
EXTENDING JUST-IN-TIME COMPILEATION FOR OP2

PROGRESS REPORT

Nathan Dunne - u1604486

Introduction

The aim of this project is to contribute to the OP2 open-source project[2], an Embedded Domain Specific Language for developing unstructured mesh based applications. It provides a programmer access to greater performance using an abstraction from hardware specific optimisations, as well as the further benefit of allowing portability between different platforms. This portability means that all programs written using the API will be able to utilise the latest hardware generation by updating only the OP2 library, a regenerating the platform specific codes.

On top of the existing work, there is space for greater performance to be gained through the use of Just-In-Time (JIT) compilation: re-compiling the source code at run-time once the inputs are known, to allow expensive repeated operations to be replaced with static constants where possible, expanding the ability of the compiler to optimise through constant propogation and other methods. This can provide significant benefit, especially over a large number of iterations as a small improvement per-iteration can compound to better offset the one-time cost of the re-compilation.

The inputs to an OP2 program will usually be the description of the unstructured mesh, as a file which is read at runtime. Since this information is not known to the ahead-of-time compiler, it is only possible to optimise based on this information "Just-In-Time", or immediately before execution.

The sequential JIT code generator has already been completed, so the primary goal is to target the NVIDIA GPU architecture using CUDA, attempting to reach a performance improvement over ahead-of-time compiled CUDA generation. Since the OP2 API itself will not be modified by the project, existing applications can be used to benchmark performance. Primarily this will be Airfoil [9], an small but industrially representative computational fluid dynamics (CFD) program, and later VOLNA[8], a non-linear shallow water equation solver. Both use the C OP2 API.

Motivations

The motivations behind this project are a desire to use and expand on knowledge of optimisations and high-performance computing gained in the second year *Advanced Computer Architecture* module, as well as an interest in compilers and code generation. It will also be personally beneficial to gain experience contributing to a real open-source project, with an existing code base to deal with. It is gratifying to know that if the project is completed to a high standard, it could be merged into the master branch and used beyond the submission for a module.

The decision to target GPUs is as a result of the increased need for parallelism in high performance computing systems and applications. GPUs are in nature aimed at graphics processing, which requires many operations to be performed in parallel, but similarly many HPC programs (including unstructured mesh apps like CFD programs) can benefit from this architecture now that they are capable of double precision.

Research

In order to begin contributing to the OP2 source, I needed to build an understanding of the existing work. There are 2 main sources for this information:

- Papers published on the OP2 website
- Analysis and Benchmarking of existing source code

While academic papers are useful to understand the context of the project and the motivations behind it [7] [6] [10]; the existing code gives a more concrete understanding of the implementation, and the form my contributions will take. Understanding the basic structure of an OP2 program, where data is divided into sets and mappings between them - and loops can be defined as a separate "kernel" of operations, is critical to making a useful contribution.

Technical Content

Prerequisites

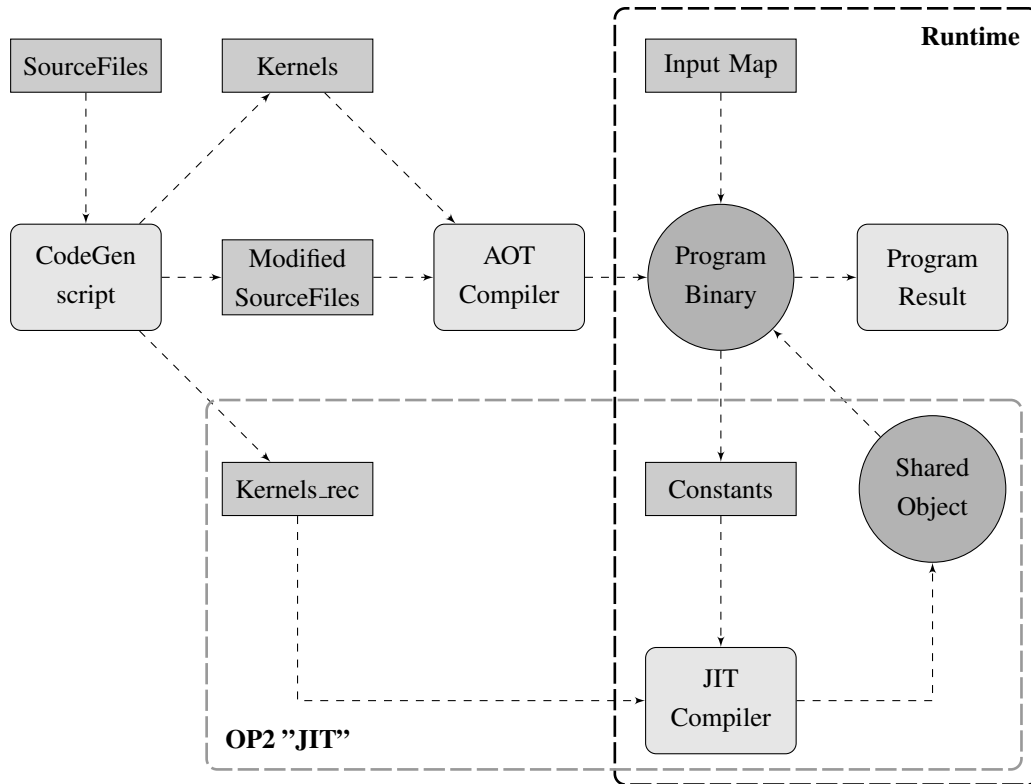
OP2 utilises two well established parallel mesh partitioning libraries: ParMetis[3] and PT-Scotch[4]. As I intend to work mainly on my personal laptop, a version of each had to be retrieved and built with the correct build arguments. The work with CUDA also required the CUDA Software Development Kit[1], and drivers for the Nvidia MX250 card I have available. Finally, the airfoil example application makes use of HDF5 file I/O operations, so a build of that library was required as well.

Tools

The main tool used in development is GitHub as a remote repository. OP2 has a number of feature branches already, with the existing JIT work done in the branch *feature/lazy-execution* of the Git repository. In order to extend this work, I've rebased the branch onto the master branch in order to benefit from the many changes committed to master since they diverged, and created a new branch *feature/jit* to track my work. Using a remote git allows for some resilience to failures on my local machine, as changes can be preserved by pushing them to the remote. It also allows easier version control, and potentially further extension to my work from other contributors in the future.

Existing Work

In gaining an understanding of the previously completed sequential JIT code generation, it became clear that "Just-In-Time Compilation" in the case of this project doesn't strictly mean the same thing it's usual definition, for example in the JVM where the intermediate representation bytecode is compiled to native machine code at run time, allowing optimisation of bytecode to take the inputs into account. In this project, the Ahead-Of-Time compiled executable issues a make command during the run of the program, so that loop functions can be generated using constants which are supplied in the input file defined as macros that can be expanded by the compiler's preprocessor. This can only be done at runtime. The regeneration of loop functions using constants should provide performance benefit over equivalent ahead of time backends, as long as the performance gain is sufficient to offset the time cost of recompiling. The process is shown below, with the operations unique to the JIT branch encircled.



To write a new code generator script, the arguments passed to existing code generator scripts for each backend:

1. masterFile: the name of the first file passed to op2.py
2. date: the date today
3. consts: master list of declared constants in source file
4. kernels: master list of loop kernels declared in source file

Where *consts* is populated from calls to *op_decl_const* in the original source files, ensuring there are no duplicates; and *kernels* comes from calls to *op_par_lllop* (also unique), where each parallel loop becomes a function with the required number of parameters, and a function definition based on the desired backend. Each loop also requires a definition of what needs to be done on each iteration, which for the example application *airfoil* is a separate header file of the same name, but this is not required.

The code generator script will iterate over the list of kernels generating the required backend inside a folder for neatness. When the executable is built it will be linked against a backend folder to link the correct definitions for each parallel loop function.

As well as the research above, the following work has been completed so far.

CUDA familiarity

Since the CUDA c++ API is not something I have used before, I experimented with writing small programs to run on a gpu, as well as referencing the *Nvidia CUDA C Programming Guide* to build my Understanding. Knowledge of how CUDA programs can solve common problems in parallelism, such as local and global shared memory, and the use of *syncthreads()* will undoubtedly be necessary to complete this project.

Hardware Resource

The graphics card in my personal laptop is sufficient to execute the generated CUDA code, and ensure correctness, however the system will be too noisy to gather representative benchmarks. To this end, I applied for and have been granted access to the Orac High Performance Computing node in the department, which will be used only for benchmarking.

Implementation

At time of writing the implementation of CUDA JIT code generation has just started. The `op2.py` script has been modified to call a `CUDA.jit` code generation file, and added to the makefile of `airfoil` to include an `airfoil.CUDA.jit` entry, however the codegeneration itself is still sequential. To complete the CUDA JIT implementation I will use the AoT CUDA implementation as a guide to speed up the process. The work is available in my *feature/jit* branch of the remote repository.

Reflection and Project Management

The timetable from the project specification (Figure 1) shows the expected level of progress at this stage. In retrospect, it is clear that the original timetable was ambitious, as the project has not run quite to time. While I am still confident that the CUDA JIT code generation can be done by the end of the autumn term (as originally agreed with my supervisor), I have reworked the timetable with a greater allocation of time given to research, and three iterations of longer implementation blocks with one week testing, rather than five rounds of 3 week blocks. Now that I have a greater understanding of the project, it seems unreasonable to suggest there would be a feature worth spending a week testing and benchmarking completed every 2 weeks of development. See Figure 2. As well as aligning with the end of the Autumn term, the modified timetable still allows for the same amount of total development time.

I would attribute the need for an extension to research time to taking an ad hoc approach to working on the project, which has resulted in not putting as much time in as I had planned, instead focussing on closer deadlines. As such I have decided that for next term I will set a weekly timetable of allocated time to work on the project.

The timetable aside, other elements of the project haven't provided any significant unexpected issues. All tools and libraries are functioning and compatible.

Ethics and Risks

As mentioned in the Project Specification, there are no obvious social or ethical issues with the project, as it will not require the collection or storage of any personal data. It is important to consider the licence under which the open source project is held, which permits redistribution of the source and binary, as long as it contains the copyright disclaimer.

I am also now aware of the Acceptable Use Policy[5] for the computing facilities provided by the SCRTP, so I will also ensure to abide by this when benchmarking using Orac.

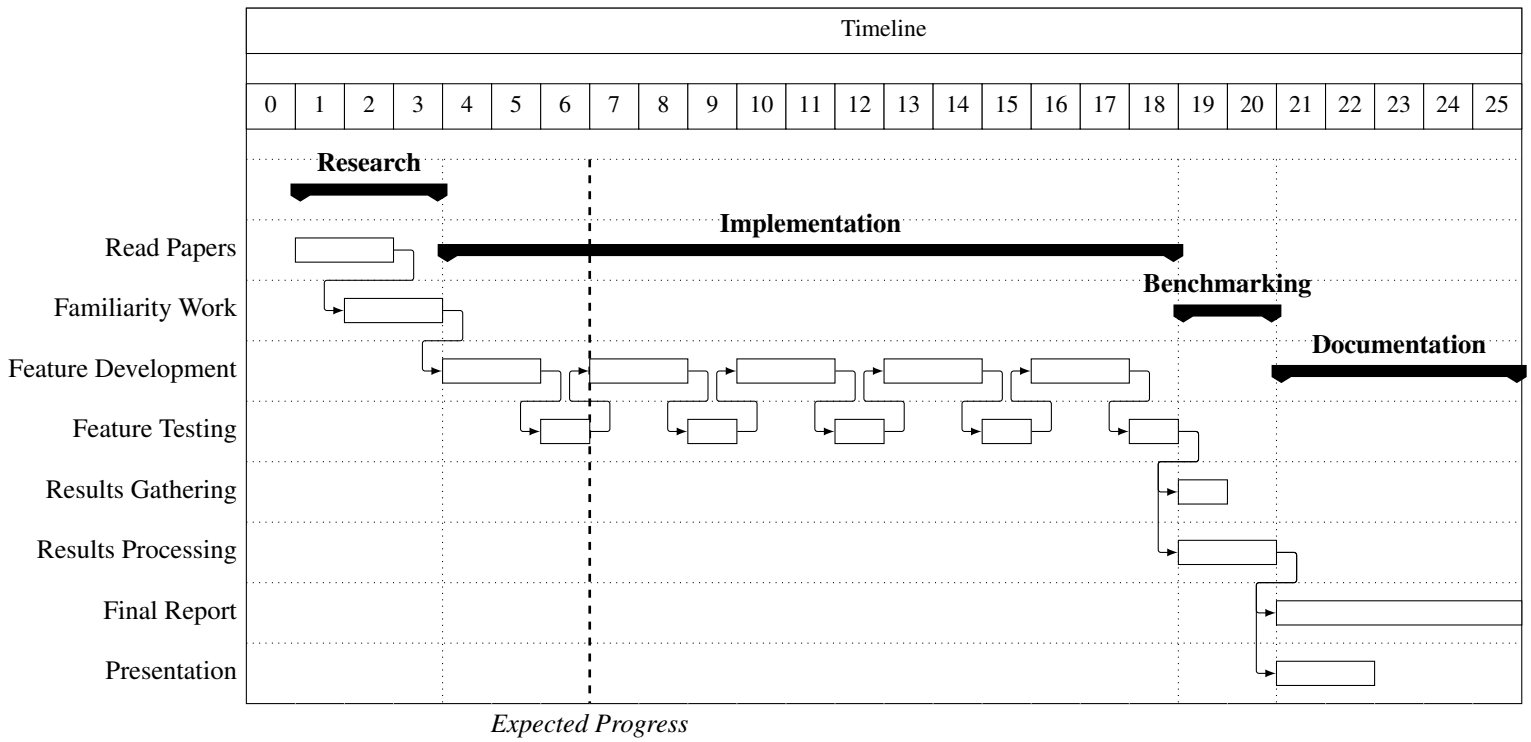


Figure 1: Original Timetable

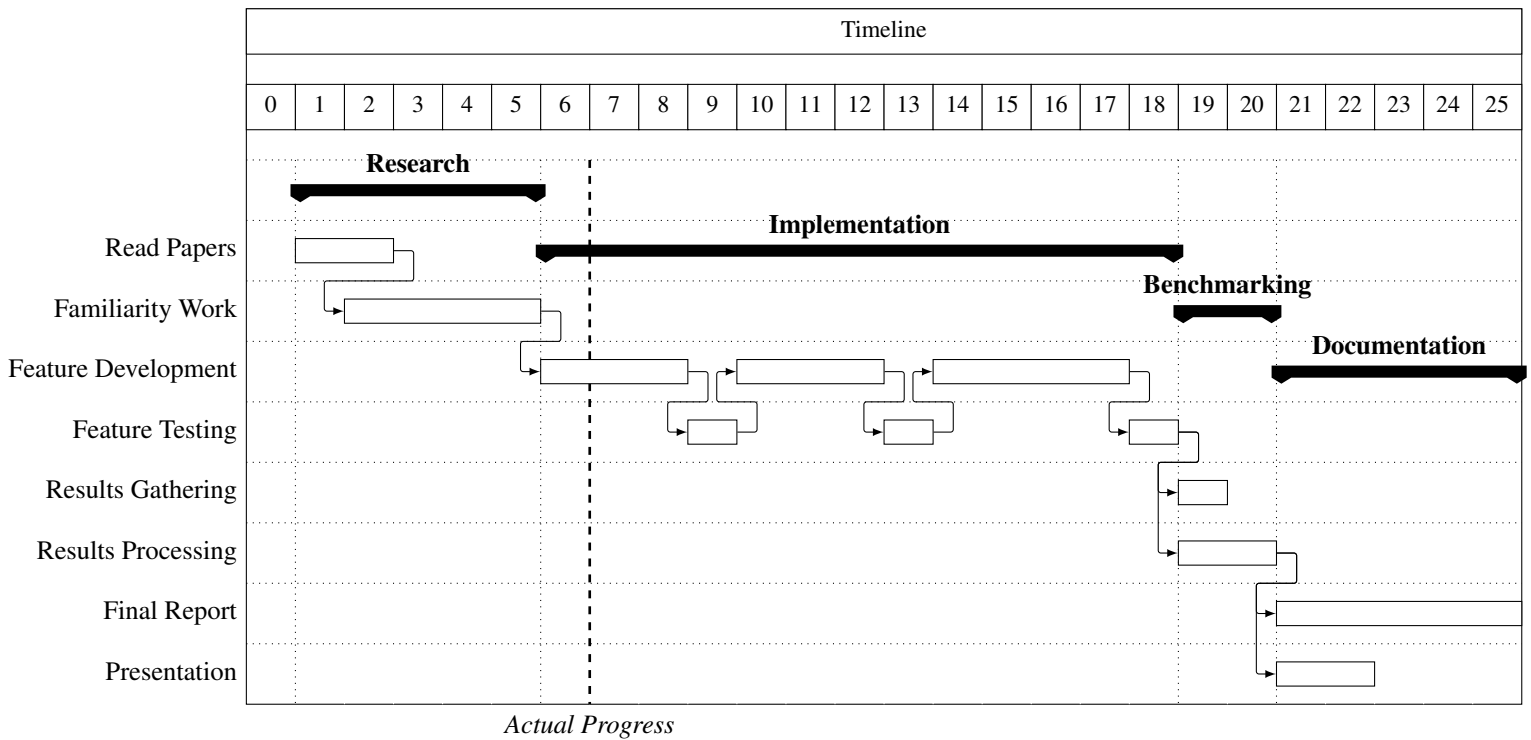


Figure 2: Modified Timetable

References

- [1] Cuda sdk. <https://developer.nvidia.com/cuda-toolkit>.
- [2] Op2-common. <https://github.com/OP-DSL/OP2-Common>.
- [3] Parmetis. <http://glaros.dtc.umn.edu/gkhome/metis/parmetis/overview>.
- [4] Pt-scotch. <https://www.labri.fr/perso/pelegrin/scotch/>.
- [5] Scientific computing rtp acceptable use policy. <https://warwick.ac.uk/research/rtp/sc/policies/aup>.
- [6] I.Z. Reguly G.R. Mudalige and M.B. Giles. Auto-vectorizing a large-scale production unstructured-mesh cfd application, 20. <https://www.oerc.ox.ac.uk/sites/default/files/uploads/profile-pages/Gihan/GRM-WPMVP.pdf>.
- [7] M.B. Giles G.R. Mudalige, I. Reguly. Op2: An active library framework for solving unstructured mesh-based applications on multi-core and many-core architectures, 20. <https://www.oerc.ox.ac.uk/sites/default/files/uploads/profile-pages/Gihan/InPar20.pdf>.
- [8] et al. Istvan Z. Reguly, Daniel Giles. The volna-op2 tsunami code, 2018. <https://www.geosci-model-dev.net/11/4621/2018/gmd-11-4621-2018.pdf>.
- [9] Istvan Reguly Mike Giles, Gihan Mudalige. Op2 airfoil example, 2012. <https://op-dsl.github.io/docs/OP2/airfoil-doc.pdf>.
- [10] I.Z. Reguly et al. Acceleration of a full-scale industrial cfd 30 application with op2, 20. <https://www.oerc.ox.ac.uk/sites/default/files/uploads/profile-pages/Gihan/OP2-Hydra.pdf>.