



Just In Time Compilation for a High-Level DSL

Nathan Dunne

1604486

3rd Year Dissertation

Supervised by Gihan Mudalige

Department of Computer Science

University of Warwick

2019–20

Abstract

TODO

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Key Words

High Performance Computing, Unstructured Mesh,
Just-In-Time Compilation

Contents

Abstract	ii
Key Words	ii
List of Figures	iv
1 Introduction	1
1.1 Background Work	1
1.2 Motivations	3
2 Related Work	3
2.1 OP2	3
2.2 JIT	3
2.3 CUDA	3
3 Specification	3
3.1 Runtime Assertions	4
3.2 System Model	4
4 Implementation	6
5 Testing	6
6 Evaluation	6
7 Future Work	6
8 Conclusion	6

Appendices	8
A example	8
Acknowledgements	9

List of Figures

1	Tri-Structured Mesh	2
2	Airfoil Tri-Unstructured Mesh	2
3	OP2 System Diagram with JIT Addition	5

1 Introduction

In the field of High Performance Computing (HPC), computers with processing power tens or hundreds of times greater than conventionally available machines are used to solve (or approximate solutions to) problems that would otherwise take an unwarrantable amount of time. Such computers have been required for some time to make use of a large degree of parallelism in order to complete with reasonable runtime: dividing work into independent subsections which can be executed simultaneously.

Many paradigms for executing parallel workloads have emerged over time: including vector instructions (SIMD), many and multi-core CPUs, clusters of interconnected computers, and General Purpose Graphical Processing Units (GPGPUs). Hardware which was originally specialised for graphical shader calculations through its very high number of processing units, allows carrying out the same operation across a very large amount of data in parallel. This hardware has been adapted in GPGPUs to perform non-specific operations that would normally have been done by the CPU.

Furthermore, a large proportion of HPC workloads involve approximating Partial Differential Equations (PDEs) to simulate complex interactions in physics problems, for example the Navier-Stokes equations for computational fluid dynamics, predicting weather patterns, or computational electro-magnetics. It is usually necessary to discretise such problems across some form of mesh, either structured (regular) or unstructured. Unstructured meshes will be discussed further in Section 1.1: Background Work.

1.1 Background Work

The OP2 library is an Open Source Domain Specific Language (DSL) which provides a high level abstraction for describing physics problems which can be abstracted to an Unstructured Mesh. Unstructured Meshes, such as Figure 2, use connectivity information to specify the mesh topology. The position of elements is highly arbitrary,

unlike structured meshes where elements follow a regular pattern (Figure 1). A particular simulation might, for example, be approximating the velocity of a fluid in each cell based on the cells around it.

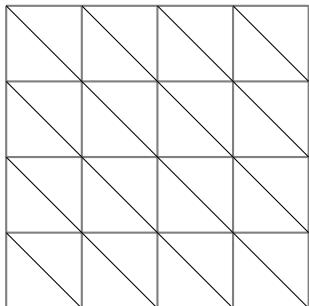


Figure 1: Tri-Structured Mesh

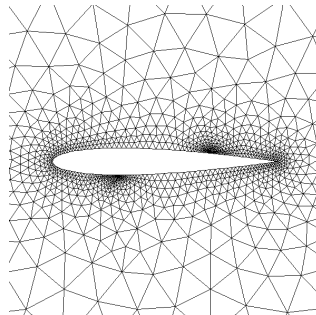


Figure 2: Airfoil Tri-Unstructured Mesh

This abstraction allows scientists and engineers to focus on the description of the problem, and separates the consideration for parallelism and efficiency into the back-end library and code generation. This is beneficial as it is unlikely that a single developer or team has the necessary expertise in both some niche area of physics with a non-trivial problem to be solved; and sufficient depth of knowledge in computer science to understand and utilise the latest generation of parallel hardware.

A further benefit is that existing solutions which make use of the OP2 library can be ported onto a new generation of hardware, by modifying only the OP2 back-end library to support the new hardware, instead of every application individually. This portability can save both time and money in development if multiple different hardware platforms are desired to be used.

The OP2 library is already able to generate optimised code for a number of platforms, including the increasingly popular NVidia CUDA for parallel programming on NVidia GPUs. However, there is always space for further benefit to be gained, and this report details an investigation into applying a new optimisation to the CUDA code generation library subsection of OP2. The optimisation is named "Just-In-Time Compilation" for its similarities to a comparable process often performed by compilers when runtime efficiency is desired.

1.2 Motivations

The idea for this project was provided by my supervisor, Dr Gihan Mudalige - an Associate Professor in the University of Warwick Computer Science Department. I selected it as it aligned with my interest in High Performance Computing, and similar experience with optimising existing codes.

Since OP2 is Open Source and freely available, the implementation I produce will become part of the library, allowing future contributors to build on my work. The project also allows me the opportunity to operate on a large codebase, where most university work is done largely within the confines of one's own code.

2 Related Work

2.1 OP2

2.2 JIT

2.3 CUDA

3 Specification

The implementation will require work in two main areas: The python code generation script, and in the OP2 library itself which is implemented in both C and Fortran. Only the C library will be modified, due to developer familiarity.

The python script will perform source-to-source translation. As input it will take the application files, which specify the structure of the program: declaring variables, and indicating where the loops parallelised by OP2 should be executed; and a kernel descriptor for each loop, which will describe the operation to be performed on each element of the set passed to the loop as a parameter. OP2 makes an important restriction that the order in which elements are processed must not affect the

final result, to within the limits of finite precision floating-point arithmetic[2, p3]. This constraint allows the code generator freedom to not consider the ordering of iterations, and select an ordering based on performance.

From this a set of valid C files must be generated, to be compiled by a normal C compiler. In the case of this project the compiler will be the NVidia C Compiler (*nvcc*), as the code generated will be include CUDA.

It is important that the resulting executable compiled from the generated code produces outputs within some tolerance of the outputs generated by executing parallel loop iterations sequentially. Correctness is always a priority over performance for any compiler.

Furthermore, it will be ensured that the OP2 API will not be altered by any modifications to the library, to ensure that all existing programs using the API are able to seamlessly use the updated version.

3.1 Runtime Assertions

As discussed in Section 2.2 on Just-In-Time Compilation, the performance gain from this optimisation technique comes from making assertions at runtime which can only be made once the input is known. The application's input will be a mesh over which to operate, which includes a large amount of data, and opens up a number of runtime optimisations. The primary target for this project is "Constant Definition": turning values specified in the input as constants into `#define` directives for the C-Preprocessor. Other possible optimisations will be discussed in Section 7, Future Work.

3.2 System Model

Figure 3 describes the new workflow of the OP2 library, with the addition of Just-In-Time compilation. As before, code generation takes the application and loop files as

by the running executable, and the functions it provides are used instead of the unoptimised versions.

4 Implementation

5 Testing

6 Evaluation

7 Future Work

8 Conclusion

References

- [1] M.B. Giles G.R. Mudalige I. Reguly. *OP2 C++ User Manual*.
https://op-dsl.github.io/docs/OP2/OP2_Users_Guide.pdf. 2013.
- [2] 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*. 2012.

Appendices

A example

Acknowledgements

TODO