

Second revision of the paper “The ESCAPE project: Energy-efficient Scalable Algorithms for Weather Prediction at Exascale”

Response to the referee 1

(grey background: text of the reviewer comment, white background: our response)

This is a second round of review of the proposed GMD paper summarising the ESCAPE project.

As the main conclusion from my earlier review of the initial version, I suggested to significantly reduce in scope and shorten the manuscript (50 pages, 30 figures). The length was clearly not the only deficiency, and the other two reviewers concurred. Yet, given the apt match between the ESCAPE project aims and the scope of GMD, the journal does seem to be a valid venue for promoting the project and allowing to give credit to its participants.

The current version is noticeably more readable, and indeed significantly shorter (28 pages, 13 figures), yet I still find that the work lacks some balance and style consistency. For instance, I doubt if the paper is a proper place to explain that one should care about index order in multi-dimensional loops or that temporary memory allocation should be done once and not within loops, or to verbosely enumerate parameters of new NVidia products (all on page 13). Similarly, I do not think there is any value in explaining to the readers that program performance benefited from removal of debug leftovers (p17/l5-6) or to promote machinelearning-technology suitability of optical processors (p20/l21-22). I recommend removing Figure 1 (if not, please reduce its resolution, currently page 5 weights over 2M, while the whole article pdf is less than 6M). Overall, please try not to limit further corrections to what is literally brought up in the reviews, but read through the whole text, and try to maintain a consistent style (e.g., passive voice vs. we-did narrative). Do not hesitate to shorten the paper even more.

We thank the reviewer for the thorough reading of the revised manuscript and the helpful review. We streamlined and shortened the presentation, and carefully re-read the whole paper while improving the balance and consistency of the presentation. Although we trimmed the original text quite substantially, we also had to add some text to address the reviewers' comments, so the revised paper is only one page shorter than its previous version. We also reduced the size of Figure 1 substantially. The entire paper has now a size of less than 2MB.

Please also make sure that the abstract and conclusions convey correctly the main lessons learnt in the project and a take-home message from the paper. For example, DSLs and optical processors are not mentioned in the abstract. In the "Conclusions and outlook" section, I suggest to limit references to ESCAPE-2 to one single paragraph (perhaps the last one).

We revised the conclusions and added a paragraph including DSLs and optical processors to the abstract.

The key remaining major concern is the necessity to deposit many of the referenced non-journal reports in persistent repositories. The current practice of solely providing links to pdfs on Google Drive with no metadata or guarantee of persistency is unacceptable (9 such references).

We published all of these documents on arxiv.org and adjusted the references in the paper accordingly.

Below, I'm listing more specific comments that might be useful.

p1/l6: what (in this context) "scientifically required" means? consider rephrasing

Scientifically required accuracy is for us the level of accuracy required to adequately represent the physical processes. We revised this sentence in the paper.

p1/l7: "algorithms should ... be ... resilient in case of ... software failure" – is it discussed, please elaborate what is meant here by software failure?

We removed the term "software".

p2/l1: "is leading" or "was leading"? (puzzling given references to ongoing ESCAPE-2)

We changed this to "was leading".

p2/l3: "under initiative FET-HPC" → "under the FET-HPC initiative"

We revised the text.

p2/l4: "goal is" or "goal was" (as above)

We decided to use "goal was" since the paper focusses on the ESCAPE 1 project which has ended.

p2/l4: "to next generation" → "to the next generation"

p2/l18: "nearly all" → "multiple"?

p2/l23: "\$US" → "US dollars"

p2/l28: "by so-called" → "by the so-called"

p3/l16: does "above" refer to resolution or location in the text, please rephrase

p3/l22: "thanks to" → "owing to"

p3/l24: "As this ... down," → "Consequently,"

We revised the text as suggested in these seven comments.

p3/l29: explain what a "highly varying kernel" is, generally it would be worth to introduce the notion of kernel in the paper

The term kernel refers to compute kernels. We explain this now in the paper.

p3/l29: what is "scientific dependency"? (data dependency?)

Scientific dependency meant the dependency due to data as well as due to physical processes. We revised the text accordingly.

p3/l29: "rather complex algorithms" - please be more specific

p3/l33: "algorithms and codes" → "algorithms and their implementations"

p3/l34: "scientifically and computationally well defined" - "well defined" reads like "well posed" and suggests that other components are not "well defined", consider rephrasing

p4/l1: add "herein introduced" before first mention of Weather & Climate dwarfs?

p4/l7: comma after 2004

p4/l3: motives → motifs

p4/l24: runnable → standalone?

p4/l1: "dwarfs created" → "dwarfs developed" or even better "List of ESCAPE dwarfs"?

p4/l4-5: "For many of the dwarfs we created so called prototypes. Each prototype implementation ..." → "Prototype implementations addressing specific hardware were developed for selected dwarfs"?

p5/l2: "and DSL..." → "and two domain-specific language (DSL) solutions: CLAW and GridTools."

p5/l2: add a final sentence mentioning usage of Fortran

p6/l12: "dwarfing" ??? (perhaps isolating)

p8/l19: comma after methods

We followed these 13 suggestions and revised the text accordingly.

p8/l25: "truncation" → "truncation error"

We revised the text. The truncation in this context is not related to the error but rather to the number of wavenumbers included in the transform and the resulting problem size.

p8/l25: "very cheap direct solve" → "efficient direct scheme"?

p8/l26: comma after paper

p8/l28: comma after Otherwise

p10/fig3 caption: "is without ocean and waves" → "was run without ocean coupling"?

p10/l11: "overlap different parts of the model" → "enable concurrent execution of different model components"?

We revised the text as suggested in these five comments.

p11/l2 and p11/l6: repetition of "parallelism exposed to the GPU", avoid duplication and explain what it means

The code needs to be restructured in such a way that the GPU is enabled to take full advantage of the inherent parallelism of the algorithm. We revised the text to make this clearer.

p11/l11: ensure listing is placed after "as follows" (or make it a figure)

We removed this listing in order shorten the paper. We believe the message of the listing is well described by the text.

p11/l11: "a factor 10", rephrase using "tenfold increase of"

p12/l3: "a FFT" → "an FFT"?

We revised the text as suggested in these two comments.

p15/l11: please remove the "Also there is still room ..." sentence

p17/l5-6: please remove the "A few percent of ..." sentence

We removed these two sentences.

p17/l13: comma after CPU

p18/l2-4: suggest starting the section with the second sentence ("The fundamental...")

p20/l1: "introducing known terms into the functions" – what terms, what functions?

We revised the text as suggested in these three comments.

p20/l15-19: explain that binary means black-white (right?), otherwise binary precision might sound confusing

Binary precision means that only bright and dark values are distinguished. We revised the text.

p20/l21-22: remove the "The inherent ability ..." sentence

We removed this sentence.

p20/l23: please mention that results from the optical processing tests will be discussed later, otherwise it seems that the subject is abruptly dropped

We revised the text.

p21/l2: "is about 15 percent" of what?

The 15 percent refer to the portion of the overall runtime. We revised the text.

p21/l9-10: please remove "There might still be a potential" sentence

We removed this sentence.

p21/l13: rephrase not to repeat "necessary"

We replaced it with "required".

p22/l2: "this section" → "preceding section"

In Latex terminology it is the preceding subsection but still the same section. We added the section number to be precise in our text.

p22/l18-19: please remove "More work on" sentence

We removed this sentence.

p22/l22: "physics equations"?, please rephrase

We revised the text.

p22/l33: "scientific correctness" of dwarfs, implementations, hardware? (please limit the use of them/their/these)

The scientific correctness referred to the implementations. We revised the text.

p23/l4-5: "The paper gives" sentence sounds like a copy-paste from an abstract

We revised this sentence.

p23/l32: Please make the "Comparing different methods..." sentence more specific

We revised the text.

general: please do not CAPITALISE Fortran – starting with Fortran 90 the all-caps name has been dropped

We changed all occurrences of Fortran.

References (in general, be consistent: capitalise only the first word of a title; use proper journal name abbreviations):

Asanovic et al. 2006: www2 → www (both work, but people seem to cite the latter)

Asanovic et al. 2009: Comm. ACM

Benard and Ginton 2019: Q. J. Royal Meteorol. Soc.; add volume, pages

Clement et al. 2018: remove "on"

Colavolpe et al. 2017: Q. J. Royal Meteorol. Soc.

Colella 2004: if impossible to locate, mention in the text after whom cited

Deconinck 2017a: use persistent repository (e.g., OpenAIRE, arXiv), NOT A GOOGLE DRIVE LINK!

Deconinck 2017b: use persistent repository (e.g., OpenAIRE, arXiv), NOT A GOOGLE DRIVE LINK!

Deconinck et al. 2017: Comput. Phys. Commun.

Douriez et al. 2018: use persistent repository (e.g., OpenAIRE, arXiv), NOT A GOOGLE DRIVE LINK!

Dziekan et al. 2019: GMDD → Geosci. Model Dev. 12, 2587-2606; doi:10.5194/gmd-12-2587-2019

Feng et al. 2012: use full title (incl. "work in progress"); capital letters in conf. name

Flamm 2018: add publisher and report number (NBER Working Paper No. 24553)

Glinton and Benard 2019: Q. J. Royal Meteorol. Soc.; add volume, pages (or remove elsewhere)

Johnston and Milthorpe 2018: use full title (incl. "Enhancing OpenCL...");

Kaltofen 2011: add quotation marks within title as in the original; use booktitle "Numerical and Symbolic Scientific Computing" instead of book series "Texts & Monographs in Symbolic Computation"; add editor names; remove Vienna

Katzav and Parker 2015: Clim. Change

Krommydas et al. 2015: remove "chun" in the surname of second author; J. Signal Process. Syst.

Kuhnlein et al. 2019: Geosci. Model Dev.

Macfaden et al 2017: Sci. Rep.

Mazauric et al. 2017a: use persistent repository (e.g., OpenAIRE, arXiv), NOT A GOOGLE DRIVE LINK!

Mazauric et al. 2017b: use persistent repository (e.g., OpenAIRE, arXiv), NOT A GOOGLE DRIVE LINK!

Mengaldo 2016: use persistent repository (e.g., OpenAIRE, arXiv), NOT A GOOGLE DRIVE LINK!

Mengaldo et al. 2018: Arch. Comput. Methods Eng.

Messer et al. 2016: Int. J. High Perform. Comput. Appl

Michalakes et al. 2015: is the "Tech. Rep. TN-484, NCAR, Boulder" correct? NCAR does not seem to list it, please use a (hopefully) more persistent url: <https://repository.library.noaa.gov/view/noaa/18654>

Muller et al. 2017: use persistent repository (e.g., OpenAIRE, arXiv), NOT A GOOGLE DRIVE LINK!

Muller et al. 2018: add volume, pages

Neumann et al. 2019: Philos. Trans. Royal Soc. A; remove 20180 148

Osuna 2018: use persistent repository (e.g., OpenAIRE, arXiv), NOT A GOOGLE DRIVE LINK!

Robinson et al. 2016: Prace → PRACE; add doi: 10.5281/zenodo.832025; remove "- Evaluations on Intel MIC?";

Schalkwijk et al. 2015: Bull. Am. Meteorol. Soc.

Schulthess et al. 2019: Comput. Sci. Eng.

Shukla et al. 2010: Bull. Am. Meteorol. Soc.

Van Bever et al. 2018: use persistent repository (e.g., OpenAIRE, arXiv), NOT A GOOGLE DRIVE LINK!

Wallemacq et al. 2018: correct author list (Wallemacq, P. and House, R.); use publisher's url instead of ResearchGate doi: <https://www.unisdr.org/we/inform/publications/61119>

Wedi et al. 2015: add doi: 10.21957/thtpwp67e

Wehner et al. 2011: J. Adv. Model. Earth Sys.

Xiao et al. 2017: add doi:10.21957/g9mjllgeq

Zheng 2018: correct volume, pages, doi: Geosci. Model Dev., 11, 3409-3426, doi:10.5194/gmd-11-3409-2018

We made all of the suggested changes and we uploaded the technical ESCAPE reports to arxiv.org.

Hope that helps.

Response to the referee 2

(grey background: text of the reviewer comment, white background: our response)

We thank the reviewer for the careful reading of the revised manuscript and the helpful comments. We accommodated all the comments in the letter or spirit. We also streamlined and shortened the presentation, and carefully re-read the whole paper with the aim to improve the presentation following guidance of both reviewers' comments. Please find the detailed point-by-point reply below.

In the abstract.

“Satisfy strict service requirements in terms of time-to-solution and energy-to-solution”. This reviewer is unaware of any service agreements for energy to solution. For any given machine, there is a time-to-solution requirement. Exploiting architectures which are more energy efficient for a given time-to-solution are explored here, but that is not the same thing. Can the authors clarify what they mean or change this sentence.

How much money can be spent on energy is limited by the budget of the institution which is paying for the supercomputer. The energy consumption has a direct impact on the decision which supercomputer the institution can afford to buy. We revised the text to make this aspect clearer.

Page 3 Line 29 “highly varying kernels”, What does this mean. Also the sentence is rather long and unwieldy and could benefit from restructuring.

Different compute kernels used in weather and climate prediction models have very different computational characteristics. We revised the text to make this clearer.

Page 4 Line 13 Motifs is a plural of Motif. English is an elastic language whose usage does change and evolve and it maybe that motives is becoming an accepted plural of motif (I have seen this usage elsewhere). It could be confused motives meaning reasons for doing something. This reviewer suggests using Motifs for clarity.

Page 10, last sentence. “The only true solution ...” this is an odd phrase, perhaps “One approach to avoiding” Is better.

We revised the text as suggested in these two comments.

Page 10/11 and figure 4. The results of Michalakes et al 2015 are used to show the amount of data that has to be communicated. Is this data from model run, or is it based on scaling from a smaller run in the paper. If it is the latter, it needs to be made explicitly clear and again, if it is the latter, is it the scaling model from the Michalakes paper?

This figure is extrapolating the results from Michalakes. We revised the text to make this clearer.

Page 11, line 11. This sentence doesn't appear to connect with anything else and should be removed.

The sentence belonged to the listing which by accident got moved to the next page. We removed the code example in order to shorten the paper. We believe that the text explains the message of the example well enough.

Page 12 figure 5.

The red open circle and green open diamond sit atop one another. This makes them hard to see and distinguish. Can this be described in the caption to make it clear. How is the operational intensity determined for a whole time-step? What was the non-optimised performance of the Kernels only and Matmult.

We revised the caption. The operational intensity was derived from measurements with nvprof as well as from counting the number of floating point operations by hand. Combining the counted floating point operations with the measured memory traffic gives us the operational intensity. Numbers for the non-optimised kernels only are not readily available.

Page 17 figure 9a

Why is the x-axis not “node valued”? i.e. 1 and 2 are separated by 1 unit but so is 20-24 and 24-30. Figure 9b has a correct scale.

We corrected Figure 9a.

Section 3.6 What is the conclusion to the work on the comparison between processors for run-time and energy consumption?

We added a statement about the conclusion of this comparison.

Page 22 Line 17

The DSL allows to perform optimisations → the DSL allows optimisations to be performed

Page 23 lines 22 and 23

Referring to papers with “like in” is rather informal language for a scientific paper.

See, or see for example is better.

Page 23 line 32

Requires to include all costs → requires all costs to be included.

We revised the text as suggested in these three comments.

The ESCAPE project: Energy-efficient Scalable Algorithms for Weather Prediction at Exascale

Andreas Müller¹, Willem Deconinck¹, Christian Kühnlein¹, Gianmarco Mengaldo¹, Michael Lange¹, Nils Wedi¹, Peter Bauer¹, Piotr K. Smolarkiewicz¹, Michail Diamantakis¹, Sarah-Jane Lock¹, Mats Hamrud¹, Sami Saarinen¹, George Mozdzyński¹, Daniel Thiemert¹, Michael Glinton², Pierre Bénard², Fabrice Voitus², Charles Colavolpe², Philippe Marguinaud², Yongjun Zheng², Joris Van Bever³, Daan Degrauwe³, Geert Smet³, Piet Termonia^{3,4}, Kristian P. Nielsen⁵, Bent H. Sass⁵, Jacob W. Poulsen⁵, Per Berg⁵, Carlos Osuna⁶, Oliver Fuhrer⁶, Valentin Clement⁷, Michael Baldauf⁸, Mike Gillard⁹, Joanna Szmelter⁹, Enda O'Brien¹⁰, Alastair McKinstry¹⁰, Oisín Robinson¹⁰, Parijat Shukla¹⁰, Michael Lysaght¹⁰, Michał Kulczewski¹¹, Miłosz Ciznicki¹¹, Wojciech Piątek¹¹, Sebastian Ciesielski¹¹, Marek Błażewicz¹¹, Krzysztof Kurowski¹¹, Marcin Procyk¹¹, Paweł Spychala¹¹, Bartosz Bosak¹¹, Zbigniew P. Piotrowski¹², Andrzej Wyszogrodzki¹², Erwan Raffin¹³, Cyril Mazaurec¹³, David Guibert¹³, Louis Douriez¹³, Xavier Vigouroux¹³, Alan Gray¹⁴, Peter Messmer¹⁴, Alexander J. Macfaden¹⁵, and Nick New¹⁵

¹European Centre for Medium-Range Weather Forecasts (ECMWF), Reading, UK

²Centre National de Recherches Météorologiques, Météo-France, Toulouse, France

³Royal Meteorological Institute (RMI), Brussels, Belgium

⁴Department of Physics and Astronomy, Ghent University, Ghent, Belgium

⁵The Danish Meteorological Institute (DMI), Copenhagen, Denmark

⁶Federal Office of Meteorology and Climatology MeteoSwiss, Zurich, Switzerland

⁷Center for Climate System Modeling, Zurich, Switzerland

⁸Deutscher Wetterdienst (DWD), Offenbach, Germany

⁹Loughborough University, Leicestershire, UK

¹⁰Irish Centre for High-End Computing (ICHEC), National University of Ireland, Galway, Ireland

¹¹Poznan Supercomputing and Networking Center (PSNC), Poznań, Poland

¹²Institute of Meteorology and Water Management - National Research Institute (IMGW-PIB), Warsaw, Poland

¹³ATOS, Bezons, France

¹⁴NVIDIA Switzerland, Zurich, Switzerland

¹⁵Optalysys Ltd., Glasshoughton, UK

Correspondence: Andreas Müller (andreas.mueller@ecmwf.int)

Abstract. In the simulation of complex multi-scale flows arising in weather and climate modelling, one of the biggest challenges is to satisfy strict service requirements in terms of time-to-solution and to satisfy budgetary constraints in terms of energy-to-solution, without compromising the accuracy and stability of the ~~calculation. These competing factors require the development of robust algorithms that can optimally exploit the targeted underlying hardware and efficiently deliver the extreme~~ computational efforts typically required in operational forecast production. These algorithms should: (i) application. These simulations require algorithms that minimise the energy footprint along with the time required to produce a solution; ~~(ii) maintain the scientifically~~, maintain the physically required level of accuracy; ~~and (iii) be numerically stable, and, are numerically stable and are~~ resilient in case of hardware ~~or software~~ failure.

The European Centre for Medium Range Weather Forecasts (ECMWF) ~~is leading a project called~~ was leading the ESCAPE (Energy-efficient SCalable Algorithms for weather Prediction on Exascale supercomputers) ~~which is project,~~ funded by Horizon 2020 (H2020) under ~~initiative~~ the FET-HPC (Future and Emerging Technologies in High Performance Computing ~~(FET-HPC)~~) initiative. The goal of ~~the ESCAPE project is~~ ESCAPE was to develop a sustainable strategy to evolve weather and climate prediction models to the next-generation computing technologies. The project partners incorporate the expertise of leading European regional forecasting consortia, university research, experienced high-performance computing centres and hardware vendors.

This paper presents an overview of ~~this~~ the ESCAPE strategy: (i) identify domain-specific key algorithmic motifs in weather prediction and climate models (Weather & Climate dwarfs), (ii) categorise them in terms of computational and communication patterns, while (iii) adapting them to different hardware architectures with alternative programming models, (iv) analyse the challenges in optimising and (v) ~~finding~~ find alternative algorithms for the same scheme. The participating weather prediction models are: IFS (Integrated Forecasting System); ALARO, a combination of AROME (Application de la Recherche à l'Opérationnel à Meso-Echelle) and ALADIN (Aire Limitée Adaptation Dynamique Développement International); and COSMO-EULAG, a combination of COSMO (Consortium for Small-scale Modeling) and EULAG (Eulerian/semi-Lagrangian fluid solver). For many of the Weather & Climate dwarfs ESCAPE provides prototype implementations on different hardware architectures (mainly Intel Skylake CPUs, NVIDIA GPUs, Intel Xeon Phi, Optalysys optical processor) with different programming models. The spectral transform dwarf represents a detailed example of the co-design cycle of an ESCAPE dwarf.

The dwarfs concept has proven to be extremely useful for rapid prototyping of alternative algorithms and their interaction with hardware; e.g. the use of a domain-specific language (DSL). Manual adaptations have lead to substantial accelerations of key algorithms in numerical weather prediction (NWP), but are not a general recipe for performance portability of complex NWP models. Existing DSLs are found to require further evolution but are promising tools for achieving the latter. Measurements of energy- and time-to-solution suggest that a future focus needs to be on exploiting simultaneous use of all available resources in hybrid CPU/GPU arrangements.

1 Introduction

Numerical weather and climate prediction capabilities represent substantial socio-economic value in ~~nearly all~~ multiple sectors of human society, namely for the mitigation of the impact of extremes, in food production, renewable energy and water management, infrastructure planning, and for finance and insurance where weather sensitive goods and services are traded. Despite significant progress achieved over past decades (Bauer et al., 2015) there are substantial shortcomings in our ability to predict, for example, weather extremes with sufficient lead time and the impact of climate change at regional or national level. Extreme weather events caused over 500 thousand casualties and over 2 trillion ~~\$US~~ US dollars economic damages in the past 20 years (Wallemacq and House, 2018). Another example is the prediction of climate change mitigation and adaptation targets. Failure to meet these targets is ranked among the leading threats to global society (World Economic Forum, 2019).

One of the key sources of model error is limited spatial (and temporal) resolution (Palmer, 2014) which implies that key physical processes that drive global circulation, like deep convection in the tropics and mesoscale eddies in the ocean, are only crudely represented in models by ~~so-called~~ the contemporary parameterisations. In addition, deficiencies in the representation of process interactions between atmosphere and ocean/sea-ice, as well as atmosphere and land, including a time-varying biosphere, are highly relevant strategic targets to improve the representation of the internal variability of the Earth system in models (Katzav and Parker, 2015). However, better spatial resolution and enhanced model complexity translate immediately into significant computational challenges (Schulthess et al., 2019), whereby spatial resolution is the most demanding because a doubling of resolution roughly translates into eight times more computations (doubling the horizontal resolution in both directions and a corresponding decrease in the time step). The critical resolution threshold at which global weather & climate models may eventually be able to overcome the bulk of the limiting deficiencies is unclear; however, O(1km) emerges as a defensible intermediate target (Shukla et al., 2010; Neumann et al., 2019).

The ESCAPE project has been motivated by the need for running Earth-system models at much higher resolution and complexity than presently available, but within the same time-critical path of daily production for weather forecasts and with the same production throughput needed for decadal/centennial climate projections as used today. This translates to computing one simulated year per wallclock day. Energy efficiency is a key requirement as power envelopes for HPC systems cannot be scaled up at the same rate as the computing demand. Obviously, this presents a substantial challenge to all aspects of computing, namely the choice and implementation of numerical methods and algorithms, and the programming models to map memory access and communication patterns onto specific hardware (Wehner et al., 2011).

The ECMWF specific strategic target is to provide skilful predictions of high-impact weather up to two weeks ahead by 2025 based on a global 5 km-scale Earth-system model ensemble, complemented by the corresponding, high-resolution data assimilation system for creating the initial conditions. This system needs sufficient model complexity in the atmosphere (physics and chemistry), oceans, sea-ice and land to represent all processes acting on such scales, and a sufficiently large ensemble size so that complex probability distributions of extreme forecast features can be sampled well enough. In terms of computational demand this is comparable to the ~~above~~-1-km resolution target for a single forecast.

The supercomputers used for numerical weather prediction (NWP) have changed dramatically over the past decades, and ECMWF's Integrated Forecasting System (IFS; Wedi et al., 2015, and references therein) has exhibited rather different levels of sustained performance ranging from 40-50% on the parallel vector-processor Cray-XMP, YMP and C90 architectures between 1984 and 1996, to 30% on Fujitsu's similar-type VPP-700 and 5000 systems between 1996 and 2002, down to 5-10% on the scalar multi-core IBM P5-7 and Cray XC30 and 40 architectures operated between 2002 and today. Despite sustained performance declining, overall FLOP performance increased exponentially thanks-owing to Moore's law, Dennard scaling and processor pricing (e.g. Flamm, 2018) so that significant upgrades of model resolution, complexity and ensemble size remained affordable in terms of capital cost and power supply. ~~As this 'natural'~~ However, as this technological evolution is slowing down, new concepts for designing algorithms and for mapping the associated computational patterns onto the available architectures - from processor to system level - are needed (Lawrence et al., 2018). Many of the algorithms used in NWP were designed well before the multi-core era ~~started~~, but even though they contain highly-tuned shared and distributed memory parallelisation, they

only achieve such limited sustained performance because of ~~poor arithmetic density~~ (i) a poor ratio of data communication /computations to computations (viz. arithmetic density) in some of the ~~highly-varying kernels~~, compute kernels (hereinafter “kernels” for brevity), (ii) sequential tasking due to strong scientific-physical/data dependencies, and ~~rather complex algorithms~~ (iii) algorithmic intricacies required to solve accurately and efficiently a multi-scale and multi-phase fluid dynamics ~~problem~~ problems on a rotating sphere ~~with sufficient accuracy~~.

If the envisioned increase in model fidelity is constrained by only marginally growing power envelopes and decelerating general purpose processor speed, then performance issues need to be addressed at the root, and a more radical redesign of the basic algorithms and ~~codes~~ their implementations used for weather prediction needs to be considered. This is why ESCAPE investigates both HPC adaptation and alternative numerical formulations ~~in scientifically and computationally well defined~~ model components, ~~of physically and computationally clearly identifiable model components - herein introduced~~ the Weather & Climate dwarfs, ~~followed by a detailed analysis of their respective computational bottlenecks, and subsequent hardware and algorithm dependent optimisations.~~

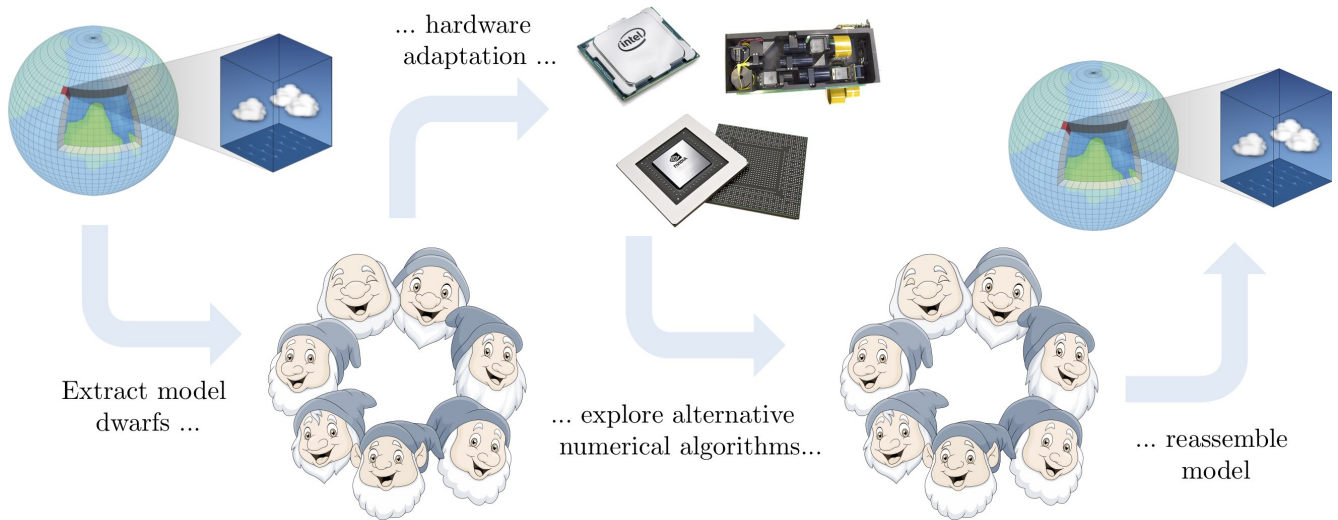
The Weather & Climate dwarf idea is introduced in Section 2. In Section 3 ~~we illustrate the usefulness~~ the usefulness of the dwarf concept is illustrated with the example of the spectral transform dwarf. The paper ends with conclusions and outlook in Section 4.

2 Weather & Climate dwarfs

2.1 Motivation

In 2004, Phillip Colella introduced the seven dwarfs of algorithms for high-end simulation in the physical sciences ~~(Colella, 2004)~~ (Colella, 2004). These were later extended to the 13 Berkeley dwarfs (Asanović et al., 2006, 2009) which are meant to represent characteristic patterns of computation and communication. These dwarfs were created to cover the characteristic computational properties of a broad range of scientific applications. These dwarfs are the basis for the OpenDwarfs benchmark suite (Feng et al., 2012; Krommydas et al., 2015; Johnston and Milthorpe, 2018) and were also applied to benchmark cloud computing in Phillips et al. (2011). In a similar fashion Kaltofen (2011) introduced the seven dwarfs of symbolic computation.

Following this idea, we categorise key algorithmic ~~motives~~ motifs specific to weather prediction and climate models and identify their specific computational and communication patterns, which in return are crucial for the entire model performance. The dwarfs thus represent domain specific mini-applications (Messer et al., 2016) which include direct input from the domain scientist, ~~together with~~ documentation, timers for profiling purposes as well as error estimates for verification purposes. In this way the dwarfs facilitate communication of ~~weather-NWP~~ weather-NWP domain specific knowledge ~~and on~~ algorithmic motifs with specialists in other domains. Different implementations of the dwarfs can be used as a first step towards optimising and adapting weather prediction and climate models to new hardware architectures and to benchmark current and future supercomputers with ~~these simple simpler~~ but relevant applications. Identifying these key algorithms ~~also thus~~ allows better collaboration between operational weather prediction centres, hardware vendors and academia ~~because they are much simpler to understand than the full models~~. The concept of dwarfs is different from the existing separation of weather and climate models into different model



Earth illustration: used under license from GraphicsRF/Shutterstock.com.
 Dwarf illustrations: used under license from Teguh Mujiono/Shutterstock.com

Figure 1. Illustration of the main idea behind the ESCAPE project. The entire model is broken down into smaller building blocks called dwarfs. These are adapted to different hardware architectures. Based on the feedback from hardware vendors and high performance computing centres alternative numerical algorithms are explored. These improvements are eventually built into the operational model.

components, such as atmosphere, land-surface and ocean for which separate dynamical core and physical parameterisation packages already exist. Instead, dwarfs define a runnable-standalone and more manageable sub-component in a hierarchy of model complexity for specific targets such as adaptation to GPUs, exploring alternative programming models, and developing performance portable domain specific languages. But dwarfs can also be used by domain scientists for developing alternative algorithms, even across model components.

The fundamental starting point of the ESCAPE project is to identify the dwarfs in the participating weather and climate models (Figure Fig. 1) and to adapt them to different hardware architectures. The knowledge gained in this adaptation is used to research alternative numerical algorithms which are better suited for those new architectures, and experiment with alternative programming models, towards improving the overall energy-efficiency of weather prediction applications.

10 2.2 List of ~~dwarfs created in~~ ESCAPE dwarfs

Table 1 gives an overview of the dwarfs defined in the ESCAPE project. These dwarfs have been chosen because they represent key algorithms that are representative of any weather and climate model, or because their computational and communication patterns represent one of the most runtime consuming parts of weather forecasting systems (Figure Fig. 2). ~~For many of the dwarfs we created so-called prototypes. Each prototype implementation for a~~ Prototype implementations addressing specific hardware (by default CPU) were developed for selected dwarfs. Each prototype implementation comes with documentation (Mengaldo, 2016; Müller et al., 2017) ~~; error-measures which allow us to quickly see if optimisations affect the results and~~

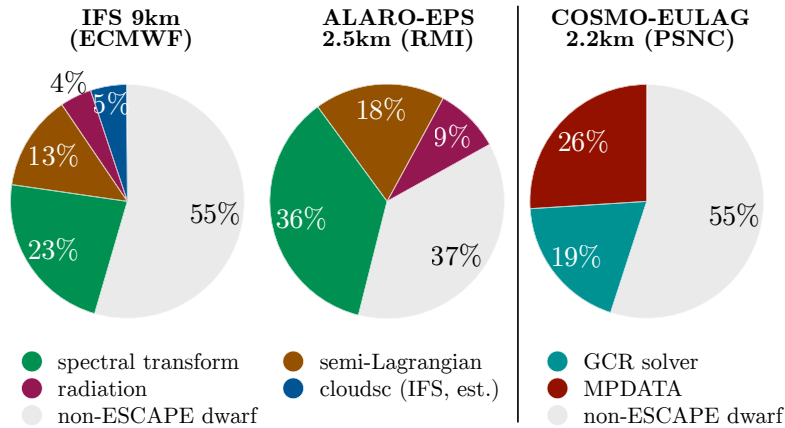


Figure 2. Portion of the forecast runtime spent in ESCAPE dwarfs for the three different models, IFS (left), ALARO ensemble prediction system (EPS) (middle) and COSMO-EULAG (right). The measurements for the IFS were taken during an operational run on 352 nodes (1408 MPI processes, 18 OpenMP threads per process). The limited area models ALARO-EPS and COSMO-EULAG used each 576 MPI processes for the simulations shown here. The IFS and ALARO-EPS are both based on the spectral transform method with latter employing 2D FFTs. For this reason the two left pie charts share one legend. COSMO-EULAG uses different methods and the pie chart has its own legend. The vertical line separates the two different legends.

~~timers to see the speedup obtained through optimisations. Table~~ as well as timers and error measures. The timers enable a quick assessment of the speedup due to the optimisations, whereas the error measures verify to what extent the optimisations affect the results. Table 1 also lists the models in which ~~we identified the dwarf~~ each dwarf was identified and from which ~~the prototype implementations~~ its prototype implementations have originated. The models include the IFS with the spectral transform method (IFS-ST), the Finite-Volume Module of the IFS (IFS-FVM), ALARO/AROME and COSMO-EULAG. In order to explore alternative discretizations in ESCAPE, ~~we created~~ a new global shallow water model called GRASS (Global Reduced A-grid Spherical-coordinate System) ~~Table~~ was developed. Table 1 further shows which dwarfs have prototypes that are based on the data structure framework Atlas ~~We use Atlas~~ (Deconinck et al., 2017), used to ease adaptation to future architectures and to avoid code duplication across prototypes. Atlas handles both the mesh generation and parallel communication aspects (Deconinck et al., 2017). In the ESCAPE project, Atlas has been improved and extended by adding support for limited area grids ~~and adding support for accelerators through~~ as well as for accelerators via GridTools (Deconinck, 2017a, b). Table 1 also shows the programming model adopted for each of the dwarf prototypes. In particular, ~~we used MPI~~ MPI was used for distributed-memory parallelism, OpenMP and OpenACC for shared-memory parallelism, and ~~DSL, that stands for domain-specific language and uses GridTools~~ two domain specific language (DSL) solutions - CLAW and GridTools - for hardware specific optimisations. Apart from DSL, all dwarfs were coded in modern Fortran.

In addition to identifying computational and communication patterns in existing models ~~we also performed research on new algorithms and approaches which have the~~ there was a search for new algorithms that show significant potential to reduce

Table 1. Overview of dwarfs identified in the ESCAPE project, short description of the computational characteristics of each dwarf and the models in which we identified the dwarf. Standalone applications including documentation, profiling and verification capabilities called prototypes have been released for many of the dwarfs. The column Atlas shows which of the dwarfs have prototypes that are based on the data structure framework Atlas (Deconinck et al., 2017). Dwarfs for which prototypes with MPI and OpenMP/OpenACC is available can be used with both programming models together as hybrid parallelisation. LAITRI is used in the full model in each MPI process individually and therefore does not include MPI parallelisation inside the dwarf. Further explanations can be found in the references given in the last column.

	computational characteristics	IFS-ST	IFS-FVM	ALARO/AROME	COSMO-EULAG	GRASS	prototype released	Atlas	MPI	OpenMP	OpenACC	DSL (Sect. 3.7)	Optalysys (Sect. 3.5)	references
Dynamics dwarfs:														
• time-stepping: spectral transform: – spherical harmonics – biFFT elliptic solver - GCR elliptic solver - multigrid HEVI	matrix multiplication, global transposition 1D and 2D FFT algorithm, global transposition iteration, global reduction interpolation for remapping nearest neighbour column communication	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	Mengaldo (2016) Mengaldo (2016) Mengaldo (2016) Müller et al. (2017) Mengaldo et al. (2018); Colavolpe et al. (2017)
• advection: semi-Lagrangian – LAITRI (3d interpolation) MPDATA	large halo communication efficient data access small halo communication, indirect addressing ²	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	Mengaldo (2016) Robinson et al. (2016) Mazaauric et al. (2017b); Osuna (2018)
• gradient computation: high order finite difference finite volume	large halo communication small halo communication, indirect addressing ²					✓		✓	✓	✓				Bénard and Glimton (2019); Glimton and Bénard (2019) Kühnlein et al. (2019)
Physics dwarfs:														
• radiation scheme: ACRANE2	many transcendental functions		✓				✓		✓	✓	✓			Müller et al. (2017); Poulsen and Berg (2017)
• cloud microphysics: CLOUDSC	single (vertical) column, many if-statements, many local variables	✓	✓				✓	✓	✓	✓	✓			Xiao et al. (2017); Clement et al. (2018)

¹ The GRASS model is currently only solving the shallow-water equations. Research on HEVI time integration methods was performed to explore future options.

² when using unstructured meshes

runtime and energy consumption ~~significantly in the future. We developed~~. This lead to the developments of a multigrid preconditioner for that reduces iterations of the Krylov-subspace solver, employed in the semi-implicit time integration in of IFS-FVM (Müller et al., 2017) ~~to reduce the number of iterations in iterative solves,~~ and a HEVI time-integration scheme with significantly improved stability (Colavolpe et al., 2017) ~~to avoid some of the global communication patterns, and connected~~
5 ~~to this, explored alternative finite difference methods~~. The latter avoids global communications, which also motivated the development of high-order finite-difference methods for reduced longitude-latitude grids on the sphere (Bénard and Glinton, 2019; Glinton and Bénard, 2019). ~~Finally, we~~ Furthermore, on the hardware side, ESCAPE explored FFTs and spherical harmonics on optical (co-)processors (Macfaden et al., 2017) to potentially scale these transforms towards higher resolutions at a fixed energy cost. An overview of the work with the optical processor can be found in ~~section ??~~ Sect. 3.5.
10 As an example of the cycle involved in identifying, dwarfing isolating, testing, adapting, optimising, and considering alternative solution procedures for a specific algorithmic motif in the ESCAPE project, ~~we shall illustrate in the next section the the~~ work on the spectral transform dwarf in is illustrated in more detail.

3 Dwarf example: spectral transform

~~We start with~~ Below is a short description of the (domain-specific) ~~spectral transform and continue with a subsection discussing~~
15 ~~computational challenges with a particular focus~~ spectral transform with a discussion of computational challenges, primarily focused on the data structures and data access patterns used in IFS. ~~After that~~ Subsequently, we present the work ~~that has been~~ done on adapting and optimising the dwarf for GPUs, CPUs and optical processors. ~~We, and~~ finish with a comparison ~~between~~ of the results obtained for the different architectures and a discussion on the sustainability of the chosen techniques.

3.1 Background

20 Each time-step of IFS is split between computations in physical space (i.e. a grid point representation) and computations with respect to spectral space (i.e. spectral coefficients at different wavenumbers). Semi-Lagrangian advection, physical parameterisations and products of terms are computed most efficiently in grid point space while horizontal gradients, semi-implicit calculations and horizontal diffusion are computed more efficiently in spectral space. The transform between these two spaces
25 is performed on the sphere with spherical harmonics, that is computing these results along longitudes in a Fast Fourier transform (FFT) and a Legendre transform (LT) along latitudes. Limited area models replace the Legendre transform with another FFT which leads to the name biFFT.

In spectral transform methods, such as the one used in IFS (Wedi et al., 2013), the specific form of the semi-implicit system facilitating large time-steps (and thus time-to-solution efficiency) is derived from subtracting a system of equations, ~~linearised~~
30 around a horizontally homogeneous reference state. The solution of this linear-system specific form is greatly accelerated by the separation of the horizontal and the vertical part, which matches the large anisotropy of horizontal to vertical grid dimensions prevalent in atmospheric and oceanic models. In spectral transform methods one uses the special property of the horizontal

Laplacian operator in spectral space on the sphere

$$\nabla^2 \psi_n^m = -\frac{n(n+1)}{a^2} \psi_n^m, \quad (1)$$

where ψ symbolises a prognostic variable, a is the Earth radius, and (n, m) are the total and zonal wavenumbers of the spectral discretisation (Wedi et al., 2013). This conveniently transforms the 3D Helmholtz problem into an array (for each zonal wavenumber) of 2D matrix operator inversions with the dimension of the vertical levels square ~~or~~ or vertical levels times the maximum truncation N_T , defined shortly, in the case of treating the Coriolis term implicitly ~~vertical levels times the maximum truncation~~, resulting in a very ~~cheap direct solve~~.

~~In this paper we focus on the efficient direct solver. The~~ computational aspects and especially ~~on~~ the data layout ~~We illustrate here are illustrated with~~ the inverse spectral transform ~~on the sphere which goes from spectral space to grid point space. The direct transform adds one numerical integration. Otherwise the only change between inverse and direct transform is that the direct transform starts with the Fourier transformation and applies the Legendre transformation afterwards. from spectral to grid-point space.~~

The inverse spectral transform begins with the spectral data $\mathbf{D}(f, i, n, m)$ which is a function of field index f (for the variable surface pressure at a single level and for wind vorticity, wind divergence and temperature at each height level), real and imaginary part i and wave numbers (zonal wave number $m = 0, \dots, N_T$ and total wave number $n = 0, \dots, N_T - m$ where ~~N_T~~ N_T is the spectral truncation). ~~Please note that we deviate here. This deviates~~ from the usual notation where total wavenumber goes from m to N_T ~~because this, because it~~ simplifies the separation between even and odd n . ~~We use here, using~~ column-major order like in ~~FORTRAN~~ Fortran, i.e. the field index f is the fastest moving index and the zonal wave number m is the slowest moving index. Typical dimensions can be seen in the operational ~~high-resolution~~ high-resolution (9km) forecast run at ECMWF: the number of fields is in this case 412 for the direct transform and 688 for the inverse transform and the number of zonal wave numbers is given by the truncation $N_T = 1279$. The number of latitudes is $2N_T + 2 = 2560$ and the number of longitudes increases linearly from 20 next to the poles to $4N_T + 20 = 5136$ next to the equator.

We take advantage of the symmetry of the Legendre polynomials for even n and anti-symmetry for odd n . The coefficients of the Legendre polynomials are pre-computed and stored in $\mathbf{P}_{e,m}(n, \phi)$ for even n and $\mathbf{P}_{o,m}(n, \phi)$ for odd n , where ϕ stands for the latitudes of our Gaussian mesh. Only latitudes on the northern hemisphere are computed. Latitudes on the southern hemisphere are reconstructed from the northern latitudes ~~as we will show later~~. In the same way ~~we split~~ the spectral data is split for each m into even part $\mathbf{D}_{e,m}(f, i, n)$ and odd part $\mathbf{D}_{o,m}(f, i, n)$. ~~We write variables over which we can parallelise our computations as indices~~ Fields (indexed f) are all independent in the parallel transform. The inverse Legendre transform is performed by computing the following matrix multiplications using BLAS:

$$\begin{aligned} \mathbf{S}_m(f, i, \phi) &= \sum_n \mathbf{D}_{e,m}(f, i, n) \cdot \mathbf{P}_{e,m}(n, \phi), \\ \mathbf{A}_m(f, i, \phi) &= \sum_n \mathbf{D}_{o,m}(f, i, n) \cdot \mathbf{P}_{o,m}(n, \phi). \end{aligned} \quad (2)$$

The resulting array for the symmetric and antisymmetric parts are now combined into the Fourier coefficients on the northern and southern hemisphere:

$$\phi > 0 : \mathbf{F}(i, m, \phi, f) = \mathbf{S}_m(f, i, \phi) + \mathbf{A}_m(f, i, \phi) ,$$

$$\phi < 0 : \mathbf{F}(i, m, \phi, f) = \mathbf{S}_m(f, i, \phi) - \mathbf{A}_m(f, i, \phi) .$$

- 5 These Fourier coefficients are finally used to compute the fields in grid point space at each latitude ϕ and longitude λ via FFT:

$$\mathbf{G}_{\phi, f}(\lambda) = \text{FFT}(\mathbf{F}_{\phi, f}(i, m)) . \quad (3)$$

3.2 Computational challenges

The computations in grid point space ~~and spectral space require all the fields~~ require all fields to be on the same computational node. The summation over the total wavenumber n in the Legendre transform (2) makes it most efficient to have all total wavenumbers (for a given field) on the same node and the Fourier transform (3) over (i, m) makes it most efficient to have all of the zonal wavenumbers m with real and imaginary part on the same node. This is only possible if the data is transposed before and after the spectral transform as well as in between Legendre and Fourier transform. These transpositions produce substantial communication which increases the contribution of the spectral transform to the overall runtime for future resolutions ([Figure Fig. 3](#)).

- 15 Simplified simulations of the MPI communications performed in ESCAPE ~~show~~ indicate that the strong scalability of the communication time for the spectral transform transpositions is better than for halo communication required by (large) semi-Lagrangian advection halos and global norm computation commonly used in semi-implicit methods (Zheng and Marguinaud, 2018).

These results indicate that halo communication will become almost as costly as the transpositions in the spectral transform method if ~~we use~~ a very large number of MPI processes is involved. An alternative which avoids transpositions and halo communication is given by the spectral element method shown in Müller et al. (2018) with explicit time integration in the horizontal direction. This leads to a very small amount of data that is communicated in each time-step because this method only communicates the values that are located along the interface between different processor domains. This method, however, requires much smaller time-steps which leads overall to an even larger communication volume ([Figure Fig. 4](#)). [Figure 4](#) is based on the model comparison presented in Michalakes et al. (2015) and does not include all of the optimisations for the spectral element method presented in Müller et al. (2018). The spectral transform results are based on the operational version of IFS and do not contain the optimisations presented in this paper. Both models have significant potential for optimisation and it is not obvious which method will have the lowest communication volume when fully optimised. ~~The only true solution, but key for both is an approach that allows reasonably large time steps. Regardless, the only way~~ to avoid waiting time during communication is to ~~overlap different parts of the model~~ enable concurrent execution of various model components such that useful computation can be done ~~while the data is communicated (Mozdzynski et al., 2015; Dziekan et al., 2019).~~ during data communication (Mozdzynski et al., 2015; Dziekan et al., 2019).

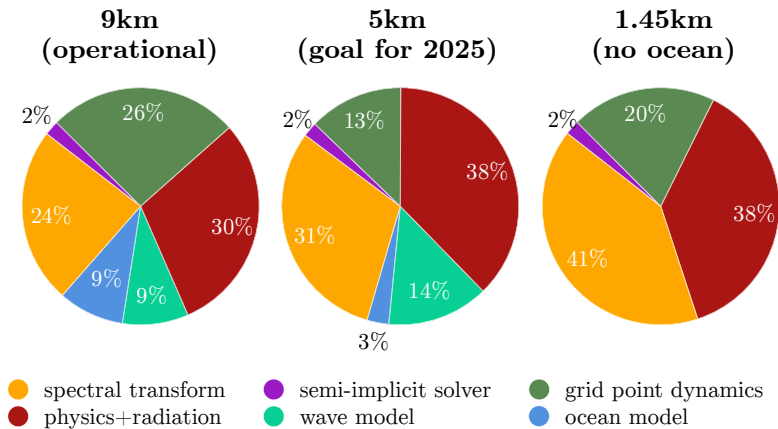


Figure 3. Cost profiles of the significant components of the IFS NWP model in percent of CPU time, at the operational 9km horizontal resolution with 137 vertical levels (left), the anticipated future horizontal resolution of 5km with 137 vertical levels (middle) and an experimental resolution of 1.45km with 62 vertical levels (right). The grid point dynamics represents the advection and gridpoint computations related to the dynamical core, semi-implicit solver represents the computations and communications internal to spectral space, spectral transform relates to the communications and computations in the transpositions from gridpoint to spectral and reverse as well as the FFT and DGEMM computations (see also spectral transform schematic below), physics+radiation relates to the cost of physical parametrisations including radiation, and finally accounting for the additional components of the wave and the ocean model. ~~The simulation at models - nonetheless excluded in the 1.45km is without ocean and wavs simulation.~~ All of these profiles have been obtained through measurements on ~~the ECMWF's~~ Cray XC40 supercomputer ~~of ECMWF~~.

3.3 GPU optimisation

For the GPU version, ~~we restructure the prototype to:~~ the prototype is restructured to: a) allow the grid-based parallelism to be fully exposed to (viz. exploited by) the GPU in a flexible manner; b) ensure that memory coalescing is achieved; and c) optimise data management. ~~We will now describe each of these in some more detail.~~

- 5 ~~The grid is a~~ To achieve performance, the two-dimensional ~~sphere,~~ grid over the sphere is augmented with a third ~~altitude~~ dimension represented by multiple fields at each point on the sphere. The updates dimension that represents multiple fields. The computations are inherently parallel across this grid, so all this parallelism ~~should can~~ be exposed to the GPU to get for maximal performance. ~~However~~ In contrast, the original implementation had a sequential loop over one of the ~~spherical~~ dimensions (at a high level in the call tree of the application). ~~We re-structured the code such that,~~ By restructuring the code for each operation,
- 10 ~~the~~ loops over the three dimensions became tightly nested, and when mapping these become tightly nested. These loops are mapped to the GPU via OpenACC directives ~~we used the "collapse" clause to instruct the compiler to collapse these to~~ collapse to a single loop, ~~such that it can map all inherent parallelism to hardware in an efficient manner. We re-structured the code as follows:-~~

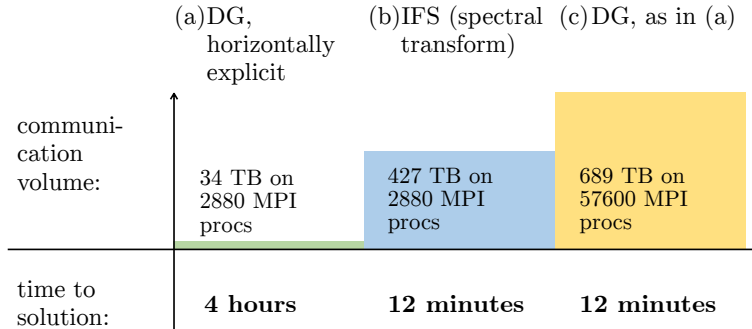


Figure 4. Overall communication volume comparing spectral element (SEM) from Müller et al. (2018) and the global spectral transform methods. The SEM requires a substantially lower amount of communication at the same number of cores, but due to the smaller timestep requires a much higher frequency of repeated communications for the given 2-day simulation. Extrapolating the number of MPI processes to achieve the same time to solution would result in a larger amount of communication for the SEM. Here we assume SEM $\Delta t = 4s$; IFS $\Delta t = 240s$; communication volume is calculated for a 48 hour forecast SEM as 290.4kBytes per MPI task and Δt ; IFS as 216mBytes per MPI task and Δt ; as a result the SEM time-to-solution is estimated as 20 x IFS based on the performance results in Michalakes et al. (2015).

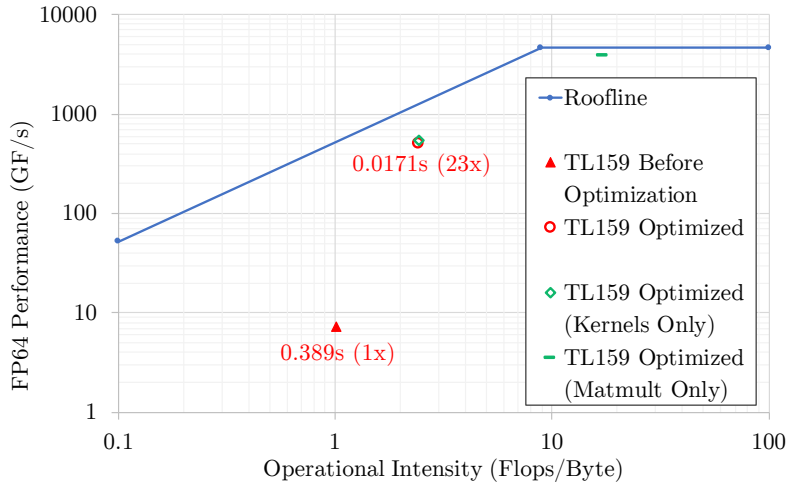


Figure 5. Roofline plot for the spectral transform dwarf at 125km resolution ($N_T = 159$) on the NVIDIA Tesla P100 GPU. The full time-step of the original prototype is represented by the solid red triangle. The corresponding time-step for the optimised prototype is represented by the red circle. Also included are partial results for kernels only (open green diamond atop the red circle) and matrix multiplication only (green dash). Each point is positioned in the plot according to its operational intensity: points under the sloping region of the roofline are limited by available memory bandwidth, and points under the horizontal region are limited by peak computational performance.

Similarly, for library calls it is important to maximise exposure of parallelism through batching computations using provided interfaces. On the GPU we perform all of the matrix multiplications in the Legendre transform (2) with a single batched call of the cuBlasDgemm library. The different matrices in (2) have different sizes because the total wavenumber goes from 0 to $N_T N_T$. To use the fully batched matrix multiplication ~~we pad each matrix~~ each matrix is padded with zeroes up to the largest size, since the library currently does not support differing sizes within a batch. This step ~~increases~~ results in a tenfold increase of the overall number of floating point operations ~~by almost a factor 10~~ but still improves the overall performance (Figure Fig. 5). ~~We perform the FFT in equation~~ The FFT in (3) is performed with the cuFFT library, ~~where we batch~~ batching over the vertical ~~dimension but multiple~~ (field) dimension. Multiple calls are still needed ~~over the spherical dimension~~ (noting for each latitude since FFTs cannot be padded in a similar way to matrix multiplications). Therefore the implementation remains suboptimal ~~here: we are still not fully exposing parallelism and there~~. There would be scope for further improvements if ~~a~~ an FFT batching interface supporting differing sizes were to become available.

~~We restructured array layouts to ensure that multiple threads on the GPU can cooperate to load chunks of data from memory in a "coalesced" manner. This would allow a high percentage of available memory throughput. This is achieved when the fastest moving index in the multidimensional array corresponds to the OpenACC loop index occurring at the innermost level in the collapsed loop nest as follows: \$!ACC parallel loop collapse(3) do k=1, ... do j=1, ... do i=1, ... array(i,j,k)=...~~

~~In FORTRAN, arrays are structured such that elements accessed by consecutive innermost indices (i in this example) are consecutive in memory. Since i is used as the innermost loop index in this collapsed loop nest, then memory coalescing is achieved.~~

Sometimes matrix transposes are necessary, but where possible these were pushed into the DGEMM library calls, which have much higher-performing implementations of transposed data accesses. There remain transpose patterns within kernels involved in transposing grid point data from column structure to latitudinal (and inverse) operations, which naturally involve transposes and are thus harder to fix through restructuring. However, ~~we optimised these~~ using the “tile” OpenACC clause, which instructs the compiler to stage the operation through multiple relatively small tiles ~~which~~, can perform the transpose operations within fast on-chip memory spaces ~~such that the~~ accesses access to global memory ~~are much~~ is more regular.

Data allocation on the GPU is expensive, as is data movement between the CPU and GPU. ~~We structured the code~~ Hence the code is structured such that the fields stay resident on the GPU for the whole timestep loop: all allocations/frees have been moved outside the timestep loop with re-use of temporary arrays, and thus all data transfer has been minimized. Arranging the model data consecutive in memory ensures that multiple threads on the GPU can cooperate to load chunks of data from memory in a "coalesced" manner. This allows a high percentage of available memory throughput.

The restructured algorithm achieves an overall speedup factor of ~~23x~~ 23x compared to the initial version which also used cuBlas and cuFFT but followed the CPU version more closely. Matrix multiplication performance is higher than the overall performance (in flops) and the operational intensity is increased into the compute-bound regime. Note that matrix multiplication is associated with $O(N^3)$ computational complexity for $O(N^2)$ memory accesses. The extra padding operations lead to larger

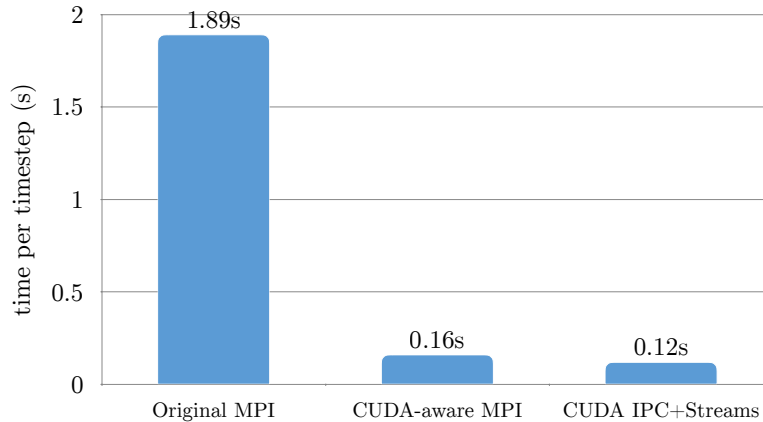


Figure 6. Computational performance of the spectral transform dwarf at 18km resolution ($N_T = 639$) on 4 NVIDIA V100 GPUs of the DGX1 with the original MPI implementation (left), CUDA-aware MPI communication (middle) and NVLink optimised communication (right). This resolution is currently used operationally for the members of the ensemble forecast at ECMWF.

N and therefore also to increased operational intensity. More details about the single GPU optimisations can be found in Mazauric et al. (2017b).

Beyond a single GPU, with multiple interconnected GPUs we see For multiple GPU arrangements there is a large benefit by using the modern NVLink interconnect and the recently announced NVSwitch from high-bandwidth interconnects (i.e. NVLink, NVSwitch) due to the high relative importance of communication for the transpositions described in section ??. Each GPU features multiple ports of high-bandwidth NVlink connections, each providing 50 GB/s of bi-directional bandwidth when using the Volta GPUs., cf. Sect. 3.2. For full bandwidth connectivity when using more than 4 GPUs we use, the NVSwitch interconnect on the DGX-2 server. The DGX-2 server has 16 Volta V100 GPUs: each with six 50 GB/s NVLink connections into the switch with routing to any of the other GPUs in the system. This allows transfers up to 300 GB/s communications between any pair of GPUs in the system, or equivalently 2.4 TB/s total throughput.

When running a single application across multiple GPUs, it is necessary to transfer data between the distinct memory spaces. Traditionally, such transfers needed to be realised via host memory and required the participation of the host CPU. Not only did this introduce additional latency, but also limited the overall bandwidth to the bandwidth offered by the PCIe bus connecting CPU and GPUs. However, modern more recent MPI implementations are CUDA-aware. This means that pointers to GPU memory can be passed directly into the MPI calls, avoiding unnecessary transfers (both in the application and in the underlying MPI implementation). This is particularly useful when using a server that features high-bandwidth NVLink connections between GPUs, in which case CUDA-aware MPI will use these links automatically. Moving our dwarf to CUDA-aware MPI gave us a speedup of 12x (Figure 12x (Fig. 6)). However, even with this optimisation the all-to-all operations remained inefficient because communication between different GPUs was not exchanged concurrently. Perfect overlap was achieved by implementing an optimised version of the all-to-all communication phase directly in CUDA using the Inter Process

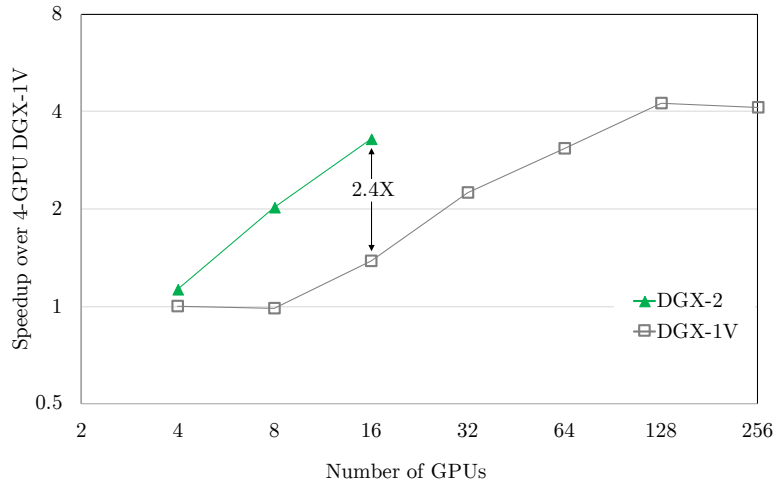


Figure 7. Computational performance of the spectral transform dwarf at 18km resolution ($N_T = 639$) on up to 16 GPUs on one DGX-2 and up to 32 DGX-1 servers connected with Infiniband. The DGX-1V uses MPI for ≥ 8 GPUs (due to lack of AlltoAll links), all others use CUDA IPC. DGX-2 results use pre-production hardware. The points were connected with lines for the purpose of improving readability.

Communication (IPC) API. Using memory handles, rather than pointers, CUDA IPC allows to share memory spaces between multiple processes, thus allowing one GPU to directly access memory on another GPU. This allowed another speedup of about 30% (Figure Fig. 6).

In Figure 7 we demonstrate how the use of Figure 7 compares a DGX-2 with NVSwitch allows significantly better scaling than the use of and the DGX-1 for a spectral transform at 18km resolution ($N_T = 639$). Note that we tune the number of MPI tasks in use: we use the The NVIDIA Multi Process Service to allow allows oversubscription of GPUs such that, e.g., the 8 GPU result on DGX-2 uses 16 MPI tasks across the 8 GPUs (i. e. 2 operating per GPU). This is because such oversubscription can sometimes. Oversubscription, a technique common to other architectures, can be beneficial to spread out any load imbalance load imbalances resulting from the spherical grid decomposition (see below) grid decomposition and hide latencies. We chose the best performing number of MPI tasks per GPU in each case.

As we increase the number of GPUs, the The scaling on DGX-1V is limited : This is because we no longer retain full connectivity and some messages must because with an increasing number of GPUs some messages go through the lower-bandwidth PCIe and QPI links and/or Infiniband when scaling across multiple servers. But on On DGX-2 with NVSwitch, all 16 GPUs have full connectivity: that is we have: i. e. the maximum peak bandwidth of 300 GB/s between each pair of GPUs in use is available. The performance is seen to scale well out scales well up to the full 16 GPUs on DGX-2, where the difference and is 2.4x faster than the result obtained with the 16 GPU (2-server) DGX-1V result is 2.4x. It can also be seen that the The speedup going from 4 to 16 GPUs on DGX-2 is 3.2x, whereas the ideal speedup would be 4x. However, initial 4x. Preliminary investigations reveal that this deviation from ideal scaling is not primarily due to communication overhead but instead due to load imbalance between the MPI tasks from the for a given spherical grid decomposition that is chosen by the

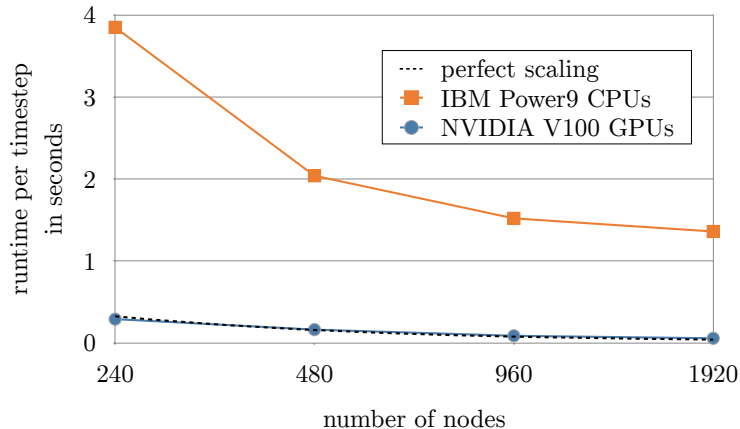


Figure 8. Strong scaling comparison between GPU optimised version and initial CPU version of the spectral transform dwarf on Summit at the Oak Ridge National Laboratory. All of these simulations are using 2.5km resolution ($N_T = 3999$) with 240 fields. The GPU version uses 6 V100 GPUs per node which leads to a maximum number of 11520 GPUs. The perfect scaling was chosen such that it goes through the GPU result at 480 nodes to illustrate the scaling efficiency of the GPU version. The points were connected with lines for the purpose of improving the readability of the plot.

~~application in each case, which would indicate that better scaling would.~~ This indicates that better scaling might be observed with a more balanced decomposition. More details about the multi-node optimisation of the spectral transform dwarf can be found in Douriez et al. (2018).

First results on the supercomputer Summit of the Oak Ridge National Laboratory are shown in [Figure Fig. 8](#). At this large scale we observe a speedup ~~of about 2x~~ between 12.5x (on 480 nodes) and 23.7x (on 1920 nodes) when comparing the optimised GPU version with the initial CPU version of the spectral transform dwarf on the same number of nodes. These simulations were run at 2.5km resolution ($N_T = 3999$) and use 240 fields. In operational application the number of fields should be larger (see [Section ??](#)). ~~We still need to optimise the memory consumption and the initialisation to test our prototype for more realistic numbers. Also there is still room for improving the setup of the environment on Summit, and exploring alternative programming models. The important message is Sect. 3.1). It has been verified~~ that the GPU optimised dwarf from ESCAPE was is capable of running on a huge number of GPUs on Summit ~~and first results suggest that we achieve with~~ a significant speedup. ~~We will do more analysis of these simulations in the future.~~

3.4 CPU optimisation

The spectral transform dwarf is based on the operational implementation used in IFS that has been continuously optimised over multiple decades. According to profiling results, it clearly appeared that the main computational intensive kernels are the FFT and matrix multiplication executed by a dedicated highly tuned library (as Intel Mathematics Kernel Library, called MKL). In support of this work ~~we looked into~~ different data scope analysis tools ~~-A comparison of the different tools is available in~~ we

explored; see Mazauric et al. (2017a). The first optimisation strategy concentrated the effort on non-intrusive optimisations which have the advantage of being portable and maintainable. Among these optimisations, the use of extensions to the x86 instruction set architecture (ISA) as SSE, AVX, AVX2, AVX-512 is notable, because it indicates how much of the source code can be vectorised by the compiler. When the compiler failed at vectorising some loops or loop nests, a deeper investigation
5 of how to use compiler directives followed. As the different instruction sets are not supported by all processors, the study proposed an intra node scalability comparison study among several available systems (at the time of benchmarking).

System tuning using Turbo frequency (TUR), Transparent Huge Page (THP), memory allocator (MAP) can be done without modifying the source code. This exposes both performance gains and interesting information on dwarf behaviour. Indeed, on ATOS BullSequana x400 supercomputer with Intel® Xeon Broadwell E5-2690v4 processors, enabling turbo offers a gain
10 equal to 11%, enabling THP gives 22%, MAP 27%, and finally the best performance (35% of performance gain) is achieved by the combination of MAP and TUR. This shows that memory management is a key point. More details about the single-node CPU optimisation can be found in Mazauric et al. (2017b).

Multi-node optimisation for CPUs focused on improving the MPI communication. The largest potential for optimisations was found to be in the preparation phase of point-to-point communications. During the preparation phase the sender side gathers
15 the local data into a contiguous buffer (Pack operation) and hands it off to the MPI library. On the receiver side, data is then scattered from a contiguous user buffer to its correct location (Unpack operation). Pack and Unpack are nearly inevitable with scattered data because Remote Direct Memory Access (RDMA) with no gather-scatter operations are known to be often less effective, notably due to the memory pinning latency (Wu et al., 2004). It also means that sender and receiver must exchange their memory layout as they may differ. The performance improvement came from reordering the loops for both pack and
20 unpack following the memory layout of the scattered buffer. This optimisation ~~decreased~~decreases the number of tests (i.e. copy or not copy) and avoids scanning memory multiple times~~which was unnecessary. We reduced this with~~, leading to a global performance gain on the whole dwarf of about 20% (Figure Fig. 9a). ~~A few percent of improvement came also from disabling MPI barriers which existed in the code to profile the MPI communications at low level. These barriers produced imbalance, the necessary synchronisations created contention and overall created a bias in the communications profile.~~ The computational
25 performance up to 30 ATOS BullSequana x400 nodes equipped with Intel® Xeon Skylake 6148 processors is shown in Figure Fig. 9b. ~~This~~The work has been performed on ATOS internal HPC resources. This optimisation can be ~~immediately directly~~ applied to the operational model IFS due to its non-intrusiveness. More details about the multi-node optimisation on CPUs can be found in Douriez et al. (2018).

The speedup of the GPU optimisations came from using highly optimised GPU libraries for batched DGEMM and FFT
30 and avoiding the transposition of temporary arrays~~which was not necessary~~. On the CPU~~we can~~, one can equally apply these improvements in the handling of temporary arrays~~as well. We used these optimisations in a newly developed~~, which lead to the development of a new serial spectral transform inside Atlas~~which is~~, used operationally at ECMWF for post-processing purposes~~since the beginning of 2019~~. Post-processing is run in serial mode due to the large number of concurrent post-processing jobs. Compared to the previous ~~operational serial transform used for post-processing we find a speedup of about 3x~~

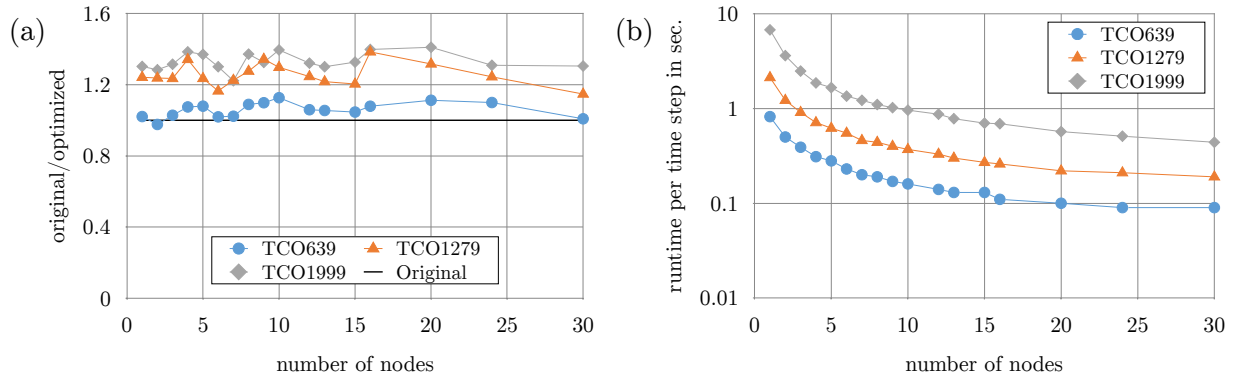


Figure 9. Performance measurements on up to 30 Intel® Xeon Skylake 6148 nodes. Subpanel (a) shows the speed-up at 18km resolution ($N_T = 639$) with regards to different optimisations of the communication preparation phase. Subpanel (b) shows the computational performance for three different resolutions of 18km ($N_T = 639$), 9km ($N_T = 1279$) and 5km ($N_T = 1999$) representative of current and future operational requirements. Note the use of log-scale. The points were connected with lines for the purpose of improving readability. All of these results have been obtained through measurements on the ATOS BullSequana x400 supercomputer of ATOS internal HPC resources.

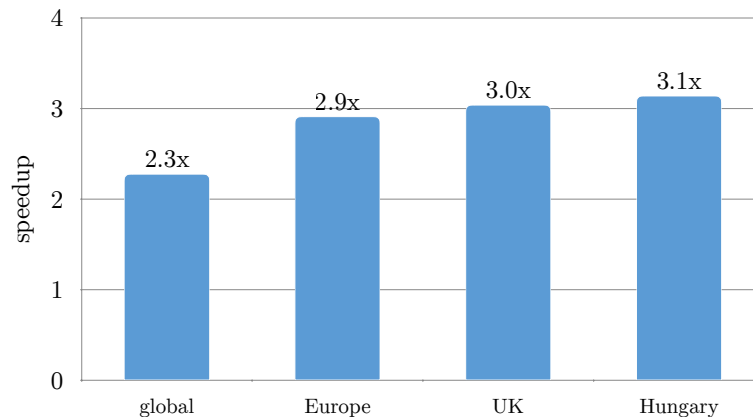


Figure 10. Speedup of the spectral transform by porting optimisations introduced in the GPU version back to the CPUs. The base line for this comparison is the current operational post-processing library used at ECMWF. This version of the spectral transform allows the computation on limited area domains. The speedup is given here for the global transform and three examples of limited domains (Europe, UK and Hungary).

(Figure-serial transform employed in ECMWF product generation, there is a speedup $2\times$ to $3\times$, in particular for limited area domains (Fig. 10).

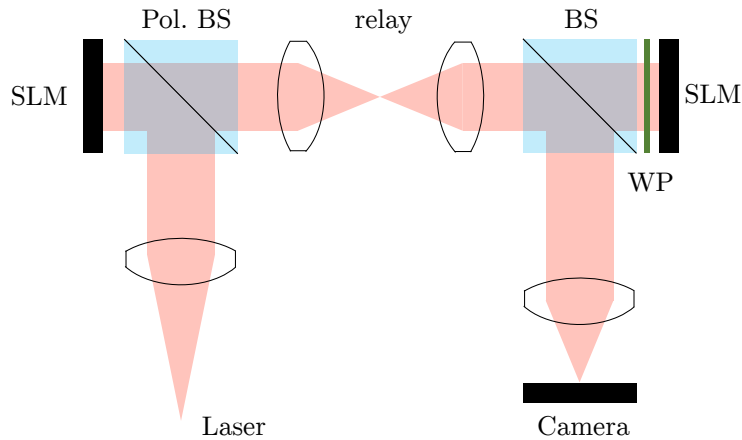


Figure 11. Illustration of the fundamental idea behind the optical processor. The laser beam is emitted on the bottom left. Two spatial light modulators (SLM) are used together to input the complex function. The system uses beamsplitters (BS) and an optical relay to image one reflective SLM onto another, followed by a lens assembly which approximates an ideal thin lens and renders the optical Fourier transform on a camera sensor. The half-waveplate (WP) before the second SLM is used to rotate linearly polarised light onto the axis of SLM action (the direction in which the refractive index switches), thus causing it to act as a phase modulator.

3.5 Optical processors

~~Optalysys have been investigating an optical implementation of the spectral transform dwarf (biFFT for limited area models as well as spherical harmonics for global models).~~ The fundamental idea behind optical processors is to encode information into a laser beam by adjusting the magnitude and phase in each point of the beam. Optalysys have been investigating an optical implementation of the spectral transform dwarf - biFFT for limited area models as well as spherical harmonics for global models. This information becomes the Fourier transform of the initial information in the focal plane of a lens. The information can be encoded into the optical beam by using spatial light modulators (SLMs) as illustrated in [Figure-Fig. 11](#). The result of the Fourier transform can be recorded by placing a camera in the focal plane of a lens. A photo of an early prototype is shown in [Figure-Fig. 12](#).

10 SLMs are optical devices with an array of pixels which can modulate an optical field as it propagates through (or is reflected by) the device. These pixels can modulate the phase (essentially applying a variable optical delay) or polarisation state of the light. Often they modulate a combination of the two. When combined with polarisers, this polarisation modulation can be converted into an amplitude modulation. Hence the modulation capability of a given SLM as a function of ‘grey level’ can be expressed by a complex vector, which describes an operating curve on the complex plane. Each pixel of the SLM is generally a
 15 1-parameter device; arbitrary complex modulation is not offered by the SLM, only some sub-set. This is one of the key issues with regards to exploiting the optical Fourier transform.

Sensor arrays - essentially common camera sensors - are used to digitise the optical field. They are in general sensitive to the intensity of the light, which is the magnitude of the amplitude squared. This poses a difficulty to sensitively measuring

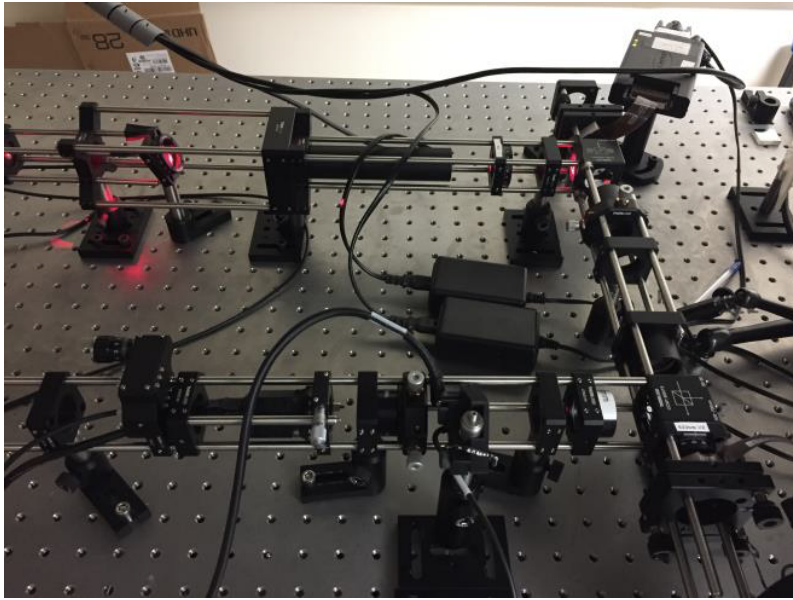


Figure 12. Photo of the first prototype of the optical processor. The final product is built into an enclosure of similar size like a GPU.

the amplitude. Moreover, ~~they~~ sensor arrays are not sensitive to optical phase. ~~We overcame this~~ The latter can be overcome with a method that ~~allowed~~ allows the optical phase to be derived from the intensity-only measurements. This ~~was done by~~ introducing known terms into the functions and is achieved by perturbing transformed functions with a specified disturbance, measuring the resulting change in the output, and deducing the phase from the properties of the Fourier transforms.

5 Each pixel of the SLM is addressable with an 8-bit value (256 levels). The SLM is not capable of independently modulating the magnitude and phase of the optical field. In the Optalysys processing system, the SLMs are configured to modulate both the amplitude and phase in a coupled manner, such that optimal correlation performance is achieved. The optical Fourier transform and all of the functions are inherently two dimensional. The propagating light beam can be thought of as a 2D function propagating and transforming along a third direction. The system is most naturally applied to 2D datasets, and many
10 problems can be mapped to an appropriate representation.

A critical aspect to realizing the potential of optical processing systems is the interface to a traditional computing platform. Bridging this gap has been a significant undertaking for Optalysys, and has resulted in the development of a custom PCIe drive board. This board interfaces to a host machine over PCIe and has direct memory access (DMA) to the system memory (RAM). It provides an interface to 4 SLMs and 2 cameras. The cameras are 4K (4096x3072). Initially, they operate at 100 Hz, but a
15 future firmware upgrade will unlock 300 Hz operation, and 600 Hz half-frame operation, dramatically increasing the potential data throughput.

There are currently two options for the SLMs. One option is using high speed binary (bright-dark) SLMs which operate at 2.4kHz. This offers correlation at binary precision. The second option is greyscale SLMs which operate at 120 Hz. This is currently the only option to reach more than binary precision. The performance of the entire processor is determined by the

part with the lowest frequency. The main bottleneck with multiple bit precision is the operating frequency of the greyscale SLM. There is currently no easy solution to increase the frequency of greyscale SLMs.

Optical processing is more appropriately applied to cases where high-throughput relatively-complex operations are the priority, with less of an emphasis on numerical precision. ~~The inherent ability of optical correlators to rapidly process convolutions naturally leads to the formation of convolution neural nets and machine learning technologies.~~ More details about the Optalysys optical processor have been published in Macfaden et al. (2017) and Mazauric et al. (2017b). The results from the optical processing tests will be discussed in the next section.

~~3.6 Comparison between processors in terms of runtime and energy consumption~~

~~We will loose a~~

3.6 Comparison between processors in terms of runtime and energy consumption

A lot of the speedup achieved by running the spectral transform on accelerators is lost if the CPUs are idle during the computation on the accelerators. ~~Also we need to take the~~ The cost of data transfer between CPU and accelerator ~~into account which needs to be considered but~~ has not been included in the speedup numbers in this section. To take full advantage of ~~the NVLink and NVSwitch we would need to run the entire simulation on a single node which requires at the currently operational resolutions more work on optimising the~~ high-bandwidth connected accelerators the entire simulation has to be run on the same node, which requires to substantially reduce the memory footprint of the model.

For the CPUs and GPUs used in this paper the overall cost is dominated by the cost of the hardware and therefore by the number of sockets/devices required to reach the desired runtime. In addition to the number of devices we also compare the energy consumption. The large number of zero operations caused in the optimised GPU version by the padding of the matrices in the Legendre transform makes it impossible to do a fair comparison between CPU and GPU by comparing metrics based on floating point operations including comparing roofline plots.

In the full operational model the 18km resolution ensemble member ($N_T = 639$) using 30 nodes on the Cray XC40 takes about 1.4s per time-step and the spectral transform component is about 15 percent (0.21s) of the overall runtime. Measurements with the dwarf on the Cray XC40 at 18km resolution ($N_T = 639$) resulted in 4.35s per timestep on a single node (4 MPI tasks, 18 threads per task), and 1.77s on 2 nodes. The energy consumption was measured at around 0.3Wh on the Cray XC40, which compares to 0.026 Wh measured on 4 V100 GPUs on a DGX1 which take 0.12s per time-step. The energy measurement on the XC40 is based on proprietary power management counters. The measurement on the V100 GPUs uses the ~~nvidia-smi~~ NVIDIA-smi monitoring tool.

Tests in ESCAPE on the latest generation of Intel Skylake CPUs have shown 0.12s per timestep using 13 Intel® Xeon Skylake 6148 nodes (connected via a fat-tree EDR network) as shown in ~~Figure Fig.~~ Fig. 9. This parallel CPU version has not seen the more radical changes which have been used in redesigning the algorithm for the parallel GPU and serial CPU version. ~~There might still be potential for more substantial optimisations in the parallel CPU version which we will explore in future research.~~

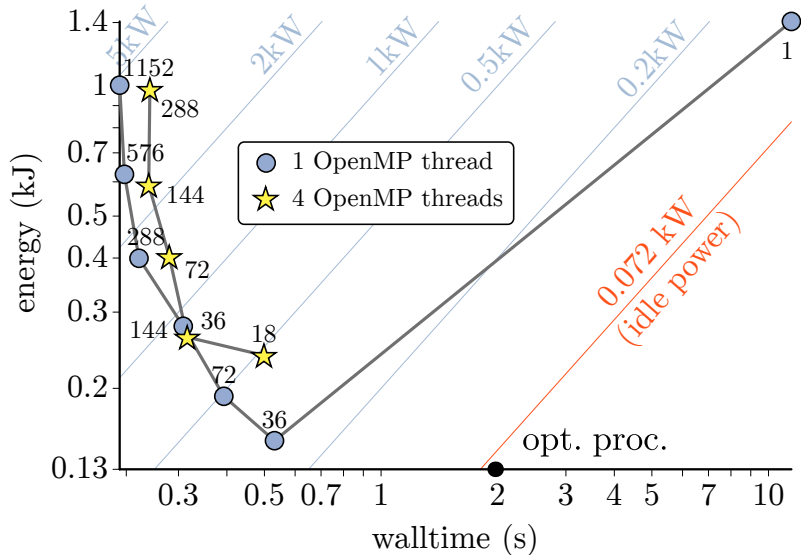


Figure 13. Log-log plot of the energy consumption vs. wall-clock time for the BiFFT dwarf and corresponding to the combination of one direct and one inverse transformation for 525 fields. Each data point is the result of averaging the outcome of two separate runs. Grey lines connect runs with the same number of OpenMP threads (1 resp. 4). Added are lines of constant power (light blue lines), including the power delivered by a node in the idle state (orange line). Indices next to each data point denote the number of MPI tasks. The black dot represents the estimate of the Optalysys optical processor when using a greyscale SLM. The performance of the optical processor at binary precision is much higher (not shown).

A comparison between CPUs and optical processor with greyscale SLM is shown in [Figure Fig. 13](#). The energy consumption of the optical processor is much lower than for the CPU. The runtime of the optical processor is larger due to the relatively slow performance of the grayscale SLM which is currently necessary to reach the precision ~~necessary~~required for NWP applications. ~~More~~The latter does not allow to exploit the superior energy-to-solution of optical processors in time-critical
5 NWP applications; more details about the comparison between CPU and optical processor can be found in Van Bever et al. (2018). In general, for hybrid arrangements of CPU and accelerators, we conclude that optimising energy-to-solution requires algorithmic changes that facilitate the simultaneous use of all the available processors.

3.7 Sustainability of code optimisation techniques

The optimisation results presented in ~~this section~~Sect. 3.6 are just an example for one of the dwarfs. Similar improvements
10 have been achieved for the advection scheme MPDATA (Douriez et al., 2018) and the radiation scheme ACRANEB2 (Poulsen and Berg, 2017). These optimisations are hardware specific and will be difficult to maintain on upcoming new architectures. As a strategy towards performance portability and sustainability we have worked in the ESCAPE project on domain specific languages through the GridTools framework. The main goal of the GridTools library is to provide a solution for weather and climate models to run one code ~~base~~ on many different architectures (portability) and achieve good performance (performance

portability). However, the main operational product of GridTools so far focused on solutions for lat-lon grid models like COSMO. The work developed in the ESCAPE project aimed at extending the DSL support for irregular grids ~~and~~, the efficient generation of backends for multiple architectures ~~and evaluating DSL useability~~.

~~These~~The DSL developments have been used to implement a portable version of the MPDATA dwarf. The DSL version
5 hides such details as the nested loops and the OpenACC directives used to specify properties of the GPU kernel and data layouts of the ~~FORTRAN~~Fortran arrays. Furthermore, the DSL allows to compose several of these operators together, which is used by the library to apply advanced performance optimisations like loop fusion or software managed caches. Comparing the ~~FORTRAN~~Fortran OpenACC kernel with the DSL version gives us a speedup of 2.1x ~~x~~ for the DSL version. This speedup could also be achieved by hand-tuned optimisation. The DSL prevents the repeated manual effort of tuning the code
10 for multiple architectures. At the same time the DSL allows ~~to perform optimisations~~ optimisations to be performed which would otherwise make the code unreadable. More details about this work including code examples on how to use the new backend to GridTools can be found in Osuna (2018). ~~More work on sustainability through domain specific languages will follow through the ESCAPE-2 project which started in October 2018.~~

In addition to GridTools ~~we also worked on using~~, ESCAPE explored domain specific languages via the CLAW DSL
15 (Clement et al., 2018) for the cloud microphysics dwarf. In particular the use of the Single Column Abstraction (SCA), where physical parameterisations are defined solely in terms of a single horizontal column, enables domain scientists to define ~~physics governing~~ equations purely in terms of vertical dependencies without needing to account for parallelisation issues. The CLAW DSL then inserts loops over the data-parallel horizontal dimension specific to the hardware architecture and programming model (OpenMP, OpenACC) via source-to-source translation, allowing multiple architectures to be targeted from a single
20 source code. A GPU implementation of the CLOUDSC dwarf generated by automated source translation tools has been used to generate similar performance results to the ones presented by Xiao et al. (2017) on K80 GPUs ~~via~~, using the OpenACC backend of the CLAW DSL.

4 Conclusions and outlook

The ESCAPE project has introduced the concept of Weather & Climate dwarfs as fundamental domain-specific building blocks
25 into the weather and climate community. ~~Their~~Dwarfs' categorisation of computational and communication patterns has been extremely useful in further breaking down the complexity of weather prediction and climate models, and advancing their adaptation to future hardware. Prototype implementations of ~~these~~the dwarfs have been used to work on ~~optimising them and using them for the purpose of benchmarking new computers. These included code optimisation and benchmarking new architectures. Dwarfs include~~ measures for verifying ~~their scientific correctness~~ scientific correctness of the implementations,
30 documentation and input from domain scientists. Our dwarfs are ~~very well suited for optimisation and benchmarking, they are~~ small enough to be fairly easy to ~~understand~~learn, and at the same time they represent a significant part of the performance of the whole weather prediction model.

The paper gives an overview of the dwarfs ~~that were~~ identified in the ESCAPE project, while illustrating ~~with the spectral transform dwarf a detailed example of the~~ in fair detail the optimisation cycle within ESCAPE. ~~In ESCAPE-2 we further identify dwarfs in other Earth System components such as sea-ice and ocean models, that are crucial for the performance of coupled applications~~ applied to the spectral transform dwarf that encapsulates historically perhaps the most fundamental algorithmic motif in operational global weather models (Williamson, 2007).

To avoid code duplication we used the data structure and mesh handling framework provided by Atlas. Atlas has been extended in ESCAPE to support limited area grids and DSL with GridTools. The participating weather models and most prototype implementations are based on ~~FORTRAN~~ Fortran and all our optimisations can be incorporated into ~~FORTRAN~~ Fortran code including the use of CUDA functions on GPUs by calling C functions from ~~FORTRAN~~. ~~We have started to incorporate the obtained code optimisations into operations which gives us a significant speedup in the spectral transform used for postprocessing~~ Fortran. Selected code optimisations are now routinely used in ECMWF's services. Besides optimisations of the existing code, improved algorithms have been developed which are specifically targeted at improving performance on large scale systems, which include a multigrid preconditioner for the elliptic solver to reduce iteration counts in iterative solvers, a HEVI time-integration scheme with significantly improved stability and alternative finite difference methods on the sphere in the context of reducing global communications across large processor counts, and include alternative solution procedures for spectral transforms at fixed energy cost, with FFTs and spherical harmonics realised on optical processors.

Code optimisations performed in the ESCAPE project targeted Intel CPUs, Intel Xeon Phi processors, NVIDIA GPUs and Optalysys optical processors. In NWP applications, the bottleneck in terms of performance is usually the memory bandwidth between processor and main memory. Having fast ~~interconnect like NVLink and NVSwitch~~ interconnects can provide a ~~massive~~ significant speedup. Having all of the processing units used by the simulation connected with such a fast interconnect is still a challenge. Using accelerators only for a small part of the code ~~destroys a lot of the benefit in terms of overall cost if the reduces their benefit if the associated~~ CPUs are idle while the accelerators perform their computations. ~~We either need to move~~ Hence, a large part of the code ~~to the accelerator (like in Fuhrer et al., 2018; Schalkwijk et al., 2015) or we need to overlap~~ has to be moved to the accelerator (cf. Fuhrer et al., 2018; Schalkwijk et al., 2015) or computations on the CPU are overlapped with computations on the accelerator ~~(like in Dziekan et al., 2019). Applying this to the participating models still requires more work. (cf. Dziekan et al., 2019).~~

Most of the work done in the ESCAPE project in terms of code optimisation focused on hardware specific optimisations. This makes the optimisations difficult to maintain on upcoming new hardware architectures. As an approach towards performance portability ~~we implemented~~ ESCAPE explored a DSL prototype for advection with MPDATA by using the GridTools framework, and a prototype of the cloud microphysics dwarf by using the CLAW DSL. DSLs are a ~~very~~ promising tool to enable good performance on multiple architectures while maintaining a single code base. However, designing a domain specific language that is user friendly and at the same time close to hand tuned performance on each architecture is challenging. ~~This work will continue in the ESCAPE-2 project which started in October 2018.~~

~~Comparing different methods requires to include all~~ When assessing the efficacy of diverse numerical methods - e.g., finite-volume vs. spectral spatial discretisation, or Eulerian vs. Lagrangian time stepping - it is important to account for all

computational costs. Spectral ~~transform~~transforms with semi-implicit, semi-Lagrangian time-integration has always been considered as being poorly suited for large supercomputers due to the high volume of communication. Our work indicates that thanks to much larger time-steps the overall communication cost is not necessarily worse than for other more local methods. Again overlapping different parts of the model such that useful computation can be done while data is communicated is a way
5 forward and needs to be high priority in future research.

When moving towards very high resolution global simulations of O(1km) or less and considering exascale computations on a variety of emerging HPC architectures, there is a continued interest and need in pursuing fundamentally different algorithmic approaches that simply do not communicate beyond a certain halo size while retaining all other favourable properties, such as removing time-step size restrictions as in the semi-Lagrangian advection with semi-implicit time-stepping(~~SISL~~), ~~and such~~
10 ~~approaches~~. Such approaches, as well as the identification of dwarfs in sea-ice and ocean model components of the Earth System, crucial for the performance of coupled applications, are further investigated in ESCAPE-2. These algorithms will utilise DSLs for performance portability.

Code availability. The data structure framework Atlas is available at <https://github.com/ecmwf/atlas> under an Apache License 2.0. The Grid-Tools framework used as a domain specific language approach for MPDATA is available at <https://github.com/GridTools/gridtools> under a
15 BSD-3-Clause license. The CLAW DSL used for the cloud microphysics dwarf is available at <https://github.com/claw-project/claw-compiler> under a BSD-2-Clause license. Model codes developed at ECMWF are the intellectual property of ECMWF and its member states, and therefore the IFS code and the IFS-FVM code are not publicly available. Access to a reduced version of the IFS code may be obtained from ECMWF under an OpenIFS license (see <http://www.ecmwf.int/en/research/projects/openifs> for further information, last access: 28 May 2019). The ALARO and ALADIN codes, along with all their related intellectual property rights, are owned by the members of the ALADIN
20 consortium and are shared with the members of the HIRLAM consortium in the frame of a cooperation agreement. This agreement allows each member of either consortium to license the shared ALADIN-HIRLAM codes to academic institutions of their home country for noncommercial research. Access to the codes of the ALADIN System can be obtained by contacting one of the member institutes or by submitting a request in the contact link below the page of the ALADIN website (<http://www.umr-cnrm.fr/aladin/>, last access: 28 May 2019) and the access will be subject to signing a standardised ALADIN-HIRLAM license agreement. The COSMO-EULAG model, along with all of its related in-
25 tellectual property rights, is owned by the members of the COSMO consortium under the cooperation agreement. This agreement allows each member of the COSMO consortium to license the COSMO-EULAG code (<http://www.cosmo-model.org/content/consortium/licencing.htm>, last access: 28 May 2019) without fee to academic institutions of their home country for noncommercial research. The code of the GRASS model is intellectual property of Météo-France and is not publicly available. The code of the ESCAPE dwarfs is intellectual property of ECMWF. A license for educational and non-commercial research can be obtained from ECMWF (see <http://www.hpc-escape.eu> for contact
30 details, last access: 28 May 2019). For the GPU optimisation of the spectral transform dwarf we used the PGI compiler version 17.10, CUDA 9.0.176 and OpenMPI 2.1.3. For the GPU optimisation of the MPDATA dwarf we used the PGI compiler version 17.10, CUDA 9.2.88 and OpenMPI 2.1.3. CUDA includes the compiler nvcc used to compile the library function wrappers, the libraries themselves and the profiling tool nvprof which we used for profiling on the GPUs. For the CPU optimisation we used Intel compilers and libraries version 2018.1.163. This includes the compilers icc and ifort and the libraries mkl and mpi. The work on the optical processor Optalysys used the MATLAB
35 software version 2017b.

Author contributions. Andreas Müller performed energy measurements for the spectral transform dwarf, implemented the optimised spectral transform for post-processing, worked on the implementation of the MPDATA dwarf, supported maintenance for all dwarf prototypes and assembled the paper. Willem Deconinck is the main developer of the data structure framework Atlas and contributed the work on improving and extending Atlas. He also was involved in creating, supporting and maintaining all of the dwarf prototypes. Christian Kühnlein and Piotr Smolarkiewicz are the main developers of IFS-FVM and contributed to the corresponding elliptic solver and MPDATA dwarfs. Gianmarco Mengaldo was involved in the creation of the dwarfs and provided improvements to early drafts of the manuscript. Michael Lange and Valentin Clement provided DSL work for the cloud microphysics dwarf. Nils Wedi and Peter Bauer supervised the entire ESCAPE project as coordinators. They were also involved in writing and improving the manuscript. Nils Wedi performed the measurements for the cost profiles in Figure 3. Michail Diamantakis implemented the semi-Lagrangian dwarf prototype. Sarah-Jane Lock supported the work on HEVI time integration methods. Mats Hamrud, Sami Saarinen and George Mozdzynski were involved in creating initial versions of the dwarf prototypes. Sami Saarinen also contributed the code for energy measurements on the Cray XC40 supercomputer. Daniel Thiemert provided project management for the ESCAPE project. Michael Ginton and Pierre Bénard implemented the GRASS model. Fabrice Voitus and Charles Colavolpe worked on HEVI time-stepping methods. Philippe Marguinaud and Yongjun Zheng performed simulations of communication costs. Joris Van Bever, Daan Degrauwe, Geert Smet and Piet Termonia provided energy and runtime measurements for ALARO, analysed the components of the RMI ensemble prediction system and implemented limited area support in Atlas. Kristian P. Nielsen, Bent H. Sass, Jacob W. Poulsen and Per Berg provided work on the radiation dwarf. Carlos Osuna and Oliver Fuhrer provided work on GridTools and were involved in adding GPU support to Atlas. Michael Baldauf contributed work on HEVI DG. Mike Gillard and Joanna Szmelter worked on improving the preconditioning for the elliptic solver. Enda O’Brien, Alastair McKinstry, Oisín Robinson, Parijat Shukla and Michael Lysaght provided the LAITRI dwarf and created initial CPU optimisations and OpenACC ports for all of the dwarf prototypes. Michał Kulczewski, Milosz Ciznicki, Wojciech Piątek, Sebastian Ciesielski, Marek Błażewicz, Krzysztof Kurowski, Marcin Procyk, Pawel Spychala and Bartosz Bosak worked on performance modelling using the roofline model which allowed a better understanding of the experimental measurements. Zbigniew Piotrowski and Andrzej Wyszogrodzki compared different time-integration methods, supported optimisation of IFS-FVM and provided energy measurements of COSMO-EULAG. Erwan Raffin, Cyril Mazauric, David Guibert, Louis Douriez and Xavier Vigouroux provided the CPU optimisations for the spectral transform and MPDATA dwarf prototypes. Alan Gray and Peter Messmer contributed the GPU optimisations for the spectral transform and MPDATA dwarf prototypes. Alexander J. Macfaden and Nick New contributed the work on using optical processors for the spectral transform.

Acknowledgements. [The authors thank three anonymous referees for their thorough reviews and many detailed suggestions that helped to improve the presentation.](#) The ESCAPE project has received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement No 671627, see also <http://www.hpc-escape.eu/>. This work has been also supported by HPC resources provided by Poznan Supercomputing and Networking Center and corresponding research activities under the MAESTRO grant number DEC-2013/08/A/ST6/00296 in National Science Centre (NCN). An award of computer time was provided by the INCITE program. This research used resources of the Oak Ridge Leadership Computing Facility, which is a DOE office of Science User Facility supported under contract DE-AC05-00OR22725.

References

- Asanović, K., Bodik, R., Catanzaro, B. C., Gebis, J. J., Husbands, P., Keutzer, K., Patterson, D. A., Plishker, W. L., Shalf, J., Williams, S. W., and Yelick, K. A.: The landscape of parallel computing research: a view from Berkeley, Tech. Rep. UCB/EECS-2006-183, EECS Department, University of California, Berkeley, <http://www2.eecs.berkeley.edu/Pubs/TechRpts/2006/EECS-2006-183.pdf>, 2006.
- 5 Asanović, K., Wawrzynek, J., Wessel, D., Yelick, K., Bodik, R., Demmel, J., Keaveny, T., Keutzer, K., Kubiawicz, J., Morgan, N., Patterson, D., and Sen, K.: A view of the parallel computing landscape, *Comm. ACM*, 52, 56, <https://doi.org/10.1145/1562764.1562783>, 2009.
- Bauer, P., Thorpe, A., and Brunet, G.: The quiet revolution of numerical weather prediction, *Nature*, 525, 47–55, <https://doi.org/10.1038/nature14956>, 2015.
- Bénard, P. and Glinton, M.: Circumventing the pole problem of reduced lat-lon grids with local schemes. Part I: analysis and model formulation, *Q. J. Royal Meteorol. Soc.*, 145, 1377–1391, <https://doi.org/10.1002/qj.3509>, 2019.
- 10 Clement, V., Ferrachat, S., Fuhrer, O., Lapillonne, X., Osuna, C. E., Pincus, R., Rood, J., and Sawyer, W.: The CLAW DSL, in: Proceedings of the Platform for Advanced Scientific Computing Conference - PASC '18, ACM Press, <https://doi.org/10.1145/3218176.3218226>, 2018.
- Colavolpe, C., Voitus, F., and Bénard, P.: RK-IMEX HEVI schemes for fully compressible atmospheric models with advection: analyses and numerical testing, *Q. J. Royal Meteorol. Soc.*, 143, 1336–1350, <https://doi.org/10.1002/qj.3008>, 2017.
- 15 Colella, P.: Defining software requirements for scientific computing, DARPA HPCS Presentation, 2004.
- Deconinck, W.: Development of Atlas, a flexible data structure framework, Tech. rep., ECMWF, <http://arxiv.org/abs/1908.06091>, 2017a.
- Deconinck, W.: Public release of Atlas under an open source license, which is accelerator enabled and has improved interoperability features, Tech. rep., ECMWF, <http://arxiv.org/abs/1908.07038>, 2017b.
- Deconinck, W., Bauer, P., Diamantakis, M., Hamrud, M., Kühnlein, C., Maciel, P., Mengaldo, G., Quintino, T., Raoult, B., Smolarkiewicz, P. K., and Wedi, N. P.: Atlas - a library for numerical weather prediction and climate modelling, *Comput. Phys. Commun.*, 220, 188–204, <https://doi.org/10.1016/j.cpc.2017.07.006>, 2017.
- 20 Douriez, L., Gray, A., Guibert, D., Messmer, P., and Raffin, E.: Performance report and optimized implementations of weather & climate dwarfs on multi-node systems, Tech. rep., ECMWF, <http://arxiv.org/abs/1908.06097>, 2018.
- Dziekan, P., Waruszewski, M., and Pawlowska, H.: University of Warsaw Lagrangian Cloud Model (UWLCM) 1.0: a modern large-eddy simulation tool for warm cloud modeling with Lagrangian microphysics, *Geosci. Model Dev.*, pp. 2587–2606, <https://doi.org/10.5194/gmd-12-2587-2019>, 2019.
- 25 Feng, W., Lin, H., Scogland, T., and Zhang, J.: OpenCL and the 13 dwarfs: a work in progress, in: Proceedings of the third joint WOSP/SIPEW International Conference on Performance Engineering - ICPE '12, ACM Press, <https://doi.org/10.1145/2188286.2188341>, 2012.
- 30 Flamm, K.: Measuring Moore’s law: evidence from price, cost, and quality indexes, National Biuro of Economic Research, Working Paper No. 24553, April 2018, <https://doi.org/10.3386/w24553>, 2018.
- Fuhrer, O., Chadha, T., Hoefler, T., Kwasniewski, G., Lapillonne, X., Leutwyler, D., Lüthi, D., Osuna, C., Schär, C., Schulthess, T. C., and Vogt, H.: Near-global climate simulation at 1km resolution: establishing a performance baseline on 4888 GPUs with COSMO 5.0, *Geosci. Model Dev.*, 11, 1665–1681, <https://doi.org/10.5194/gmd-11-1665-2018>, 2018.
- 35 Glinton, M. R. and Bénard, P.: Circumventing the pole problem of reduced lat-lon grids with local schemes. Part II: validation experiments, *Q. J. Royal Meteorol. Soc.*, 145, 1392–1405, <https://doi.org/10.1002/qj.3495>, 2019.

- Johnston, B. and Milthorpe, J.: Dwarfs on accelerators: enhancing OpenCL benchmarking for heterogeneous computing architectures, in: Proceedings of the 47th International Conference on Parallel Processing Companion - ICPP '18, ACM Press, <https://doi.org/10.1145/3229710.3229729>, 2018.
- 5 Kaltofen, E. L.: The “seven dwarfs” of symbolic computation, in: Numerical and Symbolic Scientific Computing, edited by Langer, U. and Paule, P., pp. 95–104, Springer, https://doi.org/10.1007/978-3-7091-0794-2_5, 2011.
- Katzav, J. and Parker, W. S.: The future of climate modeling, *Clim. Change*, 132, 475–487, <https://doi.org/10.1007/s10584-015-1435-x>, 2015.
- Krommydas, K., Feng, W., Antonopoulos, C. D., and Bellas, N.: OpenDwarfs: characterization of dwarf-based benchmarks on fixed and reconfigurable architectures, *J. Signal Process. Syst.*, 85, 373–392, <https://doi.org/10.1007/s11265-015-1051-z>, 2015.
- 10 Kühnlein, C., Deconinck, W., Klein, R., Malardel, S., Piotrowski, Z. P., Smolarkiewicz, P. K., Szmelter, J., and Wedi, N. P.: FVM 1.0: a nonhydrostatic finite-volume dynamical core formulation for IFS, *Geosci. Model Dev.*, 12, 651–676, <https://doi.org/10.5194/gmd-12-651-2019>, 2019.
- Lawrence, B. N., Rezny, M., Budich, R., Bauer, P., Behrens, J., Carter, M., Deconinck, W., Ford, R., Maynard, C., Mullerworth, S., Osuna, C., Porter, A., Serradell, K., Valcke, S., Wedi, N., and Wilson, S.: Crossing the chasm: how to develop weather and climate models for next generation computers?, *Geosci. Model Dev.*, 11, 1799–1821, <https://doi.org/10.5194/gmd-11-1799-2018>, 2018.
- 15 Macfaden, A. J., Gordon, G. S. D., and Wilkinson, T. D.: An optical Fourier transform coprocessor with direct phase determination, *Sci. Rep.*, 7, <https://doi.org/10.1038/s41598-017-13733-1>, 2017.
- Mazauric, C., Raffin, E., and Guibert, D.: Recommendations and specifications for data scope analysis tools, Tech. rep., ECMWF, <http://arxiv.org/abs/1908.06095>, 2017a.
- Mazauric, C., Raffin, E., Vigouroux, X., Guibert, D., Macfaden, A., Poulsen, J., Berg, P., Gray, A., and Messmer, P.: Performance report and optimized implementation of weather & climate dwarfs on GPU, MIC and Optalysys optical processor, Tech. rep., ECMWF, <http://arxiv.org/abs/1908.06096>, 2017b.
- 20 Mengaldo, G.: Batch 1: definition of several weather & climate dwarfs, Tech. rep., ECMWF, <http://arxiv.org/abs/1908.06089>, 2016.
- Mengaldo, G., Wyszogrodzki, A., Diamantakis, M., Lock, S.-J., Giraldo, F., and Wedi, N. P.: Current and emerging time-integration strategies in global numerical weather and climate prediction, *Arch. Comput. Methods Eng.*, pp. 1–22, <https://doi.org/10.1007/s11831-018-9261-8>, 25 2018.
- Messer, O. E. B., D’Azevedo, E., Hill, J., Joubert, W., Berrill, M., and Zimmer, C.: MiniApps derived from production HPC applications using multiple programming models, *Int. J. High Perform. Comput. Appl.*, 32, 582–593, <https://doi.org/10.1177/1094342016668241>, 2016.
- Michalakes, J., Govett, M., Benson, R., Black, T., Juang, H., Reinecke, A., and Skamarock, B.: AVEC Report: NGGPS Level-1 Benchmarks and software evaluation, Tech. rep., NOAA, Boulder, US, <https://repository.library.noaa.gov/view/noaa/18654>, 2015.
- 30 Mozdzyński, G., Hamrud, M., and Wedi, N. P.: A partitioned global address space implementation of the European Centre for Medium Range Weather Forecasts Integrated Forecasting System, *Int. J. High Perform. Comput. Appl.*, 29, 261–273, <https://doi.org/10.1177/1094342015576773>, 2015.
- Müller, A., Kopera, M. A., Marras, S., Wilcox, L. C., Isaac, T., and Giraldo, F. X.: Strong scaling for numerical weather prediction at petascale with the atmospheric model NUMA, *Int. J. High Perform. Comput. Appl.*, 2, 411–426, <https://doi.org/10.1177/1094342018763966>, 2018.
- 35 Müller, A., Gillard, M., Nielsen, K. P., and Piotrowski, Z.: Batch 2: definition of novel weather & climate dwarfs, Tech. rep., ECMWF, <http://arxiv.org/abs/1908.07040>, 2017.

- Neumann, P., Düben, P., Adamidis, P., Bauer, P., Brüch, M., Kornbluh, L., Klocke, D., Stevens, B., Wedi, N., and Biercamp, J.: Assessing the scales in numerical weather and climate predictions: will exascale be the rescue?, *Philos. Trans. Royal Soc. A: Mathematical, Physical and Engineering Sciences*, 377, <https://doi.org/10.1098/rsta.2018.0148>, 2019.
- Osuna, C.: Report on the performance portability demonstrated for the relevant weather & climate dwarfs, Tech. rep., ECMWF, <http://arxiv.org/abs/1908.06094>, 2018.
- 5 Palmer, T.: Climate forecasting: Build high-resolution global climate models, *Nature*, 515, 338–339, <https://doi.org/10.1038/515338a>, 2014.
- Phillips, S. C., Engen, V., and Papay, J.: Snow white clouds and the seven dwarfs, in: 2011 IEEE Third International Conference on Cloud Computing Technology and Science, IEEE, <https://doi.org/10.1109/cloudcom.2011.114>, 2011.
- Poulsen, J. W. and Berg, P.: Tuning the implementation of the radiation scheme ACRANEB2, Tech. rep., DMI report 17-22, http://www.dmi.dk/fileadmin/user_upload/Rapporter/TR/2017/SR17-22.pdf, 2017.
- 10 Robinson, O., McKinsty, A., and Lysaght, M.: Optimization of IFS subroutine LAITRI on Intel Knights Landing, Tech. rep., PRACE White Papers, <https://doi.org/10.5281/zenodo.832025>, 2016.
- Schalkwijk, J., Jonker, H. J. J., Siebesma, A. P., and Meijgaard, E. V.: Weather forecasting using GPU-based Large-Eddy Simulations, *Bull. Am. Meteorol. Soc.*, 96, 715–723, <https://doi.org/10.1175/bams-d-14-00114.1>, 2015.
- 15 Schulthess, T. C., Bauer, P., Wedi, N., Fuhrer, O., Hoefler, T., and Schär, C.: Reflecting on the goal and baseline for exascale computing: a roadmap based on weather and climate simulations, *Comput. Sci. Eng.*, 21, 30–41, <https://doi.org/10.1109/mcse.2018.2888788>, 2019.
- Shukla, J., Palmer, T. N., Hagedorn, R., Hoskins, B., Kinter, J., Marotzke, J., Miller, M., and Slingo, J.: Toward a new generation of world climate research and computing facilities, *Bull. Am. Meteorol. Soc.*, 91, 1407–1412, <https://doi.org/10.1175/2010bams2900.1>, 2010.
- Van Bever, J., McFaden, A., Piotrowski, Z., and Degrauwe, D.: Report on energy-efficiency evaluation of several NWP model configurations, Tech. rep., ECMWF, <http://arxiv.org/abs/1908.06115>, 2018.
- 20 Wallemacq, P. and House, R.: Economic losses, poverty and disasters 1998-2017, Tech. rep., <https://www.unisdr.org/we/inform/publications/61119>, 2018.
- Wedi, N. P., Hamrud, M., and Mozdzyński, G.: A fast spherical harmonics transform for global NWP and climate models, *Mon. Weather Rev.*, 141, 3450–3461, <https://doi.org/10.1175/mwr-d-13-00016.1>, 2013.
- 25 Wedi, N. P., Bauer, P., Deconinck, W., Diamantakis, M., Hamrud, M., Kühnlein, C., Malardel, S., Mogensen, K., Mozdzyński, G., and Smolarkiewicz, P.: The modelling infrastructure of the Integrated Forecasting System: recent advances and future challenges, Tech. Rep. 760, Eur. Cent. For Medium-Range Weather Forecasts, Reading, UK, <https://doi.org/10.21957/thtpwp67e>, 2015.
- Wehner, M. F., Oliker, L., Shalf, J., Donofrio, D., Drummond, L. A., Heikes, R., Kamil, S., Kono, C., Miller, N., Miura, H., Mohiyuddin, M., Randall, D., and Yang, W.-S.: Hardware/software co-design of global cloud system resolving models, *J. Adv. Model. Earth Sys.*, 3, <https://doi.org/10.1029/2011ms000073>, 2011.
- 30 Williamson, D. L.: The evolution of dynamical cores for global atmospheric models, *J. Meteor. Soc. Japan*, 85B, 241–269, <https://doi.org/10.2151/jmsj.85b.241>, 2007.
- World Economic Forum: The 2019 global risks report, <https://www.weforum.org/reports/the-global-risks-report-2019>, 2019.
- Wu, J., Wyckoff, P., and Panda, D.: High performance implementation of MPI derived datatype communication over InfiniBand, in: 18th International Parallel and Distributed Processing Symposium, 2004. Proceedings., IEEE, <https://doi.org/10.1109/ipdps.2004.1302917>, 2004.
- 35 Xiao, H., Diamantakis, M., and Saarinen, S.: An OpenACC GPU adaptation of the IFS cloud microphysics scheme, ECMWF Tech. Memo. No. 805, <https://doi.org/10.21957/g9mjjlgeq>, 2017.

Zheng, Y. and Marguinaud, P.: Simulation of the performance and scalability of MPI communications of atmospheric models running on exascale supercomputers, *Geosci. Model Dev.*, 11, 3409–3426, <https://doi.org/10.5194/gmd-11-3409-2018>, 2018.