

**T/NA086/20**  
**Code structure and coordination**

Report 2047358-TN-03

*Evaluation of Approaches to Performance  
Portability*

Steven Wright, Ben Dudson, Peter Hill, and David Dickinson

*University of York*

Gihan Mudalige

*University of Warwick*

July 28, 2021

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Application Evaluations</b>	<b>4</b>
2.1	TeaLeaf . . . . .	4
2.1.1	Performance . . . . .	4
2.1.2	Performance Portability . . . . .	5
2.2	miniFE . . . . .	8
2.2.1	Performance . . . . .	9
2.2.2	Performance Portability . . . . .	10
2.3	Laghos . . . . .	11
2.3.1	Performance . . . . .	11
2.3.2	Performance Portability . . . . .	12
2.4	CabanaPIC . . . . .	14
2.4.1	Performance . . . . .	14
2.5	VPIC . . . . .	15
2.5.1	Performance . . . . .	15
2.5.2	Performance Portability . . . . .	15
2.6	EMPIRE-PIC . . . . .	17
2.6.1	Performance . . . . .	18
2.6.2	Performance Portability . . . . .	18
<b>3</b>	<b>Conclusions</b>	<b>22</b>
3.1	Limitations . . . . .	24
	<b>References</b>	<b>24</b>

# 1 Introduction

The focus of the *code structure and coordination* work package is to establish a series of “best practices” on how to develop simulation applications for Exascale systems that are able to obtain high performance on each architecture (i.e. are performance portable) without significant manual porting efforts.

In the past decade, a large number of approaches to developing performance portable code have been developed. In this report we will begin to report on our evaluation of some of these approaches through the execution of a small number of mini-applications that implement methods similar to those likely to be required in NEPTUNE.

These applications are detailed in report 2047358-TN-02, but are summarised below for convenience:

## **TeaLeaf**

A finite difference mini-app that solves the linear heat conduction equation on a regular grid using a 5-point stencil<sup>1</sup>.

## **miniFE**

A finite element mini-app, and part of the Mantevo benchmark suite<sup>2</sup>.

## **Laghos**

A high-order curvilinear finite element scheme on an unstructured mesh<sup>3</sup>.

## **Nekbone**

A high-order spectral element application for solving the incompressible Navier-Stokes equations.<sup>4</sup>

## **CabanaPIC**

A structured PIC code built using the CoPA Cabana library for particle-based simulations<sup>5</sup>.

## **VPIC/VPIC 2.0**

A general purpose PIC code for modelling kinetic plasmas in one, two or

---

<sup>1</sup><http://uk-mac.github.io/TeaLeaf/>

<sup>2</sup><https://github.com/Mantevo/miniFE>

<sup>3</sup><https://github.com/CEED/Laghos>

<sup>4</sup><https://github.com/Nek5000/Nekbone>

<sup>5</sup><https://github.com/ECP-copa/CabanaPIC>

three dimensions, developed at Los Alamos National Laboratory<sup>6</sup>.

### **EMPIRE-PIC**

An unstructured PIC code that uses the finite-element method.

The selected applications broadly represent the algorithms of interest for the NEPTUNE project and fall in to two categories – fluid-methods and particle-methods. Within the fluid-method tranche, the applications are available implemented in a wide range of programming models, allowing us a good opportunity to evaluate the effect of programming model on the performance, and importantly the *performance portability* of that particular approach to application development. There are a relatively small number of particle-in-cell mini-applications available, and thus the selected particle-methods applications are only available implemented using Kokkos. However, this still allows us an opportunity to evaluate the appropriateness of Kokkos as a programming model for performance portable application development.

As stated previously, we will evaluate the performance portability of these applications using the metric introduced by Pennycook et al. [1], and use the visualisation techniques outlined by Sewall et al. [2].

Where possible, performance data has been taken from previously published works. Where no data exists, the data has been collected from the UK’s Tier-2 platforms, in particular Isambard’s Multi-Architecture Comparison System (MACS), ThunderX2 system and A64FX system.

As many of the applications, libraries and programming models used in this report are under active development, the data presented here is subject to change. New data is being collected all the time and analysed, and will be updated in the future where necessary. This document should therefore be considered a living document, reflecting the current state of performance portable application development focused on applications of interest for the simulation of plasma physics.

---

<sup>6</sup><https://github.com/lanl/vpic>

## 2 Application Evaluations

In this section we present performance data for a number of mini-applications, across a range of architectural platforms, using a range of different approaches to performance portability.

The applications chosen in each case are broadly representative of some of the algorithms of interest to NEPTUNE. In particular, the fluid-method based mini-apps implement algorithms that range from finite-difference (like Bout++ [3]) to high-order finite element or spectral element (like Nektar++ [4]). Similarly, the particle-methods mini-apps all implement the particle-in-cell method (like EPOCH [5]).

### 2.1 TeaLeaf

TeaLeaf is a finite difference mini-app that solves the linear heat conduction equation on a regular grid using a 5-point stencil, developed as part of the UK-MAC (UK Mini-App Consortium) project.

It has been used extensively in studying performance portability already [6, 7, 8, 9], and is available implemented using CUDA, OpenACC, OPS, RAJA, and Kokkos, among others<sup>7</sup>. The results in this section are extracted from two of these studies, namely one by Kirk et al. [7] and one by Deakin et al. [6].

In both studies, the largest test problem size (`tea_bm.5.in`) is used, a  $4000 \times 4000$  grid.

#### 2.1.1 Performance

The study by Kirk et al. executes 8 different implementation/configurations of TeaLeaf across 3 platforms, a dual Intel Broadwell system, an Intel KNL system and an NVIDIA P100 system. The raw runtime figures are presented in Figure 1. Note that in the study, some results are missing due to incompatibility (e.g. CUDA on Broadwell/KNL).

---

<sup>7</sup><http://uk-mac.github.io/TeaLeaf/>

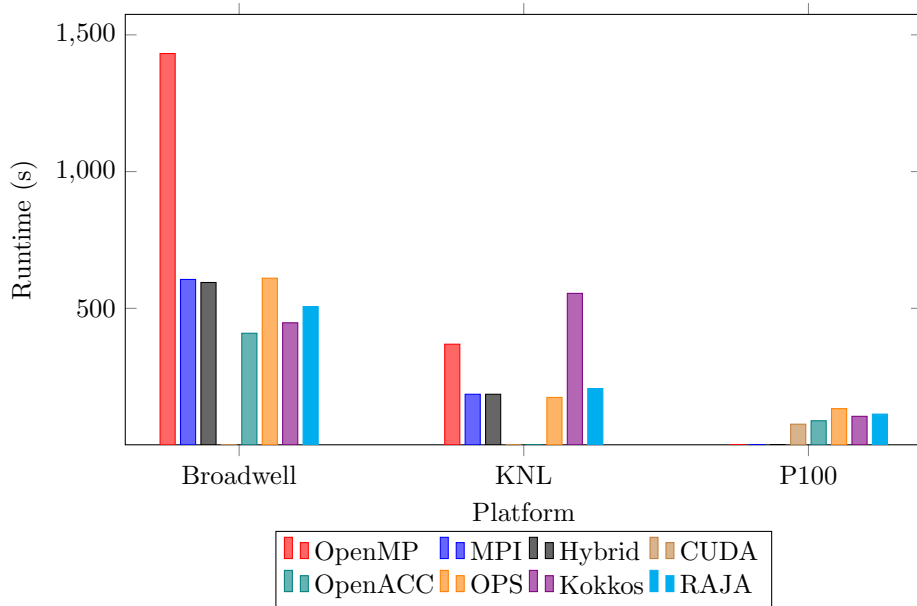


Figure 1: TeaLeaf runtime data from Kirk et al. [7]

The study by Deakin et al. is more recent, using a C-based implementation of TeaLeaf as its base. It consequently evaluates less programming models, but over a wider range of hardware, including a dual Intel Skylake system, both NVIDIA P100 and V100 systems, AMDs Naples CPU, and the Arm-based ThunderX2 platform. Runtime results are provided in Figure 2.

### 2.1.2 Performance Portability

Both studies evaluate some portable and non-portable implementations. In most cases, there is a non-portable implementation that achieves the lowest runtime, however this places a restriction on the hardware that it can target.

For study by Kirk et al. [7], Figures 3 and 4 allow us to visualise the performance portability of each approach to application development. The figures show a clear divide between portable approaches (Kokkos, OPS and RAJA), and the non-portable approaches (CUDA, OpenMP and MPI). While typically higher performance can be achieved non-portably, each of these programming models tightly binds developers to particular architectures.

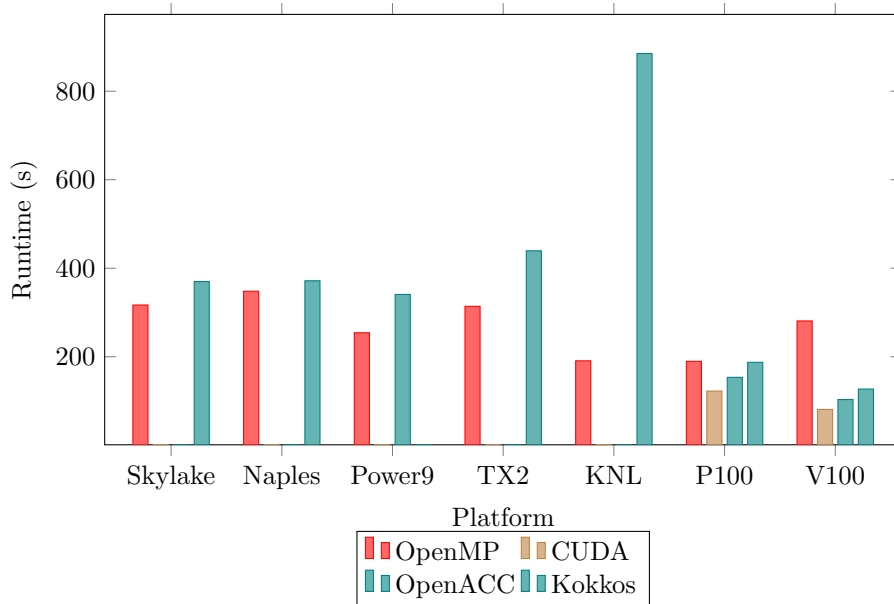


Figure 2: TeaLeaf runtime data from Deakin et al. [6]

Of the portable approaches, the best portability is achieved by RAJA, followed by Kokkos, and then the domain specific language, OPS. Referring back to Figure 1, we can see that on the Intel KNL system, the Kokkos performance is double that of other performance portable approaches, and thus skews the portability calculation. It is likely that this is the result an unidentified issue in TeaLeaf or Kokkos at the time of evaluation. Otherwise, these three programming models each achieve similar levels of performance and, importantly, portability across different architectures.

Figures 5 and 6 show the same visualisations for the data from Deakin et al. [6]. Again, the non-portable programming model (CUDA) achieves the highest performance on its target architecture. For CPU architectures OpenMP produces the highest result, and using offload directives, portability is available to GPU devices. It should be noted that to support the use of GPU devices, there are two OpenMP implementations that must be maintained (with and without offload directives), though these results are presented together here. Much like in the previous study, the performance portability of Kokkos is affected by an anomalous result on the Intel KNL platform.

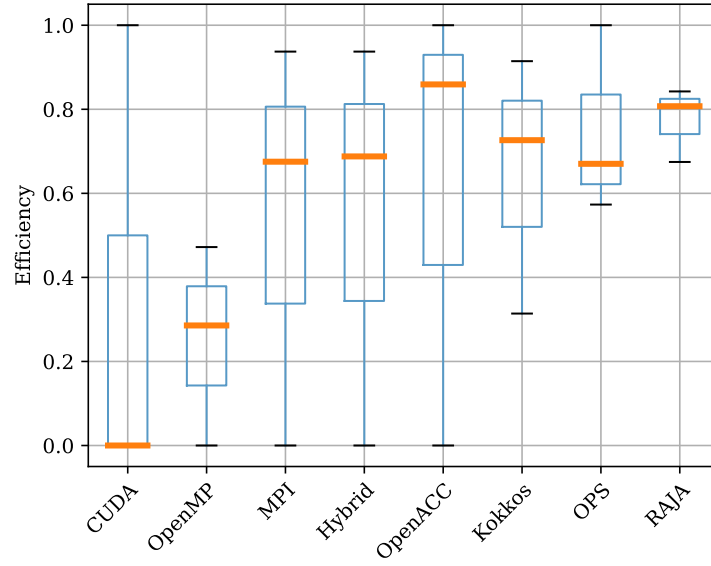


Figure 3: Box plot visualisation of performance portability from Kirk et al. [7]

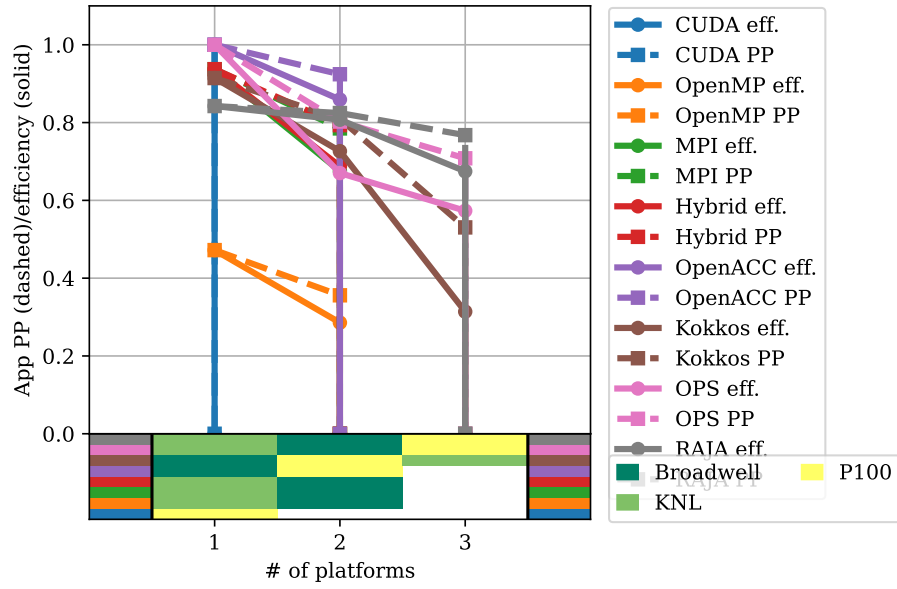


Figure 4: Cascade visualisation of performance portability from Kirk et al. [7]



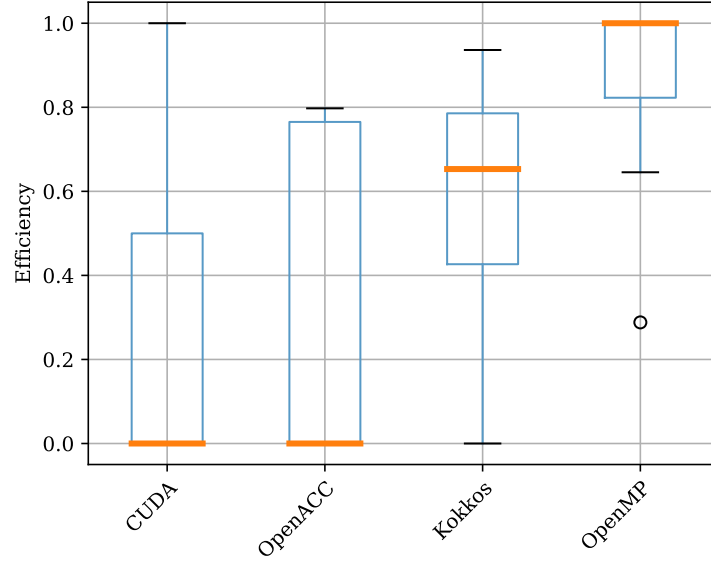


Figure 5: Box plot visualisation of performance portability from Deakin et al. [6]

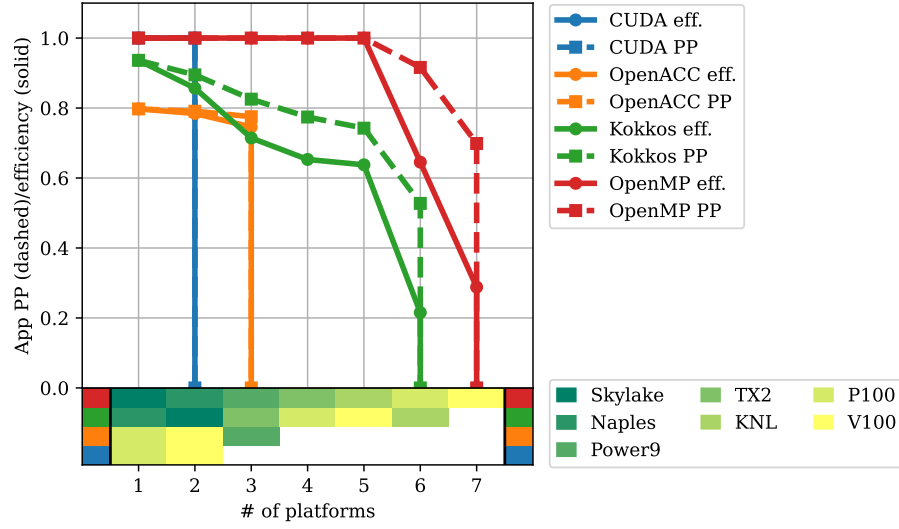


Figure 6: Cascade visualisation of performance portability from Deakin et al. [6]

## 2.2 miniFE

miniFE is a finite element mini-app, and part of the Mantevo benchmark suite [10, 11, 12, 13]. It implements an unstructured implicit finite element method

and has versions available in CUDA, Kokkos, OpenMP (3.0+ and 4.5+) and SYCL<sup>8</sup>.

While there are a number of data sources for miniFE data, many of these are limited in scope, and so to ensure consistency, all data presented in this section has been newly gathered. In all cases, a  $256 \times 256 \times 256$  problem size has been used, and all runs have been conducted on the platforms available on Isambard.

### 2.2.1 Performance

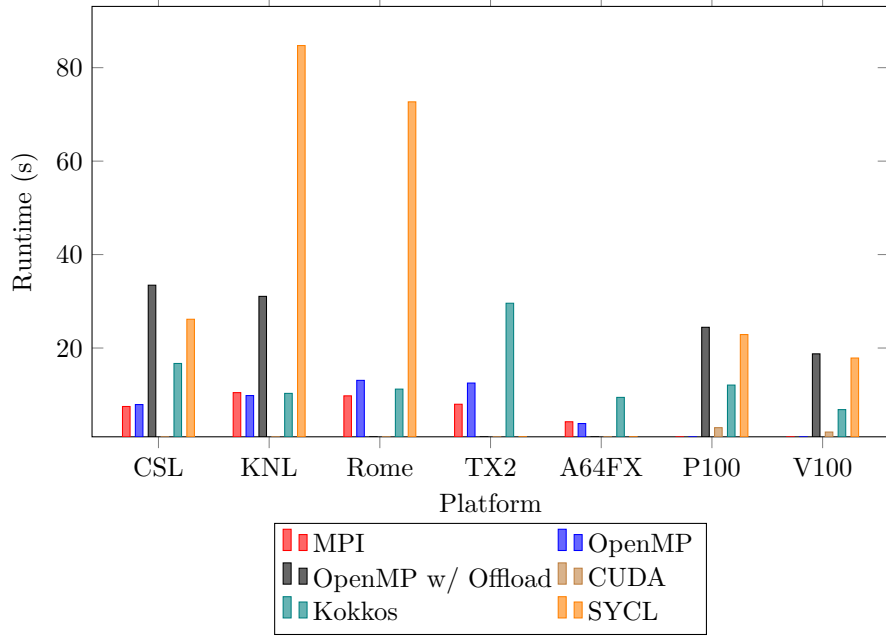


Figure 7: miniFE runtime data

The raw runtime results for these runs can be seen in Figure 7. It should be noted that the SYCL data is gathered from a miniFE port that can be found as part of the oneAPI-DirectProgramming github repository<sup>9</sup>, and is compiled using hipSYCL and GCC. Data has not yet been collected for the ARM-based systems with SYCL. The OpenMP with offload variant of miniFE runs successfully on both AMD Rome and Cavium ThunderX2 platforms, but the runtimes

<sup>8</sup><https://github.com/Mantevo/miniFE>

<sup>9</sup><https://github.com/zjin-lcf/oneAPI-DirectProgramming/tree/master/miniFE-sycl>

are several orders of magnitude greater than all other platforms (likely due to a bug in the compiled code), and so have been removed.

In many of the miniFE ports available, only the conjugate solver has been parallelised effectively, so the results presented here represent only the timing from this kernel.

### 2.2.2 Performance Portability

Figures 8 and 9 present visualisations of the performance portability of miniFE, through various approaches.

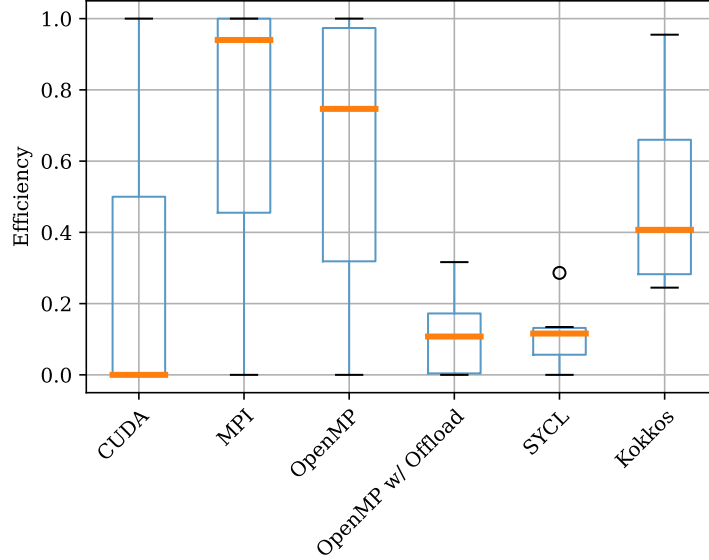


Figure 8: Box plot visualisation of performance portability of miniFE

Again, non-portable approaches (CUDA, OpenMP and MPI) often provide the best performance, but this leads to a significant restriction on the platforms that can be used. Figure 9 shows how the performance portability of miniFE evolves as more platforms are added for the Kokkos variant. While the performance is lower than native implementations, it has the advantage of being able to target every platform from a single codebase.

Likewise, both OpenMP with Offload and SYCL can target every platform,

though in some cases, we have not yet collected data (SYCL on Arm platforms), and in some cases, compiler issues are preventing reasonable performance.

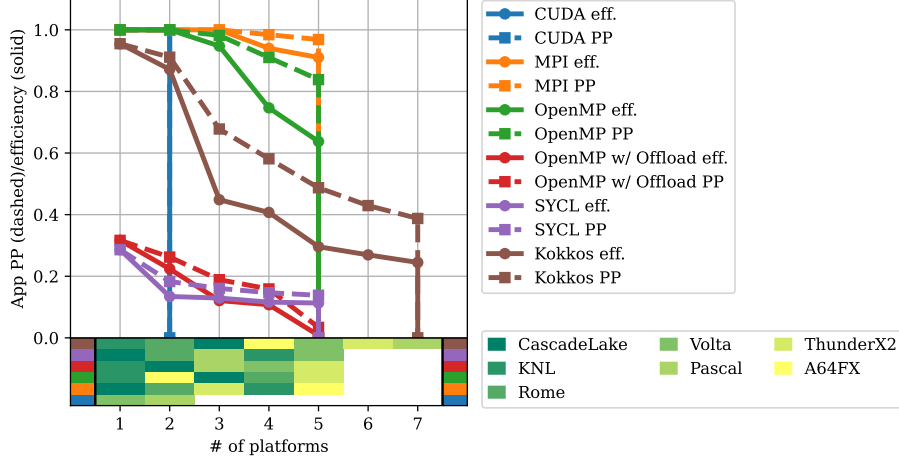


Figure 9: Cascade visualisation of performance portability of miniFE

## 2.3 Laghos

Laghos is a mini-app that is part of the ECP Proxy Applications suite [14, 15, 13]. It implements a high-order curvilinear finite element scheme on an unstructured mesh. The majority of the computation is performed by the HYPRE and MFEM libraries, and can thus use any programming model that is available for these libraries<sup>10</sup>.

The results presented below have all been collected from the Isambard platform.

### 2.3.1 Performance

Figure 10 shows the runtime for Laghos, running problem #1 (Sedov blast wave), in three dimensions, up to 1.0 second of simulated time, using partial assembly (i.e., `./laghos -p 1 -dim 3 -rs 2 -tf 1.0 -pa -f`).

<sup>10</sup><https://github.com/CEED/Laghos>

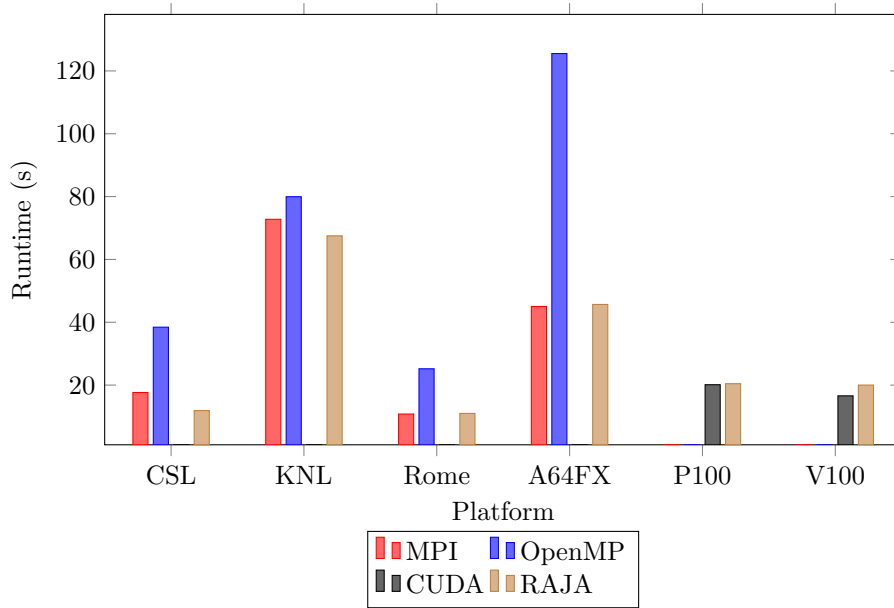


Figure 10: Laghos runtime data

### 2.3.2 Performance Portability

Portability visualisations of each implementation of Laghos are provided in Figures 11 and 12.

Again, it is clear that the highest achievable performance is often not from a portable approach. However, RAJA achieves the highest performance portability, due to being able to span each of the platforms, and provide near equal performance to the CUDA and MPI variants, and much better performance than the OpenMP implementation.

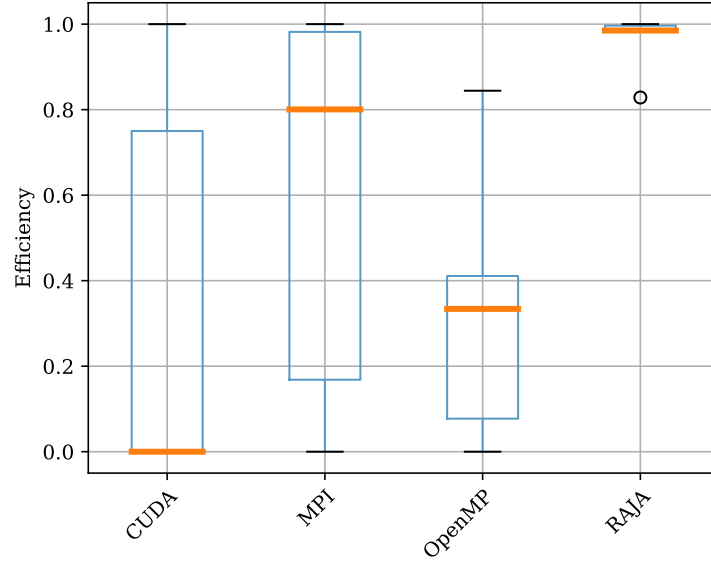


Figure 11: Box plot visualisation of performance portability of Laghos

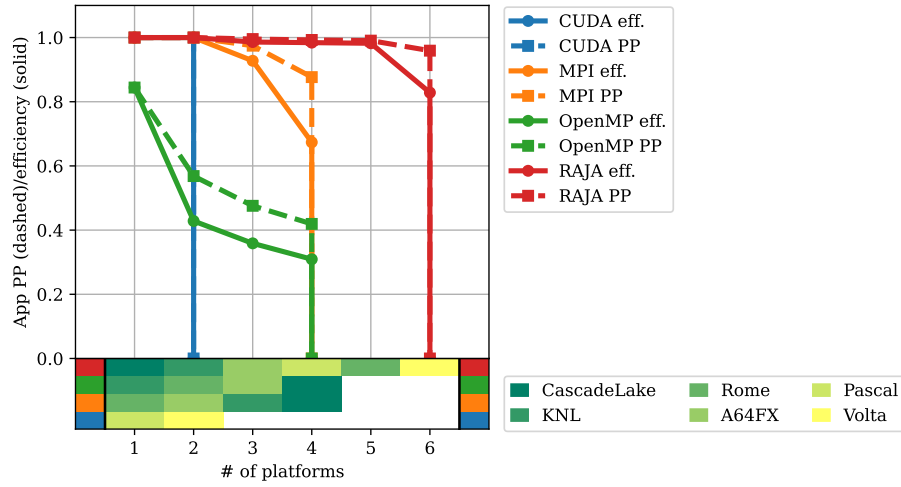


Figure 12: Cascade visualisation of performance portability of Laghos

## 2.4 CabanaPIC

CabanaPIC is a structured PIC demonstrator application built using the Co-PA/Cabana library for particle-based simulations [13]. The application uses Kokkos as its programming model for on-node parallelism and GPU use, and uses MPI for off-node parallelism<sup>11</sup>.

### 2.4.1 Performance

Since there is only a single implementation of CabanaPIC, it is not possible for us to evaluate how the programming model affects its performance portability, however, we can show how its performance changes between architectures.

Figure 13 shows the achieved runtime for CabanaPIC across four of Isambard’s platforms, running a simple 1D 2-stream problem with 6.4 million particles.

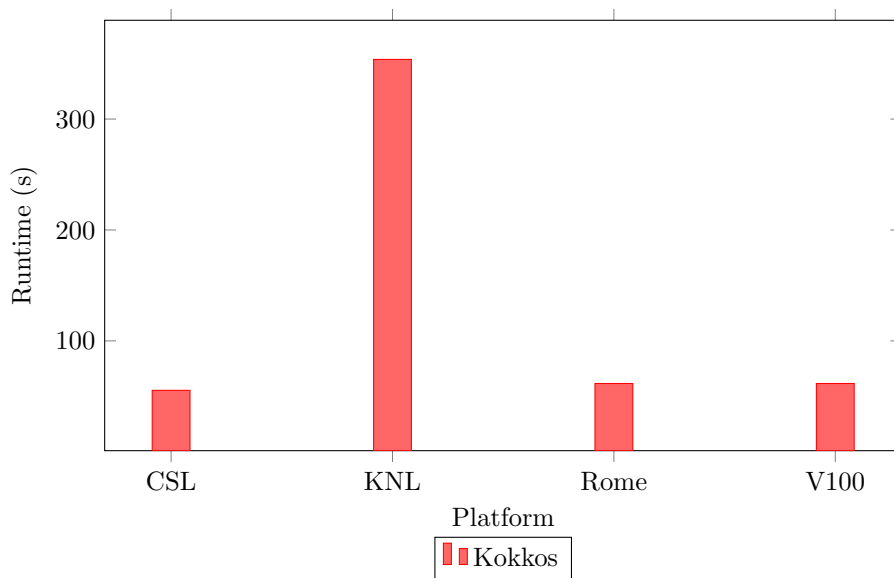


Figure 13: CabanaPIC data

Approximately equivalent performance can be seen on the CascadeLake, Rome and V100 systems. Similar to our previous Kokkos results on KNL, the runtime

<sup>11</sup><https://github.com/ECP-copa/CabanaPIC>

is significantly worse than expected, possibly indicating a Kokkos bug, or a configuration issue.

## 2.5 VPIC

Vector Particle-in-Cell (VPIC) is a general purpose PIC code for modelling kinetic plasmas in one, two or three dimensions, developed at Los Alamos National Laboratory [16]. VPIC is parallelised on-core using vector intrinsics and on-node through a choice of pthreads or OpenMP. It can additionally be executed across a cluster using MPI<sup>12</sup>.

Recently, a VPIC 2.0 [17] has been developed that adds support for heterogeneity by using Kokkos to optimise the data layout and allow execution on accelerator devices.

### 2.5.1 Performance

Figure 14 shows the runtime for the three variants of the VPIC code running on seven platforms<sup>13</sup>. This data is taken from the VPIC 2.0 study, comparing the non-vectorised, vectorised and Kokkos variants of the VPIC code. In each case, the runtime is the time taken for 500 time steps, with 66 millions particles.

### 2.5.2 Performance Portability

In terms of the performance portability of VPIC, we can see that the original and vectorised variants are only viable on the CPU architectures. Figures 15 and 16 visualise how the performance portability varies as more platforms are evaluated.

The highest performance in most cases comes from the vectorised variant of VPIC, as it achieves the best performance on all CPU platforms (except the ThunderX2, where no data is provided). However, Figure 15, when evaluating

---

<sup>12</sup><https://github.com/lanl/vpic>

<sup>13</sup>[https://globalcomputing.group/assets/pdf/sc19/SC19\\_flier\\_VPIC.pptx.pdf](https://globalcomputing.group/assets/pdf/sc19/SC19_flier_VPIC.pptx.pdf)



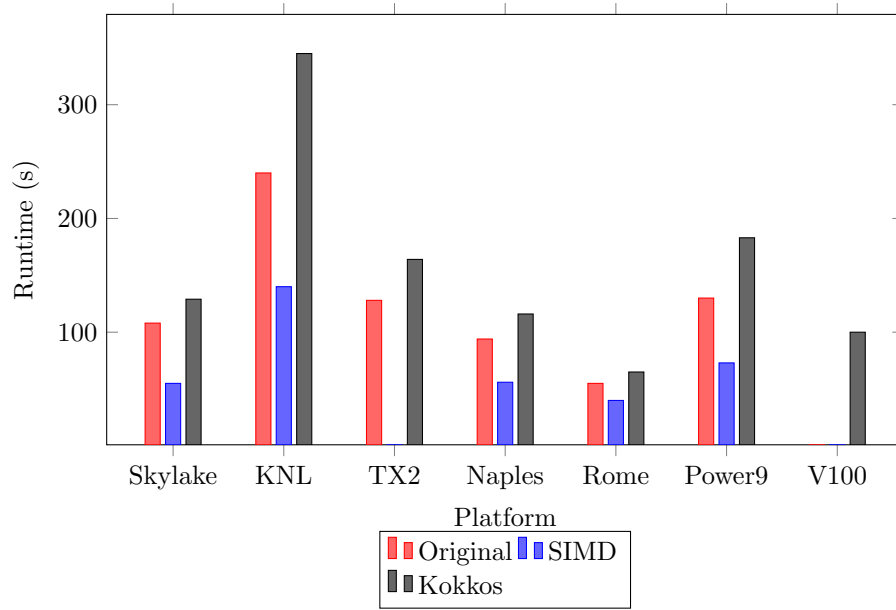


Figure 14: VPIC runtime data from Bird et al. [17]

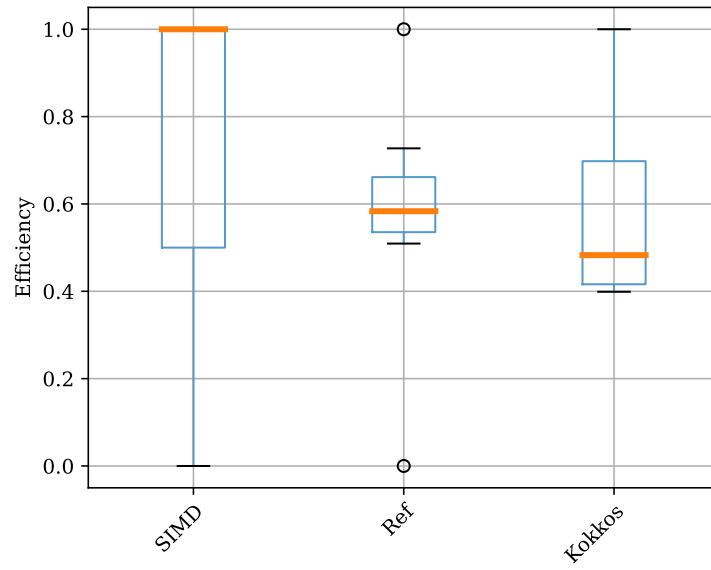


Figure 15: Box plot visualisation of performance portability of VPIC

the entire set of platforms, its performance portability would be 0, due to non-

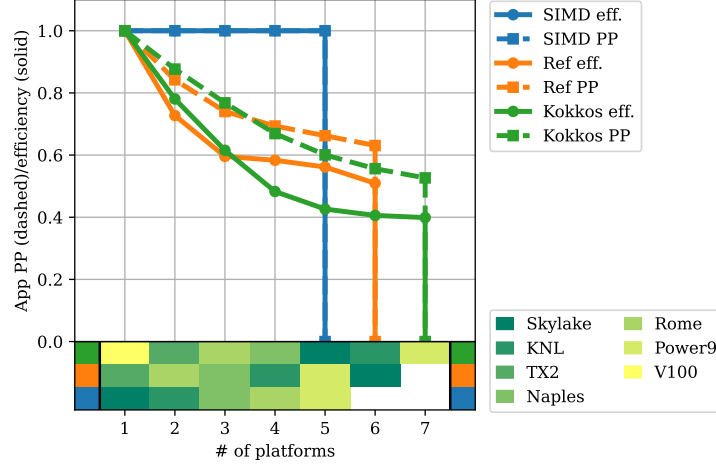


Figure 16: Cascade visualisation of performance portability of VPIC

execution on the V100 platform.

Figure 16 shows that while Kokkos performs worse than the vectorised implementation, its performance is similar the non-vectorised variant, but is also capable of execution on the V100 platform.

It should be noted that this data is from a study based on the initial implementation of VPIC using Kokkos. It is likely that these performance figures will be improved in future, potentially closing the performance gap on the vectorised implementation, while maintaining portability to heterogeneous architectures.

## 2.6 EMPIRE-PIC

EMPIRE-PIC is the particle-in-cell solver central the the ElectroMagnetic Plasma In Realistic Environments (EMPIRE) project [18]. It solves Maxwell’s equations on an unstructured grid using a finite-element method, and implements the Boris push for particle movement. EMPIRE-PIC makes extensive use of the Trilinos library, and subsequently uses Kokkos as its parallel programming model [19, 20].

### 2.6.1 Performance

The EMPIRE-PIC application is export controlled, and thus the results in this section come from the study by Bettencourt et al. [19], looking specifically at the particle kernels within EMPIRE-PIC.

Figure 17 shows the runtime of the Accelerate, Weight Fields, Move and Sort kernels within EMPIRE-PIC for an electromagnetic problem with 16 million particles (8 million H+, 8 million e-). The geometry for this problem is the tet mesh that can be seen in Figure 7 in Bettencourt et al. [19].

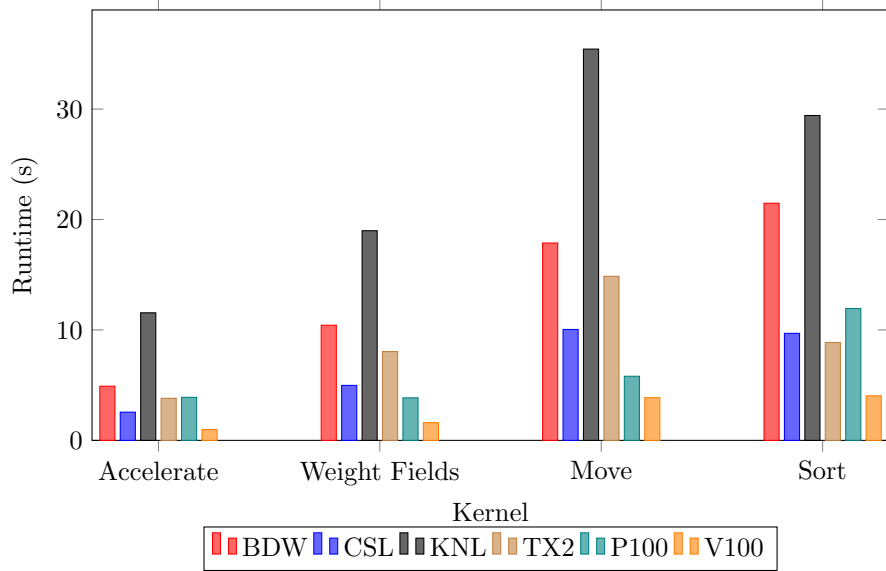


Figure 17: EMPIRE-PIC runtime data

### 2.6.2 Performance Portability

While there is only a single programming model implementation of EMPIRE-PIC, we can use the equations given in Table 2 of Bettencourt et al. [19] to calculate the FLOP/s achieved and compare this to each machines maximum performance, thus calculating the architectural efficiency. The equations presented assume the best case performance, whereby particles are evenly distributed across the domain, there is no particle migration throughout the simulation, and they are sorted at the start of the simulation. Nevertheless, they

provide a useful opportunity to analyse the performance portability of Kokkos for particle-based kernels.

Figures 18 and 19 provide visualisations of EMPIRE-PIC’s performance portability across six platforms.

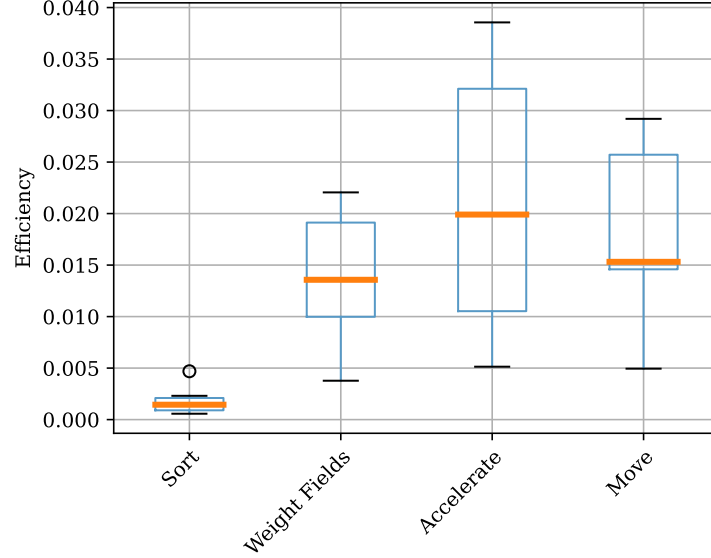


Figure 18: Box plot visualisation of performance portability for four particle kernels in EMPIRE-PIC

It is important to note that although Figure 18 shows incredibly low efficiency, this is compared to each platform’s peak performance, where a vectorised fused-multiply-add instruction must be executed each clock cycle. Achieving less than 10% of this peak performance is not unusual for a real application. In the case of the Sort kernel, the efficiency is lower still, as this is not a kernel that is bound by floating point performance.

What is clear from the performance portability visualisations is that the variance in achieved efficiency between platforms is not large, indicating that Kokkos is able to achieve a similar portion of the available performance for EMPIRE-PIC’s particle kernels. Achieved efficiency is higher on the ThunderX2 and Broadwell systems, due to less reliance on well vectorised code, and a lower available peak performance.

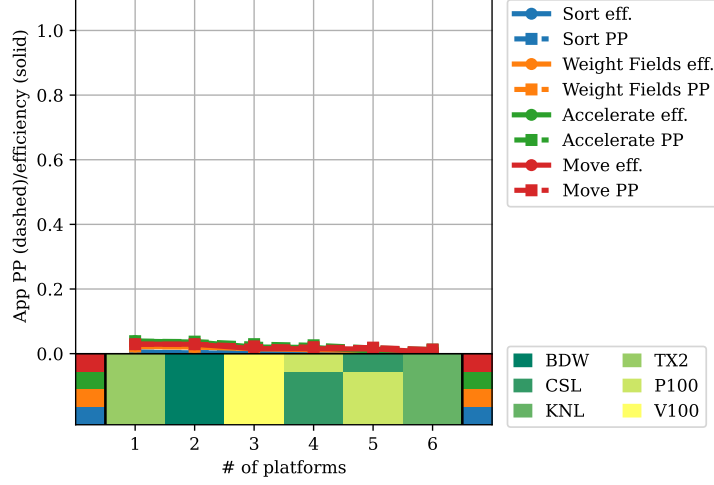


Figure 19: Cascade visualisation of performance portability for four particle kernels in EMPIRE-PIC

This data perhaps suggests that EMPIRE-PIC is not able to fully exploit the available on-core parallelism available through vectorisation. We can further support this analysis using Roofline models for the Cascade Lake, KNL, ThunderX2 and V100 systems. Figure 20 shows roofline models for these four platforms, with the four particle kernels plotted according to their arithmetic intensity and achieved FLOPs/s.

In all cases, we can see that the application is not successfully using vectorisation (and this is confirmed by compiler reports). As stated in Bettencourt et al. [19], the control flow required to handle particles crossing element boundaries leads to warp divergence on GPUs and makes achieving vectorisation difficult on CPUs. Nonetheless, on the Cascade Lake and ThunderX2 platforms, we are within an order of magnitude of the non-vectorised peak performance for the three main kernels, and we can observe that all four kernels are memory bound. For the two many-core architectures (KNL and V100), floating-point performance is further from the peak, but all kernels are likewise memory bound by available DRAM/HBM bandwidth.

Figure 20 demonstrates how vital efficient memory accesses are for achieving high performance in PIC codes, due to the relatively low arithmetic intensity of the kernels when compared to the amount of bytes that need to be moved

to and from main memory. An alternative approach to the FEM-PIC method has been explored using EMPIRE-PIC by Brown et al. [20], whereby complex particle shapes are supported using virtual particles based on quadrature rules. Using virtual particles in this manner can increase the arithmetic intensity of particle kernels without requiring significantly more data to be moved from and to main memory.

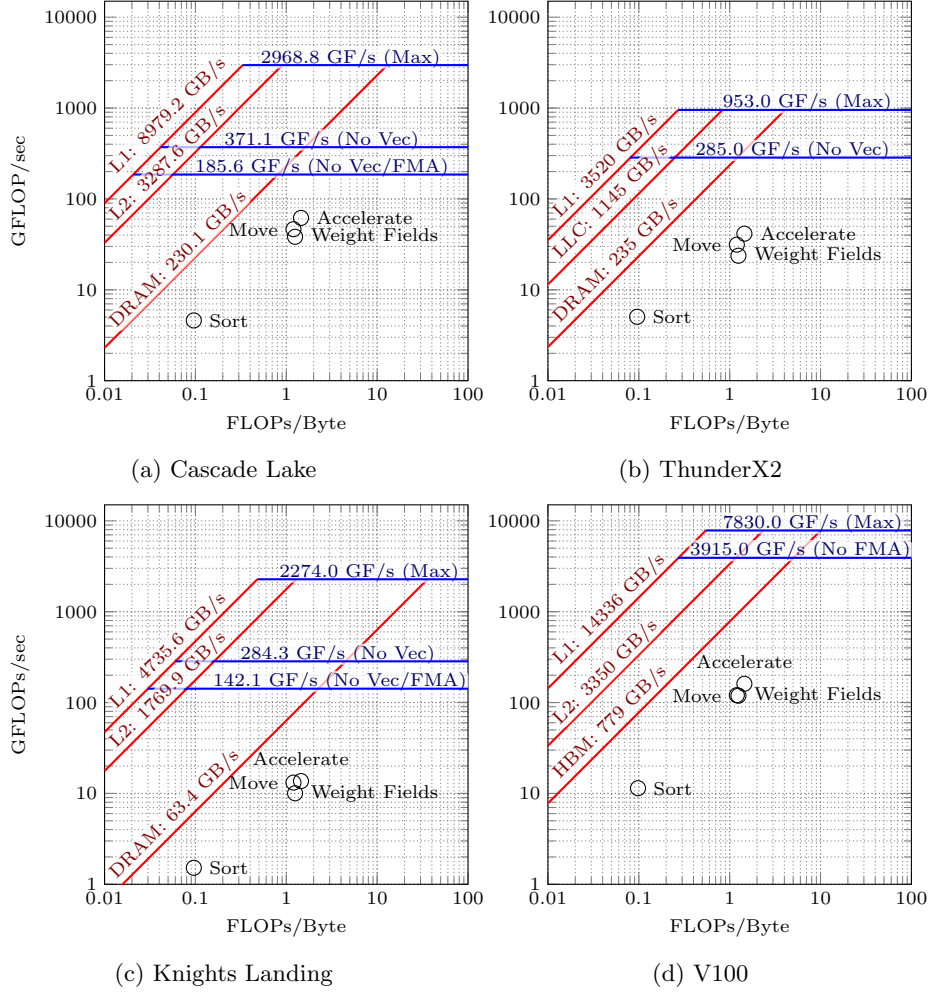


Figure 20: Roofline plots on four platforms, gathered using the Empirical Roofline Toolkit [21]

### 3 Conclusions

This report serves as a living document of the performance of applications that implement algorithms of interest to the NEPTUNE project. For each of the applications in this report, there are typically a number of alternative implementations, solving the same algorithm but using a different parallel programming model. This allows us an opportunity to assess these programming models and their appropriateness for the NEPTUNE project, with the goal of creating a set of best practices to developing plasma physics applications that are both *performant* and *portable*.

The results presented in the previous section show that in most cases, OpenMP and/or MPI provide the best performance on CPU platforms, while CUDA typically provides the best performance on NVIDIA GPUs. However, these programming models significantly affect the portability of these applications, with the former unable to use accelerators, and the latter unable to use host platforms. Developing an application that can exploit all available parallelism that is likely to be present on post-Exascale systems would therefore require developers to maintain multiple implementations of a code – potentially one for each class/generation of host or accelerator platforms.

For fluid codes, there are a number of domain specific languages (DSLs) that provide abstractions for grid-based algorithms. OPS is one such DSL targeted at structured mesh applications, and capable of code generation targeting MPI, OpenMP, OpenACC, CUDA and HIP. Our study with TeaLeaf shows that it is able to provide performance that in many cases is on par with native OpenMP and MPI, and within  $2\times$  native CUDA performance on a P100. However, such DSLs often reduce the flexibility afforded to a developer.

Besides code generation from a higher-level abstraction, GPUs can be targeted using pragma-based language extensions such as OpenMP 4.5 and OpenACC. Both offer similar functionality, but only OpenMP 4.5 allows portability between accelerator and non-accelerator platforms. However, our evaluation has shown that although OpenMP 4.5 allows us to target GPUs, different pragmas are often required to achieve sufficient performance on accelerators when compared to host systems, meaning that multiple implementations would likely need to be maintained. This is well demonstrated by our miniFE results, where the

OpenMP with offload code does successfully execute on the CPU architectures but offers significantly worse performance than OpenMP itself.

The template libraries, Kokkos and RAJA are both capable of providing full portability across all architectures, and in most cases offer good performance. The significant exception from our results is for the Intel Knights Landing platform, where Kokkos performance is typically poor. This performance gap is likely the result of a bug or memory configuration issue, but will not be investigated further due to the discontinuation of the KNL architecture. Regardless, where we are able to compare Kokkos or RAJA to a native programming model, they are typically able to achieve a runtime that is no more than 20% greater than the native programming model on CPUs and no more than 50% greater than the native programming model on GPUs, but from a single code base.

Another approach that is gaining traction is that of SYCL/DPC++. In our current benchmark set, only a single application is available implemented in SYCL (miniFE), and that implementation has been generated using Intel’s DPC++ Compatibility Toolkit. The resulting application is portable across platforms but in most cases has performance that is only slightly better than the available OpenMP 4.5 implementation. This warrants additional exploration to account for this performance difference; for such an immature programming model, it is likely that choice of compiler, and some very simple optimisations will bring performance more inline with other approaches to portability. As this project progresses, hopefully more applications will be available for evaluation, and compiler support will evolve.

For the particle methods tranche of applications, they are predominantly available using Kokkos as a parallel programming model. This does allow portable execution across all available platforms, but makes it difficult to compare performance against native implementations. In the case of VPIC, we can see that Kokkos provides performance that is inline with the original, unvectorised implementation on all platforms, and allows us to extend our platform set to include GPU devices. However, the greatest performance comes from using non-portable vector intrinsics, which in this case means maintaining an implementation for each set of vector instructions (i.e. SSE, AVX, AVX-2, AltiVec, etc.).



### 3.1 Limitations

The work presented in this report represents our initial evaluation of approaches to performance portability. We intend that this document is continually updated as new data becomes available, and as applications and implementations are developed. Currently, the data in this report contains a few limitations that we aim to rectify in future.

Firstly, due to its immaturity relative to other approaches, there are a lack of relevant fluid and particle-in-cell applications available that use the SYCL/DPC++ programming model. This means, that with the exception of miniFE, it is difficult to assess its appropriateness as an approach to performance portable application development. A recent study has by Regulý et al. shows that for a computational fluid dynamic application SYCL may be able to achieve comparable performance, though this may require different code paths for different hardware [22].

Secondly, the PIC codes assessed in this report all use the Kokkos programming model. Again, this limits our ability to reason about the appropriateness of this approach for PIC codes, but we can use the VPIC data to show that while we cannot match native, hand-vectorised performance, it can provide performance that is similar to the original implementation, and can be extended to heterogeneous architectures.

Finally, we have not currently evaluated performance on any AMD Radeon Instinct or Intel Xe hardware, due to availability of test platforms. We aim to add these platforms in the near future, when available, either through the COSMA8 system at Durham University, or through Amazon EC2 instances.

## References

- [1] S.J. Pennycook, J.D. Sewall, and V.W. Lee. Implications of a metric for performance portability. *Future Generation Computer Systems*, 92:947 – 958, 2019.
- [2] Jason Sewall, S. John Pennycook, Douglas Jacobsen, Tom Deakin, and Simon McIntosh-Smith. Interpreting and visualizing performance portabil-

- ity metrics. In *2020 IEEE/ACM International Workshop on Performance, Portability and Productivity in HPC (P3HPC)*, pages 14–24, 2020.
- [3] B D Dudson, M V Umansky, X Q Xu, P B Snyder, and H R Wilson. BOUT++: A framework for parallel plasma fluid simulations. *Computer Physics Communications*, 180:1467–1480, 2009.
  - [4] C.D. Cantwell, D. Moxey, A. Comerford, A. Bolis, G. Rocco, G. Mengaldo, D. De Grazia, S. Yakovlev, J.-E. Lombard, D. Ekelschot, B. Jordi, H. Xu, Y. Mohamied, C. Eskilsson, B. Nelson, P. Vos, C. Biotto, R.M. Kirby, and S.J. Sherwin. Nektar++: An open-source spectral/hp element framework. *Computer Physics Communications*, 192:205–219, 2015.
  - [5] T D Arber, K Bennett, C S Brady, A Lawrence-Douglas, M G Ramsay, N J Sircombe, P Gillies, R G Evans, H Schmitz, A R Bell, and C P Ridgers. Contemporary particle-in-cell approach to laser-plasma modelling. *Plasma Physics and Controlled Fusion*, 57(11):113001, sep 2015.
  - [6] Tom Deakin, Simon McIntosh-Smith, James Price, Andrei Poenaru, Patrick Atkinson, Codrin Popa, and Justin Salmon. Performance portability across diverse computer architectures. In *2019 IEEE/ACM International Workshop on Performance, Portability and Productivity in HPC (P3HPC)*, pages 1–13, 2019.
  - [7] R. O. Kirk, G. R. Mudalige, I. Z. Regulý, S. A. Wright, M. J. Martineau, and S. A. Jarvis. Achieving Performance Portability for a Heat Conduction Solver Mini-Application on Modern Multi-core Systems. In *2017 IEEE International Conference on Cluster Computing (CLUSTER)*, pages 834–841, Sep. 2017.
  - [8] Matthew Martineau, Simon McIntosh-Smith, and Wayne Gaudin. Assessing the performance portability of modern parallel programming models using tealeaf. *Concurrency and Computation: Practice and Experience*, 29(15):e4117, 2017.
  - [9] Simon McIntosh-Smith, Matthew Martineau, Tom Deakin, Grzegorz Pawelczak, Wayne Gaudin, Paul Garrett, Wei Liu, Richard Smedley-Stevenson, and David Beckingsale. TeaLeaf: A Mini-Application to Enable Design-Space Explorations for Iterative Sparse Linear Solvers. In *2017 IEEE International Conference on Cluster Computing (CLUSTER)*, pages 842–849, 2017.

- [10] Richard Frederick Barrett, Li Tang, and Sharon X. Hu. Performance and Energy Implications for Heterogeneous Computing Systems: A MiniFE Case Study. 12 2014.
- [11] Alan B. Williams. Cuda/GPU version of miniFE mini-application. 2 2012.
- [12] Meng Wu, Can Yang, Taoran Xiang, and Daning Cheng. The research and optimization of parallel finite element algorithm based on minife. *CoRR*, abs/1505.08023, 2015.
- [13] David F. Richards, Yuri Alexeev, Xavier Andrade, Ramesh Balakrishnan, Hal Finkel, Graham Fletcher, Cameron Ibrahim, Wei Jiang, Christoph Junghans, Jeremy Logan, Amanda Lund, Danylo Lykov, Robert Pavel, Vinay Ramakrishnaiah, et al. FY20 Proxy App Suite Release. Technical Report LLNL-TR-815174, Exascale Computing Project, September 2020.
- [14] J. C. Camier. Laghos summary for CTS2 benchmark. Technical Report LLNL-TR-770220, Lawrence Livermore National Laboratory, March 2019.
- [15] Robert Anderson, Julian Andrej, Andrew Barker, Jamie Bramwell, Jean-Sylvain Camier, Jakub Cervený, Veselin Dobrev, Yohann Dudouit, Aaron Fisher, Tzanio Kolev, Will Pazner, Mark Stowell, Vladimir Tomov, Ido Akkerman, Johann Dahm, David Medina, and Stefano Zampini. Mfem: A modular finite element methods library. *Computers & Mathematics with Applications*, 81:42–74, 2021. Development and Application of Open-source Software for Problems with Numerical PDEs.
- [16] K. J. Bowers, B. J. Albright, B. Bergen, L. Yin, K. J. Barker, and D. J. Kerbyson. 0.374 Pflop/s Trillion-Particle Kinetic Modeling of Laser Plasma Interaction on Roadrunner. In *Proceedings of the 2008 ACM/IEEE Conference on Supercomputing*, SC '08. IEEE Press, 2008.
- [17] Robert Bird, Nigel Tan, Scott V Luedtke, Stephen Harrell, Michela Taufer, and Brian Albright. VPIC 2.0: Next Generation Particle-in-Cell Simulations. *IEEE Transactions on Parallel and Distributed Systems*, pages 1–1, 2021.
- [18] Matthew T. Bettencourt and Sidney Shields. EMPIRE Sandia’s Next Generation Plasma Tool. Technical Report SAND2019-3233PE, Sandia National Laboratories, March 2019.

- [19] Matthew T. Bettencourt, Dominic A. S. Brown, Keith L. Cartwright, Eric C. Cyr, Christian A. Glusa, Paul T. Lin, Stan G. Moore, Duncan A. O. McGregor, Roger P. Pawlowski, Edward G. Phillips, Nathan V. Roberts, Steven A. Wright, Satheesh Maheswaran, John P. Jones, and Stephen A. Jarvis. EMPIRE-PIC: A Performance Portable Unstructured Particle-in-Cell Code. *Communications in Computational Physics*, x(x):1–37, March 2021.
- [20] Dominic A.S. Brown, Matthew T. Bettencourt, Steven A. Wright, Satheesh Maheswaran, John P. Jones, and Stephen A. Jarvis. Higher-order particle representation for particle-in-cell simulations. *Journal of Computational Physics*, 435:110255, 2021.
- [21] Yu Jung Lo, Samuel Williams, Brian Van Straalen, Terry J. Ligocki, Matthew J. Cordery, Nicholas J. Wright, Mary W. Hall, and Leonid Oliker. Roofline Model Toolkit: A Practical Tool for Architectural and Program Analysis. In Stephen A. Jarvis, Steven A. Wright, and Simon D. Hammond, editors, *High Performance Computing Systems. Performance Modeling, Benchmarking, and Simulation*, pages 129–148. Springer International Publishing, 2015.
- [22] Istvan Z. Reguly, Andrew M. B. Owenson, Archie Powell, Stephen A. Jarvis, and Gihan R. Mudalige. Under the Hood of SYCL – An Initial Performance Analysis with An Unstructured-Mesh CFD Application. In Bradford L. Chamberlain, Ana-Lucia Varbanescu, Hatem Ltaief, and Piotr Luszczek, editors, *High Performance Computing*, pages 391–410. Springer International Publishing, 2021.