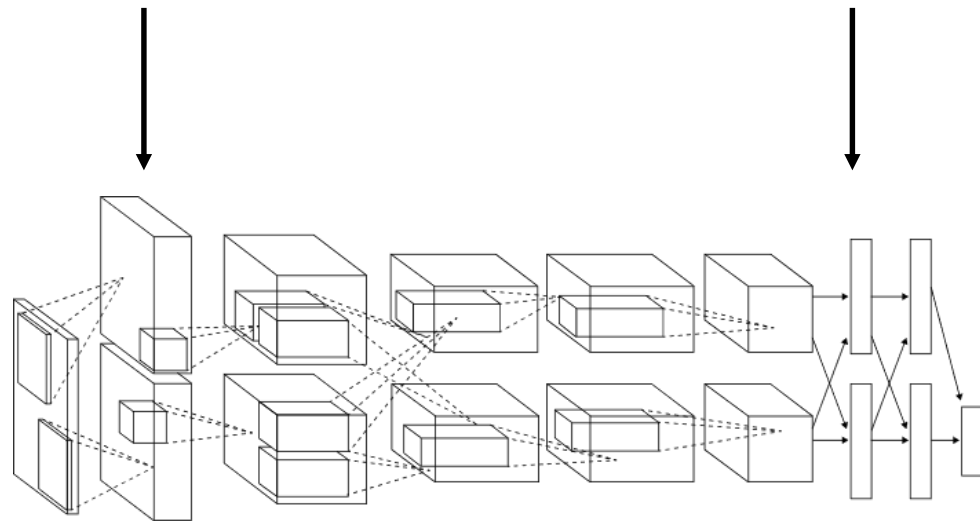# AI for Science

SEAB

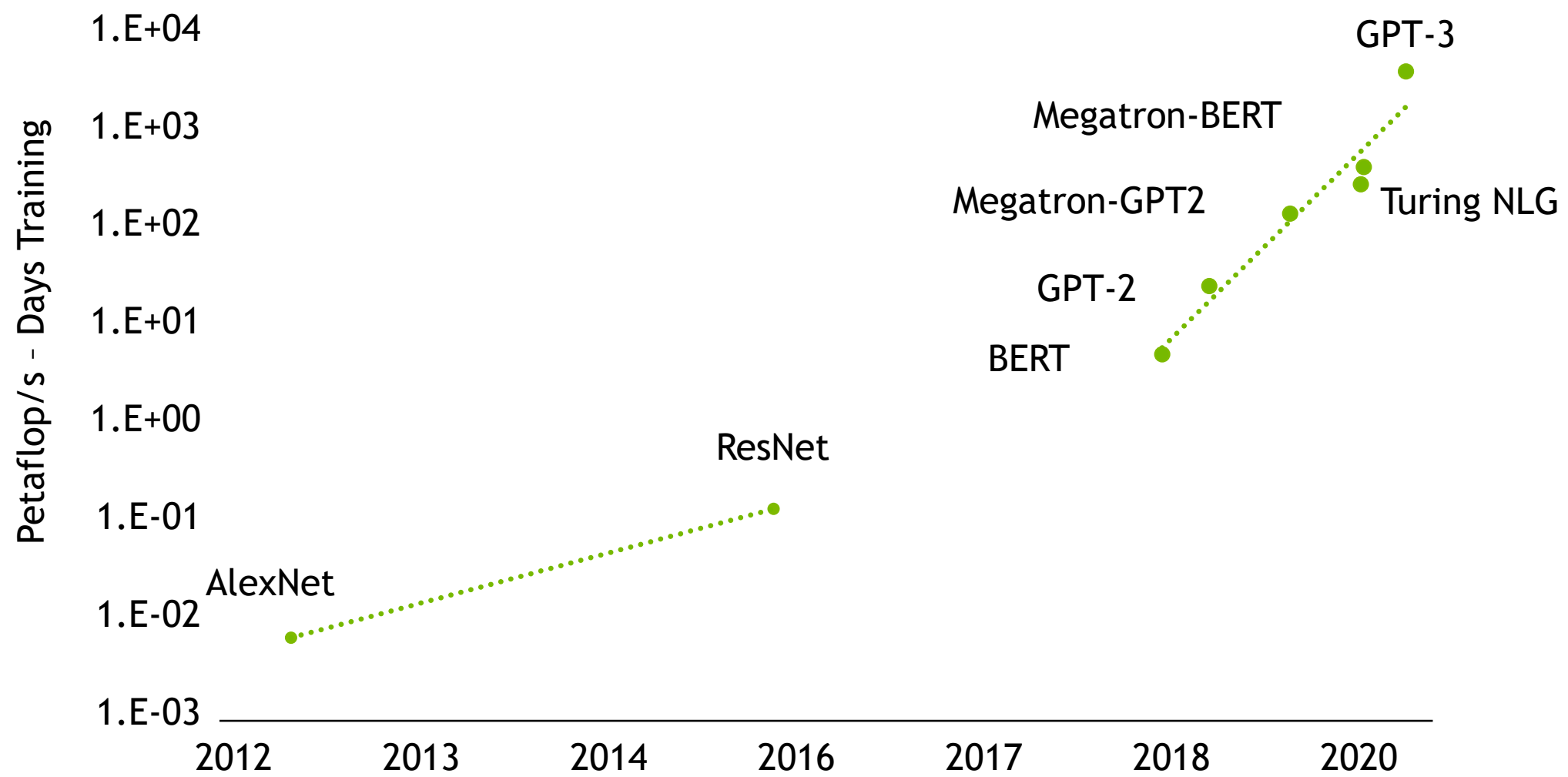October 26, 2023

**Bill Dally**

Chief Scientist and SVP of Research, NVIDIA Corporation

Adjunct Professor of CS and EE, Stanford

# Deep Learning was Enabled by GPUs

# Deep Learning is Gated by GPUs

GPT-4



Petaflop/s - Days Training

GPT-3

Megatron-BERT

Megatron-GPT2

Turing NLG

GPT-2

BERT

ResNet

AlexNet

1.E+04
1.E+03
1.E+02
1.E+01
1.E+00
1.E-01
1.E-02
1.E-03

2012   2013   2014   2016   2017   2018   2020

NVIDIA.

Gains from

Number representation
    FP32, FP16, Int8
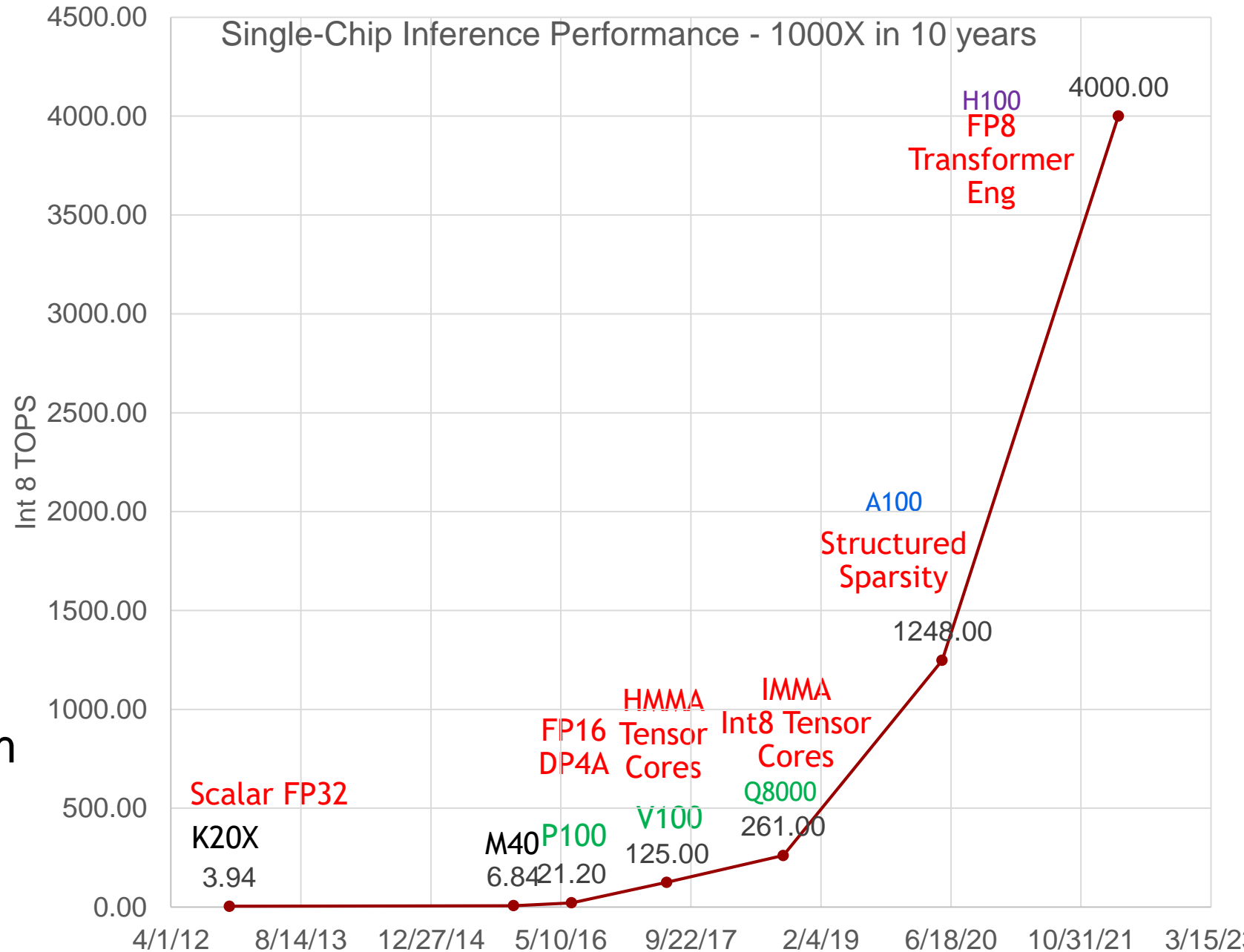    (TF32, BF16)
    16x

Complex instructions
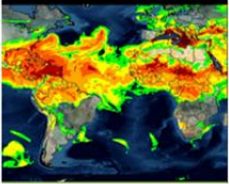    DP4, HMMA, IMMA
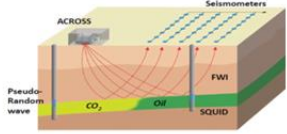    12.5x

Process
    28nm, 16nm, 7nm, 5nm
    2.5x

Sparsity – 2x

Single-Chip Inference Performance - 1000X in 10 years

Int 8 TOPS

Scalar FP32

K20X
3.94

M40
6.84

FP16
DP4A

P100
21.20

HMMA
Tensor
Cores

V100
125.00

IMMA
Int8 Tensor
Cores

Q8000

261.00

A100

Structured
Sparsity

1248.00

H100
FP8
Transformer
Eng

4000.00

# AI IN SCIENCE & ENGINEERING



**Inverse & Data Assimilation Problems**

Climate

Medical Imaging

Oil & Gas

High Energy/ Nuclear Physics

**Improved Physics & Predictions**

Radiation

Micro-mechanical Material Model

Turbulence

Molecular Dynamics

**Real Time Simulations**

Robotics

Digital Twin

Autonomous Ride & Handling

Games

**Digital Design & Manufacturing**

Heat Sink

Aerodynamics

Vias on a PCB

*Physics & Data - No Traditional Solver*

*Physics - Traditional Solver (Speed is a limitation)*

# DEEP LEARNING IN SCIENTIFIC PROBLEMS

## Modulus is an SDK to build and deploy ML/DL applications for Physical Systems

# MULTI-PHYSICS SIMULATION

## Conjugate Heat Transfer – No Training Data



**Fluid Temperature**

**Solid Temperature**

# FourCastNetv2: NVIDIA's Global Weather Simulator.

**Fully data-driven weather prediction.**

- Scope                      Global, Medium Range
- Model Type                 Full-Model AI Surrogate
- Architecture               Spherical Fourier Neural Operators
- Resolution:                25km
- Training Data:             ERA5 Reanalysis
- Initial Condition          GFS / UFS
- Speedup vs NWP             5000x
- Power Savings              $O(10^4)$

# INVERSE PROBLEM

## Finding Unknown Coefficients of a PDE: Heat Sink



**Fluid Heat Convection:**

$$0 = \nabla \cdot (D_{fluid}\nabla\theta_{fluid}) - \nabla \cdot (U\theta_{fluid}) \qquad D_{fluid} = \frac{k_{fluid}}{\rho_{fluid}c_{pfluid}}$$

**Solid Heat Conduction:**

$$0 = \nabla \cdot (k_{solid}\nabla\theta_{solid}) \qquad D_{solid} = \frac{k_{solid}}{\rho_{solid}c_{psolid}}$$

**Interface Conditions:**

$$\theta_{solid} = \theta_{fluid}$$

$$k_{solid}(N \cdot \nabla\theta_{solid}) = k_{fluid}(N \cdot \nabla\theta_{fluid})$$

**Results:**

| Property | OpenFOAM (True) | Modulus (Predicted) |
|---|---|---|
| Kinematic Viscosity $(m^2/s)$ | $1.00 \times 10^{-2}$ | $1.03 \times 10^{-2}$ |
| Thermal Diffusivity $(m^2/s)$ | $2.00 \times 10^{-3}$ | $2.19 \times 10^{-3}$ |

# Summary

- Deep learning was enabled by hardware and is paced by hardware
  - 1000x performance in last decade
- DL has a many roles in science
  - Simulation
  - Inverse problems (data interpretation)
  - Learned behavior (e.g., protein folding)
- FourCastNet
  - 5000x speed of ECMWF ICON model
  - Comparable accuracy
- LLMs can increase productivity of scientific process
- DL should be a major element of every scientist's toolbox