

Extending MKM Formats at the Statement Level

Fulya Horozal, Michael Kohlhase, and Florian Rabe

Computer Science, Jacobs University Bremen, Germany <http://kwarc.info>

Abstract. Successful representation and markup languages find a good balance between giving the user freedom of expression, enforcing the fundamental semantic invariants of the modeling framework, and allowing machine support for the underlying semantic structures. MKM formats maintain strong invariants while trying to be foundationally unconstrained, which makes the induced design problem particularly challenging.

In this situation, it is standard practice to define a minimal core language together with a scripting/macro facility for syntactic extensions that map into the core language. In practice, such extension facilities are either fully unconstrained (making invariants and machine support difficult) or limited to the object level (keeping the statement and theory levels fixed).

In this paper we develop a general methodology for extending MKM representation formats at the statement level. We show the utility (and indeed necessity) of statement-level extension by redesigning the OMDoc format into a minimal, regular core language (strict OMDoc) and an extension (pragmatic OMDoc) that maps into strict OMDoc.

1 Introduction

The development of representation languages for mathematical knowledge is one of the central concerns of the MKM community. After all, practical mathematical knowledge management consists in the manipulation of expressions in such languages. To be successful, MKM representation formats must balance multiple concerns. A format should be expressive and flexible (for depth and ease of modeling), foundationally unconstrained (for coverage), regular and minimal (for ease of implementation), and modular and web-transparent (for scalability). Finally, the format should be elegant, feel natural to mathematicians, and be easy to read and write. Needless to say that this set of requirements is over-constrained so that the design problem for MKM representation formats lies in relaxing some of the constraints to achieve a global optimum.

In languages for formalized mathematics, it is standard practice to define a minimal core language that is extended by macros, functions, or notations. For example, Isabelle [Pau94] provides a rich language of notations, abbreviations, syntax and printing translations, and a number of definitional forms. In narrative formats for mathematics, for instance, the \TeX / \LaTeX format – arguably the

most commonly used format for representing mathematical knowledge – goes a similar way, only that the core language is given by the \TeX layout primitives and the translation is realized by macro expansion and is fully under user control. This extensibility led to the profusion of user-defined \LaTeX document classes and packages that has made $\text{\TeX}/\text{\LaTeX}$ so successful.

However, the fully unconstrained nature of the extensibility makes ensuring invariants and machine support very difficult, and thus this approach is not immediately applicable to content markup formats. There, MathML3 [Aus+10] is a good example of the state of the art. It specifies a core language called “strict content MathML” that is equivalent to OpenMath [Bus+04b] and “full content MathML”. The first subset uses a minimal set of elements representing the meaning of a mathematical expression in a uniform, regular structure, while the second one tries to strike a pragmatic balance between verbosity and formality. The meaning of non-strict expressions is given by a fixed translation: the “strict content MathML translation” specified in section 4.6 of the MathML3 recommendation [Aus+10].

This language design has the advantage that only a small, regular sublanguage has to be given a mathematical meaning, but a larger vocabulary that is more intuitive to practitioners of the field can be used for actual representation. Moreover, semantic services like validation only need to be implemented for the strict subset and can be extended to the pragmatic language by translation. Ultimately, a representation format might even have multiple pragmatic front-ends geared towards different audiences. These are semantically interoperable by construction.

The work reported in this paper comes from an ongoing language design effort, where we want to redesign our OMDoc format [Koh06] into a minimal, regular core language (*strict OMDoc 2*) and an extension layer (*pragmatic OMDoc 2*) whose semantics is given by a “pragmatic-to-strict” (\mathcal{PS}) translation. While this problem is well-understood for mathematical *objects*, extension frameworks at the *statement level* seem to be restricted to the non-semantic case, e.g. the `amsthm` package for \LaTeX .

Languages for mathematics commonly permit a variety of pragmatic statements, e.g., implicit or case-based definitions, type definitions, theorems, or proof schemata. But representation frameworks for such languages do not include a generic mechanism that permits introducing arbitrary pragmatic statements — instead, a fixed set is built into the format. Among logical frameworks, Twelf/LF [PS99; HHP93] permits two statements: defined and undefined constants. Isabelle [Pau94] and Coq [BC04] permit much larger, but still fixed sets that include, for example, recursive case-based function definitions. Content markup formats like OMDoc permit similar fixed sets.

A large set of statements is desirable in a representation format in order to model the flexibility of individual languages. A large *fixed* set on the other hand is unsatisfactory because it is difficult to give a theoretical justification for fixing any specific set of statements. Moreover, it is often difficult to define the semantics of a built-in statement in a foundationally unconstrained representation

format because many pragmatic statement are only meaningful under certain foundational assumptions.

In this paper we present a general formalism for adding new pragmatic statement forms to our OMDoc format; we have picked OMDoc for familiarity and foundation-independence; any other foundational format can be extended similarly. Consider for instance the pragmatic statement of an “implicit definition”, which defines a mathematical object by describing it so accurately, that there is only one object that fits this description. For instance, the exponential function exp is defined as the (unique) solution of the differential equation $f = f'$ with $f(0) = 1$. This form of definition is extensively used in practical mathematics, so pragmatic OMDoc should offer an infrastructure for it, whereas strict OMDoc only offers “simple definitions” of the form $c := d$, where c is a new symbol and d any object. In our extension framework, the \mathcal{PS} translation provides the semantics of the implicit definition in terms of the strict definition $exp := \iota f.(f' = f \wedge f(0) = 1)$, where ι is a “definite description operator”: Given an expression A with free variable x , such that there is a unique x that makes A valid, $\iota x.A$ returns that x , otherwise $\iota x.A$ is undefined.

Note that the semantics of an implicit definition requires a definite description operator. While most areas of mathematics at least implicitly assume its existence, it should not be required in general because that would prevent the representation of systems without one. Therefore, we make these requirements explicit in a special theory that defines the new pragmatic statement and its strict semantics. This theory must be imported in order for implicit definitions to become available. Using our extension language, we can recover a large number of existing pragmatic statements as definable special cases, including many existing ones of OMDoc. Thus, when representing formal languages in OMDoc, authors have full control what pragmatic statements to permit and can define new ones in terms of existing ones.

In the next section, we will recap those parts of OMDoc that are needed in this paper. In Section 3, we define our extension language, and in Section 4, we look at particular extensions that are motivated by mathematical practice. Finally, in Section 5, we will address the question of extending the concrete syntax with pragmatic features as well.

2 MMT/OMDoc

OMDoc is a comprehensive content-based format for representing mathematical knowledge and documents. It represents mathematical knowledge at three levels: mathematical formulae at the *object level*, symbol declarations, definitions, notation definitions, axioms, theorems, and proofs at the *statement level*, and finally modular scopes at the *theory level*. Moreover, it adds an infrastructure for representing functional aspects of *mathematical documents* at the content markup level. OMDoc

Theories	Documents
Statements	
Objects	

1.2 has been successfully used as a representational basis in applications ranging from theorem prover interfaces, via knowledge based up to to eLearning systems. To allow this diversity of applications, the format has acquired a large, interconnected set of language constructs motivated by coverage and user familiarity (i.e., by pragmatic concerns) and not by minimality and orthogonality of language primitives (strict concerns).

To reconcile these language design issues for OMDoc 2, we want to separate the format into a *strict* core language and a *pragmatic* extension layer that is elaborated into strict OMDoc via a “pragmatic-to-strict” (\mathcal{PS}) translation.

For strict OMDoc we employ the foundation-independent, syntactically minimal MMT framework (see below). For pragmatic OMDoc, we aim at a language that is feature-complete with respect to OMDoc 1.2 [Koh06], but incorporates language features from other MKM formats, most notably from Isabelle/ISAR [Pau94], PVS [ORS92], and Mizar [TB85].

The MMT language was emerged from a complete redesign of the formal core¹ of OMDoc focusing on foundation-independence, scalability, modularity, while maintaining coverage of formal systems. The MMT language is described in [RK11] and implemented in [Rab08].

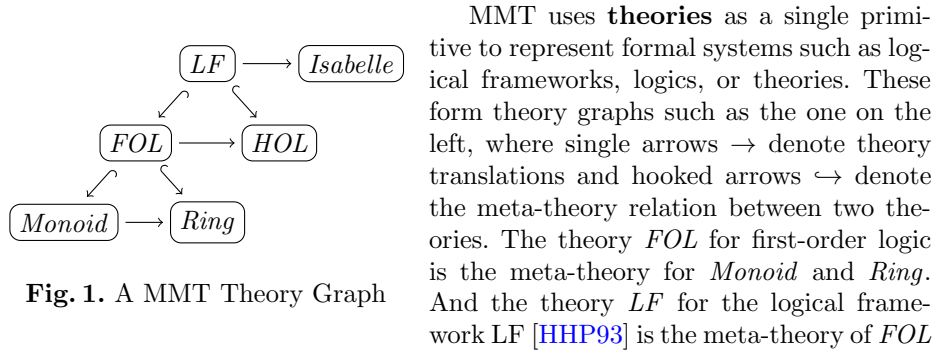


Fig. 1. A MMT Theory Graph

and *HOL* for higher-order logic. In general, we describe the theories with meta-theory *M* as *M-theories*. The importance of meta-theories in MMT is that the syntax and semantics of *M* induces the syntax and semantics of all *M*-theories. For example, if the syntax and semantics are fixed for *LF*, they determine those of *FOL* and *Monoid*.

At the statement level, MMT uses **constant** declarations as a single primitive to represent all OMDoc statement declarations. These are differentiated by the type system of the respective meta-theory. In particular, the Curry-Howard correspondence is used to represent axioms and theorems as plain constants (with special types).

In Figure 2, we show a small fragment of the MMT grammar that we need in the remainder of this paper. Meta-symbols of the BNF format are given in color.

¹ We are currently working on adding an informal (natural language) representations and a non-trivial (strict) document level to MMT, their lack does not restrict the results reported in this paper.

Modules	$G ::= (\mathbf{theory} \ T = \{\Sigma\})^*$
Theories	$\Sigma ::= \cdot \mid \Sigma, c[: E][= E] \mid \mathbf{meta} \ T$
Contexts	$\Gamma ::= \cdot \mid \Gamma, x[: E]$
Expressions	$E ::= x \mid c \mid EE^+ \mid E \Gamma.E$

Fig. 2. MMT Grammar

The module level of MMT introduces *theory declarations* $\mathbf{theory} \ T = \{\Sigma\}$. Theories Σ contain *constant declarations* $c[: E_1][= E_2]$ that introduce named atomic expressions c with optional type E_1 or definition E_2 . Moreover, each theory may declare its meta-theory T via $\mathbf{meta} \ T$.

MMT expressions are a fragment of OpenMath [Bus+04a] objects, for which we introduce a short syntax. They are formed from variables x , constants c , applications $E E_1 \dots E_n$ of functions E to a sequence of arguments E_i , and bindings $E_1 \Gamma. E_2$ that use a binder E_1 , a context Γ of bound variables, and a scope E_2 . Contexts Γ consist of variables $x[: E]$ that can optionally attribute a type E .

Example 1 (MMT-Theories). Below we give an MMT theory *Forms*, which will serve as the meta-theory of several logics introduced in this paper. It introduces all symbols needed to declare logical connectives and inference rules of a logic. The syntax and semantics of this theory are defined in terms of type theory, e.g., the logical framework LF [HHP93].

\mathbf{type} , \rightarrow , and \mathbf{lam} are untyped constants representing the primitives of type theory. \mathbf{type} represents the universe of all types, \rightarrow constructs function types $\alpha \rightarrow \beta$, and \mathbf{lam} represents the λ -binder. o is the type of logical formulas and \mathbf{proof} is a constant that assigns to each logical formula $F : o$ the type $\mathbf{proof} \ F$ of its proof.

$$\mathbf{theory} \ Forms = \{$$

$$\begin{array}{l} \mathbf{type} \\ \rightarrow \\ \mathbf{lam} \\ o : \mathbf{type} \\ \mathbf{proof} : o \rightarrow \mathbf{type} \end{array}$$

$$\}$$

3 A Framework for Language Extensions

We will now define our extension language (EL). It provides a syntactic means to define pragmatic language features and their semantics in terms of strict OMDoc.

Syntax EL adds two primitive declarations to MMT theories: *extension declarations* and *pragmatic declarations*:

$\Sigma ::= \Sigma, \mathbf{extension} \ e = \Phi$
$\mid \Sigma, \mathbf{pragmatic} \ c : \varphi$

Extension declarations **extension** $e = \Phi$ introduce a new declaration schema e that is described by Φ . Pragmatic declarations **pragmatic** $c : \varphi$ introduce new declarations that make use of a previously declared extension.

The key notion in both cases is that of *theory families*. They represent collections of theories by specifying their common syntactic shape. Intuitively, theory families arise by putting a λ -calculus on top of theory fragments Σ :

Theory Families	$\begin{array}{l} \Phi ::= \{\Sigma\} \mid \lambda x : E. \Phi \\ \varphi ::= e \mid \Phi E \end{array}$
-----------------	--

We group theory families into two non-terminal symbols as shown above: Φ is formed from theory fragments $\{\Sigma\}$ and λ -abstraction $\lambda x : E. \Phi$. And φ is formed from references to previously declared extension e and applications of parametric theory families to arguments E . This has the advantage that both Φ and φ have a very simple shape.

Any MMT theory may introduce extension declarations. However, pragmatic declarations are only legal if the used extension has been declared in the meta-theory:

Definition 1 (Legal Extension Declarations). *We say that an extension declaration **extension** $e = \lambda x_1 : E_1. \dots \lambda x_n : E_n. \{\Sigma\}$ is **legal** in an MMT theory T , if the declarations $x_1 : E_1, \dots, x_n : E_n$ and Σ are well-formed in T .*

This includes the case where Σ contains pragmatic declarations.

Definition 2 (Legal Pragmatic Declarations). *We say that a pragmatic declaration **pragmatic** $c : e E_1 \dots E_n$ is **legal** in an MMT theory T if there is a declaration **extension** $e = \lambda x_1 : E'_1. \dots \lambda x_n : E'_n. \{\Sigma\}$ in the meta-theory of T and each E_i has type E'_i .*

Example 2 (Extension Declarations). In Figure 3 we give the theory *Assertion*, which declares extensions for axiom and theorem declarations. Their semantics is defined in terms of the Curry-Howard representation of strict OMDoc.

Both extensions take a logical formula $F : o$ as a parameter. The extension *axiom* permits pragmatic declarations of the form $c : \text{axiom } F$. These abbreviate MMT constant declarations of the form $c : \text{proof } F$.

The extension *theorem* additionally takes a parameter $D : \text{proof } F$, which is a proof of F . It permits pragmatic declarations of the form $c : \text{theorem } F D$. These abbreviate MMT constant declarations of the form $c : \text{proof } F = D$.

Semantics Extension declarations do not have a semantics as such because the extension declared in M only govern what pragmatic declarations are legal in M -theories. In particular, contrary to the constant declarations in M , a model of M does not interpret the extension declarations.

The semantics of pragmatic declarations is given by elaborating them into strict declarations:

```

theory Assertion = {
  meta Forms
  extension axiom =  $\lambda F : o. \{$ 
     $c : \text{proof } F$ 
   $\}$ 
  extension theorem =  $\lambda F : o. \lambda D : \text{proof } F. \{$ 
     $c : \text{proof } F = D$ 
   $\}$ 
 $\}$ 

```

Fig. 3. An MMT Theory with Extension Declarations

Definition 3 (Pragmatic-to-Strict Translation \mathcal{PS}). A legal pragmatic declaration $\text{pragmatic } c : e \ E_1 \dots E_n$ is translated to a list of strict constant declarations

$$c.d_1 : \gamma(F_1) = \gamma(D_1), \dots, c.d_m : \gamma(F_m) = \gamma(D_m)$$

where γ substitutes every x_i with E_i and every d_j with $c.d_j$ if we have

$$\text{extension } e = \lambda x_1 : E'_1. \dots \lambda x_n : E'_n. \{d_1 : F_1 = D_1, \dots, d_m : F_m = D_m\}$$

and every expression E_i has type E'_i .

Example 3. Consider the following MMT theories in Figure 4: *HOL* includes the MMT theory *Forms* and declares a constant i as the type of individuals. It adds the usual logical connectives and quantifiers – here we only present truth (*true*) and the universal quantifier (\forall) – and introduces equality (\doteq) on expressions of type α . Then it includes *Assertion*. This gives *HOL* access to the extensions *axiom* and *theorem*.

Commutativity uses *HOL* as its meta-theory and declares a constant \circ that takes two individuals as arguments and returns an individual. It adds a pragmatic declaration named *comm* that declares the commutativity axiom for \circ using the axiom extension from *HOL*.

Commutativity' is obtained by elaborating *Commutativity* according to Definition 3.

4 Representing Extension Principles

Formal mathematical developments can be classified based on whether they follow the axiomatic or the definitional method. The former is common for logics where theories declare primitive constants and axioms. The latter is common for foundations of mathematics where a fixed theory (the foundation) is extended only by defined constants and theorems. In MMT, both the logic and the foundation are represented as a meta-theory M , and the main difference is that the definitional method does not permit undefined constants in M -theories.

However, this treatment does not capture conservative extension principles: These are meta-theorems that establish that certain extensions are acceptable

```

theory HOL = {
  meta Forms
  i      : type
  true : o
  ⋮
  ∀      : (α → o) → o
  ≐      : α → α → o
  include Assertion
}

theory Commutativity = {
  meta HOL
  o      : i → i → i
  pragmatic comm : axiom ∀x : i. ∀y : i. x o y ≐ y o x
}

theory Commutativity' = {
  meta HOL
  o      : i → i → i
  comm.c : proof ∀x : i. ∀y : i. x o y ≐ y o x
}

```

Fig. 4. A \mathcal{PS} Translation Example

even if they are not definitional. We can understand them as intermediates between axiomatic and definitional extensions: They may be axiomatic but are essentially as safe as definitional ones.

To make this argument precise, we use the following definition:

Definition 4. We call the theory family $\Phi = \lambda x_1 : E'_1. \dots \lambda x_n : E'_n. \{\Sigma\}$ **conservative** for M if for every M -theory T and all $E_1 : E'_1, \dots, E_n : E'_n$, every model of T can be extended to a model of $T, \gamma(\Sigma)$, where γ substitutes every x_i with E_i .

An extension declaration **extension** $e = \Phi$ is called **derived** if all constant declarations in Σ have a definiens; otherwise, it is called **primitive**.

Primitive extension declarations correspond to axiom declarations because they postulate that certain extensions of M are legal. The proof that they are indeed conservative is a meta-argument that must be carried out as a part of the proof that M is an adequate MMT representation of the represented formalism. Similarly, derived extension declarations correspond to theorem declarations because their conservativity follows from that of the primitive ones. More precisely: If all primitive extension principles in M are conservative, then so are all derived ones.

In the following, we will recover built-in extension statements of common representation formats as special cases of our extension declarations.

Implicit Definitions in OMDoc Implicit definitions of OMDoc 1.2 are captured using the following derived extension declaration. If the theory *ImplicitDefinitions* in Figure 5 is included into a meta-theory M , then M -theories may use implicit definitions.

Note that *ImplicitDefinitions* requires two other connectives: A description operator (ι) and a unique existential ($\exists^!$) are needed to express the meaning of an implicit definition. We deliberately assume only those two operators in order to maximize the re-usability of this theory: Using the MMT module system,


```

theory ImplicitDefinitions = {
  meta Forms
   $\exists^! : (\alpha \rightarrow o) \rightarrow o$ 
   $\iota : (\alpha \rightarrow o) \rightarrow \alpha$ 
   $\iota_{ax} : \text{proof } \exists^! x P x \rightarrow \text{proof } P (\iota P)$ 

  extension impldef =  $\lambda \alpha : \text{type}. \lambda P : \alpha \rightarrow o. \lambda m : \text{proof } \exists^! x : \alpha. P x. \{$ 
     $c : \alpha = \iota P$ 
     $c_{ax} : \text{proof } \exists^! x : \alpha. P x$ 
   $\}$ 
}

```

Fig. 5. An Extension for Implicit Definitions

any logic M in which these two operators are definable can import the theory *ImplicitDefinitions*.

More specifically, *ImplicitDefinitions* introduces the definite description operator as a new binding operator (ι), and describes its meaning by the axiom $\exists^! x P(x) \Rightarrow P(\iota P)$ formulated in ι_{ax} for any predicate P on α . The extension *impldef* permits pragmatic declarations of the form $f : \text{impldef } \alpha P m$, which defines f as the unique object which makes the property P valid. This leads to the well-defined condition that there is indeed such a unique object, which is discharged by the proof m . The pragmatic-to-strict translation from Section 3 translates the pragmatic declaration $f : \text{impldef } \alpha P m$ to the strict constant declarations $f.c : \alpha = \iota P$ and $f.c_{ax} : \text{proof } \exists^! x : \alpha P x$.

Mizar-Style Functor Definitions The Mizar language [TB85] provides a wide (but fixed) variety of special statements, most of which can be understood as conservative extension principles for first-order logic. A comprehensive list of the corresponding extension declarations can be found in [IKR11]. We will only consider one example in Figure 6.

```

theory FunctorDefinitions = {
  meta Forms
  extension functor =  $\lambda \alpha : \text{type}. \lambda \beta : \text{type}. \lambda \text{means} : \alpha \rightarrow \beta \rightarrow o. \{$ 
     $f : \alpha \rightarrow \beta$ 
     $f_{thm} : \text{means } x (f x)$ 
   $\}$ 
}

```

Fig. 6. An Extension for Mizar-Style Functor Definitions

The theory *FunctorDefinitions* describes Mizar-style implicit definition of a unary functor. This is different from the one above because it uses a primitive extension declaration that is well-known to be conservative. In Mizar, the axiom f_{thm} is called the definitional theorems induced by the implicit definition. Using

the extension *functor*, one can introduce pragmatic declarations of the form **pragmatic** $c : \text{functor } A B P$ that declare functors c from A to B that are defined by the property P .

Flexary Extensions The above two examples become substantially more powerful if they are extended to implicit definitions of functions of arbitrary arity. This is in fact supported by our extension language by using an LF-based logical framework with term and type sequences and natural numbers as parameters of extension. We omit the formal details of this framework here for simplicity and refer to [Hor12] instead. We only give one example in Figure 7 that demonstrates the potential.

```

theory CaseBasedDefinitions = {
  meta Forms
   $\wedge : o \rightarrow o \rightarrow o$ 
   $\Rightarrow : o \rightarrow o \rightarrow o$ 
   $\forall : (\alpha \rightarrow o) \rightarrow o$ 
  extension casedef =  $\lambda n : \mathbb{N}. \lambda \alpha : \text{type}. \lambda \beta : \text{type}. \lambda c : (\alpha \rightarrow o)^n.$ 
     $\lambda d : (\alpha \rightarrow \beta)^n. \lambda \rho : \text{proof } \forall x : \alpha. \bigvee^! [c_i x]_{i=1}^n. \{$ 
       $f : \alpha \rightarrow \beta$ 
       $ax : \text{proof } \forall x : \alpha. \wedge [c_i x \Rightarrow (f x) = (d_i x)]_{i=1}^n$ 
     $\}$ 
 $\}$ 

```

Fig. 7. An Extension for Case-Based Definitions

The theory *CaseBasedDefinitions* introduces an extension that describes the case-based definition of a unary function f from α to β that is defined using n different cases where each case is guarded by the predicate c_i together with the respective definiens d_i . Such a definition is well-defined if for all $x \in \alpha$ exactly one out of the $c_i x$ is true, which we abbreviate using an n -ary $\bigvee^!$ connective. Note that the formulation of this proof uses a special sequence constructor: for example, $[c_i x]_{i=1}^n$ simplifies to the sequence $c_1 x, \dots, c_n x$.

The pragmatic declaration **pragmatic** $f : \text{casedef } n \alpha \beta c_1 \dots c_n d_1 \dots d_n \rho$ corresponds to the following function definition:

$$f(x) = \begin{cases} d_1(x) & \text{if } c_1(x) \\ \vdots & \vdots \\ d_n(x) & \text{if } c_n(x) \end{cases}$$

HOL-Style Type Definitions Due to the presence of λ -abstraction and a description operator in HOL [Chu40], a lot of common extension principles become derivable in HOL, in particular, implicit definitions.

But there is one primitive definition principle that is commonly accepted in HOL-based formalizations of the definitional method: A Gordon/HOL type

definition [Gor88] introduces a new type that is axiomatized to be isomorphic to a subtype of an existing type. This cannot be expressed as a derivable extension because HOL does not use subtyping.

```

theory Types = {
  meta Forms
   $\forall : (\alpha \rightarrow o) \rightarrow o$ 
   $\exists : (\alpha \rightarrow o) \rightarrow o$ 
   $\doteq : (\alpha \rightarrow \alpha) \rightarrow o$ 
  extension typedef =  $\lambda\alpha : \text{type} \cdot \lambda A : \alpha \rightarrow o. \lambda P : \text{proof} \exists x : \alpha. A x. \{$ 
    T           : type
    Rep         :  $T \rightarrow \alpha$ 
    Abs         :  $\alpha \rightarrow T$ 
    Rep'        :  $\text{proof } \forall x : T. A (\text{Rep } x)$ 
    Rep_inverse :  $\text{proof } \forall x : T. \text{Abs } (\text{Rep } x) \doteq x$ 
    Abs_inverse :  $\text{proof } \forall x : \alpha. A x \Rightarrow \text{Rep } (\text{Abs } x) \doteq x$ 
  }
}

```

Fig. 8. An Extension for HOL-Style Type Definitions

The theory *Types* in Figure 8, formalizes this extension principle. Our symbol names follow the implementation of this definition principle in Isabelle/HOL [NPW02]. Pragmatic declarations of the form **pragmatic** $t : \text{typedef } \alpha A P$ introduce a new non-empty type t isomorphic to the predicate A over α . Since all HOL-types must be non-empty, a proof P of the non-emptiness of A must be supplied. More precisely, it is translated to the following strict constant declarations:

- $t.T : \text{type}$ is the new type that is being defined,
- $t.Rep : t.T \rightarrow \alpha$ is an injection from the new type $t.T$ to α ,
- $t.Abs : \alpha \rightarrow t.T$ is injection from α to the new type $t.T$,
- $t.Rep'$ states that the property A holds for any term of type $t.T$,
- $t.Rep_inverse$ states that the injection of any element of type $t.T$ to α and back is equal to itself,
- $t.Abs_inverse$ states that if an element satisfies A , then injecting it to $t.T$ and back is equal to itself.

HOL-based proof assistants implement the type definition principle as a built-in statement. They also often provide further built-in statements for other definition principles that become derivable in the presence of type definitions, e.g., a definition principle for record types. For example, in Isabelle/HOL [NPW02], HOL is formalized in the Pure logic underlying the logical framework Isabelle [Pau94]. But because the type definition principle is not expressible in Pure, it is implemented as a primitive Isabelle feature that is only active in Isabelle/HOL.

5 Syntax Extensions and Surface Languages

Our definitions from Section 3 permit pragmatic *abstract* syntax, which is elaborated into strict abstract syntax. For human-oriented representations, it is desirable to complement this with similar extensions of pragmatic *concrete* syntax. While the pragmatic-to-strict translation at the abstract syntax level is usually non-trivial and therefore not invertible, the corresponding translation at the concrete syntax level should be compositional and bidirectional.

5.1 OMDoc Concrete Syntax for EL Declarations

First we extend OMDoc with concrete syntax that exactly mirrors the abstract syntax from Section 3. The declaration **extension** $e = \lambda x_1 : E_1. \dots \lambda x_n : E_n. \{\Sigma\}$ is written as

```
<extension name="e">
  <parameter name="x1"> $E_1$ </parameter>
  ⋮
  <parameter name="xn"> $E_n$ </parameter>
  <theory>
     $\Sigma$ 
  </theory>
</extension>
```

Here we use the box notation \boxed{A} to gloss the XML representation of an entity A given in abstract syntax.

Similarly, the pragmatic declaration **pragmatic** $c : e E_1 \dots E_n$ is written as

```
<pragmatic name="c" extension="⟨M⟩?e">
   $E_1 \dots E_n$ 
</pragmatic>
```

Here $\langle M \rangle$ is the meta-theory in which e is declared so that $\langle M \rangle ? e$ is the MMT URI of the extension.

Example 4. For the implicit definitions discussed in Section 3, we use the extension *impldef* from Figure 5, which we assume has namespace URI $\langle U \rangle$. If ρ is a proof of unique existence for an f such that $f' = f \wedge f(0) = 1$, then the exponential function is defined in XML by

```
<pragmatic name="exp" extension="⟨U⟩?ImplicitDefinitions?impldef">
   $\lambda f. f' = f \wedge f(0) = 1$   $\rho$ 
</pragmatic>
```

5.2 Pragmatic Surface Syntax

OMDoc is mainly a machine-oriented interoperability format, which is not intended for human consumption. Therefore, the EL-isomorphic syntax introduced is sufficient in principle – at least for the formal subset of OMDoc we have discussed so far.

OMDoc is largely written in the form of “surface languages” – domain-specific languages that can be written effectively and transformed to OMDoc in an automated process. For the formal subset of OMDoc, we use a MMT-inspired superset of the Twelf/LF [PS99; HHP93] syntax, and for informal OMDoc we use $\text{\S}\text{\TeX}$ [Koh08], a semantic extension of \TeX / \LaTeX .

For many purposes like learning the surface language or styling OMDoc documents, **pragmatic surface syntax**, i.e., a surface syntax that is closer to the notational conventions of the respective domain, has great practical advantages. It is possible to support, i.e., generate and parse, pragmatic surface syntax by using the macro/scripting framework associated with most representation formats.

For instance, we regained the XML syntax familiar from OMDoc 1.2 via notation definitions that transform between **pragmatic** elements and the corresponding OMDoc 1.2 syntax (see Appendix A). For Twelf/LF, we would extend the module system preprocessor, and for Isabelle we would extend the SML-based syntax/parsing subsystem. We have also extended $\text{\S}\text{\TeX}$ as an example of a semi-formal surface language (see Appendix B). Here we used the macro facility of \TeX as the computational engine. We conjecture that most practical surface languages for MKM can be extended similarly.

6 Conclusion & Future Work

In this paper, we proposed a general statement-level extension mechanism for MKM formats powered by the notion of theory families. Starting with MMT as a core language, we are able to express most of the pragmatic language features of OMDoc 1.2 as instances of our new extension primitive. Moreover, we can recover extension principles employed in languages for formalized mathematics including the statements employed for conservative extensions in Isabelle/HOL and Mizar. We have also described a principle how to introduce corresponding pragmatic concrete syntax. We exemplified this by defining a general mechanism how extension statements automatically introduce user-friendly concrete syntax in the $\text{\S}\text{\TeX}$ semantic authoring language.

The elegance and utility of the extension language is enhanced by the modularity of the OMDoc 2 framework, whose meta-theories provide the natural place to declare extensions: the scoping rules of the MMT module system supply the justification and intended visibility of statement-level extensions. In our examples, the Isabelle/HOL and Mizar extensions come from their meta-logics, which are formalized in MMT.

For full coverage of OMDoc 1.2, we still need to capture abstract data types and proofs; the difficulties in this endeavor lie not in the extension framework

but in the design of suitable meta-logics that justify them. For OMDoc-style proofs, the $\bar{\lambda}\mu\tilde{\mu}$ -calculus has been identified as suitable [ASC06], but remains to be encoded in MMT. For abstract data types we need a λ -calculus that can reflect signatures into (inductive) data types; the third author is currently working on this.

The fact that pragmatic extensions are declared in meta-theories points towards the idea that OMDoc metadata and the corresponding metadata ontologies [LK09] are actually meta-theories as well (albeit at a somewhat different level); we plan to work out this correspondence for OMDoc 2.

Finally, we observe that we can go even further and interpret the feature of definitions that is primitive in MMT as pragmatic extensions of an even more foundational system. Then definitions $c : E = E'$ become pragmatic notations for a declaration $c : E$ and an axiom $c = E'$, where $=$ is an extension symbol introduced in a meta-theory for equality. Typing can be handled similarly. This would also permit introducing other modifiers in declarations such as $<$: for subtype declarations.

References

- [ASC06] Serge Autexier and Claudio Sacerdoti Coen. “A Formal Correspondence Between OMDoc with Alternative Proofs and the $\bar{\lambda}\mu\tilde{\mu}$ -calculus”. In: *Mathematical Knowledge Management (MKM)*. Ed. by Jon Borwein and William M. Farmer. LNAI 4108. Springer Verlag, 2006, pp. 67–81.
- [Aus+10] Ron Ausbrooks et al. *Mathematical Markup Language (MathML) Version 3.0*. W3C Recommendation. World Wide Web Consortium (W3C), 2010. URL: <http://www.w3.org/TR/MathML3>.
- [BC04] Y. Bertot and P. Castéran. *Coq’Art: The Calculus of Inductive Constructions*. Springer, 2004.
- [Bus+04a] S. Buswell et al. *The Open Math Standard, Version 2.0*. Tech. rep. See <http://www.openmath.org/standard/om20>. The Open Math Society, 2004.
- [Bus+04b] Stephen Buswell et al. *The Open Math Standard, Version 2.0*. Tech. rep. The OpenMath Society, 2004. URL: <http://www.openmath.org/standard/om20>.
- [Chu40] A. Church. “A Formulation of the Simple Theory of Types”. In: *Journal of Symbolic Logic* 5.1 (1940), pp. 56–68.
- [Gor88] M. Gordon. “HOL: A Proof Generating System for Higher-Order Logic”. In: *VLSI Specification, Verification and Synthesis*. Ed. by G. Birtwistle and P. Subrahmanyam. Kluwer-Academic Publishers, 1988, pp. 73–128.
- [HHP93] R. Harper, F. Honsell, and G. Plotkin. “A framework for defining logics”. In: *Journal of the Association for Computing Machinery* 40.1 (1993), pp. 143–184.

- [Hor12] F. Horozal. “Logic Translations with Declaration Patterns”. <https://svn.kwarc.info/repos/fhorozal/pubs/patterns.pdf>. 2012.
- [IKR11] M. Iancu, M. Kohlhase, and F. Rabe. *Translating the Mizar Mathematical Library into OMDoc format*. Tech. rep. KWARC Report-01/11. Jacobs University Bremen, 2011.
- [KMR08] M. Kohlhase, C. Müller, and F. Rabe. “Notations for Living Mathematical Documents”. In: *Mathematical Knowledge Management*. Ed. by S. Autexier et al. Vol. 5144. Lecture Notes in Computer Science. Springer, 2008, pp. 504–519.
- [Koh06] Michael Kohlhase. *OMDOC – An open markup format for mathematical documents [Version 1.2]*. LNAI 4180. Springer Verlag, Aug. 2006. URL: <http://omdoc.org/pubs/omdoc1.2.pdf>.
- [Koh08] Michael Kohlhase. “Using L^AT_EX as a Semantic Markup Format”. In: *Mathematics in Computer Science 2.2* (2008), pp. 279–304. URL: <https://svn.kwarc.info/repos/stex/doc/mcs08/stex.pdf>.
- [LK09] Christoph Lange and Michael Kohlhase. “A Mathematical Approach to Ontology Authoring and Documentation”. In: *MKM/-Calcuemus Proceedings*. Ed. by Jacques Carette et al. LNAI 5625. Springer Verlag, July 2009, pp. 389–404. ISBN: 978-3-642-02613-3. URL: <http://kwarc.info/kohlhase/papers/mkm09-omdoc4onto.pdf>.
- [NPW02] T. Nipkow, L. Paulson, and M. Wenzel. *Isabelle/HOL — A Proof Assistant for Higher-Order Logic*. Springer, 2002.
- [ORS92] S. Owre, J. Rushby, and N. Shankar. “PVS: A Prototype Verification System”. In: *11th International Conference on Automated Deduction (CADE)*. Ed. by D. Kapur. Springer, 1992, pp. 748–752.
- [Pau94] L. Paulson. *Isabelle: A Generic Theorem Prover*. Vol. 828. Lecture Notes in Computer Science. Springer, 1994.
- [PS99] F. Pfenning and C. Schürmann. “System Description: Twelf - A Meta-Logical Framework for Deductive Systems”. In: *Lecture Notes in Computer Science* 1632 (1999), pp. 202–206.
- [Rab08] F. Rabe. *The MMT System*. see <https://trac.kwarc.info/MMT/>. 2008.
- [RK11] F. Rabe and M. Kohlhase. “A Scalable Module System”. see <http://arxiv.org/abs/1105.0548>. 2011.
- [TB85] A. Trybulec and H. Blair. “Computer Assisted Reasoning with MIZAR”. In: *Proceedings of the 9th International Joint Conference on Artificial Intelligence*. Ed. by A. Joshi. 1985, pp. 26–28.

A Pragmatic OMDoc Syntax

For implicit definitions we have the notation definition in listing 1.1 (slightly simplified elemental XML notation). Recall from [KMR08] that notation definitions are essentially transformation rules that use named `expr` elements as variables in the heads and `render` for the recursive calls of the transformation.

Listing 1.1. A Notation Definition for implicit Definitions

```
<notation format="OMDoc-pragmatic">
  <prototype>
    <pragmatic name="c" extension="⟨V⟩?ImplicitDefinitions?impldef">
      <expr name="prop"/>
      <expr name="proof"/>
    </pragmatic>
  </prototype>
  <rendering>
    <definition type="implicit" name="c" exunique="c.exu">
      <render name="prop"/>
    </definition>
    <assertion name="c.exu" type="theorem">
      <OMA>
        <OMS cd="ImplicitDefinitions" name="\implDef"/>
        <render name="prop"/>
      </OMA>
    </assertion>
    <proof for="c.exu"><render name="proof"/></proof>
  </rendering>
</notation>
```

Note that the “direction” of the notation definition may be confusing at first. But it becomes natural: Here the OMDoc 1.2 syntax can be regarded as a surface language into which the **statement** elements are presented into. But note that the notation definitions can also be “turned around” for parsing. We are currently working on parsing technology that integrates this idea.

B Pragmatic Concrete Syntax in \LaTeX

Now we show how to extend one of the OMDoc surface languages (formats intended for human authoring and transformation to OMDoc) with pragmatic concrete syntax. We have extended the \LaTeX format [Koh08] with an extension language, which we present here as an example how the ideas presented in this paper can be extended to semi-formal MKM formats. The next listing shows the extensions for implicit definitions. The `\extension` macro takes the name of the extension as a first required argument and then uses the second and third arguments to define an environment of that name in the usual way. The main problem in \TeX -based formats is argument passing; in our extension language

we use two ways. “Small” arguments are passed as key/value pairs; we use the `\extkey` macro to declare keys. “Large” arguments are naturally packaged in environments of their own which are declared as “continuations” by the `cont` key in the optional argument of **extension**.

Listing 1.2. Extending \LaTeX by Implicit Definitions

```
\begin{module}[id=ImplicitDefinitions]
\ldots
\extkey{idef}{name}
\extkey{idef}{title}
\extension{idef}
{\begin{def}[\namarg{idef}{title}]\label{idef.\namarg{idef}{name}}}}
{\end{def}}
\extkey{exu}{for}
\extension[cont=idef]{exu}
{\begin{thm}[Existence and Uniqueness for \ref{\namarg{exu}{for}}]
\label{\namarg{exu}{for}.exu}}
{\end{thm}}}
\extkey{exup}{for}
\extension[cont=idef]{exup}
{\begin{proof} (for \ref{\namarg{exu}{for}})}
{\end{proof}}}
\end{module}
```

Note the structural similarity of these extension declarations to the rendering element in Listing 1.1. Given this we can use the three induced environments once we have imported the module above.

```
\importmodule{ImplicitDefinitions}
\begin{idef}[name=exp,title=Exponential Function]
  The exponential function is that function  $f$  with  $f'=f$  and  $f(0)=1$ .
\end{idef}
\begin{exu}[for=exp]
  There is exactly one function  $f$  with  $f'=f$  and  $f(0)=1$ .
\end{exu}
\begin{exup}[for=exp]
  \ldots
\end{exup}
```

Note that the OMDoc generated from the can either be in the pragmatic concrete syntax from Section A or directly in the concrete syntax from Section 5.1.