

# Introduction to GPU Computing



Hans Henrik Brandenborg Sørensen  
DTU Computing Center

DTU Compute

<hhbs@dtu.dk>



$$f(x+\Delta x) = \sum_{i=0}^{\infty} \frac{(\Delta x)^i}{i!} f^{(i)}(x)$$

$\int_a^b \epsilon^{\theta} + \Omega^{\sqrt{17}} \int \delta e^{i\pi} =$   
 $\Theta^{\infty} = \{2.7182818284$   
 $\Sigma^{\gg}, \sum!$

# Practicalities

- Teacher(s) - Week 3
  - Hans Henrik B. Sørensen
  - Sebastian, Martin, Jonas
- Lectures / GPU Exercises – Mon-Tue
  - Learning by doing!
  - Complementary to Assignment 3 (should not be handed in)
- Assignment #3: GPU Matrix Multiplication and GPU Poisson Problem – Wed-Fri
  - Collaborate on developing the code + report
  - Fill out “responsibilities” in the addendum like week 1+2

# Learning objectives for week 3

- This week we focus on learning
  - Explain high-performance GPUs
  - Interplay of computer components like CPU, GPU, caches, and memory
  - Write parallel programs with OpenMP / Offloading
  - Write efficient programs for multi-core processors and many-core GPUs / Multi-GPUs
- Pre-requisites:
  - Proficiency in C/C++
  - Access to a so-called CUDA enabled GPU (e.g. DCC or home-pc with NVIDIA CC.>=8.0).

# Overview for week 3

- GPU computing – Mon
  - Introduction to GPU computing
  - OpenMP offload basics
  - OpenMP offload data mapping
- Performance tuning of GPU applications – Tue
  - Introduction to GPU tuning techniques
  - OpenMP offload memory accessing
  - OpenMP offload profiling tools / advanced topics
- Assignment 3 – Wed

# OpenMP books for week 3

OpenMP API 5.2 | C/C++ C/C++ content | Fortran or For Fortran content | [n.n.n] Sections in 5.2. spec | [n.n.n] Sections in 5.1. spec | See Clause info on pg. 9 | Page 1

**OpenMP**  
openmp.org

## Getting Started

Navigating this reference guide

Directives and Constructs | 1 C/C++ Variables | 14  
Clauses | 1 C/C++ Environment Variables | 14  
Runtime Library Routines | 10 Using OpenMP Tools | 16

An Examples Document and a link to a GitHub repository with code samples is at [link.openmp.org/examples51](https://link.openmp.org/examples51).

## OpenMP Examples Document

Fortran: Fortran directives are formed with comments in free form and fixed form sources (codes).

## Directives and Constructs

OpenMP constructs consist of a directive and, if defined in the syntax, an associated structured block that follows.

- OpenMP directives except SIMD and any declarative directive may not appear in Fortran PURE procedures.
- structured-block is a construct or block of executable statements with a single entry at the top and a single exit at the bottom.
- strictly-structured-block is a structured block that is a Fortran BLOCK construct.
- loosely-structured-block is a structured block that isn't strictly structured and doesn't start with a Fortran BLOCK construct.
- omp-integer-expression is of a C/C++ scalar int type or Fortran scalar integer type.
- omp-logical-expression is a C/C++ scalar expression or Fortran logical expression.

### Data environment directives

threadprivate [§ 2.1] [2.1.3]

Specifies that variables are replicated, with each thread having its own copy. Each copy of a threadprivate variable is initialized once prior to the first reference to that copy.

C/C++  
Fortran

!\$omp threadprivate [list]

list:

C/C++ A comma-separated list of file scope, namespace scope, or static block-scope variables that do not have incomplete types.

Fortran A comma-separated list of named variables and named common blocks. Common block names must appear between slashes.

### declare reduction [§ 2.1] [2.1.3]

Declares a reduction operator that can be used in a reduction, in\_reduction, or task\_reduction clause.

C/C++  
Fortran

!\$pragma omp reduction [operator : type-list : combiner] [initializer-clause]

!\$omp declare reduction & [operator : type-list : combiner] [initializer-clause]

combines:

C/C++ An expression

For An assignment statement or a subroutine name followed by an argument list.

initializer-clause: initializer [initializer-exp]

initializer-exp: omp\_parallel\_initializer or function-name [argument-list]

reduction-identifier:

C/C++ A base language identifier (for), or an identifier (for C++) or one of the following operators: \*, /, %, &, &&, ||

For A base language identifier, user-defined operator, or one of the following operators: +, -, .and., .or., .eqv., .neqv.; or one of the following intrinsic procedure names: max, min, land, lor, ieor.

C/C++ hypothesis-list: A list of type names

For type-list: A list of type specifiers that must not be CLASS(\*) or abstract type.

### Memory management directives

Memory spaces [§ 2.1.13]

Predefined memory spaces represent storage resources for storage and retrieval of variables.

Memory space: Standard allocation intent

omp\_default, mem\_space Default storage

omp\_align, mem\_space Large capacity

omp\_cyclic, mem\_space Variables with constant values

omp\_low\_bw, mem\_space High bandwidth

omp\_low\_lt, mem\_space Low latency

when [context-identifier-specification]: [directive-variant]

Conditionally select a directive variant.

otherwise [directive-variant]

Conditionally select a directive variant. otherwise was named default in previous versions.

Continued

OMP1121-01-OMP52



# Our suggestion (if you get hooked)

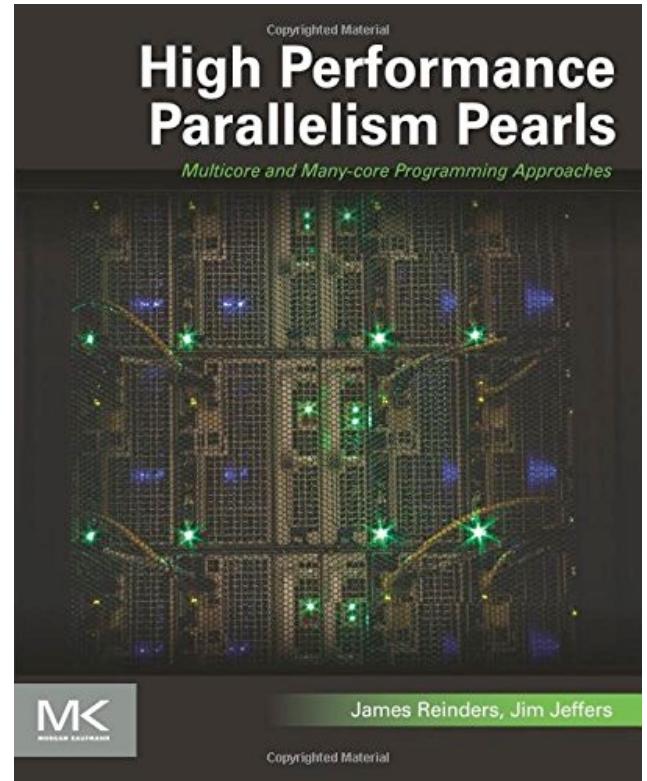


The screenshot shows the "HPC SDK DOCUMENTATION" page for the NVIDIA HPC SDK Version 22.11. The left sidebar contains links to various documentation sections: Accelerated Computing, HPC Compilers, Programming Models, CUDA Fortran Programming, OpenACC Getting Started, Other Documentation, Fortran CUDA Interfaces, Using OpenACC with MPI Tutorial, Support Services, Quick Start Guide, Terms & Conditions, and Archives. The main content area is titled "NVIDIA HPC SDK Version 22.11 Documentation" and "HPC Compilers". It includes links to the "HPC Compilers User's Guide" (describing how to use the HPC Fortran, C, and C++ compilers and program development tools), the "HPC Compiler Reference Manual" (including reference information for using the NVIDIA HPC compilers and program development tools), and sections on "Programming Models" (C++ Parallel Algorithms, CUDA Fortran Programming Guide, OpenACC Getting Started Guide), and "Other Documentation" (NVIDIA Fortran CUDA Library Interfaces, Using OpenACC with MPI Tutorial).

<https://docs.nvidia.com/hpc-sdk/compilers/index.html>

<https://www.openmp.org/spec-html/5.1/openmp.html#openmpch2.html>

<https://stackoverflow.com/questions/tagged/openmp+gpu>

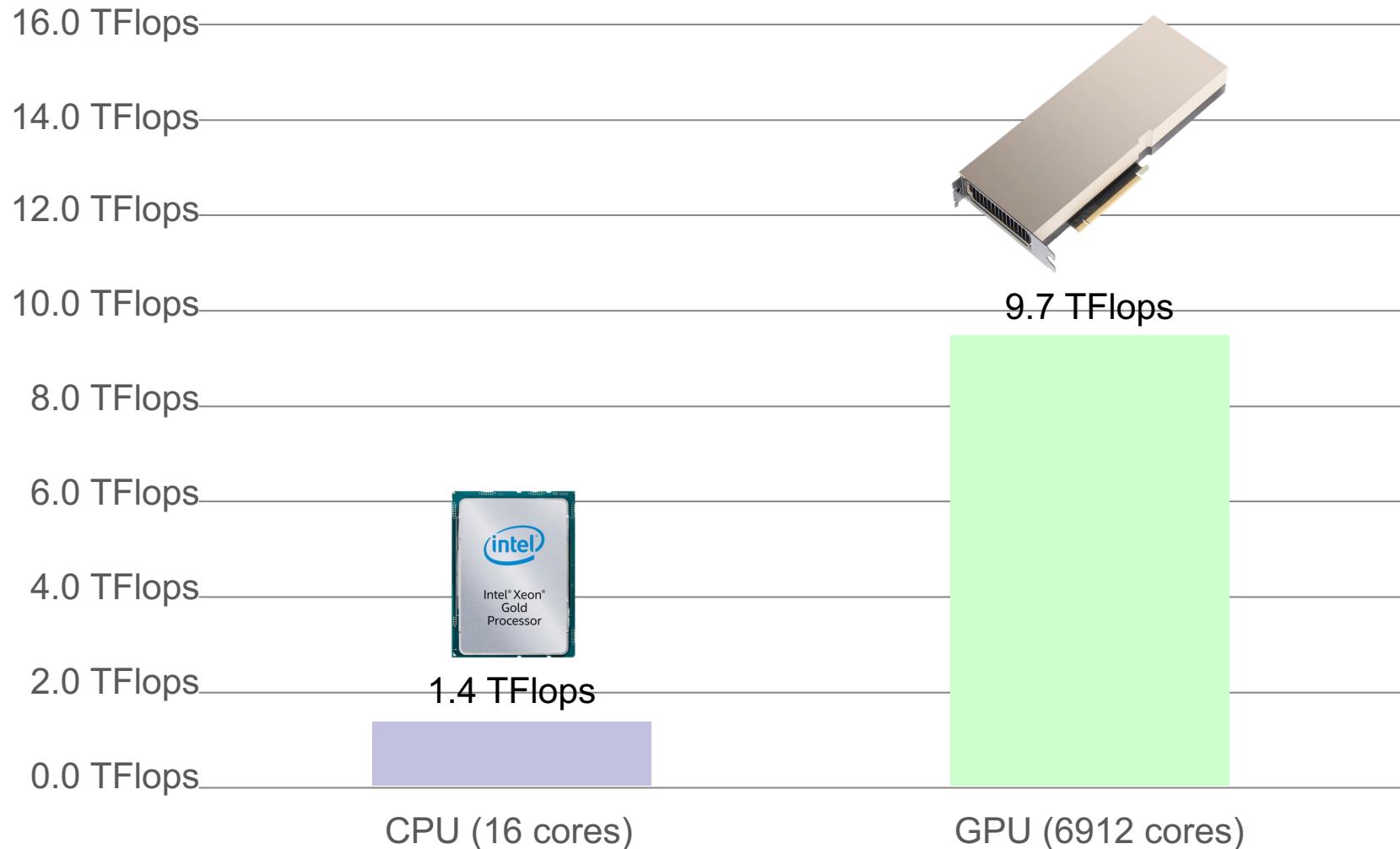


# Introduction to GPU computing

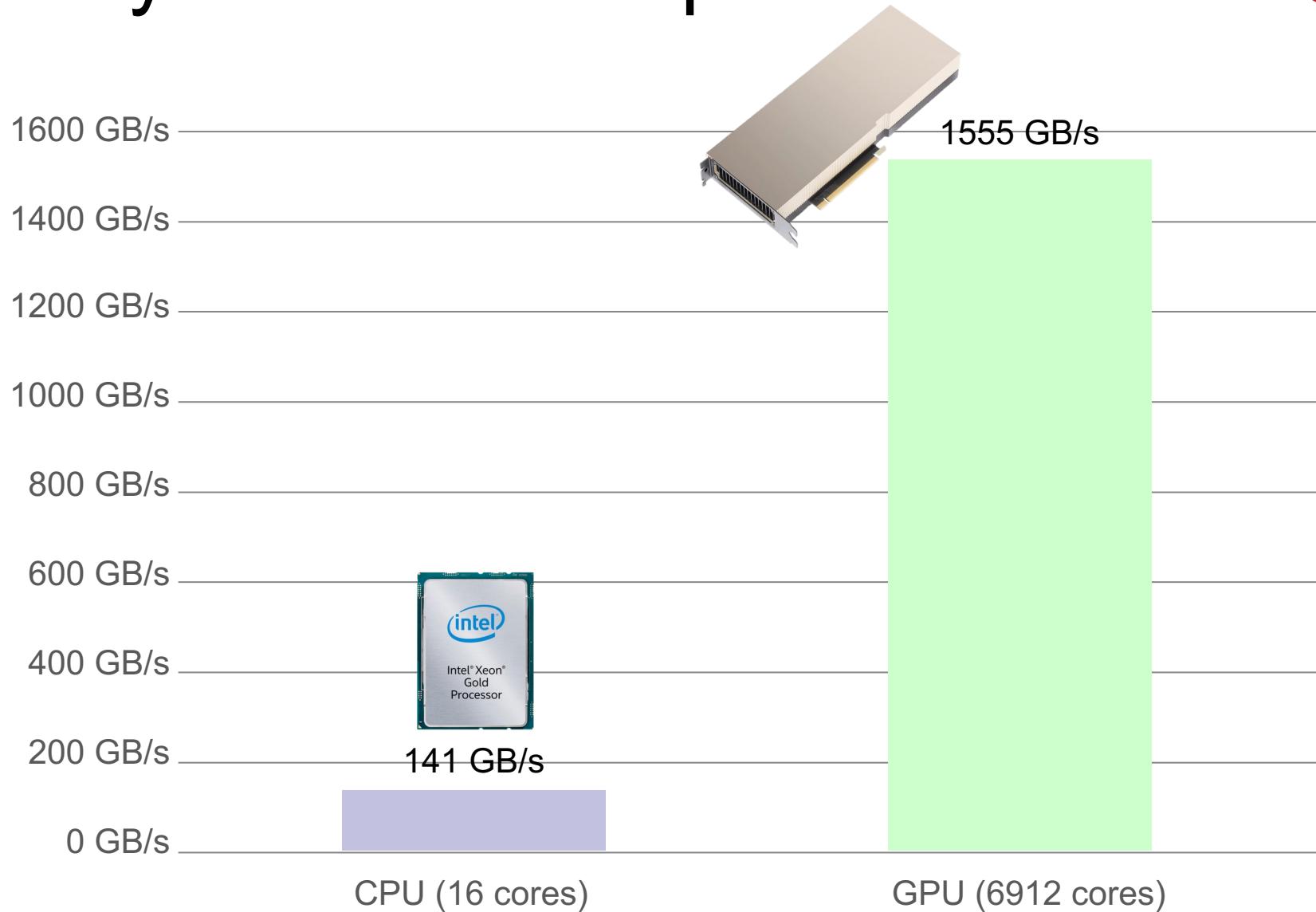
# Outline

- Why are GPUs important in HPC?
- Why are GPUs different from CPUs?
- GPUs as accelerators
- GPU hardware for 02614 course
- Motivation slides

# Why are GPUs important in HPC?



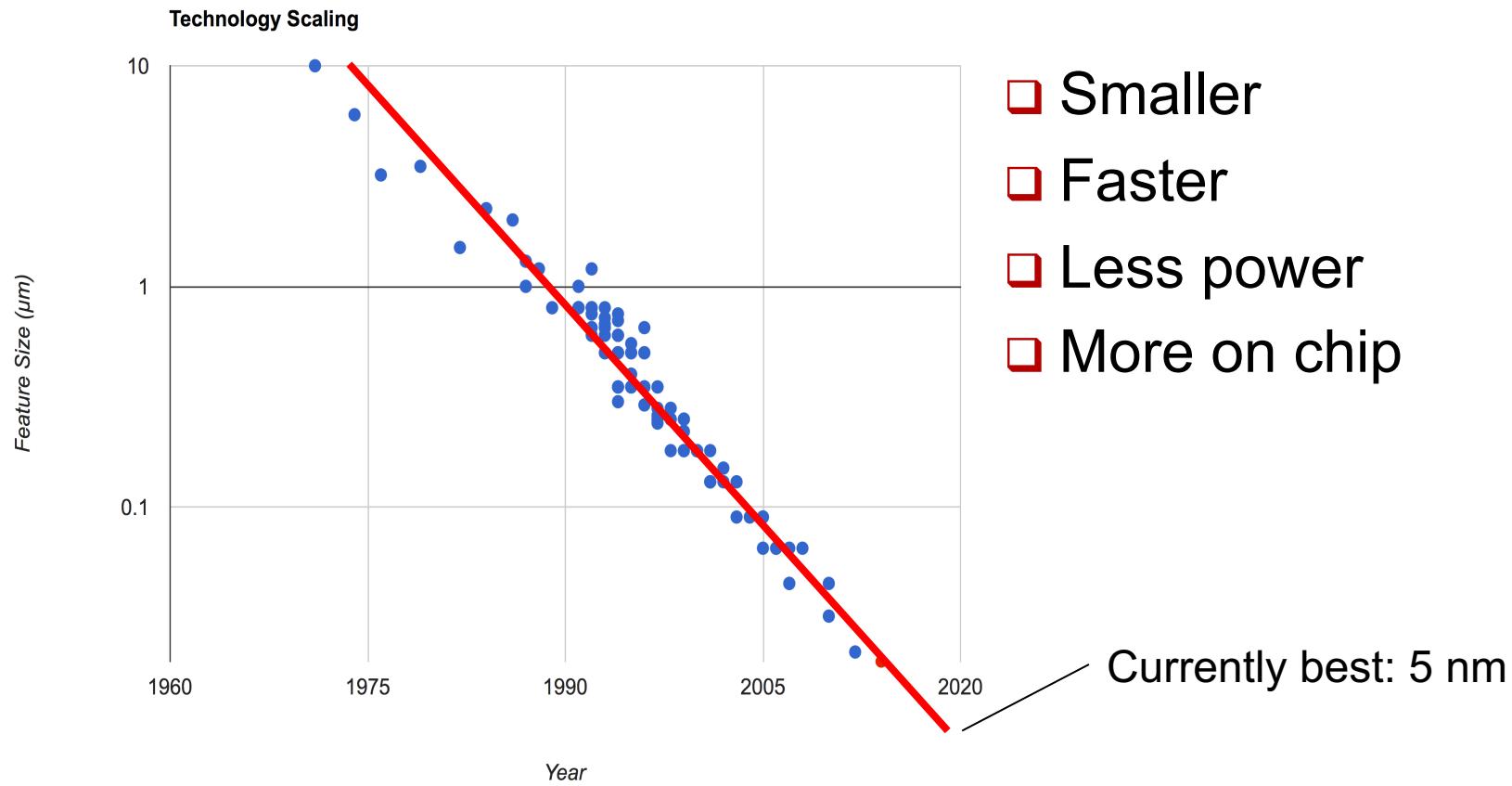
# Why are GPUs important in HPC?



# Why are GPUs different from CPUs?

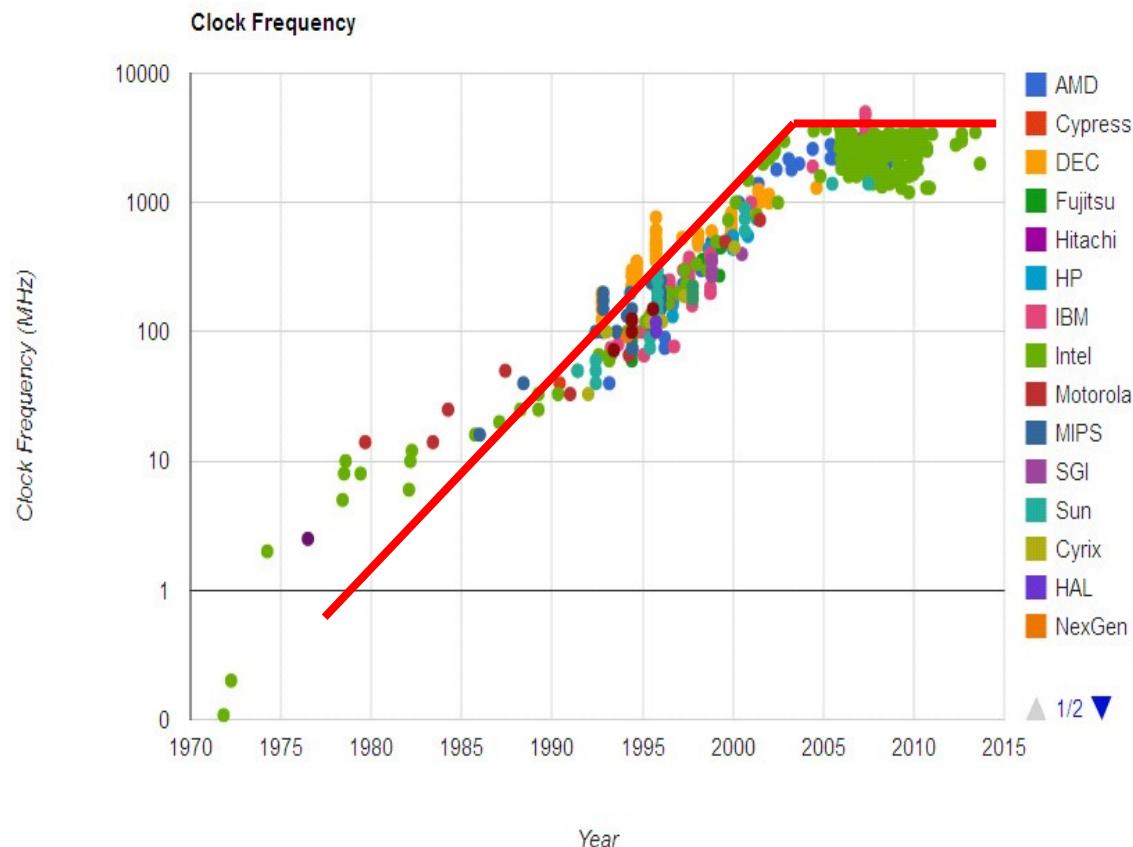
# Recap from week 1: Good news

- Transistor size over the years; still decreases



# Recap from week 1: Bad news

- Clock speed over the years; stagnated since 2005



- Increasing over many years
- However, over the last decade the clock speeds have essentially remained constant

# Why not increase clock speed?

- Power = Frequency x Voltage<sup>2</sup>
- Heat produced depends on power



# Why not increase clock speed?

- Power = Frequency x Voltage<sup>2</sup>
- Heat produced depends on power
- “Nard scaling”
  - Reduce transistor voltage to counter higher clock speed
  - Broke down in 2005 because of weakened current in the wires (need to distinguish 0 and 1)



# Why not increase clock speed?

- Power = Frequency x Voltage<sup>2</sup>
- Heat produced depends on power
- “Nard scaling”
  - Reduce transistor voltage to counter higher clock speed
  - Broke down in 2005 because of weakened current in the wires (need to distinguish 0 and 1)
- What matters today: # operations per Watt!
- Trade-off favors slower simpler processors
  - More operations per watt
  - Frequency kept low



# Building a modern processor

- What is the goal?

# Building a modern processor

- What is the goal? Roughly two choices...

A.

Latency  
(time to complete a task)

[e.g. seconds]

B.

Throughput  
(tasks completed per unit time)

[e.g. jobs/hour]

# Building a modern processor

- What is the goal? Roughly two choices...

A.

Latency  
(time to complete a task)

[e.g. seconds]

B.

Throughput  
(tasks completed per unit time)

[e.g. jobs/hour]

- Unfortunately these goals are not always aligned



# CPUs target latency (traditionally)



## ■ Usual task of traditional CPUs

### □ Desktop applications / OS

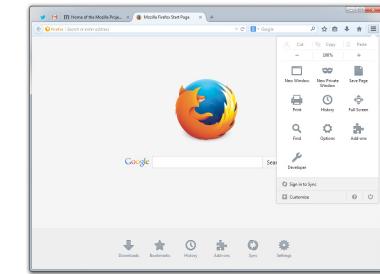
- Lightly threaded
- Lots of branches
- Lots of (indirect) memory accesses



Mac OS



Windows 10



## ■ CPUs try to minimize the time to complete a particular task – often to support user interaction!

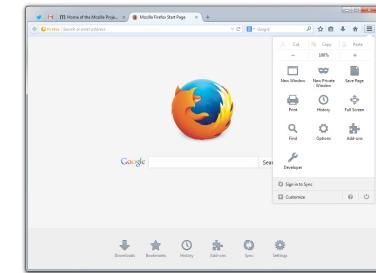
# CPUs target latency (traditionally)



## ■ Usual task of traditional CPUs

- Desktop applications / OS

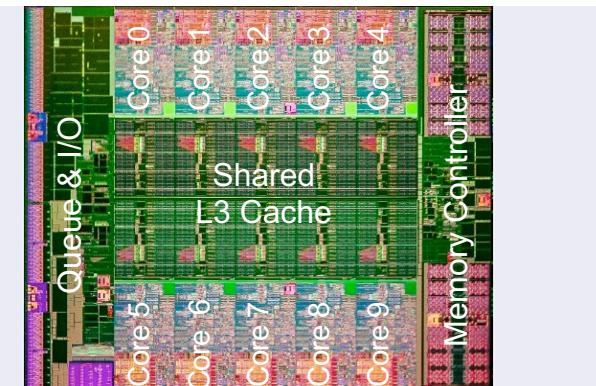
- Lightly threaded
- Lots of branches
- Lots of (indirect) memory accesses



## ■ CPUs try to minimize the time to complete a particular task – often to support user interaction!

## ■ Complex control hardware

- + Flexibility in performance
- + Lightly parallel
- - Expensive in terms of power



# GPUs target throughput

- GPUs are designed to **compute pixels** – fast!
  - Rendering video games in real-time
  - Play HD movies on smart phones
  - Render visual effects for movies...
- More concerned about the number of pixels per second than the latency of any particular pixel!



# GPUs target throughput

- GPUs are designed to **compute pixels** – fast!

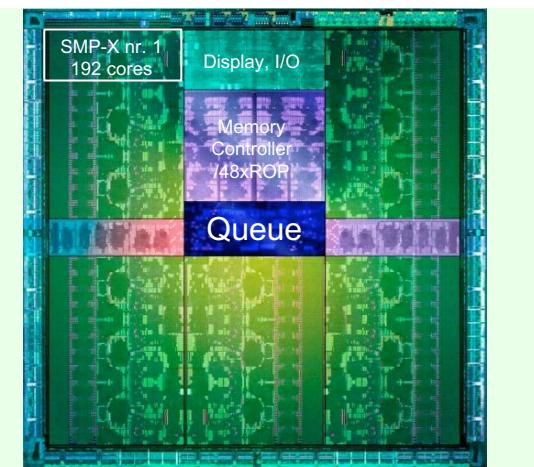
- Rendering video games in real-time
  - Play HD movies on smart phones
  - Render visual effects for movies...



- More concerned about the number of pixels per second than the latency of any particular pixel!

- Simpler control hardware

- + More transistors for computation
  - + Power efficient
  - - Highly parallel
  - - More restrictive in performance



# Hardware hierarchy

- CPU (typical)
  - Processing Unit (L3 cache)
  - Core (L2 cache / L1 cache / Instr. cache)

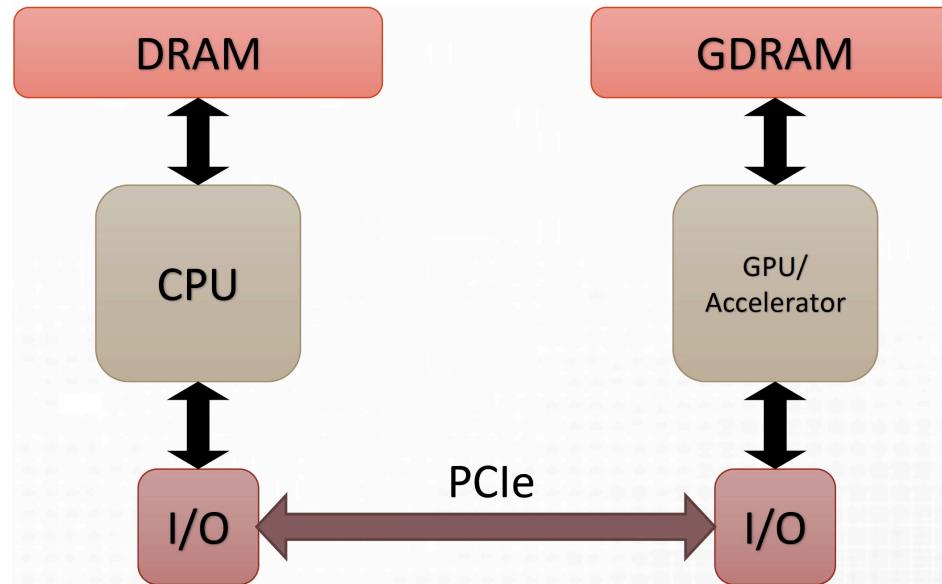
# Hardware hierarchy

- CPU (typical)
  - Processing Unit (L3 cache)
  - Core (L2 cache / L1 cache / Instr. cache)
- GPU (Nvidia)
  - Processing Unit (L2 cache)
  - GPC – Graphics Processing Cluster
  - TPC – Texture Processing Cluster
  - SM – Streaming Multiprocessor (64 “cores” / L1 cache)
  - Processing block (Instr. cache)
  - “Core”

# GPUs as accelerators

# GPUs as accelerators

- Problem: Still require OS, I/O, and scheduling
- Solution: “Hybrid system”
  - CPU provides management
  - Accelerators such as GPUs provide compute power



# Types of accelerators

## ■ GPUs

- HPC high-end versions – Tesla branch
- DP downgraded versions – Titan branch
- Gamer versions – RTX branch



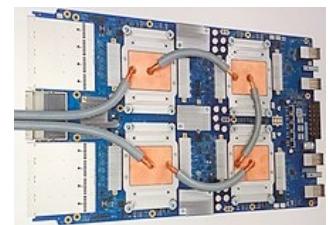
## ■ Custom many-core processors

- Japan / China



## ■ TPUs (Tensor Processing Unit)

- AI accelerator – application-specific integrated circuit (ASIC) developed by Google



# Accelerators in Top500



Rank	System	Cores	Rmax (PFlop/s)	Rpeak (PFlop/s)	Power (kW)
1	<b>Frontier</b> - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, <b>AMD Instinct MI250X</b> , Slingshot-11, HPE, DOE/SC/Oak Ridge National Laboratory, United States	8,730,112	1,102.00	1,685.65	21,100
2	<b>Supercomputer Fugaku</b> - Supercomputer Fugaku, A64FX 48C 2.2GHz, Tofu interconnect D, Fujitsu RIKEN Center for Computational Science, Japan	7,630,848	442.01	537.21	29,899
3	<b>LUMI</b> - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, <b>AMD Instinct MI250X</b> , Slingshot-11, HPE, EuroHPC/CSC, Finland	2,220,288	309.10	428.70	6,016
4	<b>Leonardo</b> - BullSequana XH2000, Xeon Platinum 8358 32C 2.6GHz, <b>NVIDIA A100 SXM4 64 GB</b> , Quad-rail NVIDIA HDR100 Infiniband, Atos EuroHPC/CINECA, Italy	1,463,616	174.70	255.75	5,610
5	<b>Summit</b> - IBM Power System AC922, IBM POWER9 22C 3.07GHz, <b>NVIDIA Volta GV100</b> , Dual-rail Mellanox EDR Infiniband, IBM, DOE/SC/Oak Ridge National Laboratory, United States	2,414,592	148.60	200.79	10,096

1st exascale

AMD Instinct  
(MI250X)

AMD Instinct  
(MI250X)

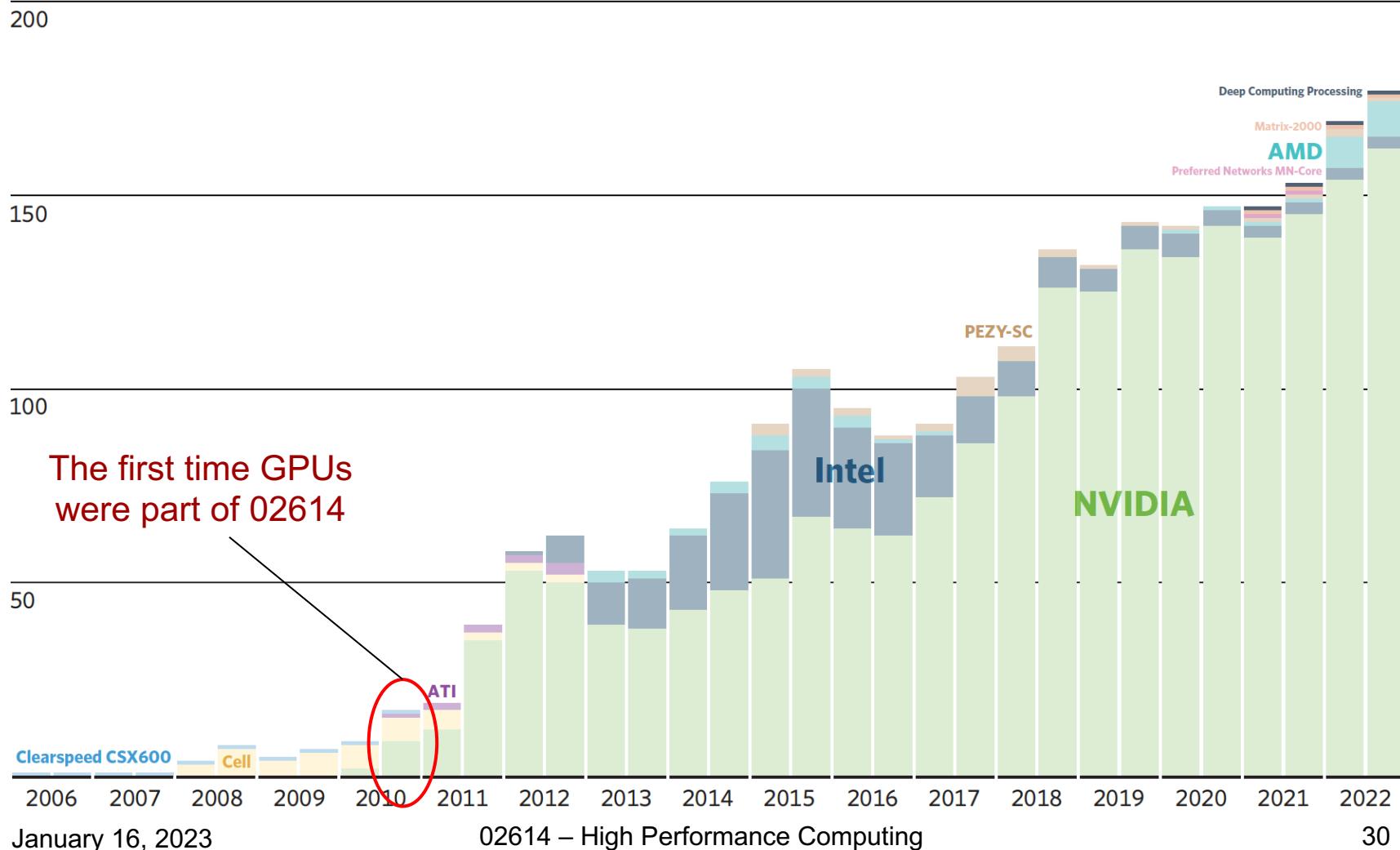
Nvidia Tesla  
(Ampere)

Nvidia Tesla  
(Volta)

# Accelerators in Top500



## ACCELERATORS/CO-PROCESSORS



# Nvidia GPU architectures

- Five generations of high-performance GPUs:

	# GPUs	Name	Year	Architecture	CUDA cap.	CUDA cores	Clock MHz	Mem GiB	SP peak GFlops	DP peak GFlops	Peak GB/s
Kepler	5	Tesla K40c	2013	GK110B (Kepler)	3.5	2880	745 / 875	11.17	4291 / 5040	1430 / 1680	288
	8	Tesla K80c (dual)	2014	GK210 (Kepler)	3.7	2496	562 / 875	11.17	2796 / 4368	932 / 1456	240
Pascal	8	*TITAN X	2016	GP102 (Pascal)	6.1	3584	1417 / 1531	11.90	10157 / 10974	317.4 / 342.9	480
	22	Tesla V100	2017	GV100 (Volta)	7.0	5120	1380	15.75	14131	7065	898
Volta	12	Tesla V100-SXM2	2018	GV100 (Volta)	7.0	5120	1530	31.72	15667	7833	898
	6	Tesla A100-PCIE	2020	GA100 (Ampere)	8.0	6912	1410	39.59	19492	9746	1555
Hopper	-	Tesla H100-SXM5	2022	GH100 (Hopper)	9.0	8448	1650	79.18	38984	19492	3110

Source: [http://www.hpc.dtu.dk/?page\\_id=2129](http://www.hpc.dtu.dk/?page_id=2129)

#Cores,  
Mem, and  
GB/s keep  
going up!

Clock freq.  
level off!

Peak  
increases 2x  
every ~3  
years!

# Trends in HPC due to GPUs



- Improvements at individual computer node level are greatest
  - Less data transfer
  - Heterogeneous computing
  - Avoiding MPI

# Trends in HPC due to GPUs

- Improvements at individual computer node level are greatest
  - Less data transfer
  - Heterogeneous computing
  - Avoiding MPI
- “Super-nodes”: Nvidia DGX-1
  - Eight tightly linked high-end GPUs
  - 40,960 cores / 960 TFlops in 1 node



# Trends in HPC due to GPUs

- Improvements at individual computer node level are greatest
  - Less data transfer
  - Heterogeneous computing
  - Avoiding MPI
- “Super-nodes”: Nvidia DGX-1
  - Eight tightly linked high-end GPUs
  - 40,960 cores / 960 TFlops in 1 node
- Communication costs are increasing
  - Synchronization-reducing algorithms
  - Communication lower-bound algorithms

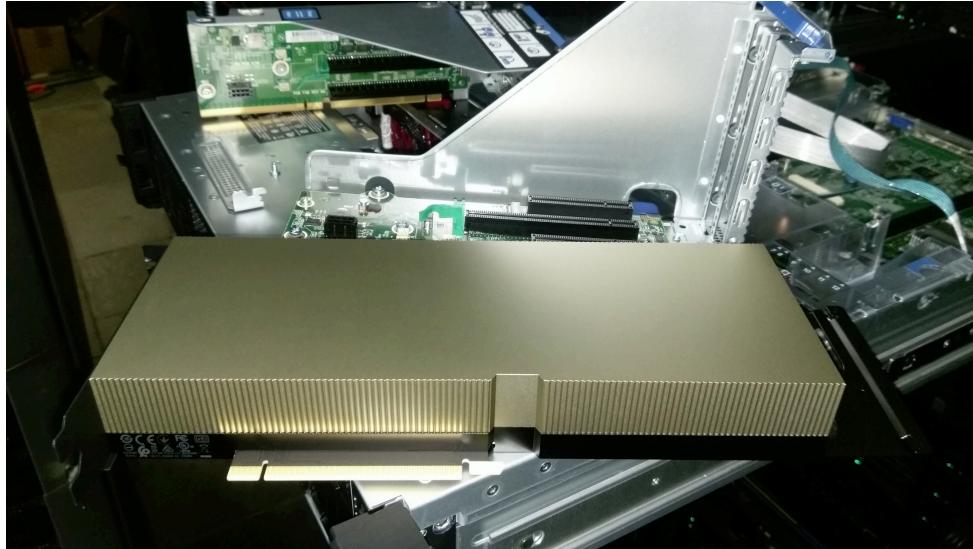


# GPU hardware for this course

# Ampere cluster for 02614

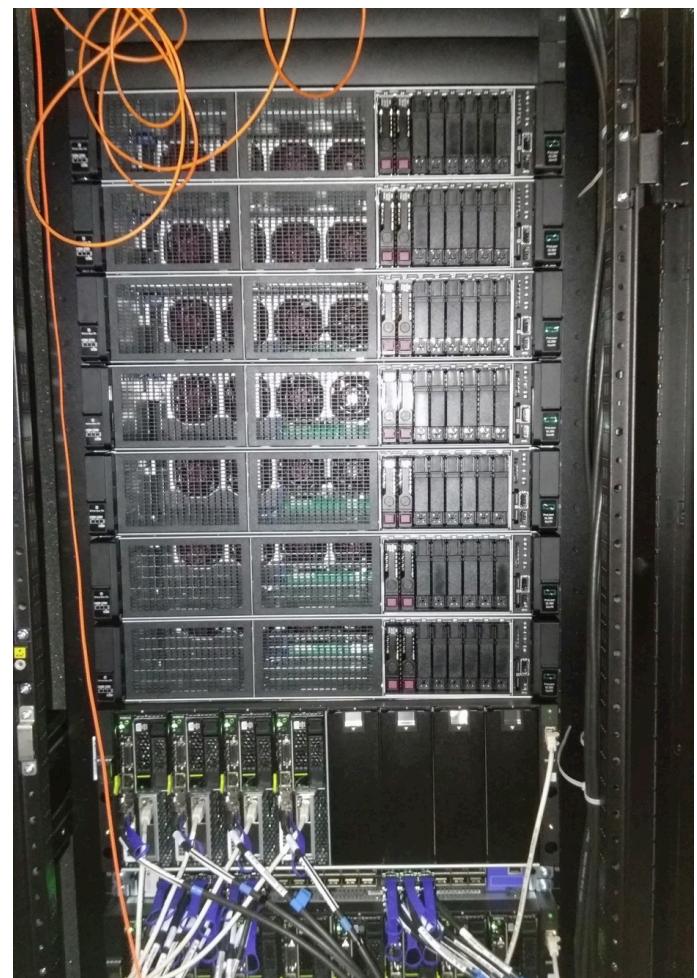


Delivered: December 23, 2020



1 HPE node:

2 x Intel Xeon Gold 6126 CPU @ 2.90GHz (16 cores)  
2 x Tesla A100-PCIE-40GB (6912 cores)  
768 GB DDR4 @ 2666MHz



6 HPE nodes:

Theoretical peak: 133,2 TFlops (FP64)

# Ampere cluster for 02614

A100-PCIE-40GB

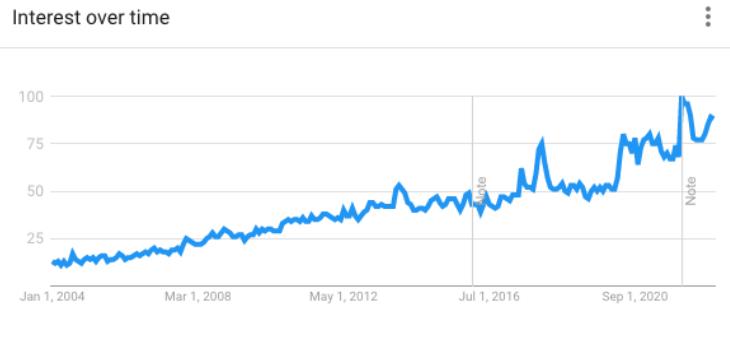


- 108 streaming multiprocessors (SM) – 40 GB memory
- Multi-user system: GPUs are in default mode (shared)
- Batch system: Script sets GPUs to exclusive mode

# Google trends: GPUs still going up



## GPU



NVIDIA DGX-1  
AI Supercomputer

ORDER NOW

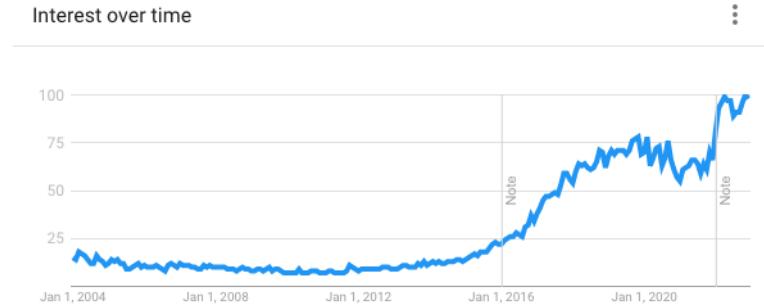
THE WORLD'S FIRST AI SUPERCOMPUTER IN A BOX

Get faster training, larger models, and more accurate results on [deep learning](#) with the NVIDIA® DGX-1™. This is the world's first purpose-built system for deep learning and AI accelerated analytics,

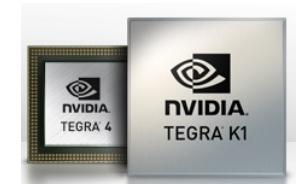
## CPU



## Machine learning



TEGRA K1—THE WORLD'S MOST ADVANCED MOBILE PROCESSOR



# Motivation Jan 2020

**Valgte søgekriterier**

GPU

NULSTIL →

GEM SØGNING →

**Geografi** ▾

**Arbejdsmønster** ▾

**Ansættelsesvilkår** ▾

**Ansættelsens varighed** ▾

**Arbejdstid** ▾

**3 jobopslag** Sorter efter: Bestre match → VIS PÅ KORT →

**Porting CUDA and CPU code from Windows Nvidia P4/T4 to Linux Nvidia Jetson Worksome ApS**  
... windows Visual Studio C++ code to Linux and at the same time moving cuda code from P4/T4 GPU to Jetson Please apply with your rate and an estimate on the project cost/length ...

Bemærk! Jobannoncen åbner i en ny fane

Indrykket: 20. november 2019 - Storbritannien  
Deltid (5 - 36 timer ugentligt)  
Ansøgningsfrist: 15. januar 2020

Tip en ven

Id: 5076183

**Programmør til High Performance Computing (HPC) Forsvarets Efterretningstjeneste**  
...medarbejdere, hvor din arbejdsgave primært vil være at programmere symmetriske multiprocessorsystemer (fx GPU).  
Som en del af indhentningssektoren kommer du til at arbejde...

Bemærk! Jobannoncen åbner i en ny fane

Indrykket: 19. december 2019 - 2100 København Ø  
Fuldtid  
Ansøgningsfrist: 12. januar 2020

Tip en ven  Ruteplan

Id: 5090019

**PhD fellow in Computer Science KU - SCIENCE - DATALOGISK INSTITUT - UP1**  
...probability distributions and sampling from them; generating high-performance vectorized multicore or GPU code; and applications to deep learning, Bayesian inference and probabilistic programming.<>...

Bemærk! Jobannoncen åbner i en ny fane

Indrykket: 10. december 2019 - 2100 København Ø  
Fuldtid  
Ansøgningsfrist: 15. januar 2020

Tip en ven  Ruteplan

Id: 5085539

**3 jobopslag**

Lukas Christian Høghøj, *Large-scale modelling on GPUs with OpenACC*, 2019

Patrick Møller Jensen and Julian Thomas Reckeweg Olsen, *GPU beamforming*, 2019

Konstantinos Gkanos, *Interactive, real-time room acoustic simulations*, 2019

Gandalf Saxe and Oisin D. Kiær, *Low Energy Transfer Orbits to Mars using Evolution Strategies*, 2019

Mia Sandra Nicole Siemon, *Comparison of GPU programming models*, 2019

Mathias Sorgenfri Lorenz, *Scaling analysis of a multi-GPU Poisson solver*, 2018

Tim Felle Olsen and Mathias Sorgenfri Lorenz, *Large-scale computations on modern GPUs*, 2018

Nick Clausen, *Leveraging the GPU architecture of embedded systems for model predictive control problems*, 2018

...

# Motivation Jan 2022 - 2023

[PhD scholarship in Cardiac Vector Flow Ultrasound Imaging](#)

[Danmarks Tekniske Universitet](#), Kongens Lyngby

A 3-year PhD scholarship is available at the Center for Fast Ultrasound Imaging (CFU), Department of Health Technology (DTU Health Tech), DTU and at the Department of Cardiology, Rigshospitalet from March 2022. It is conducted in collaboration between the two departments, and the overall purpose of the project is to further develop fast vector flow-based imaging of the heart and make it applicable for clinical use.

This PhD project develops a method for fast vector flow imaging on a Verasonics research scanner and determines its performance and applicability to cardiac diagnostics.



Tjek jobglæden:

★★★★★ 391 evalueringer

INDRYKKET: 17-12-2021

**Apply for this job**

Apply no later than 21 January 2022

Apply for the job at DTU Health Tech by completing the following form.

**Apply online**




International Workshop on OpenMP  
 ↳ IWOMP 2022: [OpenMP in a Modern World: From Multi-device Support to Meta Programming](#) pp 81–93 | [Cite as](#)

## Feasibility Studies in Multi-GPU Target Offloading

[Anton Rydahl](#)✉, [Mathias Gammelmark](#) & [Sven Karlsson](#)

Conference paper | [First Online: 20 September 2022](#)

178 Accesses

Part of the [Lecture Notes in Computer Science](#) book series (LNCS, volume 13527)

### Abstract

Many of the largest supercomputers are based on heterogeneous architectures with multiple general-purpose graphics processing units (GPGPUs) per compute node. While many APIs

[Download book PDF](#)



[Download book EPUB](#)



[Sections](#)

[Figures](#)

[References](#)

[Abstract](#)

[Introduction](#)

[Background](#)

[Design Patterns](#)

[https://link.springer.com/chapter/10\\_1007/978-3-031-15922-0\\_6#Sec9](https://link.springer.com/chapter/10_1007/978-3-031-15922-0_6#Sec9)

# End of lecture