

# Performance Analysis and Assessment for Computational Research Software—A Workbook

Tobias Weinzierl and Thomas Flynn

January 6, 2026



# Preface

## April 2025

We organised a series of performance analysis workshops in spring 2025, funded through the Digital Research Infrastructure (DRI) programme under the umbrella of our HAI-End project. These workshops differed from previous installments in that we did not focus on particular tools. Instead, we conducted brainstorming and discussion exercises centered around a fundamental question: How do you begin performance assessment when you know little to nothing about your code?

While some tools such as Intel's Application Performance Snapshots (APS) or Linaro's MAP provide high-level overviews, we approached this challenge not from a tools perspective but from an observational point of view. We asked: What are the essential properties we must understand before diving deeper into code analysis?

This inquiry led to a high-level, top-down approach, documented through mindmaps, flowcharts, and case studies. The raising approach is fundamentally different from the existing POP methodology, which collects fine-grained metrics and successively aggregates them into high-level metrics and flaws. It also differs from approaches offered by and supported through tools like Scalasca, which require a certain level of code maturity and system knowledge from the user.

Both direction of travel are important and very valuable, and we think they might be more sophisticated and insightful compared to what we came up with. However, our technique is very simple and does require next to no specialised software and tools knowledge. We developed a straightforward, accessible starting point to initiating performance analysis. This first version of the document summarises these foundational ideas.

Some of the sections are not yet written, and a lot more are yet to come. However, we decided to share them early, to be able to gather feedback and polish the ideas. This is only the first version of a living document.

Take home: - Target group has to be clearer: in-depth or application-domain focused  
- Tool information density plus training material is not appropriate for latter group -  
Future workshops have to cover both domains separately

- Application people want to do performance engineering but don't know how to do (and that they have to do) assessment first - They pick the first tool and then stick to it and hence often miss out on big picture

- Success model: pair up tool expert with domain expert - Tool training is there for

assessor, but not necessarily for the domain expert who wants to do performance engineering - This contradicts my previous assumption that tools are there predominantly for performance engineer

- Maybe the assessments have to come first - Success stories are key (like the ones POP has delivered)

- A lot of tools are too difficult to use; counterpart is MAQAO or APS or Perf Snapshot

- Pitch assessment gurus directly for tool workshops - Create for others a tutorial "how to assess your code's performance"

Thomas Flynn, Tobias Weinzierl  
Durham, January 6, 2026

# Contents

|            |  |           |
|------------|--|-----------|
| <b>I</b>   | <b>Introduction</b>                            | <b>7</b>  |
| <b>1</b>   | <b>How to read</b>                             | <b>9</b>  |
| <b>2</b>   | <b>Terminology</b>                             | <b>13</b> |
| 2.1        | What we do . . . . .                           | 13        |
| 2.2        | Code development workflows . . . . .           | 14        |
| 2.3        | Fundamental HPC and assessment terms . . . . . | 16        |
| 2.4        | Bottom-up vs. top-down assessment . . . . .    | 19        |
| <b>3</b>   | <b>Preparation</b>                             | <b>21</b> |
| 3.1        | Set up the benchmarks . . . . .                | 21        |
| 3.2        | Working environment . . . . .                  | 23        |
| 3.3        | Compiler setup . . . . .                       | 23        |
| 3.4        | Understand the code's complexity . . . . .     | 25        |
| 3.5        | I/O . . . . .                                  | 25        |
| 3.6        | Machine details . . . . .                      | 26        |
| 3.7        | Labelling of code parts . . . . .              | 27        |
| 3.8        | Existing optimisations . . . . .               | 28        |
| <b>II</b>  | <b>Performance Assessment</b>                  | <b>31</b> |
| <b>4</b>   | <b>High-level first glance</b>                 | <b>33</b> |
| 4.1        | The assessment rubrics . . . . .               | 34        |
| 4.2        | Core performance . . . . .                     | 35        |
| 4.3        | Intra-node (node-level) performance . . . . .  | 38        |
| 4.4        | Inter-node performance . . . . .               | 43        |
| 4.5        | GPU performance . . . . .                      | 46        |
| 4.6        | I/O performance . . . . .                      | 47        |
| 4.7        | Localisation . . . . .                         | 48        |
| <b>III</b> | <b>Case Studies</b>                            | <b>51</b> |
| <b>5</b>   | <b>High-level assessments</b>                  | <b>53</b> |

|           |   |           |
|-----------|---|-----------|
| <b>IV</b> | <b>Discussion</b>                                     | <b>55</b> |
| <b>6</b>  | <b>Summary and outlook</b>                            | <b>57</b> |
| 6.1       | Recap . . . . .                                       | 57        |
| 6.2       | Outlook . . . . .                                     | 58        |
| <b>V</b>  | <b>Appendix</b>                                       | <b>61</b> |
| <b>A</b>  | <b>Acknowledgements</b>                               | <b>63</b> |
| A.1       | Funding councils and supporting initiatives . . . . . | 63        |
| A.2       | People . . . . .                                      | 63        |
| A.3       | Partners . . . . .                                    | 64        |
| <b>B</b>  | <b>Submission checklist</b>                           | <b>65</b> |
| B.1       | Code submission . . . . .                             | 65        |
| B.2       | Code of conduct . . . . .                             | 68        |
| B.3       | Outcomes . . . . .                                    | 68        |

## **Part I**

# **Introduction**





# Chapter 1

## How to read

As with any workbook — i.e. a book that’s there to guide actual work not to teach a big subject — there are numerous ways to read the text. You can read the manuscript chapter by chapter, but this is unlikely to be the most fun experience. Below, we collect some ideas on how to substructure the reading:

### Context

This document is part of an ongoing body work of designing and implementing a structured performance analysis methodology with the future aim of building a performance assessment service.

The work presented here is a precursor to a far more detailed text which will include a structured methodology for conducting performance assessments. In this iteration we introduce: relevant terminology for High-Performance Computing (HPC) and performance analysis; preparing for an assessment; and finally how to begin an assessment from a high-level.

### Current state of manuscript

Both Chapter 2 and 3 are predominantly complete though with notable exceptions such as Section 3.4 which are yet to be completed, and are labelled as ‘t.b.d’.

Within Chapter 4, the first three sections are largely finished and have been written in conjunction with implementing these methodologies during recent performance analysis workshops. Section 4.4 has been started but the ending of this section still needs significant work. We want to implement more of this ‘inter-node’ assessment before making recommendations that are unfounded. Sections 4.5 and 4.6 are yet to be worked through, though Section 4.7 is mostly complete.

Clearly sections labelled as ‘complete’ may also need to be extended as further work is done, as the content may become insufficient. Note however that some forward references may link to topics to come in future versions of this manuscript, e.g. discussions of tools, performance assessment service, etc.

Finally, Chapter 6 will not be part of the final document (as with this section and the previous) but recaps the current progress of this document whilst also discussing the content to come in the next iteration of this document.

## First-time assessment

We recommend to run through the preparatory remarks in Chapter 3 first and to complete all the checkboxes (dos) there. Once completed, it is time to run through the high-level assessment remarks in Chapter 4. With most codes, it is highly unlikely that a code will only show one kind of flaw in the high-level analysis.

We propose to run through the detailed analysis step by step afterwards. Per step, it makes sense to consult the tools summary (Chapter ??) and notably to try out different things. Where appropriate, the reader might prefer to read through one of the case studies in Part ??, just to get some inspiration on what tools and approaches to try out and how to report outcomes.

## Assessment

If you have already completed numerous assessments, then the manuscript is mainly a reference document. You'd start with the high-level assessment in Chapter 4, before you dig into the various subchapters depending on the initial assessment outcome. The case studies are not of particular relevance (anymore), as you have already completed your own case studies.

## Feedback

The quintessential objective of any assessment is to give feedback to developers. Once they make changes to their source code, it is important to re-start the assessment from scratch. There is always this temptation to return to the particular characteristics of interest that have led to a report, to see how the code changes by the development team have altered these characteristics, and then to continue from there.

It is important to reassess the “flaw” metrics. However, this should only be a validation that something has changed. After that, we have to start from scratch, as performance optimisation is a complicated process where one change might introduce flaws, improvement or character changes in a completely different part of the code or affect a completely different flavour of the performance. If we alter a particular code part that makes the scaling deteriorate for example, this might have a severe knock on effect on the core performance. If we improve the core performance, the GPU offloading might all of a sudden run into data transfer issues. There are numerous examples for these side-effects.

## How it works

- ⊞ Conduct a performance assessment which both covers all aspects of the code (cmp. Chapter 4) plus goes in-depth all the way through for all flaws identified.
- ⊞ Hand this completed assessment back to the developers to address the flaws.
- ⊞ Restart the assessment from scratch top-down, i.e. starting with the high-level assessment covering all code runtime nuances.

**How it does not work**

- ⊞ Mix code optimisation and performance assessment workflows.
- ⊞ Restart the assessment from the place where you stopped when you reported a flaw and continue digging through the assessment decision metrics without redoing a high-level, overarching assessment.



## Chapter 2

# Terminology

Performance and efficiency, i.e. the speed of a code, are critical non-functional properties of scientific codes. Measuring, assessing and interpreting them is challenging; without even considering tuning. Most annoyingly, there is no standardised way how to analyse and assess performance. There is no vanilla workflow. Every team, every project, every scientist create their own analysis workflow or process, and this processes often changes from paper to paper or experiment to experiment.

Despite the vagueness in the metrics and the processes, there are some well-defined, common terms in the community that we should use. There are also some key concepts and some jargon that we use throughout the write-up. So we better introduce them first.

### 2.1 What we do

From hereon, we rely on a few key assumptions:

1. We assume that the person studying the code has *limited knowledge* of the code's internals, the application domain, and the algorithms used.
2. We evaluate *how* a code performs and whether certain characteristics deserve closer examination.
3. We are interested in the code's behavior on *one particular machine* for *one particular experimental setup* or a closely defined set of experiments.
4. Whenever we explain *why* a code performs as it does, this explanation comes from a machine's point of view: It does not argue about the algorithm or science case.

**Doctrine 1** Performance Assessment *documents code behaviour from an outsider's point of view. It captures what we see (benchmarking).* Performance Analysis *explains why a code behaves the way it does.* Performance Engineering *improves the performance.*

The three steps typically build up on each other. Performance assessment is a little bit like a health screening in medicine. We run a couple of standard benchmarks (similar to blood pressure, oxygen saturation and pulse for example). If some metrics are not in the region where we want them to be, we issue further tests, i.e. we switch to the analysis part.

The skills we need first and foremost are to know what tests we should conduct or commission, and we have to be able to read the outcomes. If the test metrics are not what they should be, we eventually transfer the patient (code) to a specialist and run further tests. Based upon these tests, we finally start the treatment.

#### Reader's guide

This text starts with a performance assessment before it switches to some performance analysis. It does not cover aspects of performance engineering.

Our working assumption is that performance engineering requires inevitably domain or code expertise which we cannot/might not have as generalists.

**Assessment crime 1** *Developers tend to skip the assessment part and sometimes the analysis.*

They want to get their hands dirty and improve the code. However, this is deeply unscientific and unreasonable without a proper assessment and triggers lots of follow-up complications. Systematic work on performance always follows the sequence assessment–analysis–improvement (engineering).

#### How it works

- ▣ Ignore everything you know nothing about the code and focus on the data you measure.
- ▣ Stop any development while you assess and analyse.

#### How it does not work

- ▣ Start with the attitude “I want to fix a flaw”.
- ▣ Skip assessment steps, as you already know what’s going on.

## 2.2 Code development workflows

#### Reader's guide

Feel free to skip the remainder of this chapter if you are a first-time reader.

Systematic scientific computing typically implements an iterative approach: First, we measure or benchmark how a code performs and use these data to inform our performance assessment, i.e. our reporting. Second, we conduct performance analysis to

explain why the code exhibits the observed behaviour. Finally, this analysis feeds into program tuning or optimisation. After that, we repeat the workflow.

If we take our list of activities and integrate it within the definition of the performance analysis workflow, it becomes clear that the tuning step is out-of-scope here. Another related technique is out-of-scope, too: The most advanced papers in scientific computing, particularly within high-performance computing, deliver a performance model; essentially a pen-and-paper assessment of what a code should be able to deliver. This is impossible if we adhere to the principles above: having limited knowledge and only explaining observed behavior rather than relating it to external factors (such as input data or algorithmic choices) or algorithmic context. We commit to a data-driven approach. Obviously, such an approach might complement an analytic, theoretical assessment of code, but we recommend considering this as an entirely separate cup of tea.

We will inevitably face difficulties if we truly do not know the code we are evaluating: Explaining certain effects is so much easier with some (rudimentary) understanding of the code. However, there are three compelling reasons to treat code as a (logical) black box: First, performance assessment is often delegated to non-developers: colleagues from computing centers, centralised HPC specialists, or external service providers. If we want to discuss performance analysis methodology, it is counterproductive to assume in-depth code knowledge. Many people conducting assessments cannot be experts in every code they examine. Second, requiring code expertise within assessment teams would imply that core developers are always available. This might be an ideal in the spirit of eXtreme Programming [1], where the customer must be accessible at all times. It is unrealistic in science, where PIs are busy with various tasks, developers frequently leave projects—as they graduate, accept new positions, or are reassigned—and where code components are often written by individuals and are so complex that only a few people understand the core parts. Finally (and most importantly), knowledge about a codebase tends to bias the assessment.

Developers inevitably make assumptions about reasons for code behaviour when working on performance analysis. (Senior) Developers tend to “defend” their code, explaining any flaws as natural consequences of their domain’s challenging nature, their awareness of future requirements, or their in-depth knowledge of other setups not covered by the present configuration. Simultaneously, we have frequently observed junior developers downplaying or challenging existing code—likely to prove themselves by pointing out “how things really should work”—and consequently interpreting every piece of data through this lens. This list is not comprehensive. The whole team dynamics dimension deserves further discussion (in Section ??). For the time being, we conclude that having emotionally invested developers involved in early assessment is counterproductive due to their inherent bias.

**Work in progress/todo** The link to team dynamics is broken as this part is work in progress and not yet written.

**Assessment crime 2** *People conduct assessments as code experts. Due to this, they deliver biased assessments since they think they already know the code’s performance characteristics and its flaws without an objective assessment of its “real” behaviour.*

For a high-quality assessment, it helps to step back and (pretend to) know nothing. This way, we avoid entering performance analysis with bias. Adhere strictly to the observational perspective. Even if it is your own code, try to detach yourself and pretend you have no prior knowledge. Abstract.

#### How it works

- ▣ Pretend you know nothing about the code. Blend out your expert knowledge and start on a fresh sheet.
- ▣ Ignore what others, notably developers pretend to know or have found. This makes your assessment biased.
- ▣ Force yourself not to do any tuning before the assessment and analysis are complete (within reason).

#### How it does not work

- ▣ Start with the attitude “I know what this code does and where there are weaknesses”.
- ▣ Involve someone who’s emotionally involved in the code’s development into the reporting. They will always defend their brain child.

## 2.3 Fundamental HPC and assessment terms

We try to keep the HPC jargon in this write-up to a minimum. You might argue that this is a bad idea as we want to run scientifically sound assessments. However, we have to acknowledge that most scientific code today is written by people without a computer science background, and that most scientists do not have an extensive HPC training. It is therefore not only a matter of accessibility of the performance analysis and assessment itself to try to avoid too much jargon—if we use too much of it, we run the risk that our outputs, i.e. the assessments themselves, cannot be read and interpreted appropriately by our target audience. After all, most of the assessments are not done for ourselves, but for other scientists.

At the same time, performance analysis tools are written by computer scientists, software stacks are written by computer scientists, and computer hardware is built by computer scientists. It is hence clear that there’s a lot of jargon buried in any conversation around performance. We have to find a healthy middle ground to make precise, in-depth assessments and analyses, without immersing ourselves in technical fuzz.

**Some of the most important performance analysis terms** When we speak of *profiling*, we mean recording and presenting (accumulated) metrics. The classic example is measuring how much time we spend in a function or waiting for a message to arrive. *Tracing* is the counterpart to profiling. Here, we record performance data of individual events. We are interested in the sequence of operations occurring within the system. Profiling focuses on the effect of program execution while tracing focuses on the source



of behaviour (Figure 2.1). Many developers confuse these definitions, which we took from the VI-HPS community.

Tracing can certainly be used to generate profiles, provided the recorded performance data encompasses all quantities required for a particular metric. Conversely, the reverse problem is ill-defined: it is generally impossible to reconstruct a trace from a profile without specialized knowledge of the code’s internal workings. One might immediately ask why we don’t trace all the time. The answer is equally simple: It is too expensive. Tracing generates enormous amounts of data. Therefore, all real-world tracing relies heavily on filtering, where you determine beforehand which types of events to trace, which quantities to measure, and which program parts to monitor. Consequently, we revise our introductory statement: most traces are *not* well-suited for deriving program metrics, as they are, by definition, incomplete.

To create a profile, we distinguish between two techniques: We can *sample* program execution. Essentially, we run the code, pause it at regular intervals, and each time examine our metric of interest at that point. Over time, we gather a good impression of which functions are called most frequently, how many MFlops/s we achieve, and so forth. Some people refer to this process as monitoring. As with all statistical methods, the data quality improves the longer an application runs. However, the risk of missing important but brief effects remains.

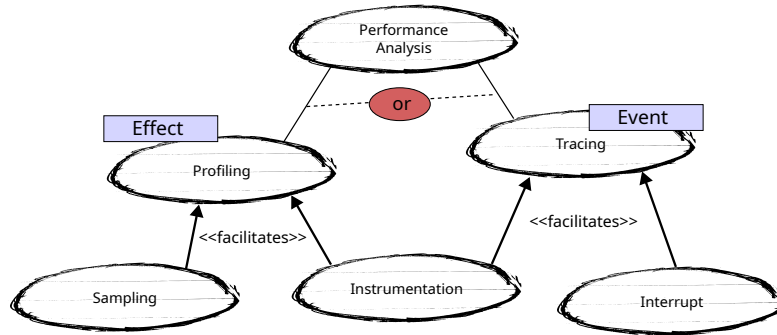


Figure 2.1: Overview of some fundamental performance analysis concepts.

Alternatively, we can *instrument* an application. Instrumentation involves inserting measurement code into our actual code which then measures quantities of interest whenever these insertion points are reached. We can insert either at compile time, or in a postprocessing step after the compile, or alter the code as we go. Instrumentation yields very accurate data (provided we instrument all relevant code sections), but it can become prohibitively expensive.

Most functions in our applications call other functions. Therefore, we must always clearly state whether any metric we present (e.g., runtime) is *inclusive*, comprising all contributions from called routines, or *exclusive*, representing only the function itself.

When tracking metrics, measuring time and call frequency is straightforward: we either measure system time or simply count occurrences. Other metrics are more complicated to obtain. Fortunately, modern systems offer performance counters. They

are special registers within the chip that automatically count certain events. For example, they can be configured to count every floating-point operation or every main memory access, allowing us to sample them or read their values upon completion. Performance counter registers are expensive to manufacture, so they are limited in number, and not every vendor tracks all possible events. Consequently, some analyses might not be possible on every system, and complex analyses might require multiple code executions, as only a few event types can be measured per run.

As we know about the existence of hardware counters, it is clear how sampling and tracing can obtain machine data beyond mere timing and call counts: They simply monitor (read out) the counters. Through this, they can make statements about the number of floating point operations completed, the number of cache misses, and so forth. For many of these concepts, we can invert the control flow: Rather than polling a counter or quantity, we can define *interrupts*, i.e. situations in which the program flow is stopped and we call our profiler or tracer. Interrupts can be realised through hardware or software. Their beauty is that the performance analysis sits somewhere outside of the core program logic—most of the time, our code runs totally unaware of the analysis in the background.

**Hardware jargon** We revise the bare minimum of hardware jargon here to get started, and then refer to further resources to revise details.

Supercomputers consist of many different individual computers linked through a high-performance interconnect (network). These computers are called *nodes*.

Each node is a full-grown computer and hence hosts many individual *cores*. All the cores on a node share their memory (but cores from different nodes do typically not share memory). Technically, many vendors call their cores *hardware threads*, as a core is a technical building block and might host multiple threads. We use hardware thread and core as synonyms.

On larger systems, you do not use the compute nodes directly. Instead, you log into a *login node*, you compile your code there, and then you hand your code over to a *scheduler*. The most popular scheduler today is Slurm. Slurm takes your request, i.e. how many nodes you want to have and for how long, as well as all the other requests from all other users and puzzles out a fair and efficient way to assign the resources. For performance analysis, it makes sense to investigate if your system of choice offers a test or benchmarking queue with a high priority. These queues usually offer only a few nodes and restrict the maximum application runtime. In return, jobs handed over to these queues are scheduled with high priority.

Compilation should be on the login node, as the actual compute nodes might not provide all the required software, and many production nodes are also shielded from the Internet, i.e. even updating your repository is a pain. Benchmarking should not be done on the login nodes, as you share them with all other users. In contrast, you can tell the scheduler that you want to have nodes *exclusive*, so no other user interferes with your measurements. This is particularly important for all data stemming from performance counters, as these counters are a hardware feature that we read out, i.e. our profiler or tracer cannot distinguish if a hardware counter's data stems from your code or another code if multiple users coexist on one machine.

**Core performance quantities** The runtime of a code is written here as  $t(p)$  where the  $p$  is either the number of nodes or number of cores. It depends on the context. The *speedup* of a program is defined as

$$S(p) = \frac{t(1)}{t(p)}, \quad (2.1)$$

which is typically smaller than or equal to one and tells us how much faster we made an application by using  $p$  compute units (cores or nodes) rather than only one.

The *parallel efficiency* is given by

$$E(p) = \frac{S(p)}{p}. \quad (2.2)$$

It quantifies to which degree adding further compute resources was worth it.

## 2.4 Bottom-up vs. top-down assessment

### Reader's guide

It might be better to skip this section and to read it after the first high-level assessment is completed.

**Work in progress/todo** This part is yet to be written and will basically point out how the POP methodology differs from what we do here.



## Chapter 3

# Preparation

### Reader's guide

This section is absolutely key before you start an assessment!

In this chapter, we discuss various steps that we might want to do *before* we dive into any benchmarking or assessment. It summarises activities that ensure that we are well-prepared to run assessments, but it also comprises steps that ensure that we understand the code sufficiently to come up with unbiased and valid statements.

There is a checklist in [Appendix B](#) which summarises some of the key information you might want to ask from a code before you start any benchmarking or analysis. This checklist is even more high level than the present discussion, but might provide an idea what minimal information an assessor requires about a code before we get started. It is kind of a customer-facing summary of this chapter's content.

### 3.1 Set up the benchmarks

Before we start a performance assessment, we need to agree on what type of measurements to perform: There has to be one well-defined benchmark. If we are given a set of setups and try to compare them, this will make our task complicated. We run the risk of comparing different characteristics produced by different code parts.

- This one benchmark has to run without any user interaction. In the normal case, the sponsor has to provide a script which builds and runs the whole benchmark “as a black-box”.
- It has to be very easy to scale that benchmark up and down, i.e. to make the problem solved slightly bigger or smaller.
- The benchmark has to be a simulation run which is representative of a real-world run.

- The benchmark does not run for too long, and it does not terminate too early. “Good” runs are often 3–5 minutes, so we get a good coverage of performance characteristics, but we do not wait for ages.

None of the properties are surprising: We need a well-defined toy setup, and we have to be able to modify its size to make statements on its scalability. In this context, it is important that we know exactly how changes of input sizes affect the compute workload: We have to know or have to be told the computational complexity of the underlying algorithms. If we have a pure Monte Carlo setup, then doing 100 times more trials should last around 100 times longer. If we have a parabolic Partial Differential Equation (PDE) solver with explicit time stepping and a regular grid, we have to know that the problem sizes increases by a factor of  $2^d$  if we half the mesh spacing. However, these algorithms also have to reduce the time step size, so if we measure over a fixed time span, we actually will have four times more time steps and therefore increase the workload by a factor of  $2^{d+2}$  once we half the mesh size. Finally, if we run an Adaptive Mesh Refinement (AMR) code for a hyperbolic setup with explicit time stepping or a Smoothed-Particle Hydrodynamics (SPH) code, we need some kind of output that tells us exactly how many particles or mesh cells we have and what time step sizes we have picked. In this case, time step sizes go down linearly with mesh spacing. At the end of the day, we are interested in the *throughput* of the system, i.e. the time we need to update one quantity of interest. It is this metric that we will start from in the next steps.

However, if domain experts tell us afterwards, i.e. after we have completed our assessment—they will always do that afterwards when they are unhappy with the outcome—that this particular run is not characteristic for larger runs, as it lacks a feature or does not really use a particular routine heavily, our whole assessment exercise is corrupted. And with it the whole assessment will be a big disappointment. It is their job to pick a setup that represents what we are interested in and is nevertheless reasonably small.

Having a “small” setup that completes in reasonable time is important since some performance analysis tools induce significant overheads, i.e. make the code significantly slower. If our vanilla run already needs more than an hour, we will never get an answer on time. Furthermore, a lot of tools gather significant amounts of data per second. Long-running simulations thus run the risk that our analysis tools fail to handle the amount of data collected.

#### Checklist

- 1 Your benchmark compiles as a black box, i.e. you can give it to another person together with a README file and there might be a few instructions that they have to copy-n-paste into their terminal, but other than that there’s no other preparatory work to be completed.
- 2 Your benchmark runs without any user interaction.
- 3 Your benchmark finishes within less than ten minutes on a single core.<sup>a</sup>

<sup>a</sup>To analyse different aspects of the code, we do not need to rely on the same benchmark. For

example, when conducting scaling analyses, there may be performance characteristics that are only revealed at high core counts. Hence, for some codes we may set up our benchmark to have runtimes exceeding 10 minutes at low core counts, such that the problem size is large enough to see high core count scaling effects. However, at high core counts we would expect to see this benchmark runtime sit within this 10 minute upper bound.

No matter how we create this setup, we have to be able to run all experiments with absolutely minimal user interaction. If we have to run through a complex pipeline for each individual experiment, we will not be able to measure anything productively.

#### How it works

- ▣ Ask the code owner (“assessment customer”) to provide you with a ready-to-use benchmark.
- ▣ Validate that this benchmark meets the criteria outlined in this section. Otherwise, return it immediately.

#### How it does not work

- ▣ Try to make a code run and design a benchmark yourself. This is the job of a domain/code expert.
- ▣ Expect insight from benchmarks that somebody without domain expertise has chosen.

## 3.2 Working environment

**Work in progress/todo** Yet has to be written, but will cover

- Access to some queues with quick turnaround times (recommended);
- A couple of standard assessment tools. We refer to them later (forward link). Reader recommendations: Study them as you go along.
- Exclusive access to nodes.
- Some higher user privileges (advanced).
- A good relationship to your sysad.

## 3.3 Compiler setup

Before we start, it is important to double-check whether we have used the correct compiler settings. It makes no sense to benchmark a code which is not producing tailored code for our particularly machinery (therefore it is important to benchmark on the right system), and it makes no sense to benchmark a code that is not using the latest compiler features.

First of all, we work with the the most aggressive optimisation that our compiler can offer. The minimum choice is to compile with `-O3`, although some other compilers might offer further granularity. Also consider using `-ffast-math` which enables aggressive floating-point optimisation.

Second, we create code that is tailored towards our particular machinery. The correct setting for this depends once again on your compiler. With LLVM and GNU, you can use flags starting with `-mXXXX` to specify machinery-specific optimisations, i.e. you can exactly say for which processor generation, instruction set, ... you want for the compiler to create binary code. If your compile node is of the same type as your test node, then simply use `-xhost` and the compiler will pick the highest instruction set available. The flag is called `-mhost` with GNU.

**Work in progress/todo** A question which was raised at RSECon25: when applying aggressive optimisation, do we include validation checks on the code? I.e., to ensure that the validity of the code still passess when potentially using approximations in the higher level optimisation flags.

#### Intel compiler

With the Intel compiler, you sometimes can pick between the generic LLVM (community) versions of a feature and one optimised by Intel itself (with the latter obviously performing better on bespoke Intel architecture). For example, `-fopenmp` enables the general OpenMP implementation whereas `-fiopenmp` gives you an in-house version of the same.

And while `-xhost` uses Intel-specific optimisations, `-mhost` will use LLVM's optimisation stack, i.e. Intel has replaced some optimisation steps with bespoke variants, which are enabled if and only if you go down the `-xhost` route. Be aware that Intel's optimisation passes typically will be disabled if you use the Intel toolchain on other vendor's hardware, i.e. you'll get unoptimised code.

Finally, we might want to add debug symbols. These instruct the compiler to insert additional information into the source code which tools can then later use to draw inference from, e.g. which source code line a certain instruction stems from. Usually, you cannot know from a single instruction in a machine code, where this comes from in the original code. With the debug systems, we inject exactly this information into the executable. Therefore, we should add `-g`. Some compilers support more detailed information (e.g. using `-g3`). It requires studying the compiler's feedback whether these flags disable optimisations, in which case we might have to balance insight vs. best-case speed.

#### Checklist

- 1 Your benchmark is compiled with debug symbols (`-g`).
- 2 Your benchmark is compiled with the highest optimisation level (`-O3` plus



a hardware-specific instruction set such as picked by `-xhost`).

#### MAQAO

If you run MAQAO over your code, the tool will report back to you what optimisations you might have missed out throughout the compilation. The top-level rubric *Global* provides you with information on the optimisation chosen, but it also makes recommendations which options might lead to better runtime.

The nice thing about MAQAO is that the tool can handle situations where you translate different translation units with different options. Also, the recommendation aspect goes beyond pure reporting and is helpful.

#### How it works

- ⊞ Use the most aggressive compiler optimisation that preserve the code semantics (aka still return valid results) and nevertheless are tailored towards your assessment machinery.
- ⊞ Check what optimisations the users' build system uses by default. Often, this is a low-hanging fruit.

#### How it does not work

- ⊞ Use a low optimisation level and machine-generic optimisations and then draw conclusions on code behaviour on a particular machine.

### 3.4 Understand the code's complexity

#### Reader's guide

The complexity discussion is mainly relevant when we conduct an MPI (weak scaling) assessment. In many cases, we might be able to skip it initially.

**Work in progress/todo** This part is missing but we need some discussion of the  $O(N^2)$  vs.  $O(N)$  behaviour that we see frequently.

### 3.5 I/O

**Work in progress/todo** This section is only partially complete. We cover the terminal outputs but miss out on the "real" I/O.

**Assessment crime 3** *The code writes a lot of information to the terminal.*

String operations are excessively expensive.

**Checklist**

- 1 Your benchmark does **not** write (a lot of) files.
- 2 Your benchmark does **not** write a lot of information to the terminal.

In this context, it is important that there's no I/O. We don't want to study complex outputs; we are busy enough with handling the output of performance analysis tools.

**Assessment crime 4** *The code reads or write to files excessively.*

More importantly, if our code writes excessive amounts of data, we almost certainly profile the I/O efficiency rather than the runtime behaviour. Unless we want to explicitly study I/O, we are well-advised to switch off any file outputs, database outputs, but also excessive file reads. A log file alone can become problematic when these files become suddenly huge or are updated frequently. In the best case, a performance assessment benchmark writes and reads barely any data at all, i.e. it largely hard-coded and silent.

**How it works**

- ▣ Disable I/O, exclude reading and writing files from any assessment as far as possible (and as long as you don't want to benchmark exactly this I/O).
- ▣ Remove terminal dumps. If the developers want to preserve the output (e.g. for validation or debugging), suggest to embed them into macros which can be compiled out of the code prior to production runs.

**How it does not work**

- ▣ Start with a high-level assessment straightaway. In particular any shared memory measurement will be corrupted by excessive terminal dumps.

## 3.6 Machine details

Modern supercomputers are very complex and it is virtually impossible to memorise all the hardware details of a testbed given the zoo of processor types, generations and flavours. While it is possible to consult your compute centre's webpages or to consult your system administrator, it might be more reasonable to obtain that data directly from the machine.

**Assessment crime 5** *Never grab the specification from the login node. Instead, ask for an interactive terminal on a compute node and get the data there. Many supercomputers' login nodes are slightly different to the compute nodes, i.e. slightly different make, more memory, more cores, . . . So you might end up with the wrong hardware specification if you analyse a login node.*

**Linux command line**

On every Linux system, information about the CPU on the system is held in a central file. You can read it out by typing `cat /proc/cpuinfo` into the terminal.

**Likwid**

If you have Likwid installed, you can invoke `likwid-topology` to obtain in-depth information about your CPU.

**MAQAO**

In MAQAO, you find some topology information in the menu rubric *Topology*.

**How it works**

- ⊞ Gather the exact spec of your testbed under realistic load. If this information is already available (e.g. due to previous assessment exercises), it makes sense to reuse it.

**How it does not work**

- ⊞ Throw advanced tools onto your code and expect that they report the theoretical machine capabilities themselves.

## 3.7 Labelling of code parts

**Reader's guide**

You might want to skip this section initially and return to it later, once you need information about what certain code parts do.

Software typically contains quite a lot of boilerplate code: predominantly administrative routines that enable calculations but don't deliver actual science flops. For several follow-on exercises, it will be very useful to know which code parts deliver science and which code parts do not directly add scientific value. This is domain knowledge. Using it contradicts, to some degree, our ambition to run a performance assessment in a black-box manner.

However, we do not ask an assessor to understand what a code does. We require the assessor to know which code parts contribute towards the science case. For this reason, it is absolutely sufficient to hand a first profile back to a domain expert and to ask them to label those routines that do the actual calculations.

As an example: In a linear algebra code, the actual matrix-vector products, scalar products, . . . all deliver science. The allocations for temporary block matrices, the exchange routines for rows of the matrices, any indexing over the sparsity pattern, and so forth do not contribute calculations. They are absolutely essential, but they do not give us the MFlops/s.

As your analysis progresses, you might run into situations where functions suddenly appear to be hot that you have not seen before. In this case, you can kind of safely assume that these functions do *not* contribute to the science directly. If they didn't show up for a small-scale, few-core run, then they likely do not feed directly into the result.

#### gprof

t.b.d.

#### VTune

Start a multithreaded run and mark the actual computational phase in the timeline at the bottom. Switch to the *Bottom up* view. *Zoom in and filter by selection* will narrow down the profile to the region of interest. Sort by *Effective Time* and change to *Show Data As . . . Percentage*. Make notes of all the entries above roughly 1% (rule of thumb) and/or track the top 10 routines.

#### How it works

- ⊞ Upon your very first profiling, ask a code expert which code parts (functions) actually perform calculations, i.e. contribute towards the science.

#### How it does not work

- ⊞ Assume that all routines deliver science (even though you may assume that all are necessary).

## 3.8 Existing optimisations

#### Reader's guide

Playing around with (user-driven) optimisations can become tricky and might require the assessor to communicate with developers (if you are not the same person). For an initial assessment, it might be very reasonable to skip this section.

It is rare that we work with a code that hasn't gone through a series of optimisation steps already. Most developers think at one point about how to make a realisation fast. However, optimisation steps can backfire:

First, many code developers have no systematic training in HPC and/or do not apply rigorous performance analysis. Consequently, their realisation does not tackle hotspots. In the worst case, the resulting optimisation steps make the code more complicated than it actually has to be.

Second, many scientific codes are old. They contain legacy code parts. If these code parts had been optimised, the optimisation steps might not be valid for today's computer architectures anymore. For example, compute time used to be expensive and performance engineers hence used to introduce precomputed data tables in the past. Today, memory bandwidth is likely the most precious resource. Having extensive lookup tables can actually impose additional stress on the memory interconnect, and a recomputation of values might be advantageous.

Finally, we have to be frank with our community: A lot of people who consider themselves to be outstanding computational scientists actually have a very limited understanding of computer architecture. That's not a problem per se, as long as people are aware of their lack of knowledge. It is problematic when we see developers introducing or pushing for "optimisations" which actually are not supported by data or insight, but rather individual "genius" and "expert knowledge".

**Assessment crime 6** *We work with a code that the developers consider to be "highly optimised". Due to these optimisations, we start from a code base that is actually overly complicated (thus hiding the "real" performance flaws) and might actual underperform due to these optimisation steps.*

Many scientific codes have not been co-designed as HPC codes. They have been written first and foremost with functional requirements ("the science") in mind. The HPC came later. This is reflected in the code design: Optimisations are realised as add-ons wrapping around the core science (yet sometimes, unfortunately, also penetrating it). Therefore, we always have to keep in mind that performance has been seen as a second-class citizen in most development projects. By the time a project has reached a certain maturity and people suffer from performance flaws, they will ask for a performance assessment (and then want to be told that they have already done the heavy lifting or that their science is so challenging and completely different to everything else that it is very hard to get good performance). They won't ask for it prior to that.

If possible, it is helpful to disable all of these historic optimisations. We can then start with the performance assessment on a green field and re-inject the optimisations one by one as we go along. In an ideal world, they should make the code faster, but this has to be validated!

#### How it works

- ⊠ Disable existing user-driven tuning options.
- ⊠ Independently validate that they actually make the code faster.

#### How it does not work

- ⊠ Assume historical performance optimisation steps to pay off on current hardware.

- ⊖ Rely on domain scientists that their optimisation has been driven by measurements (formal performance analysis) and an in-depth insight into the machine.

# **Part II**

## **Performance Assessment**





## Chapter 4

# High-level first glance

There are various ways or strategies to obtain a first high-level overview and how to get started if we want to understand a code's performance better. Every team, every project, every scientist have their own approach, and it often changes from project to project. In this chapter, we propose a workflow running along few very simple principles:

1. We study the performance characteristics for *different machine part/aspects separately*.
2. We use a simple *blackbox performance metric* per dimension to find out first if there are issues or not.
3. *We do not use any tools* (yet).
4. We are not interested in explanations of *why* we see any performance flaws. We simply document *what flaws* we observe.

Per dimension, we later break all performance flaws found down into subrubrics, and dig deeper and deeper into the code. Along this way, we construct working hypotheses of why the code performs poorly. But this comes later. For the time being, we focus in a bird's eye view, which is an analysis which we can conduct, for example, with a simple electronic spreadsheet, a few tools and mere timings.

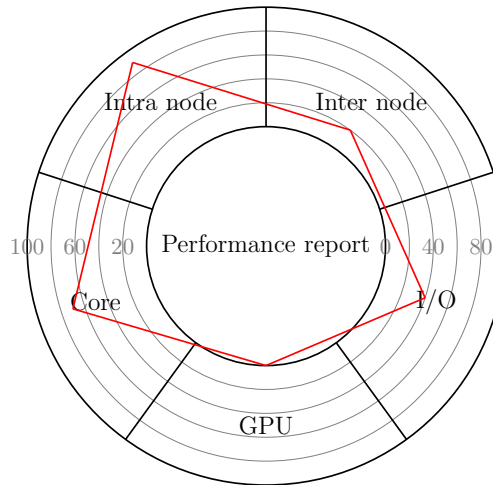
**Assessment crime 7** *Modern performance analysis tools (e.g. Intel Performance Application Snapshot (APS) and Linaro MAP) and methodologies offer a high-level overview as kind of an automatic service. This is a very valuable feature, as it provides you with a quick overview of your code's characteristics. While we appreciate these opportunities, we are however not sure if this is really the right starting point for performance assessment. In our work approach, we do not use one single run and its characterisation to get started, but we use many different runs (with different configurations) to zoom into particular feature dimensions such as intra-node scaling and see if there are issues there. Instead of one global overview, we obtain a set of metrics which we combine into an application profile. We do not rely on a **one run uncovers it all** approach.*

## 4.1 The assessment rubrics

We classify codes along five different rubrics:

1. core performance,
2. intra node,
3. inter node,
4. GPUs, and
5. input/output.

For each of these rubrics, we initially aim for a single three-valued: All is fine (green), there may be issues (yellow) or there are obvious flaws (red). In line with our overall presentation approach, we use thresholds of 60% and 80% to distinguish these classes. That is, if the GPU performance is somewhere close to 100%, all is fine. However, if it drops below 80%, we flag it as yellow. Once the metrics we define for GPUs underrun 60%, we claim that there are definitely major issues. It is convenient to present the overall data via a spiderweb diagram (below for a code that has I/O issues, doesn't scale particularly well between nodes, has no GPU yet reasonably good core performance and node scaling):



The idea to look at the five different performance dimensions separately is natural yet not does come along without shortcomings. Many performance issues within a code affect multiple dimensions. Think about a poor GPU utilisation that materialises in bad load balancing as well. When we dig down recursively into finer and finer metrics later on, we have to take such cross-influence into account. Eventually, we might have to return to our highest level of assessment again and again.

We propose to use different benchmarks for the individual rubrics. Each benchmark will follow the recommendations from Section 3.1, but we will tailor the benchmark

characteristics to the effect we want to study, such that the effects of interest are most pronounced while others are less dominant. Notably, we will disable I/O for all studies unless we study I/O itself (cmp. Section 3.5). This implies that the final picture might not reflect the real behaviour of a science run.

As we analyse the dimensions of performance characteristics independently, the order does not really make a difference. In theory. In practice, we have to start somewhere. We have personally had the experience that it is typically favourable to start with the node-level analysis (intra-node): If the individual cores are not used properly, MPI often does not make that much of a difference and we struggle to keep GPUs busy. Often, the next dimension is naturally identified once we look at the node-level performance, i.e. it tells us what to analyse next.

**Work in progress/todo** I think the previous paragraph might need to be updated, as we are typically beginning with core analysis now...

#### How it works

- ⊞ Ignore most things you know about tools to avoid that you are not able to see obvious shortcomings.
- ⊞ When you report on your findings, focus exclusively on what you observe. Avoid any attempt to interpret things.

#### How it does not work

- ⊞ Throw a huge, powerful tool onto your code and expect it to flag problems immediately.
- ⊞ Start to measure without a clear plan of what you want to measure.

## 4.2 Core performance

For the core-level analysis, we examine how efficiently your code harnesses the chip's microarchitecture. Does your code leverage the most advanced instruction sets available? Are we maybe even compute-bound? Does it properly utilise the memory interconnect? Are we memory-bound? Today's chips come equipped with powerful vector computing capabilities and sophisticated cache and streaming features designed to move data smoothly between memory and compute registers. Our analysis evaluates how effectively your implementation takes advantage of these capabilities.

To capture core behaviour at a high-level we simply want to measure the code's peak performance relative to the maximum of the hardware. Tools can tempt us to reach immediately for a roofline model, or potentially even more detailed metrics around cache misses, vectorisation ratios, etc. These topics are relevant to a core performance analysis, but we want to start high-level or we otherwise risk focusing on an individual metric which obfuscates or distracts our analysis. By focusing simply on peak performance, our initial core performance investigation should sit us somewhere in the 'traffic light' table below.

### 4.2.1 Benchmark preparation

Before beginning our data collection we want to ensure that our data is a reliable and accurate representation of the peak, single-core performance of the code.

**Assessment crime 8** *We assess a setup that fits completely into one of the caches.*

**Assessment crime 9** *We study a code that is not translated with the most aggressive compiler optimisation for the particular architecture that we study.*

Both properties are easy to check: Measure the total memory footprint of your code and ensure that the working set exceeds the largest cache you find on your system (cmp. remarks on caches below); techniques for making these measurements are given in Section 3.1. Further to that, check your compiler’s manual and ensure you use a bespoke code version for your target system—something that the `-march` or `-xarch` as well as the highest compiler optimisations (`-O3`) for example give you; greater details for these checks are given in Section 3.3.

To some degree, the argument holds the other way round for the cores: Modern chips are often not capable of cooling all of their cores if all of them run at full compute capacity. They therefore downclock the system when they encounter lots of compute-intensive operations. As we only assess scalar peak performance in this first exercise, this effect should not arise or be negligible. You might however want to double-check.

### 4.2.2 Data Collection

Once we have prepared our benchmark we are now ready to begin making measurements, and the data we need for this metric are:

- The maximum single-core, peak performance of the hardware
- The obtained single-core, peak performance of our executable

To make these measurements, we must turn to a performance tool. Multiple performance tools are able to make these measurements, and we list below basic instructions of how to perform these measurements with some of our ‘go to’ tools.

#### Likwid

To measure the peak performance of a hardware example with Likwid it is only a few lines

```
likwid-bench -t peakflops -w S0:16kB:1
```

This command pins the `peakflops` microbenchmark to a single core on the first socket (labelled `S0` here), to ensure maximum thread affinity. We use a convention advised by the Likwid developers in which the `peakflops` microbenchmark is performed with an input dataset half the size of the L1 cache []. For example, in the commands above we use a 16kB input which is half the size of a 32kB L1 cache (which we can find by running the `likwid-topology` command). As we

scale the `peakflops` microbenchmark across a node, we scale up this data input size by the number of cores, e.g., if want to run this analysis across 128 cores we use an input of  $128 \times 16\text{kB}$ , giving

```
likwid-bench -t peakflops -w N:128*16kB:128
```

though we stick to just single-core performance for now.

The instruction set usage is something we can evaluate per core. Depending on the character of the code, we might want to pick a more advanced instruction set (e.g. `peakflops_sse` or `peakflops_avx_fma`). As long as we argue about the overall code performance, we find the plain peak performance to be absolutely sufficient. Being within 80% of that one for the total program execution is typically a good sign.

Once we have the thresholds, the code's performance itself can again be used by multiple tools. Again, `likwid` is our personal favourite:

```
likwid-perfctr -f -C 0 -g FLOPS_DP ./mycode
```

which provides you with the information relevant, although the exact names of the groups might differ from system to system.

#### Amplifier

t.b.d.

#### Linaro MAP

t.b.d.

### 4.2.3 Evaluation

Once we know the maximum performance and the obtained performance, we can compute their ratio, i.e., the normalised peak  $C_{\text{peak}}$ . This metric is a little bit weird (and we will revise it later), as we typically normalise against the scalar peak performance. This is surprising given that all mainstream CPUs support vectorisation and hence can deliver several times more than that. However, we empirically found that starting with a comparison against the baseline is a good starting point: Few codes spend all of their runtime exclusively in very compute-heavy routines. Therefore, there will always be code parts which do not exploit vectorisation. If those that are compute-heavy do so and hence lift the average total peak into the regime of the scalar peak, this is already an overall efficient code.

**Performance metric (rule of thumb) 1** *The delivered MFlops/s normalised against the scalar peak is a good high-level indicator for the compute efficiency.*

| Peak FLOPS proportion            | Description  |
|----------------------------------|--|
| $C_{\text{peak}} \geq 0.8$       | Core compute performance is good.  |
| $0.8 > C_{\text{peak}} \geq 0.6$ | Core compute performance is not effectively using the hardware to its maximum potential. |
| $0.6 > C_{\text{peak}}$          | Core compute performance is poor.  |

If  $C_{\text{peak}} \geq 0.8$ , the code uses the compute capabilities (very) efficiently. If  $0.6 \leq C_{\text{peak}} < 0.8$ , the code does not use the compute units particularly efficiently and we have to investigate further.

#### How it works

- ▣ Run a standard benchmark to determine a realistic maximum throughput (bandwidth) and peak performance for a single core.
- ▣ Assess your code to find out if it makes sufficient use of the available bandwidth and compute units.

#### How it does not work

- ▣ Apply some performance counters as a black box without a focus on bandwidth and peak performance.
- ▣ Jump straight for a tool with gives significant data on vectorisation, cache misses, etc.
- ▣ Immediately dig into any performance flaw without evaluating the other four rubrics first.

### 4.3 Intra-node (node-level) performance

On an individual node, our first goal has to be that we use all the hardware threads efficiently. Some people prefer the term cores of threads. The term intra-node hence means all those cores sharing one memory. Hereby it does not matter if we have a code using MPI or a genuine multithreading language: We assess the runtime characteristics for a code physically sharing its memory. It might be written logically in a distributed memory fashion.

Our method of choice to argue about node-level performance is strong scaling. In principle, we ask ourselves “what would happen if we had only one core on our node” (we use the term core and hardware thread as synonyms from hereon). If our code needs 80 cores on one node, then we would expect that it only needs 40 cores on two, and 20 cores on four. This is obviously an artificial question to ask: Why would we not use all the node? Despite the fact that we have the cores available anyway and therefore might want to use them, finding an answer to the “what if” questions helps us to uncover inefficiencies. How much performance do we leave on the way as we use more and more cores? How much more computation should we be able to squeeze out of the machine?

The underlying performance model is called Amdahl’s law, and we will return to this law in a future chapter. For the time being, we can abstract from any performance model or law, and simply focus on the scaling behaviour over cores, as we have defined it in Section 2.3: We alter the number of cores available to our code until we have reached all the cores available, i.e. use the whole machine, and track the time to solution dependence on core counts.

### 4.3.1 Benchmark preparation

**Work in progress/todo** Do we insist upon a single problem setup? I don’t think so, but do we go with a rule of thumb that our memory use should be greater than 20% of our node memory. For example, we could

- **Core** - this benchmark should exceed the cache size
- **Intra-node** - this should consist of atleast 20% of memory, so on single core this could take up to say an hour, so long as the runtime decreases reasonably for the higher core counts

Throughout this assessment, we use one fixed problem setup. The setup has to be chosen such that we can compute it on a single core if necessary. However, it also has to be reasonably big such that it makes sense to exploit all threads of a node.

**Assessment crime 10** *We start an intra-node assessment with a code which only uses a tiny fraction of the main memory.*

We may assume that a node and its memory are somehow balanced. If we don’t use a reasonable amount of this memory for a test run, then there is simply not enough work to do to utilise all cores. We can barely blame the code. As a rule of thumb, we should have a memory footprint which is at least 20% of the total RAM available on the whole node.

### 4.3.2 Data collection

For the data collection, we first run the simulation without any parallelisation, which is, trivially, a single-core run. This gives us  $t_{\text{serial}}$ . After that, we switch to the parallel code version and vary the core count  $p_{\text{intra}}$ . This gives us  $t_{\text{intra}}(p_{\text{intra}})$ .

**Assessment crime 11** *Use a code base that is relatively inefficient on a single core and use this as baseline for any scaling claims.*

It is important that we question the code developers on different code variants prior to the assessment: If they offer, for example, quite sophisticated task-based parallelism for shared memory, we still should calibrate all runtime against a serial code without these tasks. That is a fair benchmark that avoids any tasking overhead where it is not needed.

Consequently,  $t_{\text{serial}} \neq t_{\text{intra}}(1)$  since we use different executables—once compiled with support for shared memory and once without. The two measurements will differ

and we can assume  $t_{\text{serial}} \leq t_{\text{intra}}(1)$ . If the code is not available without any parallelisation, use  $t_{\text{serial}} = t_{\text{intra}}(1)$ . After that, we collect the remaining  $t_{\text{intra}}(p_{\text{intra}})$  with the parallel code variant.

**Assessment crime 12** *We do not use the compute resources exclusively. For example in SLURM, clusters often have exclusive queues, or the `--exclusive` flag can be passed.*

### OpenMP

In OpenMP, we can simply set `OMP_NUM_THREADS` to achieve this. After that, we run our code and should automatically get the correct core count. Depending on your system's configuration, you might want to set some additional command line variables such as

```
export OMP_PROC_BIND=close
```

or similar to ensure that the operating system does not swap your threads around, i.e. decide at one point to migrate a thread to a different core. This never gives good performance. Be careful as well with setting the OpenMP threads manually while using another core specification technique such as Slurm's configs or tasksets (both below): If you do so, you might force all the OpenMP threads onto a few cores only, i.e. accidentally over- or undersubscribing the testbed.

### Slurm

Using Slurm, we do not have to set environment variables ourselves, instead Slurm manages resources and so we can request a quantity of cores with the `-c` flag, e.g.,

```
#!/bin/bash
```

```
#SBATCH -c 2
```

```
./my_exec
```

or equivalently use `#SBATCH --cpus-per-task 2`.

### Linux command line

The one thing that works always is using

```
taskset -c 0-3 ./myexec
```



to employ, e.g., the first four cores only. However, you should be careful when you use this in combination with Slurm or OpenMP which also make choices on cores. You might end up with competing or contradictory settings and end up with a suboptimal system usage.

### MPI

If we have a code base which employs MPI only, we set the number of ranks per node through the `mpirun` call or the scheduler (Slurm) environment.

For hybrid MPI+X codes, we have a certain freedom: We can use an arbitrary number of MPI ranks—indeed our analysis also works for pure MPI codes as long as we stick to one node—and use threads on top such that we exploit all  $p_{\text{threads}}$  hardware threads of interest. Picking one rank and four threads is, logically, equivalent to picking two ranks and two threads each. It is part of the assessment to play around with these degrees of freedom. Once again: If we spot sudden runtime jumps (yellow dots in example sketch) and if these jumps coincide with a NUMA domain or socket, then this might be a strong indicator that we should use MPI instead of threading. However, one might argue that this is already part of an in-depth analysis rather than a high-level overview.

We recommend that for a code of this form: when working with a researcher or developer, it best to request a typical configuration that they may use. They may be able to advise on their experiences, and from there we can then perturb about this configuration.

### Likwid

Likwid has the in-built `likwid-pin` tool for handling thread-affinity for a specified number of cores. Cores can be requested explicitly, i.e., specific NUMA domains, sockets, etc. can be requested as each core has a unique identifier (as with `taskset`). This gives the analyst greater flexibility in the choice of cores across the node, whilst also being flexible to incorporate with other Likwid tools.

T.B.D.

If you have the opportunity to mask out the initialisation time and to focus only on the “real” runtime of your benchmark code, use it at this point. Notably when you work with very short execution times, the initialisation overhead might otherwise dominate your measurements.

### 4.3.3 Evaluation

Our roadmap to determine how good the shared memory behaviour is, i.e. our strategy to come up with a single number, is straightforward. First, we calculate the code’s parallel efficiency. We use a slightly modified version as compared to Equations (2.1) and (2.2), where we calibrate the runtimes against the serial code runtime:

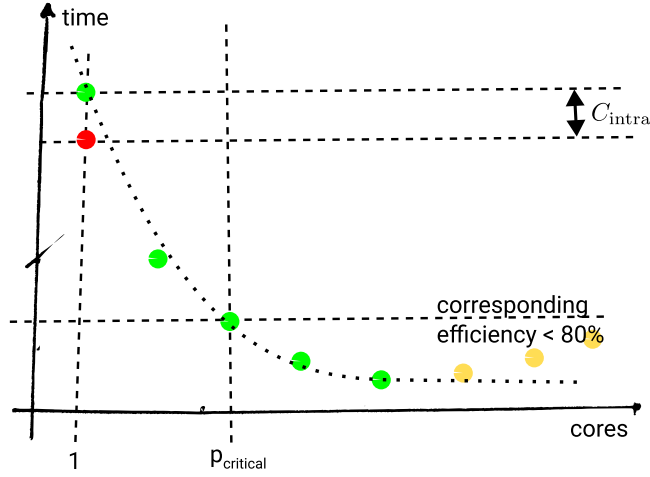


Figure 4.1: A scatter plot tracking the (strong) scaling data feeding into our intra-node analysis. The red node is the code’s runtime without any parallelisation enabled, while the green and yellow measurements stem from the parallelised code over various cores. For our high-level analysis, we are interested at which point (core count) the resulting parallel efficiency falls below an acceptable threshold.

$$E_{\text{intra}}(p_{\text{intra}}) = \frac{t_{\text{serial}}}{t(p_{\text{intra}}) \cdot p_{\text{intra}}}$$

is normalised against the serial runtime without any parallelism support. We track how much we sacrifice by introducing the parallelism in the first place, which is typically a better indicator for inter-node performance flaws than the textbook efficiency definition.

Second, we determine the number of cores for which our parallel efficiency falls below the thresholds 80% or 60% respectively (Figure 4.1):

$$\begin{aligned} p_{\text{critical,intra}}^{80\%} &= \arg \max_{p_{\text{intra}}} \{E_{\text{intra}}(p_{\text{intra}}) \geq 0.8\} \\ p_{\text{critical,intra}}^{60\%} &= \arg \max_{p_{\text{intra}}} \{E_{\text{intra}}(p_{\text{intra}}) \geq 0.6\}. \end{aligned}$$

It is not unheard of that codes yield data where the runtime starts to raise again once we exceed a certain core count (yellow dots in Figure 4.1). This means that the speedup deteriorates completely. The above metric tries to flag early if we run into such a regime or also into one that leads to a stagnation of the runtime and hence to an efficiency decline.

Finally, we introduce our performance metric. It normalises the critical thresholds against the maximum number of cores or threads  $p_{\text{max,intra}}$  that we have on our system.

**Performance metric (rule of thumb) 2** *We evaluate*

$$C_{intra}^{80\%} = \frac{p_{critical,intra}^{80\%}}{p_{max,intra}} \quad \text{and} \quad C_{intra}^{60\%} = \frac{p_{critical,intra}^{60\%}}{p_{max,intra}}$$

to determine how well the code scales over a given node.

We use the two metrics to determine if we consider the code to have an intra-node flaw.

| Efficiency  | Description  |
|---|--|
| $C_{intra}^{80\%} \geq 0.8$                               | Shared memory scaling is good.   |
| $C_{intra}^{80\%} < 0.8 \wedge C_{intra}^{60\%} \geq 0.6$ | Shared memory scaling is not particularly good, and we notably might run into problems with the next generation of chips where the core count will increase. |
| otherwise   | Shared memory scaling is poor.   |

#### How it works

- ⊞ Compile your code without multithreading/parallel support and compare its performance to a code without any parallelisation to determine the overhead of the parallel implementation.
- ⊞ Run a strong scaling study and compare the efficiency to prescribed thresholds how many cores you want to the system to use reasonably efficiently.

#### How it does not work

- ⊞ Skip intra-node scalability studies, as the code is “only” MPI. If a code is only MPI, you should still do strong scaling tests on a single node, where MPI physically uses shared memory.
- ⊞ Perform a sole strong scalability study without taking the implementation penalty into account.
- ⊞ Jump straight into a further “in-depth” analysis without evaluating the other high-level metrics.

## 4.4 Inter-node performance

**Work in progress/todo** This section is in early development, but the key takeaways currently are to perform a weak scaling analysis. Similar to the intra-node rubric, we again monitor parallel efficiency and note when it drops below 80% and 60%.

Inter-node behaviour means that the code runs over multiple physically separated nodes. In our case, this implies that it uses MPI. We work with a programming model which treats memory physically and logically as separate. There are software solutions that pretend to the code to have one large shared memory. Such virtually shared memory solutions would fall into this assessment category as well, so long as we use more than one physical processor.

Our method of choice to argue about inter-node performance is weak scaling.

### 4.4.1 Benchmark preparation

For weak scaling, we increase the problem size as we add additional nodes such that the approximate work load for each node remains relatively constant (in the ideal case).

**Assessment crime 13** *We scale up the problem such that the individual work per nodes becomes very small.*

As a rule of thumb, we can simply track the memory footprint per rank and ensure that it stays within a certain regime such as 60% of the available memory per node.

When preparing a benchmark for the inter-node, weak scaling analysis a key point is raised for the performance analyst: we must have some understanding of how to scale the problem size.

### 4.4.2 Data collection

The data collection is very similar to the strong scaling studies in Section 4.3: We track the cost of the computation (time-to-solution) over multiple nodes for each problem size. For the bigger setups, we won't be able to evaluate the problem on only one node, but that's, in general, not necessary. If we have a problem of size  $N$  on one node, a problem of size  $2N$  does only have to be evaluated on two nodes. There is no need to run this one on two nodes as well.

When preparing the benchmark we discussed that performance analyst must understand how to scale up the problem. However, when collecting data there are also key questions of how to scale up the resources.

For example, if we take an example of a hybrid parallel MPI+OpenMP code. Our intra-node analysis may have found that this code is efficient over 32 threads on one node (where we only define 1 rank on this node). How are we to scale up these resources? Do we define 32 threads per node and simply assign 1 MPI rank for each node? We could therefore run 2 nodes with 32 threads each, we cannot simply say that this is more performant than 4 nodes with 16 threads. We see this as already delving into a deeper analysis of efficient resource use and so would advise the analyst at this stage to continue behaving somewhat 'ignorant'.

When defining the resources for the weak-scaling, the analyst could again return to the code's developers and ask: what are the typical assignments of MPI ranks, threads, etc. for their code. For example, many multi-node GPU codes define an MPI rank per GPU. It is possible that the code's developers have some 'best practice' which the analyst could use and test. Otherwise, we simply advise scaling up the best performing intra-node resource definition.

### 4.4.3 Evaluation

Our roadmap to determine how good the distributed-memory behaviour is, i.e. our strategy to come up with a single number, is straightforward. First, we calculate the code's parallel efficiency. In the context intra-node scaling, we define parallel efficiency in terms of the serial runtime,  $t_{\text{serial}}$ , which is preferably a 'true' serial run, i.e., not just

a single core run but a run in which any parallel libraries are explicitly turned off. We here define our weak-scaling parallel efficiency as

$$E_{\text{inter}}(p_{\text{inter}}) = \frac{t(1)}{t(p_{\text{inter}})},$$

in which  $p_{\text{inter}}$  represents the number of nodes. We avoid defining inter-node parallel efficiency in terms of a serial run, as typically a single-node run will still be running in parallel on this node, e.g., using OpenMP multithreading across this node.

Using this inter-node parallel efficiency we can then define the node counts above which the efficiency drops below 80% and 60%

$$\begin{aligned} p_{\text{critical,inter}}^{80\%} &= \arg \max_{p_{\text{inter}}} \{E_{\text{inter}}(p_{\text{inter}}) \geq 0.8\} \\ p_{\text{critical,inter}}^{60\%} &= \arg \max_{p_{\text{inter}}} \{E_{\text{inter}}(p_{\text{inter}}) \geq 0.6\}. \end{aligned}$$

We then use these node counts to define inter-node metrics which are node count thresholds for 80% and 60% parallel efficiency normalised against some measure of the maximum number of nodes available.

**Performance metric (rule of thumb) 3** *We evaluate*

$$C_{\text{inter}}^{80\%} = \frac{p_{\text{critical,inter}}^{80\%}}{p_{\text{max,inter}}} \quad \text{and} \quad C_{\text{inter}}^{60\%} = \frac{p_{\text{critical,inter}}^{60\%}}{p_{\text{max,inter}}}$$

*to determine how well the code scales over a given node.*

**Work in progress/todo** I think we need a better definition of  $p_{\text{max,inter}}$  in this context. This is a less intuitive value than in the intra-node analysis. We either use the maximum number of nodes available on the system, or some upper bound set by the resource allocation, or even scale of the algorithm, etc.

| Efficiency  | Description  |
|---|--|
| $C_{\text{inter}}^{80\%} \geq 0.8$                                      | Distributed memory scaling is good.                  |
| $C_{\text{inter}}^{80\%} < 0.8 \wedge C_{\text{inter}}^{60\%} \geq 0.6$ | Distributed memory scaling is not particularly good. |
| otherwise   | Distributed memory scaling is poor.                  |

**How it works**

▣ t.b.d.

**How it does not work**

▣ t.b.d.

## 4.5 GPU performance

**Work in progress/todo** This section is still in its early stages of development, but the key takeaway is that: for a code with GPU offloading we simply want to measure the GPU utilisation/occupation. These terms can be used interchangeably, and so here we mean what proportion of the GPU's hardware does the code approximately use.

### 4.5.1 Benchmark preparation

t.b.d.

### 4.5.2 Data Collection

t.b.d.

Intel VTune

t.b.d.

AMD rocprofiler-system

t.b.d.

Nvidia Nsight

Nvidia Nsight Systems  
Nvidia Nsight Compute t.b.d.

### 4.5.3 Evaluation

**Performance metric (rule of thumb) 4** *The proportion of the GPU hardware used by the computational kernels. This is typically referred to as GPU occupancy.*

| GPU Occupation           | Description  |
|--------------------------|--|
| $C_{GPU} \geq 0.8$       | The GPU is well occupied with work.  |
| $0.6 \leq C_{GPU} < 0.8$ | The GPU is not being occupied with work well, this could become a significant issue with other GPUs including newer generations. |
| $C_{GPU} < 0.6$          | The GPU is not being significantly under utilised.   |

#### How it works

⊠ t.b.d.

**How it does not work**

☐ t.b.d.

**4.6 I/O performance**

**Work in progress/todo** This section is also still in its early stages of development, but the key takeaway is: we gauge the ‘high-level’ I/O performance of a benchmark by the proportion of runtime it spends in reads and writes.

**4.6.1 Benchmark preparation**

When preparing the benchmark for core, intra-node, inter-node and GPU analysis, we advised turning off I/O as this can make some of the measurements unreliable. However, I/O performance is crucial for some code’s and so here we advise turning on some example of typical I/O load, which again can be information provided by the code’s developers.

**4.6.2 Data Collection**

t.b.d.

**Darshan**

For a dynamic runtime analysis we setup Darshan with

```
export DARSHAN_DIR=/apps/developers/tools/darshan/3.3.1/1/
```

```
intel-2021.4-intelmpi-2021.6 # path to Darshan
export DARSHAN_LOGDIR=./darshan_logs
mpiexec -genv LD_PRELOAD ${DARSHAN_DIR}/lib/libdarshan.so
-np 4 ./my-exe
```

the logging environment variable name can be different so make sure to run

```
darshan-config --log-path
```

to find the name on your system. Once this has created a \*.darshan binary this can analysed via

```
darshan-parser --total my_log.darshan
```

The relevant metrics from this output are in the form

```
STDIO_F_*_TIME: cumulative time spent in different
types of functions.
```

Therefore depending on different forms of I/O — e.g., POSIX, MPIIO, STDIO, H5F, H5D, PNETCDF\_FILE, LUSTRE, DFS, DAOS — we need to sum together the read and write times.

**VTune**

t.b.d.

### 4.6.3 Evaluation

**Performance metric (rule of thumb) 5** *We evaluate*

$$C_{I/O} = 1 - \frac{t_{read} + t_{write}}{t_{total}}$$

where  $t_{read}$  and  $t_{write}$  are the times the code spends in reading and writing data, respectively, whilst  $t_{total}$  is simply the total runtime (inclusive of  $t_{read}$  and  $t_{write}$ ).

| I/O Proportion           | Description  |
|--------------------------|--|
| $C_{I/O} \geq 0.8$       | I/O does not dominate runtime.                     |
| $0.6 \leq C_{I/O} < 0.8$ | I/O forms a significant part of the total runtime. |
| $C_{I/O} < 0.6$          | I/O dominates the runtime.                         |

#### How it works

⊠ t.b.d.

#### How it does not work

⊠ t.b.d.

## 4.7 Localisation

So far, we have assessed our code as one big black box.

- We did not ask *which part of the code* causes a problem.
- We did not ask *when* it arises throughout the execution.
- We did not ask *where*, i.e. on which rank (if we run over multiple nodes) a problem arises.
- We distinguished, to some degree, *what* we want to study in our code, i.e. which quantity of interest.



Before we continue to assess the code and dive deeper into details, it makes sense to gather the information from the first three bullet points above or at least to discuss various approaches for gathering said information. Hereby, we must consider the type of flaw we are investigating: Single-core flaws imply that the “where” question does not play a significant role. Inter-node flaws might suggest that the “where” question is the primary concern to address.

Finally, the most frequent quantity of interest (metric) will be time spent in routines. In this case, the routines consuming most time are also called *hotspots* of a code. However, there are many other metrics that affect runtime, and we might want to study them separately—even though all of them eventually affect the time-to-solution<sup>1</sup>. Not all tools can provide answers for all metrics. We therefore must carefully select our repertoire of tools.

**gprof**

t.b.d.

**Intel VTune**

t.b.d.

### How it works

- ⊞ Once you have a first high-level overview, run a profiler to confirm that effect is either localised or uniformly found over all compute units.
- ⊞ If the flaw is found on few units only or only for a certain compute time or within one function which is not responsible for the majority of the runtime, study the flaw for this program phase, compute unit, function separately from all others.
- ⊞ Upon your very first profiling, ask a code expert which code parts (functions) actually perform calculations (cmp. Section 3.7).

### How it does not work

- ⊞ Jump straight into some traces without spending time to think what metrics and program parts are relevant to answer your hypotheses regarding performance flaws.
- ⊞ Draw conclusions from global observations on all code parts and compute resources (cores or ranks).

<sup>1</sup> Some codes are interested in other metrics such as energy, but they are still rare.



# **Part III**

## **Case Studies**



## Chapter 5

# High-level assessments

**Work in progress/todo** In this chapter we aim to include example high-level performance assessments applied to a handful of code examples. These high-level assessments are reports which are generated from a Jupyter notebook (currently in development).



# **Part IV**

## **Discussion**





## Chapter 6

# Summary and outlook

This document is still early into its preparation. However, it demonstrates some key points around our main focus of building a performance analysis methodology that seeks to avoid jumping straight for a performance tool. Instead we want a methodology that starts high-level and only reaches for tool, when appropriate. Rather than being lead by the tool, we want to be lead *towards* a tool by the methodology.

### 6.1 Recap

Though this iteration of the document is just a brief introduction, it covers some of the necessities for establishing a rigorous performance assessment service.

We believe that beginning with an overview of terminology in Chapter 2, including the fundamentals of HPC, is a vital starting point. Many users and developers come to HPC without a technical background in computing, and it is important now more than ever to set out an agreed set of terms for this performance analysis methodology. Beyond acting as a glossary of HPC and performance terms, Chapter 2 also establishes some of the key concepts of how we are to conduct a performance assessment such as treating a code as a blackbox and noting measurements of runtime performance objectively, without jumping quickly to deductions.

Preparation for a performance assessment is key to reproducibility. Chapter 3 lays out some details of our ‘experimental reproducibility’, i.e. before carrying out a performance assessment, we need to establish the fundamentals of setting up our benchmark. Once the benchmark example has been agreed upon, setting up the compiler and other variables such as I/O, etc. is a crucial part of how to make sure that varying, e.g. the number of cores gives reliable data. It should be noted that within this preparation, we may also see performance gains in tuning compiler options. This preparation stage is a key point of the performance assessment in which engagement with the service user is crucial. In order to get a useful benchmark example which is indicative of a real simulation run, the domain specialist needs to set this up.

This current document finishes with a high-level approach to an assessment in Chapter 4. One reason for this is we find often without this initial high-level approach,

HPC users will immediately jump to a tool which can give a myriad of data and metrics. This sometimes sends an HPC user down a wrong path as they may not be following a structure methodology to rigorously analyse performance; or, we have even found that some people are put off by the overwhelming data and the user gives up on the analysis as they get a sensation of: “where do I begin?”.

The high-level approach can be a stepping stone into a more detailed analysis as it simply breaks down performance into five topics: core, intra-node, inter-node, GPU, I/O. For some tool developers, the tools direct the user towards an almost ‘Matryoshka doll’ model of performance, i.e. core performance sits within intra-node performance which sits within inter-node performance. We believe this to be an over simplification, and in reality, these performance topics can never be truly de-coupled from one another. The high-level approach presented here is not seeking to focus the analyst down one particular avenue or another, it simply gives an overview of the performance landscape for a particular code.

## 6.2 Outlook

Having established that this document is not the finished work, it is worthwhile detailing what is to come in future versions. In particular: 1. what is left out of the current methodology; 2. what is to be discussed about the implementation of this methodology into a performance assessment service.

Firstly, performance topic decision trees. We believe there is value in beginning with the high-level performance assessment, and we indeed recommend beginning with this for any performance assessment. However, this is by definition not a complete analysis in and of itself. Through a series of performance analysis workshops we have begun to design a series of decision trees.

Each performance topic will have a decision tree, i.e. having conducted the high-level assessment the analyst will have to make a decision as to which topic to begin with and the decision trees will direct the analyst through a series of binary outcomes to isolate individual performance issues. For example, for intra-node performance we can begin by asking if the code is dominantly parallelised, if not then our performance engineering should be directed towards increasing the scope of parallelisation (though again we note that said engineering falls outside the scope of this work). If there is significant parallelism, then we can ask if there is good load-balancing. If there is good load balance then we can explore synchronisation, if not we turn to checking if there are enough tasks, etc. as we work down the tree. We plan to have decision trees for each performance topic with good progress made on core, intra-node and inter-node with our focus next turning to GPU and I/O.

Secondly, once we have introduced the decision trees we can describe our entire performance assessment methodology including case-studies. As we work through the creation of this manuscript, we are employing this methodology on existing codes whilst also engaging with other users/developers through performance analysis workshops. We include points on building the performance analysis results into a finished report alongside these case studies.

Lastly, we will give details as to the implementation of the assessment service from

an organisational stand point. Beyond setting out the assessment methodology, there are also a number of points around how do you actually build a functional assessment service from the interactions with scientists and developers through to the interface of the service.

One component of the service interface is planned to be a reproducibility analysis, i.e. if an HPC developer is deciding to engage with a performance assessment then it is worth auditing the reproducibility of their software. For example, is the code documented, using version control, implementing CI/CD pipeline, etc. A performance assessment can mean a developer is already giving their time to stop and reflect on their code and so introducing some time to also implement some good software practices can have two effects: 1. it helps the sustainability of their software going forward; 2. it helps the job of the analyst if the building and running of the software is well documented and hosted on a well maintained repository.

In summary, this document lays out just beginning of delivering a performance assessment, and the future iteration of this document will expand upon this methodology whilst also discussing the infrastructure of how to implement this methodology as a functional performance assessment service.



**Part V**

**Appendix**



# Appendix A

## Acknowledgements

### A.1 Funding councils and supporting initiatives

This work was supported by the Science and Technology Facilities Council [grant number UKRI/ST/B000293/1]. This work also was supported by the Engineering and Physical Sciences Research Council [grant number UKRI1801]. The underlying projects HAI-End and SHAREing (Skills Hub for Accelerated Research Environments Inspiring the Next Generation) is part of the cross-UKRI Digital Research Infrastructure initiative.



### A.2 People

This manuscript is the result of many meetings, workshops and conversations, and many people chipped in ideas and remarks that eventually helped us to improve the text. Special thanks are due to

- Brian Wylie
- Ben Clark who knocked down the first prototype of the accompanying Jupyter notebook in only one week of his school placement.
- Eva Fernandez Amez who created all the upskilling videos that we created whenever we spotted topics at the workshops which deserve a brief revision.

### A.3 Partners

We thank the companies AMD, Intel and Linaro for their continuing support of these activities by sending delegates to our workshops and contributing to our document.

Durham University is a proud partner of VI-HPS (<https://www.vi-hps.org/>) and appreciates the input of our colleagues and friends of the research consortium. Without their input, none of this would have been possible.



## Appendix B

# Submission checklist

This chapter hosts a set of checkpoints that groups might introduce *before* they start to conduct assessments and analyses. They form kind of a common sense contract between the code owner and the person conducting a review and ensure that some typical hiccoughs and problems are avoided. As it has this legal flair, we call the parties involved the code *owner* vs. the code *analyst*, even though these roles in practice might be blurred and taken up by code developers themselves. The checklists complement the aspects and checklists discussed in Chapter 3. While Chapter 3 is written from an analysis/assessment point of view and hence targets predominantly the analyst, the checklists here are more owner-focused and notably highlight the interaction between analysts and owners.

### B.1 Code submission

Before we accept a piece of code for the assessment and analysis, it is important to identify what we expect from the owner submitting the code. This covers first of all the artefact itself:

| Checklist (Code) |  |
|------------------|--|
| 1                | There is a central archive with all files required to run the code. The archive is self-contained.   |
| 2                | This archive is frozen, i.e. no one changes the object of an assessment while the assessment is going on.  |
| 3                | There is a well-defined, characteristic run of the code which the analyst can use to study performance of the software. This benchmark should not take 'too long'; for further details please see Section 3.1. |

The second point is obviously once written down, but deserves to be highlighted: If an object of study changes as we try to assess it, it is virtually impossible to get the big picture right. Obviously, continuous benchmarking and performance analysis can and

should be part of professional scientific software development. However, if we try to conduct an objective assessment, we have to work with a fixed object of study. That means, we should freeze the development for the time being or at least work on a stable (release) branch. For example, the owner can create a *tag* within the repository. Tags are often used to annotate specific code release versions, but here we can also tag the repository at a point for assessment. If a code does not (yet) support all dimensions of an assessment (cmp. Chapter 4) or some dimensions are not yet stable, this is the time to flag this, so we can exclude them from the analysis.

Submitting a characteristic run means that the owner describes how we can start such a run on the reference system ourselves. This is the litmus test: If we cannot (even) start a baseline run, then the code/submission is obviously not yet in a shape to be assessed and we have to stop immediately.

The attribute “characteristic” requires the owner to doublecheck if the setting handed over makes sense. Both parties want the outcomes of the benchmarking and analysis to be meaningful. We assume that the owner wants to learn something about the code from an objective analyst, while the analyst wants to know right from the start that their work is relevant. We want to avoid a situation where we are told afterwards that the analysis’ results are meaningless because “in the real world, we would use a completely different setup”. This has happened to us before.

Further to the mere artefact submission, we also require more detailed information from the owner. We begin by outlining the required documentation that we expect from the owner:

| Checklist (Documentation) |   |
|---------------------------|---|
| 1                         | Software dependencies, e.g., libraries such as FFTW, LINPACK, etc. preferably with versions if known. |
| 2                         | Compilers, preferably with version if known, along with the typical compile flags.                    |
| 3                         | Parallel libraries/models, e.g., MPI, hybrid parallelism with MPI+X, CUDA, OpenMP, etc.               |
| 4                         | Detailed build instructions.  |

This required documentation we see as necessary for an analyst to start by cloning the repository, and almost without thinking, produce a correct executable for the “characteristic” code run. Though we treat the code as a black box, the parallel model is requested to guide the analyst to which aspects of a software’s performance can be analysed, i.e., clearly a code which supports MPI only cannot have its GPU performance analysed.

With this information an analyst should be able to build the executable and so next we request typical resource allocations for the software such that the analyst can begin to define job scripts and plan out the assessment. The requested resource allocation information is listed as:

**Checklist (Resources)**

- 1 The owner should, if known, state some example runtimes of the characteristic benchmark along with the hardware examples and configurations that were used.
- 2 The owner should submit example configurations of the parallel resources used. For example, if CUDA+MPI is used then how many ranks are assigned per GPU? If using OpenMP+MPI then how many threads are typically assigned per rank.

This information is to direct the analyst towards ‘typical’ job and resource allocations. For example, some hybrid parallel MPI+OpenMP codes use quite specific allocations of ranks per node and threads per rank, and it is not feasible to explore all combinations. Hence, requesting an example job setup, preferably in the form a job script, can direct the analyst and make the performance report more relevant to the owner. Furthermore, knowledge of typical runtimes along with job setups can be useful for the analyst to sanity check results and support that the benchmark configuration is in keeping with the guidelines touched on in Section 3.1.

Having built the executable and designed the relevant job scripts, we finally request information on benchmark configuration:

**Checklist (Configuration)**

- 1 The owner should describe how to vary the benchmark, i.e., which flags can be varied including typical value ranges for these parameters.
- 2 The owner should supply details of how the benchmark scales with the input parameters.
- 3 The owner should describe how, if possible, the I/O of the benchmark can be switched off or on.

The questions on configuration — at least for the high-level assessment — are mainly for a weak scaling analysis, because this is likely the only point in which the analyst is straying away from the example benchmark. Therefore, the analyst needs to understand how to scale the benchmark both from an interface standpoint of how do we vary the input parameters, but also how does the algorithm scale. Put simply: how do we vary the input parameters to, e.g., double the problem size?

The final question on I/O is that when not studying I/O performance, it can be beneficial for the analyst to be able to switch this off and avoid I/O skewing performance measurements.

All of the requested information above is to aid the analyst to be able to perform our high-level assessment as just “turning the handle”. However, technical documentation of a code is not all that is required for a performance assessment service; professional conduct and managing expectations.

## B.2 Code of conduct

Above we discussed the importance of defining characteristic benchmarks to avoid (from experience) performance reports being disregarded as not indicative of “real world” runs. Along the lines of “this has happened to us before”, we here write down some simple principles how to behave regarding the analysis:

| Checklist |  |
|-----------|--|
| 1         | The owner is available to the analyst for questions around the code while the assessment project is going on, and try to answer questions re code usage and setups quickly and constructively. |
| 2         | The owner will refrain from challenging the assessment methodology.  |
| 3         | The owner commits to acknowledge all.  |

The first item can be subject of debate.

**Work in progress/todo** Greater context needs to be added here.

## B.3 Outcomes

It is important to sort out what happens with any outcomes before the assessment starts. The two items below can be subject of debate, but they align with principles of open, fair and reproducible science:

| Checklist |   |
|-----------|---|
| 1         | The owner grants the analyst the right to publish their review.   |
| 2         | The analyst commits to an objective assessments which refrains from “blaming” a code for certain behaviour. |

This sounds defensive and it is defensive, but we have observed numerous situations where developers reacted annoyed up to angry when they read assessments. Scientists tend to develop emotional attachment to their codes, and are not amused if an independent analysis highlights shortcomings. It is important to acknowledge that this is exactly the purpose of an analysis. It is not there to confirm how exceptionally great some piece of code is already, but should inform future development.

# Bibliography

- [1] Kent. Beck. *Extreme programming explained : embrace change*. Addison-Wesley, Reading, Mass, 1999.

# Index

- AMD rocprofiler-system, [46](#)
- Amplifier, [37](#)
- Core, [18](#)
- Darshan, [47](#)
- Efficiency, [41](#)
- gprof, [28](#), [49](#)
- Hardware counters, [17](#)
- Hardware thread, [18](#)
- Intel compiler, [24](#)
- Intel VTune, [46](#), [49](#)
- Likwid, [27](#), [36](#), [41](#)
- Linaro MAP, [37](#)
- Linux command line, [27](#), [40](#)
- MAQAO, [25](#), [27](#)
- MFlops/s, [37](#)
- MPI, [41](#)
- Node, [18](#)
- Nvidia Nsight, [46](#)
- OpenMP, [40](#)
- Parallel efficiency, [19](#)
- Peak performance, [37](#)
- Profiling, [16](#)
  - Exclusive, [17](#)
  - Inclusive, [17](#)
  - Instrumentation, [17](#)
  - Sampling, [17](#)
- Requirements
  - Non-functional, [13](#)
  - Slurm, [40](#)
  - Speedup, [19](#), [41](#)
  - Tracing, [16](#)
    - Filtering, [17](#)
  - VTune, [28](#), [48](#)