

# Retour GTC 2021 : GPU Technical Conference

Pierre Aubert

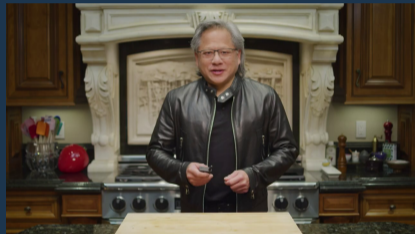


GTC 2021 <https://gtc21.event.nvidia.com/>

This is not possible to summarize **1600** talks !

Only focus on :

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

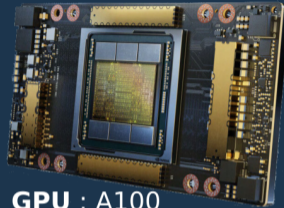


Keynotes video link



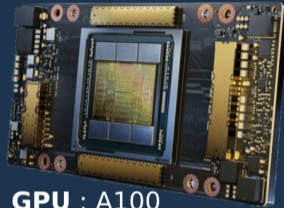
**GTC 2020 :**  
**- NVidia bought Mellanox**

**GTC 2020 :**  
- **NVidia** bought **Mellanox**



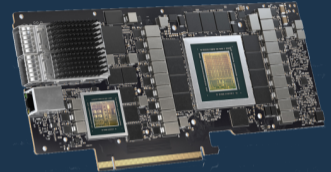
**GPU : A100**  
**40 GB and 80 GB**

**GTC 2020 :**  
- **NVidia** bought **Mellanox**



**GPU : A100**  
**40 GB and 80 GB**

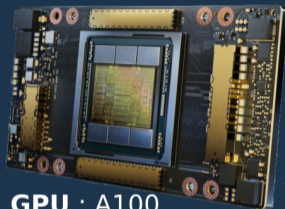
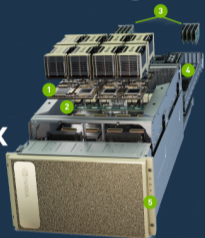
**DPU : Bluefield 2**



GTC 2020 :

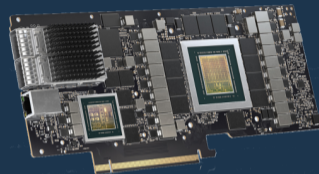
- **NVidia** bought **Mellanox**

**DGX**



**GPU** : A100  
**40 GB** and **80 GB**

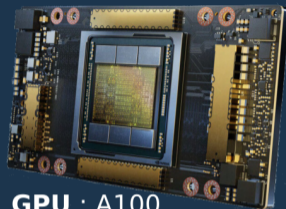
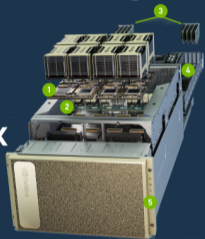
**DPU** : Bluefield 2



GTC 2020 :

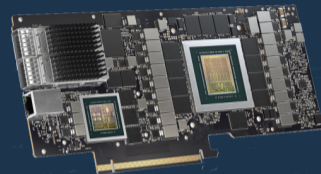
- **NVidia** bought **Mellanox**

**DGX**



**GPU** : A100  
**40 GB** and **80 GB**

**DPU** : Bluefield 2



GTC 2021 :

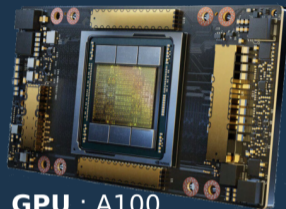
- **NVidia** bought **ARM**



GTC 2020 :

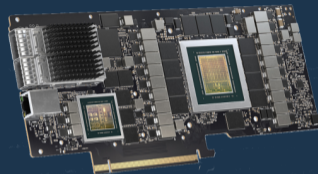
- **NVidia** bought **Mellanox**

DGX



GPU : A100  
40 GB and 80 GB

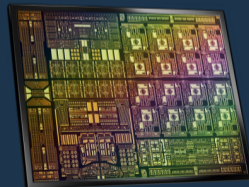
DPU : Bluefield 2



GTC 2021 :

- **NVidia** bought **ARM**

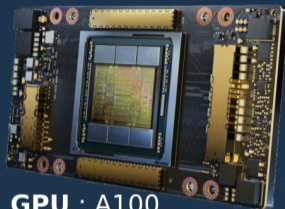
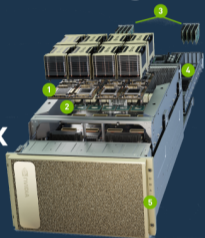
DPU : Bluefield 3



GTC 2020 :

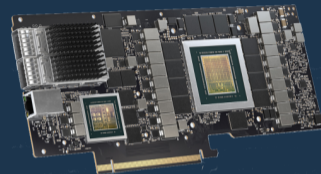
- **NVidia** bought **Mellanox**

**DGX**



**GPU** : A100  
40 GB and 80 GB

**DPU** : Bluefield 2



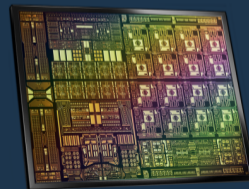
GTC 2021 :

- **NVidia** bought **ARM**

**DGX**

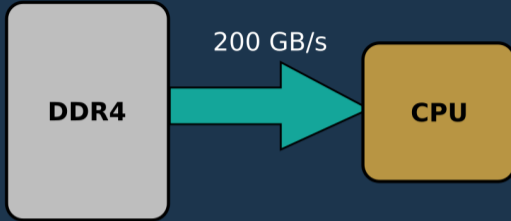


**DPU** : Bluefield 3

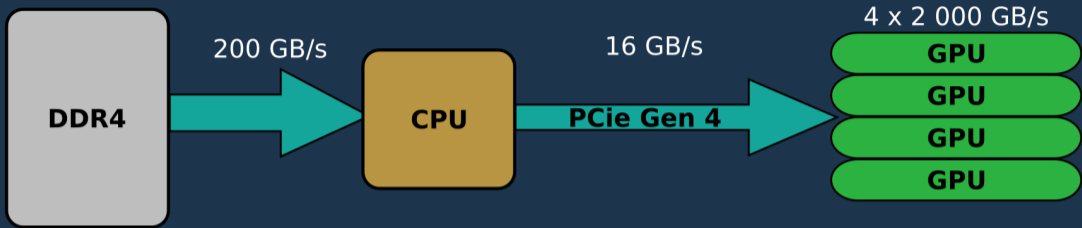


**DDR4**

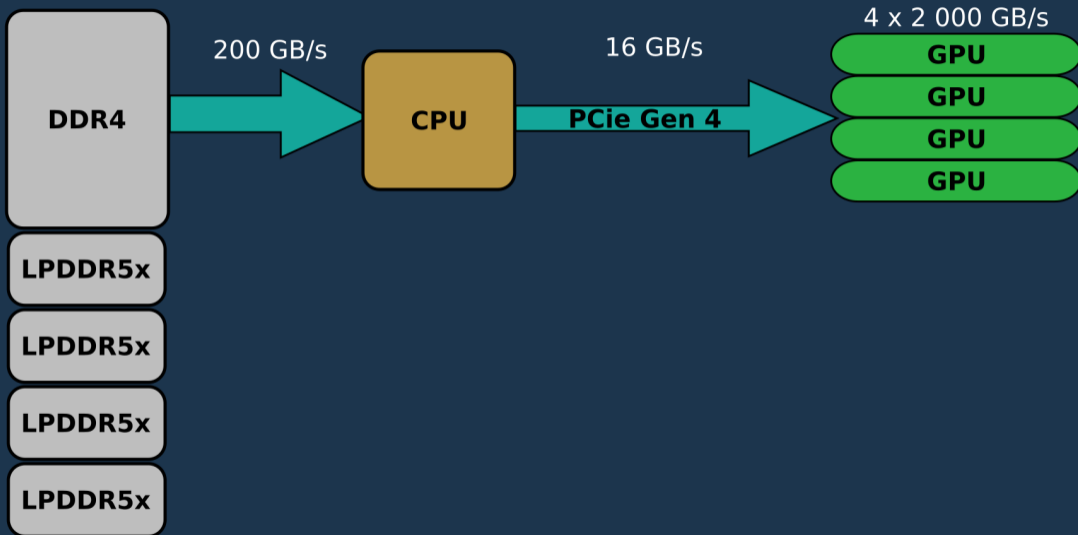
# Hardware evolution (NVIDIA)



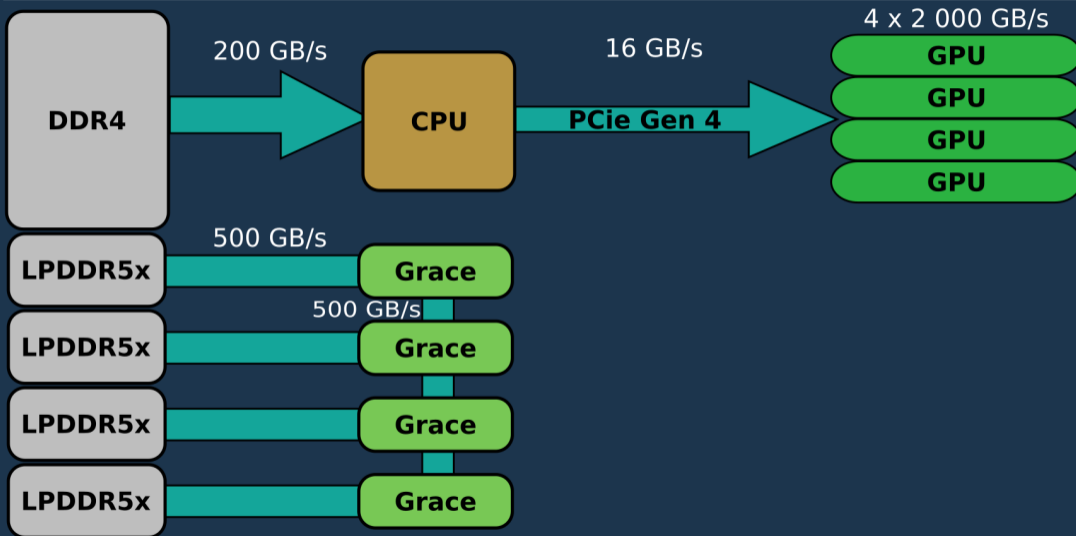
# Hardware evolution (NVIDIA)



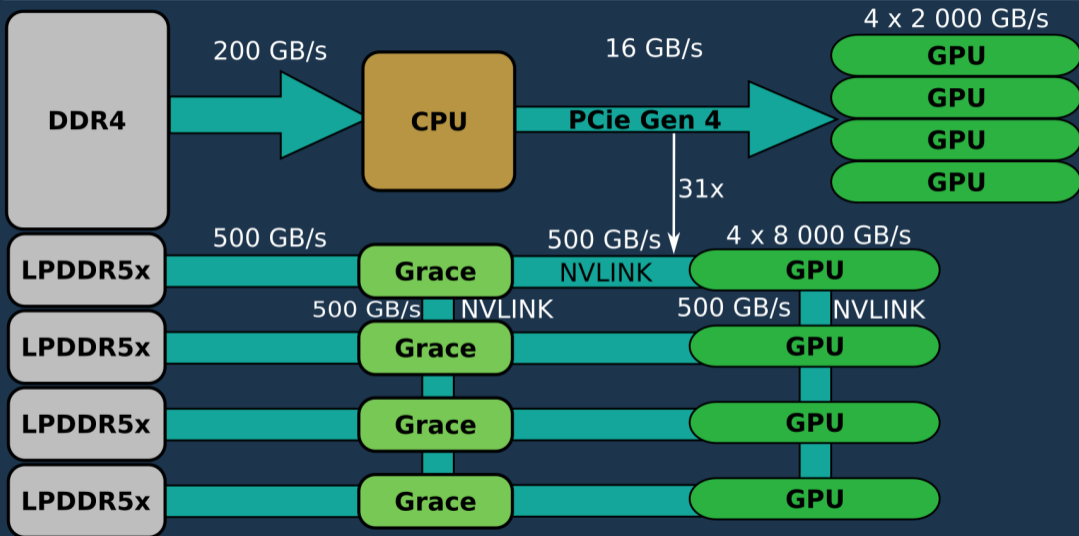
# Hardware evolution (NVIDIA)



# Hardware evolution (NVIDIA)

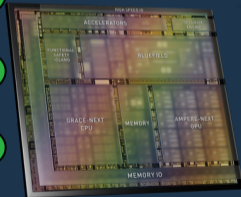
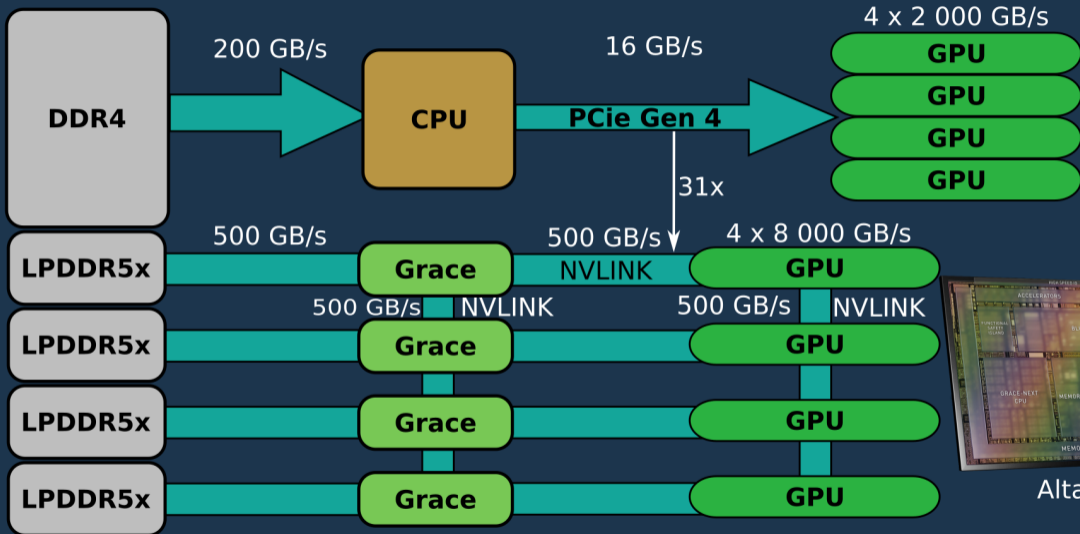


# Hardware evolution (Nvidia)



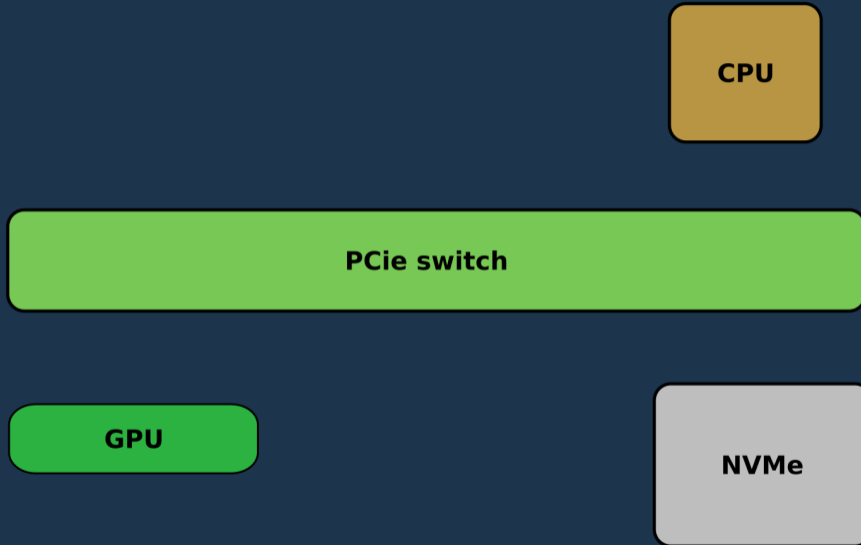


# Hardware evolution (Nvidia)

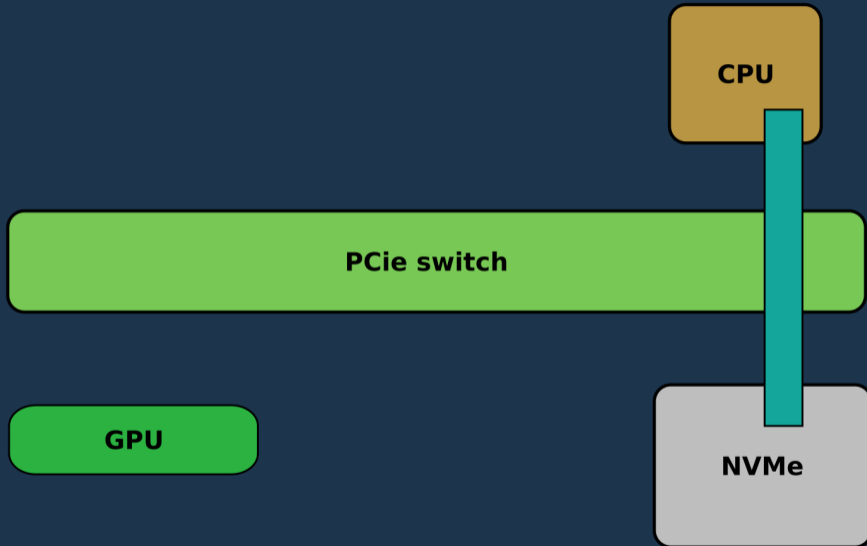


Altan

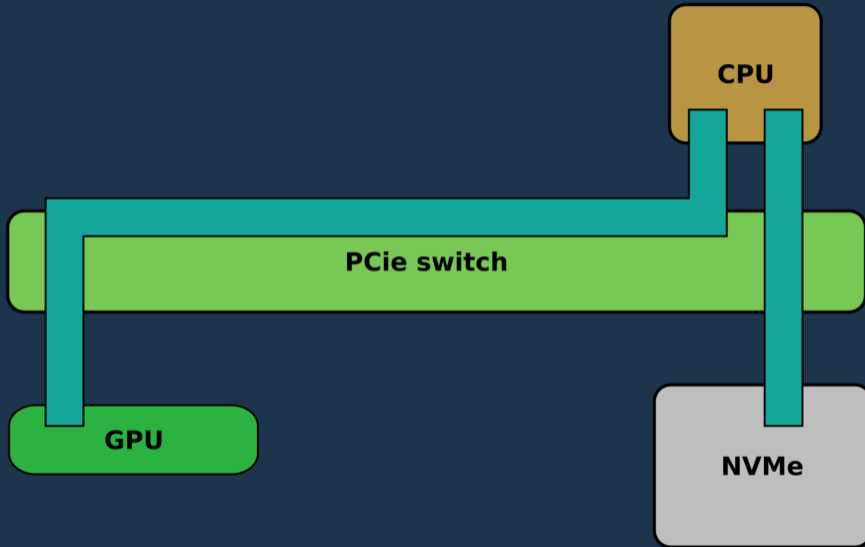
# Hardware evolution (NVIDIA)



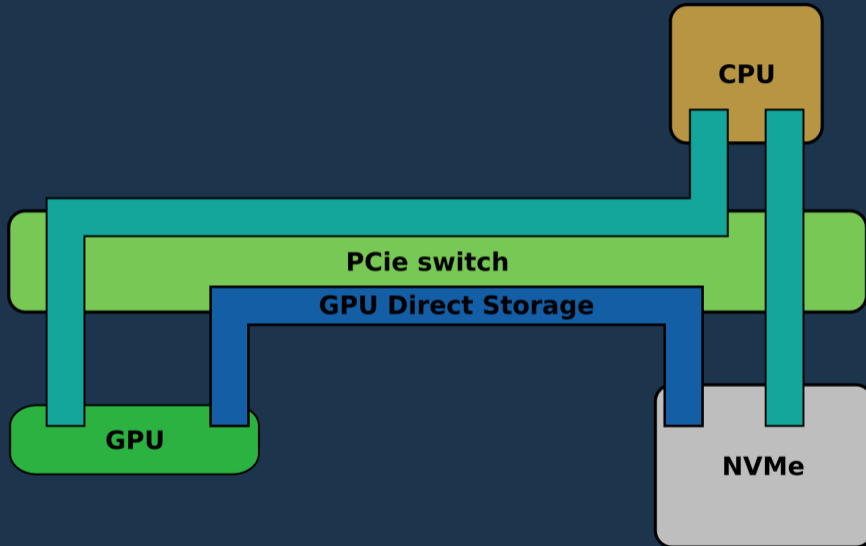
# Hardware evolution (NVIDIA)



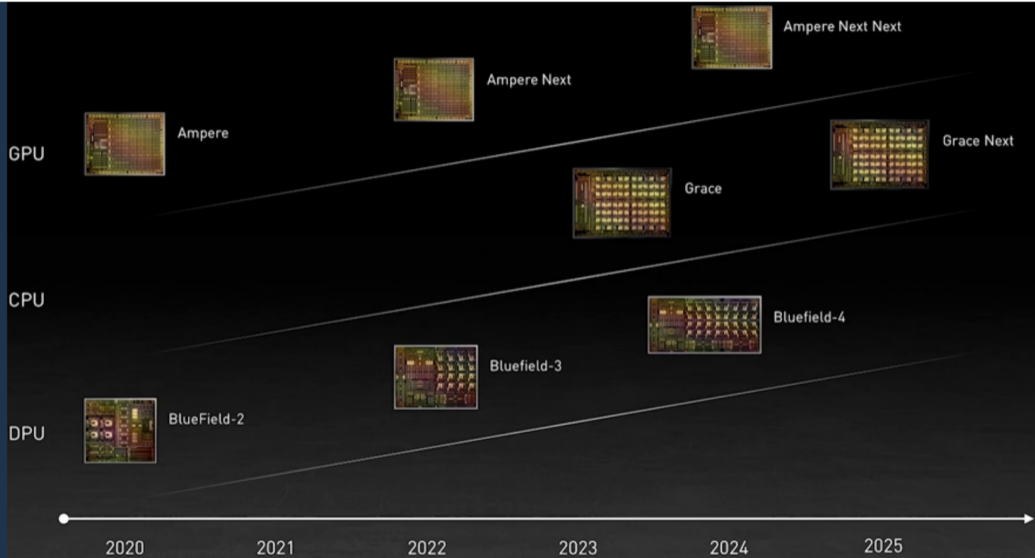
# Hardware evolution (NVIDIA)



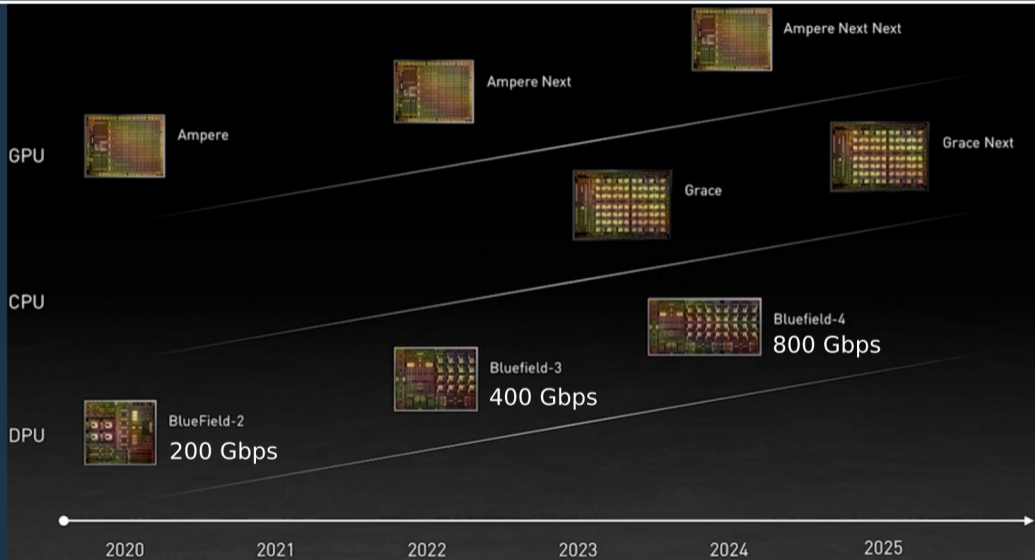
# Hardware evolution (NVIDIA)



# Hardware evolution (Nvidia)



# Hardware evolution (Nvidia)







Today

Today

nvc++

C++17 : **std::transform**  
Fortran : **do concurrent**

Today

nvc++

*-stdpart  
target GPU*

C++17 : **std::transform**  
Fortran : **do concurrent**

Today

nvc++

*-stdpart  
target GPU*

C++17 : **std::transform**

Fortran : **do concurrent**

**Compute Capabilities :**

- Ampere : 8
- Turing : 7
- Volta : 7
- Pascal : 6

Today

nvcc++

*-stdpart  
target GPU*

C++17 : **std::transform**  
Fortran : **do concurrent**

**Compute Capabilities :**

- Ampere : 8
- Turing : 7
- Volta : 7
- Pascal : 6

nvcc

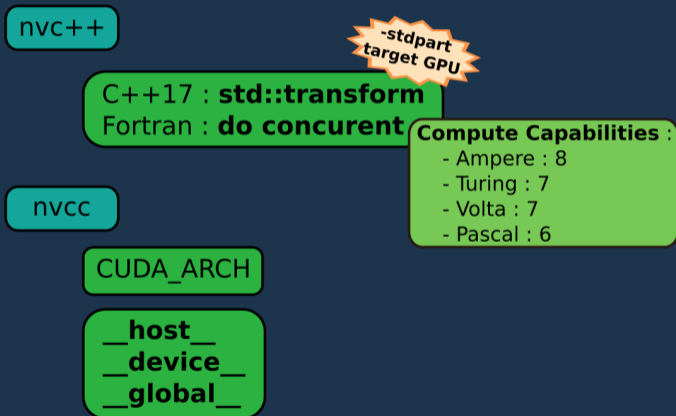
CUDA\_ARCH

**\_\_host\_\_  
\_\_device\_\_  
\_\_global\_\_**

# NVCC, NVC++ evolution (NVIDIA)

Today

Future (2022-2023)



# NVCC, NVC++ evolution (NVIDIA)

Today

Future (2022-2023)

nvc++

nvc++

*-stdpart  
target GPU*

C++17 : **std::transform**  
Fortran : **do concurrent**

**Compute Capabilities :**  
- Ampere : 8  
- Turing : 7  
- Volta : 7  
- Pascal : 6

nvcc

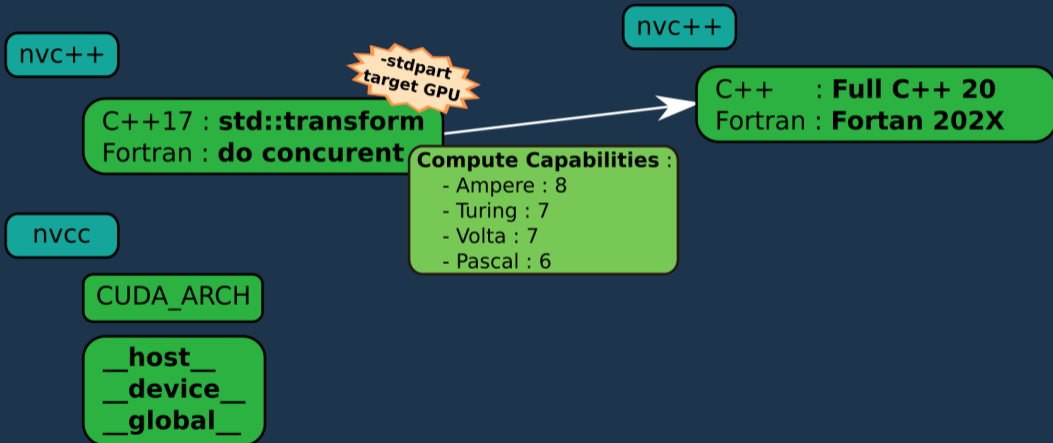
CUDA\_ARCH

\_\_host\_\_  
\_\_device\_\_  
\_\_global\_\_

# NVCC, NVC++ evolution (Nvidia)

Today

Future (2022-2023)





# NVCC, NVC++ evolution (NVIDIA)

Today

Future (2022-2023)

nvc++

C++17 : **std::transform**  
Fortran : **do concurrent**

**-stdpart  
target GPU**

**Compute Capabilities :**  
- Ampere : 8  
- Turing : 7  
- Volta : 7  
- Pascal : 6

nvcc

CUDA\_ARCH

\_\_host\_\_  
\_\_device\_\_  
\_\_global\_\_

nvc++

C++ : **Full C++ 20**  
Fortran : **Fortran 202X**

**Stencil :**  
- Basic  
- Customizable



# NVCC, NVC++ evolution (Nvidia)

Today

Future (2022-2023)

nvc++

C++17 : **std::transform**  
Fortran : **do concurrent**

**-stdpart  
target GPU**

**Compute Capabilities :**  
- Ampere : 8  
- Turing : 7  
- Volta : 7  
- Pascal : 6

nvcc

CUDA\_ARCH

**\_\_host\_\_  
\_\_device\_\_  
\_\_global\_\_**

nvc++

C++ : **Full C++ 20**  
Fortran : **Fortran 202X**

Automatic **Unified Memory**

**Stencil :**  
- Basic  
- Customizable

# NVCC, NVC++ evolution (Nvidia)

Today

Future (2022-2023)

nvc++

C++17 : **std::transform**  
Fortran : **do concurrent**

nvcc

CUDA\_ARCH

\_\_host\_\_  
\_\_device\_\_  
\_\_global\_\_

**-stdpart  
target GPU**

**Compute Capabilities :**  
- Ampere : 8  
- Turing : 7  
- Volta : 7  
- Pascal : 6

nvc++

C++ : **Full C++ 20**  
Fortran : **Fortran 202X**

Automatic **Unified Memory**

Automatic **Host/Device**  
determination

**Stencil :**  
- Basic  
- Customizable

# NVCC, NVC++ evolution (Nvidia)

Today

Future (2022-2023)

nvc++

C++17 : **std::transform**  
Fortran : **do concurrent**

nvcc

CUDA\_ARCH

\_\_host\_\_  
\_\_device\_\_  
\_\_global\_\_

**-stdpart  
target GPU**

**Compute Capabilities :**  
- Ampere : 8  
- Turing : 7  
- Volta : 7  
- Pascal : 6

nvc++

C++ : **Full C++ 20**  
Fortran : **Fortran 202X**

Automatic **Unified Memory**

Automatic **Host/Device**  
determination

**Stencil :**  
- Basic  
- Customizable

**No more :**  
\_\_host\_\_  
\_\_device\_\_

# NVCC, NVC++ evolution (Nvidia)

Today

Future (2022-2023)

nvc++

C++17 : **std::transform**  
Fortran : **do concurrent**

**-stdpart  
target GPU**

**Compute Capabilities :**  
- Ampere : 8  
- Turing : 7  
- Volta : 7  
- Pascal : 6

nvcc

CUDA\_ARCH

**\_\_host\_\_  
\_\_device\_\_  
\_\_global\_\_**

nvc++

C++ : **Full C++ 20**  
Fortran : **Fortran 202X**

Automatic **Unified Memory**

Automatic **Host/Device**  
determination

**if target(description)**

**Stencil :**  
- Basic  
- Customizable

**No more :**  
\_\_host\_\_  
\_\_device\_\_



# NVCC, NVC++ evolution (Nvidia)

Today

Future (2022-2023)

nvc++

C++17 : **std::transform**  
Fortran : **do concurrent**

**-stdpart  
target GPU**

**Compute Capabilities :**  
- Ampere : 8  
- Turing : 7  
- Volta : 7  
- Pascal : 6

nvcc

CUDA\_ARCH

**\_\_host\_\_  
\_\_device\_\_  
\_\_global\_\_**

nvc++

C++ : **Full C++ 20**  
Fortran : **Fortran 202X**

**Stencil :**  
- Basic  
- Customizable

Automatic **Unified Memory**

Automatic **Host/Device**  
determination

**No more :  
\_\_host\_\_  
\_\_device\_\_**

**if target**(description)

Still (for specific optimisations) :  
**\_\_host\_\_  
\_\_device\_\_  
\_\_global\_\_**



```

d = 2.5 * ceil(transpose(a)) + 3.0 * abs(transpose(b))
d = 2.5 * ceil(transpose(a)) + 3.0 * abs(b)
d = reshape(a, shape=[ni, nj, nk])
d = reshape(a, shape=[ni, nk, nj])
d = 2.5 * sqrt(reshape(a, shape=[ni, nk, nj], order=[1, 3, 2]))
d = alpha * conjg(reshape(a, shape=[ni, nk, nj], order=[1, 3, 2]))
d = reshape(a, shape=[ni, nk, nj], order=[1, 3, 2])
d = reshape(a, shape=[nk, ni, nj], order=[2, 3, 1])
d = reshape(a, shape=[ni*nj, nk])
d = reshape(a, shape=[nk, ni*nj], order=[2, 1])
d = reshape(a, shape=[64, 2, 16, 16, 64], order=[5, 2, 3, 4, 1])
d = abs(reshape(a, shape=[64, 2, 16, 16, 64], order=[5, 2, 3, 4, 1]))
c = matmul(a, b)
c = matmul(transpose(a), b)
c = matmul(reshape(a, shape=[m, k], order=[2, 1]), b)
c = matmul(a, transpose(b))
c = matmul(a, reshape(b, shape=[k, n], order=[2, 1]))

```

```

c = matmul(transpose(a), transpose(b))
c = matmul(transpose(a), reshape(b, shape=[k, n], order=[2, 1]))
d = spread(a, dim=3, ncopies=nk)
d = spread(a, dim=1, ncopies=ni)
d = spread(a, dim=2, ncopies=nx)
d = alpha * abs(spread(a, dim=2, ncopies=nx))
d = alpha * spread(a, dim=2, ncopies=nx)
d = abs(spread(a, dim=2, ncopies=nx))
d = transpose(a)
d = alpha * transpose(a)
d = alpha * ceil(transpose(a))
d = alpha * conjg(transpose(a))
c = c + matmul(a, b)
c = c - matmul(a, b)
c = c + alpha * matmul(a, b)
d = alpha * matmul(a, b) + c
d = alpha * matmul(a, b) + beta * c

```

**Current** default situation

**New** default situation



**Current** default situation

**Scalar**



**New** default situation

**Current** default situation

**Scalar**



**Single-threaded**



**New** default situation

**Current** default situation

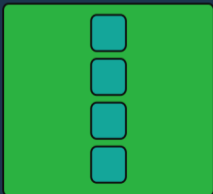
**Scalar**



**Single-threaded**



Execution **in order**



**New** default situation

# Paradigms evolution

**Current** default situation

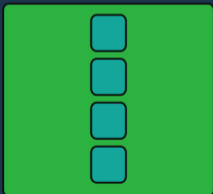
**Scalar**



**Single-threaded**



Execution **in order**



Recent machine :

- 64 cores
  - 16 SIMD units
- uses < 0.1% of computing power

**New** default situation

# Paradigms evolution

**Current** default situation

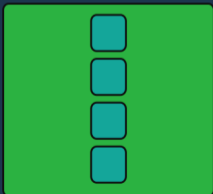
**Scalar**



**Single-threaded**



Execution **in order**



Recent machine :  
- 64 cores  
- 16 SIMD units  
uses < 0.1% of  
computing power

**New** default situation

**Vectorized**



# Paradigms evolution

**Current** default situation

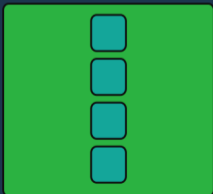
**Scalar**



**Single-threaded**



Execution **in order**



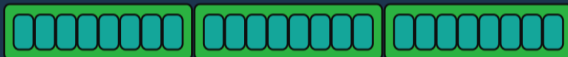
Recent machine :  
- 64 cores  
- 16 SIMD units  
uses < 0.1% of  
computing power

**New** default situation

**Vectorized**



**Multi-threaded**



# Paradigms evolution

**Current** default situation

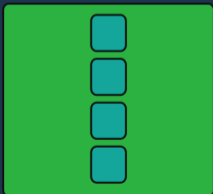
**Scalar**



**Single-threaded**



Execution **in order**



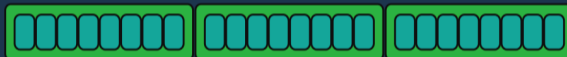
Recent machine :  
- 64 cores  
- 16 SIMD units  
uses < 0.1% of  
computing power

**New** default situation

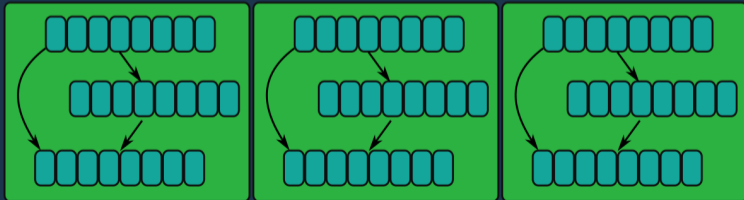
**Vectorized**



**Multi-threaded**



Execution **out of order**



# Paradigms evolution

NVC++  
2022-2023

**Current** default situation

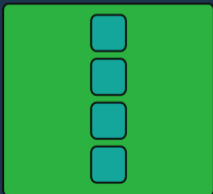
**Scalar**



**Single-threaded**



Execution **in order**



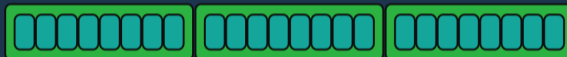
Recent machine :  
- 64 cores  
- 16 SIMD units  
uses < 0.1% of  
computing power

**New** default situation

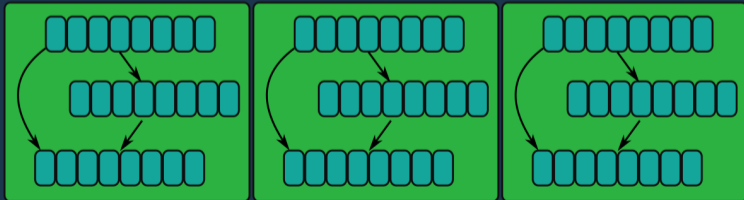
**Vectorized**



**Multi-threaded**



Execution **out of order**



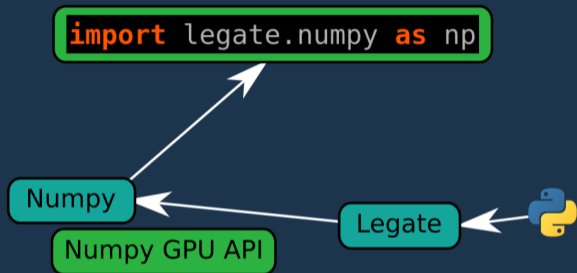


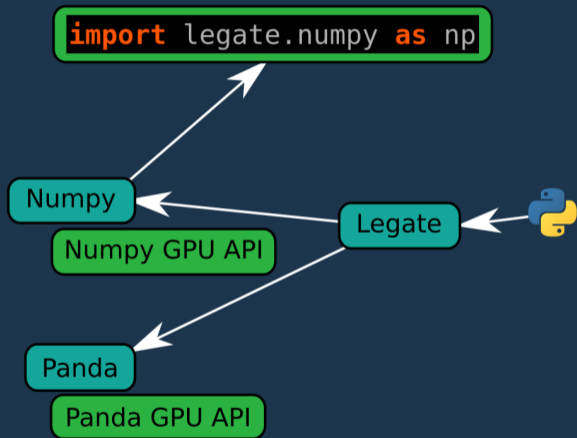


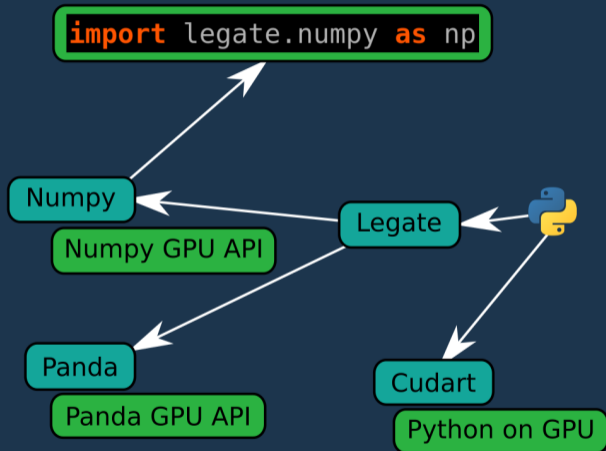




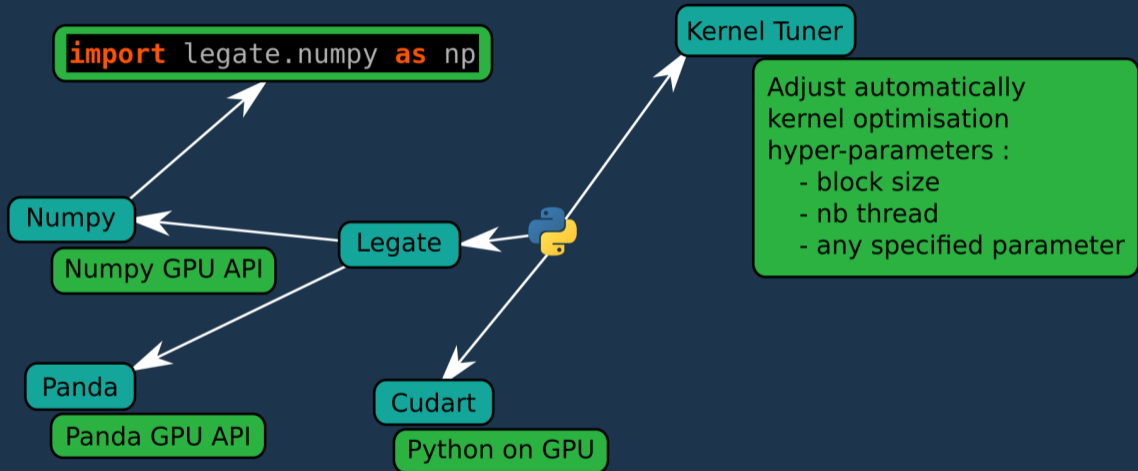






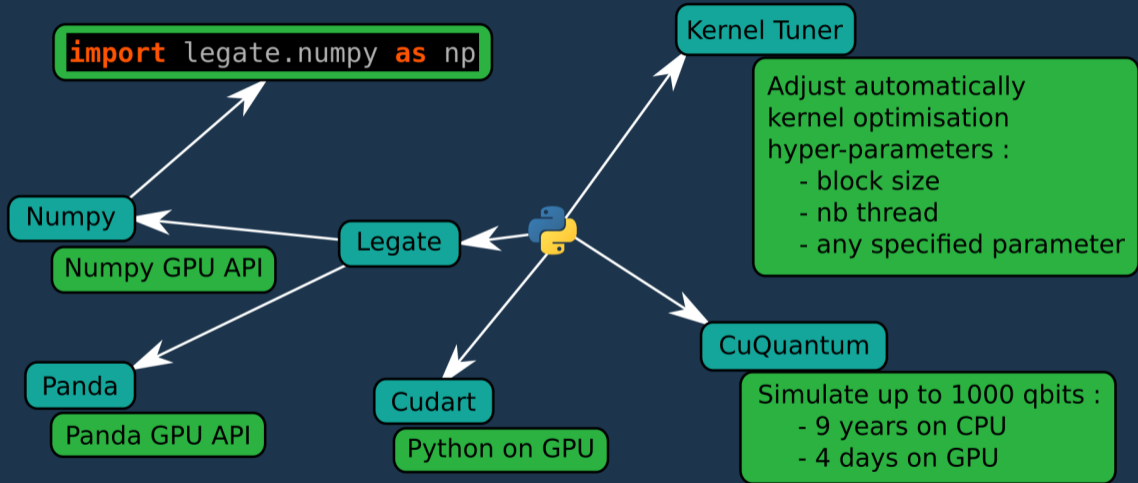


# Python (ecosystem) evolution (NVIDIA)



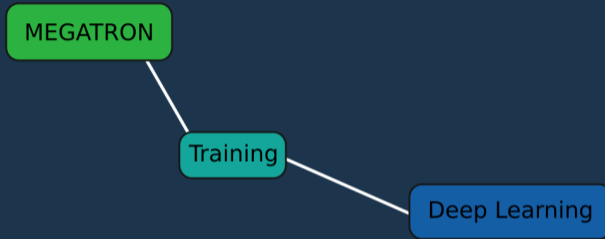


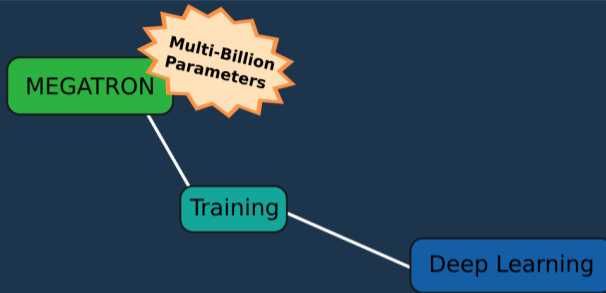
# Python (ecosystem) evolution (NVIDIA)

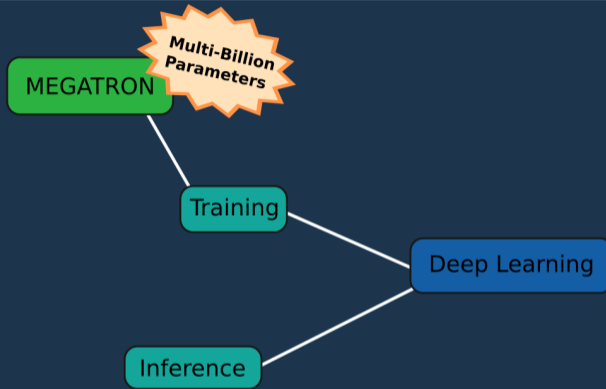


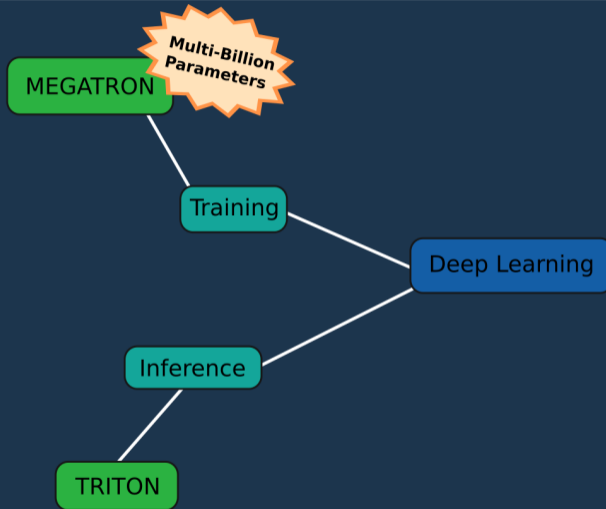
Deep Learning



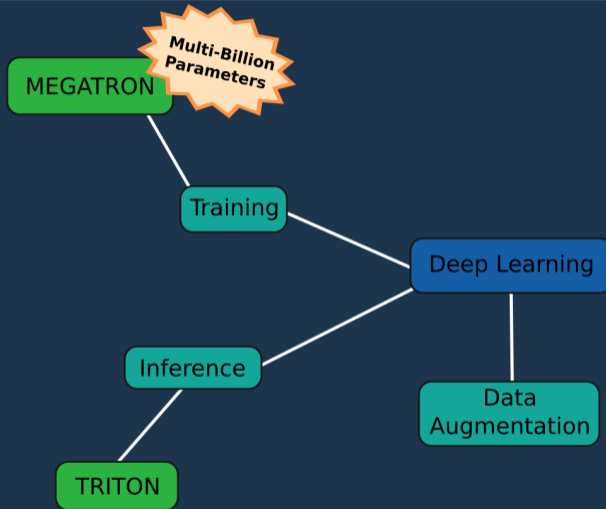






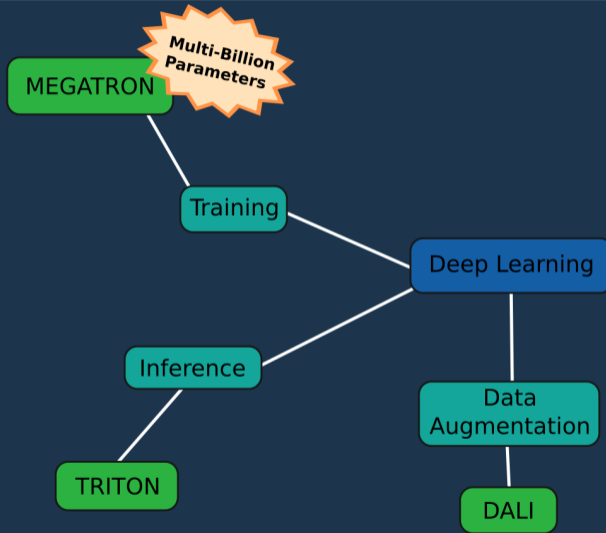


# Deep Learning evolution

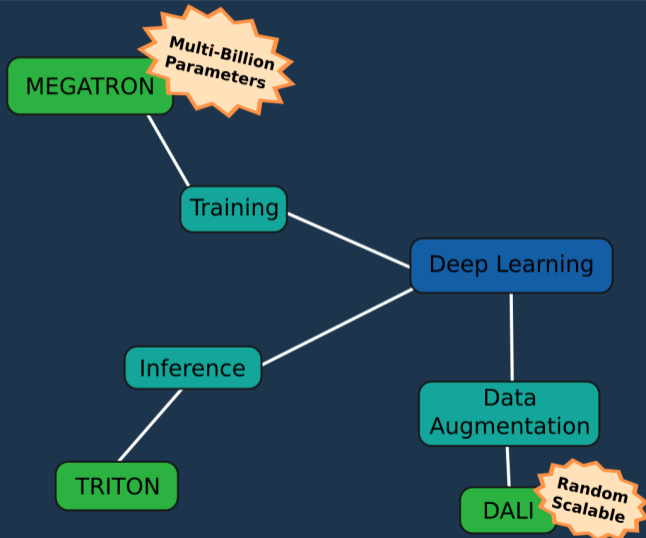




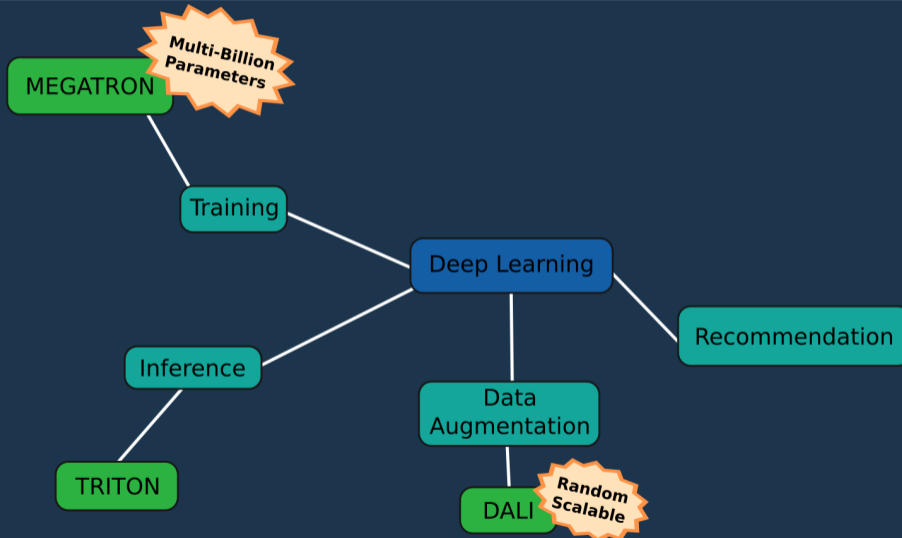
# Deep Learning evolution



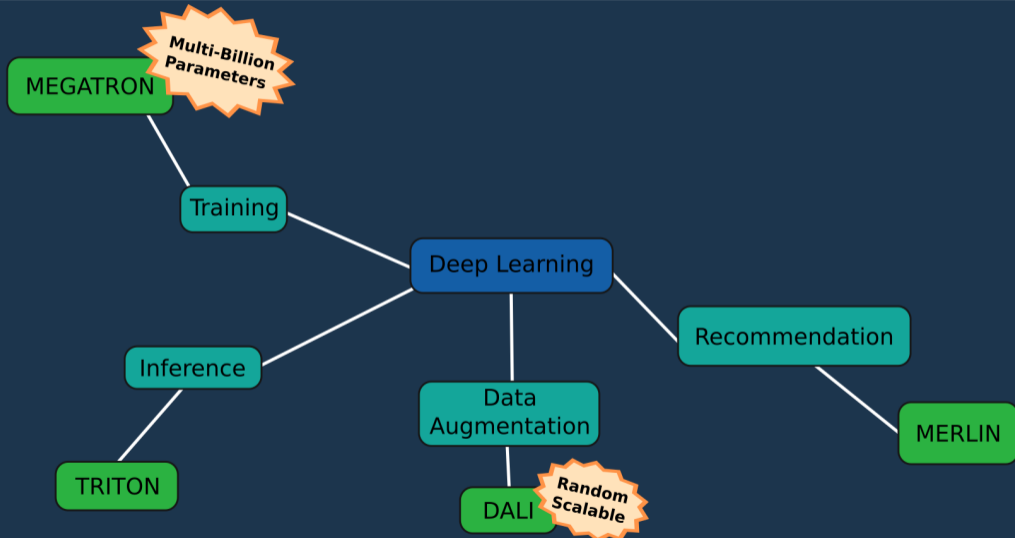
# Deep Learning evolution



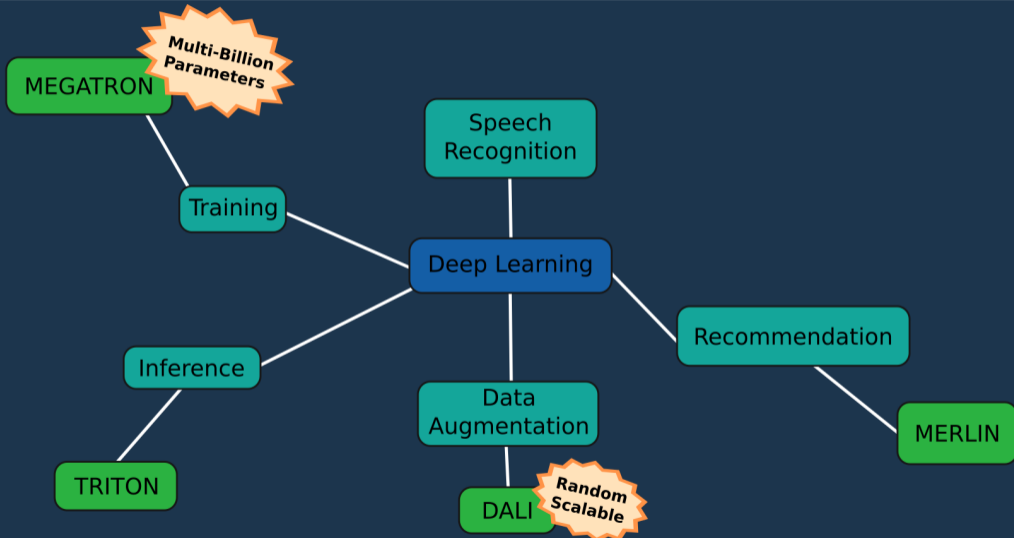
# Deep Learning evolution



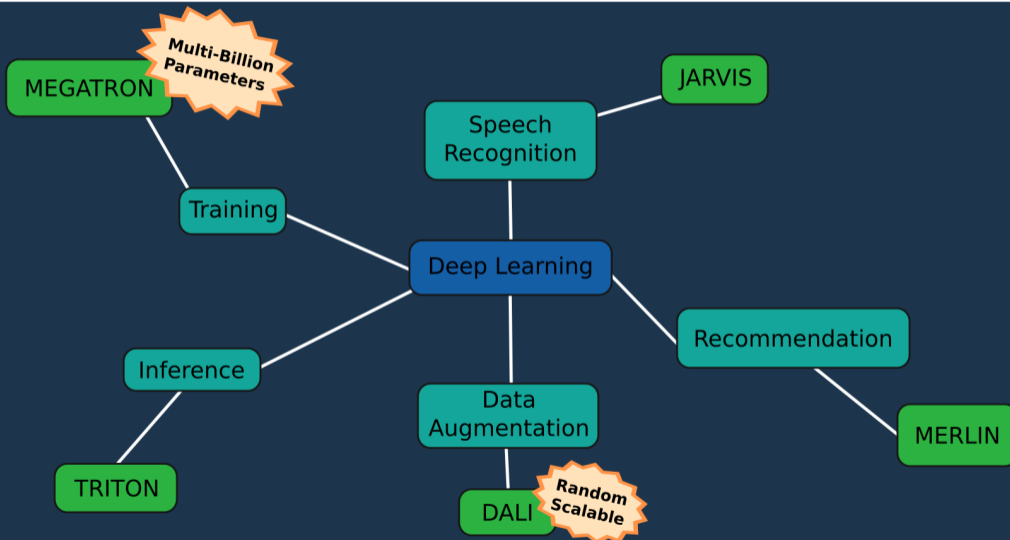
# Deep Learning evolution



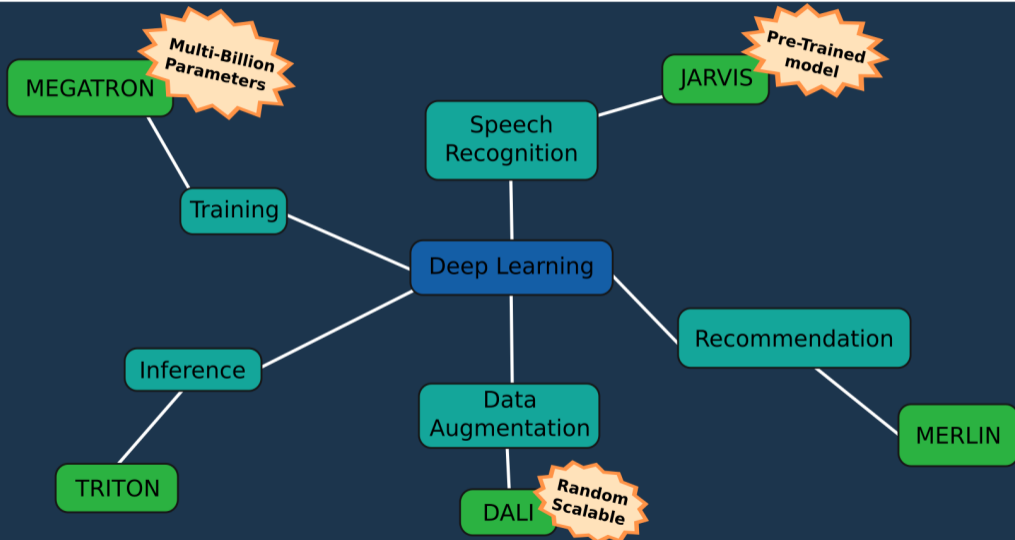
# Deep Learning evolution



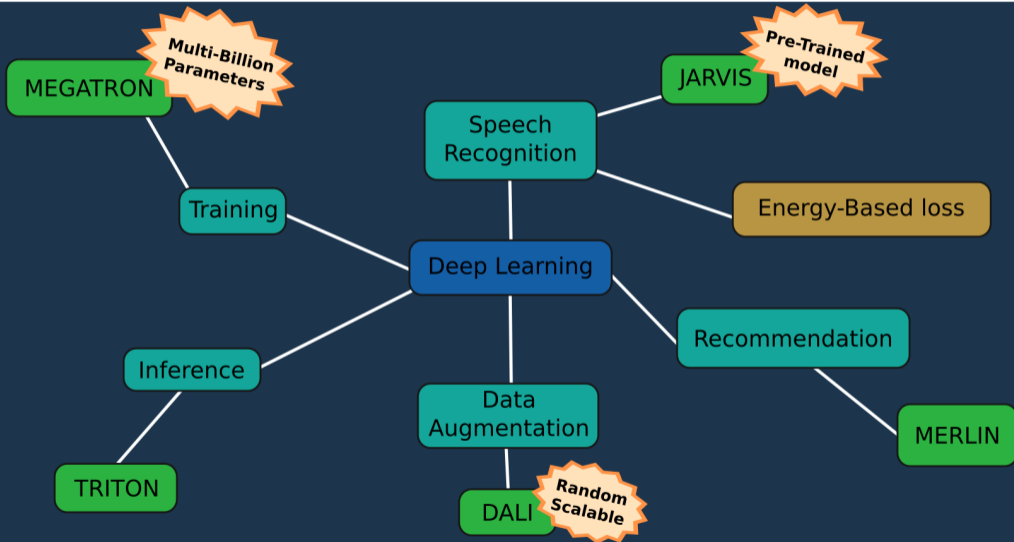
# Deep Learning evolution



# Deep Learning evolution

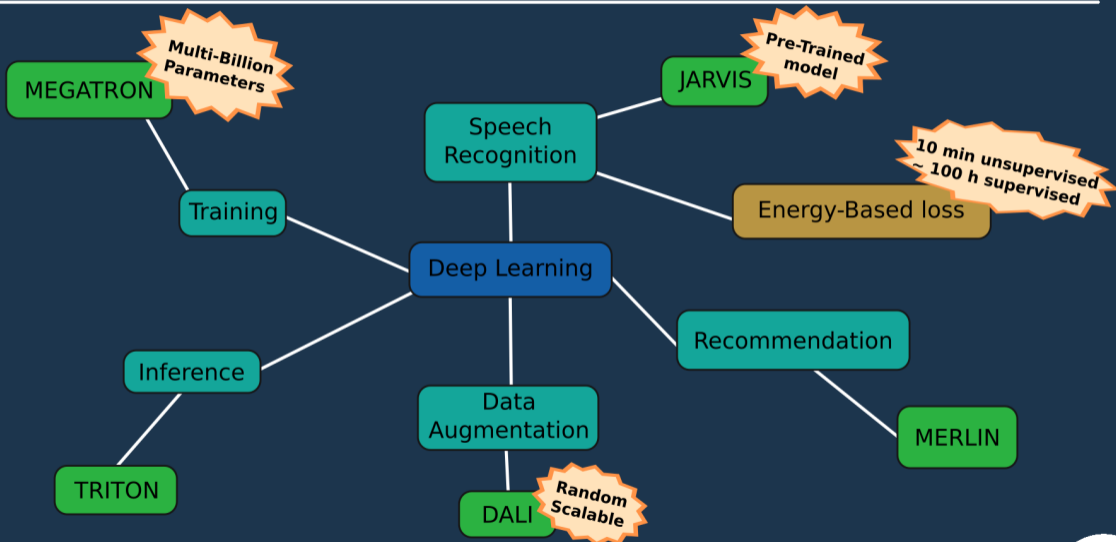


# Deep Learning evolution

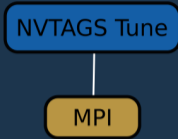


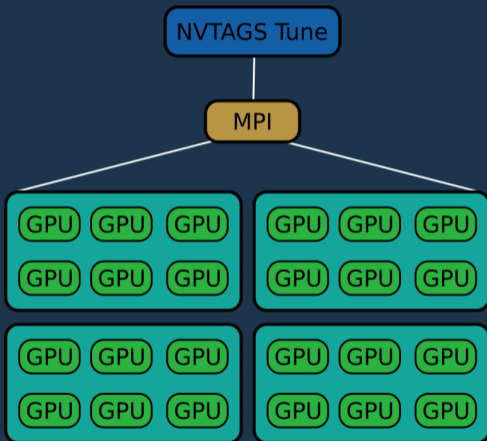


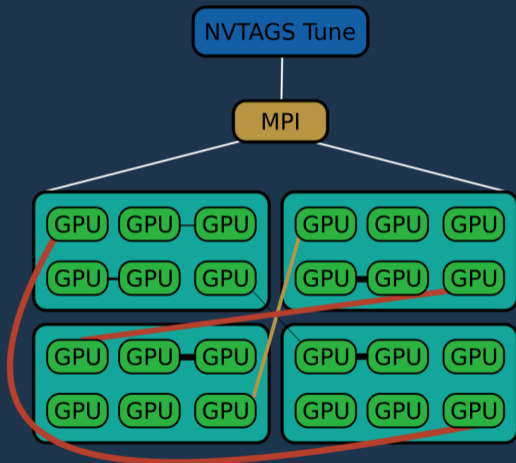
# Deep Learning evolution

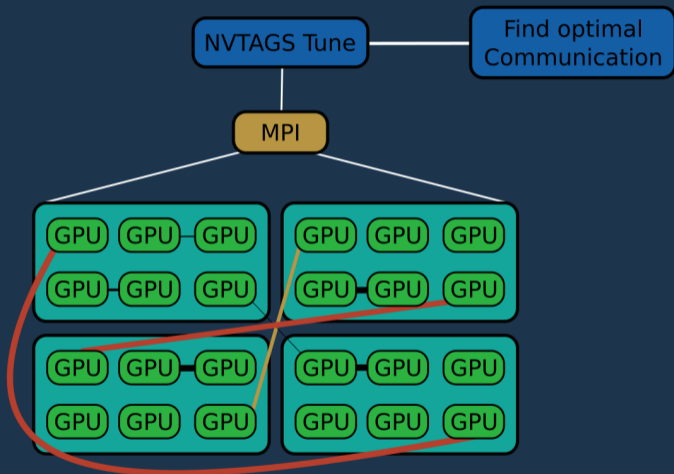


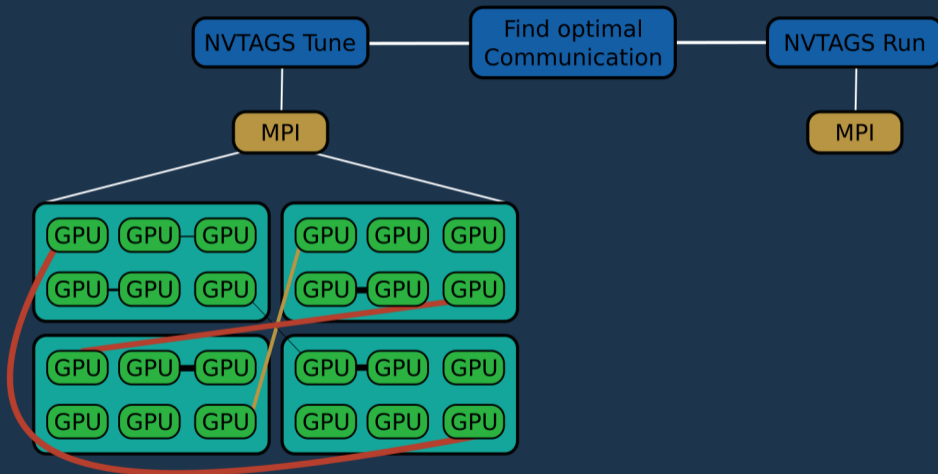


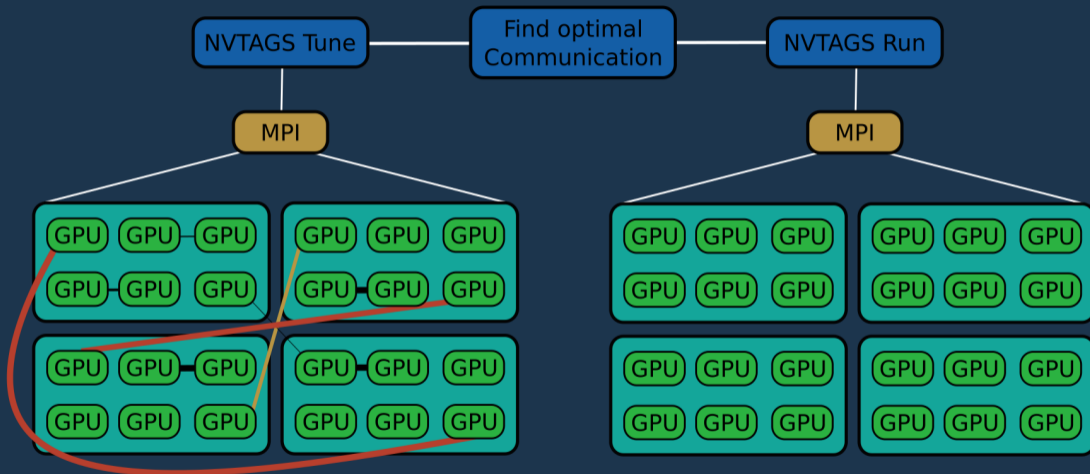




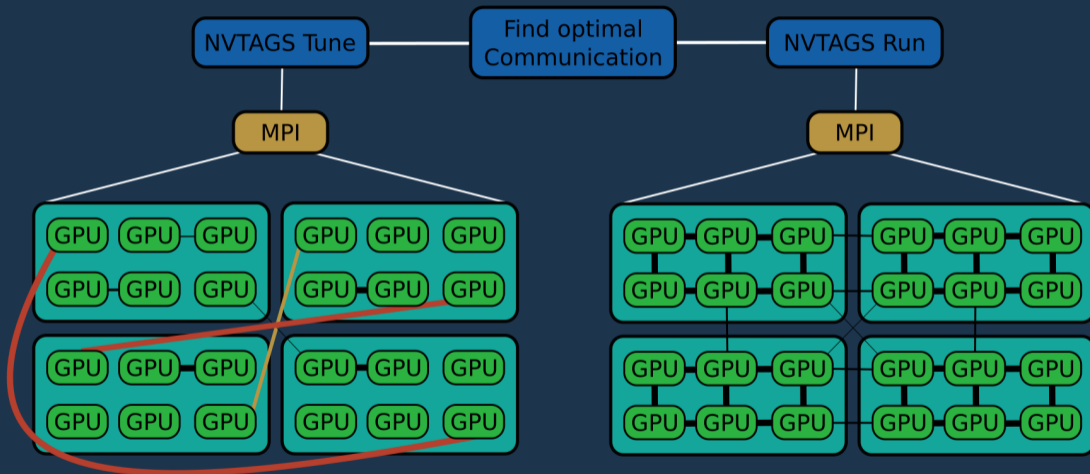


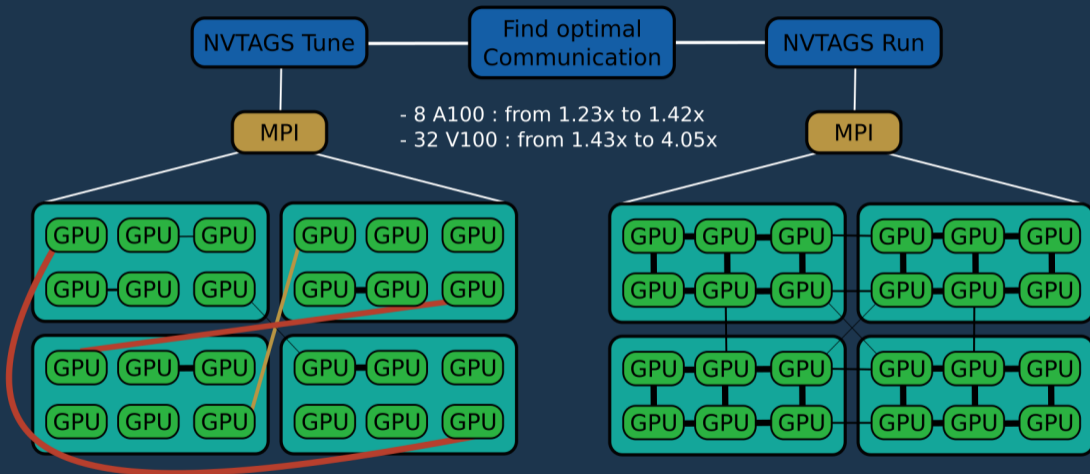


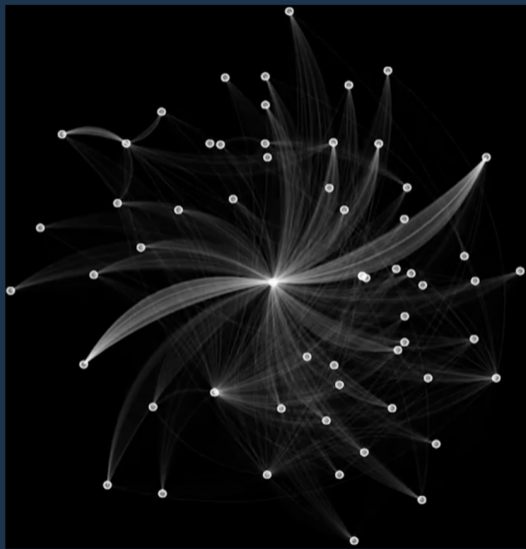




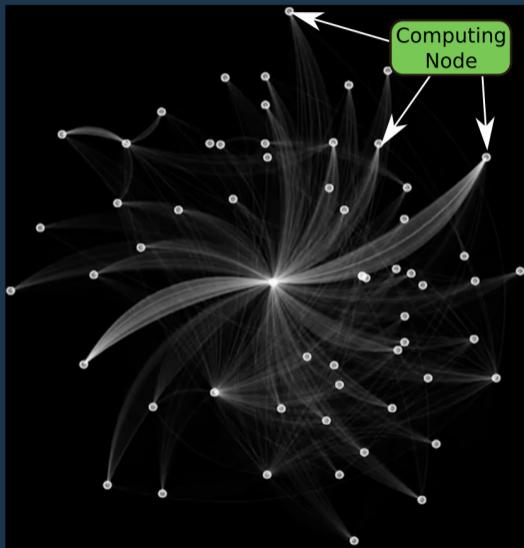




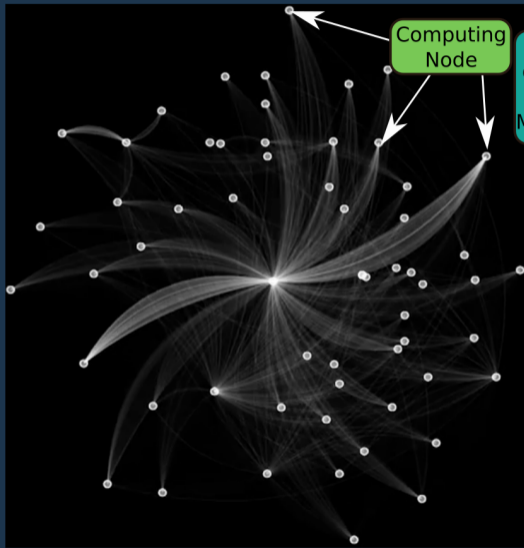




# Morpheus (NVIDIA)

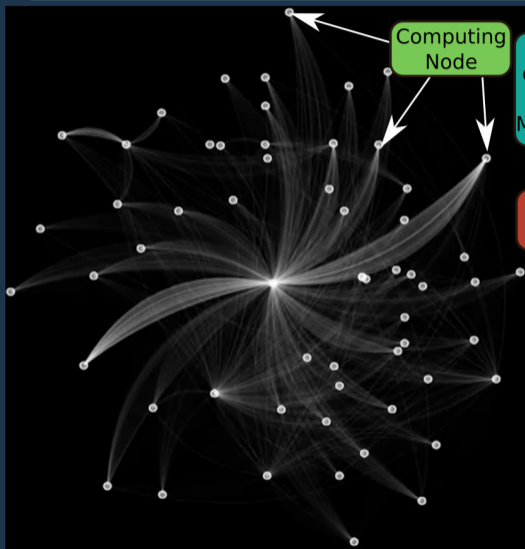


# Morpheus (Nvidia)



Deep Analysis  
of all exchanged  
packets with  
Machine-Learning

# Morpheus (NVIDIA)

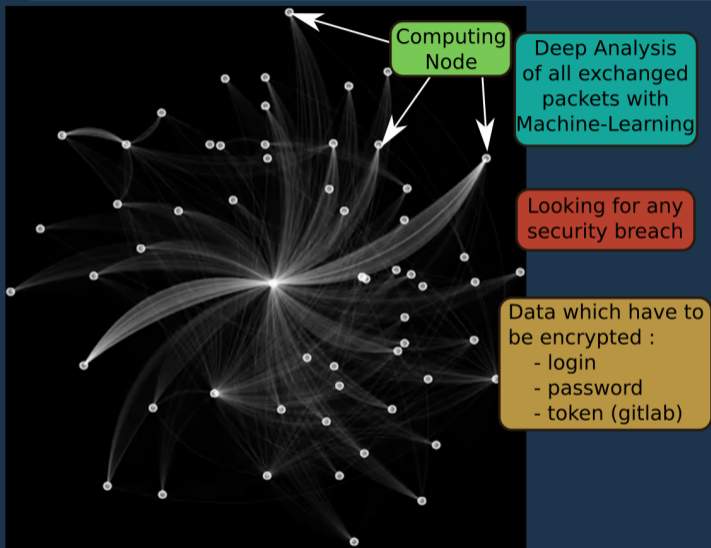


Computing Node

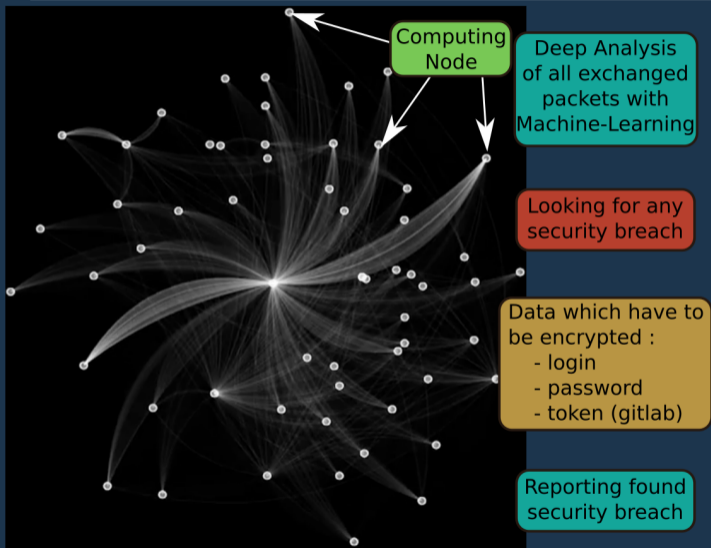
Deep Analysis of all exchanged packets with Machine-Learning

Looking for any security breach

# Morpheus (Nvidia)

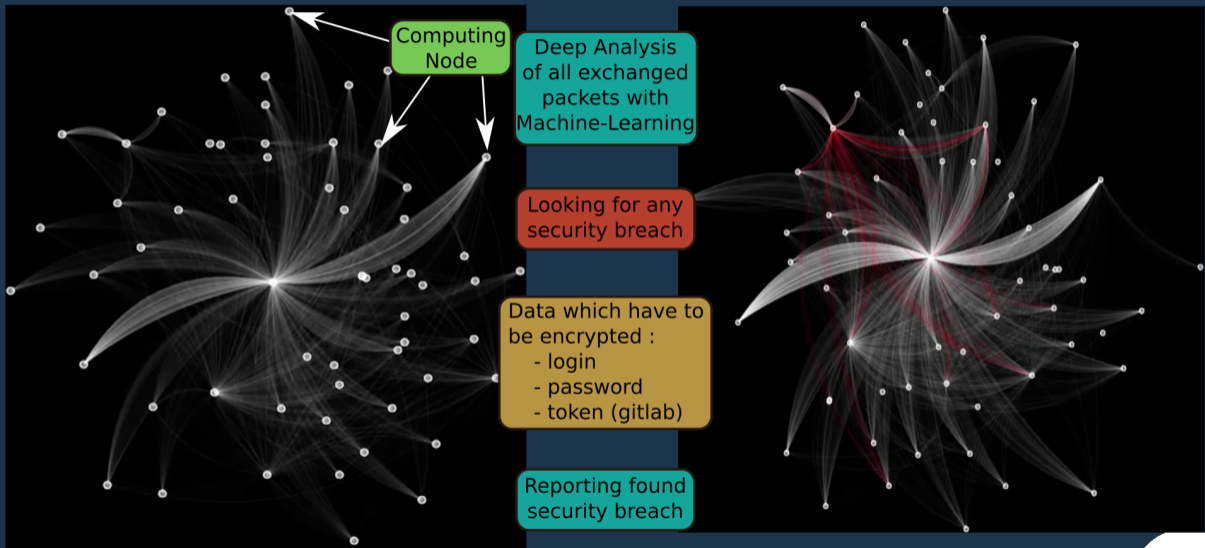


# Morpheus (Nvidia)

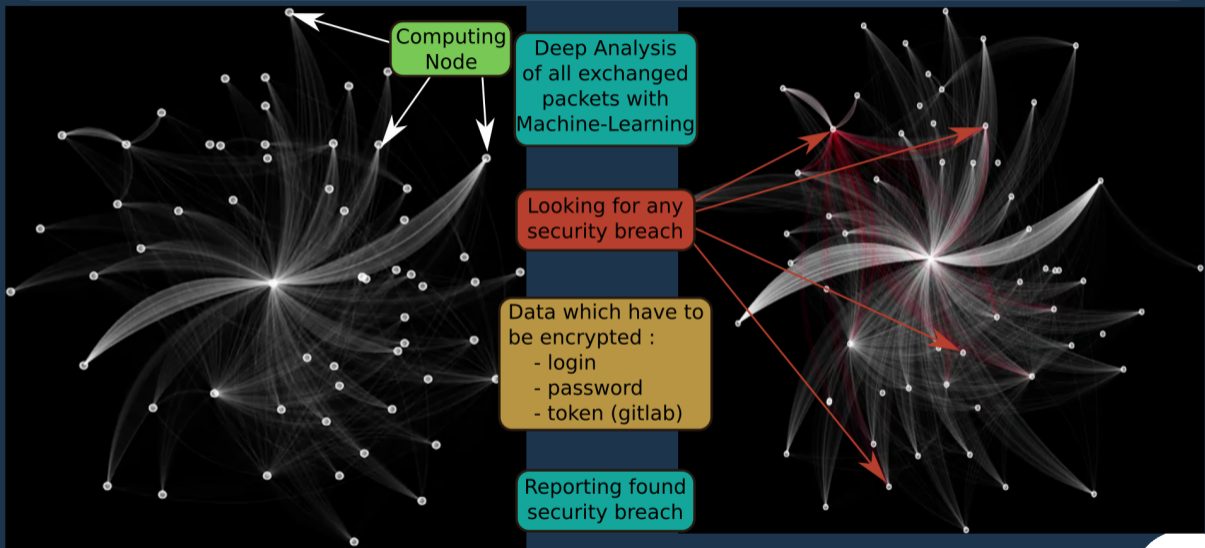




# Morpheus (Nvidia)



# Morpheus (Nvidia)



# Energy consumption (ATOS)

2010 : 0.23 GFlops/W

2020 : 25 GFlops/W (Juelich supercomputer - 3 744 A100), free cooling

2010 : 0.23 GFlops/W

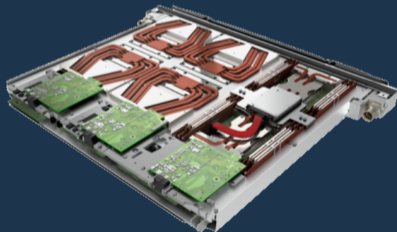
2020 : 25 GFlops/W (Juelich supercomputer - 3 744 A100), free cooling

- Classic Data Center : 100 % computing, 100 % cooling
- Water cooling Data Center : 100 % computing, 50 % cooling
- Direct liquid cooling Data Center : 100 % computing, 5 % cooling

2010 : 0.23 GFlops/W

2020 : 25 GFlops/W (Juelich supercomputer - 3 744 A100), free cooling

- Classic Data Center : 100 % computing, 100 % cooling
- Water cooling Data Center : 100 % computing, 50 % cooling
- Direct liquid cooling Data Center : 100 % computing, 5 % cooling



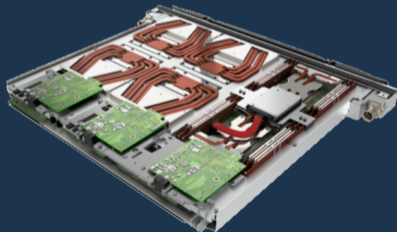
Bull-SeqwanaX1125

# Energy consumption (ATOS)

2010 : 0.23 GFlops/W

2020 : 25 GFlops/W (Juelich supercomputer - 3 744 A100), free cooling

- Classic Data Center : 100 % computing, 100 % cooling
- Water cooling Data Center : 100 % computing, 50 % cooling
- Direct liquid cooling Data Center : 100 % computing, 5 % cooling



Bull-SeqwanaX1125



Everything is cooled :

- CPU
- GPU
- switches

# Future of HPC

Programming  
Languages

Software

Hardware

Machine  
Learning

# Future of HPC

Programming  
Languages

Software

Complexify

Hardware

Machine  
Learning



# Future of HPC

Programming  
Languages

Software

Complexify

Hardware

Can increase  
gap between  
peak and measured  
performances

Machine  
Learning

# Future of HPC

Programming  
Languages

Software

Reduce gap between  
peak and measured  
performance

Can increase  
gap between  
peak and measured  
performances

Complexify

Hardware

Machine  
Learning

# Future of HPC

Programming  
Languages

Software

Code  
Generation

Reduce gap between  
peak and measured  
performance

Can increase  
gap between  
peak and measured  
performances

Complexify

Hardware

Machine  
Learning

# Future of HPC

Programming  
Languages

Software

Code  
Generation

- Handle hardwares
- Aggressive optimisation
- Fasten development

Reduce gap between  
peak and measured  
performance

Complexify

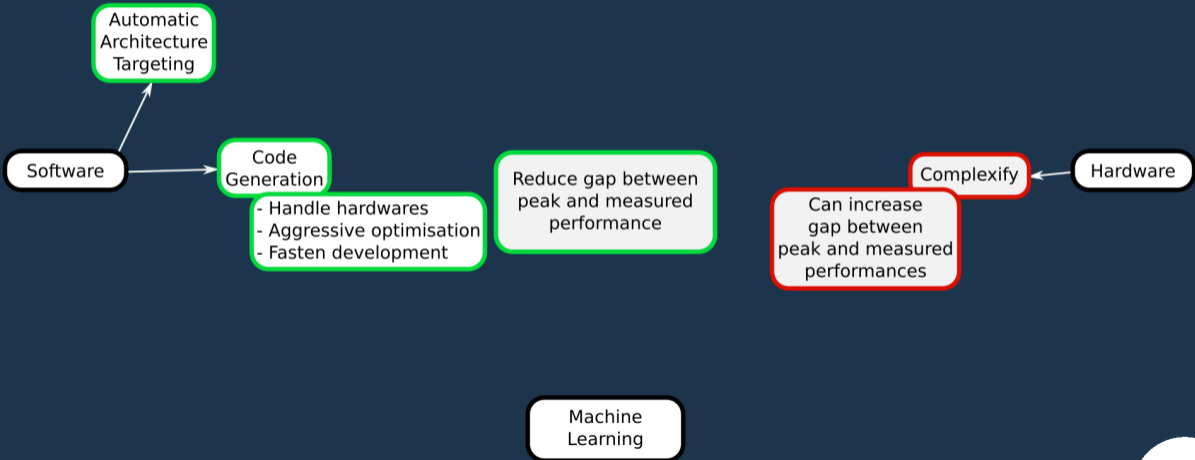
Can increase  
gap between  
peak and measured  
performances

Hardware

Machine  
Learning

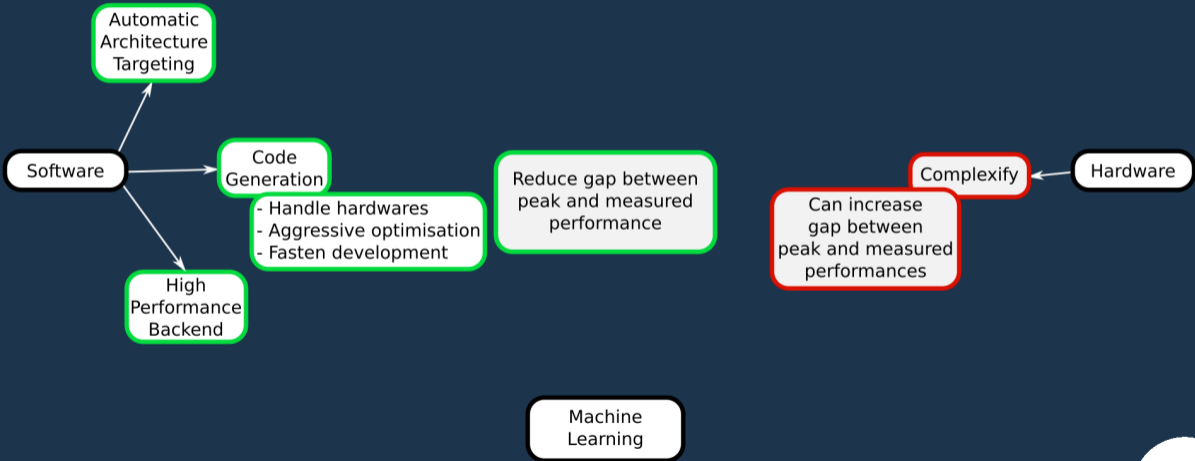
# Future of HPC

Programming Languages



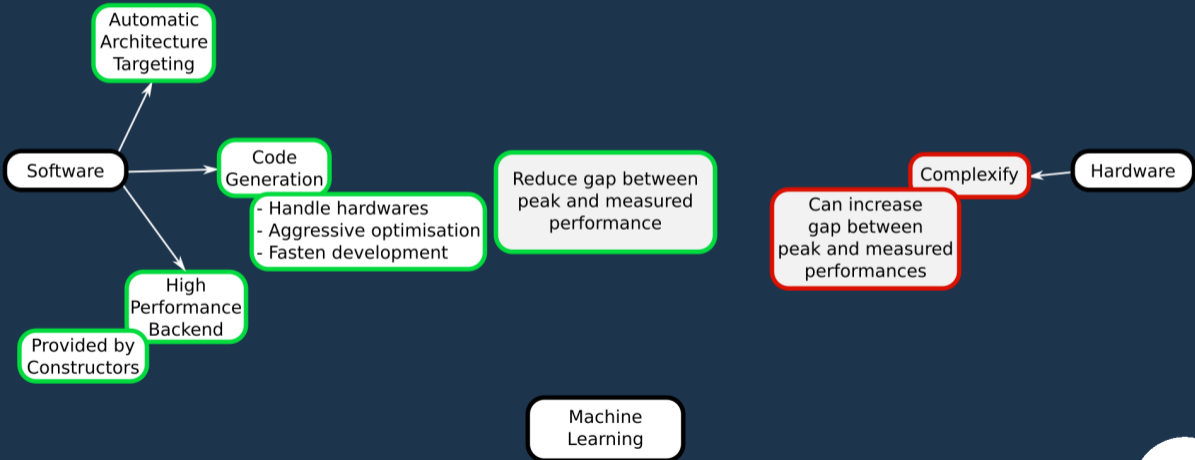
# Future of HPC

Programming Languages



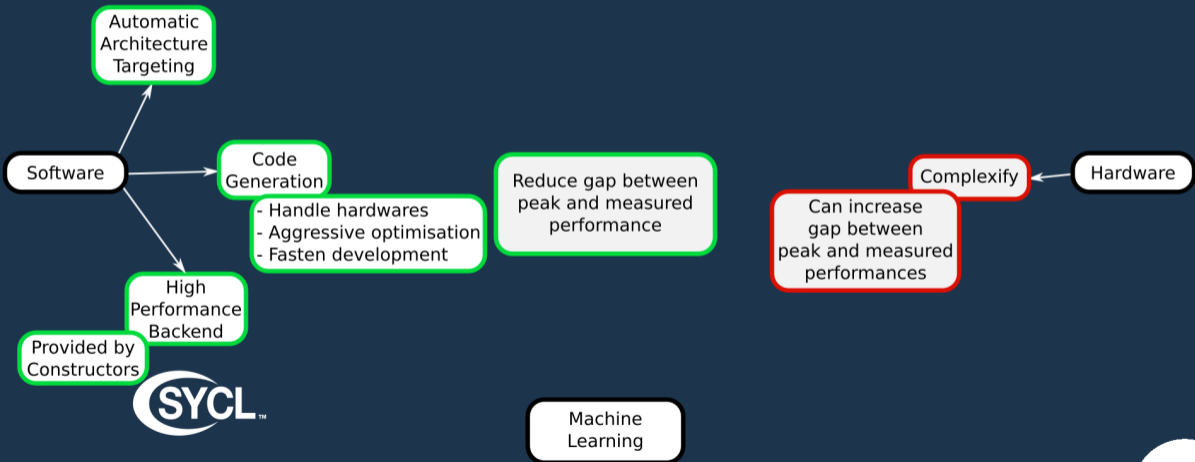
# Future of HPC

Programming Languages



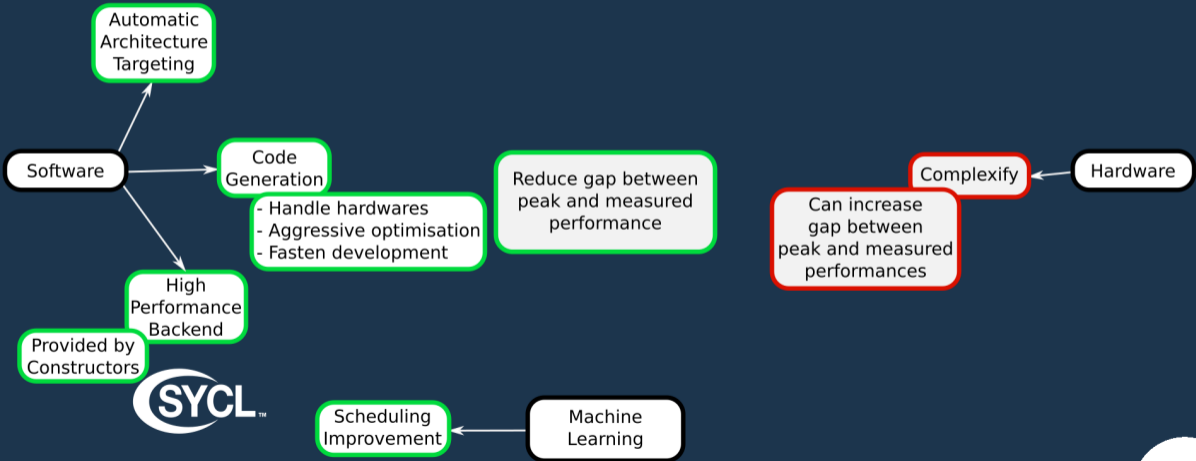
# Future of HPC

Programming Languages

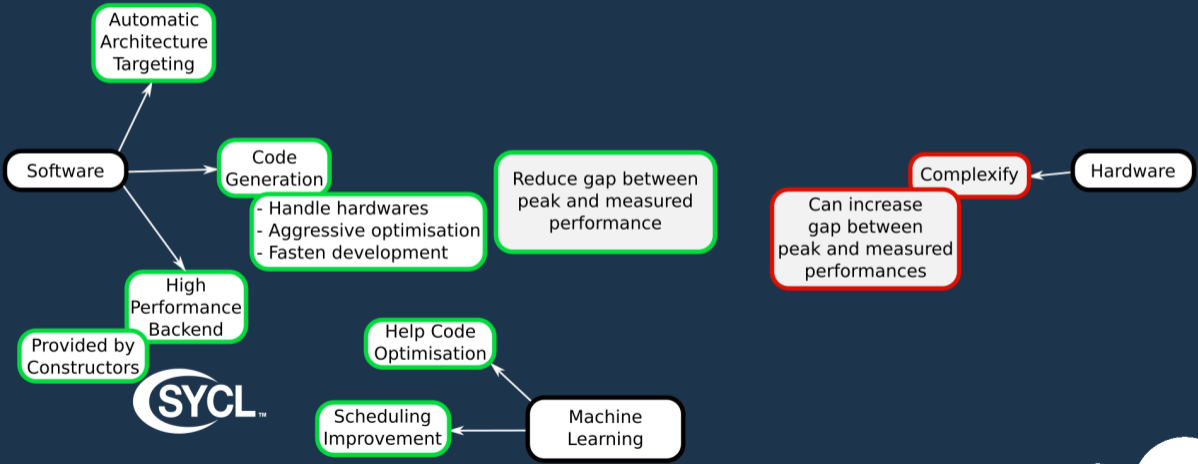




Programming Languages

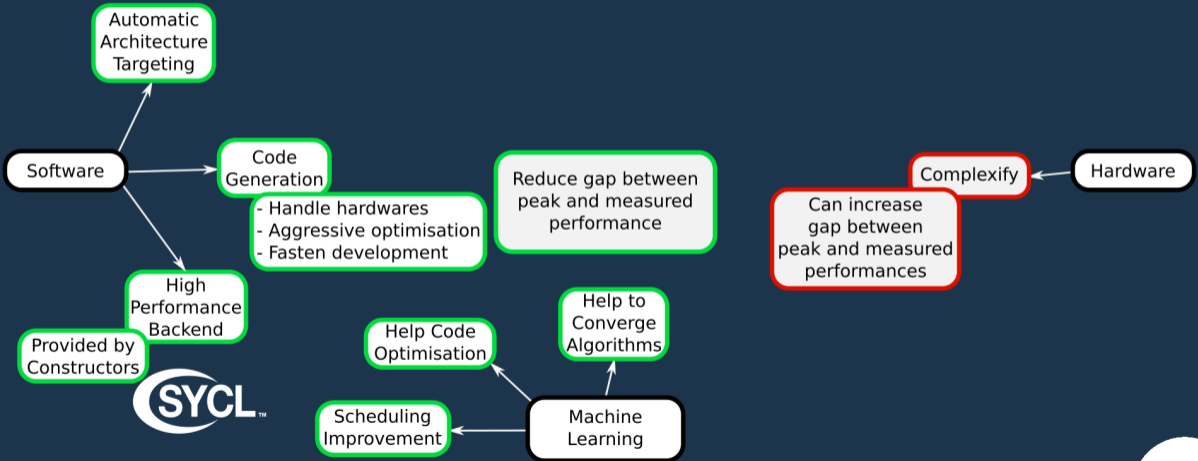


Programming Languages

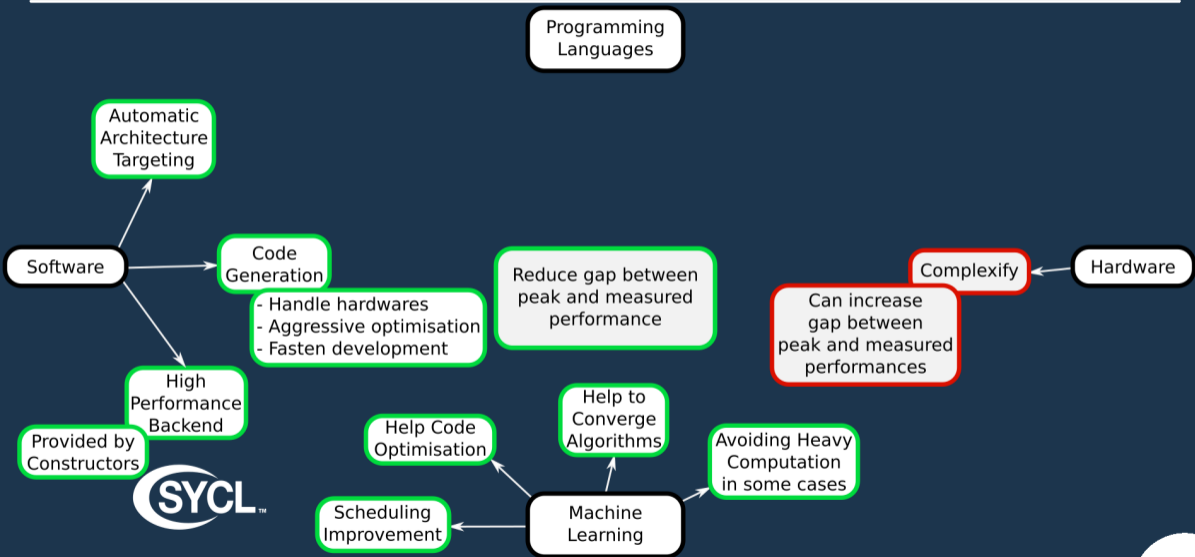


# Future of HPC

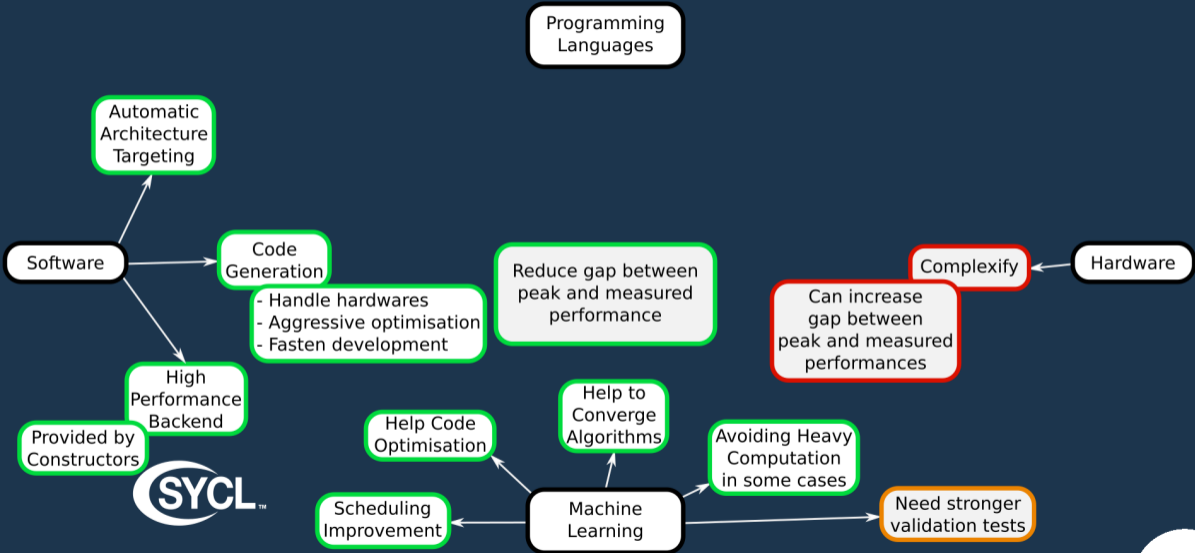
Programming Languages



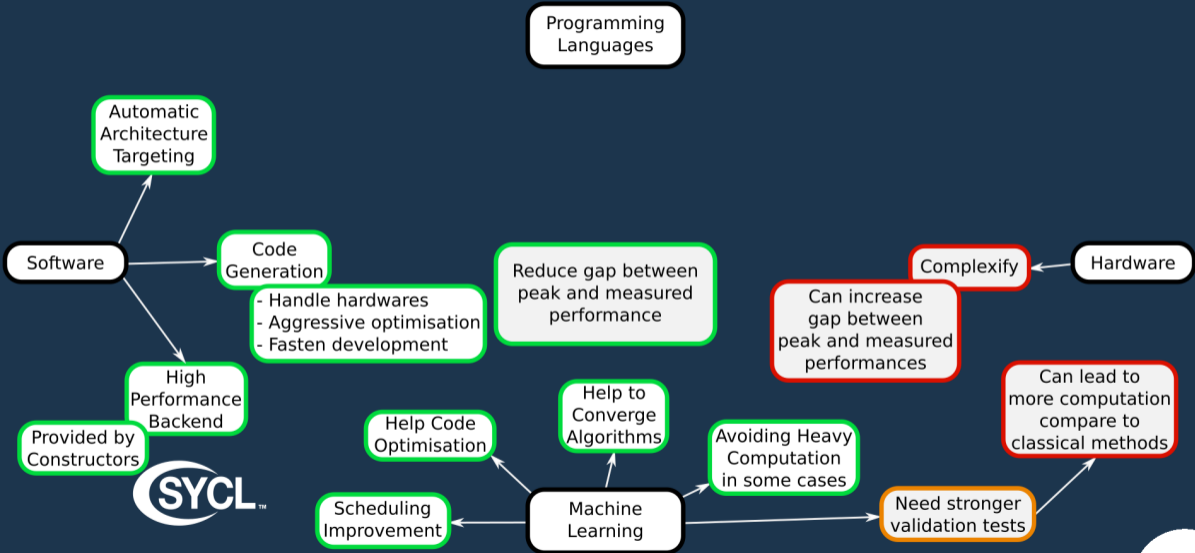
# Future of HPC



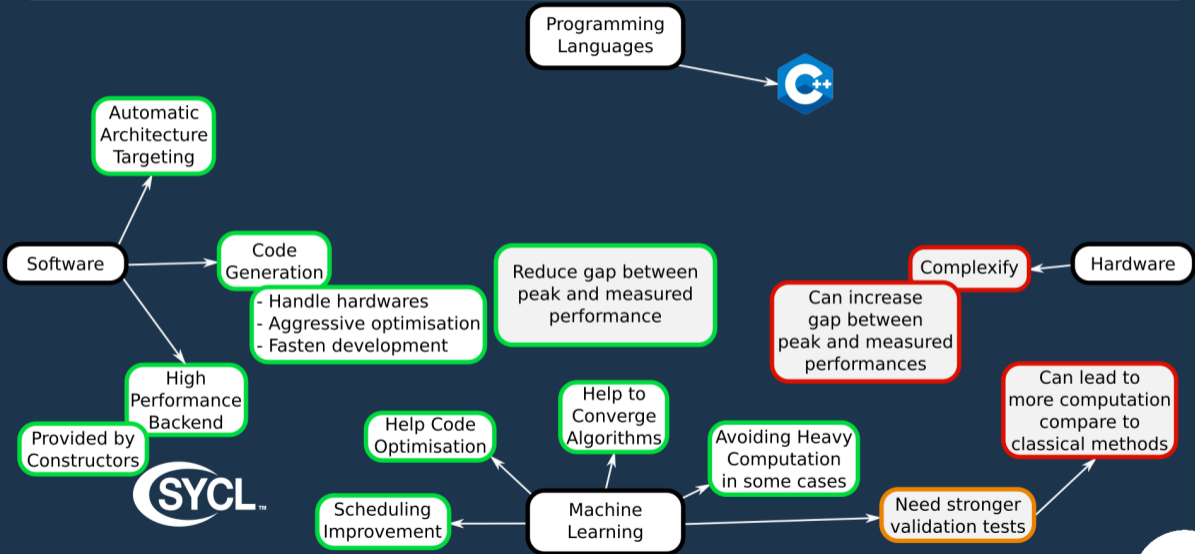
# Future of HPC



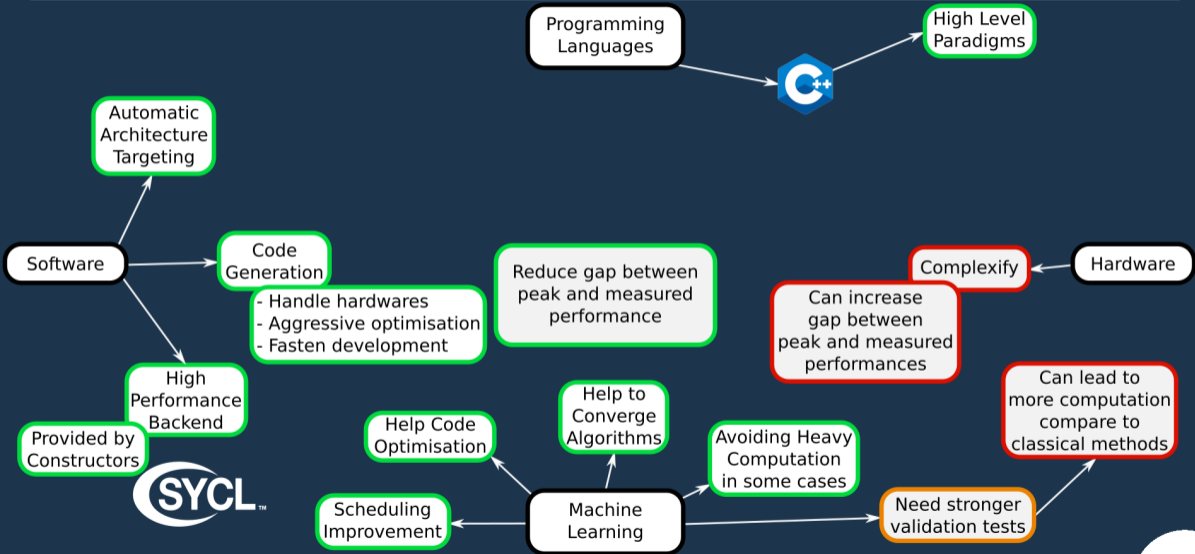
# Future of HPC



# Future of HPC

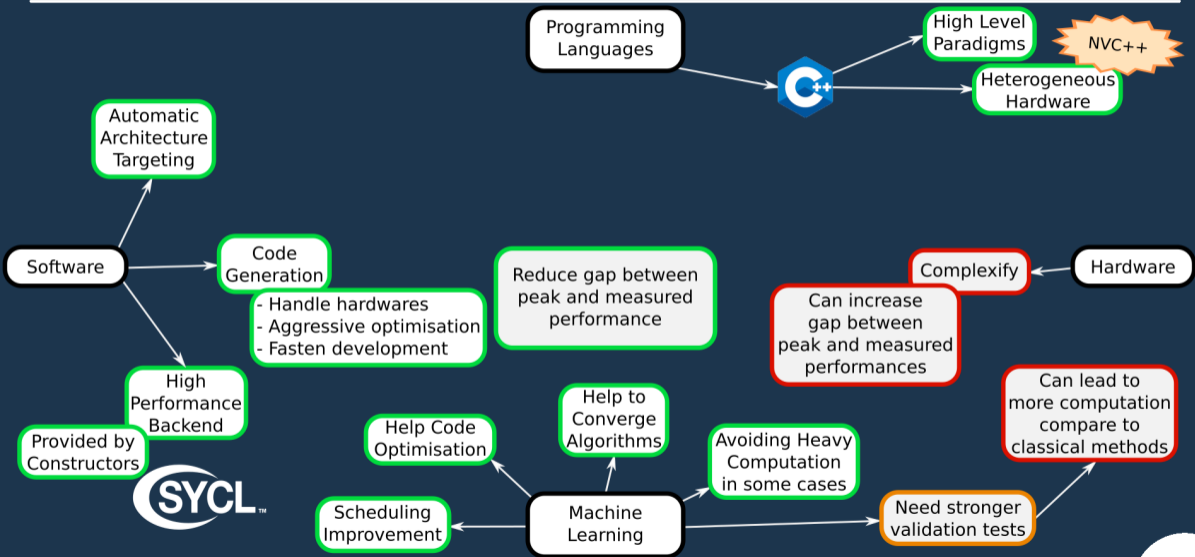


# Future of HPC

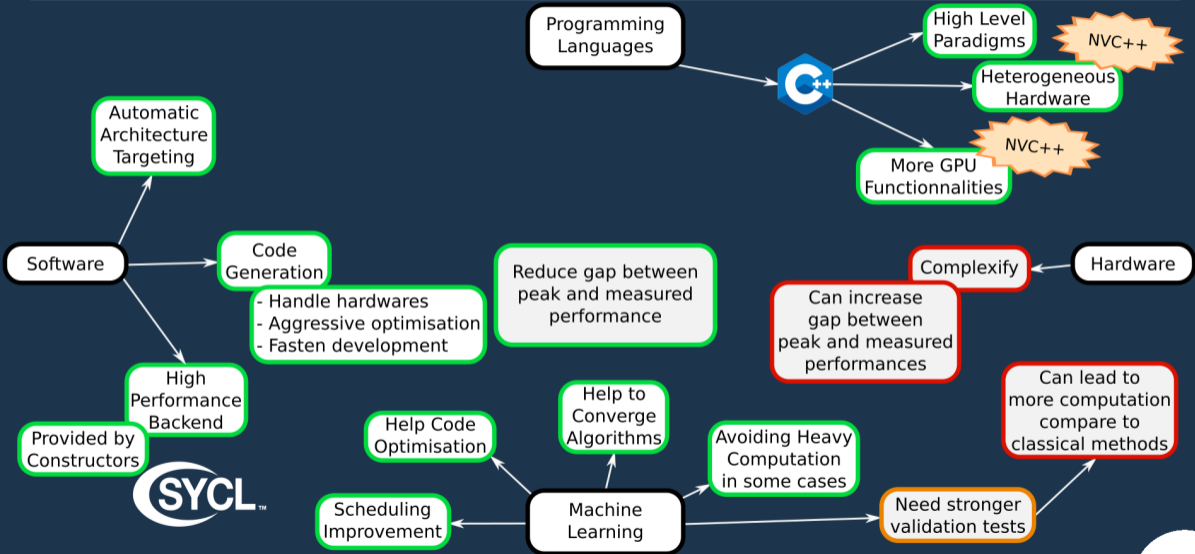




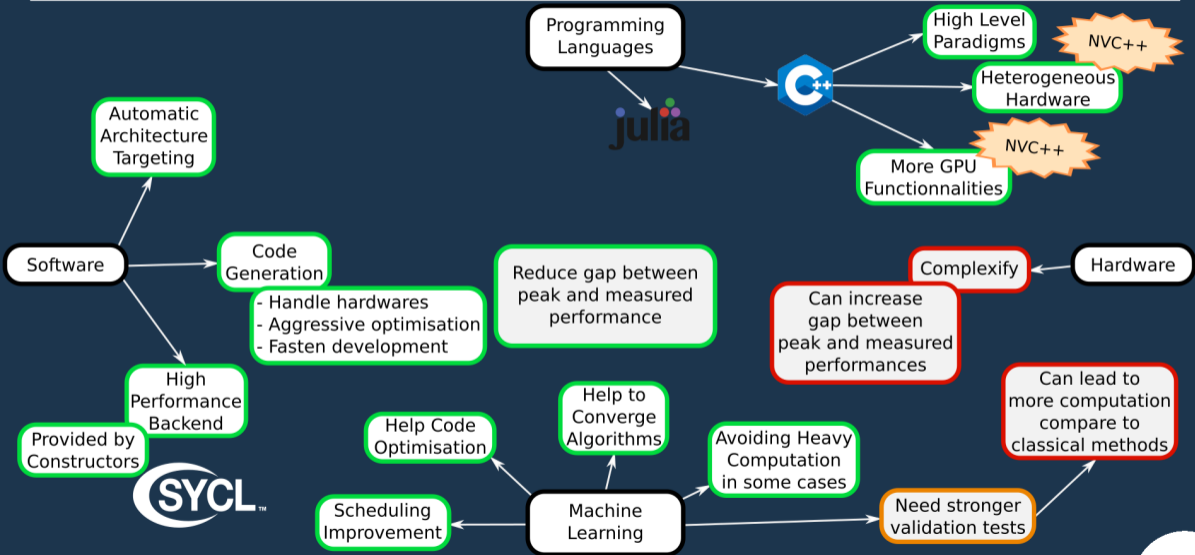
# Future of HPC



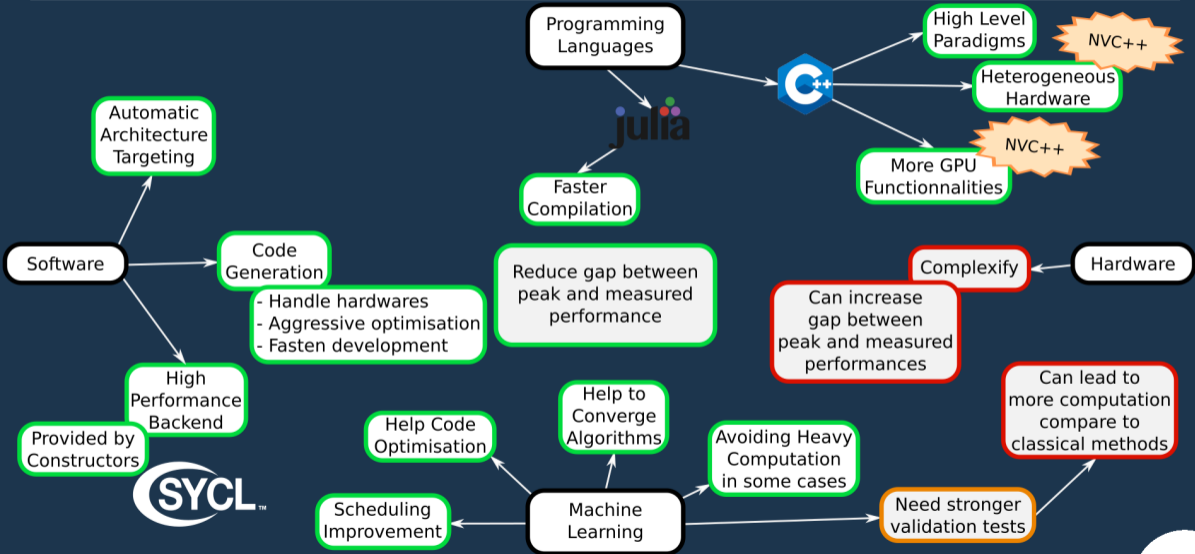
# Future of HPC



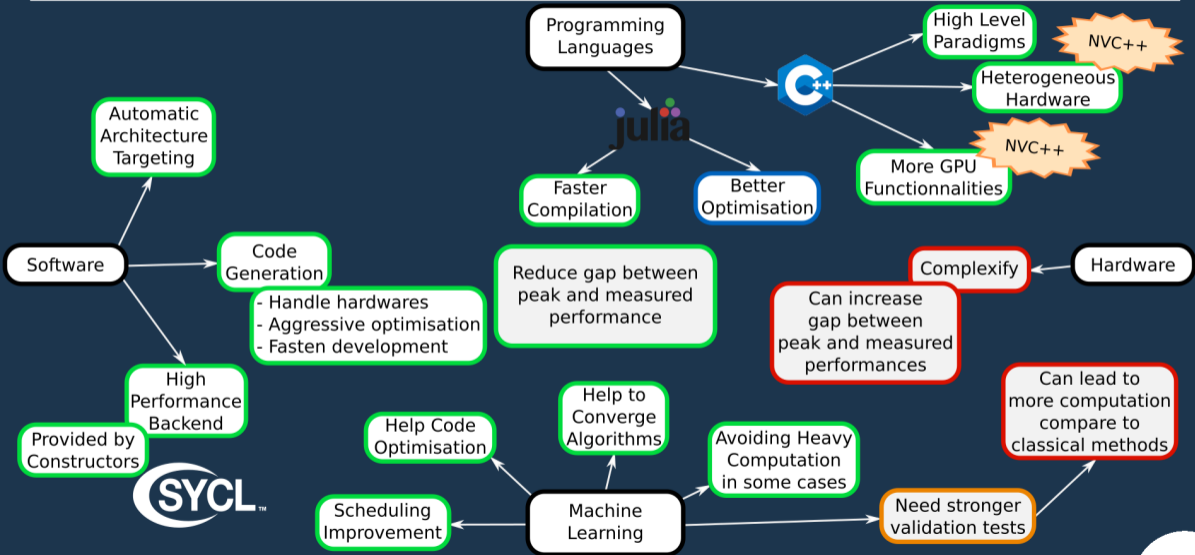
# Future of HPC



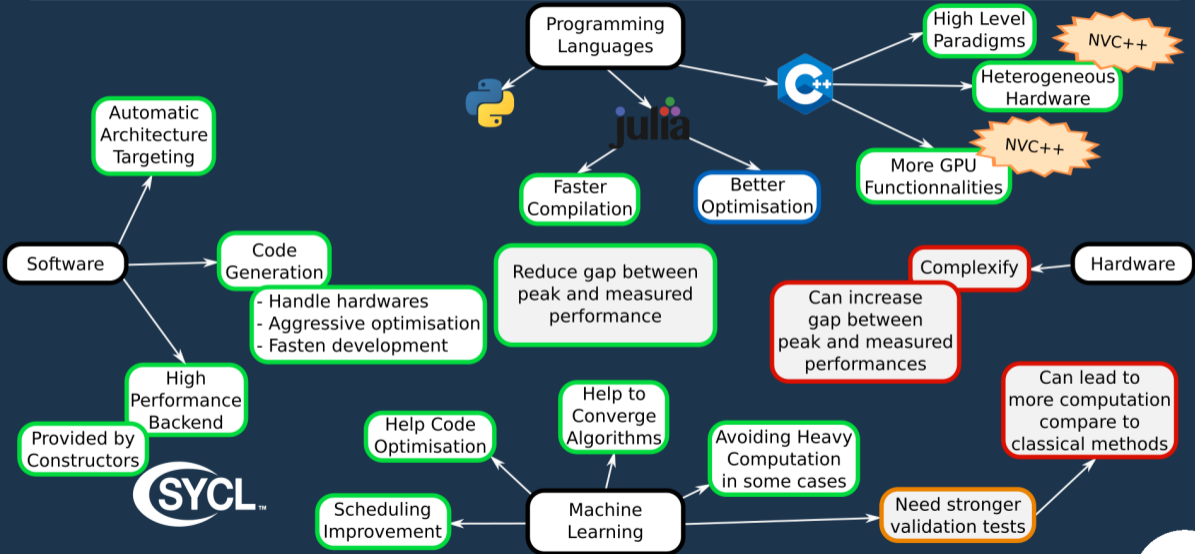
# Future of HPC



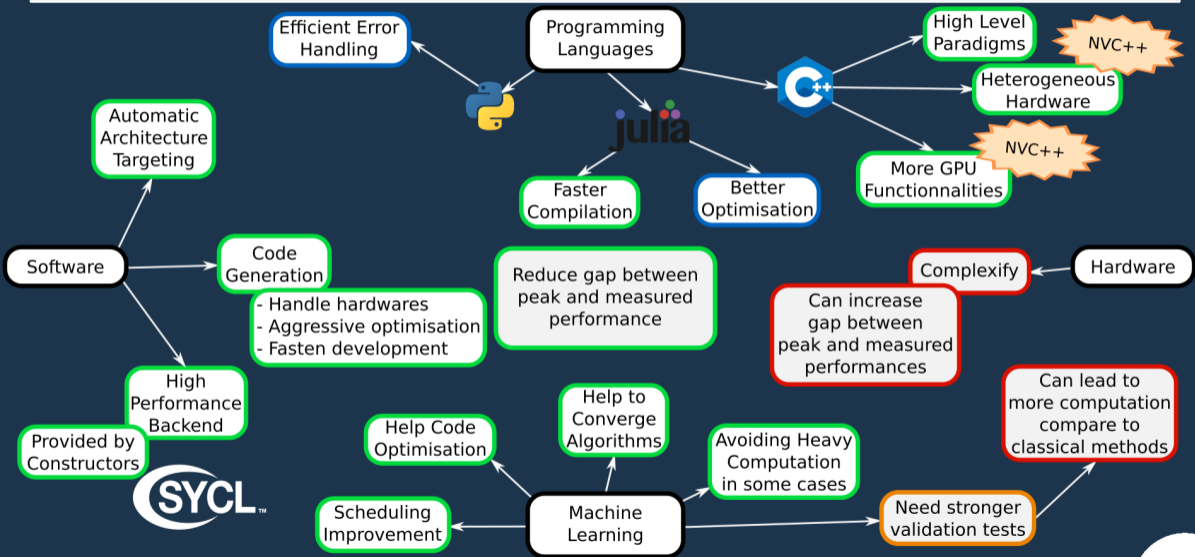
# Future of HPC



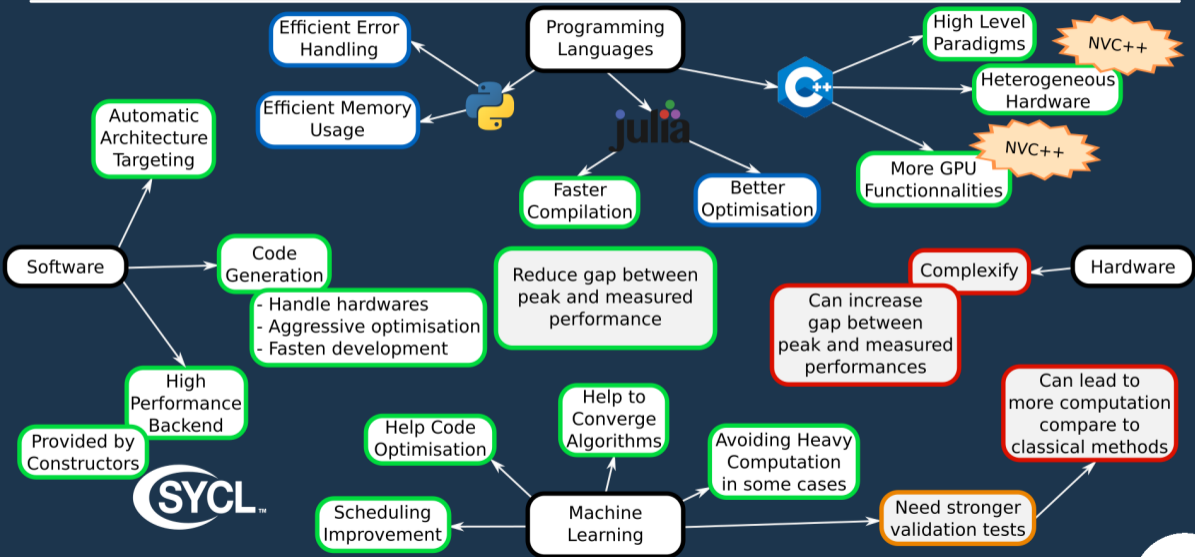
# Future of HPC



# Future of HPC

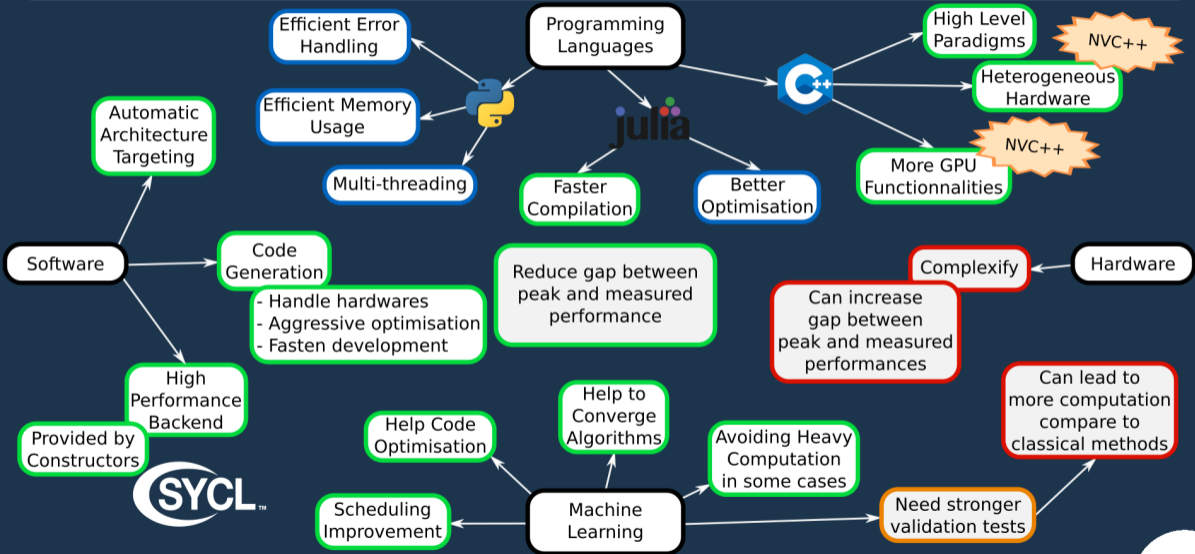


# Future of HPC

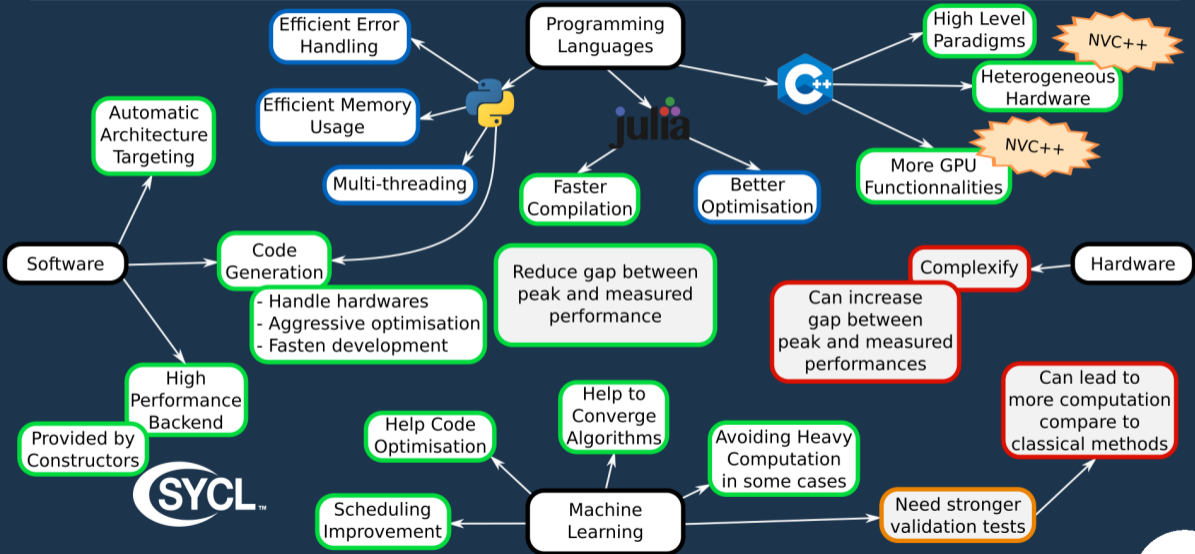




# Future of HPC



# Future of HPC



# Future of compilation stack

Goal

Time



# Future of compilation stack

Goal

Hardware

Time

# Future of compilation stack

Goal

Compiler

Hardware

Time

# Future of compilation stack

Goal

Language

Compiler

Hardware

Time

Goal

Libraries

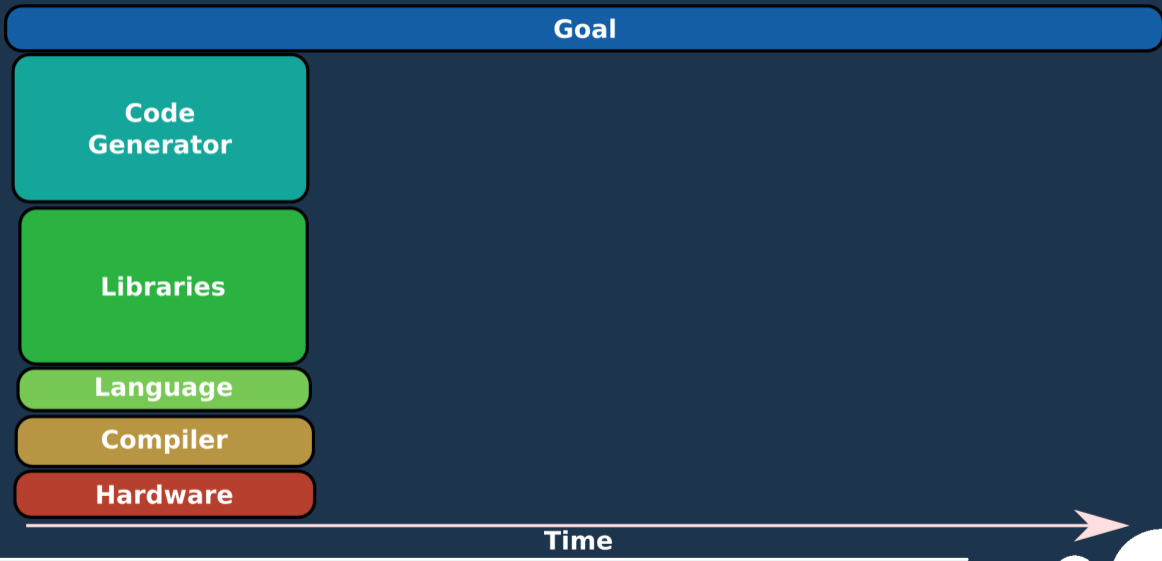
Language

Compiler

Hardware

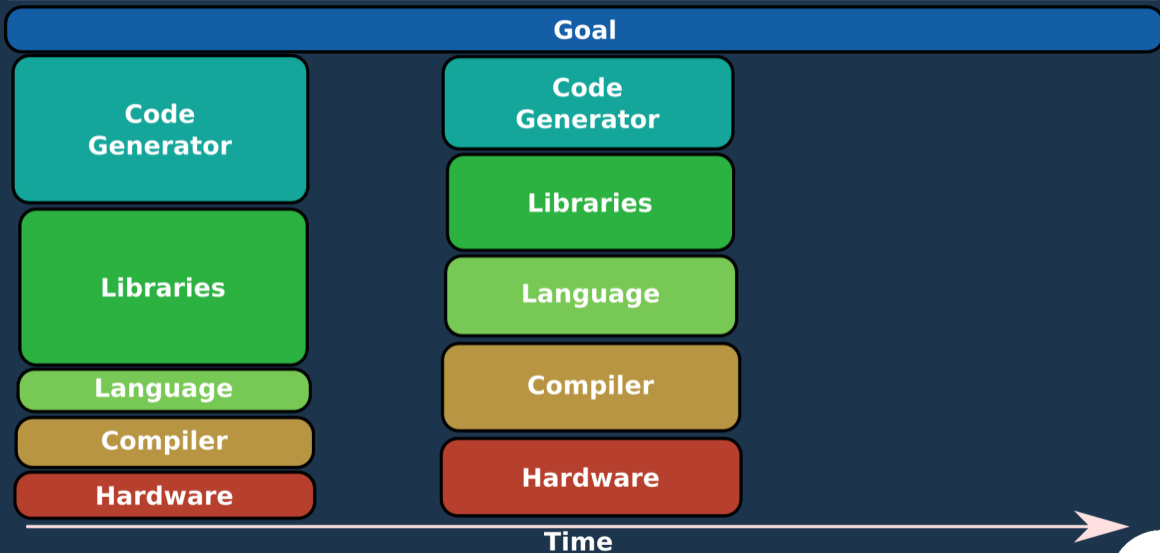
Time

# Future of compilation stack

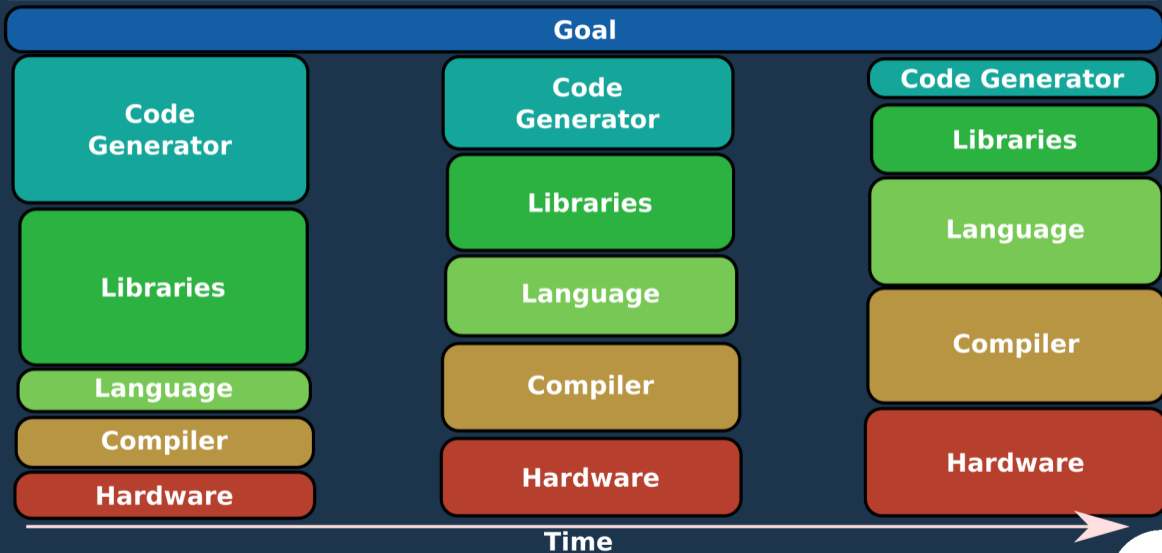




# Future of compilation stack



# Future of compilation stack



- ▶ **Advancing Exascale: Faster, Smarter, and Greener! (Presented by Atos) [SS32805]**
  - ▶ From 0.23 GFlops/W in 2010 to 25 GFlops/W 2020 (Juelich supercomputer - 3 744 A100), free cooling
- ▶ **Overcoming HPC Application Communication Bottlenecks with Intelligent and Automatic Resource Selection [S31547]**
  - ▶ NVTAGS : speed up from 1.42x to 4x just with jobs placement
  - ▶ NVTAGS doc
- ▶ **Accelerating GPU-Enabled HPC and Data Science Applications with On-the-Fly Compression [S31664]**
- ▶ **Optimizing Communication on GPU-Based HPC Systems for Dask and cuML Using MVAPICH2-GDR [S31627]**
- ▶ **Simplify Integrating Large-Scale Hybrid CPU and GPU Systems to Accelerate HPC Workloads (Presented by QCT) [SS32967]**

- ▶ **Go Beyond HPC: GPU Direct Storage, Parallel File Systems, and More (Presented by NetApp) [SS33181]**
  - ▶ Up to 17 GB/s in NVME read
- ▶ **Advancing Exascale: Faster, Smarter, and Greener! (Presented by Atos) [SS32805]**
  - ▶ **Classic Data Center** : 100 % computing, 100 % cooling
  - ▶ **Water cooling Data Center** : 100 % computing, 50 % cooling
  - ▶ **Direct liquid cooling Data Center** : 100 % computing, 5 % cooling

- ▶ **CUDA Code Optimization and Auto-Tuning Made Easy [S31429]**
  - ▶ Kernel Tuner code generator
- ▶ **RTCore for Compute: Exploiting Computational Patterns Using NVIDIA RTX [S31809]**
  - ▶ Particles simulation using RTX (noisy first 5 min)
- ▶ **What's New in OptiX [S31736] 55min**
- ▶ **Accelerating FFT toward Exascale Computing [P31661]**
- ▶ **Optimizing Lossless Compression Algorithms on the GPU [S32401]**
- ▶ **Accelerating GPU-Enabled HPC and Data Science Applications with On-the-Fly Compression [S31664]**
- ▶ **Fast and Simple Hash Tables [S31466] 38min**
- ▶ **GPU Accelerated Signal and Sensor Processing from Prototype to Deployment**

- ▶ **Inside NVC++ and NVFORTRAN [S31358] 33min**
  - ▶ C++ 17 and Fortran 202X on GPU with nvc++
- ▶ **A Deep Dive into the Latest HPC Software [S31286]**
  - ▶ Pragma for OpenACC
- ▶ **CUDA: New Features and Beyond**
  - ▶ CUDA Stream with Asynchronous allocation
  - ▶ Parallel functions loading

- ▶ **The Energy-Based view of Self-Supervised Learning (Yann Lecun)**
  - ▶ Same results with 10 min **unsupervised** compare to 100 h of **supervised** for speech recognition
- ▶ **NVIDIA DALI: GPU-Powered Data Preprocessing**
  - ▶ Up to 96% of performance peak with **Triton**, and **DALI** even if data are compressed
- ▶ **GPU Performance for Recommendation DNNs**
- ▶ **The Challenges in Hardware Aware Inference Optimization for Deep Learning**
- ▶ **Fast and Furious A GPU Inference Acceleration for State-of-the-Art YOLO v4 Model**
- ▶ **Starting a Deep Learning team with DGX Station A100 and Cnvr.io Presented by Microway Inc.**

- ▶ **Morpheus**
  - ▶ Spots sensitive data exchange (login, password, token)
- ▶ **Morpheus: AI Inferencing for Cybersecurity Pipelines [S33291]**
- ▶ **Building a GPU-Accelerated Cyber Flyaway Kit with NVIDIA Morpheus (Presented by Booz Allen Hamilton) [S32077]**
- ▶ **Cybersecurity is a Data Problem: Using NVIDIA Morpheus to Increase Network Visibility and Aid in Identifying Potential Network Problems [S32761]**



GTC 2021 notes from Reprises