



**Create and Explore a graph from Agda
Definitions to analyse projects and speed up
compilation**

Dissertation

Julian Camilo Vidal Polania

School of Computer Science

College of Engineering and Physical Sciences

University of Birmingham

2024-25

Abstract

Usually 100-300 words stating the salient points of the report. It should help your reader to decide whether the report is relevant to her or his interests.

I will mention:

- How I got the s-expressions and how they were parsed
- The uses of the graph and the queries you can make
- How much it can speed up compilation and how

Acknowledgements

Abbreviations

ACB

Apple Banana Carrot

Contents

Abstract	ii
Acknowledgements	iii
Abbreviations	iv
List of Figures	vii
List of Tables	viii
1 Introduction	1
1.1 Background	1
1.2 Motivation	1
1.3 Problem Statement	2
1.4 Objectives	2
1.5 Project Structure	3
2 Literature Review	4
2.1 Related works	4
2.2 Parallel Compilation	5
2.3 Conclusion	5
3 Legal, Social, Ethical and Professional Issues	6
4 System Requirements	7
5 Design	8
6 Implementation	9
7 Testing and Success Measurement	10
8 Project Management	11

9	Evaluation	12
10	Conclusion	13
11	Appendix	16

List of Figures

List of Tables

CHAPTER 1

Introduction

1.1 Background

When you first encounter a massive code base, containing years of development from different people all adding their own code with their own style and thought process, it can be overwhelming and difficult to parse through. On most code-bases you can ignore the fundamental functions that underpin a given library, but in Agda it is quite useful to have knowledge of the underlying definitions of the project.

Agda is a functional programming language that can be used as a proof assistant, following the propositions as types logic system [2]. Since proofs in Agda can be created almost entirely from scratch, defining the most basic and fundamental types and operations it means that projects will quickly balloon in size. While most of the time it is possible to use higher level definitions when developing a proof, knowing the how the used types are fundamentally defined can be quite beneficial.

1.2 Motivation

Due to the size of these projects and how complex type definitions can be it becomes difficult to fully grasp what the relationships between all the definitions are. This creates the need for a tool that can analyze Agda projects and give the user an easy to use interface to query the relationships between definitions, in order to better understand the proofs. A tool that would allow the user to explore these relationships would be of great benefit.

Also, large Agda projects tend to deal with long compilation times. Agda is a proof assistant, during "compilation", it type checks the entire project which can take a 5 minutes or more depending on the computer and project. Agda will keep track of what modules have already been compiled, so working on high level modules isn't an issue as Agda will only type check files that have changed in

Add more
citations
to articles
and different
sources

the project and high level modules don't have dependents that would need to be re-checked. In contrast, low level files with many dependents, such as a module defining the natural numbers, will cause a significant portion of the project to re-compile. This effect will be particularly pronounced during refactoring where significant changes will be made to the bottom of the project, causing Agda to re-check the whole project making most of the time spent refactoring, spent waiting.

Agda type checks the project sequentially, going through the dependencies of each modules, this is done to maintain safety as a module won't be compiled until its dependencies have. A tool that analyzes the relationship between modules and finds modules that are safe to compile in parallel could lead to significant time savings.

1.3 Problem Statement

Agda doesn't provide a built-in tool to extract the definitions from a project, but there is an s-expression extractor Agda version by Andrej Bauer [3] that will allow a better view into Agda's internal representation. S-expressions aren't well suited data structure for the queries on definition relationships, mostly contain syntax information that can be discarded and each module has an s-expressions file. A data structure that lends itself well towards queries of relationship between definitions, that only contains information about the definitions and stores all the definitions from a project, such as a graph. A tool that could explore the dependency graph of the definitions, would allow for a faster understanding of a large code base.

While Agda doesn't have a tool to extract definitions from projects, it does have a tool to create a dependency graph from the modules. For compilation time improvement, only the modules are of interest as such an effective representation is already available. The challenge is what is the most effective way to traverse a dependency graph, concurrently, while maintaining safety (i.e. don't compile modules and its dependencies at the same time).

1.4 Objectives

This objective of this project is to create a dependency graph of definitions from an Agda project by using the s-expression representation of modules using the s-expression extractor [3]. Then create build a CLI that will create the dependency graph from any Agda project and give the user access to commands that query the dependency graph. These queries will let the user explore the graph, in order to gain insight into the Agda project. As there aren't currently any tools that allows you to query the definitions of an Agda project.

The objective is also to build another CLI tool that will create the module dependency graph from an Agda project and generate the makefile with the optimal order that modules should be compiled to reduce compile time.

1.5 Project Structure

First, there will be an explanation of what s-expressions are and how they are being read. These s-expressions use tags to name different symbols in a module, only the tags that contain definitions are important. Then, show how the s-expressions are structured, the way that definitions can be found and getting their dependencies. Using the information gained for the s-expressions, they will be turned into a graph that can be easily queried for relationships. Some example queries will be demonstrated, to underline their utility and implementation.

Secondly, by analyzing different algorithms for traversing a module dependency graph there will be an exploration on how parallelism can improve compile time of an Agda project. Results from different algorithms and projects will be shown.

CHAPTER 2

Literature Review

There are many ways to represent code using graphs, each with their own utility and purpose. There are Abstract Syntax Tree (AST), Control Flow Graph (CFG), Data Flow Graph (DFG), Program Dependence Graph (PDG), and Code Property Graph (CPG). These representations encode all the behaviour and properties of a project, they can have many uses such as for code vulnerability detection [8]. These representation also allows for static analysis, where code isn't ran but the structure is analyzed for software validation [7]. For this project, these graphs encode far more information that is require, as for the dependency graph only the relationship between definitions is of value.

Add litera-
ture about
code decay

2.1 Related works

The Language Server Protocol (LSP) is used by IDEs, such as Visual Studio Code and IntelliJ, to provide features like goto definition and goto references. The LSP sits between the code editor and the language server, it is the language server that analyzes the code structure to support the features [9]. Language server has the functionality to implement this project, but, they are meant to be used while editing code and the language server isn't made to be used with custom queries. Examples of LSP's are Jedi [12] and Agda's language server [1], while they contain the functionality needed they are only meant to work on the current open file, not a whole project.

Ctags [13] is a tool that indexes all the symbols in a project, this is helpful to get all the definitions from an Agda project. But it doesn't capture the relationship between the symbols it finds.

Graph Buddy is an interactive code dependency browsing and visualization tool [5]. It takes large Java codebases turns them into Semantic Code Graphs (SCG), and creates a visualization of this graph. This graphs shows dependencies between modules, classes and methods which helps the developer better understand the project and tackle the pervasive issue of code decay [4]. Graph

Add section
explaining
dependency
graph

Buddy is integrated as to an IDE as a plugin, where the user can seamlessly explore the visualization.

There is also a tool [6] that visualizes the call graph of a coding project, this allows for better developer experience. It works in real-time, while the user is editing the code the graph will automatically update to show the changes.

2.2 Parallel Compilation

Type-checking is a computationally expensive task that hinders the workflow of a developer, this has lead to work to parallelize type-checking algorithms or make them incremental. Parallel type-checking aims to type check different parts of a project at the same time, while incremental type-checking aims to allow the developer to type check a small change in the project without having to type-check the whole project again.

An example is the work by Newton et al [11] which seeks to parallelize type-checking with Haskell. Also, the work by Zwaan et al. [15] using scope graphs for incremental type-checking.

However, this project doesn't aim to optimize the type checking algorithm themselves, rather, find independent modules that can be type checked together. This aligns more close with the following paper that explores reducing FPGA compile time by changing from a monolithic compilation style to compiling separate blocks in parallel [14].

Fundamentally the compilation problem in this project is a scheduling problem, which is NP-complete[10]. This means that there many algorithms that attempt to tackle this problem by following different assumptions as shown by Yves Robert [10].

2.3 Conclusion

There are many tools that analyze and visualize the overall structure of programming projects. Allow for static code analysis, where software can be validated and developers can have an easier time exploring a project. However, most of the tools are not easy to query by the user, meant to be used while editing a file and are reserved towards more popular languages like Java. A tool that is able to read an entire Agda project and gives the user access to the underlying graph is still missing.

Slow compilers is a common problem, hurting developer experience. Due to the monolithic nature of compilers, parallelization becomes a route to follow when optimizing compilation time. While parallelization can be applied to the type-checking algorithm itself, this project looks to type check modules in parallel while the type checking algorithm remains the same. This problem is closer to a scheduling problem, where the goal is to find the optimal way to assign tasks to multiple machines to reduce completion time. There is no tool in Agda that attempts to apply an scheduling algorithm to the type-checking of modules, which could lead to significant speed ups.

Add section explaining DAGs

Add section explaining scheduling problem

Add section explaining s-expressions

CHAPTER 3

Legal, Social, Ethical and Professional Issues

The project uses several open-source libraries, such as NetworkX and sexpdata. It is important to comply with the licenses of these libraries, such as MIT and BSD to avoid legal issues. Without compliance this tool can't be made publicly available. This project does not handle sensitive user data, any future additions that might involve user data must comply with data protection regulations.

The project has potential for positive impact on developer productivity by giving tools to understand large codebases quickly and reducing down time during refactoring. It is important for this tool to be accessible and inclusive to a wide range of users. Since, the tool runs on a terminal the user can choose the interface that suits them best. It is also important to provide documentation and helpful error messages to help users unfamiliar with Python to use this tool.

Transparency and accountability are critical for a project that will analyze the personal and professional projects of developers. The tool must give accurate results when querying definitions and creation compilation order. Misleading output can introduce errors during development which is unethical and counter productive.

This project must follow the professional standards set by BCS Code of Conduct. These standards ensure that software is reliable, secure and ethical. The tool will evolve over time and it is critical for it to be well-tested, to avoid introducing errors.

CHAPTER 4

System Requirements

This section should detail your understanding of what you are planning to create. The section should aim to break down the overarching aims of the work into clear, measurable requirements that can be used in the evaluation of the project. This is why we often function of functional and non-functional requirements.

- Explain what the project should do and how it should do it
- Functional requirements
 - Create graph from agda project
 - Should allow for asked queries (the ones mentioned in the Mathodon thread)
 - Reduce compilation time by 10%?
- Non-Functional requirements
 - Installation shouldn't be difficult
 - Queries should be done in under 2 seconds
 - Creation of the graph shouldn't take more than 10 minutes

CHAPTER 5

Design

Communicating how you think about the composition of your system and how it works. You might detail the ways in which the overall system will be broken down into subsystems. Detail should then be provided on the design of each of these subsystems.

- A system that turns agda into s-expressions
- A system that reads s-expressions and turns it into a graph
- A system that queries that graph
- A way to store the graph for future use

CHAPTER 6

Implementation

This section should discuss how you went about developing a system that was consistent with your design to meet your stated requirements. The implementation of subsystems should be accurately documented, with any implementation difficulties being acknowledged. The Design and Implementation sections can be grouped in the Project Report, if these are tightly coupled. Likely omitted for the Project Proposal, though should mention your proposed implementation technologies somewhere.

- Explain the s-expressions extractor
- How they are loaded into python
- How they are stored for future use
- How the compilation uses graph dot

CHAPTER 7

Testing and Success Measurement

As well as documenting system testing, this section should also describe any unit testing or integration testing performed. If you are not familiar with unit, integration or system testing then it would be a good idea to investigate these notions and consider about how they relate to your project. This section might also detail any performance, reliability or usability testing performed, with quantification, i.e., numeric measurements, being used wherever possible. All those points are systems focused through. If you're doing something that's research focused or more of a social / analytical study then you need to think about how you'll measure success.

- Measure compilation time
- Measure time to build graph
- Measure time to make queries

CHAPTER 8

Project Management

What is the timeline for your project? What software development methodology will you use? Can you identify the major milestones? These types of questions are important for project planning and should be addressed directly.

- Gantt chart
- Weekly updates with supervisor
- Getting the s-expression extractor was a major milestone
- Using python instead of clojure was a significant improvement

CHAPTER 9

Evaluation

Usually you evaluate your project with regard to the functional and non-functional requirements you set out in the earlier chapter. This doesn't necessary mean that your project was successful but, if these requirements were appropriately specified, you it's likely that your project was successful. You might be reiterating some points from the Testing and Success Measurement chapter in your Project Report.

- Mention compilation time reduction, and improvements
- Mention how some queries take too long but were all implemented
- Mention how trees take some time but can be re-used instantaneously
- Creates graph from any agda project

CHAPTER 10

Conclusion

Generally speaking, section can take different forms - as a minimum you would normally provide a brief summary of your project work and a discussion of possible future work. You may also wish to reiterate the main outcomes of your project and give some idea of how you think the ideas dealt with by your project relate to real-world situations, etc.. For the Proposal, don't mention future work but do summarise your document.

- Compilation time could be further reduced
- A clever way of removing cycles could lead to faster queries
- Many useless nodes that are part of Agda's backend but not necessary for the user.
- Work on making the whole process more automatic and easily distributable

Bibliography

- [1] agda. Github - agda/agda-language-server, Dec 2024. Language Server for Agda.
- [2] Agda Developers. Agda documentation. <https://agda.readthedocs.io/en/latest/overview.html>.
- [3] Andrej Bauer. Agda S-expression Extractor. <https://github.com/andrejbauer/agda/tree/release-2.7.0.1-sexp?tab=readme-ov-file>.
- [4] A. Bandi, B. J. Williams, and E. B. Allen. Empirical evidence of code decay: A systematic mapping study. In *2013 20th Working Conference on Reverse Engineering (WCRE)*, pages 341–350, 2013.
- [5] K. Borowski, B. Balis, and T. Orzechowski. Graph buddy — an interactive code dependency browsing and visualization tool. In *2022 Working Conference on Software Visualization (VISSOFT)*, pages 152–156, 2022.
- [6] I.-A. CSÁSZÁR and R. R. SLAVESCU. Interactive call graph generation for software projects. In *2020 IEEE 16th International Conference on Intelligent Computer Communication and Processing (ICCP)*, pages 51–58, 2020.
- [7] R. Fairley. Tutorial: Static analysis and dynamic testing of computer software. *Computer*, 11(4):14–23, 1978.
- [8] L. T. Gia Lac, N. Cao Cuong, N. D. Hoang Son, V.-H. Pham, and P. T. Duy. An empirical study on the impact of graph representations for code vulnerability detection using graph learning. In *2024 International Conference on Multimedia Analysis and Pattern Recognition (MAPR)*, pages 1–6, Aug 2024.
- [9] N. Gunasinghe and N. Marcus. *Language Server Protocol and Implementation: Supporting Language-Smart Editing and Programming Tools*. Apress, 2021.

- [10] Y.-K. Kwok and I. Ahmad. Static scheduling algorithms for allocating directed task graphs to multiprocessors. *ACM Comput. Surv.*, 31(4):406–471, Dec. 1999.
- [11] R. R. Newton, O. S. Ağacan, P. Fogg, and S. Tobin-Hochstadt. Parallel type-checking with haskell using saturating lvars and stream generators. *SIGPLAN Not.*, 51(8), Feb. 2016.
- [12] pappasam. Github - pappasam/jedi-language-server, Dec 2024. A Python language server exclusively for Jedi.
- [13] universal ctags. Github - universal-ctags/ctags, Dec 2023.
- [14] Y. Xiao, D. Park, A. Butt, H. Giesen, Z. Han, R. Ding, N. Magnezi, R. Rubin, and A. DeHon. Reducing fpga compile time with separate compilation for fpga building blocks. In *2019 International Conference on Field-Programmable Technology (ICFPT)*, pages 153–161, 2019.
- [15] A. Zwaan, H. van Antwerpen, and E. Visser. Incremental type-checking for free: using scope graphs to derive incremental type-checkers. *Proc. ACM Program. Lang.*, 6(OOPSLA2), Oct. 2022.

CHAPTER 11

Appendix

Some information, for example program listings, is useful to include within the report for completeness, but would distract the reader from the flow of the discussion if it were included within the body of the document. Short extracts from major programs may be included to illustrate points but full program listings should only ever be placed within an appendix. Remember that the point of appendices is to make your report more readable. I don't expect to see apprentices in any Project Proposals.