

# A GLIMPSE ON GPU IN ASTROPHYSICS

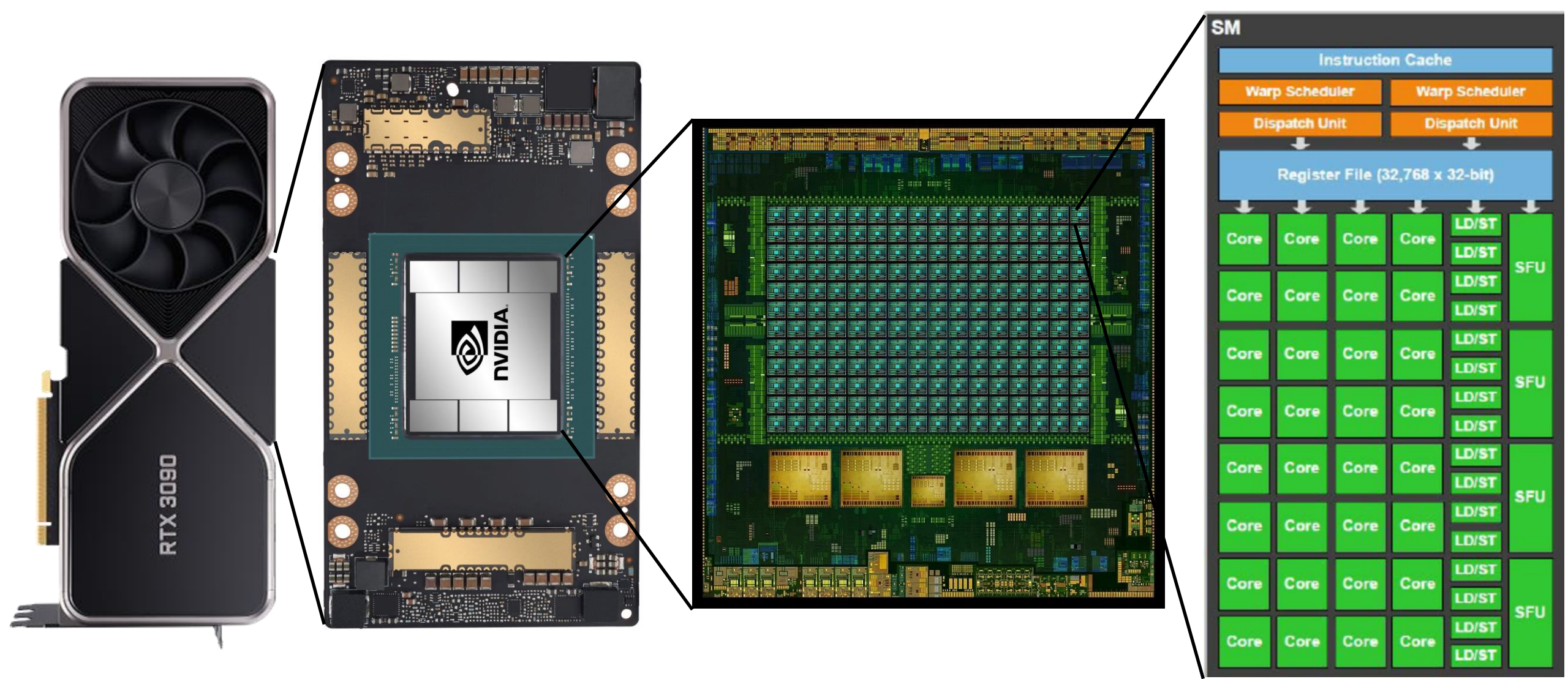
---

Lile Wang (王力乐)

KIAA/PKU

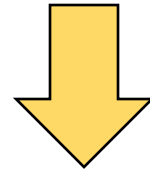






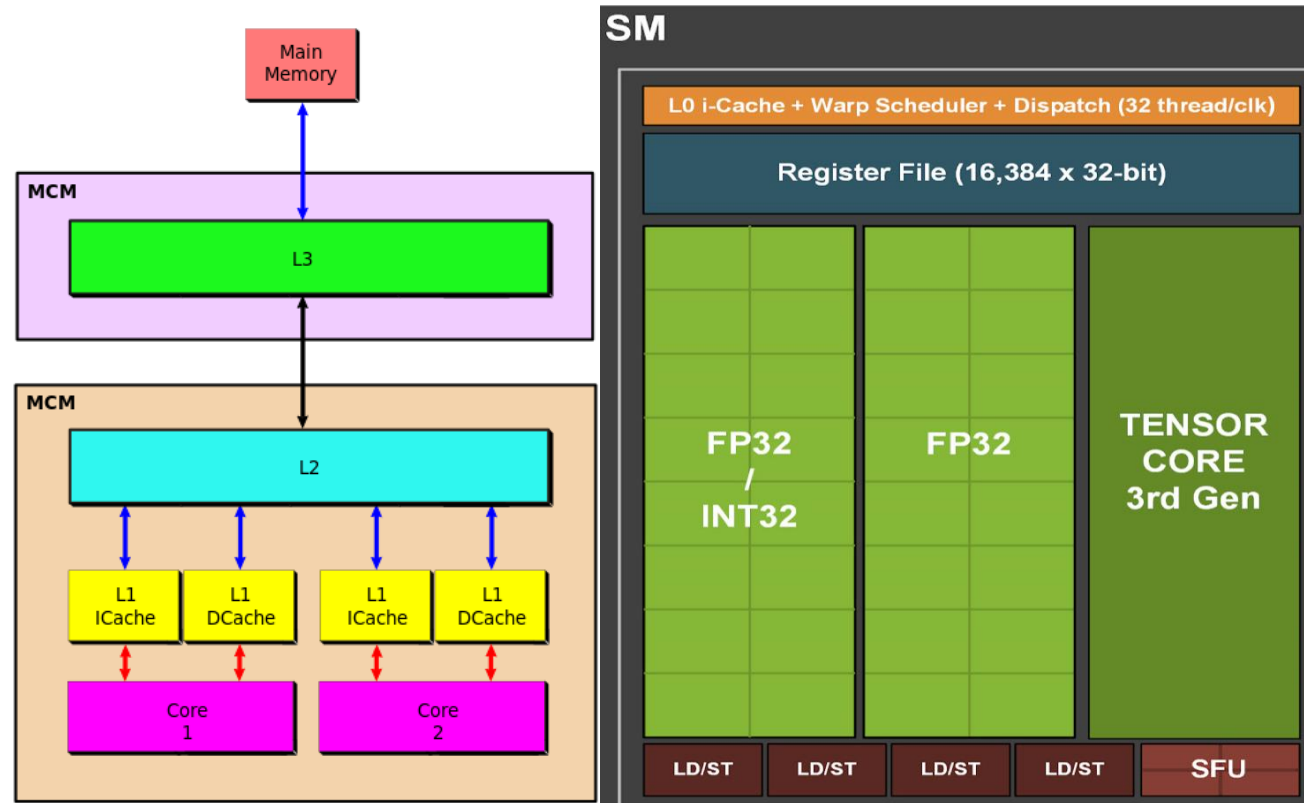
What is a Graphics Processing Unit?

- GPUs: Jet liners  
SIMD (before Turing)  
High latency mem access  
High throughput  
Lower (?) frequency  
Worse versatility



Slow, but Massive

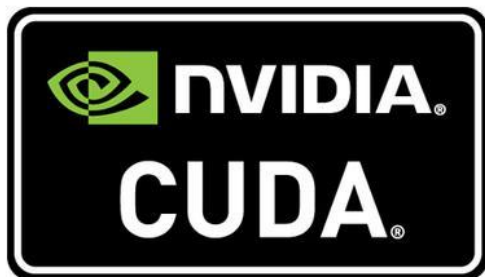
- Typical CPU (AMD Ryzen 7 5800X):
- 8 Cores / CPU
- Register: 64 bits / Unit
- L1 Cache: 64 KB / Core
- L2 Cache: 512 KB / Core
- L3 Cache: 32 MB overall



Typical GPU (NVIDIA RTX 3080):

- 82 SM / GPU
- 16 Units / SM
- Shared I-Cache (!!!)
- Register: 64 kB / SM
- L1 Cache: 128 KB / SM (64 KB / SM  $\leq$  Turing)
- L2 Cache: 18 MB overall

# Advantages and Disadvantages



Map  
(映射)

- Distribute computation tasks over all computation units properly

Cache  
(缓存)

- Minimize global graphics memory access, eliminate unnecessary ones



Stream  
(执行流)

- Coordinate computation and data transfer, avoid waiting

# Basic Guidelines



## CUDA

- Mature, Reliable (?), Powerful, Modernized (C++-20)
- Limited to NVIDIA



## HIP/ROCm

- Almost IDENTICAL to CUDA, e.g. cudaMemcpy vs hipMemcpy
- HIP-CPU: Use MT CPU to “emulate” GPUs (Useful!)



## OpenCL

- Deceased support from major manufacturers (and decreasing)
- cHIP-STAR: OpenCL + SPIR-V --> HIP



## OneAPI

- Premature: programming model, system quality...
- cHIP-STAR: OneAPI --> HIP



## MUSA

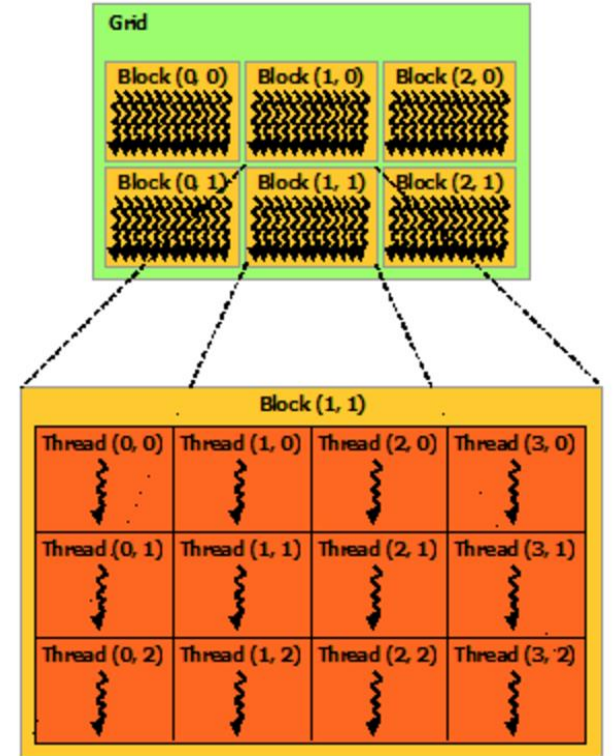
- Mostly identical to CUDA
- Premature support for computational purposes

# Major Programming Models

```
#include <iostream>
```

```
__global__ void kernel( float * a, float * b, float * tgt )  
{  
    const int idx = threadIdx.x + blockDim.x * blockIdx.x;  
    c[ idx ] = a[ idx ] + b[ idx ];  
    return;  
}
```

```
int main( )  
{  
    // Memory allocation, initialization  
    kernel <<< n_block, n_thread >>> ( a, b, c );  
    // Resource releasing and garbage collection  
    return 0;  
}
```



# Major Programming Models for GPUs



Instruction:

$res[?] = a[?] + b[?];$

Block 0  
w/ N Threads

...

Block m  
w/ N  
Threads

Thread 0

Thread 1

...

Thread (N-1)

Thread n

? -> 0

? -> 1

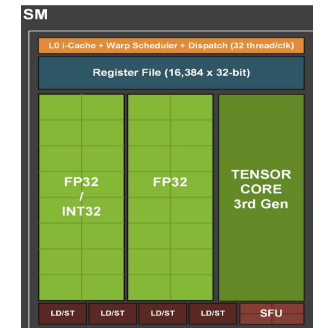
? -> (N-1)

? -> m \* N + n

```
template < typename f_T > __global__ void kernel  
( f_T * res, const f_T * a, const f_T * b )  
{  
    const int idx = blockIdx.x * blockDim.x + threadIdx.x;  
    res[ idx ] = a[ idx ] + b[ idx ];  
    return;  
}
```



...



# Naive Mapping: Uniform Task Dispatching

```

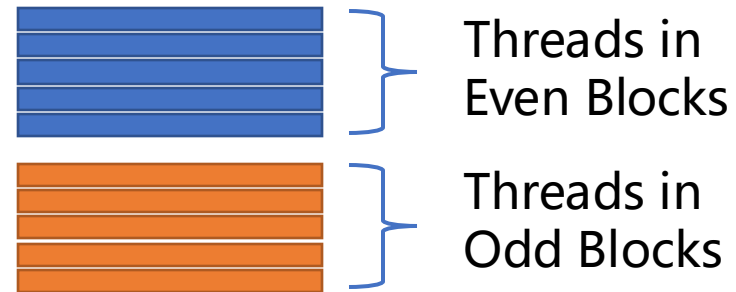
_global__ void kernel( ... )
{
    const int idx = blockIdx.x * blockDim.x + threadIdx.x;
    if( idx % 2 == 0 )
    {
        // Operations for even indices;
    }
    else
    {
        // Operations for odd indices;
    }
    return;
}

```

```

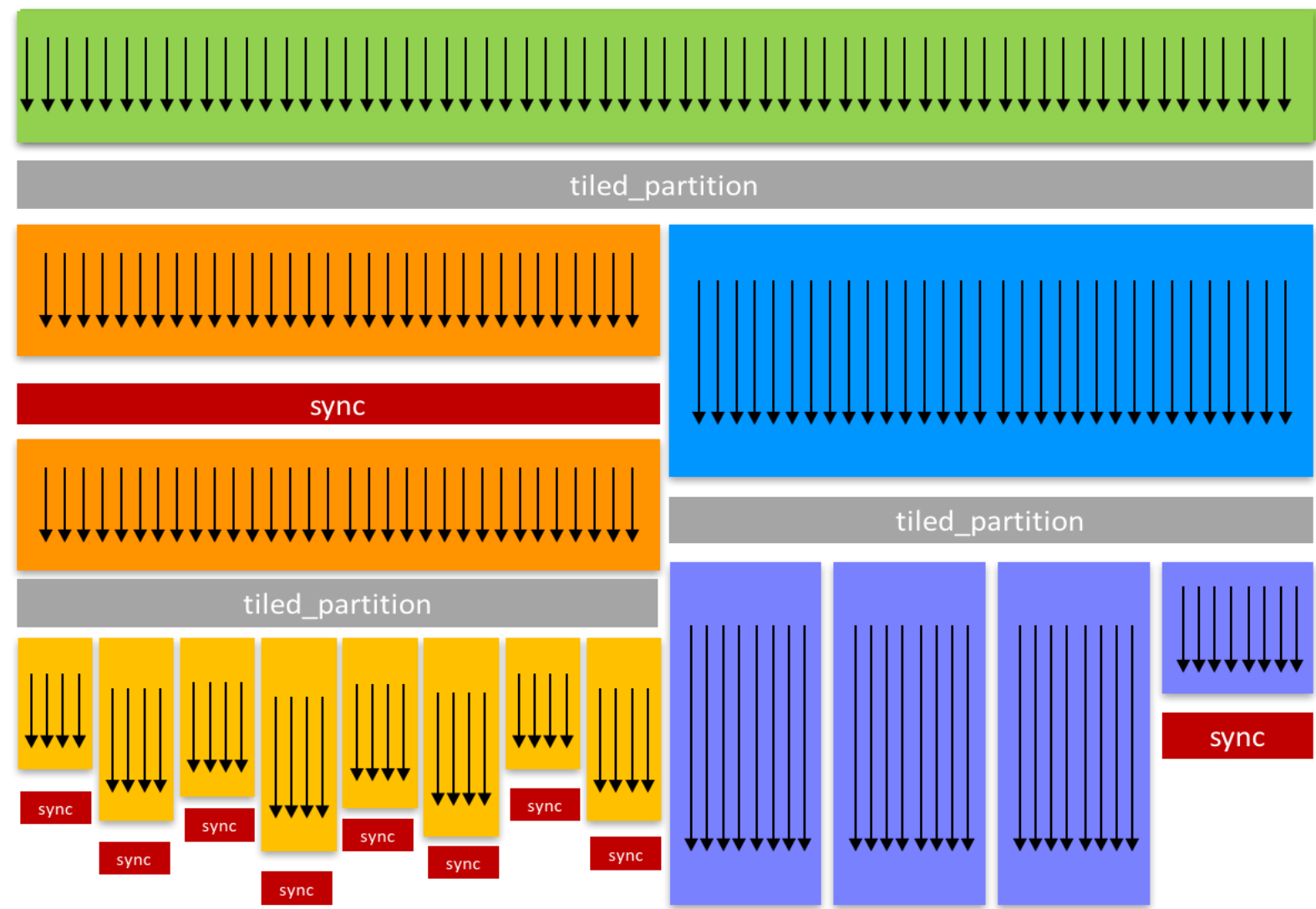
_global__ void kernel( ... )
{
    const int idx_0 = blockIdx.x * blockDim.x + threadIdx.x;
    if( blockIdx.x & 1 )
    {
        const int idx = 2 * idx_0;
        // Operations for even indices;
    }
    else
    {
        const int idx = 2 * idx_0 + 1;
        // Operations for odd indices;
    }
    return;
}

```

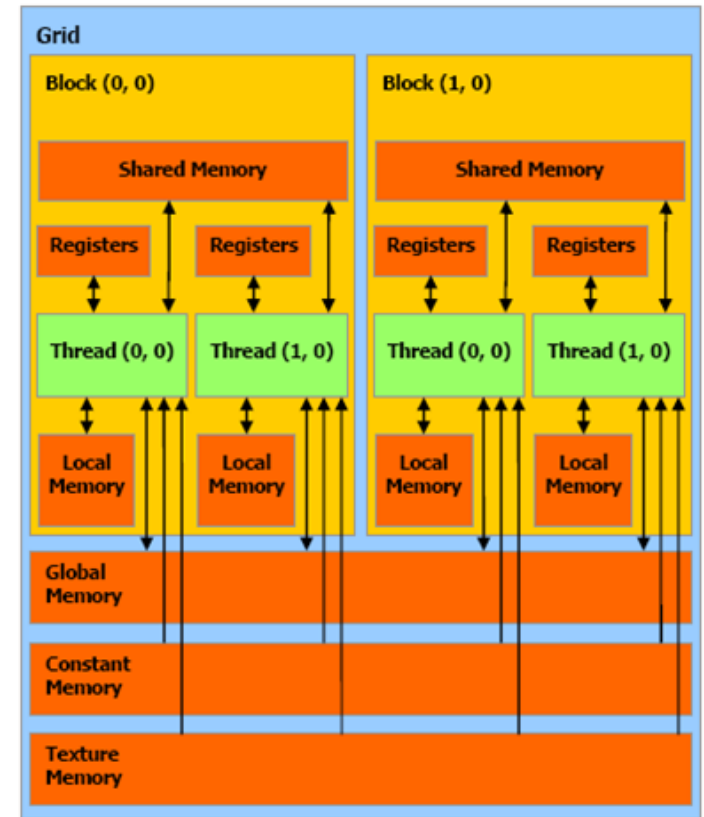
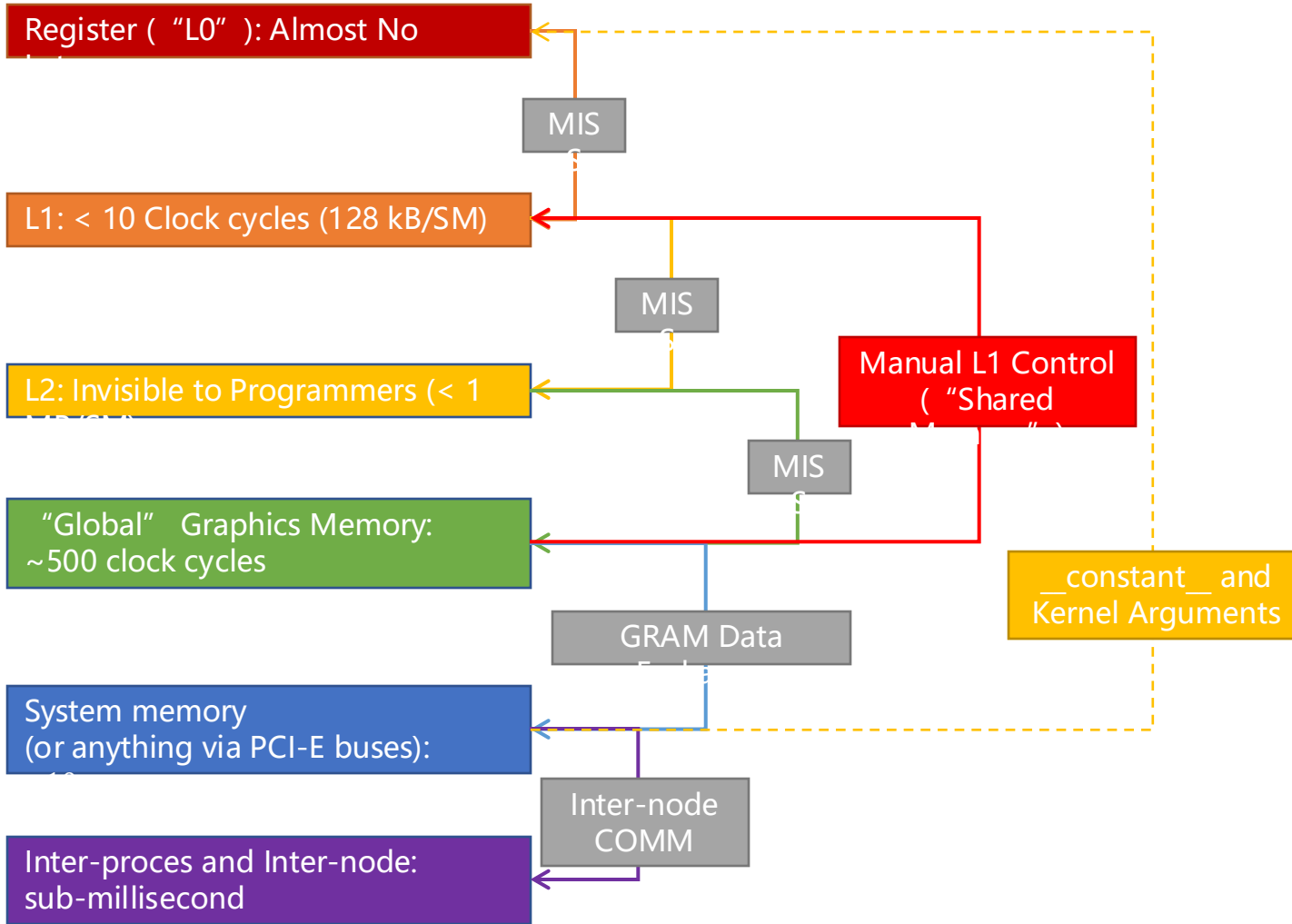


Less Naive Mapping: Dispatch More Wisely

- Finer or broader
- Single-thread granularity
- Why is it possible? (hint: Cache)
- CUDA ONLY (for now)



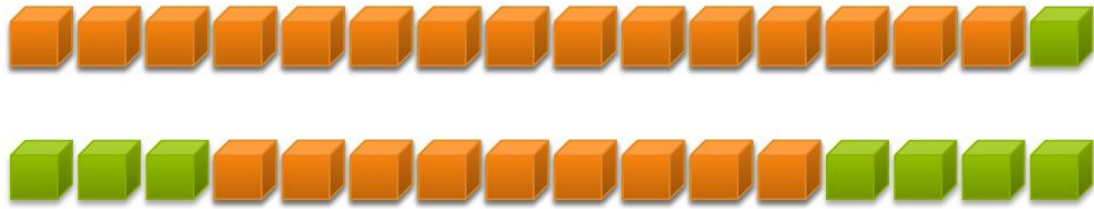
# Finer Thread Manipulation: Cooperative Group



# Cache: Ubiquitous Structures

## Example: Convolution, etc. in 1D

- Assuming a logical radius of 3:  
Each data point used 3 times



## Example: Minimum Timestep

- Binary reduction for the minimum

```
__global__ void kernel_shared_static( ... )
{
    __shared__ float var[ 64 ];
    // Operations...
    return;
}

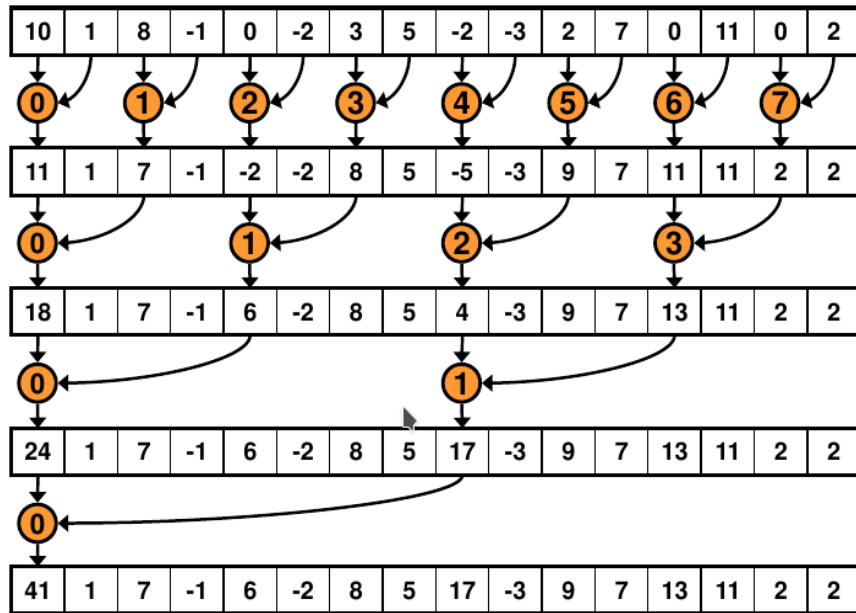
__global__ void kernel_shared_dynamic( ... )
{
    extern __shared__ float var[ ];
    // Operations...
    return;
}

int main( )
{
    // ...
    kernel_shared_static <<< n_bl, n_th, 0 >>> ( ... );
    kernel_shared_dynamic <<< n_bl, n_th, n_sh >>> ( ... );
    // ....
}
```

Shared Memory: L1 Cache with Manual Ctrl

# Example: Reduction of Arrays

- Reduction for the sum at  $O(\log N)$
- Avoid congesting global atomicAdd



```
__global__ void kernel_shared_static( ... )
{
    __shared__ float var[ 64 ];
    // Operations...
    return;
}

__global__ void kernel_shared_dynamic( ... )
{
    extern __shared__ float var[ ];
    // Operations...
    return;
}

int main( )
{
    // ...
    kernel_shared_static <<< n_bl, n_th, 0 >>> ( ... );
    kernel_shared_dynamic <<< n_bl, n_th, n_sh >>> ( ... );
    // ....
}
```

Shared Memory: L1 Cache with Manual Ctrl

## Communicate

- The best way that threads **communicate** with in a block

## Race Conditions

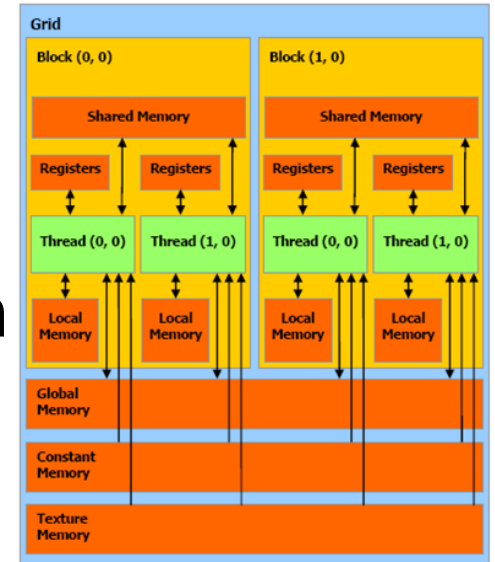
- Well-defined shared-memory model: Pay extra attention **race conditions!**

## Local Dynamic Memory

- Shared memory is the only way of reserving dynamic memory **locally**

## Reduce Usage

- **Reduce** shared memory usage when possible for more context on an SM

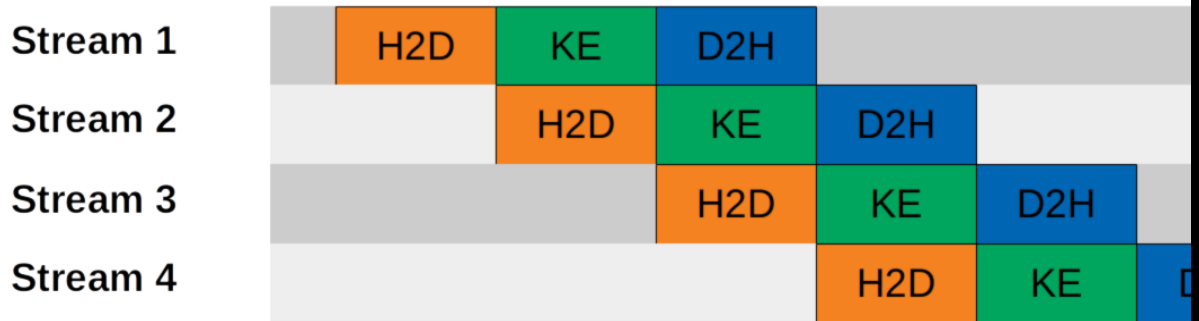


# More About Shared Memory

## Serial Model



## Concurrent Model



```
__global__ void kernel( ... )  
{  
    // Operations...  
    return;  
}  
  
int main( )  
{  
    // ...  
    cudaMemcpyAsync( p_device[ 0 ], p_host[ 0 ], n_bytes,  
                    cudaMemcpyDefault, stream_0 );  
    kernel <<< n_bl, n_th, n_sh, stream_0 >>> ( ... );  
  
    cudaMemcpyAsync( p_device[ 1 ], p_host[ 1 ], n_bytes,  
                    cudaMemcpyDefault, stream_1 );  
    kernel <<< n_bl, n_th, n_sh, stream_1 >>> ( ... );  
    // ....  
  
    cudaStreamSynchronize( stream_0 );  
    cudaStreamSynchronize( stream_1 );  
  
    return 0;  
}
```

Time

# Streams: Higher-Level Coordination



## Pinned Memory

- Unpinned memory may be paged out
- Use `cudaMallocHost` and `cudaFreeHost`

## Distributing and Dispatching Tasks Properly

- If memory exchange is necessary, hide it behind computation
- Divide the computation operations when necessary

## More Possibilities

- “Virtual memory” for the GPUs
- Inter-process or Inter-node communications



# More About Streams

- Spatial discretization: Control volume
- Temporal discretization: Explicit(!)

- In a “control volume” , the integral of conserved quantity

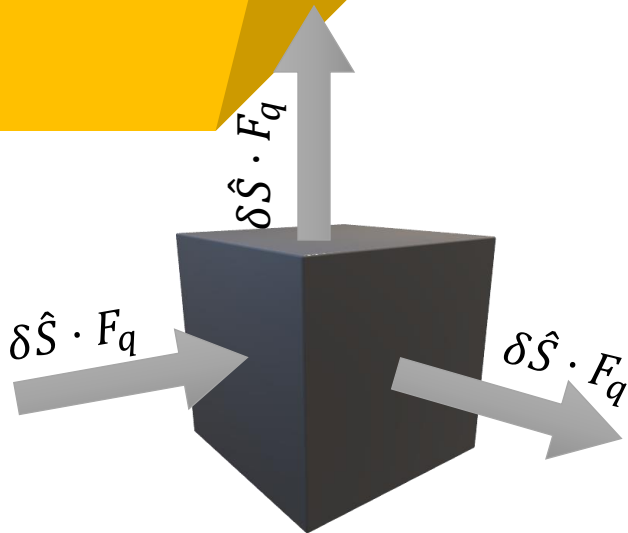
$$\Delta \int q \, dV = \sum_{\text{surface}} \int dt \int d\hat{S} \cdot F_q$$

- Finite volume discretization:

$$\Delta \int q \, dV \simeq \delta V \Delta q$$

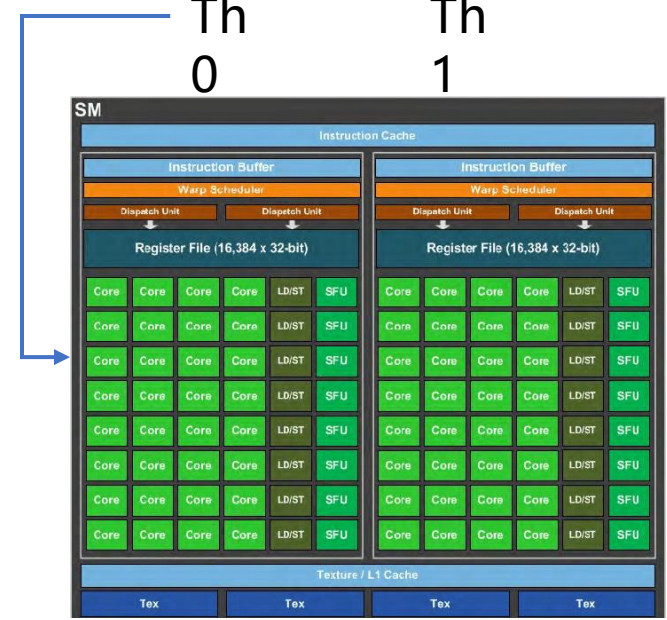
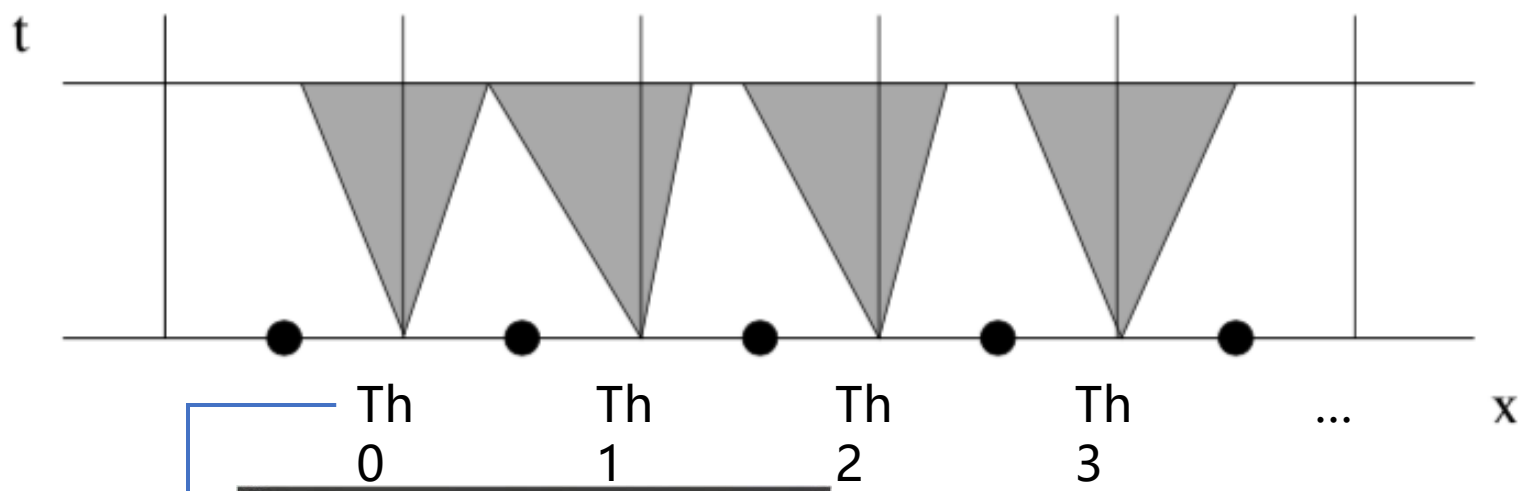
$$\int d\hat{S} \cdot F_q \simeq \delta \hat{S} \cdot F_q$$

$$\int dt \rightarrow \text{“Initial Value Problem”}$$



# Finite Volume Method

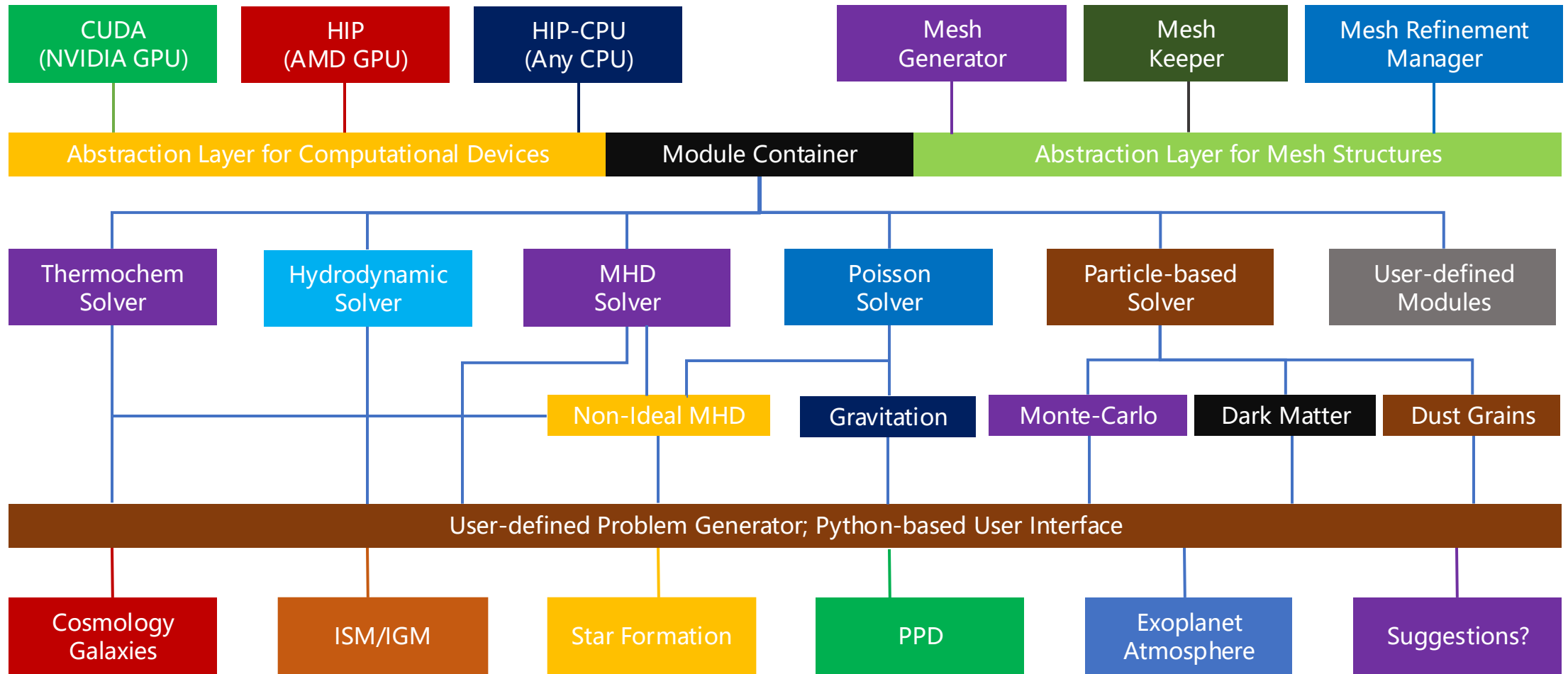
- Godunov method naturally suitable for parallelization
- GPU: The de-facto computational device for this
- However: No adequate implements



✈ CPU 集群  
 ✈ GPU 集群



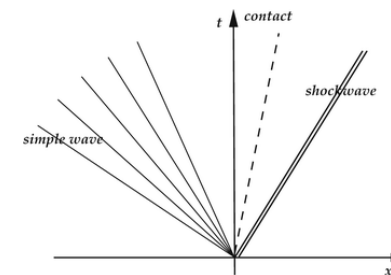
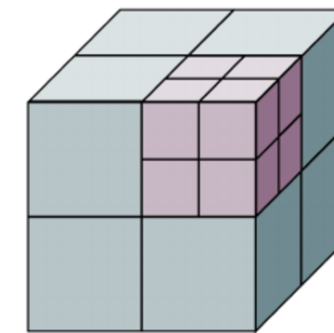
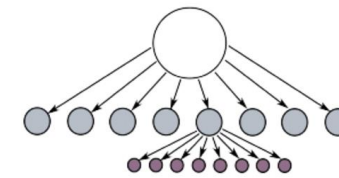
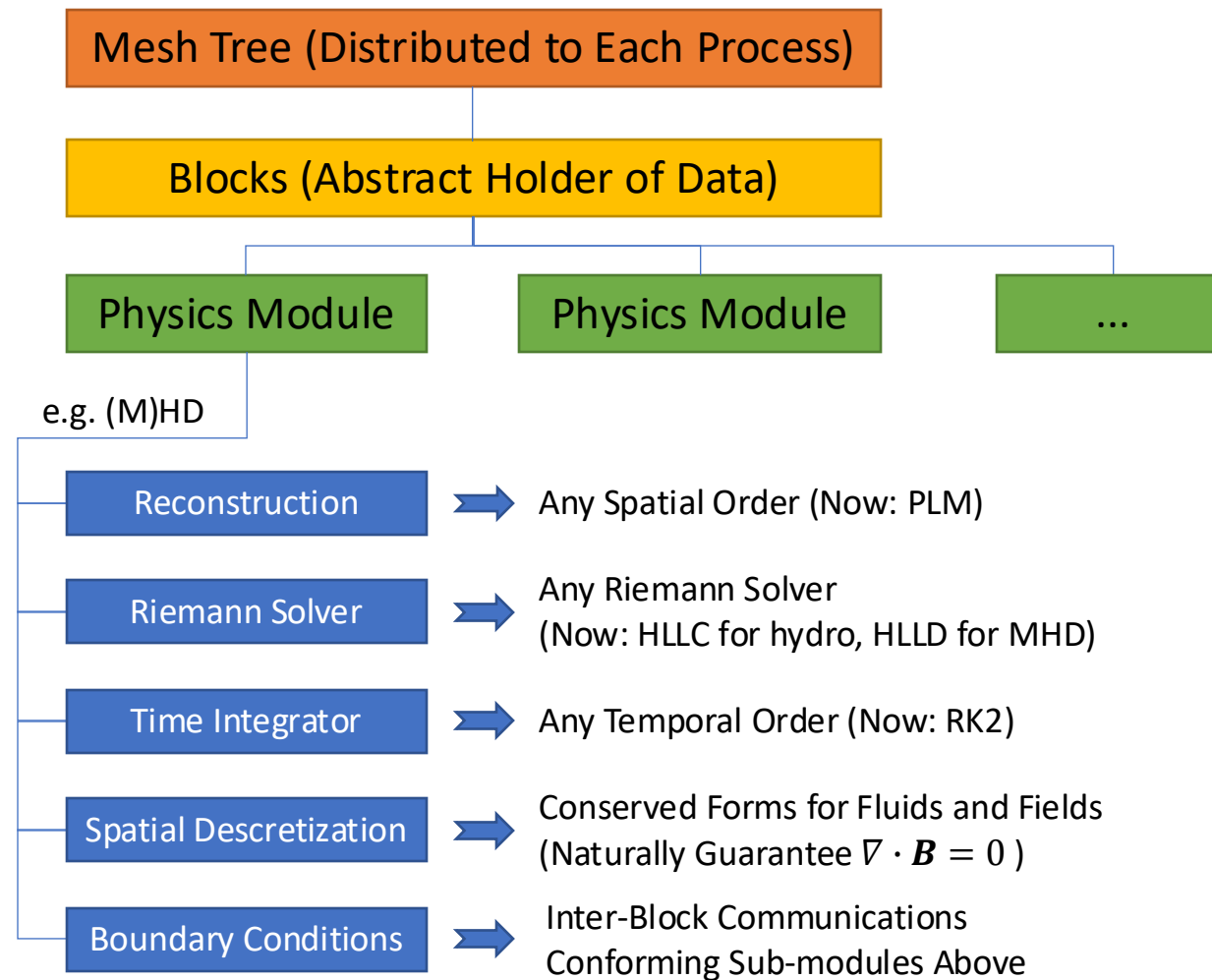
# Use GPUs for Godunov Method



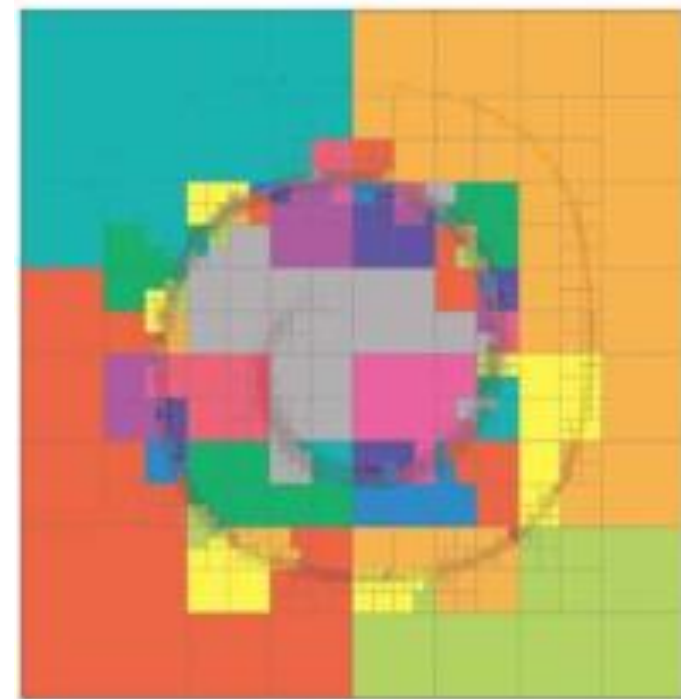
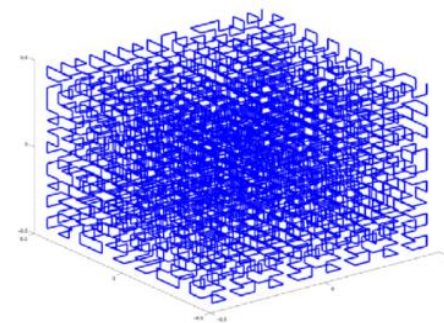
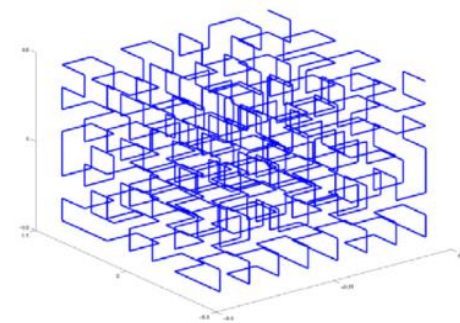
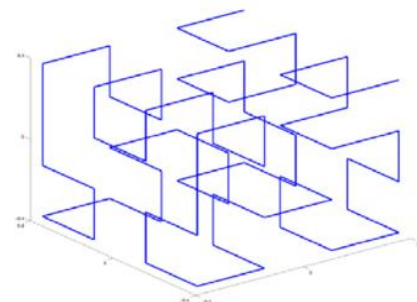
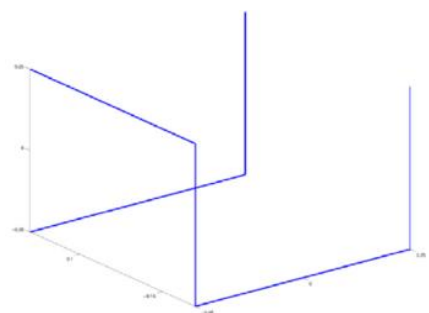
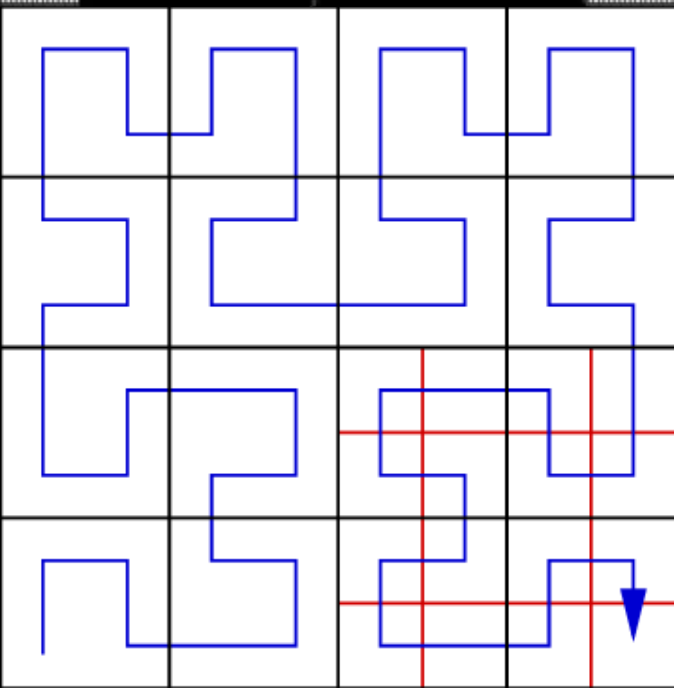
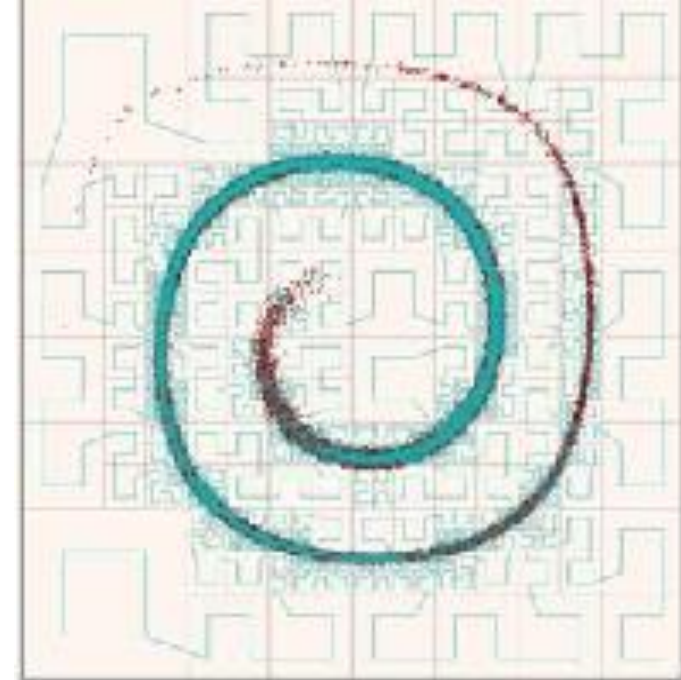
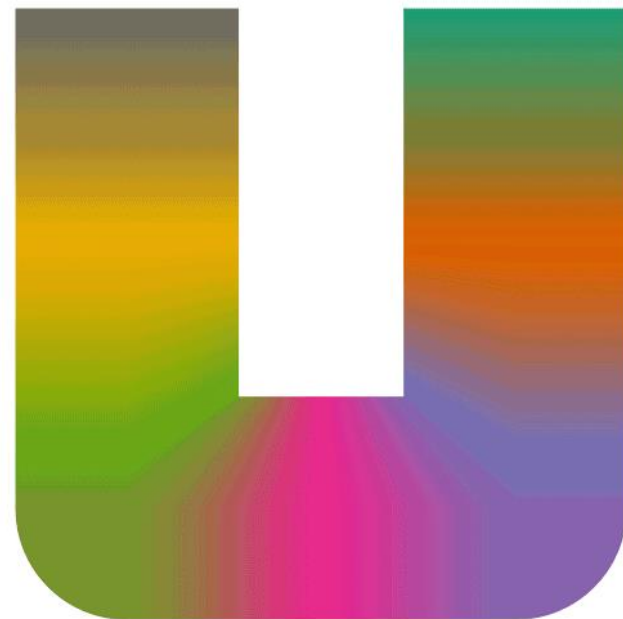
# New Heterogeneous Code: Structures



- **Baseline: Mesh Management**
  - Mesh Tree: **Unaware** of Data
  - Blocks: **Not Specific** to Modules
  - Modules: **Not Specific** to Algorithms
- Physics Modules "Tell" the Mesh Manager about What They Need
  - Easy to Emplace New Modules, Sub-modules, Algorithms
- Abstracted Interfaces
  - With Hardware (GPU, CPU, etc.)
  - With Communication Layer (MPI)
- Optimizations are Specific!
  - Can (and Should) be Accomplished on Sub-module Level
  - Can be Trimmed for Special Needs
- Computer Science Requirements: Type Calculus based on Templates



# Why KRATOS is Generic and Fast



- NVIDIA
- AMD
- RISC-V
- ARM(!)
  - Phone
  - Tablet

Device 0 [Radeon RX 7900 XTX] PCIe GEN 4@16x RX: N/A TX: N/A  
 GPU 2541MHz MEM 1249MHz TEMP 58°C FAN 18% POW 302 / 303 W  
 GPU[|||||||||||||||||||||||||||||||||100%] MEM[|||||] 3.054Gi/23.98

PID	USER	DEV	TYPE	GPU	GPU MEM	CPU	HOST MEM	Command
9640	lilew	0	Graphic	0%	0MiB	0%	5MiB	nvtop

F2 Setup F6 Sort F9 Kill F10 Quit F12 Save Config

```

cycle = 00050, t = 5.775195e-01, dt = 1.560404e-02, Speed = 4.019526e+08
cycle = 00060, t = 7.413505e-01, dt = 1.706026e-02, Speed = 4.015043e+08
cycle = 00070, t = 9.188548e-01, dt = 1.835290e-02, Speed = 4.017142e+08
cycle = 00080, t = 1.108971e+00, dt = 1.956057e-02, Speed = 4.020546e+08
cycle = 00090, t = 1.289782e+00, dt = 1.702364e-02, Speed = 4.021649e+08
cycle = 00100, t = 1.451583e+00, dt = 1.562206e-02, Speed = 4.016577e+08
cycle = 00110, t = 1.602745e+00, dt = 1.467418e-02, Speed = 4.023116e+08
cycle = 00120, t = 1.746775e+00, dt = 1.417207e-02, Speed = 4.013577e+08
cycle = 00130, t = 1.887855e+00, dt = 1.407335e-02, Speed = 4.017169e+08
cycle = 00140, t = 2.029158e+00, dt = 1.420160e-02, Speed = 4.012160e+08
cycle = 00150, t = 2.172597e+00, dt = 1.449955e-02, Speed = 4.009746e+08
cycle = 00160, t = 2.319522e+00, dt = 1.492584e-02, Speed = 4.014976e+08
cycle = 00170, t = 2.470988e+00, dt = 1.539055e-02, Speed = 4.012055e+08
cycle = 00180, t = 2.627619e+00, dt = 1.593827e-02, Speed = 4.013733e+08
cycle = 00190, t = 2.789773e+00, dt = 1.650582e-02, Speed = 4.015471e+08
cycle = 00200, t = 2.958111e+00, dt = 1.713032e-02, Speed = 4.015099e+08
cycle = 00210, t = 3.133088e+00, dt = 1.779585e-02, Speed = 4.015091e+08
cycle = 00220, t = 3.315028e+00, dt = 1.849798e-02, Speed = 4.008501e+08
cycle = 00230, t = 3.504164e+00, dt = 1.922949e-02, Speed = 4.010758e+08
cycle = 00240, t = 3.700743e+00, dt = 1.999533e-02, Speed = 4.008514e+08
  
```

1:w1.p4 1:zsh\* "desktop" 20:57 07-Dec-23

Device 0 [NVIDIA GeForce RTX 3080 Ti] PCIe GEN 1@ 8x RX: 0.000 KiB/s TX  
 GPU 210MHz MEM 405MHz TEMP 42°C FAN 0% POW 17 / 350 W  
 GPU[|||||] MEM[|||] 0%] MEM[|||] 1.06

Device 1 [NVIDIA GeForce RTX 3090] PCIe GEN 4@ 8x RX: 1.638 GiB/s TX  
 GPU 1770MHz MEM 9501MHz TEMP 56°C FAN 30% POW 356 / 350 W  
 GPU[|||||||||||||||||||||||||||||||||100%] MEM[|||||] 3.65

PID	USER	DEV	TYPE	GPU	GPU MEM	CPU	HOST MEM	Command
2751532	lilew	1	Compute	98%	3408MiB	14%	103%	./kratos
954	root	0	Graphic	0%	411MiB	3%	0%	/usr/lib
2738497	lilew	0	Graphic	0%	86MiB	1%	0%	/usr/lib
2497029	lilew	0	Graphic	0%	64MiB	1%	0%	/usr/bin
2688062	lilew	0	Graphic	0%	56MiB	0%	0%	/proc/se

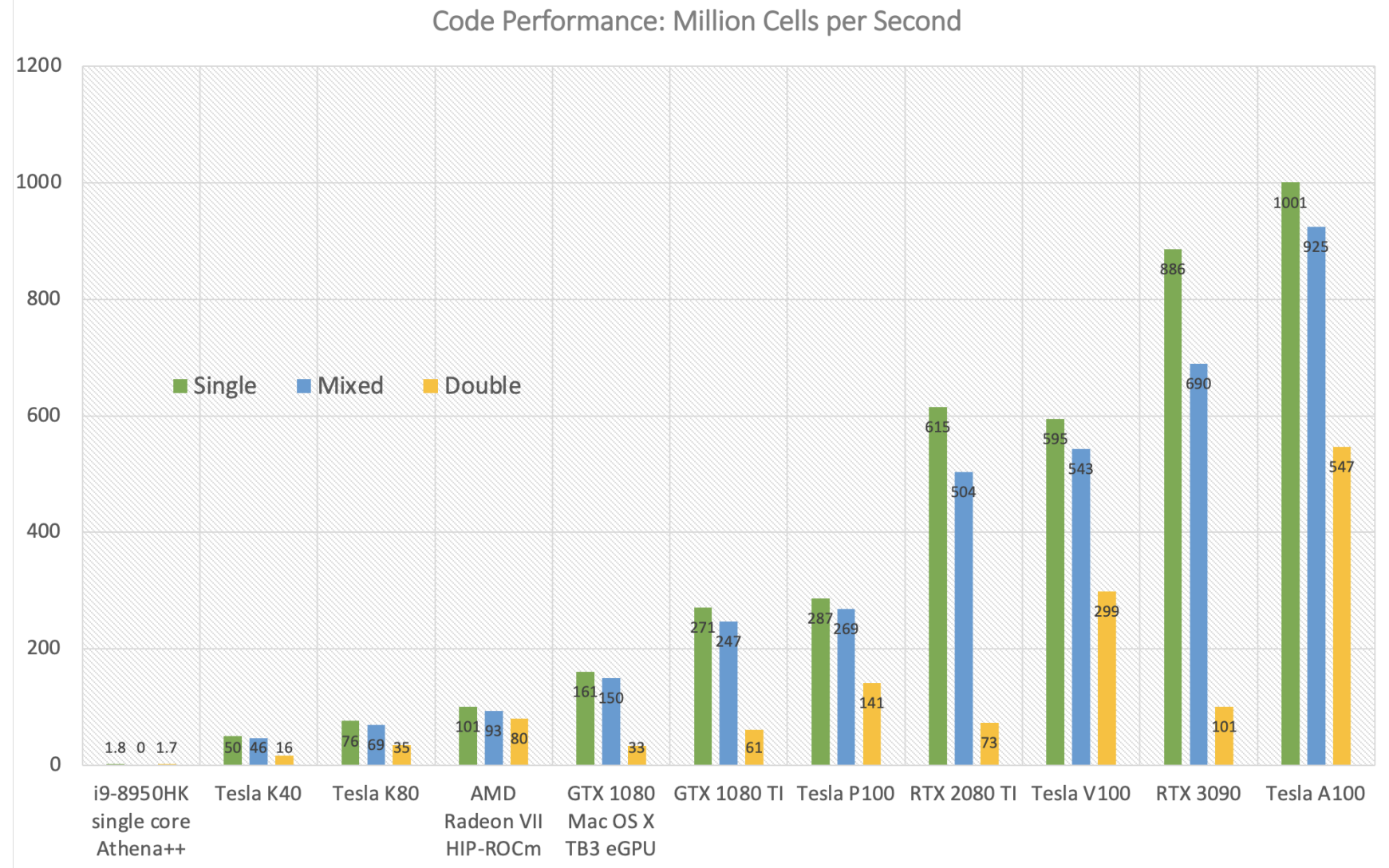
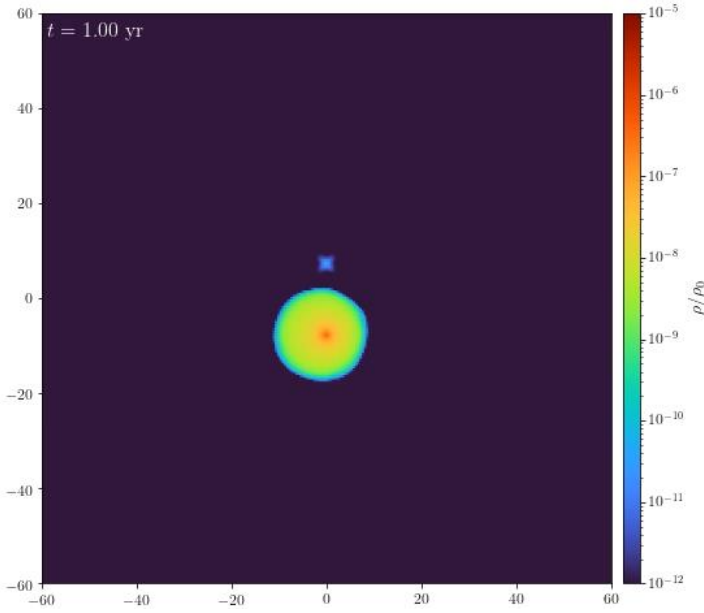
F2 Setup F6 Sort F9 Kill F10 Quit F12 Save Config

```

cycle = 00010, t = 8.758848e-02, dt = 8.441104e-03, Speed = 3.874580e+08
cycle = 00020, t = 1.821420e-01, dt = 1.046634e-02, Speed = 3.893777e+08
cycle = 00030, t = 2.962584e-01, dt = 1.226030e-02, Speed = 3.895929e+08
cycle = 00040, t = 4.286711e-01, dt = 1.401186e-02, Speed = 3.887565e+08
cycle = 00050, t = 5.775201e-01, dt = 1.560409e-02, Speed = 3.898962e+08
cycle = 00060, t = 7.413522e-01, dt = 1.706045e-02, Speed = 3.862590e+08
cycle = 00070, t = 9.188587e-01, dt = 1.835320e-02, Speed = 3.897396e+08
cycle = 00080, t = 1.108978e+00, dt = 1.956005e-02, Speed = 3.866890e+08
cycle = 00090, t = 1.289779e+00, dt = 1.702265e-02, Speed = 3.868385e+08
cycle = 00100, t = 1.451571e+00, dt = 1.562119e-02, Speed = 3.903923e+08
cycle = 00110, t = 1.602727e+00, dt = 1.467377e-02, Speed = 3.861778e+08
cycle = 00120, t = 1.746754e+00, dt = 1.417186e-02, Speed = 3.871235e+08
cycle = 00130, t = 1.887832e+00, dt = 1.407331e-02, Speed = 3.890181e+08
cycle = 00140, t = 2.029135e+00, dt = 1.420174e-02, Speed = 3.859638e+08
cycle = 00150, t = 2.172576e+00, dt = 1.450043e-02, Speed = 3.846060e+08
cycle = 00160, t = 2.319506e+00, dt = 1.492592e-02, Speed = 3.883941e+08
cycle = 00170, t = 2.470984e+00, dt = 1.539114e-02, Speed = 3.863917e+08
cycle = 00180, t = 2.627631e+00, dt = 1.593971e-02, Speed = 3.852376e+08
cycle = 00190, t = 2.789810e+00, dt = 1.650766e-02, Speed = 3.869391e+08
cycle = 00200, t = 2.958181e+00, dt = 1.713276e-02, Speed = 3.877465e+08
  
```

1:w1.p4 1:zsh\* "desktop" 20:57 07-Dec-23

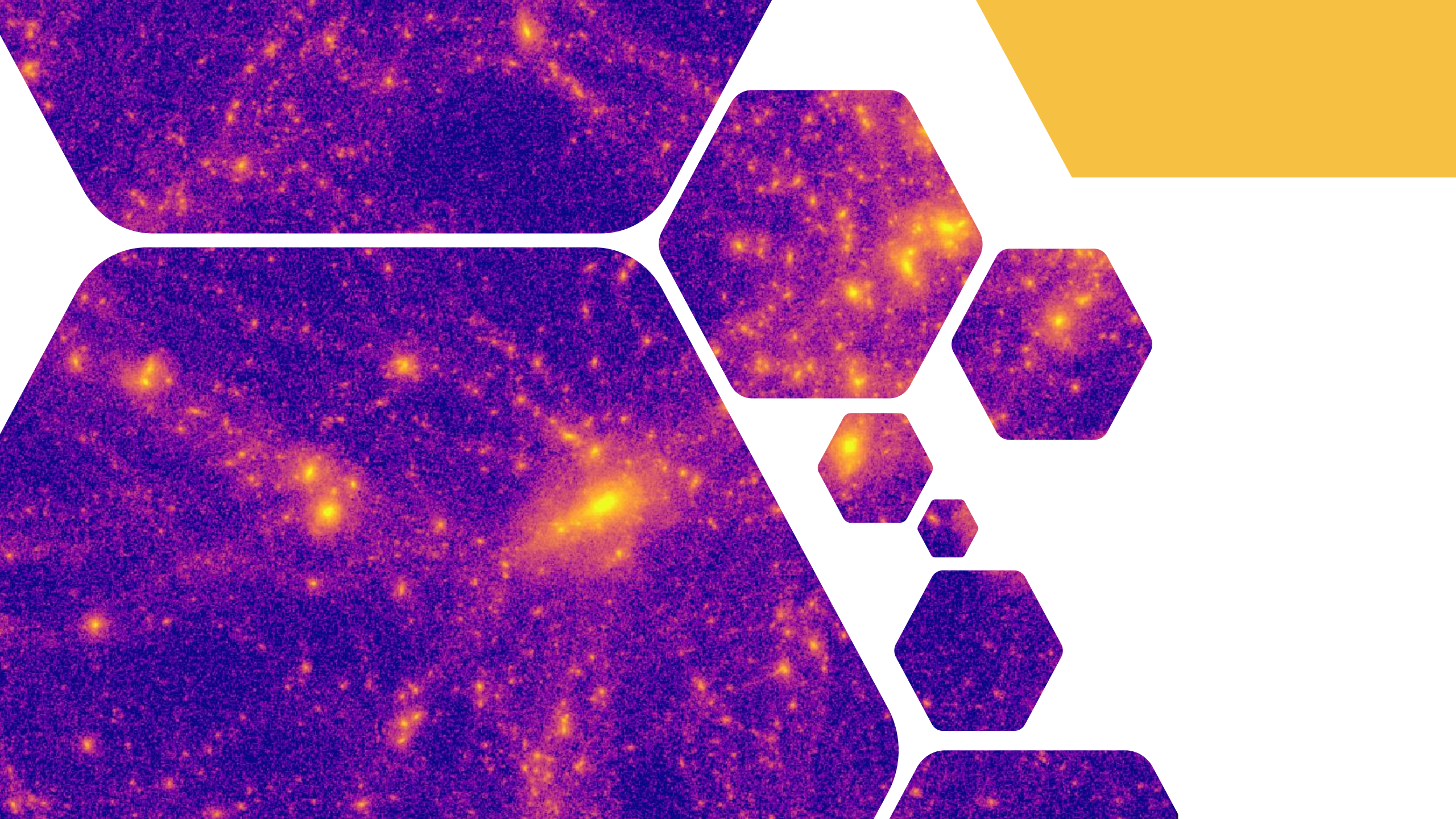
- Test case: Binary Outflows
- Single-card:  $10^9$  cells/s on A100
- Strong Scaling: 99.5% on 4xA100

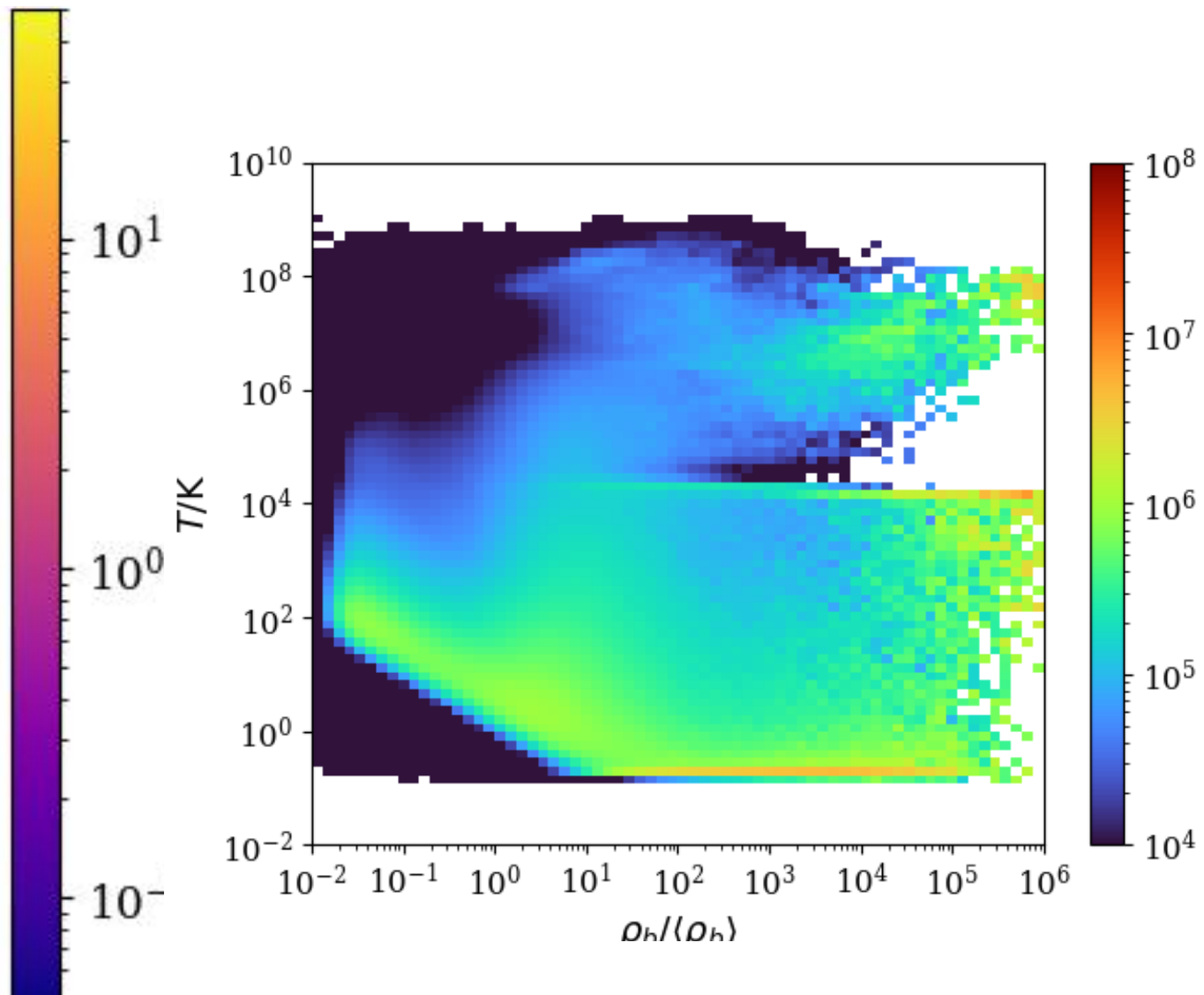
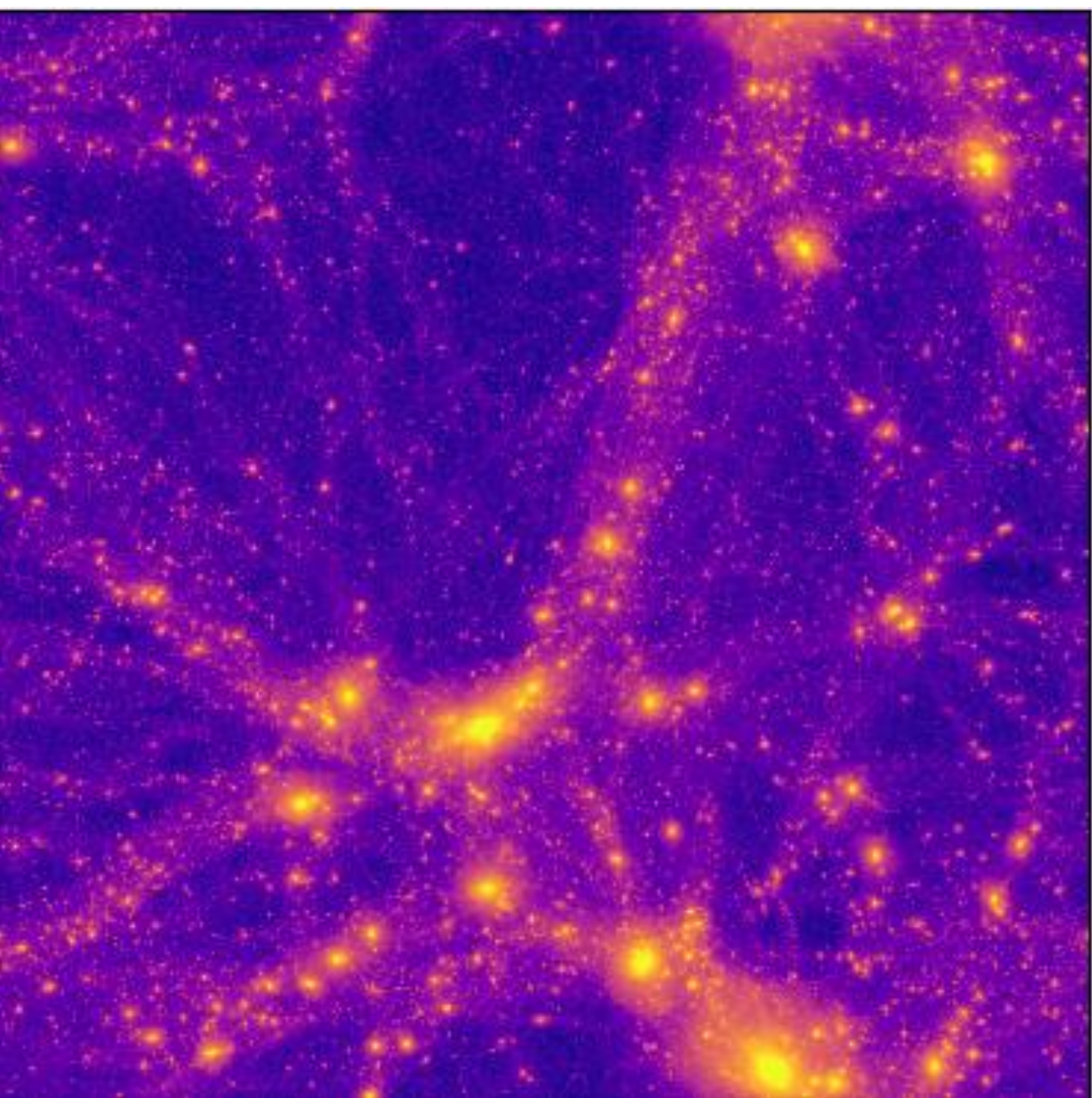


# Usages and Tests: Hydrodynamics & MHD



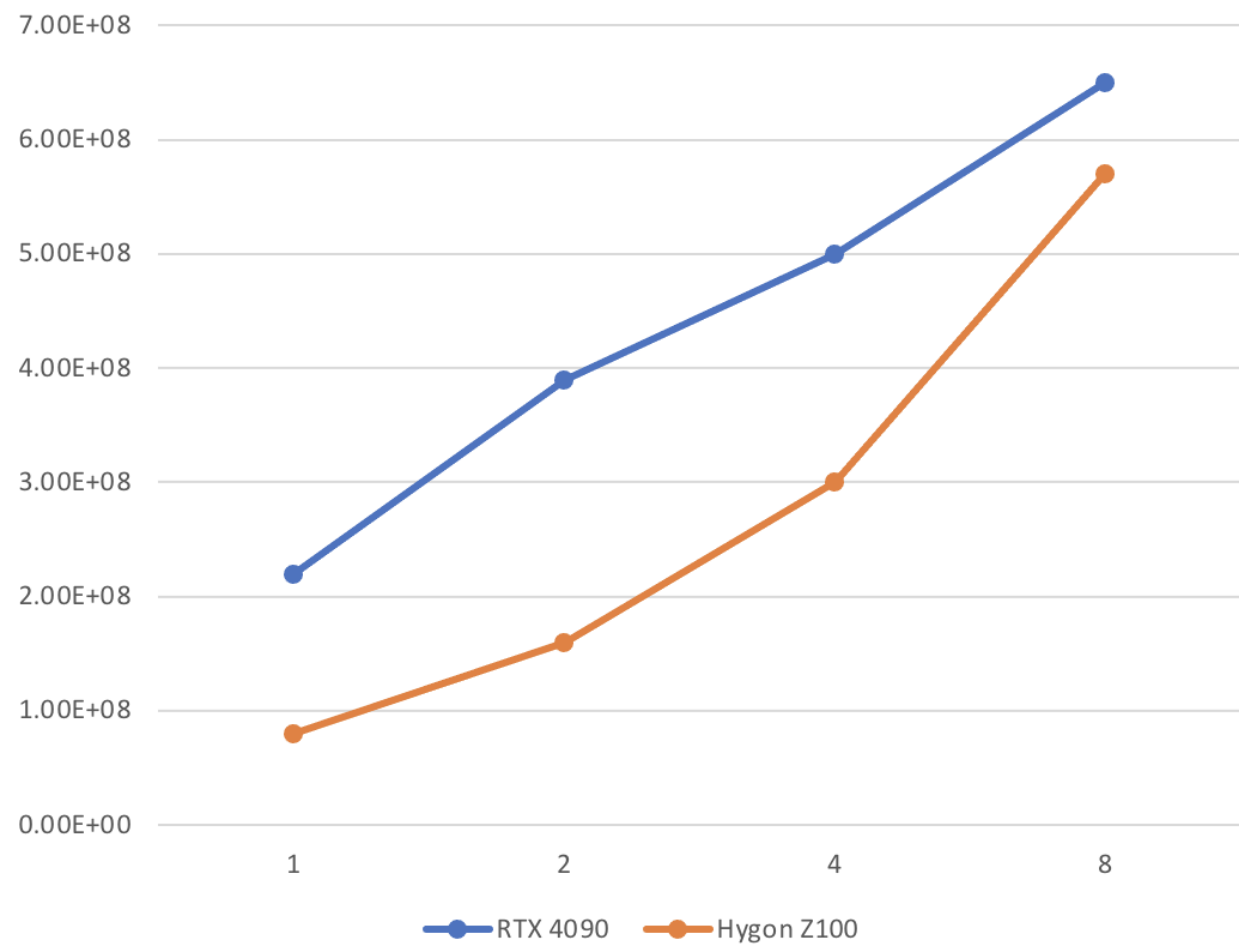




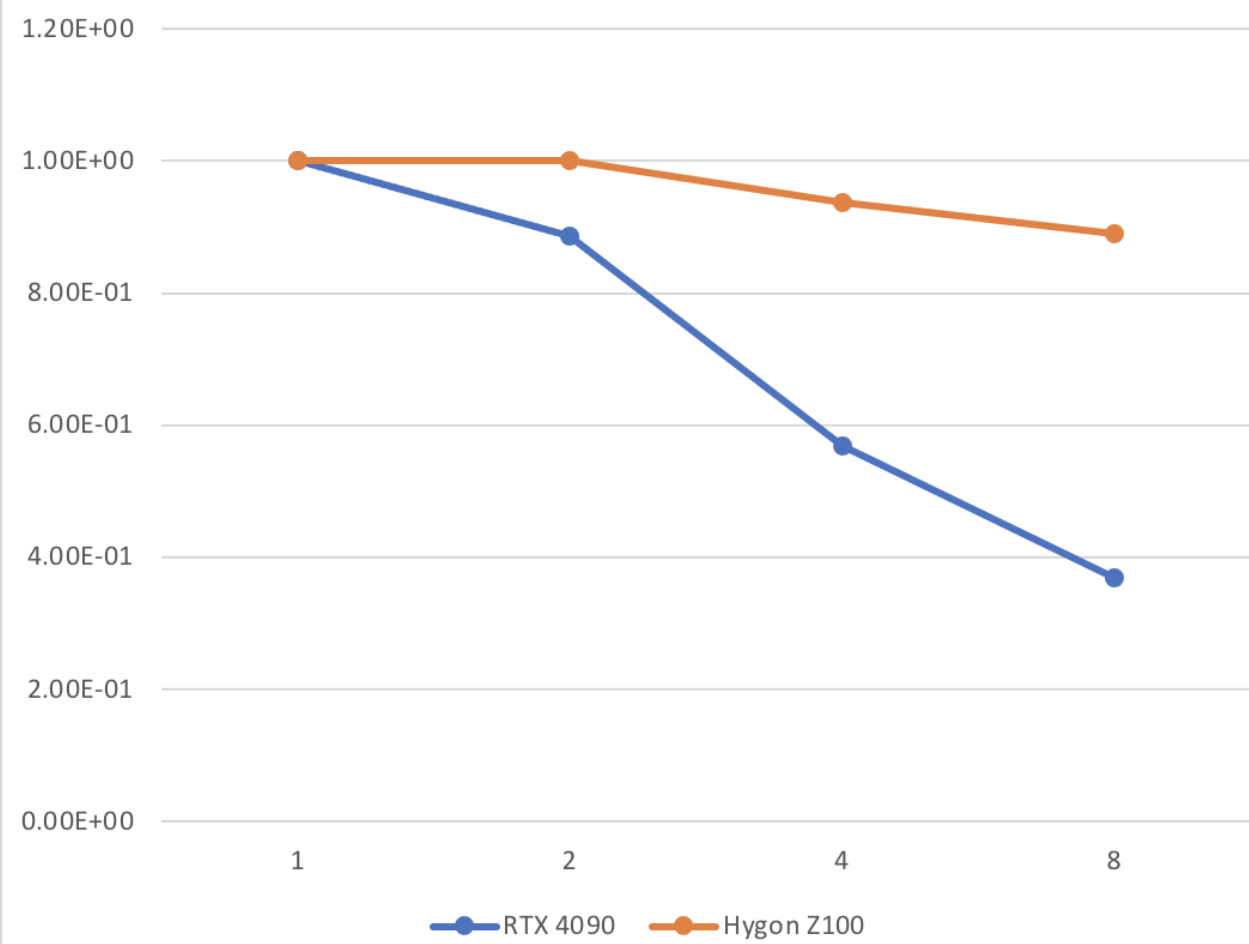


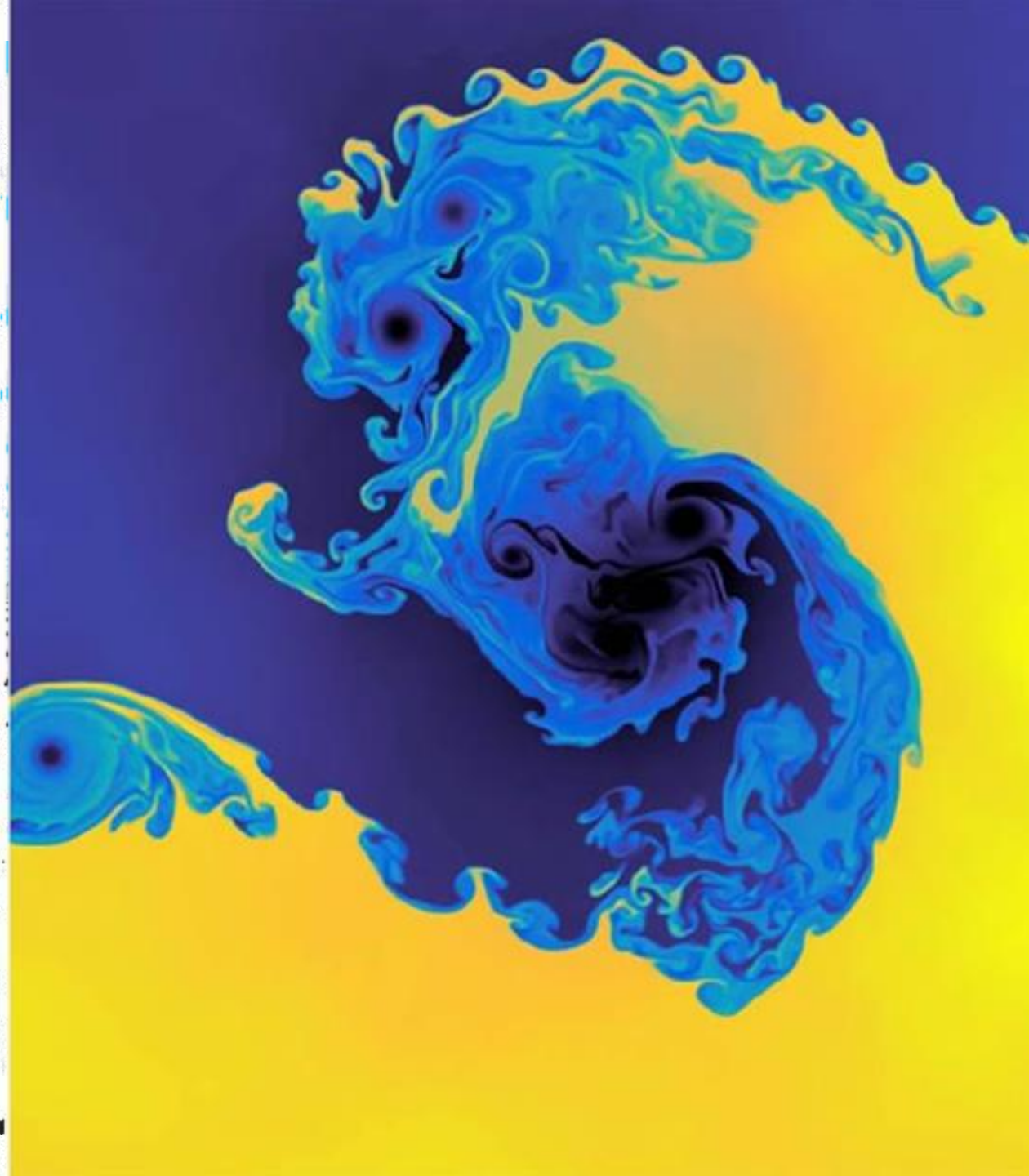
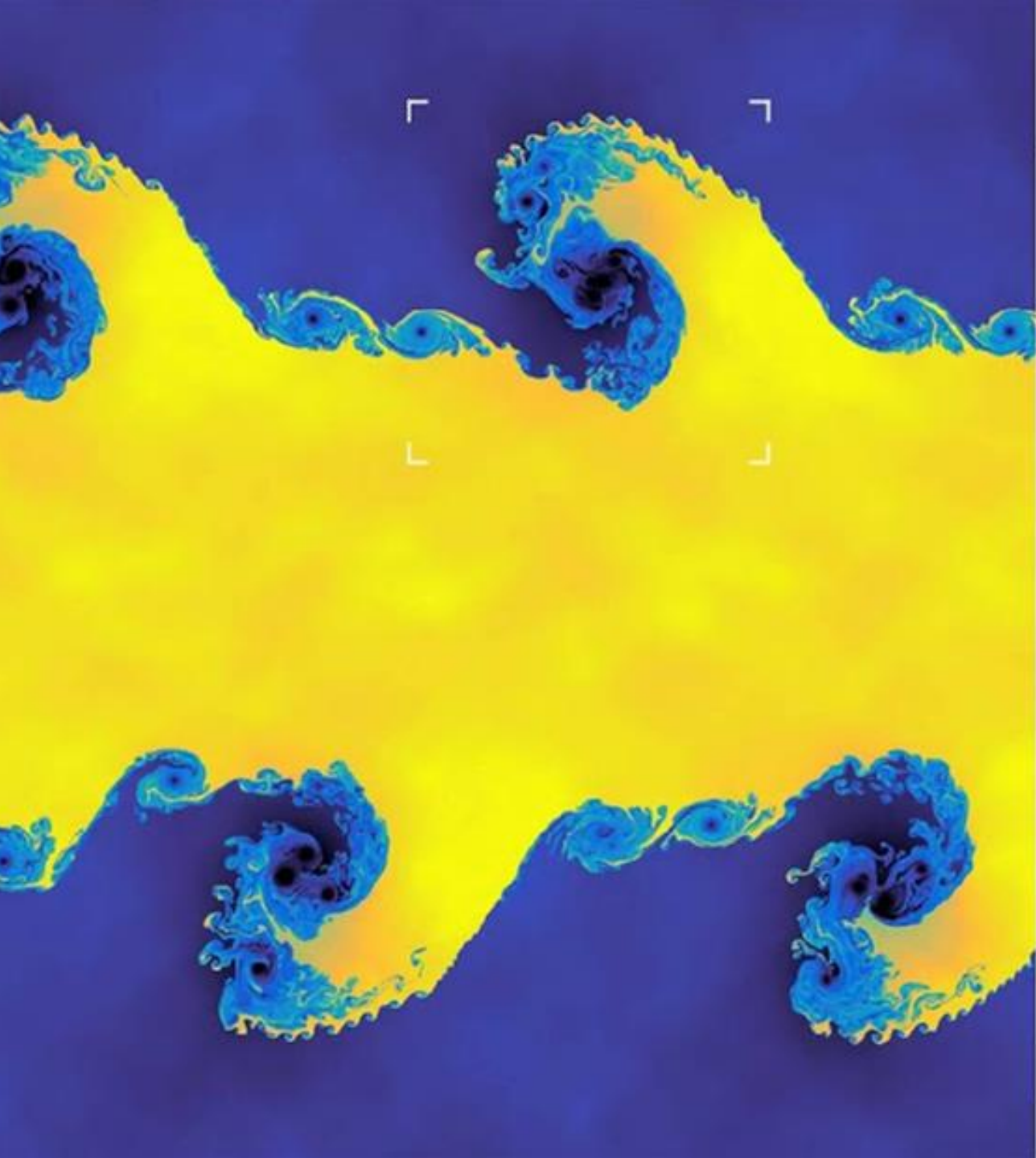
- MUSIC generated initial conditions  
<https://www-n.oca.eu/ohahn/MUSIC/>
- Particle module for DM  
(1 or 2 order LF)
- VL2 mixed-precision hydro with dual-energy formalism  
(Why VL2?)
- Multi-grid Poisson solver
- How fast is fast?

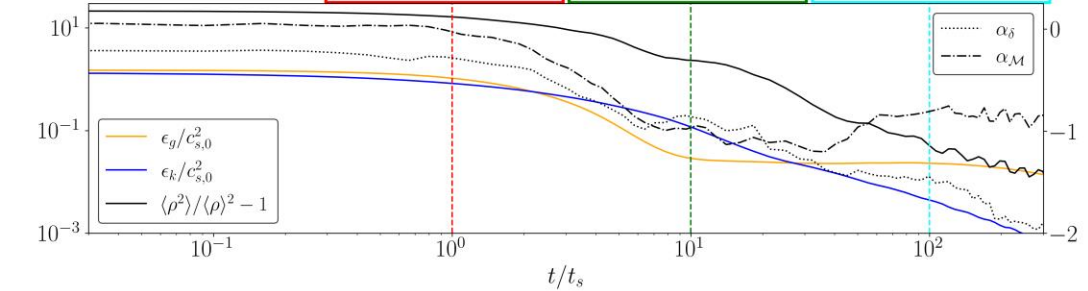
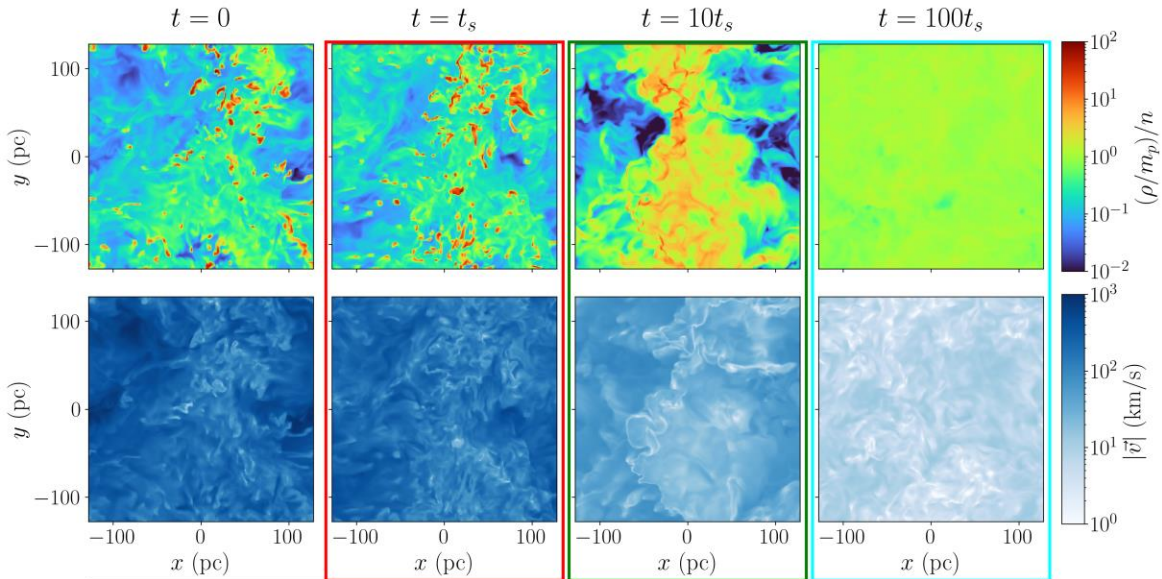
Cosmlogy: Speed



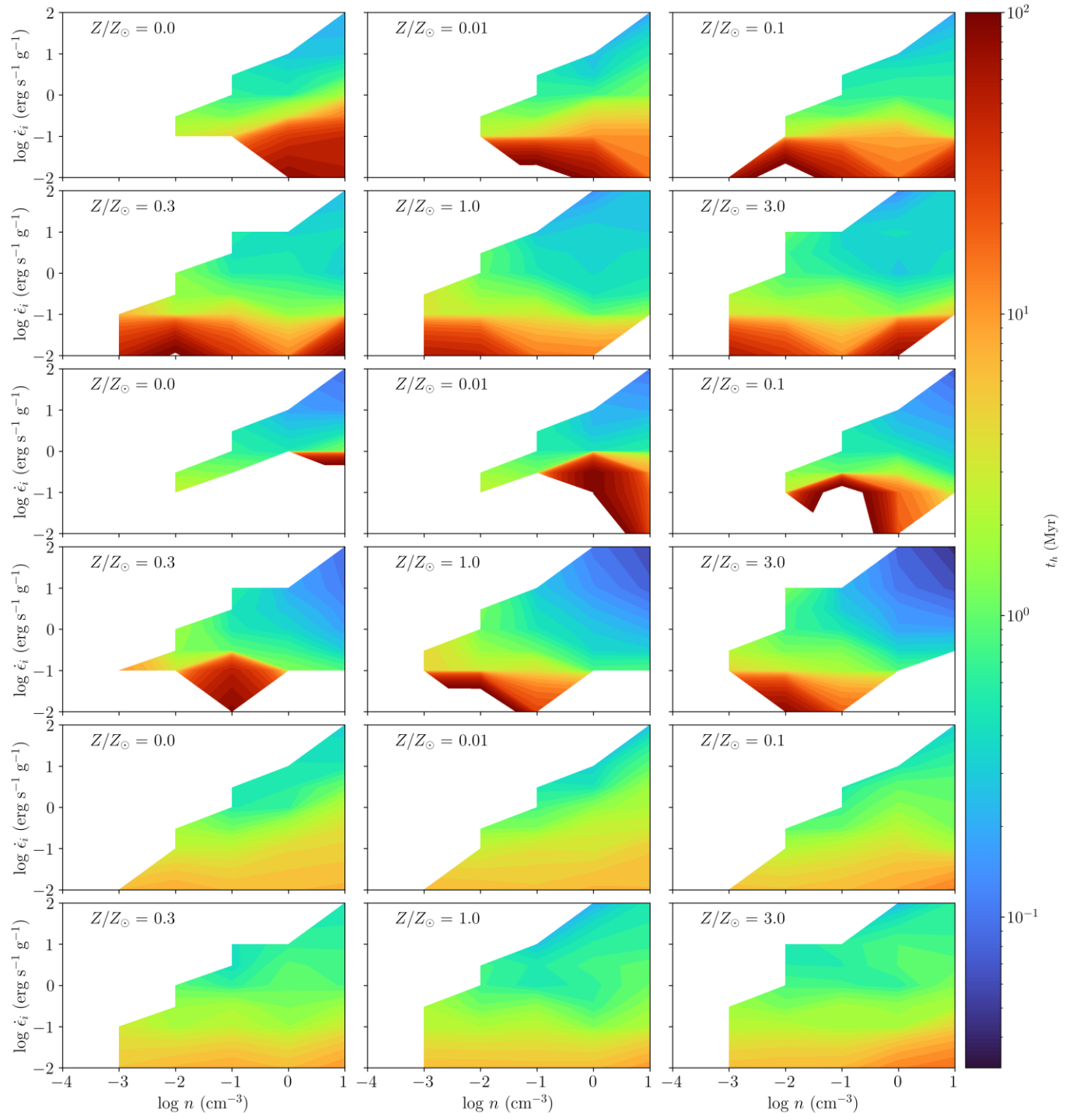
Cosmology: Weak Scaling



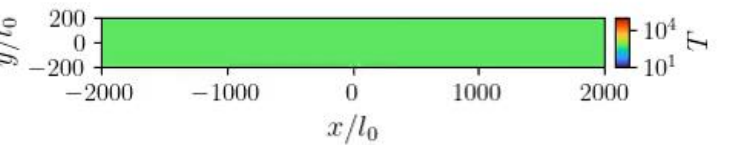
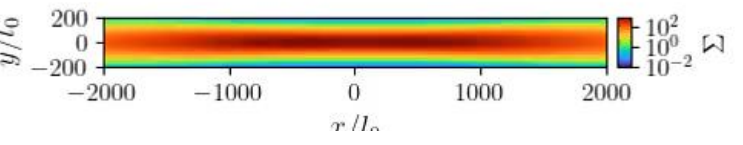
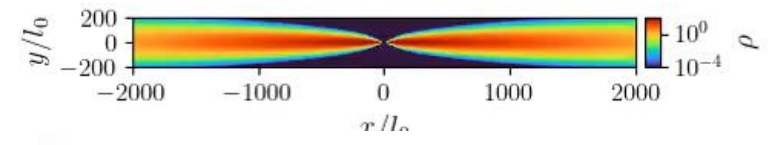
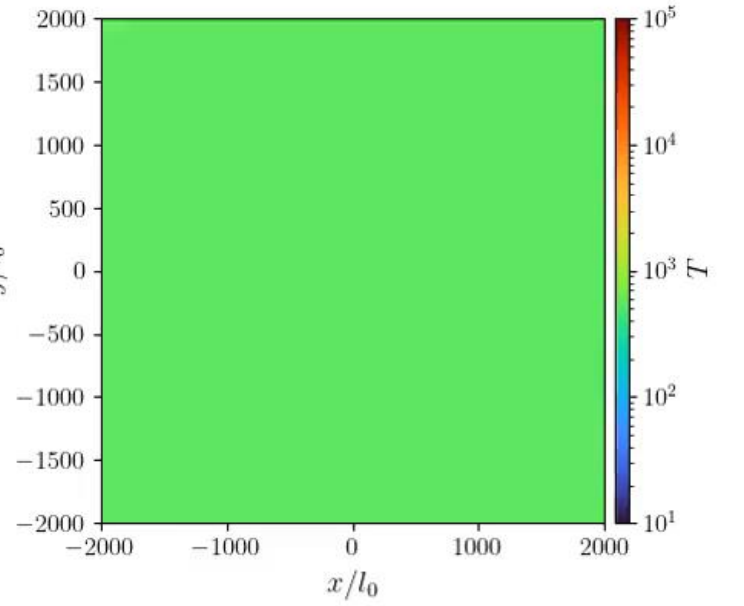
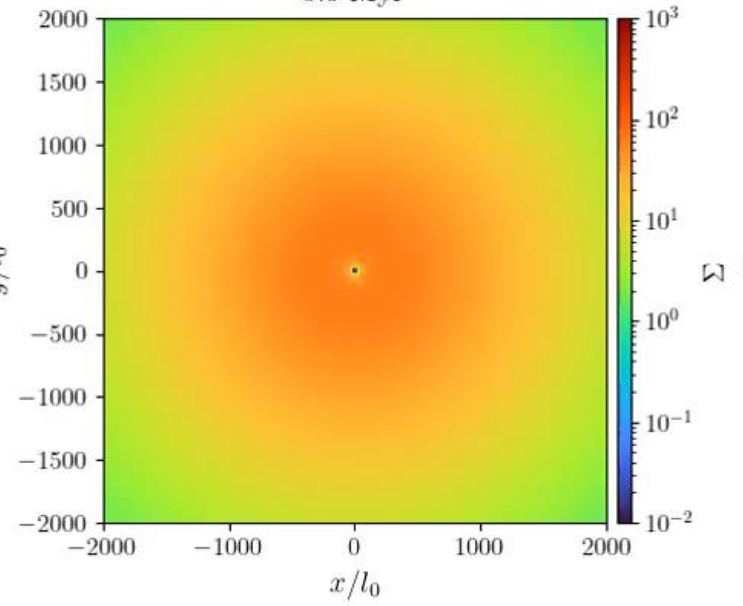
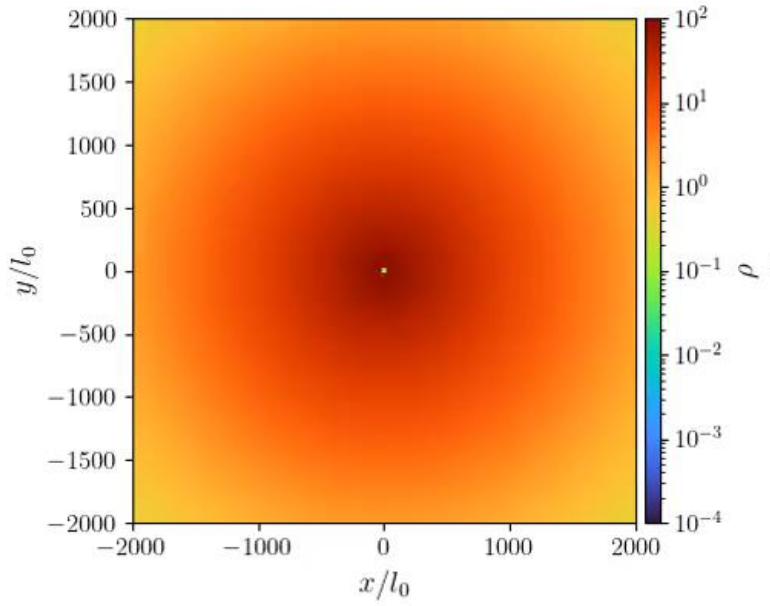




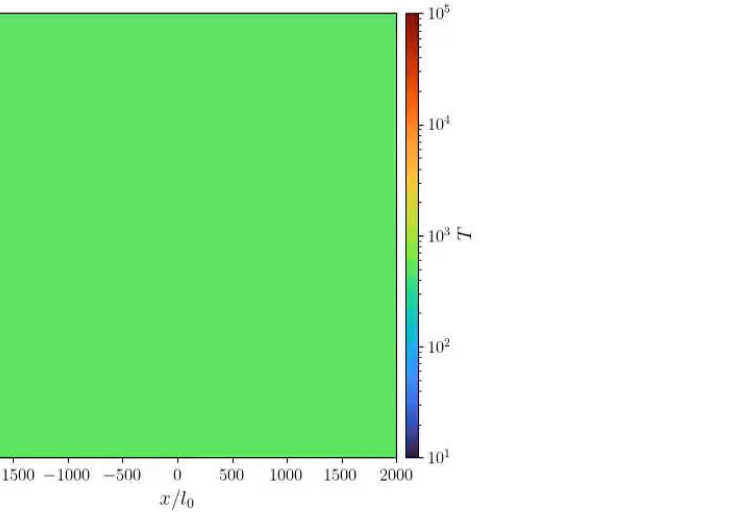
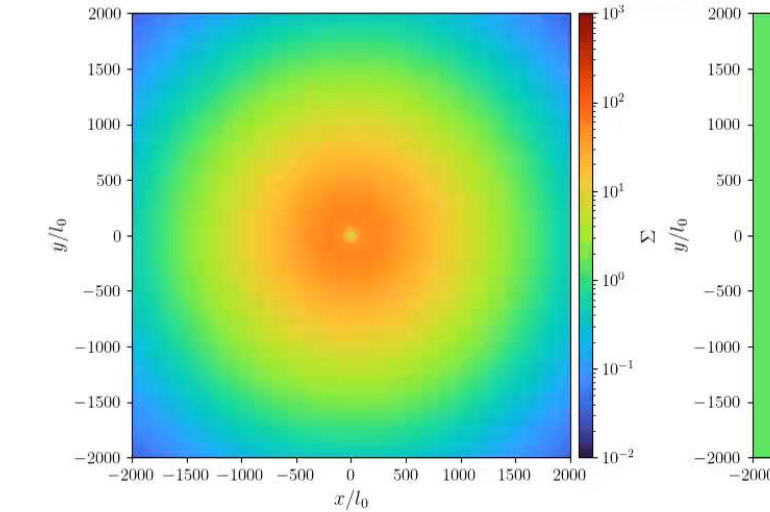
- 252 sims
- 1.5 days
- On this guy
- arXiv:2406.18920



0.0 Myr



0.0 Myr

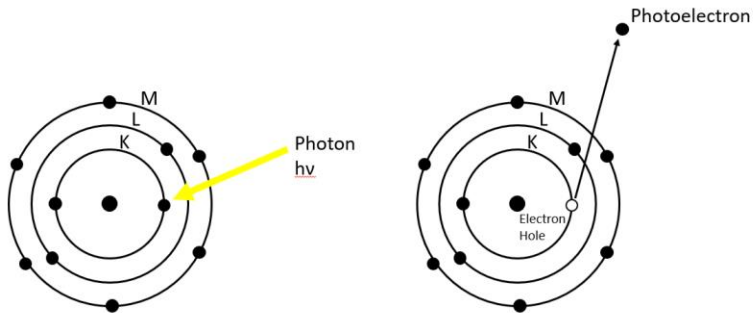
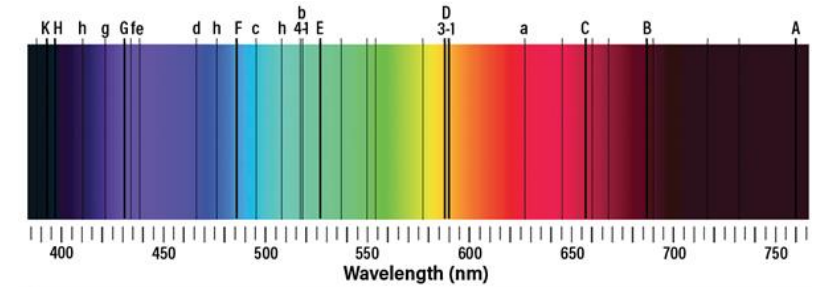
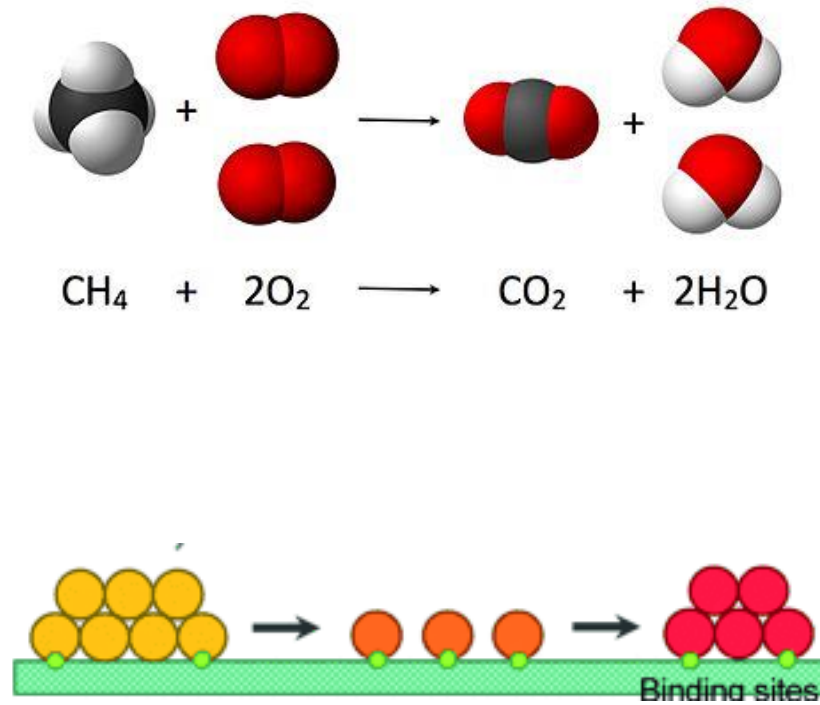
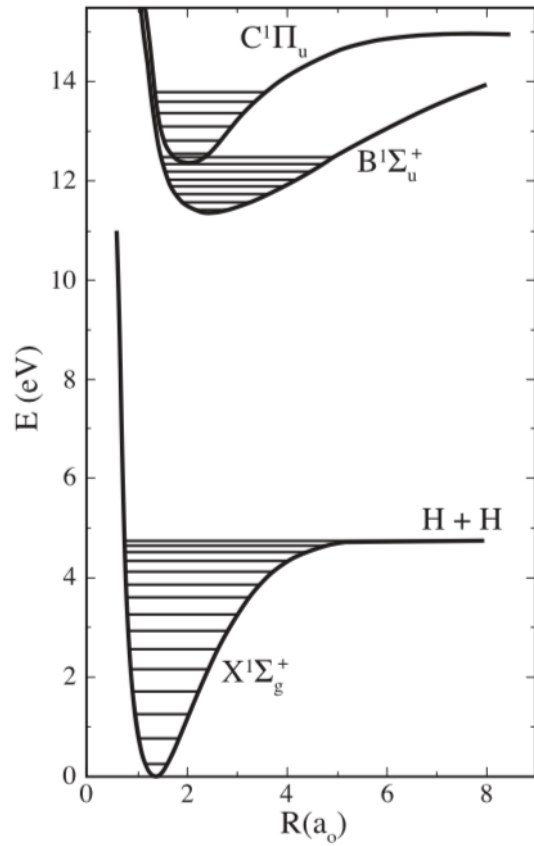


Lv, Wang, Ho et al., in prep.

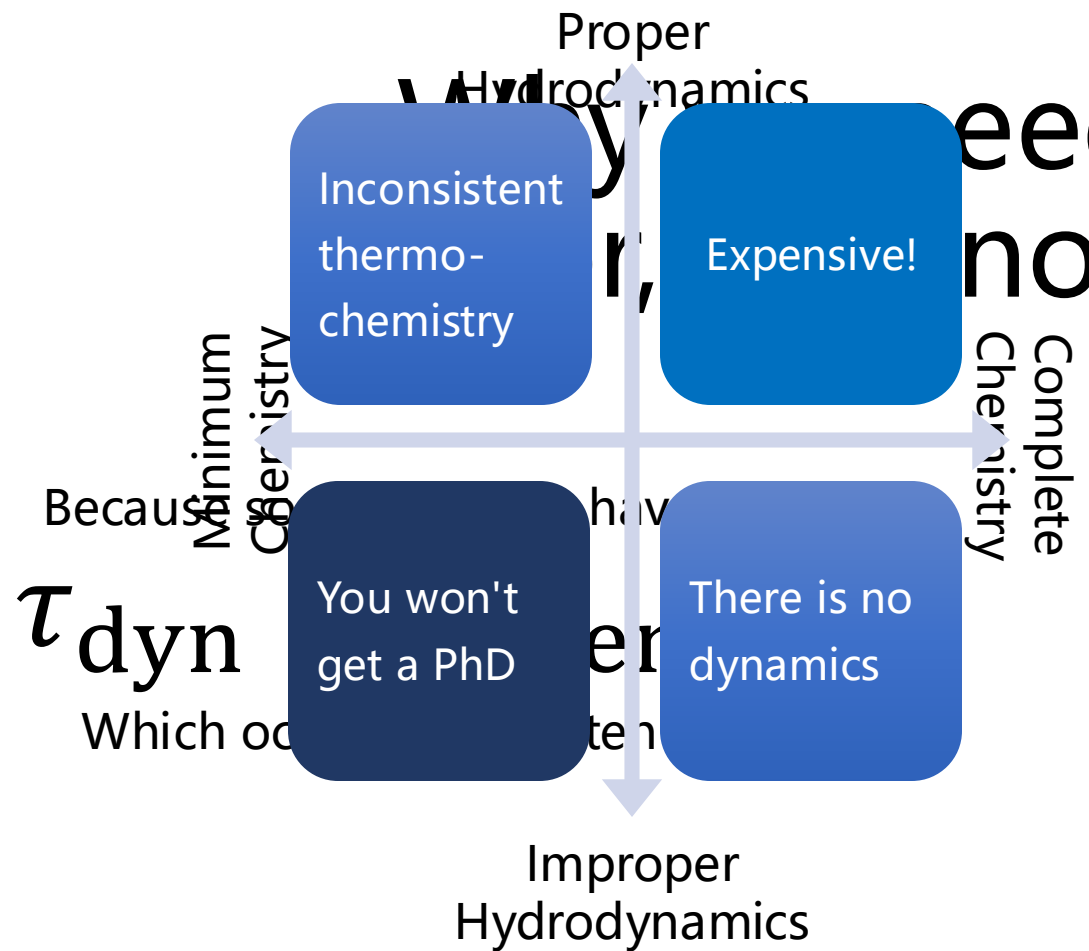
The slide features a central white diamond shape with a thin white border, set against a light gray background. The word "Microphysics" is centered within the diamond. The background is decorated with four overlapping diamond shapes in the corners: yellow in the top-left and bottom-right, and blue in the top-right and bottom-left. The diamonds are semi-transparent and overlap each other and the central diamond.

# Microphysics





Chemistry ( "Microphysics" ) is Necessary



Need it  
not equilibrium?



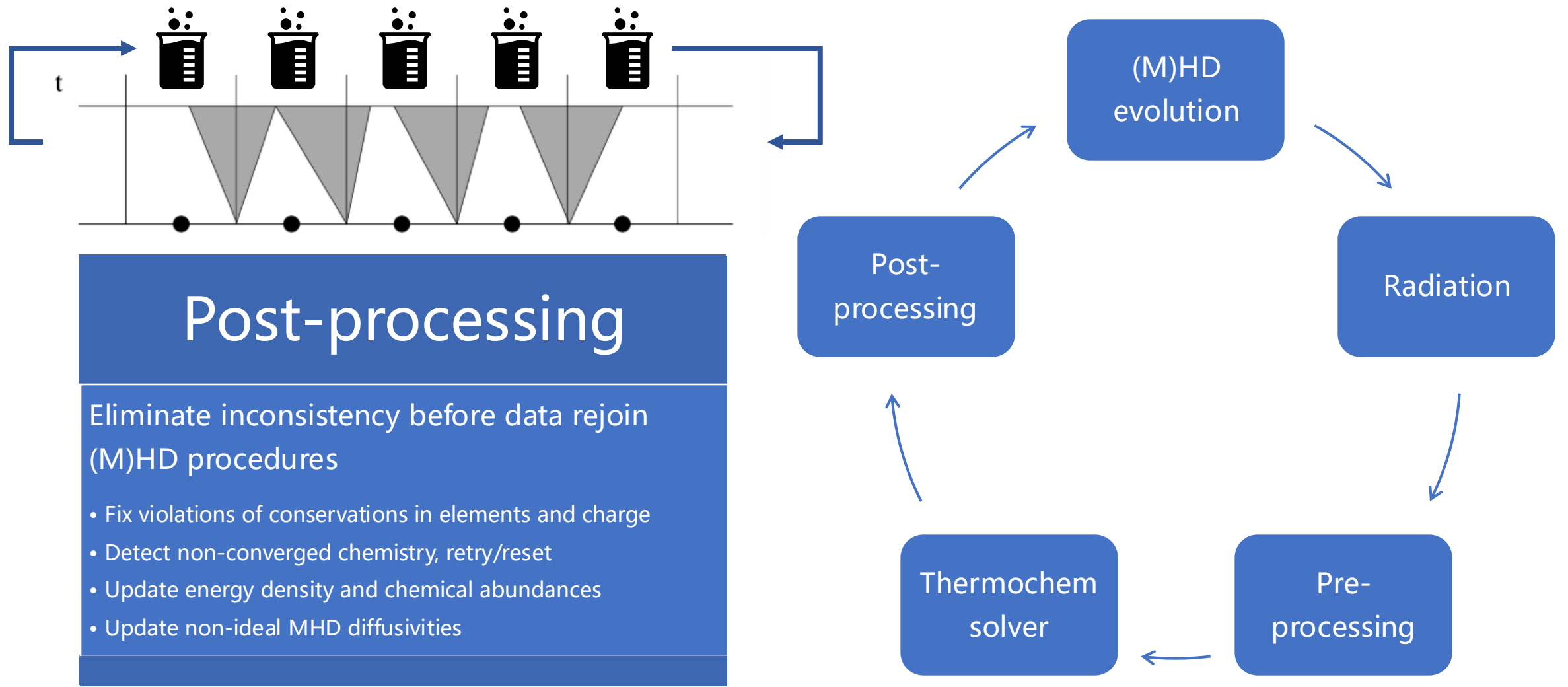
$\tau_{dyn}$   
Which occurs

Because

have

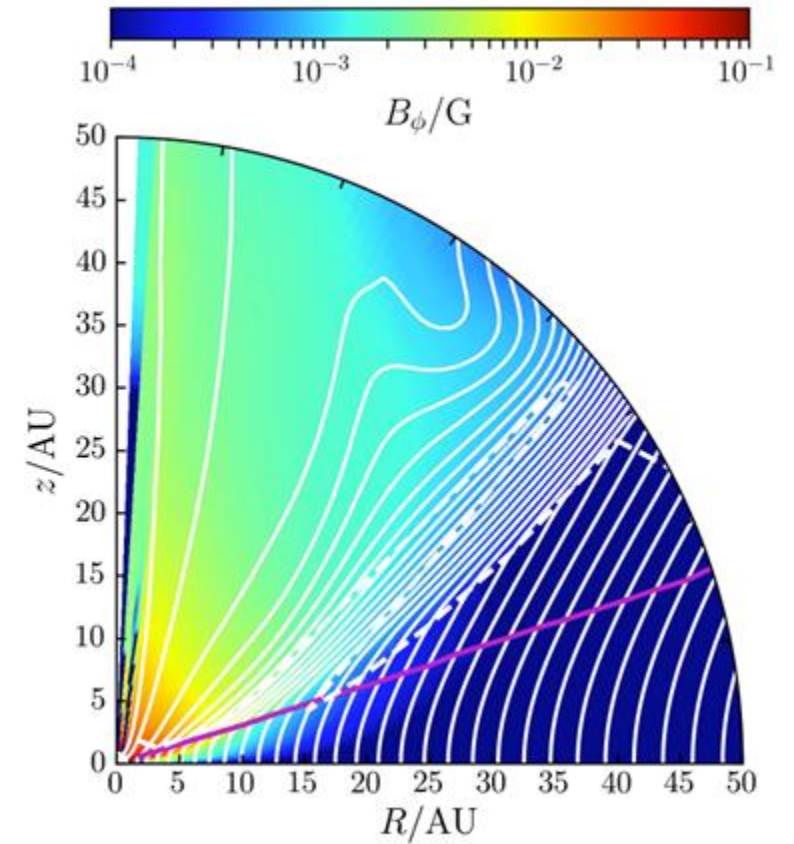
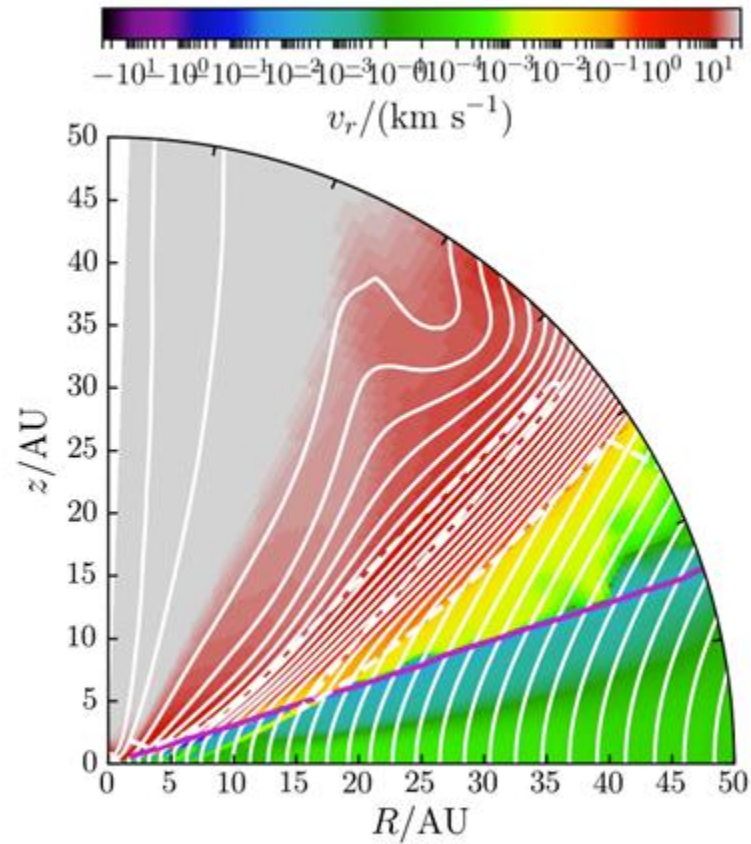
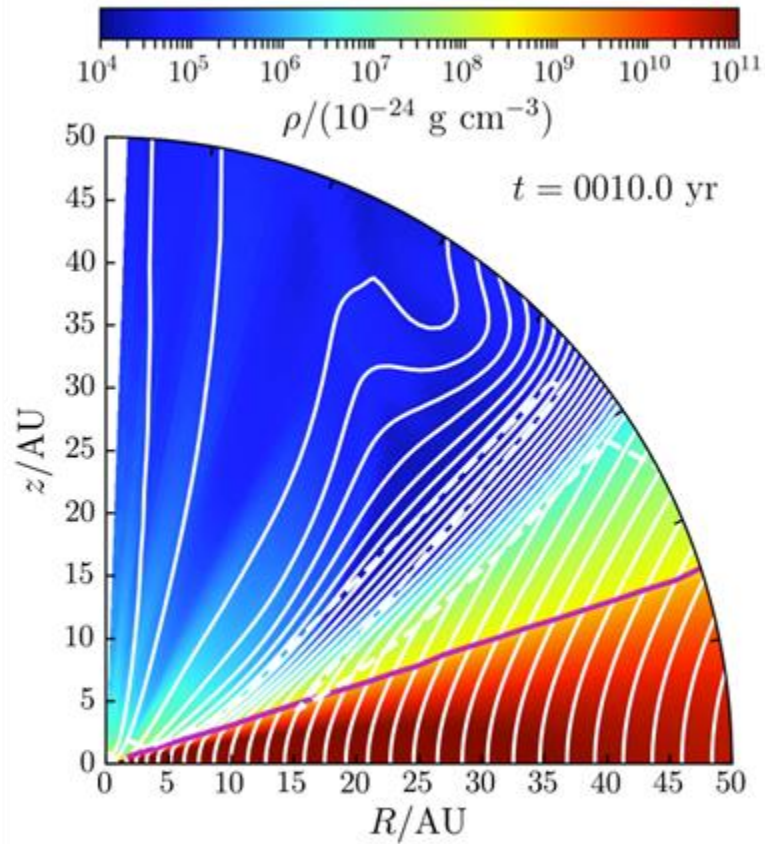
er

ten

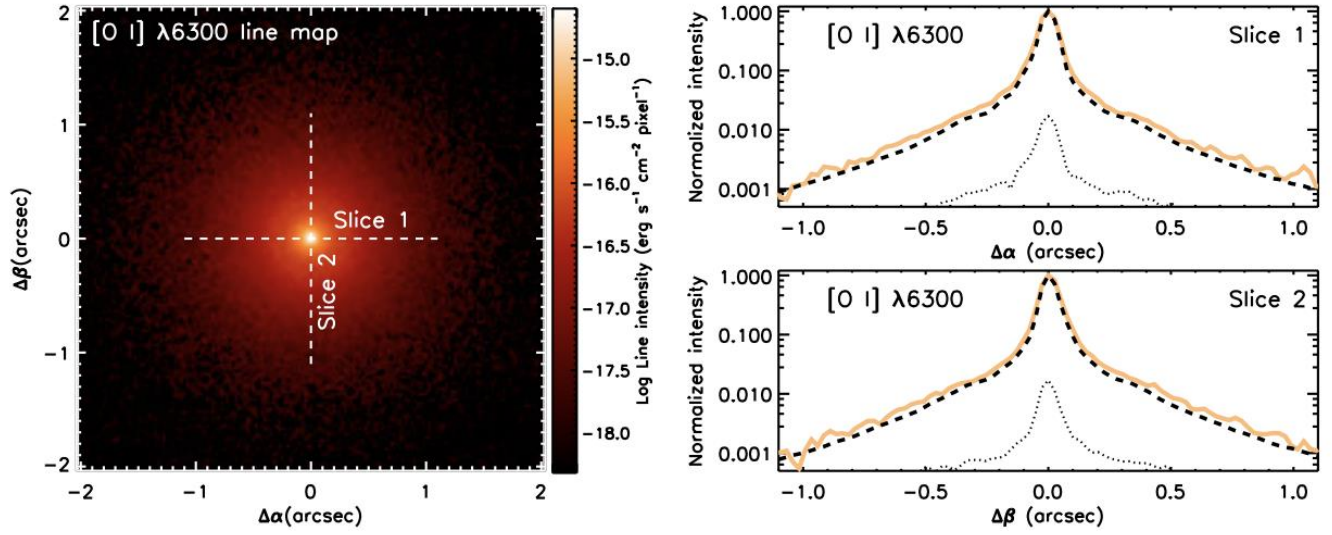
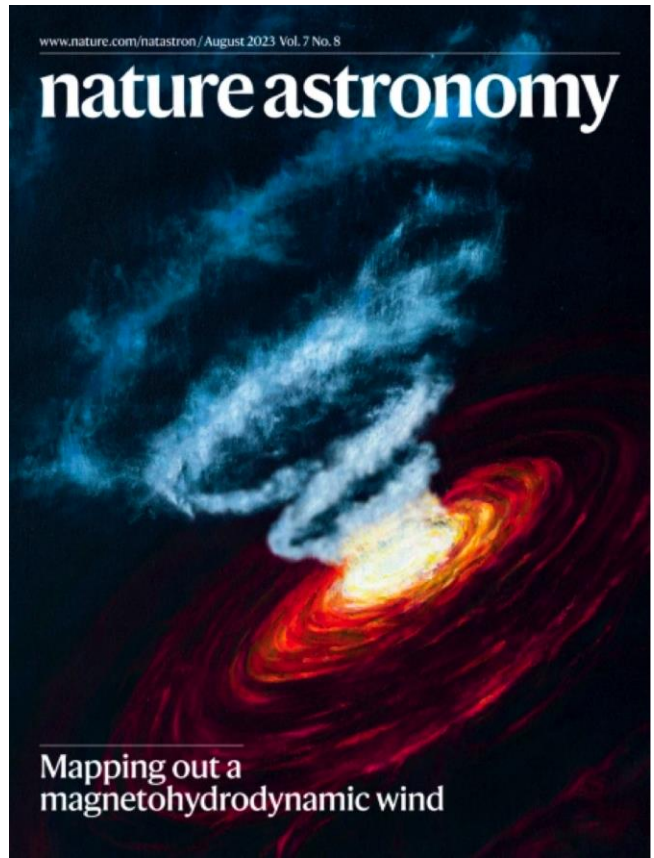
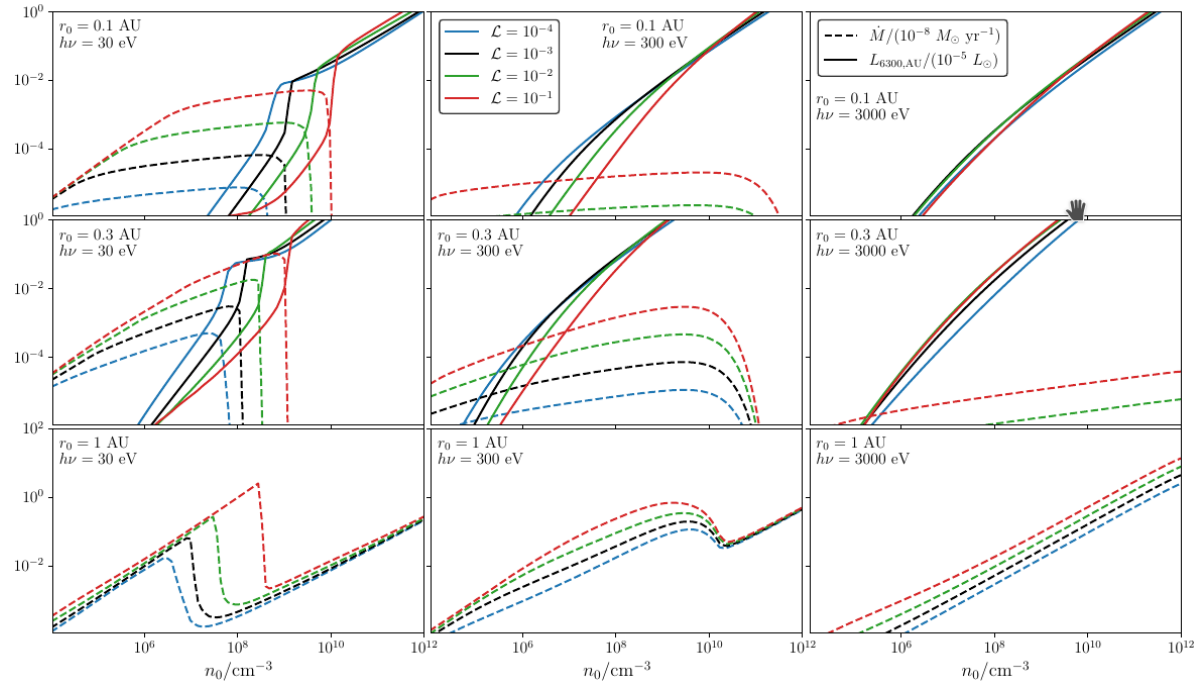
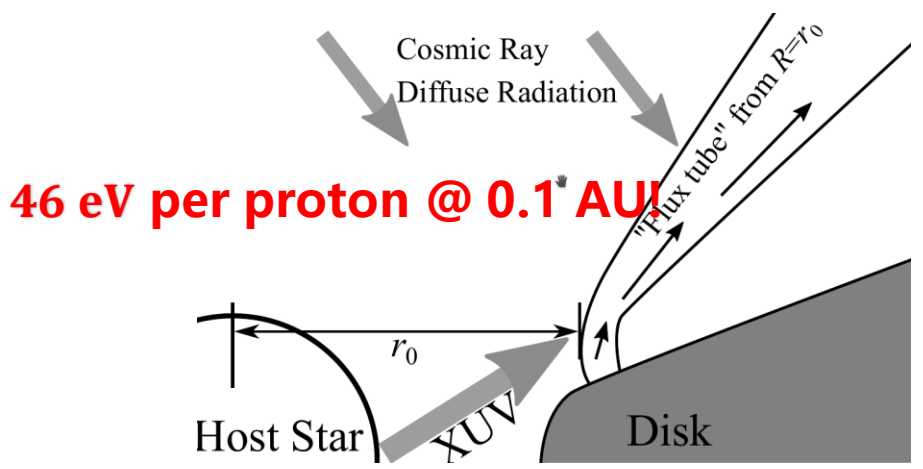


# Realtime GPU Microphysics

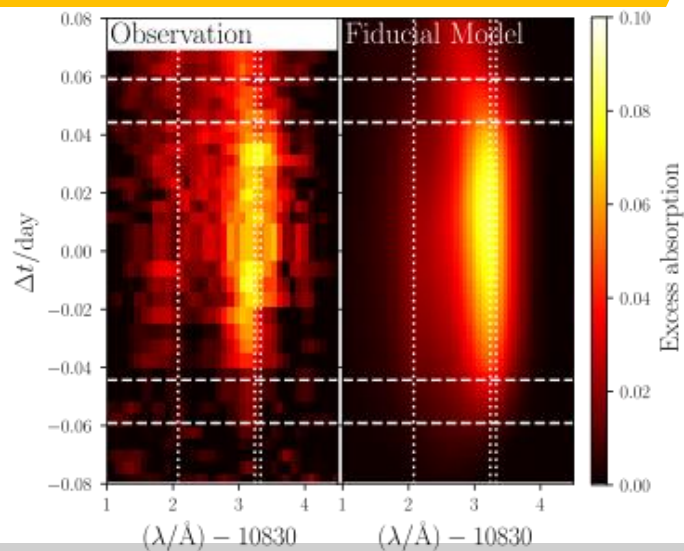
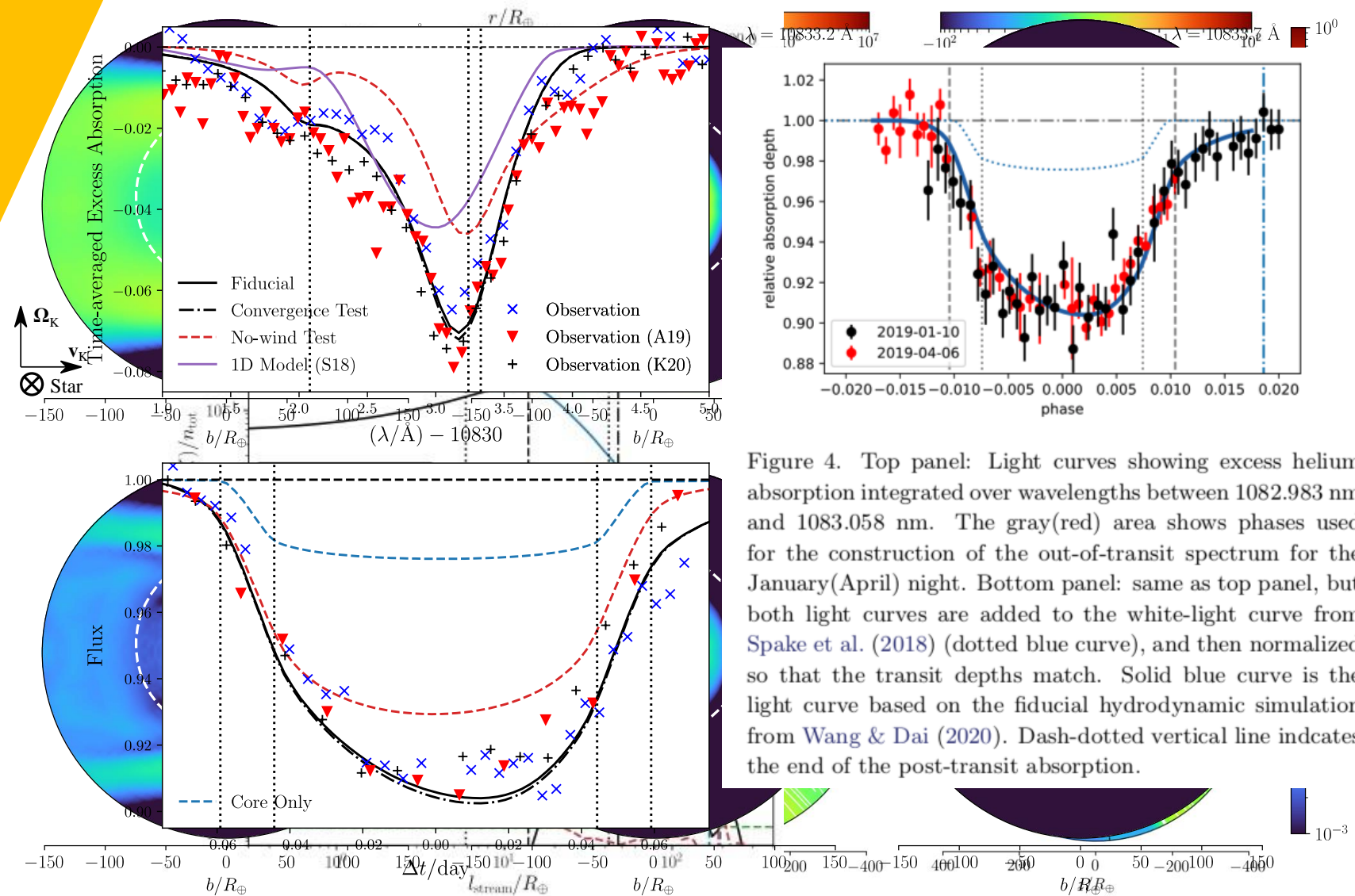
L. Wang  
PhD Thesis  
(2018)



# Wind-driven Accretion: Consistent Microphysics



- $\tau_{\text{dyn}} \sim 1$  hr,  $\tau(\text{He } 2^3S) \sim 1$  hr
- He I 10830 Å absorption: ionized outflows
- Full hydro simulations for observed spectra

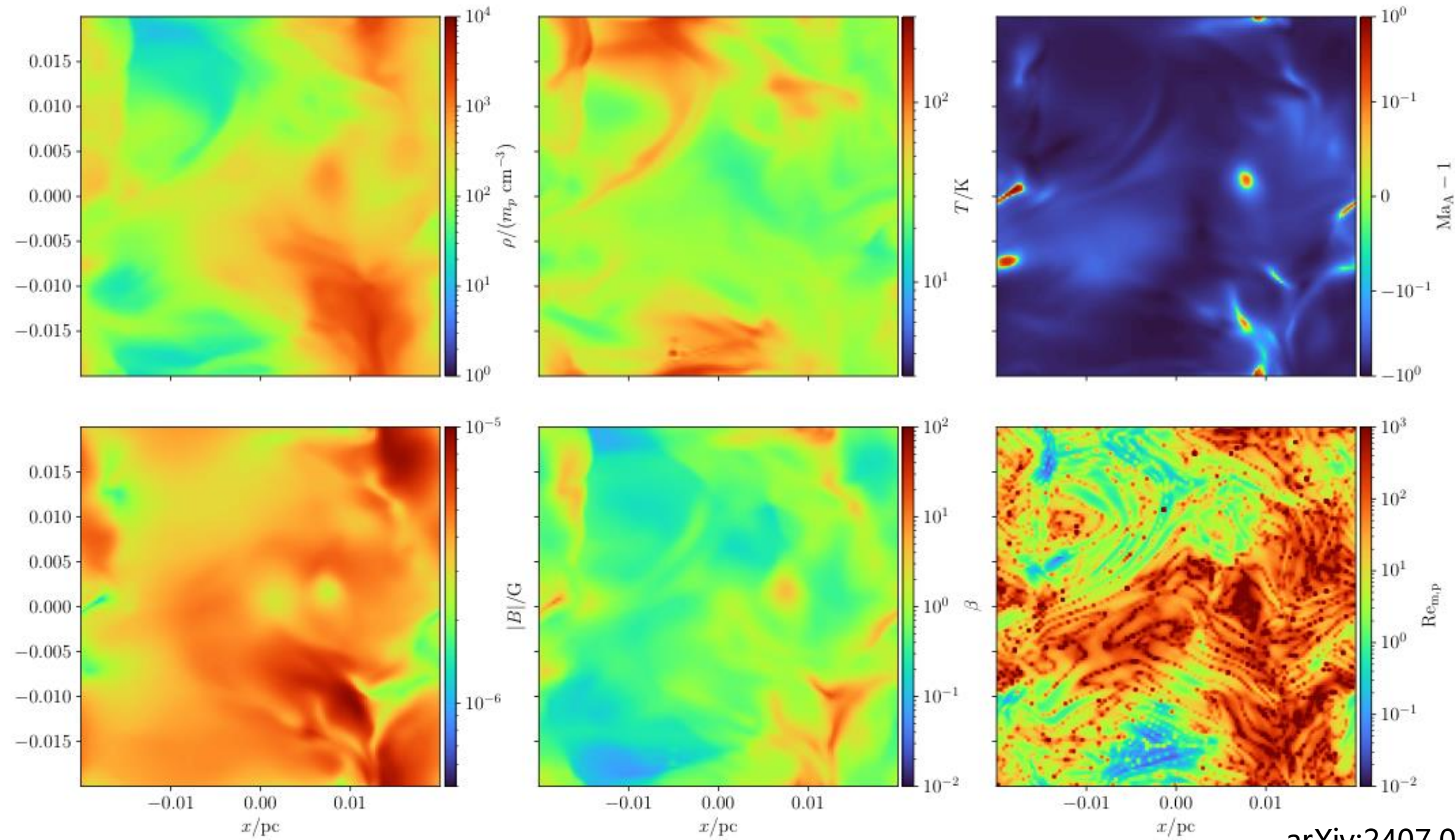


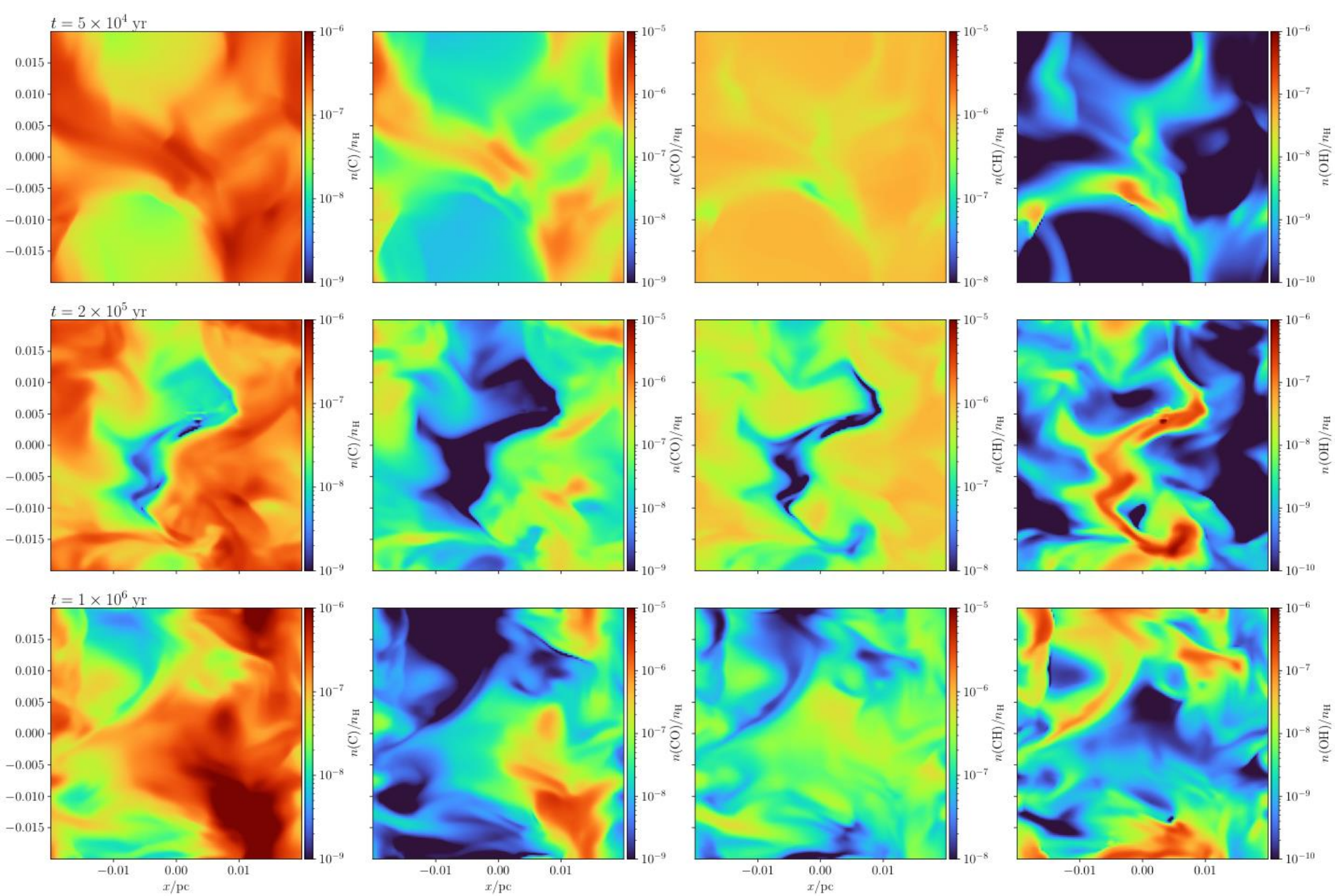
Direct Evidence of Dispersal:

A

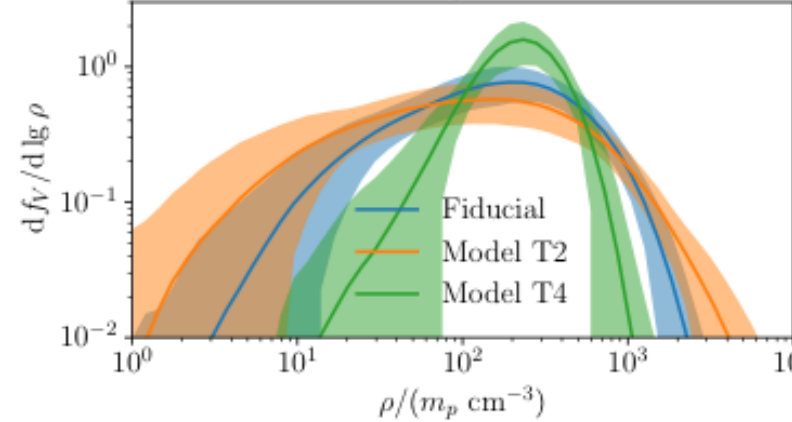
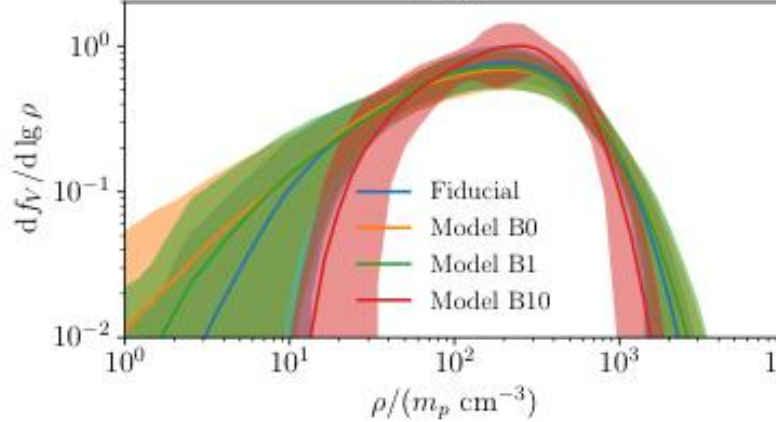
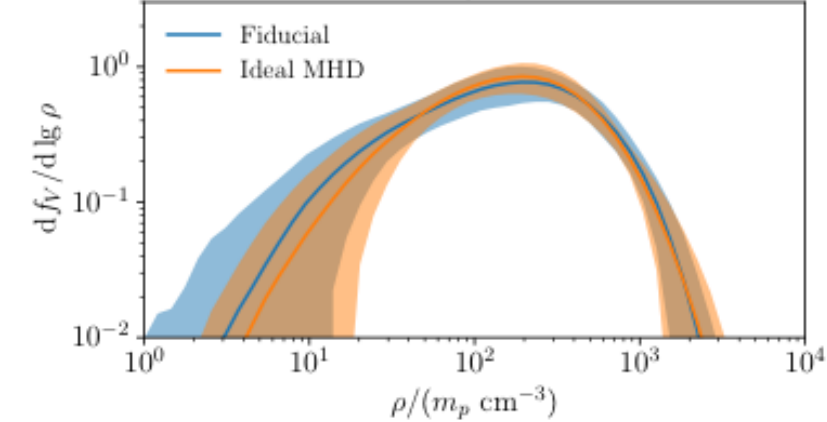
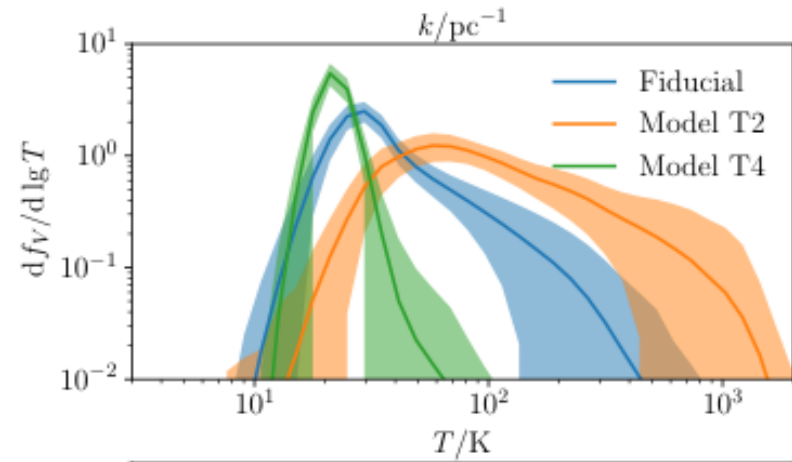
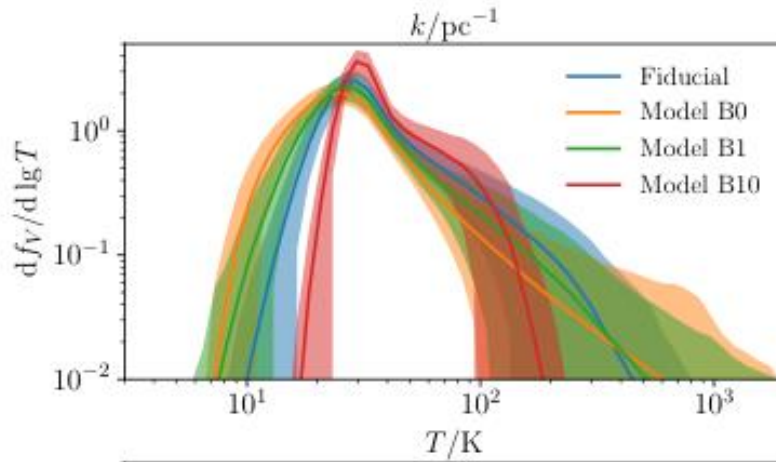
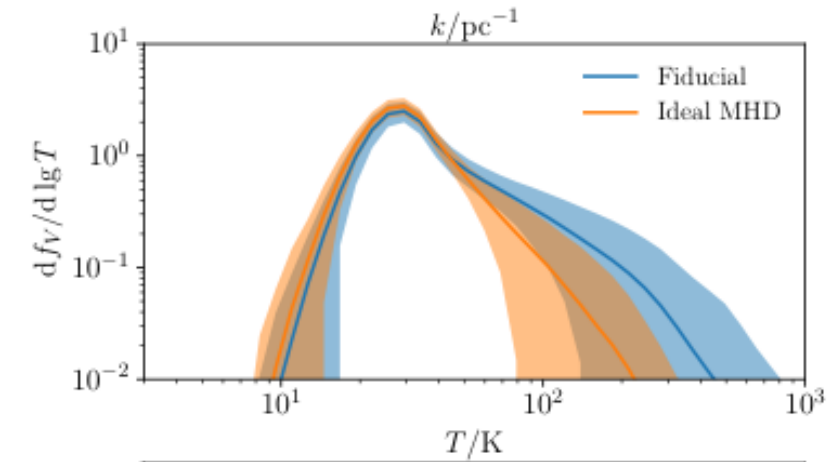
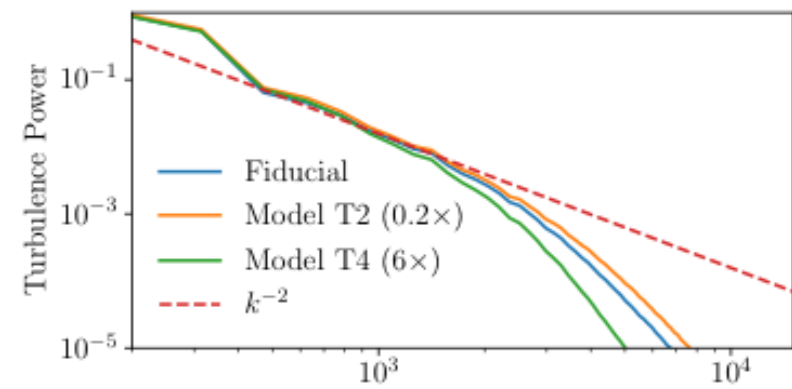
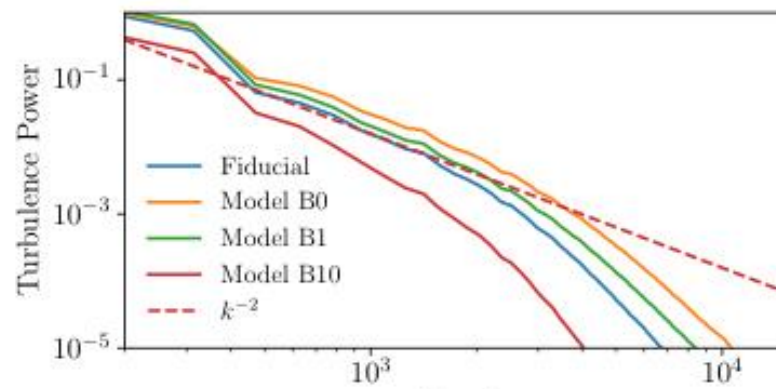
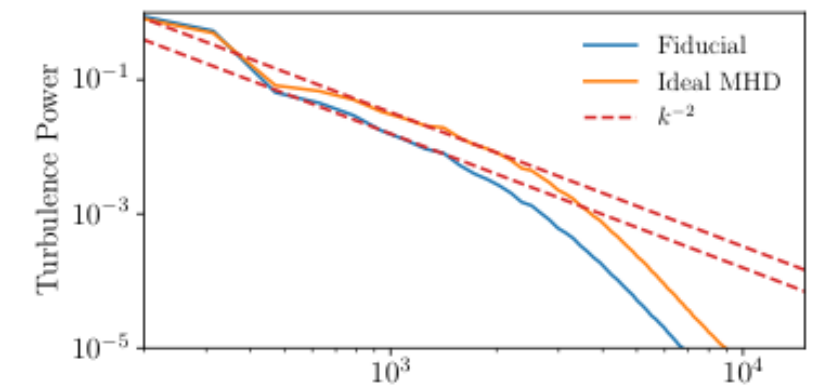
He I 10830

Wang & Dai  
2021 ApJ, 914,  
99









- Any CUDA tutorials--Simply search "CUDA tutorial" with Google/Bing. Specifically, please search these topics:
  - Stream (graphs), Shared Memory, Multi-thread Algorithms
  - cuda-gdb, cuda-memcheck
  - Template metaprogramming of C++
- Documents of CUDA and HIP/ROCm
- <https://github.com/ROCm-Developer-Tools/HIP-CPU>
- Write your own GPU project and play with it!



Search-Engine-Oriented Programming

Questions?

科学吹牛

Aerodynamics of a cow.

