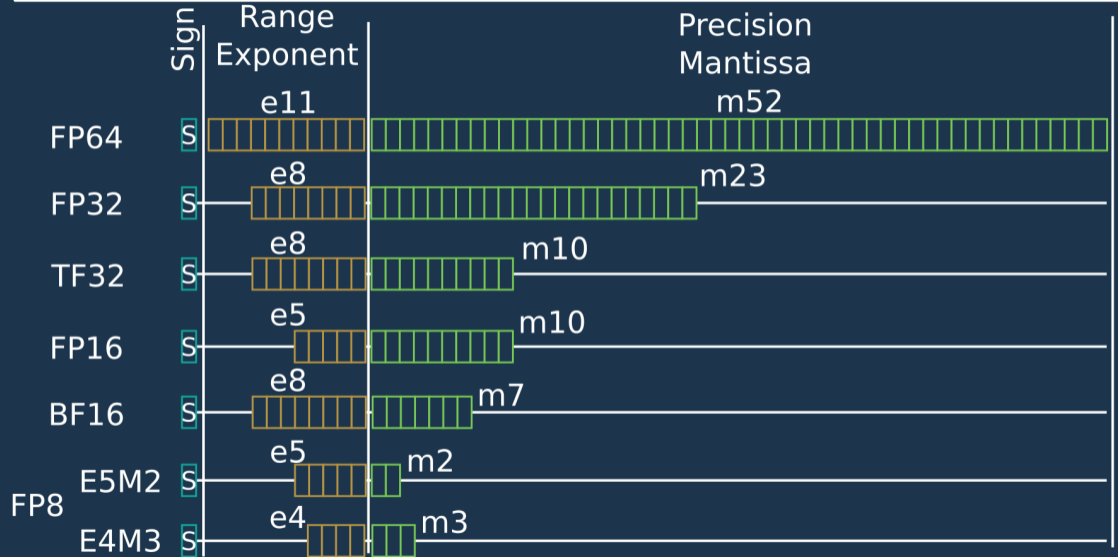
Computing Precision



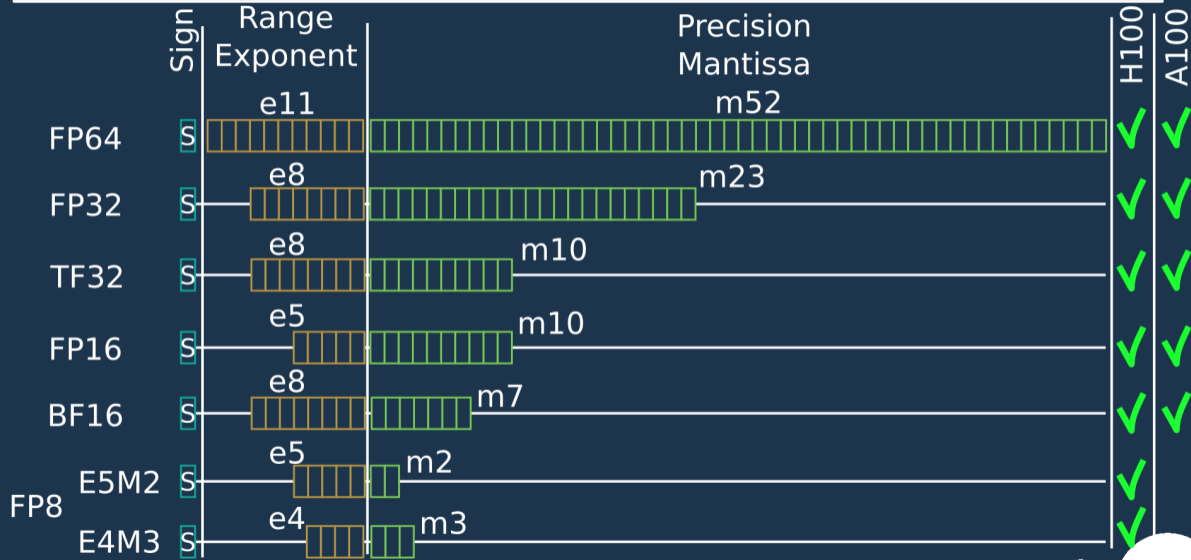
Computing Precision



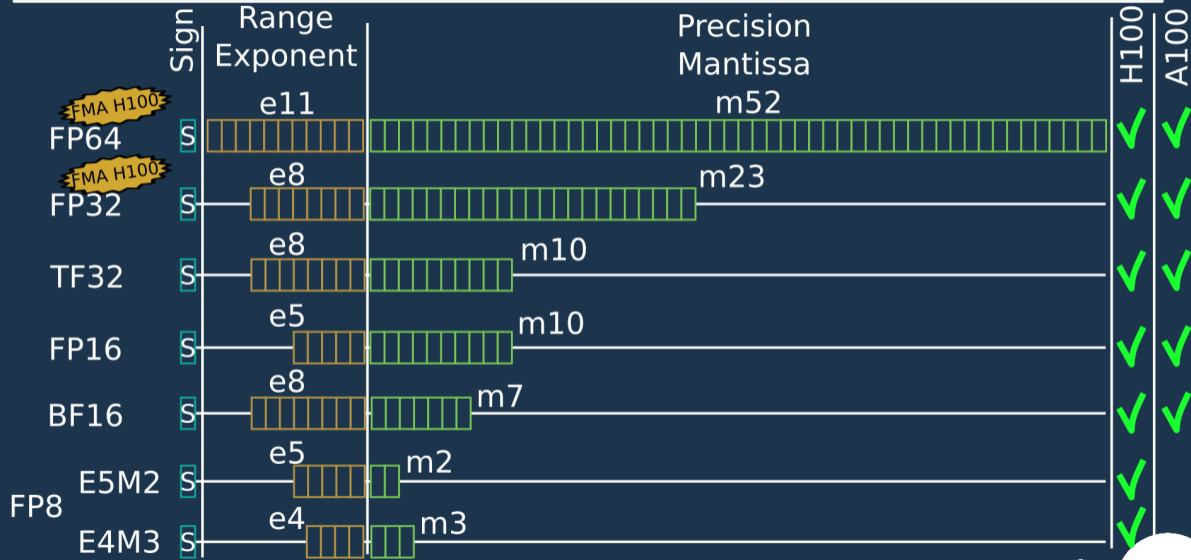
Computing Precision



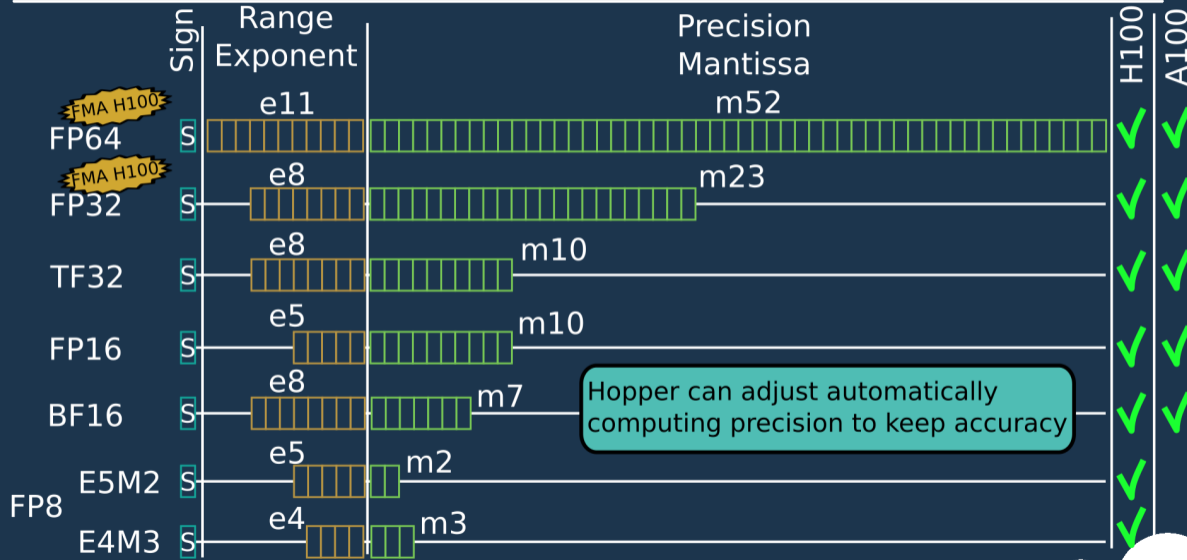
Computing Precision



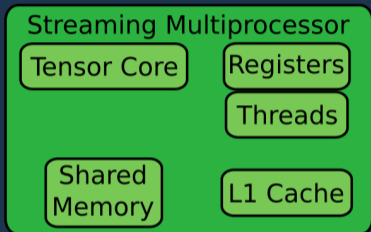
Computing Precision



Computing Precision



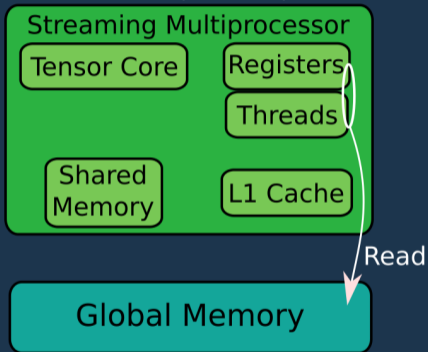
Tensor Memory Acceleration (TMA)



Thread acceleration

Tensor Memory
Acceleration (TMA)

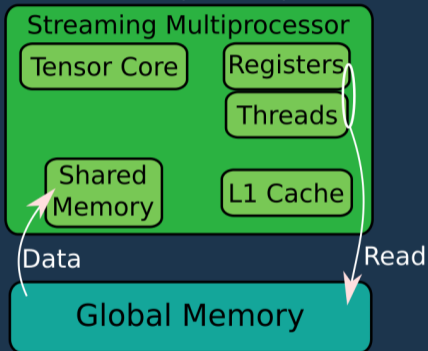
A100 (no TMA)



Thread acceleration

Tensor Memory
Acceleration (TMA)

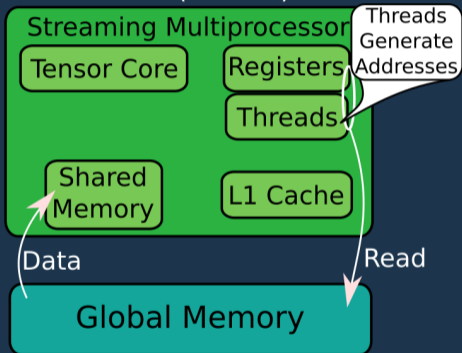
A100 (no TMA)



Thread acceleration

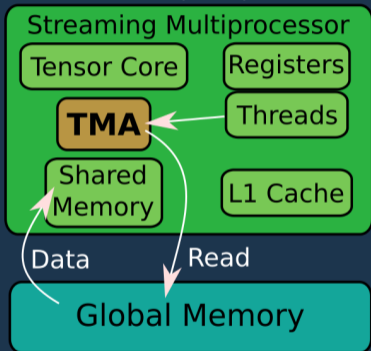
Tensor Memory
Acceleration (TMA)

A100 (no TMA)



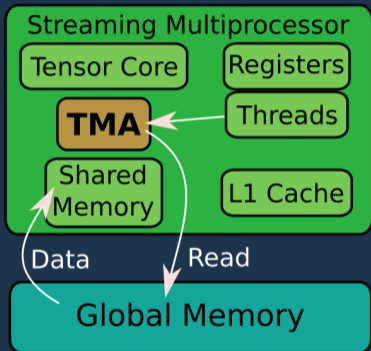
Tensor Memory
Acceleration (TMA)

H100 (TMA)



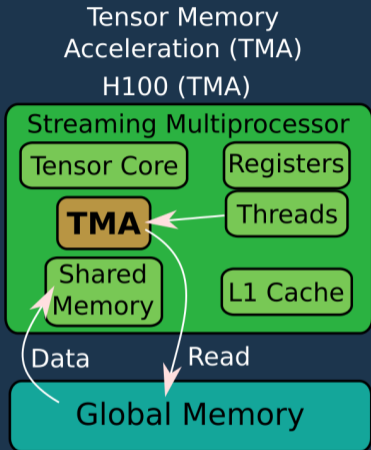
Thread acceleration

Tensor Memory
Acceleration (TMA)
H100 (TMA)

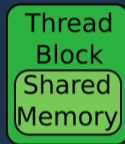
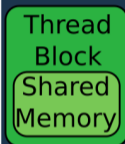


Distributed
shared memory

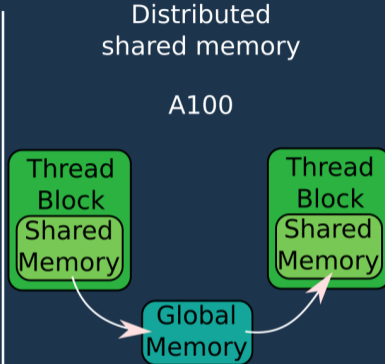
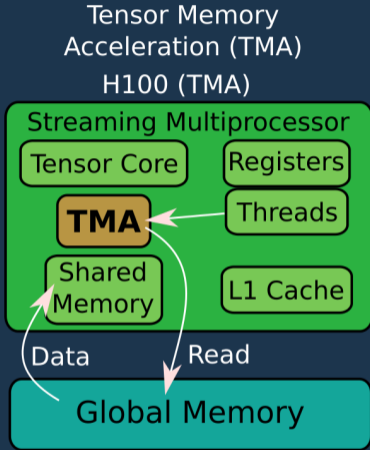
Thread acceleration



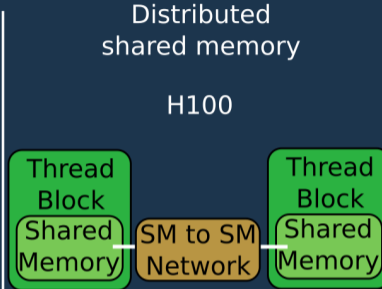
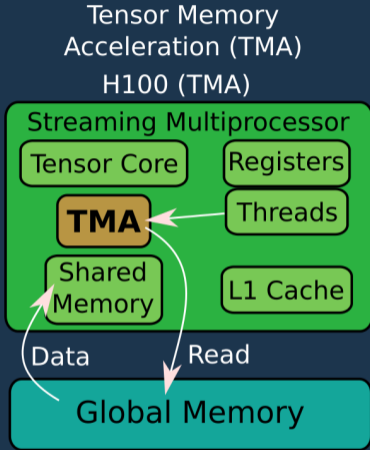
Distributed shared memory



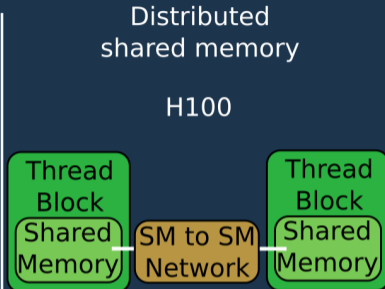
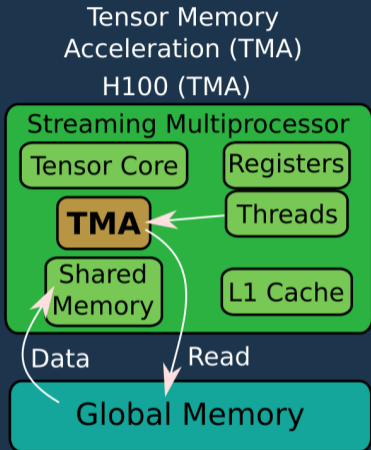
Thread acceleration



Thread acceleration

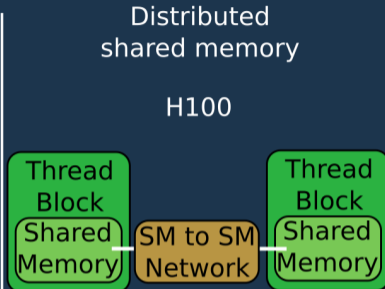
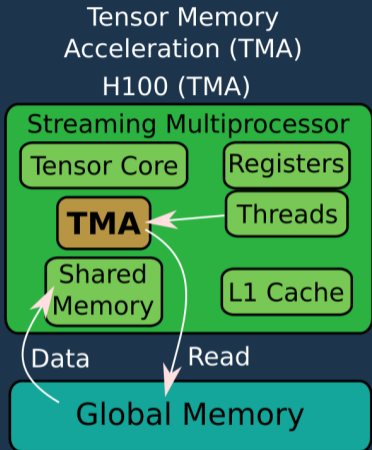


Thread acceleration



Asynchronous
Transaction Barrier

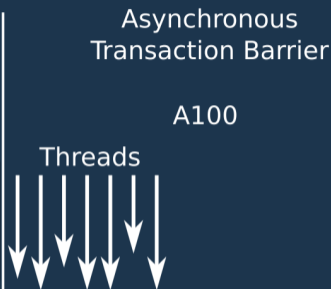
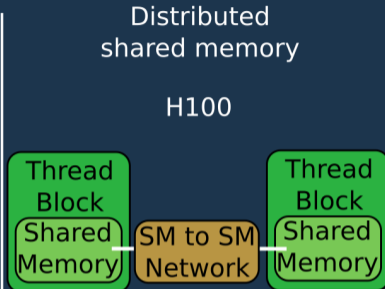
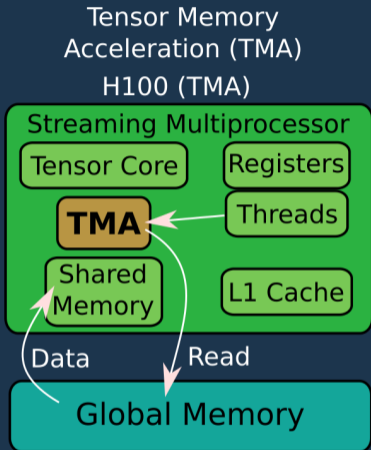
Thread acceleration



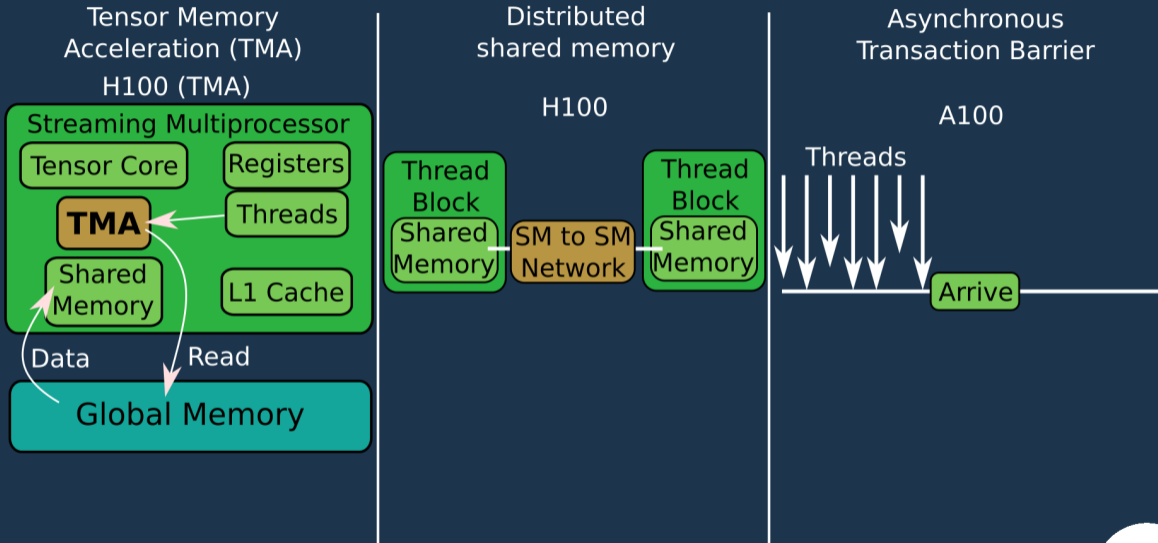
Asynchronous Transaction Barrier

A100

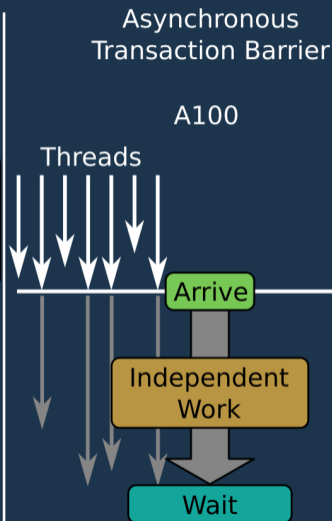
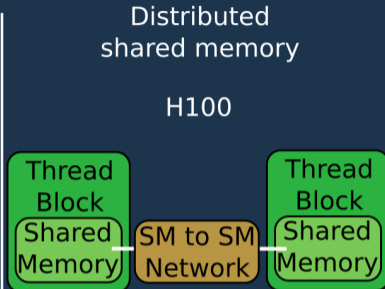
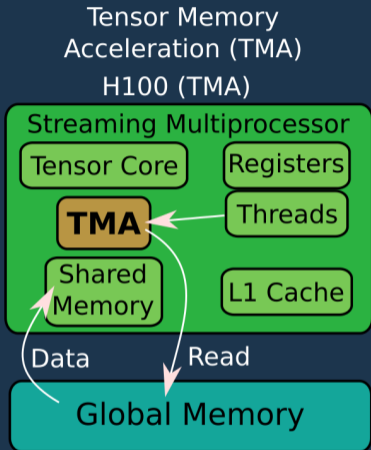
Thread acceleration



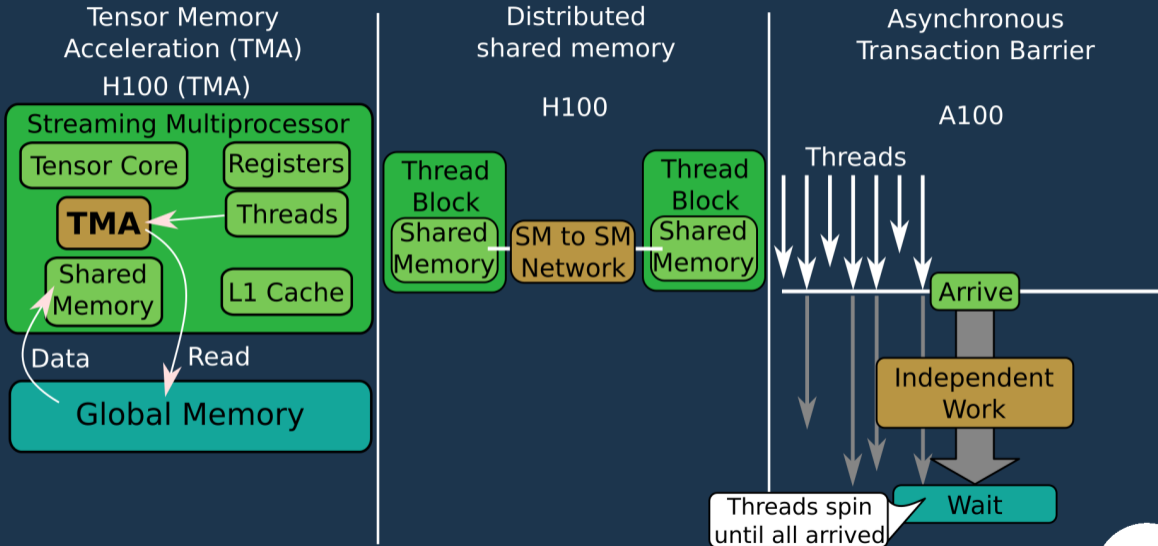
Thread acceleration



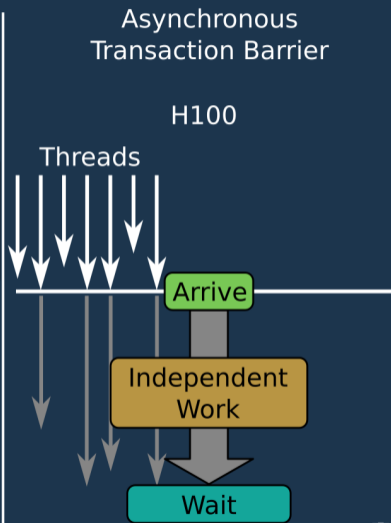
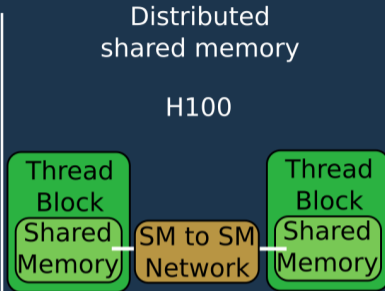
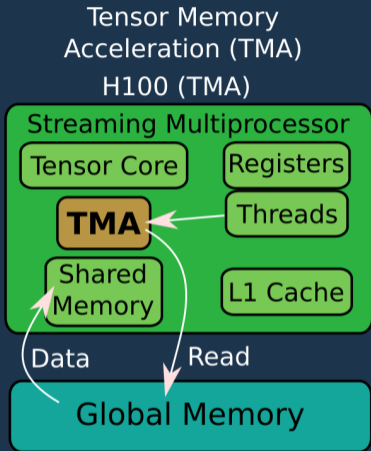
Thread acceleration



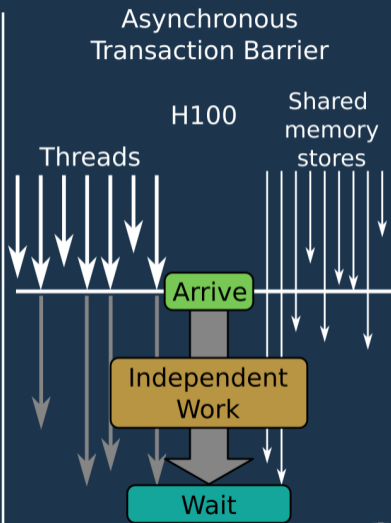
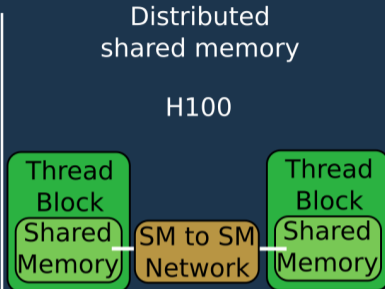
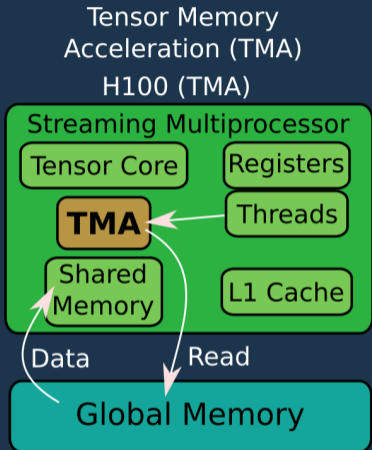
Thread acceleration



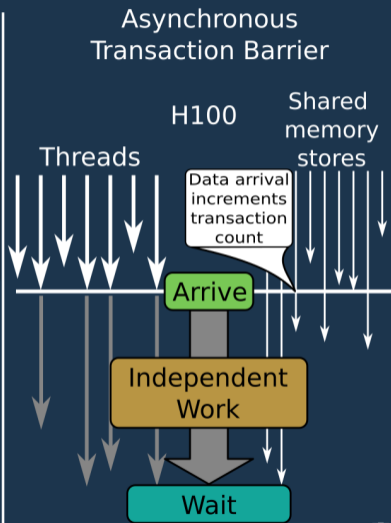
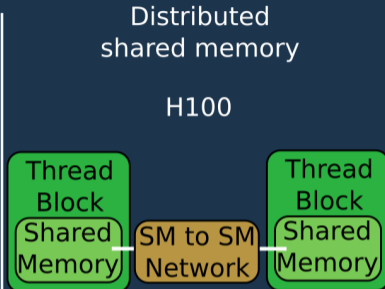
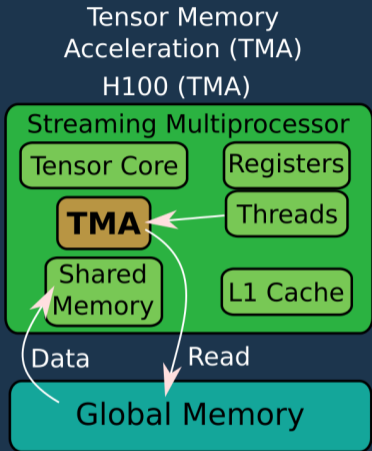
Thread acceleration



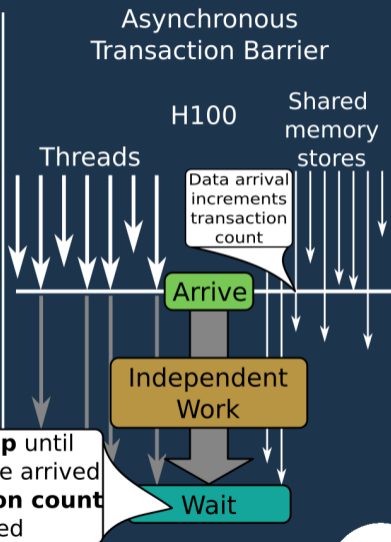
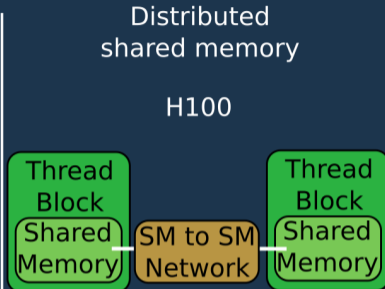
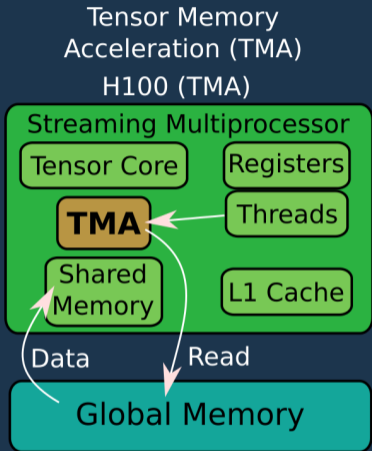
Thread acceleration



Thread acceleration



Thread acceleration



New intrinsic functions

New intrinsic functions

Dynamic Programming Acceleration (DPX) : recursive algorithms

Dynamic Programming Acceleration (DPX) : recursive algorithms

Floyd-Warshall Algorithm

Route optimization



Dynamic Programming Acceleration (DPX) : recursive algorithms

Floyd-Warshall Algorithm

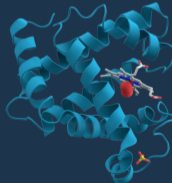
Route optimization



Smith-Waterman and
Needleman-Wunsch
Algorithm

Pattern Matching

- DNA sequence alignment
- Protein classification
- Protein folding



New intrinsic functions

Dynamic Programming Acceleration (DPX) : recursive algorithms

Floyd-Warshall Algorithm

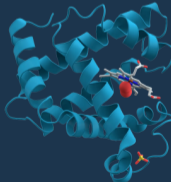
Route optimization



Smith-Waterman and
Needleman-Wunsch
Algorithm

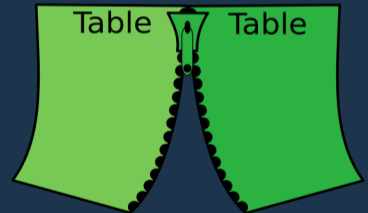
Pattern Matching

- DNA sequence alignment
- Protein classification
- Protein folding

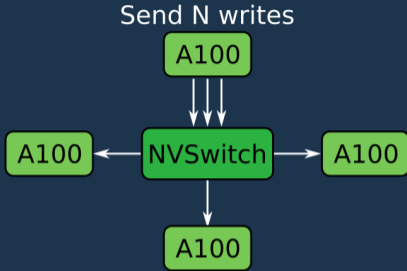


SQL Query optimisation

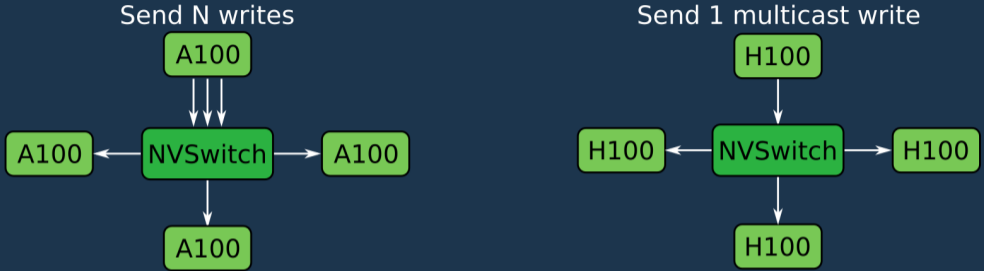
Join in optimal order



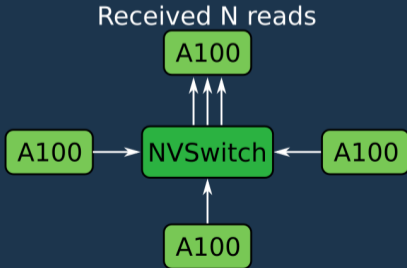
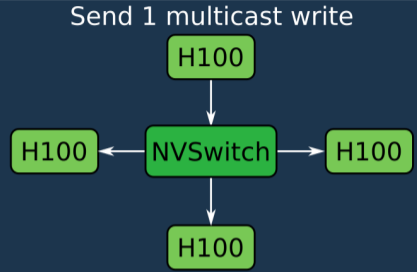
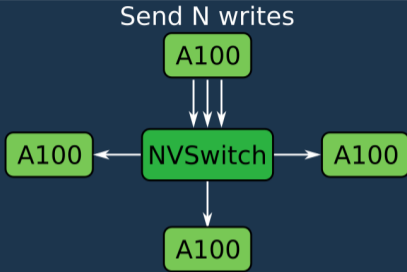
NVLink/NVSwitch acceleration



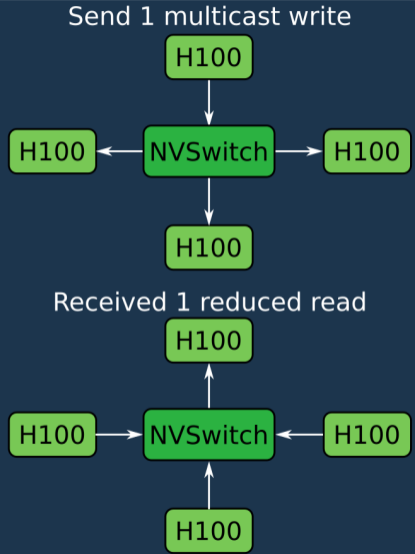
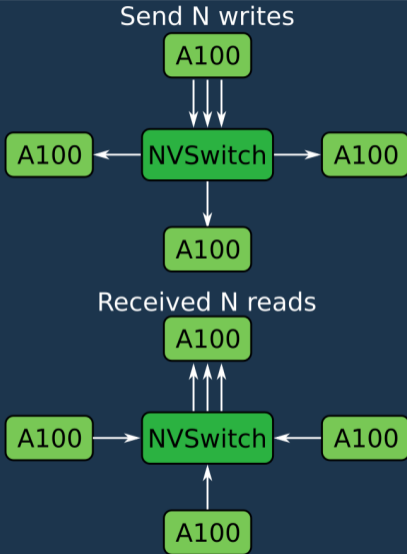
NVLink/NVSwitch acceleration



NVLink/NVSwitch acceleration



NVLink/NVSwitch acceleration



Thread block cluster

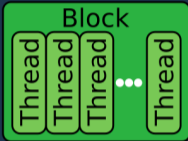
A100 (108 SMs)

A100 (108 SMs)

Thread

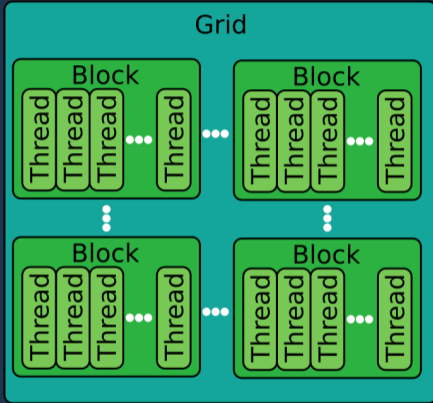
Thread block cluster

A100 (108 SMs)



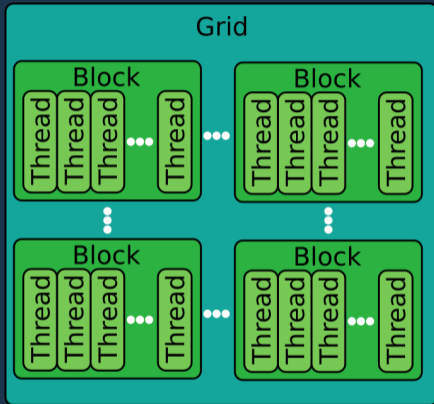
Thread block cluster

A100 (108 SMs)



Thread block cluster

A100 (108 SMs)

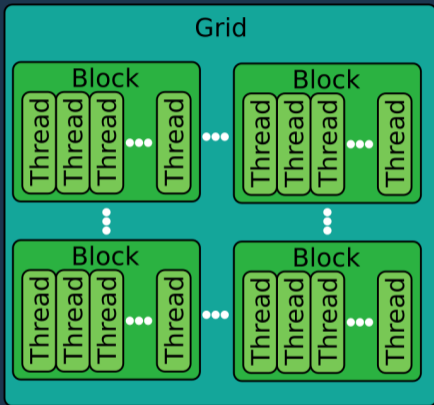


1 Block ~ 1 SM < 1% A100

Thread block cluster

A100 (108 SMs)

H100 (132 SMs)

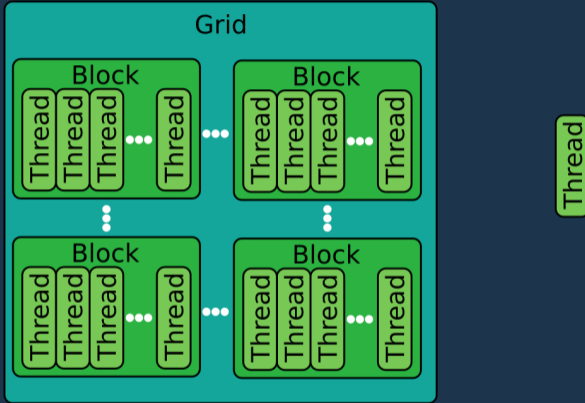


1 Block ~ 1 SM < 1% A100

Thread block cluster

A100 (108 SMs)

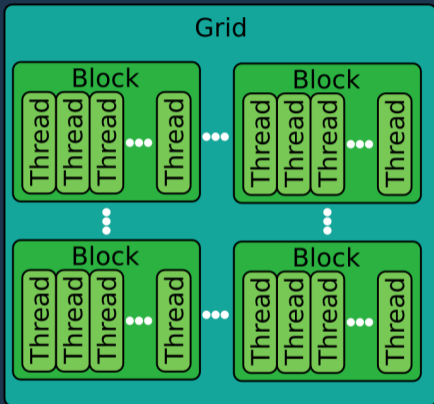
H100 (132 SMs)



1 Block ~ 1 SM < 1% A100

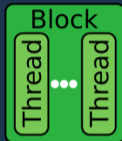
Thread block cluster

A100 (108 SMs)



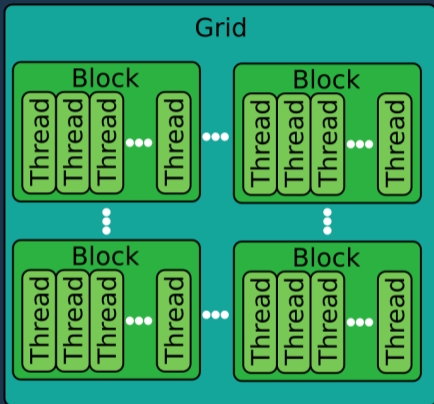
1 Block ~ 1 SM < 1% A100

H100 (132 SMs)



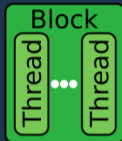
Thread block cluster

A100 (108 SMs)



1 Block ~ 1 SM < 1% A100

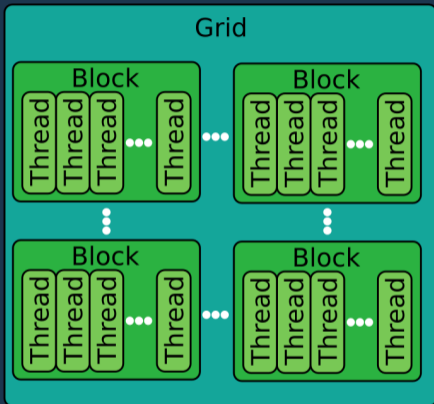
H100 (132 SMs)



1 Block ~ 1 SM < 0.8 % H100

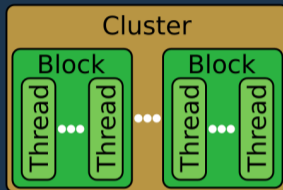
Thread block cluster

A100 (108 SMs)



1 Block ~ 1 SM < 1% A100

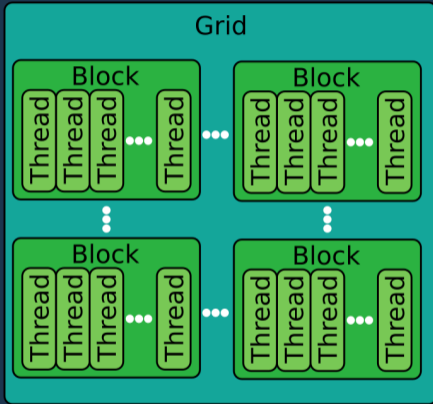
H100 (132 SMs)



1 Block ~ 1 SM < 0.8 % H100

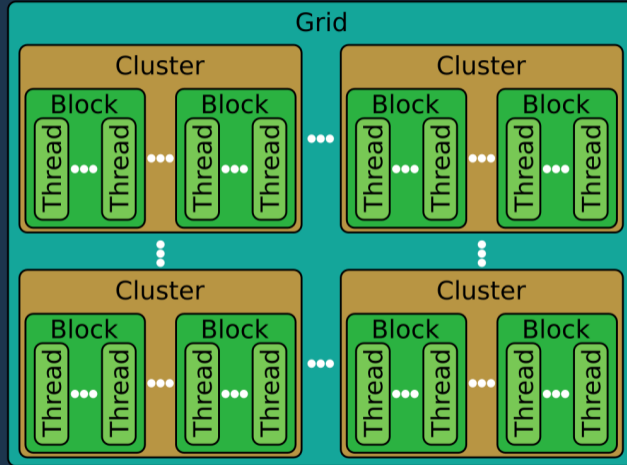
Thread block cluster

A100 (108 SMs)



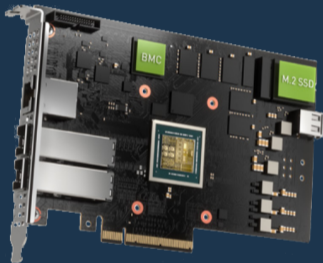
1 Block ~ 1 SM < 1% A100

H100 (132 SMs)



1 Block ~ 1 SM < 0.8 % H100

DPU : Data Processing Unit

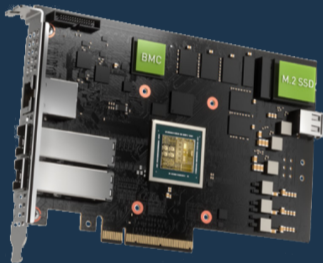


Bluefield 3 DPU – 2x 200Gb/s
FHHL form factor

Bluefield 3

DPU : Data Processing Unit

1, 2, 4 ports with up to 400Gb/s connectivity
(Ethernet or NDR InfiniBand)

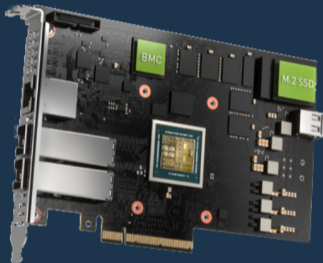


Bluefield 3 DPU – 2x 200Gb/s
FHHL form factor

Bluefield 3

DPU : Data Processing Unit

1, 2, 4 ports with up to 400Gb/s connectivity
(Ethernet or NDR InfiniBand)

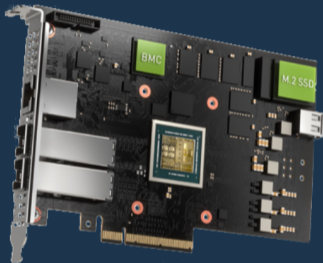


16 GB on-board DDR5 memory

Bluefield 3 DPU – 2x 200Gb/s
FHHL form factor

Bluefield 3

DPU : Data Processing Unit



Bluefield 3 DPU – 2x 200Gb/s
FHHL form factor

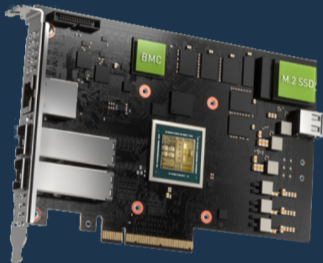
1, 2, 4 ports with up to 400Gb/s connectivity
(Ethernet or NDR InfiniBand)

16 GB on-board DDR5 memory

1 GbE out-of-band management port

Bluefield 3

DPU : Data Processing Unit



Bluefield 3 DPU – 2x 200Gb/s
FHHL form factor

1, 2, 4 ports with up to 400Gb/s connectivity
(Ethernet or NDR InfiniBand)

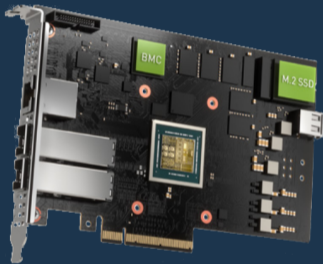
16 GB on-board DDR5 memory

1 GbE out-of-band management port

M.2 / U.2 connectors options for
direct attached storage

Bluefield 3

DPU : Data Processing Unit



Bluefield 3 DPU – 2x 200Gb/s
FHHL form factor

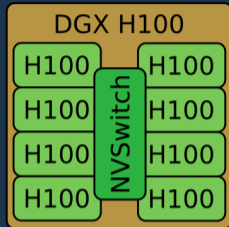
1, 2, 4 ports with up to 400Gb/s connectivity
(Ethernet or NDR InfiniBand)

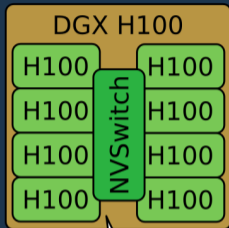
16 GB on-board DDR5 memory

1 GbE out-of-band management port

M.2 / U.2 connectors options for
direct attached storage

Form factors: HHHL, FHHL



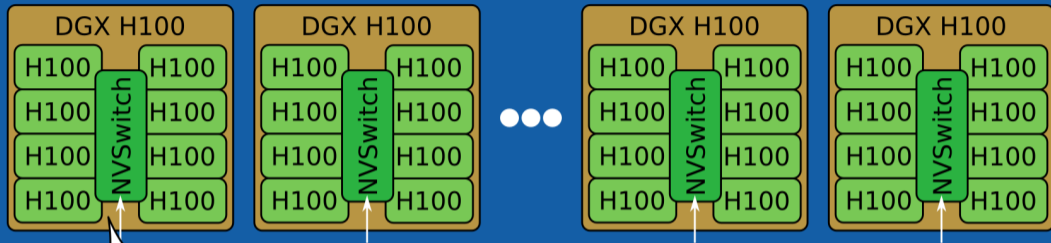


DGX H100 :

- 8 H100
- 32 PFlop (AI performance)
- 640 GB HBM3 Memory
- 24 TB/s memory Bandwidth

DGX H100 / DGX POD H100

DGX POD H100

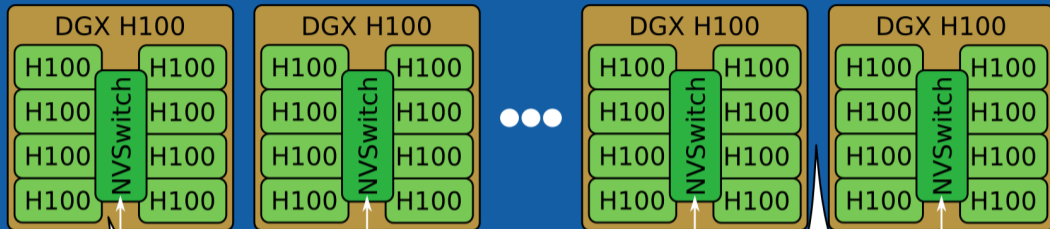


DGX H100 :

- 8 H100
- 32 PFlop (AI performance)
- 640 GB HBM3 Memory
- 24 TB/s memory Bandwidth

DGX H100 / DGX POD H100

DGX POD H100



DGX H100 :

- 8 H100
- 32 PFlop (AI performance)
- 640 GB HBM3 Memory
- 24 TB/s memory Bandwidth

DGX POD H100 :

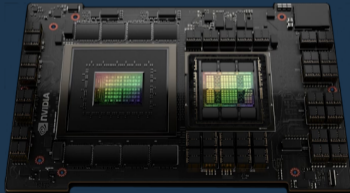
- 32 DGX H100
- 1 EFlop (AI performance)
- 20 TB HBM3 Memory
- 70 TB/s memory Bandwidth
- 192 TFlops of Sharp in-Network compute



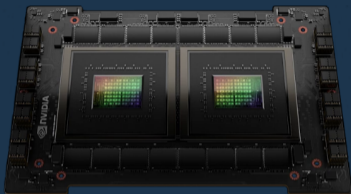
Grace Supership



Grace Supership

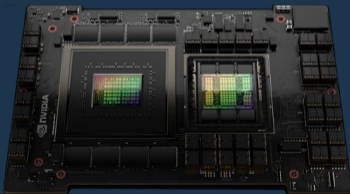


Grace Hopper



Grace Supership

1906 - 1992



Grace Hopper

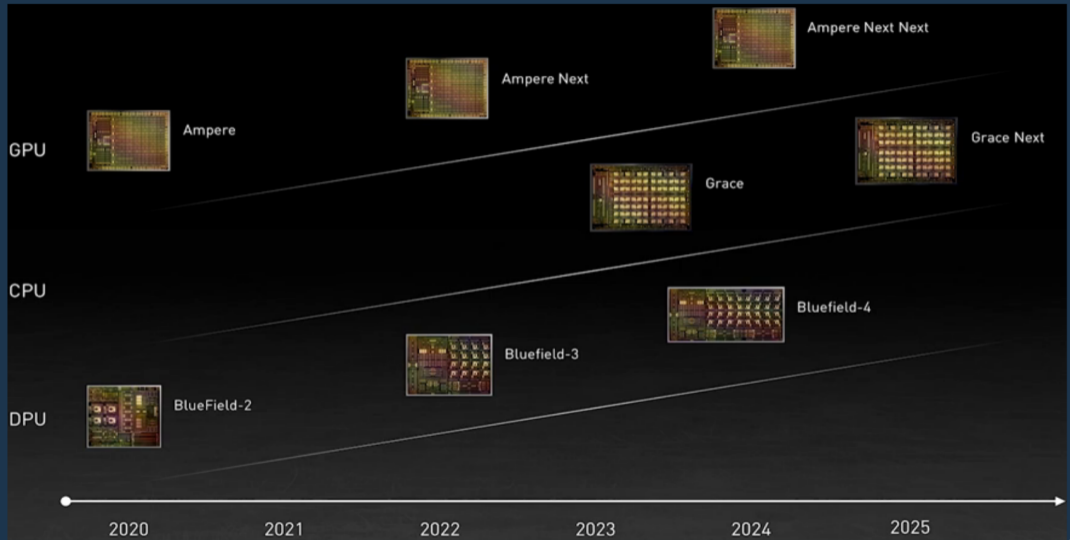


Grace architecture 2023

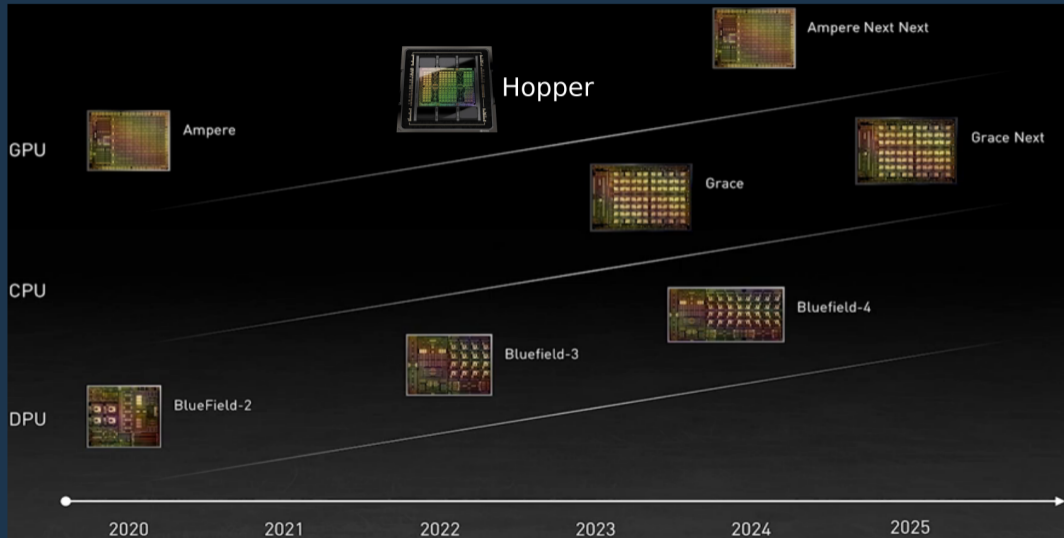
CX7 : 400 Gbit/s



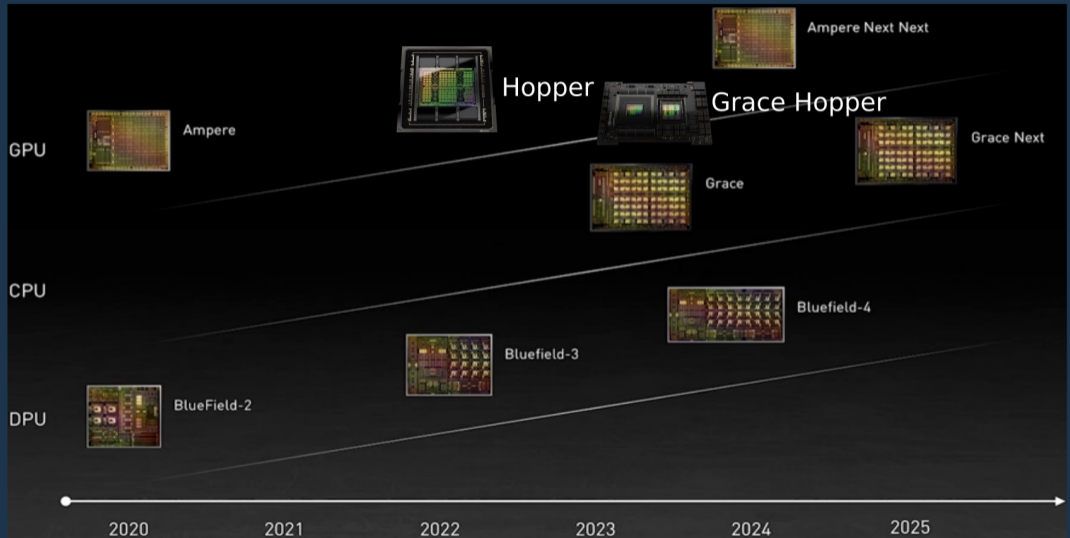
NVidia roadmap



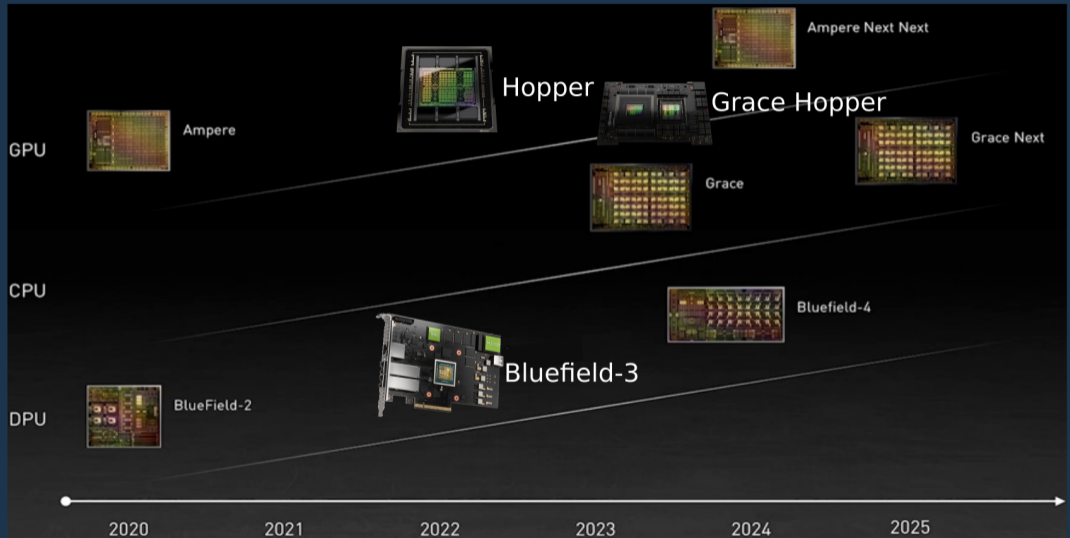
NVidia roadmap



NVidia roadmap



NVidia roadmap



Fermat "Storage Attached Computing"



- Storage : **Flash**
- Computing : **ARM + FPGA**

Fermat "Storage Attached Computing"



- Storage : **Flash**
- Computing : **ARM + FPGA**

Extract & Label atmospheric and river feature from **8 TB** Timed Raw **NOAA Satellite**

Dell - Xeon Server :

- **2** Intel Xeon Gold 6240 @ 2.6 GHz
- **18** cores, **36** threads, 24.75 MB cache each
- **384 GB** DDR4 @ 2933 MHz
- Storage **RAID 5**, **16 x 7.67 TB** SSD
- OS : CentOS Linux release 7.6.1810

Time to complete : **44 minutes**

Fermat :

- **9** FPGAs
- **192 GB** DDR4 DRAM
- **8 TB** NAND Flash

Time to complete : **22 seconds (118x)**

NVC++

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurrency

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVC++

NVCC

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVCC

Autodetection GPU

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVCC

Autodetection **GPU**

128 bits integers supported since version 11.6

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVCC

Autodetection GPU

128 bits integers supported since version 11.6

Link Time Optimization (LTO)

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVCC

Autodetection **GPU**

128 bits integers supported since version 11.6

Link Time Optimization (LTO)

Parallel compilation for different compute capabilities (-d -thread)

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurrency

C++23 :

- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :

- Executors
- Linear algebra

NVCC

Autodetection **GPU**

128 bits integers supported since version 11.6

Link Time Optimization (LTO)

Parallel compilation for different compute capabilities (-d -thread)

Implementation detail (-d -v)
(ex: a+b)

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVCC

Autodetection **GPU**

128 bits integers supported since version 11.6

Link Time Optimization (LTO)

Parallel compilation for different compute capabilities (-d -thread)

Implementation detail (-d -v)
(ex: a+b)

More Parallel compilation coming



NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVCC

Autodetection **GPU**

128 bits integers supported since version 11.6

Link Time Optimization (LTO)

Parallel compilation for different compute capabilities (-d -thread)

Implementation detail (-d -v)
(ex: a+b)

More Parallel compilation coming



Ongoing migration **LLVM 7**



NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVCC

Autodetection **GPU**

128 bits integers supported since version 11.6

Link Time Optimization (LTO)

Parallel compilation for different compute capabilities (-d -thread)

Implementation detail (-d -v)
(ex: a+b)

More Parallel compilation coming



Ongoing migration **LLVM 7**



C++20 :



- Module but no cuda call

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVCC

Autodetection **GPU**

128 bits integers supported since version 11.6

Link Time Optimization (LTO)

Parallel compilation for different compute capabilities (-d -thread)

Implementation detail (-d -v)
(ex: a+b)

More Parallel compilation coming



Ongoing migration **LLVM 7**



C++20 :



- Module but no cuda call

NVFortran

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVCC

Autodetection **GPU**

128 bits integers supported since version 11.6

Link Time Optimization (LTO)

Parallel compilation for different compute capabilities (-d -thread)

Implementation detail (-d -v)
(ex: a+b)

More Parallel compilation coming



Ongoing migration **LLVM 7**



C++20 :



- Module but no cuda call

NVFortran

Uses **nvTENSOR**

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVCC

Autodetection **GPU**

128 bits integers supported since version 11.6

Link Time Optimization (LTO)

Parallel compilation for different compute capabilities (-d -thread)

Implementation detail (-d -v)
(ex: a+b)

More Parallel compilation coming



Ongoing migration **LLVM 7**



C++20 :



- Module but no cuda call

NVFortran

Uses **nvTENSOR**

Tensor Core available

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurrency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVCC

Autodetection **GPU**

128 bits integers supported since version 11.6

Link Time Optimization (LTO)

Parallel compilation for different compute capabilities (-d -thread)

Implementation detail (-d -v) (ex: a+b)

More Parallel compilation coming



Ongoing migration **LLVM 7**



C++20 :



- Module but no cuda call

NVFortran

Uses **nvTENSOR**

Tensor Core available

Do Concurrent :

- No function call
- Reduction not available yet

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurrency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVCC

Autodetection GPU

128 bits integers supported since version 11.6

Link Time Optimization (LTO)

Parallel compilation for different compute capabilities (-d -thread)

Implementation detail (-d -v)
(ex: a+b)

More Parallel compilation coming



Ongoing migration **LLVM 7**



C++20 :



- Module but no cuda call

NVFortran

Uses **nvTENSOR**

Tensor Core available

Do Concurrent :

- No function call
- Reduction not available yet

Same Perf as **OpenACC**

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurrency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVCC

Autodetection **GPU**

128 bits integers supported since version 11.6

Link Time Optimization (LTO)

Parallel compilation for different compute capabilities (-d -thread)

Implementation detail (-d -v) (ex: a+b)

More Parallel compilation coming



Ongoing migration **LLVM 7**



C++20 :



- Module but no cuda call

NVFortran

Uses **nvTENSOR**

Tensor Core available

Do Concurrent :

- No function call
- Reduction not available yet

Same Perf as **OpenACC**

Cuda Fortran not portable

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurrency

C++23 :



- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :



- Executors
- Linear algebra

NVCC

Autodetection GPU

128 bits integers supported since version 11.6

Link Time Optimization (LTO)

Parallel compilation for different compute capabilities (-d -thread)

Implementation detail (-d -v) (ex: a+b)

More Parallel compilation coming



Ongoing migration **LLVM 7**



C++20 :



- Module but no cuda call

NVFortran

Uses **nvTENSOR**

Tensor Core available

Do Concurrent :

- No function call
- Reduction not available yet

Same Perf as **OpenACC**

Cuda Fortran not portable

For now **NVidia** only :



- **Intel, Cray** ongoing

NVC++

C++17 :

- Combine **std::transform** and **std::reduce** in the same call

C++20 :

- Parallelism Thread concurrency

C++23 :

- Range based parallel algorithm (mdspan, mdarray)
- Multi-Dim Array abstractions

C++2X :

- Executors
- Linear algebra

NVCC

Autodetection GPU

128 bits integers supported since version 11.6

Link Time Optimization (LTO)

Parallel compilation for different compute capabilities (-d -thread)

Implementation detail (-d -v) (ex: a+b)

More Parallel compilation coming 

Ongoing migration **LLVM 7** 

C++20 : 

- Module but no cuda call

NVFortran

Uses **nvTENSOR**

Tensor Core available

Do Concurrent :

- No function call
- Reduction not available yet

Same Perf as **OpenACC**

Cuda Fortran not portable

For now **NVidia** only : 

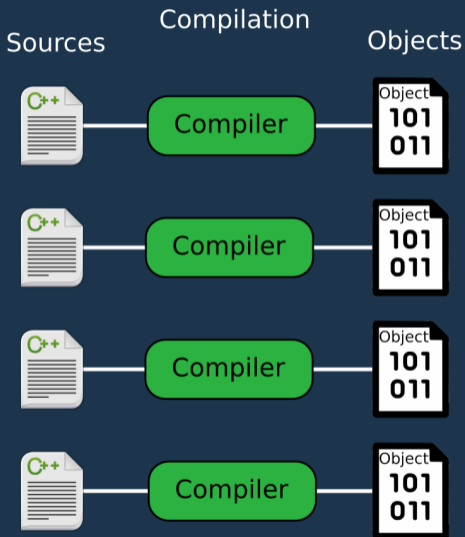
- **Intel, Cray** ongoing

Fortran Forever :

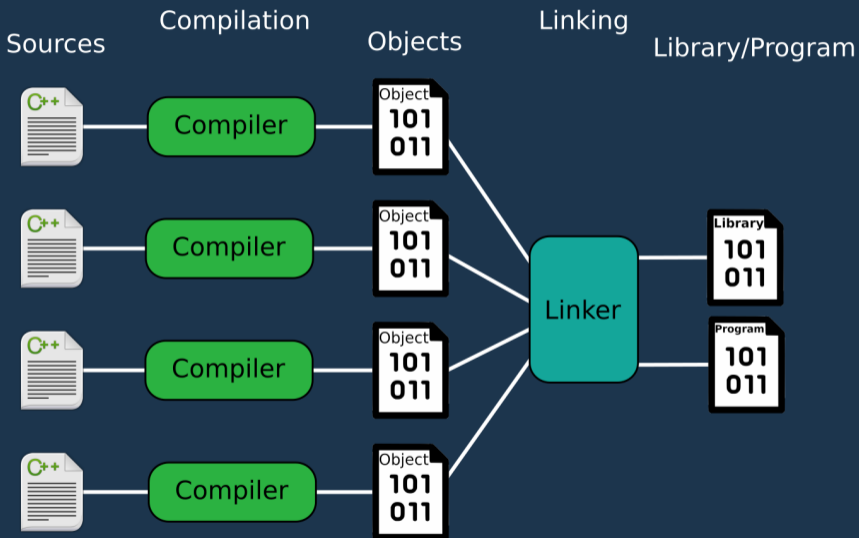
- Modern Fortran **>= F18**

Sources

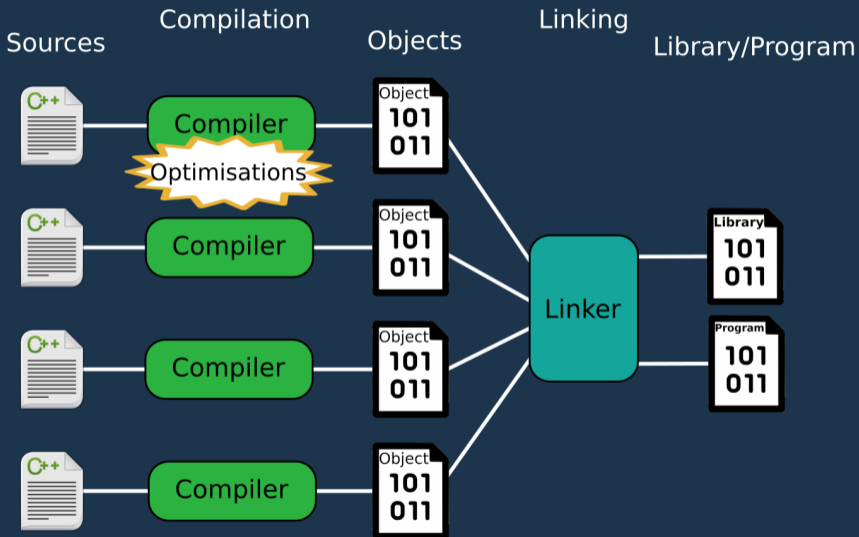




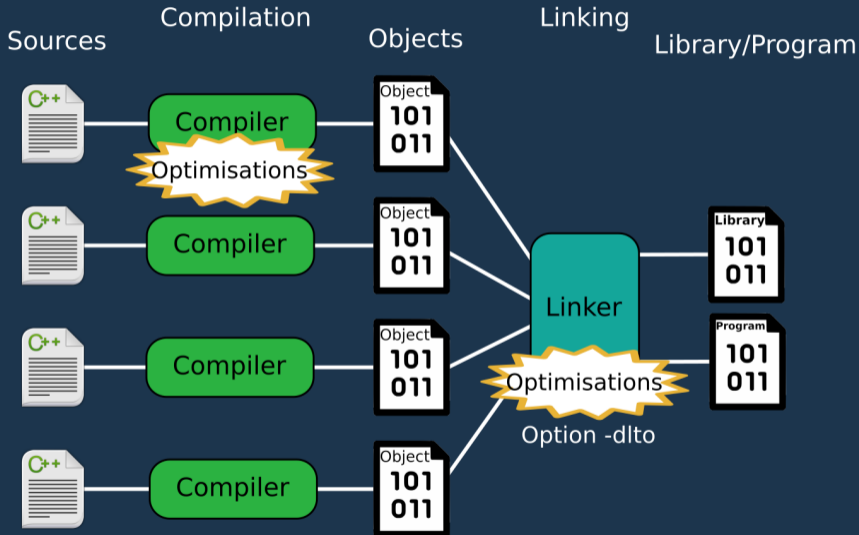
Link Time Optimization



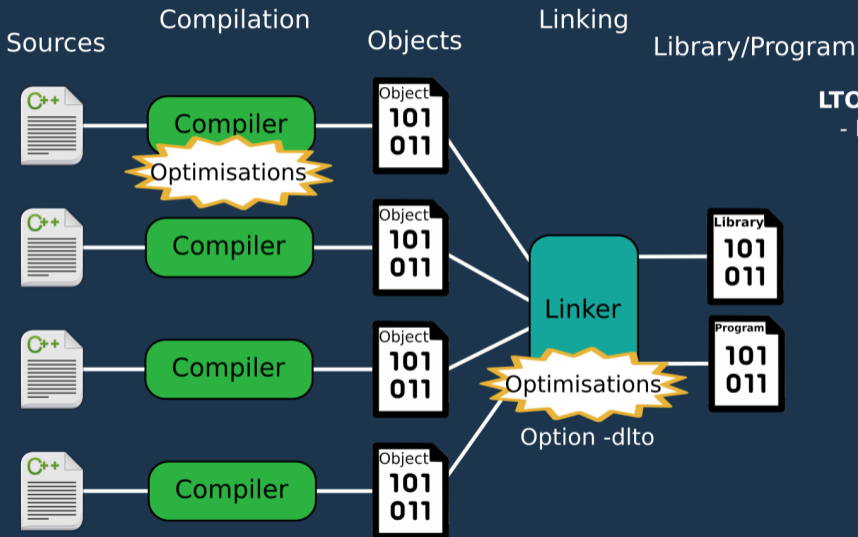
Link Time Optimization



Link Time Optimization

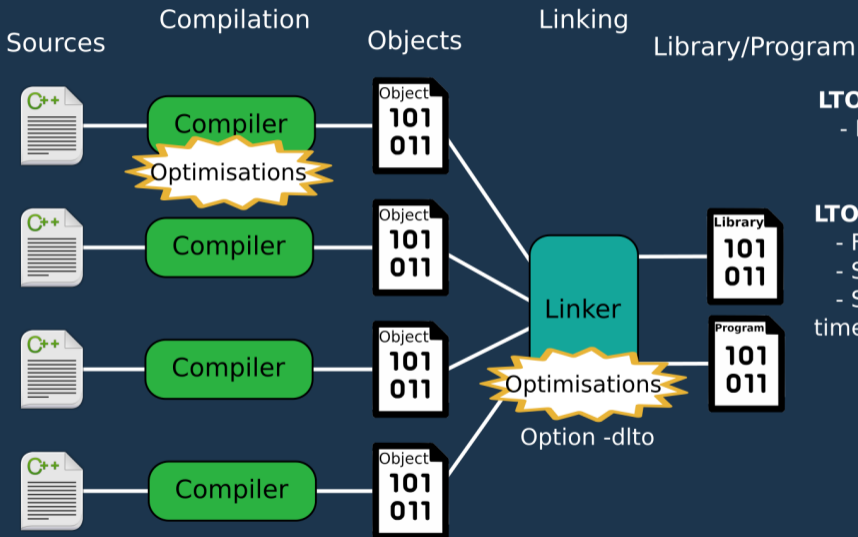


Link Time Optimization



LTO Optimisations :
- Remove dead code
(even kernel)

Link Time Optimization



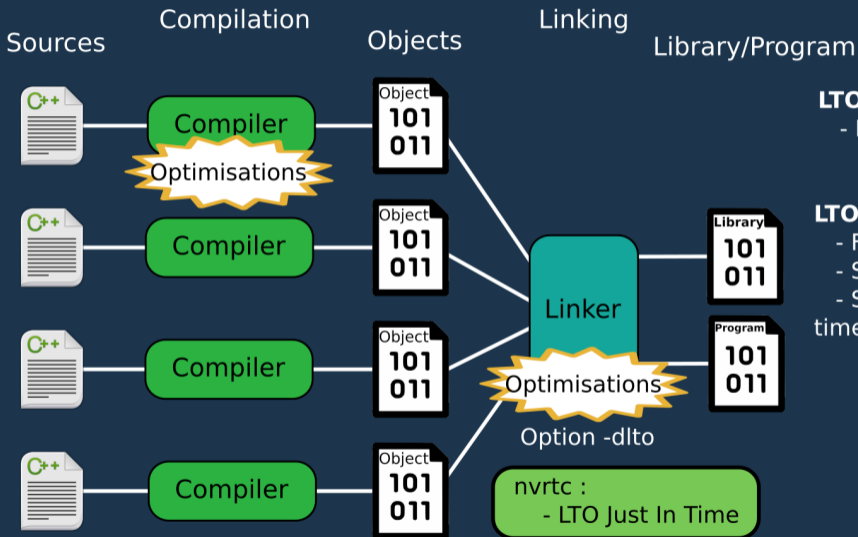
LTO Optimisations :

- Remove dead code (even kernel)

LTO :

- Faster compilation
- Slower linking
- Speed up execution time from 5 to 91 %

Link Time Optimization



LTO Optimisations :

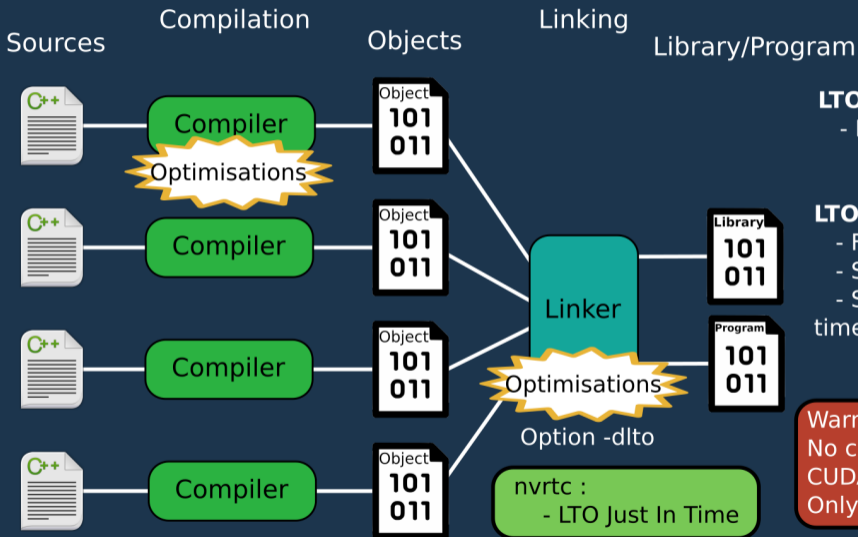
- Remove dead code (even kernel)

LTO :

- Faster compilation
- Slower linking
- Speed up execution time from 5 to 91 %

nvrvc :
- LTO Just In Time

Link Time Optimization



LTO Optimisations :

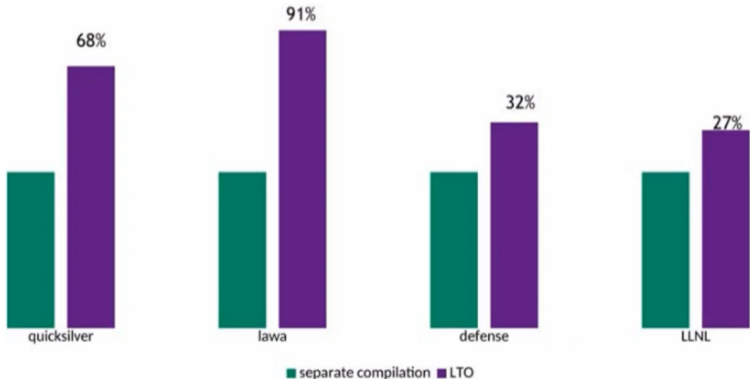
- Remove dead code (even kernel)

LTO :

- Faster compilation
- Slower linking
- Speed up execution time from 5 to 91 %

Warning :
No compatibility on
CUDA minor versions 11.X
Only from CUDA 12.0

Compilation Elapsed Time



gram

LTO Optimisations :

- Remove dead code (even kernel)

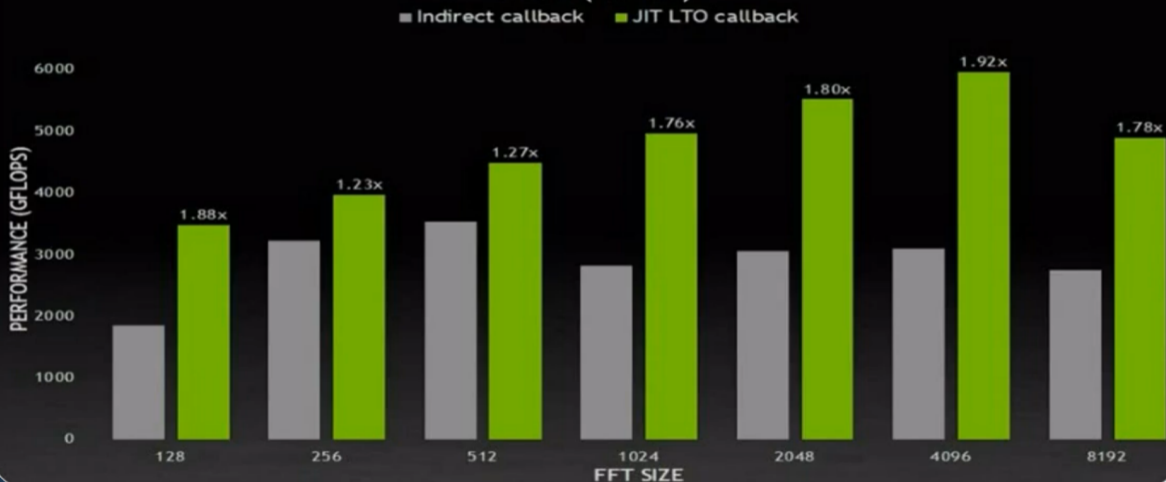
LTO :

- Faster compilation
- Slower linking
- Speed up execution time from 5 to 91 %

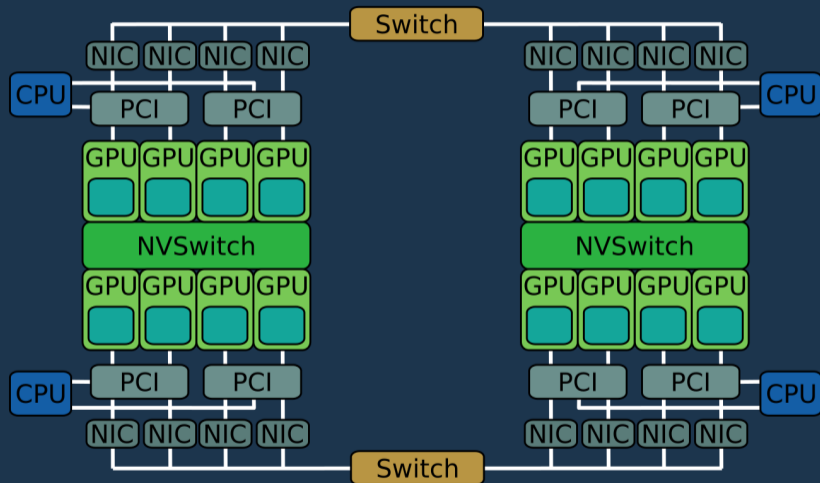
Warning :
No compatibility on
CUDA minor versions 11.X
Only from CUDA 12.0

Link Time Optimization

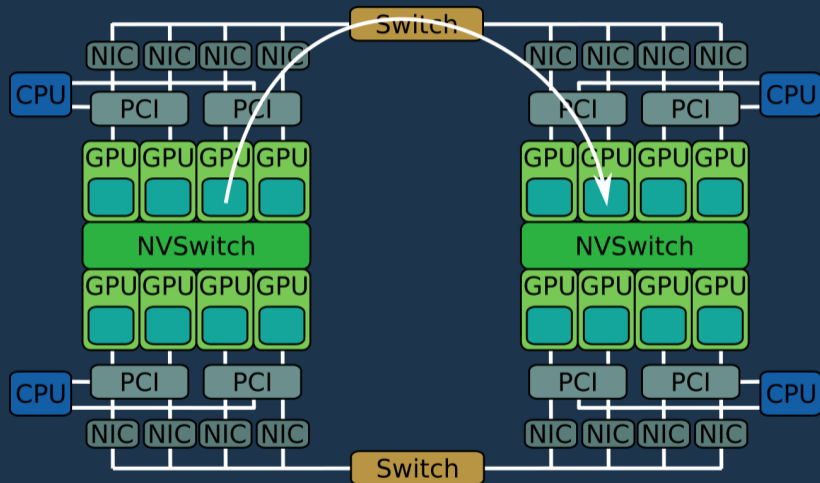
INDIRECT CALLBACK VS JIT LTO CALLBACK
A100 (40 GB)



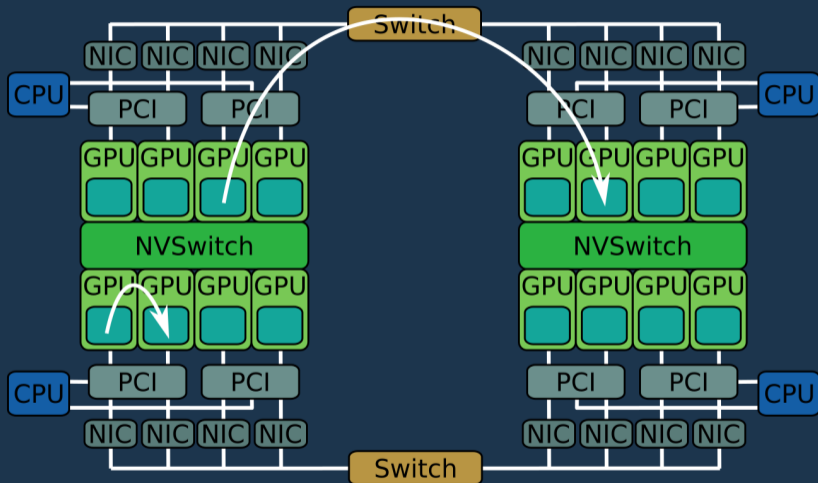
Communication : NVShmem



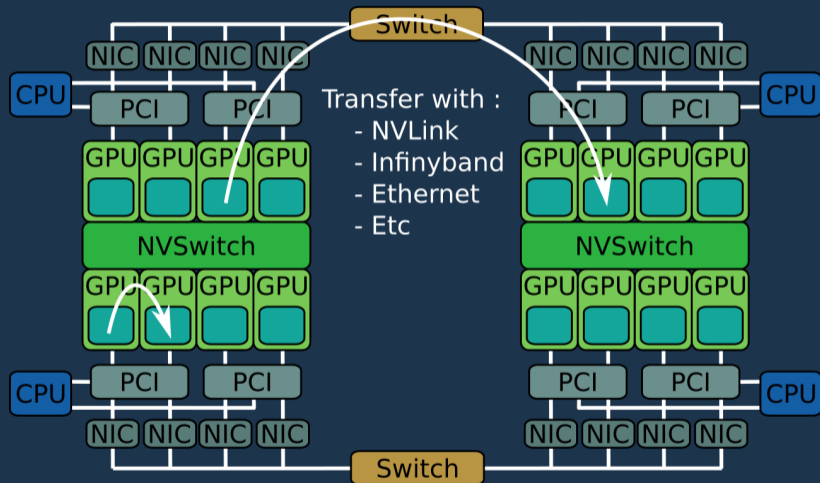
Communication : NVShmem



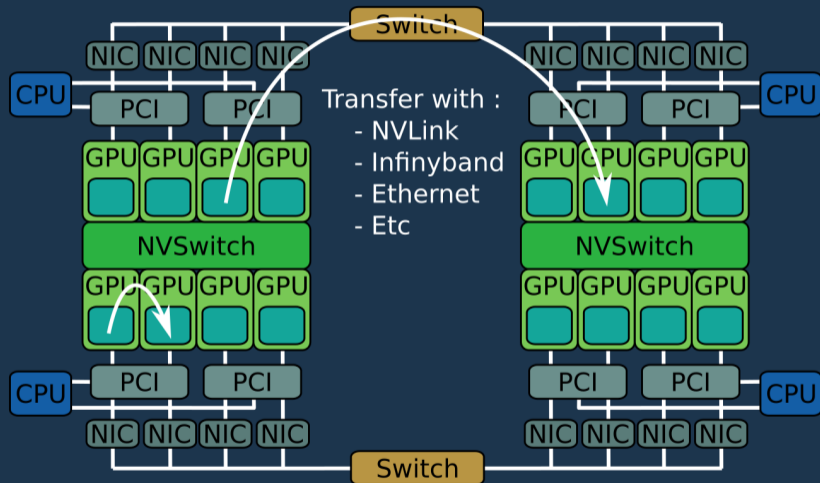
Communication : NVShmem



Communication : NVShmem



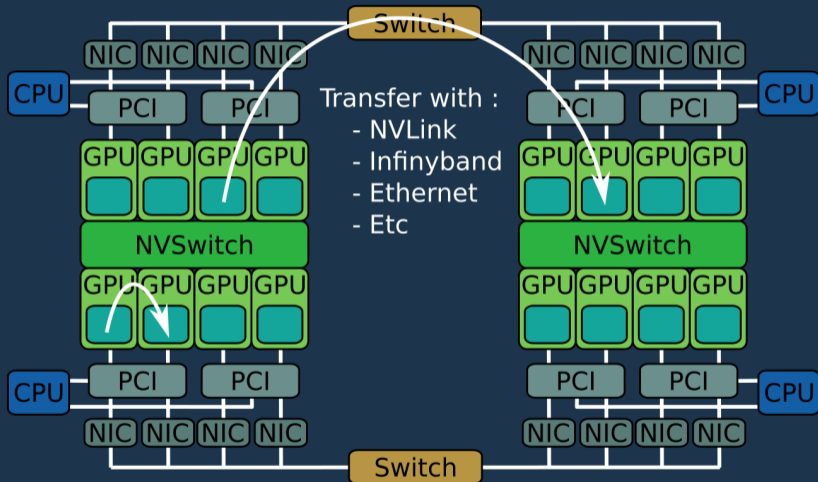
Based on OpenShmem



Communication : NVShmem

Based on OpenShmem

Scalable point to point communication

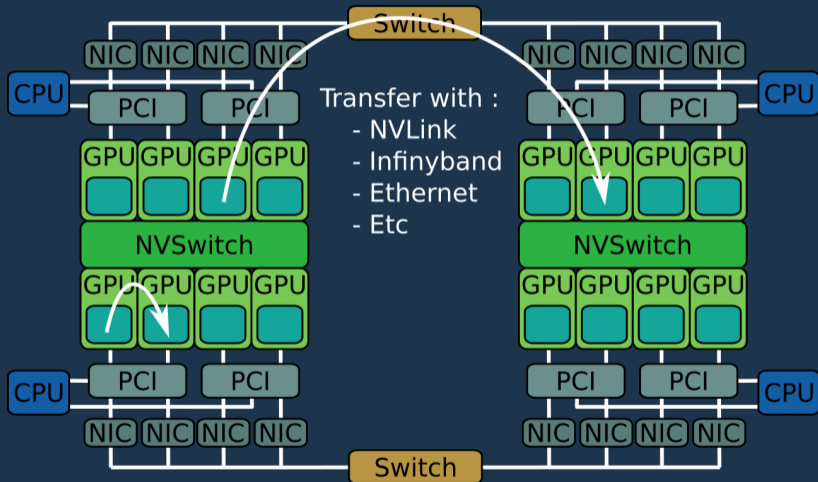


Communication : NVShmem

Based on OpenShmem

Scalable point to point communication

Symmetric allocations



Communication : NVShmem

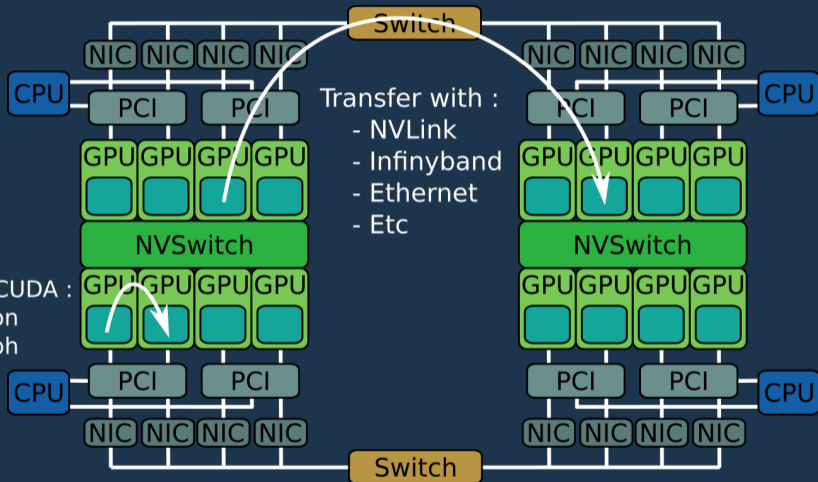
Based on OpenShmem

Scalable point to point communication

Symmetric allocations

Communication integrated with CUDA :

- GPU Kernel-Initiated Operation
- Operation CUDA/Stream graph
- CPU Initiated operations



Communication : NVShmem

Based on OpenShmem

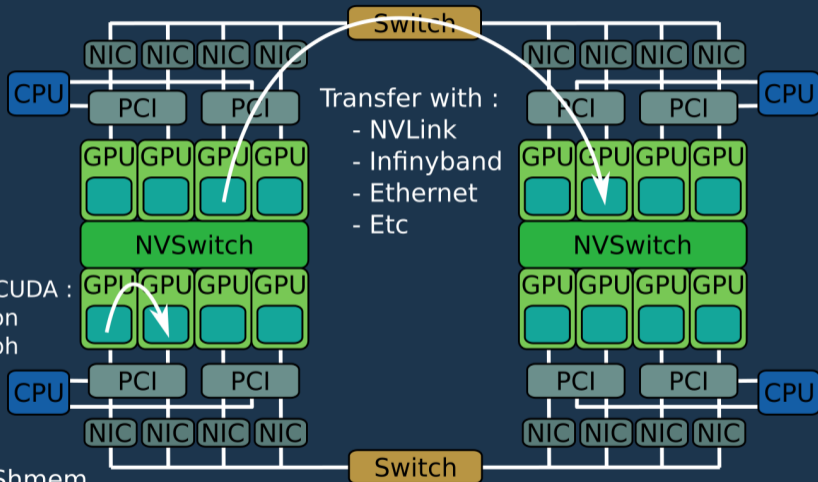
Scalable point to point communication

Symmetric allocations

Communication integrated with CUDA :

- GPU Kernel-Initiated Operation
- Operation CUDA/Stream graph
- CPU Initiated operations

MagnumIO is based on NVShmem



Communication : MagnumIO

RTX

HPC

RAPIDS

AI:
RIVA,
MERLIN

CLARA

METRO

DRIVE

ISAAC

AERIAL

CUDA-X

AI

DA

Graph

ML

DL Training

DL Inference

HPC

Linear Algebra

Parallel Algorithms

Signal Processing

Image Processing

...

CUDA

MAGNUM IO

NETWORK IO

STORAGE IO

IN-NETWORK COMPUTE

IO MANAGEMENT

- GPU Direct RDMA
- MOFED
- NCCL
- NVidia ASAP - Accelerating Switch and Packet Processing
- NVidia HPC-X
- NVShmem

CLARA

METRO

DRIVE

ISAAC

AERIAL

CUDA-X

ML

DL Training

DL Inference

AI

HPC

Parallel Algorithms

Signal Processing

Image Processing

...

CUDA

MAGNUM IO

NETWORK IO

STORAGE IO

IN-NETWORK COMPUTE

IO MANAGEMENT

Communication : MagnumIO

- GPU Direct RDMA
- MOFED
- NCCL
- NVidia ASAP - Accelerating Switch and Packet Processing
- NVidia HPC-X
- NVShmem

CLARA

METRO

DRIVE

ISAAC

AERIAL

CUDA-X

ML

DL Training

DL Inference

HPC

Parallel Algorithms

Signal Processing

Image Processing

...

- GPU Direct Storage
- SNAP

CUDA

MAGNUM IO

NETWORK IO

STORAGE IO

IN-NETWORK COMPUTE

IO MANAGEMENT

Communication : MagnumIO

- GPU Direct RDMA
- MOFED
- NCCL
- NVidia ASAP - Accelerating Switch and Packet Processing
- NVidia HPC-X
- NVShmem

- Hardware tag matching
- NVidia Scalable Hierarchical Aggregation and Reduction Protocol (SHARP)

- GPU Direct Storage
- SNAP

CLARA

METRO

DRIVE

ISAAC

AERIAL

CUDA-X

ML

Inference

Parallel Algorithms

Signal

Image Processing

CUDA

MAGNUM IO

NETWORK IO

STORAGE IO

IN-NETWORK COMPUTE

IO MANAGEMENT

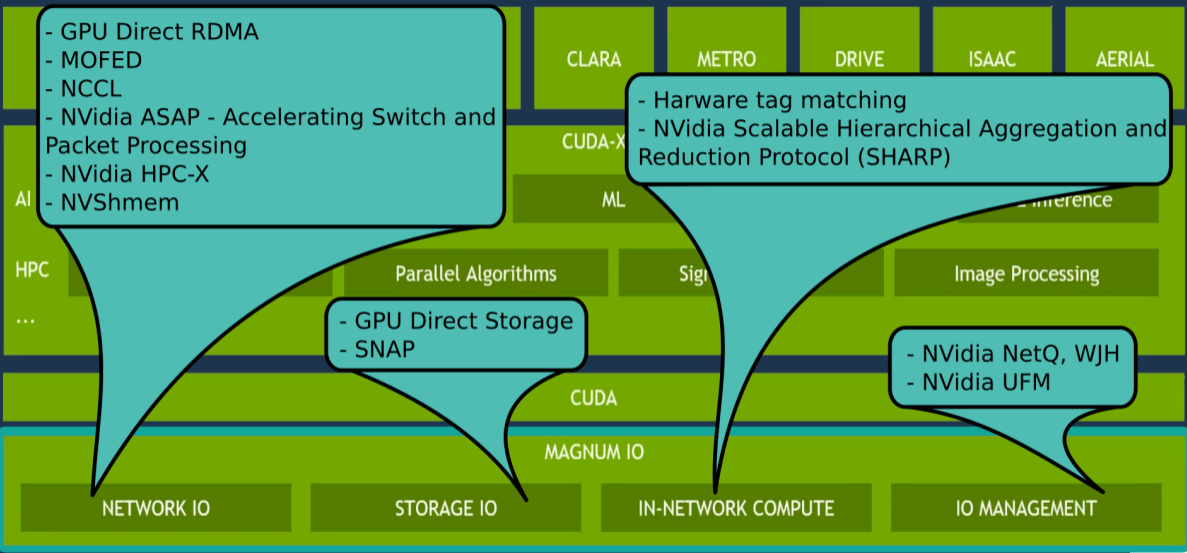
Communication : MagnumIO

- GPU Direct RDMA
- MOFED
- NCCL
- NVidia ASAP - Accelerating Switch and Packet Processing
- NVidia HPC-X
- NVShmem

- Hardware tag matching
- NVidia Scalable Hierarchical Aggregation and Reduction Protocol (SHARP)

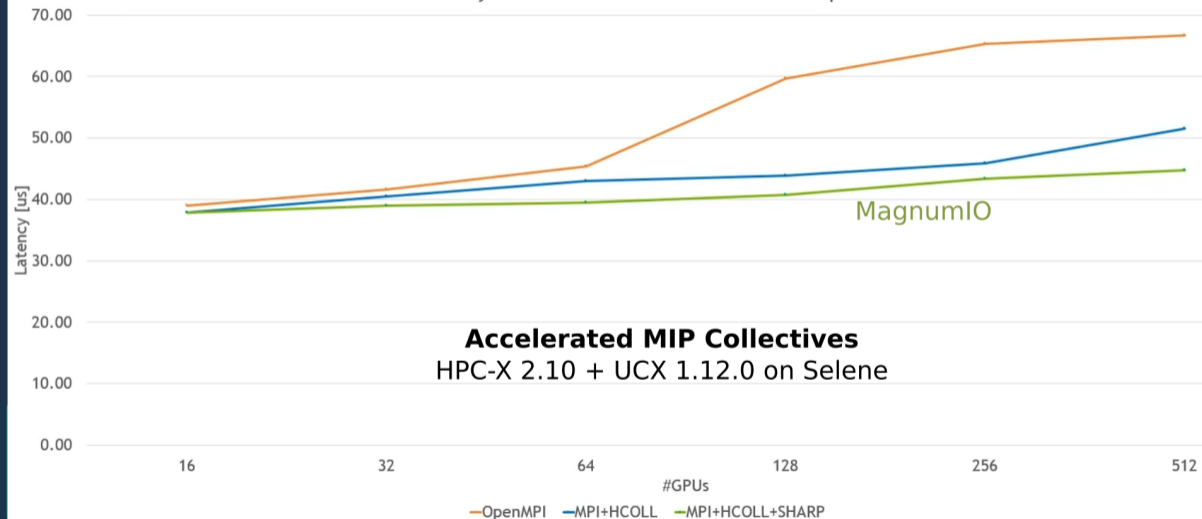
- GPU Direct Storage
- SNAP

- NVidia NetQ, WJH
- NVidia UFM

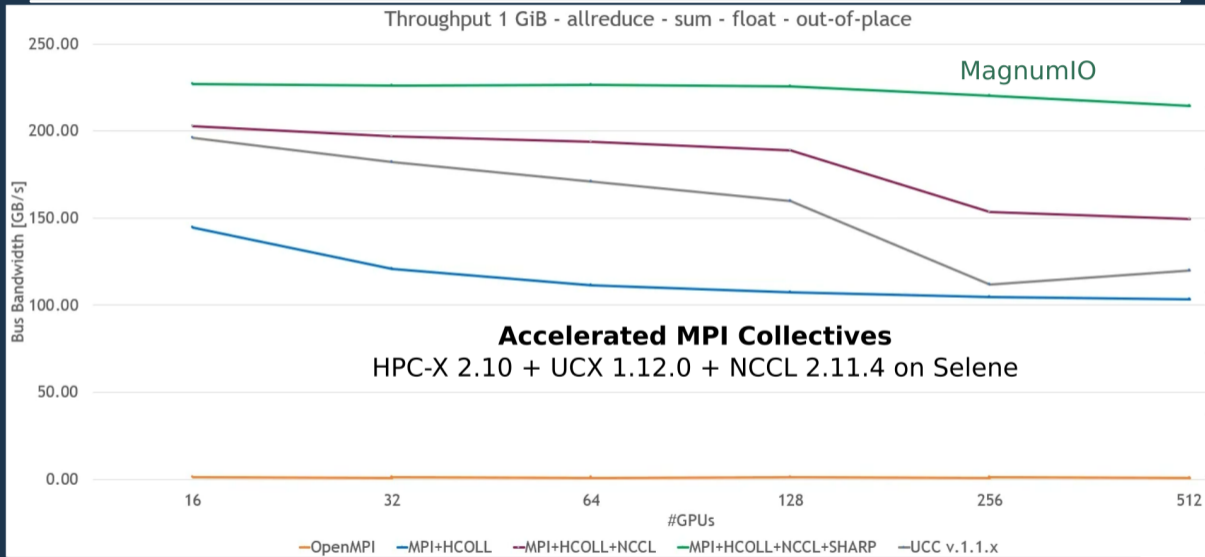


Communication : MagnumIO

Latency 8b - allreduce - sum - float - out-of-place

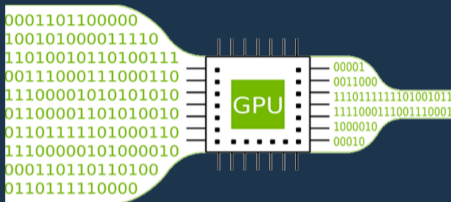


Communication : MagnumIO



Data Compression : NVComp

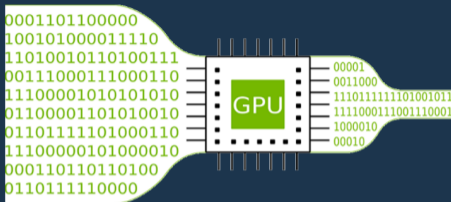
C/C++ standard calls



Compression / decompression algorithms :

- LZ4
- Snappy
- DGeflate
- Cascaded
- Bitcomp
- ANS (Asymmetric Numeral System)

C/C++ standard calls



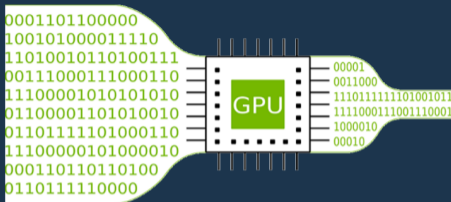
Compression / decompression algorithms :

- LZ4
- Snappy
- DGeflate
- Cascaded
- Bitcomp
- ANS (Asymetric Numeral System)

Release 2.2 :

- Can be used with **GPU Direct Storage**
- High Level API simplified
- No kernel call available

C/C++ standard calls



Compression / decompression algorithms :

- LZ4
- Snappy
- DGeflate
- Cascaded
- Bitcomp
- ANS (Asymetric Numeral System)

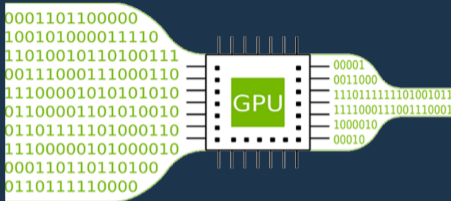
Release 2.2 :

- Can be used with **GPU Direct Storage**
- High Level API simplified
- No kernel call available

Some advices :

- To be used on files of **few MB**
- With **chunks** > 8kB
- Decompressed data must be accessible by the GPU

C/C++ standard calls



Compression / decompression algorithms :

- LZ4
- Snappy
- DGeflate
- Cascaded
- Bitcomp
- ANS (Asymmetric Numeral System)

Release 2.2 :

- Can be used with **GPU Direct Storage**
- High Level API simplified
- No kernel call available

Some advices :

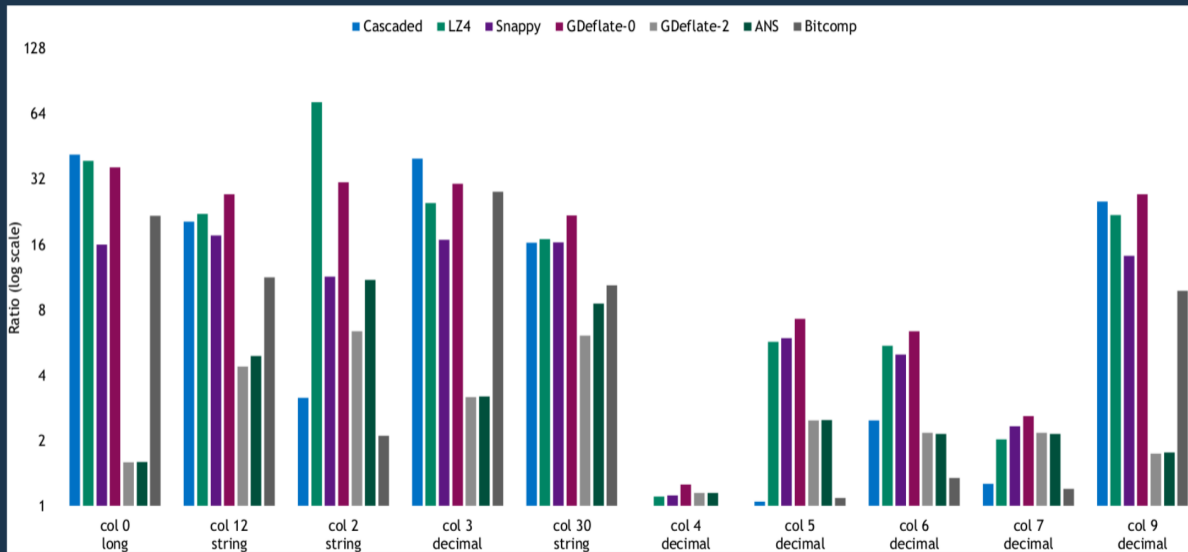
- To be used on files of **few MB**
- With **chunks** > 8kB
- Decompressed data must be accessible by the GPU

News in future release 2.3 :

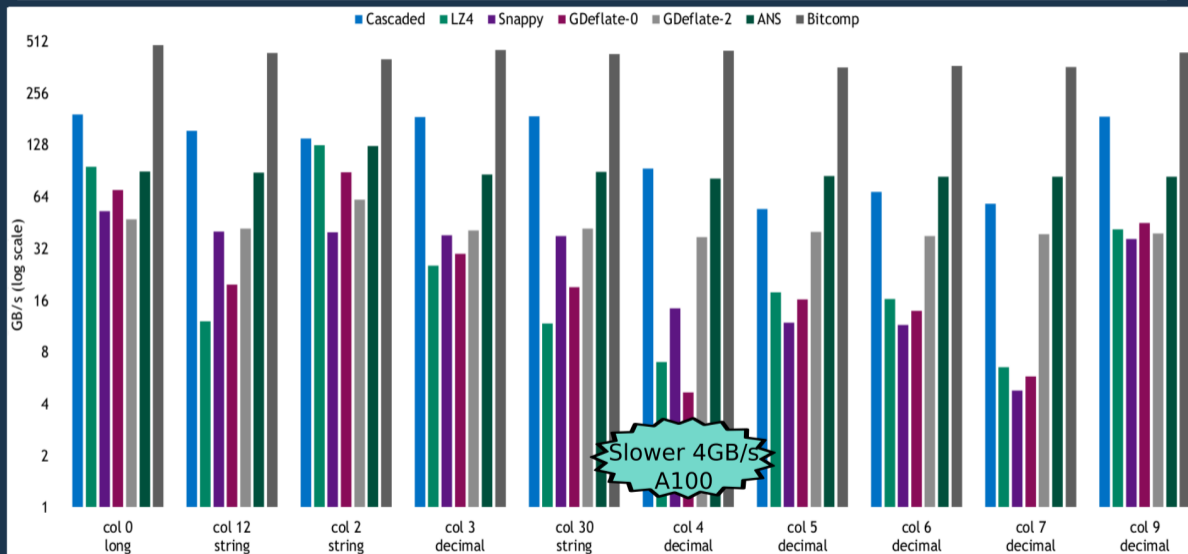


- Add **Zstd** decompression
- Add **BitShuffle** (**float** compression)

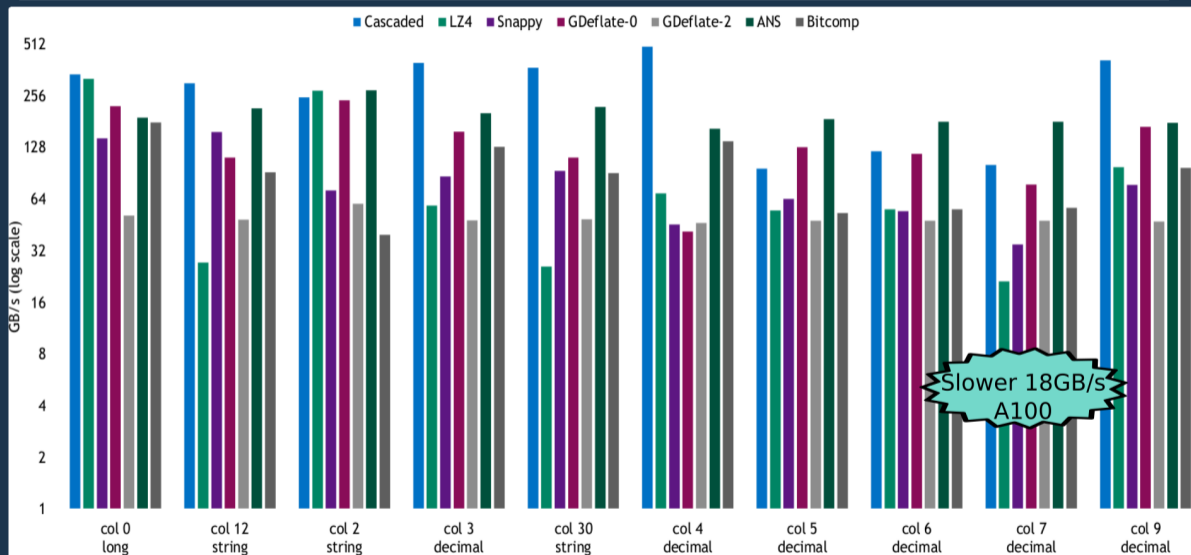
Data Compression : NVComp



Data Compression : NVComp



Data Compression : NVComp



Deep Learning

- Speed up 3x to 6x
- Can use FP16 + FP32

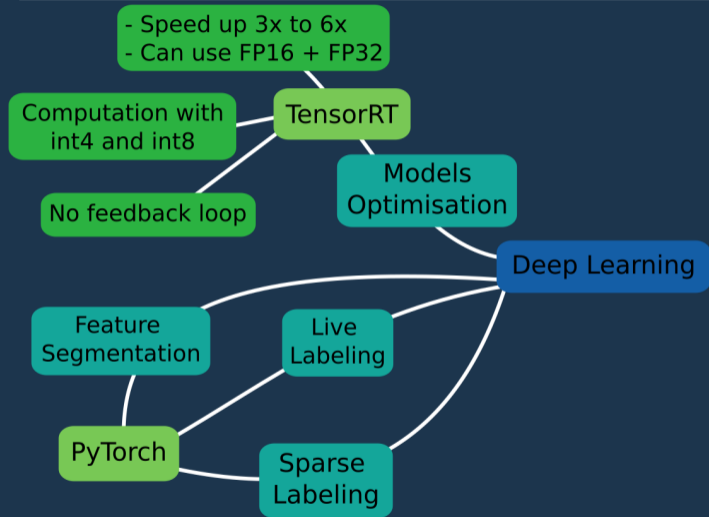
Computation with
int4 and int8

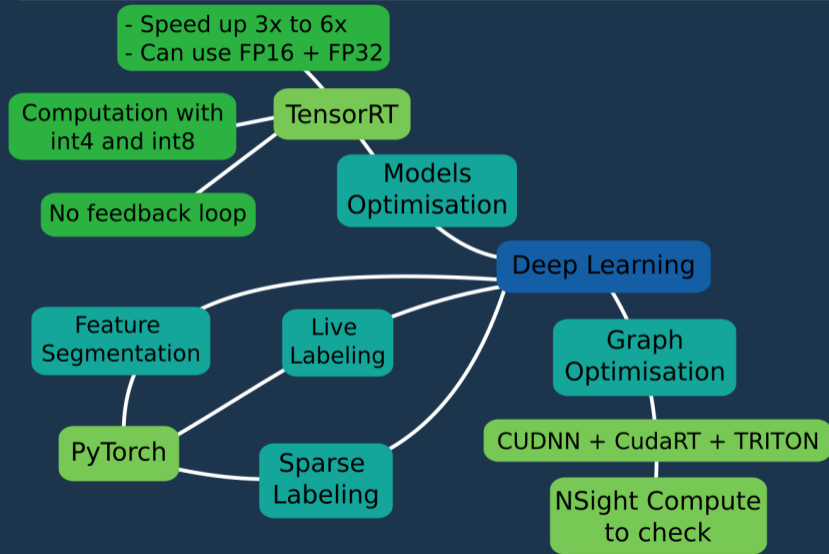
TensorRT

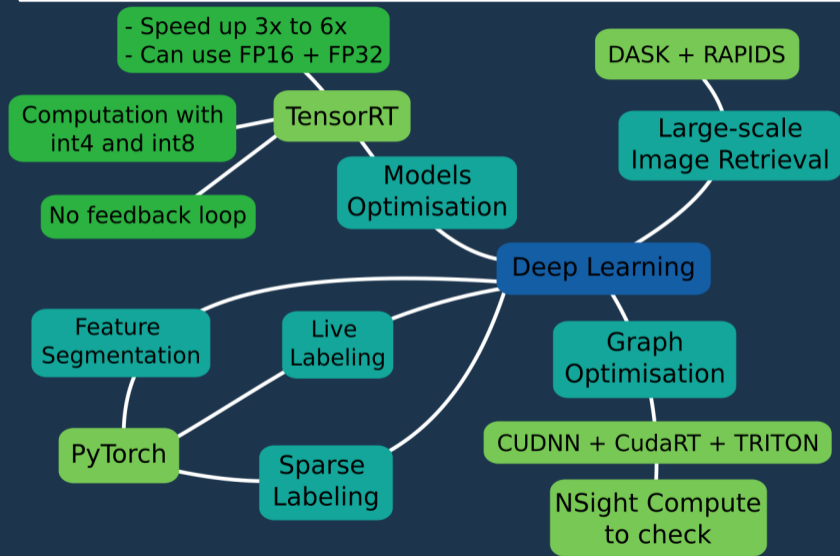
No feedback loop

Models
Optimisation

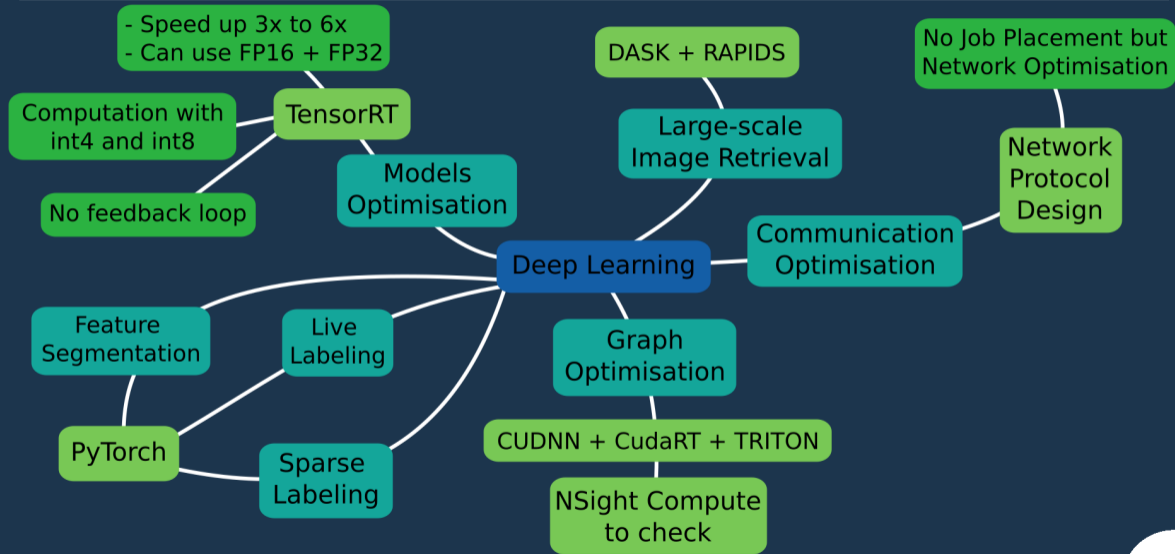
Deep Learning



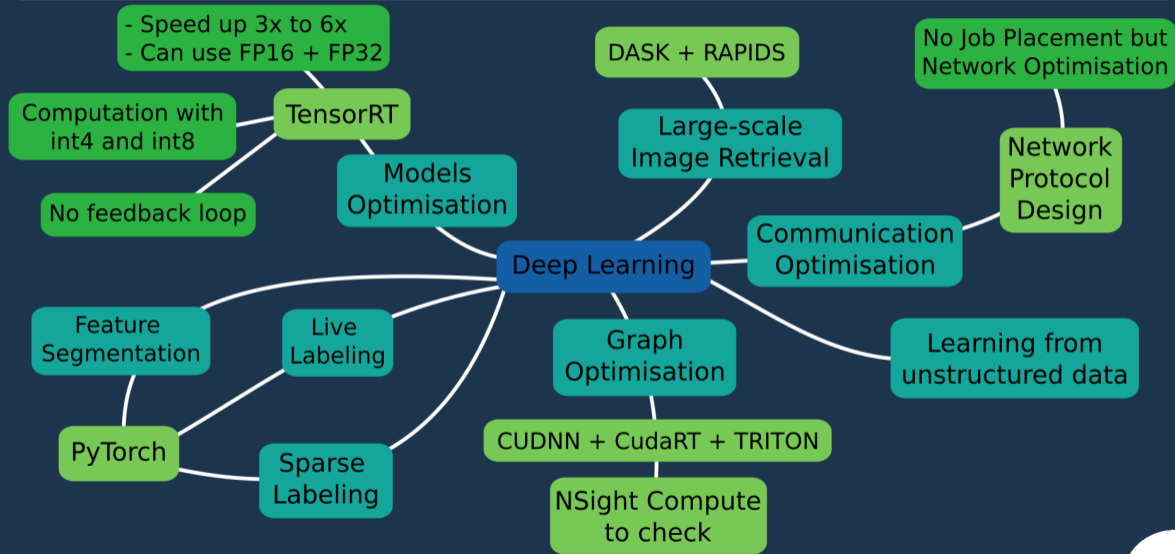




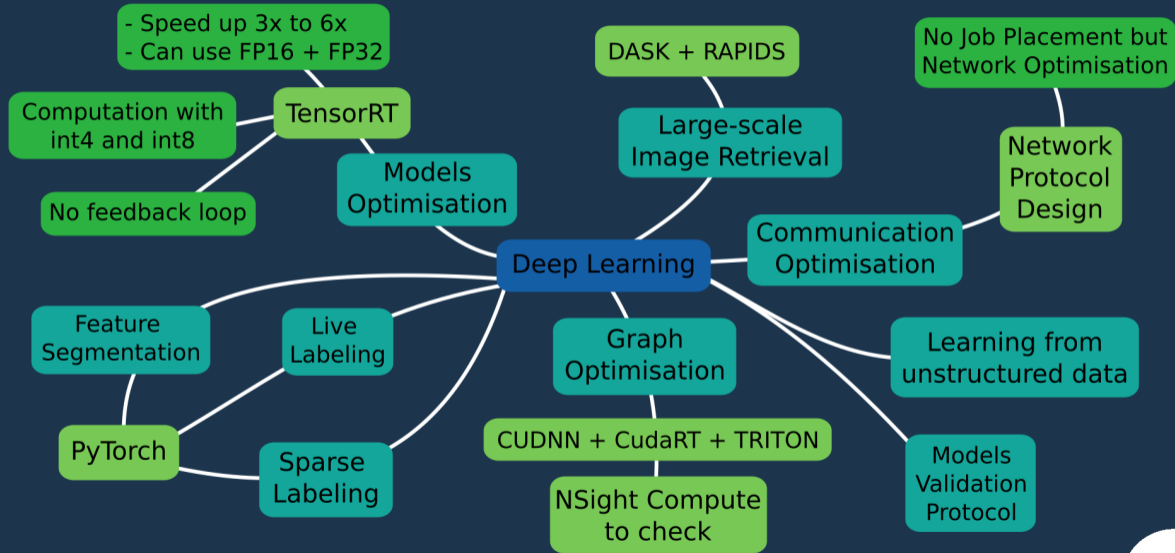
Deep Learning



Deep Learning



Deep Learning



GTC 2022 notes from Reprises website

Sébastien Gignoux : l'architecture Fermat (see 11h-11h15 of first day)