# LUMI

Feedback from the use of the largest supercomputer in Europe
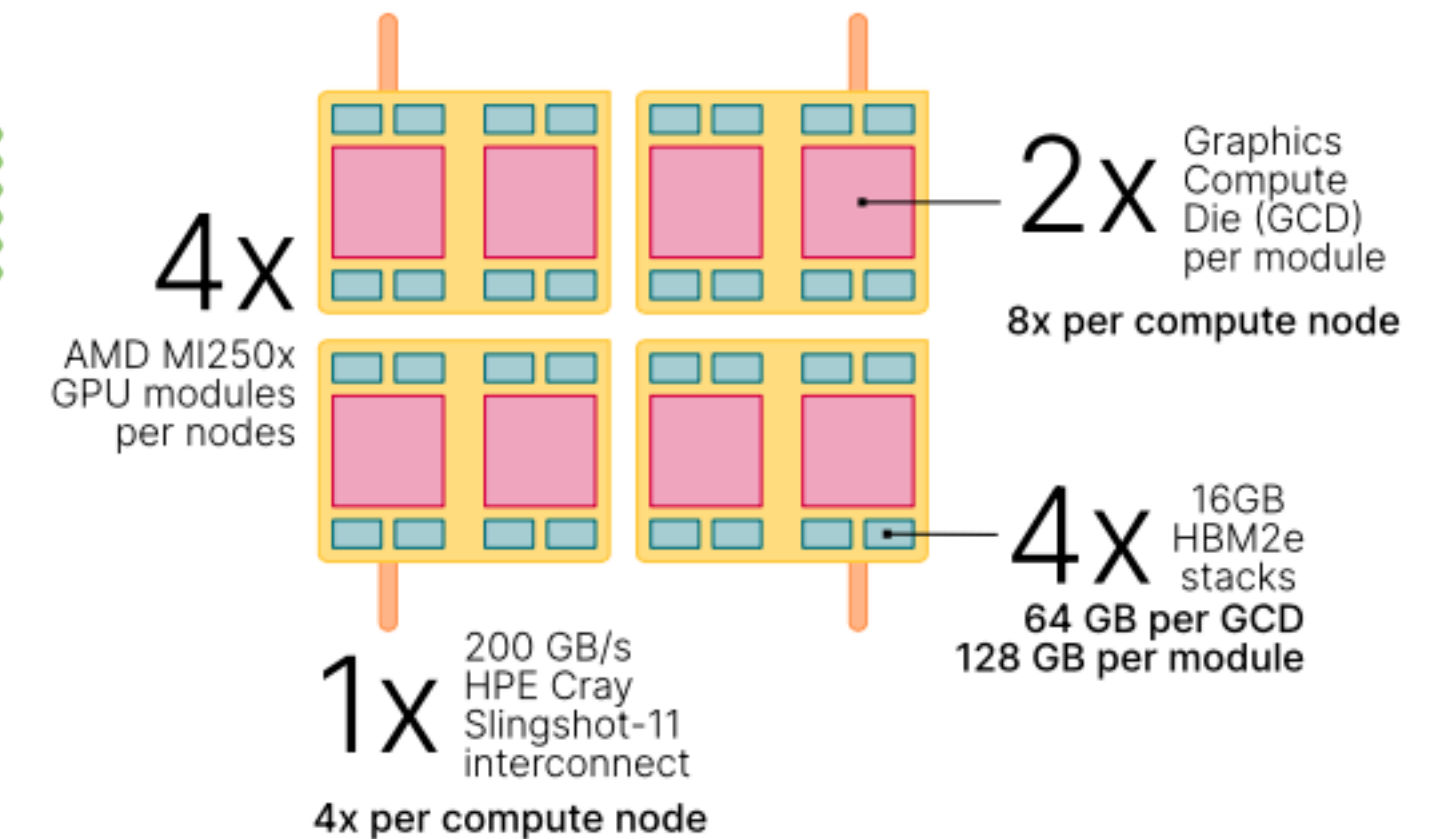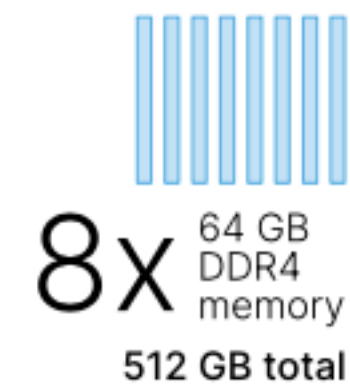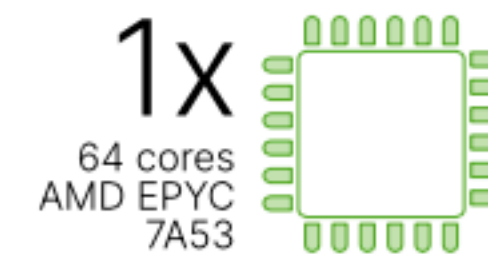
Geoff Lesur

# What is LUMI?

- 2978 nodes with
  - 4 AMD MI250x GPUs
  - a single 64 cores AMD EPYC "Trento" CPU.
- The aggregated HPL Linpack performance of LUMI-G is 379.70 PFlop/s.



*Overview of a LUMI-G compute node*

~ 9x Adastra GPU @ CINES

# Accessing LUMI



eurohpc-ju.europa.eu

- Granted an allocation for the « extreme scale access mode »

- 3,484,000 GPU h allocated on 1 year

# The project « turbulent disks »

## In collaboration with H. Latter (Univ. of Cambridge, UK)



| Vertical Shear Instability | Convective Overstability | Zombie Vortex Instability |
|---|---|---|
| $\tau_{\rm cool}\Omega_{\rm K} \ll 1$ | $\tau_{\rm cool}\Omega_{\rm K} \sim 1$ | $\tau_{\rm cool}\Omega_{\rm K} \gg 1$ |
| $q \neq 0$ | $-1 < p/q < 1/(\gamma-1)$ | $|z| \gtrsim \sqrt{\gamma/(\gamma-1)}H$ |
| $\alpha_{\rm SS} \sim 10^{-4}$ | $\alpha_{\rm SS} \sim 10^{-3}$ | $\alpha_{\rm SS} \sim (10^{-5}\text{–}10^{-4})^{\dagger}$ |
| Outcome: turbulence & vortices | | |

[Lesur+ 2023, PPVII]

- Full disk?
- Survival of corrugation modes?
- Dust settling?
- Vortices?

- Full disk?
- Impact of vertical stratification?
- Competition with VSI?

- Compressibility effects?
- Survival with curvature?
- Propagation to the midplane?

# The challenge

# #1: system stability

```
Gravity: central mass gravitational potential ENABLED with M=1
TimeIntegrator: using 3rd Order (RK3) integrator.
TimeIntegrator: Using adaptive dt with CFL=0.8 .
TimeIntegrator: will stop after 43.5 hours.
Main: Creating initial conditions.
Main: Cycling Time Integrator...
TimeIntegrator:              time |    cycle |      time step | cell (updates/s) | MPI overhead (%)
TimeIntegrator:     0.000000e+00 |        0 |   1.000000e-06 |              N/A |              N/A
TimeIntegrator:     1.593742e-05 |       10 |   2.593742e-06 |     1.228411e+11 |         6.229152
TimeIntegrator:     5.727500e-05 |       20 |   6.727500e-06 |     1.263376e+11 |         6.696973
TimeIntegrator:     1.644940e-04 |       30 |   1.744940e-05 |     1.260329e+11 |         6.868502
TimeIntegrator:     4.425926e-04 |       40 |   4.525926e-05 |     1.260741e+11 |         6.713010
TimeIntegrator:     1.163909e-03 |       50 |   1.173909e-04 |     1.258320e+11 |         6.739630
TimeIntegrator:     3.034816e-03 |       60 |   3.044816e-04 |     1.251714e+11 |         6.986361
TimeIntegrator:     6.519082e-03 |       70 |   3.555275e-04 |     1.245450e+11 |         7.566431
TimeIntegrator:     1.007368e-02 |       80 |   3.553797e-04 |     1.238352e+11 |         8.051875
TimeIntegrator:     1.362683e-02 |       90 |   3.552368e-04 |     1.236095e+11 |         8.408736
TimeIntegrator:     1.717857e-02 |      100 |   3.550985e-04 |     1.230930e+11 |         8.999555
TimeIntegrator:     2.072898e-02 |      110 |   3.549738e-04 |     1.191618e+11 |        10.278076
...
TimeIntegrator:     6.680007e-02 |      240 |   3.539673e-04 |     8.560018e+10 |        31.998791
TimeIntegrator:     7.033955e-02 |      250 |   3.539236e-04 |     8.455433e+10 |        32.354674
TimeIntegrator:     7.387860e-02 |      260 |   3.538831e-04 |     8.536546e+10 |        31.871857
TimeIntegrator:     7.741726e-02 |      270 |   3.538455e-04 |     8.451818e+10 |        32.560332
...
TimeIntegrator:     3.213042e-01 |      960 |   3.533728e-04 |     7.226361e+10 |        39.876328
TimeIntegrator:     3.248379e-01 |      970 |   3.533726e-04 |     7.318458e+10 |        39.556451
TimeIntegrator:     3.283717e-01 |      980 |   3.533725e-04 |     7.281833e+10 |        39.854619
TimeIntegrator:     3.319054e-01 |      990 |   3.533724e-04 |     7.167976e+10 |        40.422609
Main: Reached maximum number of integration cycles.
Main: Reached t=0.335439
Main: Completed in 10 minutes 53 seconds and 1000 cycles
Main: Perfs are 8.174090e+10 cell updates/second
MPI overhead represents 34% of total run time.
```

# #1: system stability (cont'd)

## process #1

```
Gravity: central mass gravitational potential ENABLED with M=1
TimeIntegrator: using 3rd Order (RK3) integrator.
TimeIntegrator: Using adaptive dt with CFL=0.8 .
TimeIntegrator: will stop after 43.5 hours.
Main: Creating initial conditions.
Main: Cycling Time Integrator...
TimeIntegrator:          time |  cycle |    time step | cell (updates/s) | MPI overhead (%)
TimeIntegrator:  0.000000e+00 |      0 | 1.000000e-06 |              N/A |             N/A
TimeIntegrator:  1.593742e-05 |     10 | 2.593742e-06 |     7.951065e+09 |        3.249777
TimeIntegrator:  5.727500e-05 |     20 | 6.727500e-06 |     8.028158e+09 |        3.289534
TimeIntegrator:  1.644940e-04 |     30 | 1.744940e-05 |     8.014591e+09 |        3.362284
TimeIntegrator:  4.425926e-04 |     40 | 4.525926e-05 |     7.998629e+09 |        3.436400
TimeIntegrator:  1.163909e-03 |     50 | 1.173909e-04 |     7.964717e+09 |        3.557654
TimeIntegrator:  3.034816e-03 |     60 | 3.044816e-04 |     7.903785e+09 |        3.754029
TimeIntegrator:  7.327936e-03 |     70 | 4.867873e-04 |     7.863393e+09 |        3.629722
TimeIntegrator:  1.219500e-02 |     80 | 4.866076e-04 |     7.798375e+09 |        3.673882
TimeIntegrator:  1.706028e-02 |     90 | 4.864313e-04 |     7.740056e+09 |        3.972529
TimeIntegrator:  2.192384e-02 |    100 | 4.862645e-04 |     7.717967e+09 |        4.076347
TimeIntegrator:  2.678575e-02 |    110 | 4.861024e-04 |     7.671914e+09 |        4.058943
TimeIntegrator:  3.164608e-02 |    120 | 4.859492e-04 |     7.693121e+09 |        4.226520
TimeIntegrator:  3.650490e-02 |    130 | 4.858028e-04 |     7.691725e+09 |        4.318757
TimeIntegrator:  4.136232e-02 |    140 | 4.856696e-04 |     7.680140e+09 |        4.480481
TimeIntegrator:  4.621844e-02 |    150 | 4.855419e-04 |     7.635158e+09 |        4.922653
TimeIntegrator:  5.107329e-02 |    160 | 4.854182e-04 |     7.480662e+09 |        5.866868
TimeIntegrator:  5.592695e-02 |    170 | 4.853023e-04 |     7.055460e+09 |        8.341483
TimeIntegrator:  6.077947e-02 |    180 | 4.851931e-04 |     6.830506e+09 |        9.758574
TimeIntegrator:  6.563094e-02 |    190 | 4.850902e-04 |     6.774158e+09 |        9.520462
TimeIntegrator:  7.048139e-02 |    200 | 4.849914e-04 |     6.675399e+09 |       10.142318
TimeIntegrator:  7.533088e-02 |    210 | 4.848983e-04 |     6.639014e+09 |        8.720803
TimeIntegrator:  8.017947e-02 |    220 | 4.848107e-04 |     6.534640e+09 |        9.950846
TimeIntegrator:  8.502719e-02 |    230 | 4.847271e-04 |     6.529307e+09 |        9.881086
TimeIntegrator:  8.987411e-02 |    240 | 4.846491e-04 |     6.414888e+09 |        9.832200
TimeIntegrator:  9.472026e-02 |    250 | 4.845747e-04 |     6.354311e+09 |       10.337835
TimeIntegrator:  9.956569e-02 |    260 | 4.845042e-04 |     6.309454e+09 |       10.386557
TimeIntegrator:  1.044104e-01 |    270 | 4.844372e-04 |     6.315063e+09 |       10.402304
TimeIntegrator:  1.092545e-01 |    280 | 4.843735e-04 |     6.211119e+09 |       11.243196
TimeIntegrator:  1.140980e-01 |    290 | 4.843146e-04 |     6.175291e+09 |       10.965431
TimeIntegrator:  1.189409e-01 |    300 | 4.842579e-04 |     6.113392e+09 |       11.152301
TimeIntegrator:  1.237832e-01 |    310 | 4.842048e-04 |     6.023583e+09 |       10.558754
TimeIntegrator:  1.286250e-01 |    320 | 4.841539e-04 |     5.978660e+09 |       11.492593
TimeIntegrator:  1.334663e-01 |    330 | 4.841063e-04 |     5.992310e+09 |       11.478762
TimeIntegrator:  1.383072e-01 |    340 | 4.840615e-04 |     5.931508e+09 |       11.957992
TimeIntegrator:  1.431476e-01 |    350 | 4.840187e-04 |     5.882077e+09 |       11.923003
TimeIntegrator:  1.479876e-01 |    360 | 4.839782e-04 |     5.926814e+09 |       11.963534
TimeIntegrator:  1.528272e-01 |    370 | 4.839398e-04 |     5.821199e+09 |       12.286236
TimeIntegrator:  1.576665e-01 |    380 | 4.839036e-04 |     5.874892e+09 |       11.783085
```

## process #24

```
Gravity: central mass gravitational potential ENABLED with M=1
TimeIntegrator: using 3rd Order (RK3) integrator.
TimeIntegrator: Using adaptive dt with CFL=0.8 .
TimeIntegrator: will stop after 43.5 hours.
Main: Creating initial conditions.
Main: Cycling Time Integrator...
TimeIntegrator:          time |  cycle |    time step | cell (updates/s) | MPI overhead (%)
TimeIntegrator:  0.000000e+00 |      0 | 1.000000e-06 |              N/A |             N/A
TimeIntegrator:  1.593742e-05 |     10 | 2.593742e-06 |     7.951023e+09 |        3.904353
TimeIntegrator:  5.727500e-05 |     20 | 6.727500e-06 |     8.028158e+09 |        4.072326
TimeIntegrator:  1.644940e-04 |     30 | 1.744940e-05 |     8.008421e+09 |        4.212089
TimeIntegrator:  4.425926e-04 |     40 | 4.525926e-05 |     8.004786e+09 |        4.183434
TimeIntegrator:  1.163909e-03 |     50 | 1.173909e-04 |     7.961734e+09 |        4.307905
TimeIntegrator:  3.034816e-03 |     60 | 3.044816e-04 |     7.893845e+09 |        3.920318
TimeIntegrator:  7.327936e-03 |     70 | 4.867873e-04 |     7.868658e+09 |        3.739326
TimeIntegrator:  1.219500e-02 |     80 | 4.866076e-04 |     7.793467e+09 |        3.770076
TimeIntegrator:  1.706028e-02 |     90 | 4.864313e-04 |     7.739504e+09 |        3.993650
TimeIntegrator:  2.192384e-02 |    100 | 4.862645e-04 |     7.708094e+09 |        4.231982
TimeIntegrator:  2.678575e-02 |    110 | 4.861024e-04 |     7.680370e+09 |        4.114648
TimeIntegrator:  3.164608e-02 |    120 | 4.859492e-04 |     7.686638e+09 |        4.273643
TimeIntegrator:  3.650490e-02 |    130 | 4.858028e-04 |     7.699252e+09 |        4.157207
TimeIntegrator:  4.136232e-02 |    140 | 4.856696e-04 |     7.682005e+09 |        4.177367
TimeIntegrator:  4.621844e-02 |    150 | 4.855419e-04 |     7.624046e+09 |        3.774363
TimeIntegrator:  5.107329e-02 |    160 | 4.854182e-04 |     7.430682e+09 |        3.920086
TimeIntegrator:  5.592695e-02 |    170 | 4.853023e-04 |     7.076315e+09 |        4.172325
TimeIntegrator:  6.077947e-02 |    180 | 4.851931e-04 |     6.801539e+09 |        3.103568
TimeIntegrator:  6.563094e-02 |    190 | 4.850902e-04 |     6.748423e+09 |        4.771410
TimeIntegrator:  7.048139e-02 |    200 | 4.849914e-04 |     6.691333e+09 |        6.146956
TimeIntegrator:  7.533088e-02 |    210 | 4.848983e-04 |     6.617100e+09 |        3.880033
TimeIntegrator:  8.017947e-02 |    220 | 4.848107e-04 |     6.554727e+09 |        3.714286
TimeIntegrator:  8.502719e-02 |    230 | 4.847271e-04 |     6.506650e+09 |        2.458089
TimeIntegrator:  8.987411e-02 |    240 | 4.846491e-04 |     6.414570e+09 |        2.240765
TimeIntegrator:  9.472026e-02 |    250 | 4.845747e-04 |     6.345239e+09 |        2.172663
TimeIntegrator:  9.956569e-02 |    260 | 4.845042e-04 |     6.308356e+09 |        2.252992
TimeIntegrator:  1.044104e-01 |    270 | 4.844372e-04 |     6.317087e+09 |        2.679004
TimeIntegrator:  1.092545e-01 |    280 | 4.843735e-04 |     6.210179e+09 |        2.672542
TimeIntegrator:  1.140980e-01 |    290 | 4.843146e-04 |     6.158299e+09 |        2.735128
TimeIntegrator:  1.189409e-01 |    300 | 4.842579e-04 |     6.117275e+09 |        2.465897
TimeIntegrator:  1.237832e-01 |    310 | 4.842048e-04 |     6.013434e+09 |        2.136093
TimeIntegrator:  1.286250e-01 |    320 | 4.841539e-04 |     5.980344e+09 |        2.037715
TimeIntegrator:  1.334663e-01 |    330 | 4.841063e-04 |     5.979734e+09 |        1.836741
TimeIntegrator:  1.383072e-01 |    340 | 4.840615e-04 |     5.933302e+09 |        2.202197
TimeIntegrator:  1.431476e-01 |    350 | 4.840187e-04 |     5.885458e+09 |        1.841658
TimeIntegrator:  1.479876e-01 |    360 | 4.839782e-04 |     5.929866e+09 |        2.268691
TimeIntegrator:  1.528272e-01 |    370 | 4.839398e-04 |     5.818551e+09 |        4.272906
TimeIntegrator:  1.576665e-01 |    380 | 4.839036e-04 |     5.878025e+09 |        2.621614
```

# #1: system stability (cont'd)

# #1: system stability (end)

- On some GPUs, the cooling was inefficient (growth of « biology » in the pipes = algae)

- To avoid overheating, these GPUs were slowed down automatically by the system.

- The MPI processes running on those « hot » GPUs were always lagging behind

- All the other processes had to wait for the hot one to finish their task

➡ MPI overhead was increasing on every processes but the hot ones

Addition of an MPI balance diagnostic in Idefix, to quickly identify these problems.

# #2 file size

```
-rw-rw---- 1 lesurg l-ipag 2.2T Mar 29 22:13 data.0017.vtk
-rw-rw---- 1 lesurg l-ipag 2.2T Apr  6 00:10 data.0018.vtk
-rw-rw---- 1 lesurg l-ipag 2.2T Apr 23 04:52 data.0019.vtk
-rw-rw---- 1 lesurg l-ipag 2.2T Apr 25 17:49 data.0020.vtk
-rw-rw---- 1 lesurg l-ipag 2.2T May  4 22:06 data.0021.vtk
-rw-rw---- 1 lesurg l-ipag 2.2T May 11 23:29 data.0022.vtk
-rw-rw---- 1 lesurg l-ipag 2.2T May 29 17:21 data.0023.vtk
-rw-rw---- 1 lesurg l-ipag 2.2T Jun 11 06:51 data.0024.vtk
-rw-rw---- 1 lesurg l-ipag 2.2T Jun 23 06:52 data.0025.vtk
```

- VTK files are typically above 1TB
  (NB: idefix is efficient: ~320s to write one of the above VTK file)

- How to open such a large file? (Paraview won't, unless with a lot or RAM)

- We don't have the ressources to directly load these files…

Solution: on-the-fly VTK slices (i.e. 2D slices of the 3D domain). See this afternoon

# #3 node failure

```
MPICH ERROR [Rank 3284] [job id 6471201.0] [Mon Mar 25 16:22:20 2024] [nid006342] - Abort(1009397903) (rank 3284 in comm
error in PMPI_Waitall: Other MPI error, error stack:
PMPI_Waitall(378)..............: MPI_Waitall(count=2, req_array=0x14bf148, status_array=0x7ffe3211d520) failed
MPIR_Waitall(167)..............:
MPIR_Waitall_impl(51)..........:
MPID_Progress_wait(201)........:
MPIDI_Progress_test(97)........:
MPIDI_OFI_handle_cq_error(1067): OFI poll failed (ofi_events.c:1069:MPIDI_OFI_handle_cq_error:Input/output error - UNDELI
...
srun: error: Node failure on nid005989
srun: Force Terminated StepId=6471201.0
slurmstepd: error: *** JOB 6471201 ON nid005280 CANCELLED AT 2024-03-25T16:29:07 DUE TO NODE FAILURE, SEE SLURMCTLD LOG
***
```

… this happens **a lot**

# #3 node failure (cont'd)

After 2 months of intense exchange with LUMI support:

*I have been discussing failure rate with the SysAdmins team and it is*
*approximately 3 random GPU node failures per day observed (there are almost*
*3000 GPU nodes in total). This rough statistics gives at least some*
*understanding of expected level of reliability or mean time to failure. I have*
*been also again instructed these failures do not share any common*
*characteristics. In other words, current diagnosis is "bad luck".*

In other words, a job running on 1000 GPUs for 24h *is expected to fail*

Solution(s):
- self-erasing dump files (but not satisfactory: significant I/O overheads)
- Idefix should rely on a resilience library (e.g. kokkos resilience): is it enough?

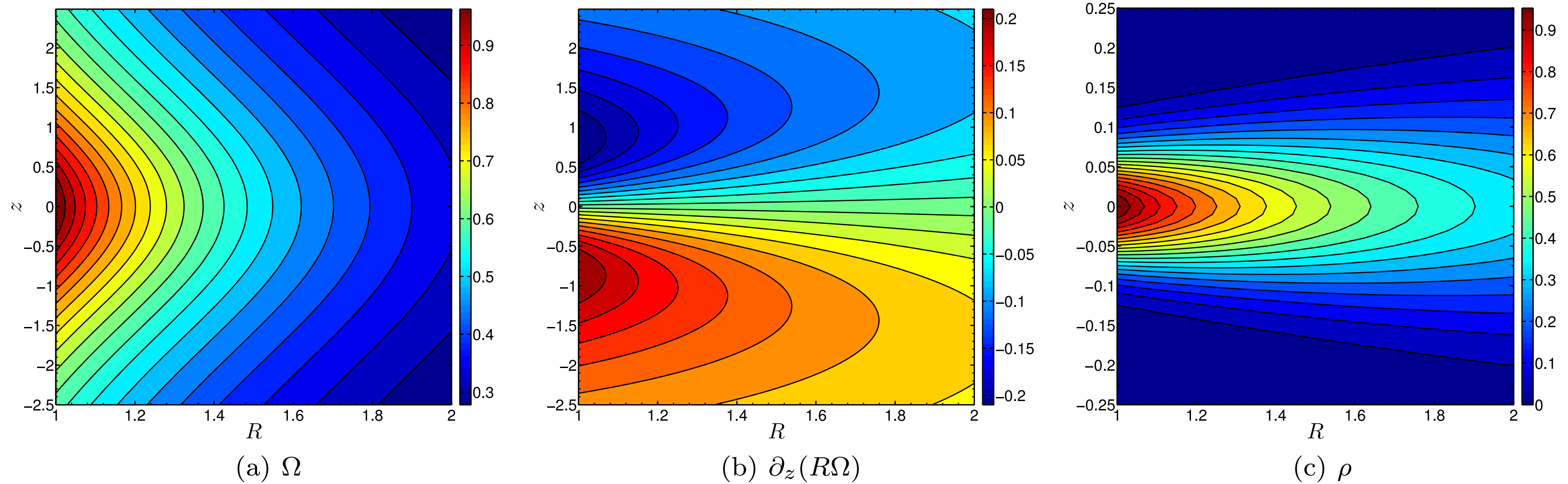# VSI

# VSI=vertical shear instability



(a) $\Omega$           (b) $\partial_z(R\Omega)$           (c) $\rho$

**Figure 1.** Basic state for the locally isothermal disc with $q = -1$, $p = -1.5$ and $c_0 = 0.05$. The left-hand panel shows a contour plot of $\Omega$ on the $(R, z)$ plane. The middle panel is a similar contour plot, but this shows the magnitude of the vertical shear $\partial_z(R\Omega)$, which has a maximum at $|z| \sim 1$ (whereas the scaleheight at the inner radial boundary is 0.05). The right-hand panel shows the density $\rho$.

[Barker & Latter 2015]

$$T \propto R^{-q}$$

linear growth rate:      $\sigma \approx |\partial_z(R\Omega)| \sim \epsilon|q|\Omega,$
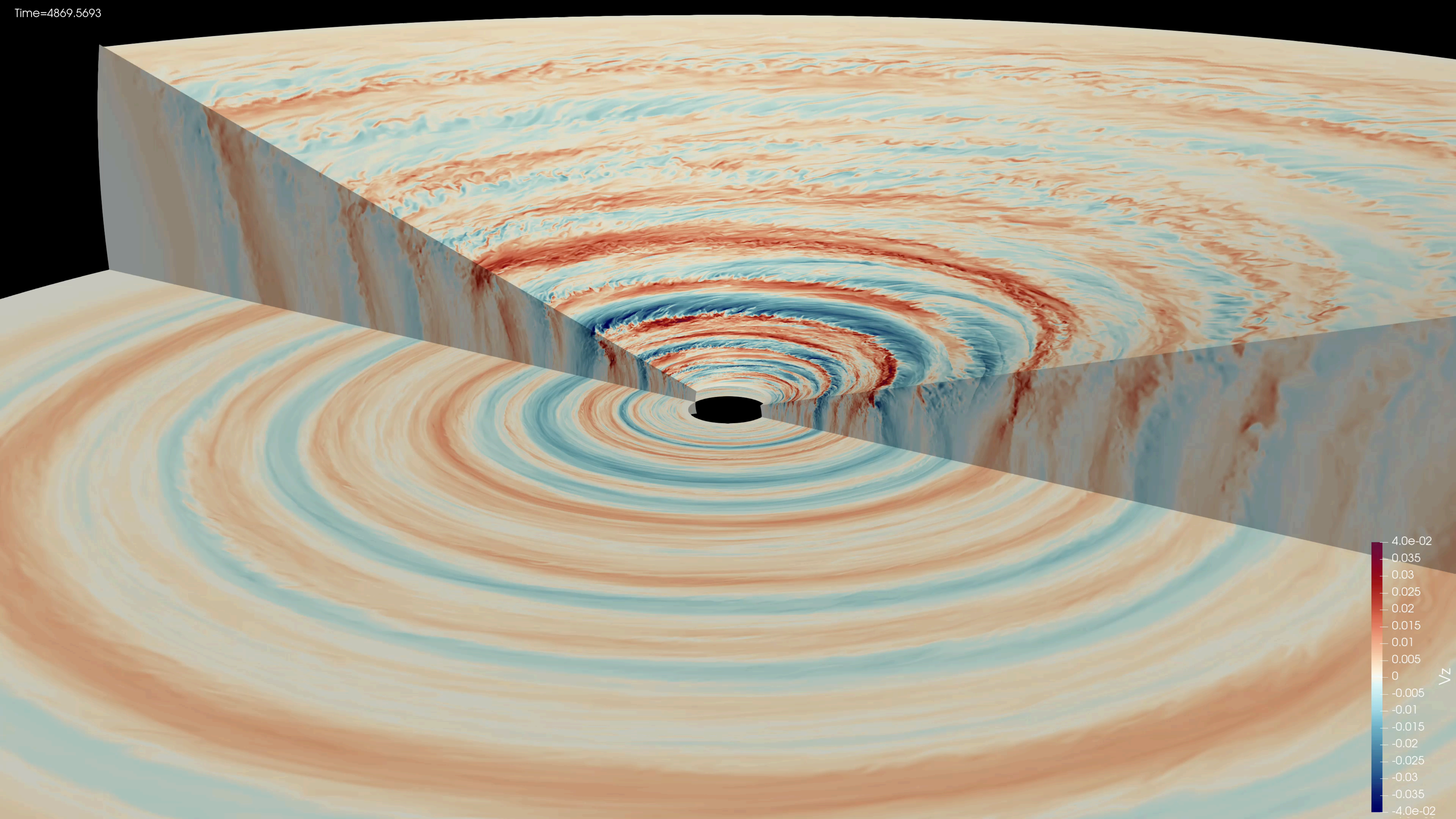
# A high resolution VSI run

- Idefix code on AMD Mi250 GPUs (LUMI, Finland)

- 70pts/H, 3 directions

- High order reconstruction (LimO3)

- Fargo scheme

- Large aspect ratio ($R_{\mathrm{out}}/R_{\mathrm{in}}$=25)
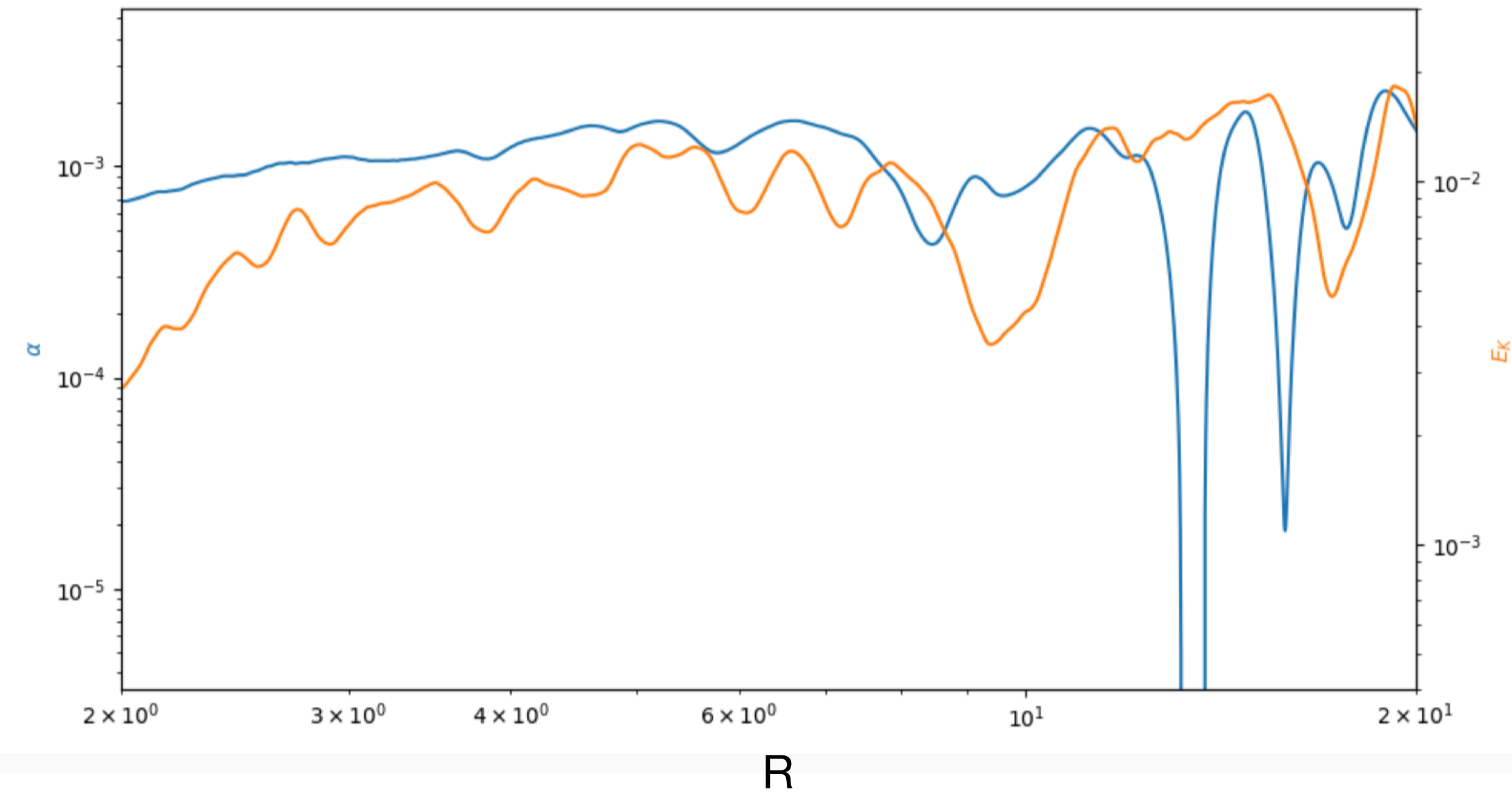
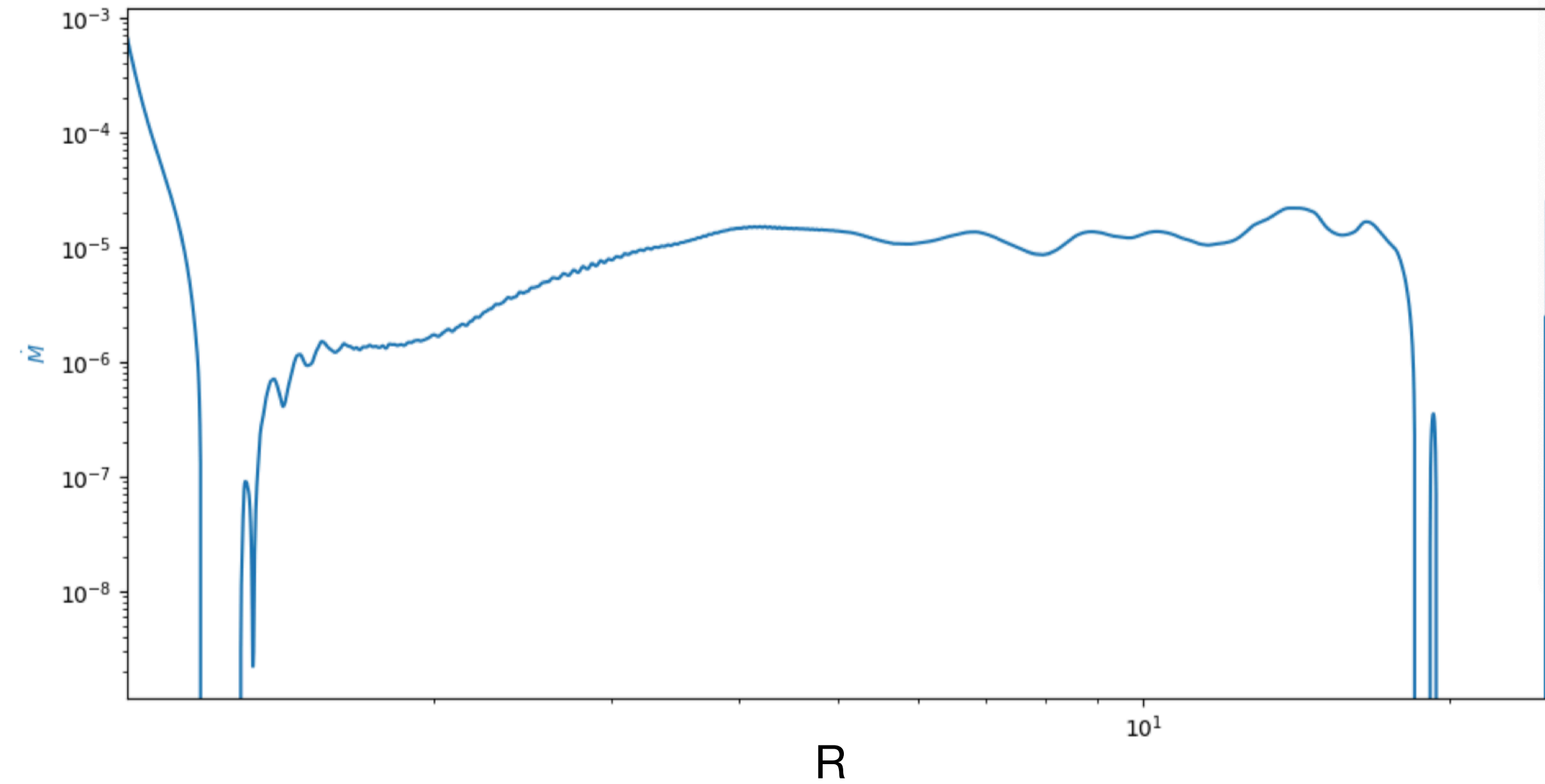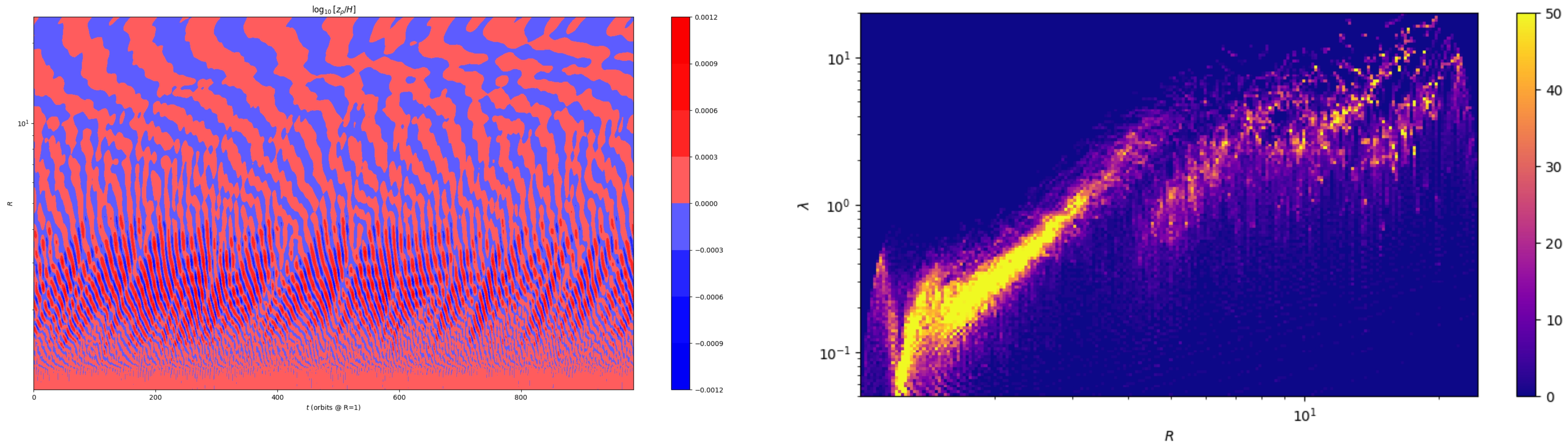- 4 dust size (pressure less fluid)

Time=4869.5693

Vz

4.0e-02
0.035
0.03
0.025
0.02
0.015
0.01
0.005
0
-0.005
-0.01
-0.015
-0.02
-0.025
-0.03
-0.035
-4.0e-02

# Average properties

$\alpha$ and kinetic energy

Accretion rate

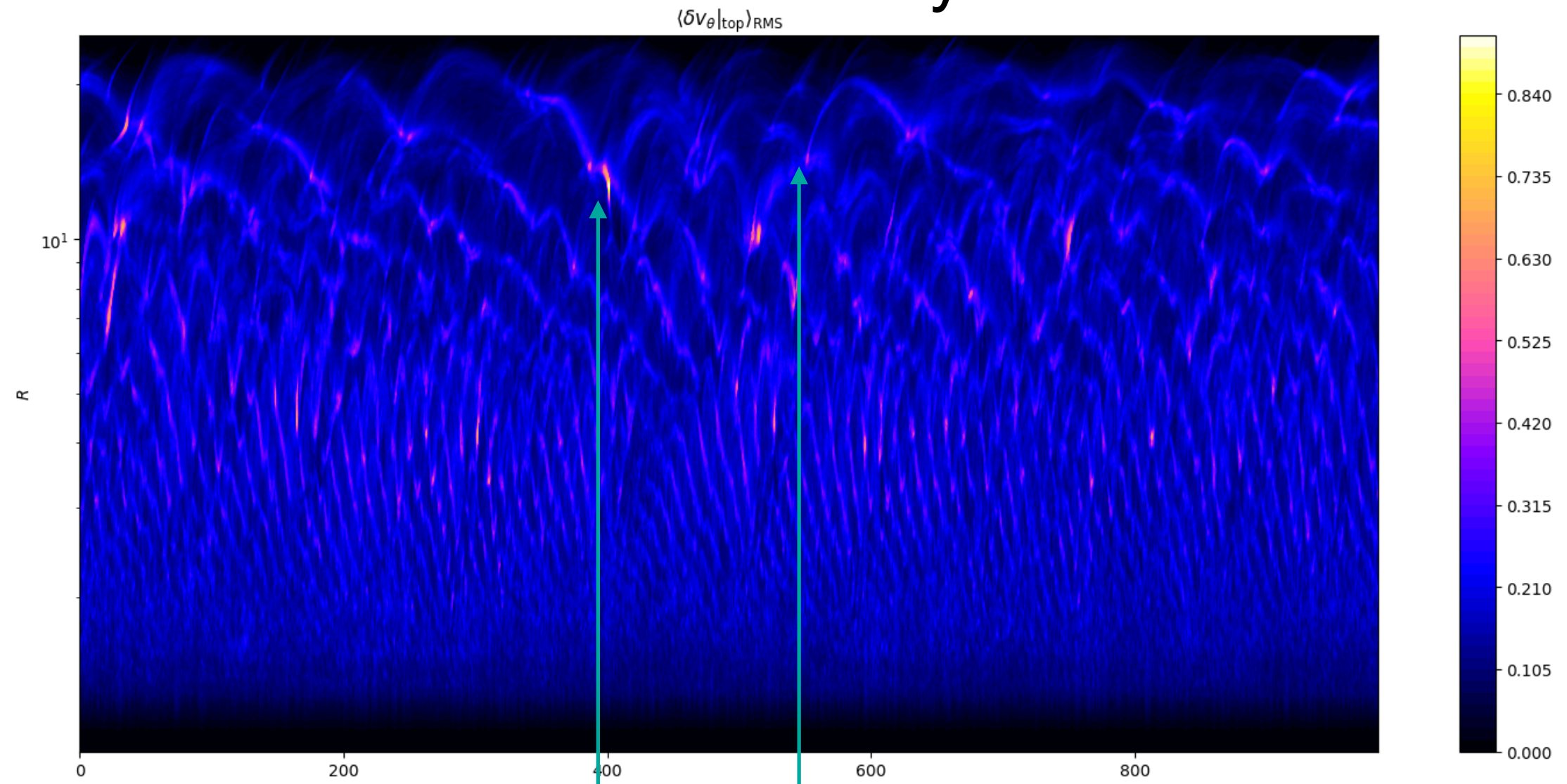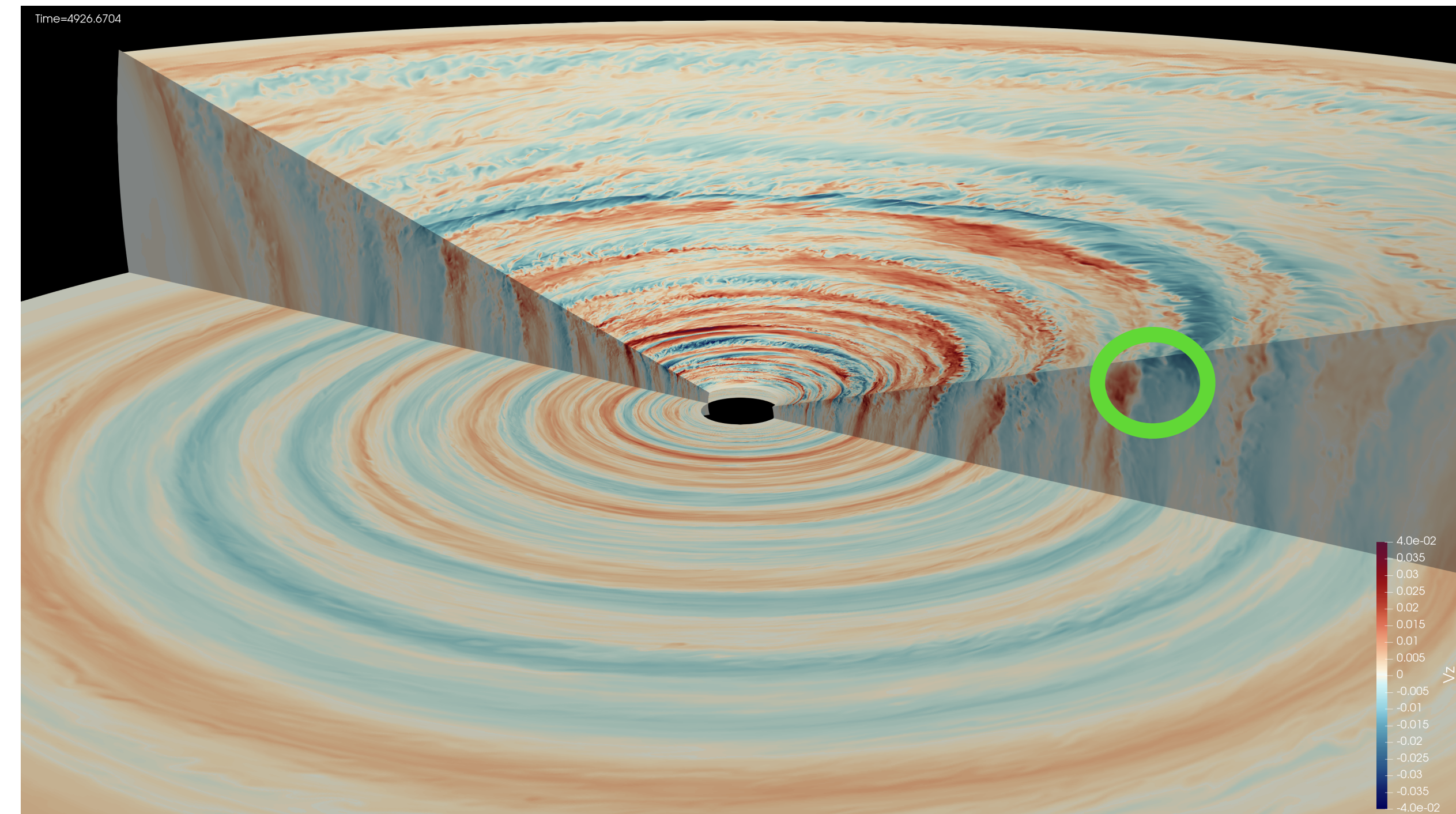# pattern

```
                                    linscale = 1,
                                    vmin =-0.1,
                                    vmax = 0.1),)
plt.yscale('log')
plt.xlabel('$t$ (orbits @ R=1)')
plt.ylabel('$R$')
plt.title(r'$\log_{10}[z_\rho/H]$')
plt.colorbar(c)
```

<matplotlib.colorbar.Colorbar at 0x152e0c7b1390>





<Figure size 640x480 with 0 Axes>

## Time averaged (over the last 100 orbits)
- We recover large-scale mean

```
fld.shape
```

(5563, 2576)

```
tavg=p.t[-1]-100*2*np.pi
idx=np.argwhere(p.t>tavg)[:,0]
```

```
Sigma=p.r[:,None]*p.rho
Sigma_avg=np.mean(Sigma[:,idx],axis=1)
Sigma_org=Sigma[:,0]
```

## Turbulence

# Turbulence

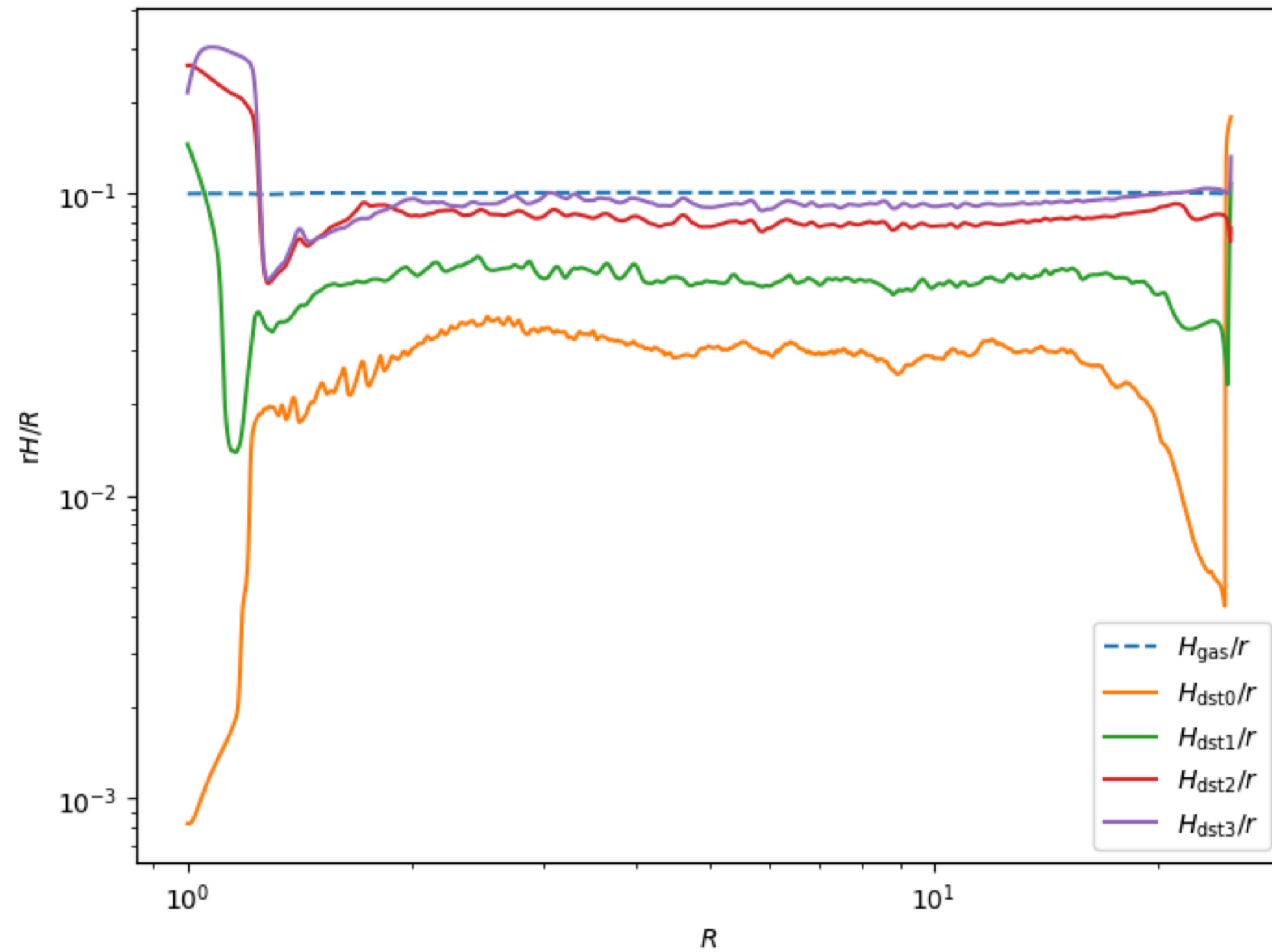## RMS turbulent velocity @ 3H



## Mean vertical velocity @3H

# Dust settling



Settling under-estimated in the innermost regions.

No radial dependence?

# Conclusion

- Idefix runs efficiently on pre-exascale European machines. However:

  - These machines are very sensitive to heating problems, leading to reduced performances and imbalanced computations

  - The failure rate is high, impacting significantly the runs when using more than 100 nodes (+1000GPUs)

  - File size is huge

- But eventually, it works :-)