

# The Coq Library as a Theory Graph

Dennis Müller, Florian Rabe, and Claudio Sacerdoti Coen<sup>3</sup>

<sup>1</sup> University Erlangen-Nuremberg

<sup>2</sup> LRI Paris

<sup>3</sup> University of Bologna

**Abstract.** Representing proof assistant libraries in a way that allows further processing in other systems is becoming increasingly important. It is a critical missing link for integrating proof assistants both with each other or with peripheral tools such as IDEs or proof checkers. Such representations cannot be generated from library source files because they lack semantic enrichment (inferred types, etc.) and only the original proof assistant is able to process them. But even when using the proof assistant’s internal data structures, the complexities of logic, implementation, and library still make this very difficult.

We describe one such representation, namely for the library of Coq, using OMDoc theory graphs as the target format. Coq is arguably the most formidable of all proof assistant libraries to tackle, and our work makes a significant step forward.

On the theoretical side, our main contribution is a translation of the Coq module system into theory graphs. This greatly reduces the complexity of the library as the more arcane module system features are eliminated while preserving most of the structure. On the practical side, our main contribution is an implementation of this translation. It takes the entire Coq library, which is split over hundreds of decentralized repositories, and produces easily-reusable OMDoc files as output.

## 1 Introduction and Related Work

*Motivation* A critical bottleneck in interactive theorem proving is data sharing, both between proof assistants and between proof assistants and related tools. The general situation is in stark contrast to the global push for FAIR data practices [oFD18], i.e., findable, accessible, interoperable, and reusable sharing of data. Currently, for example, any reuse of a library must go through the respective proof assistant. Thus, any novel idea is typically limited to the implementation framework and data flows provided by the proof assistant; and out-of-the-box experiments that by-pass the proof assistant are expensive, often prohibitively so. This limitation is particularly relevant as proof assistants are becoming more mature and many challenges are shifting from prover design to library management tasks like refactoring, reuse, search, and user interfaces.

For multiple reasons, Coq is the most formidable target for library sharing among all proof assistants: Firstly, the logic of Coq is arguably the most complex among all major proof assistants. This applies not only to the core logic, but also

to the processing pipeline from user-visible source to kernel representation and the library-building features of module system, sections, etc. Secondly, the code base of the Coq system has grown to a point where it is very hard for outsiders to navigate it. Thirdly, Coq has been so successful that its library is now so vast that it is non-trivial to even understand its extent — it is now split over several hundred repositories with non-uniform build processes.

*Contribution* Our overall goal is making the Coq library easier to access, search, interoperate with, and reuse. Even though we make a substantial first step, comprehensive progress on this front will take years. Concretely, in this paper, we translate the high-level structure of the Coq library that is visible to the Coq kernel (including modules and sections but not records or type classes) into the MMT language [RK13] for theory graphs. Theory graphs are an attractive target format because they allow preserving most of the library structure while being significantly simpler. Besides being designed to maintain theory graphs, MMT also provides a flexible logical framework that allows us to define the logical syntax of Coq. Thus, our MMT theories include all information in the Coq kernel including universes, inductive types, proof terms, and termination of recursive functions.

We translate all 49 Coq packages that are distributed via `opam` (a package manager originally designed for ocaml software) and that compile with the latest version of Coq (8.9.0). These comprise more than 383,500 logical declarations, making ours one of the largest proof assistant library translations ever and the largest for Coq.

*Related Work* Multiple ad hoc exports between different proof assistant libraries have been achieved. The general design of instrumenting the proof assistant kernel to export the library as a trace of checking it was first applied in [OS06]. This has proved to be the only feasible design, and all major later exports including ours employed variants of it. For example, Coq was the target format in [KW10].

Exports specifically into MMT were achieved for Mizar in [IKRU13], HOL Light in [KR14], PVS in [KMOR17], and very recently for Isabelle in not-yet published work. This overall line of research was presented in [KR16].

Similarly, to the MMT research, proof assistant libraries have been exported into Dedukti [BCH12]. Coquine is the tool used for translating Coq to Dedukti, and while MMT exports focus on preserving the original structure, Coquine focuses on reverifying proofs. Unlike our logic definition, Coquine includes a formalization of the typing rules in order to type check the export. In order to make this feasible, this translation eliminates several features of the Coq logic so that the typing rules become simpler. Our export, on the contrary, makes the dual trade-off, covering the entire logic at the expense of representing the typing rules. Concretely, the original version [BB12] covered most of the small standard library, using simplifications like collapsing the universe hierarchy. A later reimplemention [Ass15] used a more faithful representation of the logic. But it still omitted several language features such as modules, functors and universe polymorphism and therefore could not translate a significant part of the library.

[CK18] develops a translation of Coq into first-order logic in order to apply automated provers as hammers. Like Coqine, it only covers a subset of the language. It can in principle be used to translate the entire library, but that would have very limited use: Indeed, due to the nature of this application, it is actually *beneficial* to ignore all modular structure and even to not care about soundness, for example to collapse all universes.

*Overview* The structure of our paper follows the three major parts of building a library: the core logical language, the language for library building, and the library itself. Sect. 3 describes the representation of the Coq logics (As we will see, Coq technically provides a few slightly different logics.) in MMT. This results in a few manually written MMT theories. Sect. 4 specifies how the library language features of Coq can be mapped to MMT theories and morphisms. And Sect. 5 describes the implementation that translates the Coq library into MMT. We recap the relevant preliminaries about Coq and MMT in Sect. 2, and we discuss limitations and future work in Sect. 6.

*Acknowledgments* The authors were supported by DFG grant RA-18723-1 OAF and EU grant Horizon 2020 ERI 676541 OpenDreamKit. Tom Wiesing helped with scripting for automatically creating the necessary `mathhub.info` repositories and pushing to them.

## 2 Preliminaries

### 2.1 Coq

We give only an extremely dense explanation of the Coq language and refer to Appendix A<sup>4</sup> and [Coq15] for details. We use a **grammar** for the *abstract* syntax seen by the Coq kernel (Fig. 1) because that is the one that our translation works with. Even though we do not have space to describe all aspects of the translation in detail, we give an almost entire grammar here in order to document most of the language features we translate. A slightly more comprehensive grammar is presented in the companion paper [Sac19]: the omitted features do not pose additional problems to the translation to MMT and are omitted to simplify the presentation.

The Coq library is organized hierarchically into (from high to low) packages, nested directories, files, nested modules, and nested sections, where “nested” means that multiple levels of the same kind may be nested. Modules and sections are optional, i.e., logical declarations can occur in files, modules, or sections. When forming base logic expressions  $E$ , universes  $U$ , module expressions  $M$ , and module type expressions  $T$ , declarations can be referred to by **qualified identifiers**  $e$ ,  $u$ ,  $m$ , resp.  $t$  formed from

1. The root identifier that is defined by the Coq package. Typically, but not necessarily, every package introduces a single unique root identifier.

---

<sup>4</sup> Available online at <https://kwarc.info/people/dmueller/pubs/CoqImport.pdf>

```

decl ::= — base logic declarations
      e@{y*} : E [ := E ]
      | Universe u
      | Constraint  $u(< | \leq | =)u$ 
      | (Inductive | CoInductive) (e@{y*} : E := (e@{y*} : E)*)*
      — section declarations and variables in sections
      | Section  $\underline{s} := \text{decl}^*$ 
      | Variable  $x : E$ 
      | Polymorphic Universe  $y$ 
      | Polymorphic Constraint  $(u | y)(< | \leq | =)(u | y)$ 
      — module (type) declarations
      | Module Type  $\underline{m} (\underline{m} : T)^* <: T^* := (T | \text{decl}^*)$ 
      | Module  $\underline{t} (\underline{m} : T)^* [ : T ] <: T^* [ := (M | \text{decl}^*) ]$ 
E ::= — base logic expressions
    e@{U*} | x | Prop | Set | Type@{U} |  $\Pi x : E.E$  |  $\lambda x : E.E$  | E E
    | Match e E E E* | (Fix | CoFix)  $\mathbb{N}(\underline{e} : E := E)^*$  | let  $x : E := E$  in E
    | E.N | (E : E)
U ::= — universes
    u | y | max U U | succ U
T ::= — module type expressions
    [!]  $t m^* | T$  with  $e' := E$  |  $T$  with  $m' := M$ 
M ::= — module expressions
    [!]  $m m^*$ 
x, y ::= variables for term, universe respectively
e, u, m, t, s ::= qualified identifiers of expressions, universes, modules, module types, sections
e', m' ::= relative qualified identifiers of expressions, modules
e, u, m, t, s ::= fresh (unqualified) identifiers

```

**Fig. 1.** Coq Kernel Grammar

2. One identifier each for every directory or file that leads from the package root to the respective Coq source file.
3. One identifier each (possibly none) for every nested module (type) inside that source file that contains the declarations.
4. The unqualified name e, u, m, resp. t.

Note that section identifiers do not contribute to qualified names: the declarations inside a section are indistinguishable from the ones declared outside the section. Relative qualified names are always relative to a module type, i.e. they are missing the root identifiers and the directory identifiers.

**Expressions** are the usual  $\lambda$ -calculus expressions with dependent products  $\Pi x : \text{term.term}$  (used to type  $\lambda$ -abstractions), let binder, **let...in**, sorts **Prop**, **Set**, **Type@{U}** (used to type types), casts  $(E : E)$ , (co)inductive types with primitive pattern-matching, (co)fixpoints definitions (i.e. terminating recursive functions) and record projections  $(E.N)$ . Notably, Coq maintains a partially or-

dered **universe hierarchy** (a directed acyclic graph) with consistency constraints of the form  $U(< | \leq)U'$ .

**Module types and modules** are the main mechanism for grouping base logic declarations. Public identifiers introduced in modules are available outside the module via qualified identifiers, whereas module types only serve to specify what declarations are public in a module. We say that a module  $M$  *conforms* to the module type  $T$  if

- $M$  declares every constant name that is declared in  $T$  and does so with the same type,
- $M$  declares every module name  $\underline{m}$  that is declared in  $T$  and does so with a module type that conforms to the module type (= the set of public declarations) of  $\underline{m}$  in  $T$ ,
- if any such name has a definiens in  $T$ , it has the same definiens in  $M$ .

Conformation of a module *type* to a module type is defined accordingly.

Both modules and module types may be defined in two different ways: *algebraically* as an abbreviation for a module (type) expression (the definiens), or *interactively* as a new module (type) given by a list of declarations (the body). Every module (type) expression can be elaborated into a list of declarations so that algebraic module (type) declarations can be seen as abbreviations of interactive ones. A module may also be abstract, i.e., have neither body nor definiens. The  $<$ : and  $:$  operators may be used to attach conformation conditions to a module (type), and we will explain their semantics in Sect. 4.2.

Module (type)s can be abstracted over a list of module bindings  $\underline{m} : T$ , which may be used in the definiens/body. When the list is not empty the module (type) is called a *functor* (type). A functor must be typed with a functor type that has the same list of module bindings. Conformation induces a notion of subtyping between module and functor types. Coq treats functor types contravariantly and allows for higher-order functors. However, from our experiments, it seems that this feature is never used in any of the libraries we exported from Coq.

Module (type) expressions can be obtained by functor application, whose semantics is defined by  $\beta$ -reduction in the usual way, unless “!” annotations are used. According to complex rules that we will ignore in the rest of the paper for lack of space, the “!” annotations performs  $\beta$ -reduction and then triggers the replacement of constants defined in the actual functor argument with their definiens. Finally, the **with** operator adds a definition to an abstract field in a module type.

**Sections** may be used to subdivide files and module (type)s. These are similar to module functors except that they abstract over base logic declarations, which are interspersed in the body and marked by the **Variable** and **Polymorphic** keywords. The section itself has no semantics: outside the section, all normal declarations are  $\lambda\Pi$ -abstracted over all **Variable/Polymorphic** declarations.

## 2.2 MMT

MMT aims at being a universal representation language for formal systems. Its syntax was designed carefully to combine simplicity and expressivity. Fig. 2 gives the fragment needed for Coq, and we refer to [RK13,Rab17] for details.

```

decl ::= Theory  $l =^{[E]}$  decl*
      | Morph  $l : E \rightarrow E =^{[E]}$  decl*
      | include  $E$ 
      |  $l [ : E ] [ = E ]$ 
      | Rule Scala object
 $E$  ::=  $g \mid g?l \mid x \mid g?l((x [ : E ] [ = E ])^*, E^*)$ 
 $g$  ::=  $URI?l \mid g/l$ 
 $l$  ::= local identifiers

```

Fig. 2. MMT Grammar

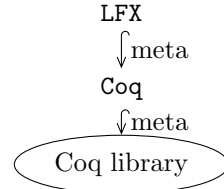
An expression **Theory**  $l =^{[E]}$   $B$  introduces a named **theory**  $l$  with body  $B$  and (optional) meta-theory  $E$ . In the simplest case, nested theories create a tree of declarations, whose leafs are **constant** declarations  $l : E_1 = E_2$  introducing a local identifier  $l$  with optional type  $E_1$  and definiens  $E_2$ .

Named theories have global identifiers  $g$  of the form  $g = NS?l_1 / \dots / l_n$  where  $NS$  is the namespace URI assigned by the containing source file, and the  $l_i$  are the local identifier of the nested theories (i.e.  $l_n$  is the local identifier of the theory itself,  $l_{n-1}$  is the local identifier of the containing theory etc.).

Every constant has a unique URI of the form  $g?l$  where  $g$  is the global identifier of the containing theory and  $l$  is the local identifier of the constant. In an expression  $E$ , every theory or constant is always referenced via its global identifier.

Every theory  $g$  induces a set of **expressions**  $E$  formed from identifiers, bound variables  $x$ , and composed expressions  $g?l(C, E_1, \dots, E_n)$ , where  $C$  is a (possibly empty) list of variables  $x[( : E')][[ = E'']]$  considered bound in the subexpressions  $E_i$ . The semantics of expressions is signaled by the chosen constructor  $g?l$ , which is usually a constant declared in the meta-theory.

Meta-theories yield the language-independence of MMT: The meta-theory  $L$  of a theory  $t$  defines the language in which  $t$  is written. For example, the diagram on the right indicates the form of the theory graph we built in this paper: all theories in the Coq library will be translated to MMT theories with meta-theory **Coq**, which in turn has meta-theory **LFX**. LFX is an extension of the logical framework LF [HHP93] that is strong enough to define the Coq logic in a theory **Coq**. The semantics of LFX itself is obtained by declaring **rules**: these are Scala objects that are injected dynamically into the MMT kernel. We refer to [Rab18] for the general mechanism and the definition of LF. Most importantly, LFX declares



1. 5 constants for forming the composed expressions `type`,  $A \rightarrow B$ ,  $\Pi x : A. B$ ,  $\lambda x : A. t$  and function applications  $f(a_1, \dots, a_n)$ ;
2. syntax rules that render the composed expression `LFX?lambda(x : A, t)` as  $\lambda x : A. t$ ; and
3. about 10 typing rules for the LF type system.

The MMT theory graphs arise from adding **morphisms**  $m : s \rightarrow t =^M B$ . Here  $m$  is a morphism that maps the meta-theory of  $s$  to  $t$ , and the body  $B$  must contain a defined constant  $g?c = E$  for every  $s$ -constant  $g?c$  and some  $t$ -expression  $E$ . Then the homomorphic extension of  $m$  maps any  $s$ -expression to a  $t$ -expression in a way that, critically, preserves all judgments, e.g., if  $\vdash_s e : E'$ , then  $\vdash_t m(e) : m(E')$ . Such morphisms have theorem-flavor and are used to represent language translations and refinement or interpretation theorems. Alternatively, include declarations **include**  $s$  in a theory  $t$  are used to create morphisms  $s \rightarrow t$  that hold by definition: their semantics is that all  $s$ -constants are also visible to  $t$ , which implies that the identity map of constants is a theory morphism. Include morphisms are used to represent inheritance and extension relations between theories and are depicted in diagrams as  $s \hookrightarrow t$ . Expressions over a theory  $t$  can use all constants declared in  $t$ , in the meta-theory of  $t$ , a parent theory of  $t$ , or in theories included into  $t$ .

### 3 Defining the Coq Base Logics in a Logical Framework

We define an MMT theory `Coq` that defines the base logic of Coq. This theory will occur as the meta-theory of all MMT theories generated from files in the Coq library (except when flags are used, see below). The theory `Coq` is available at <https://gl.mathhub.info/Coq/foundation/blob/master/source/Coq.mmt>. We briefly describe it in the sequel and refer to Appendix B<sup>5</sup> for details.

*Expressions* Due to lack of space, we only present the encoding of the PTS fragment of Coq, i.e. we omit `let...in`, projections, inductive types, pattern matching and (co)fixpoints. The encoding is quite straightforward.

```

Theory Coq =LFX
  univ : type   max : univ → univ → univ   succ : univ → univ
  expr : type   Prop : expr   Set : expr   Type : univ → expr
  Π : expr → (expr → expr) → expr   λ : expr → (expr → expr) → expr
  app : expr → expr → expr
  hastype : expr → expr → type
  exprOfType : expr → type = λe : expr. {x : expr | hastype x e}

```

Our representation of the syntax is a Curry-style encoding, in which all expressions have the same LF type and the binary typing judgment between expressions is formalized by separate judgment `hastype`. We do not give any typing

<sup>5</sup> Available online at <https://kwarc.info/people/dmueller/pubs/CoqImport.pdf>

rules here, but they could now be added in a straightforward way (except of course that Coq’s typing rules are very complex and doing so is correspondingly time consuming). There are alternative Church-style encodings, where a Coq expression of type  $E$  is represented as an LF-expression of type `exprOfType E` for an operator `exprOfType : expr → type`. These would be preferable because they allow declaring Coq identifiers as, e.g., `zero : exprOfType Nat` instead of erasing their type by using a declaration `zero : expr`. This is also why they are used in Coqine [BB12,Ass15] to formalize the calculus of constructions in Dedukti. However, Church encodings introduce so much representational overhead that they would make the translation of the entire Coq library infeasible. To gain the best of both worlds, we use predicate subtyping to define the `exprOfType E` as the subtype of `expr` containing those  $x$  for which the judgment `hastype x E` holds.

With these declarations in place, we can for example translate the definition of a universe-polymorphic identity function

$$\text{id@}\{y\} : \Pi A : \text{Type@}\{y\}. A \rightarrow A := \lambda A : \text{Type@}\{y\}. \lambda x : A. x$$

of Coq to the following MMT definition over the theory `Coq`

$$\begin{aligned} \text{id} : \Pi y : \text{univ}. \text{exprOfType} (\Pi (\text{Type } y) (\lambda A. A \rightarrow A)) \\ = \lambda y. \lambda (\text{Type } y) (\lambda A. \lambda A (\lambda x. x)) \end{aligned}$$

This captures all relevant information of the Coq definition with minimal representational overhead. Note how Coq’s  $\Pi$  and  $\lambda$ -binding are represented using LF higher-order abstract syntax, whereas universe polymorphism is represented directly using LF’s binders.

*Logic Variants* Maybe surprisingly, Coq does not actually use a single logic: it offers flags that allow choosing variants of the type theory. Two flags are of particular importance as they are required by some of the libraries:

- `-impredicative-set` changes the typing rule of the dependent function space  $\Pi$  so that the type of functions that takes in input an inhabitant of a large universe and return a set is still a (small) set instead of being a larger type; the flag is inconsistent when assumed together with any axiom of choice and classical logic
- `-type-in-type` squashes all universe except `Prop` and `Set` into the single universe `Type`. The resulting inconsistency `Type : Type` is acceptable and useful in some applications, e.g., those that focus on computation rather than deduction and need the possibility to write non terminating functions.

All variants can be formalized similarly using slightly different typing rules. As we omit the typing rules anyway, we simply create theories `Coq`, `ImpredicativeCoq`, and `InconsistentCoq`, all of which include `CoqSyntax` and then contain placeholder comments for the typing rules. When extracting the library from Coq, we record the flags used to compile each Coq file. Depending on that information, we choose one of the above three theories as the meta-theory of the MMT-theory that is the translation of that file.



## 4 Representing the Coq Structuring Language in MMT

### 4.1 Overview

Coq	MMT
package	namespace
directory	namespace
file	file that declares a theory
module type	theory
module	theory
visibility of a module $m$ to $p$	inclusion morphism $m \hookrightarrow p$
module typing $M : T$	morphism $T \rightarrow M$
module conformation $M <: T$	morphism $T \rightarrow M$
module type conformation $T <: T'$	morphism $T' \rightarrow T$
section	theory
variable in a section	constant
any base logic declaration	constant

**Fig. 3.** Overview of the Translation

Fig. 3 gives an overview of our translation. Above the file level, our translation preserves the structure of Coq exactly: every directory or file in a Coq package is translated to a corresponding directory resp. file in an MMT archive. Therefore, all Coq directories result in MMT namespace URIs.

A **qualified identifier** consisting of root  $r$ , directories  $d_1, \dots, d_r$ , file name  $f$  with extension  $v$ , modules (types)  $m_1, \dots, m_s$ , and name  $n$  is translated to the MMT URI `coq : /r/d1/.../dr?f/m1/.../ms?n`. Note that the MMT URI makes clear, which parts of the qualified identifier are directory, file, or module (type) names without having to dereference any part of the URI.

The only subtlety here is that we translate every Coq source file to a theory. Effectively, we treat every Coq file  $f.v$  in directory  $D$  like the body of a module of name  $f$ ; and we translate it to an OMDoc file  $f.omdoc$  containing exactly one MMT theory with URI  $D?f$ .

If we translated files to namespaces instead of theories, the above MMT URI would be `coq : /r/d1/.../dr/f?/m1/.../ms?n`. We would have preferred this, but it is not possible: In Coq, base logic declarations may occur directly in files whereas MMT constants may only occur inside theories. Thus, we have to wrap every Coq source file into an MMT theory. This is inconsequential except that we have to add corresponding include declarations in MMT: for every file  $f'$  that is referenced in a file  $f$ , the resulting MMT theory must include the MMT theory of  $f'$ . Fortunately, this information is anyway stored by Coq so that this is no problem.

In the sequel, we write  $\bar{i}$  for the MMT translation of the Coq item  $i$  except that, if  $i$  is a Coq identifier, we write  $i$  in MMT as well if no confusion is possible.

We omit the translation of base logic expressions and refer to Appendix C<sup>6</sup> for the translation of sections.

## 4.2 Modules and Module Types

We translate all files and module (type)s to MMT theories. Thus, the parent  $p$  of every module (type), which is either a file or a module (type), is always translated to an MMT theory; and every module (type) with parent  $p$  is translated to an MMT theory nested into  $\bar{p}$ . Overall, Coq’s tree of nested module (type) and constant declarations is translated to an isomorphic tree of nested MMT theories and constants, augmented with the theory morphisms induced by module type conformation (explained below).

We first consider the non-functor case and generalize to functors in Sect. 4.3. An algebraic module (type) is translated by first computing its explicit representation as an interactive one, according to the meta-theory of Coq. In this way we only have to consider the interactive case and we lift from the user the burden of understanding the intricacies of algebraic module (type) resolution (e.g. the complex semantics given by “!” annotations, or the issue of generativity for functors application).

*Module and Module Types as Theories* So let us consider an interactive module type **Module Type**  $t <: T_1 \dots T_n := B$ . We translate it to an MMT theory  $\bar{t}$  whose body arises by declaration-wise translation of the declarations in  $B$ . However, we have to treat universes specially because Coq maintains them globally: all universes and constraint declared in  $B$  are not part of  $\bar{t}$  and instead treated as if they had been declared at the beginning of the containing source file. We discuss the treatment of  $<$ : attributions below.

A module is translated in exactly the same way as a module type. The semantic difference between modules and module types is that a module  $m$ , once closed, exports all declarations in its body to its parent  $p$ . We capture this difference in MMT exactly by additionally generating an include declaration **include**  $\bar{m}$  after the theory  $\bar{m}$ , which makes  $\bar{m}$  available to  $\bar{p}$ .

It may be tempting to alternatively translate module types  $t$  to theories  $\bar{t}$  and a module  $m : t$  to a theory morphism  $\bar{m} : \bar{t} \rightarrow \bar{p}$ . This would elegantly capture how every module is an implementation of the module type by providing definitions for the abstract declarations in  $t$ . But that is not possible because Coq allows abstract fields even in modules, and such modules would not induce MMT theory morphisms. A maybe surprisingly example is the following, which is well-typed in Coq:

```

Module Type  $s := e : \text{False}$ 
Module  $m : s := e' : \text{False}, e : \text{False} := e'$ 
 $x : \text{False} := m.e$ 

```

Here the abstract declaration of  $e'$  in the module  $m$  is allowed even though it is used to implement the interface  $s$  of  $m$ .

<sup>6</sup> Available online at <https://kwarc.info/people/dmueller/pubs/CoqImport.pdf>

*Conformation as a Theory Morphism* Now we translate the attributions  $<: T_i$  on a module (type) and the attributions  $: T$  on a module. Our translation does not distinguish modules and module types, and if multiple attributions  $<: T_i$  are present, they are translated individually. Thus, we only need to consider the cases  $m <: T$  and  $m : T$ . In both cases, our translation consists of a morphism  $\overline{m}^*$  from  $\overline{T}$  to  $\overline{m}$  that witnesses that  $m$  conforms to  $T$ .

Inspecting the grammar, we see that  $T$  has normal form  $(t\ m_1 \dots m_r)\ \mathbf{with}\ k_1 := K_1 \dots \mathbf{with}\ k_s := K_s$ , where  $k := K$  unifies the cases of constant and module instantiations. As we will see below, if  $t$  is a module type functor (i.e., if  $r! = 0$ ), its module parameters  $x_1 : T_1, \dots, x_r : T_r$  are translated as if  $t$  were not functorial and the  $x_i : T_i$  were abstract modules in the body of  $t$ . Accordingly, we treat  $T$  in the same way as  $t\ \mathbf{with}\ x_1 := m_1 \dots \mathbf{with}\ x_r := m_r\ \mathbf{with}\ k_1 := K_1 \dots \mathbf{with}\ k_s := K_s$ , and therefore we can restrict attention to the case  $r = 0$ .

We know that  $t$  is translated to a theory  $\overline{t}$ , and if  $T$  is well-typed, recursively translating the  $K_i$  already yields a partial theory morphism  $\varphi$  from  $t$  to  $\overline{p}$ . Because  $\overline{m}$  is a theory nested into  $\overline{p}$ ,  $\varphi$  is also a morphism into  $\overline{m}$ . It remains to extend  $\varphi$  with assignments for the remaining declarations of  $\overline{t}$ . Now we observe that if  $m$  conforms to  $T$  in Coq, we obtain a well-typed MMT theory morphism  $\overline{m}^*$  by extending  $\varphi$  with assignments  $\overline{t}?k := \overline{m}?k$  for every such name  $k$ . (The converse is also true if we add typing rules to Coq that adequately capture the typing relation of base logic expressions.)

For  $<:$  attributions, this is all we have to do. But  $\overline{t} \xrightarrow{\overline{m}^*} \overline{m} \xrightarrow{r} \overline{m.\mathbf{impl}}$  an attribution  $m : T$  is stronger than an attribution  $m <: T$ . It additionally restricts the interface of  $m$  to what is declared in  $T$ . Therefore, we have to do a little bit more in the case  $m : T$  as shown on the right:

1. We rename the theory  $\overline{m}$  to  $\overline{m.\mathbf{impl}}$ .
2. We create a second theory  $\overline{m}$  that is a copy of  $\overline{t}$  where all qualified names use  $m$  in place of  $t$ .
3. We create a morphism  $r : \overline{m} \rightarrow \overline{m.\mathbf{impl}}$  that maps every name of  $\overline{m}$  to itself.
4. We create the renaming morphism  $\overline{m}^*$  with codomain  $\mathfrak{m}$  in the same way as for the case  $m <: T$ . The morphism just performs the renaming since  $\overline{m}$  and  $\overline{t}$  only differs on names.

The  $:$  attributions of Coq are peculiar because  $\overline{m.\mathbf{impl}}$  can never be referenced again — the morphism  $r$  can be seen as a dead end of the theory graph. In fact, trying to understand this part of the translation made us realize the following curiosity about the Coq module system. Consider the well-typed example on the right, where we use indentation for scoping. The attribution  $m : s$  in the module type  $t$  hides the definition of the field  $f$  in the module  $m$ . Because that definition is never considered again, the module  $n$  can supply a different definition for  $f$  later on.

```

Module Type s := f : Nat
Module Type t :=
  Module m : s := f : Nat := 0
Module n : t :=
  Module m := f : Nat := 1

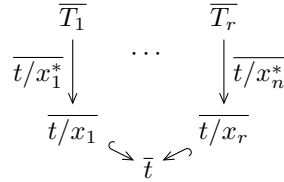
```

Indeed, the Coq kernel imperatively throws away  $\overline{m.\text{impl}}$  after checking it. When  $f$  declares a logical axiom instead of a type like **Nat**, the behaviour is somewhat more intuitive: if we only care that a definition (i.e., proof) exists, it is fine to give two different ones in different places. But this treatment is markedly different from analogous features of other languages: In object-oriented programming,  $n$  would not be allowed to redefine  $m$  because the definition of  $f$  is still inherited even if remains inaccessible. Similarly, in theory graphs with hiding [MAH06,CHK<sup>+</sup>12], the model  $n$  of  $t$  would be required to implement  $e$  in a way that is consistent with the hidden definition in  $t$ .

### 4.3 Functors

*Declaring Functors* In many ways the parameters of an interactive module (type) can be treated in the same way as the declarations in its body  $B$ . Indeed, the declaration **Module (Type)**  $t(x_1 : T_1, \dots, m_r : T_r) := B$  is well-typed iff **Module (Type)**  $t := \text{Module } x_1 : T_1, \dots, \text{Module } m_r : T_r, B$  is.

This motivates what we call the *covariant* translation of functors, which we employ: parameters of interactive modules or module types are translated as if they occurred as abstract module declarations at the beginning of the body. Thus, the two variants of  $t$  above are translated to the exact same MMT theory. The resulting diagram is shown on the right. Note that the theories  $\overline{t/x_i}$  are nested into  $\overline{t}$



and additionally included into  $\overline{t}$ . We also add metadata to the declarations of the  $\overline{t/x_i}$  to record the fact that they used to be functor parameters. Algebraic module (type) functors and  $<$ : and  $:$  attributions are handled in the same way as in Sect. 4.2.

Technically, only a *contravariant* translation that translates functors to theory-valued functions would capture the semantics of functors adequately. For example, the covariant translation of **Module**  $m(x_1 : T_1, \dots, m_r : T_r) := B$  results in the same diagram as above with an additional inclusion morphism from  $\overline{m}$  into the parent  $p$  of  $m$ . Thus, the theories  $\overline{t/x_i}$  become falsely included into  $\overline{p}$ . More formally:

1. the covariant translation preserves well-typedness only if the library does not rely on Coq's contravariant rule for functor type subtyping, which is the case for all the libraries exported so far.
2. the covariant translation does not reflect well-typedness.

However, considering that the Coq library is already well-typed and that the covariant translation is so much simpler, that is sufficient for many practical applications.

*Applying Functors* Coq functor application may be partial and curried. Thus, it is sufficient to restrict attention to  $r = 1$ . So consider a module type declaration **Module Type**  $t(x : T) := B$  and a module  $m : T$ . We have to define the translation of the module type  $t(m)$ , whose semantics is determined in Coq by substituting  $m$  for  $x$  in  $B$ . We want the translation to be compositional, i.e., defined in terms of the theory  $\bar{t}$  arising from the translation of  $t$  and the morphism  $\bar{m}^* : \bar{t} \rightarrow \bar{m}$  arising from the translation of  $m$ , as in the diagram on the right. As defined above, the functor  $t$  is translated to a theory  $\bar{t}$  with a nested theory  $\overline{t/x}$  that conforms to  $\bar{T}$  as well as an include of  $\bar{t/x}$ . Let  $p$  be the Coq file or module (type) in which  $t(m)$  is well-typed; thus,  $\bar{p}$  is a theory that includes  $\bar{m}$ .

$$\begin{array}{ccccc}
 \bar{T} & \xrightarrow{\overline{t/x^*}} & \overline{t/x} & \hookrightarrow & \bar{t} \\
 & & \downarrow \bar{m}^* & & \downarrow \\
 & & \bar{m} & & \\
 & & \downarrow & & \\
 \bar{p} & \hookrightarrow & & & \overline{t(m)}
 \end{array}$$

This situation is well-known in MMT theory graphs: to translate  $t(m)$ , we have to create a new theory nested into  $\bar{p}$  such that the diagram is a pushout. The canonical choice of a pushout [Rab17,CMR16] amounts to copying over all declarations in  $\bar{t}$  except for replacing all occurrences of  $x$  by the homomorphic translation along  $\bar{m}^*$ . This yields the same theory as translating the flattening of  $t(m)$ .

## 5 Translating the Coq Library

Our translation is implemented in two steps. Firstly, Coq is instrumented via kernel hooks to export the internal kernel data structures into Coq-near gzipped XML files. This part of the translation is described in detail in [Sac19]. Secondly, we read these XML files into MMT and translate them to MMT data structure, which we then write out to disk as OMDoc files. (Actually, we use xz-compressed OMDoc files because the uncompressed files would be too large.)

This separation into a Coq-export and an MMT-import may seem inefficient. But this design has proved very successful in the past [IKRU13,KR14,KMOR17]. Moreover, it allows separating the formidable practical task of exporting anything from the theoretical task of specifying the translation.

Notably, the whole export of the 49 opam packages for Coq libraries that currently compile with Coq 8.9.0 (recently released) comprises about 1.3 million XML files totaling 224.7 GB (interestingly, merely counting the number and sizes of XML files takes around 15 minutes). More packages will be translated in the future as soon as they are ported from previous Coq versions. Translating to MMT only the Coq standard library takes about 22 hours on a standard laptop, converting 15.4GB of (uncompressed) XML into 28.9MB of (compressed) OMDoc. This reduction is not only due to a high compressibility of the OMDoc, but also reflects the fact that every declaration in Coq corresponds to multiple XML files with partially redundant information.

## 6 Conclusion

*Evaluation* Our translation covers entirely the syntax of the Coq language and it preserves typing and soundness, with the exception of higher order functors, functor declarations and contravariant functor subtyping. The latter three features do not seem to occur in the 49 libraries that have an opam package which is up-to-date with the last Coq release. As more libraries become available, we will have to verify that our covariant functor translation is still adequate.

Moreover, we are confident that, if and when future work yields a complete formalization of the Coq typing rules in an LF-like logical framework, our translation will be in a format suitable for rechecking the entire library — with the obvious caveat that such a rechecking would face even more serious scalability issues than we had to overcome so far.

An obvious way to verify that the exported information is sound and complete for type-checking would be to implement an importer for Coq itself or, in alternative, for an independent verifier for the logic of Coq, like the one implemented in Dedukti or the one developed in the HELM/MoWGLI projects [APS+03] and later incorporated into Matita 0.5.x series [ARCT11]. In both cases one would need to develop a translation from MMT theories to the modular constructs of the language, which requires more research. For example, no translation of MMT theories and theory morphisms into modules, module types and functors is currently known.

We would also like to stress that independent verification is not the aim of our effort: the main point of exporting the library of Coq to MMT is to allow independent services over them, like queries, discovery of alignments with libraries of other tools or training machine learning advisers that can drive hammers. Most of these services can be implemented even if the typing information is incomplete or even unsound (e.g. if all universes are squashed to a single universe, making the logic inconsistent).

*Limitations and Future Work* Our translation starts with the Coq kernel data structures and is thus inherently limited to the structure seen by the kernel. Therefore, record types and type classes are presented just as inductive types, that is the way they are elaborated before passing them to the kernel. This is unfortunate as recent Coq developments, most importantly the Mathematical Components project [GGMR09], have made heavy use of records to represent theory graph-like structuring and an unelaborated representation would be more informative to the user and to reasoning tools.

In fact, even sections are not visible to the kernel, and we were able to include them because we were able to reconstruct the section structure during our translation. We expect that similar efforts may allow for including record types and canonical structures in the theory graph in the future and we plan to start working on that next.

Many libraries avoid module and functors and achieve modularity using other more recent features of Coq that are invisible at the kernel level, like type classes. Moreover type classes, canonical structures, coercions, etc. are necessary

information to extend a library because they explain how the various mathematical notions are meant to be used. While the already cited services that we plan to provide do not depend on them, importing the library in another system to build on top of it surely does. Therefore a future challenge will be to find system independent generalizations and representations of such constructs, which will be necessary to incorporate them into a logic and system independent tool like MMT.

Our formal representation of Coq declarations includes the types of all constants and variables, but we use a single type in the logical framework for all Coq expressions. As we explain in Sect. 3, we consider a typed representation of expressions infeasible at this point. Our representation does not include the typing rules for the expression, but this is not due to a principal limitation: it is possible to add these rules to let MMT type-check the library. But formalizing the rules of the Coq type system is in itself a major challenge, and representing the details of, e.g., Coq’s treatment of pattern matching or sort polymorphism may even require innovations in logical framework design.

## References

- APS<sup>+</sup>03. Andrea Asperti, Luca Padovani, Claudio Sacerdoti Coen, Ferruccio Guidi, and Irene Schena. Mathematical Knowledge Management in HELM. *Ann. Math. Artif. Intell.*, 38(1-3):27–46, 2003.
- ARCT11. Andrea Asperti, Wilmer Ricciotti, Claudio Sacerdoti Coen, and Enrico Tassi. The Matita Interactive Theorem Prover. In *Automated Deduction - CADE-23*, volume 6803 of *Lecture Notes in Computer Science*, pages 64–69. Springer, 2011.
- Ass15. A. Assaf. *A framework for defining computational higher-order logics*. PhD thesis, École Polytechnique, 2015.
- BB12. M. Boespflug and G. Burel. CoqInE: Translating the Calculus of Inductive Constructions into the lambda Pi-calculus Modulo. In D. Pichardie and T. Weber, editors, *Proof Exchange for Theorem Proving*, 2012.
- BCH12. M. Boespflug, Q. Carbonneaux, and O. Hermant. The  $\lambda\Pi$ -calculus modulo as a universal proof language. In D. Pichardie and T. Weber, editors, *Proceedings of PxTP2012: Proof Exchange for Theorem Proving*, pages 28–43, 2012.
- CHK<sup>+</sup>12. M. Codescu, F. Horozal, M. Kohlhase, T. Mossakowski, and F. Rabe. A Proof Theoretic Interpretation of Model Theoretic Hiding. In T. Mossakowski and H. Kreowski, editors, *Recent Trends in Algebraic Development Techniques 2010*, pages 118–138. Springer, 2012.
- CK18. L. Czajka and C. Kaliszyk. Hammer for coq: Automation for dependent type theory. *Journal of Automated Reasoning*, 61(1-4):423–453, 2018.
- CMR16. M. Codescu, T. Mossakowski, and F. Rabe. Selecting Colimits for Parameterisation and Networks of Specifications. In M. Roggenbach and P. James, editors, *Workshop on Algebraic Development Techniques*, 2016.
- Coq15. Coq Development Team. The Coq Proof Assistant: Reference Manual. Technical report, INRIA, 2015.
- GGMR09. F. Garillot, G. Gonthier, A. Mahboubi, and L. Rideau. Packaging mathematical structures. In S. Berghofer, T. Nipkow, C. Urban, and M. Wenzel,

- editors, *Theorem Proving in Higher Order Logics*, pages 327–342. Springer, 2009.
- HHP93. R. Harper, F. Honsell, and G. Plotkin. A framework for defining logics. *Journal of the Association for Computing Machinery*, 40(1):143–184, 1993.
- IKRU13. M. Iancu, M. Kohlhase, F. Rabe, and J. Urban. The Mizar Mathematical Library in OMDoc: Translation and Applications. *Journal of Automated Reasoning*, 50(2):191–202, 2013.
- KMOR17. M. Kohlhase, D. Müller, S. Owre, and F. Rabe. Making PVS Accessible to Generic Services by Interpretation in a Universal Format. In M. Ayala-Rincon and C. Muñoz, editors, *Interactive Theorem Proving*, pages 319–335. Springer, 2017.
- KR14. C. Kaliszyk and F. Rabe. Towards Knowledge Management for HOL Light. In S. Watt, J. Davenport, A. Sexton, P. Sojka, and J. Urban, editors, *Intelligent Computer Mathematics*, pages 357–372. Springer, 2014.
- KR16. M. Kohlhase and F. Rabe. QED Reloaded: Towards a Pluralistic Formal Library of Mathematical Knowledge. *Journal of Formalized Reasoning*, 9(1):201–234, 2016.
- KW10. C. Keller and B. Werner. Importing HOL Light into Coq. In M. Kaufmann and L. Paulson, editors, *Interactive Theorem Proving*, pages 307–322. Springer, 2010.
- MAH06. T. Mossakowski, S. Autexier, and D. Hutter. Development graphs - Proof management for structured specifications. *J. Log. Algebr. Program*, 67(1–2):114–145, 2006.
- oFD18. European Commission Expert Group on FAIR Data. Turning fair into reality, 2018. URL: <https://doi.org/10.2777/1524>.
- OS06. S. Obua and S. Skalberg. Importing HOL into Isabelle/HOL. In N. Shankar and U. Furbach, editors, *Automated Reasoning*, volume 4130. Springer, 2006.
- Rab17. F. Rabe. How to Identify, Translate, and Combine Logics? *Journal of Logic and Computation*, 27(6):1753–1798, 2017.
- Rab18. F. Rabe. A Modular Type Reconstruction Algorithm. *ACM Transactions on Computational Logic*, 19(4):1–43, 2018.
- RK13. F. Rabe and M. Kohlhase. A Scalable Module System. *Information and Computation*, 230(1):1–54, 2013.
- Sac19. Claudio Sacerdoti Coen. A plugin to export Coq libraries to XML. In *12th International Conference on Intelligent Computer Mathematics (CICM 19)*, Lecture Notes in Artificial Intelligence, 2019.