



Funded by the
European Union

Co-ordinated by  **ECMWF**

ESCAPE



Energy-efficient Scalable Algorithms for Weather Prediction at Exascale

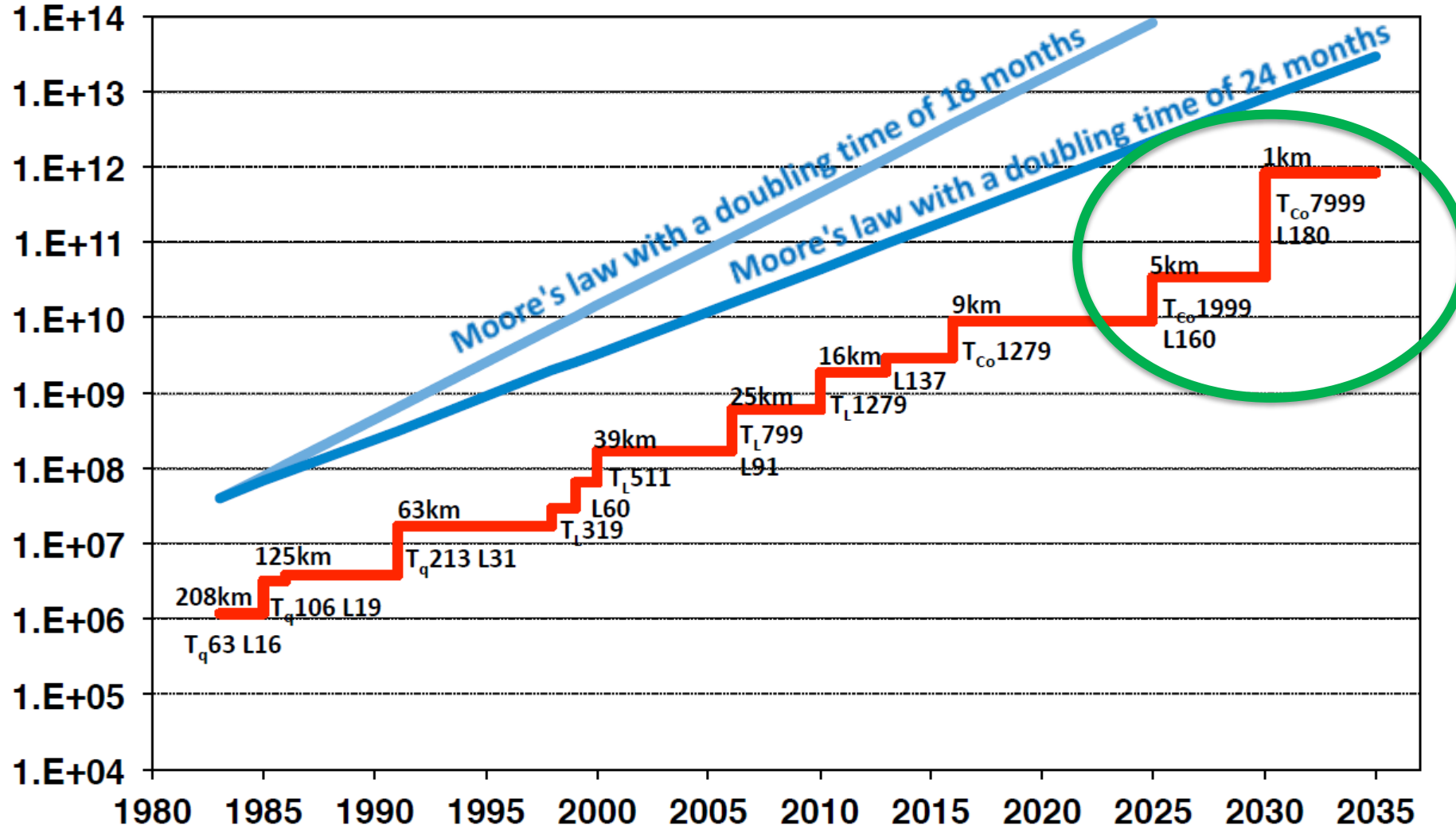
Nils P. Wedi, Peter Bauer, A. Mueller, W. Deconinck, ESCAPE project partners
European Centre for Medium-Range Weather Forecasts (ECMWF)



Outline

- Motivation: Emerging constraints for ensemble-based assimilation and forecasts of Weather & Climate with increasing complexity
- An intermediate goal: globally uniform weather & climate modelling at 1 km horizontal resolution
- *ESCAPE(-2) stands for*
 - Pioneering approaches for refactoring society critical legacy codes
 - Energy-efficient accelerator use in global weather & climate prediction
 - Co-development of novel mathematical algorithms & hardware adaptation
 - Defining and encapsulating the fundamental algorithmic building blocks ("Weather and Climate Dwarfs")
 - Reviewing the need for precision
 - Pioneering algorithm development with hardware adaptation using DSL toolchains
 - A HPCW benchmark and cross-disciplinary Verification, Validation, and Uncertainty Quantification (VVUQ)
 - Application resilience

Computational power drives spatial resolution



Gap of sustained and peak performance

Steepness of gradient from 10km to 1km

(Schulthess et al, 2018)

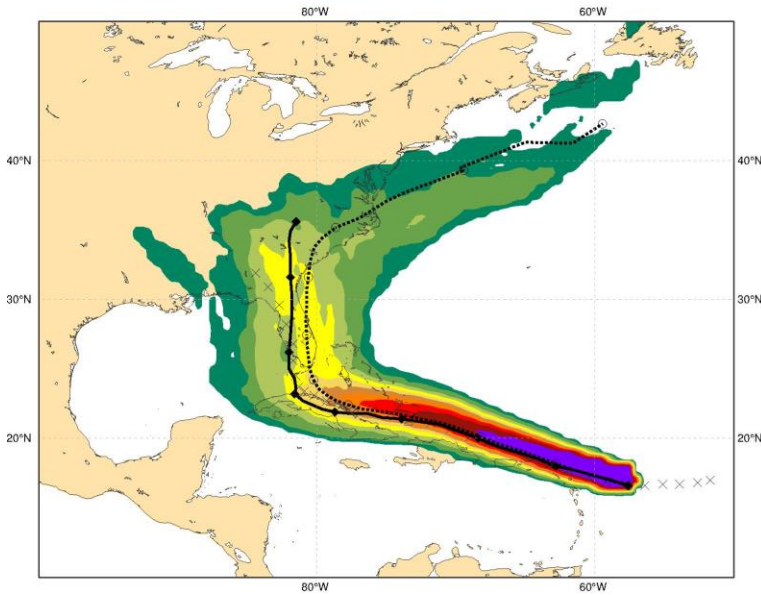
ECMWF's progress in degrees of freedom
(levels x grid columns x prognostic variables)

Hurricane IRMA 18km vs 5km ensemble

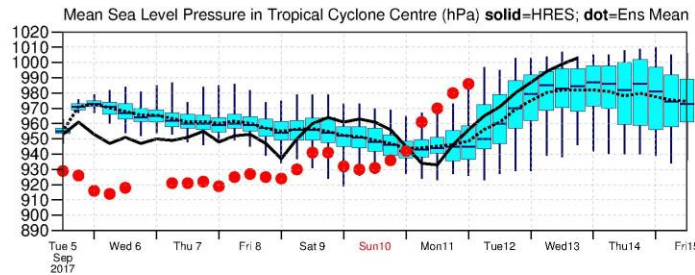
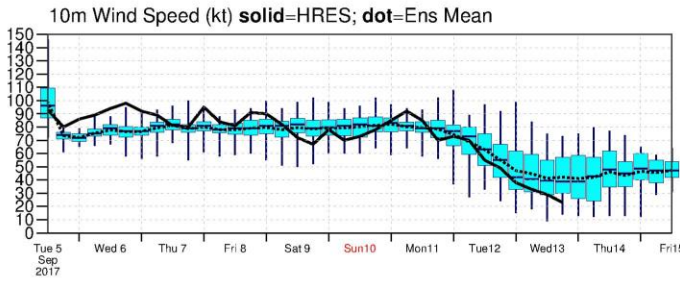
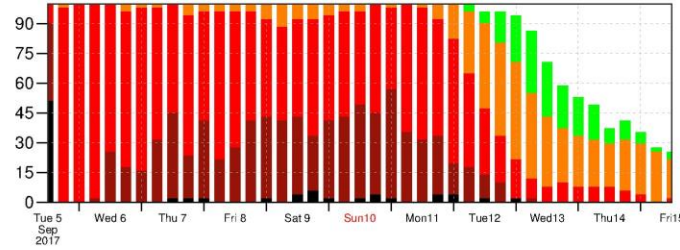
Date 20170905 12 UTC @ ECMWF

Probability that **IRMA** will pass within 120 km radius during the next 240 hours
tracks: **solid**=HRES; **dot**=Ens Mean [reported minimum central pressure (hPa) 929]

5-10 10-20 20-30 30-40 40-50 50-60 60-70 70-80 80-90 > 90%



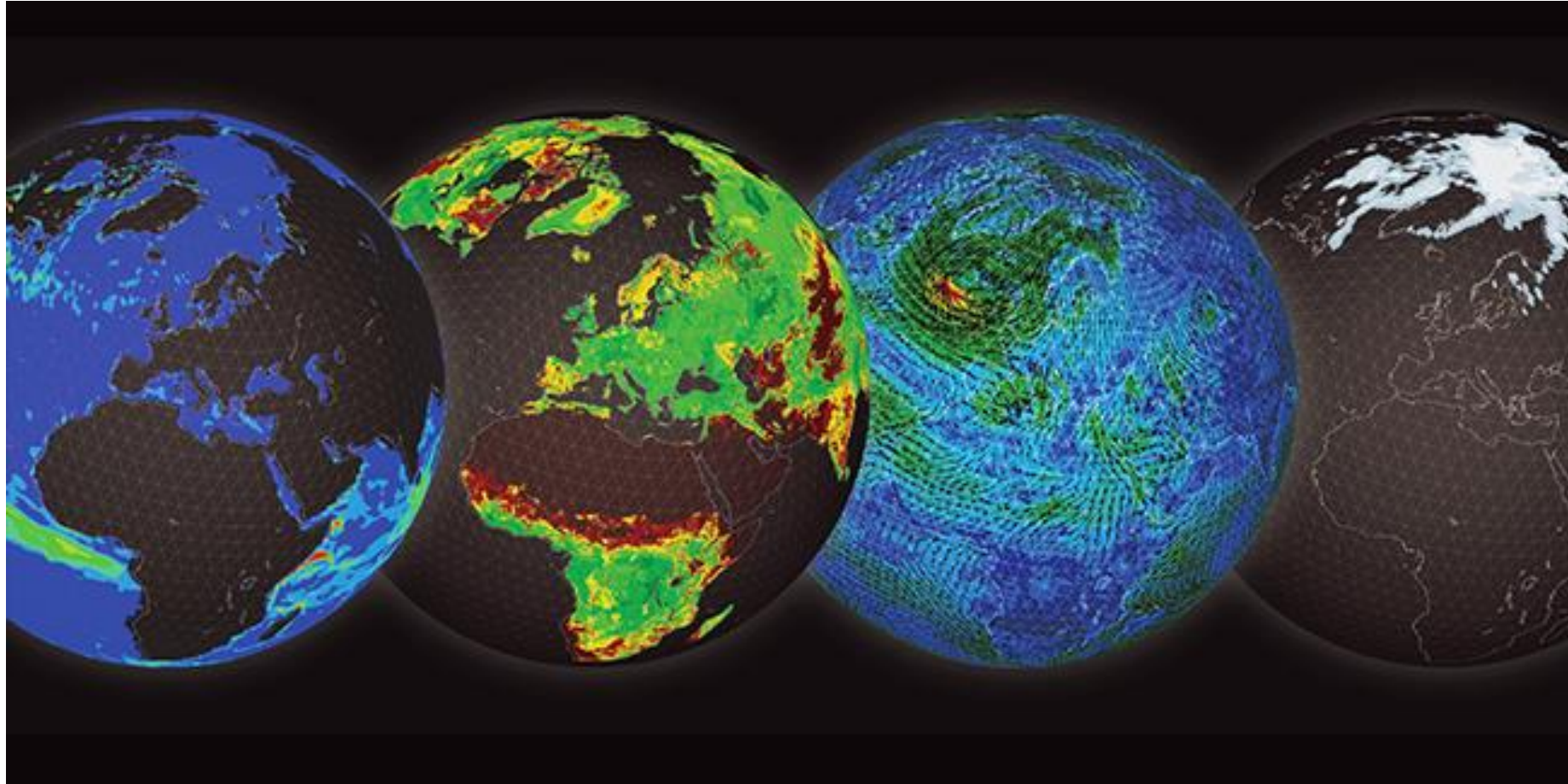
Probability (%) of Tropical Cyclone Intensity falling in each category
TD [up to 33] TS [34-63] HR1 [64-82] HR2 [83-95] HR3 [> 95 kt]



ECMWF Strategy 2025
a 5km ensemble ...

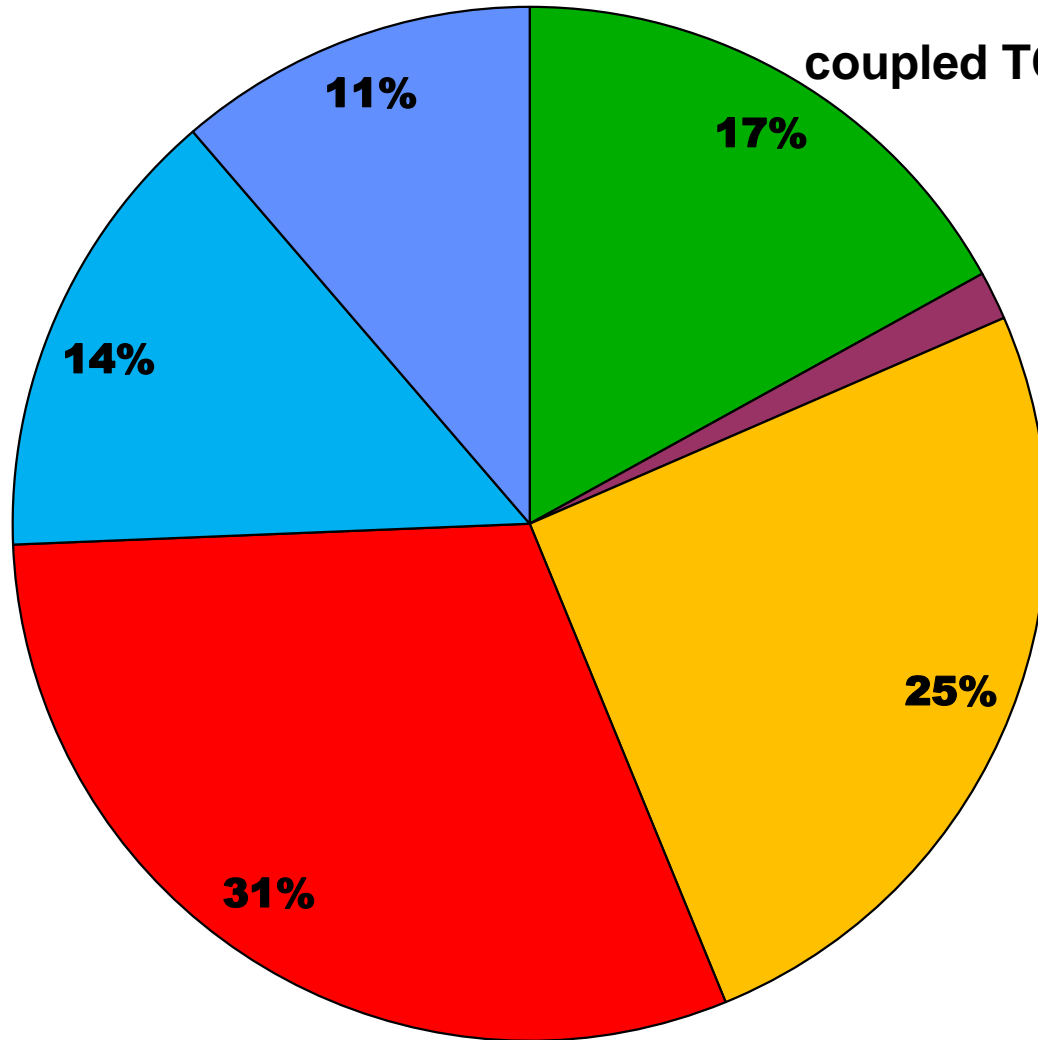
S. Lang & L. Magnusson

Ocean – Land – Atmosphere – Sea ice



Where do we spend the time ? Cycle 45r1

■ GP_DYNAMICS ■ SI_SOLVER ■ SP_TRANSFORMS ■ PHYSICS+RAD ■ WAVEMODEL ■ OCEANMODEL



coupled TCo1279 L137 (~9km operational) run

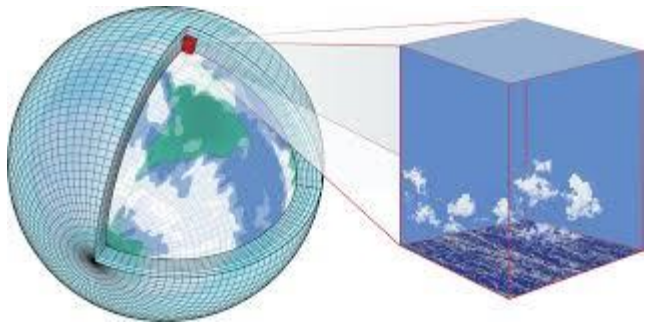
Single electrical group:
~52 minutes wallclock time
(single electrical group==384 nodes)

1408 MPI tasks x 18 threads
290 FC/day

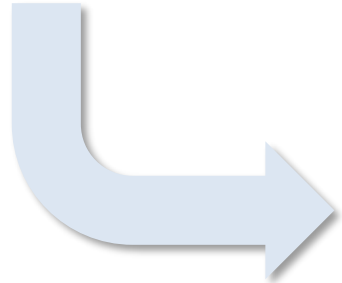
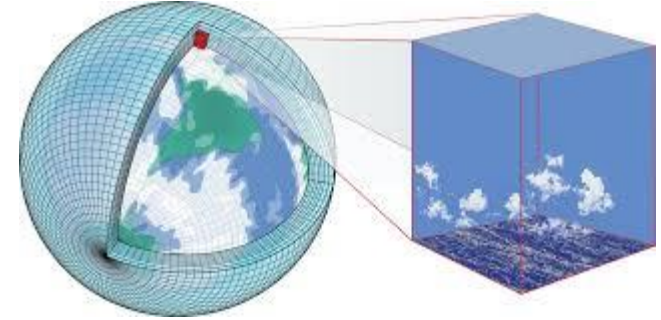
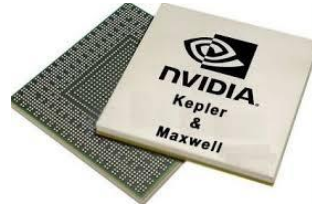
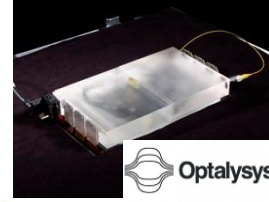


Weather & Climate Dwarfs

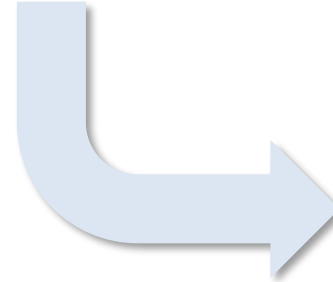
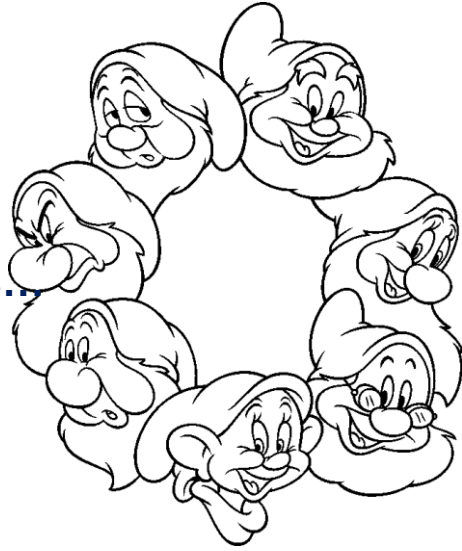
(hpc-escape.eu)



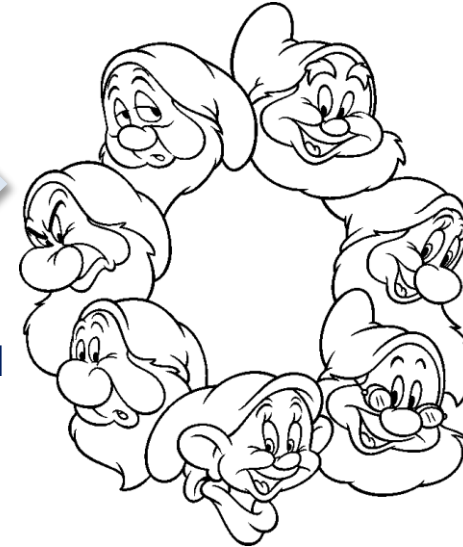
... hardware adaptation ...



Extract model dwarfs...



... explore alternative numerical algorithms ...

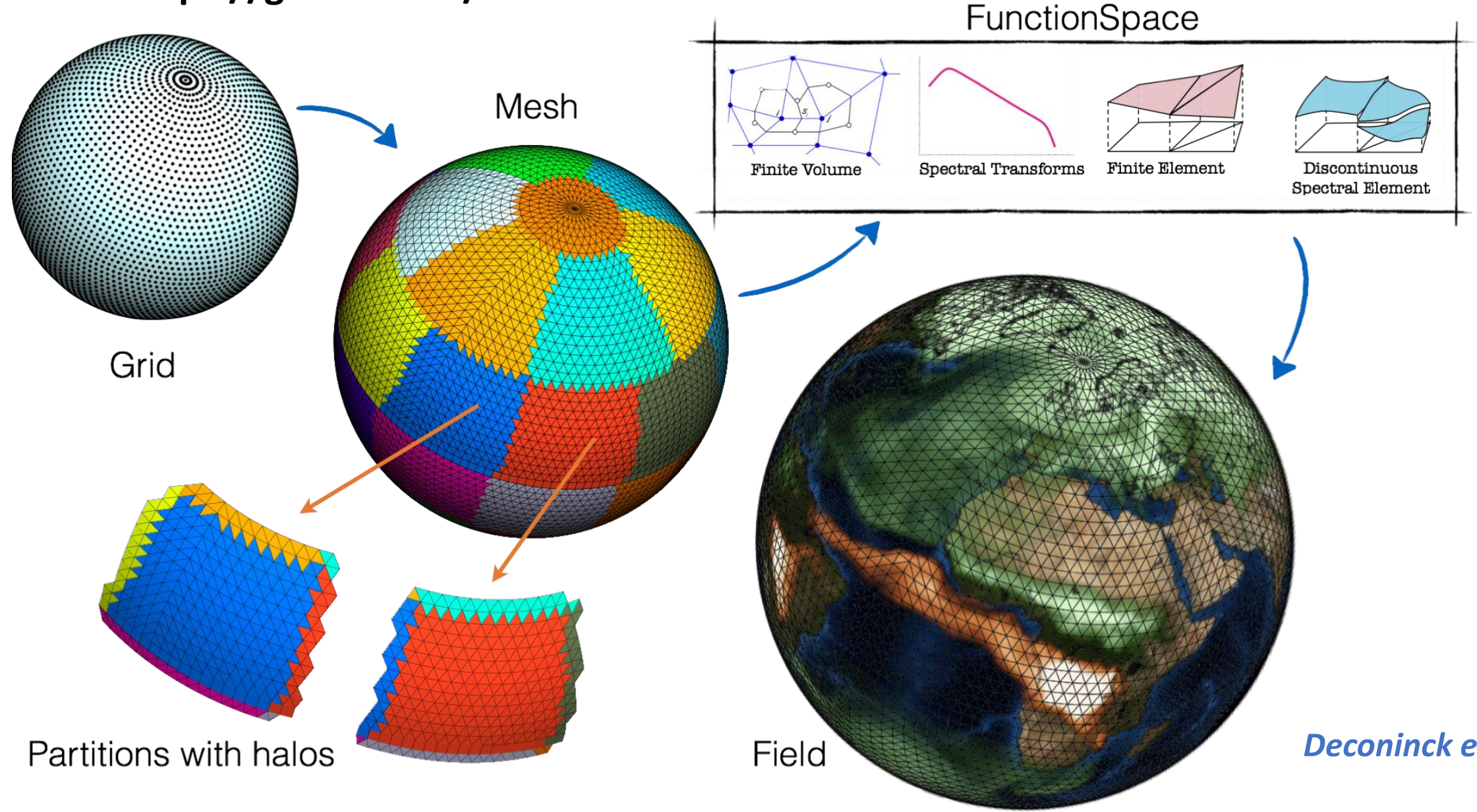


... reassemble model and benchmark



Atlas: a library for NWP and climate modelling

<https://github.com/ecmwf>



Deconinck et al. 2017



$$F(\Psi_L, \Psi_R, U) \equiv [U]^+ \Psi_L + [U]^- \Psi_R \quad (3a)$$

$$U \equiv \frac{u\delta t}{\delta x}, \quad [U]^+ \equiv 0.5(U + |U|), \quad [U]^- \equiv 0.5(U - |U|)$$

Advection (MPDATA)



```
template <uint_t Color> struct upwind_flux {
using flux = accessor<0, enumtype::inout, icosahedron_topology_t>;
using pD =
    in_accessor<1, icosahedron_topology_t::vertices>;
using vn = in_accessor<2, icosahedron_topology_t::vertices>;

typedef boost::mpl::vector<flux, pD, vn> arg_list;

template <typename Evaluation> static void Do(Evaluation &eval) {
    constexpr auto neighbors_offsets =
        connectivity<edges, vertices, Color>::neighbors_offsets;
    constexpr auto ip0 = neighbors_offsets[0];
    constexpr auto ip1 = neighbors_offsets[1];

    float_type pos = math::max(eval(vn()), (float_type)0);
    float_type neg = math::min(eval(vn()), (float_type)0);

    eval(flux()) = eval(pos * pD(ip0) + neg * pD(ip1));
}
};
```

```
ite_upwind_flux(this, pflux, pD, pVn)
=>, intent(inout) :: this
t(out) :: pflux(:, :)
t(in) :: pVn(:, :), pD(:, :)
s, zneg
jges
jvels
je, jlev, ip1, ip2

debug('compute_upwind_flux')

s%dimensions%nb_edges
s%dimensions%nb_levels

DO SCHEDULE(STATIC) PRIVATE(jedge, jlev, ip1, ip2, zpos, zneg)
    _edges
    >de(1, jedge)
    >de(2, jedge)
    _levels
        = max(0._wp, pVn(jlev, jedge))
        = min(0._wp, pVn(jlev, jedge))
    jedge) = pD(jlev, ip1)*zpos + pD(jlev, ip2)*zneg
enddo
!$OMP END PARALLEL DO
end subroutine compute_upwind_flux
```

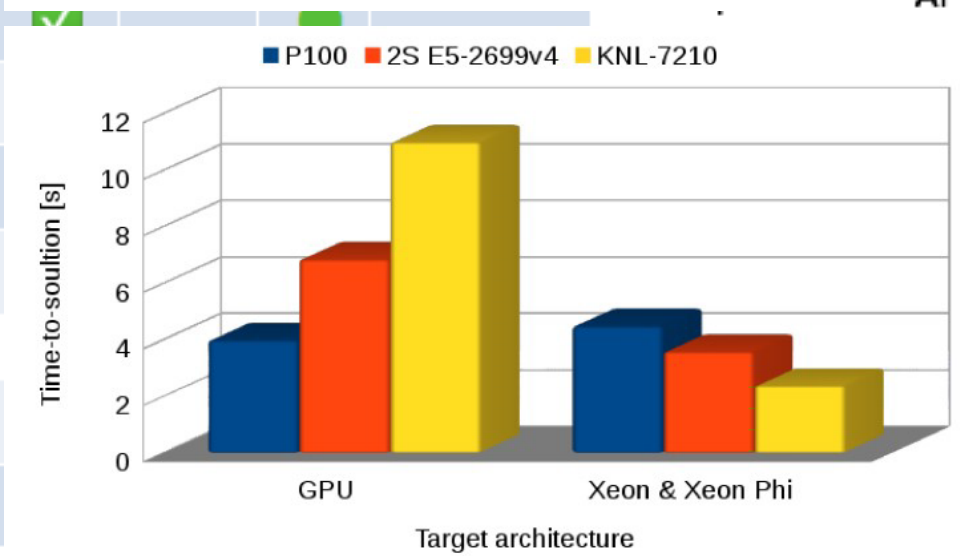
Complementary skills of CLAW, GridTools (MeteoSwiss) and Atlas (ECMWF)



Dwarf	prototype implemented	documented	based on Atlas	MPI	Open MP	Open ACC	DSL	Optalysys
D - spectral transform - SH	✓	✓	✓	✓	✓	✓		
D - spectral transform - biFFT	✓	✓		✓	✓	✓		✓
D - advection - MPDATA	✓	✓	✓	✓	✓	✓	✓	
D - advection - semi-Lagrangian	✓	✓	✓	✓				
D - elliptic solver - GCR	✓	✓	✓	✓				
P - cloud microphysics - CloudSC	✓	✓		✓				
P - radiation scheme - ACRANEB2	✓	👷	👷	✓				
I - LAIRI (3d interpol. algorithm)	✓	✓						
planned next:								
D - advection - discontinuousGalerkin	●	●	●	●				
D - elliptic solver - multigridPrecon	●	●	●	●				

✓: first version running
 👷: in progress
 ●: planned

Comparison of software optimized for GPU and Xeon processors



Poulsen & Berg (2017)

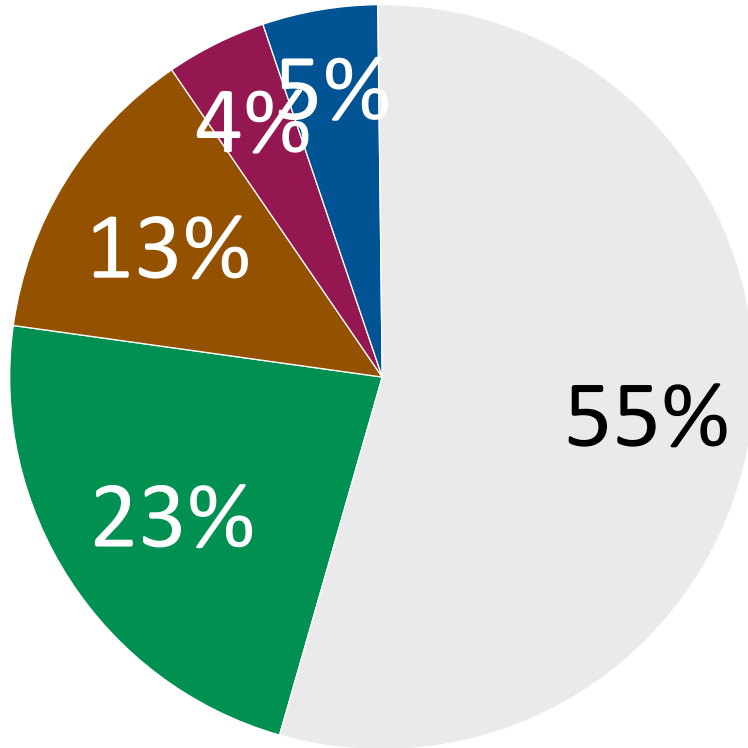


Funded by the European Union

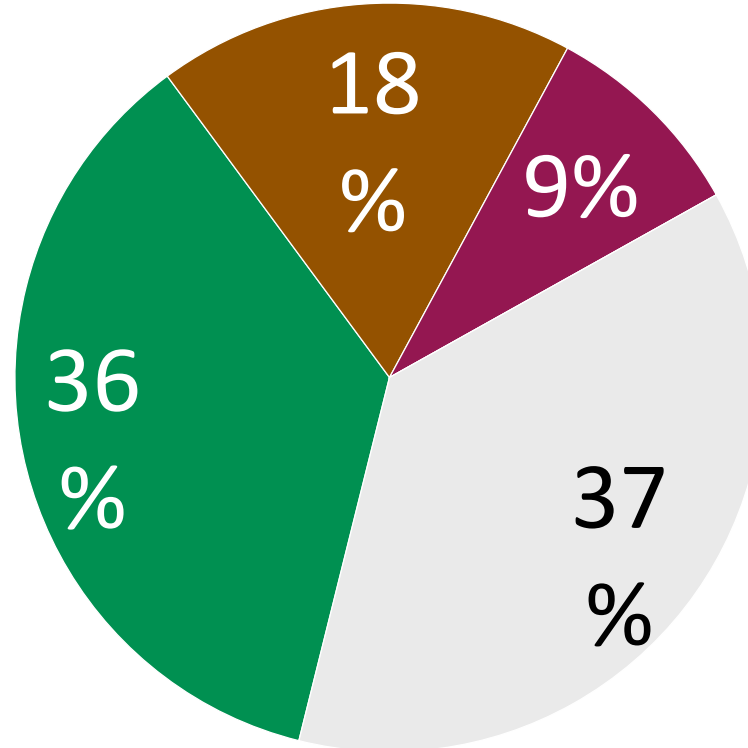
dwarfs

ESCAPE

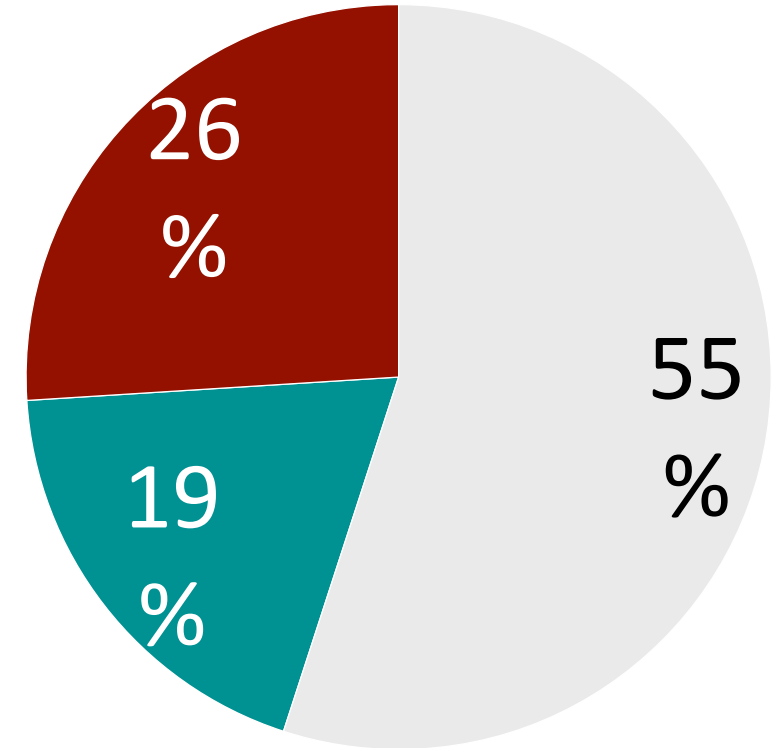
IFS 9km (ECMWF)



ALARO-EPS 2.5km (RMI)



COSMO-EULAG 2.2km (PSNC)

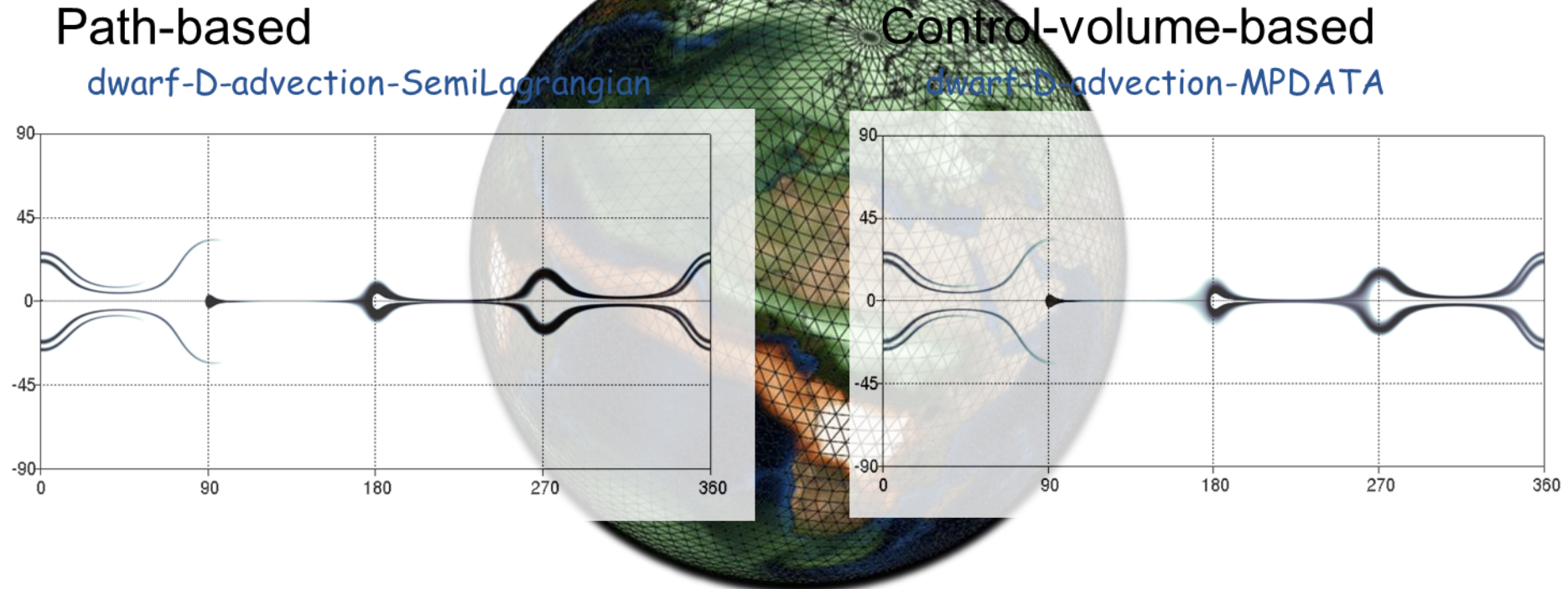


■ spectral transform

■ GCR solver

Advection

Rossby-Haurwitz test case after 7 days



Atlas library support for both prototype implementations

Moist baroclinic instability with FVM and spectral-transform IFS (ST) using large-scale condensation and diagnostic precipitation:

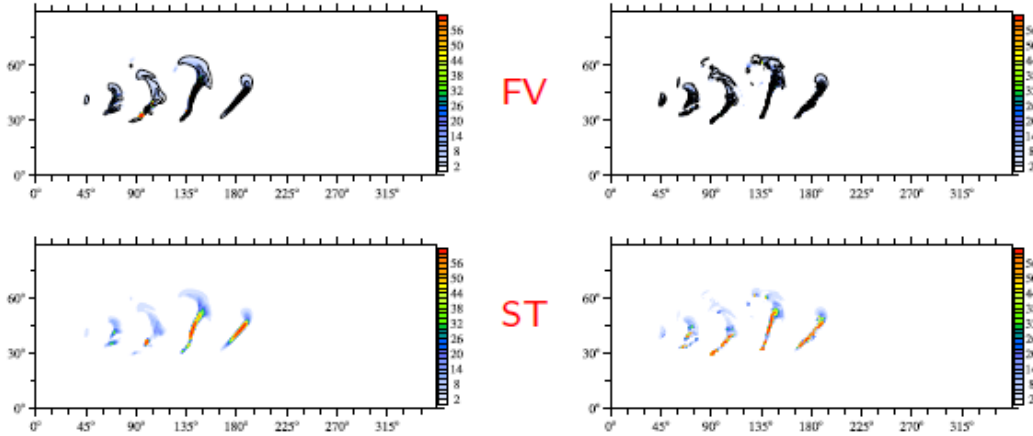
Alternative dynamical core choices on the same grid with the same physics!

Christian Kuehnlein

Precipitation (mm/day) at day 10

O160/TCo159, $\Delta_h \approx 62$ km

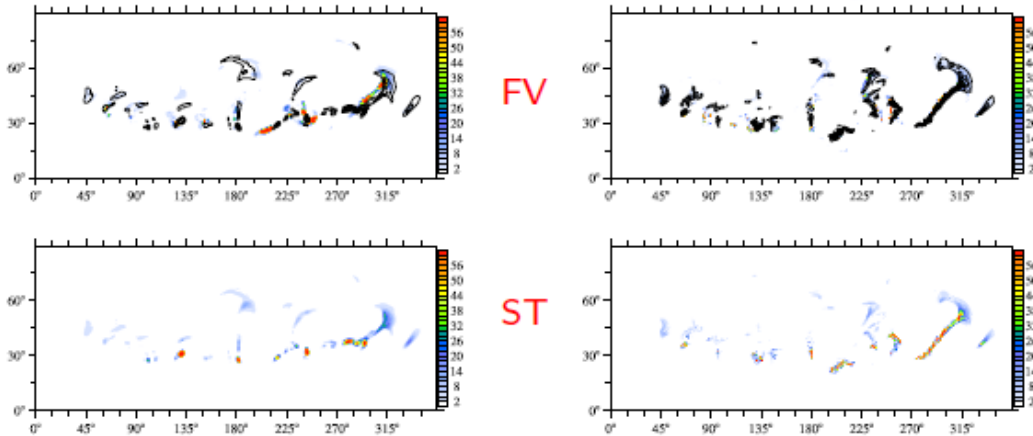
O640/TCo639, $\Delta_h \approx 18$ km



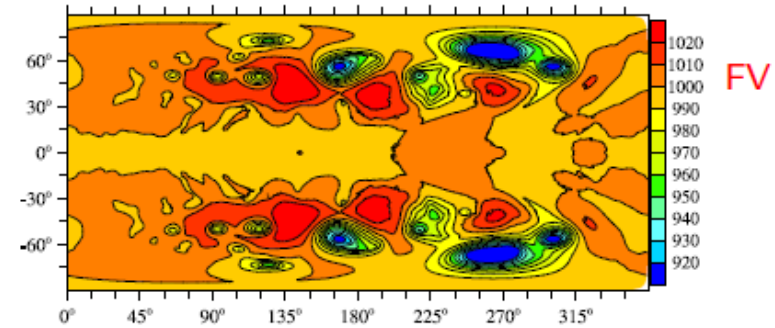
Precipitation (mm/day) at day 15

O160/TCo159, $\Delta_h \approx 62$ km

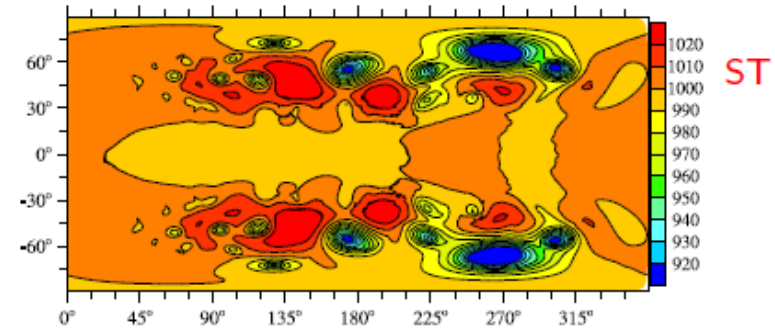
O640/TCo639, $\Delta_h \approx 18$ km



Surface pressure O640/TCo639, $\Delta_h \approx 18$ km, day 15



Finite volume



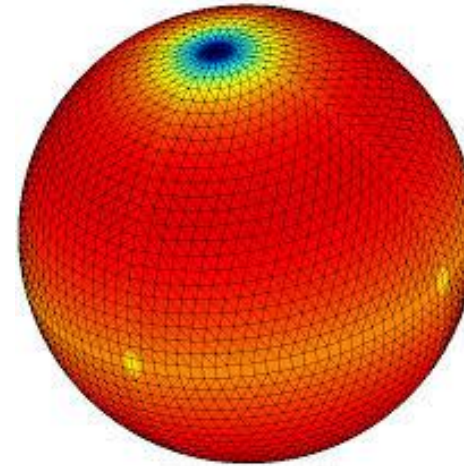
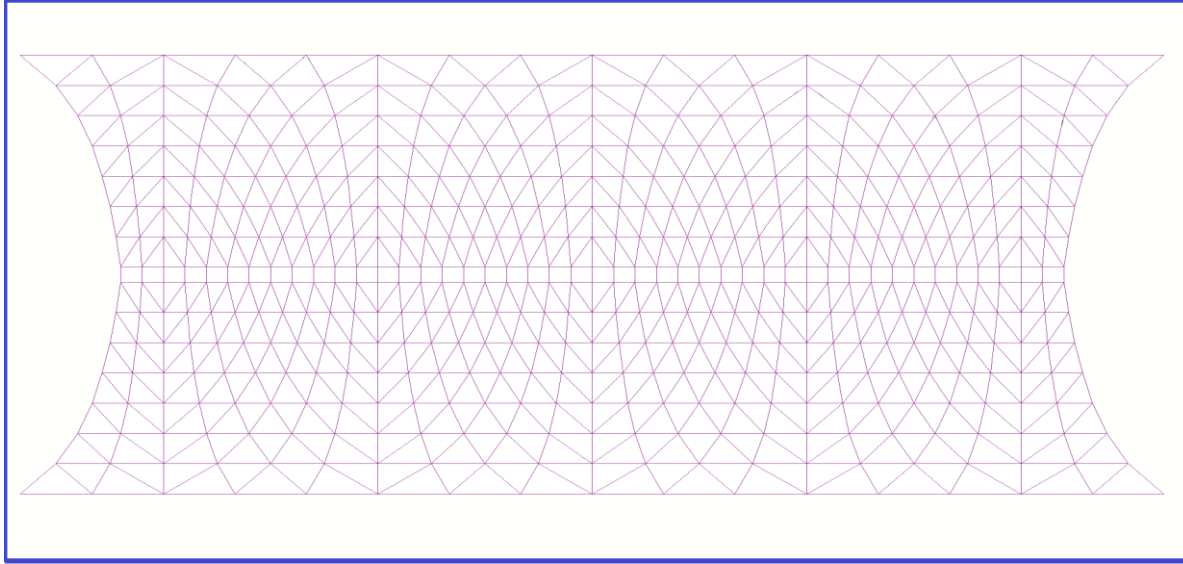
Spectral transform

Another alternative: higher-order finite-difference development by *M. Glinton & P. Bénard*

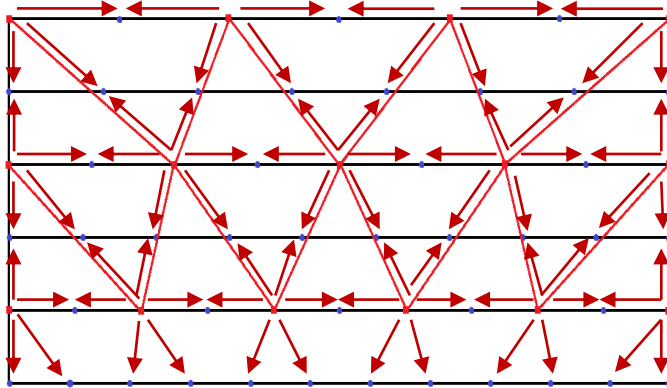
Bespoke Krylov solvers



Parallel restriction, prolongation and Atlas mesh generation

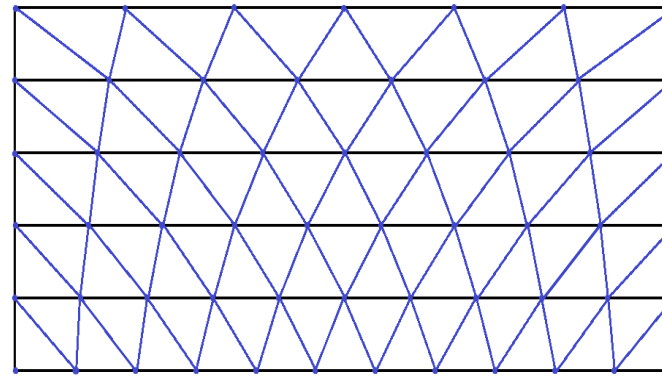


Distributed & Atlas



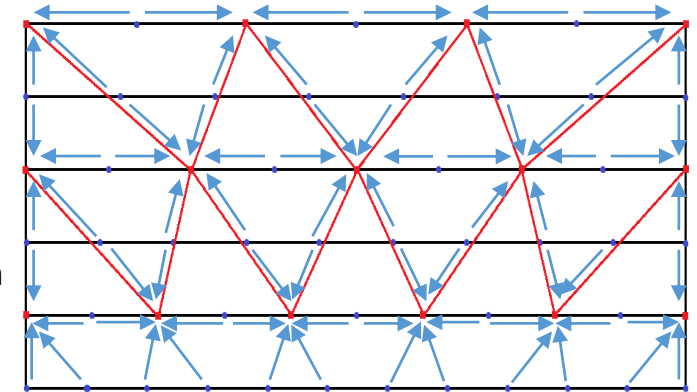
Coarse mesh

interpolation



Fine mesh

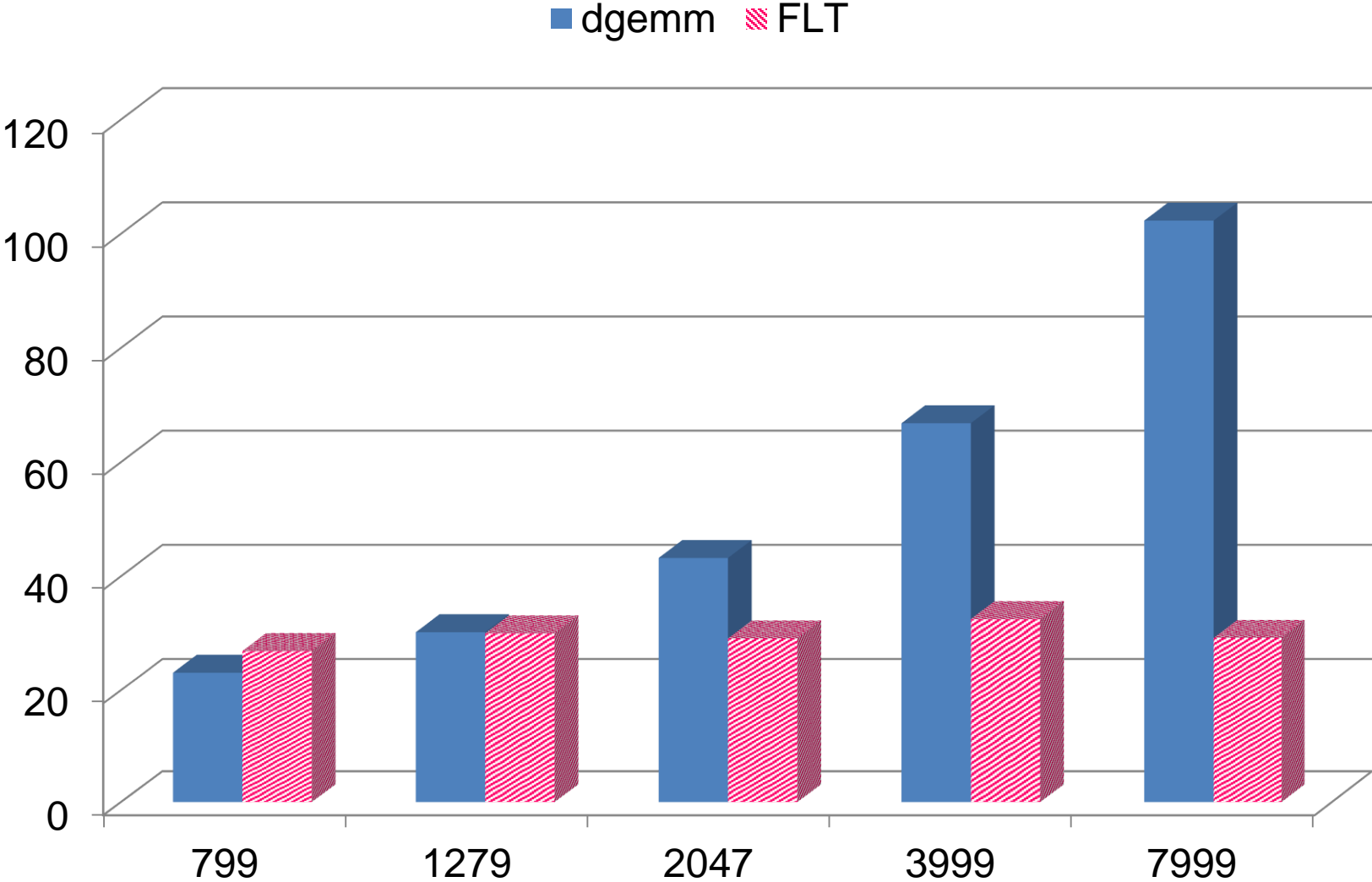
restriction



Coarse mesh

Fast Legendre Transform

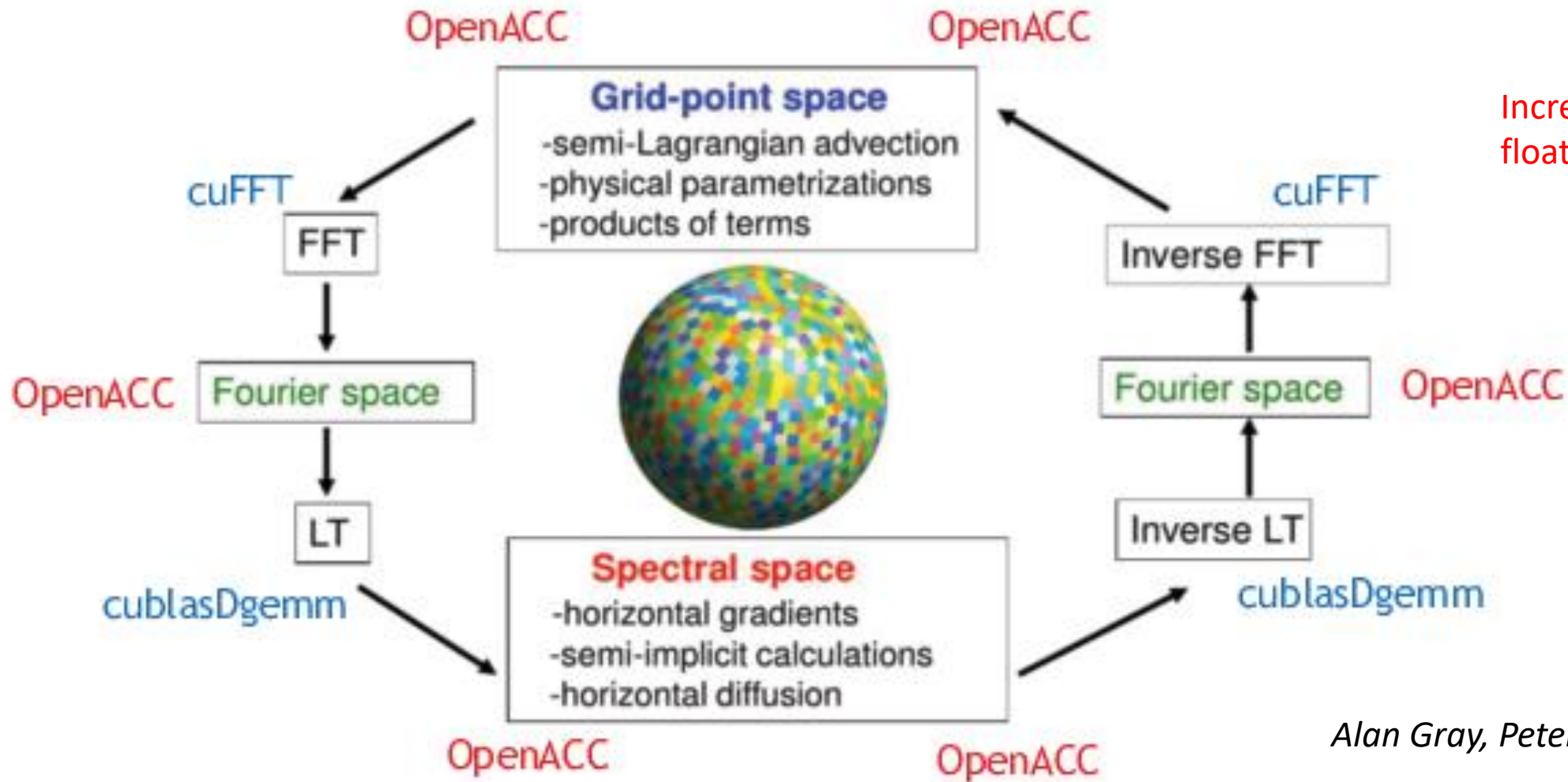
Reduces number of floating point operations



Number of floating point operations for direct or inverse spectral transforms of a single field, scaled by $N^2 \log^3 N$



Schematic description of the *spectral transform method* in the ECMWF IFS model



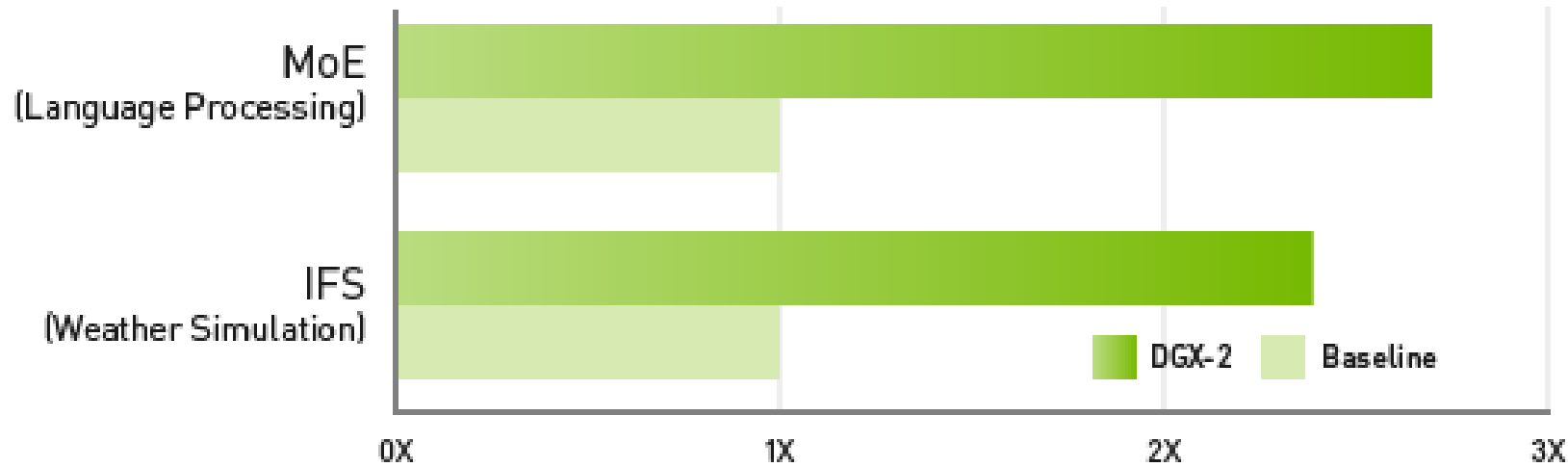
Increases number of floating point operations

Alan Gray, Peter Messmer, NVIDIA



Will Deep Learning influence algorithmic choices for weather & climate ?

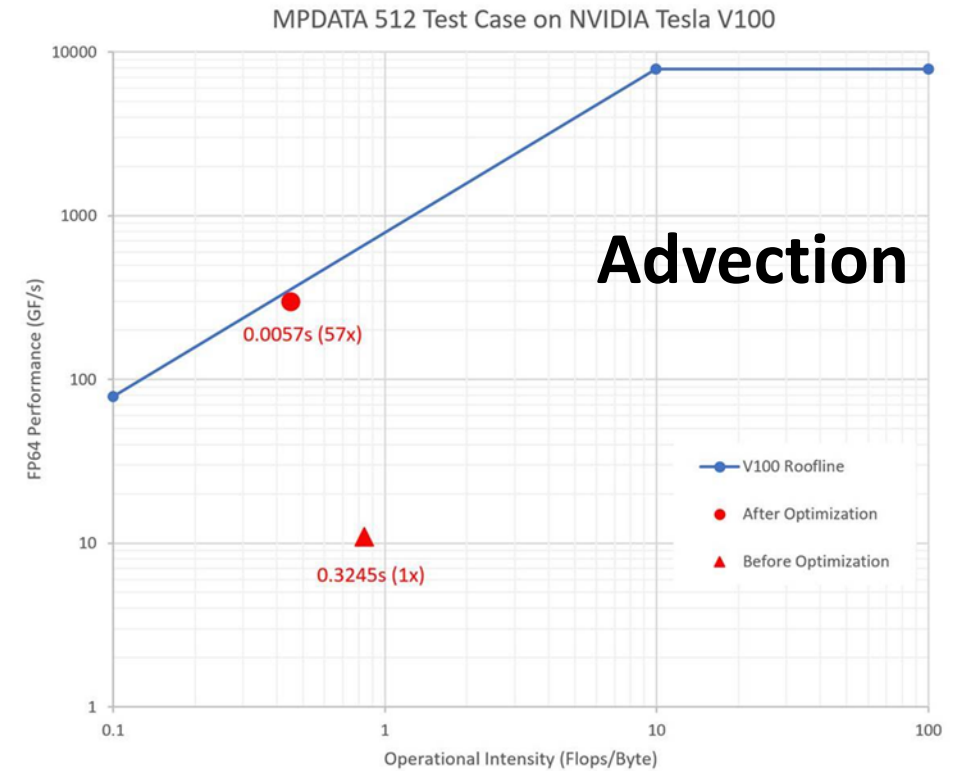
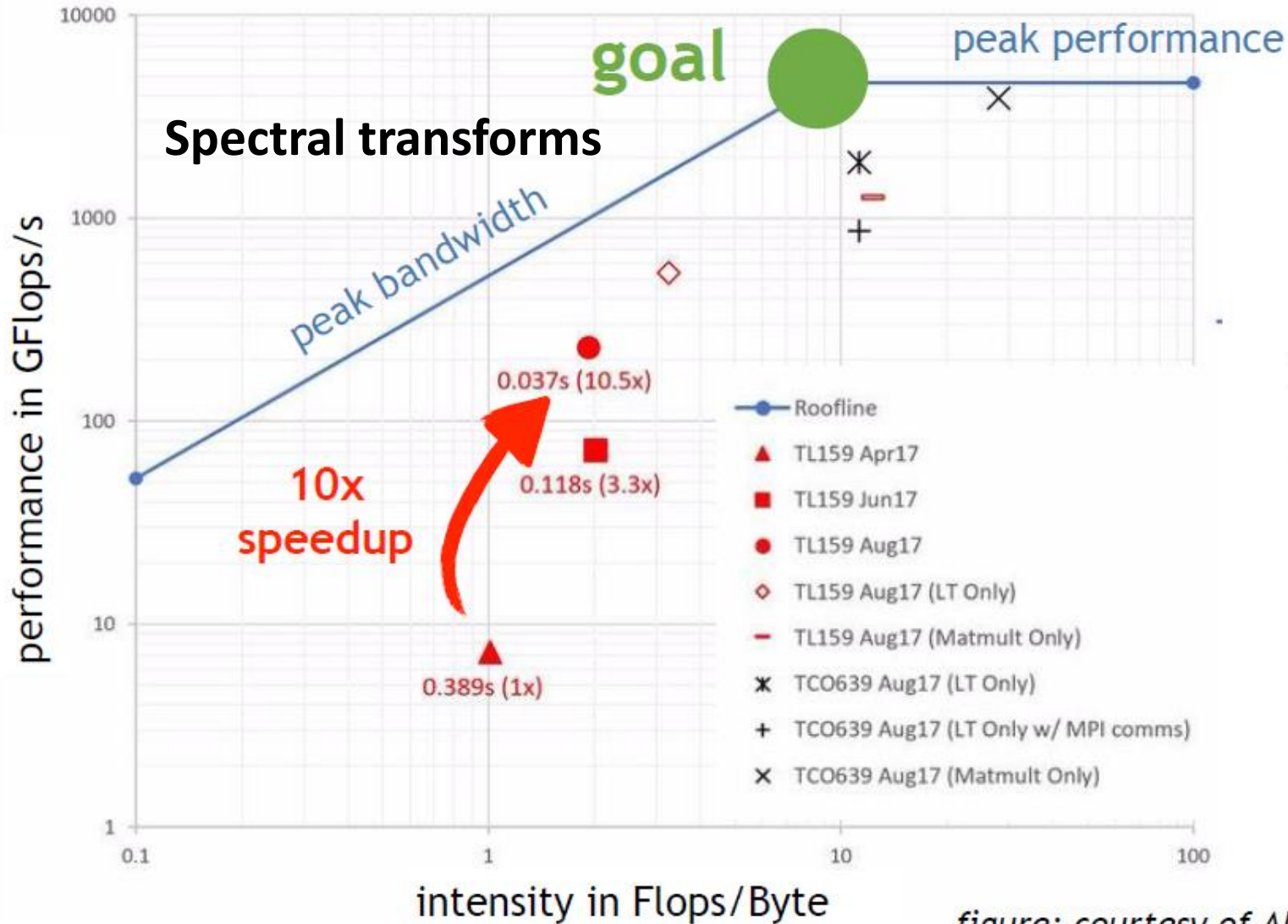
NVSwitch Delivers a >2X Speedup for Deep Learning and HPC*



System Configs: Each of the two DGX-1 servers have dual-socket Xeon E5 2690v4 Processor, 8 x V100 GPUs; servers connected via a 4 EDR (100Gb) InfiniBand connections. DGX-2 server has dual-socket Xeon Scalable Processor Platinum 8168 Processors, 16 x Tesla V100 GPUs.

See talk by A. Gray

<https://news.developer.nvidia.com/nvswitch-leveraging-nvlink-to-maximum-effect/>



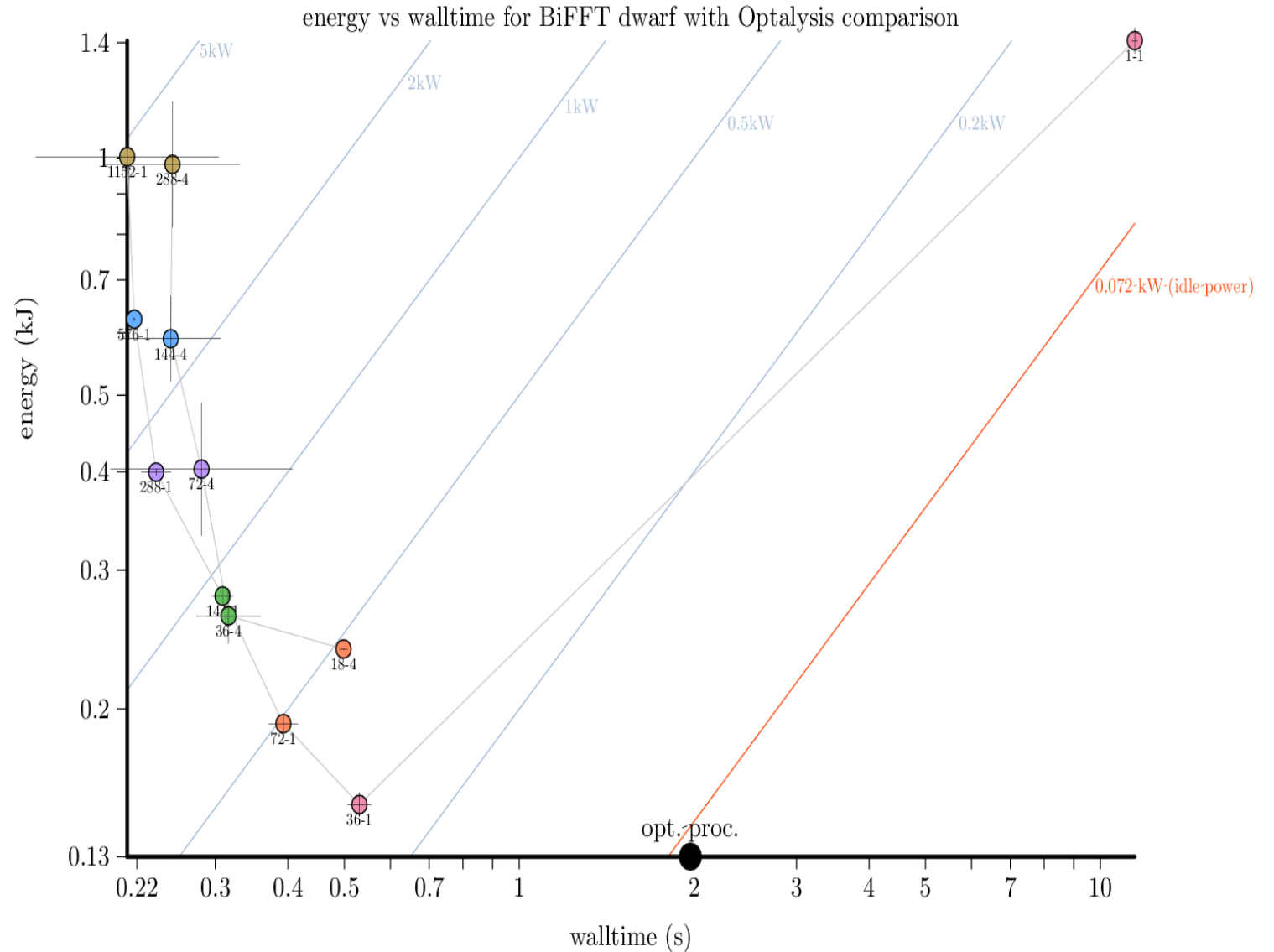
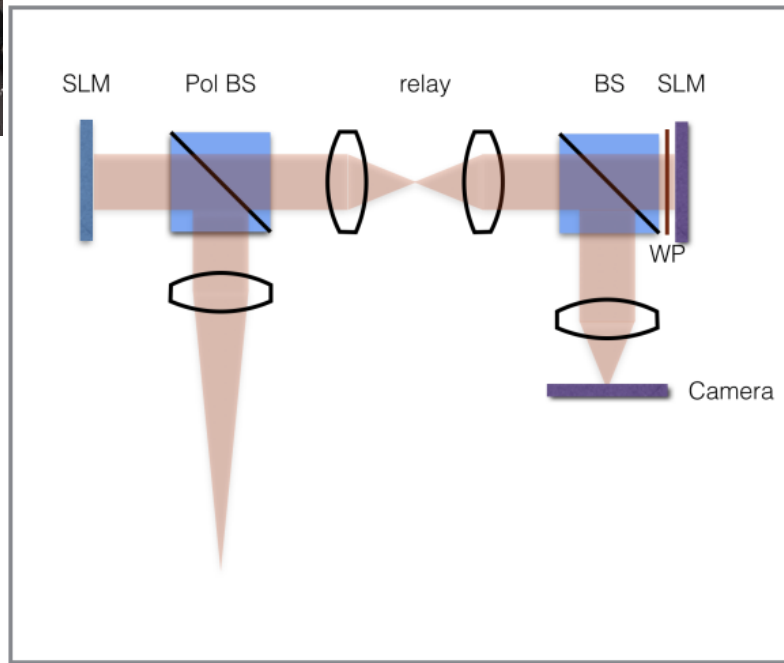
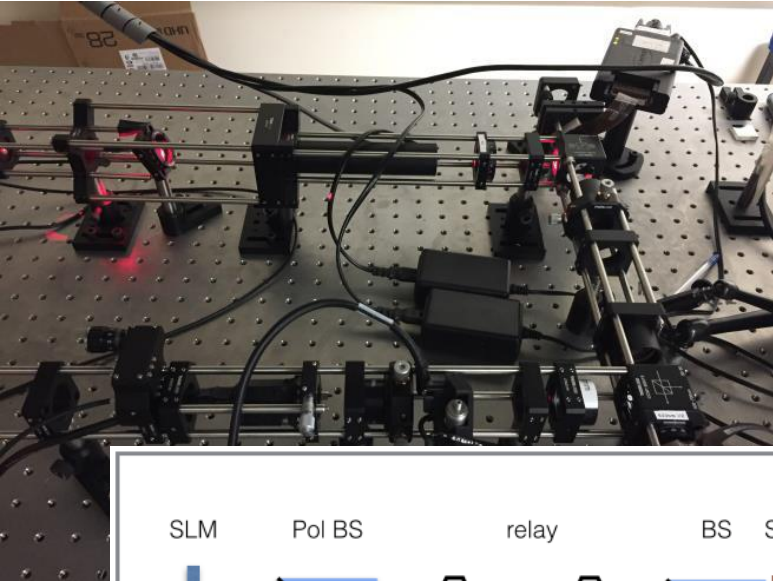
GPU speed-up

figure: courtesy of Alan Gray, Peter Messmer (NVIDIA)



Funded by the European Union

Optalysys: optical processor for spectral transform (biFFT and spherical harmonics) at speed of light ?

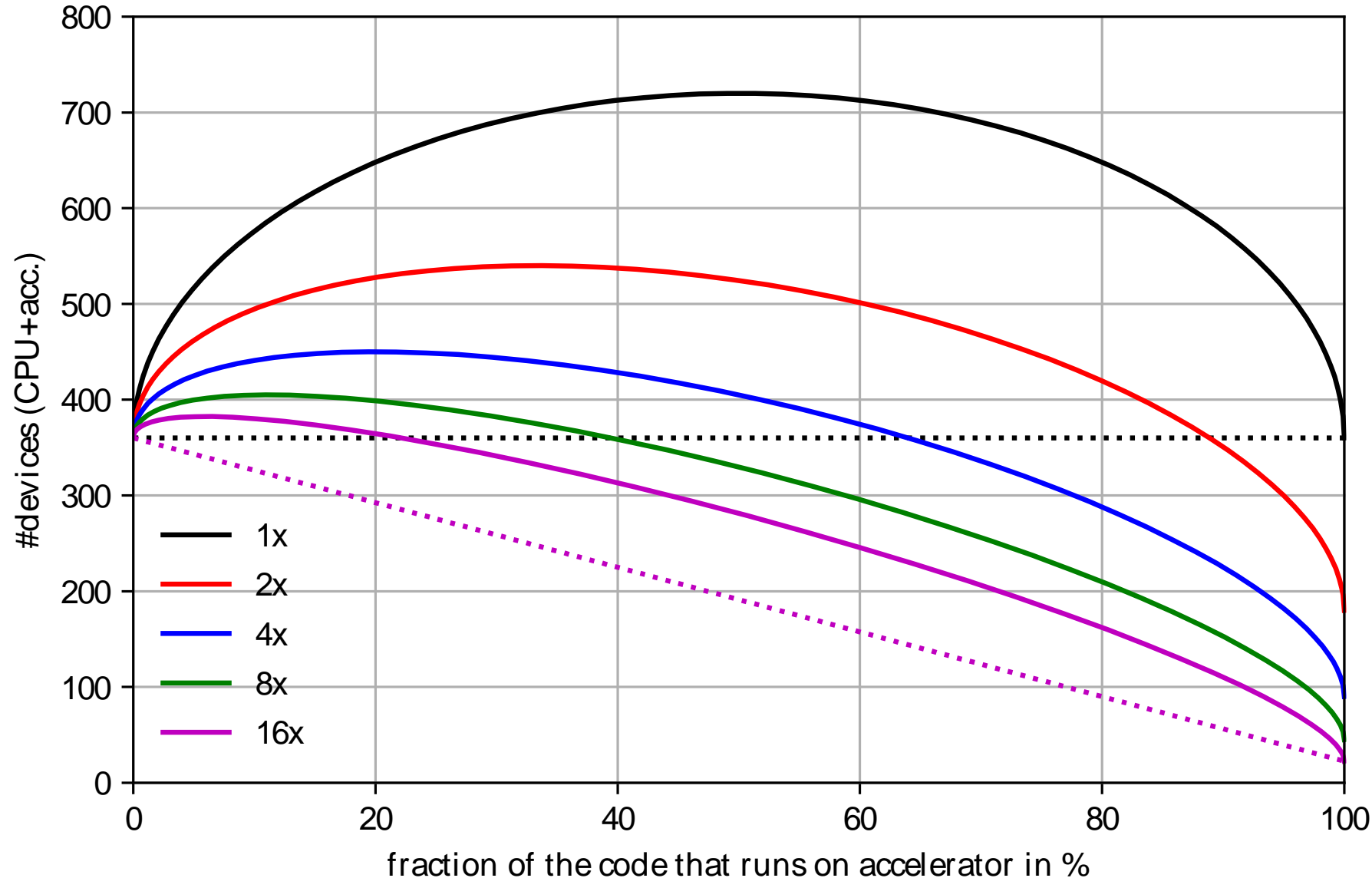




Funded by the European Union

Benefit of accelerators – theoretical model number of devices (acc.+CPU) at MétéoFrance

ESCAPE



Assumes:
sequential execution
perfect scalability

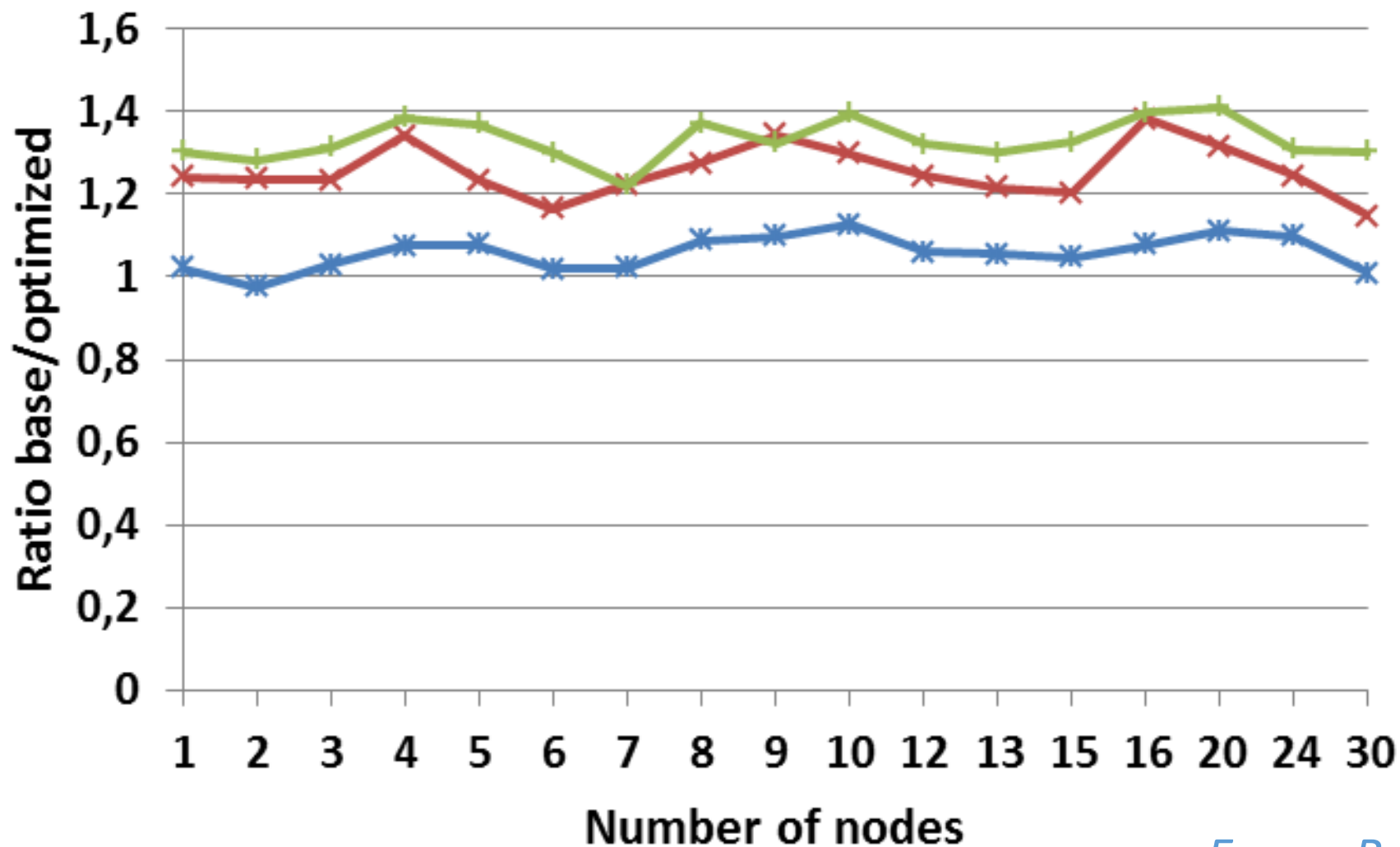
*Philippe Marginaud,
Andreas Mueller*



Funded by the European Union

Spectral transform optimisation by Atos/Bull

ESCAPE



transposition in Fourier transform:
2-3x speedup

- Barrier
- Pack/Unpack
- Both

measurements for the plot: SKX, similar results on Cray at ECMWF

Erwan Raffin

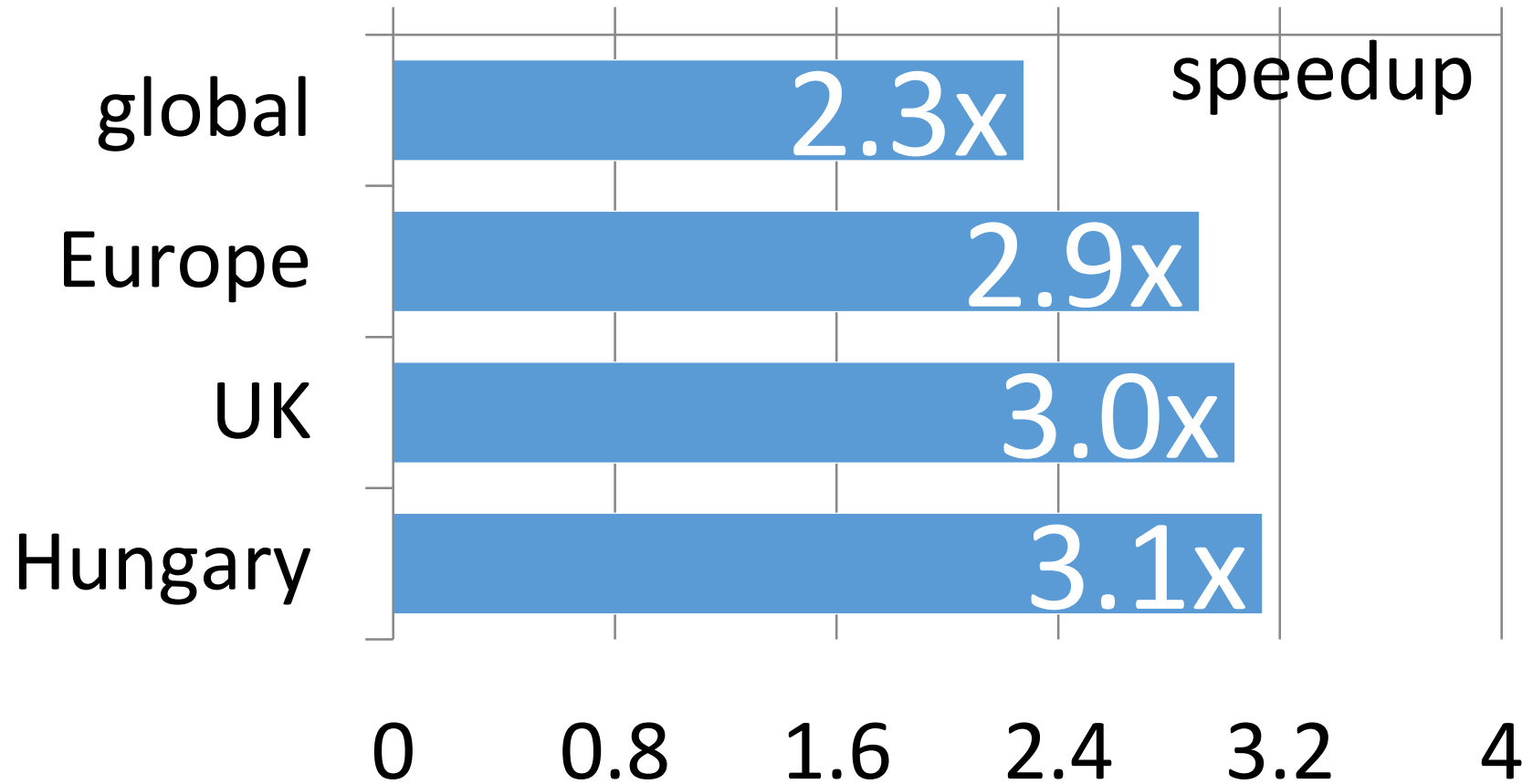


Funded by the European Union

Optimisations in IFS on CPUs postprocessing of spectral data at 9km resolution

ESCAPE

Single CPU



speedup

Andreas Mueller

speedup compared to current operational transform used for postprocessing

Geosci. Model Dev., 11, 1665–1681, 2018
<https://doi.org/10.5194/gmd-11-1665-2018>
© Author(s) 2018. This work is distributed under
the Creative Commons Attribution 4.0 License.



Geoscientific
Model Development

Open Access

Near-global climate simulation at 1 km resolution: establishing a performance baseline on 4888 GPUs with COSMO 5.0

Oliver Fuhrer¹, Tarun Chadha², Torsten Hoefler³, Grzegorz Kwasniewski³, Xavier Lapillonne¹, David Leutwyler⁴, Daniel Lüthi⁴, Carlos Osuna¹, Christoph Schär⁴, Thomas C. Schulthess^{5,6}, and Hannes Vogt⁶

¹Federal Institute of Meteorology and Climatology, MeteoSwiss, Zurich, Switzerland

²ITS Research Informatics, ETH Zurich, Switzerland

³Scalable Parallel Computing Lab, ETH Zurich, Switzerland

⁴Institute for Atmospheric and Climate Science, ETH Zurich, Switzerland

⁵Institute for Theoretical Physics, ETH Zurich, Switzerland

⁶Swiss National Supercomputing Centre, CSCS, Lugano, Switzerland

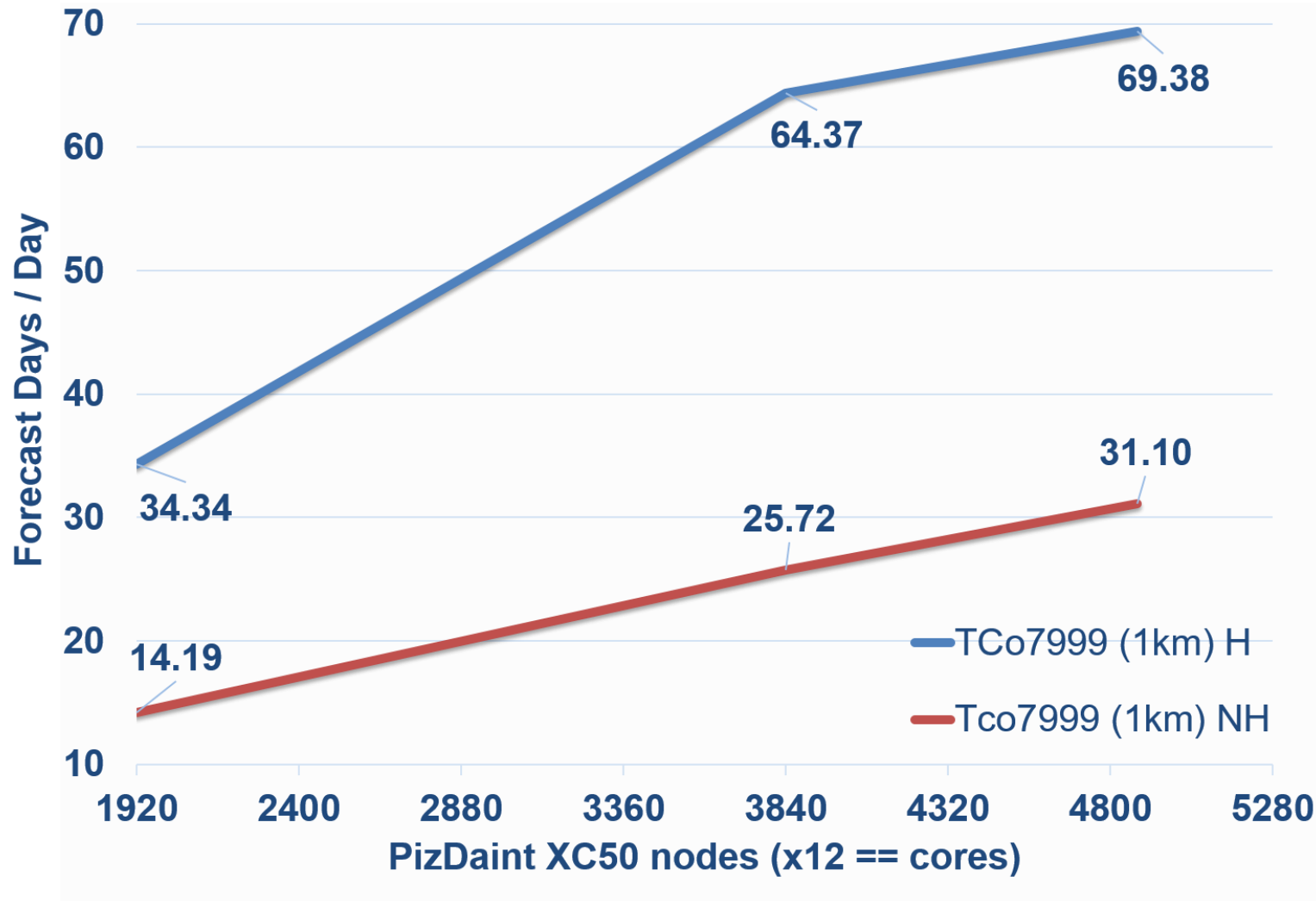
Correspondence: Oliver Fuhrer (oliver.fuhrer@meteoswiss.ch)

Received: 16 September 2017 – Discussion started: 5 October 2017

Revised: 7 February 2018 – Accepted: 8 February 2018 – Published: 2 May 2018

IFS 1km: strong scaling on PizDaint

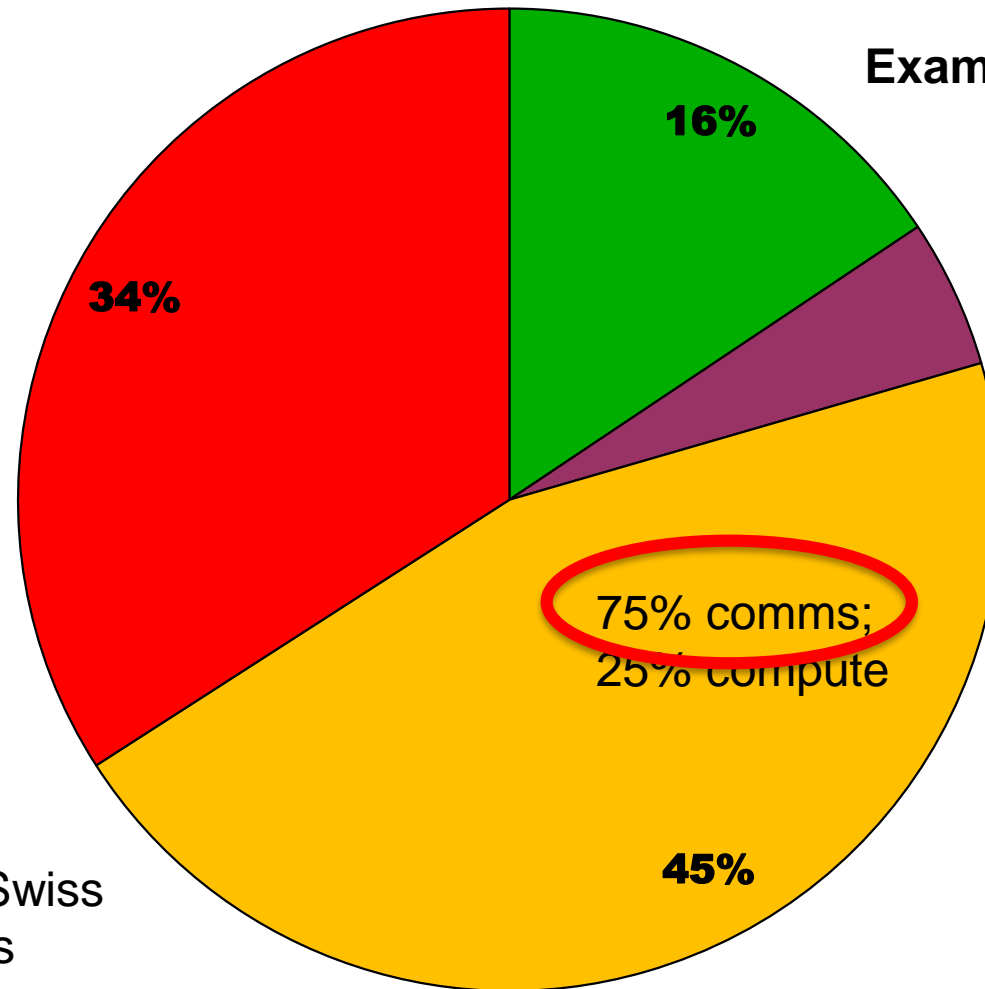
Goal ~1 year / Day



Result of algorithmic changes and single precision

Many thanks to
Thomas Schulthess &
Maria Grazia Giuffreda !

The cost profile of a 1.25km IFS atmosphere simulation on Piz Daint (CPU only)



Example: TCo7999 L62 (~1.25km)

4880 MPI tasks x 12 threads

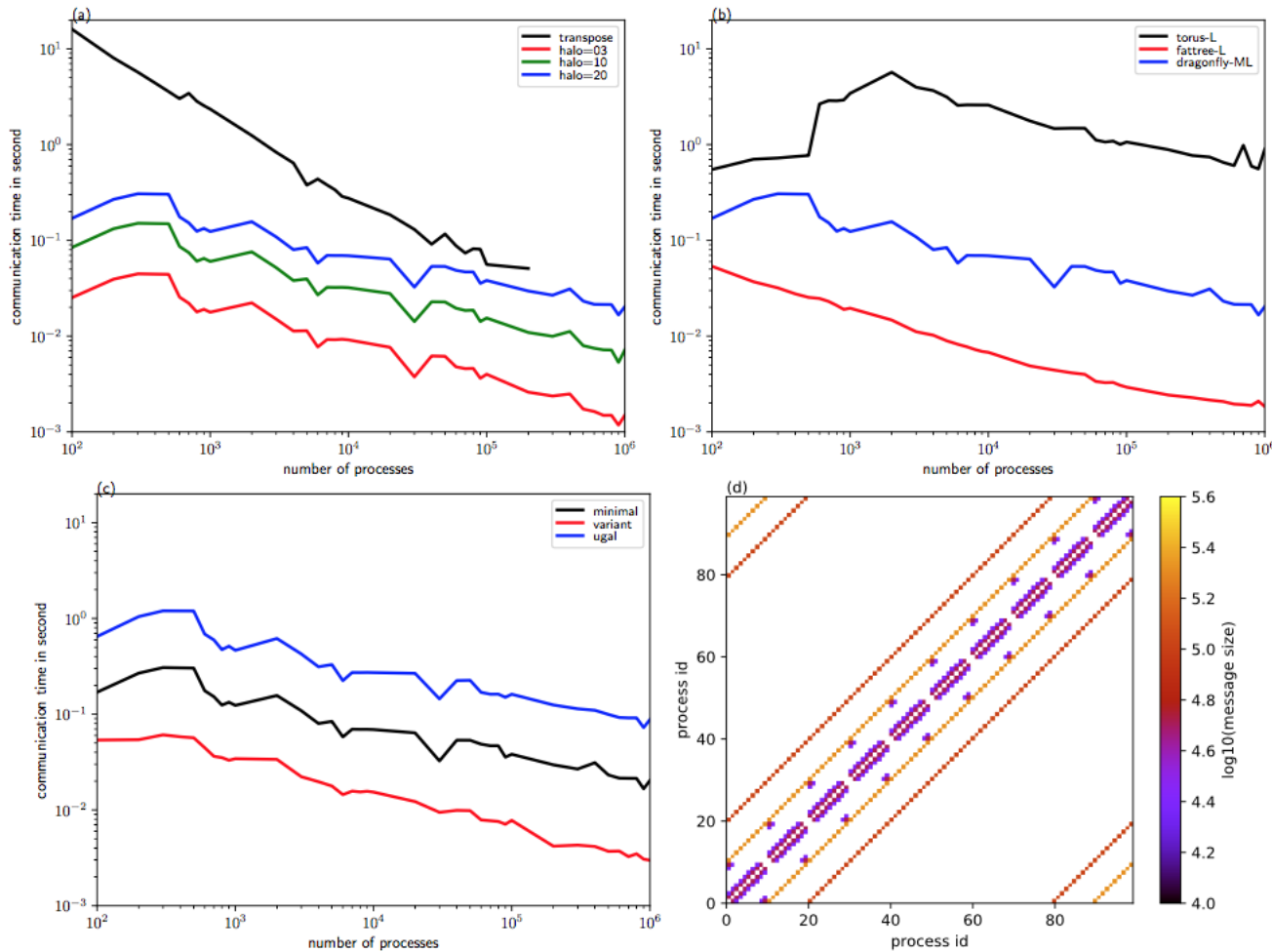
69 FC/day ~ 0.19 SYPD

single precision / FLT

~85.21 MWh / SY

Based on the Piz Daint, Swiss Cray XC50 Haswell, Aries interconnect, ~5000 nodes total

Simulating performance and scalability of MPI communications

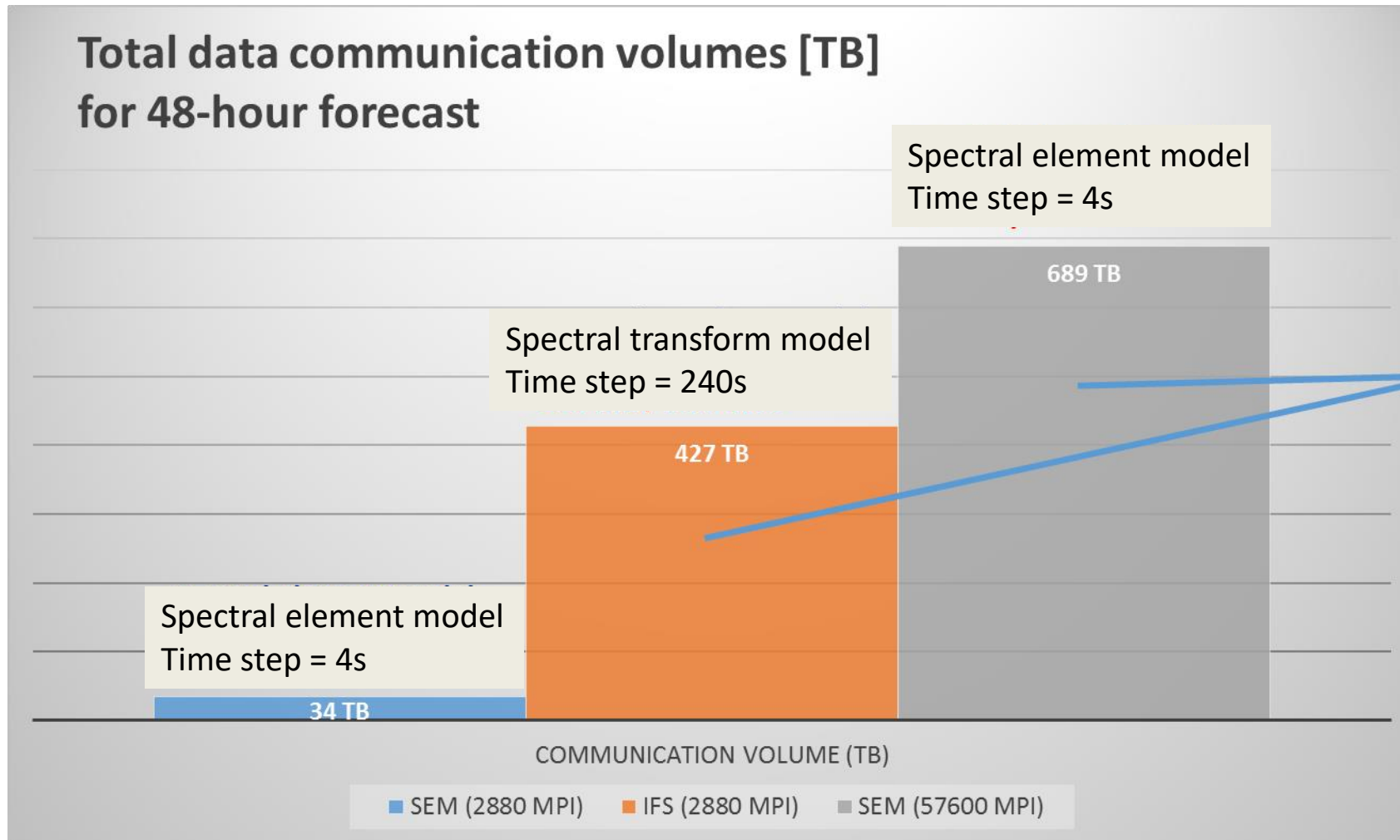


Communication time as a function of halo size, topology & routing algorithm, compared to spectral transpositions

Zheng and Marginaud (GMD, 2018)

MPI collectives at least across a subset of nodes are required in NWP & climate!

Communication is bad – small time steps are worse



Same time to solution!
Energy efficiency?

Data movement x100 (x1000)
more expensive than
computations in time (energy)!

[Shalf et al. 2011]



Funded by the
European Union

Co-ordinated by  **ECMWF**

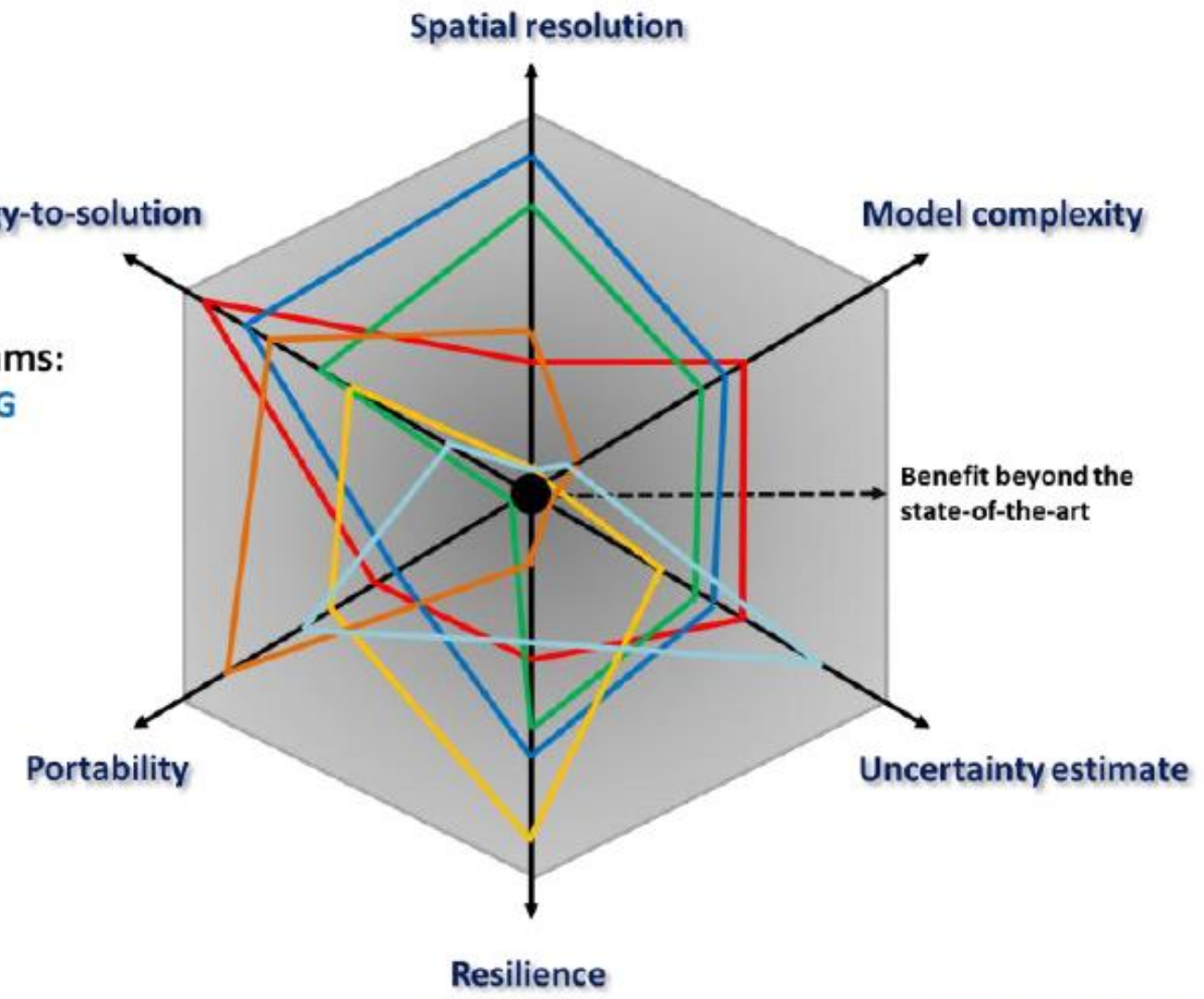
ESCAPE 2

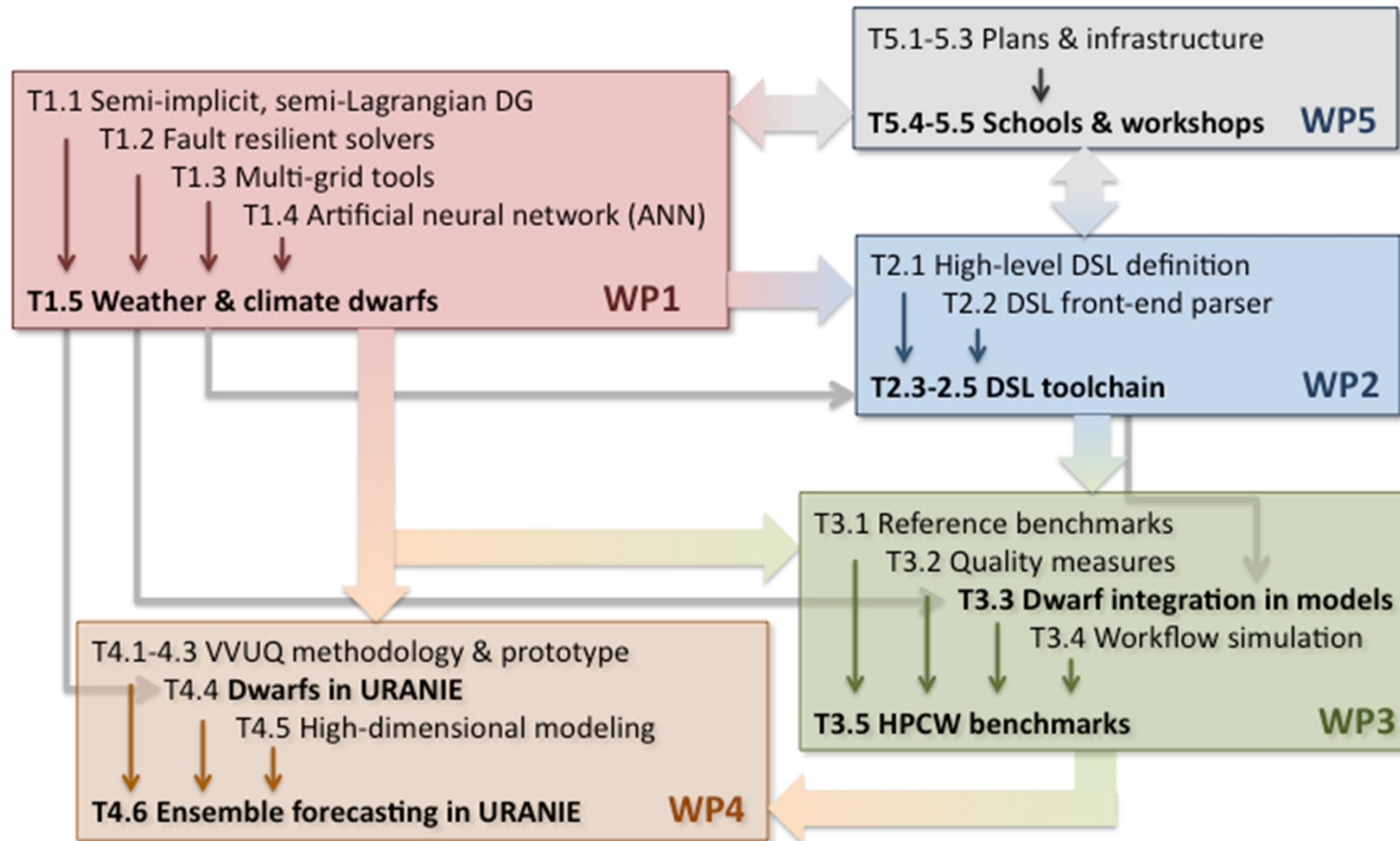




Mathematical methods and algorithms:

- Semi-implicit, semi-Lagrangian CG/DG
 - Hierarchical multigrid tools
 - Fault resilient solver
 - Artificial neural networks
- and:
- DSL toolchain
 - Ensemble based URANIE
- State-of-the-art





ESCAPE Summary



- Numerical weather prediction & climate needs sustained efforts to evolve together with emerging computing opportunities
- *ESCAPE and ESCAPE-2 will deliver*
 - Pioneering approaches for refactoring society critical legacy codes
 - Weather & climate dwarfs
 - Energy-efficient accelerator use in global weather & climate prediction
 - Scrutiny of the need for precision
 - Co-development of novel mathematical algorithms & hardware adaptation
 - Pioneering mathematical algorithm development with hardware adaptation using DSL toolchains
 - A HPCW benchmark and cross-disciplinary Verification, Validation, and Uncertainty Quantification (VVUQ)
 - Application resilience





48h forecast ~9km

Take the “Turing test” of climate & weather modelling (T. Palmer)

<http://gigapan.com/gigapans/206287>

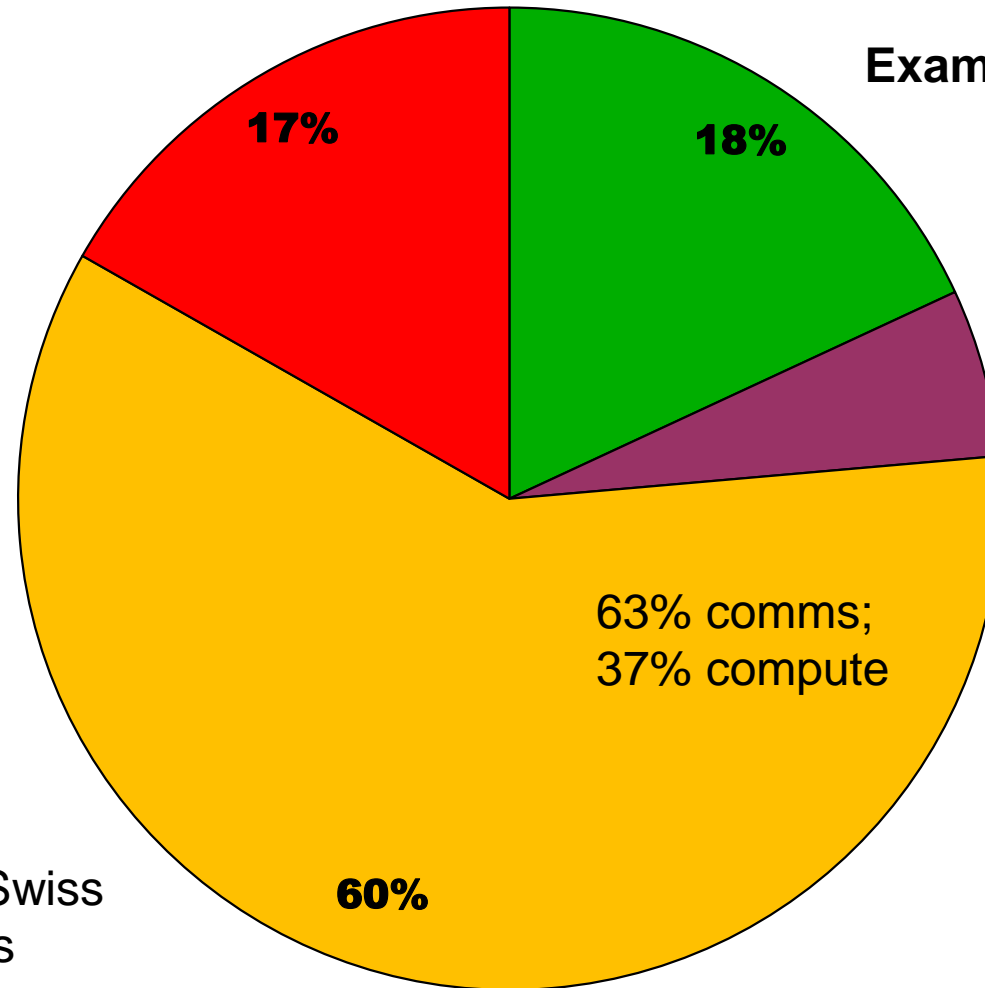
48h forecast ~1km

Additional slides

The cost profile of a 1.25km (non-hydrostatic) IFS atmosphere simulation Piz Daint



Example: TCo7999 L62 (~1.25km)



4880 MPI tasks x 12 threads

32 FC/day ~ 0.088 SYPD

single precision / FLT

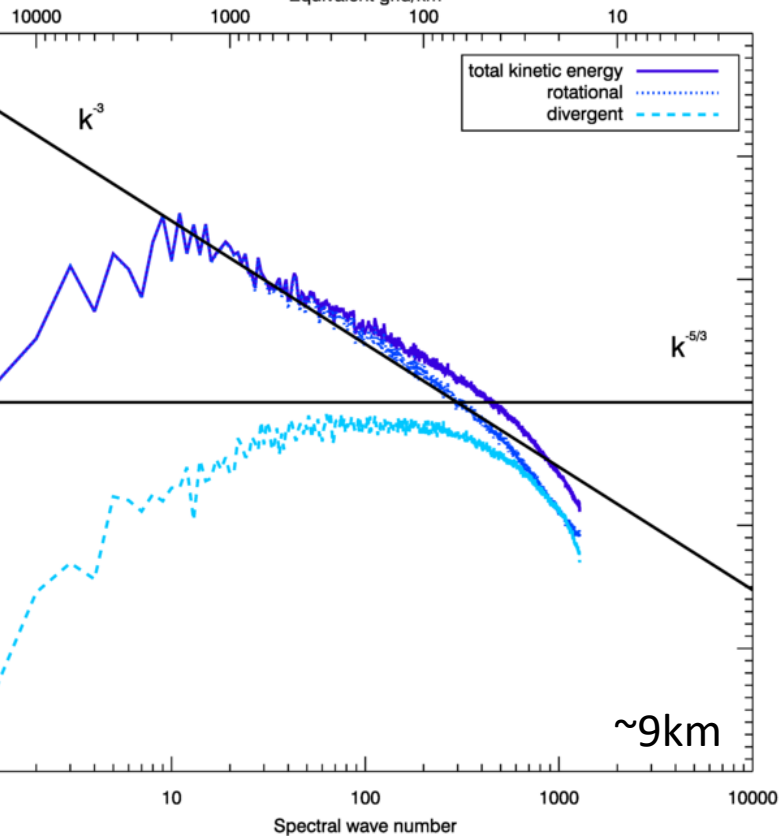
~191.74 MWh / SY

Based on the Piz Daint, Swiss Cray XC50 Haswell, Aries interconnect, ~5000 nodes total

Global KE - Spectra ~500hPa

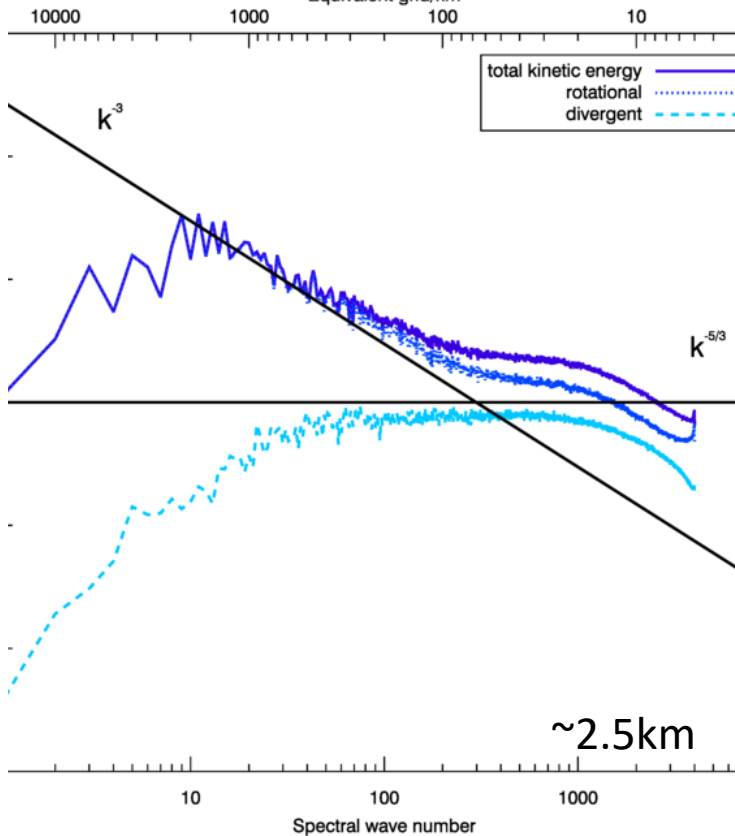
Horizontal kinetic energy spectra

Equivalent grid/km



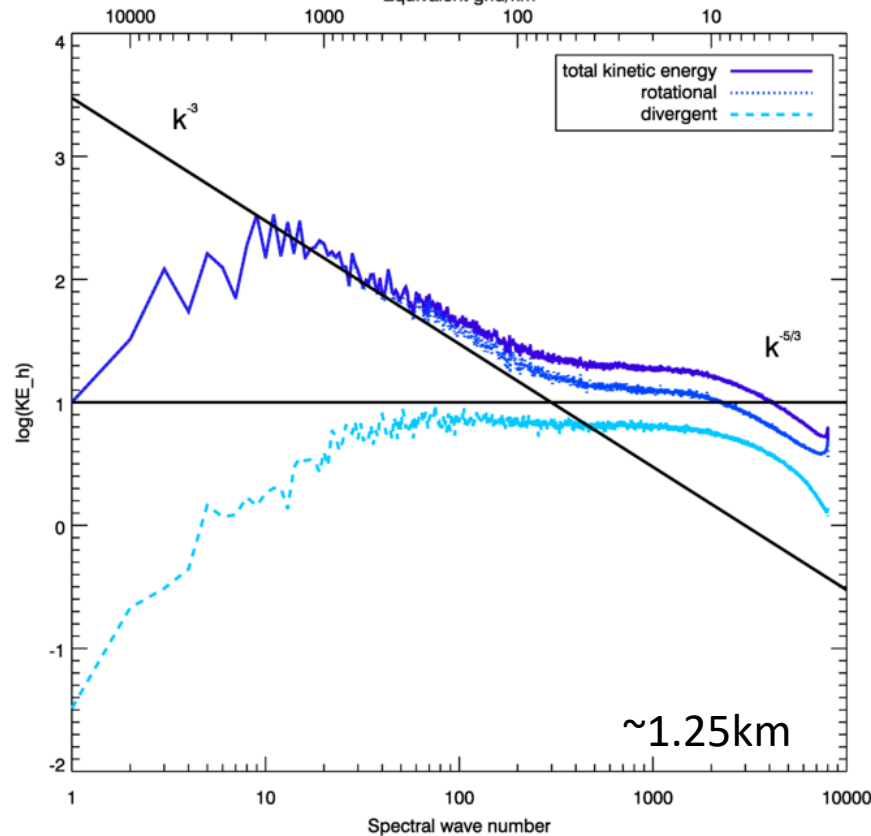
Horizontal kinetic energy spectra

Equivalent grid/km



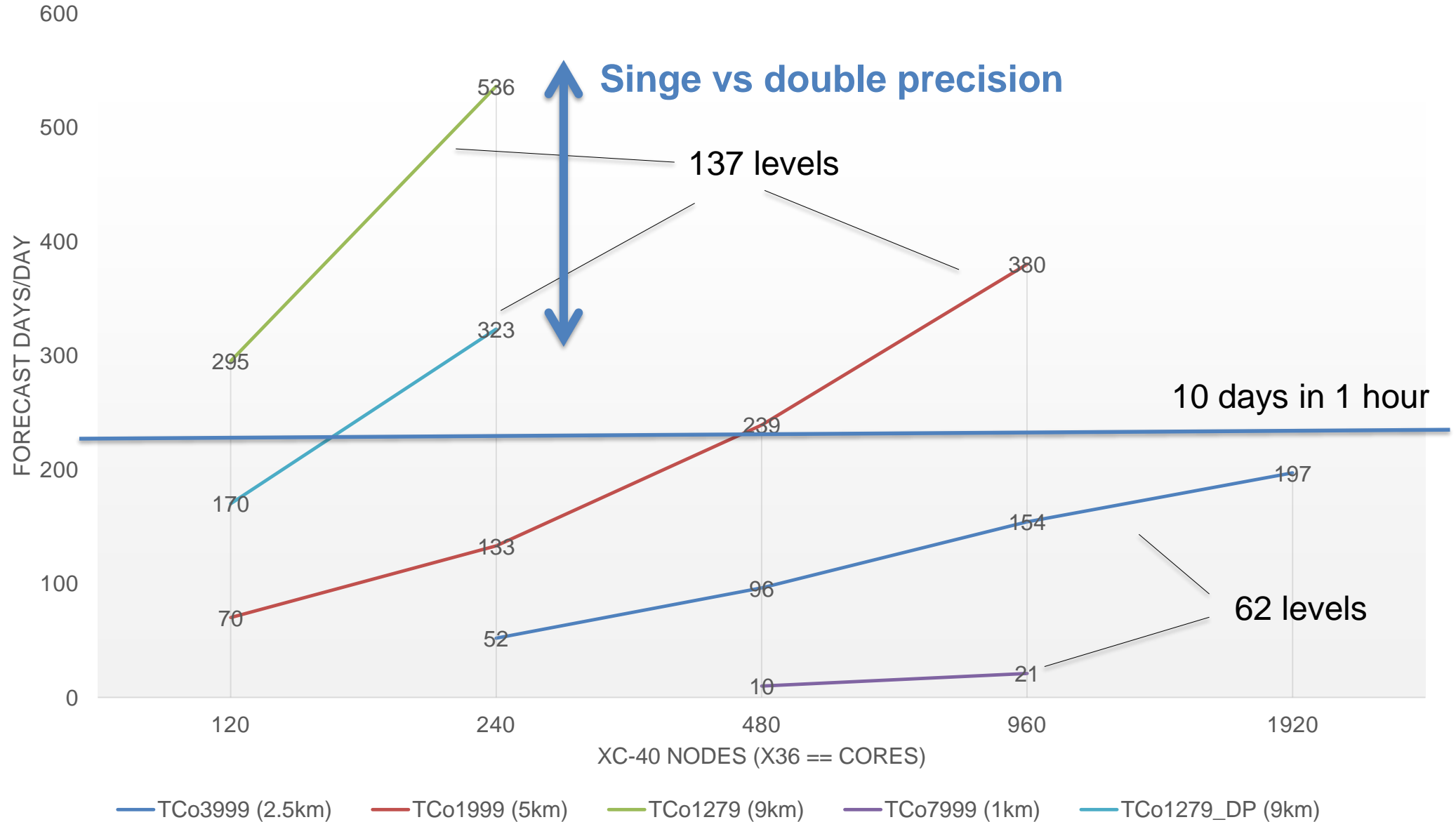
Horizontal kinetic energy spectra

Equivalent grid/km



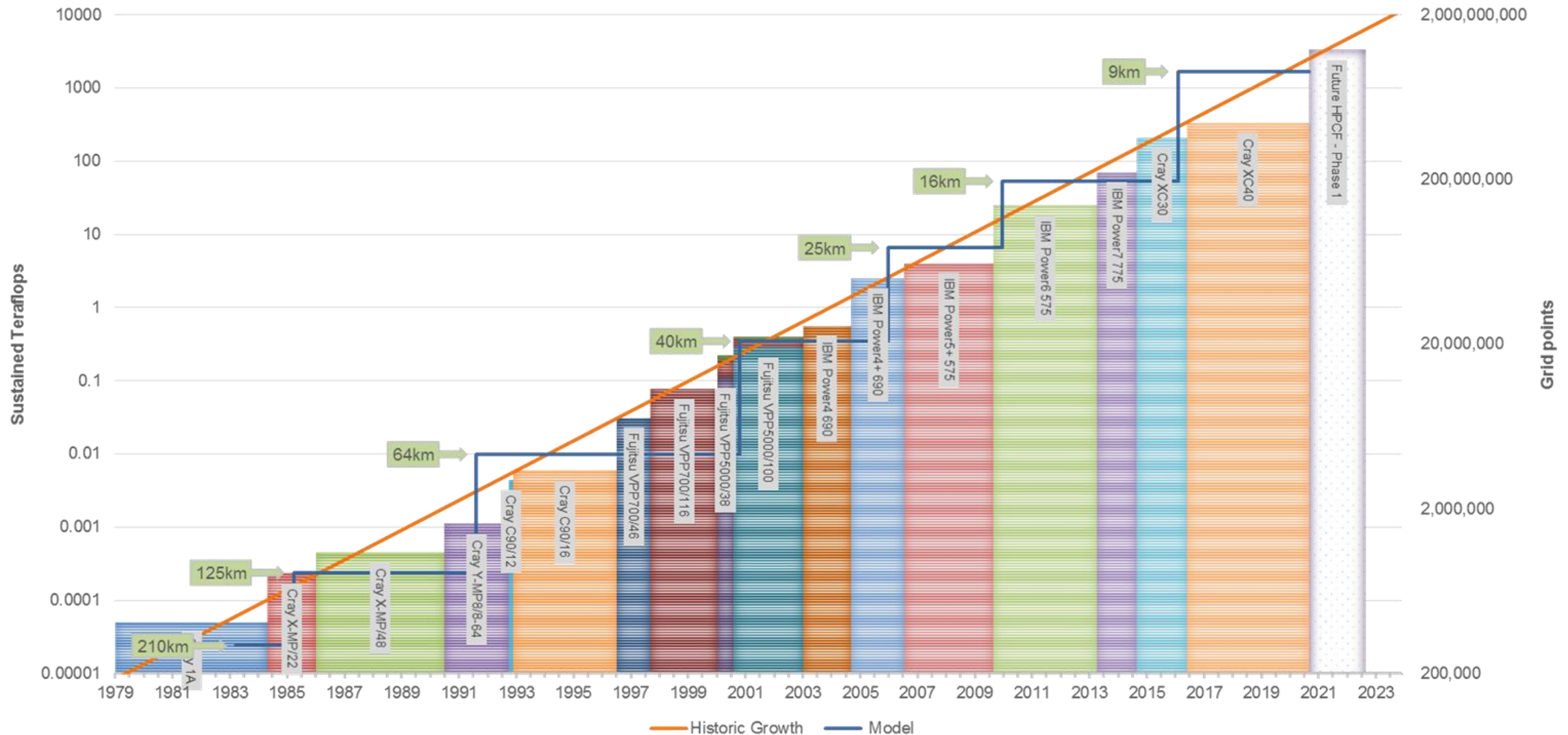
Resolve rather than parametrize much of the crucial vertical transport of momentum and heat

IFS single precision performance – Atmosphere only (no I/O)



(Vana, Dueben et al 2017)

Sustained HPC performance



Ensemble of assimilations and forecasts

