



Oak Ridge User Conference

Technical Challenges for Ab-initio Simulations of Low and High-Pressure Turbines

R. Pichler, R. D. Sandberg, V. Michelassi and G. Laskowski

ACKNOWLEDGEMENTS

- ▶ ORNL Directors Grant - access to a GPU system
- ▶ John Levesque and Tom Edwards from CRAY COE - initial porting
- ▶ Jeff Larkin - Nvidia

OUTLINE

Code structure and algorithms

CPU performance and scaling

ACC porting

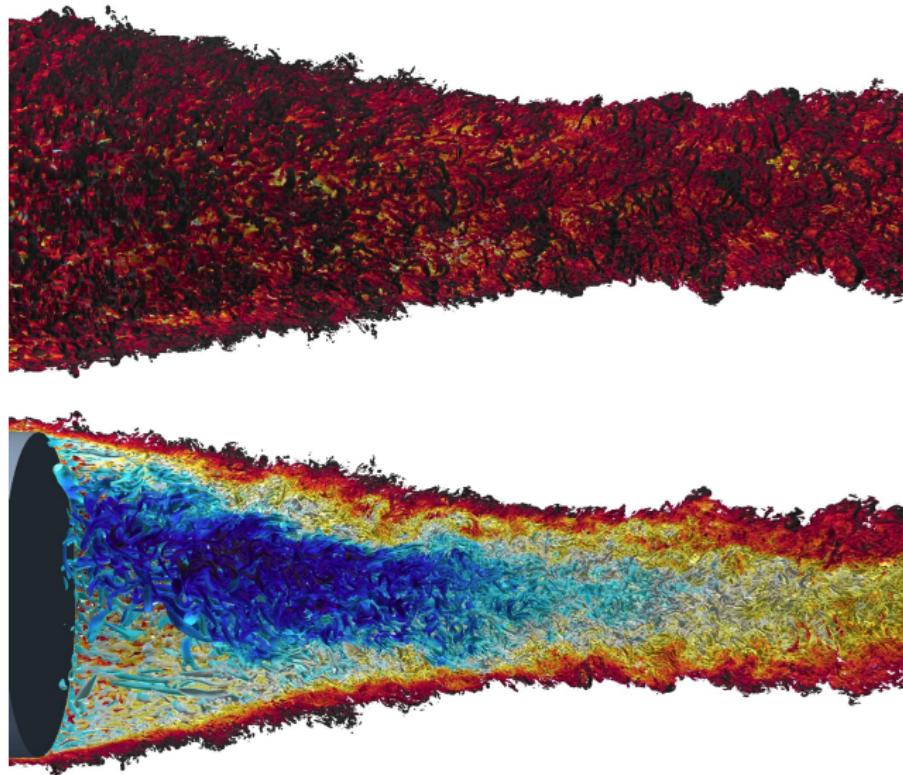
Porting of the Computationally heavy Part

Limiting data transfer

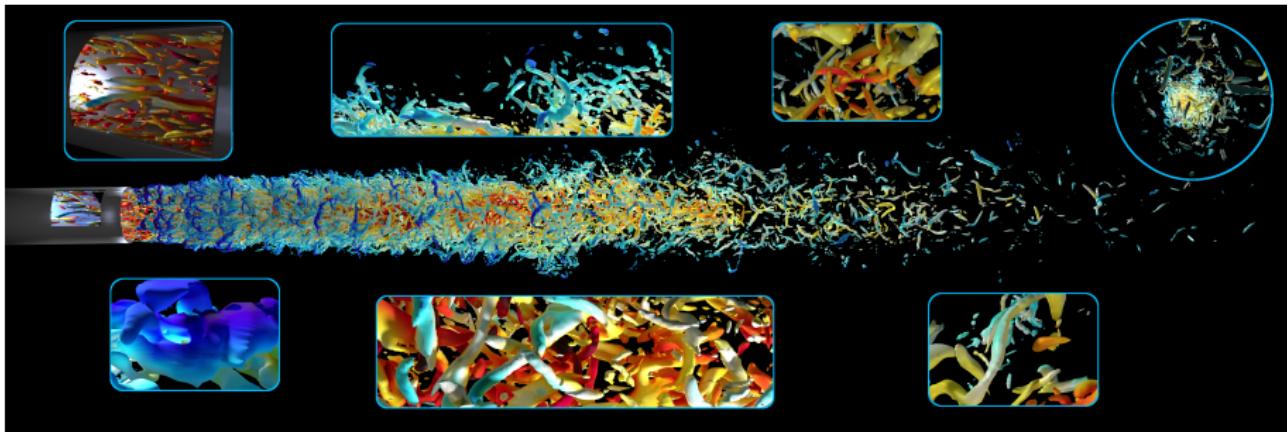
Additional important tweaks

Conclusions and Future

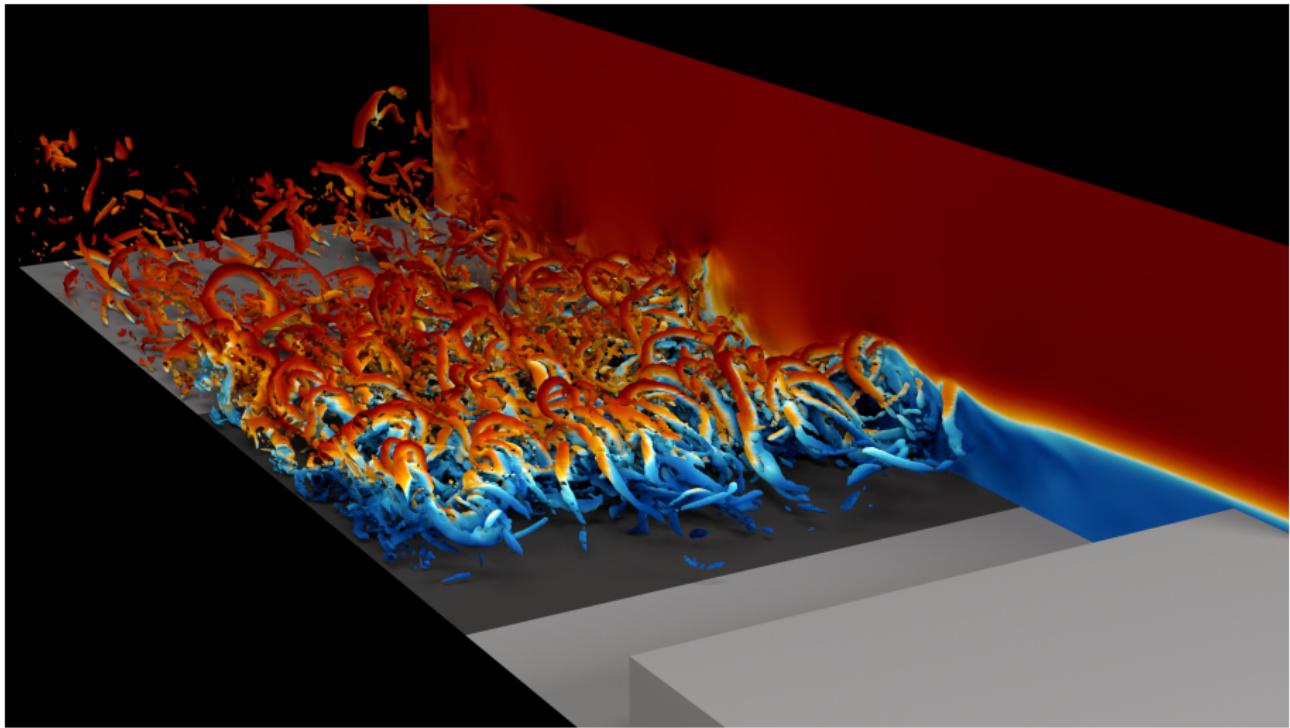
RESEARCH AREAS



RESEARCH AREAS



RESEARCH AREAS



HIPSTAR

High Performance Solver for Turbulence and Aeroacoustic Research

HIPSTAR

High Performance Solver for Turbulence and Aeroacoustic Research
Solves unsteady 3d compressible Flow equations - continuum

- ▶ Conservation of mass
- ▶ Conservation of momentum (3x)
- ▶ Conservation of energy

HIPSTAR

High Performance Solver for Turbulence and Aeroacoustic Research
Solves unsteady 3d compressible Flow equations - continuum

- ▶ Conservation of mass
- ▶ Conservation of momentum (3x)
- ▶ Conservation of energy
- ▶ Euler Frame of Reference
- ▶ neglect gravity

$$\frac{\partial}{\partial t} (\rho u_i) + \frac{\partial}{\partial x_j} (\rho u_j u_i) = -\frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} (t_{ji}) \quad (1)$$

$$\frac{\partial}{\partial t} (\rho \vec{u}) + \rho (\vec{u} \cdot \nabla) \vec{u} = -\nabla p + \nabla \cdot \underline{t} \quad (2)$$

NUMERICAL APPROACH

$$\frac{\partial}{\partial t} (\rho u_i) + \frac{\partial}{\partial x_j} (\rho u_j u_i) = -\frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} (t_{ji}) \quad (3)$$

rewrite as

$$\frac{\partial}{\partial t} (\rho u_i) = \underbrace{-\frac{\partial}{\partial x_j} (\rho u_j u_i) - \frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} (t_{ji})}_{\text{RHS}} \quad (4)$$

NUMERICAL APPROACH

$$\frac{\partial}{\partial t} (\rho u_i) + \frac{\partial}{\partial x_j} (\rho u_j u_i) = -\frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} (t_{ji}) \quad (3)$$

rewrite as

$$\frac{\partial}{\partial t} (\rho u_i) = \underbrace{-\frac{\partial}{\partial x_j} (\rho u_j u_i) - \frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} (t_{ji})}_{\text{RHS}} \quad (4)$$

integrate numerically w.r.t time using explicit Runge-Kutta method

$$\rho u_i|_{n+1} = \rho u_i|_n + \Delta t \sum_{k=1}^5 b_k \mathbf{RHS}|_n \quad (5)$$

NUMERICAL APPROACH

$$\frac{\partial}{\partial t} (\rho u_i) + \frac{\partial}{\partial x_j} (\rho u_j u_i) = -\frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} (t_{ji}) \quad (3)$$

rewrite as

$$\frac{\partial}{\partial t} (\rho u_i) = \underbrace{-\frac{\partial}{\partial x_j} (\rho u_j u_i) - \frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} (t_{ji})}_{\text{RHS}} \quad (4)$$

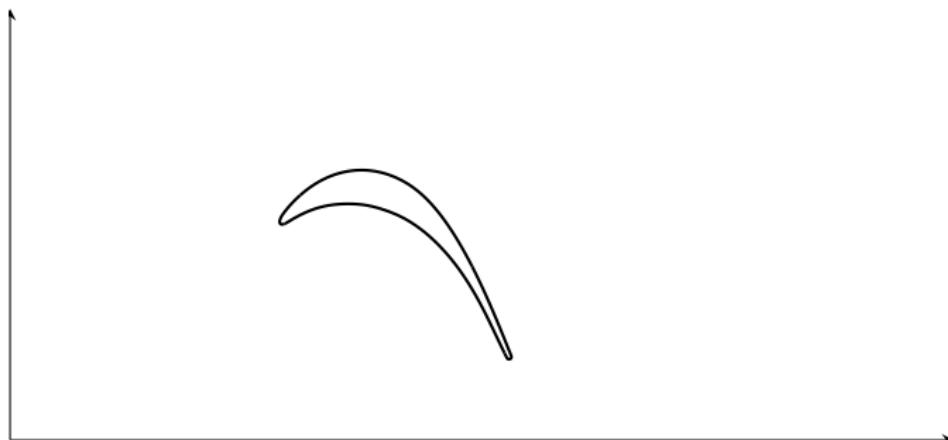
integrate numerically w.r.t time using explicit Runge-Kutta method

$$\rho u_i|_{n+1} = \rho u_i|_n + \Delta t \sum_{k=1}^5 b_k \mathbf{RHS}|_n \quad (5)$$

spatial derivative discretized by 4th order finite differences and Fourier method

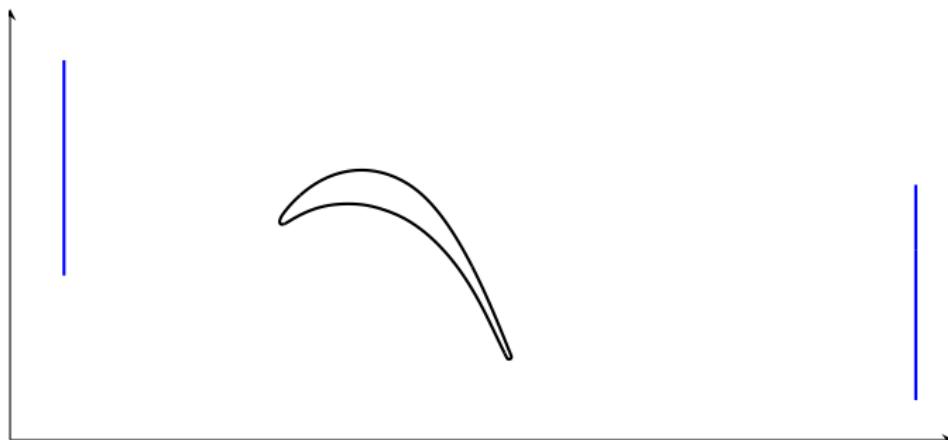
PARALLELIZATION

- discretize domain



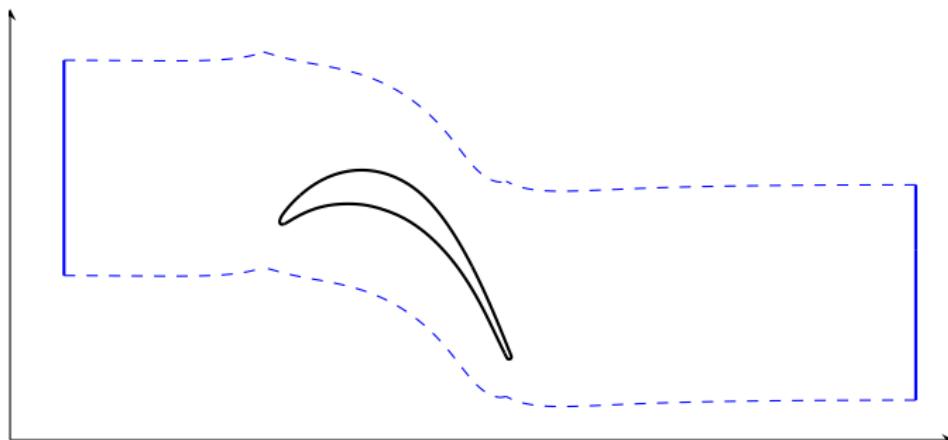
PARALLELIZATION

- discretize domain



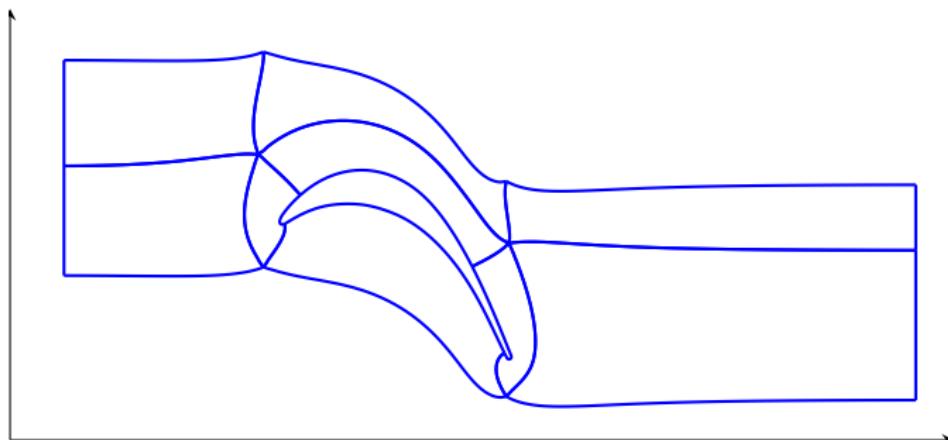
PARALLELIZATION

- discretize domain



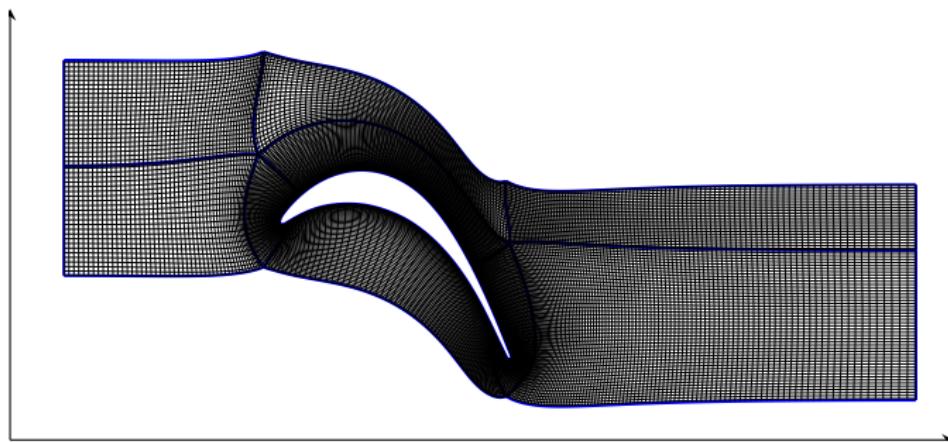
PARALLELIZATION

- ▶ discretize domain
- ▶ multi block grid for domain



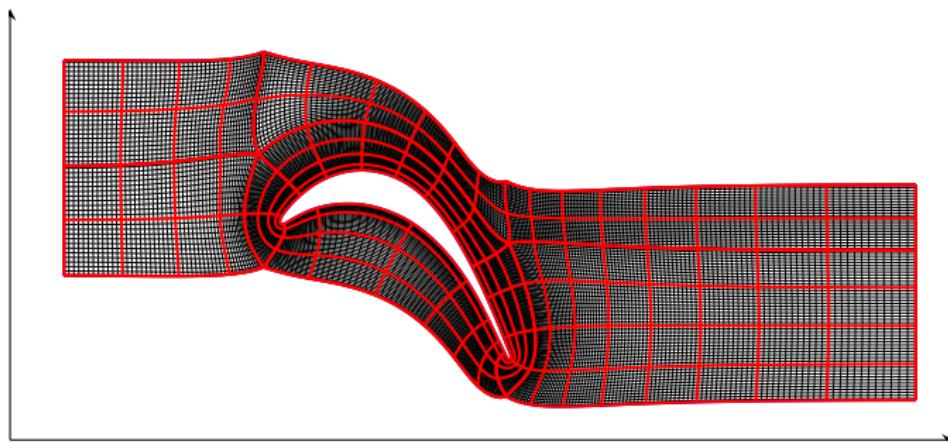
PARALLELIZATION

- ▶ discretize domain
- ▶ multi block grid for domain



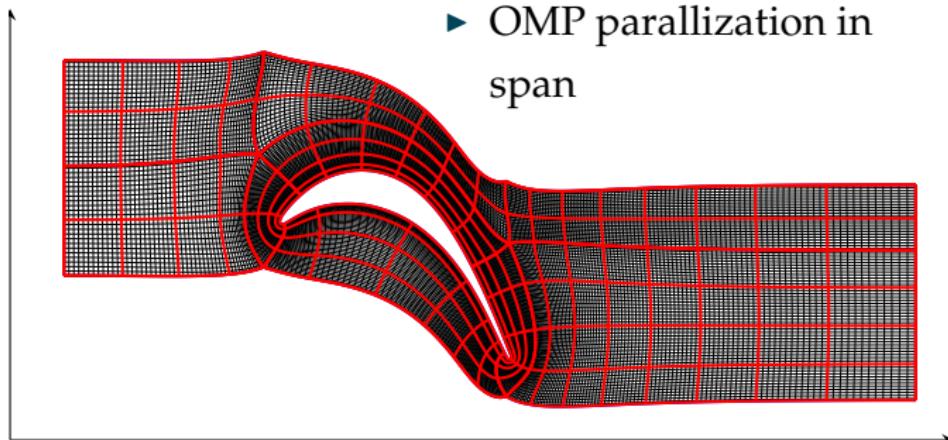
PARALLELIZATION

- ▶ discretize domain
- ▶ multi block grid for domain
- ▶ decompose grid using MPI

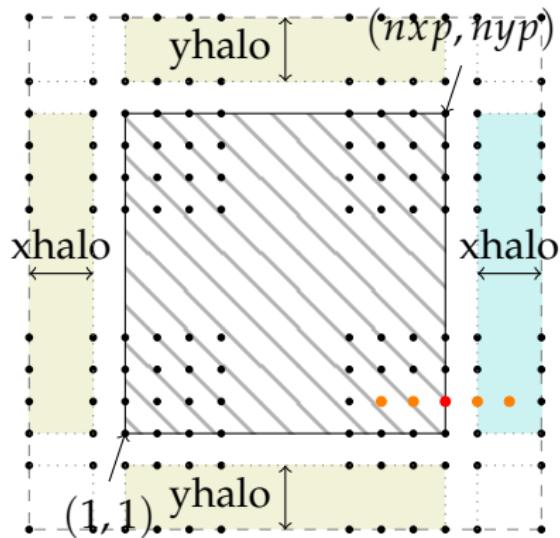


PARALLELIZATION

- ▶ discretize domain
- ▶ multi block grid for domain
- ▶ decompose grid using MPI
- ▶ OMP parallelization in span

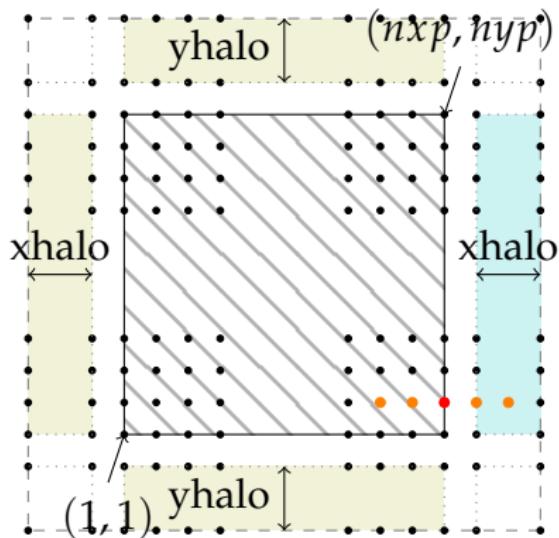


DOMAIN DECOMPOSITION



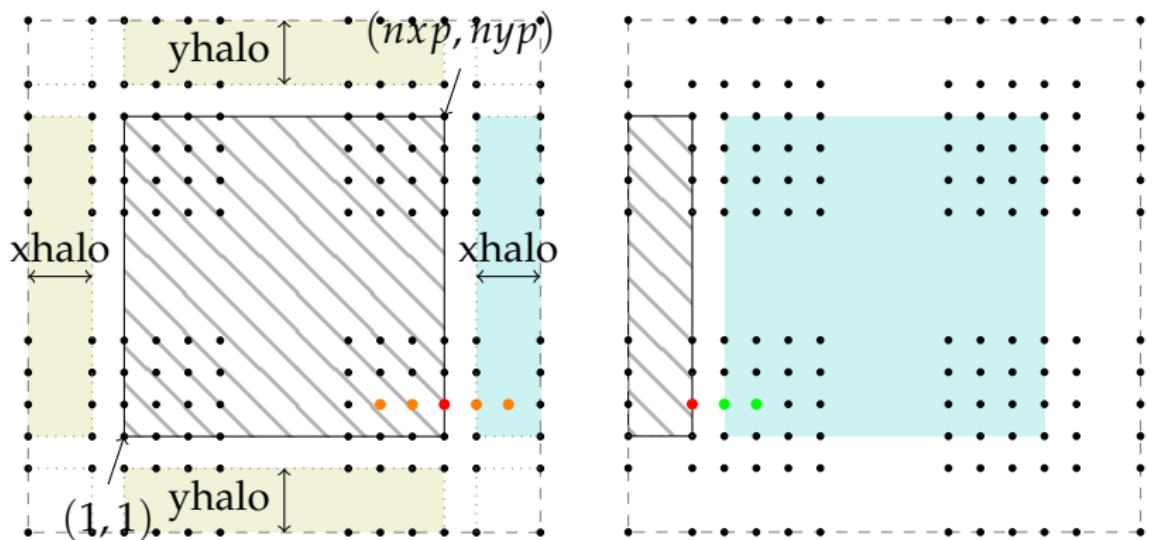
DOMAIN DECOMPOSITION

$$\frac{\partial}{\partial t} (\rho u_i) = \underbrace{-\frac{\partial}{\partial x_j} (\rho u_j u_i) - \frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} (t_{ji})}_{\text{RHS}}$$



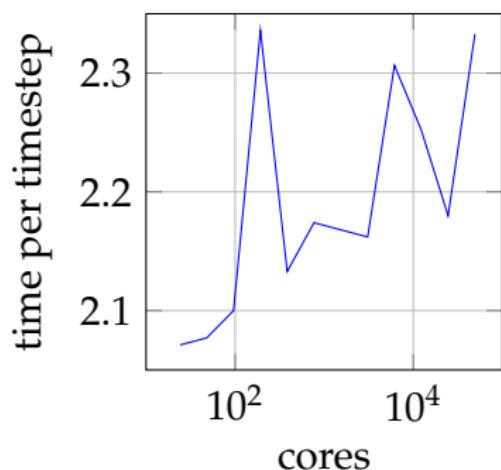
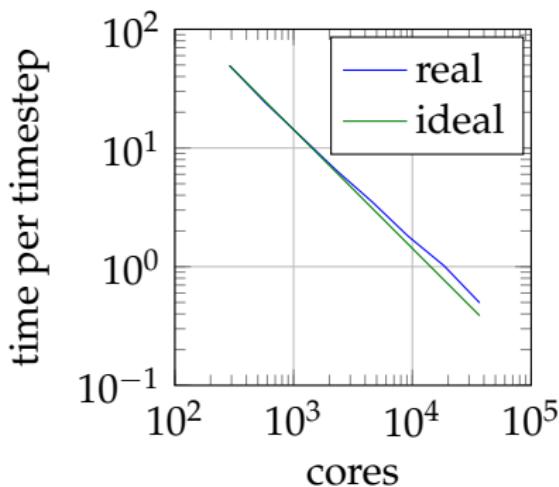
DOMAIN DECOMPOSITION

$$\frac{\partial}{\partial t} (\rho u_i) = \underbrace{-\frac{\partial}{\partial x_j} (\rho u_j u_i) - \frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} (t_{ji})}_{\text{RHS}}$$



CPU SCALING - PLANE JET

- Archer - UK national facility XC30
- strong scaling: speed up of 102 vs. 128
- weak scaling: 24 to 49152 cores



CODE STRUCTURE - CPU PROFILE

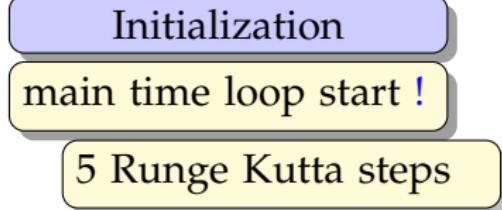
Initialization

main time loop start !

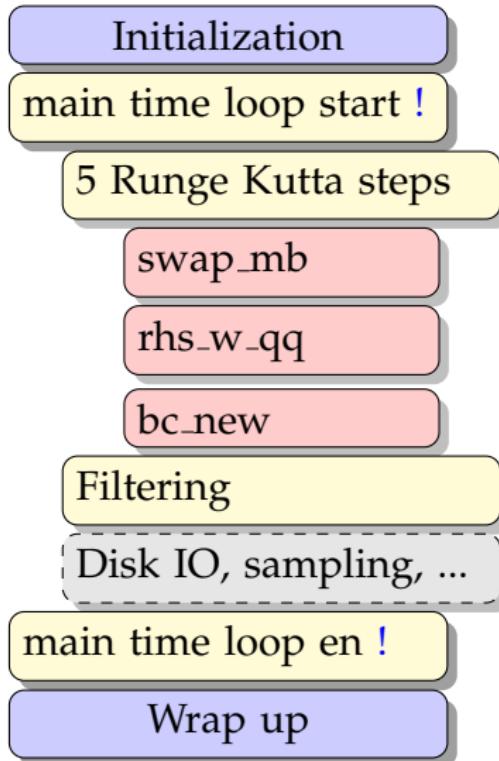
main time loop en !

Wrap up

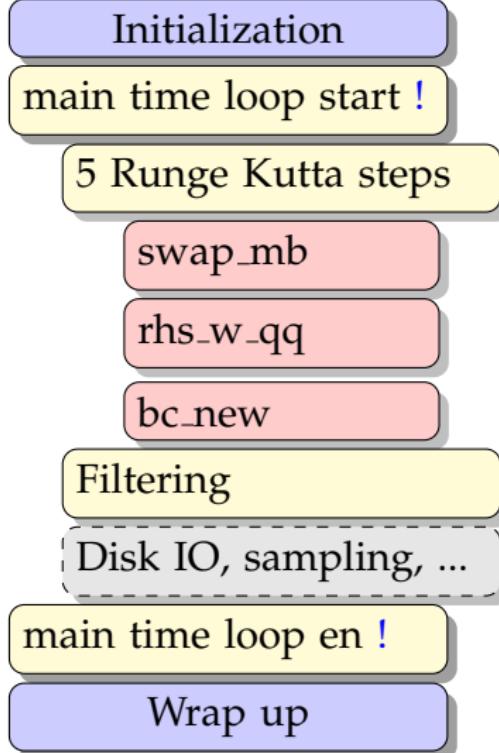
CODE STRUCTURE - CPU PROFILE



CODE STRUCTURE - CPU PROFILE

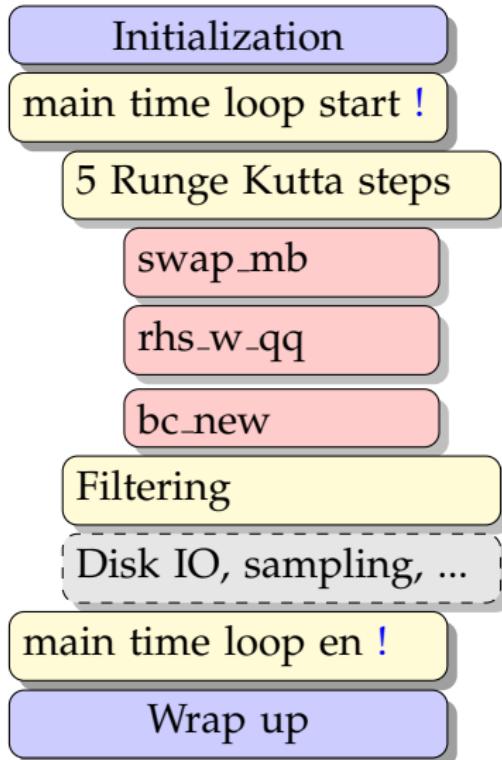


CODE STRUCTURE - CPU PROFILE

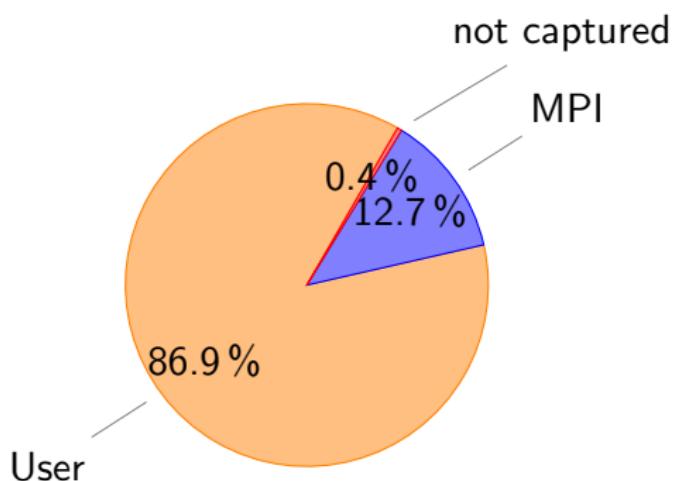


- ▶ Piz Daint - Xeon E5-2670 8C
2.600GHz
- ▶ CrayPAT - trace MPI and user functions in main loop

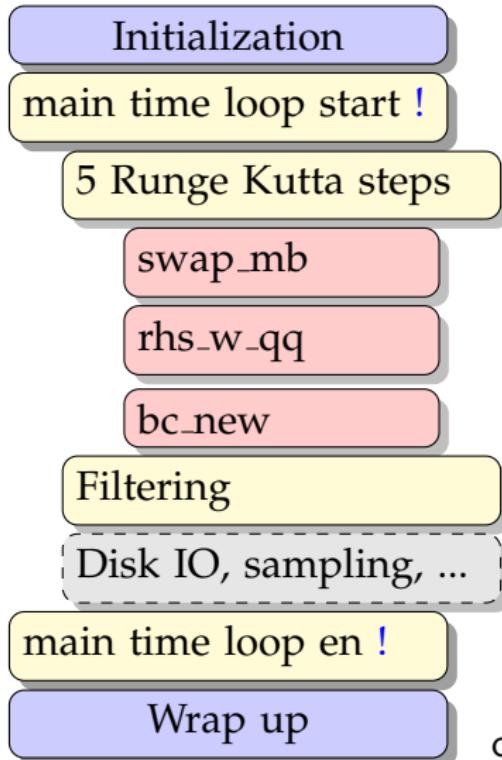
CODE STRUCTURE - CPU PROFILE



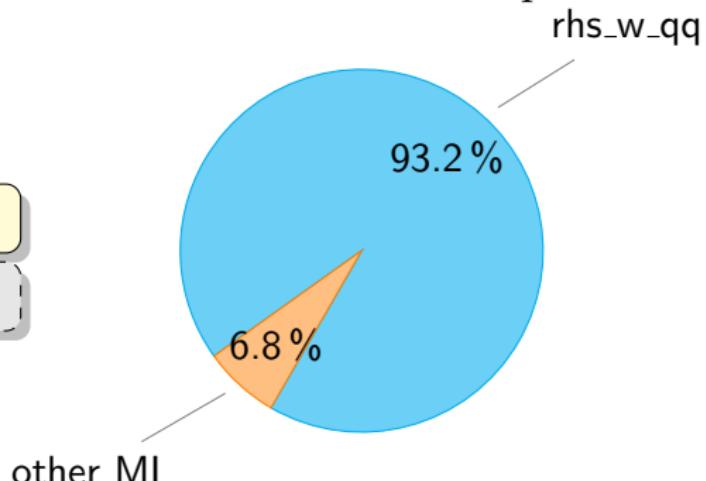
- ▶ Piz Daint - Xeon E5-2670 8C 2.600GHz
- ▶ CrayPAT - trace MPI and user functions in main loop



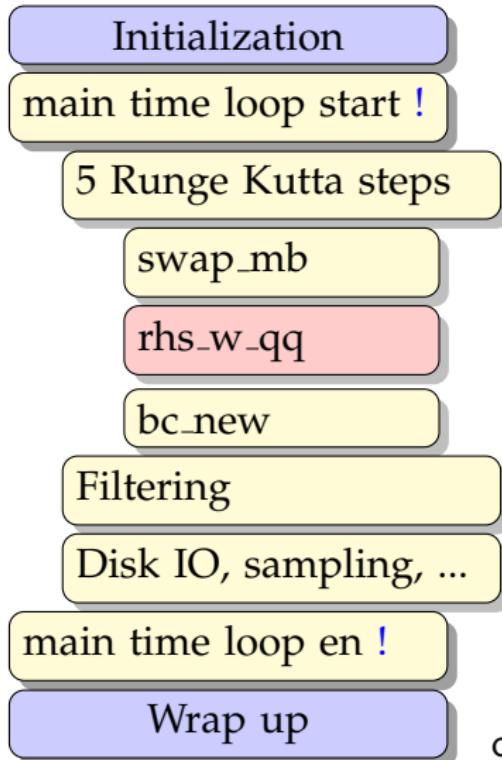
CODE STRUCTURE - CPU PROFILE



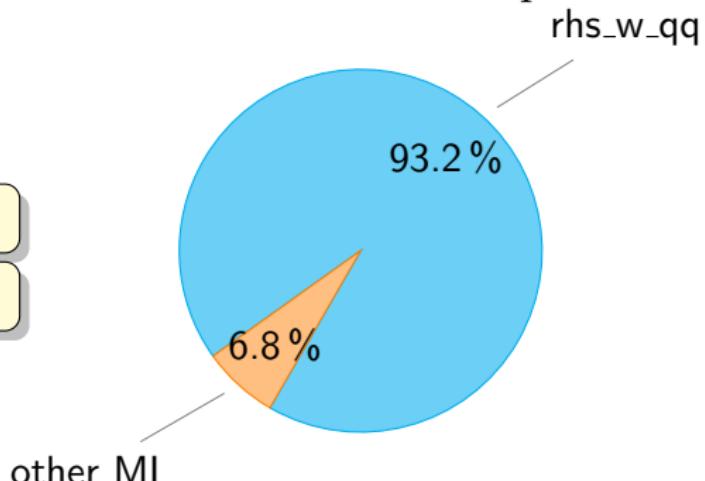
- ▶ Piz Daint - Xeon E5-2670 8C 2.600GHz
- ▶ CrayPAT - trace MPI and user functions in main loop



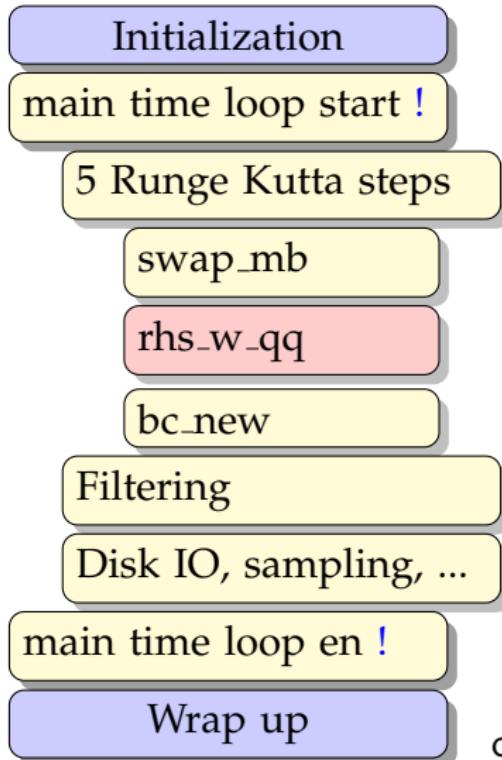
CODE STRUCTURE - CPU PROFILE



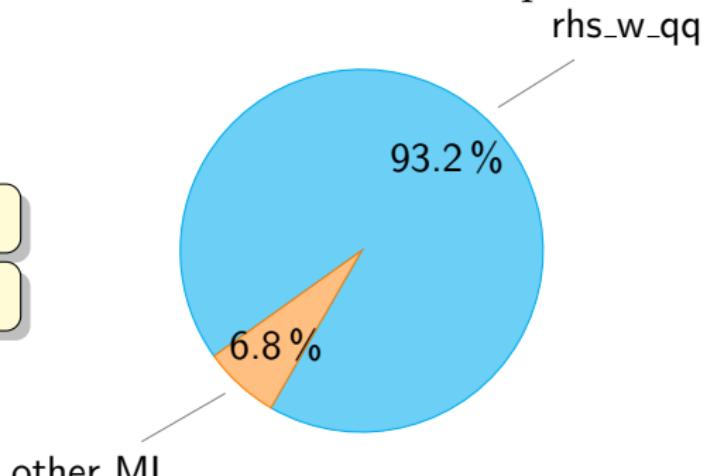
- ▶ Piz Daint - Xeon E5-2670 8C 2.600GHz
- ▶ CrayPAT - trace MPI and user functions in main loop



CODE STRUCTURE - CPU PROFILE



- ▶ Piz Daint - Xeon E5-2670 8C 2.600GHz
- ▶ CrayPAT - trace MPI and user functions in main loop



1. step: port rhs_w_qq ⇒ data transfer

STEP 1 - COMPUTATIONALLY HEAVY PART

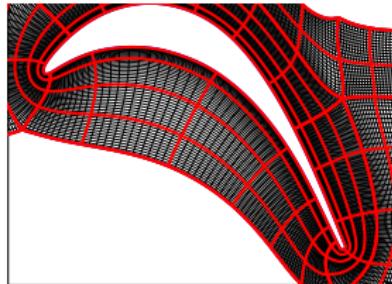
- ▶ “Only” 30% of the code have to be ported
- ▶ done by John Levesque and Tom Edwards (Cray)

STEP 1 - COMPUTATIONALLY HEAVY PART

- ▶ “Only” 30% of the code have to be ported
- ▶ done by John Levesque and Tom Edwards (Cray)
- ▶ add ACC statements to loops
- ▶ define data regions ⇒ CPU-GPU data transfer

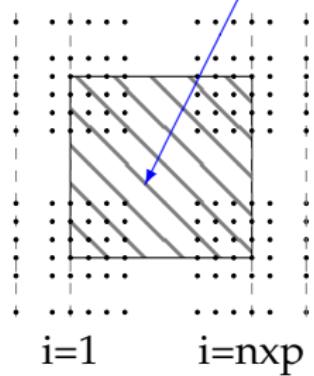
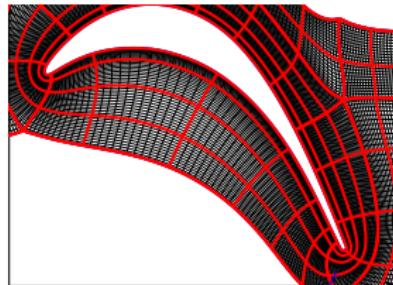
STEP 1 - COMPUTATIONALLY HEAVY PART

$$\mathbf{RHS} = -\frac{\partial}{\partial x_j} (\rho u_j u_i) - \frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} (t_{ji})$$



STEP 1 - COMPUTATIONALLY HEAVY PART

$$\mathbf{RHS} = -\frac{\partial}{\partial x_j} (\rho u_j u_i) - \frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} (t_{ji})$$



STEP 1 - COMPUTATIONALLY HEAVY PART

```
do k=nzp_start ,nzp_end
```

$$\mathbf{RHS} = -\frac{\partial}{\partial x_j} (\rho u_j u_i) - \frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} (t_{ji})$$

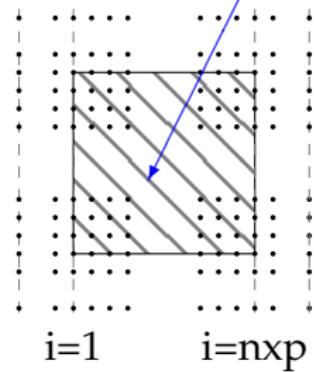
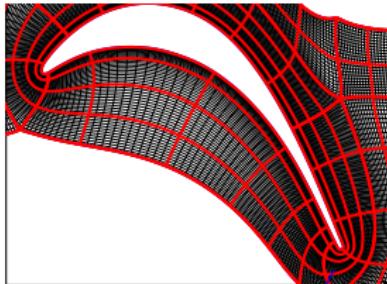
```
do j=jbegwp ,jendwp
```

```
  do i=ibegwp ,iendwp  
    ddw4(i,j,1) = qph(i,j,k,2)*r_et(i,j)  
  end do  
end do
```

```
call ddx1(ddw4(1-xhalo,1-yhalo,1)  
&           ,ddw1(1-xhalo,1-yhalo,1) ,1)
```

```
do j=1,nyp
```

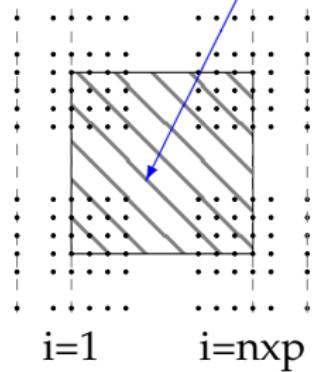
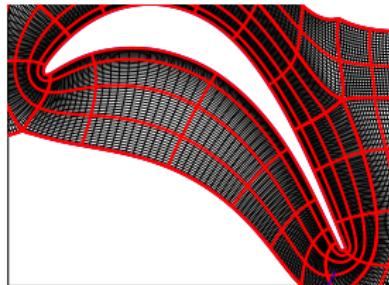
```
  do i=1,nxp  
    q(i,j,k,2)= q(i,j,k,2)+ddw1(i,j,1)  
  end do  
end do  
end do
```



STEP 1 - COMPUTATIONALLY HEAVY PART

$$\mathbf{RHS} = -\frac{\partial}{\partial x_j} (\rho u_j u_i) - \frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} (t_{ji})$$

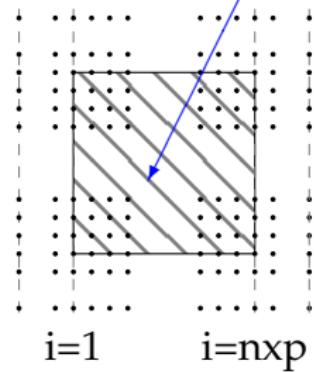
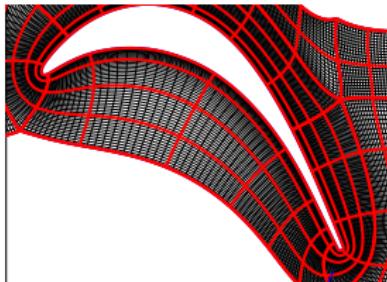
```
do k=nzp_start ,nzp_end  
  
  do j=jbegwp ,jendwp  
  
    do i=ibegwp ,iendwp  
      ddw4(i,j,k) = qph(i,j,k,2)*r_et(i,j)  
    end do  
  end do  
  end do  
  
  call ddx1_all(ddw4(1-xhalo,1-yhalo,1)  
&           ,ddw1(1-xhalo,1-yhalo,1) ,nzp)  
  
  do k=nzp_start ,nzp_end  
  
    do j=1,nyp  
  
      do i=1,nxp  
        q(i,j,k,2)= q(i,j,k,2)+ddw1(i,j,k)  
      end do  
    end do  
  end do
```



STEP 1 - COMPUTATIONALLY HEAVY PART

```
!$acc parallel loop gang
    do k=nzp_start ,nzp.end
!$acc loop worker
    do j=jbegwp ,jendwp
!$acc loop vector
    do i=ibegwp ,iendwp
        ddw4(i,j,k) = qph(i,j,k,2)*r_et(i,j)
    end do
    end do
!$acc end parallel loop
    call ddx1i_all(ddw4(1-xhalo,1-yhalo,1)
&           ,ddw1(1-xhalo,1-yhalo,1) ,nzp)
!$acc parallel loop gang
    do k=nzp_start ,nzp.end
!$acc loop worker
    do j=1,nyp
!$acc loop vector
    do i=1,nxp
        q(i,j,k,2)= q(i,j,k,2)+ddw1(i,j,k)
    end do
    end do
    end do
!$acc end parallel loop
```

$$\mathbf{RHS} = -\frac{\partial}{\partial x_j} (\rho u_j u_i) - \frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} (t_{ji})$$



STEP 1 - COMPUTATIONALLY HEAVY PART

- ▶ compare runtime of
 - ▶ hybrid open MPI/OMP with 1 MPI rank per node (16 OMP threads)
 - ▶ ACC version using 1 MPI rank per node

STEP 1 - COMPUTATIONALLY HEAVY PART

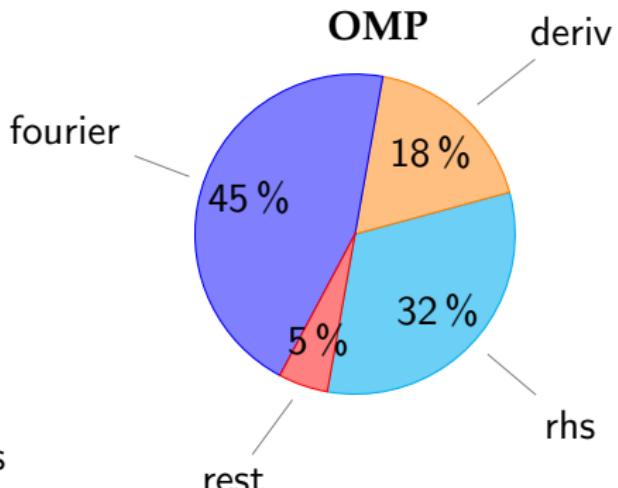
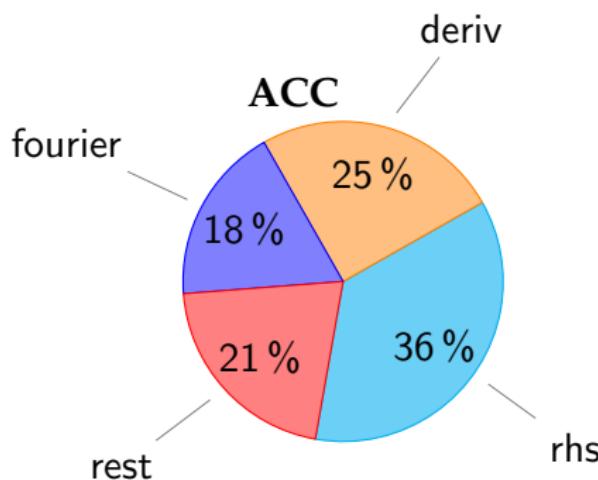
- ▶ compare runtime of
 - ▶ hybrid open MPI/OMP with 1 MPI rank per node (16 OMP threads)
 - ▶ ACC version using 1 MPI rank per node
- ▶ overall speedup of 18% (GPU vs. CPU)

STEP 1 - COMPUTATIONALLY HEAVY PART

- ▶ compare runtime of
 - ▶ hybrid open MPI/OMP with 1 MPI rank per node (16 OMP threads)
 - ▶ ACC version using 1 MPI rank per node
- ▶ overall speedup of 18% (GPU vs. CPU)
- ▶ data transfer between GPU and CPU \approx 30% of runtime

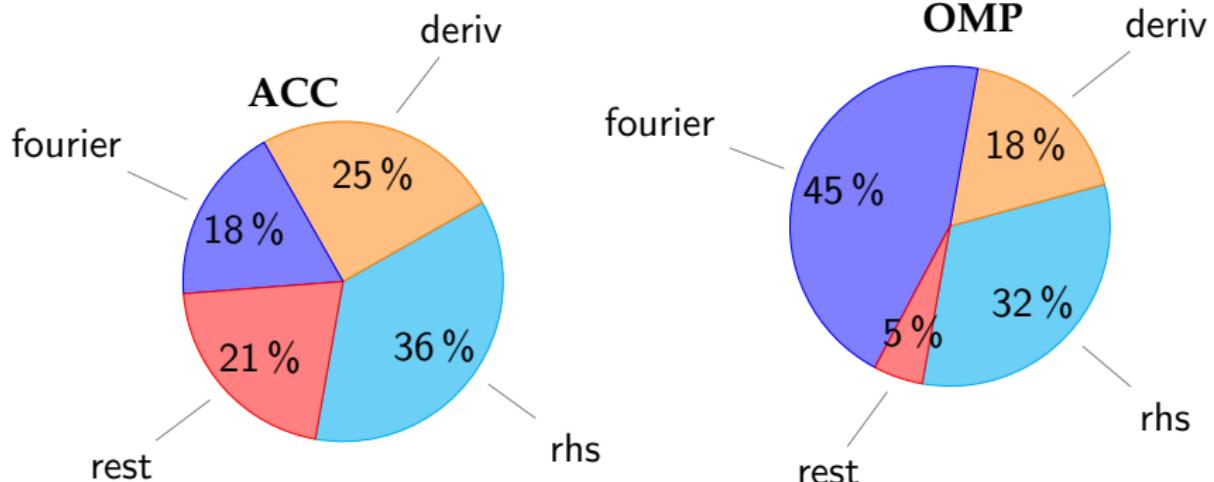
STEP 1 - COMPUTATIONALLY HEAVY PART

- *rest* ⇒ CPU code runs slower (1 vs 16 OMP threads)



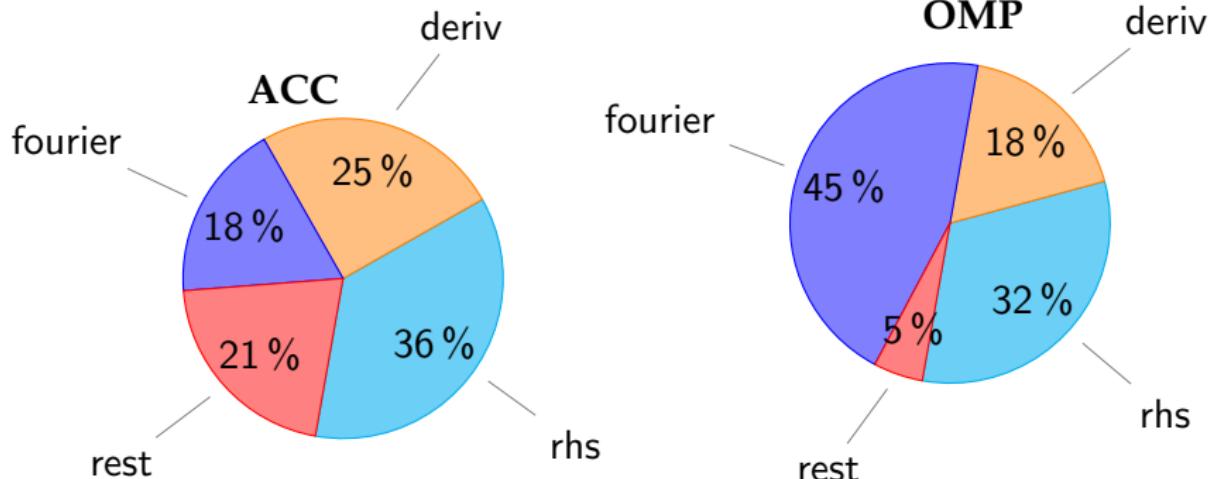
STEP 1 - COMPUTATIONALLY HEAVY PART

- *rest* ⇒ CPU code runs slower (1 vs 16 OMP threads)
- FFTW (cufft) are significantly faster

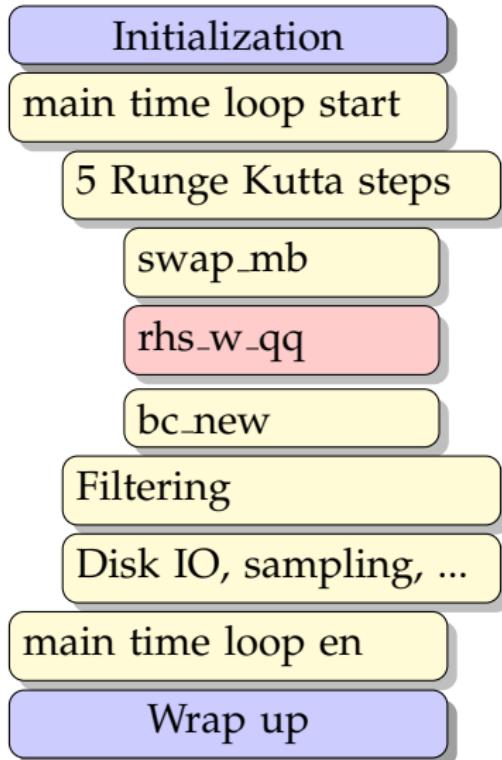


STEP 1 - COMPUTATIONALLY HEAVY PART

- *rest* ⇒ CPU code runs slower (1 vs 16 OMP threads)
- FFTW (cufft) are significantly faster
- *deriv* and *rhs* similar relative fraction

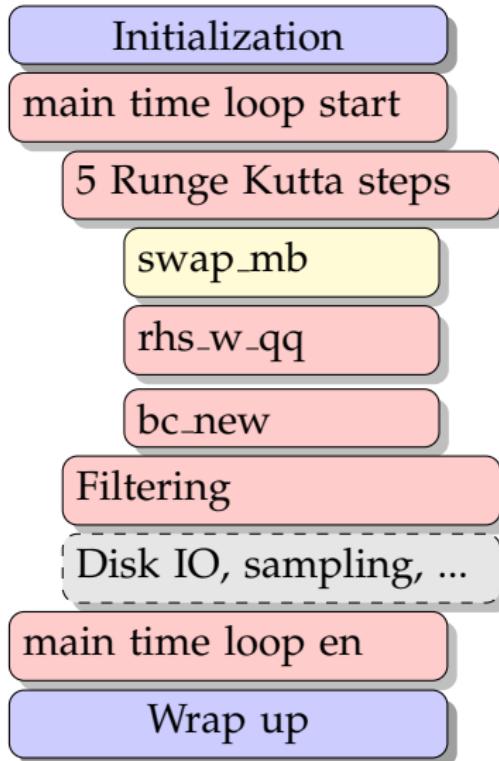


STEP 2 - LIMITTING DATA TRANSFER



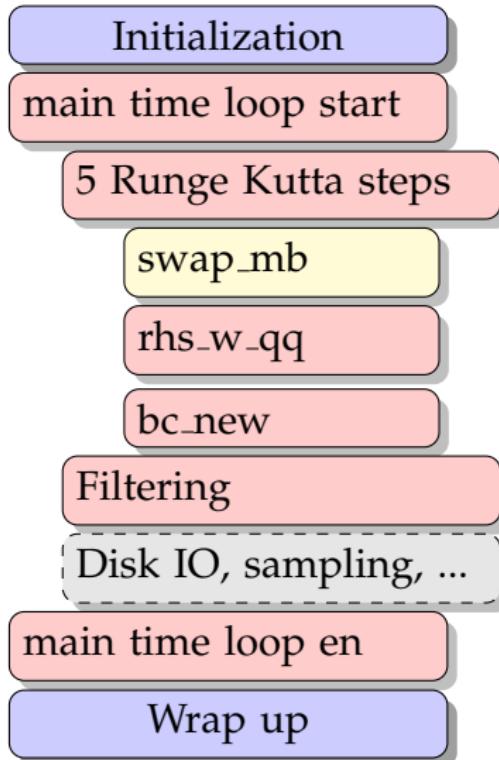
- ▶ 30% of walltime spent in CPU-GPU data transfer
- ▶ move everything to GPU except for

STEP 2 - LIMITTING DATA TRANSFER



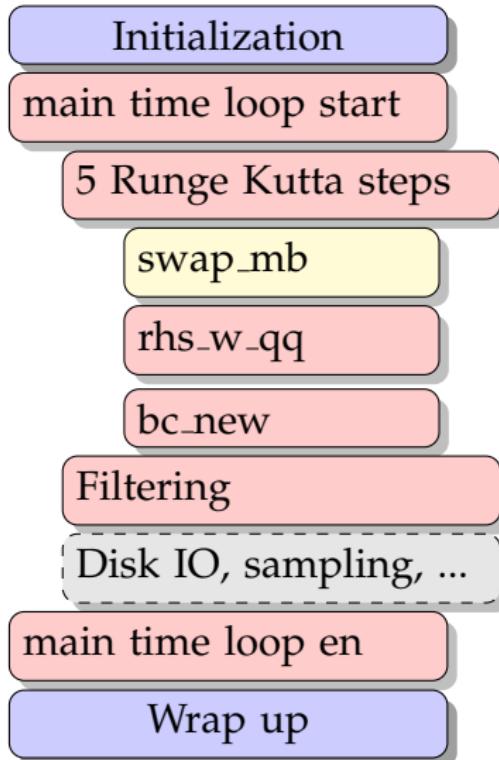
- ▶ 30% of walltime spent in CPU-GPU data transfer
- ▶ move everything to GPU except for
 - ▶ swapping - MPI communication

STEP 2 - LIMITTING DATA TRANSFER



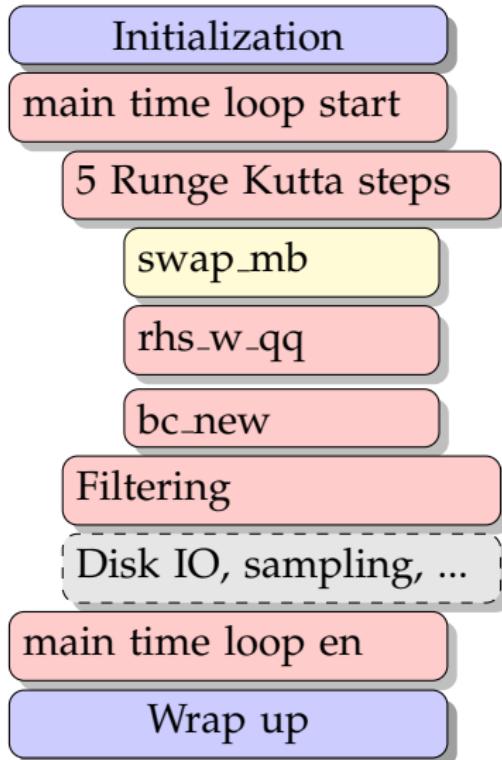
- ▶ 30% of walltime spent in CPU-GPU data transfer
- ▶ move everything to GPU except for
 - ▶ swapping - MPI communication
 - ▶ Disk IO - MPIIO

STEP 2 - LIMITTING DATA TRANSFER



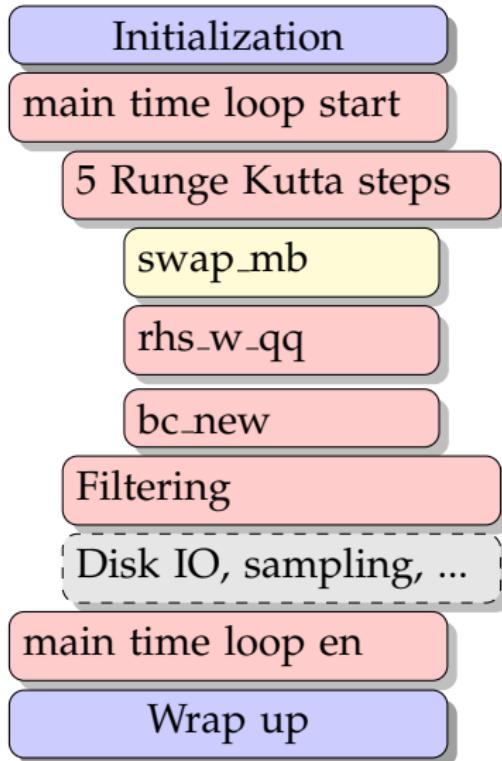
- ▶ 30% of walltime spent in CPU-GPU data transfer
- ▶ move everything to GPU except for
 - ▶ swapping - MPI communication
 - ▶ Disk IO - MPIIO
 - ▶ Initialization, Wrap up

STEP 2 - LIMITTING DATA TRANSFER



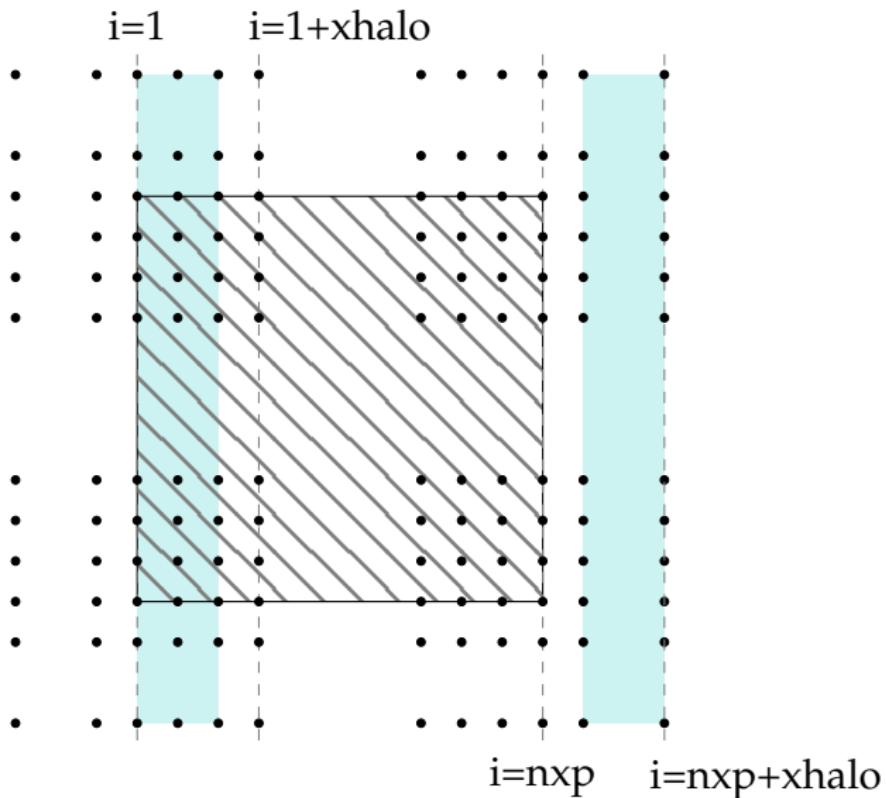
- ▶ 30% of walltime spent in CPU-GPU data transfer
- ▶ move everything to GPU except for
 - ▶ swapping - MPI communication
 - ▶ Disk IO - MPIIO
 - ▶ Initialization, Wrap up
- ▶ only halos are transferred between GPU and CPU
- ▶ about 70% of code have to be ported

STEP 2 - LIMITTING DATA TRANSFER



- ▶ 30% of walltime spent in CPU-GPU data transfer
- ▶ move everything to GPU except for
 - ▶ swapping - MPI communication
 - ▶ Disk IO - MPIIO
 - ▶ Initialization, Wrap up
- ▶ only halos are transferred between GPU and CPU
- ▶ about 70% of code have to be ported
- ▶ **runtime increase by 7x**

STEP 2 - LIMITING DATA TRANSFER



STEP 2 - LIMITING DATA TRANSFER

```
!$acc parallel loop private(1)
do n=1,nv
  do k=1,kpoints
    do j=1-yhalo ,nyp+yhalo
      do i=1,xhalo
        1 = i + (j+yhalo-1)*xhalo
&          + (k-1)*(nyp+2*yhalo)*xhalo
&          + (n-1)*(kpoints)*(nyp+2*yhalo)*xhalo
        ahoxm(1) = a(i,j,k,n)
      enddo
    enddo
  enddo
enddo
```

- ▶ base version: > 1600s

STEP 2 - LIMITING DATA TRANSFER

```
!$acc parallel loop private(1) collapse(2)
    do n=1,nv
        do k=1,kpoints
!$acc loop vector
        do j=1-yhalo ,nyp+yhalo
!dir$ unroll
        do i=1,xhalo
            1 = i + (j+yhalo-1)*xhalo
&            + (k-1)*(nyp+2*yhalo)*xhalo
&            + (n-1)*(kpoints)*(nyp+2*yhalo)*xhalo
            ahoxm(1) = a(i,j,k,n)
        enddo
        enddo
        enddo
    enddo
```

- ▶ base version: > 1600s
- ▶ final version: $\approx 12s$

STEP 2 - LIMITING DATA TRANSFER

```
!$acc parallel loop private(1) collapse(2)
    do n=1,nv
        do k=1,kpoints
!$acc loop vector
        do j=1-yhalo ,nyp+yhalo
!dir$ unroll
        do i=1,xhalo
            1 = i + (j+yhalo-1)*xhalo
&            + (k-1)*(nyp+2*yhalo)*xhalo
&            + (n-1)*(kpoints)*(nyp+2*yhalo)*xhalo
            ahoxm(1) = a(i,j,k,n)
        enddo
        enddo
        enddo
    enddo
```

- ▶ base version: > 1600s
- ▶ final version: $\approx 12s$
- ▶ $\approx 30 - 40\%$ speed up w.r.t. CPU version

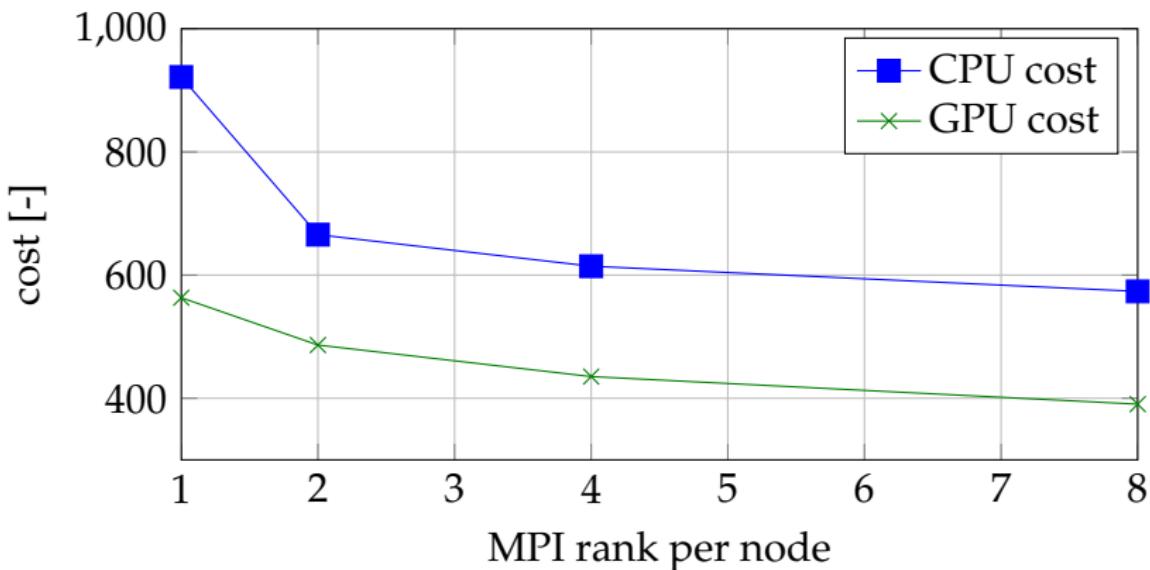
STEP 2 - LIMITING DATA TRANSFER

```
!$acc parallel loop private(1) collapse(2)
    do n=1,nv
        do k=1,kpoints
!$acc loop vector
        do j=1-yhalo ,nyp+yhalo
!dir$ unroll
        do i=1,xhalo
            1 = i + (j+yhalo-1)*xhalo
&            + (k-1)*(nyp+2*yhalo)*xhalo
&            + (n-1)*(kpoints)*(nyp+2*yhalo)*xhalo
            ahoxm(1) = a(i,j,k,n)
        enddo
        enddo
        enddo
    enddo
```

- ▶ base version: > 1600s
- ▶ final version: $\approx 12s$
- ▶ $\approx 30 - 40\%$ speed up w.r.t. CPU version
- ▶ packing on CPU fastest if using $l=l+1$
- ▶ GPU faster as CPU

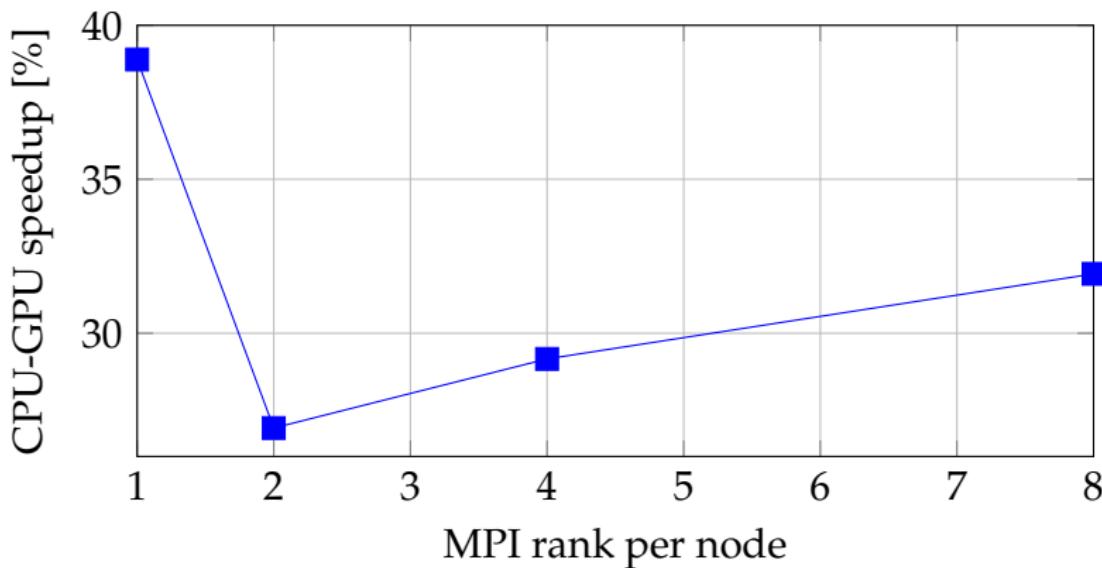
STEP 2 - SPEED-UP

- ▶ strong scaling
- ▶ Using 64 - 512 nodes



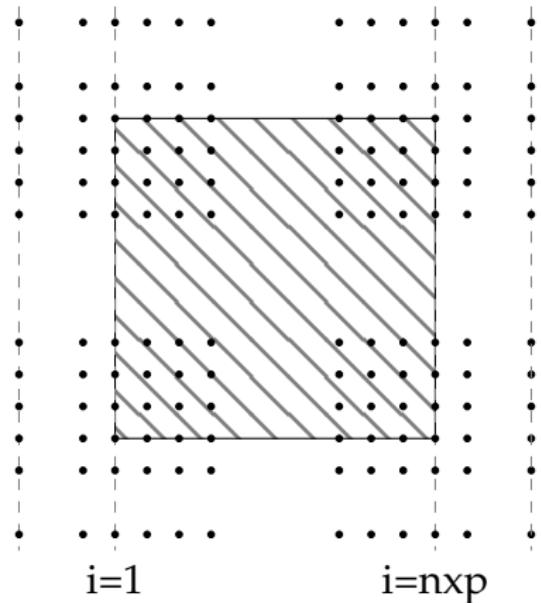
STEP 2 - SPEED-UP

- ▶ strong scaling
- ▶ Using 64 - 512 nodes
- ▶ speed up around 30%



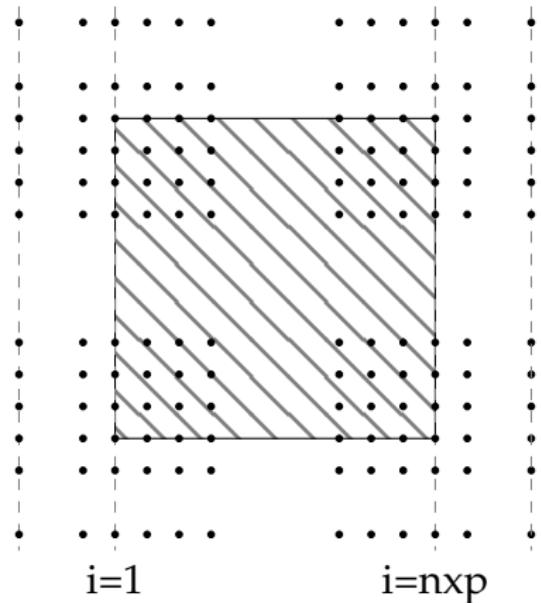
STEP 3 - ADDITIONAL TWEAKS

- ▶ mostly innermost loop length:
 $nxp + 2 \times xhalo$
- ▶ generate grid such that
innermost loop length is 32,64
or 128



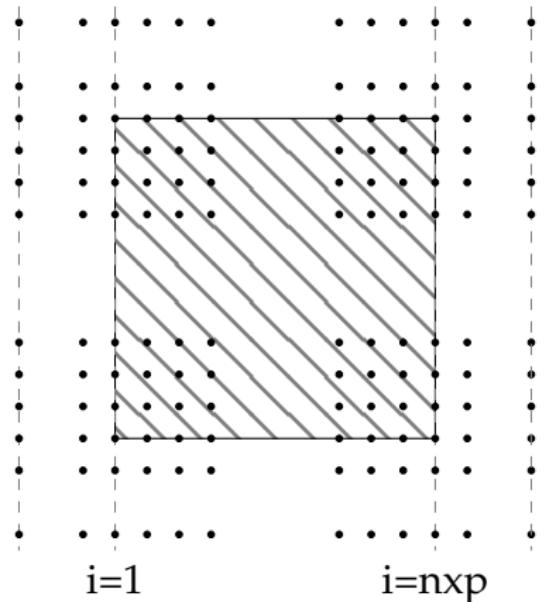
STEP 3 - ADDITIONAL TWEAKS

- ▶ mostly innermost loop length:
 $nxp + 2 \times xhalo$
- ▶ generate grid such that
innermost loop length is 32,64
or 128
- ▶ set `vector_length()` accordingly
- ▶ **set `nxp,nyp,nzp` as parameter**



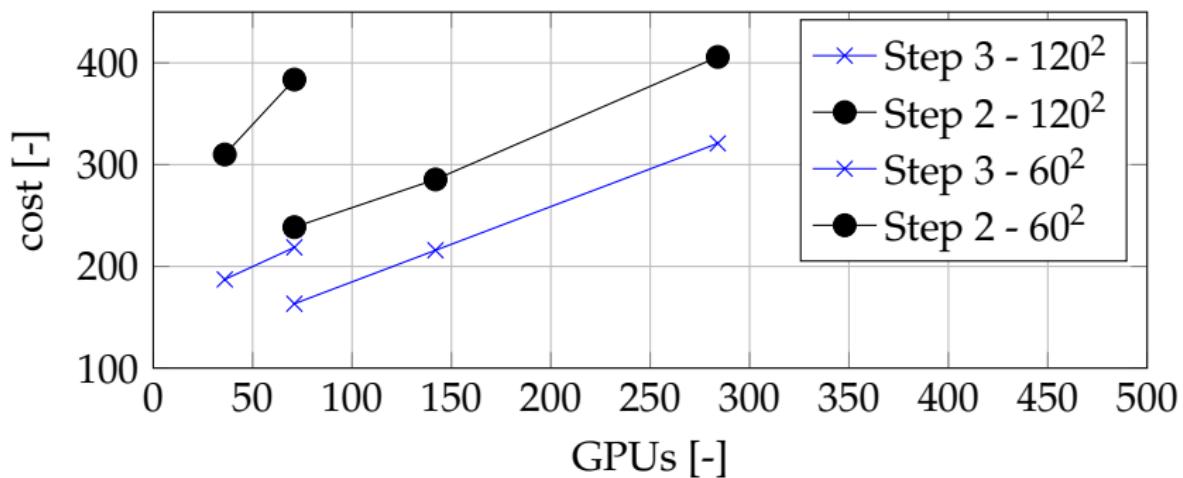
STEP 3 - ADDITIONAL TWEAKS

- ▶ mostly innermost loop length:
 $nxp + 2 \times xhalo$
- ▶ generate grid such that
innermost loop length is 32,64
or 128
- ▶ set `vector_length()` accordingly
- ▶ **set `nxp,nyp,nzp` as parameter**
- ▶ **less flexible but significant
performance improvement**



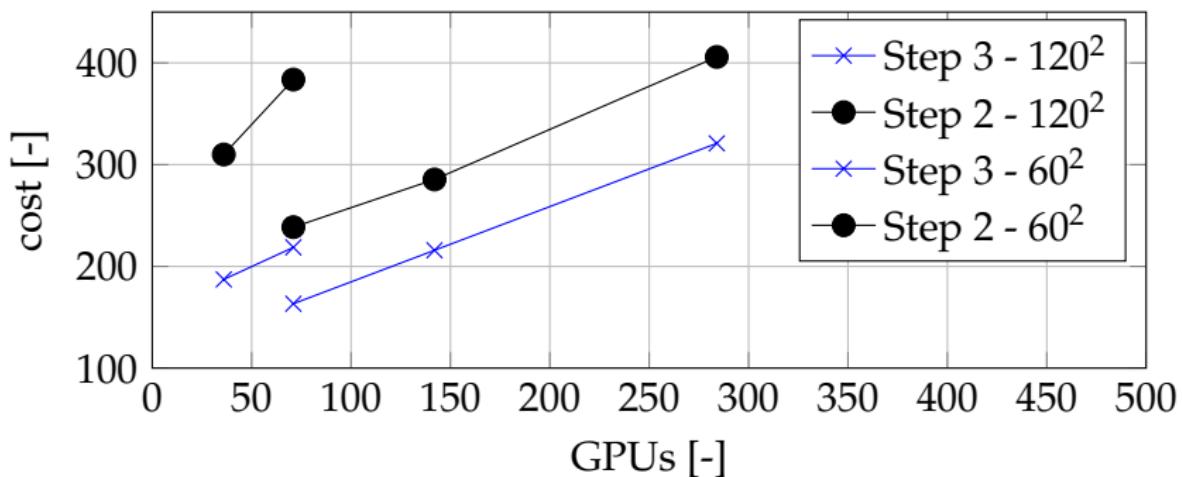
STEP 3 - SPEED UP

- ▶ strong scaling
- ▶ almost 2x speed w.r.t. Step 2



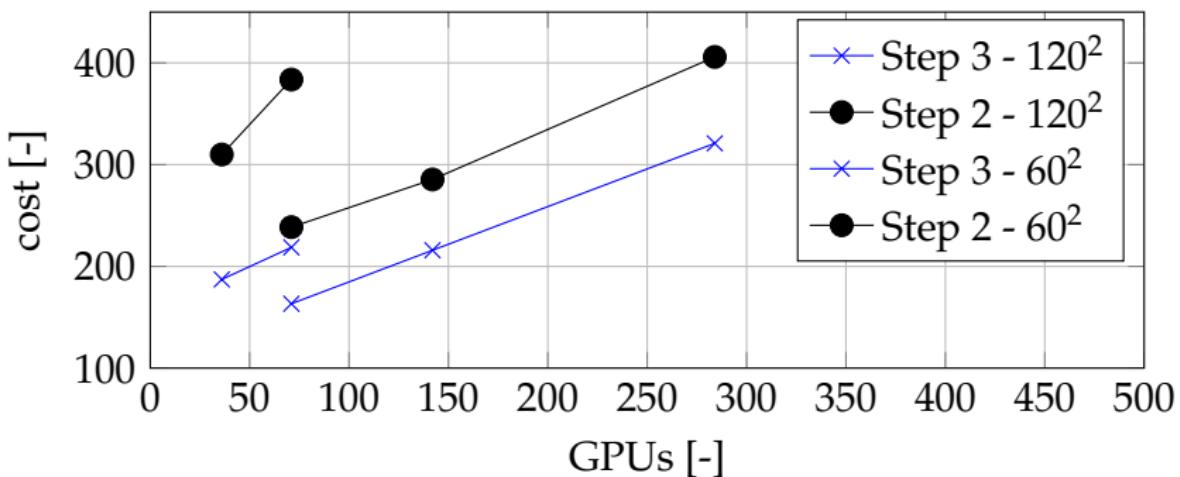
STEP 3 - SPEED UP

- ▶ strong scaling
- ▶ almost 2x speed w.r.t. Step 2
- ▶ “best performance” by placing 4 MPI ranks
- ▶ try 28^2 (vector_length(32))



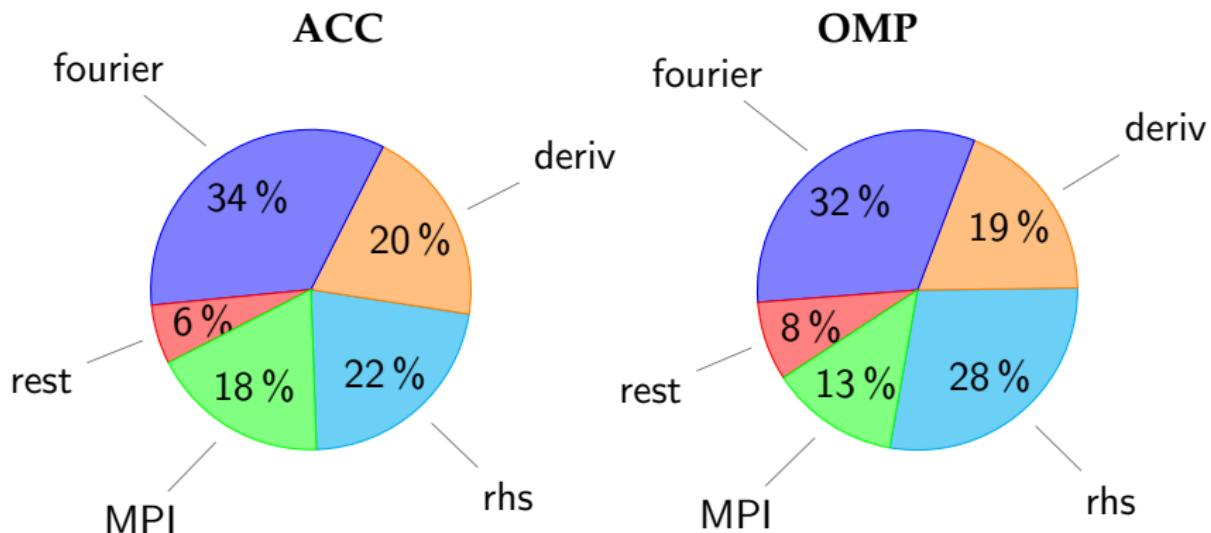
STEP 3 - SPEED UP

- ▶ strong scaling
- ▶ almost 2x speed w.r.t. Step 2
- ▶ “best performance” by placing 4 MPI ranks
- ▶ try 28^2 (`vector_length(32)`)
- ▶ speed up w.r.t. CPU $\approx 2x$



STEP 3 - PROFILES

- ▶ Piz Daint
- ▶ CrayPAT only works for 1 MPI rank per node
- ▶ similar distribution of work according to routines



CONCLUSIONS AND FUTURE WORK

- ▶ HiPSTAR achieved more than 2x speed up (1 GPU vs. 1 CPU)
- ▶ GPU code faster but less flexible

CONCLUSIONS AND FUTURE WORK

- ▶ HiPSTAR achieved more than 2x speed up (1 GPU vs. 1 CPU)
- ▶ GPU code faster but less flexible
- ▶ CPU does not perform number crunching \Rightarrow different ratio of CPU to GPU

CONCLUSIONS AND FUTURE WORK

- ▶ HiPSTAR achieved more than 2x speed up (1 GPU vs. 1 CPU)
- ▶ GPU code faster but less flexible
- ▶ CPU does not perform number crunching ⇒ different ratio of CPU to GPU

ACC Hackaton CSCS

- ▶ correct problems with PGI compilers
- ▶ port rest of the code: FSI, Adjoint method
- ▶ cooperate with experts

CONCLUSIONS AND FUTURE WORK

- ▶ HiPSTAR achieved more than 2x speed up (1 GPU vs. 1 CPU)
- ▶ GPU code faster but less flexible
- ▶ CPU does not perform number crunching ⇒ different ratio of CPU to GPU

ACC Hackaton CSCS

- ▶ correct problems with PGI compilers
- ▶ port rest of the code: FSI, Adjoint method
- ▶ cooperate with experts

Longer term (Summit)

- ▶ Multi GPU environment - which route to go?
- ▶ Are the algorithms optimized for CPU as well the best ones for GPU?