

Seismic Imaging at Scale: The Future of HPC and AI - A NVIDIA Perspective

João Paulo Navarro - Sr. Solutions Architect

jnavarro@nvidia.com

What Is NVIDIA Energy Samples

GPU Accelerated Building Blocks For Expediting Exploration Algorithms,
Maximizing Performance Lowering Operational Costs For Energy Customers

Customer/ISV Applications

Quickly evaluate and integrate GPU accelerated techniques into applications in the cloud and on-prem

Energy SDK

Energy sped-of-light CUDA-first sample algorithms for on-prem and cloud such as **FWI**, **RTM**, Kirchhoff, ResSim, imaging, etc.

Compression and upscaling techniques for handling and pipelining subsurface, surface, edge, and far edge data

Physics ML & AI models to support digital transformation and energy transition

HPC SDK

RAPIDS

PhysicsNeMo

Integrate the best of HPC SDK, RAPIDS, PhysicsNeMo, and GPU Direct Storage (GDS)



Leverage the GPU architectural features and tested on various system configs for on-prem / cloud deployment

Seismic Imaging

Use Cases

- Full Waveform Inversion (FWI)
- Reverse Time Migration (RTM)
- Kirchhoff Depth Migration
- Surface-related multiple elimination (SRME)

Challenges

Massive workloads that require significant memory for processing and interpretation

Solution

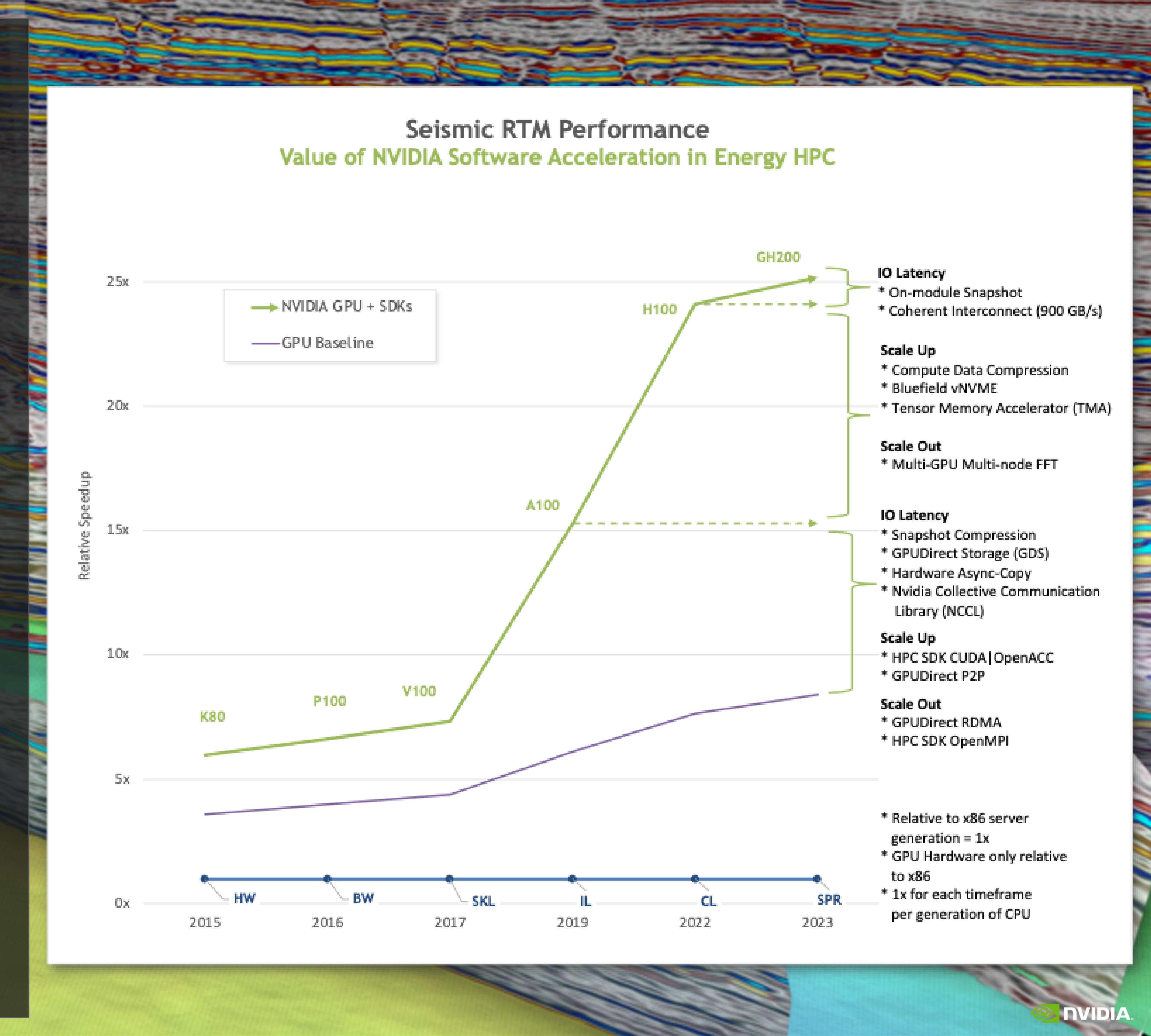
- Internal code development built to be optimized on the GPU
- Solutions provided under NDA to partners, continued development and optimization with their development teams

NVIDIA Solution Stack

- Hardware: A100, H100, H200, GH200
- Software: Energy SDK, HPC SDK

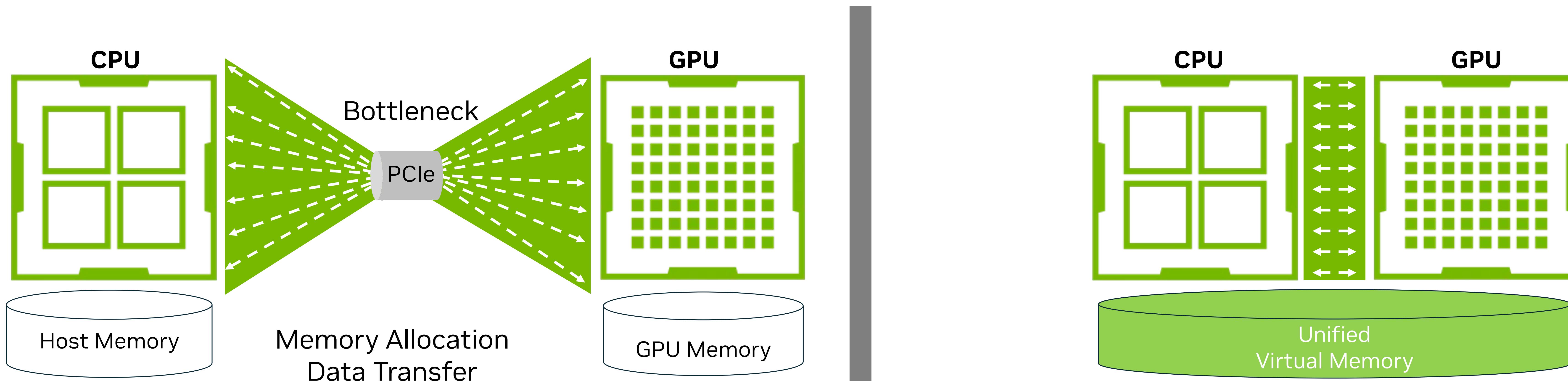
Outcomes

- 9x relative speedup of seismic + GPU
- 25x relative speedup of seismic + GPU + SDK
- Software-defined solution allows for significant acceleration of seismic interpretation
- Key ISVs include SLB, Shearwater, AspenTech (former Emerson)



Challenges with Traditional Accelerated Systems

Bridging the GAP: Unifying CPU and GPU Memory



- ❑ PCIe bottlenecks CPU-GPU communication
- ❑ GPU can't access CPU memory directly
- ❑ Complex Programmability
- ❑ Low CPU Performance per Watt

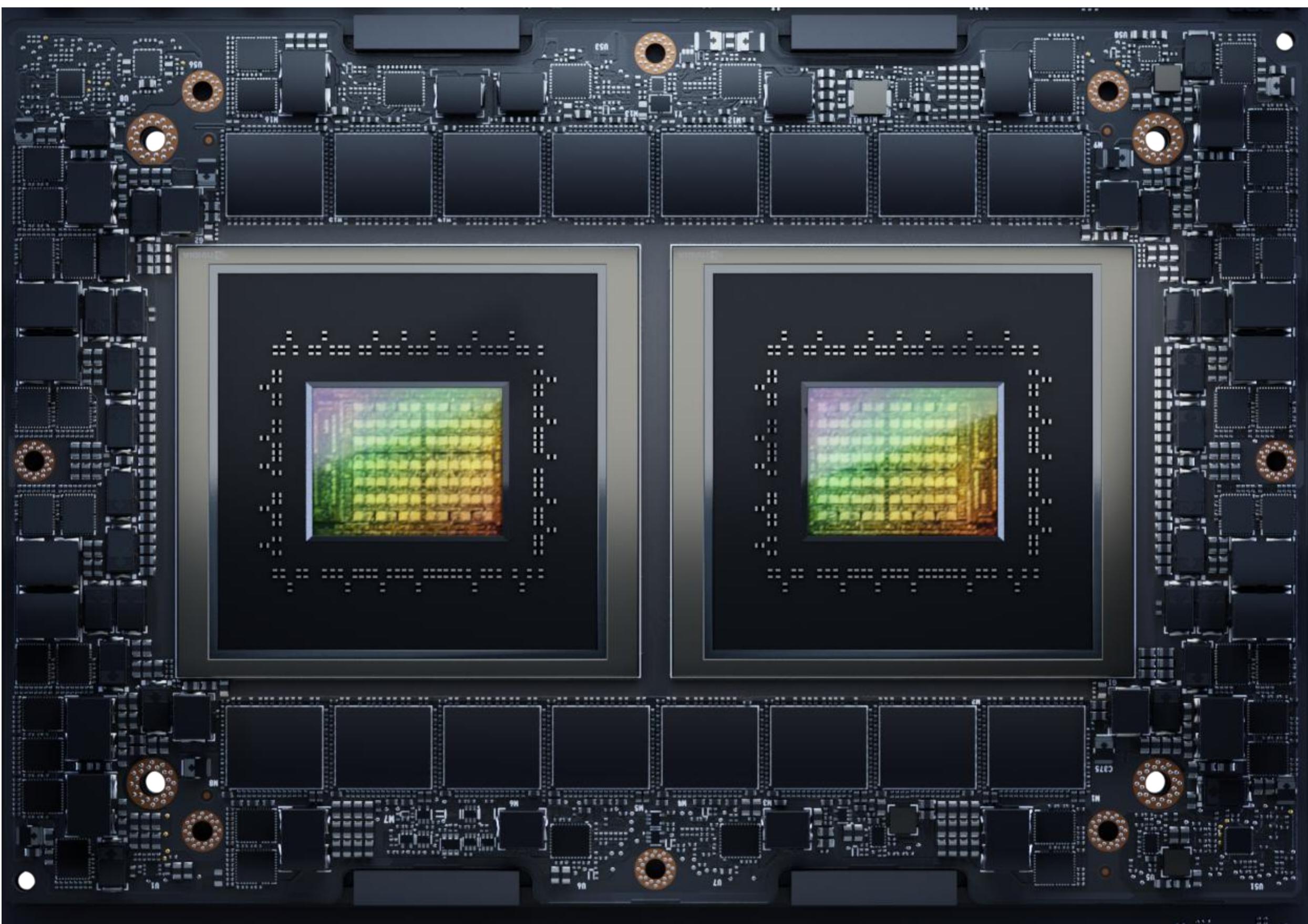
Legacy PCIe
Architecture

We need a new chip
architecture:

- ✓ Accelerate CPU-GPU data transfers
- ✓ Coherent CPU and GPU cache memory
- ✓ Simplify Programmability
- ✓ Deliver more performance per Watt

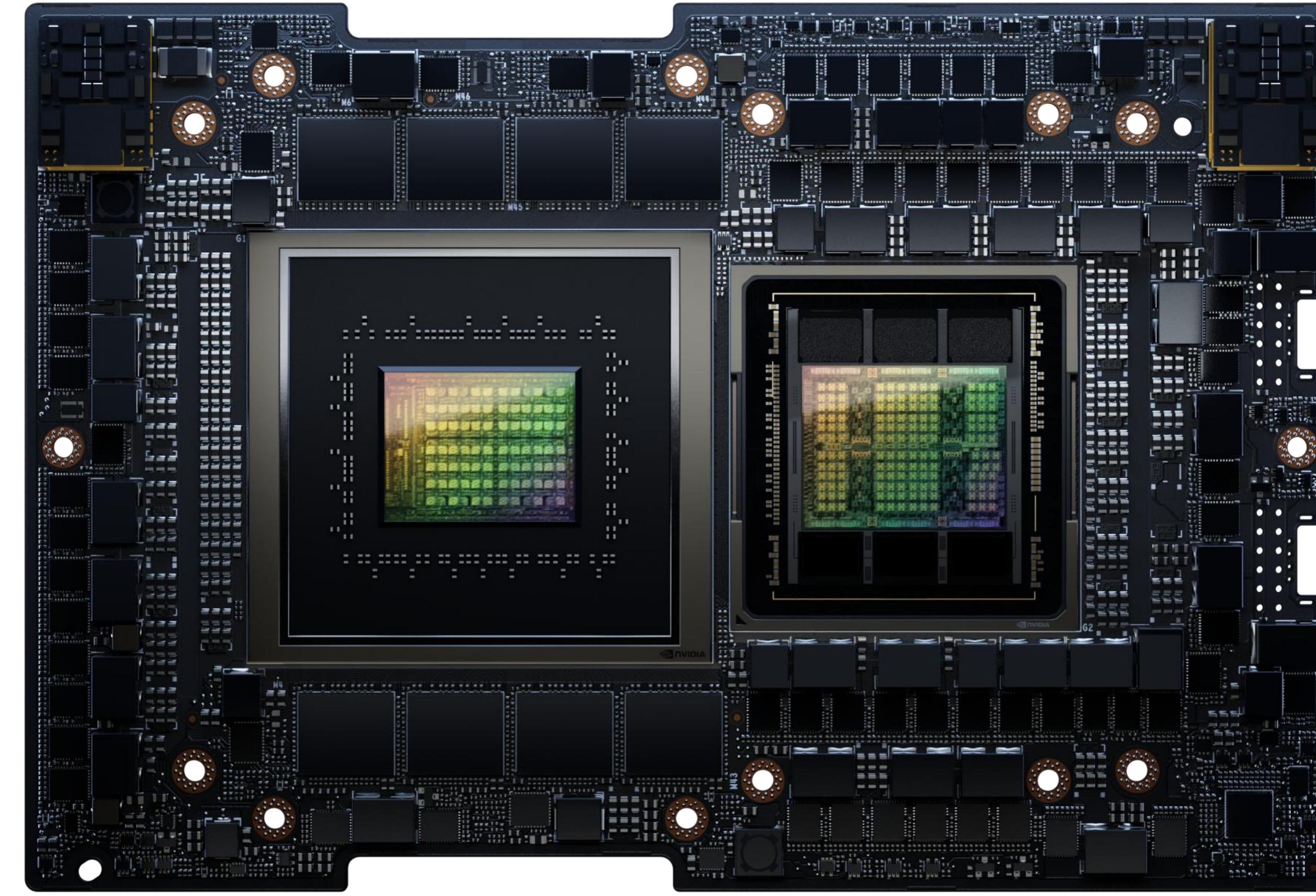
NVIDIA Grace for Cloud, AI and HPC Infrastructure

Grace CPU Superchip
CPU Computing



CPU-based applications where absolute performance, energy efficiency, and data center density matter, such as scientific computing, data analytics, enterprise and hyperscale computing applications

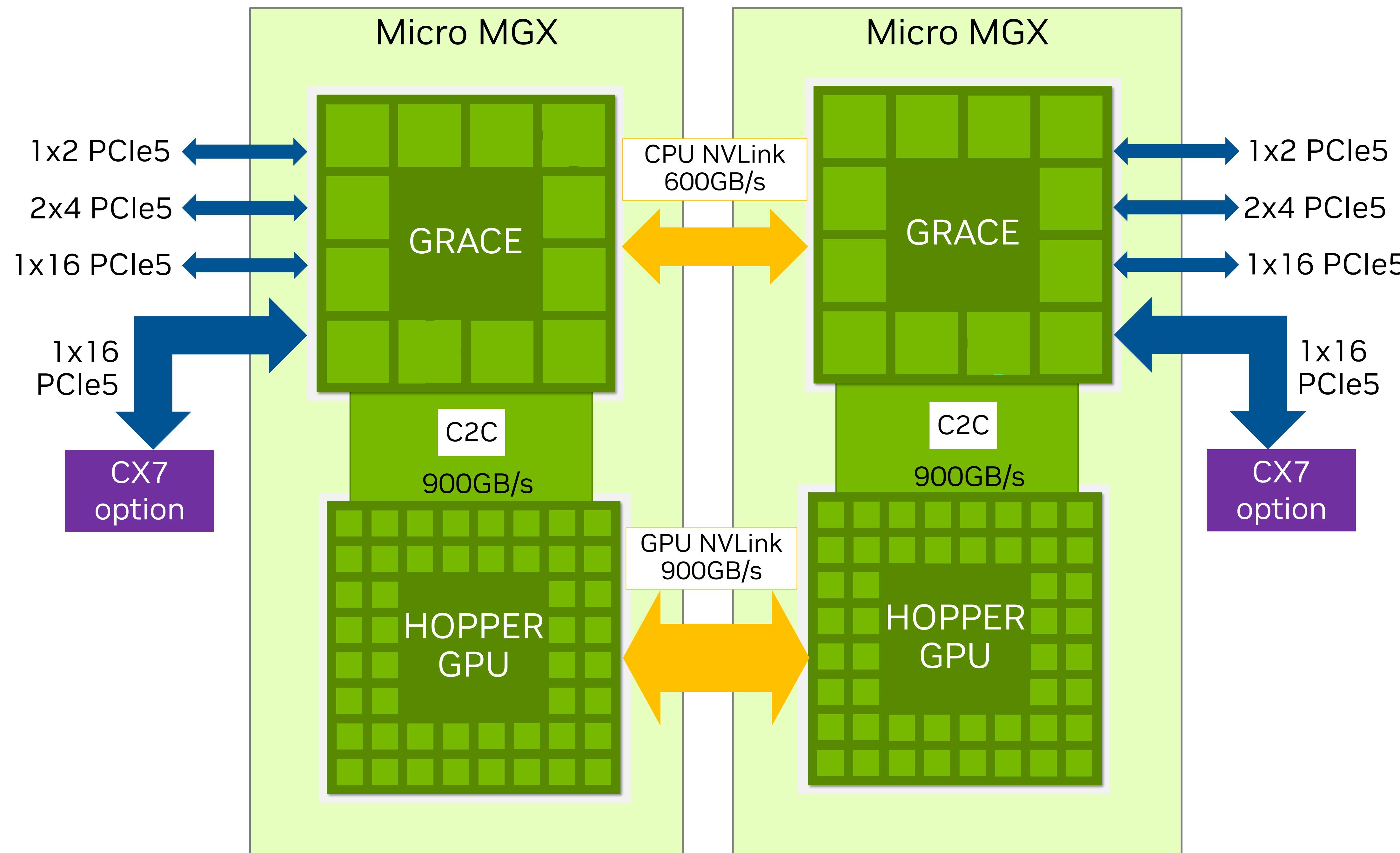
GH200 Grace Hopper Superchip
Large Scale AI & HPC



Accelerated applications where CPU performance and system memory size and bandwidth are critical; tightly coupled CPU & GPU for flagship AI & HPC. Most versatile compute platform for scale out.

GH200 NVL2 Node Architecture

Single Large GPU for Era of Accelerated Computing and Gen AI



- Single node, single OS
- Unified address space across 2 CPU + 2 GPU
- Each GPU is cache coherent to the dual CPU
- Option for up to 4x CX7
- 8 PF of AI performance
- 10 TB/s of GPU memory bandwidth and up to 1.2TB total fast memory
- Best performance for compute and memory intensive workloads
 - Mainstream LLM
 - Data Processing
 - RAG

How to Go Further With the Same Chip?

GPU: A Processor Made of Processors

Multiple Layers of Specialized Cores Handle Diverse Data Types

H100 Chip Overview



NVIDIA H100 SXM5 ¹	
Peak FP64 ¹	30 TFLOPS
Peak FP64 Tensor Core ¹	60 TFLOPS
Peak FP32 ¹	60 TFLOPS
Peak FP16 ¹	120 TFLOPS
Peak BF16 ¹	120 TFLOPS
Peak TF32 Tensor Core ¹	500 TFLOPS 1000 TFLOPS ²
Peak FP16 Tensor Core ¹	1000 TFLOPS 2000 TFLOPS ²
Peak BF16 Tensor Core ¹	1000 TFLOPS 2000 TFLOPS ²
Peak FP8 Tensor Core ¹	2000 TFLOPS 4000 TFLOPS ²
Peak INT8 Tensor Core ¹	2000 TOPS 4000 TOPS ²

Table 1. NVIDIA H100 Tensor Core GPU preliminary performance specs

Mixed-Precision Wave Propagation

Physics vs. Limitations of FP16

Discretization of the Elastic Wave Equation

• $v_p = 3500 \text{ m.s}^{-1}$
• $v_s = 2000 \text{ m.s}^{-1}$
• $\rho = 2000 \text{ kg.m}^{-3}$
• $\Delta t = 0.001 \text{ s}$

$$v_x = v_x + \frac{\Delta t}{\rho} (D_x \sigma_{xx} + D_y \sigma_{xy} + D_z \sigma_{xz})$$

$$\sigma_{xx} = \sigma_{xx} + \Delta t (\lambda + 2\mu) D_x v_x + \Delta t \lambda (D_y v_y + D_z v_z) + \Delta t s_{xx}$$

$24 \times 10^6 \text{ Pa.s}$ (overflow in FP16)

$5 \times 10^{-7} \text{ s.m}^3.\text{kg}^{-1}$ (subnormal value for FP16)

Data Type	Exponent	Mantissa	Min (normal)	Min (subnormal)	Max
FP16	5	10	6.104×10^{-5}	5.961×10^{-8}	$65,504$
FP32	8	23	1.176×10^{-38}	1.401×10^{-45}	3.403×10^{38}

Rescaling the Wave Equation

Compute normalized quantities

$$S = \max(s_{ij})$$

Maximum of source signal

$$C = \max(C_{ij})$$

Maximum of stiffness coefficient

Scale wavefields, sources, and material parameters

$$\sigma'_{ij} = \frac{1}{\Delta t S} \sigma_{ij}$$

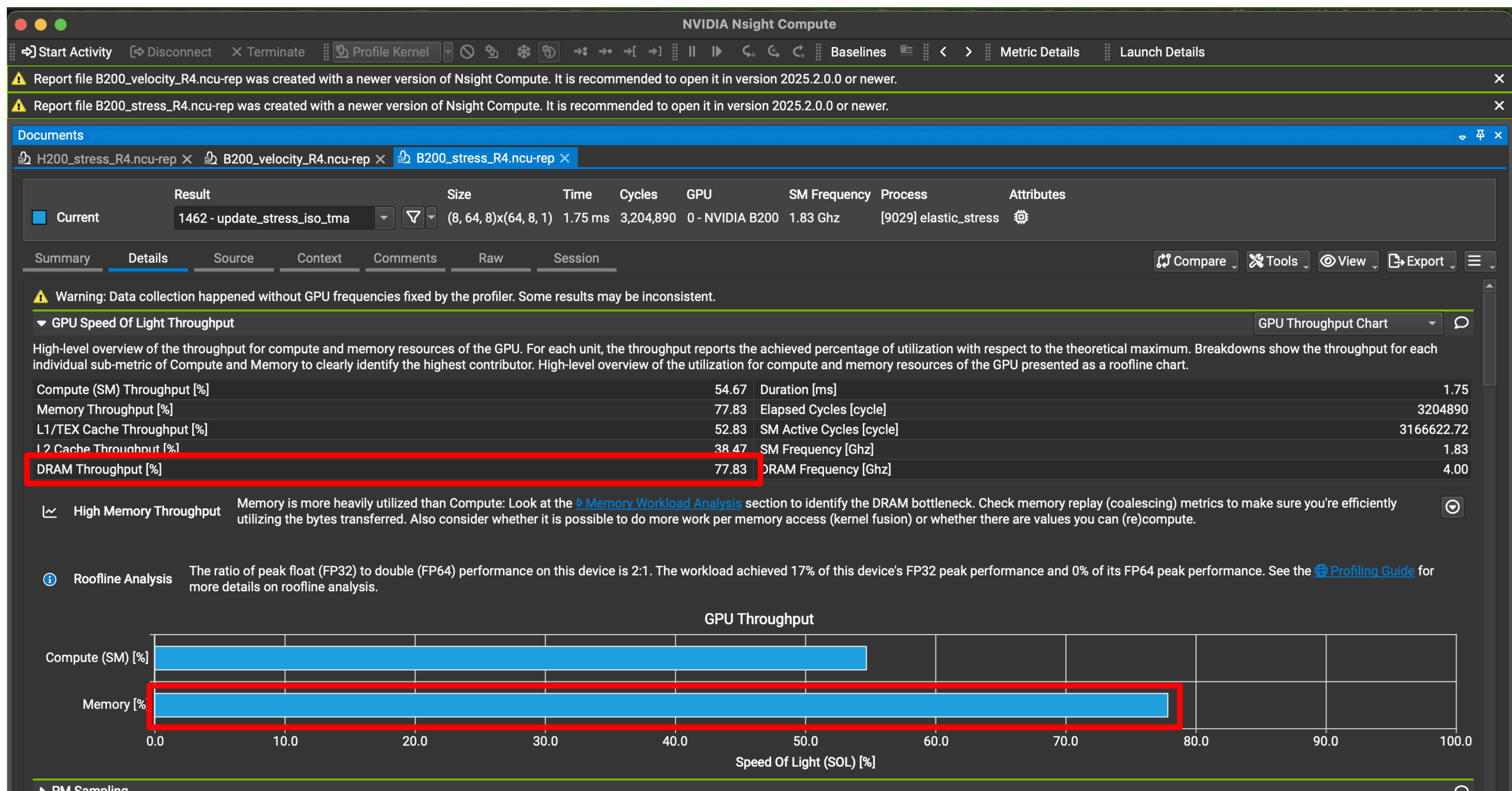
$$v'_{ij} = \frac{C}{S} v_{ij}$$

$$s'_{ij} = \frac{1}{\Delta t S} s_{ij}$$

$$C'_{ij} = \frac{1}{\Delta t C} C_{ij}$$

$$\rho'_{ij} = \frac{1}{\Delta t / C} \rho$$

Motivation and Strategy (Memory Bound)

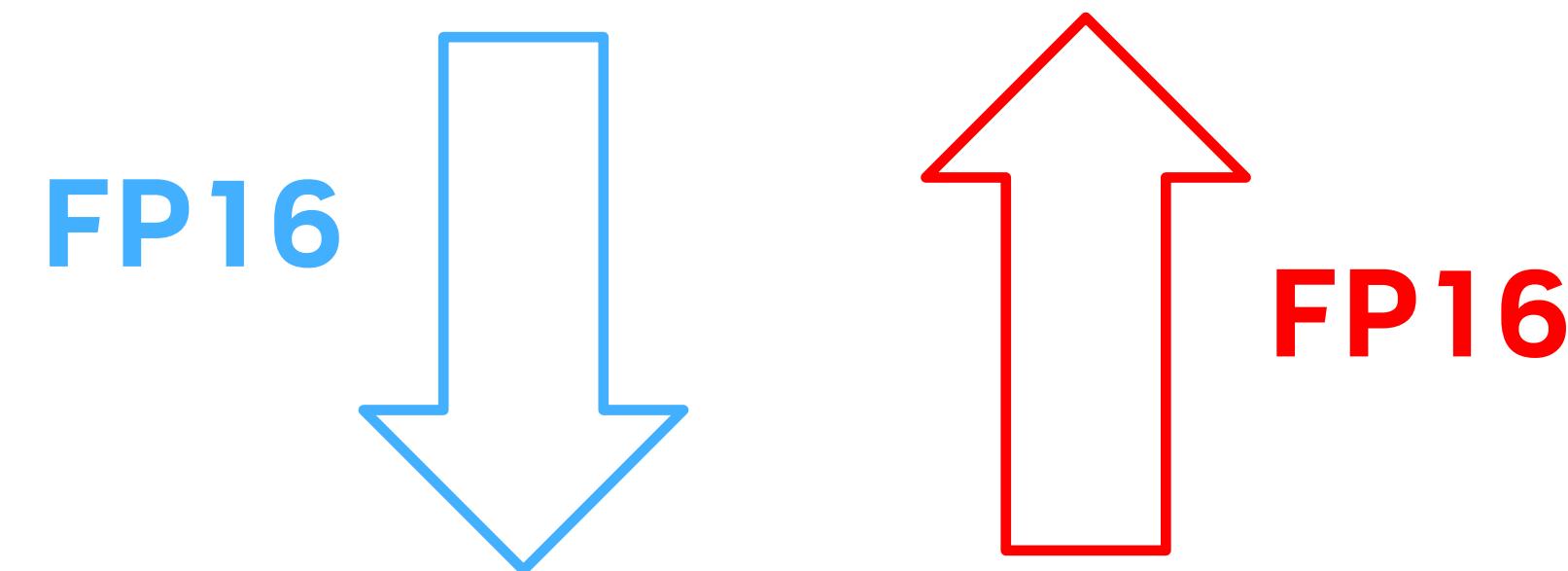


⚠️ Elastic kernels are **Bandwidth Bound!**

GPU global memory / L2 cache

Wavefields (FP16)

Elastic models*



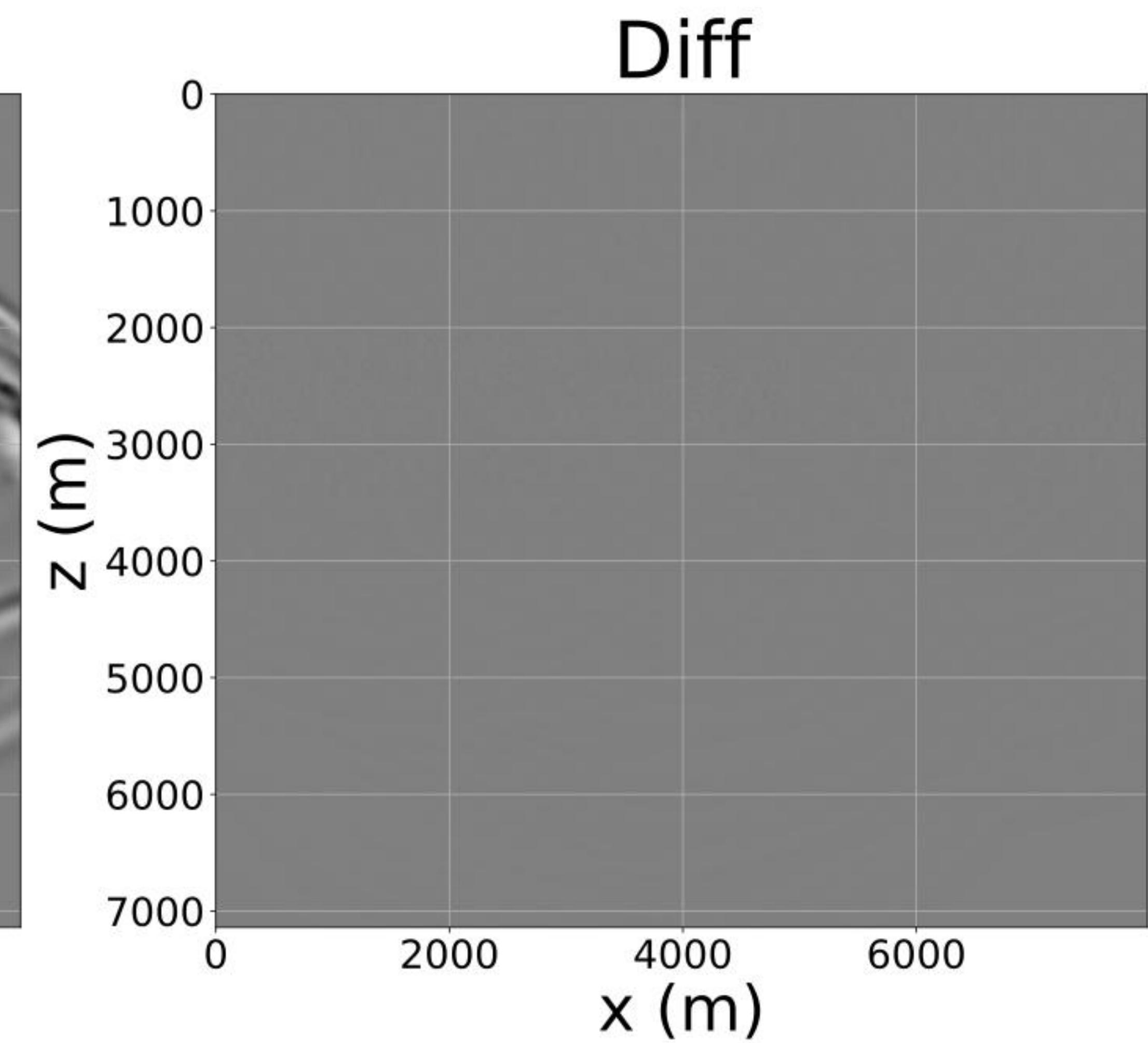
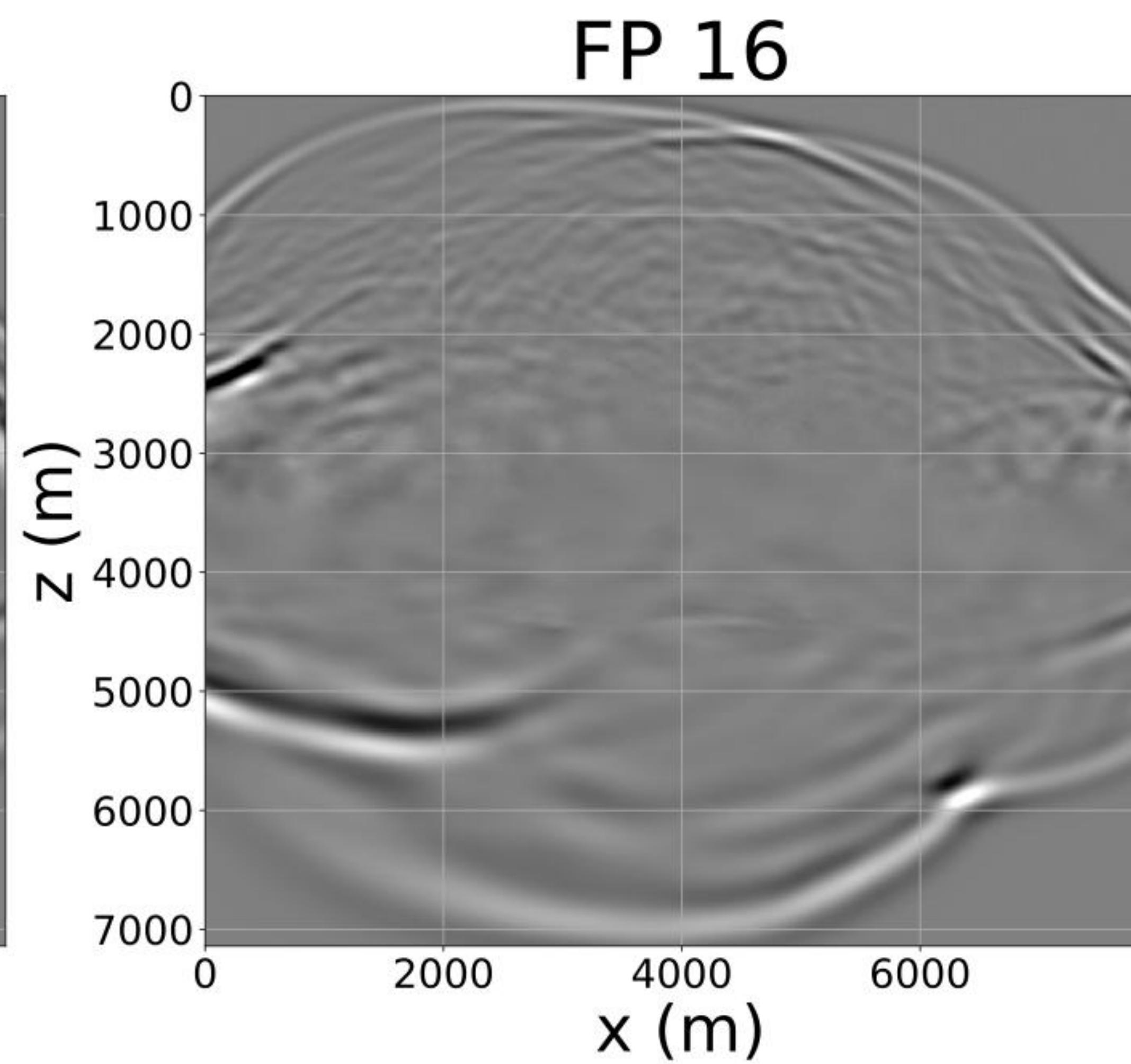
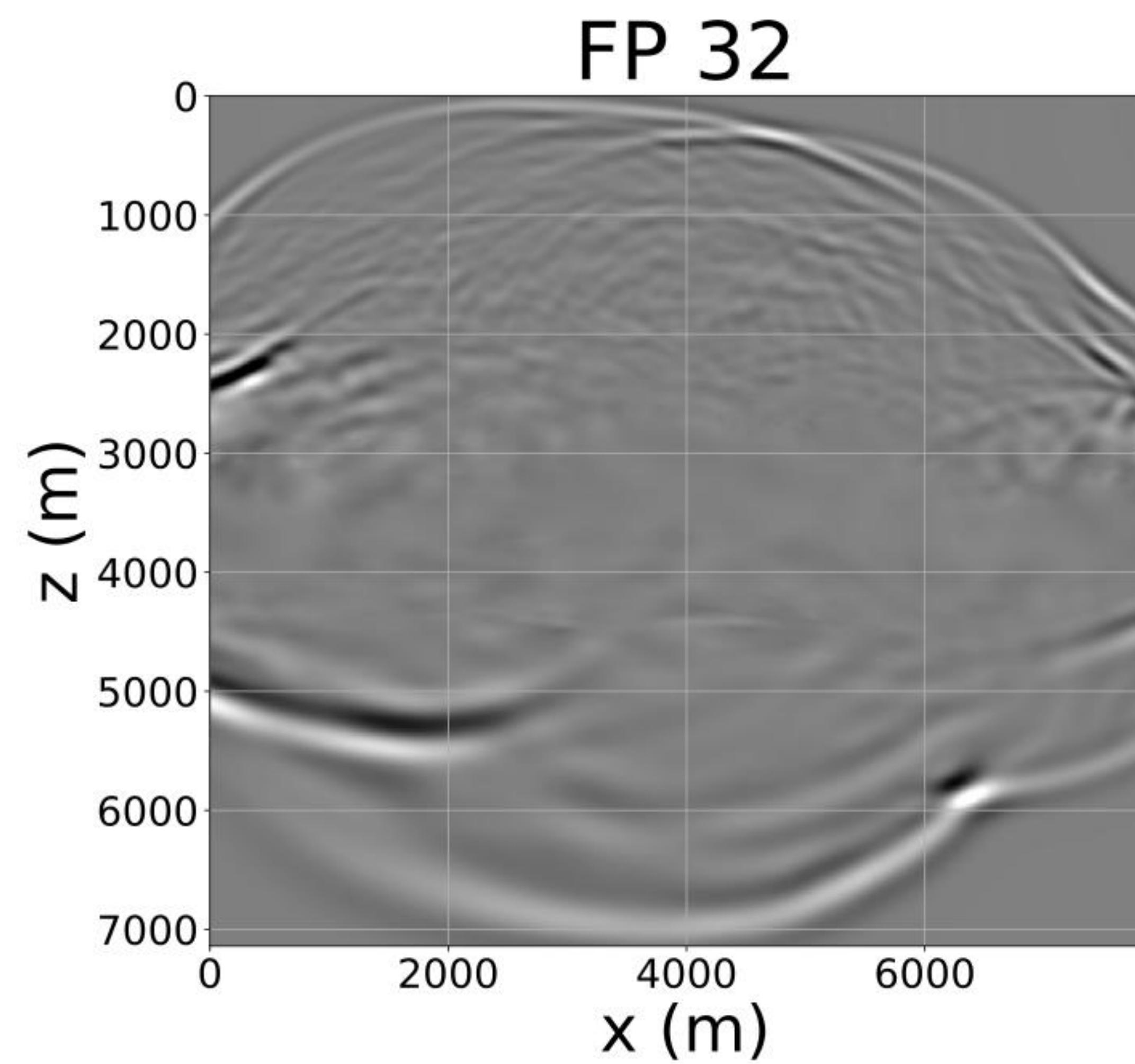
- Load wavefields and elastic models*
- Cast wavefields and elastic models* to FP32
- Update wavefield using FP32 arithmetic
- Store results in FP16

Kernel

* Elastic models may be stored in FP32 or compressed and stored in lower precision

Experiment: Heterogeneous VTI

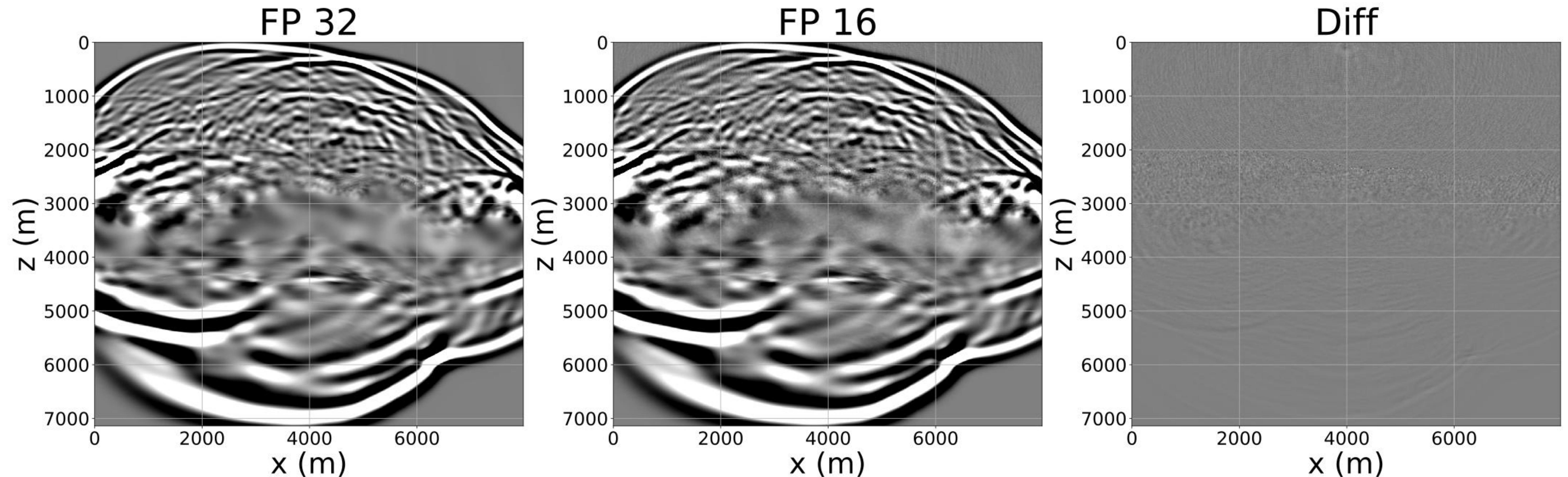
Results and Validation



Pressure snapshots
Clip=10%

Experiment: Heterogeneous VTI

Results and Validation

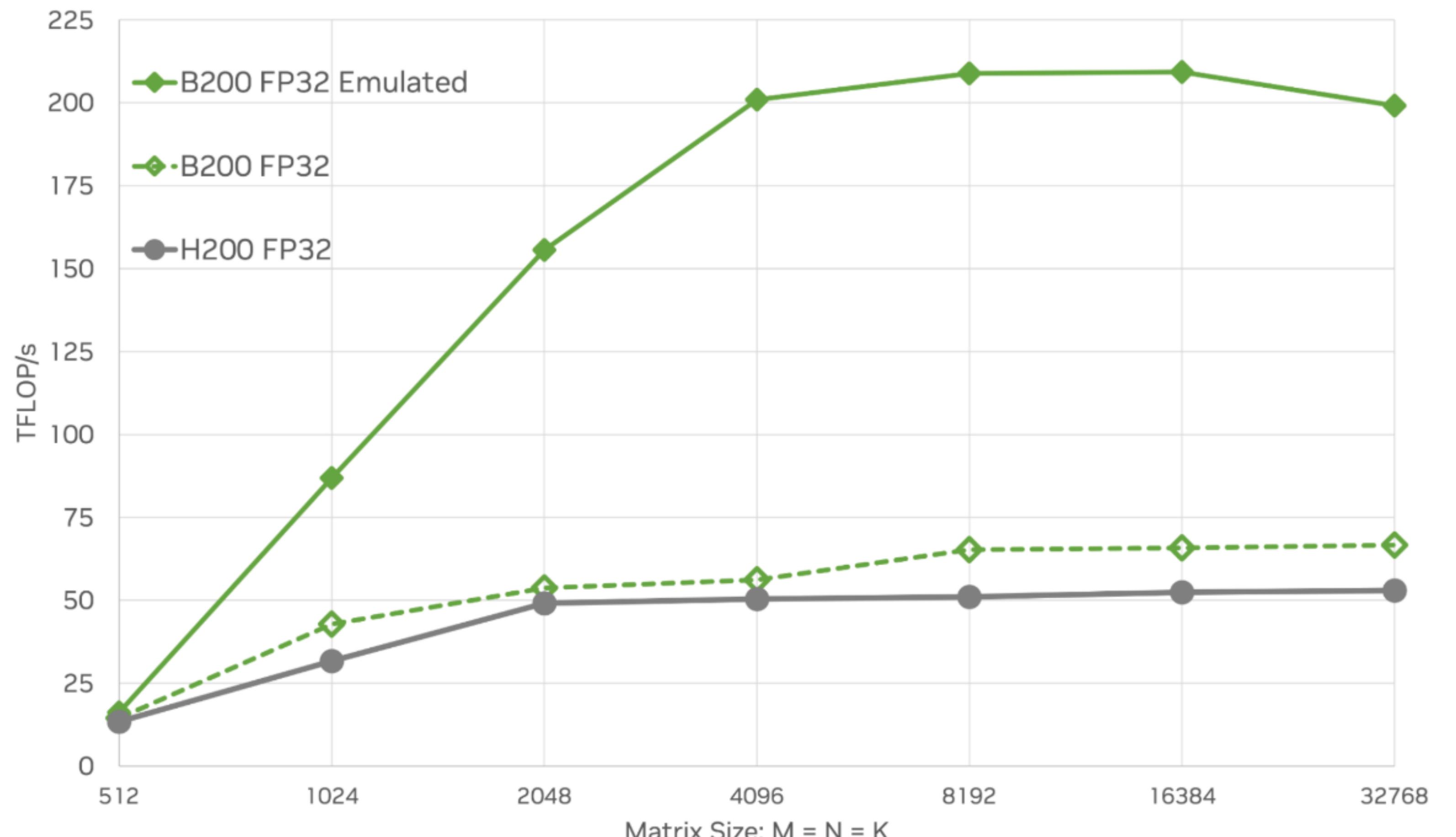


Pressure snapshots
Clip=1%

- We maintained the geophysical accuracy required for modeling and FWI gradients
 - Reducing the memory footprint **by 50%** and unlocking a potential 2x speedup

Floating Point Emulation Methods

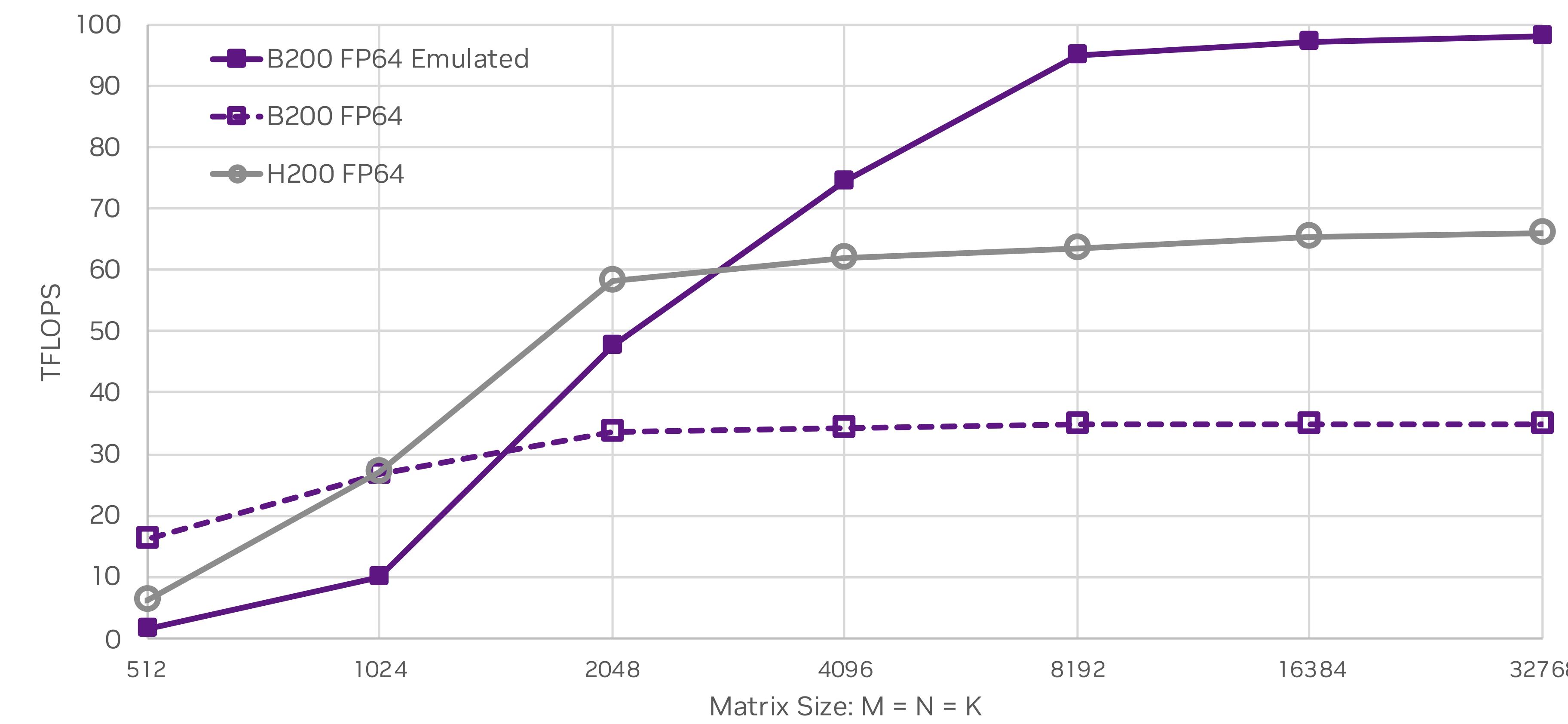
Floating Point Emulation Methods



Performance (TFLOP/s) of emulated FP32 compared to native FP32 on B200 and H200 GPUs

- **FP32 using BF16 Tensor Cores¹**

- Released with **CUDA 12.9 in cuBLAS** for Blackwell GPUs
- Being tested for accuracy and performance impact in:
 - Weather, quantum circuit and condensed matter simulations
 - Dense Linear Algebra (QR, LU)
- **Uses 9 inner matmuls in BF16 (BF16x9)**



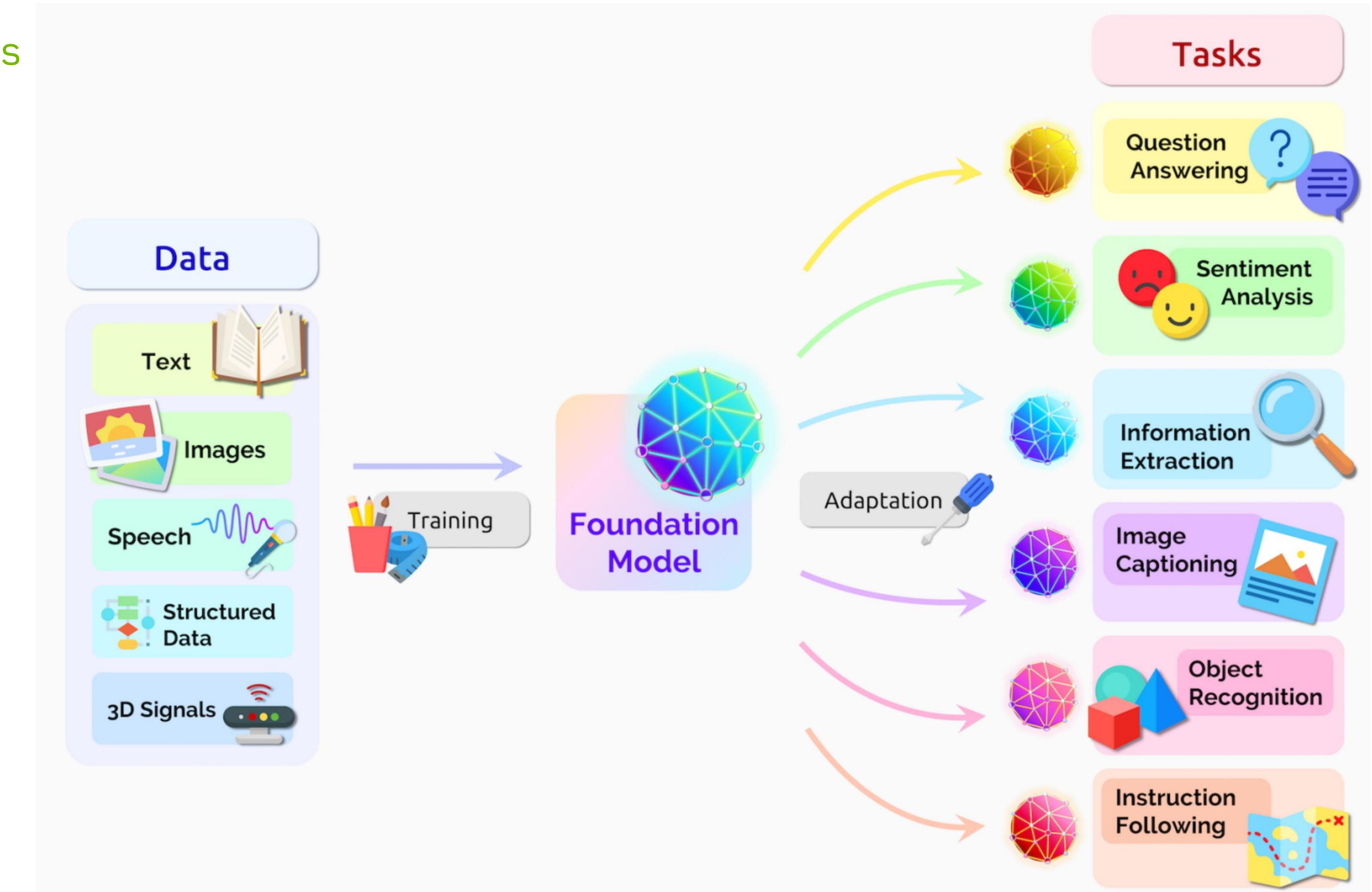
FP64 Emulation using INT8 tensor cores (s=7 or 56 bits)

AI Perspective

What are Foundation Models

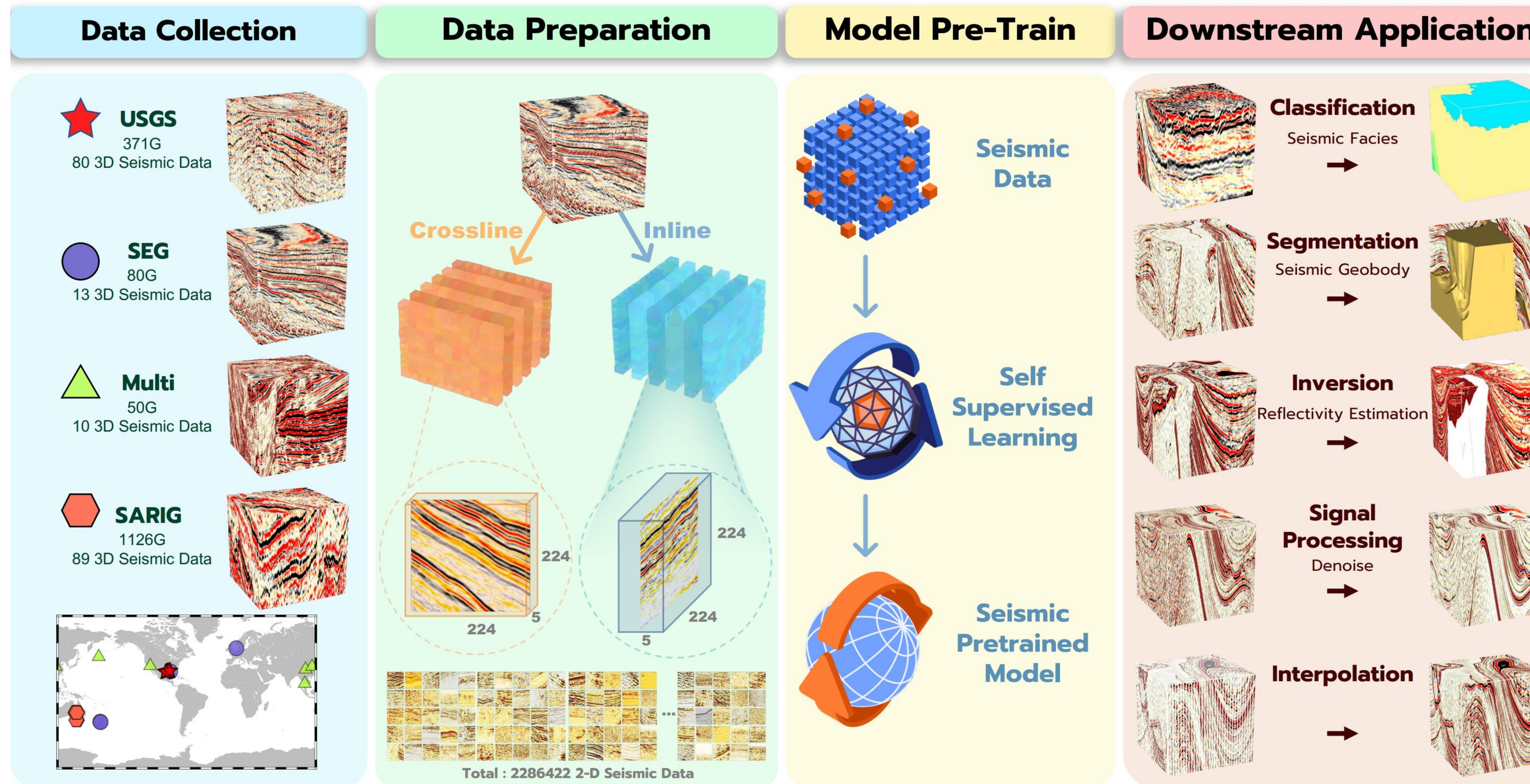
General purpose machine learning models

- Trained on massive unlabeled datasets to handle a wide variety of jobs from translating text to analyzing seismic images.
- General purpose adaptable models that can be fine-tuned for diverse tasks like computer vision, natural language processing with minimal task-specific training.
- Scalable and reusable architectures which leverage learned patterns to perform new tasks efficiently, reducing the need for extensive retraining or custom model development.



Source: [What Are Foundation Models? | NVIDIA Blogs](#)

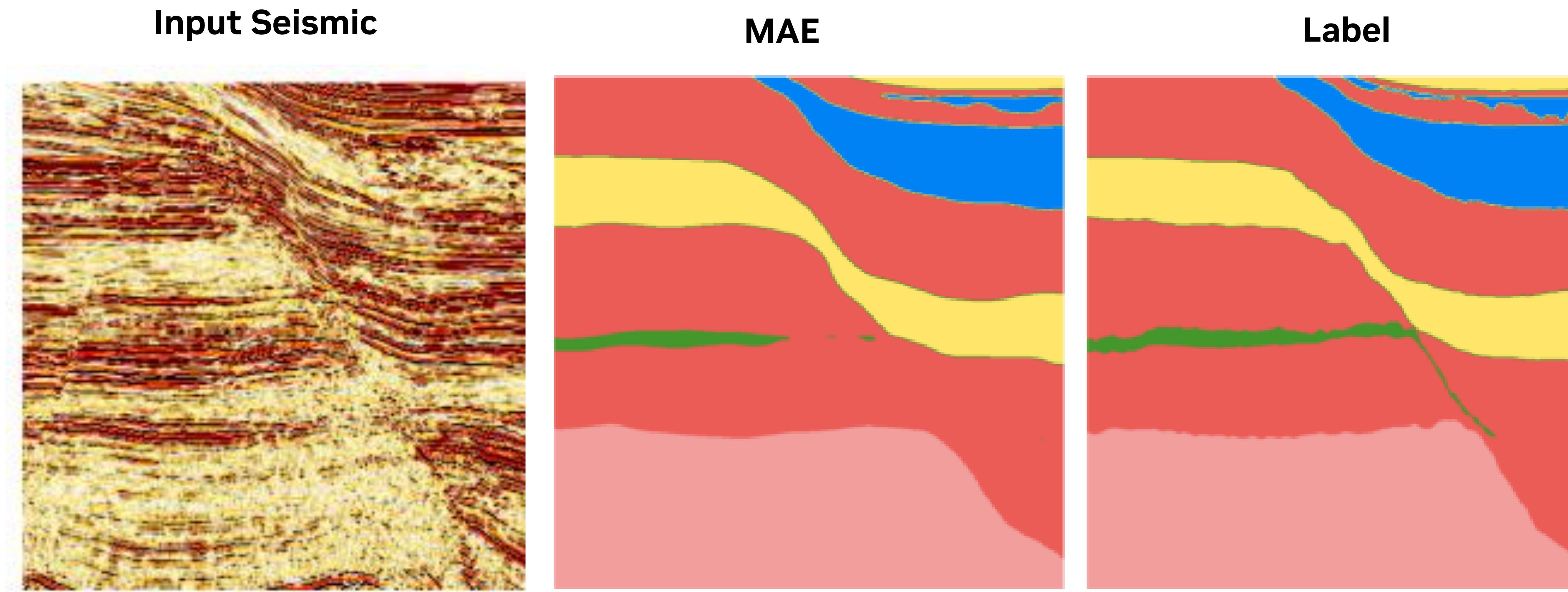
Seismic Foundation Model (SFM)



SFM represents the first successful large-scale foundational model trained exclusively within the seismic domain [[reference](#)].

Masked Auto-Encoder fine-tuned on 200 samples

Trained with self-supervised learning on 2.2 million seismic images utilizing NVIDIA TAO



Facies classification

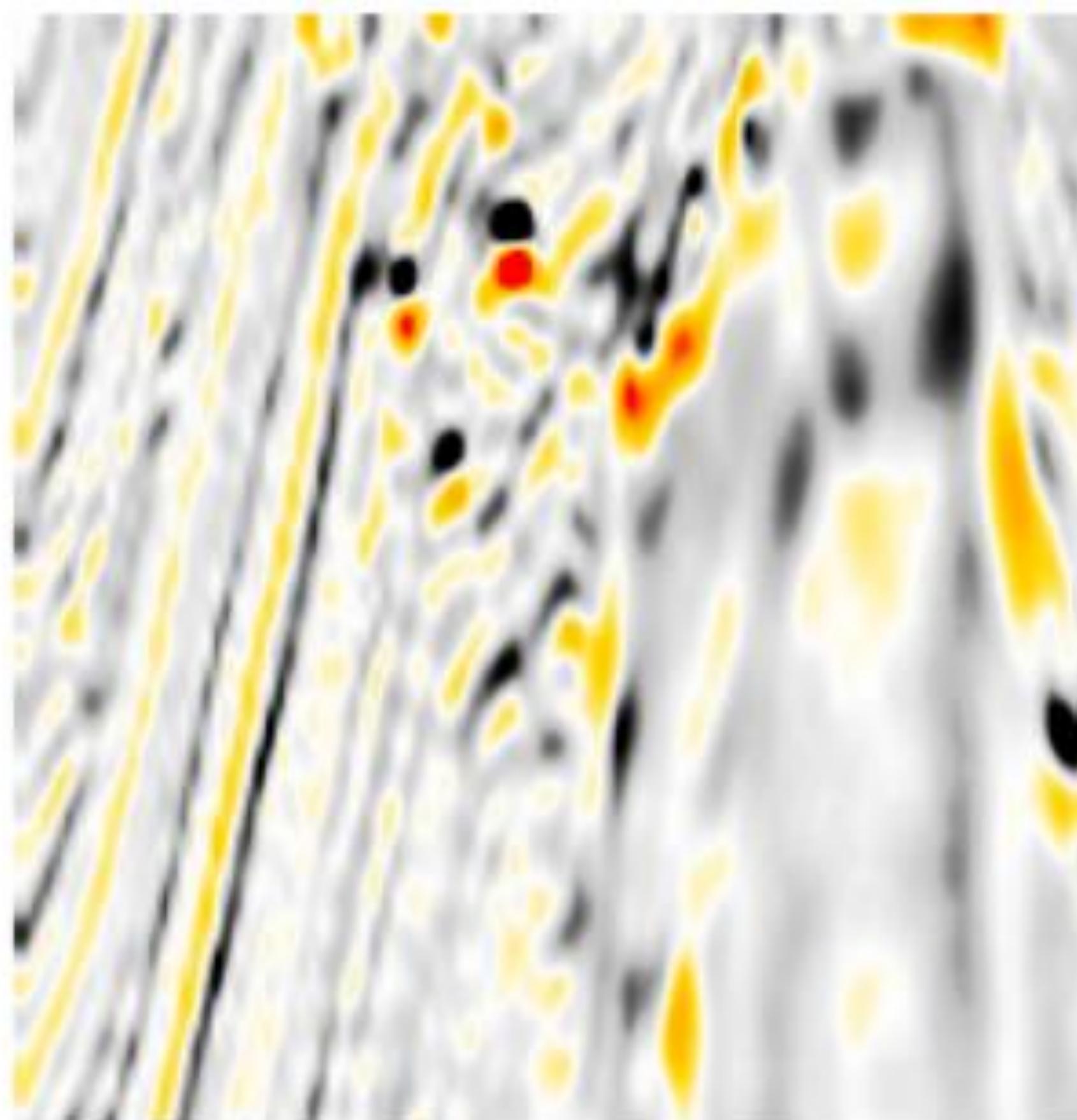
(Results reproduced from [Sheng et al., 2025](#))

“Seismic Foundation Model (SFM): a new generation deep learning model in geophysics

Masked Auto-Encoder fine-tuned on 3000 samples

Trained with self-supervised learning on 2.2 million seismic images utilizing NVIDIA TAO

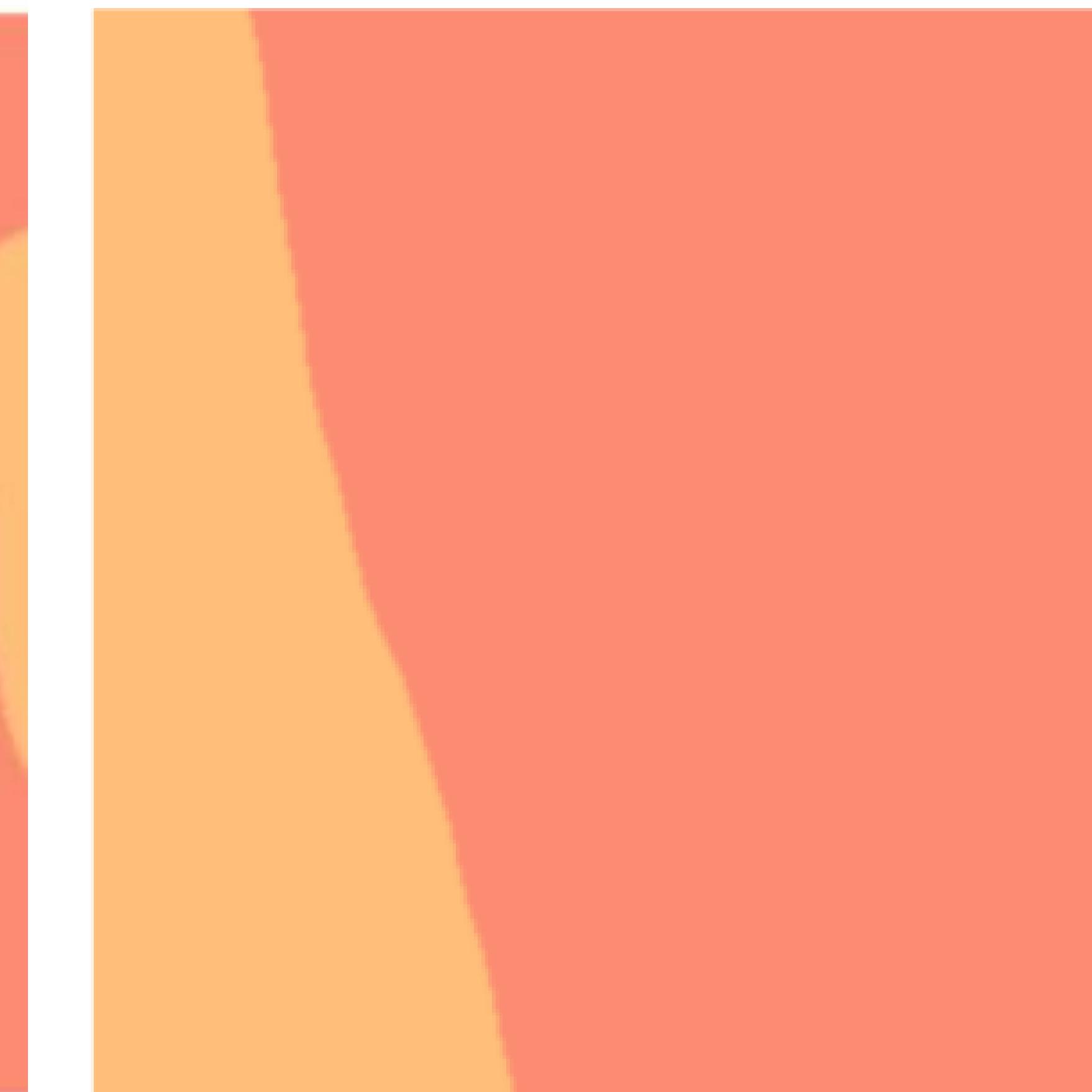
Input Seismic



U-Net



MAE



Label



Salt Segmentation

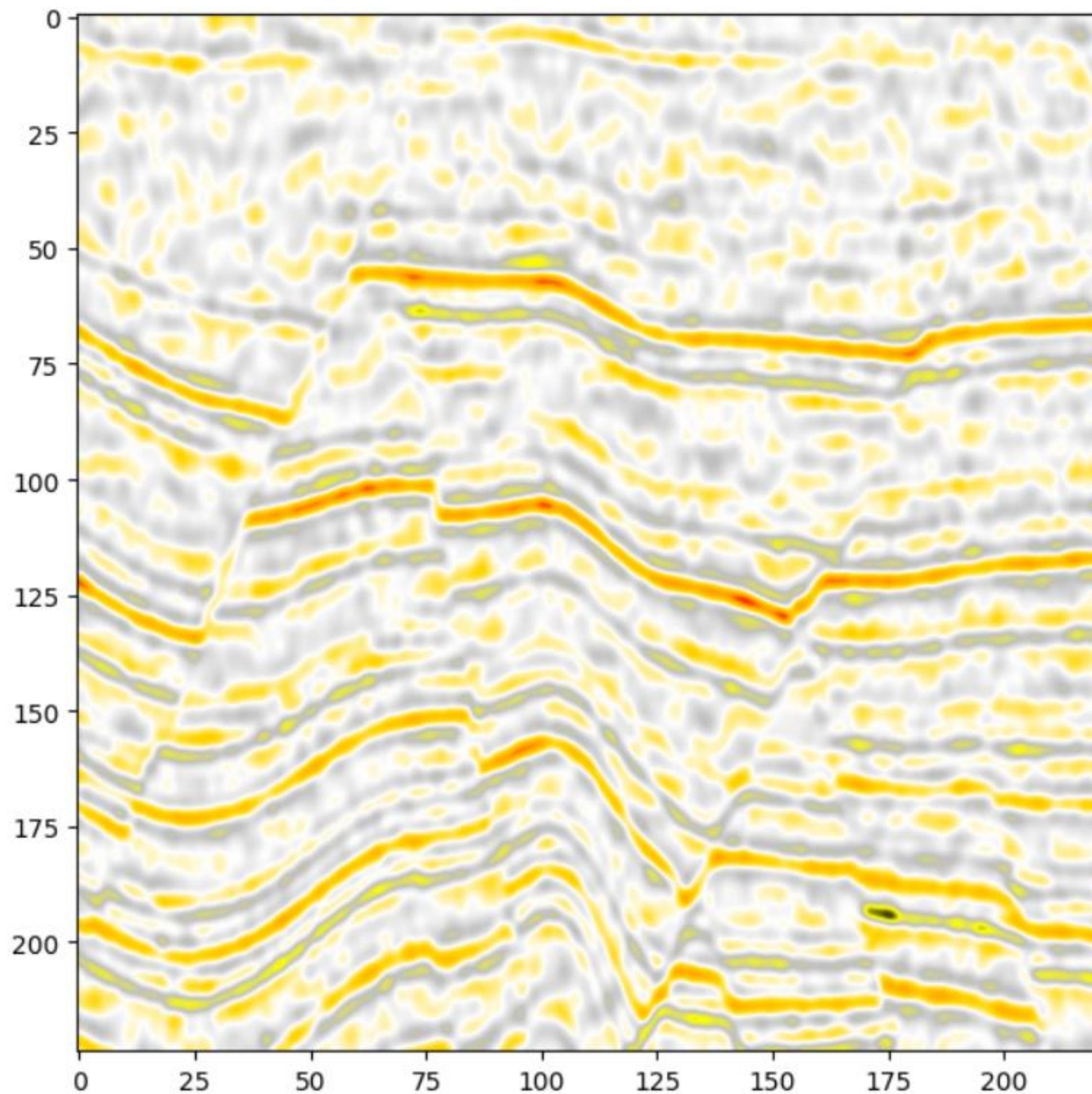
(Results reproduced from [Sheng et al., 2025](#))

“Seismic Foundation Model (SFM): a new generation deep learning model in geophysics

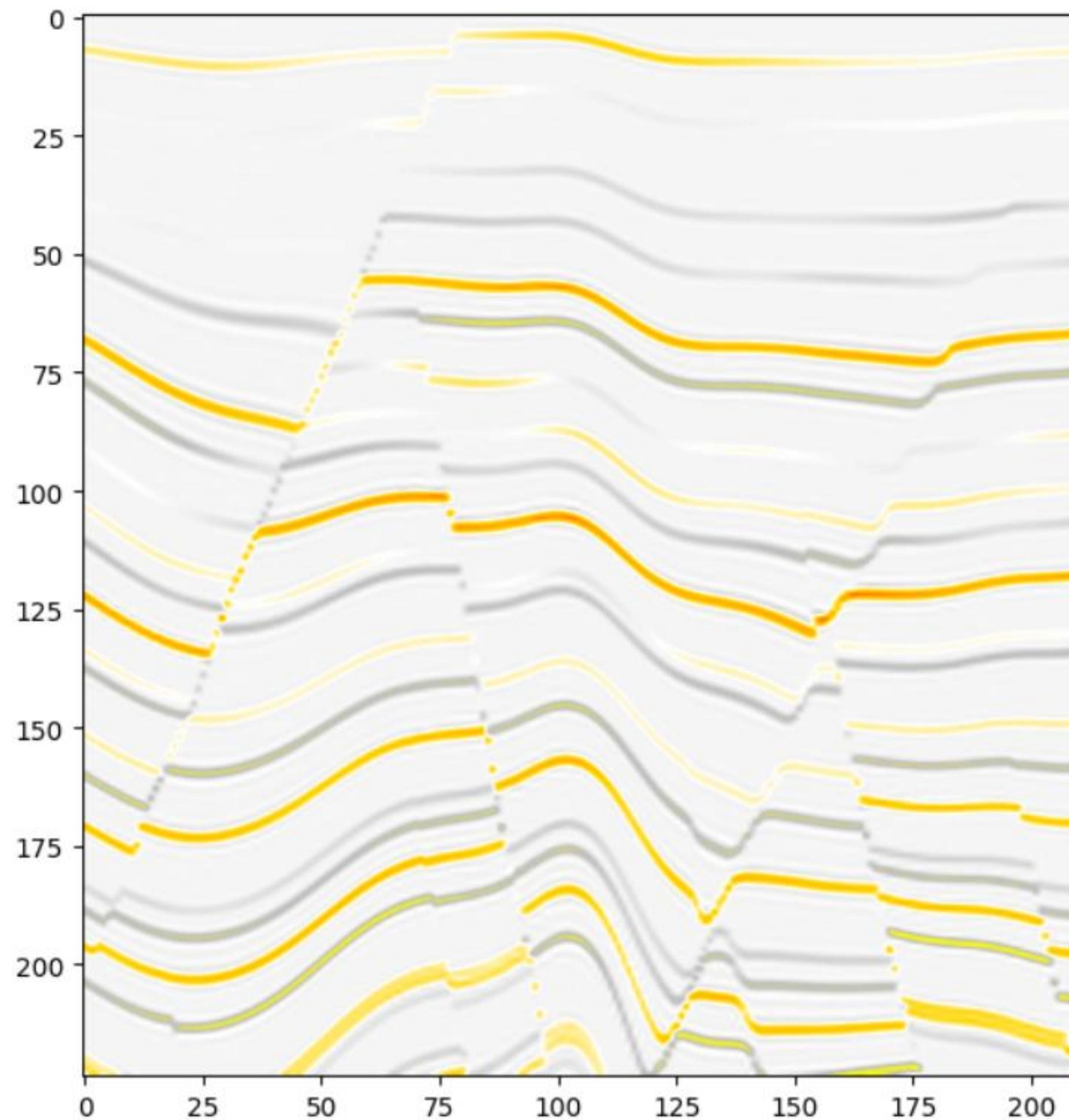
Masked Auto-Encoder fine-tuned on 1000 samples

Trained with self-supervised learning on 2.2 million seismic images utilizing NVIDIA TAO

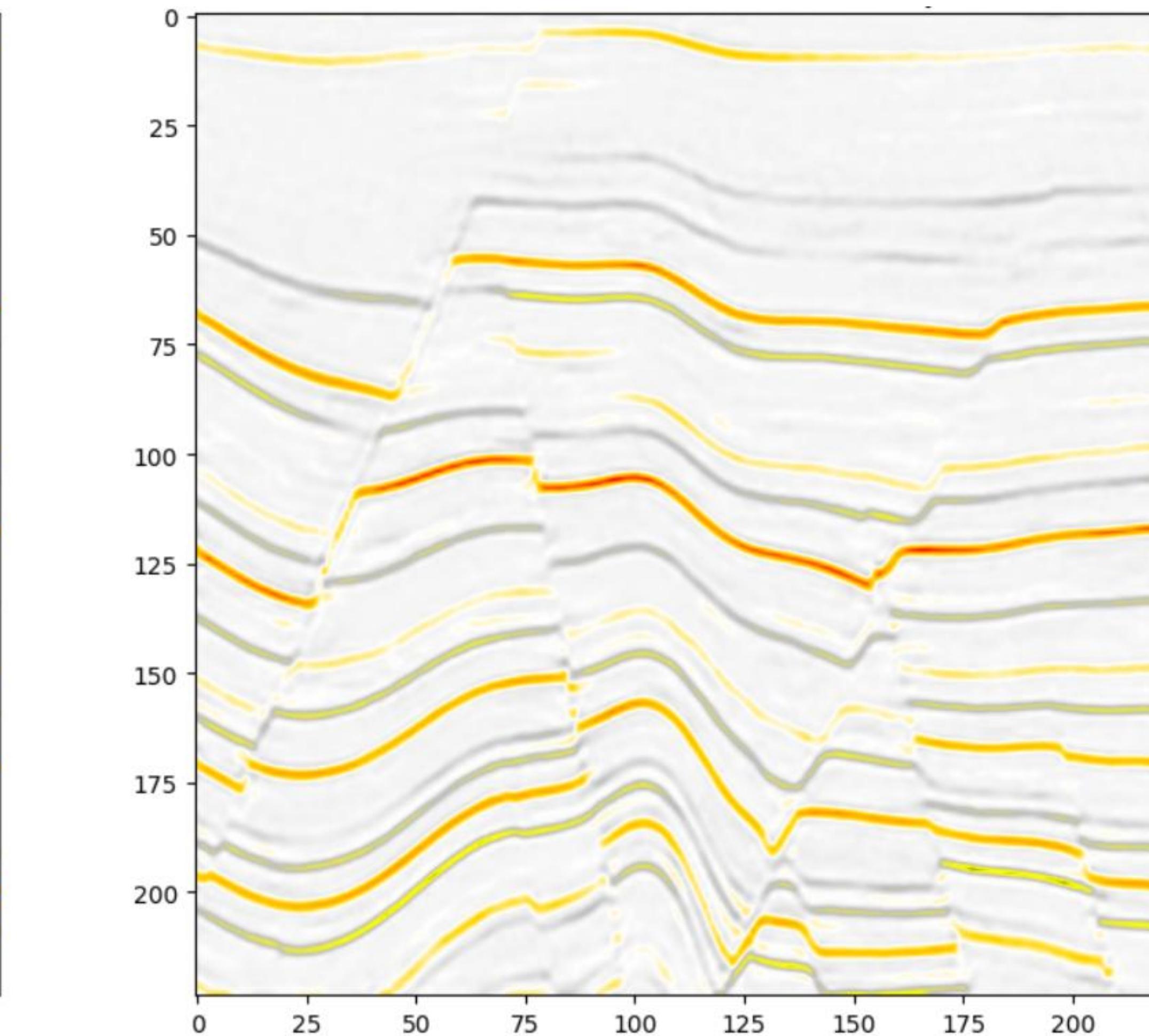
Input Seismic



MAE



Label



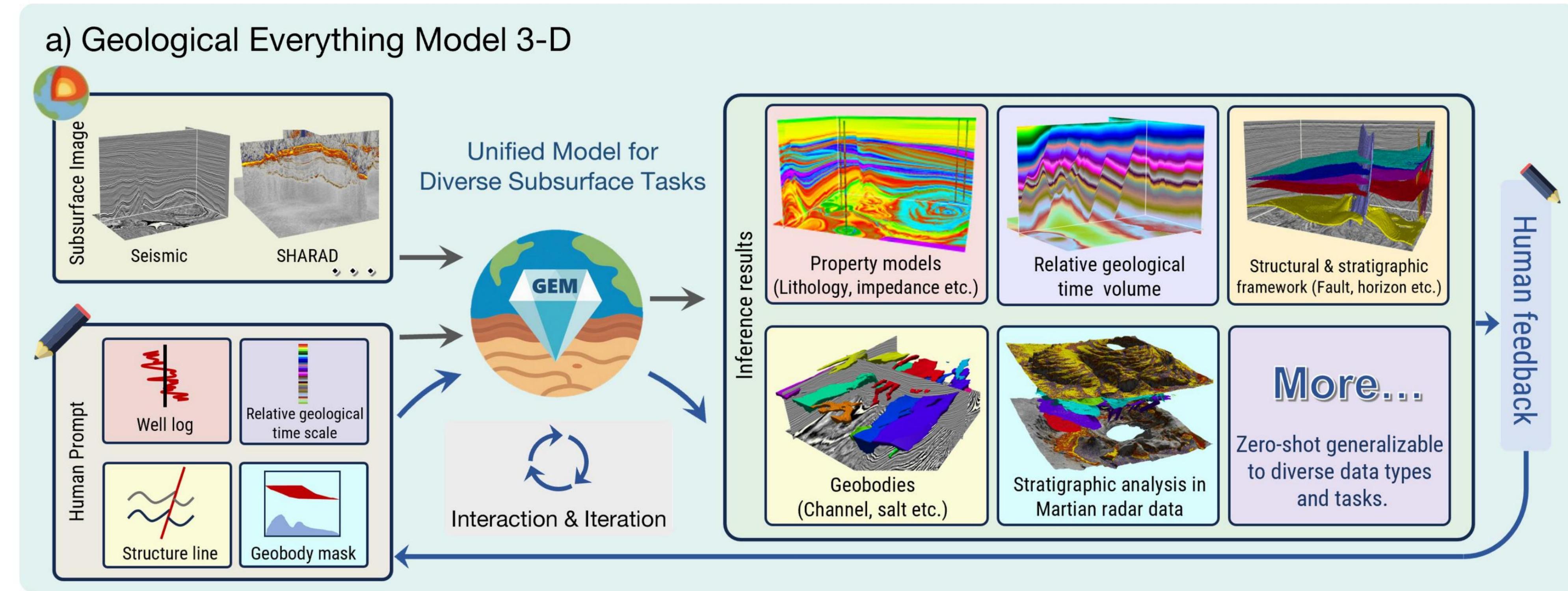
Reflectivity inversion

(Results reproduced from [Sheng et al., 2025](#))

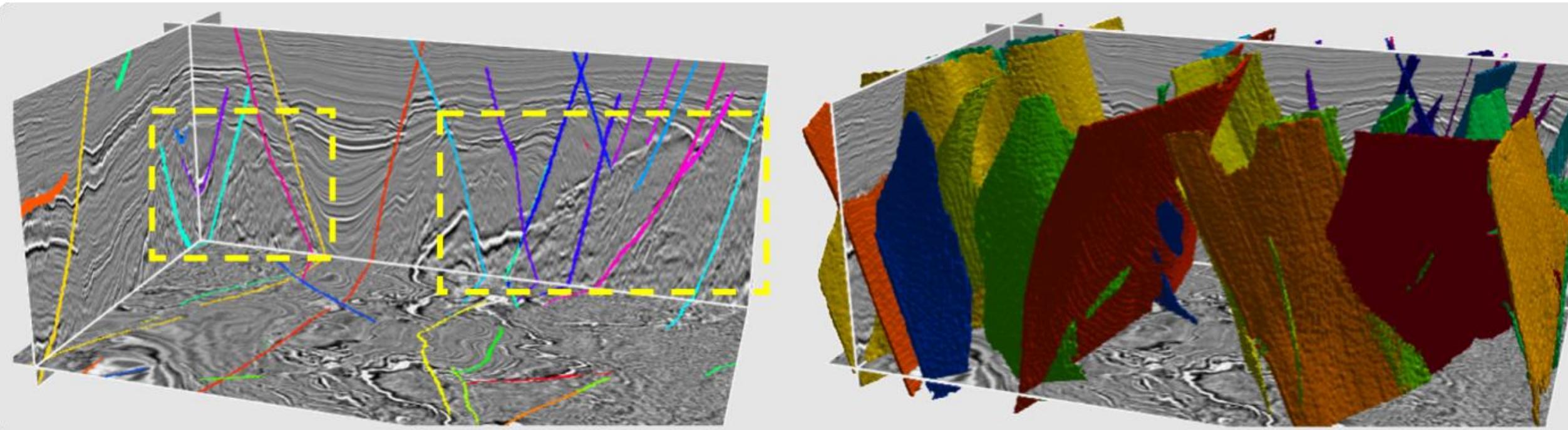
“Seismic Foundation Model (SFM): a new generation deep learning model in geophysics

Geological Everything Model 3D

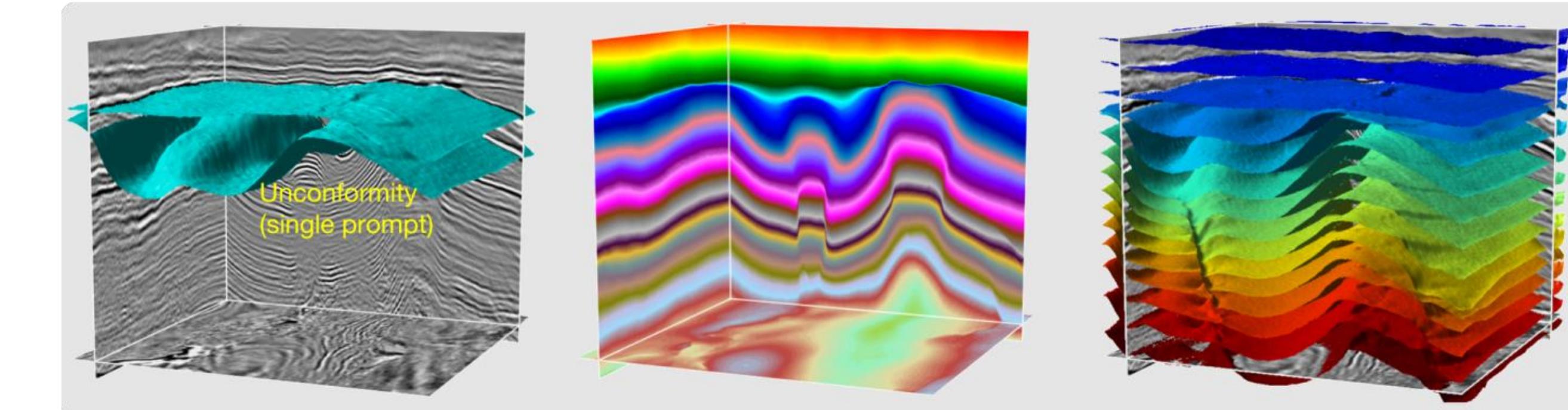
A Promptable Foundation Model for Unified and Zero-Shot Subsurface Understanding



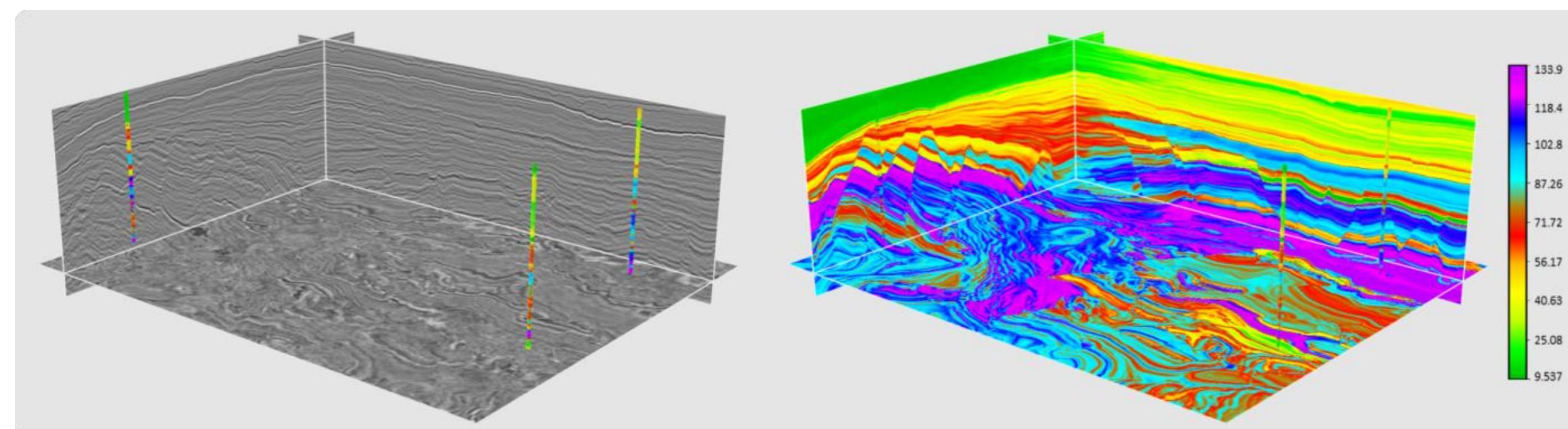
Baiyun Fault



Delft Unconformity and RGT



Poseidon GammaRay



Try Diffusion Model for Full-Waveform Inversion Using NVIDIA PhysicsNeMo

NVIDIA PhysicsNeMo Framework latest ▾

Overview

Getting Started

- System Requirements
- Install Guide

User Guide

- Training Recipe
- Logging and Checkpointing
- Model Architectures
- PhysicsNeMo Distributed
- Physics-guided
- Performance
- Data Curation
- Model Evaluation and Inference
- Symbolic Abstractions

Examples

- PhysicsNeMo Examples Catalog

Library Documentation

PhysicsNeMo

- API Reference

Examples

- Introductory Examples
- CFD Examples
- Weather and Climate

Diffusion Model for Full-Waveform Inversion (FWI)

Problem Overview

In the context of geophysics, Full Waveform Inversion (FWI) is a seismic imaging technique that reconstructs subsurface properties, also called velocity model, by fitting the recorded seismic waveform. It underpins a range of applications, including:

- Hydro-carbon exploration and production, where an accurate velocity model guides drilling decisions.
- CO₂ storage, ensuring the integrity of underground reservoirs used for carbon capture and sequestration.
- Global and regional seismology, helping characterise tectonic processes and earthquakes.
- Analogous elastic/acoustic imaging modalities such as medical ultrasound and non-destructive testing.

The present example is tailored to the elastic wave equation in the context of hydro-carbon exploration, but the same framework can be applied to other wave equations and applications.

The following introduces a few key concepts that are essential to FWI in the context of hydro-carbon exploration:

- Velocity model $\mathbf{x}(r) = [V_p, V_s, \rho]$ – a 3-D image over coordinates $r = (z, x, y)$, where z is the depth, and x and y are the surface coordinates. The P-wave velocity is denoted by V_p , the S-wave velocity by V_s , and ρ is the density. The velocity model spans several kilometres and is discretised at metre-scale resolution.

https://docs.nvidia.com/physicsnemo/latest/physicsnemo/examples/geophysics/diffusion_fwi/README.html

