

# nb3dffft

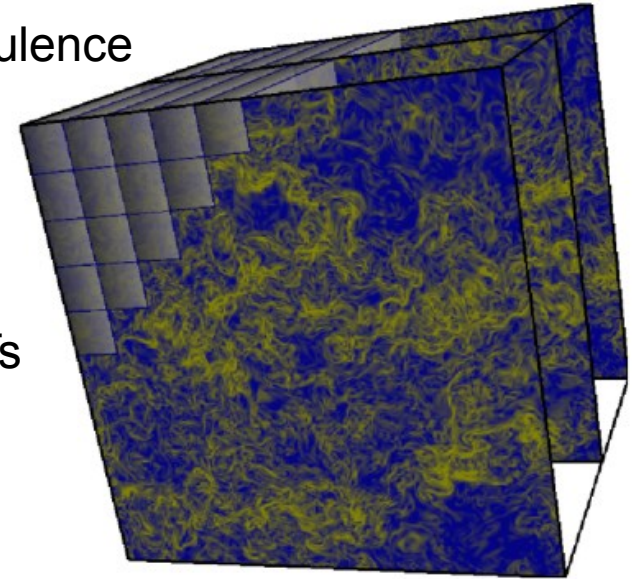
Jens Henrik Göbbert

Jülich Supercomputing Centre, Forschungszentrum Jülich GmbH, Germany

# Motivation

## Pseudo-Spectral Application

- isotropic homogeneous decaying/forced turbulence
  - turbulent channel flow
  - turbulent premixed flames
  - scalar mixing
- 
- 80% of the compute time required by 3d-FFTs
    - all-to-all communication
    - 2x data redistribution among processes



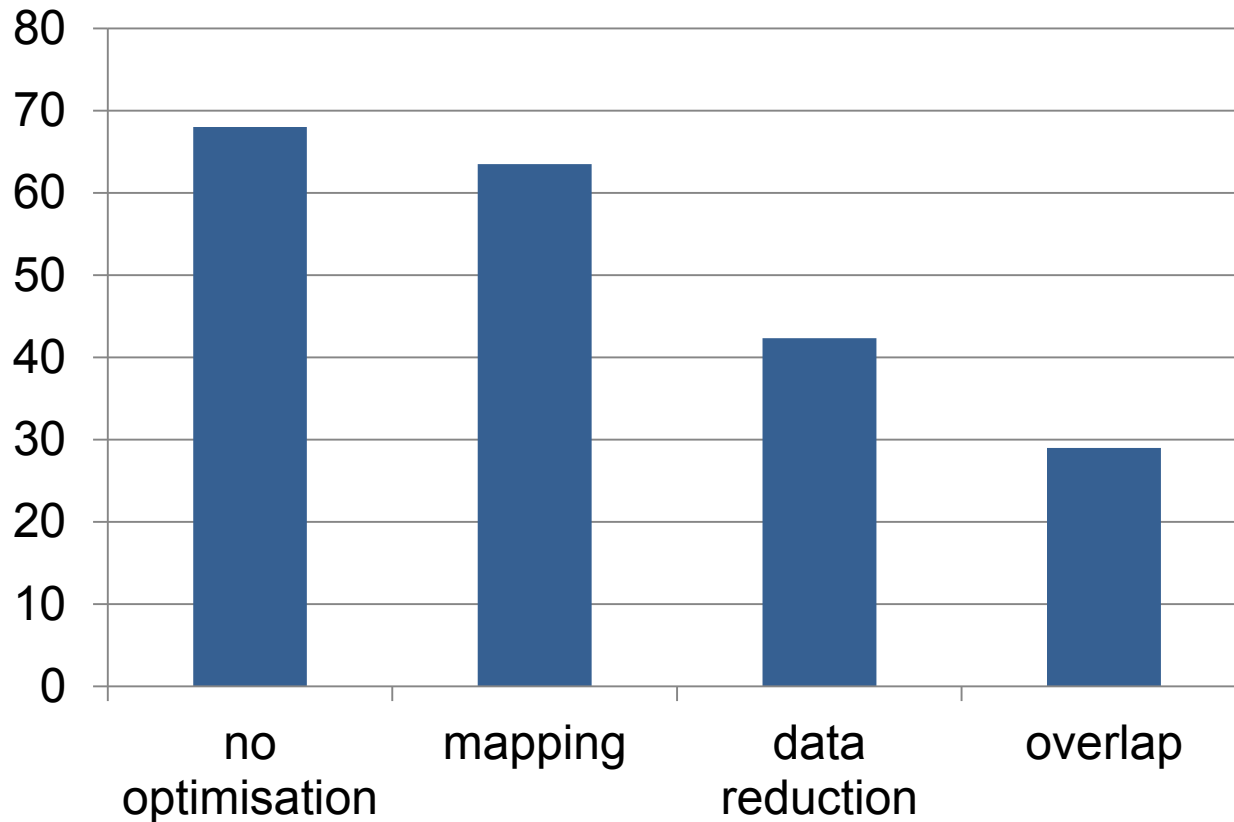
development of the improved 3d-FFT library ,nb3dfft‘

- based on P3DFFT (thanks to Dmitry Pekurovsky)
- data reduction by including the filter process (dealiasing)
- overlapping communication and computation of multiple 3D-FFTs

# Why we wrote our own 3D-FFT

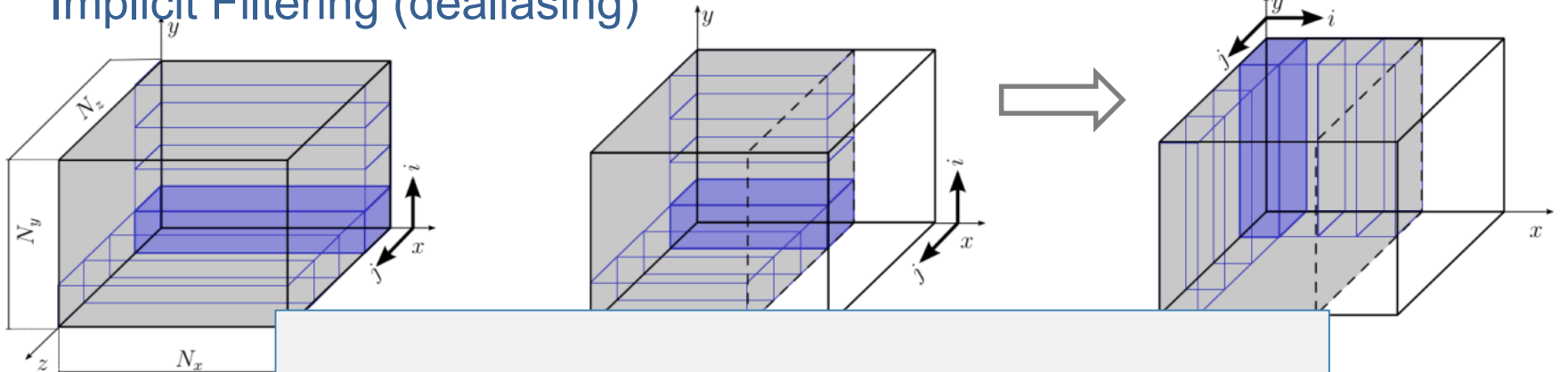
It is all about saving compute time

## Million core/h required to compute $6144^3$ project



# nb3dfft Optimizations

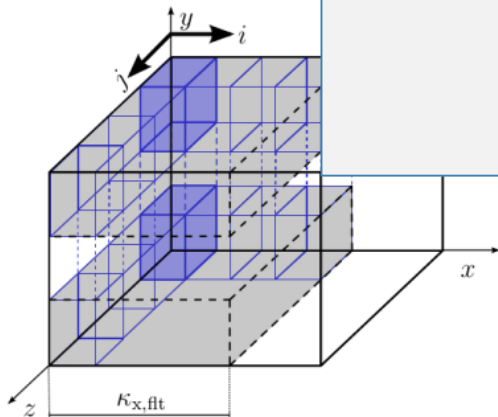
## Implicit Filtering (dealiasing)



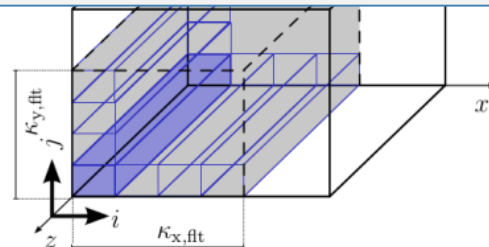
(a) 1d-FFTs applied in  $x$  direction

**70.4% data reduction**  
**44.4% less data to send**  
**44.4% less 1d-FFTs to compute**

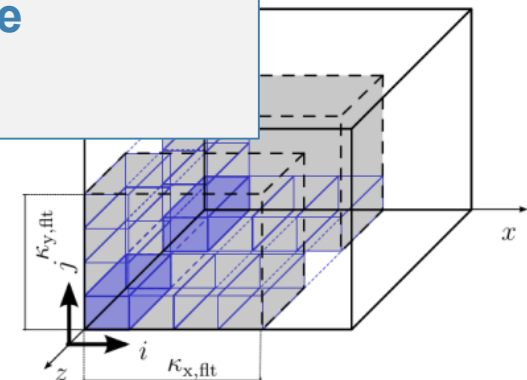
and applying 1d-FFT in  $z$  direction



(d) cut-off filtering in  $y$  while packing for the transpose



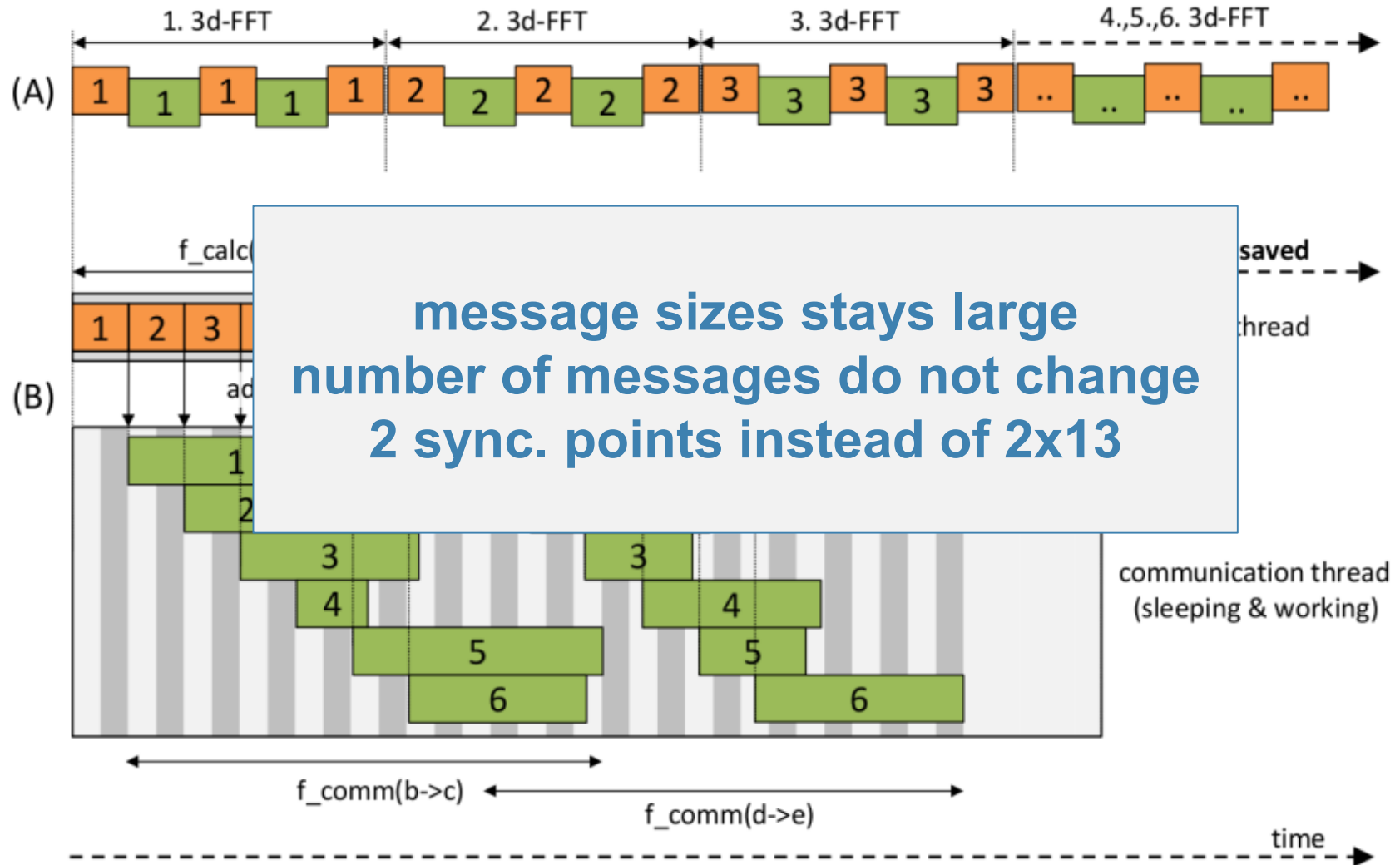
(e) unpacking and applying 1d-FFT in  $z$  direction



(f) cut-off filtering in  $z$  considering access patterns

# nb3dffft Optimizations

Overlapping communication & computation



# Overlapping Communication & Computation

worker + scheduler thread

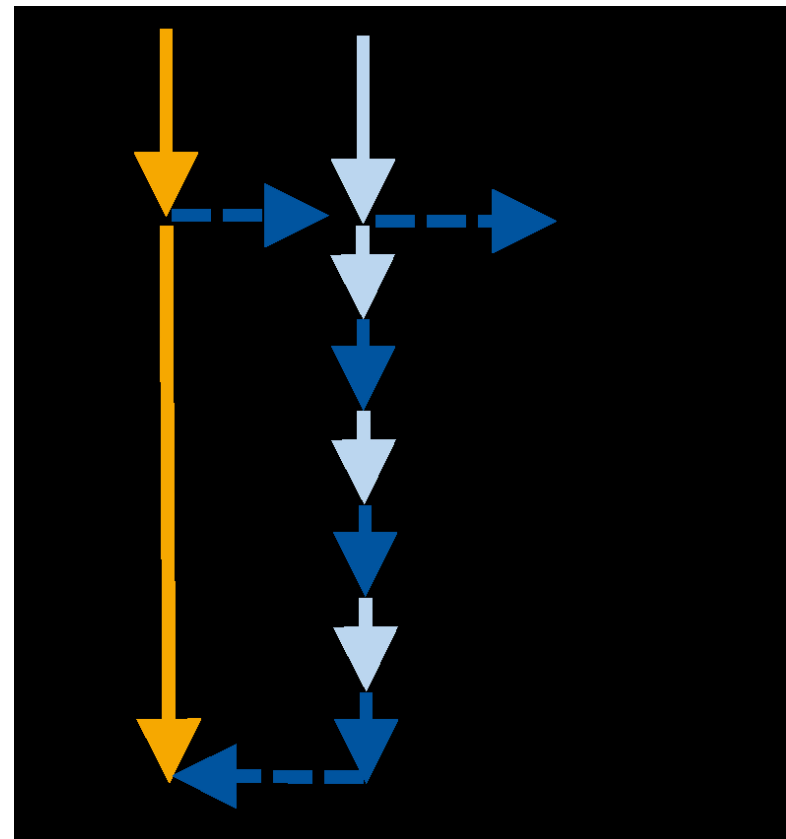
nb3dffft relies on the MPI 3.0 for

- non-blocking collective operations and OpenMP for starting
- scheduler + worker threads.

The **worker** passes any non-blocking MPI operation to the **scheduler**.

The **scheduler** executes the MPI commands and iteratively calls MPI\_TEST to proceed.

worker thread      scheduler thread



# Overlapping Communication & Computation

## application optimizations

```

1: procedure ITERATE VELOCITY ( $n \rightarrow [n+1]$ ), PASSIVE SCALAR ( $[n-1] \rightarrow n$ )
2:   input (velocity):  $\hat{u}_{1|2|3}^n, \hat{G}_{4|5|6}^{n-1}, \hat{f}_u^{n-1}$ 
3:   input (passive scalar):  $\hat{\phi}_{17}^{n-1}, J_{19}^{n-1}, \hat{J}_{18}^{n-2}, \hat{u}_{16}^{n-1}$ 
4:   for all wave numbers ( $\kappa$ ) do:                                ▷ compute vorticity
5:      $\hat{\omega}_{x,10}^n \leftarrow \kappa_z \hat{u}_{y,2}^n - \kappa_y \hat{u}_{z,3}^n$ 
6:      $\hat{\omega}_{y,11}^n \leftarrow \kappa_x \hat{u}_{z,3}^n - \kappa_z \hat{u}_{x,1}^n$ 
7:      $\hat{\omega}_{z,12}^n \leftarrow \kappa_y \hat{u}_{x,1}^n - \kappa_x \hat{u}_{y,2}^n$ 
8:
9:
10:
11:  for all wave numbers ( $\kappa$ ) do:                                ▷ advance passive scalar
12:     $\hat{\phi}_{17}^n \leftarrow F_{\text{adv}}(\hat{\phi}_{17}^{n-1}, \hat{J}_{19}^{n-1}, \hat{J}_{18}^{n-2}, \hat{u}_{16}^{n-1})$ 
13:     $\hat{J}_{19}^{n-1} \stackrel{\text{zero copy}}{=} \hat{J}_{18}^{n-2}$                                 ▷ swap memory
14:    for all grid points ( $N$ ) do:                                ▷ compute  $G^n = \omega^n \times u^n$ 
15:       $G_{x,13}^n \leftarrow u_{y,8}^n \cdot \omega_{z,12}^n - u_{z,9}^n \cdot \omega_{y,11}^n$ 
16:       $G_{y,14}^n \leftarrow u_{z,9}^n \cdot \omega_{x,10}^n - u_{x,7}^n \cdot \omega_{z,12}^n$ 
17:       $G_{z,15}^n \leftarrow u_{x,7}^n \cdot \omega_{y,11}^n - u_{y,8}^n \cdot \omega_{x,10}^n$ 
18:       $\widehat{\nabla} \hat{\phi}_{10|11|12}^n \leftarrow \hat{\phi}_{17}^n$                                 ▷ 3x compute derivatives
19:
20:
21:   $J_{19}^n \leftarrow u_{7|8|9}^n \cdot \nabla \phi_{20|21|22}^n$                                 ▷ compute convective term
22:   $\hat{G}_{13|14|15}^n \stackrel{\text{zero copy}}{=} \hat{G}_{4|5|6}^{n-1}$                                 ▷ swap memory
23:   $\hat{u}_{16}^n \stackrel{\text{zero copy}}{=} \hat{u}_2^n$                                 ▷ swap memory
24:   $\hat{f}_u^n \leftarrow F(\hat{f}_u^{n-1})$                                 ▷ compute forcing energy
25:  for wave numbers ( $|\kappa_{ijk}| < k_f = 2\sqrt{2}$ ) do:            ▷ add forcing energy
26:     $\hat{u}_{1|2|3}^{\text{force}} \leftarrow F_{\text{force}}(\hat{u}_{1|16|3}^n, \hat{f}_u^n)$ 
27:  for all wave numbers ( $\kappa$ ) do:                                ▷ advance velocity
28:     $\hat{u}_{1|2|3}^{\text{step1}} \leftarrow F_{\text{step1}}(\hat{u}_{1|2|3}^{\text{force}}, \hat{G}_{4|5|6}^n, \hat{G}_{13|14|15}^{n-1})$ 
29:  for all wave numbers ( $\kappa$ ) do:                                ▷ apply projection tensor
30:     $\hat{u}_{1|2|3}^{n+1} \leftarrow F_{\text{step2}}(\hat{u}_{1|2|3}^{\text{step1}})$ 
31:     $\langle \hat{u}_{1|2|3}^{n+1} \rangle = 0$                                 ▷ set mean velocity to zero
32:  output (velocity):  $\hat{u}_{1|2|3}^{n+1}, \hat{G}_{4|5|6}^n, u_{7|8|9}^n, \hat{u}_{16}^n$ 
33:  output (passive scalar):  $\hat{\phi}_{17}^n, J_{19}^n, \hat{J}_{18}^{n-1}$ 

```

# nb3dffft Optimizations

## Overlapping communication & computation

- message sizes stay large
- number of messages do not change
- 2 synchronization points instead of 2x13

### Performance increase by overlapping comm & comp

	<b>1024<sup>3</sup></b>	<b>2048<sup>3</sup></b>	<b>4096<sup>3</sup></b>	<b>6144<sup>3</sup></b>
1024 nodes	15.50%	21.60%	-	-
2048 nodes	-	26.63%	29.63%	-
4096 nodes	-	17.70%	19.78%	-
8192 nodes	-	11.64%	22.69%	17.53%



Getting pseudo-spectral codes for DNS **scale on large systems** is challenging

- 3d-FFTs require lots of MPI\_alltoall() calls, which cannot be avoided
- for best performance key functions have to be rewritten

3d-FFT can be optimized, if it is not thought as a single function call, but as part of the algorithm

- Moving the required **filter operation** of the main algorithm into the 3d-FFT reduces the data to be sent by 44.4% and the 1d-FFTs by 44.4%
- Calling multiple 3d-FFT for different fields at the same time allows **overlapping of communication and computation**  
-> ~20% performance gain

**nb3dfft** is available under GPL License at

<https://gitlab.version.fz-juelich.de/goebbert/nb3dfft>

# Appendix

nb3dff

# nb3dffft Optimizations

## Implicit Filtering

### Performance Measurements

#### RWTH Compute Cluster

- 672 MPI processes on 56 nodes
- domain size  $2048^3$

#### JUQUEEN

- 16.384 MPI processes on 8192 cores
- domain size  $4096^3$

