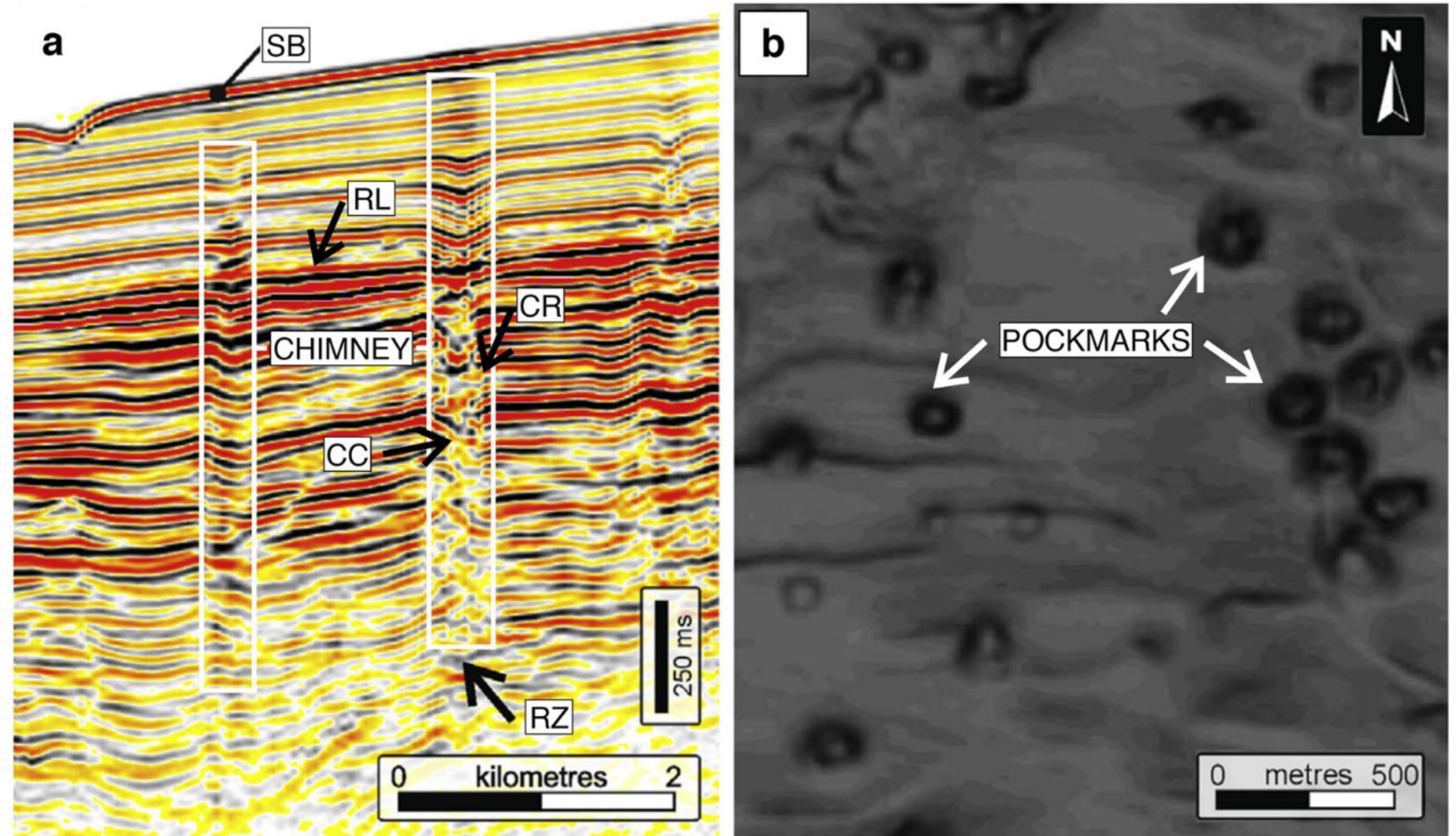
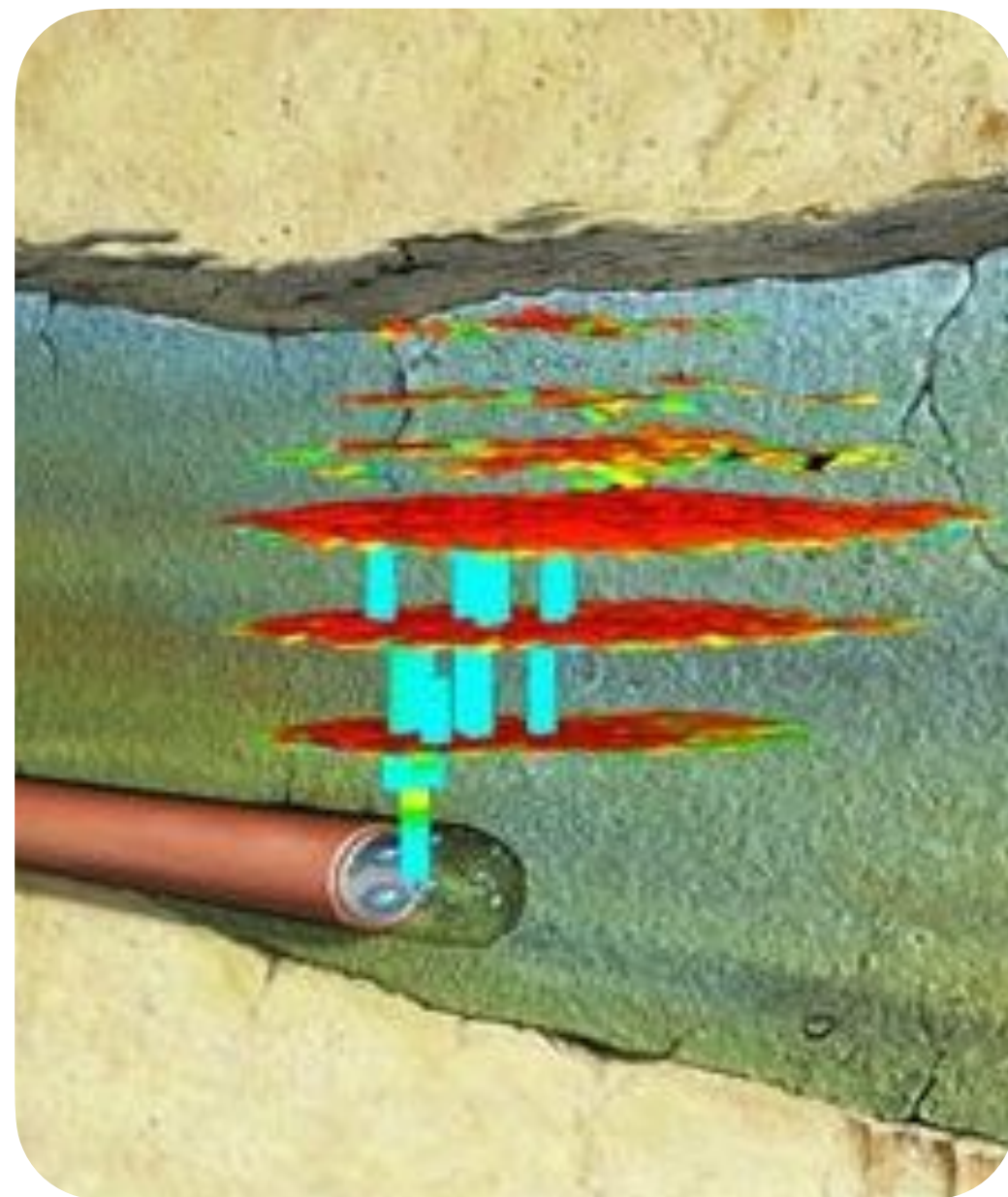


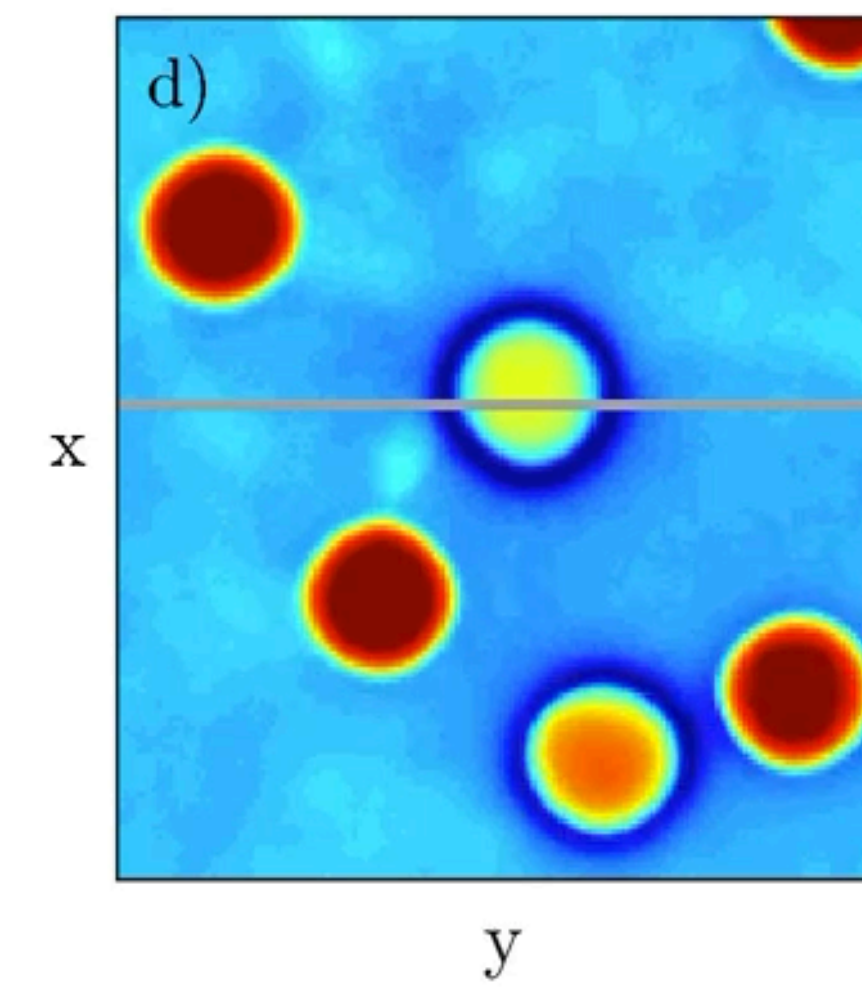
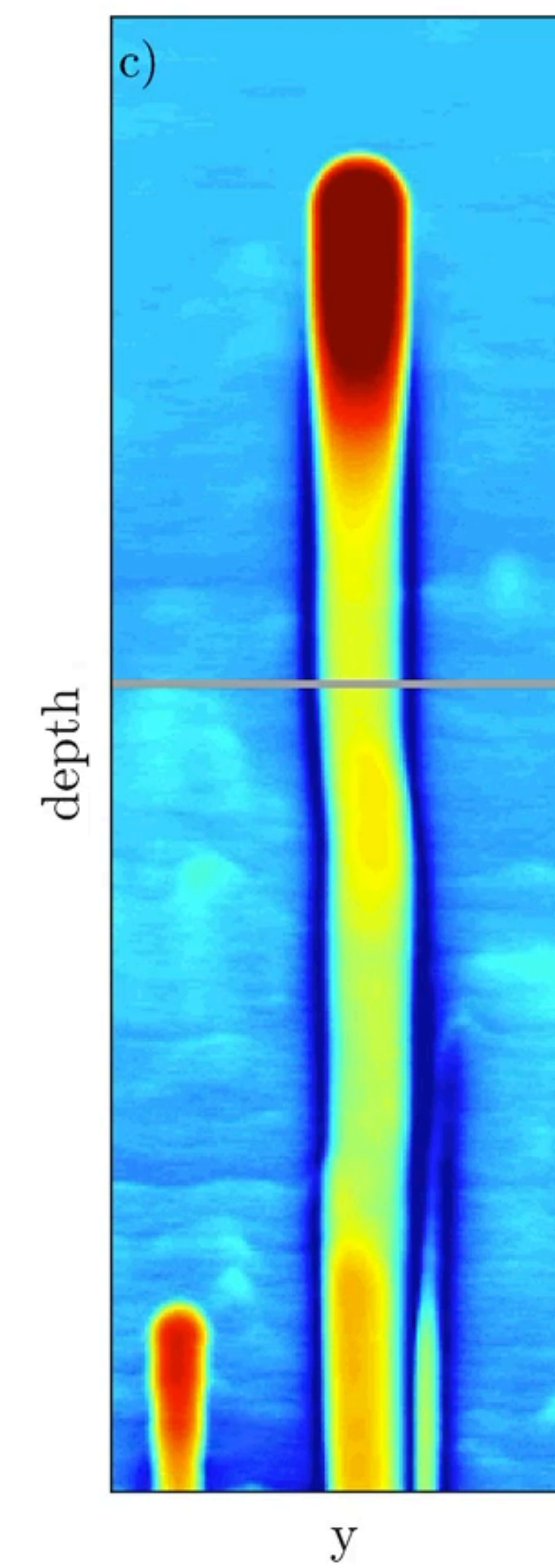
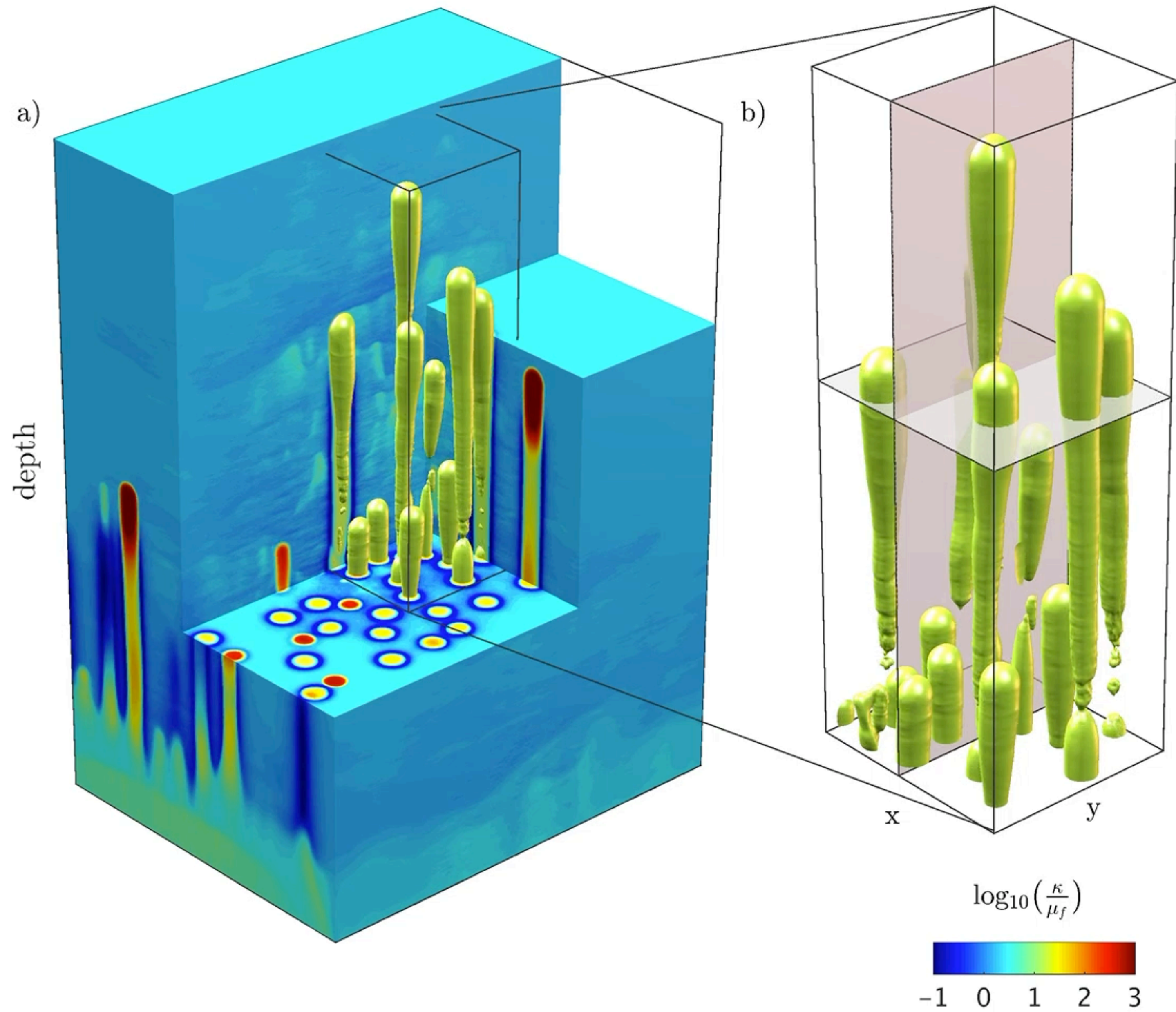
PDEs meet GPUs

Using Julia to tackle the hardware limit

GPU computing

In earth sciences





Agenda

PDEs meet GPUs

1. Performance limiters

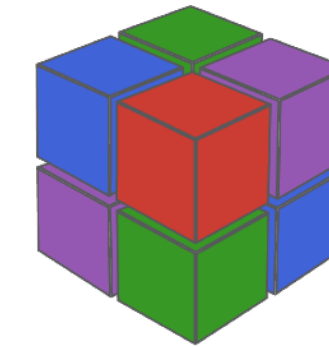
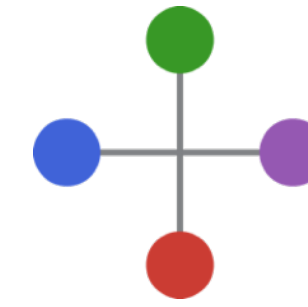
Hardware-aware methods and discretisations

2. GPUs and Julia

Portability, performance and backend abstraction

3. GPUs and Julia in HPC

Scaling from laptops to supercomputers

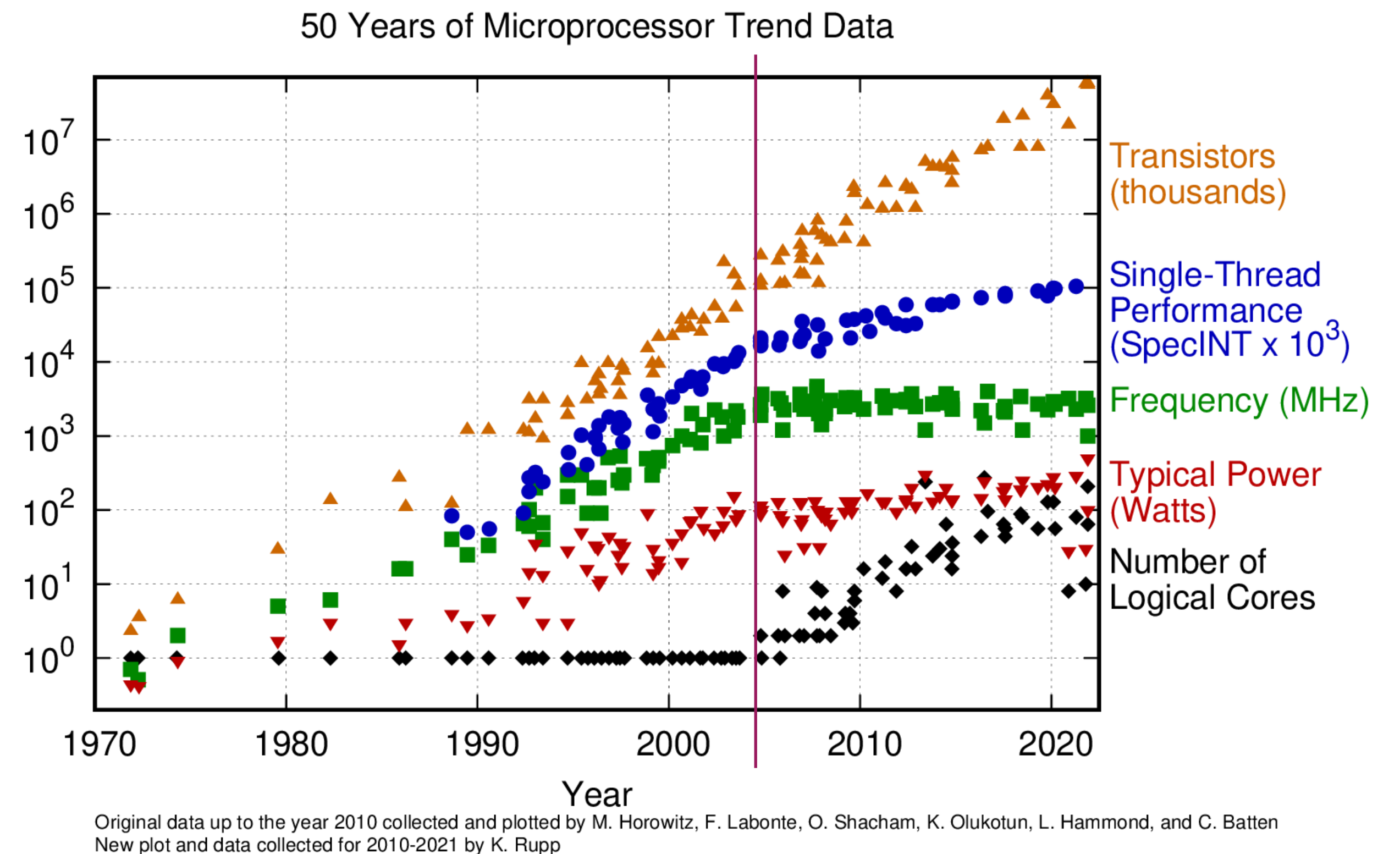


Performance limiters

Parallel computing

The memory wall

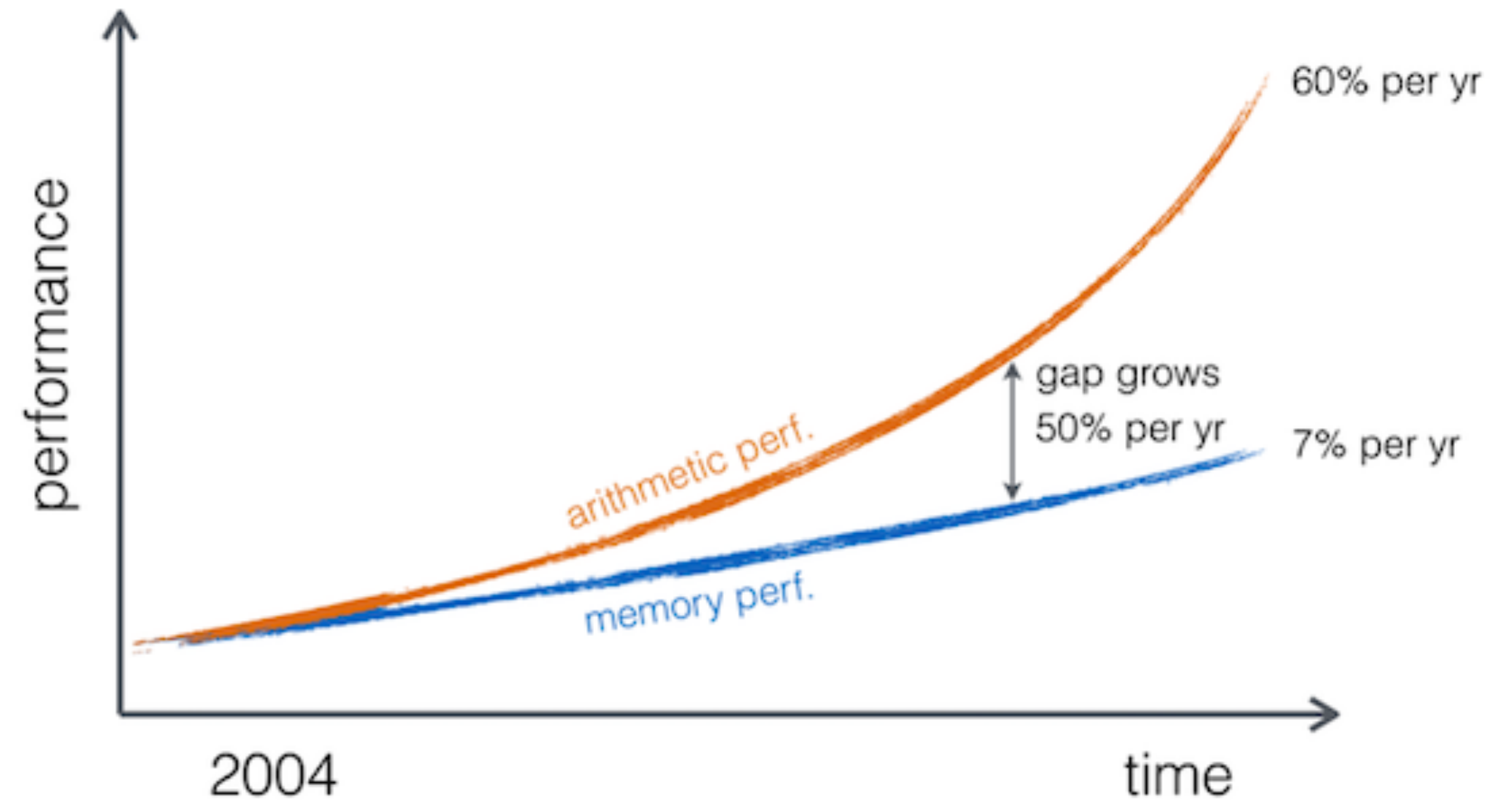
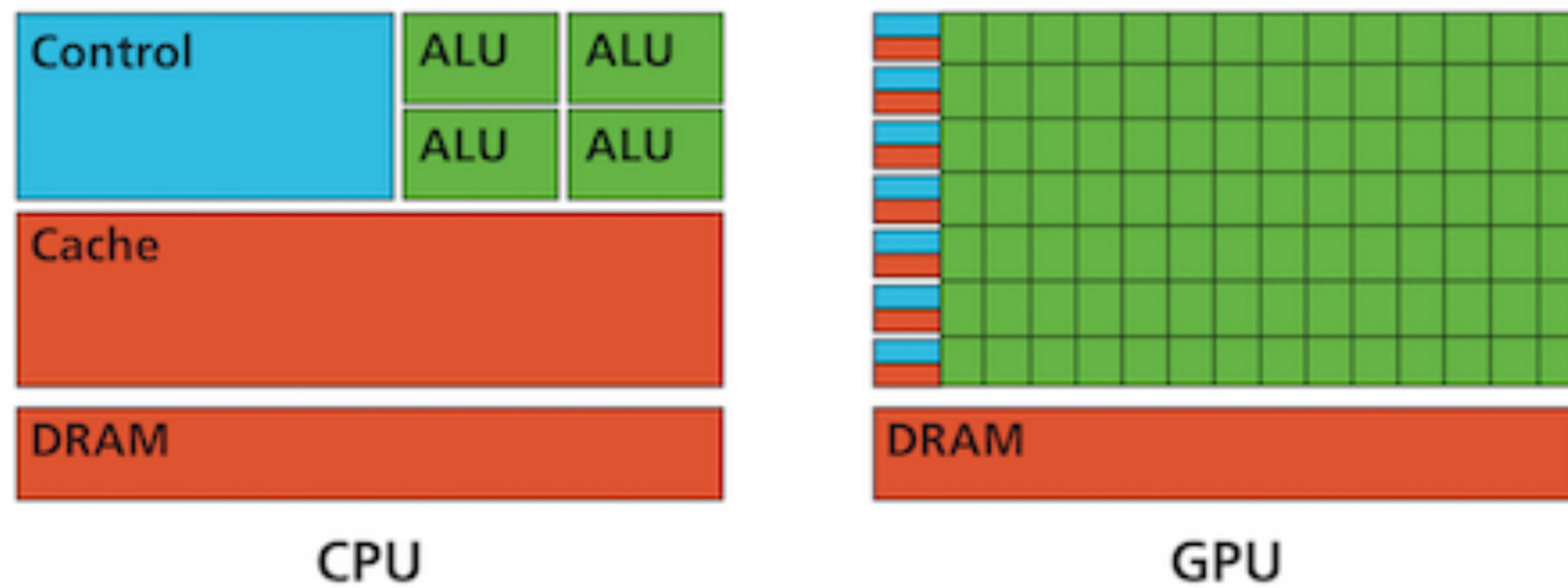
- Hitting the memory wall in 2004.
- Recent CPUs and GPUs have many cores.
- They use their parallelism to hide latency (i.e. overlapping execution times).
- Multi-core CPUs and GPUs share similar challenges.



GPU vs CPUs

Program hardware accelerators

- GPUs are massively parallel devices
- SIMD machine (programmed using threads - SPMD)



Further increases the FLOPs vs Bytes gap

Hardware limiters

CPU and GPU capabilities

- Recent GPUs and CPUs are memory bound (can do much more compute vs accessing numbers from RAM).
- Imbalance: we can do >50 FLOPs per number accessed from main memory.

Imbalance:

$$\frac{\text{Computation peak performance [TFLOP/s]}}{\text{Memory access peak performance [TB/s]}} \times \text{size of a number [Bytes]}$$

Device	TFLOP/s (FP64)	Memory BW TB/s	Imbalance
Nvidia GH200	34	4	$34/4 \times 8 = 68$
Nvidia A100	9.7	1.55	$9.7 / 1.55 \times 8 = 50$
AMD EPYC 7282 "Rome" 16 cores	0.7	0.085	$0.7 / 0.085 \times 8 = 65$

FLOPs are for free in memory bound regimes

Scientific applications

Memory bandwidth limited

- Most algorithms require only a few FLOPs ...
- Compared to the amount of bytes accessed from main memory.
- FLOPs metric is not the most adequate for reporting application performance.

First derivative example $\partial A / \partial x$:

The “cost” of an isolated evaluation of a finite-difference first derivative (e.g. computing a flux q):

$$q = -D \frac{\partial A}{\partial x}$$

Which in discrete form reads

$$q[ix] = -D * (A[ix+1] - A[ix]) / dx$$

1 read + 1 write => 2 x 8 = 16 bytes transferred

1 (fused) muladd => 1 (or 2) FLOPs

assuming D , ∂x are scalars and q , A are Float64 arrays read from main memory

Performance metric

The effective memory throughput

- Need a memory throughput-based performance evaluation metric: T_{eff} [GB/s].
- The upper-bound of T_{eff} is T_{peak}
- Defining T_{eff} we assume:
 - We evaluate an iterative stencil solver.
 - The problem size is much larger than the cache size.

The effective memory access A_{eff} [GB] is the sum of:

- Twice the memory footprint of unknown fields D_u (that depend on their own history)
- Known fields D_k that do not change during iterations

$$A_{\text{eff}} = 2D_u + D_k$$

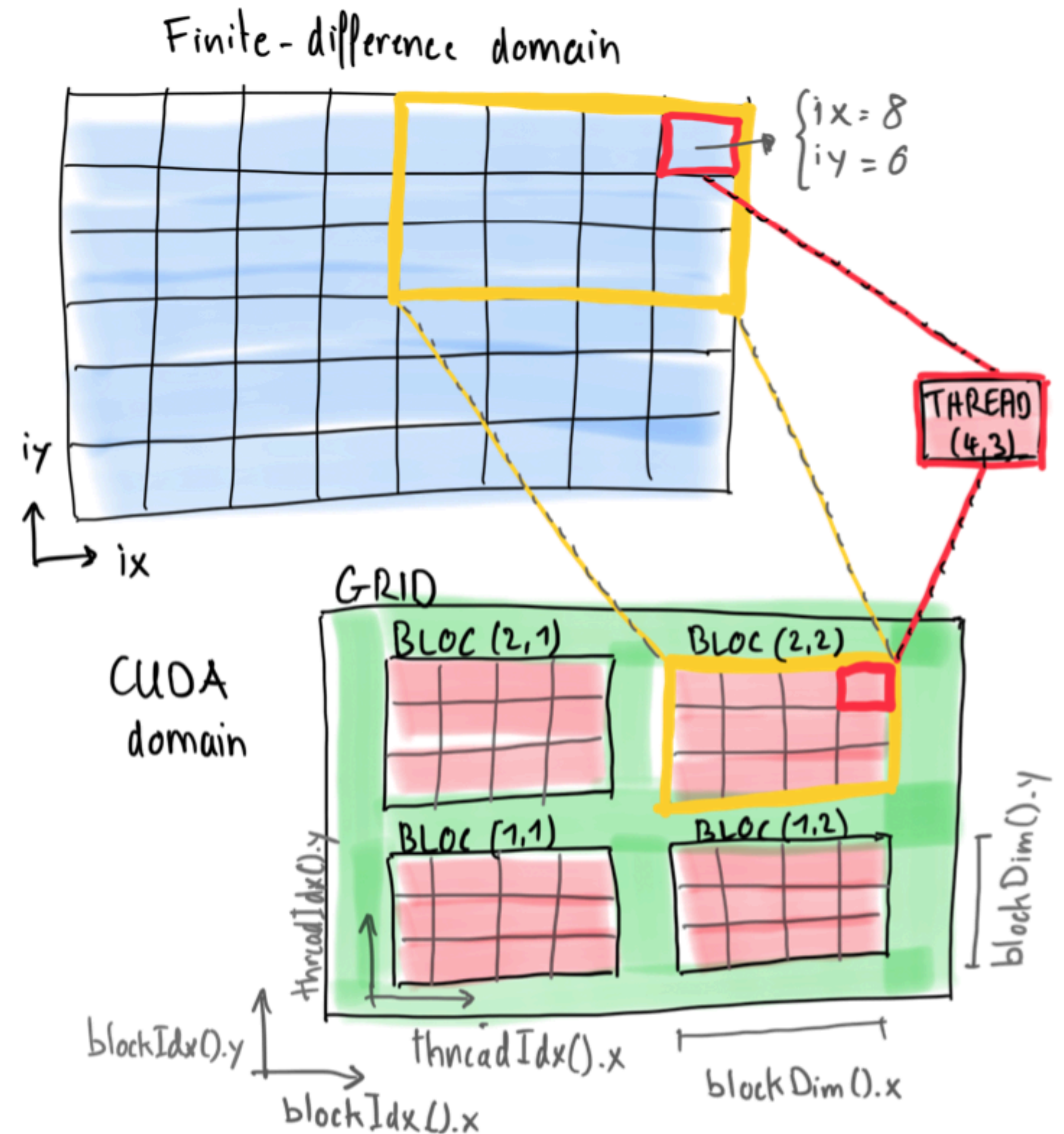
The effective memory access divided by the time per iteration t_{it} [sec] defines the effective memory throughput [GB/s]:

$$T_{\text{eff}} = \frac{A_{\text{eff}}}{t_{\text{it}}}$$

Matching the grids

Hardware-aware discretisation

- GPU architecture matching the computational domain.
- Finite-differences naturally map the underlying hardware layout.
- Use lower-order methods to limit memory traffic.
- Ensure accuracy by high spatial and temporal resolution.



Pseudo-transient Matrix-free solvers

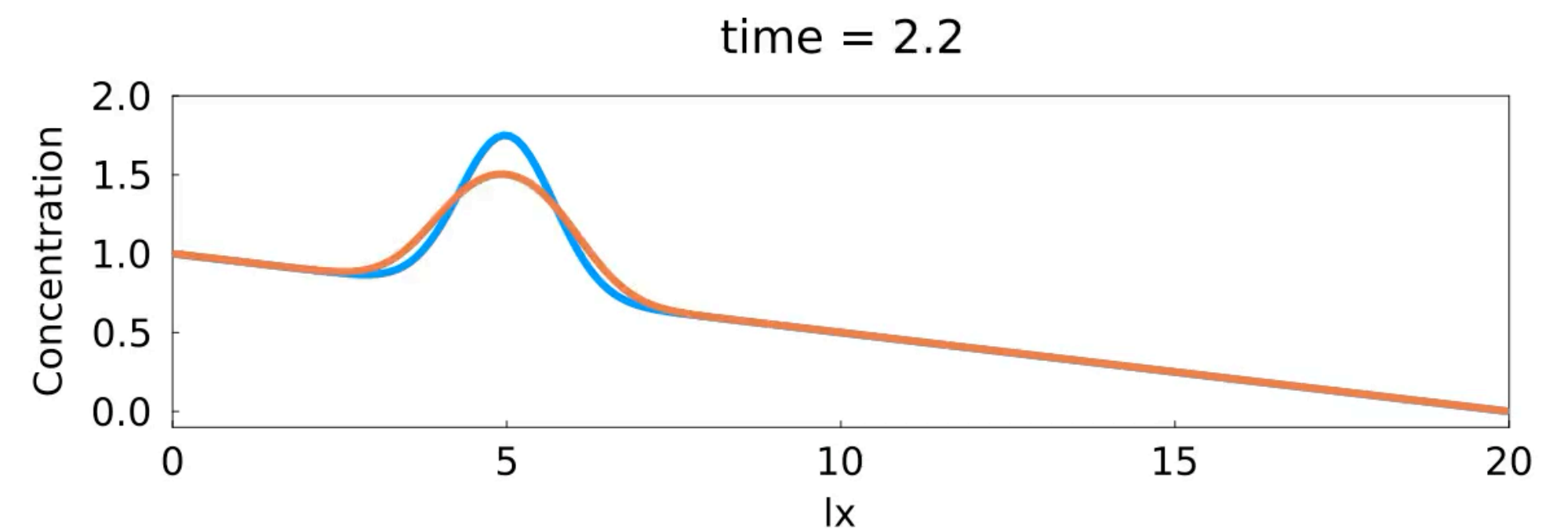
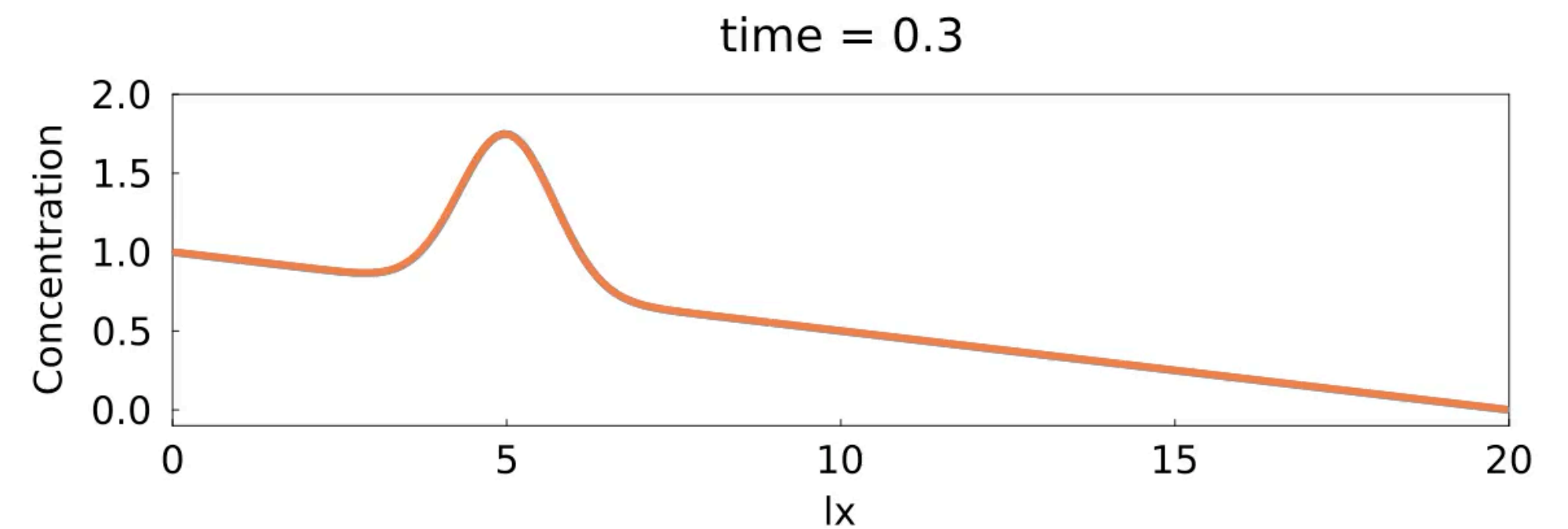
- Fast and implicit elliptic solvers.
- Analogy to transient physics with damping.
- Local algorithms, no global reductions.

$$\frac{\partial C}{\partial \tau} = \frac{\partial^2 C}{\partial x^2} \quad \rho \frac{\partial^2 C}{\partial \tau^2} + \frac{\partial C}{\partial \tau} = \frac{\partial^2 C}{\partial x^2}$$

- Challenge on finding the iteration parameter ρ (estimating $\lambda_{\min}, \lambda_{\max}$)

<https://doi.org/10.5194/gmd-15-5757-2022>

<https://doi.org/10.5194/egusphere-2025-5641>



<https://pde-on-gpu.vaw.ethz.ch/lecture3/>

GPUs and Julia

Julia GPU

Portability and performance

- Supports array and kernel programming.
- Integration with vendor libs such as cu/rocBLAS, cu/rocFFT, cu/rocSPARSE and more ...
- Integration with dev tools like Nsight, CUPTI, rocprof.
- Support various toolkit versions with automatic installation through JLLs.

```
using CUDA
const dim = 100_000_000
const a = 3.1415

x = CUDA.ones(dim)
y = CUDA.ones(dim)
z = CUDA.zeros(dim)

# (a) SAXPY via high-level broadcasting
CUDA.@sync z .= a .* x .+ y

# (b) SAXPY via CUBLAS
CUDA.@sync CUBLAS.axpy!(dim, a, x, y)


# (c) SAXPY via CUDA kernel
function saxpy_gpu_kernel!(z, a, x, y)
    i = (blockIdx().x - 1) * blockDim().x +
        threadIdx().x
    if i <= length(z)
        @inbounds z[i] = a * x[i] + y[i]
    end
    return nothing
end

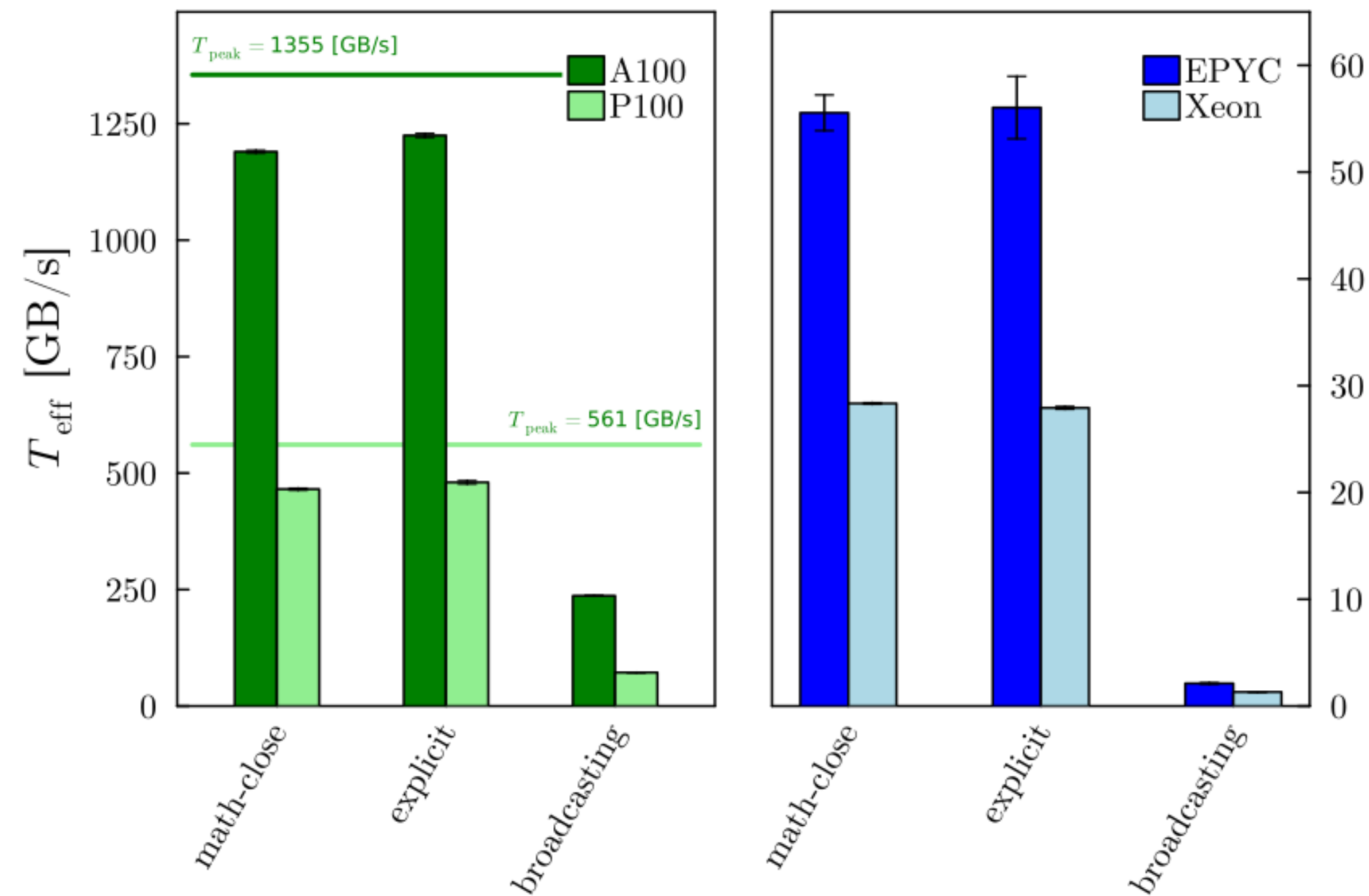
# launch configuration
nthreads = 1024
nblocks = cld(dim, nthreads)

# execute the kernel
CUDA.@sync @cuda(
    threads = nthreads,
    blocks = nblocks,
    saxpy_gpu_kernel!(z, a, x, y)
)
```

Backend abstraction

Stencil computations

 ParallelStencil.jl enables parallel stencil computations on GPUs & CPUs.



<https://github.com/omlins/ParallelStencil.jl>

Benchmark application: 3D heat diffusion solver

```

1 using ParallelStencil
2 using ParallelStencil.FiniteDifferences3D
3 @init_parallel_stencil(CUDA, Float64, 3)
4
5 @parallel memopt=true optvars=T function step!(
6     T2, T, Ci, lam, dt, _dx, _dy, _dz)
7     @inn(T2) = @inn(T) + dt*(
8         lam*@inn(Ci)*(@d2_xi(T)*_dx^2 +
9             @d2_yi(T)*_dy^2 +
10            @d2_zi(T)*_dz^2 ) )
11     return
12 end
13
14 function diffusion3D()
15     # Physics
16     lam      = 1.0           #Thermal conductivity
17     c0       = 2.0           #Heat capacity
18     lx=ly=lz = 1.0           #Domain length x|y|z
19
20     # Numerics
21     nx=ny=nz = 512           #Nb gridpoints x|y|z
22     nt       = 100          #Nb time steps
23     dx       = lx/(nx-1)    #Space step in x
24     dy       = ly/(ny-1)    #Space step in y
25     dz       = lz/(nz-1)    #Space step in z
26     _dx, _dy, _dz = 1.0/dx, 1.0/dy, 1.0/dz
27
28     # Initial conditions
29     T = @ones(nx,ny,nz).*1.7 #Temperature
30     T2 = copy(T)             #Temperature (2nd)
31     Ci = @ones(nx,ny,nz)./c0 #1/Heat capacity
32
33     # Time loop
34     dt = min(dx^2,dy^2,dz^2)/lam/maximum(Ci)/6.1
35     for it = 1:nt
36         @parallel memopt=true step!(
37             T2, T, Ci, lam, dt, _dx, _dy, _dz)
38             T, T2 = T2, T
39         end
40     end
41 end
42
43 diffusion3D()

```

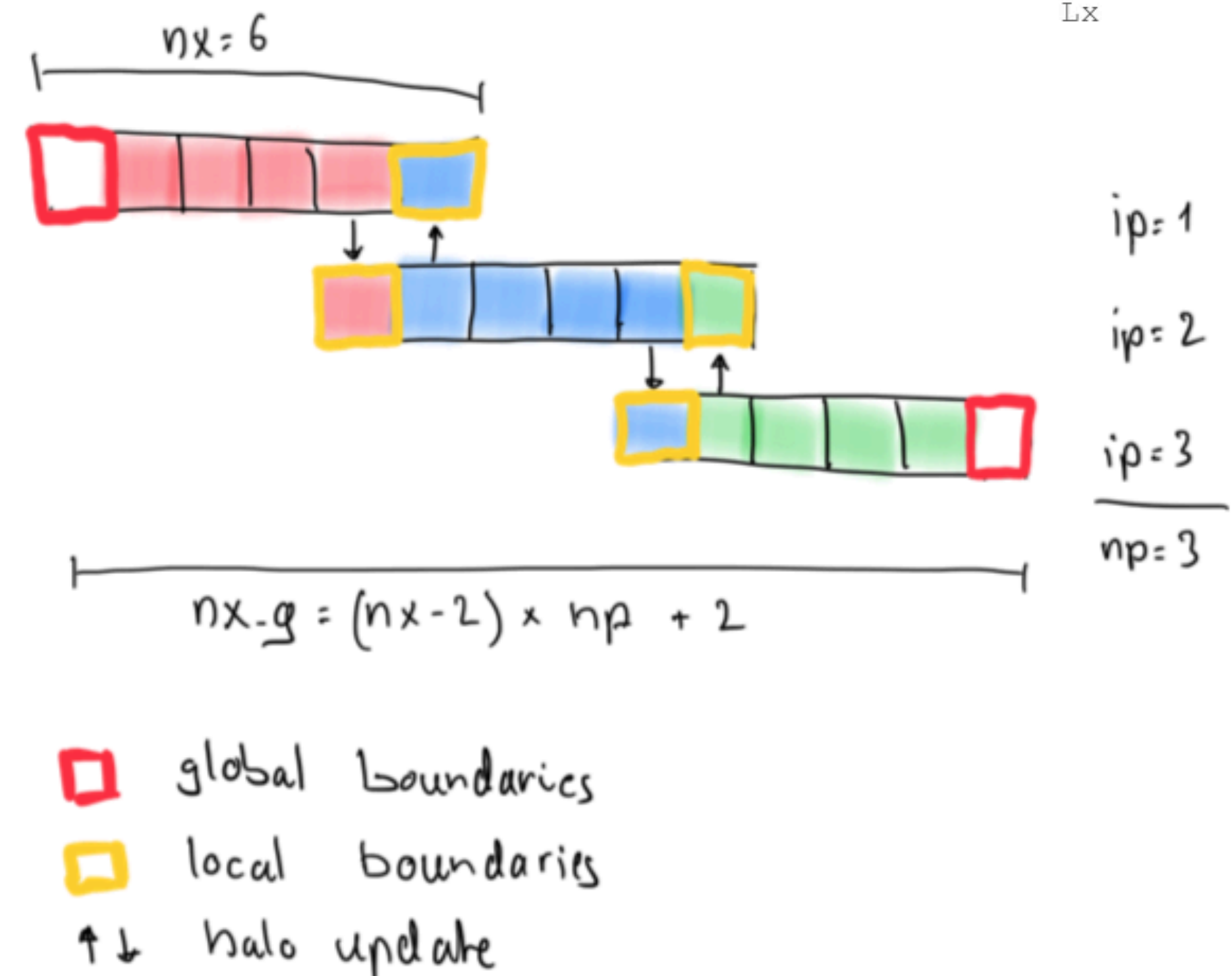
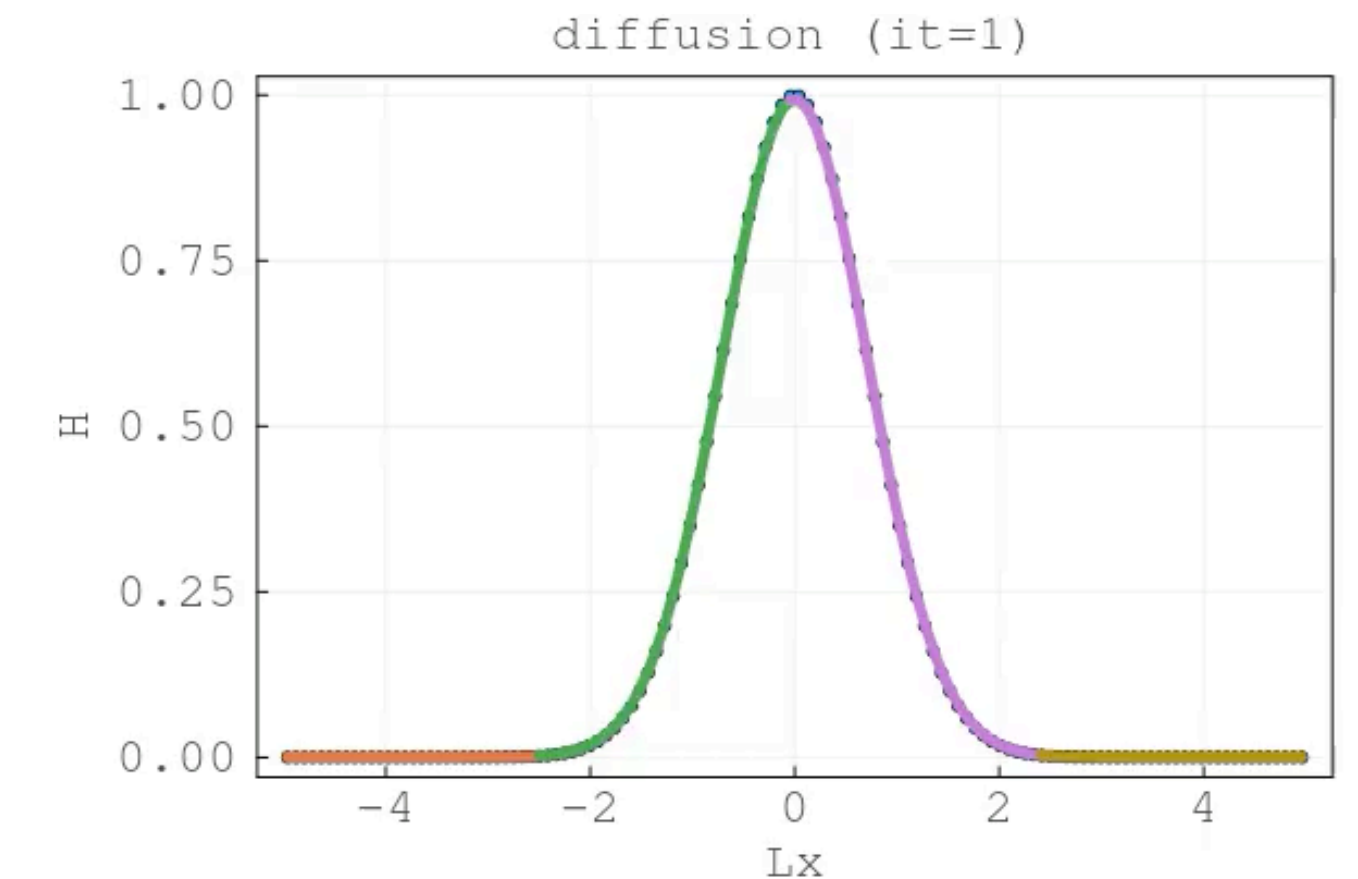
<https://doi.org/10.21105/jcon.00138>

GPUs and Julia in HPC

Multi-GPUs

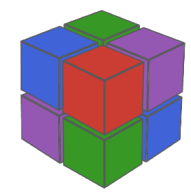
Julia and distributed memory computing

- Leveraging more than a single GPU
- Distributed memory computing
- GPU-aware MPI
- Implicit global grid
- Halo exchange
- Latency hiding (communication — computation overlap)



Distributed computing

Julia and distributed memory computing



ImplicitGlobalGrid.jl

- Relies on Julia MPI wrapper MPI.jl (could also use NCCL or other transport library).
- Only 3 additional lines needed to massively parallelise the previous example! ←
- Supports GPU-aware MPI (CUDA, ROCm).



<https://github.com/eth-cscs/ImplicitGlobalGrid.jl>

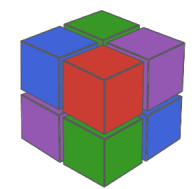
Benchmark application: 3D heat diffusion solver

```
1 using ImplicitGlobalGrid
2 using ParallelStencil
3 using ParallelStencil.FiniteDifferences3D
4 @init_parallel_stencil(CUDA, Float64, 3)
5
6 @parallel function step!(T2,T,Ci,lam,dt,dx,dy,dz)
7     @inn(T2) = @inn(T) + dt*(
8         lam*@inn(Ci)*(@d2_xi(T)/dx^2 +
9             @d2_yi(T)/dy^2 +
10                @d2_zi(T)/dz^2 ) )
11     return
12 end
13
14 function diffusion3D()
15     # Physics
16     lam      = 1.0           #Thermal conductivity
17     c0       = 2.0           #Heat capacity
18     lx=ly=lz = 1.0           #Domain length x|y|z
19
20     # Numerics
21     nx=ny=nz = 512           #Nb gridpoints x|y|z
22     nt       = 100           #Nb time steps
23     me,      = init_global_grid(nx, ny, nz) ←
24     dx       = lx/(nx_g()-1) #Space step in x
25     dy       = ly/(ny_g()-1) #Space step in y
26     dz       = lz/(nz_g()-1) #Space step in z
27
28     # Initial conditions
29     T = @ones(nx,ny,nz).*1.7 #Temperature
30     T2 = copy(T)             #Temperature (2nd)
31     Ci = @ones(nx,ny,nz)./c0 #1/Heat capacity
32
33     # Time loop
34     dt = min(dx^2,dy^2,dz^2)/lam/maximum(Ci)/6.1
35     for it = 1:nt
36         @hide_communication (16, 2, 2) begin
37             @parallel step!(T2,T,Ci,lam,dt,dx,dy,dz)
38             update_halo!(T2) ←
39         end
40         T, T2 = T2, T
41     end
42
43     finalize_global_grid() ←
44 end
45
46 diffusion3D()
```

<https://doi.org/10.21105/jcon.00137>

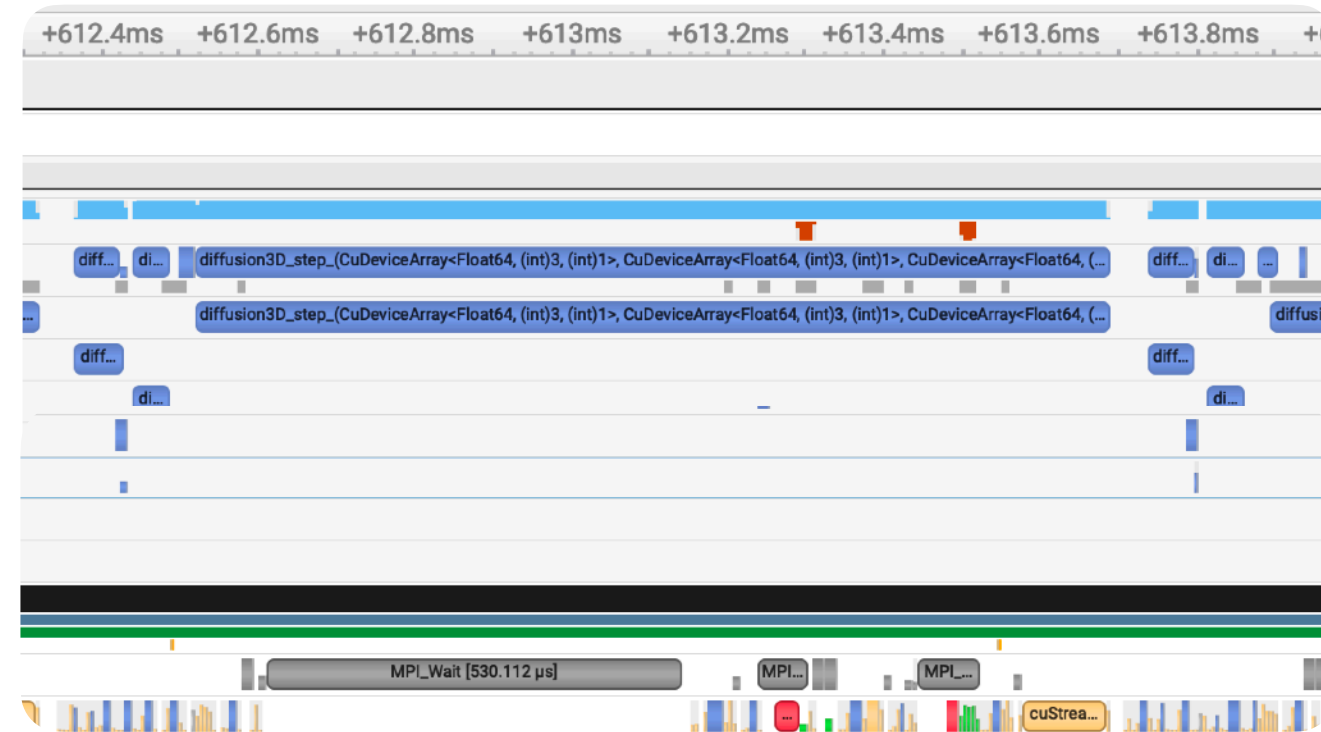
Julia on Nvidia GPUs

Scaling on ALPS



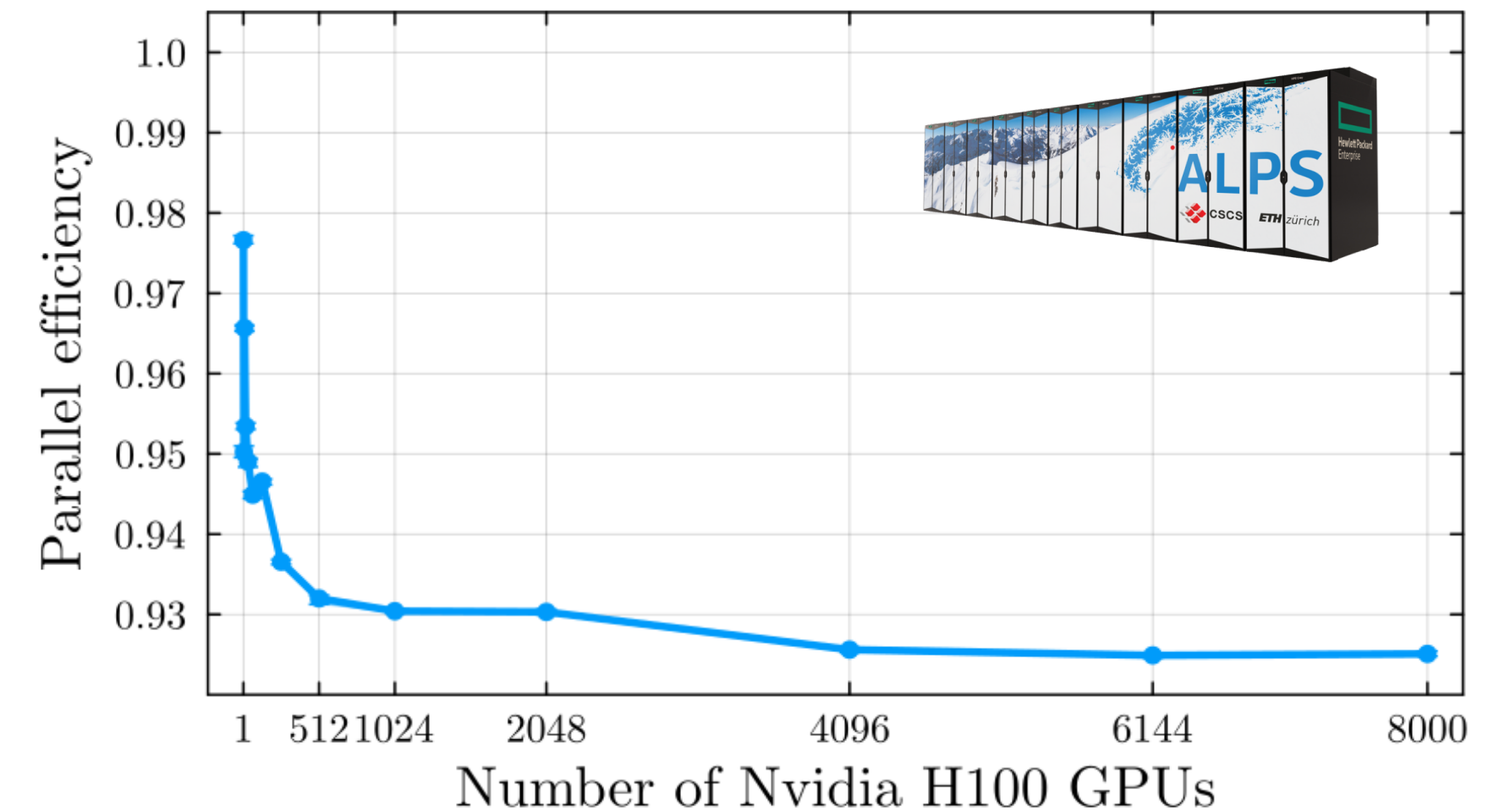
ImplicitGlobalGrid.jl

- Parallel efficiency close to 1
- Weak scaling on CSCS supercomputers
ALPS (Nvidia GH200) & previous Piz
Daint (Nvidia P100)
- Communication
latency hiding

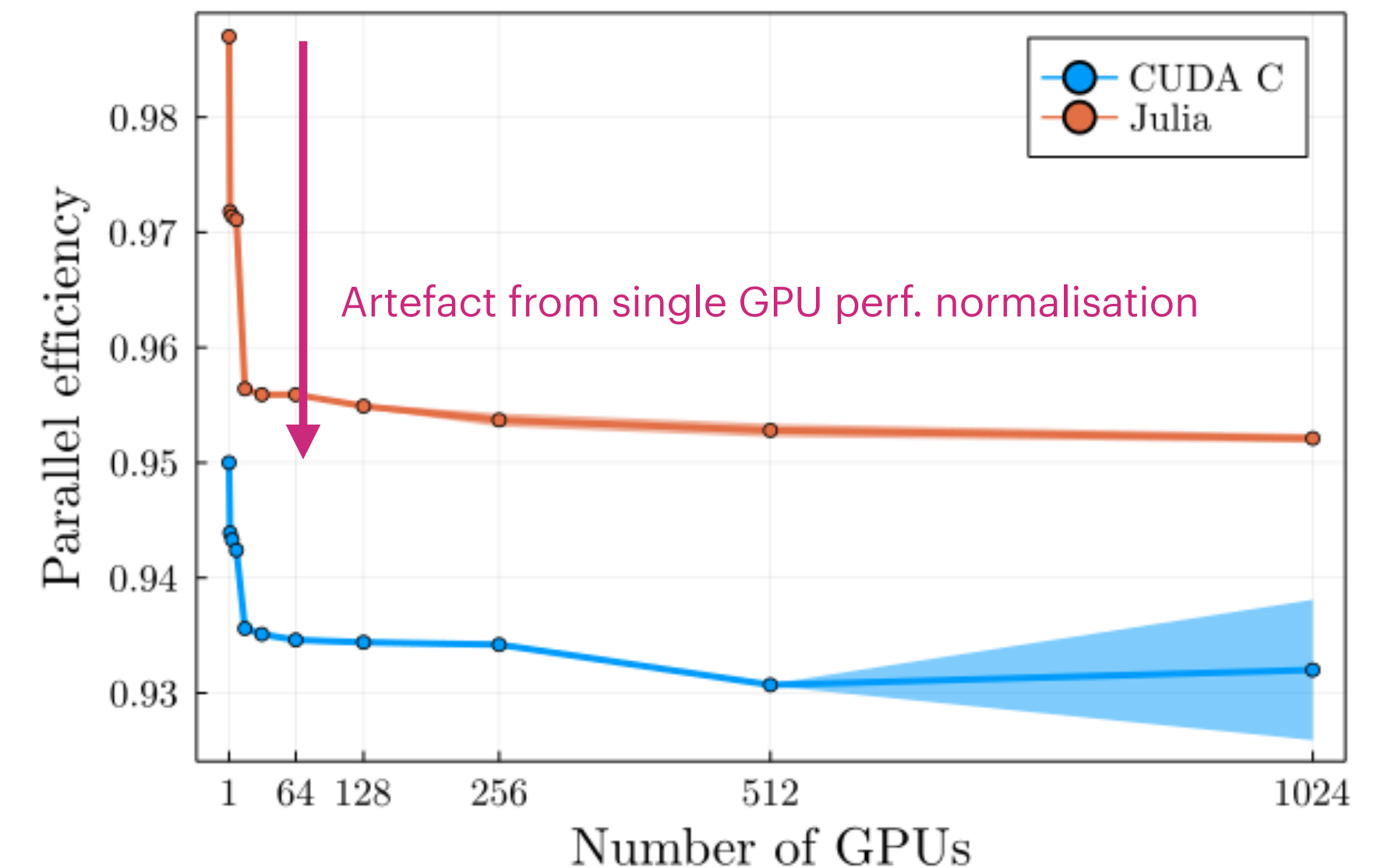


<https://github.com/eth-cscs/ImplicitGlobalGrid.jl>

Benchmark application: 3D heat diffusion solver

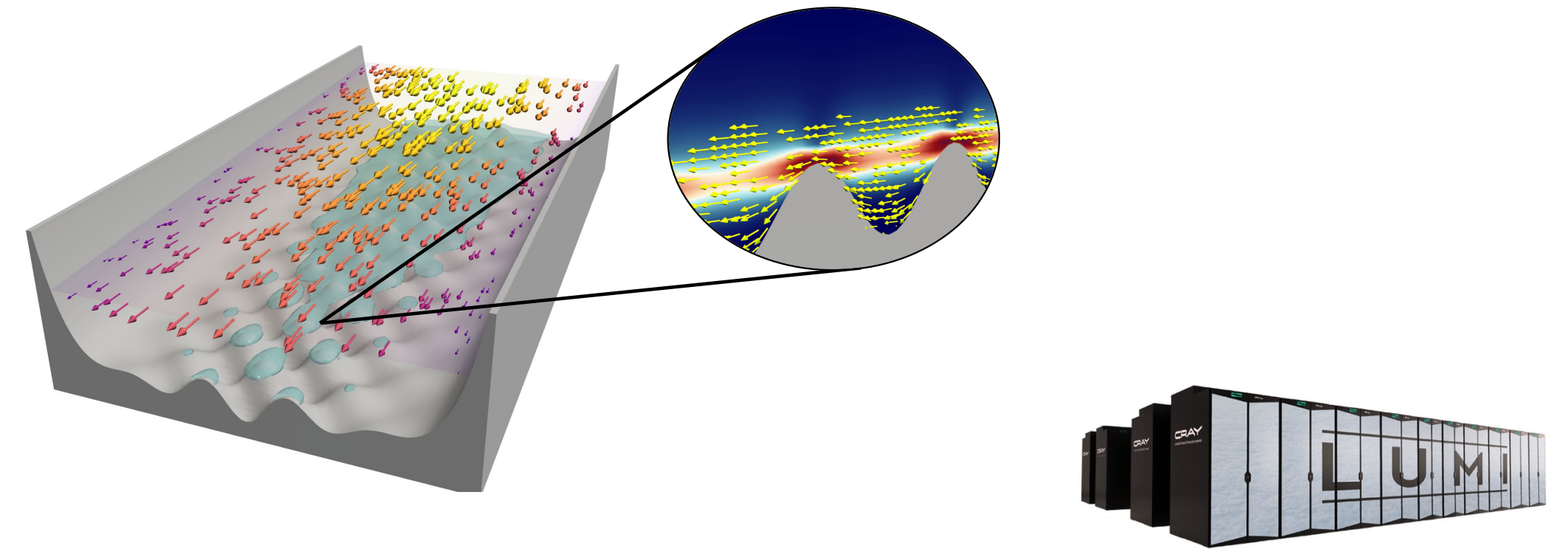



Real world application: 3D poro-visco-elastic two-phase flow solver

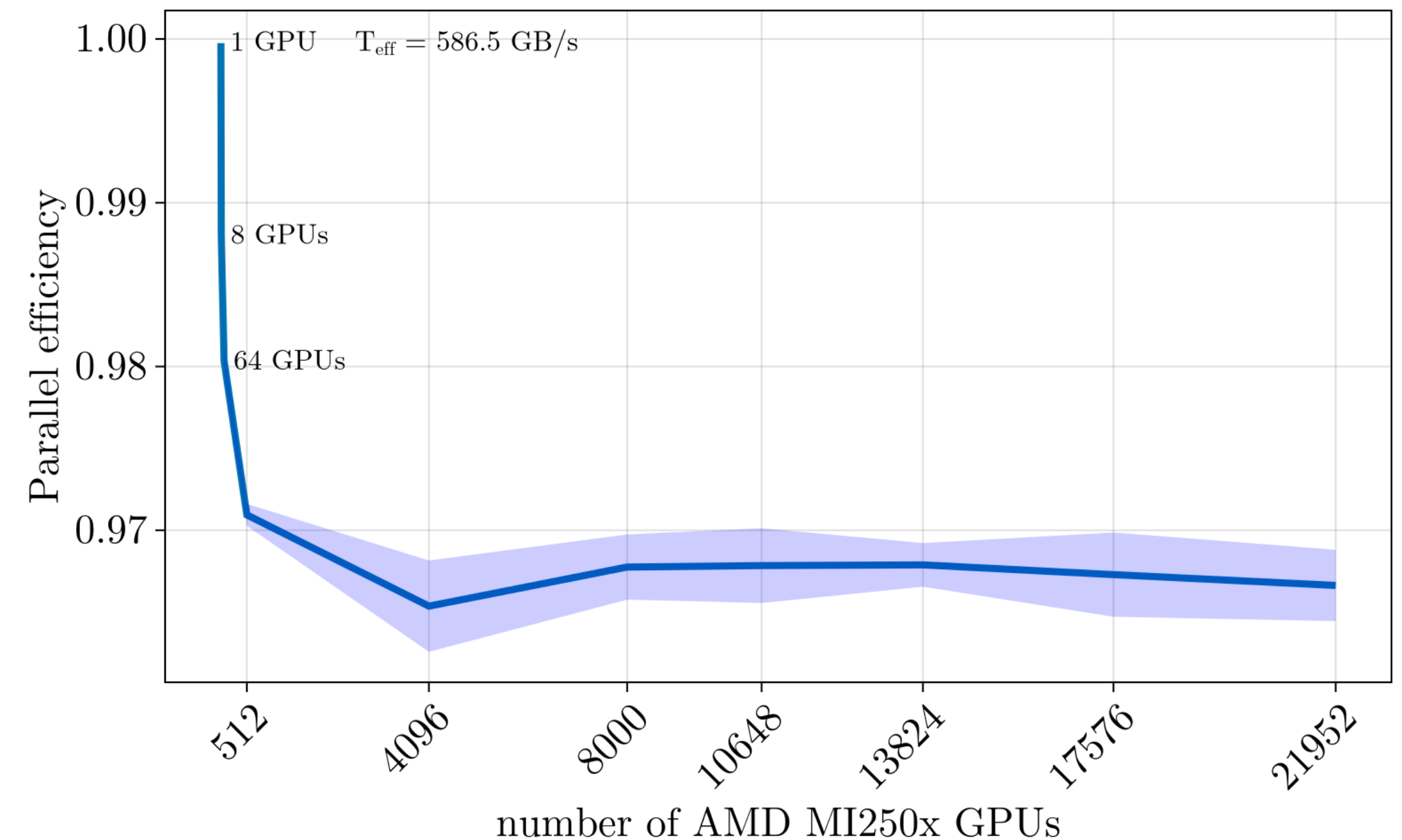


Julia on AMD GPUs

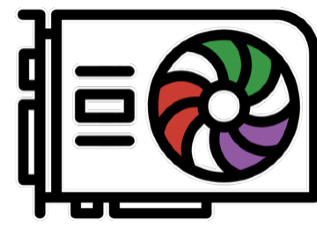
Scaling on LUMI



- 97% of parallel efficiency in weak scaling on 21'952 AMD MI250x GPUs.
- **Challenge:** AMDGPU.jl requirements to use task local approach.
-  Chmy.jl — Overlapping communication and computation using cooperative multitasking in Julia.
- Each task gets local GPU context and stream.

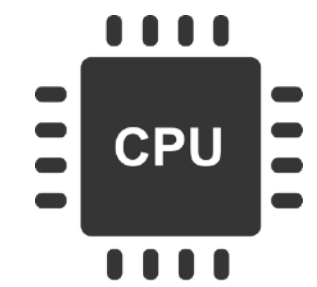


Outlook



Backends

CPUs vs GPUs vs TPUs



CPUs: ARM, x86, ...

GPUs: AMD, Nvidia, Apple, ...

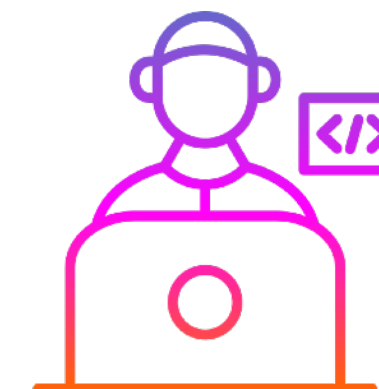


Design

Prototyping vs Production

Modularity

Abstractions

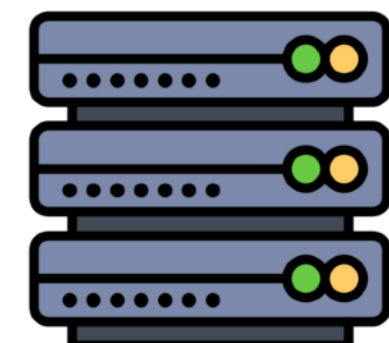


Users

Learning vs Advanced

Teaching vs Research

Timescales (PhD, project)



Compute infrastructure

Laptop, Workstation

University clusters

Supercomputers



Code conciseness

Monster codes: terrible to maintain

Worse; rewrite?

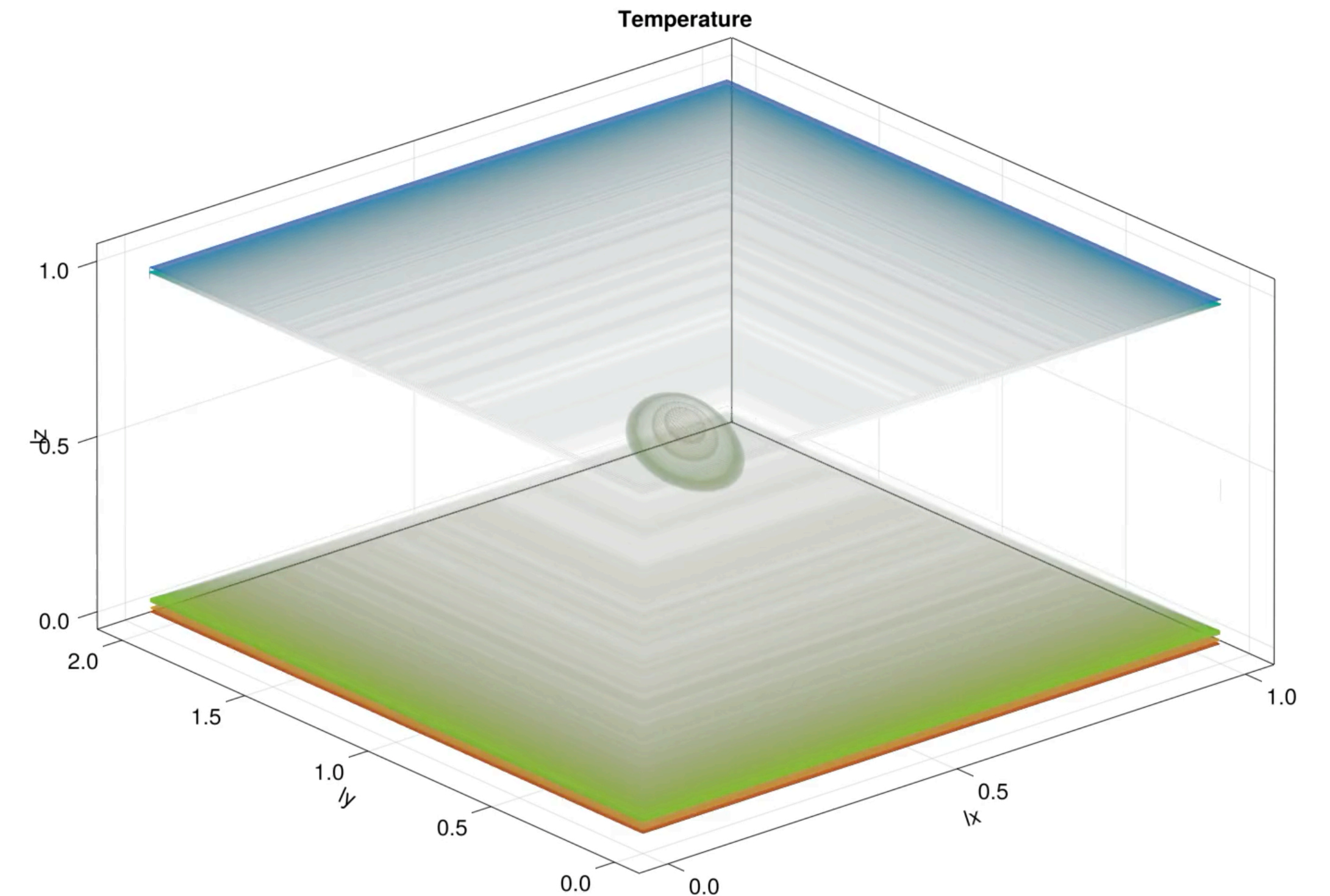
Education

Lateral spreading of supercomputing

- ETHZ course on PDEs, GPUs & HPC with Julia

The screenshot shows a web browser window with the URL pde-on-gpu.vaw.ethz.ch. The page title is "Solving PDEs in parallel on GPUs with Julia". The main content includes a welcome message, a note that the 2024 edition starts on Tuesday, Sept. 17, 12h45, and course information stating it covers state-of-the-art methods in modern parallel GPU computing. A small 2D heatmap titled "Temperature" is shown, with axes ranging from -20 to 20 and a color scale from -50 to 50. The page also features a navigation menu on the left with sections like "Welcome", "Logistics", "Homeworks", "Software install", "Extras", and "Part 1 - Introduction".

<https://pde-on-gpu.vaw.ethz.ch/>



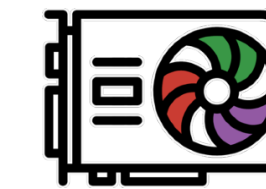
22 ETH Learning and Teaching Journal, Vol 4, No 1, 2023

**Teaching supercomputing and software engineering skills
to science and engineering students**

<https://learningteaching.ethz.ch/index.php/lt-eth/article/view/235>

∂ GPU4GEO Differentiable multi-physics solvers

for extreme-scale geophysics simulations on GPUs



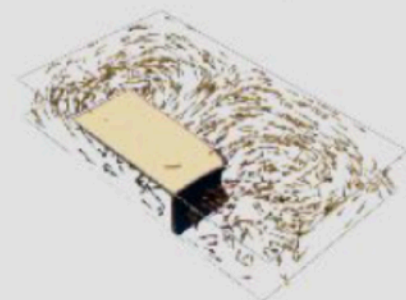
Gradient-based optimisation

GPU4GEO – forward modelling

Flagship Applications

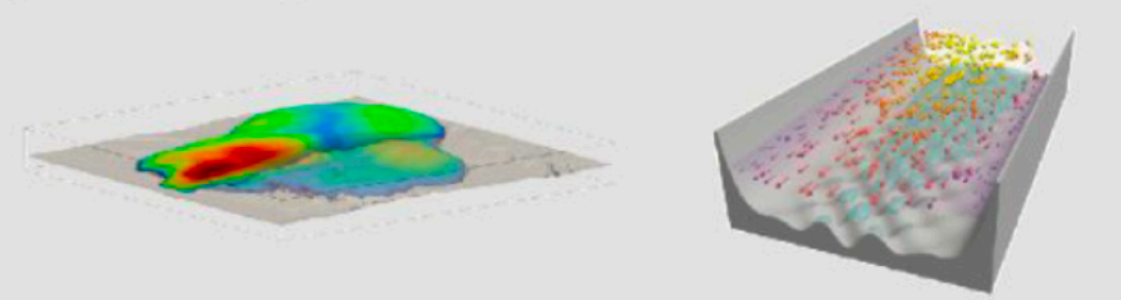
Geodynamics

JustRelax.jl



Ice flow

FastIce.jl



Forward PT solvers

$$\frac{\partial u}{\partial \tau} + R(u) = 0$$

Forward model outputs

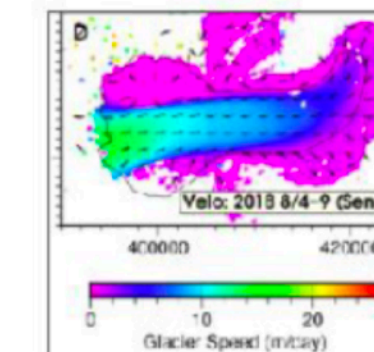
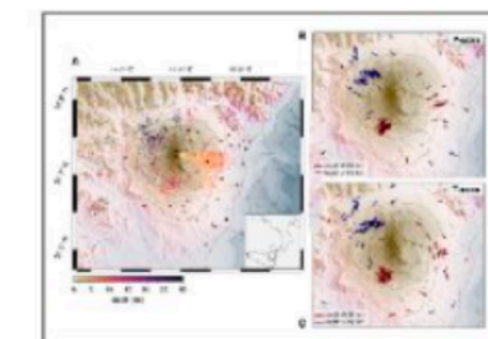
- Velocities
- Stresses
- Geometry evolution

Autodiff (AD)

∂ GPU4GEO – sensitivity and inverse modelling

Data

- In-situ field observations
- Remote sensing products
- Meteorological time series



Enzyme.jl

Adjoint PT solvers

$$\frac{\partial \psi}{\partial \tau} + \left[\frac{\partial R(u)}{\partial u} \right]^T \psi = \frac{\partial J}{\partial u}$$

Thank you
For your attention



<https://gpu4geo.org/>

<https://juliagpu.org/>



<https://github.com/JuliaGeodynamics>

<https://github.com/PTsolvers>

ludovic.raess@unil.ch



luraess