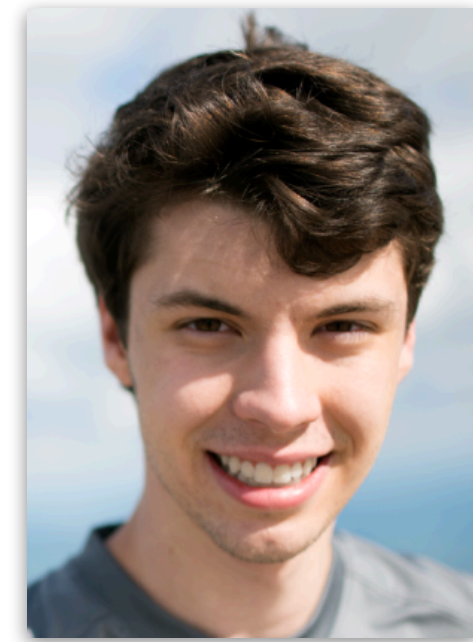
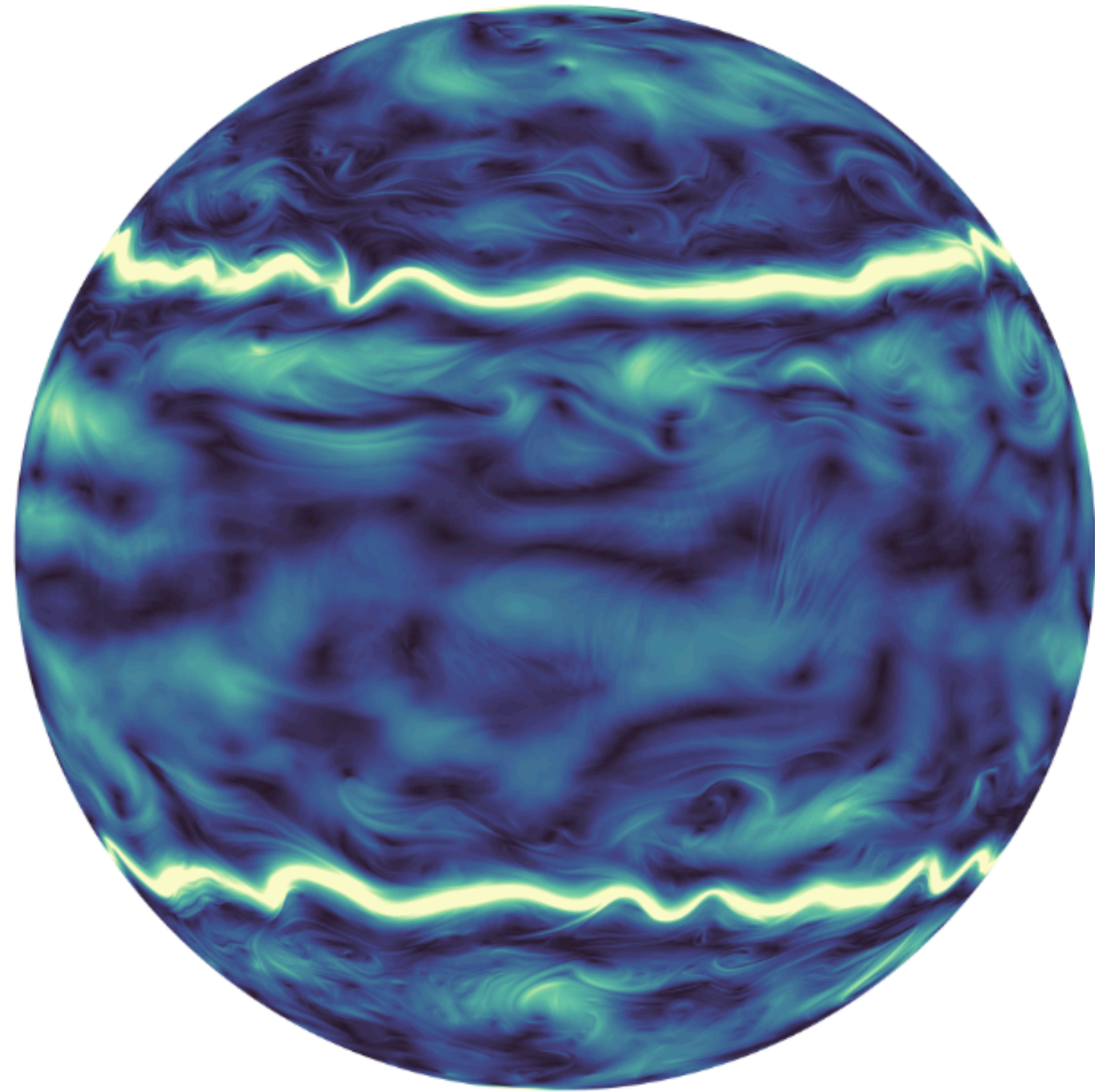
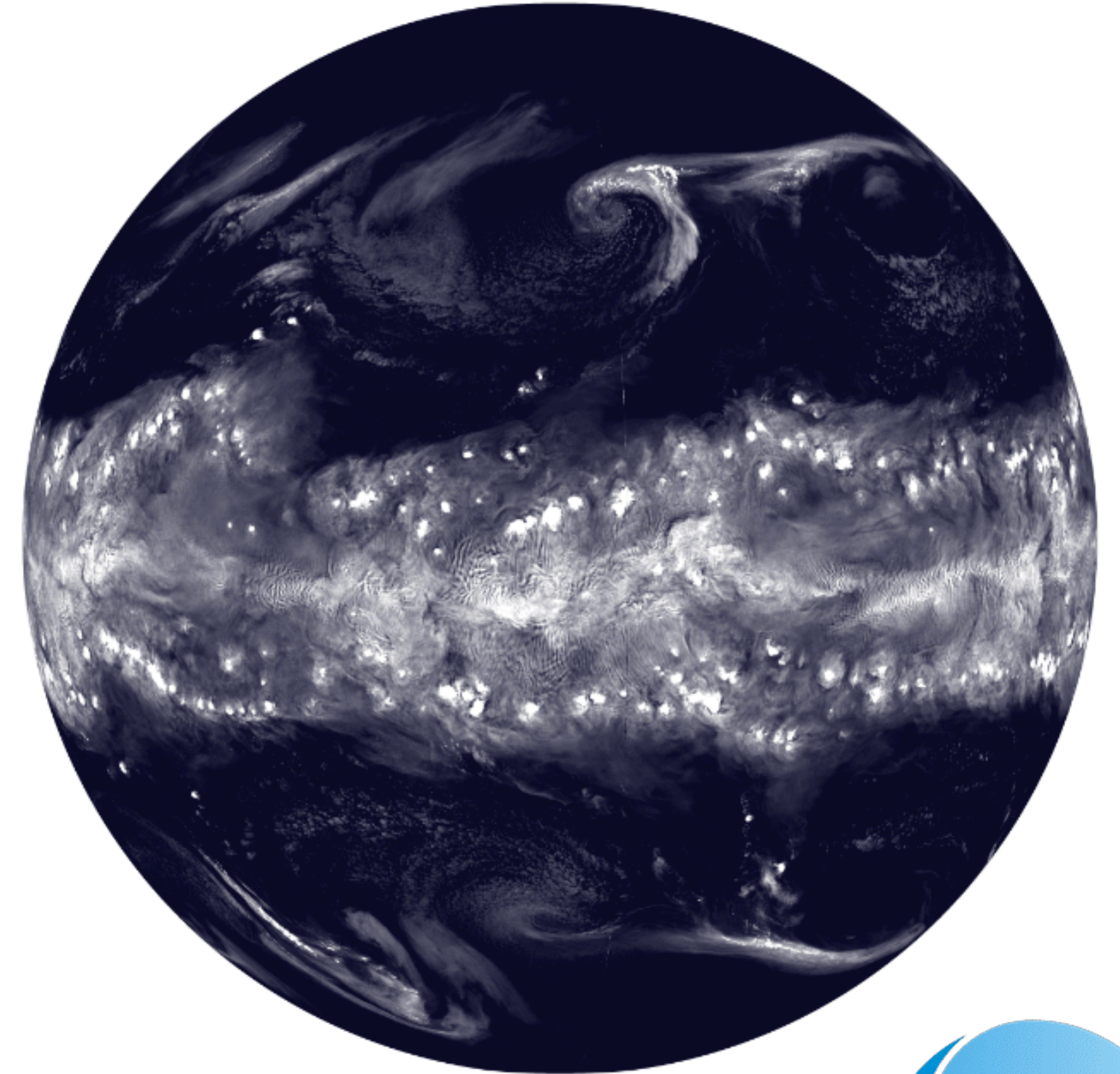


# Making Waves in the Cloud: A Paradigm Shift for Scientific Computing through Compiler Technology



William S. Moses

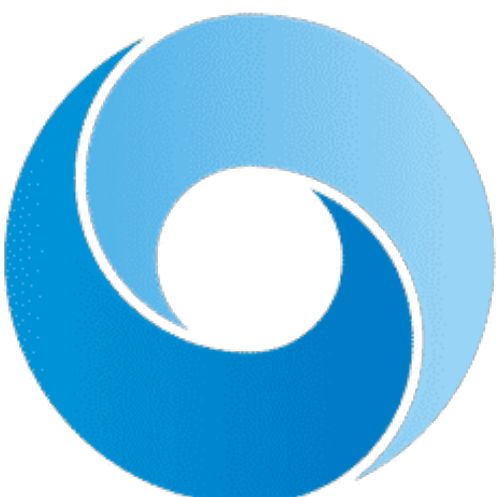


[wsmoses@illinois.edu](mailto:wsmoses@illinois.edu)

University of Cambridge Seminar  
April 21, 2026



UNIVERSITY OF  
**ILLINOIS**  
URBANA-CHAMPAIGN



William S. Moses<sup>†§</sup>, Mosè Giordano<sup>★</sup>, Avik Pal<sup>‡</sup>, Gregory Wagner<sup>‡¶</sup>, Daniel Kyte-Zable<sup>□¶</sup>, Ivan R Ivanov, Paul Berg<sup>∇</sup>, Johannes Blaschke, Roman Lee<sup>♡</sup>, Jules Merckx<sup>△</sup>, Arpit Jaiswal<sup>◆</sup>, Joseph Kump<sup>#</sup>, Patrick Heimbach<sup>#</sup>, Rosie Zou<sup>§</sup>, Tongfei Guo<sup>§</sup>, Son Vu<sup>bb</sup>, Sergio Sanchez-Ramirez<sup>◇</sup>, Simone Silvestri, Nora Loose<sup>♣</sup>, Ivan Ho, Vimarsh Sathia<sup>†</sup>, Jan Hueckelheim<sup>♠</sup>, Johannes De Fine Licht, Kevin Gleason<sup>§</sup>, Ludovic Räss<sup>◇◇</sup>, Gabriel Baraldi, Dhruv Apte<sup>#</sup>, Tobias Bischoff<sup>¶</sup>, Lorenzo Chelini<sup>◆</sup>, Jacques Pienaar<sup>§</sup>, Gaetan Lounes, Navid Constantinou, William R. Magro<sup>§</sup>, Michel Schanen<sup>♠</sup>, Alexis Montoison<sup>♠</sup>, Alan Edelman<sup>‡</sup>, Samarth Narang, Tobias Grosser, Keno Fischer<sup>‡</sup>, Valentin Churavy<sup>#</sup>, Robert Hundt<sup>§</sup>,  
Albert Cohen<sup>§</sup>, Oleksandr Zinenko<sup>bb \*</sup>

UIUC<sup>†</sup>, Google<sup>§</sup>, UCL<sup>★</sup>, MIT<sup>‡</sup>, NVIDIA<sup>◆</sup>, UT Austin<sup>#</sup>, Brium<sup>bb</sup>, [C]Worthy<sup>♣</sup>, BSC<sup>◇</sup>, Argonne National Laboratory<sup>♠</sup>, Lawrence Berkeley National Laboratory<sup>♡</sup>, Cambridge<sup>b</sup>, JuliaHub<sup>‡</sup>, Johannes Gutenberg University, Mainz & University of Augsburg<sup>#</sup>, BFH<sup>∇</sup>, Ghent University<sup>△</sup>, Aeolus Labs<sup>¶</sup>, Brown University<sup>□</sup>, University of Lausanne<sup>◇◇</sup>

# Computing Hardware is No Longer For Everybody

---

# Computing Hardware is No Longer For Everybody

---

## **NVIDIA Puts Grace Blackwell on Every Desk and at Every AI Developer's Fingertips**

NVIDIA Project DIGITS With New GB10 Superchip Debuts as World's Smallest AI Supercomputer Capable of Running 200B-Parameter Models

# Computing Hardware is No Longer For Everybody

---

## Exclusive: Meta begins testing its first in-house AI training chip

By Katie Paul and [Krystal Hu](#)

March 11, 2025 2:37 PM GMT+1 · Updated March 11, 2025



[1/2] The logo of Meta Platforms' business group is seen in Brussels, Belgium December 6, 2022. REUTERS/Yves Herman/File Photo [Purchase Licensing Rights](#)



## NVIDIA Puts Grace Blackwell on Every Desk and at Every AI Developer's Fingertips

NVIDIA Project DIGITS With New GB10 Superchip Debuts as World's Smallest AI Supercomputer Capable of Running 200B-Parameter Models

# Computing Hardware is No Longer For Everybody

## Exclusive: Meta begins testing its first in-house AI training chip

By Katie Paul and [Krystal Hu](#)

March 11, 2025 2:37 PM GMT+1 · Updated March 11, 2025



[1/2] The logo of Meta Platforms' business group is seen in Brussels, Belgium December 6, 2022. REUTERS/Yves Herman/File Photo [Purchase Licensing Rights](#)



ANTHROPIC

Claude ▾ API ▾ Solutions ▾ Research ▾ Commitments ▾ Learn ▾ News [Try Claude](#)

Product

### Claude 3.5 Haiku on AWS Trainium2 and model distillation in Amazon Bedrock

Dec 3, 2024 · 3 min read

## NVIDIA Puts Grace Blackwell on Every Desk and at Every AI Developer's Fingertips

NVIDIA Project DIGITS With New GB10 Superchip Debuts as World's Smallest AI Supercomputer Capable of Running 200B-Parameter Models

# Computing Hardware is No Longer For Everybody

## Exclusive: Meta begins testing its first in-house AI training chip

By Katie Paul and [Krystal Hu](#)

March 11, 2025 2:37 PM GMT+1 · Updated March 11, 2025



[1/2] The logo of Meta Platforms' business group is seen in Brussels, Belgium December 6, 2022. REUTERS/Yves Herman/File Photo [Purchase Licensing Rights](#)



## NVIDIA Puts Grace Blackwell on Every Desk and at Every AI Developer's Fingertips

NVIDIA Project DIGITS With New GB10 Superchip Debuts as World's Smallest AI Supercomputer Capable of Running 200B-Parameter Models

ANTHROPIC

Claude ▾ API ▾ Solutions ▾ Research ▾ Commitments ▾ Learn ▾ News

Try Claude

Product

### Claude 3.5 Haiku on AWS Trainium2 and model distillation in Amazon Bedrock

Dec 3, 2024 · 3 min read

## Elon Musk's xAI is reportedly trying to borrow \$12,000,000,000 for even more Nvidia GPUs, an impulse all PC gamers can truly understand

News By [Andy Edser](#) published 23 July 2025

I've checked down the back of the sofa, and I'm not sure I can cover it.



# Computing Hardware is No Longer For Everybody

## Exclusive: Meta begins testing its first in-house AI training chip

By Katie Paul and [Krystal Hu](#)

March 11, 2025 2:37 PM GMT+1 · Updated March 11, 2025



[1/2] The logo of Meta Platforms' business group is seen in Brussels, Belgium December 6, 2022. REUTERS/Yves Herman/File Photo [Purchase Licensing Rights](#)

## NVIDIA Puts Grace Blackwell on Every Desk and at Every AI Developer's Fingertips

NVIDIA Project DIGITS With New GB10 Superchip Debuts as World's Smallest AI Supercomputer Capable of Running 200B-Parameter Models

ANTHROPIC

Claude ▾ API ▾ Solutions ▾ Research ▾ Commitments ▾ Learn ▾ News [Try Claude](#)

Product

### Claude 3.5 Haiku on AWS Tr

Dec 3, 2024 · 3 min read

## OpenAI's Sam Altman is dreaming of running 100 million GPUs in the future - model distillation in Amazon 100x more than it plans to run by December 2025

News By [Efosa Udimwen](#) published July 26, 2025

OpenAI scale-up will give its investors something to think about



## Elon Musk's xAI is reportedly trying to borrow \$12,000,000,000 for even more Nvidia GPUs, an impulse all PC gamers can truly understand

News By [Andy Edser](#) published 23 July 2025

I've checked down the back of the sofa, and I'm not sure I can cover it.



# Computing Hardware is No Longer For Everybody

## Exclusive: Meta begins testing its first in-house AI training chip

By Katie Paul and Krystal Hu

March 11, 2025 2:37 PM GMT+1 · Updated March 11, 2025



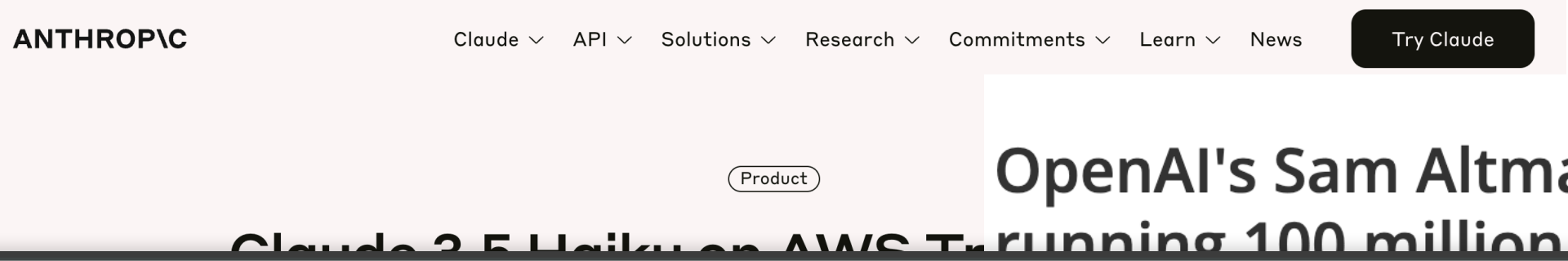
[1/2] The logo of Meta Platforms' business group is seen in Brussels, Belgium. Photo by Herman/Photo: Purchase Licensing Rights

### Ironwood: The first Google TPU for the age of inference

- When scaled to 9,216 chips per pod for a total of 42.5 Exaflops, Ironwood supports more than 24x the compute power of the world's largest supercomputer – El Capitan – which offers just 1.7 Exaflops per pod. Ironwood delivers the massive parallel processing power necessary for the most demanding AI workloads, such as super large size dense LLM or MoE models with thinking capabilities for training and inference. Each individual chip boasts peak compute of 4,614 TFLOPs. This represents a monumental leap in AI capability. Ironwood's memory and network architecture ensures that the right data is always available to support peak performance at this massive scale.

## NVIDIA Puts Grace Blackwell on Every Desk and at Every AI Developer's Fingertips

NVIDIA Project DIGITS With New GB10 Superchip Debuts as World's Smallest AI Supercomputer Capable of Running 200B-Parameter Models



## OpenAI's Sam Altman is dreaming of running 100 million GPUs in the future - by December

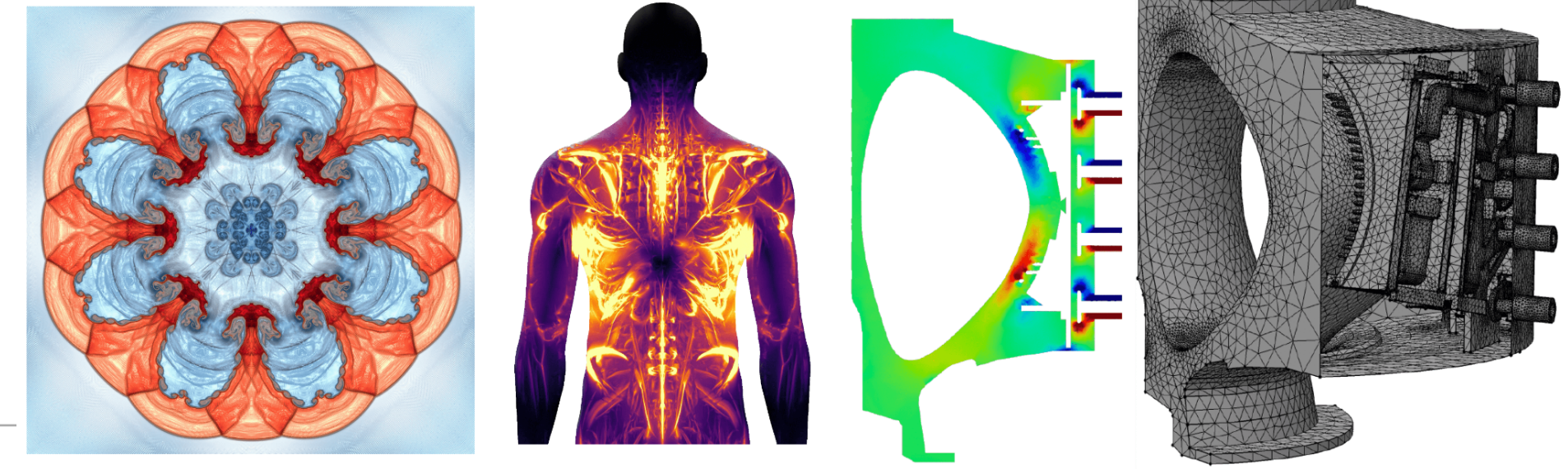
## more NVIDIA GPUs, an impulse all PC gamers can truly understand

News By Andy Edser published 23 July 2025

I've checked down the back of the sofa, and I'm not sure I can cover it.

Comments (2)

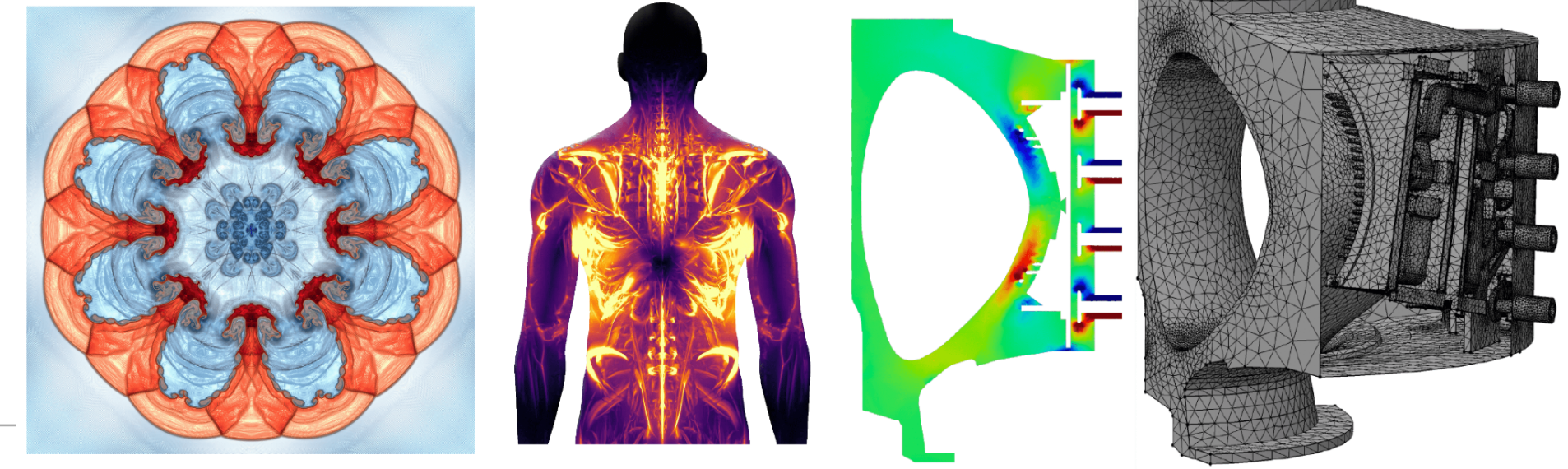
# Lingua Franca of Scientific Computing



- Scientists do not write TPU\* code

```
__global__  
void AddNodeForcesFromElems_kernel( Index_t numNode,  
Index_t padded_numNode,  
const Int_t* nodeElemCount,  
const Int_t* nodeElemStart,  
const Index_t* nodeElemCornerList,  
const Real_t* fx_elem,  
const Real_t* fy_elem,  
const Real_t* fz_elem,  
Real_t* fx_node,  
Real_t* fy_node,  
Real_t* fz_node,  
const Int_t num_threads)  
{  
  int tid=blockDim.x*blockIdx.x+threadIdx.x;  
  if (tid < num_threads)  
  {  
    Index_t g_i = tid;  
    Int_t count=nodeElemCount[g_i];  
    Int_t start=nodeElemStart[g_i];  
    Real_t fx,fy,fz;  
    fx=fy=fz=Real_t(0.0);  
  
    for (int j=0;j<count;j++)  
    {  
      Index_t pos=nodeElemCornerList[start+j]; // Uncoalesced access here  
      fx += fx_elem[pos];  
      fy += fy_elem[pos];  
      fz += fz_elem[pos];  
    }  
  
    fx_node[g_i]=fx;  
    fy_node[g_i]=fy;  
    fz_node[g_i]=fz;  
  }  
}
```

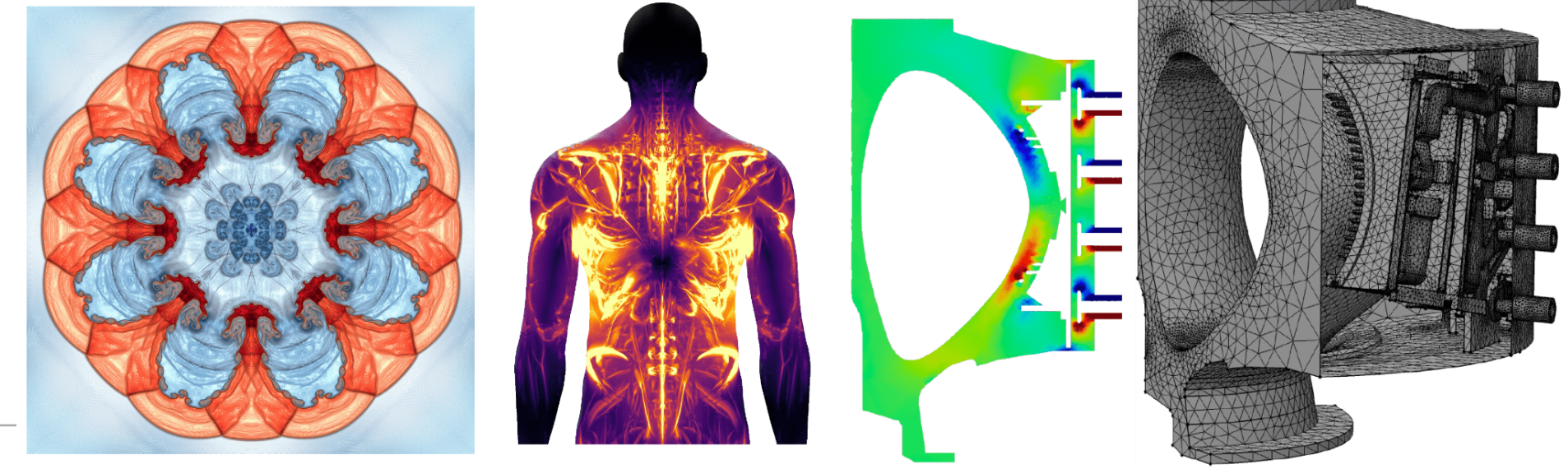
# Lingua Franca of Scientific Computing



- Scientists do not write TPU\* code
  - BIG (MFEM library alone is 737K LOC)

```
__global__  
void AddNodeForcesFromElems_kernel( Index_t numNode,  
Index_t padded_numNode,  
const Int_t* nodeElemCount,  
const Int_t* nodeElemStart,  
const Index_t* nodeElemCornerList,  
const Real_t* fx_elem,  
const Real_t* fy_elem,  
const Real_t* fz_elem,  
Real_t* fx_node,  
Real_t* fy_node,  
Real_t* fz_node,  
const Int_t num_threads)  
{  
    int tid=blockDim.x*blockIdx.x+threadIdx.x;  
    if (tid < num_threads)  
    {  
        Index_t g_i = tid;  
        Int_t count=nodeElemCount[g_i];  
        Int_t start=nodeElemStart[g_i];  
        Real_t fx,fy,fz;  
        fx=fy=fz=Real_t(0.0);  
  
        for (int j=0;j<count;j++)  
        {  
            Index_t pos=nodeElemCornerList[start+j]; // Uncoalesced access here  
            fx += fx_elem[pos];  
            fy += fy_elem[pos];  
            fz += fz_elem[pos];  
        }  
  
        fx_node[g_i]=fx;  
        fy_node[g_i]=fy;  
        fz_node[g_i]=fz;  
    }  
}
```

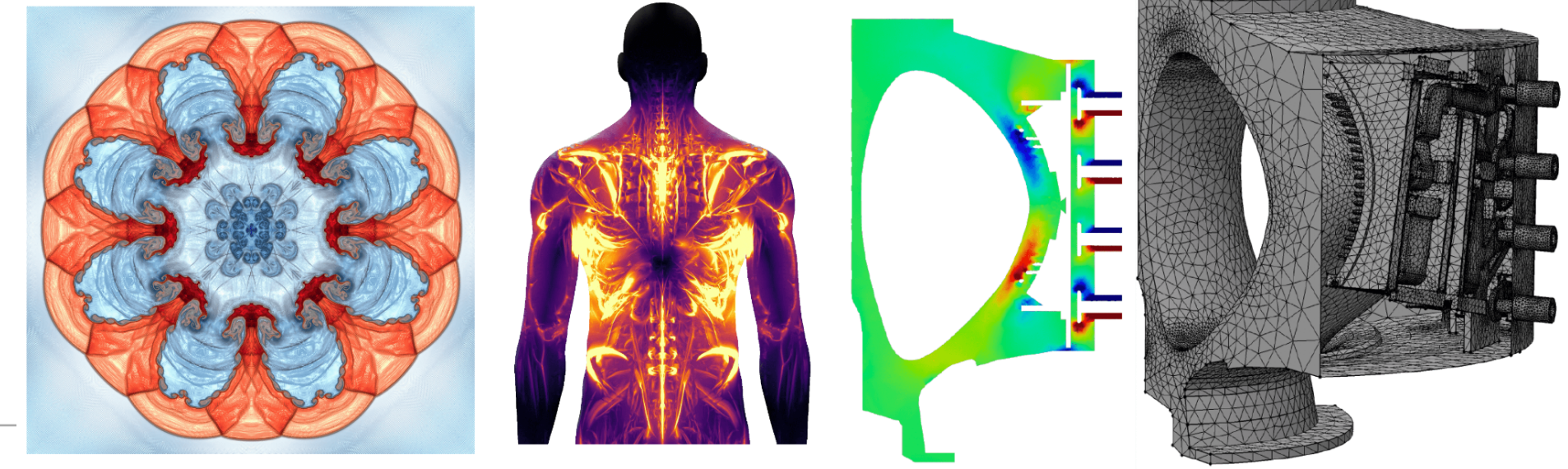
# Lingua Franca of Scientific Computing



- Scientists do not write TPU\* code
  - BIG (MFEM library alone is 737K LOC)
  - Templated

```
__global__  
void AddNodeForcesFromElems_kernel( Index_t numNode,  
                                   Index_t padded_numNode,  
                                   const Int_t* nodeElemCount,  
                                   const Int_t* nodeElemStart,  
                                   const Index_t* nodeElemCornerList,  
                                   const Real_t* fx_elem,  
                                   const Real_t* fy_elem,  
                                   const Real_t* fz_elem,  
                                   Real_t* fx_node,  
                                   Real_t* fy_node,  
                                   Real_t* fz_node,  
                                   const Int_t num_threads)  
{  
    int tid=blockDim.x*blockIdx.x+threadIdx.x;  
    if (tid < num_threads)  
    {  
        Index_t g_i = tid;  
        Int_t count=nodeElemCount[g_i];  
        Int_t start=nodeElemStart[g_i];  
        Real_t fx,fy,fz;  
        fx=fy=fz=Real_t(0.0);  
  
        for (int j=0;j<count;j++)  
        {  
            Index_t pos=nodeElemCornerList[start+j]; // Uncoalesced access here  
            fx += fx_elem[pos];  
            fy += fy_elem[pos];  
            fz += fz_elem[pos];  
        }  
  
        fx_node[g_i]=fx;  
        fy_node[g_i]=fy;  
        fz_node[g_i]=fz;  
    }  
}
```

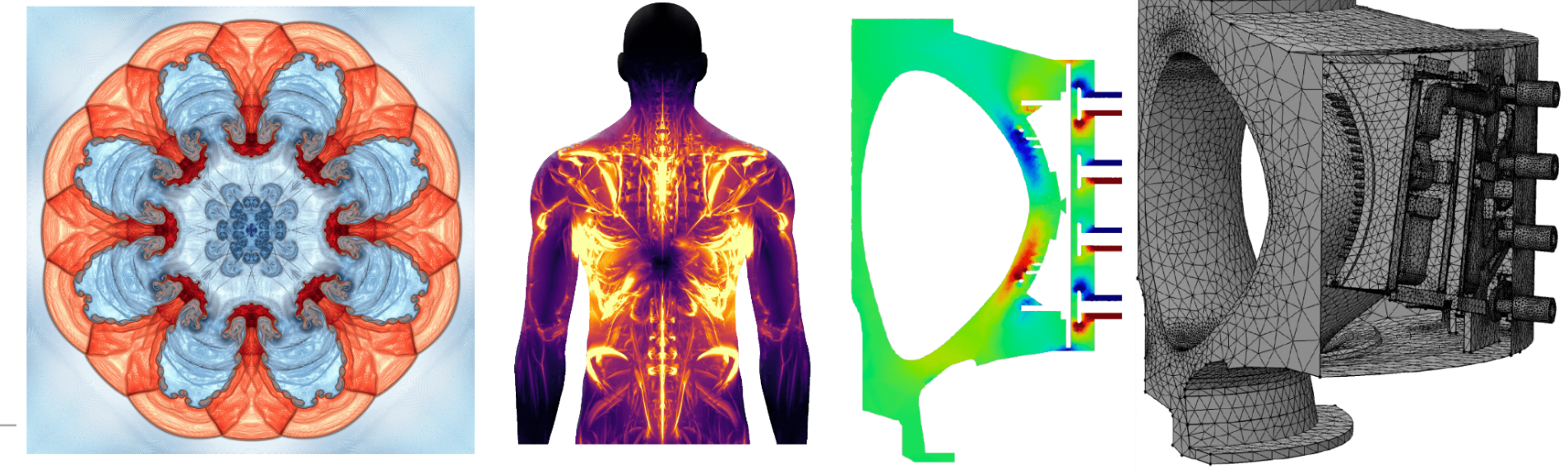
# Lingua Franca of Scientific Computing



- Scientists do not write TPU\* code
  - BIG (MFEM library alone is 737K LOC)
  - Templated
  - Not in Python

```
__global__  
void AddNodeForcesFromElems_kernel( Index_t numNode,  
Index_t padded_numNode,  
const Int_t* nodeElemCount,  
const Int_t* nodeElemStart,  
const Index_t* nodeElemCornerList,  
const Real_t* fx_elem,  
const Real_t* fy_elem,  
const Real_t* fz_elem,  
Real_t* fx_node,  
Real_t* fy_node,  
Real_t* fz_node,  
const Int_t num_threads)  
{  
    int tid=blockDim.x*blockIdx.x+threadIdx.x;  
    if (tid < num_threads)  
    {  
        Index_t g_i = tid;  
        Int_t count=nodeElemCount[g_i];  
        Int_t start=nodeElemStart[g_i];  
        Real_t fx,fy,fz;  
        fx=fy=fz=Real_t(0.0);  
  
        for (int j=0;j<count;j++)  
        {  
            Index_t pos=nodeElemCornerList[start+j]; // Uncoalesced access here  
            fx += fx_elem[pos];  
            fy += fy_elem[pos];  
            fz += fz_elem[pos];  
        }  
  
        fx_node[g_i]=fx;  
        fy_node[g_i]=fy;  
        fz_node[g_i]=fz;  
    }  
}
```

# Lingua Franca of Scientific Computing



- Scientists do not write TPU\* code
  - BIG (MFEM library alone is 737K LOC)
  - Templated
  - Not in Python
  - Sometimes\* in CUDA

```
template <>
struct RajaCuWrap<3>
{
    template <const int BLCK = MFEM_CUDA_BLOCKS, typename DBODY>
    static void run(const int N, DBODY &&d_body,
                   const int X, const int Y, const int Z, const int G)
    {
        RajaCuWrap3D(N, d_body, X, Y, Z, G);
    }
};
```

```
__global__
void AddNodeForcesFromElems_kernel( Index_t numNode,
                                     Index_t padded_numNode,
                                     const Int_t* nodeElemCount,
                                     const Int_t* nodeElemStart,
                                     const Index_t* nodeElemCornerList,
                                     const Real_t* fx_elem,
                                     const Real_t* fy_elem,
                                     const Real_t* fz_elem,
                                     Real_t* fx_node,
                                     Real_t* fy_node,
                                     Real_t* fz_node,
                                     const Int_t num_threads)
{
    int tid=blockDim.x*blockIdx.x+threadIdx.x;
    if (tid < num_threads)
    {
        Index_t g_i = tid;
        Int_t count=nodeElemCount[g_i];
        Int_t start=nodeElemStart[g_i];
        Real_t fx,fy,fz;
        fx=fy=fz=Real_t(0.0);

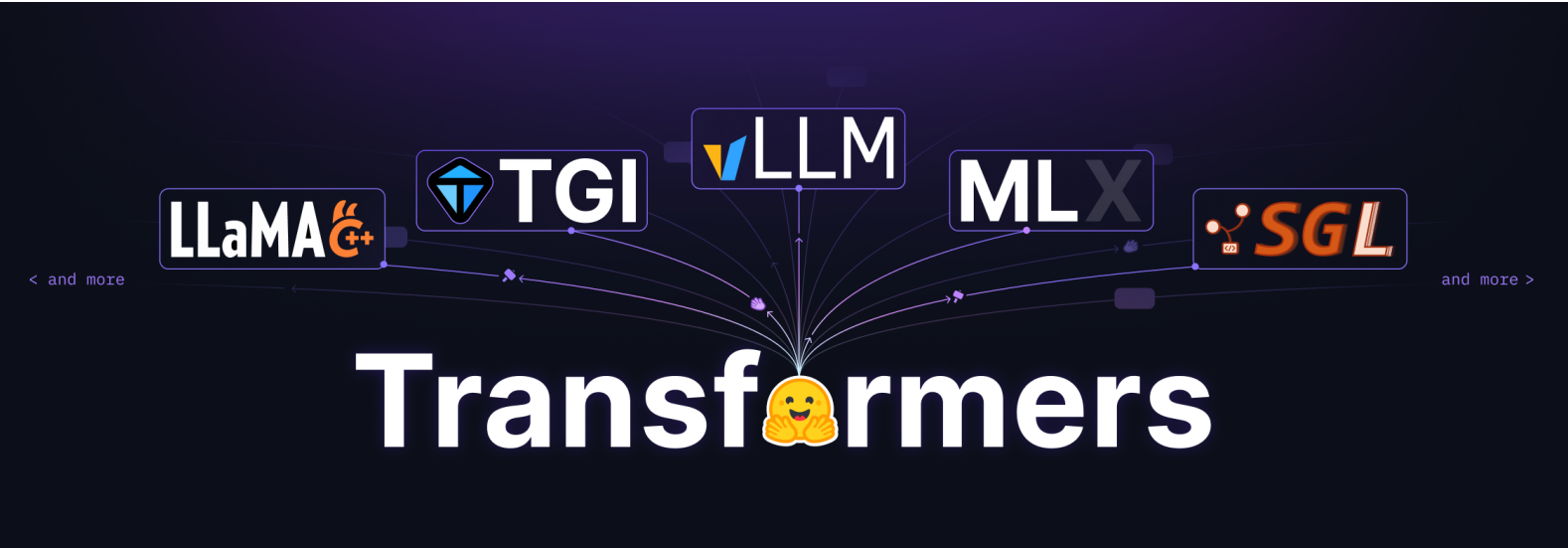
        for (int j=0;j<count;j++)
        {
            Index_t pos=nodeElemCornerList[start+j]; // Uncoalesced access here
            fx += fx_elem[pos];
            fy += fy_elem[pos];
            fz += fz_elem[pos];
        }

        fx_node[g_i]=fx;
        fy_node[g_i]=fy;
        fz_node[g_i]=fz;
    }
}
```

# How do we write ML Accelerator code now?

---

# How do we write ML Accelerator code now?



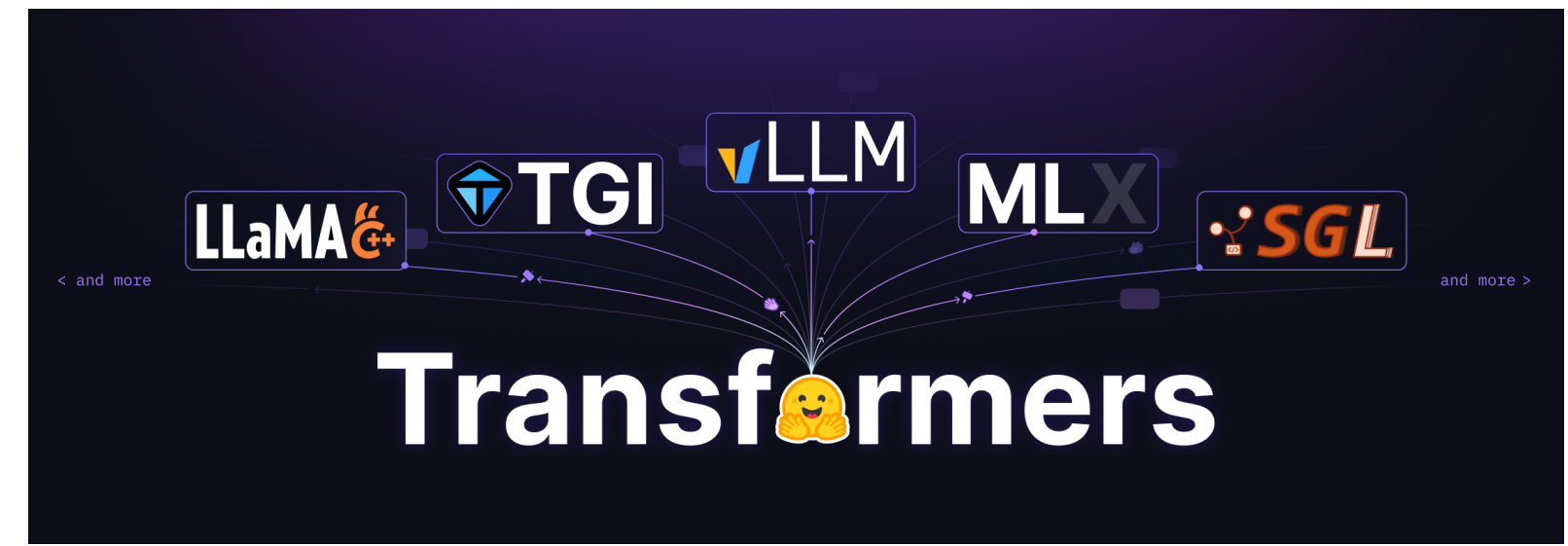
**Stable Diffusion**

Stable Diffusion was made possible thanks to a collaboration with [Stability AI](#) and [Runway](#) and builds upon our previous work:

[High-Resolution Image Synthesis with Latent Diffusion Models](#)  
[Robin Rombach\\*](#), [Andreas Blattmann\\*](#), [Dominik Lorenz](#), [Patrick Esser](#), [Björn Ommer](#)  
[CVPR '22 Oral](#) | [GitHub](#) | [arXiv](#) | [Project page](#)

A horizontal row of five small images generated by Stable Diffusion. From left to right: a person in a dark suit riding a motorcycle through a forest; an astronaut in a white suit playing a piano; a white unicorn in a field; a dog wearing a brown leather jacket and a hat; and a brown bear wearing a space helmet and a life-support pack.


# How do we write ML Accelerator code now?



### Stable Diffusion

Stable Diffusion was made possible thanks to a collaboration with [Stability AI](#) and [Runway](#) and builds upon our previous work:

[High-Resolution Image Synthesis with Latent Diffusion Models](#)  
[Robin Rombach\\*](#), [Andreas Blattmann\\*](#), [Dominik Lorenz](#), [Patrick Esser](#), [Björn Ommer](#)  
[CVPR '22 Oral](#) | [GitHub](#) | [arXiv](#) | [Project page](#)




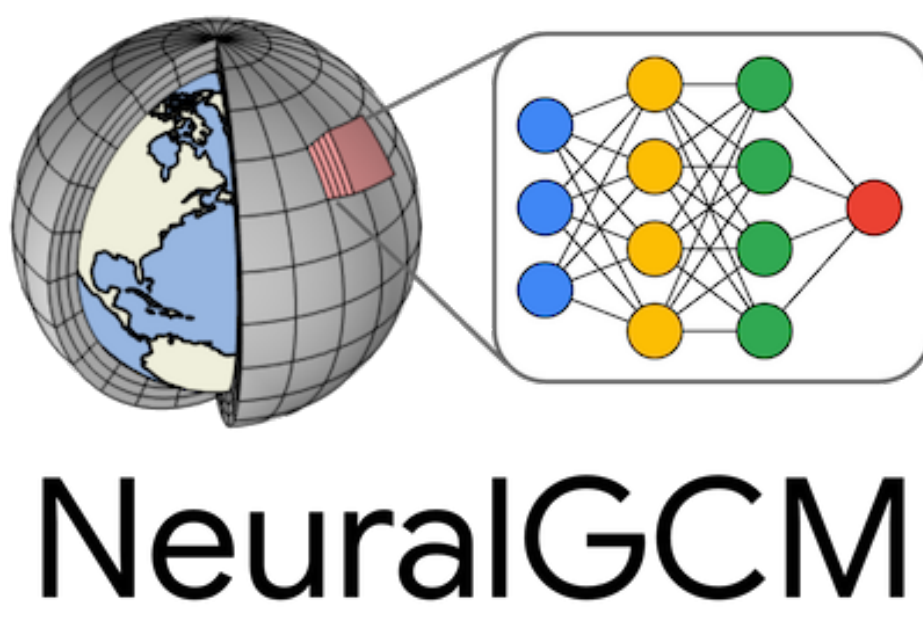
### JAX, M.D.

Accelerated, Differentiable, Molecular Dynamics

[Quickstart](#) | [Reference docs](#) | [Paper](#) | [NeurIPS 2020](#)

Build passing DOI [10.5281/zenodo.14220247](#) pypi [v0.2.8](#) license [Apache 2.0](#)

Molecular dynamics is a workhorse of modern computational condensed matter physics. It is frequently used to simulate materials to observe how small scale interactions can give rise to complex large-scale phenomenology. Most molecular dynamics packages (e.g. HOOMD Blue or LAMMPS) are complicated, specialized pieces of code



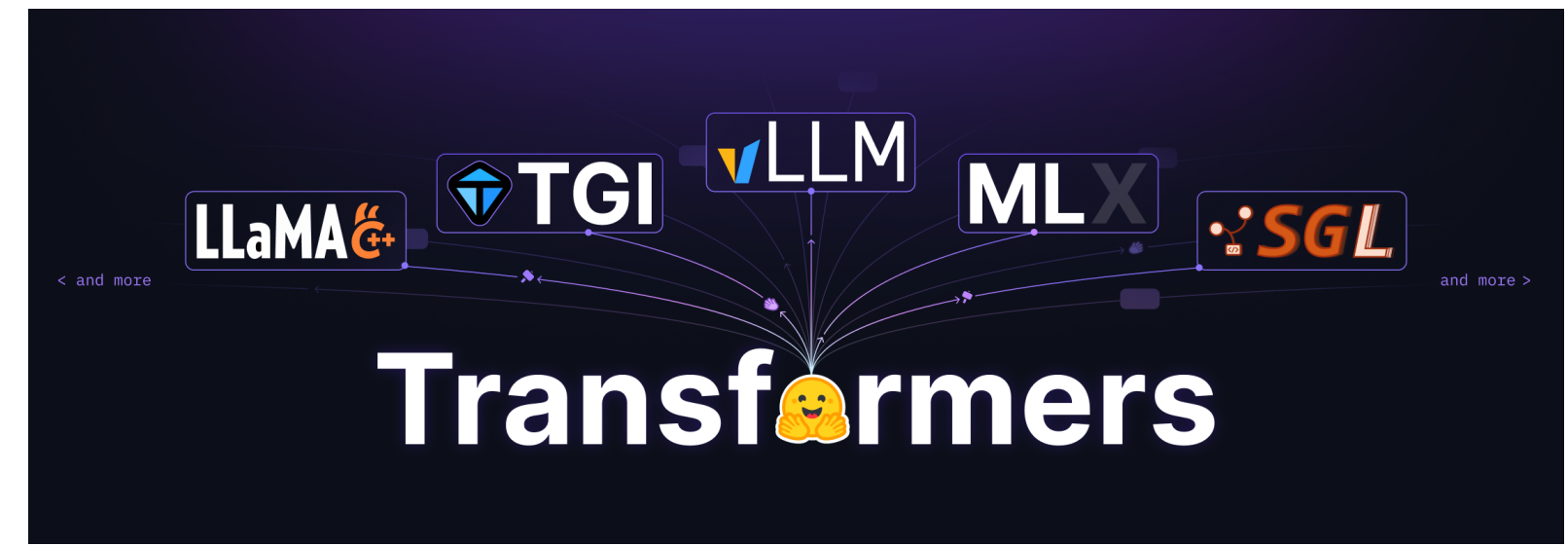
### jaxspec

PYPI [v0.3.0](#) PYTHON [>=3.10,<3.13](#) DOCS passing COVERAGE 94% SLACK

⚠️ jaxspec is still in early release: expect bugs, breaking API changes, undocumented features and lack of functionalities

jaxspec is an X-ray spectral fitting library built in pure Python. It can currently load an X-ray spectrum (in the OGIP standard), define a spectral model from the implemented components, and calculate the best parameters using state-of-the-art Bayesian approaches. It is built on top of JAX to provide just-in-time compilation and automatic differentiation of the spectral models, enabling the use of sampling algorithm such as NUTS.


# How do we write ML Accelerator code now?



### Stable Diffusion

Stable Diffusion was made possible thanks to a collaboration with [Stability AI](#) and [Runway](#) and builds upon our previous work:

[High-Resolution Image Synthesis with Latent Diffusion Models](#)  
[Robin Rombach\\*](#), [Andreas Blattmann\\*](#), [Dominik Lorenz](#), [Patrick Esser](#), [Björn Ommer](#)  
[CVPR '22 Oral](#) | [GitHub](#) | [arXiv](#) | [Project page](#)




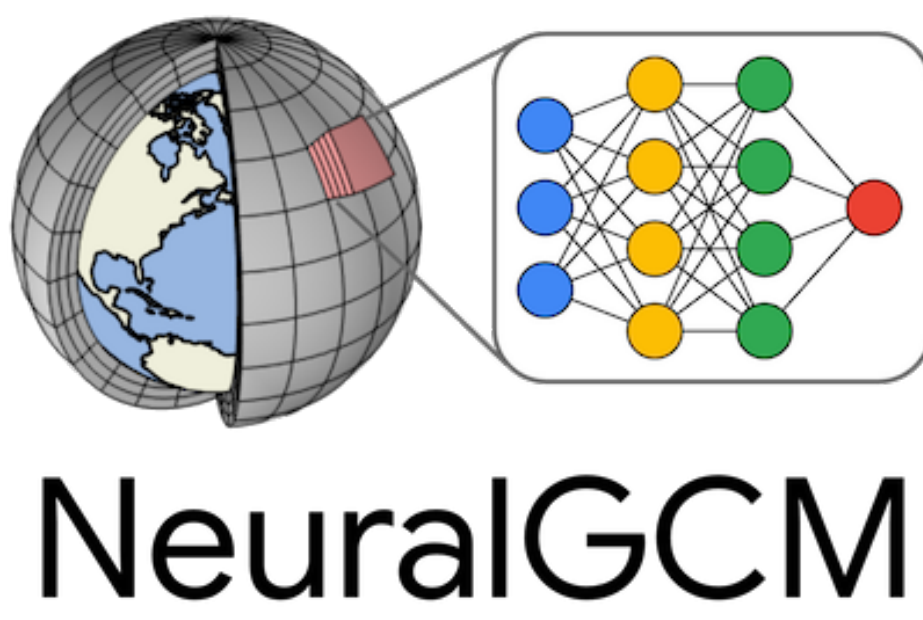
### JAX, M.D.

Accelerated, Differentiable, Molecular Dynamics

[Quickstart](#) | [Reference docs](#) | [Paper](#) | [NeurIPS 2020](#)

Build passing DOI [10.5281/zenodo.14220247](#) pypi [v0.2.8](#) license [Apache 2.0](#)

Molecular dynamics is a workhorse of modern computational condensed matter physics. It is frequently used to simulate materials to observe how small scale interactions can give rise to complex large-scale phenomenology. Most molecular dynamics packages (e.g. HOOMD Blue or LAMMPS) are complicated, specialized pieces of code



### jaxspec

PYPI v0.3.0 PYTHON >=3.10,<3.13 DOCS PASSING COVERAGE 94% SLACK

⚠️ jaxspec is still in early release: expect bugs, breaking API changes, undocumented features and lack of functionalities

jaxspec is an X-ray spectral fitting library built in pure Python. It can currently load an X-ray spectrum (in the OGIP standard), define a spectral model from the implemented components, and calculate the best parameters using state-of-the-art Bayesian approaches. It is built on top of JAX to provide just-in-time compilation and automatic differentiation of the spectral models, enabling the use of sampling algorithm such as NUTS.

## Rewrite it in JAX/PyTorch!

# The Exascale Computing Project (ECP) ECP by the Numbers

The ECP ran from 2016–2024 and was the largest software research, development, and deployment project managed to date by the US Department of Energy (DOE). The \$1.8 billion project was a joint effort by the DOE Office of Science and the National Nuclear Security Administration that funded nearly 2,800 multidisciplinary individuals over the lifetime of the project to uplift the high-performance computing community toward capable exascale platforms, software, and application codes. The outcome was the delivery of an exascale computing ecosystem to provide breakthrough solutions that address future challenges in energy assurance, economic competitiveness, healthcare, and scientific discovery, as well as growing security threats. The ECP exascale ecosystem includes DOE mission-critical application codes, the underlying supporting software technologies, and mechanisms for their deployment and integration.

ECP was a grand convergence of advances in modeling and simulation, software tools and libraries, data analytics, machine learning, and artificial intelligence in support of delivering the world's first capable exascale ecosystem.

The payoff is here: exascale computing is revolutionizing nearly every domain of science.

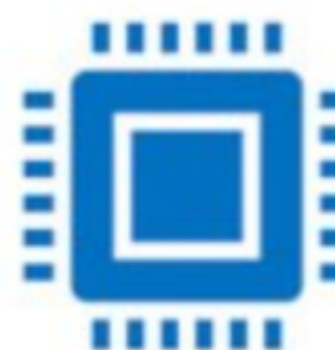
Created to develop the nation's first capable exascale computing ecosystem, this unprecedented DOE research, development, and deployment project has already made a huge impact on computational science:



**2,800 collaborators** funded to develop exascale applications, software, and hardware.



**Game-changing results** in a broad spectrum of science and engineering application areas.



**2 different GPU architectures** now proven to work with exascale environments.



**First and only** open-source scientific software stack developed for scalability and available across all HPC platforms, including cloud computing.

# The Exascale Computing Project (ECP) ECP by the Numbers

The ECP ran from **2016–2024** and was the largest software research, development, and deployment project managed to date by the US Department of Energy (DOE). The **\$1.8 billion** project was a joint effort by the DOE Office of Science and the National Nuclear Security Administration that funded nearly 2,800 multidisciplinary individuals over the lifetime of the project to uplift the high-performance computing community toward capable exascale platforms, software, and application codes. The outcome was the delivery of an exascale computing ecosystem to provide breakthrough solutions that address future challenges in energy assurance, economic competitiveness, healthcare, and scientific discovery, as well as growing security threats. The ECP exascale ecosystem includes DOE mission-critical application codes, the underlying supporting software technologies, and mechanisms for their deployment and integration.

ECP was a grand convergence of advances in modeling and simulation, software tools and libraries, data analytics, machine learning, and artificial intelligence in support of delivering the world's first capable exascale ecosystem.

The payoff is here: exascale computing is revolutionizing nearly every domain of science.

Created to develop the nation's first capable exascale computing ecosystem, this unprecedented DOE research, development, and deployment project has already made a huge impact on computational science:



**2,800 collaborators** funded to develop exascale applications, software, and hardware.



**Game-changing results** in a broad spectrum of science and engineering application areas.



**2 different GPU architectures** now proven to work with exascale environments.



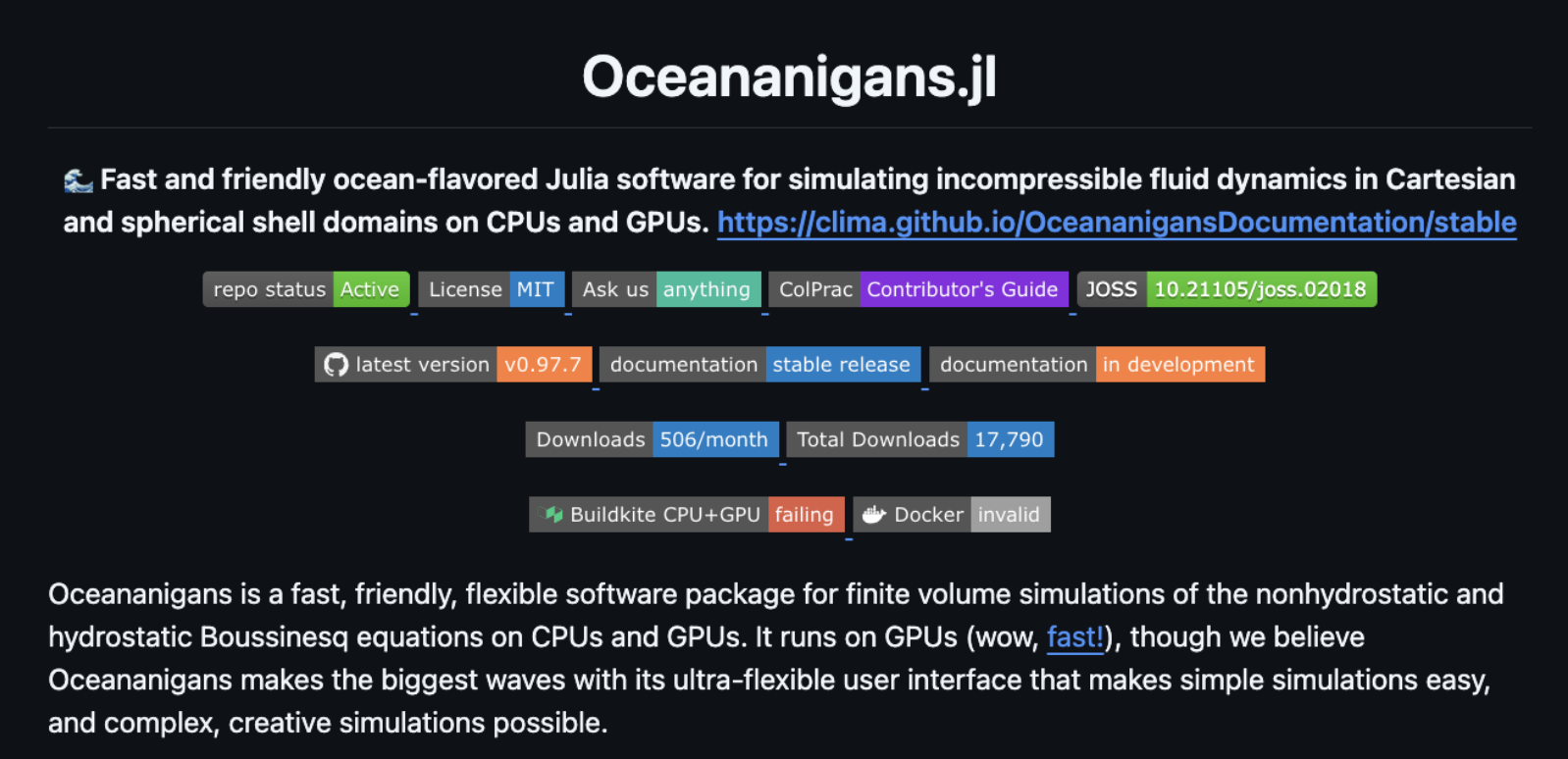
**First and only** open-source scientific software stack developed for scalability and available across all HPC platforms, including cloud computing.

# Looking More Deeply at Scientific Code

```
function stencil_kernel(y, x)
    i = threadIdx().x + (blockIdx().x - 1) * blockDim().x
    if i <= length(x) - 2
        y[i] = x[i] - 2 * x[i + 1] + x[i + 2]
    end
end

function model(...)
    @cuda threads=... blocks=... stencil_kernel(y, x)
    @cuda threads=... blocks=... stencil_kernel(x, y)
end
```

> 277 such kernels



**Oceananigans.jl**

Fast and friendly ocean-flavored Julia software for simulating incompressible fluid dynamics in Cartesian and spherical shell domains on CPUs and GPUs. <https://clima.github.io/OceananigansDocumentation/stable>

repo status **Active** License **MIT** Ask us **anything** ColPrac **Contributor's Guide** JOSS **10.21105/joss.02018**

latest version **v0.97.7** documentation **stable release** documentation **in development**

Downloads **506/month** Total Downloads **17,790**

Buildkite CPU+GPU **failing** Docker **invalid**

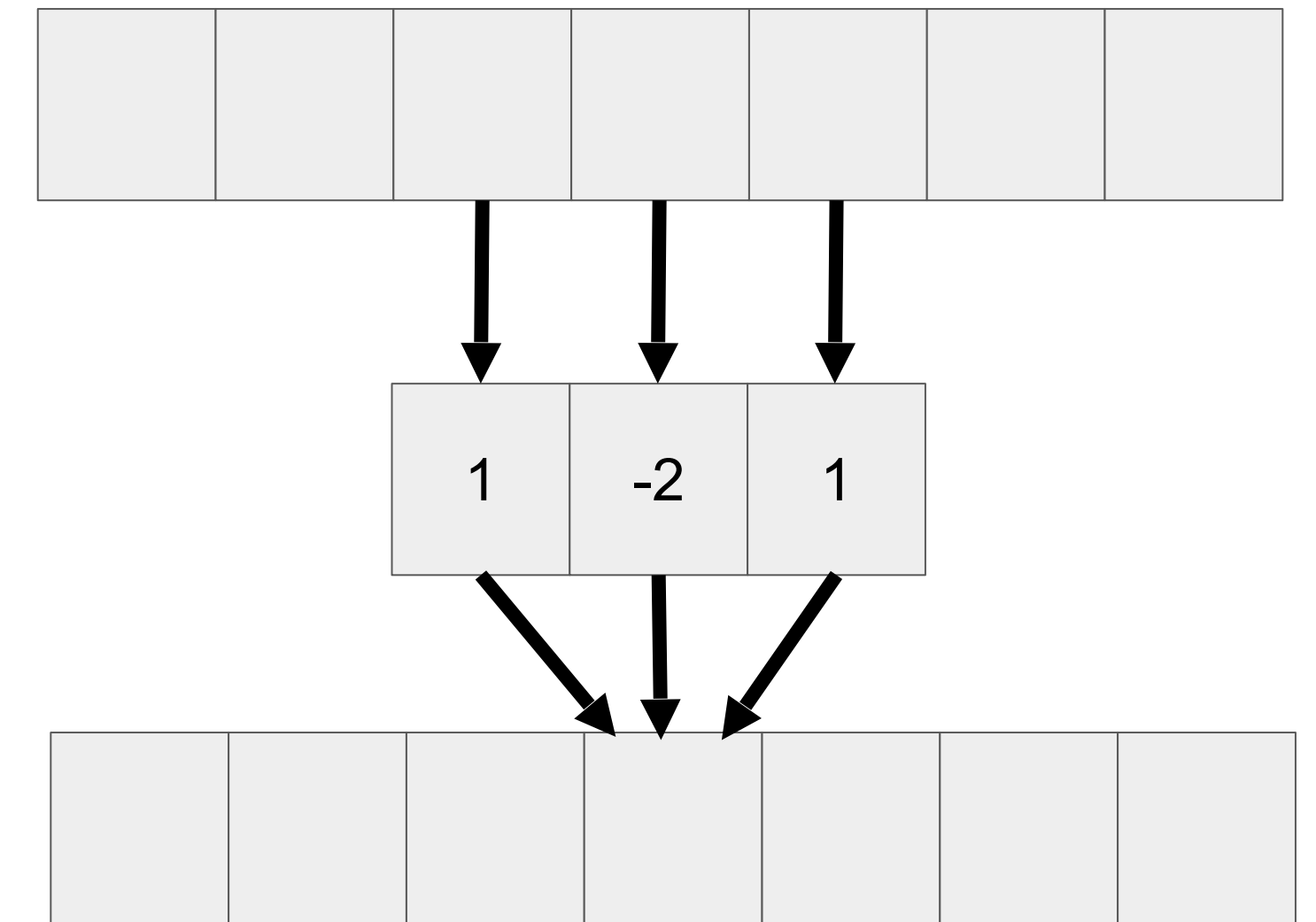
Oceananigans is a fast, friendly, flexible software package for finite volume simulations of the nonhydrostatic and hydrostatic Boussinesq equations on CPUs and GPUs. It runs on GPUs (wow, **fast!**), though we believe Oceananigans makes the biggest waves with its ultra-flexible user interface that makes simple simulations easy, and complex, creative simulations possible.

# Looking More Deeply at Scientific Code

```
function stencil_kernel(y, x)
    i = threadIdx().x + (blockIdx().x - 1) * blockDim().x
    if i <= length(x) - 2
        y[i] = x[i] - 2 * x[i + 1] + x[i + 2]
    end
end

function model(...)
    @cuda threads=... blocks=... stencil_kernel(y, x)
    @cuda threads=... blocks=... stencil_kernel(x, y)
end
```

> 277 such kernels



# CUDA to Accelerator IR (StableHLO)

- New framework for raising and optimizing the structure within existing kernels to stablehlo!
  - 1) Compile Kernels to LLVM
  - 2) Raise the underlying structure in MLIR
  - 3) Multi-dimensionalize it into tensor operators
  - 4) Optimize
- Compiled single-node CUDA version of code to execute on thousands of distributed TPUs and GPUs

```
function stencil_kernel(y, x)
  i = threadIdx().x + (blockIdx().x - 1) * blockDim().x
  if i <= length(x) - 2
    y[i] = x[i] - 2 * x[i+1] + x[i+2]
  end
end

function model(...)
  @cuda threads=... blocks=... stencil_kernel(y, x)
  @cuda threads=... blocks=... stencil_kernel(x, y)
end
```

Compilation

```
define void @julia_difference_kernel_890({}* %y, {*} %x) {
top:
  %3 = call i32 @llvm.nvvm.read.ptx.sreg.tid.x()
  %4 = add nuw nsw i32 %3, 1
  ...
  br i1 %.not, label %common.ret, label %L31
}
```

Raising

```
func.func @kernel(%y : memref<100xf64>, %x : memref<100xf64>) {
  affine.parallel %arg1 = 0 to 100 {
    %x1 = affine.load %x[%arg1]
    %x2 = affine.load %x[%arg1 + 1]
    ...
    affine.store %sum, %y[%arg1]
  }
}
```

Multi-Dimensionalization

```
%x1 = stablehlo.slice %x[1:98]
%x2 = stablehlo.slice %x[2:99]
%mul = stablehlo.multiply %x2, tensor<2.0>
%add = stablehlo.add %x1, %mu
```

Optimization

```
res = stablehlo.convolve %x, tensor<[1.0, -4.0, 6.0, -4.0, 1.0]>
```

# GPU Programming via LLVM

- Mainstream compilers do not have a high-level representation of parallelism, making optimization difficult or impossible
- This is accentuated for GPU programs where the kernel is kept in a separate module & synchronization is a barrier to optimization.

```
__global__ void normalize(int *out, int* in, int n) {  
    int tid = blockIdx.x;  
    if (tid < n)  
        out[tid] = in[tid] / sum(in, n);  
}  
  
void launch(int *out, int* in, int n) {  
    normalize<<<n>>(d_out, d_in, n);  
}
```

Host Code

```
target triple = "x86_64-unknown-linux-gnu"  
  
define void @_Z6launchPiS_i(i32* %out,  
                           i32* %in,  
                           i32 %n) {  
    call i32 @pushCallConfiguration(...)  
    call i32 @cudaLaunch(@_device_stub, ...)  
    ret void  
}
```

Device Code

```
target triple = "nvptx"  
  
define void @_Z9normalize(i32* %out,  
                        i32* %in, i32 %n) {  
    %4 = call i32 @llvm.tid.x()  
    %5 = icmp slt i32 %4, %n  
    br i1 %5, label %6, label %13  
  
6:  
    %8 = getelementptr i32, i32* %in, i32 %4  
    %9 = load i32, i32* %8, align 4  
    %10 = call i32 @_Z3sumPii(i32* %in, i32 %n)  
    %11 = sdiv i32 %9, %10  
    %12 = getelementptr i32, i32* %out, i32 %4  
    store i32 %11, i32* %12, align 4  
    br label %13  
  
13:  
    ret void  
}
```

# GPU Programming via MLIR

- Preserve Host & Device code through frontend  
(Clang Plugin for C++, JIT Package for Julia, etc)
- Enables optimization between caller and kernel
- Enable parallelism-specific optimization

```
__global__ void normalize(int *out, int *in, int n) {  
    int tid = blockIdx.x;  
    if (tid < n)  
        out[tid] = in[tid] / sum(in, n);  
}  
  
void launch(int *out, int* in, int n) {  
    normalize<<<n>>>(d_out, d_in, n);  
}
```

```
func @_Z6launch(%out: memref<?xi32>,  
               %in: memref<?xi32>, %n: i32) {  
    %c1 = constant 1 : index  
    %c0 = constant 0 : index  
  
    parallel (%tid) = (%c0) to (%n) step (%c1) {  
        %2 = load %in[%tid]  
        %sum = call @_Z3sumPii(%in, %n)  
        %4 = divsi %2, %sum : i32  
        store %4, %out[%tid]  
        yield  
    }  
    return  
}
```



# GPU Programming via MLIR

- Preserve Host & Device code through frontend  
(Clang Plugin for C++, JIT Package for Julia, etc)
- Enables optimization between caller and kernel
- Enable parallelism-specific optimization

```
__global__ void normalize(int *out, int *in, int n) {  
  int tid = blockIdx.x;  
  if (tid < n)  
    out[tid] = in[tid] / sum(in, n);  
}  
  
void launch(int *out, int* in, int n) {  
  normalize<<<n>>>(d_out, d_in, n);  
}
```

```
func @_Z6launch(%out: memref<?xi32>,  
               %in: memref<?xi32>, %n: i32) {  
  %c1 = constant 1 : index  
  %c0 = constant 0 : index  
  %sum = call @_Z3sumPii(%in, %n)  
  parallel (%tid) = (%c0) to (%n) step (%c1) {  
    %2 = load %in[%tid]  
  
    %4 = divsi %2, %sum : i32  
    store %4, %out[%tid]  
    yield  
  }  
  return  
}
```

# GPU Programming via MLIR

```
func @launch(%h_out : memref<?xf32>, %h_in : memref<?xf32>, %n : i64) {  
  parallel.for (%gx, %gy, %gz) = (0, 0, 0) to (grid.x, grid.y, grid.z) {  
    %shared_val = memref.alloc : memref<f32>  
    parallel.for (%tx, %ty, %tz) = (0, 0, 0) to (blk.x, blk.y, blk.z) {  
      if %tx == 0 {  
        store ..., %shared_val[] : memref<f32>  
      }  
      polygeist.barrier(%tx, %ty, %tz)  
      ...  
    }  
  }  
}
```

# Synchronization via Memory

- Synchronization (`sync_threads`) ensures all threads within a block finish executing `codeA` before executing `codeB`
- The desired synchronization behavior can be reproduced by defining `sync_threads` to have the union of the memory semantics of the code before and after the sync.
- This prevents code motion of instructions which require the synchronization for correctness, but permits other code motion (e.g. index computation).

```
codeA(fib(idx));  
sync_threads;  
codeB(fib(idx));
```



```
off = fib(idx);  
codeA(off);  
sync_threads;  
codeB(off);
```

# Synchronization via Memory

- High-level synchronization representation enables new optimizations, like sync elimination.
- A synchronize instruction is not needed if the set of read/writes before the sync don't conflict with the read/writes after the sync.

```
__global__ void bpnnp_layerforward(...) {
    __shared__ float node[HEIGHT];
    __shared__ float weights[HEIGHT][WIDTH];

    if ( tx == 0 )
        node[ty] = input[index_in] ;

    // Unnecessary Barrier #1
    // None of the read/writes below the sync
    // (weights, hidden)
    // intersect with the read/writes above the sync
    // (node, input)
    __syncthreads();


    // Unnecessary Store #1
    weights[ty][tx] = hidden[index];

    __syncthreads();

    // Unnecessary Load #1
    weights[ty][tx] = weights[ty][tx] * node[ty];
    ...
}
```

# Synchronization via Memory

- Here
- Another



## High-Performance GPU-to-CPU Transpilation and Optimization via High-Level Parallel Constructs

William S. Moses  
wmoses@mit.edu  
MIT CSAIL  
United States

Ivan R. Ivanov  
ivanov@m.titech.ac.jp  
Tokyo Tech  
Japan

Jens Domke  
jens.domke@riken.jp  
RIKEN  
Japan

Toshio Endo  
endo@is.titech.ac.jp  
Tokyo Tech  
Japan

Johannes Doerfert  
jdoerfert@llnl.gov  
LLNL  
United States

Oleksandr Zinenko  
zinenko@google.com  
Google  
France

**Abstract**  
While parallelism remains the main source of performance, architectural implementations and programming models change with each new hardware generation, often leading to costly application re-engineering. Most tools for performance portability require manual and costly application porting to yet another programming model.

We propose an alternative approach that automatically translates programs written in one programming model (CUDA), into another (CPU threads) based on Polygeist/MLIR. Our approach includes a representation of parallel constructs that allows conventional compiler transformations to apply transparently and without modification and enables parallelism-specific optimizations. We evaluate our framework by transpiling and optimizing the CUDA Rodinia benchmark suite for a multi-core CPU and achieve a 58% geometric speedup over handwritten OpenMP code. Further, we show how CUDA kernels from PyTorch can efficiently run and scale on the CPU-only Supercomputer Fugaku without user intervention. Our PyTorch compatibility layer making use of transpiled CUDA PyTorch kernels outperforms the PyTorch CPU native backend by 2.7x.

**CCS Concepts:** • Software and its engineering → Compilers; • Theory of computation → Parallel computing models.

**Keywords:** Polygeist, MLIR, CUDA, Barrier Synchronization

**ACM Reference Format:**  
William S. Moses, Ivan R. Ivanov, Jens Domke, Toshio Endo, Johannes Doerfert, and Oleksandr Zinenko. 2023. High-Performance GPU-to-CPU Transpilation and Optimization via High-Level Parallel Constructs. In *The 28th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming (PPoPP '23)*, February 25–March 1, 2023, Montreal, QC, Canada. ACM, New York, NY, USA, 16 pages. <https://doi.org/10.1145/3572848.3577475>

**1 Introduction**  
Despite x86 CPUs and NVidia GPUs remaining primary platforms for computation, customized and emerging architectures play an important role in the computing landscape. A custom version of an ARM CPU, A64FX, is even used in one of the top supercomputers Fugaku [49] where its high-bandwidth memory is expected to compete with that of GPUs. However, these architectures are often overlooked by efficiency-oriented frameworks and libraries. For example, PyTorch [44] targeting Intel's oneDNN [28] backend expectedly underperforms on ARM due to architecture differences and even Fujitsu's customized oneDNN [20] does not yield competitive performance on some kernels. Such situations call for performance portability.

Many non-library approaches for performance portability have been proposed and include language extensions (e.g., OpenCL [14], OpenACC [26]), parallel programming frameworks (e.g., Kokkos [3]), domain-specific languages (e.g., SPIRAL [17], Halide [47] or Tensor Comprehensions [64]). All of these approaches still require legacy applications to be ported, and sometimes entirely rewritten, due to differences in the language, or the underlying programming model.

We explore an alternative approach based on a fully automated compiler that takes code in one programming model (CUDA) and produces a binary targeting another one (CPU threads). While GPU-to-CPU translation has been explored in the past [9, 23, 58], it was rarely able to produce efficient code. In fact, optimizations for CPUs and even generic compiler transforms, such as common sub-expression elimination or loop-invariant code motion, are hindered by the lack of analyzable representations of parallel constructs inside the compiler [39]. As representations of parallelism within a mainstream compiler have only recently begun to

- 27% speedup on real code, 2.7x on PyTorch cross compilation!

## Retargeting and Respecializing GPU Workloads for Performance Portability

```
__void bpn_layerforward(...) {  
    __float node[HEIGHT];  
    __float weights[HEIGHT][WIDTH];
```

```
    == 0 )  
    ty] = input[index_in] ;
```

```
    // Necessary Barrier #1  
    // of the read/writes below the sync  
    // (ights, hidden)  
    // intersect with the read/writes above the sync  
    // (de, input)  
    threads();
```

```
    // Necessary Store #1  
    [ty][tx] = hidden[index];
```

```
    __syncthreads();
```

```
    // Unnecessary Load #1
```

```
    weights[ty][tx] = weights[ty][tx] * node[ty];
```

```
    ...  
}
```

# Synchronization via Memory

- A unified representation of parallelism enables programs in one parallel architecture (e.g. CUDA) to be compiled to another (e.g. historically OpenMP, now TPUs)
- Some backends do not have block synchronization
- Lower a top-level synchronization by distributing the parallel for loop around the sync, and interchanging control flow

```
parallel_for %i = 0 to N {  
  codeA(%i);  
  sync_threads;  
  codeB(%i);  
}
```



```
parallel_for %i = 0 to N {  
  codeA(%i);  
}  
parallel_for %i = 0 to N {  
  codeB(%i);  
}
```

# Synchronization via Memory

- A unified representation of parallelism enables programs in one parallel architecture (e.g. CUDA) to be compiled to another (e.g. historically OpenMP, now TPUs)
- Some backends do not have block synchronization
- Lower a top-level synchronization by distributing the parallel for loop around the sync, and interchanging control flow

```
parallel_for %i = 0 to N {  
  for %j = ... {  
    codeB1(%i, %j);  
    sync_threads;  
    codeB2(%i, %j);  
  }  
}
```



```
for %j = ... {  
  parallel_for %i = 0 to N {  
    codeB1(%i, %j);  
    sync_threads;  
    codeB2(%i, %j);  
  }  
}
```

# LLVM to StableHLO

---

LLVM/NVVM Dialect

```
llvm.call @_nv_fabsf(%arg0)
llvm.br
```

Arith + Control Flow

```
%0 = math.abs %arg0
cf.br
```

SCF (While)

```
scf.while %arg = %c0 {
  %arg < %c10
} do {
  ...
}
```

SCF (For)

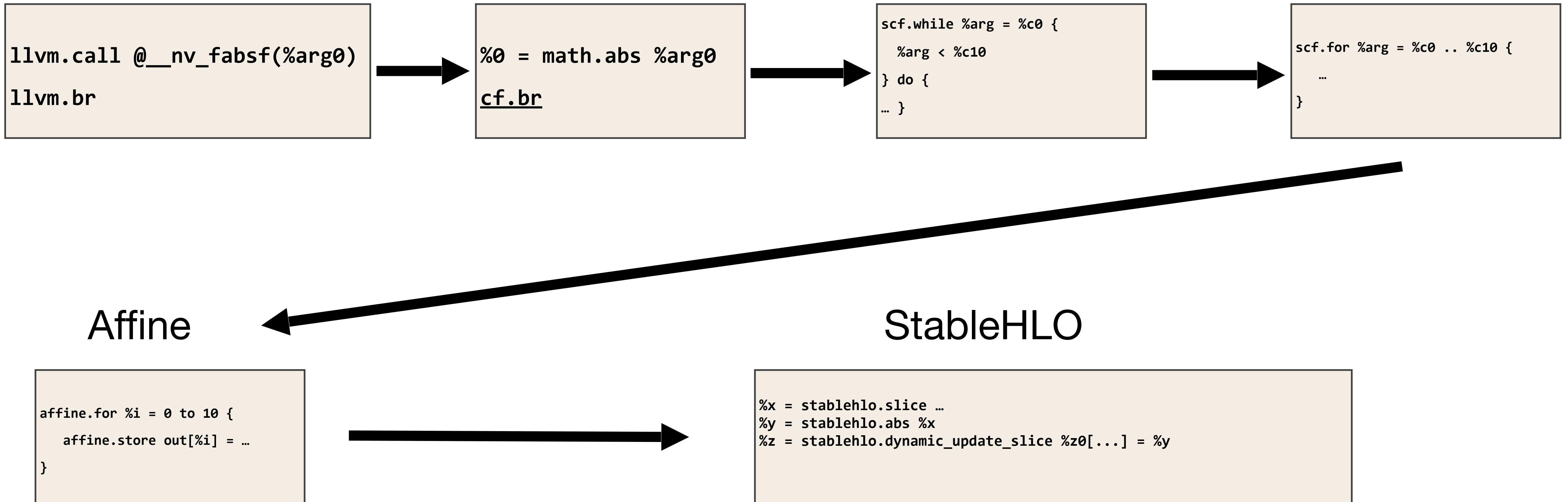
```
scf.for %arg = %c0 .. %c10 {
  ...
}
```

Affine

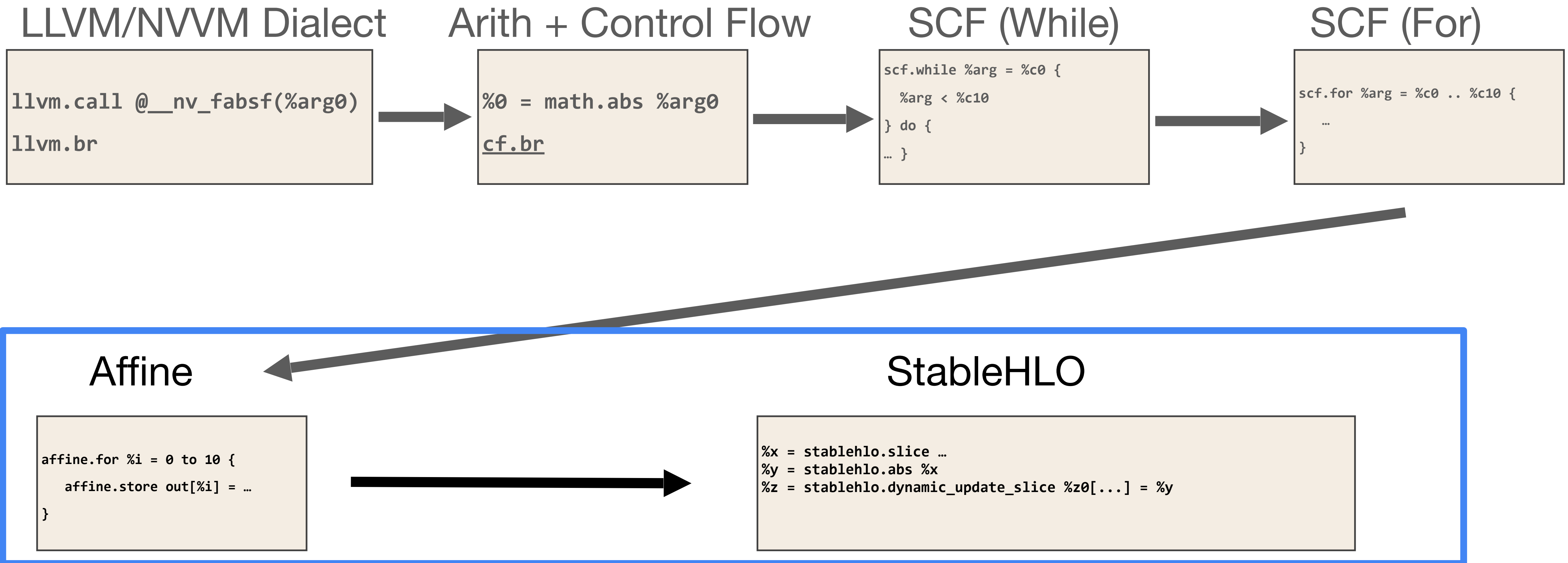
```
affine.for %i = 0 to 10 {
  affine.store out[%i] = ...
}
```

StableHLO

```
%x = stablehlo.slice ...
%y = stablehlo.abs %x
%z = stablehlo.dynamic_update_slice %z0[...] = %y
```



# LLVM to StableHLO



# Affine to StableHLO

---

- Represent *permissive, device-agnostic parallelism*
  - Legal to re-order and interchange instructions
  - One execution (lock-step), runs all of A1, then all of A2, etc
  - Lets us form efficient tensor (stablehlo) versions of kernels

```
parallel.for (%tx, %ty, %tz) = (0,0,0) to (5,7,9){  
    %A1 = load x[%tx, %ty, %tz]  
  
    %A2 = sin(%A1)  
  
    store y[%tx, %ty, %tz] = %A2  
  
    ...  
}
```

# Affine to StableHLO

---

- Represent *permissive, device-agnostic parallelism*
  - Legal to re-order and interchange instructions
  - One execution (lock-step), runs all of A1, then all of A2, etc
  - Lets us form efficient tensor (stablehlo) versions of kernels

```
%A1 = stablehlo.slice %x[0:5, 0:7, 0:9]
parallel.for (%tx, %ty, %tz) = (0,0,0) to (5,7,9){
    %A2 = sin(%A1)
    store y[%tx, %ty, %tz] = %A2
    ...
}
```

# Affine to StableHLO

---

- Represent *permissive, device-agnostic parallelism*
  - Legal to re-order and interchange instructions
  - One execution (lock-step), runs all of A1, then all of A2, etc
  - Lets us form efficient tensor (stablehlo) versions of kernels

```
%A1 = stablehlo.slice %x[0:5, 0:7, 0:9]
%A2 = stablehlo.sine %A1
parallel.for (%tx, %ty, %tz) = (0,0,0) to (5,7,9){
  store y[%tx, %ty, %tz] = %A2
  ...
}
```

# Affine to StableHLO

---

- Represent *permissive, device-agnostic parallelism*
  - Legal to re-order and interchange instructions
  - One execution (lock-step), runs all of A1, then all of A2, etc
  - Lets us form efficient tensor (stablehlo) versions of kernels

```
%A1 = stablehlo.slice %x[0:5, 0:7, 0:9]
%A2 = stablehlo.sine %A1
%Y2 = stablehlo.dynamic_update_slice
      %Y[0:5, 0:7, 0:9], %A2
parallel.for (%tx, %ty, %tz) = (0,0,0) to (5,7,9){
  ...
}
```

# StableHLO ... to better StableHLO

- The direct vectorization of the code works, but may not be efficient.
- We will lost the convolution!
- Perform tensor-level optimizations on stablehlo to recover and optimize higher-level structures

```
%x1 = stablehlo.slice %x[1:98]
%x2 = stablehlo.slice %x[2:99]
%mul = stablehlo.multiply %x2, tensor<2.0>
%add = stablehlo.add %x1, %mu
...
```

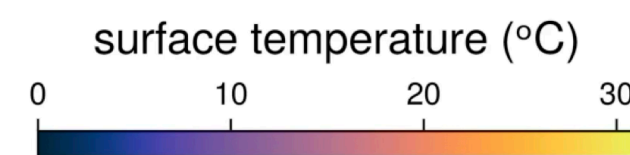
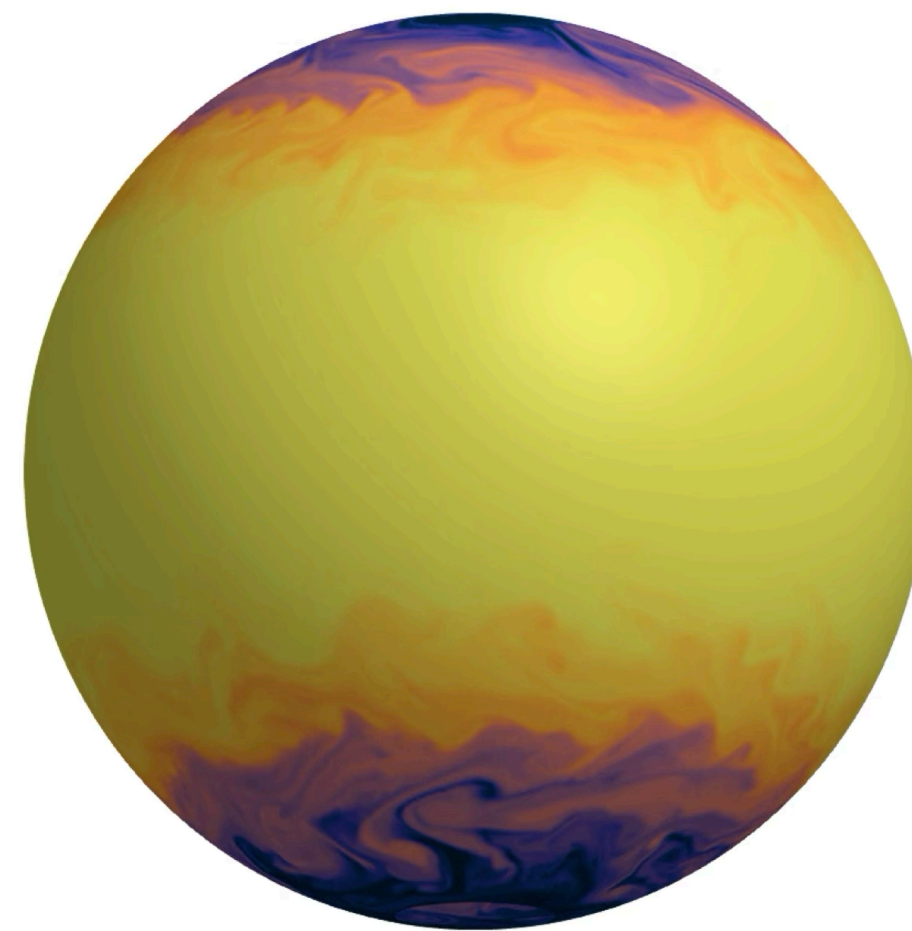
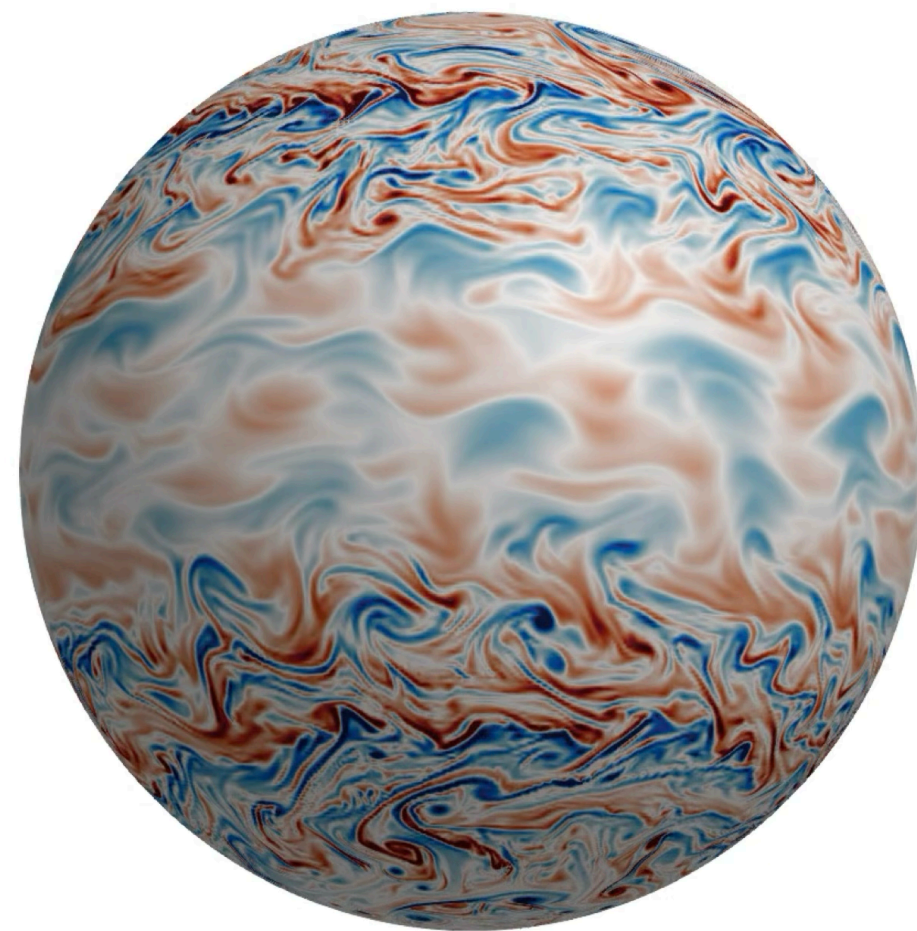
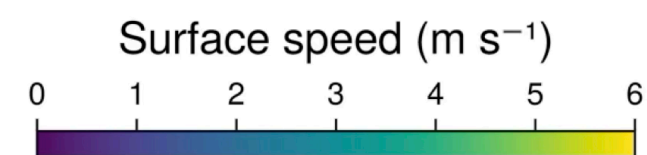
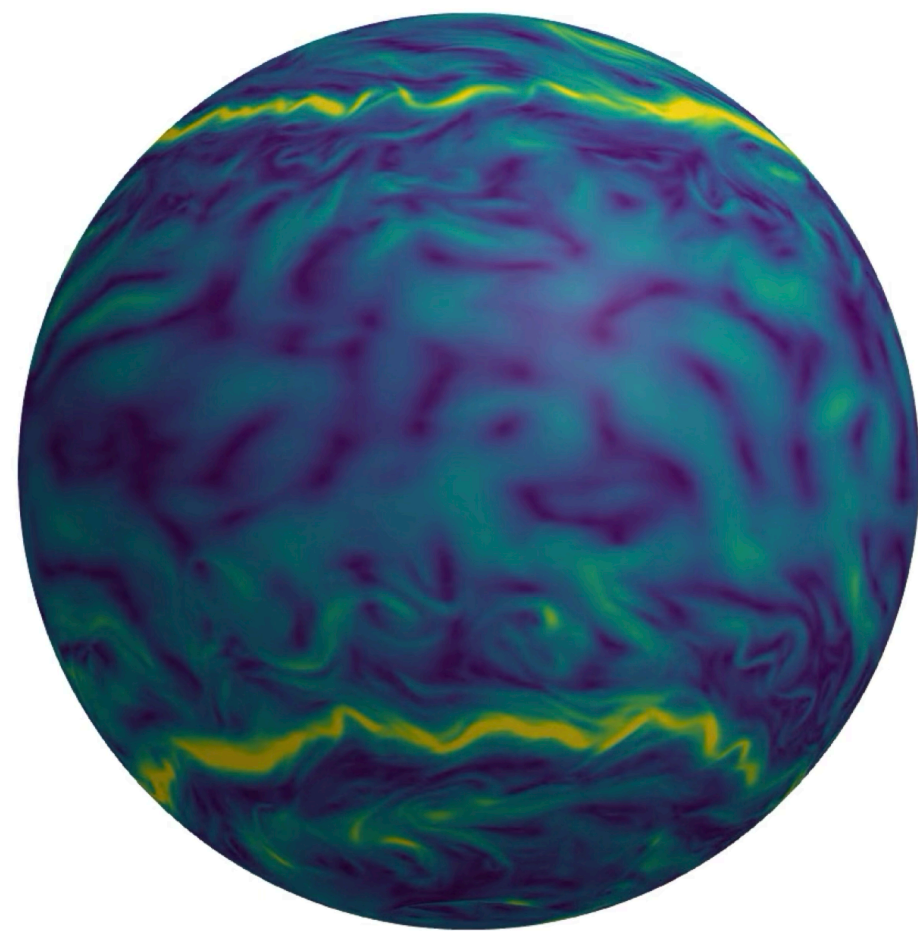
```
%y = stablehlo.convolve %x, tensor<[1.0, -2.0, 1.0]>
%z = stablehlo.convolve %y, tensor<[1.0, -2.0, 1.0]>
```

```
%z = stablehlo.convolve %x, tensor<[1.0, -4.0, 6.0, -4.0, 1.0]>
```



# CUDA to Accelerator IR (StableHLO)

165.0 days



```
function stencil_kernel(y, x)
  i = threadIdx().x + (blockIdx().x - 1) * blockDim().x
  if i <= length(x) - 2
    y[i] = x[i] - 2 * x[i+1] + x[i+2]
  end
end

function model(...)
  @cuda threads=... blocks=... stencil_kernel(y, x)
  @cuda threads=... blocks=... stencil_kernel(x, y)
end
```

Compilation

```
define void @julia_difference_kernel_890({}* %y, {}* %x) {
top:
  %3 = call i32 @llvm.nvvm.read.ptx.sreg.tid.x()
  %4 = add nsw i32 %3, 1
  ...
  br i1 %.not, label %common.ret, label %L31
}
```

Raising

```
func.func @kernel(%y : memref<100xf64>, %x : memref<100xf64>) {
  affine.parallel %arg1 = 0 to 100 {
    %x1 = affine.load %x[%arg1]
    %x2 = affine.load %x[%arg1 + 1]
    ...
    affine.store %sum, %y[%arg1]
  }
}
```

Multi-Dimensionalization

```
%x1 = stablehlo.slice %x[1:98]
%x2 = stablehlo.slice %x[2:99]
%mul = stablehlo.multiply %x2, tensor<2.0>
%add = stablehlo.add %x1, %mul
```

Optimization

```
res = stablehlo.convolve %x, tensor<[1.0, -4.0, 6.0, -4.0, 1.0]>
```

# Single Node Code => Multi Node Code

---

- Each tensor and operation was automatically sharded across the cluster
- Creates suboptimal shadings for intermediate values

```
%left = stablehlo.slice %x[1:2] : tensor<100x100> -> tensor<1x100>  
%x1 = stablehlo.dynamic_update_slice %x[0:1] = %left : tensor<100x100>
```

- Creates redundant communications

```
%left1 = enzymexla.rotate_left %x, 1  
%left2 = enzymexla.rotate_left %x, 2  
%left3 = enzymexla.rotate_left %x, 3  
%result = %left1 + %left2 + %left3
```

- Built novel distributed operations and optimizations

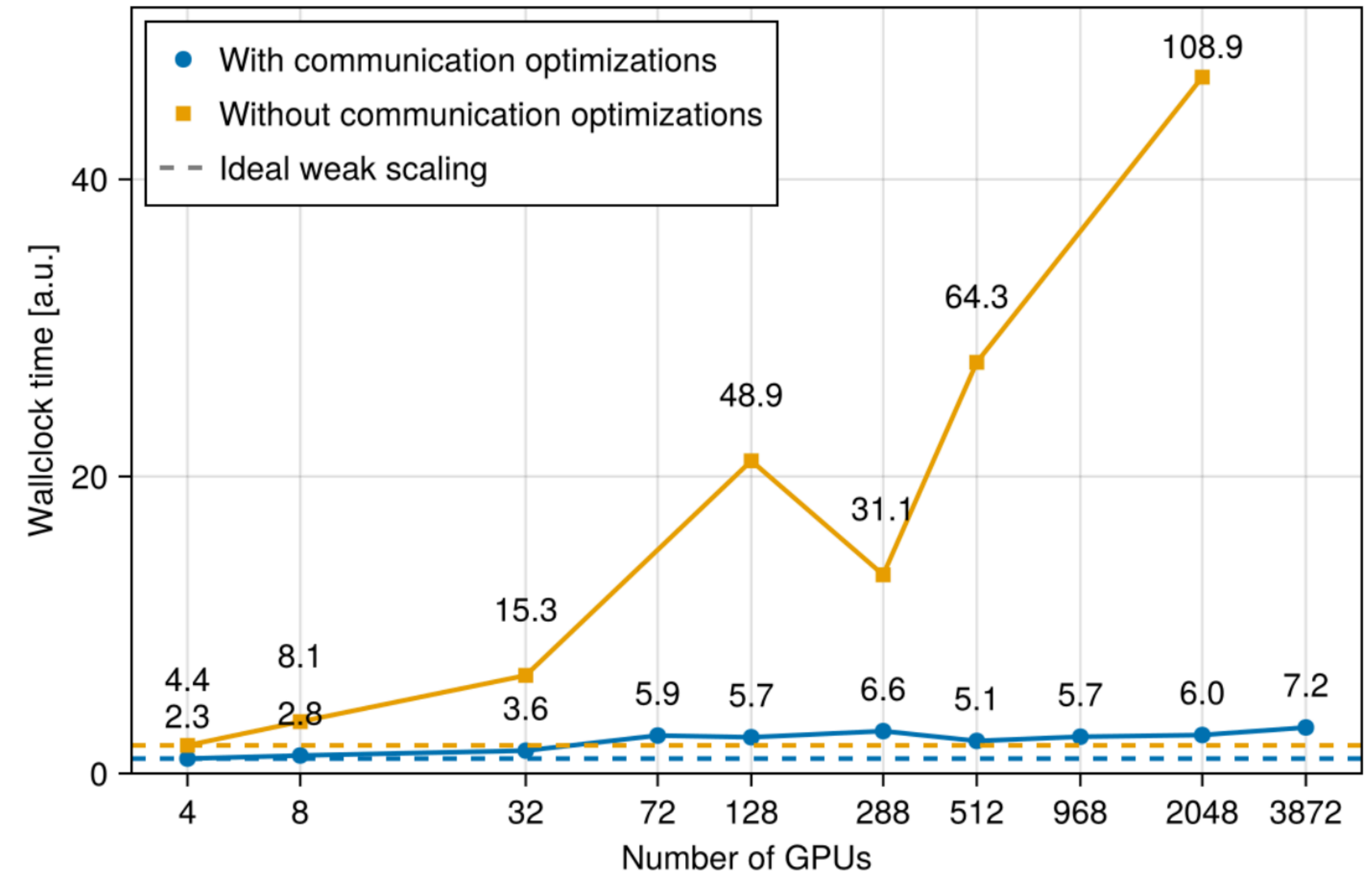
# Distributed Optimizations

- Remove redundant communications (akin to automatic discovery of manual halo exchange)

```
# Perform all rotates in a single communication [extending the halo]
%left1, %left2, %left3 = enzymexla.multi_rotate_left %x, 3

%result = %left1 + %left2 + %left3
```

- Fuse resharding operations into consumers to avoid unnecessary sends
- Follow sub tensor provenance & create novel communication partial-CSE to avoid (partially) redundant sends



# Accelerator-Specific Optimizations

---

- ML accelerators contain the majority of their flops in specialized matrix cores. Detection of convolutions are useful, but stencils contain more general patterns than a pure stencil
  - $2a - b \Rightarrow [2, -1]^T * [a, b]$
  - Generalization of Im2Col
- Accelerators often don't contain native f64 (or even f32) operations.
  - Automatic rewriting of higher level floats in terms of multiple limbs of smaller floats to preserve accuracy while maintaining performance.
  - Even with the additional work, converting f32->2xbf16 was 10x faster than f32 due to bf16-specific hardware

# Performance Results

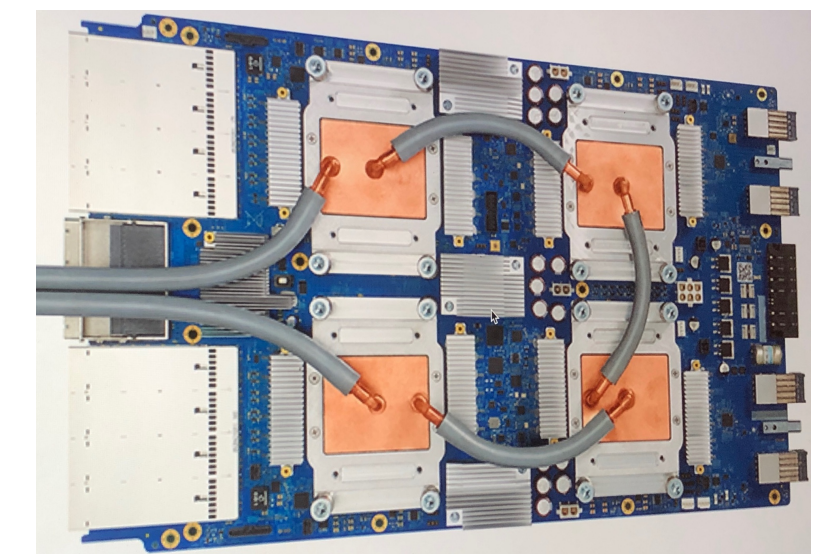
- Successfully ran single-node Oceanangians.jl on thousands of distributed accelerators
  - Alps (2688 nodes x 4 NVIDIA GH200 GPUs)
  - Perlmutter (1536 nodes x 4 NVIDIA A100 GPUs)
  - 8,192 Google TPUs v7 (4614 F8 TFLOPS each)
- Good Single-Node Perf (CPU)
  - Vanilla Model: 272.0seconds
  - Tensor Optimis: 11.5seconds



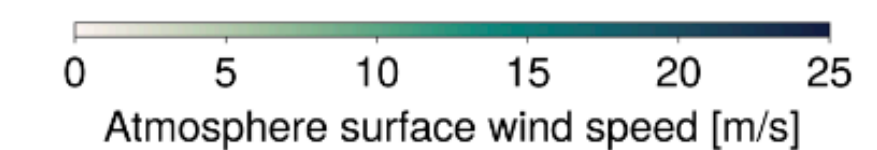
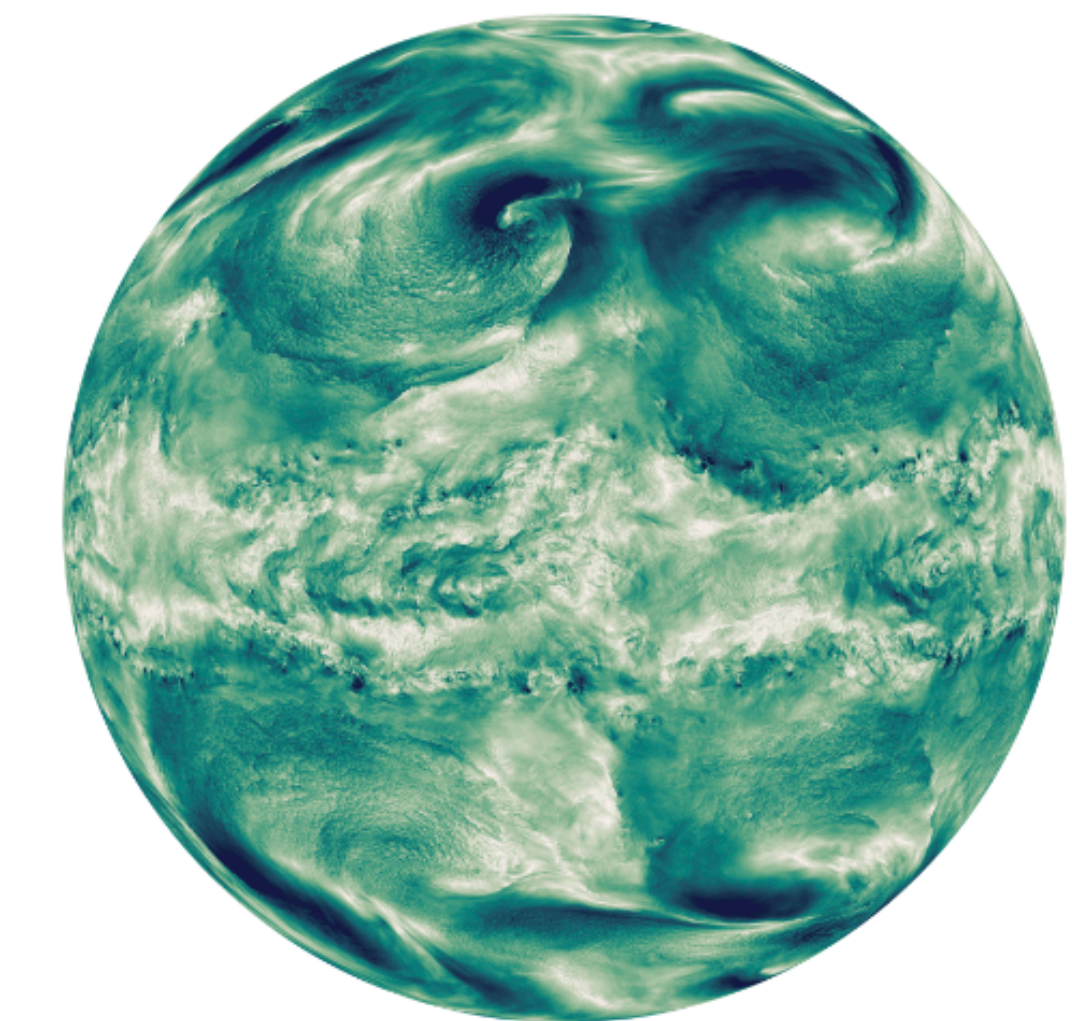
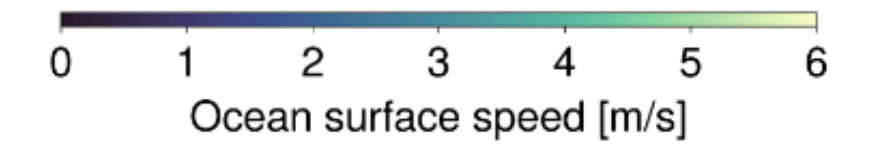
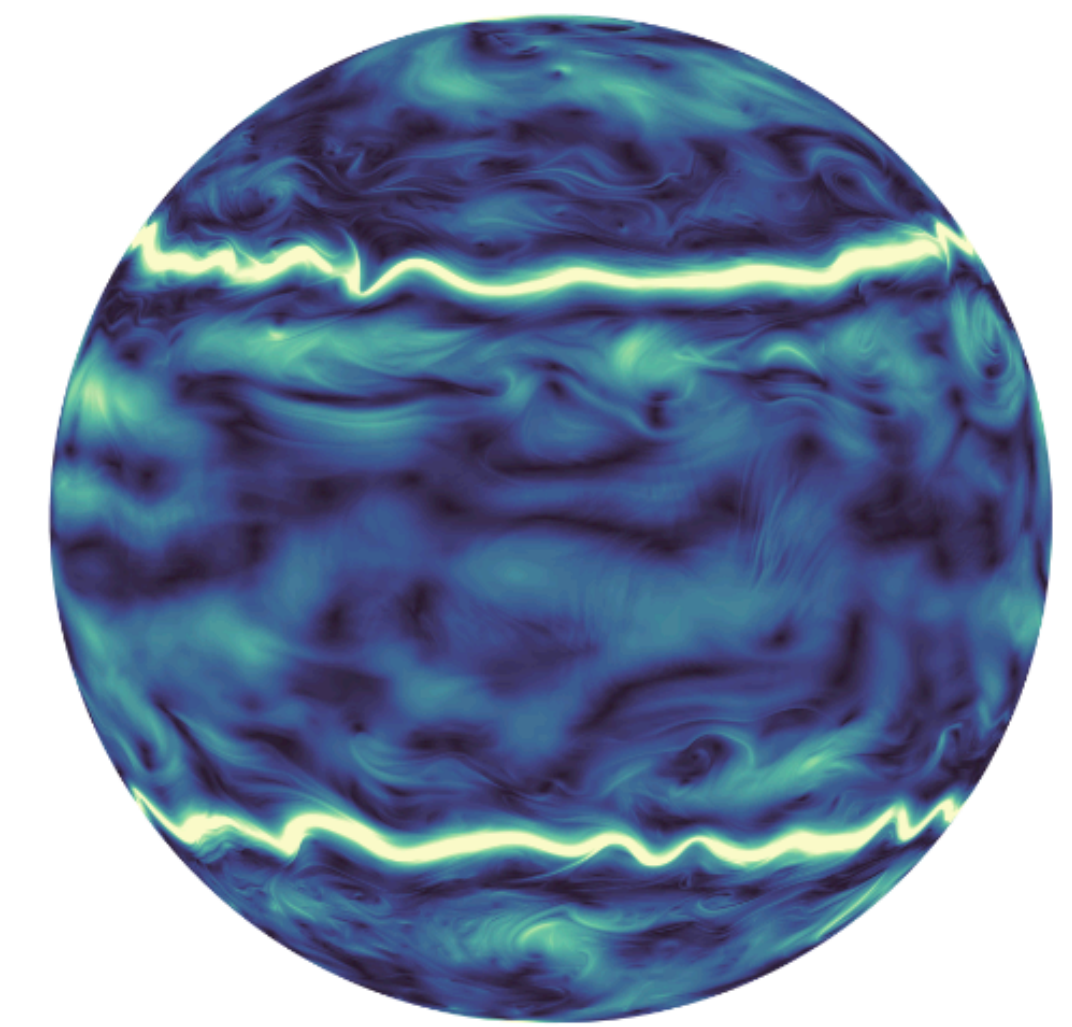
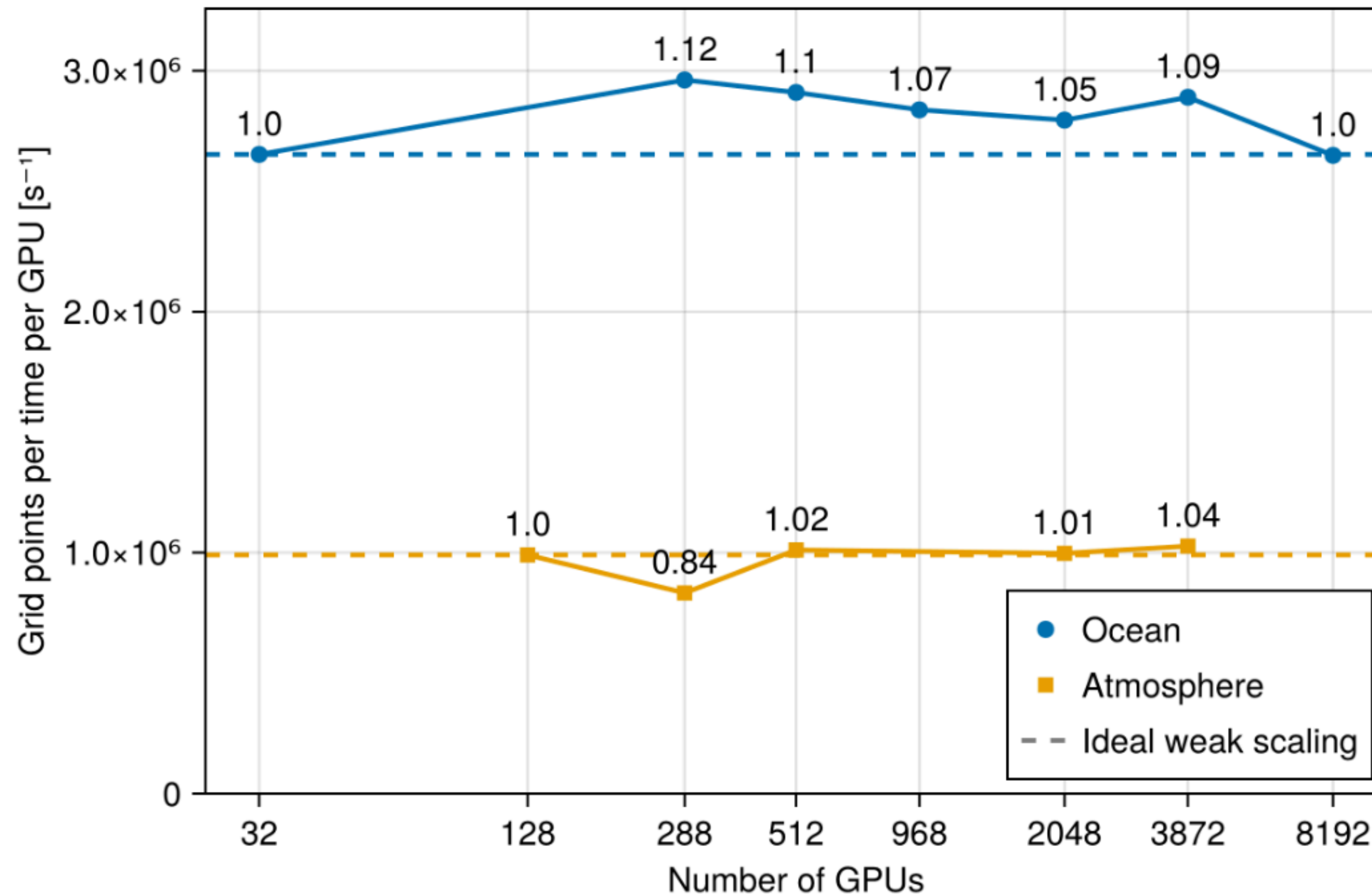
**NVIDIA**®



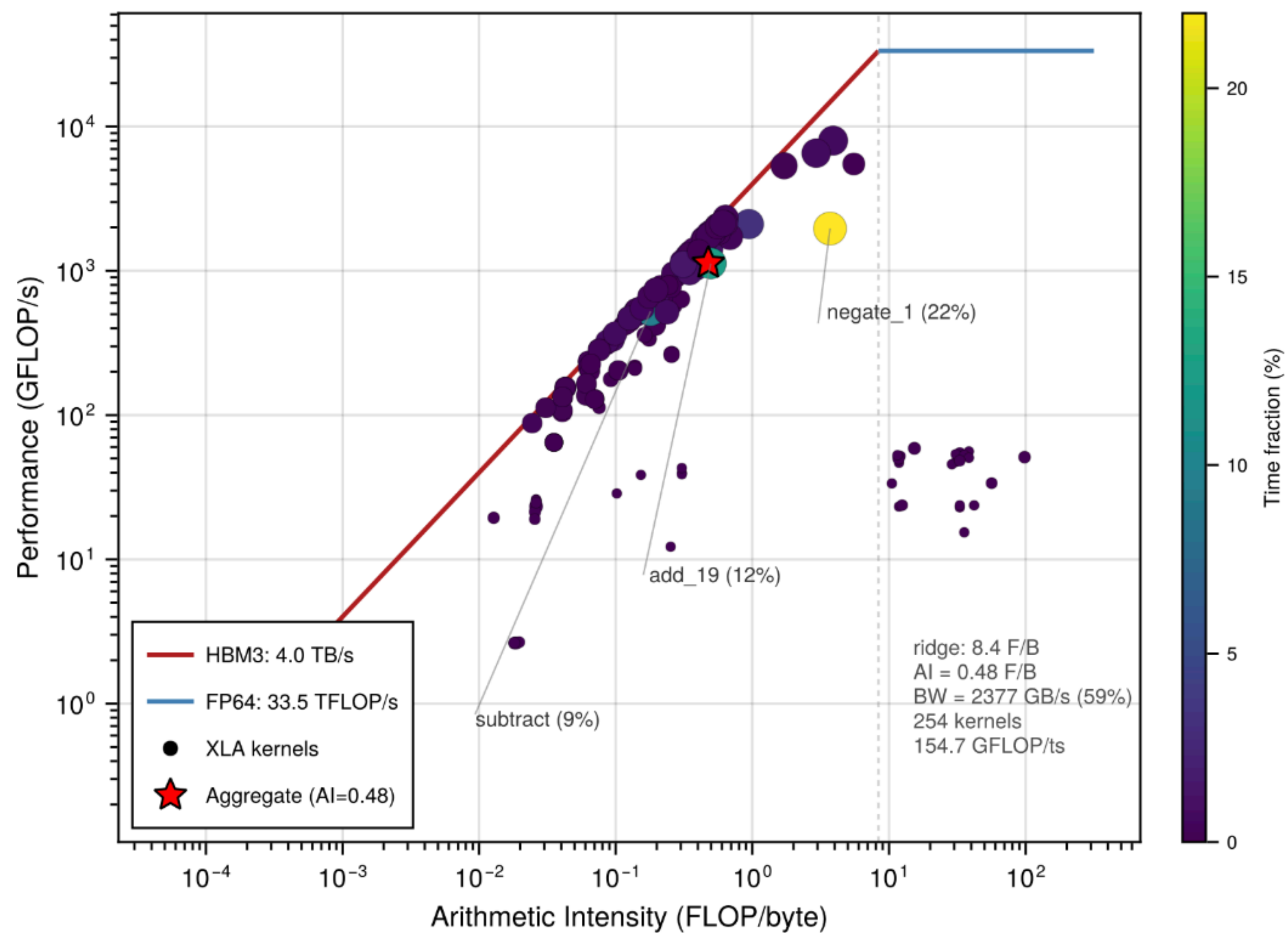
Google Cloud



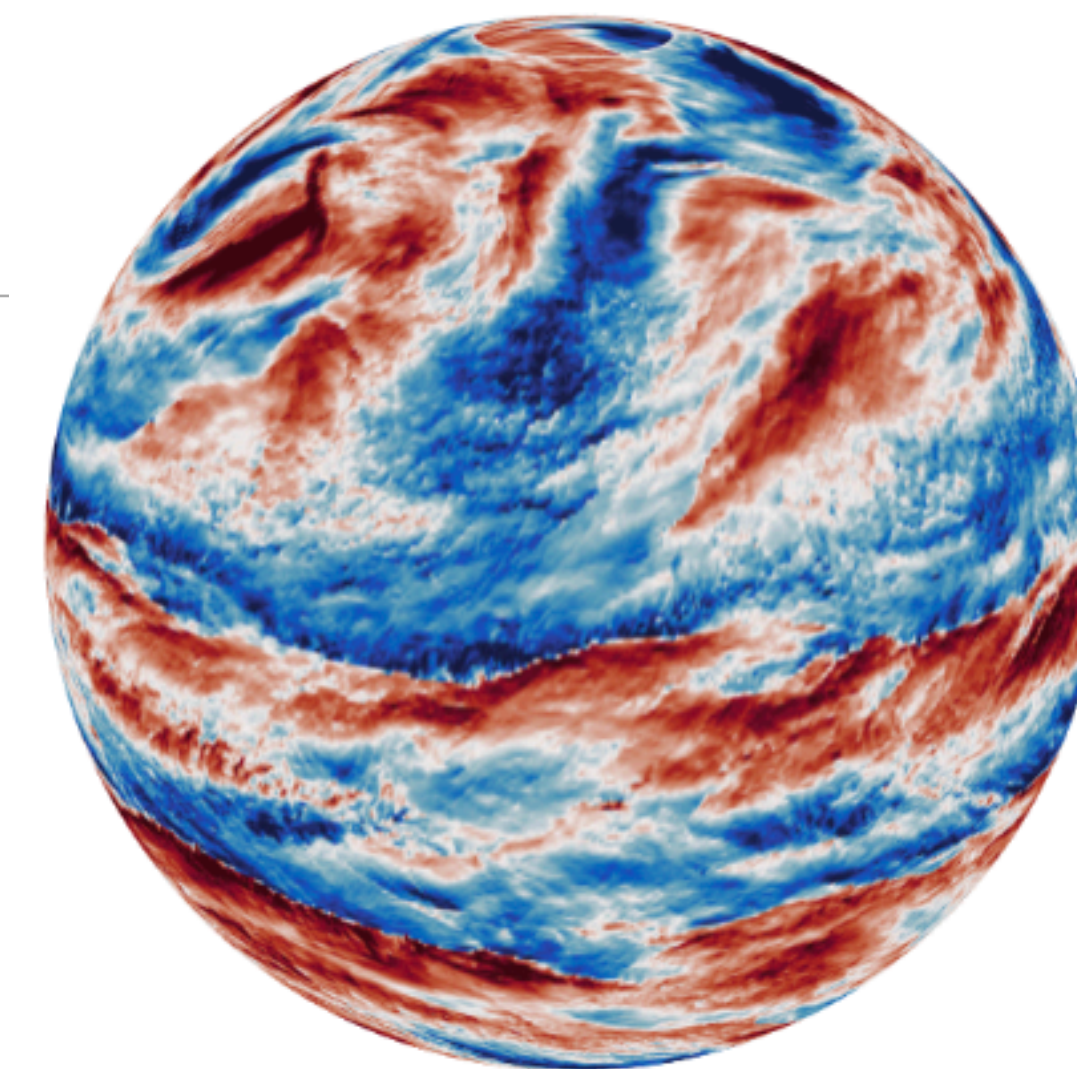
# Performance Results (Alps / NVIDIA GPU)



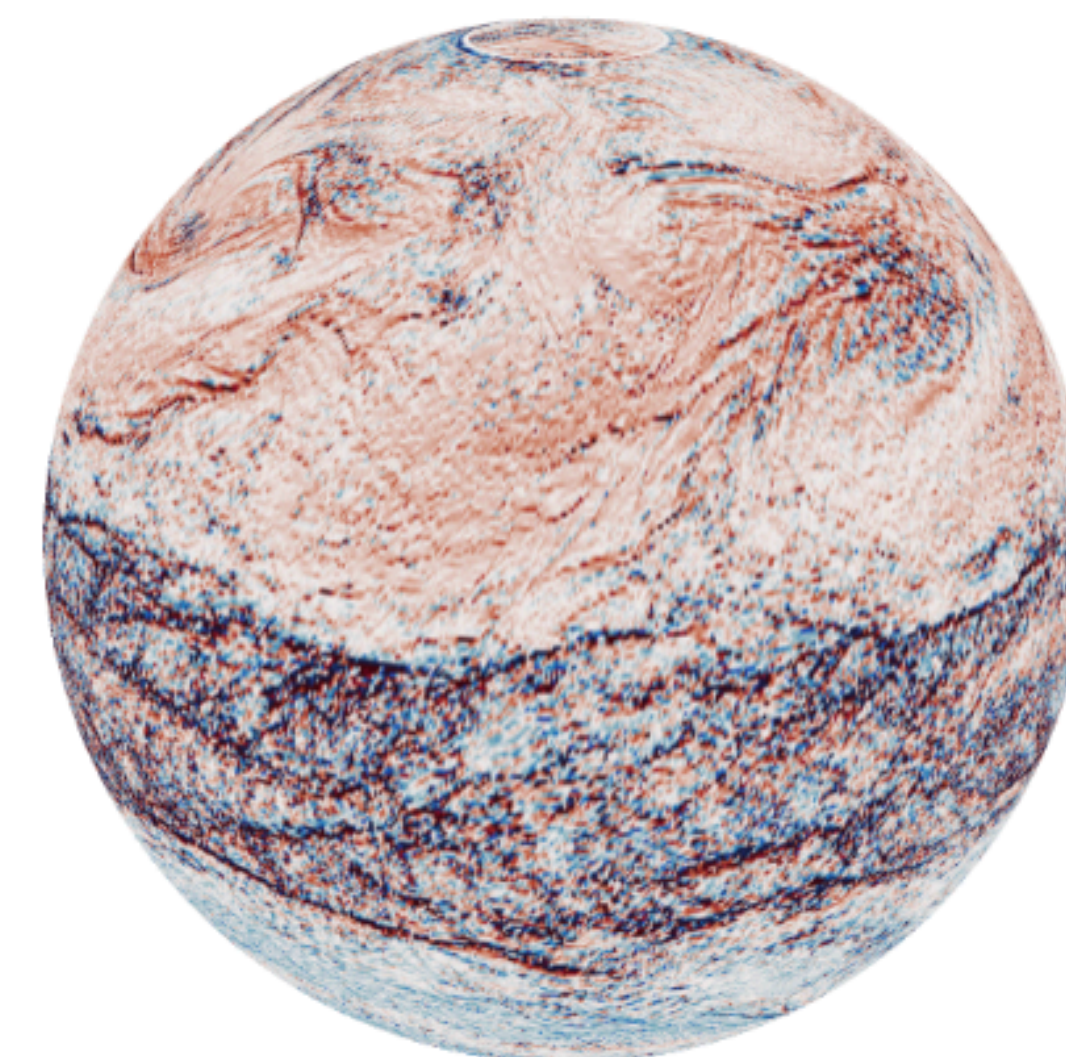
# Performance Results (Alps / NVIDIA GPU)



Meridional momentum,  $\rho v$

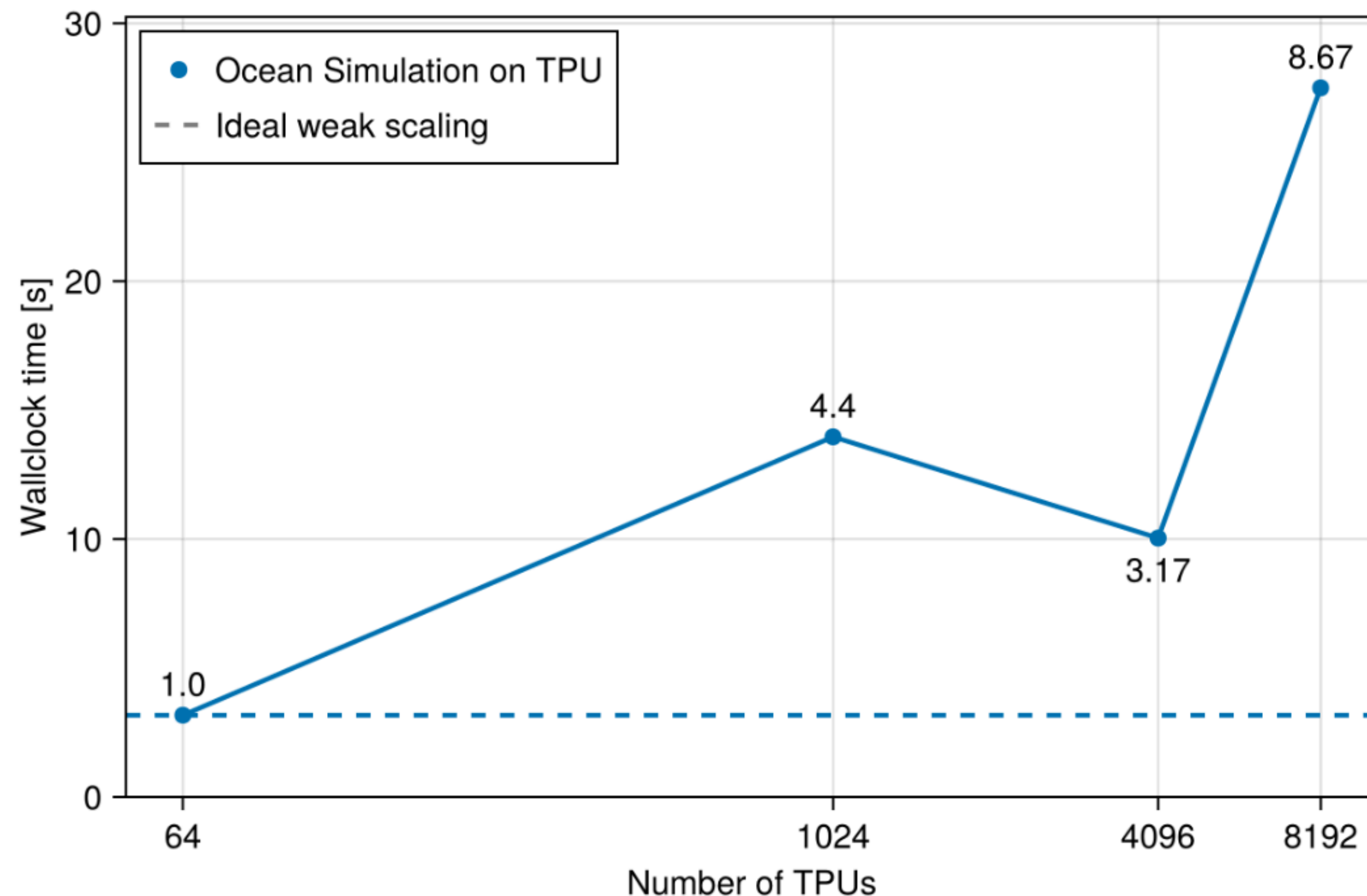


Sensitivity  $\partial J / \partial(\rho v_0)$   
 $J = V^{-1} \int (\rho \theta)^2 dV$



# Performance Results (Google Cloud / TPUs)

- At peak utilized of 97% of the TPU cluster, containing 36 FP8 exaFLOP!
- 1.493 PiB of HBM to store fields!
- 8x128 registers created some additional comms existing patterns don't yet handle (in progress)



Operation	Percent of Execution
Concatenate	39.04%
Reduce-Window	35.01%
Loop-Fusion 1	19.71%
Data Formatting	2.89%
Slice	1.59%
X64Combine	0.88%
Collective-Permute	0.48%

# Conclusions

---

- Computing hardware is increasingly moving to domain-specific accelerators, leaving existing scientific workloads in the dust
- New tool to extract the existing accelerator-friendly tensor operators written in existing parallel code and run them on distributed accelerators
- Opens the door for moving workloads to where you want to run them, without needing to re-engineer them
- Works generically on LLVM code, with explicit frontends for C++ ([github.com/EnzymeAD/Reactant](https://github.com/EnzymeAD/Reactant)) and Julia ([github.com/EnzymeAD/Reactant.jl](https://github.com/EnzymeAD/Reactant.jl))

