

Dhruv C. Makwana

Exploring the structure of mathematical theories using graph databases

Computer Science Tripos, Part II

Trinity College

May 10, 2017

Proforma

Name:	Dhruv C. Makwana
College:	Trinity College
Project Title:	Exploring the structure of mathematical theories using graph databases
Examination:	Computer Science Tripos, Part II, 2016–2017
Word Count:	11,014
Project Originator:	Dr. Timothy G. Griffin
Supervisor:	Dr. Timothy G. Griffin

Original Aims of the Project

This project aims to (a) represent Coq libraries as Neo4j (graph) databases and (b) create a library of Neo4j tools with the goal of highlighting the structure of and relationship between proof-objects. It also aims to have more features but still perform comparably to existing tools which help a user understand a Coq library. As an extension, this project aims to work with the Odd Order Theorem (part of Mathematical Components Coq library).

Work Completed

This project represents Coq libraries as Neo4j (graph) databases and provides a library of Neo4j tools that highlight the structure of and relationship between proof-objects. It performs comparably to existing tools which help a user understand a Coq library. This project works with the Odd Order Theorem (part of Mathematical Components Coq library) extension task.

Special Difficulties

None.

Declaration

I, Dhruv C. Makwana of Trinity College, being a candidate for Part II of the Computer Science Tripos, hereby declare that this dissertation and the work described in it are my own work, unaided except as may be specified below, and that the dissertation does not contain material that has already been used to any substantial extent for a comparable purpose.

Signed: Dhruv C. Makwana

Date: 10/05/2017

Contents

1	Introduction	1
1.1	Motivation	2
1.2	A Database Approach	2
1.3	Aims of this Project	3
1.4	Sample of Coq	4
1.5	Summary	4
2	Preparation	5
2.1	Project Planning	6
2.2	Requirements Analysis	6
2.3	Technologies Used	7
2.4	Starting Point	7
2.5	Coq Proof-Assistant	8
2.6	Existing Tools for Coq	8
2.7	Neo4j	10
2.8	Existing Tools for Neo4j	12
2.9	Summary	16
3	Implementation	17
3.1	Coq object-files to CSV	18
3.2	CSV to Neo4j	22
3.3	Query Library	23
3.4	Project Related	24
3.5	Dead-ends	25
3.6	Summary	27
4	Evaluation	29

4.1	Features	30
4.2	Performance	31
4.3	Library of Queries	33
4.4	Summary	42
5	Conclusions	43
5.1	In Hindsight	44
5.2	Future Work	44
	Bibliography	45
	A Full Model	47
	B Timings	49
	C Sample Code	51
	D Project proposal	59

1 | Introduction

1.1	Motivation	2
1.2	A Database Approach	2
1.3	Aims of this Project	3
1.4	Sample of Coq	4
1.5	Summary	4

Coq is an interactive proof-assistant ([The Coq development team, 2004](#)); it allows users to formalise mathematical theories into machine-checked proofs. Coq libraries, which encode mathematical theories, are difficult to understand. In this dissertation, I will describe a new tool I implemented for Coq users that aimed to (a) represent Coq libraries as Neo4j (graph) databases and (b) provide a library of queries with the goal of highlighting the structure of and relationship between proof-objects. It achieves (b) by using network-analysis techniques usually associated with social-networks.

1.1 Motivation

Coq proof-scripts are notoriously difficult to understand. Not only do these proof-scripts encode a mathematical theory, which can be difficult to understand in and of itself, they serve as verbose instructions on how to create proof-term whose type matches a certain proposition, rather than a statement of *why* something is true. There are currently no tools which help with the challenge of gaining a *high-level* understanding of a large Coq library.

An example of a large Coq library is [Gonthier et al. 2013](#). It translates (into Coq) two dense books on the Feit-Thompson Odd Order Theorem ([Peterfalvi 2000](#), [Bender et al. 1994](#)), which contain thousands of proofs and definitions. It is part of a larger, general trend of formalising substantial bodies of mathematics and marks the first time that we have a new way representing mathematics, different from the usual combination of formulas and natural-language prose. Such a turning point provides an opportunity to explore the novel representations and analyses possible.

However, it is still the case that when a user approaches a large Coq library, they are left to understand and keep track of several aspects of the library (such as implicit assumptions, previously defined results and the types and conventions behind any notation) by themselves. There is very little opportunity to consider and compare different approaches for arriving at a result (i.e. number of assumptions, number of steps, some notion of the importance such as number of uses).

I obtained such an opportunity by using a query-based approach: by creating a tool that expresses Coq libraries as databases. Then I was able query and analyse them to gain insights like those mentioned above.

1.2 A Database Approach

Mathematical theories are highly interconnected structures of definitions and proofs. As such, relational or a document-oriented models are not adequate representations for answering questions such as ‘What depends on this lemma and how many such things are there?’ or ‘What are the components of this definition?’.

Even a static, graph-based approach on its own is insufficient. Simply outputting a dependency graph is ineffective for all but the smallest of libraries. Figure 1.1 shows one such dependency graph for a medium-sized Coq library.

Graph databases have been developed to deal with highly connected data sets and path-oriented queries. That is, graph databases are optimised for computing transitive-closure and related queries, which pose a huge challenge for traditional, relational databases. Due to the emergence of massive data sets thanks to large-scale social and advertising networks, new techniques for representing and analysing such data have been developed.

For this project, I used the Neo4j ([Neo4j](#)) graph database system (and its expres-

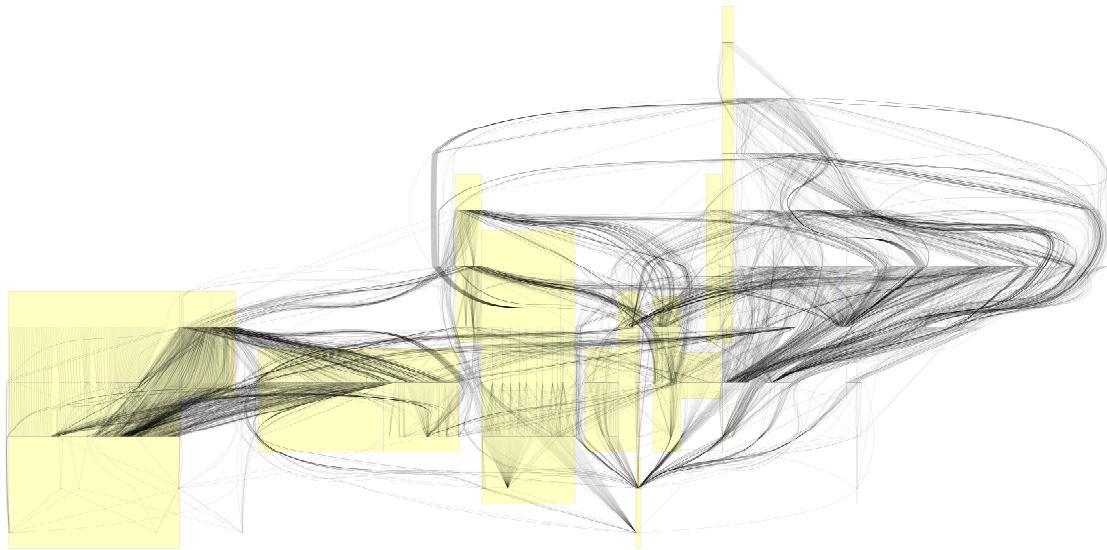


Figure 1.1 – Static graph of [CAS](#) Coq library. For all but the smallest of libraries, simply outputting a dependency graph is an ineffective way of understanding them. For example, if a user wanted to know whether or not X was used by Y, with a picture, they would have to visually trace backwards from Y to every node it depended on and check if X was one of them. In a relational database, such a query would not be possible with standard relational algebra (only with a highly-inefficient recursive-join); in a document database, a user would have to explicitly construct relations and joins at the application level. With graph databases, a user could simply enter a query (like the one in Table 4.3).

sive query language *Cypher*) to create a tool for Coq users which applies such network-analysis techniques to Coq libraries and help understand them.

1.3 Aims of this Project

In this project, I aimed to write a tool for Coq users that:

- represents Coq libraries as Neo4j graph databases. This involves
 - exploring and choosing the correct model (schema)
 - converting and extending existing code to output CSVs (comma-separated values)
 - writing new programs to extract extra information (omitted from other, existing tools)
 - writing new programs to automate database creation.
- provides a library of Neo4j queries, intended
 - to highlight the structure of and relationship between proof-objects
 - to provide several network-analysis techniques.

Note that the aims centre around the *mechanics and details* of creating a model for Coq proof-objects and applying network analysis techniques to it, *not creating a polished, user-friendly product*. Since *usability was not a focus of this project*, I did not conduct a user-study whilst evaluating this tool.

1.4 Sample of Coq

For readers unfamiliar with Coq, the following small example will clarify the starting point of this project.

```
(* A generic type-variable *)
Variable T : Set.

(* An inductive definition of a slist, similar to ML datatypes *)
Inductive slist : Set :=
| Empty : slist
| Sing : T -> slist
| Concat : slist -> slist -> slist.

(* Transforming slists into regular by pattern matching. *)
Fixpoint flatten (xs : slist) : list T :=
  match xs with
  | Empty => nil
  | Sing x => cons x nil
  | Concat ys zs => flatten ys ++ flatten zs
  end.

(* We prove flatten distributes over concatenation. *)
Theorem flatten_dist_concat : forall ls1 ls2 : slist,
  flatten (Concat ls1 ls2) = flatten ls1 ++ flatten ls2.
  reflexivity.
Qed.
```

1.5 Summary

First, I explained why Coq libraries are difficult to understand. Then, I said that the trend of formalising large bodies of mathematics such as the Feit-Thompson Odd Order Theorem (Peterfalvi 2000, Bender et al. 1994) presents an opportunity to develop new representations and analyses of mathematical theories. After that, I outlined why databases, specifically graph databases, were a promising approach to this new opportunity. Finally, I listed the aims of this project.

2 | Preparation

2.1	Project Planning	6
2.2	Requirements Analysis	6
2.3	Technologies Used	7
2.4	Starting Point	7
2.5	Coq Proof-Assistant	8
2.6	Existing Tools for Coq	8
2.6.1	Coqdoc	9
2.6.2	Coqdep	9
2.6.3	CoqSerAPI	9
2.6.4	dpdgraph	9
2.6.5	Comparison	10
2.7	Neo4j	10
2.8	Existing Tools for Neo4j	12
2.8.1	APOC: Awesome Procedures on Cypher	12
2.8.2	igraph	12
2.8.3	igraph Algorithms	13
2.8.4	visNetwork	14
2.8.5	R	15
2.9	Summary	16

In this chapter, I will describe how I planned this project. I will start by mentioning the development methodology and tools I used. Following that, I will present and explain the requirements I used to guide the implementation of my project. From there, I will state my choice of technologies for this project. I will then go on to elaborate my starting point for this project. This will consist of a brief commentary on Coq, a comparison of existing tools that aim to help a user understand a Coq library and a description of Neo4j and some of its relevant plugins.

2.1 Project Planning

This project had two distinct phases: deciding how to model Coq proof-objects (an exploratory phase) and applying network-analysis techniques (a technical phase). For both, I chose a spiral software development model: think of an idea, modify the code, propagate the necessary changes, evaluate the end-result and repeat. This allowed me to experiment with ideas flexibly and easily.

I found Git (git-scm.com) and GitHub (github.com/dc-mak) invaluable during this project; they allowed me to easily track, revert and review changes. I could safely test new features on different branches before merging them and store a copy of my work in multiple places. As it became apparent that precisely specified versioning, build-dependencies and tests were useful in spotting errors early, I added GitHub extensions such as [Travis-CI](#) for automated continuous-integration builds.

2.2 Requirements Analysis

I developed and brought together several distinct components for *modelling* and translating (from Coq to the chosen model), displaying and *interacting* (Neo4j/Cypher) and *computing* and analysing (plugins). Below is a list of required features I used throughout development to guide and provide context for implementation decisions.

- **Modelling:** The model should

M1 include as much relevant data as possible. Here, relevant means useful to understanding a large library, but not so much so as to obfuscate any information or make learning how to use the project more difficult.

M2 be flexible to work with and easy to translate. One could imagine different front-ends for interacting with and computing data from the model.

M3 strike a balance between size and pre-computing too much data. Figuring out which pieces of data can be reconstructed later and which are beneficial to compute during modelling will be a matter of experimentation and weighing up ease of implementation versus ease of later processing.

- **Interacting:** Interacting with the model should

I1 primarily, allow users to understand the data. The following two points follow from this principal goal.

I2 support both graphical and textual modes of use. Small queries and novice users are likely to benefit from the presence of a well-designed

GUI. However, larger queries requiring more computation and flexibility will benefit from a traditional, shell-like interface.

I3 be interactive and extensible. A static presentation of data, even in a GUI, would fail to make full use of graph-databases and the ability to query, in whatever way the user desires, information dynamically.

- **Computing:** Working with the model’s data should

C1 be enabled by a core library of good defaults. Certain, common functions should be ready ‘out-of-the-box’ and provide users all they need to get started.

C2 allow the user to add their own functions. It is not possible to imagine and implement all the functionality users may desire and so it would be very useful to have a way to extend this project to suit their needs.

2.3 Technologies Used

Choice of implementation languages was, although an important decision, almost completely dictated by the programs at the core of the project (Coq and Neo4j).

Coq and its plugins are written in OCaml (ocaml.org); I continued using it for two reasons. First, it is almost always wiser to work with and modify existing systems (for example, integrating with the build-system) and more representative of real-world work. Second, as a functional language, OCaml benefits from a strong, static (and inferred) type-system (which allows for easy experimentation and greater confidence in correctness). Hence, using OCaml for other parts of the tool which need not necessarily be in OCaml (for example, the `dpt2` utility) was a welcome and easy decision. OCaml has several other benefits too, such as inductively-defined datatypes (useful for working with Coq’s grammar) and good editing tools.

Similarly, Neo4j and its plugins are (usually) written in Java, but several languages are supported for the latter, both by Neo4j officially and by the community. As I will explain in Section 2.7, I found R to be the most suitable for achieving this project’s goals.

2.4 Starting Point

Both Coq and Neo4j have a rich ecosystem of tools and libraries built for them. I examined the current landscape in detail to determine the software available and see if I could use any of them as a starting point.

2.5 Coq Proof-Assistant

The Coq proof-assistant, implemented in OCaml, can be viewed as both a system of logic – in which case it is a realisation of the *Calculus of Inductive Constructions* – and as a *dependently-typed* programming language. Its expressiveness and automation facilities make it well-suited to *large-scale* developments.

On the logical side, Coq lays claim to feats of modern software projects such as the Four-Colour Theorem ([Gonthier, 2008](#)) (60,000 lines) and the aforementioned *Feit-Thompson* theorem (approximately 170,000 lines, 15,000 definitions and 4,200 theorems).

On the programming language side, Coq has served as the basis for many equally impressive projects. The *CompCert Verified C Compiler* ([Leroy, 2012](#)) demonstrates the practical applications of theorem-proving and dependently-typed programming by implementing and proving correct an optimising compiler for the C programming language. *DeepSpec* ([Pierce, 2016](#)), a recently announced meta-project, aims to integrate several large projects such as *CertiKOS* (operating system kernels), *Kami* (hardware), *Vellvm* (verifying LLVM) and many more in the hope of providing complete, *end-to-end* verification of real-world systems.

Coq is a *huge* project, developed by INRIA (France), and its size and complexity are best experienced through untangling the source code for oneself. Just for the implementation of the system (not including the standard library), Coq features approximately 3 major ASTs (Abstract Syntax Trees), 6 transformations between them, 3000 visible types, 9000 APIs and 521 implementation files containing 228,000 lines of dense, functional OCaml.

However, most of this massive project is sporadically (and tersely) documented. Even after I received some guidance via the Coq developers' mailing-list, I spent several hours browsing the source code trying to understand how the system worked. Although I had some prior familiarity with *using* Coq (as an introduction to tactical theorem-proving and dependently-typed programming), it was not useful for understanding the internals beyond context and how to compile and use programs and libraries. However, it did serve as invaluable insight for designing the model during the implementation phase.

2.6 Existing Tools for Coq

There are many tools for Coq that, like my project, *aim to help a user understand a library*. I studied several to learn their approaches and analyse their strengths and weakness. What follows is a detailed account of each tool and why it did not meet this project's aims and requirements.

2.6.1 Coqdoc

Coqdoc is a documentation tool for Coq projects, included as part of the Coq system. It can output to raw text, HTML, L^AT_EX and a few other formats to help a user navigate and understand a Coq library. Although it supports prettifying code with syntax highlighting and Unicode characters, its most relevant feature was its hyperlinking: potentially useful for building dependency graphs.

However, the whole tool works on an entirely *lexical* level, with no formal parsing or understanding of the code structure. Hence, since coqdoc could not meet any of the modelling requirements (completeness M1, flexibility M2 and size/pre-computation M3) I did not use it.

2.6.2 Coqdep

Coqdep is a utility included with Coq that helps users understand a Coq library's *module-level* dependencies by tracking **Require** and **Import** statements. Although on first impressions, this tool seemed to offer more flexibility than coqdoc, it was even more restrictive: it simply searches for keywords (such as **Require** or **Import** for Coq and **open** or dot-notation module usage for OCaml) per file and outputs them accordingly. As with coqdoc (and for the same reasons), I did not use coqdep either.

2.6.3 CoqSerAPI

Coq Serialized (S-expression) API (github.com/ejgallego/coq-serapi) is a new library and communication protocol aiming to make low-level interactions easier using OCaml datatypes and s-expressions, particularly for tool/IDE developers. It has a starting point for gathering some statistics on proof-objects in a project. While this is likely to be useful in the future, it is still far from complete and is more geared towards interactive *construction* (via a tool/IDE) rather than *analysis*. As such, tracking dependencies (critical to the modelling requirements) is not possible.

2.6.4 dpdgraph

dpdgraph (github.com/Karmaki/coq-dpdgraph) is a project which helps users understand the dependencies between proof-objects in a Coq library. It does so by extracting information from compiled Coq object-files to a .dpd file. It includes two example tools: **dpd2dot** (for producing a .dot file for static visualisation) and **dpdusage** (for finding unused definitions). Its developers intended it to be a starting point for other tools to build upon.

Although lots of information such as notation, the relationship between constructors and the types they construct, proof tactics, the precise kind of an object (for

example, fixpoint, class, lemma, theorem, etc.) and which module an object belongs to was missing, I thought it unlikely that the information was not present in the compiled object files (since they are necessary for term-construction and type-checking).

So, assuming the data was already present in those files, but simply *ignored or unused*, I focused on understanding and augmenting dpdgraph to add the missing pieces to the model and output it to CSVs.

2.6.5 Comparison

I have just discussed the many existing tools for Coq that *aim to help a user understand a library*. Table 2.1 summarises the main features of each. The features chosen reflect the strengths of each tool and are justified and elaborated upon below.

Bundled with Coq, coqdoc produces **hyperlinked source code**, meaning details such as **precise kinds** of a proof-object, **constructors of a type** and **type-signatures** are immediately visible; hence those five dimensions were included. Also included in a Coq system is coqdep: a tool for modelling **module-level dependencies** that can output a .dot file to present the dependencies **graphically**.

The Coq Serialised (s-expression) API is an **interactive** IDE communication protocol with facilities for there gathering of some basic **statistics**. Here, interactive is used to mean that information is not presented all at-once, *statically*, but can instead be queried dynamically at run-time. An example of a static display is Figure 1.1; it shows a medium-sized Coq library as output by dpdgraph.

In principle, coqdep and dpdgraph can also support some degree of interactivity, with support from other tools (which translate .dot files to interactive JavaScript), although this is rarely done. Lastly, dpdgraph models **object dependencies** well, with some scope for distinguishing precise kinds and displaying information graphically.

2.7 Neo4j

Neo4j is a graph database system implemented in Java. Traditional, relational database theory and systems are designed with the goal of storing and manipulating information in the form of *tables*. As such, working with highly interconnected data (such as a social-network graph) is best tackled with the alternative approach of *graph databases*.

Briefly, a (directed) *graph* is defined as $G = (V, E)$ where V is a set of vertices or *nodes* and $E \subseteq V \times V$ is a set of edges or *relations* between two nodes. A *graph database* is an OLTP (online transaction processing, meaning operated upon live, as data is processed) database management system with CRUD (create, read, update and delete) operations acting on a graph data model. Relations are there-

	Source Code	Hyperlinks	Precise Kinds	Constr. & Types	Type-Sig.	Module depend.	Graphical rep.	Interactivity	Statistics	Object depend.
Coqdoc	■	■	■	■	■	■	■	■	■	■
Coqdep	■	■	■	■	■	■	■	■	■	■
CoqSerAPI	■	■	■	■	■	■	■	■	■	■
dpgraph	■	■	■	■	■	■	■	■	■	■

■ Has feature ■ Does not have feature ■ Can be extended to support it

Table 2.1 – Comparison of Features. There are many tools for Coq that, like my project, *aim to help a user understand a library*. This table summarises the main features of each. The features chosen reflect the strengths of each tool considered. Detailed commentary is provided in Subsection 2.6.5.

fore promoted to first-class citizens and (like nodes) can also be manipulated and analysed.

Neo4j supports both graphical and textual modes of use and is easily extensible (through Cypher plugins and several language-specific bindings and libraries). It meets all the interaction requirements of helping users to understand data because it is flexible in its use and extensible in its capabilities. It even includes a tool to import CSV files containing nodes and edges into a new database. This allowed me to focus on extracting as much information as possible and expressing it in a simple format.

Neo4j also includes an interactive, graphical interface, accessible through an ordinary web-browser. As can be seen in Figure 2.1, the tool offers

- an overview of the current labels, relations and properties in the database
- interactive, syntax-highlighted input
- graphical representation of query result (with options to view it as rows like a shell, or raw JSON text results) with profiling information along the bottom
- easy access to favourite queries and scripts (the star on the left)
- easy access to documentation and system information (the book on the left)
- other features such as browser sync, settings and the ‘about’ section.

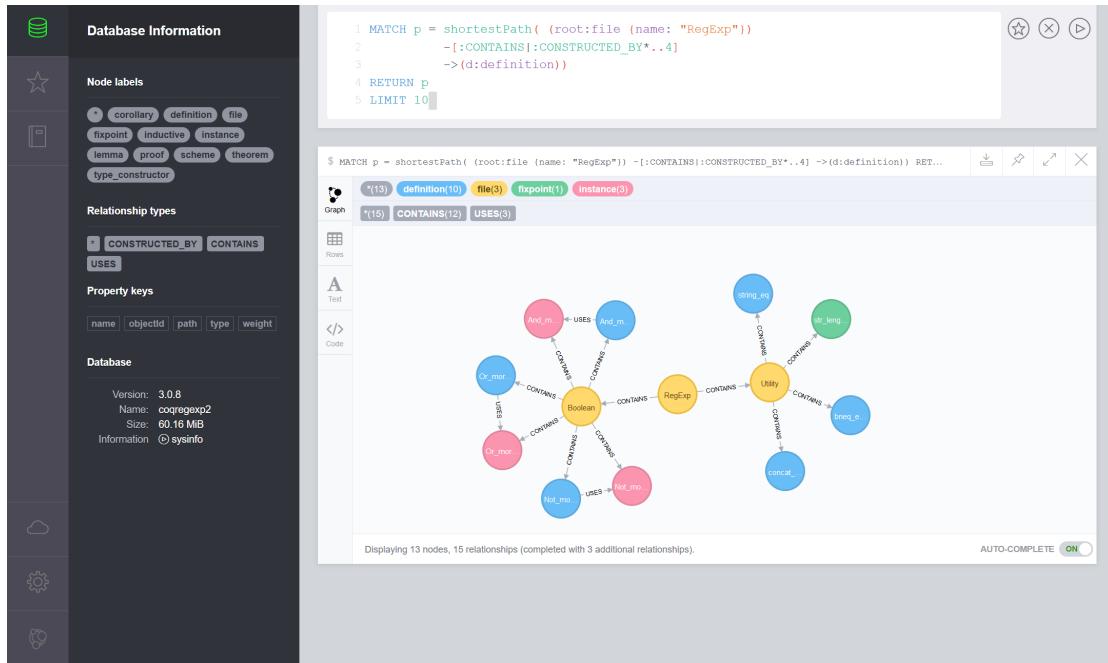


Figure 2.1 – Neo4j Interactive Browser. See Section 2.7 for a full list of features.

2.8 Existing Tools for Neo4j

Neo4j features rich integration with many languages, libraries and tools. Of those, I found the following to be the most relevant and useful tools for meeting this project requirements.

2.8.1 APOC: Awesome Procedures on Cypher

Awesome Procedures on Cypher, or *APOC* for short, is a community-maintained Java plugin featuring several network-analysis algorithms, callable from within Cypher itself. Although there are other Cypher extension libraries (such as MazeRunner), APOC is easy to install, well-documented, up-to-date and the most comprehensive, and therefore the obvious choice as a foundation.

Thus, APOC helps step towards meeting the *interaction* requirements for this project by being easy to understand, flexible to use and extensible; it even goes part-way towards meeting the *computation* requirements.

2.8.2 igraph

APOC's focus is on interacting with and combining different sorts and sources of data and so it lacks graph analysis functionality *beyond* the basics. The fact that it is implemented in Java adds to its limitations: it is not well-suited to more intense analyses over large graphs of libraries and is insufficient to *fully* meet the *computation* requirements of this project.

For such tasks, `igraph` is ideal: it is described on its website as a *collection of network analysis tools, with the emphasis on efficiency, portability and ease of use*. `igraph` offers lots of high-performance graph algorithms (it is written in C/C++, with bindings for R and Python). Some of these algorithms and their uses are described in Subsection 2.8.3.

Although `igraph` is not as easy to interact with as APOC, the extra capabilities it provided were indispensable towards achieving the *computation* requirements of a core library of good defaults.

2.8.3 `igraph` Algorithms

`igraph` offers a comprehensive set of network-analysis algorithms. Becoming familiar with these was a challenge and I spent several hours reading Newman's *Networks* (Newman, 2008) in order to understand and use them correctly. Network-analysis typically centres around two broad classes of measures: centrality and community detection; both of which are described next.

I Centrality

Centrality measures offer a way to characterise a node's *importance*.

Betweenness centrality is a measure based on shortest paths. For each node v , the fraction of shortest paths (from source s to target t , σ_{st}) which pass through it $\sigma_{st}(v)$ is its betweenness centrality. When applied to mathematical theories, how “unavoidable” a given object is for the results which mention it (Freeman, 1977).

$$g(v) = \sum_{s \neq v \neq t} \frac{\sigma_{st}(v)}{\sigma_{st}} \quad (2.1)$$

Closeness centrality is a measure also based on shortest paths. For each node, the sum of the length of all shortest paths to every other node is its closeness centrality. In a directed, dependency graph of mathematical theories, this corresponds to how “foundational” a node is. It is typically calculated as the reciprocal of *farness* $\sum_y \sigma_{xy}$, scaled by the number of nodes in the graph (N) to allow comparisons with graphs of different sizes (Bavelas, 1950).

$$C(x) = \frac{N}{\sum_y \sigma_{xy}} \quad (2.2)$$

PageRank is a variant of **eigenvector** centrality. For an adjacency matrix \mathbf{A} , the v^{th} component of the eigenvector \mathbf{x} (whose entries must all be non-negative, corresponding to the greatest eigenvalue λ) is the v^{th} node's *relative, eigenvector* centrality. Normalising the eigenvector provides the *absolute eigenvector* centralities (Newman, 2008).

The principle behind such a measure is that connections to high-ranking nodes contribute more to a node's rank than connections to low-ranking nodes. As such, for mathematical theories, it is like a *weighted in-degree/number-of-uses* which takes into account the importance of the nodes using a given type, proof or definition.

PageRank is similar; we want a vector \mathbf{r} that satisfies equation 2.3 (Page et al., 1999) instead of the equation $\mathbf{Ar} = \lambda\mathbf{r}$ (for number of nodes N , random-jump probability $1 - d$ and stochastic adjacency matrix \mathbf{L}). Hence, it can be viewed as the probability of a randomly-perusing mathematician coming across a given proof or definition on a first reading.

$$\mathbf{r} = \frac{(1 - d)}{N} \mathbf{1} + d\mathbf{Lr} \quad (2.3)$$

II Community Detection

Complex networks can exhibit community structure; that is, the graph can be (roughly) divided into sparsely-connected, dense groups. Although mathematical theories are often divided into sections, chapters and books, the following algorithms provide scope for re-evaluating these groupings.

Label propagation is a simple, near linear time method for determining which community a node belongs to. Each node starts with a unique label, after which, on successive iterations, it adopts the label held by most of its neighbours, until a consensus is reached. This whole procedure is repeated a few times and an aggregate result constitutes the output (Raghavan et al., 2007).

Edge betweenness – like betweenness centrality – is also based on shortest paths, with the idea that edges separating communities are likely to have high edge betweenness (since all shortest paths must pass through them). Successively removing the edge with the greatest betweenness value and re-computing over the remaining edges will result in a rooted tree, a hierarchical map (called a dendrogram) where the root represents the whole graph and the leaves represent individual nodes (Newman and Girvan, 2004).

Modularity is a measure of how well network can be divided. Formally, it is the fraction of edges that fall within a given grouping (across the whole graph) minus the expected number of those which could have fallen within the group by chance (and so is a real number between -0.5 and 1). Calculating this in an optimal manner is an NP-complete problem, and so I used a fast and greedy version of the algorithm instead (Clauset et al., 2004).

2.8.4 visNetwork

There exist *several* visualisation programs for Neo4j; however, many are for commercial, industrial use and offer the features/complexity (and pricing) to match.

All tools that offer live visualisation with built-in Cypher query execution (for example, KeyLines, TomSawyer, Linkurious) are proprietary, require a fee to use and offer more granularity than I needed. Offline (and open-source) solutions (which require data to be exported in some manner before visualisation) such as Gephi or Alchemy.js also offer many features, but at the cost of a steep learning curve.

Ultimately, I chose [visNetwork](#), which is an R library exporting to JavaScript because it was simple to use and easy to integrate with the rest of this project.

2.8.5 R

R is a statistics-oriented programming language, part of the Free Software Foundation's GNU project. It is relevant for this project because it offers an easy way to tie together Neo4j (through official bindings), igraph and visNetwork.

To take advantage of this convenience, I had to learn R for this project. This was not too difficult because R is a well-documented, relatively easy language to pick-up.

2.9 Summary

In this chapter, I gave a detailed account of how I planned this project. At the start, I mentioned the choice of development methodology (spiral) and the development tools I used (Git, GitHub and Travis-CI). Following that, I presented and explained the requirements this project should meet for modelling, interacting with and analysing Coq libraries. From there, I stated how the choice of technologies was dictated by the decision to use Coq and Neo4j.

I then elaborated on my starting point for this project. I gave a description of Coq, its uses and explained why (because of its complexity and poor documentation) it is a difficult system to work with. I also compared existing tools that aim to help a user understand a Coq library (coqdoc, coqdep, CoqSerAPI and dpdgraph) against this project's requirements and showed that although none satisfied all the requirements, dpdgraph provided a platform (albeit limited) to build upon.

Similarly, I gave a description of Neo4j and some of its plugins (APOC, igraph and visNetwork) and explained their features, advantages and disadvantages in relation to which of the project requirements they met. I discussed how I needed to learn the R programming language to tie together Neo4j, igraph and visNetwork in order to meet this project's interacting and computing requirements.

3 | Implementation

3.1	Coq object-files to CSV	18
3.1.1	Algorithm	18
3.1.2	Modelling	19
3.1.3	Translation	21
3.2	CSV to Neo4j	22
3.3	Query Library	23
3.3.1	Java Library	23
3.3.2	R Library	23
3.4	Project Related	24
3.4.1	Testing	24
3.4.2	Continuous-Integration Builds	24
3.4.3	Tooling	24
3.5	Dead-ends	25
3.5.1	Coqdoc	25
3.5.2	Coq source-files to CSV	25
3.6	Summary	27

In this chapter, I will describe how I implemented this project. What follows is an account of the programs I wrote, problems I encountered, solutions I implemented and tests I conducted. Figure 3.1 shows an overview of this project’s components and how they fit together.

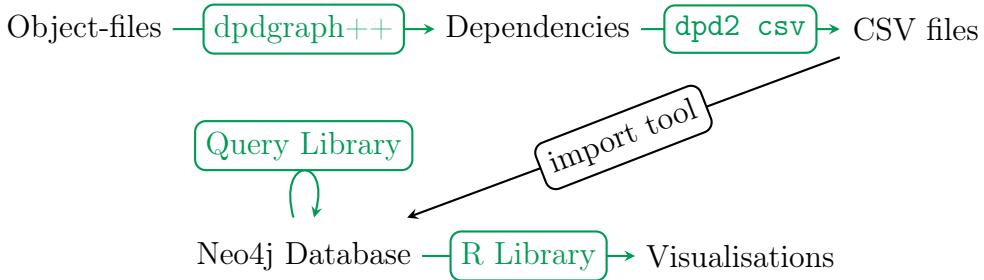


Figure 3.1 – Structure of Project. Nodes show the format of the information. Edges show the transformations. Edges in green show my contributions.

3.1 Coq object-files to CSV

This section of implementation corresponds to ‘dpdgraph++’ and ‘dpd2 csv’ on Figure 3.1: modelling the data contained in and the structure of a Coq library as CSVs. First, I will briefly describe the algorithm I used to construct a dependency graph constructed from a compiled Coq library. Then, I will elaborate upon how I modelled the Coq library by going through the attributes of each node and edge in the dependency graph. Finally, I will describe the format of the .dpd file output, which contains the graph constructed from the compiled Coq library.

3.1.1 Algorithm

This is a high-level overview of the algorithm I used in ‘dpdgraph++’ to construct a dependency graph from a compiled Coq library.

Input: A Coq script containing a list of files in the library (I wrote a shell-script to generate this file automatically) and the path to the compiled Coq library.

Output: A .dpd file containing a description of a graph whose nodes represent Coq proof-objects and whose edges represent various relations between the proof-objects.

1. For each module, for each proof-object in the module:
 - (a) if the proof-object is not in the set of nodes, add it.
 - (b) for each module in the chain of modules from the proof-object’s parent to the root of the library:
 - i. if the module is not in the set of nodes, add it.
 - ii. if there is not an edge between a module and its child in the set of edges, add it.
 - (c) collect the proof-object’s dependencies by recursing over its AST.
 - (d) for each dependency:
 - i. if the dependency is not in the set of nodes, add it.

- ii. if there is not an edge between the proof-object and its dependency in the set of edges, add it.
 - iii. otherwise, retrieve the edge and increment its *weight attribute* (representing the number of uses)
2. For each node in the set of nodes, output all its attributes (attributes described in Subsection 3.1.2; output format described in Subsection 3.1.3).
 3. For each edge in the set of edges, output all its attributes (source, destination and weight; output format described in Subsection 3.1.3).

3.1.2 Modelling

Although I used dpdgraph as a starting point, I rewrote all the attribute-assigning code because the meaning of the attributes it initially collected was not obvious and the terminology it used did not correspond to recognisable Coq constructs.

I added the following attributes and improved the dependency-collection for inductive types and their constructors (as I will describe in III Types and Constructors on page 20).

I Precise Kinds

A proof-object can be one of the following things, corresponding to an AST term: `module`, `class`, `type_constructor`, `inductive_type`, `definition`, `assumption` or `proof`.

Optionally, some terms have more precise terminology, for distinguishing different constructs. For example, when writing Coq, there is no `Proof` keyword; instead `Theorem`, `Lemma`, `Fact`, `Remark`, `Property`, `Proposition` and `Corollary` all are *synonyms* for proofs. Full details can be found in Appendix A.

So, to model this categorisation, I assigned each node a ‘kind’ *label*: labels are strings used to group nodes into subsets; since a node can belong to more than one subset, it can have more than one label assigned to it. To model the more precise distinctions, I optionally assigned some nodes a ‘subkind’ label.

II Recursive Modules

Every proof-object is contained within a module and every module is either contained in a file, is a file or is a directory. To model this inclusion relation, analyse module-level dependencies (like `coqdep`) and compare the results of community detection algorithms with the module structure of a library, I augmented ‘dpdgraph++’ to include modules in the graph. Using the Coq API, I could get the fully-qualified path (path from the library’s root) of a proof-object. Using a variant datatype, I expanded the type of a node in the graph to include modules and propagated the changes throughout the project.

However, because the Coq API returned module paths as strings, modules were in the model as a flat structure: modules could be related to objects but not to other modules. I fixed this by inferring (splitting the fully-qualified path) and adding all the “ancestors” of a module with the correct relations (parent as source, child as destination, repeatedly up to the root module).

III Types and Constructors

One of the most glaring omissions from dpdgraph’s initial model was the inability to relate a type to its constructor(s). To fix this, I improved the dependency-collection code (see Appendix C).

Expanding an AST term for type-constructors shows which type it constructed (note that types have no information about which constructors construct them). Since dependencies were constructed in a depth-first manner *down* the AST, I had to store the type and type-constructor relation reversed but output it in the correct order.

To do so, I compared each pair of nodes to see if it (a) was a type and a constructor and if so, (b) the constructor’s fully-qualified type matched the fully-qualified name of the type. If both these criteria were met, I swapped the direction of the edge output.

IV Types

Another unmissable omission from dpdgraph’s initial model was the inability to see the type of a proof-object. Type theory is central to a Coq user’s work and being able to include them in the model, would, along with kinds, subkinds and modules, help towards meeting the modelling requirement M1 of including as much relevant data as possible.

Coq’s type-checking algorithm is complex. To replicate it within ‘dpdgraph++’, I followed the functions called for the Coq command `Check <expression>` (for printing the type of a given expression). This led to the algorithm for *getting* the type, which I then implemented.

A subtlety I had to overcome was *using* the output in .dpd and CSV files. I replaced newlines, quotation marks, and commas with hash signs, single-quote marks and underscores respectively, because the former were used to delimit data in .dpd and CSV files and would have otherwise caused errors when they were being parsed into subsequent programs.

V Relations

Modelling relations was the most interesting aspect of deciding how to represent information. I wanted to keep consistent the notion of expanding a node to see more details: if a user is looking at an object, they should be able to expand the

object to see what the object depends on; if a user is looking at a module, they should be able to expand the module to see the objects contained within that module; if a user is looking at a type, they should be able to expand the type and see its constructors.

I modelled this idea with the following relations: `(src)-[:USES]->(dst)` for dependencies, `(type)-[:CONSTRUCTED_BY]->(constr)` for types, and `(module)-[:CONTAINS]->(object)` for modules.

I had two problems whilst implementing these relations. First was finding and matching types and constructors (details of which are described in III Types and Constructors on page 20).

Second was balancing expressiveness against simplicity. I considered relations of the following format, `X_USES_Y` for kinds X and Y. Although this was useful for fewer kinds, I decided its specificity when subkinds are included in the model made the model too large and complex (on the order of n^2 relations). Since Cypher allows pattern-matching and filtering based on kinds and subkinds anyway, I chose the simplified model presented above.

3.1.3 Translation

This subsection corresponds to ‘dpd2 csv’ on Figure 3.1. Once a model is constructed (in the form of a graph) by ‘dpdgraph++’ and output to a .dpd file, it is translated by the dpd2 tool to a CSV file for use by Neo4j’s import tool to create a database. I will now present an overview of the .dpd and CSV formats, as well as the dpd2 tool itself.

I dpd Format

‘dpdgraph++’ outputs the graph representing the model as a .dpd file with the following format for nodes (one per line):

`N: <id> "<name>" [<property>=<value>];`

for example (full type elided for brevity),

`N: 76 "matches" [type="... -> bool", subkind=fixpoint, kind=definition, path="RegExp.Definitions",];`

and likewise for edges:

`E: <src id> <dst id> [<property>=<value>];`

for example,

`E: 150 145 [type=CONTAINS, weight=1,];`

II CSV

Once it is output as a `.dpd` file, a model can be translated to various other formats. For example, `dpdgraph` includes a tool to output a `.dot` file (a format used extensively for *visualising* graphs by many tools) from a `.dpd` file. This is what I used to generate Figure 1.1.

However, for the purpose of this project, I translated the `.dpd` file (using the `dpd2` tool described below) to *two* CSV files: one for nodes and one for edges, for use with Neo4j’s import tools with the following headers.

```
objectId:ID(Object), name, kind:LABEL, subkind:LABEL, path, type
```

Here we see `name`, `path` and `type` declared as properties, `kind` and `subkind` declared as labels and the `objectId` field declared as a unique identifier (or in relational terms, a key) for the nodes (in this schema, called ‘Objects’) in the graph.

```
:START_ID(Object), :END_ID(Object), weight:int, :TYPE
```

Similarly, here we see relations (between ‘Objects’ as declared previously) named according to the value under the `:TYPE` column (for example, `CONTAINS`, `USES` or `CONSTRUCTED_BY`), each with an integer property, `weight`.

Using a CSV format makes adding extra properties, labels and relations very straightforward. This means it is easy to integrate other tools with this project and so extend this project with new features.

III dpd2 Tool

Initially, this tool started out as the `dpd2dot` utility (bundled with `dpdgraph`). I refactored it into a more general, `dpd2` tool which could accept the file-type (`.dot` or CSV) as a command-line argument.

Whilst using this tool, I discovered it had a bug. By default, `dpd2` attempts to remove reflexive and transitive dependencies (i.e. $a \rightarrow a$ and removing $a \rightarrow c$ if $a \rightarrow b$ and $b \rightarrow^+ c$) in a depth-first manner. For large graphs, doing so is stack-intensive, and thus causes a stack-overflow error. Since this “feature” was unnecessary (because it *removed* useful information) and appeared time-consuming to fix, I passed the `-keep-trans` flag on subsequent uses to avoid the issue altogether.

3.2 CSV to Neo4j

The ‘import tool’ in Figure 3.1 is a command-line program that is included with Neo4j. Given a target directory, a CSV file containing the nodes of the graph and a separate CSV file containing the edges of the graph, the import tool constructs a database. As explained in II CSV, the headers determine the labels, properties and IDs of nodes as well as the sources, destinations, properties and types of edges.

3.3 Query Library

This section of the implementation corresponds to the ‘Query Library’ and the ‘R Library’ shown in Figure 3.1. I will now outline the plugins I worked with to meet the interaction and computation requirements (page 6).

3.3.1 Java Library

APOC (Awesome Procedures on Cypher) provides many convenient procedures. Of note are:

- functions for constructing and examining a *meta-graph* showing which labels and relations are in the database and how they are connected. They allow a user to see an overview of the model.
- a few key graph algorithms, such as: node and path expansion, spanning tree, Dijkstra’s shortest paths, A* search, label propagation (for community detection) and centrality measures (betweenness, closeness and PageRank).
- functions for regular-expressions and mathematics, useful for selecting nodes based on statistics (for example, PageRank) or selecting nodes based on their path property.

Each of these can be called directly from within Cypher; for example, writing `CALL apoc.meta.graph()` would return the meta-graph of the database.

3.3.2 R Library

For the ‘R Library’, (to meet requirement C1 for a core set of good defaults) I wrote several example programs which automatically processed data from a new database, stored the information again for later use (to avoid re-computation) and output appropriate visualisations. Since processing could take on the order of minutes, I also provided status updates (informing the user of the tasks being executed), execution times for each task, as well as progress-bars where possible and relevant (for example, committing a transaction to the database). See Appendix C for details.

I Visualisation

There are many interesting ways to visualise the plethora of data that translating Coq libraries (representing mathematical theories) into graph databases produces.

Surprisingly, igraph was able to help with visualising graphs as well analysing them. Whilst visNetwork – with its own JavaScript, force-directed, physics rendering – produced more aesthetically pleasing results for smaller graphs, the webpages

it output for larger graphs took intolerably long to render inside a typical web-browser (Firefox/Opera). Thanks to an (experimental) integration with igraph (specifically, igraph’s layout mechanisms), I could pre-compute graph layouts in fast, native C/C++ *before* rendering graphs in a web-browser.

3.4 Project Related

During implementation, I learnt several skills and lessons about correct project management. Small things, such as grep-ing a code base or keeping track of time and a log of work done, proved to be useful. However, to ensure this project ran smoothly on a larger-scale, I focused more on the following areas.

3.4.1 Testing

For most of this project, I conducted testing manually by inspecting the output. Though this was tedious, it was the only way to do so when the model was undergoing continual development. Once I fixed the model, I shifted my focus to visualisation and used automated tests. These tests were particularly useful when I had to duplicate this project because for compatibility with libraries I was using to test this project (see Subsection 3.4.2).

I used some of my older Coq proof-scripts – problems I solved when learning Coq – to check output on a small scale (where I knew every single aspect of the library, and testing turnaround was quick). I used Coq’s Standard Library as a large-scale stress-test, to ensure all constructs were translated correctly. When I did find problems, I could usually trace them back using output from debug statements I had embedded into ‘dpdgraph++’.

3.4.2 Continuous-Integration Builds

Although this project never failed to build locally, there were occasionally problems when trying to build it on the project supervisor’s machine. The problem was compounded when it became apparent that smaller libraries (such as CoqRegExp and the solved problems) relied on a version of Coq (8.5.2) older than the one Mathematical Components (needed for this project’s moon-shot, the Odd Order Theorem) relied on (8.6). I set up Travis-CI for continuous-integration builds, made version dependencies precise and explicit, and removed much of the inconvenience and uncertainty surrounding builds on other machines.

3.4.3 Tooling

At first, I ran the project half on Windows and half on a Linux VM (virtual machine), using shared folders. I had a few reasons for this: I had already set

up Coq and Neo4j on Windows, both are easier to interact with in a graphical environment and starting up a VM just for some experimentation took too long.

Eventually, this set-up became confusing and time-consuming and I made the leap to running the project fully on a Linux VM. Nevertheless, I had serious issues when working on R integrations: running a Java database inside a Linux VM was insufferably slow and switching databases was not easy. I solved this problem by setting up SSH reverse-port-forwarding (to let R inside the Linux VM connect to a Neo4j instance running directly on Windows) for decent performance.

To be the most productive during longer sessions of work, I also set up editor integrations for OCaml, dramatically reducing the edit-compile cycle (especially for a strong, statically-typed programming language as OCaml). At this point, Coq's use of non-standard OCaml features, extensions and build-systems became particularly frustrating. For example, I spent a considerable time was spent untangling the Makefile inherited from dpdgraph to have cleaner, out-of-source builds and reduce the mess in the current working directory, all to no avail as I realised how complex of the build-system for Coq plugins is.

3.5 Dead-ends

I have now finished describing the whole project; all parts of Figure 3.1. I will now explain why I did not further pursue a couple of seemingly obvious, alternative strategies to modelling Coq libraries. Both centre around the same theme: analysing the Coq source-files directly instead of compiling them and then analysing the object-files.

3.5.1 Coqdoc

I tried to modify coqdoc's output into a useful format but this did not prove fruitful because the purely lexical tokenisation coqdoc does cannot infer or preserve as much information as full parsing.

3.5.2 Coq source-files to CSV

So, why not try parsing the source-files directly? This is an appealing idea, because with my current approach, there is still some information being lost during compilation: there is an *apparent* absence or duplication of some modules and a lack of notation and tactics.

The former arises from the use of *functors*: modules that take other modules as arguments (used to abstract over arguments, tactics, definitions and proofs). A concrete example (from the Coq Standard Library) is the theory of total orders, minimums and maximums, which is applied to naturals, integers and rational numbers (as well as any other ordered types). Such functors cannot be compiled

unless fully applied (explaining the *absence* of some modules); when fully applied, they essentially copy their *structure* into each instance (explaining the *duplication* of some module names and structures).

I Exploring Solutions

Several issues stop parsing from being a pragmatic solution, least of which are the size of the AST and complexity of parsing Coq files (Coq uses non-standard OCaml tools to implement an *extensible grammar*, making the whole situation quite ugly to work with).

1. *Modules and functors are represented identically* in the AST, with the former as a special case of the latter, making it very difficult to distinguish between the two on an AST-level. By far, this was the biggest roadblock.
2. *Module types, or signatures*, must be incorporated into the model for functors to make sense. Modules and signatures share a many-to-many relation: a signature can be satisfied by multiple modules and a module can satisfy multiple signatures. To express this correctly would require *signature-matching*, a notoriously difficult task (and the reason why few languages support ML-style modules).
3. Further issues involve resolving objects into a global namespace and knowing which compiler flags were given during compilation (to match physical directories to logical modules), all of which complicate matching and merging with the compiled proof-objects.

II Resolution

I made some progress in tackling these issues, but it took too long, so I could not justify spending more time on it. Functors are not used in many projects, (especially given the popularity and ease of use of *type-classes* for expressing generalisations).

Furthermore, including functors risked violating requirement M1 (by retaining information that could be difficult to understand or use). Additionally, given that names and structures of instantiated modules are duplicated, it could be possible to *reconstruct* generalisations *representing functors* once the database is created (hence respecting requirement M3 by not pre-computing too much data).

As such, extracting information directly from Coq source-files did not contribute to the overall project, but presents clear ways forward for future work.

3.6 Summary

I just presented an in-depth account of the programs I wrote (dpdgraph++, dpd2, Query & R libraries), problems I encountered (recursive modules, relating types and constructors, incorporating type-signatures, CSV translation, dependency and version tracking, impenetrable build-systems and parsing information directly from Coq source-files), solutions I implemented and tests I conducted. Throughout, I referenced the project requirements to justify important decisions.

4 | Evaluation

4.1 Features	30
4.2 Performance	31
4.2.1 Setup	31
4.2.2 Inefficiencies	31
4.2.3 Graph Analysis & Visualisation	33
4.3 Library of Queries	33
4.3.1 Small: CoqRegExp	33
4.3.2 Large: Odd Order Theorem	35
4.4 Summary	42

In this chapter, I will show that this project meets its aims (as listed in Section 1.3) by presenting evidence for this claim. I will do so by (a) comparing this project with existing tools that aim to help a user understand a Coq library and (b) providing sample output and explaining the insights they provide.

```

MATCH (a)-[:USES]->(b),
      (src:module)-[:CONTAINS]->(a),
      (dst:module)-[:CONTAINS]->(b)
WHERE src.objectId <> dst.objectId
CREATE UNIQUE (src)-[r:DEPENDS_ON]->(dst)
SET r.weight = coalesce(r.weight, 0) + 1
RETURN r

```

Listing 1 – Query to set Module Dependencies

4.1 Features

I talked through most of the first aim (how I represent Coq libraries as Neo4j databases) in the Implementation chapter. To assess the capabilities of this project, and thereby the suitability of my chosen model, I will compare all the programs listed in Subsection 2.6, Existing Tools for Coq with this project against the features listed there (features chosen to reflect strengths of each tool considered).

This project *can be extended* to support **linking to source code** by modifying either the model, database or JavaScript visualisations to link to the relevant webpages output by coqdoc. So, a user could switch between a graphical overview and a detailed inspection at will.

Whenever a node is visible, it can be expanded to see the nodes it depends on, so in that sense, it supports **hyperlinks**, though, unlike hyperlinks, such expansion is done in place, thus retaining the *context* of its use.

Thanks to the *kind* and *subkind* labels, this project supports **precise kinds**. Also, any **type** can, via the **CONSTRUCTED_BY** relation, be expanded to see its **constructors**.

Due to the type property, each object's **type-signature** is also visible. Crucially, it is a fully-expanded type-signature, making explicit any assumptions introduced (perhaps hundreds of lines prior or in a different file) into the environment.

Module dependencies are set with the query in Listing 1 using the **CONTAINS** relation. **Interactivity** is achieved through the Neo4j browser interface and the JavaScript visualisations; **statistics** are achieved through the library of queries; **graphical representations** by both.

An important limitation of CoqSerAPI is that its **statistics** are (at the time of writing) simply three counters, whereas this project offers many sophisticated graph metrics and the ability (through a queriable database) to gain *any* sort of information a user is interested in.

And finally, **object dependencies** are at the heart of this project: by querying a Neo4j graph database, a user can understand and manipulate this relation in a

much more flexible and scalable manner than any visualisation can manage.

Therefore, this project either supports, or can be extended to support, every feature supported by other tools. This project also supports additional features not in other tools: structural queries, network-analysis algorithms and interactive visualisations.

4.2 Performance

Now, I will compare this project’s execution time to most of the tools from the previous section. This project is slightly slower than other tools (Figure 4.1) for *translating* Coq libraries because it has more features and is more flexible than them. However, it is much faster for *analysing and visualising* than `dpd2dot`.

4.2.1 Setup

To evaluate timings, I used the Coq (8.6) Standard Library, due to its sheer size (564 modules, 5823 definitions, 23,892 proofs). For `coqdoc` and `coqdep`, I modified Coq’s Makefiles to measure execution time using bash’s `time` command. At the time of writing, CoqSerAPI’s statistics were not fully/usably implemented, so I did not include it in this comparison. For `dpdgraph`, I took separate measurements for outputting a `.dpd` file and converting that file to a `.dot` format. I took a similar approach for this project, so that the comparison was as fair as possible. Details of experimental setup can be found in Appendix B.

4.2.2 Inefficiencies

Setting up a graph database from scratch can take some time. Assuming a Coq project is already compiled, the following steps need to take place:

1. generating file with a list of all the modules to be examined (15 ms),
2. compile that file using the Coq compiler (12.6 s),
3. convert the output `.dpd` file to CSV files (72.7 s),
4. create a Neo4j database from those CSV files (12.8 s)

When steps 1 and 2 are taken on their own, we see that the changes to accommodate a more detailed model *only resulted in a 25% slowdown*, which is acceptable given this is a stress-test and the difference is of the order of seconds. Creating a database from CSV files takes a similar amount of time, also acceptable given the size of the graph (31,088 nodes and 850,434 edges).

So, the real bottleneck is step 3: converting `.dpd` files to CSV. During execution, `dpd2` reads in a 25MB `.dpd` file and outputs two CSV files of size 7MB (nodes)

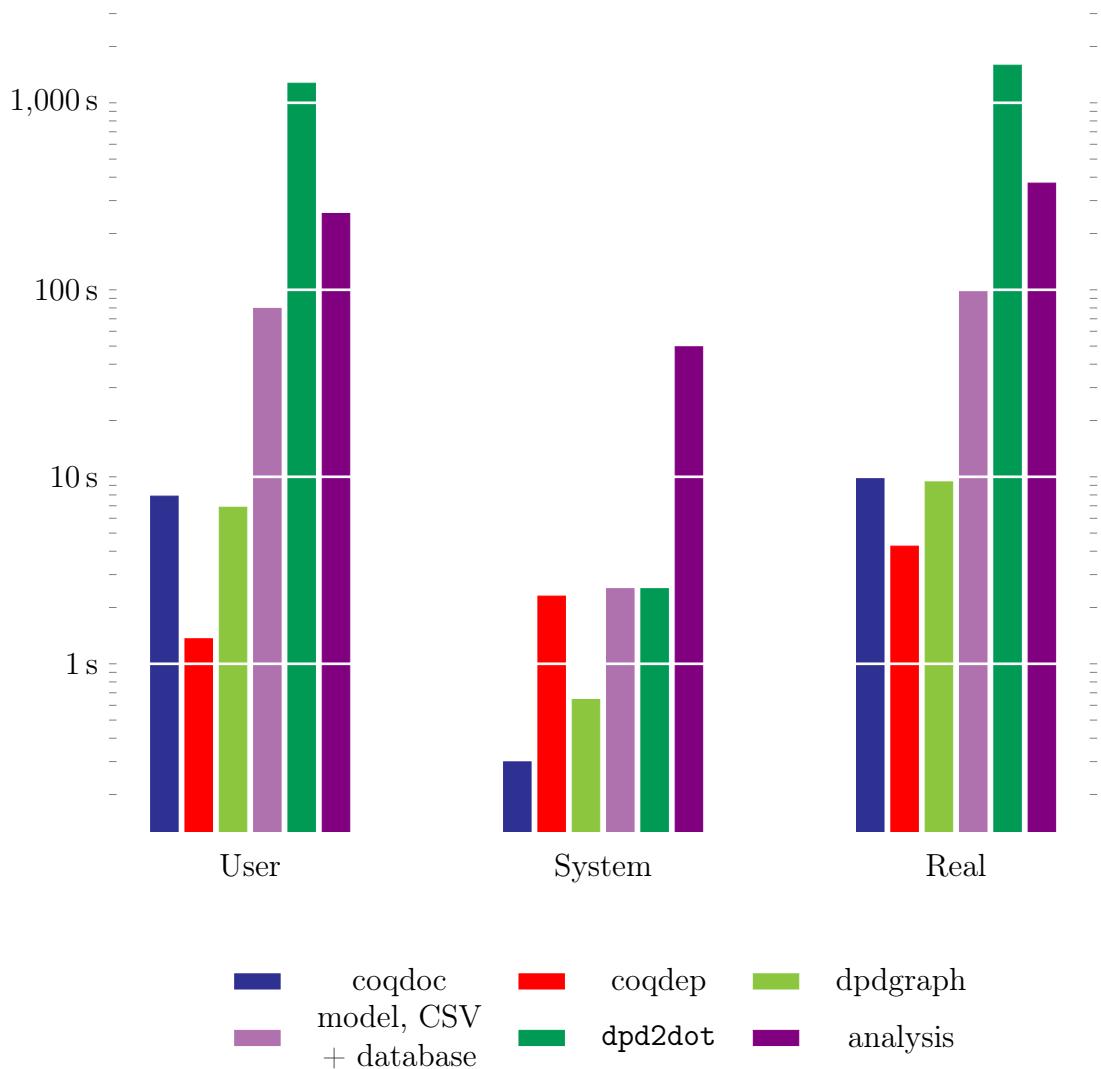


Figure 4.1 – Comparison of Execution Times. Note that the data is presented on a *logarithmic* scale. Details of experimental setup can be found in Appendix B. We see coqdoc takes very little time to run, and coqdep even less (which is to be expected considering their purely lexical approach). Somewhat surprisingly, dpdgraph runs just as quickly as coqdoc; its increase in system time can be explained by the 14MB .dpd file output by dpdgraph. There is an order-of-magnitude slowdown with this project when *translating* a Coq library (98s versus 9.4s); a more detailed examination in Subsection 4.2.2 explains precisely *what* occurs during database creation and where inefficiencies lie. However, it is *much* more efficient to analyse and produce visualisations for a Coq library once it has been translated into a graph database (6 minutes versus 27 minutes)

and 8MB (edges), so IO is likely to be a factor, as is reconstructing the graph in memory. I output the model to a `.dpd` to make it easier to extend this project with other tools. It is likely that outputting a CSV directly would have resulted in being able to bypass this phase altogether, but from a software-engineering point-of-view, the trade-off there is increased coupling for faster execution.

4.2.3 Graph Analysis & Visualisation

Once a graph is created, the last step is to *use* the data by analysing and visualising it. Here, this project shows a significant improvement over `dpd2dot`.

`dpd2dot`'s execution time dwarfs the analyses carried out by this project. Whereas `dpd2dot` took *27 minutes*, to convert a 13MB `.dpd` file to a 24MB dot file, in about one-quarter of the time, an R script ran (a) PageRank and closeness centrality algorithms over all proofs and definitions in the database and (b) output *8* different 9MB visualisations of the data. Analyses took less than a minute; visualisations took from 20 to 90 seconds each.

It should be noted that `dpd2dot` does not do graph *layout*: it just splits the graph into sub-graphs (based on modules) and assigns a colour and a label to each node (based on their properties). Converting the `.dot` file to a viewable format (e.g. a scalable vector-graphic or SVG) is up to another tool (that being said, I cancelled the command `dot -Tsvg` to produce an SVG after it failed terminate within a few *hours*).

4.3 Library of Queries

I talked through how I provide several-network analysis techniques in the Implementation chapter. To assess whether this tool highlights the structure of and relationship between proof-objects (the second aim), I will now show the output of the library of queries on the small case of a Coq Regular-Expression library and on the large case of the project's moon-shot, the Odd Order Theorem.

All visualisations (except Figure 4.2) show only definitions (shown as triangles, reminiscent of the \triangleq symbol sometimes used for definitions) and proofs (shown as squares, reminiscent of the end-of-proof \square symbol), except for Figures 4.4, 4.6 and 4.7 which also include modules (as circles).

4.3.1 Small: CoqRegExp

I will now show, how, without studying any code, a Coq user can use this project to understand the structure of the CoqRegExp library.

```
MATCH (obj) RETURN LABELS(obj), count(*) AS total
ORDER BY total DESC LIMIT 5
```

LABELS(obj)	total
proof, lemma	79
definition	17
proof, theorem	16
definition, instance	14
type_constructor	8

Table 4.1 – Top 5, most common kinds of proof-objects in CoqRegExp with their frequency and the query used to obtain them above the table.

```
MATCH (n) WITH collect(n) AS nodes
CALL apoc.algo.pageRank(nodes) YIELD node, score
RETURN node.name, LABELS(node), node.path, score
ORDER BY score DESC LIMIT 5
```

node.name	node.path	LABELS(node)	score
RegExp	RegExp.Definitions	inductive_type, inductive	7.12518
matches	RegExp.Definitions	definition, fixpoint	3.12445
Or	RegExp.Definitions	type_constructor	2.35997
re_eq	RegExp.Definitions	definition	2.30212
Cat	RegExp.Definitions	type_constructor	2.04202

Table 4.2 – Top 5 proof-objects by PageRank in CoqRegExp, with the modules they are in, their kinds, their PageRank values, and the query used to obtain them above the table.

I Neo4j & APOC

As visible by the examples in Tables 4.1 and 4.2, and Figure 4.2, a user can make arbitrary queries on the CoqRegExp library database. The first gives an indication of the most common kinds of proof-objects: mostly lemmas, with definitions and proofs following. The second uses an APOC procedure for PageRank over all proof-objects to quantify the importance of each. The definition of a regular-expression and the definition of what it means for a regular-expression to match a string top the table, with other, fundamental theory components/proof-objects following. The third also uses an APOC procedure to produce a visual representation of the schema of the database.

II Visualisation

Intuitively, in Figures 4.3 and 4.4, the light blue cluster represents *executable functions*; the orange cluster represents *proofs of correctness*; the light purple

```
CALL apoc.meta.graph
```

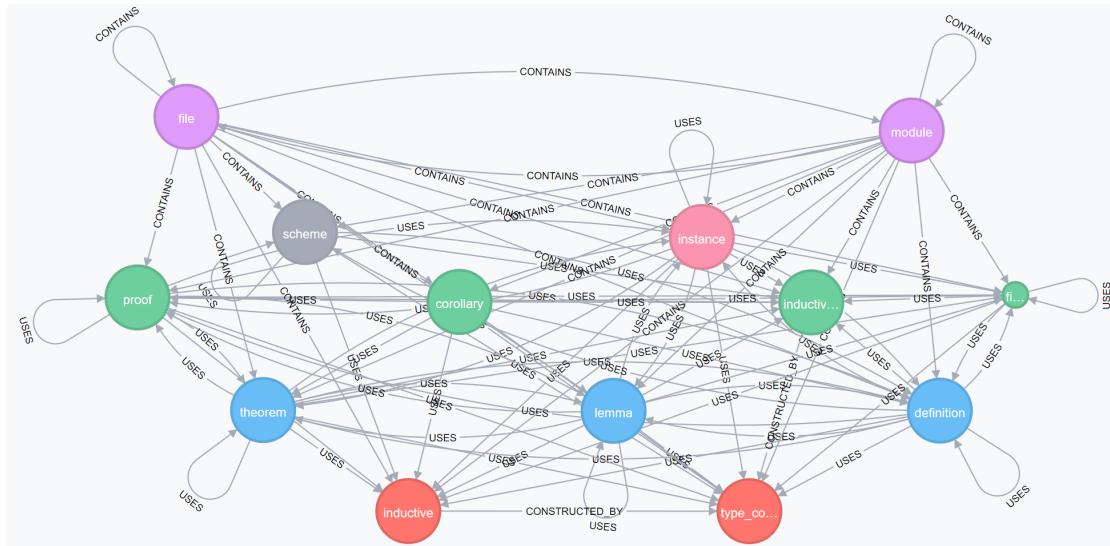


Figure 4.2 – Meta-graph of CoqRegEx representing the labels (kinds/subkinds) and relations (USES, CONTAINS, CONSTRUCTED_BY) between them. Above the figure is the Cypher query which calls the APOC procedure to produce this. At the bottom, red nodes represent types and constructors; above them, blue nodes represent definitions; above them, green nodes represent proofs. The grey node represents a scheme; the pink node represents an instance and the two purple nodes at the top represent files and modules.

cluster represents proofs using the definition of `string length`; the light green cluster represents proofs involving `converting strings to regular-expressions`; the violet cluster represents proofs related to `nullable strings`.

The node at the centre of the light blue cluster is a function which computes whether a given regular-expression matches a given string; the node at the centre of the orange cluster is a definition which defines what it means for two regular-expressions to be equal. These nodes represent two proof-objects that have high PageRank values, ‘matches’ and ‘re_eq’ respectively.

4.3.2 Large: Odd Order Theorem

This Coq library closely follows the structure of the source material it encodes: Peterfalvi (Peterfalvi, 2000) and Bender & Glauberman (Bender et al., 1994). Each section (chapter) in the original books is a file/module in the Coq library; each definition/lemma/corollary/etc. corresponds to the same in the books. Following the convention in the Coq library, ‘Bender & Glauberman’ will henceforth be abbreviated to ‘BG’ and ‘Peterfalvi’ to ‘PF’. For brevity, ‘Odd Order Theorem’ will also be abbreviated, to ‘OOT’.

I will now show, how, without studying any code, or reading BG or PF, a Coq user can use this project to understand the structure of BG and PF.

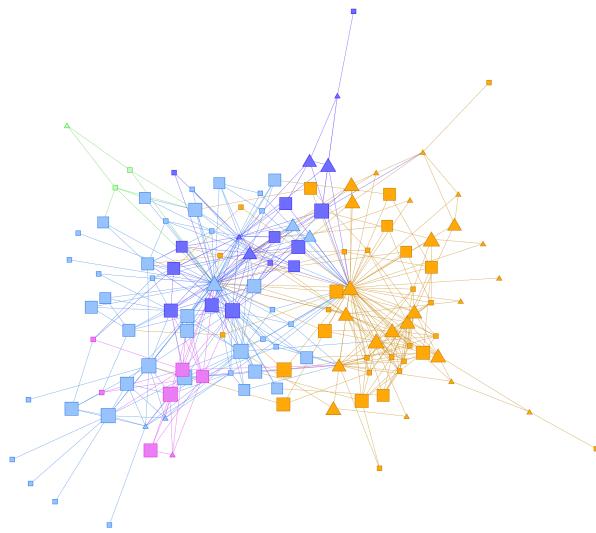


Figure 4.3 – Force-directed visualisation of definitions and proofs in CoqRegExp. Colours assigned by modularity clustering largely correspond to the way a human might group nodes visually: two major clusters with a few, smaller clusters (see II Visualisation for an interpretation of what each colour corresponds to). The size of the nodes corresponds to betweenness centrality scores (split up into 10 logarithmically equal-width buckets). Edges represent the USES relation; edge-directions (source-uses-destination) are omitted for clarity since they generally point toward the centre of a cluster.

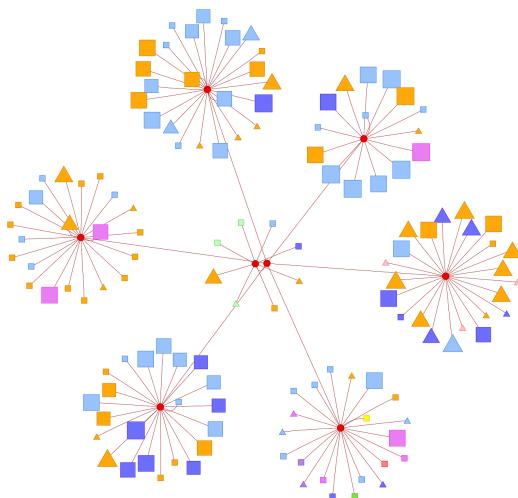


Figure 4.4 – Force-directed visualisation of definitions, proofs and modules in CoqRegExp. Setup is like Figure 4.3, except edges represent the CONTAINS relation. Each module tends to have two main parts: those relating to **executable functions** those relating to **proofs of correctness**, occasionally accompanied by a few definitions/proofs on **nullability** or **string length**.

```

MATCH (q:proof)
WHERE NOT (:proof {path: "mathcomp.odd_order.PFsection14",
                     name: "Feit_Thompson"})
      - [:USES*] ->(q:proof)
RETURN q.name, replace(q.path, "mathcomp.odd_order.", "") AS path

```

q.name	path
pcore_Fcore	BGsection15
main	stripped_odd_order_theorem...
ell_sigma_leq_2	BGsection14
Ptype_trans	BGsection14
P1type_trans	BGsection14

Table 4.3 – Five (of 89), proofs in the OOT Coq library that do not ultimately lead to the proof of the Feit-Thompson OOT. Module names (the ‘path’ property) have been shortened to remove redundant information. The ‘stripped_odd_order_theorem’ is a self-contained proof relying on only basic Coq features and is not part of BG or PF.

I Neo4j

By modelling Coq libraries as Neo4j graph databases, we can answer powerful questions. For example: there are 1064 proofs in BG and PF, *does every one of them ultimately lead to the proof of the Feit-Thompson OOT?* In general, how many ‘dead-ends’ – proofs, definitions or types – are there in these two books, and where are they? These questions have been answered in Tables 4.3 and 4.4 respectively.

II Visualisation

New to this section, edges are coloured. Unless a figure is stated as being ‘flipped’, the colour of an edge is the same colour as its source (user); otherwise, it is the same colour as its destination (used).

In Figure 4.5, we see that modularity clustering over proofs and definitions connected by the USES relation distinguishes *more* groups than force-directed visualisation does, despite appearing to be a highly interconnected theory.

Figure 4.6 shows visually what one might expect intuitively: *proofs and definitions within the same module tend to belong to the same group*. It explains why modularity clustering could distinguish more groups than force-directed visualisation could. However, it does not explain why the groups *cross* module boundaries and represent *multiple* modules.

Figure 4.7 explains why the groups cross module boundaries: they demarcate different parts/phases of BG and PF. Both the figure and the list below show this Coq library is not quite the linear chain of dependencies ranges one may expect.

```

MATCH (q) WHERE NOT((q:module) OR ()-[:USES]->(q))
RETURN LABELS(q),
       REPLACE(q.path, "mathcomp.odd_order.", "") AS path,
       COUNT(*) AS total
ORDER BY total DESC LIMIT 7
  
```

LABELS(q)	path	total
definition, scheme	stripped_odd_order_theorem	27
proof, lemma	BGsection16	6
proof, lemma	BGsection15	5
proof, lemma	PFsection8	5
proof, lemma	BGsection10	4
proof, remark	BGsection14	4
proof, lemma	PFsection5	4

Table 4.4 – Top 7 kinds of proof-objects in the OOT Coq library which are never used again (of which there are 107), grouped by module and ordered by frequency. Module names (the ‘path’ property) have been shortened to remove redundant information. The ‘stripped_odd_order_theorem’ is a self-contained proof of the entire OOT relying only on basic Coq features and is not part of BG or PF. Many unused results are simply called ‘lemma’ instead of the more descriptive and conventional ‘corollary’ or ‘remark’.

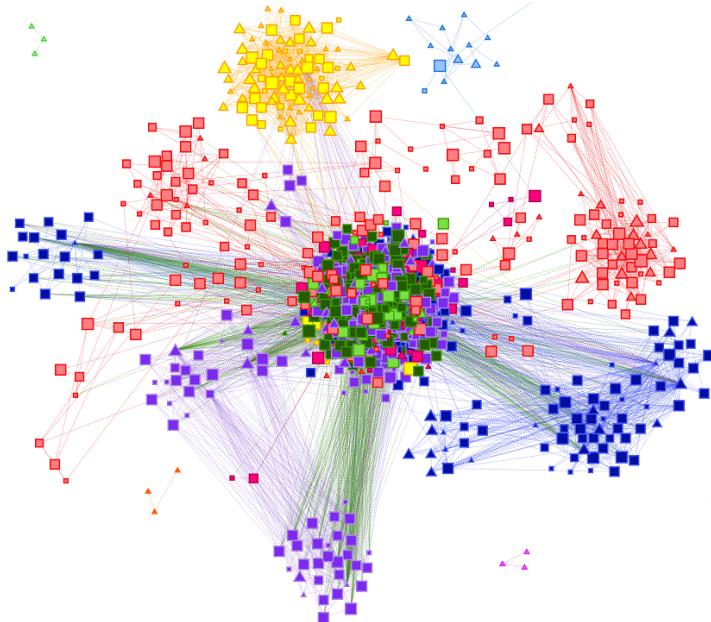


Figure 4.5 – Force-directed visualisation of definitions and proofs in the OOT Coq library (some nodes omitted for clarity). Setup is the same as that of Figure 4.3. Observe that modularity clustering can distinguish between nodes in the centre as belonging to different groups. This distinction is more evident in Figures 4.6, 4.7, 4.8 and 4.9; it is explained in II Visualisation.

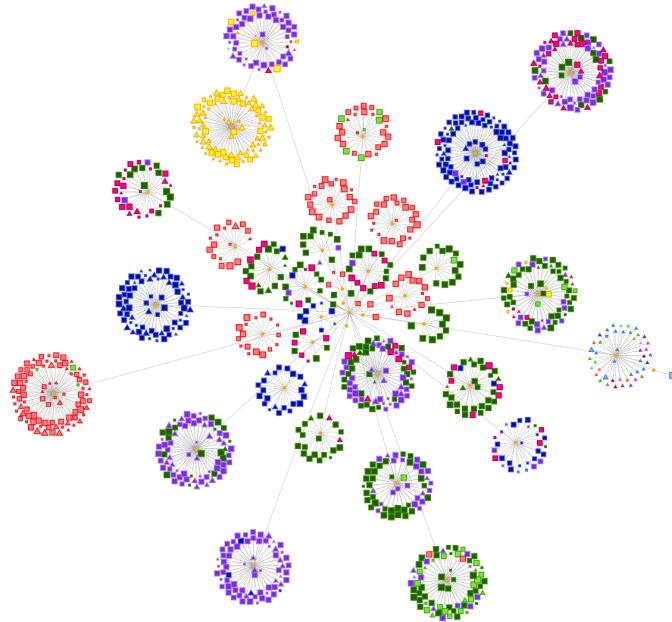


Figure 4.6 – Force-directed visualisation of definitions, proofs and modules in the OOT Coq library. Setup is the same as that of Figure 4.4. The module in yellow – a self-contained result by the name of ‘CyclicTlisoReflexion’ – is the only one *nested* inside another (PF 3). Unlike CoqRegExp, here, proofs and definitions in each module tend to belong to the same group. Note that groupings *cross* module boundaries; this is elaborated upon in II Visualisation.

Instead, the overlapping ranges of these groups suggest that the material is heavily interconnected and that there are other ways of approaching the material other than the linear presentation of BG and PF.

- pink: BG 1-6 and Appendices A/B/C
- dark green: BG 7-13 and PF 9-10/12-14
- purple: BG 14/16 and PF 3-4/8
- yellow: CyclicTlisoReflexion in PF 3
- magenta: BG 15 and PF 11
- dark blue: PF 1-2/5-7.
- light blue: self-contained proof of the OOT.

Figures 4.8 and 4.9 show two extremes of how one may approach the material.

1. Figure 4.8 shows the hierarchical, Sugiyama layout. A consequence of this algorithm is that nodes are placed closest to their first use, usually done as a heuristic to minimise edge crossings. Intuitively, one can think of this approach as the lazy student who only studies a particular definition/proof *just before* it is needed by some other part of the text (or like call-by-need evaluation in a non-strict programming language).

2. Figure 4.9 shows that swapping the direction of the edges does not simply flip the layout. In this vertically reflected image, nodes are placed *lower/closest to the nodes using them*. Intuitively, one can think of this approach as the eager student who studies *any* definitions/proofs that will be required by *any* later result *as early as possible* (or like call-by-value evaluation in a programming language).

One could conjecture that Figure 4.8 shows how mathematics is developed (a narrow foundation, top-down refinement, creating proofs/definitions as and when needed, several independent lines of thought) and Figure 4.9 shows how mathematics is presented (a broad base of all the groundwork that will be needed first, followed by a bottom-up, rapidly developing, linear argument which uses the base extensively). The similarities between Figures 4.7 and 4.9 are evidence for this interpretation.

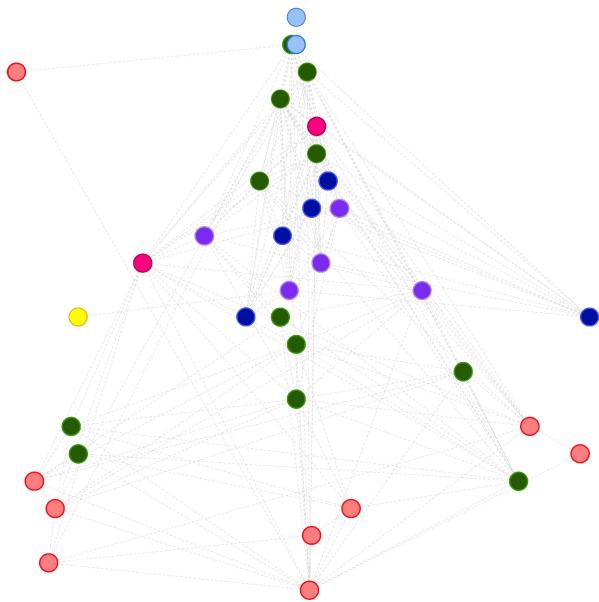


Figure 4.7 – Sugiyama (hierarchical) layout (the direction of edges always points downwards, above DEPENDS_ON below) of **modules** in the OOT Coq library. Modules are coloured according the most frequent (mode) colour of all definitions and proofs in them. Edges represent the DEPENDS_ON relation between modules defined in Listing 1. This visually represents the chapter dependencies in BG and PF, starting with BG 1 at the bottom and ending with PF 14 (containing the actual proof of the Feit-Thompson OOT) at the top.

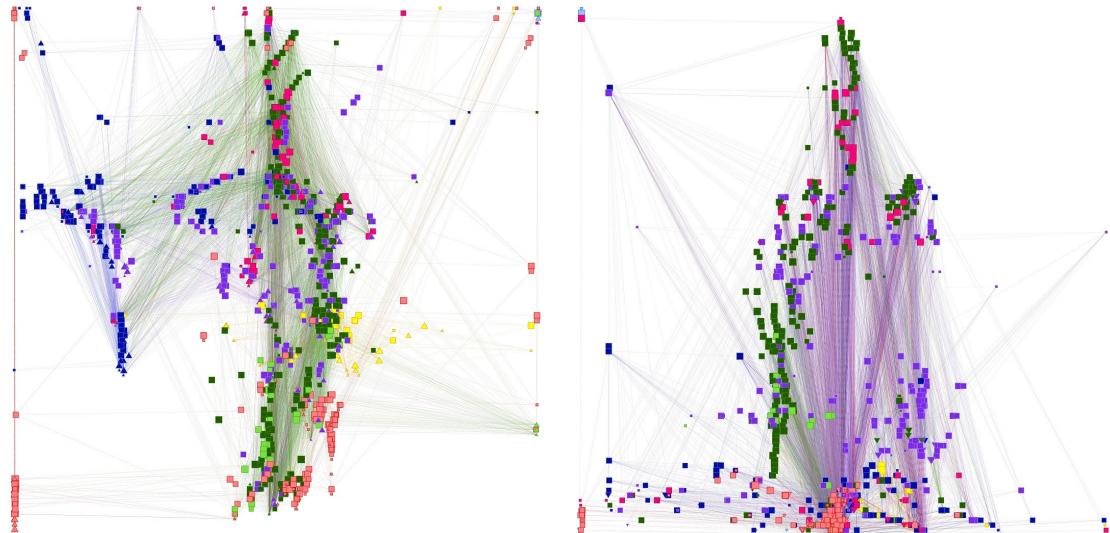


Figure 4.8 – Sugiyama (hierarchical) layout (the direction of edges point downwards) of **definitions and proofs** in the OOT Coq library. Colours are assigned by modularity clustering. Edges represent the USES relation: above USES below. A by-product of this layout is that nodes are placed higher/closest to their *users*.

Figure 4.9 – Setup is the same as Figure 4.9, except that (a) the edge directions are *reversed* and (b) the image is reflected vertically. It is still the case that above USES below. However, doing this has the effect of placing a node lower/closest to the nodes *using* it. An interpretation of this effect is provided in II Visualisation.

4.4 Summary

In this chapter, I evaluated how well my project achieved its two major aims (Section 1.3).

1. To evaluate the first aim of choosing the correct model (as part of representing Coq libraries as Neo4j graph databases), I compared this project’s features and performance against other, existing tools that aim to help a user understand a Coq library.
 - I showed that this project either supports, or can be extended to support, every feature supported by other tools. This project also supports additional features not in other tools: structural queries, network-analysis algorithms and interactive visualisations.
 - This project is slightly slower than other tools (Figure 4.1) for *translating* Coq libraries because it has more features and is more flexible than them. However, it is much faster for *analysing and visualising* than `dpd2dot`.
2. To evaluate the second aim of analysing and highlighting the structure of Coq libraries, I examined this project’s output on two examples: CoqRegExp (a small case) and OOT (a large case). It is not possible to gain the following insights using other tools.
 - Analysing the output for CoqRegExp revealed the following insights: lemmas were the most common kind of proof-object; the definitions of a regular-expression and what it means for a regular-expression to match a string were among the most important proof-objects (by PageRank values); each module and the library as a whole is primarily split into two groups ([executable functions](#) and [proofs of correctness](#));
 - Analysing the output for OOT revealed the following insights: there are many proof-objects in BG and PF that do not lead to the proof of the Feit-Thompson OOT (about 4-6 every section, most of them lemmas); clustering proof-objects shows grouping at the module-level; these groups demarcate different parts/phases of BG and PF; this phase demarcation is also reflected in the *reversed dependencies* at the proof-object level; the differences obtained by reversing proof-object dependencies could be a visual representation of the difference between *developing* mathematics (top-down, as needed) versus *presenting* mathematics (bottom-up, eagerly).

5 | Conclusions

5.1 In Hindsight	44
5.2 Future Work	44

I have just described a tool that allows a user to understand a Coq library more effectively than existing tools. It does so by translating Coq libraries to Neo4j graph databases. I gave examples of new questions that can now be answered by querying the library and showed how network-analysis algorithms and visualisations can provide new insights into the structure of a Coq library. I used a variety of technologies, (Coq, Neo4j, OCaml and R) tools (APOC, RNeo4j, igraph, visNetwork) and algorithms (centrality, clustering and other statistics) in this project. Its feature set, performance and output on real-world Coq libraries show that this project has potential to be useful to the Coq community.

5.1 In Hindsight

I learnt *a lot* during this project. Working with existing systems *and* creating something novel at the same time was difficult to juggle mentally. I found that I became better at learning how to use different frameworks and managing my time effectively as the project progressed. Initially, I was intrigued and fascinated by the project idea, suggested by the project supervisor. After I had completed preliminary preparations, the novelty, difficulty and usefulness of the concept had taken root.

However, Coq's lack of good, structured documentation made implementing the project frustrating. It was easy to get stuck in the details and spend hours staring at the entire Coq compiler code base just to figure out how to do something. Sometimes, having several options for how to proceed but no clear way of comparing between them made implementation difficult. It was rarely *writing* the code which was the issue, but *knowing what to aim for*. It was in those moments that having an experienced, focused and clear supervisor was the most useful.

I found project management and organisation easy and immensely helpful in providing a useful structure for guiding implementation. A good grasp of Unix, Git and programming languages proved to be essential for executing this project as smoothly as possible. Evaluating the project gave me the chance to switch from *developing* to *using* and exploring what was possible with the project in its final form. With the benefit of hindsight, the only thing I would do differently would be to make sure my development workflow (Makefiles, Linux VM, Git, editors, continuous-integration builds) was set up from the start – it would have been easier and useful for longer.

5.2 Future Work

Immediate extensions to this project could come from fleshing out the model further, to include tactics, notation and other aspects of the Coq system. More ambitious extensions could involve somehow presenting some information to a user *as they are working* on something, or using data science techniques to analyse the output model as a means to providing new insights or a more helpful compiler which can suggest how to (re)structure a project. Although very useful and easily interpretable in the context of mathematical theories, this project's core concept could also be applied to any programming language and its libraries, to provide new ways of becoming familiar with and understanding increasingly complex software projects.

Bibliography

- Alex Bavelas. Communication patterns in task-oriented groups. *The Journal of the Acoustical Society of America*, 22(6):725–730, 1950.
- Helmut Bender, George Glauberman, and Walter Carlip. *Local analysis for the odd order theorem*, volume 188. Cambridge University Press, 1994.
- Aaron Clauset, Mark EJ Newman, and Cristopher Moore. Finding community structure in very large networks. *Physical review E*, 70(6):066111, 2004.
- Linton C. Freeman. A set of measures of centrality based on betweenness. *Sociometry*, 40(1):35–41, 1977. ISSN 00380431. URL <http://www.jstor.org/stable/3033543>.
- Georges Gonthier. Formal proof—the four-color theorem. *Notices of the AMS*, 55(11):1382–1393, 2008.
- Georges Gonthier, Andrea Asperti, Jeremy Avigad, Yves Bertot, Cyril Cohen, François Garillot, Stéphane Le Roux, Assia Mahboubi, Russell O’Connor, Sidi Ould Biha, et al. A machine-checked proof of the odd order theorem. In *International Conference on Interactive Theorem Proving*, pages 163–179. Springer, 2013.
- Georges Gonthier, Assia Mahboubi, and Enrico Tassi. A Small Scale Reflection Extension for the Coq system. Research Report RR-6455, Inria Saclay Ile de France, 2015. URL <https://hal.inria.fr/inria-00258384>.
- Xavier Leroy. The compcert c verified compiler. *Documentation and user’s manual. INRIA Paris-Rocquencourt*, 2012.
- The Coq development team. *The Coq proof assistant reference manual*. LogiCal Project, 2004. URL <http://coq.inria.fr>. Version 8.0.
- Neo4j. Neo4j. neo4j.com. Accessed: 13/10/2016.
- Mark EJ Newman. The mathematics of networks. *The new palgrave encyclopedia of economics*, 2(2008):1–12, 2008.
- Mark EJ Newman and Michelle Girvan. Finding and evaluating community structure in networks. *Physical review E*, 69(2):026113, 2004.

Tobias Nipkow, Markus Wenzel, and Lawrence C. Paulson. *Isabelle/HOL: A Proof Assistant for Higher-order Logic*. Springer-Verlag, Berlin, Heidelberg, 2002. ISBN 3-540-43376-7.

Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab, 1999.

Thomas Peterfalvi. *Character theory for the odd order theorem*, volume 272. Cambridge University Press, 2000.

Benjamin C Pierce. The science of deep specification (keynote). In *Companion Proceedings of the 2016 ACM SIGPLAN International Conference on Systems, Programming, Languages and Applications: Software for Humanity*, pages 1–1. ACM, 2016.

Usha Nandini Raghavan, Réka Albert, and Soundar Kumara. Near linear time algorithm to detect community structures in large-scale networks. *Physical review E*, 76(3):036106, 2007.

A | Full Model

Here is the full translation of the Coq AST to kinds and subkinds.

```
1 let kind_of_constref =
2 let open Decl_kinds in function
3 | IsDefinition def -> ("definition", Some (match def with
4   | Definition -> "definition"
5   | Coercion -> "coercion"
6   | SubClass -> "subclass"
7   | CanonicalStructure -> "canonical_structure"
8   | Example -> "example"
9   | Fixpoint -> "fixpoint"
10  | CoFixpoint -> "cofixpoint"
11  | Scheme -> "scheme"
12  | StructureComponent -> "projection"
13  | IdentityCoercion -> "coercion"
14  | Instance -> "instance"
15  | Method -> "method"))
16 | IsAssumption a ->
17  ("assumption", Some (match a with
18  | Definitional -> "definitional"
19  | Logical -> "logical"
20  | Conjectural -> "conjectural")))
21 | IsProof th ->
22  ("proof", Some (match th with
23  | Theorem -> "theorem"
24  | Lemma -> "lemma"
25  | Fact -> "fact"
26  | Remark -> "remark"
27  | Property -> "property"
28  | Proposition -> "proposition"
29  | Corollary -> "corollary"))
30
31 let kind_of_ind ind =
32 let (mib,oib) = Inductive.lookup_mind_specif (Global.env ()) ind in
33 if mib.Declarations.mind_record <> None then
34  let open Decl_kinds in
```

```
35 begin match mib.Declarations.mind_finite with
36 | Finite -> "recursive_inductive"
37 | BiFinite -> "recursive"
38 | CoFinite -> "corecursive"
39 end
40 else
41 let open Decl_kinds in
42 begin match mib.Declarations.mind_finite with
43 | Finite -> "inductive"
44 | BiFinite -> "variant"
45 | CoFinite -> "coinductive"
46 end
47
48 let get_constr_type typ =
49   Names.KerName.to_string (Names.MutInd.user typ)
50
51 let kind_of_gref gref =
52 if Typeclasses.is_class gref then
53   ("class", None)
54 else
55   match gref with
56   | Globnames.ConstRef cst ->
57     kind_of_constref (Decls.constant_kind cst)
58
59   | Globnames.ConstructRef ((typ, _), _) ->
60     ("type_constructor", None)
61
62   | Globnames.IndRef ind ->
63     ("inductive_type", Some (kind_of_ind ind))
64
65   | Globnames.VarRef _ ->
66     assert false
67
68
69 let kind_of_obj = function
70 | G.Node.Gref gref ->
71   kind_of_gref gref
72 | G.Node.Module modpath ->
73   ("module", match modpath with
74   | Names.ModPath.MPbound _ -> Some "bound"
75   | Names.ModPath.MPdot _ -> Some "module"
76   | Names.ModPath.MPfile _ -> Some "file")
```

B | Timings

Below is the table of measurements I conducted of different tools that aim to help a user understand a Coq library on a Surface Pro 3 (Intel Haswell i7-4650U 1.7-3GHz, 8GB RAM, 512GB SSD) running Fedora 24 inside VirtualBox 5.1.16 on Windows 10 (Creator's Update). I rebooted the VM between measurements to reduce the effects of caching.

Timing	Test 1	Test 2	Test 3	Test 4	Test 5	Mean	Std. Dev.
DOC: U	7886	7885	7879	7875	7879	7880.8	4.6
DOC: S	307	303	288	295	300	298.6	7.4
DOC: R	10018	9997	9997	10017	10007	10007.2	10.3
DEP: U	1363	1365	1357	1359	1363	1361.4	3.3
DEP: S	2307	2299	2314	2292	2297	2301.8	8.7
DEP: R	4263	4254	4257	4257	4256	4257.4	3.4
DPD: U	6869	6868	6872	6878	6876	6872.6	4.3
DPD: S	654	631	639	646	650	644	9.1
DPD: R	9396	9381	9390	9404	9382	9390.6	9.7
CRT: U	79544	79529	79510	79526	79517	79525.2	12.9
CRT: S	2522	2518	2528	2522	2523	2522.6	3.6
CRT: R	98146	98173	98165	98156	98159	98159.8	10.1
2DT: U	1272235	1272244	1272264	1272250	1272239	1272246.4	11.3
2DT: S	2517	2537	2519	2517	2515	2521	9.1
2DT: R	1588302	1588305	1588304	1588305	1588299	1588303	2.5
ANL: U	255915	255909	255922	255914	255927	255917.4	7.1
ANL: S	49568	49565	49567	49565	49572	49567.4	2.9
ANL: R	371877	371884	371876	371890	371886	371882.6	6.0

Table B.1 – Table of execution time measurements of tools that aim to help a user understand a Coq library. All timings are in milliseconds. DOC=coqdoc, DEP=coqdep, DPD=dpd2graph, CRT=creation, 2DT=dpd2dot, ANL=analysis, U=user, S=system, R=real. Mean and sample standard deviation rounded to the nearest decimal point.

C | Sample Code

Here is the dependency collection code (recursing down the AST).

```
1 let collect_dependance gref =
2   match gref with
3   | Globnames.VarRef _ -> assert false
4
5   | Globnames.ConstRef cst ->
6     let cb = Environ.lookup_constant cst (Global.env()) in
7     let cl = match Global.body_of_constant_body cb with
8       Some e -> [e]
9     | None -> [] in
10    let cl = match cb.Declarations.const_type with
11      | Declarations.RegularArity t -> t::cl
12      | Declarations.TemplateArity _ -> cl in
13    List.fold_right collect_long_names cl Data.empty
14
15   | Globnames.IndRef i ->
16     let _, indbody = Global.lookup_inductive i in
17     let ca = indbody.Declarations.mind_user_lc in
18     Array.fold_right collect_long_names ca Data.empty
19
20   | Globnames.ConstructRef (i,_) ->
21     let _, indbody = Global.lookup_inductive i in
22     let ca = indbody.Declarations.mind_user_lc in
23     (* So a constructor and its type are linked, BUT WRONG WAY AROUND *)
24     add_inductive i (Array.fold_right collect_long_names ca Data.empty)
```

Here is a small excerpt of how edges are formatted and corrected.

```
1 let out_edge fmt _g e =
2
3   (* incorporate src & dst types, flip if constructor & ind & match
4    * TO FIX WRONG WAY DEPENDENCY LINK FROM SEARCHDEPEND.ML4 *)
5   let src, dst, rel_type =
6     let matches typ n =
```

```
7   let dirname, name = G.Node.split_name (G.Node.obj n) in
8   get_constr_type typ = dirname ^ "." ^ name in
9
10  let src, dst = G.Edge.src e, G.Edge.dst e in
11
12  match G.Node.obj src, G.Node.obj dst with
13  | G.Node.Gref (Globnames.ConstructRef ((typ, _), _)),
14    G.Node.Gref (Globnames.IndRef _) when matches typ dst ->
15    dst, src, "CONSTRUCTED_BY"
16  | G.Node.Module _, G.Node.Module _ ->
17    dst, src, "CONTAINS"
18  | G.Node.Module _, G.Node.Gref _ ->
19    src, dst, "CONTAINS"
20  | _, _ ->
21    src, dst, "USES" in
22
23  let edge_attribs =
24    [ ("type", rel_type) ; ("weight", string_of_int (G.Edge.nb_use e))] in
25
26  (* NOTE: Flipped src and dst - src USED_BY dst <=> dst USES src*)
27  Format.printf fmt "E: %d %d [%a];@."
28    (G.Node.id dst)
29    (G.Node.id src)
30  pp_attribs edge_attribs
```

Here is the analysis of the CoqRegEx library.

```
1 suppressMessages(library(RNeo4j))
2 suppressMessages(library(igraph))
3 library(tictoc)
4
5 # Connect to DB instance
6 cat("Connecting to database... "); tic()
7 graph <- startGraph(
8   "http://localhost:7474/db/data/",
9   username="neo4j",
10  password="Neo4j")
11 time <- toc(quiet=TRUE); time <- time$toc - time$tic
12 cat(sprintf("done. (%.2fs)\n", time))
13
14 # Grab nodes...
15 cat("Getting nodes (definitions and proofs) "); tic()
16 nodes <- cypher(graph,
17   MATCH (obj:definition),
18
19   RETURN obj.objectId AS id,
```

```
20         obj.name AS label,
21         obj.path AS title,
22         \"triangle\" AS shape,
23
24     UNION
25     MATCH (obj:proof),
26
27     RETURN obj.objectId AS id,
28             obj.name AS label,
29             obj.path AS title,
30             \"square\" AS shape,
31     )
32
33 # and edges.
34 cat("and edges... ")
35 edges <- cypher(graph, "
36     MATCH (src)-[edge]->(dst)
37     WHERE (src:definition OR src:proof) AND (dst:definition OR dst:proof)
38     RETURN src.objectId AS from, dst.objectId AS to, edge.weight AS weight")
39 time <- toc(quiet=TRUE); time <- time$toc - time$tic
40 cat(sprintf("done. (%.2fs)\n", time))
41
42 # Constructing graph
43 cat("Constructing igraph... "); tic()
44 d_ig <- graph_from_data_frame(edges, directed=TRUE, nodes)
45 ig <- graph_from_data_frame(edges, directed=FALSE, nodes)
46 time <- toc(quiet=TRUE); time <- time$toc - time$tic
47 cat(sprintf("done. (%.2fs)\n", time))
48
49 # Calculate metrics on nodes
50 node_metric <- function(algorithm, graph) {
51     cat(sprintf("%s ... ", algorithm)); tic()
52
53     result <- switch(algorithm,
54                         PageRank = page_rank(graph)$vector,
55                         Betweenness = betweenness(graph),
56                         Closeness = closeness(graph))
57
58     time <- toc(quiet=TRUE); time <- time$toc - time$tic
59     cat(sprintf("done. (%.2fs)\n", time))
60
61     return(result)
62 }
63
64 # Compute PageRank (proofs and definitions)
65 nodes$pagerank <- node_metric("PageRank", d_ig)
```



```
112         "definition_proof_modularity",
113         "definition_proof_edge_betweenness",
114         "definition_proof_label_prop)); tic()
115 set =
116 MATCH (obj { objectId :.toInt({OBJID}) })
117 SET obj.definition_proof_pagerank = toFloat({PGR}),
118     obj.definition_proof_betweenness = toFloat({BTW}),
119     obj.definition_proof_closeness = toFloat({CLOSE}),
120     obj.definition_proof_modularity = toInt({MODGROUP}),
121     obj.definition_proof_edge_betweenness = toInt({EBTW}),
122     obj.definition_proof_label_prop = toInt({LBL_PRP})
123 "
124 transaction <- newTransaction(graph)
125 progressBar <- txtProgressBar(min=0,
126                                 max=nrow(nodes),
127                                 char='=',
128                                 width=80,
129                                 style=3)
130 for (i in 1:nrow(nodes)) {
131   row <- nodes[i,]
132   appendCypher(transaction,
133                 set,
134                 OBJID=row$id,
135                 PGR=row$pagerank,
136                 BTW=row$betweenness,
137                 CLOSE=row$closeness,
138                 MODGROUP=row$modularity,
139                 EBTW=row$edge_betweenness,
140                 LBL_PRP=row$label_prop)
141   setTxtProgressBar(progressBar, i)
142 }
143 close(progressBar)
144 cat("Committing transaction... ")
145 commit(transaction)
146 time <- toc(quiet=TRUE); time <- time$toc - time$tic
147 cat(sprintf("done. (%.2fs)\n", time))
```

Here is the visualisation of the CoqRegEx library.

```
1 suppressMessages(library(RNeo4j))
2 suppressMessages(library(igraph))
3 suppressMessages(library(visNetwork))
4 library(tictoc)
5
6 # Connect to DB instance
7 cat("Connecting to database... "); tic()
```

```

8 graph <- startGraph(
9   "http://localhost:7474/db/data/",
10  username="neo4j",
11  password="Neo4j")
12 time <- toc(quiet=TRUE); time <- time$toc - time$tic
13 cat(sprintf("done. (%.2fs)\n", time))
14
15 # Grab nodes...
16 cat("Getting nodes (definitions and proofs) "); tic()
17 nodes <- cypher(graph, "
18   MATCH (obj:definition),
19
20   RETURN obj.objectId AS id,
21         obj.path + '.' + obj.name AS title,
22         \"triangle\" AS shape,
23         obj.definition_proof_pagerank AS pagerank,
24         obj.definition_proof_betweenness AS betweenness,
25         obj.definition_proof_closeness AS closeness,
26         obj.definition_proof_edge_betweenness AS edge_betweenness,
27         obj.definition_proof_label_prop AS label_prop,
28         obj.definition_proof_modularity AS modularity
29   UNION
30   MATCH (obj:proof),
31
32   RETURN obj.objectId AS id,
33         obj.path + '.' + obj.name AS title,
34         \"square\" AS shape,
35         obj.definition_proof_pagerank AS pagerank,
36         obj.definition_proof_betweenness AS betweenness,
37         obj.definition_proof_closeness AS closeness,
38         obj.definition_proof_edge_betweenness AS edge_betweenness,
39         obj.definition_proof_label_prop AS label_prop,
40         obj.definition_proof_modularity AS modularity")
41
42 # and edges.
43 cat("and edges... "); tic()
44 edges <- cypher(graph, "
45   MATCH (src)-[edge]->(dst)
46   WHERE (src:definition OR src:proof) AND (dst:definition OR dst:proof)
47   RETURN src.objectId AS from, dst.objectId AS to, edge.weight AS weight")
48 time <- toc(quiet=TRUE); time <- time$toc - time$tic
49 cat(sprintf("done. (%.2fs)\n", time))
50
51 # Output visualisations
52 visualise <- function(nodes, edges, filename, layout_opts,
53                         edge_opts, skipIgraph=FALSE) {

```

```
54     cat(sprintf("Outputting %s... ", filename)); tic()
55
56     g <- visNetwork(nodes, edges, width="1600px", height="1600px") %>%
57         visInteraction(navigationButtons=TRUE,
58                         dragNodes=FALSE, zoomView=FALSE)
59
60     g <- do.call(visEdges, append(list(g), edge_opts))
61
62     if (!skipIgraph) {
63         g <- do.call(visIgraphLayout, append(list(g), layout_opts))
64     } else {
65         visLayout(g, randomSeed=layout_opts$randomSeed)
66     }
67
68     visSave(g, file = filename)
69
70     time <- toc(quiet=TRUE); time <- time$toc - time$tic
71     cat(sprintf("done. (%.2fs)\n", time))
72     return(g)
73 }
74
75 # Bucket
76 bucket <- function(x) { return(cut(log(0.001+x), 10, labels=FALSE)) }
77
78 # Non DrL layouts
79 other_layout <- list(randomSeed=1492)
80
81 # Edge options
82 edge_opts <- list(arrows="middle", color=list(opacity=1), dashes=FALSE)
83
84 # Modularity (FR)
85 nodes$value <- bucket(nodes$betweenness)
86 nodes$group <- nodes$modularity
87 other_layout$layout <- "layout_with_fr"
88 edge_opts <- list(color=list(opacity=1), dashes=FALSE)
89 visualise(nodes, edges, "direct.html", other_layout, edge_opts)
90
91 # Grid graph
92 nodes$value <- cut(nodes$betweenness, 10, labels=FALSE)
93 nodes$group <- nodes$modularity
94 other_layout$layout <- "layout_on_grid"
95 edge_opts <- list(arrows="middle", color=list(opacity=0.7), dashes=TRUE)
96 visualise(nodes, edges, "grid.html", other_layout, edge_opts)
97
98 # Hierarchical graph
99 nodes$value <- bucket(nodes$betweenness)
```

```
100 nodes$group <- nodes$modularity
101 other_layout$layout <- "layout_with_sugiyama"
102 edge_opts <- list(color=list(opacity=1), dashes=TRUE)
103 visualise(nodes, edges, "hierarchical.html", other_layout, edge_opts)
104
105
106 # Construct a hierarchical network
107 cat("Getting nodes (modules) "); tic()
108 modules <- cypher(graph, "
109   MATCH (obj:module)
110   RETURN obj.objectId AS id,
111       coalesce(obj.path + '.', '') + obj.name AS title,
112       \"circle\" AS shape,
113       NULL AS pagerank,
114       NULL AS betweenness,
115       NULL AS edge_betweenness,
116       NULL AS closeness,
117       NULL AS modularity,
118       NULL AS label_prop,
119       9 AS value")
120
121 # and edges.
122 cat("and edges... ")
123 contains <- cypher(graph, "
124   MATCH (src)-[edge:CONTAINS]->(dst)
125   WHERE (src:module) AND (dst:module OR dst:definition OR dst:proof)
126   RETURN src.objectId AS from, dst.objectId AS to,
127   )
128 time <- toc(quiet=TRUE); time <- time$toc - time$tic
129 cat(sprintf("done. (%.2fs)\n", time))
130
131 # Modularity (FR, modules)
132 nodes$value <- bucket(nodes$betweenness)
133 nodes$group <- nodes$modularity
134 modules$group <- max(nodes$group)+1
135 other_layout$randomSeed=440
136 edge_opts <- list(color=list(opacity=1), dashes=FALSE)
137 visualise(rbind(nodes, modules),
138             contains,
139             "module.html",
140             other_layout,
141             edge_opts,
142             skipIgraph=TRUE)
```

D | Project proposal

Computer Science Tripos – Part II – Project Proposal

Exploring the structure of mathematical theories
using graph databases

Dhruv C. Makwana, Trinity College

Originator: Dr. Timothy G. Griffin

Project Supervisor: Dr. Timothy G. Griffin

Directors of Studies: Dr. Frank Stajano & Dr. Sean B. Holden

Project Overseers: Dr. David J. Greaves & Prof. John Daugman

Introduction and Description of the Work

This project aims to (a) represent Coq libraries as Neo4j (graph) databases and (b) create a library of Neo4j queries with the goal of highlighting the structure and relationship between the representations of the proof-objects.

Mathematics textbooks aimed at professionals/researchers follow a well-established rhythm: define some constructions and some properties on them and prove theorems on both, with lemmas, corollaries and notation interspersed throughout. Such a presentation is concise but limiting: it is linear; it forces the reader to keep track of dependencies such as implicit assumptions, previously defined results and the types and conventions behind any notation used; and it offers little opportunity to consider and compare different approaches for arriving at a result (i.e. number of assumptions, number of steps, some notion of the importance of a result such as number of uses by later results).

With the increasing popularity of interactive theorem-provers such as Coq ([The Coq development team, 2004](#)) and Isabelle ([Nipkow et al., 2002](#)), many mathematical theories (such as the formidably large Feit-Thompson Odd Order Theorem [Peterfalvi 2000, Bender et al. 1994](#)) have been ([\(Gonthier et al., 2013\)](#)) or

are being translated and formalised into machine-checked proof-scripts. However, these proof-scripts on their own inherit the same disadvantages as the aforementioned textbooks, as well as some new ones: they are usually more verbose and explicit and are primarily designed for automation/computation than readability. The former (usually out of necessity to convey to the computer the intended meaning) leads to unnecessary “noise” in the proof and the latter departs from the vocabulary or flow a natural-language presentation may have.

The database world is currently experiencing a tremendous explosion of creativity with the emergence of new data models and new ways of representing and querying large data sets. *Graph databases* have been developed to deal with highly connected data sets and path-oriented queries. That is, graph databases are optimised for computing transitive-closure and related queries, which pose a huge challenge for traditional, relational databases.

A graph-based approach to the representation and exploration of the structure of proof-objects would be a far more natural expression of the complex relationships (i.e. chains of dependencies) involved in constructing mathematical theories. Questions such as “What depends on this lemma and how many such things are there?” or “What are the components of this definition?” could thus be expressed concisely (questions which are not even expressible with standard relational databases systems such as SQL). A popular graph database, Neo4j ([Neo4j](#)) with an expressive query language *Cypher* will be used for this project.

Resources Required

Software

Several components of software will be required for executing this project, all of which are available for free online.

For using the proof-scripts, the Coq proof assistant will be required, as well as the Proof General proof assistant ([proofgeneral.github.io/](#)) for the Emacs ([www.gnu.org/software/emacs/](#)) text-editor.

For writing the plug-in to access Coq proof-objects, the parser and associated modules in the source code will be required ([github.com/coq/coq](#)) written in the OCaml programming language ([ocaml.org](#)) with the OCaml’s Package Manager OPAM ([opam.ocaml.org](#)).

For building the library of (Cypher) queries, Neo4j Community Edition will be used.

Hardware

Implementation and testing will be done on both Windows 10 and a Linux Virtual Machine as appropriate and convenient on a Surface Pro 3 (Intel Haswell i7-4650U

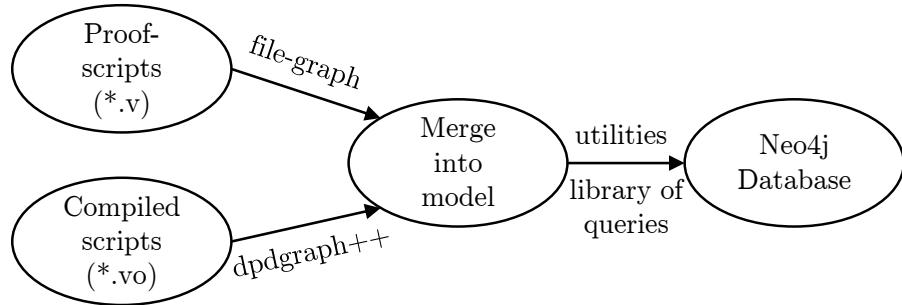


Figure D.1 – System Components

1.7-3GHz, 8GB RAM, 512GB SSD) with a personal GitHub account and physical backup drive (Seagate 1TB) making hourly backups using Windows' File History.

Starting Point

Some existing tools offer part of the solutions: these will be used and combined as appropriate. A large part of the project will rely on my knowledge of OCaml and Coq usage and internals.

Coq-dpdgraph (github.com/Karmaki/coq-dpdgraph) is a tool which analyses dependencies between *compiled* Coq proofs. As such, desirable information about notation, tactics, definitions and the relationship between a type and its constructors is lost.

Coqdep is a utility included with Coq which analyses dependencies *at the module level* by tracking `Require` and `Import` statements.

Coq SerAPI (github.com/ejgallego/coq-serapis) is a work-in-progress library and communication protocol for Coq designed to make low-level interaction with Coq easier, especially for IDEs. It has a starting point for gathering some statistics of proof-objects in a project.

All of these tools have the same disadvantage: they present information statically, with no way to query and interact with the information available.

Substance and Structure of the Project

The project will have three major parts, as shown in Figure D.1.

Processing Compiled Files

First, using coq-dpdgraph as a starting point, a tool which expresses a compiled proof-script as CSV files (shown as “dpdgraph++” in the diagram). Finding what information can and should be extracted will be an iterative process. Although coq-dpdgraph is functional, it is very basic with no way of even relating the relationship between a (co-)inductive type and its constructors, hence much work is to be done to even come close to utilising the full potential of compiled proof-scripts.

Processing Source Code Directly

Second, using Coq’s sophisticated extensible-parser, to parse, gather and convert to CSV files the desirable but missing information coq-dpdgraph does not extract (shown as “file-graph” in the diagram). An interesting feature of Coq’s parser is that it allows new constructs and notation to be defined: this is used heavily in some projects and therefore poses a great challenge for simply understanding and using the parser effectively.

Extraction and Analysis Tools

Lastly, writing utilities to automate analysis of Coq files and importing them into Neo4j and libraries of queries to run on imported data in Neo4j. Since it is not known what sort of data can be extracted and what will be useful or interesting to know, modelling the data – in this case the structure and objects of a mathematical proof – will be a non-trivial task which will be tackled iteratively.

Extensions

Extensions for this project will come from the process of adapting the project to be compatible with SSReflect ([Gonthier et al., 2015](#)), part of the Mathematical Components set of tools for Coq. These set of tools use low-level hooks in the Coq plugin system to significantly alter the specification and computation of proofs. As such, although they allow for large-scale projects to be formalised more easily, they are non-standard and would thus be very difficult to support fully.

Success Criteria

Alongside a planned and written dissertation describing the work done, the following criteria will be used to evaluate the success of this project:

1. A schema of attributes and relations for each proof-object is defined.
2. Programs which convert proof-scripts and compiled proofs to CSV files are implemented.
3. A library of queries in order to manipulate and explore the proof-objects is implemented.
4. These new sets of tools are shown to have more capabilities and perform comparably to existing tools for exploring mathematical theories.

Timetable and milestones

Date	Milestone
21-10-2016	Complete Project Proposal
04-11-2016	Finish a prototype compiled-to-CSV tool. Get familiar with Neo4j Cypher. Understand how to use the Coq parser.
18-11-2016	Refine compiled-to-CSV tool: tests and documentation. Explore queries possible and start the library. Begin work on translating Coq constructs from proof-scripts.
02-12-2016	Finish a prototype script-to-CSV tool.
16-12-2016	Test and document script-to-CSV tool.
30-12-2016	Begin work on integrating tools into one workflow.
13-01-2017	Stabilise and document whole project so far. Prepare presentation for CoqPL Conference.
27-01-2017	Look at SSReflect and evaluate changes to be made.
10-02-2017	Incorporate changes from feedback/new features.
24-02-2017	Test and document the new features.
10-03-2017	Write Introduction, Preparation and Implementation chapters.
24-03-2017	Fix bugs/unexpected problems.
07-04-2017	Write Evaluation and Conclusion chapters.
21-04-2017	Fix bugs/unexpected problems.
05-05-2017	Complete Dissertation (references, bibliography, appendix, formatting).
