

**Report on OpenDreamKit deliverable D6.5**  
**GAP/SAGE/LMFDB Interface Theories and**  
**Alignment in OMDoc/MMT for System**  
**Interoperability**

John Cremona, Michael Kohlhase, David Lowry-Duda,  
Dennis Müller, Markus Pfeiffer, Florian Rabe, Nicolas  
M. Thiéry, Tom Wiesing

**ABSTRACT.** There is a large ecosystem of mathematical software systems and knowledge bases. Individually, these are optimized for particular domains and functionalities, and together they cover many needs of practical and theoretical mathematics. However, each system specializes on one particular area, and it remains very difficult to solve problems that need to involve multiple systems. Some integrations exist, but they are ad-hoc and have scalability and maintainability issues. In particular, there is not yet an interoperability layer that combines the various systems into a virtual research environment (VRE) for mathematics.

The OpenDreamKit project aims at building a toolkit for such VREs. It suggests using a central system-agnostic formalization of mathematics (Math-in-the-Middle, MitM) as a mathematical pivot point for semantic-preserving translations in the needed interoperability layer. In this report, we report on a series of case studies that instantiates the MitM paradigm with the systems GAP, SAGEMATH, LMFDB, and SINGULAR to perform distributed computation in group, ring, and number theory.

Our work involves massive practical efforts, including a novel formalization of computational group theory, improvements and extensions of the involved software systems and knowledge bases, an extension of the underlying knowledge management system to cope with large theories, and a novel mediating system that sits at the center of a star-shaped integration layout between mathematical software systems and knowledge bases.

Together with deliverable report **D6.8**, this report describes the implementation and initial evaluation of the MitM integration and interoperability paradigm initially envisioned in deliverables **D6.2** and **D6.3**. The MitM paradigm constitutes the core development goal of **WP6** and the curated content described in this report enables running non-trivial integrations. In the future we expect further extending the reach of the integration in terms of both systems covered as well as knowledge available in the central MitM ontology.



|   |   |
|---|---|
| Due on  | 31/08/2017  |
| Delivered on  | 05/09/2018  |
| Lead  | Friedrich-Alexander Universität Erlangen/Nürnberg (FAU) |
| Progress on and finalization of this deliverable has been tracked publicly at:<br><a href="https://github.com/OpenDreamKit/OpenDreamKit/issues/139">https://github.com/OpenDreamKit/OpenDreamKit/issues/139</a> |   |

## Contents

|   |    |
|---|----|
| 1. Introduction   | 4  |
| 2. The Math-in-the-Middle Approach  | 7  |
| 2.1. The MitM Ontology  | 7  |
| 2.2. OMDoc/MMT as a Knowledge Representation Format                       | 8  |
| 2.3. Specifying System Dialects via System API Theories                   | 8  |
| 2.4. MitM-based Distributed Computation                                   | 9  |
| 2.5. Mathematical Knowledge Bases   | 10 |
| 3. The MitM Ontology  | 11 |
| 3.1. The MitM Foundation  | 11 |
| 3.2. Computational Group Theory   | 12 |
| 3.3. L-Functions and Modular Forms  | 13 |
| 4. Concrete Encodings of MitM Objects                                     | 14 |
| 5. Integrating Computation Systems with MitM: GAP, SAGEMATH, and SINGULAR | 17 |
| 5.1. SAGEMATH   | 17 |
| 5.2. GAP  | 19 |
| 5.3. SINGULAR   | 19 |
| 5.4. Alignments   | 20 |
| 6. Integrating Databases with MitM: The LMFDB Case Study                  | 22 |
| 6.1. LMFDB Overview   | 22 |
| 6.2. LMFDB as a Set of Virtual Theories                                   | 23 |
| 6.3. Ascribing Encodings in Schema Theories                               | 25 |
| 6.4. Translating Queries  | 26 |
| 7. MitM-Based Distributed Computation                                     | 28 |
| 8. Conclusion   | 30 |
| Bibliography  | 33 |

## 1. Introduction

**Motivation.** There is a large and vibrant ecosystem of open-source software systems for mathematics. These range from calculators, which perform simple computations, via mathematical databases, which curate collections of mathematical objects, to powerful modeling tools and computer algebra systems (CAS).

These systems can be very specific, often focusing on a narrow area of mathematics. For example, among databases, the “Online Encyclopedia of Integer Sequences” (OEIS) focuses on sequences over  $\mathbb{Z}$  and their properties, and the “L-Functions and Modular Forms Database” (LMFDB) [Cre16; LMFDB] on objects in number theory pertaining to Langland’s program. Among CAS, GAP [GAP] excels at discrete algebra with a focus on group theory, SINGULAR [SNG] focuses on polynomial computations with special emphasis on commutative and non-commutative algebra, algebraic geometry, and singularity theory. Finally, SAGEMATH [Sage] aims to be a general purpose software for computational pure mathematics by integrating many systems including the aforementioned ones, together with a large body of code for the SAGEMATH library itself written in Python.

For a mathematician, however, (a user, which we call Jane) the systems themselves are not relevant. Instead, she only cares about being able to solve problems. Because it is typically not possible to solve a mathematical problem within a single system, Jane has to work with multiple systems and combine the results to reach a solution. Currently there is very little tool support for this in practice; so Jane has to isolate sub-problems that the respective systems are amenable to, formulate them in the respective input language, collect intermediate results, and reformulate them for the next system — a tedious and error-prone process at best, a significant impediment to scientific progress at worst. Solutions for some situations certainly exist, which can help get Jane unstuck, but these are ad-hoc and only for specific often-used system combinations. Moreover, each of these ad hoc solutions requires a lot of maintenance and scales badly to multi-system integration. To add insult to injury, the knowledge bases Jane would like to use — ranging from Wikipedia to theorem prover libraries — are usually only accessible via the restricted API of a dedicated web information system or the low-level API to the raw database content. What we would want is a “programmable, mathematical API” which would give access to the knowledge-bases programmatically via their mathematical constructions and properties.

**Related Work.** There are essentially two approaches to tackling the problem of system interoperability and integration. Firstly, **ad-hoc-integration** uses a programming language for translating object data structures between systems, possibly using shared memory. The most prominent proponent is the SAGEMATH system, which uses Python as the glue language and various Python-to-X bridges for Master/Slave-type integration. Here, the “preservation of semantics” property that is so important for a sound system integration is implicitly embedded in the glue code and has to be painstakingly maintained by the community over system releases. Recently, the OSCAR project [OC] was started to achieve a similar but much deeper integration of a fixed set of four computer algebra systems using Julia as the glue language.

Secondly, **middleware-integration** uses a dedicated formal language and protocol for the standardized representation and communication of objects. Outside the domain of mathematics, this has been successfully applied in, e.g., CORBA [Gro], which uses object-oriented modeling and stub generation for the exchange of business objects, REST-based Web Services [WSDL07; Mit03] that exchange XML or JSON-based objects, or more recently Protocol Buffers [PB], which model structured data as JSON-like text files. Middleware integration excels at integrating disparate systems across machine, programming language, and even operating system borders. It works best if the respective object models are compatible, which the respective communities try to achieve by standardizing domain ontologies and schemata. This schema standardization

and mediation between differing object models turns out to be the main problem involved in interoperability.

However, in mathematics, domain model standardization is inherently difficult: The set of potential object *types* is dynamically extensible and thus itself part of the standardization problem. Moreover, essentially no object type admits canonical representations for its objects so that different-but-equivalent constructions must be supported in practice. Middleware integration for mathematics was attempted using the OpenMath [Bus+04] and content MathML [Aus+10] languages and the SCSCP protocol [Fre+]. Their underlying object models are isomorphic and represent mathematical objects as syntax trees with binding starting from symbols, variables, and literals. Extensibility is achieved by using an open-ended set of symbols that are introduced and specified in special documents called content dictionaries (CDs). Due to this extensibility, any interoperability that preserves the meaning of the communicated objects hinges on the availability of standardized CDs, which have to be manually written by the respective community. In the past, the main problem of such math middleware approaches has been the lack of high-impact CDs that both fully specified the semantics of symbols and were supported by concrete systems in a way that respected these specifications.

**Towards a Virtual Research Environment for Mathematics.** A goal of the OpenDreamKit project is tackling the integration problem for mathematical systems systematically by building virtual research environments (VRE) on top of the existing systems. This requires a joint user interface — the OpenDreamKit project adopts Jupyter [Jup] and active documents [Koh+11] — and an interoperability layer that allows passing problems and results between the disparate systems. For the latter, OpenDreamKit explores both ad-hoc (via Sage) and (in this deliverable) middleware integration.

We have dubbed OpenDreamKit’s middleware approach **Math-in-the-Middle (MitM)** [Deh+16]. It provides an interoperability framework based on a central, system-independent **ontology** of mathematical knowledge and on system **API theories** that specify the interfaces of the various systems. Here the system API theories of a systems  $S$  are system-*near* in the sense that they introduce symbols that correspond exactly to the functionality implemented by  $S$ . But they are system-independent in the sense that they are separate from the source code of  $S$ . Thus, they can be seen as a formal and thus machine-actionable documentation of the interface of  $S$ .

We use the OMDoc/MMT [Koh06; RK13] language as the representation format for both the ontology and the system API theories. And we use **alignments** [Mül+17b; Mül+17a] to describe the relations between the symbols declared in the ontology and those in the API theories. This allows building math middleware on top of the MMT system [MMT].

**Contribution.** In this report we instantiate the MitM paradigm in two concrete case studies. In the first one, we show distributed computation involving the GAP, SAGEMATH, and SINGULAR systems. In the second one, we show the integration of the mathematical knowledge base LMFDB into MitM-based computation.

For the CAS case study, we will use the following running example from computational group theory: Jane wants to experiment with invariant theory of finite groups. She works in the polynomial ring  $R = \mathbb{Z}[X_1, \dots, X_n]$ , and wants to construct an ideal  $I$  in this ring that is fixed by a group  $G \leq S_n$  acting on the variables, linking properties of the group to properties of  $I$  and the quotient of  $R$  by  $I$ .

To construct an ideal that is invariant under the group action, it is natural to pick some polynomial  $p$  from  $R$  and consider the ideal  $I$  of  $R$  that is generated by all elements of the orbit  $O = \text{Orbit}(G, R, p) \subseteq R$ . For effective further computation with  $I$ , she needs a Gröbner base of  $I$ .

Jane is a SAGEMATH user and wants to receive the result in SAGEMATH, but she wants to use GAP's orbit algorithm and SINGULAR's Gröbner base algorithm, which she knows to be very efficient. For the sake of example, we will work with  $n = 4$ ,  $G = D_4$  (the dihedral group), and  $p = 3 \cdot X_1 + 2 \cdot X_2$ , but our results apply to arbitrary values. Notably, the this dihedral group is called  $D_4$  in SAGEMATH but  $D_8$  in GAP due to differing conventions in different mathematical communities — a small example of the obstacles to system interoperability that MitM tackles.

For the LMFDB case study, Joanna (a colleague of Jane's) wants to investigate the number fields which are generated by the coefficients of Hilbert modular forms (HMFs). For many totally real number fields  $F$  of low degree, and for many levels  $\mathcal{N}$  for each field, the LMFDB contains information about all HMFs of level  $\mathcal{N}$  (of parallel weight 2 and trivial character). Each of these HMFs is an eigenform for the Hecke algebra with eigenvalues generating a number field; the same number field contains the coefficients of the standard Fourier expansion of the HMF, which are expressible in terms of the eigenvalues. In the LMFDB, each HMF's Hecke field  $K$  is stored by means of a defining polynomial which has been obtained as a by-product of the computation of the HMF itself, and is in no way canonical or minimal, making study of these fields difficult and — even in simple cases — obscure. For example, the Hecke field  $K = \mathbb{Q}(\sqrt{2})$  may occur for more than one HMF, defined by the polynomial  $x^2 - 2$  for some and by the polynomial  $x^2 - 2x - 1$  for others. Hence Jane would like to be able to extract these defining polynomials from the LMFDB, use them to define number fields in SAGEMATH, find simpler polynomials defining the same fields, and study their arithmetic properties (for example, their class numbers). To this end, some of the Hecke fields may themselves be in the LMFDB's collection of number fields, in which case the information about them which Jane needs is already computed and stored, but this will in most cases be hidden since the defining polynomials used in the HMF database will often not be the one stored in the number fields database.

**Overview.** In Section 2, we recap the MitM paradigm. In Section 3, we describe the MitM ontology. While the ontology describes the *abstract* syntax of mathematical objects, Section 4 introduces a codec framework for describing their *concrete* syntax in different systems. In Section 5 and 6, we describe the integration of the computation systems GAP, SAGEMATH, and SINGULAR and the LMFDB databases with the MitM architecture. In Section 7, we present the resulting virtual research environment built on these systems in action. Section 8 concludes the report.

## 2. The Math-in-the-Middle Approach

Figure 1 shows the basic MitM design. We want to make the systems  $A$  to  $H$  with system dialects  $a$  to  $h$  interoperable. A P2P translation regime ( $n(n-1)$  translations between  $n$  systems) is already intractable for the systems in the OpenDreamKit project (more than a dozen). Alternatively, an “industry standard” regime, where one system dialect is declared as the standard is infeasible because no system dialect subsumes all others — not to mention the political problems such a standardization would induce. Instead, MitM uses a central mathematical ontology that provides an independent mediating language, via which all participating systems are aligned. All mathematical knowledge shared between the systems and exposed to the high-level VRE user is expressed using the vocabulary of this ontology. Crucially, while every system dialect makes implementation-driven, system-specific design choices, the MitM ontology can remain close to the knowledge published in the mathematical literature, which already serves as an informal interoperability layer.

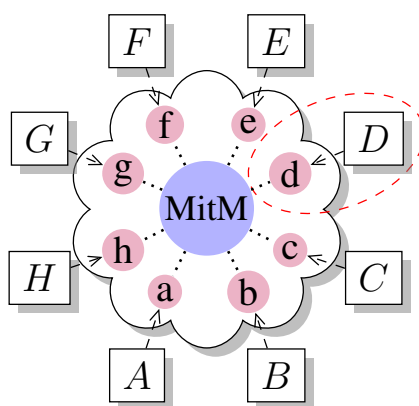


FIGURE 1. MitM Paradigm

The following sections describe the three components of the MitM paradigm in more detail.

### 2.1. The MitM Ontology

In the center, we have the **MitM Ontology**, which is a formalization of the mathematical knowledge behind the systems  $A$  to  $H$ . As a formalization framework, it uses the OMDOC/MMT format, which was designed with this specific application in mind. We do not go into the details of OMDOC/MMT here (but see Section 2.2) – for our purposes, it suffices to assume that an OMDOC/MMT theory graph formalizes a language for mathematical objects as a set of typed symbols with a (formal or informal) specification of their semantics. For example, the MitM-symbol `PolynomialRing` takes a ring  $r$  of coefficients and a number  $n$  of variables and returns the ring  $r[X_1, \dots, X_n]$  of polynomials.

Note that the purpose of the MitM ontology is not the formal verification of mathematical theorems (as for most existing formalizations such as [Gon+13] of group theory), but to act as a pivot point for integrating systems. This means that it can be much nearer to the informal but rigorous presentation of mathematical knowledge in the literature. While each system dialect makes compromises and optimizations needed for a particular application domain, the MitM ontology follows the existing and already informally standardized mathematical knowledge and can thus serve as a standard interface layer between systems.

Importantly, the MitM ontology does not have to include any definitions<sup>1</sup> or proofs – it only has to declare the types of all relevant symbols and state (but not prove) the relevant theorems. This makes it possible for users like Jane to extend the MitM ontology quickly whereas extending formalizations usually requires extensive efforts by specialists.

<sup>1</sup>Of course, definitions are one possible way to specify the semantics of MitM-symbols.



## 2.2. OMDoc/MMT as a Knowledge Representation Format

The realization of the MitM approach crucially depends on the information architecture of the OMDoc/MMT language [Koh06; RK13] and its implementation in the MMT system [Rab13; MMT].

In OMDoc/MMT knowledge is organized in **theories**, which contain information about mathematical concepts and objects in the form of **declarations**. Theories are organized into an “object-oriented” inheritance structure via **inclusions** and **structures** (for controlled multiple inheritance), which is augmented via truth-preserving mappings between theories called **views**, which allow to relate concepts of pre-existing theories and transport theorems between these. Inclusions, structures, and views impose a graph structure on the represented mathematical knowledge, called a **theory graph**.

We observe that even very large mathematical knowledge spaces about abstract mathematical domains can be represented by small, but densely connected, theory graphs, if we make all inherited material explicit in a process called **flattening**. The OMDoc/MMT language provides systematic names (MMT URIs) for all objects, properties, and relations in the induced knowledge space, and given the represented theory graph, the MMT system can compute them on demand.

Generally, knowledge in a knowledge space given by a theory graph loaded by the MMT system can be accessed by either giving it’s MMT URI, or by uniquely describing it via a set of conditions. To achieve the latter, MMT has a Query Language called QMT [Rab12], which allows even complex conditions to be specified. Currently, the MMT system loads the theory graph into main memory at startup and interleaves incremental flattening and query evaluation operations on the MMT data structures until the result has been produced.

In this report we show that the MitM approach, its OMDoc/MMT-based realization, and distribution via the SCSCP protocol are sufficient for

- distributed, federated computation between multiple computer algebra systems (Sage, GAP, and Singular; see Section 7),
- providing programmatic, mathematics-level system interfaces for mathematical knowledge bases (LMFDB; see Section 6),

and that the MitM ontology of abstract group theory can be represented in OMDoc/MMT efficiently. This setup is effective because

- the knowledge spaces behind abstract and computational mathematics can be represented in theory graphs very space-efficiently: The compression factors between a knowledge space and its theory graph – we call it the **TG factor** – exceeds two orders of magnitude even for small domains.
- only small parts of the knowledge space are traversed for a given computation.

## 2.3. Specifying System Dialects via System API Theories

**System Dialects.** It is unavoidable that each system induces its own language for mathematical objects. This is the cause of much incompatibility because even subtle differences make naive integration impossible. Moreover, due to the difficulty of the involved mathematics and the effort of maintaining the implementations, such differences are aplenty.

Fortunately, we can at least easily abstract from the user-facing surface syntax of these languages: scalable interoperability can anyway only be achieved by acting on the internal data structures of the systems. Thus, only the much simpler internal abstract syntax needs to be considered. In particular the OpenMath/cMathML model for objects-as-formulae is sufficient. As OMDoc/MMT uses OpenMath at the level of domain objects and statements/declarations, it is a very good fit; but as we have seen in the discussion in the introduction, we need a set of content dictionaries that cover the domain of applicability of the respective system. OMDoc/MMT represents content dictionaries, as theories and dependencies and translations between content dictionaries – a feature that OpenMath does not have – as theory morphisms. We call a graph



of OMDoc/MMT content dictionaries that define the input language of a system (the abstract syntax trees encoded as OpenMath objects whose symbols are declared in these CDs) **system API theories**.

The symbols that build the abstract syntax trees can be split into two kinds: **constructors** build primitive objects without involving computation, and **operations** compute objects from other objects (including predicates, which we see as operations that return booleans). For purposes of interoperability it is desirable to abstract from this distinction and consider both as typed symbols. This abstraction is important because systems often disagree on the choice of constructors. Thus, we can represent the interfaces of the systems  $A$  to  $H$  as OMDoc/MMT theory graphs  $a$  to  $h$  that declare the constructors and operations (but omit all implementations of the operations) of the respective system.

Given the theory graph  $a$  representing the system dialect of  $A$ , we can express all objects in the language of system  $A$  as OMDoc/MMT objects using the symbols of  $a$ . We refer to these objects as  $A$ -objects. It is conceptually straightforward to write (or even automatically generate) the theory graph  $a$  and to implement a serializer and parser for  $A$ -objects as a part of  $A$ .<sup>2</sup> This is because no consideration of interoperability and thus no communication with the developers of other systems is needed.

The MitM approach stipulates that interface theories and interface views are maintained and released together with the respective systems, whereas the core MitM ontology represents the mathematical scope of the VRE and is maintained with it. In fact in many ways, the core MitM ontology is the conceptual essence of the mathematical VRE.

**Alignments with the Ontology.** The above reduces the interoperability problem to relating each system dialect to the MitM ontology. Each system dialect overlaps with the language of the ontology, but no system implements all ontology symbols and every system implements idiosyncratic operations that are not useful as a part of the ontology. Therefore, some system dialect symbols are related to corresponding symbols in the MitM ontology. We use these symbols of the MitM ontology as an intermediate representation to bridge between any two systems, e.g., by translating  $A$ -objects to the corresponding ontology objects and then those to the corresponding  $B$ -objects. This can be done using OMDoc/MMT **alignments** (see Section 5.4)

## 2.4. MitM-based Distributed Computation

The final missing piece for a system interoperability layer for a VRE toolkit is a practical way of transporting (OpenMath) objects between systems. This requires two steps.

Firstly, if the system dialects and alignments are known, we can automatically translate  $A$ -objects to  $B$ -objects in two steps:  $A$  to ontology and ontology to  $B$ . This two-step translation has been implemented in [Mül+17a] based on the MMT system [Rab13; MMT], which implements the OMDoc/MMT format along with logical and knowledge management algorithms.

Secondly, each system  $A$  has to be able to serialize/parse  $A$ -objects and to send them to/receive them from MMT. In the OpenDreamKit project we use the OpenMath SCSCP (Symbolic Computation Software Composability) protocol [Fre+] for that. SCSCP is essentially a distributed remote-procedure-call system based on OpenMath, which is itself the fragment of OMDoc/MMT used for representing objects. It is straightforward to extend a parser/serializer for  $A$ -objects to an SCSCP clients/server by implementing the SCSCP protocol on top of, e.g., sockets or using an existing SCSCP library.

<sup>2</sup>However, as we see below, this may still be surprisingly difficult in practice.

## 2.5. Mathematical Knowledge Bases

But the OpenDreamKit VRE must also include mathematical data sources, which curate the objects, models, examples, and counterexamples of mathematics. There are various large-scale sources of mathematical knowledge. These include

- generic information systems like Wikipedia,
- collections of informal but rigorous mathematical documents – e.g. research libraries, publisher’s “digital libraries”, or the Cornell preprint arXiv,
- literature information systems like zbMATH or MathSciNet,
- databases of mathematical objects – like the GAP group libraries, the Online Encyclopedia of Integer sequences (OEIS [[Slo03](#); [OEIS](#)]), and the L-Functions and Modular Forms Database (LMFDB [[Cre16](#); [LMFDB](#)]),
- fully formal theorem prover libraries like those of Mizar, Coq, PVS, and the HOL systems.

We will use the term **mathematical knowledge bases** to refer to them collectively and restrict ourselves to those that are available digitally. They are very useful in mathematical research, applications, and education. Commonly these systems are only accessible via a dedicated web interface that allows humans to query or browse the databases. A programmatic interface, if it exists at all, is usually system-specific, to use it, users need to be familiar both with the mathematical background and internal structure of the system in question. No predominant standard exists, and these interfaces usually only expose the low-level raw database content. We claim that mathematicians and other scientists desire a “programmatic, mathematical API” that gives access to the knowledge-bases programmatically via their mathematical constructions and properties. We focus on addressing this problem in the OpenDreamKit project (see Section 6.1)

For our implementation we interpret mathematical knowledge bases as OMDOC/MMT theory graphs – modular, flexi-formal representations of mathematical objects, their properties, and relations. This embedding gives us a common conceptual framework to handle different knowledge sources, and the modular and heterogeneous nature of OMDOC/MMT theory graph can be used to reconcile differing ontological commitments of the knowledge sources with in this conceptual framework.

But knowledge sources like the LMFDB or the OEIS contain millions of mathematical objects. For such knowledge sources, the classical MMT system is not yet suitable:

- the knowledge space corresponding to the data base content cannot be compressed by “general mathematical principles” like inheritance. Indeed, redundant information is already largely eliminated by the data base schema and the “business logic” of the information system it feeds.
- typically large parts of the knowledge space need to be traversed to obtain the intended results to queries.

Therefore, we extend the concept of OMDOC/MMT theories – which carry the implicit assumption of containing only a small number of declarations (see [[FGT92](#)] for a discussion) – to **virtual theories**, which can have an unlimited (possibly infinite) number of declarations. To contrast the intended uses we will call the classical OMDOC/MMT theories **concrete theories**. In practice, a virtual theory is represented by concrete approximations: OMDOC/MMT works with a concrete theory, whose size changes dynamically as a suitable backend infrastructure generates declarations on demand. Thus, from the system perspective, virtual theories behave just like concrete theories but without the assumption that all declarations are loaded at once; instead, declarations are loaded lazily.

We also update the knowledge management algorithms in the MMT system so that they can directly deal with the databases underlying the knowledge bases. Here we provide a systematic solution for encoding/decoding between low-level representations in standard databases and high-level mathematical representations.

### 3. The MitM Ontology

Jane’s use case involves groups and actions, polynomials, rings and ideals, and Gröbner bases and Joanna’s in addition elliptic curves, Hecke algebras, and modular forms, all of which must be formalized in the MitM ontology. First we discuss the MitM Foundation — i.e. the language used for specifying the ontology — in Section 3.1. Then we describe two ontology fragments as examples in Section 3.2 and 3.3. The first is a manual formalization for computational group theory (CGT) that is needed in particular for Jane’s use case. The latter is an informal ontology fragment for the part of number theory covered by Langland’s program; it is needed for integrating and generated from the LMFDB system (see Section 6).

#### 3.1. The MitM Foundation

OMDOC/MMT formalizations must be relative to a foundational logic, which is itself formalized in OMDoc/MMT. As a foundation for all formalizations in MitM [Mit], we use a polymorphic dependently typed  $\lambda$ -calculus with two universes type and kind (roughly analogous to sets and proper classes in set theory) and subtyping. It provides dependent function types  $\{a:A\}B(a)$ , representing the type of all functions mapping an argument  $a:A$  to some element of type  $B(a)$ . If  $B$  does not depend on the argument  $a$ , we obtain the simple function type  $A \rightarrow B$ .

For formulas, we use a type prop and a higher order logic where quantifiers range over any type. We furthermore follow the judgments-as-type paradigm by declaring a function  $\vdash:\text{prop} \rightarrow \text{type}$  mapping propositions to the **type of their proofs**, which allows us to declare proof rules as functions mapping proofs (of the premises) to a proof (of the conclusion).

The judgment  $A <: B$  expresses that  $A$  is a subtype of  $B$ . We use power types (the type of subtypes of a type) and predicate subtyping  $\{a:A \mid P(a)\}$ . The latter makes type-checking undecidable, but that is necessary for natural formalizations in many areas of mathematics.

```
theory Loop : base:?Logic =

  theory loop_theory : base:?Logic =
    | include ?Unital/unital_theory |
    | inverse : U → U | # 1-1 prec 24 |
    | axiom_inverse : ⊢ prop_inverse op (inverse) e |
  loop = Mod loop_theory | role abbreviation |

  inverseOf : {G : loop} dom G → dom G | # 2-1 prec 25 |
  | = [G] [a] (G.inverse) a |

theory Group : base:?Logic =

  theory group_theory : base:?Logic =
    | include ?Monoid/monoid_theory |
    | include ?Loop/loop_theory |
  group = Mod group_theory | role abbreviation |

  automorphisms : group → ℕ |
  cyclic : group → prop |
  degree : group → ℕ |
  order : group → ℕ |
```

FIGURE 2. MitM Ontology Fragment

A critical question is what additional types to provide. To be practical at all, this must include basic types for aggregation (e.g., records) and collection (e.g., lists). Based on a survey of practical needs, we compiled such a set in OpenDreamKit deliverable **D6.8**. [D6.818]

Moreover, we need an open-ended set of types for mathematical structures such as fields and rings. Here we have developed a novel type constructor [MRK18] that turns any MMT theory into the corresponding dependent record type. This combines the benefits of MMT’s axiomatic theories and a flexible type system. Concretely, for a theory  $T$ , the type  $\text{Mod } T$  is the record type of models of  $T$ . For example, in Figure 2 we define groups in terms of an MMT theory `group_theory` that declares the operations and axioms in a way that corresponds to the mathematical definition of groups. Then `group = Mod group_theory` is the type of all models of that theory, i.e., the type of groups. Any element  $g : \text{group}$  thus represents an actual group, whose operations and axioms can be accessed via record field projections (e.g. `g.inverse` yields the inverse operation of  $g$  (imported from loops via the include `?Loop/loop_theory-statement`). Since axioms are turned into record type fields as well, actually constructing a record of type `group` corresponds to proving that the field universe and the operations provided in the record do in fact form a group.

### 3.2. Computational Group Theory

Relative to the foundation, we can formalize the actual mathematical knowledge we are interested in. As a running example, we describe a formalization of computational group theory (CGT), which is inspired by the corresponding implementation in GAP. It uses several different levels of abstraction – currently *abstract*, *representation*, *implementation*, and *concrete*. From our experience, we expect this pattern to be applicable across computational algebra, possibly with additional levels of abstraction. The left box in Figure 3 shows the levels and their relation to the constructors and operations of GAP.

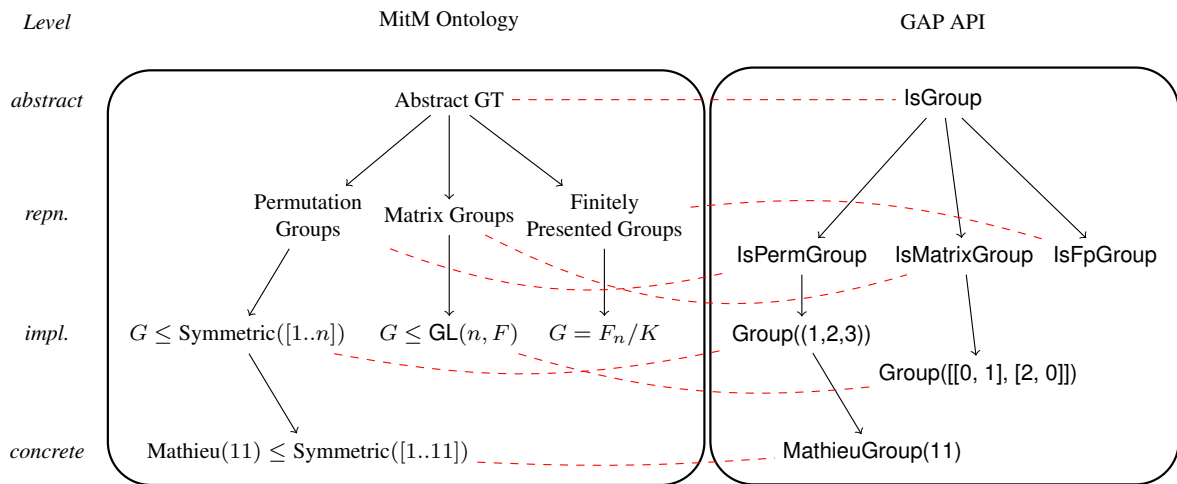


FIGURE 3. Alignments between the MitM Ontology and the GAP API

**Abstract Level.** This contains the abstract, axiomatic theory of *Groups*: the group axioms, generating sets, homomorphisms, group actions, stabilisers, and orbits. This also easily leads into definitions of centralisers – i.e. stabilisers of elements under conjugation – and normalisers – i.e. stabilisers of subgroups under conjugation, stabiliser chains, Sylow- $p$  subgroups, Hall subgroups, and many other concepts.

OMDOC/MMT also allows expressing that there are different equivalent definitions of a concept: We defined group actions in two ways and used *views* to express their equivalence.

**Representation Level.** Abstract groups are represented in different ways as concrete objects suitable for computation: as groups of permutations, groups of matrices, finitely presented groups, algebraic constructions of groups, or using polycyclic presentations.

Additionally, mathematicians often compute with canonical representatives of an isomorphism class of groups: When group theorists talk about the “Dihedral group of order 8”, they often have a particular representation in mind, for example as a group that acts on the square by rotations and reflections. In GAP this group would be represented as a group of permutations, or by a polycyclic presentation.

Many representations arise naturally from *group actions*: If we are considering symmetry in a setting where we want to apply group theory, we start with a group action, for example a group acting on a graph by permuting its vertices.

The universal tool to bridge the gap between groups, representations and canonical representatives are *group homomorphisms*, particularly embeddings and isomorphisms, which are used extensively in GAP. This is reflected in our approach.

**Implementation Level.** At this level we encode implementation details: Permutation groups in GAP are considered as finite subgroups of the group  $S_{\mathbb{N}+}$ , and defined by providing a set of generating permutations. GAP then computes a stabiliser chain for a group that was defined this way, and naturally considers the group to be a subgroup of  $S_{[1..n]}$ , where  $n$  is the largest point moved.

**Concrete Level.** It is at the concrete level where the computation happens: while the higher levels are suitable for mathematical deduction and inference, this level is where GAP (or any other system providing computational group theory) does its work. If a group (or a group action) has been constructed by giving generators through MitM, GAP can now compute the size of the group, its isomorphism type, and perform all the other operations that are available via the GAP system dialect.

### 3.3. L-Functions and Modular Forms

Since 2011, the “L-Functions and Modular Forms Database” (LMFDB) [Cre16; LMFDB] documents the system functionality in a system of “knowls” (knowledge items), which contain short informal definitions of the mathematical concepts and information about their relation to others. They were originally designed to make the LMFDB user interface self-contained and self-documenting. There are currently 945 knowls in LMFDB, which describe about as many mathematical concepts (some knowls introduce multiple concepts, some are non-mathematical).

Figure 4 shows the knowl for a “Hecke Algebra”, where another knowl for the “weight” of a Bianchi modular form is displayed inline (the typical mode of interaction with knowls). In LMFDB, the set of knowls can be downloaded from the system API as JSON where the text is encoded as HTML with  $\LaTeX$  formulae and knowl references as custom markup.

As knowls are informal but have some structure, we can represent them elegantly in OMDOC/MMT using the  $\S\TeX$ -encoding [Koh08; sTeX], a variant of  $\TeX/\LaTeX$  that uses special  $\TeX$  macros for marking up OMDOC/MMT structure. We have implemented a converter that generates  $\S\TeX$  theories from LMFDB knowls.





LMFDB [Knowledge](#) → [Hecke algebra](#) [Feedback](#) · [Hide Menu](#)

## Hecke algebra

[show](#) · [mf.bianchi.hecke\\_algebra](#) · [all knows](#) · [up](#) · [search:](#)  [go](#)

Let  $K$  be an imaginary quadratic field and  $\mathcal{M}_k(\Gamma)$  the space of [Bianchi modular forms](#) of [weight](#)  $k$  and [level](#)  $\Gamma = \Gamma_0(\mathcal{N})$  for some integral ideal  $\mathcal{N}$  of  $\mathcal{O}_K$ .

[Weight of a Bianchi modular form](#)

The **weight** of a Bianchi modular form is the positive integer  $k$  occurring in the transformation formula. Also, Bianchi modular forms are vector-valued functions  $F: \mathcal{H}_3 \rightarrow \mathbb{C}^{k+1}$  where  $\mathcal{H}_3$  is Hyperbolic 3-space, so for example a Bianchi modular form of weight 2 takes values in  $\mathbb{C}^3$ .

[permalink](#)

The (level  $\mathcal{N}$ ) **Hecke algebra**  $\mathbb{T}$  is a commutative algebra of linear **Hecke operators** acting on  $\mathcal{M}_k(\Gamma)$ . It preserves the cuspidal subspace  $\mathcal{S}_k(\mathcal{N})$  of Bianchi cusp forms and the [new subspace](#)  $\mathcal{S}_k(\mathcal{N})^{\text{new}}$ , and also preserves the rational and integral structures on these.

The Hecke algebra is generated by operators  $T_p$  for primes  $p$  of  $K$ . Each  $T_p$  is self-adjoint on  $\mathcal{S}_k(\mathcal{N})$  and hence has eigenvalues which are totally real algebraic integers. The new space  $\mathcal{S}_k(\mathcal{N})^{\text{new}}$  has a basis consisting of forms which are eigenforms for all Hecke operators.

**Authors:**

- [John Cremona](#)

**Knowl status:**

- Review status: beta
- Last edited by John Cremona on 2017-07-14 13:02:15.370000

FIGURE 4. Knowls in the LMFDB User Interface

## 4. Concrete Encodings of MitM Objects

When integrating systems with the star-shaped MitM architecture, some translation of concrete formats is necessary. While this is not surprising, it leads to an important difference between the integration of computation systems and knowledge bases: the former but not the latter include a programming environment that provides all necessary infrastructure for implementing the reformatting. Therefore, to integrate with databases, it is convenient to standardize some encodings that translate between high-level datatypes in the MitM ontology and concrete representations that can be sent to and received from databases. Even though the concepts and implementations are universally applicable, we will use the LMFDB as a concrete motivating setting.

This is particularly critical as the databases used for the scalable physical storage of large datasets usually offer only very simple data structures. For example, a JSON database (as underlies LMFDB) offers only limited-precision integers, boolean, strings, lists, and records as primitive objects. An SQL database offers only records of basic objects. Neither provides a type system. Consequently, the objects stored in the database are very different from the sophisticated mathematical objects expected by the mathematical software systems in the OpenDreamKit VRE toolkit.

Therefore, databases like LMFDB must encode this complex mathematical objects as simple database objects. Consider, for example, the **degree** of an elliptic curve (as we will in Section 6). Its *semantic* type is  $\mathbb{Z}$ , but its *physical* type in LMFDB is IEEE 754 a mixture of 64-bit floating point numbers and strings: integers that exceeds  $2^{53} - 1$  are stored as JSON strings containing the corresponding decimal representation.

To formally specify these encodings codecs, we introduce a new OMDoc/MMT theory *Codecs* as a part of the MitM ontology. Our codecs are indexed by semantic types: the type constructor *codec* maps a semantic type to a new type of codecs for it. For instance, the

object `StandardInt` of type `codec  $\mathbb{Z}$`  is a codec that translates between LMFDB’s idiosyncratic float/string-representation and MitM’s integers. Note that there can be multiple different codecs for the same semantic type. For example, `IntAsArray` encodes integers  $x$  as lists of 64-bit integers consisting of the digits of  $x$  with respect to base  $2^{64}$ . Figure 5 shows a collection of atomic codecs useful in the LMFDB context.

| Codecs                      |  |   |
|-----------------------------|--|---|
| codec                       | : type $\rightarrow$ type                      |   |
| <code>StandardPos</code>    | : <code>codec <math>\mathbb{Z}^+</math></code> | JSON number if small enough,<br>else JSON string of decimal expansion |
| <code>StandardNat</code>    | : <code>codec <math>\mathbb{N}</math></code>   |   |
| <code>StandardInt</code>    | : <code>codec <math>\mathbb{Z}</math></code>   |   |
| <code>IntAsArray</code>     | : <code>codec <math>\mathbb{Z}</math></code>   | JSON List of Numbers  |
| <code>IntAsString</code>    | : <code>codec <math>\mathbb{Z}</math></code>   | JSON String of decimal expansion                                      |
| <code>StandardBool</code>   | : <code>codec <math>\mathbb{B}</math></code>   | JSON Booleans   |
| <code>BoolAsInt</code>      | : <code>codec <math>\mathbb{B}</math></code>   | JSON Numbers 0 or 1   |
| <code>StandardString</code> | : <code>codec <math>\mathbb{S}</math></code>   | JSON Strings  |

FIGURE 5. Some Codecs specified in MMT ( $\mathbb{N}$ ,  $\mathbb{Z}$ ,  $\mathbb{Z}^+$  are as usual,  $\mathbb{B}$  are booleans, and  $\mathbb{S}$  are Unicode strings)

We do not (and do not have to) define the actual encoding/decoding functions in OMDOC/MMT. It is more important to identify the codecs needed in practice, introduce names for them, and spell out their semantics. Then it is straightforward to implement them in any other programming language used interfacing with LMFDB.

Concretely, we have implemented them in Scala, the language underlying the MMT system. Additionally, the `Codecs` theory annotates each codec declaration with a reference to the Scala class implementing the codec. That way, MMT can run the encoding/decoding functions of the codec.

The above is only sufficient for atomic semantic types, which typically correspond to one (or more) atomic codecs. Consider now the field `isogeny_matrix` of elliptic curves. The semantic representation of one possible value (namely for the curve 11a1 ) of this field is the matrix on the right.

$$M = \begin{pmatrix} 1 & 5 & 25 \\ 5 & 1 & 5 \\ 25 & 5 & 1 \end{pmatrix}$$

The semantic type operator `Matrix` takes one type argument (the element type, integers in this case) and two value arguments (the dimensions, 3 and 3 in this case) and constructs the respective matrix type. In principle, one could give a codec for each matrix type that comes up in a database schema. But a much more elegant solution is to specify **codec operators** in analogy to type operators. A codec operator for a type operator with  $k$  type and  $l$  value arguments, takes  $k$  codec and  $l$  value arguments. For example, a codec operator for matrices takes a codec  $C : \text{codec } E$  for the element type  $E$  and the dimensions  $m$  and  $n$  and returns a codec of type `codec (Matrix  $E$   $m$   $n$ )`.

| Codecs (continued)          |   |  |
|-----------------------------|---|--|
| <code>StandardList</code>   | : <code>{<math>T</math>} codec <math>T \rightarrow</math> codec List(<math>T</math>)</code>               | JSON list, recursively coding each element of the list |
| <code>StandardVector</code> | : <code>{<math>T, n</math>} codec <math>T \rightarrow</math> codec Vector(<math>n, T</math>)</code>       | JSON list of fixed length $n$                          |
| <code>StandardMatrix</code> | : <code>{<math>T, n, m</math>} codec <math>T \rightarrow</math> codec Matrix(<math>n, m, T</math>)</code> | JSON list of $n$ lists of length $m$                   |

FIGURE 6. Some MMT Codec Operators for LMFDB; compare with Figure 5.

Like codecs, codec operators are represented in MMT in two ways: as declarations inside the theory `Codecs` (see Figure 6 for a list and Figure 9 in Section 6.1 for the general setting) and



as a corresponding Scala function that maps codecs to codecs. When reading the declarations, note that we make use of the dependent function types of the MitM foundation: curly brackets denote dependent function arguments, i.e., arguments that may occur in later argument types and the result type.

With these declarations, we recover the LMFDB encoding of isogeny matrices by applying the codec operator `StandardMatrix`, which encodes matrices as lists of lists, to the codec `StandardInt` and the dimension 3 and 3. The resulting codec

`StandardMatrix( $\mathbb{Z}$ , 3, 3, StandardInt)`

encodes the above matrix as `[[1,5,25],[5,1,5],[25,5,1]]`.

## 5. Integrating Computation Systems with MitM: GAP, SAGEMATH, and SINGULAR

We now show how we produce OMDOC/MMT theory graphs that specify the system dialects of GAP, SINGULAR, and SAGEMATH. The three systems are sufficiently different that we can consider the development presented in this section a meaningful case study in the methodology and difficulty of exposing the APIs of real-world systems as of formally described system dialects.

In each case, we had to overcome major implementation difficulties and invest significant manpower. In fact, even the serialization of internal abstract syntax trees as OMDOC/MMT objects proved difficult, for different system-specific reasons. In the following, we summarize these efforts.

### 5.1. SAGEMATH

We first consider our previous work [Deh+16] regarding a direct (i.e., without MitM) integration of SAGEMATH and GAP. Here SAGEMATH's native interface to GAP is upgraded from the **handle paradigm** to the **semantic handles** paradigm. In the former, when a system  $A$  delegates a calculation to a system  $B$ , the result  $r$  of the calculation is not converted to a native  $A$  object (unless it is of some basic type); instead  $B$  just returns a handle  $h$  (i.e., some kind of reference) to the  $B$ -object  $r$ . Later,  $A$  can run further calculations with  $r$  by passing it as argument to functions or methods implemented by  $B$ . Additionally, with a **semantic** handle,  $h$  behaves in  $A$  as if it was a native  $A$  object. In other words, one adapts the API satisfied by  $r$  in  $B$  to match the API for the same kind of objects in  $A$ . For example, the method call `h.cardinality()` on a SAGEMATH handle  $h$  to a GAP group  $G$  triggers in GAP the corresponding function call `Size(G)`.

This approach avoids the overhead of back and forth conversions between  $A$  and  $B$  and enables the manipulation of  $B$ -objects from  $A$  even if they have no native representation in  $A$ . However, if these  $B$ -objects need to be acted on by native operations of  $A$  or other systems (as in Jane's scenario), we actually have to convert the objects  $r$  between  $A$  and  $B$ .

**5.1.1. SAGEMATH API Theories.** In [Deh+16] we describe the extraction of some of SAGEMATH's API from its **categories**. This exploited the mathematical knowledge explicitly embedded in the code to cover a fairly large area of mathematics (hundreds of kinds of algebraic structures such as groups, algebras, fields, ...), with little additional efforts or need to curate the output. This extraction did not cover the constructors, knowledge about which is critical for (de)serialization, nor other areas of mathematics (graph theory, elliptic curves, ...) where SAGEMATH developers currently do not use categories (usually because the involved hierarchies of abstract classes are shallow and easily maintained by hand).

To extract more APIs, we took the following approach:

- (1) We constructed a list of typical SAGEMATH objects.
- (2) We used introspection to analyze those objects, crawling recursively through their hierarchy of classes to extract constructors and available methods together with some mathematical knowledge.

At this stage, the list of objects was crafted by hand to cover Jane's scenarios and some others. In a later stage, we plan to take advantage of one of SAGEMATH's coding standards: every concrete type must be instantiated at least once in SAGEMATH's tests and the instance passed through a generic test suite that runs sanity checks for its advertised properties (e.g. associativity, ...). Therefore, by a simple instrumentation of SAGEMATH's test framework, we could run our exporter on a fairly complete collection of SAGEMATH objects.

The process remains brittle and the export will eventually require much curation:

- The signature of methods is incomplete: it specifies the number and names of the arguments, but only the type of the first argument.

- For constructors, the type of all the arguments is known, but only for the specific call that led to the construction of the introspected object.
- There is no distinction between mathematically relevant methods and purely technical ones like data structure manipulation helpers.
- The export is very large and seems of limited use without alignments with the MitM ontology. At this stage we do not foresee much opportunities to produce such alignments other than manually (but cf. [Mül+17b] for a machine-learning based approach).

Nonetheless, we consider this an important first step toward fully automatic extraction of the SAGEMATH API. Moreover, we expect further improvements by code annotations in SAGEMATH (e.g., the ongoing porting of SAGEMATH from PYTHON 2 to PYTHON 3 will enable **gradual typing**, which we hope to become widely adopted by the community) or using type inference in SAGEMATH and/or MitM.

5.1.2. *Serialization and Deserialization.* Because SAGEMATH is based on PYTHON, it benefits from its native serialization support. For example, the dihedral group  $D_4$  is serialized as a binary string, which encodes the following straight line program to be executed upon deserialization:

```
pg_unreduce = unpickle_global('sage.structure.unique_representation', 'unreduce')
pg_DihedralGroup =
    unpickle_global('sage.groups.perm_gps.permgroup_named', 'DihedralGroup')
pg_make_integer = unpickle_global('sage.rings.integer', 'make_integer')
pg_unreduce(pg_DihedralGroup, (pg_make_integer('4'),), {})
```

The first three lines recover the constructors for integers and for dihedral groups from SAGEMATH's library. The last line applies them to construct successively the integer 4 and  $D_4$ .

Up to concrete syntax, this serialization is already close to the desired SAGEMATH system dialect. We can therefore extend PYTHON's native (de)serializer to use OMDoc/MMT as an alternative serialization format (using the PYTHON OpenMath interface library [POMa] developed in the OpenDreamKit project). The following shows the corresponding OpenMath syntax tree in Python

LISTING 1. The dihedral group  $D_4$  in OpenMath Syntax

```
OMApplication(
    elem=OMSymbol(name='DihedralGroup',
                  cd='sage.groups.perm_gps.permgroup_named', cdbase='http://python.org'),
    arguments=[OMApplication(
        elem=OMSymbol(name='make_integer',
                      cd='sage.rings.integer', cdbase='http://python.org'),
        arguments=[OMBytes(bytes='4')]))])
```

and XML respectively:

LISTING 2. The dihedral group  $D_4$  in XML Format

```
<OMA xmlns="http://www.openmath.org/OpenMath">
  <OMS name="DihedralGroup"
    cd="sage.groups.perm_gps.permgroup_named" cdbase="http://python.org"/>
  <OMA>
    <OMS name="make_integer" cd="sage.rings.integer" cdbase="http://python.org"/>
    <OMI>4</OMI>
  </OMA>
</OMA>
```

This approach has the additional advantage of benefiting from future optimizations implemented in PYTHON's serialization, like structure sharing for identical subexpressions.

Still, systematically expanding OMDoc/MMT serialization to the *entire* SAGEMATH library requires significant manpower and can only be a long-term goal. To increase community support, our design elegantly decouples the problem into *i*) instrumenting the serialization to generate OMDoc/MMT as an alternative target format, and *ii*) structural improvements of the serialization that benefit SAGEMATH in general.

In particular, our serialization of SAGEMATH objects is **by construction** rather than **by representation**, i.e., we serialize the constructor call that was used to build an object instead of the low-level PYTHON representation of the resulting object. This is important to hide implementation details and allow for straightforward alignments. From the origin, the SAGEMATH community has internally promoted good support for serialization as this is a fundamental building block for communication between parallel processes, databases, etc. Thus, it already values serialization by construction as superior because it is usually more concise and more robust under changes to SAGEMATH. Therefore, independent of the purposes of this report, we expect a synergy with the SAGEMATH community toward improving serialization.

## 5.2. GAP

In [Deh+16], we already described our general approach to extract APIs from the GAP system. We have now improved on this work considerably.

Firstly, we improved the MitM foundation so that the primitives of GAP's type system can be expressed in the MitM ontology.<sup>3</sup> GAP's type system heavily uses subtyping: **filters** express finer and finer subtypes of the universal type `IsObject`. Moreover, an object in GAP can learn about its properties, meaning its type is refined at runtime: a group can learn that it is Abelian or nilpotent and change its type accordingly.

Secondly, we devised and implemented a special treatment of GAP's constructors during serialization. As GAP only has a weak notion of object construction, we achieved this by manually identifying and annotating all functions that create objects in the GAP code base and then instrumenting them to store which arguments they were called with. With the constructor annotation in place, it is possible to have GAP represent any object in a running session as either a primitive type (integers, permutations, transformations, lists, floats, strings), or as a constructor applied to a list of arguments.

The instrumentation itself is minimal – 57 lines of GAP code, plus 100 lines for serializing and parsing. The main – and indeed considerable – challenge was to identify the constructors and their arguments. In GAP, objects are created by calling the function `Objectify` with a type and some arguments. Hence we analyzed all call-sites to this function and some light inference of the enclosing function. This amounted to 665 call sites in the GAP library and an additional 1664 in the standard package distribution. The instrumentation will be released as part of a future version of GAP, making GAP fully MitM capable.

As a major positive side-effect of our work, this instrumentation led to general improvements of the type infrastructure in GAP. For example, it enables static type analysis, which can be used to optimize the dynamic method dispatch and thus hopefully lead to efficiency gains in the system.

## 5.3. SINGULAR

As we only need a very small part of SINGULAR for our case study, we were able to use the existing OpenMath content dictionaries for polynomials [OMCP] as the SINGULAR system dialect. These are part of a standard group of content dictionaries that describe (some) mathematical objects at a high level of abstraction to be universally applicable. OMDoc/MMT understands OpenMath, i.e., it can use these content dictionaries as OMDoc/MMT theories.

<sup>3</sup>In the future MMT might even serve as an external type-checker for GAP.

Building on the OpenMath toolkits for OpenMath phrasebooks [POMa] and SCSCP communication [POMb] in PYTHON – which were developed for SAGEMATH in the OpenDreamKit project, we wrapped SINGULAR in a thin layer of PYTHON code that provides SCSCP communication. This work was undertaken by a student as part of a summer internship in about a week without prior expert knowledge of the system.

Since then, Sebastian Gutsche has tested the feasibility of the MitM approach by testing the effort of adding a set of system API theories for the SINGULAR system. He has generated the full set necessary to achieve a comprehensive coverage of the SINGULAR dialect by heuristically parsing definitions, call patterns, and comment strings in the Singular C code. The pure coding exercise took him – a very experienced Singular user and part-time developer – three days. This shows that the joining costs involved in adding a CAS is small, in this case smaller than developing the SAGEMATH glue code had been.

#### 5.4. Alignments

Finally we have to curate the **alignments** between the system dialects and the MitM ontology.

To make the systems interoperable and translate objects and expressions, it is crucial to inform the middleware which symbols in the respective system API theories represent which mathematical concepts from the ontology. Then (in the simplest case) translation reduces to simply substituting all symbols in an expression. A major difficulty for middleware integration is correctly handling the subtleties in more complicated cases.

To work this out in detail, [Mül+17b] introduced OMDOC/MMT **alignments**. Technically, these are pairs of OMDOC/MMT symbol identifiers decorated by a set of key-value pairs. The alignments of  $a$ -symbols with the MitM ontology determine which  $A$ -objects correspond to MitM-objects.

The alignment of  $a$ -symbols to ontology symbols must be spelled out manually. But this is usually straightforward and easy even for inexperienced users. For example, the following line aligns GAP’s symbol `IsCyclic` (in the file `lib/grp.gd`) with the corresponding symbol `cyclic` in the MitM ontology.

```
gap:/lib?grp?IsCyclic mitm:/smglom/algebra?group?cyclic direction="both" type="VRE"
```

Here the two key-value pairs at the end are used to signify that this alignment is part of a group of alignments called “VRE” and can be used for translations in both directions.

Thus, we can reduce the problem of interfacing  $n$  systems to *i*) curating the MitM ontology for the joint mathematical domain, *ii*) generating  $n$  theory graphs for the system dialects, *iii*) maintaining  $n$  collections of alignments with the MitM ontology.

Alignments form an independent part of the MitM interoperability infrastructure. They obey a separate development schedule: The MitM ontology is developed by the community as a whole as the understanding of a mathematical domain changes. The system dialects are released together with the systems according to their respective development cycle. The alignments bridge between them and have to mediate these cycles. The alignments are currently produced and curated manually. In the future, we will also consider automatically extracting alignments from API theories. This is possible if system developers annotate their functions with the corresponding URIs in the MitM ontology, which has already been done in an ad-hoc nature in the SAGEMATH library for similar translations.

Unfortunately, in general, even when two symbols represent the same mathematical operation, there may be subtle differences between implementations that require translations, e.g., symbols can differ in the order of arguments, the choice of default values for omitted arguments, etc. In those instances, we have two options. If the translation is not interesting or difficult, we use the advanced alignments introduced by [Mül+17b]. For example, this is indicated for the representation of polynomials as lists of coefficients starting with the highest vs. starting with the lowest coefficient. Otherwise, we formalize *both* variants as separate symbols in the MitM

ontology and formalize the translations in MMT and prove their correctness. If two variants are present in the ontology, a symbol in a system API is aligned with the most similar one in the ontology. For example, this is indicated for the representation of polynomials as dense vs. sparse lists of coefficients. In practice, the choice between these two options is not clear-cut and must often be made pragmatically.

**Example.** Consider again the dihedral group in SAGEMATH from Listings 1 and 2 in Section 5.1.1. We want to translate it to GAP. First we align the SAGEMATH constructor for dihedral groups (`sage:?sage.groups.perm_gps.permgroup_named?DihedralGroup`) with a corresponding symbol in the MitM ontology (say, for instance `mitm:?group?dihedralGroup`). Now when the MMT system receives an object from SAGEMATH, it uses these alignments to replace all SAGEMATH symbols by their system-neutral MitM analog. Similar alignments induce the replacement of `make_integer(4)` with the plain OpenMath integer 4. Second, to translate from MitM to GAP, we make a corresponding alignment and traverse it from MitM to GAP.

Concretely, we could align `mitm:?Groups?dihedralGroup` with the constructor for dihedral groups in GAP, namely `gap:/grp?basic?DihedralGroupCons`. However, as mentioned in Section 1, the definitions of dihedral groups in SAGEMATH and GAP differ with respect to their arguments — the group  $D_4$  in SAGEMATH is called  $D_8$  in GAP. So with this second alignment, MMT would falsely translate SAGEMATH's  $D_4$  to GAP's  $D_4$ .

We could bridge this discrepancy with a complex alignment that performs the easy translation. However, both conventions for naming dihedral groups are wide-spread in practice, and users may deem this translation interesting enough to make explicit. Thus, we can alternatively declare two different symbols in the MitM ontology: `mitm:?group?dihedralGroupSmall` and `mitm:?group?dihedralGroupLarge` corresponding to the two different ways to construct dihedral groups from natural numbers. Now we align the SAGEMATH symbol with the former and the GAP symbol with the latter using the following alignments:

LISTING 3. Alignments for Dihedral Groups

```
sage:?sage.groups.perm_gps.permgroup_named?DihedralGroup
  mitm:?group?dihedralGroupSmall direction="both" type="VRE"
gap:/grp?basic?DihedralGroupCons
  mitm:?group?dihedralGroupLarge direction="both" type="VRE"
```

Furthermore, we formalize the translation rule in MitM that maps `dihedralGroupSmall( $n$ )` to `dihedralGroupLarge( $2n$ )` and vice versa. Then translating from SAGEMATH to GAP is a five-step process as illustrated in Figure 7 consisting of (from top to bottom):

- (1) SAGEMATH exports the object `DihedralGroup(4)` to OpenMath using its system dialect.
- (2) MMT applies alignment to obtain the corresponding OpenMath object for MitM.
- (3) MMT applies the above translation rule. (In general, none or multiple rules may have to applied.)
- (4) MMT applies alignment to obtain the corresponding OpenMath in GAP's system dialect.
- (5) GAP imports the object to obtain `DihedralGroupCons(8)`.



|                                |   |
|--------------------------------|---|
| SAGEMATH                       | DihedralGroup(4)  |
| OM (SAGEMATH)                  | <pre> &lt;OMA cdbase="http://python.org"&gt;   &lt;OMS name="DihedralGroup"     cd="sage.groups.perm_gps.permgroup_named"/&gt;   &lt;OMA&gt;     &lt;OMS name="make_integer" cd="sage.rings.integer"/&gt;     &lt;OMI&gt;4&lt;/OMI&gt;   &lt;/OMA&gt; &lt;/OMA&gt; </pre> |
| OM (MitM)                      | <pre> &lt;OMA cdbase='http://mathhub.info/MitM/Core'&gt;   &lt;OMS name='dihedralGroupSmall' cd='group'/&gt;   &lt;OMI&gt;4&lt;/OMI&gt; &lt;/OMA&gt; </pre>   |
| OM (MitM, after applying view) | <pre> &lt;OMA cdbase='http://mathhub.info/MitM/Core'&gt;   &lt;OMS name='dihedralGroupLarge' cd='group'/&gt;   &lt;OMI&gt;8&lt;/OMI&gt; &lt;/OMA&gt; </pre>   |
| OM (GAP)                       | <pre> &lt;OMA cdbase='http://www.gap-system.org/grp'&gt;   &lt;OMS name='DihedralGroupCons' cd='basic'/&gt;   &lt;OMI&gt;8&lt;/OMI&gt; &lt;/OMA&gt; </pre>  |
| GAP                            | DihedralGroupCons(8)  |

FIGURE 7. Stepwise Translation of Dihedral Group  $D_4$  from SAGEMATH to GAP using the alignments from Listing 3

## 6. Integrating Databases with MitM: The LMFDB Case Study

The mathematical software systems to be integrated via the MitM approach have so far been computation-oriented, e.g., computer algebra systems. Their API theories typically declare types and functions on these types (the latter including constants seen as nullary functions). Even though database systems differ drastically from these in many respects, they are very similar at the MitM level: a mathematical database defines

- some types: each table's schema is essentially one type definition,
- many constants: each table entry is one constant of the corresponding type.

Thus, we can reuse many of the same concepts. In particular, the API theories must contain definitions of the database schemas.

Apart from standard software engineering tasks, this leaves three conceptual problems we had to solve:

**P1** Turn the database schemas and tables into OMDoc/MMT theories and declarations.

**P2** Lift data in *physical* representation (as records of the underlying database) to OMDoc/MMT object in *semantic* representation.

**P3** Translate semantic queries to queries about physical representations so that they can be executed directly on the database without loading the entire theory into MMT.

We deal with **P1** in Section 6.2, with **P2** in Section 6.3, and with **P3** in Section 6.4.

In this report, we focus on LMFDB as an example, of which we give an overview in Section 6.1. However, our methods are general enough to apply to many other mathematical databases such as OEIS, or findstat.

### 6.1. LMFDB Overview

The “L-Functions and Modular Forms Database” (LMFDB [LMF]) is a large database, storing among other mathematical objects millions of L-Functions, modular forms and curves, along with their properties. Technically, it uses a MongoDB database with a Python web frontend.



LMFDB has several sub-databases, e.g., for elliptic curves or transitive groups. Within each of these, every object is stored as a single JSON record. Figure 8 shows an example: each property of this JSON object corresponds to a property of the underlying mathematical object. For example, the `degree` property — here 1 — of the JSON objects corresponds to the degree of modular parametrization of the underlying elliptic curve.

```
{
  "degree": 1,
  "x-coordinates_of_integral_points": "[5,16]",
  "isogeny_matrix": [[1,5,25],[5,1,5],[25,5,1]],
  "label": "11a1",
  "_id": "ObjectId('4f71d4304d47869291435e6e')",
  ...
}
```

FIGURE 8. Part of an elliptic curve in LMFDB (some fields omitted for brevity)

Other properties are more complex: the value of the `isogeny_matrix` property is a list of lists representing a matrix. This disconnect between JSON encoding and mathematical meaning can become much more severe, e.g., the `x-coordinates_of_integral_points` field is semantically a list of integers but (due to the sizes limits on integers) is encoded as a string.

The LMFDB API [Lmf] exposes a querying interface that can be used either by humans via the web or programmatically via JSON-based GET requests over HTTP. Queries must name the sub-database to be queried and consist of a set of key-value pairs that correspond to an SQL `where` clause. However, while LMFDB offers a programmable API for accessing its contents, this API sits at the level of the underlying MongoDB, and not the level of mathematical objects. For example, to retrieve all Abelian objects in the subdatabase of transitive groups, we expect to use the key-value pair `commutative = true`. However, these values need to be encoded to be understood by MongoDB. We need to realize that the database schema actually uses the key `ab` for commutativity, that it has boolean values, and that the schema encodes `true` as 1. Thus, the actual query to send is <http://www.lmfdb.org/api/transitivegroups/groups/?ab=1>.

In this example, all steps are relatively straightforward. But in general, e.g. when searching for all elliptic curves with a specific isogeny matrix, this not only requires good familiarity with the mathematical background but also with the system internals of the particular LMFDB sub-database; a skill set commonly found in neither research programmers nor average mathematicians.

## 6.2. LMFDB as a Set of Virtual Theories

The set of constants in a database table – while finite – can be arbitrarily large. In particular, all LMFDB tables<sup>4</sup> are just finite subsets of infinite sets, whose size is not limited by mathematical specifications but by computational power: the database holds all objects that users have computed so far and grows constantly as more objects are computed. LMFDB tables usually include a naming system that defines unique identifiers (which are used as the database keys) for

<sup>4</sup>Technically, until July 2018, LMFDB was implemented using MongoDB and comprises a set of sets (each one called a database) of JSON objects. MongoDB allows each JSON object in a collection to be different (with a different schema), though in practice almost all objects in each collection had the same schema apart from some missing data components. The schema for each collection had to be documented elsewhere, in an inventory, which since 2017 has been itself stored as a database within the LMFDB. During 2018, however, work has been ongoing to migrate the LMFDB to use PostgreSQL (with a fixed schema for each table) as the underlying database, without any change to the external API. In both cases, due to the conventions used, we can understand the LMFDB conceptually as a set of tables of a relational database, keeping in mind that every row is a tuple of arbitrary JSON objects.

these objects, and these identifiers are predetermined even for those objects that have not been computed yet.

Thus, it is not practical to fix a set of concrete API theories. Instead, the API theories must be split into two parts: for each database table, we need

- a concrete theory called the **schema theory** that defines the schema and other relevant information about the type of objects in the table and
- a virtual theory called the **database theory** that contains one definition for each value of that type (using the LMFDB identifier as the name of the defined constant).

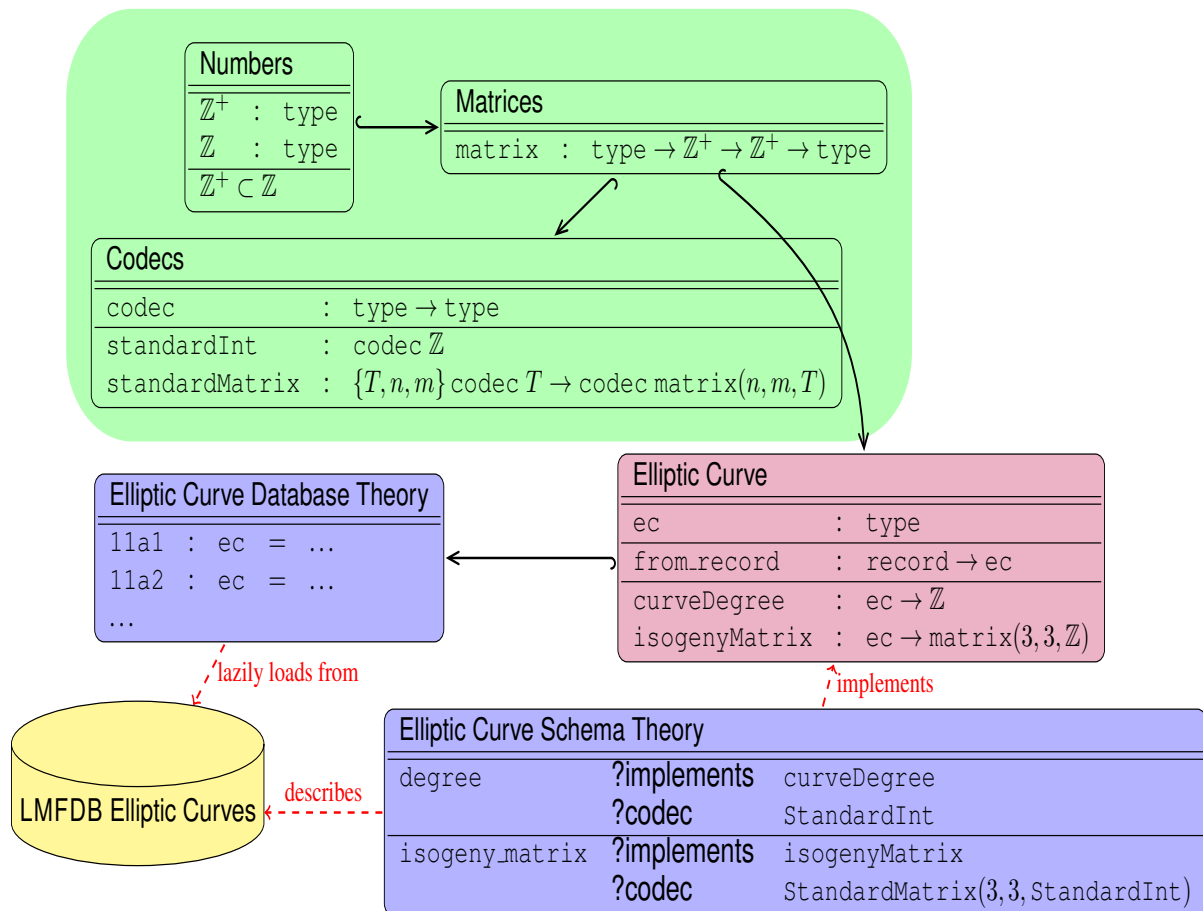


FIGURE 9. Virtual theory for LMFDB elliptic curves (some declarations omitted)

Conceptually, it is straightforward how to implement each LMFDB table as a virtual MMT theory  $V$ : we use an initially empty concrete theory  $C$ , and whenever an identifier  $\text{id}$  of  $V$  is requested, MMT dynamically adds the corresponding declaration of  $\text{id}$  to  $C$ . Because MMT already abstracts from the physical realizations of persistent storage, we only have to implement a new storage instance that connects to LMFDB, retrieves the JSON object with identifier  $\text{id}$ , and turns it into an OMDoc/MMT declaration.

A sketch of our overall solution is given in Figure 9. The relevant parts of the MitM ontology comprise simple ones like numbers and matrices (in green) and LMFDB-specific ones like elliptic curves (in red). The remaining theories (in blue) form the LMFDB API theories: the schema theory and the database theory, which we describe below.

LMFDB's original technical realization, using MongoDB, did not require formalizing the schema of each table. Instead, the tables were generated systematically and therefore followed an implicit schema that could — in principle — be obtained from the documentation or reverse-engineered from the tables. Until 2017 the documentation of these implicit schemas was created

and maintained manually by LMFDB developers, and as a result was incomplete and frequently out of date. During 2017-2018 a new LMFDB Inventory was created, taking as its starting point both the manually prepared inventory (which contained human definitions and explanations of the content of each data field) and a new dynamically created schema obtained by analysis of the data in each collection. This process, which was both necessitated by the requirements of the Math-in-the-Middle approach and made possible in practice through the provision of research software engineers funded by ODK, revealed numerous inconsistencies in the LMFDB data which developers have since been able to address. Moreover, having this detailed schema for each collection in the MongoDB databases also fed in to the migration process, expected to be complete by August 2018, in which the MongoDB free-format collections are being replaced by PostgreSQL tables, each of which has a completely specified and formalized schema.

Note that (and here LMFDB critically differs from, e.g., the OEIS), the mathematical definitions and concepts involved in the LMFDB data and tables are extremely deep so that reverse-engineering the associated schemas from the data itself is only possible in practice with the aid of experts. As the first such table for which a formal schema was to be created, before the development of the new comprehensive Inventory, we chose one for which the existing documentation was most complete, and which originated with one of the current authors (John Cremona), who sat down with the Math-in-the-Middle team at an ODK workshop to formalize the corresponding schema in OMDoc/MMT. In the following, we will use this table as a running example. Our methods extend immediately to any other table once its schema has been formalized.

Our formalization models elliptic curves in a very simple fashion by using an abstract type `ec`. The constructor `from_record` takes an MMT record and returns an elliptic curve. Properties of elliptic curves are formalized as functions out of this type. We list only two here as examples: the `degree`, an integer, and the `isogeny_matrix`, a  $3 \times 3$  matrix of integers. We omit the relevant axioms, which are not essential for our purposes here. As usual for the MitM approach, the model of elliptic curves does not rely on LMFDB, nor any other system, so that we can integrate other knowledge sources about elliptic curves or to future versions of the LMFDB with changed structure.

### 6.3. Ascribing Encodings in Schema Theories

If we ignore encoding issues, schema theories are straightforward: they contain one declaration of the same name for each field within an LMFDB record. This specifies only the semantic type of each field and does not relate it to the MitM formalization. To handle the encoding as a physical type, we annotate each declaration with the codec that the databases for the values of that field. Moreover, to connect the schema theories to the MitM formalization, we additionally annotate each field with the corresponding property of elliptic curves from the MitM theory. We can now understand the last unexplained parts of Fig. 9. `?implements` is the symbol used to annotate the metadatum, which MitM property a schema field corresponds to. And `?codec` similarly annotates the codec to each field.

For example, the `degree` field implements the `curveDegree` property in the elliptic curve theory and uses the `StandardInt` codec. Thus, the schema theories determine the entire relation between semantic and physical objects.

The database theory is a virtual theory and contains one declaration per LMFDB record. Given the URI of an object in the respective database, our MMT backend for LMFDB first retrieves the appropriate record from LMFDB – in the case of `11a1` this corresponds to retrieving the JSON found in Figure 8. Then, for each field, it uses the annotated codec (which is an OMDoc/MMT expression) to build an actual codec (as a runnable Scala function) and runs its decoding function. Next, it passes the resulting record to the `from_record` constructor,

which yields an elliptic curve in the MitM theories. Finally, this elliptic curve is added as a new declaration in the database theory.

This example already shows that the virtual theories framework — while conceptually slightly complex — is declarative in the sense that it can be configured by supplying schema theories. This task only requires knowledge about the underlying mathematics and how it is represented in the MitM ontology, the encodings of the respective mathematical types in the basic data types of the underlying database, and the data base schema. In particular no knowledge of the MMT internals or programming skills are needed and the prerequisite knowledge is exactly what LMFDB contributors have when they set up the particular sub-database.

#### 6.4. Translating Queries

Recall that MMT has a general-purpose Query Language called QMT [Rab12], which allows users to find knowledge subject to even complex conditions. We continue by briefly addressing **P3**: query translation; for a complete discussion we refer the interested reader to [Wie17].

In practice, most queries involving virtual theories so far have a shape similar to the one that LMFDB supports: Finding all objects within a single sub-database for which a specific field has a specific value. As an example, consider again the query of finding all Abelian transitive groups. QMT has an MMT-powered surface syntax, which can be used to express this query as:

```
x in (related to ( literal 'lmfdb:db/transitivegroups?group ' by (object declares))
| holds x (x commutative x == true)
```

The example consists of two parts, first we find all objects declared in the theory interface theory `lmfdb:db/transitivegroups?group` (line 1), and then we restrict this set of results to all those for which the `commutative` property is `true` (line 2). Notice that this the example shown here is the formal equivalent of the LMFDB query shown in Section 6.1. The key difference is that this query does not require knowing the record structure of LMFDB— apart from knowing the proper sub-db, instead it only relies on knowing the mathematical semantics (commutativity) of the query in question.

Recall that to evaluate a query prior to the introduction of virtual theories, the MMT system loaded the theory graph into main memory and then interleaved incremental flattening and query evaluation operations on the MMT data structures until a result had been produced. But it is infeasible to first load all potentially relevant data into memory, and only then proceed with evaluation. This would require loading a copy of LMFDB into main memory, something that virtual theories were designed to avoid.

The low-level API of LMFDB and similar system provides a new approach for making queries towards virtual theories. First, the MMT query is translated into a system-specific information-retrieval language — in the case of LMFDB this is currently a MongoDB-based syntax. Next, this translated query is sent to the external API. Upon receiving the results, these are translated back into OMDoc/MMT with the help of already existing functionality in the appropriate virtual theory backend.

This leaves just one problem unsolved — translating queries into the system-specific API. Note that it is insufficient to simply translate queries as a whole: On the one hand a general QMT query may or may not involve a virtual theory, on the other hand, it may also involve several (unrelated) virtual theories. This makes it necessary to filter out queries involving virtual theories, and assign them to a specific backend, and then translate only these parts.

Achieving this automatically is a non-trivial problem. Queries are inductive in nature, and one could attempt to intercept each of the intermediate results. However, this would require a check on each intermediate result to first determine if it comes from a virtual theory or not, and then potentially switching the entire evaluation strategy, leading to a computationally expensive implementation.

Instead of intercepting each result, we extended the Query Language to allow users to annotate sub-queries for evaluation with a specific virtual theory backend. This allows the system to immediately know which parts of a query have to be evaluated in MMT memory, and which have to be translated and sent to an external system. This turns the example above into:

```
use "lmfdb" for {*
  x in (related to ( literal `lmfdb:db/transitivegroups?group )
    by (object declares)) | holds x (x commutative x == true)
*}
```

Here, we have simply wrapped the entire query with a `use lmfdb` statement, indicating the query should be evaluated using LMFDB.

The encoding of this specific query can be achieved using codecs. The query corresponds to the URL <http://www.lmfdb.org/api/transitivegroups/groups/?ab=1>. Next, the LMFDB API returns a set of JSON objects corresponding to all Abelian transitive groups. These can then be decoded into OMDOC/MMT objects using the procedure described in Section 6.3, i.e. for each field we look up the corresponding codec and use it to deconstruct the field, eventually creating an MMT record. Afterwards, these OMDOC/MMT terms can then be passed to the user as a result to the query.

### 7. MitM-Based Distributed Computation

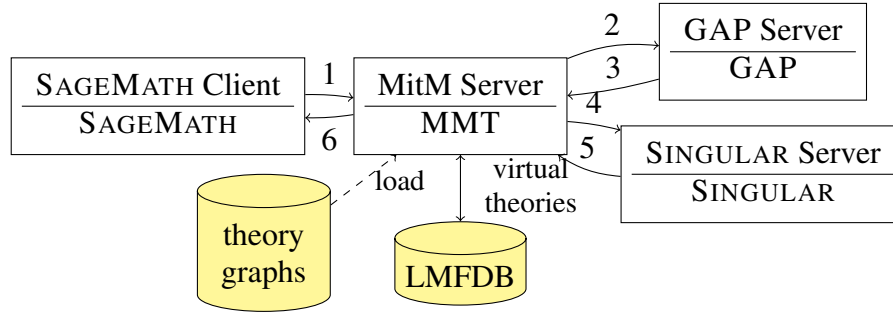


FIGURE 10. MitM Interaction in Jane's Use Case

Figure 10 shows the overall architecture with an MitM server as the central mediator. All arrows represent the transfer of OMDOC/MMT objects via SCSCP. Critically, the MitM server also maintains the alignments and uses them to convert between system dialects.

We have extended the MMT system [Rab13] with an SCSCP server/client so that it can receive/send objects from/to computation systems. For the GAP server, we built on pre-existing SCSCP support. To obtain an SCSCP server for SINGULAR, which does not have native SCSCP support, we wrapped SINGULAR in a Python script that includes the pycscsp library [POMb]. In SAGEMATH, we directly programmed the client interface to the MitM server.

The resulting system forms the nucleus of the OpenDreamKit interoperability layer. It can already delegate computations between the four participating systems as long as the exchanged objects are covered by the MitM ontology, the alignments, and the formalizations of the system dialects.

**Jane's Use Case.** Initially, Jane has already built in SAGEMATH the ring  $R = \mathbb{Z}[X_1, X_2, X_3, X_4]$ , the group  $G = D_4$ , the action  $A$  of  $G$  on  $R$  that permutes the variables, and the polynomial  $p = 3 \cdot X_1 + 2 \cdot X_2$ . She now calls

```
MitM.Singular(MitM.Gap.orbit(G, A, p)).Ideal().Groebner().sage()
```

which results in the following steps (the numbers on the edges of the graph of Figure 10 indicate the order of communications when processing Jane's use case):

- (1) Jane uses SAGEMATH to call the MitM server with the command above, which includes both the computation to be performed and information about which system to use at which step.
- (2) The MitM server translates `MitM.Gap.orbit(G, A, p)` to the GAP system dialect and sends it to GAP.
- (3) GAP returns the orbit:

$$O = [3X_1 + 2X_2, 2X_3 + 3X_4, 3X_2 + 2X_3, 3X_3 + 2X_4, \\ 2X_2 + 3X_3, 3X_1 + 2X_4, 2X_1 + 3X_4, 2X_1 + 3X_2].$$

- (4) The MitM server translates `MitM.Singular(O).Ideal().Groebner()` to the SINGULAR system dialect and sends it to SINGULAR.
- (5) SINGULAR returns the Gröbner base  $B$ .
- (6) The MitM server translates  $B$  to the SAGEMATH system dialect and sends it to SAGEMATH, where the result is shown to Jane.

$$B = [X_1 - X_4, X_2 - X_4, X_3 - X_4, 5X_4].$$



**Alternative Use Case.** Suppose Jon, one of Jane's colleagues, prefers working in GAP, and he wants to compute the Galois group of the rational polynomial  $p = x^5 - 2$ . He discovers the GAP package `radroot`, which promises this functionality, but unfortunately the package does not work for this polynomial and thus GAP alone cannot solve Jon's problem.

Jon hears from Jane that he should use SAGEMATH, because she knows it can compute Galois groups. So, from GAP, he calls

```
G := MitM("Sage", "GaloisGroup", p)
```

which gives him the desired Galois group as a GAP permutation group. Having heard of Jane's experiments, he can further run her orbit and Gröbner basis calculation starting from this new group, without leaving his favorite computing environment.

Finally, Jon, being a proficient GAP user, also knows that he can now install a **method** in GAP by calling

```
InstallMethod(GaloisGroup, "for a polynomial", [IsUnivariatePolynomial],
  p -> MitM("Sage", "GaloisGroup", p))
```

that will compute the Galois group of any rational polynomial transparently for him whenever he calls `GaloisGroup` for a rational polynomial in GAP. And thus (at the price of using multiple systems) a significant part of the 1800-line `radroot` package can be replaced by a few lines in GAP, taking advantage of the work of the SAGEMATH community and participating in any future improvements of SAGEMATH. In fact, Sage itself delegates to the PARI system – another one of the OpenDreamKit systems – for this computation. So in the future GAP might directly delegate to PARI instead, bypassing the need of iterated translations.



## 8. Conclusion

**Summary.** The MitM approach envisions integrating mathematical software systems based on formalizations of the underlying mathematical knowledge that provide the center of a star-shaped communication architecture. The evaluation and application of this approach requires several major practical investments:

- the curation of the MitM Ontology,
- the generation of formal specifications of APIs for concrete systems,
- the alignment of the ontology and the system APIs,
- the implementation of a MitM server that uses alignments to translate between systems,
- communication modules in each concrete system that communicate with the MitM server (using the SCSCP protocol).

In this report, we have described substantial progress on each aspect. Due to the enormity of the overall goal, we have focused on concrete case studies that have driven this progress. Concretely, we focused on group theory and modular forms as well as a few adjacent areas of mathematics to draw integration examples from. Moreover, we picked four concrete systems to integrate: the computation systems SAGEMATH, GAP, and SINGULAR, and the mathematical database collection LMFDB.

Our efforts led to several foundational innovations that were needed to apply MitM in practice. These included the notion of virtual theories and mathematical codecs that relate large collections of mathematical objects to high-level formalizations of mathematics. Virtual theories provide feasibly small windows to large databases. And codecs perform systematic syntax translations between mathematical languages and general purpose database languages.

**Evaluation.** Our case studies show that MitM-based integration is an achievable goal. Critically, the MitM-based approach to interoperability of data sources and systems keeps systems and data sources “as is”. Only their APIs are documented in a machine-actionable way that can be utilized for remote procedure calls, content format mediation, and service discovery. As a consequence, interaction between systems is very flexible and requires only relatively small investments for each system’s developer community. Delegation-based workflows can either be programmed directly or embedded into the interaction language of the mathematical software systems. Moreover, this can be done in a way that requires only minimal changes to users’ existing work flows.

In comparison to middleware systems for distributed objects computing like CORBA, or ProtoBuffers, the MitM paradigm solves the “semantics preservation problem” by supplying the MitM ontology as a “specification layer” and the OMDOC/MMT framework as a joint meaning space, which is infinite in principle and only limited by the coverage of the MitM ontology.

The MitM paradigm is more general than the ad-hoc integration approach in used in SAGEMATH, in particular, as already our two use cases show, there is no dedicated master system — the system that the user feels at home in. In Jane’s use case it is still SAGEMATH, in Jon’s it is GAP, and in Joanna’s it is LMFDB.

Incidentally, we can understand the MitM paradigm as an extension of the OpenMath approach: MitM uses the OpenMath encoding for transporting objects and the OpenMath-based SCSCP transport protocol, which the OpenMath society endorsed as an OpenMath standard in 2017 prompted by the OpenDreamKit project. Here the innovation of the MitM effort in **WP6** is

- (1) using OMDOC/MMT as a much more expressive language for content dictionaries that can accommodate both formal and informal representations of the background knowledge,
- (2) using theory morphisms that specify meaning-preserving relations between theories to modularize ontologies,

- (3) the provisioning of the MitM ontology and thousands of system API theories, which constitute high-impact content dictionaries for OpenMath,
- (4) working use cases that demonstrate the feasibility the decades-old “OpenMath Vision” of distributed, meaning-preserving computation in mathematics.

Compared to ad-hoc translations, MitM-based interoperability is relatively inefficient as objects have to be serialized into (possibly large) OMDOC/MMT objects, transferred via SCSCP to MMT, parsed, translated into another system dialect, serialized and transferred, and parsed again. On the other hand, instead of implementing and maintaining  $n^2$  translations, we only have to establish and maintain  $n$  collections of system APIs and their alignments to the MitM ontology. This makes the management of interoperability much more tractable:

- (1) The MitM ontology is developed and maintained as a shared resource by the community. We expect it to be well-maintained, since it can directly be used as a documentation of the functionality of the respective systems.
- (2) All the workflows are star-shaped: instead of requiring expert knowledge in two systems — a rare commodity even in open-source projects, and even for the system experts involved in this report — and keeping up with their changes, the MitM approach only needs expertise and change management for single systems.

All in all, these translate into a “business model” for MitM-based cooperation in terms of the necessary investment and achievable results, which is based on the well-known *network effects*: the joining costs are in the size of the respective system, whereas the rewards — i.e. the functionality available by delegation — is in the size of the network.

**Future Work.** Of course, MitM is very young and only backed by a minuscule developer community (the participants of **WP6** in the OpenDreamKit), so tool support, documentation and community buy-in are just starting to develop. The most obvious path of future work is in scaling up the approach in two orthogonal direction: write ontologies for more areas of mathematics and connect more systems and databases to the MitM ecosystem. Both may require substantial efforts but due to our initial developments the marginal costs can now be quite affordable, as the Singular case study shows (see Section 5.3). In the long run, we expect that the MitM-generated qualitative boost to mathematicians’ productivity will these costs.

Concrete candidates for additional systems include PARI GP and the remaining databases of LMFDB. Moreover, we are working on integrating the Online Encyclopedia of Integer Sequences (OEIS [Slo03; OEIS]) — a difficult undertaking due to the weak structure it enforces on its data. In [LK16] we have already (heuristically) formalized OEIS contents in OMDOC/MMT; the next step will be to come up with appropriate codecs based on this basis and develop schema theories for OEIS.

Moreover, the new OSCAR project [OC], which aims at a high-efficiency integration of the computer algebra systems GAP, SINGULAR, Polymake, and Antic using shared-memory communication and Julia as a glue language, is closely related to OpenDreamKit. The most promising synergy between OpenDreamKit and OSCAR, which we have recently started exploring, is to jointly use and further develop the MitM ontology as a documentation and integration facility across all systems.

The network effect described above can be enhanced further by technical refinements. For instance, if we annotate alignments with a “priority” value that specifies how canonically/efficiently/powerfully a given system implements a given MitM operation, then we can let the MMT mediator automatically choose a suitable target system for a requested computation (as opposed to our current setup where Jane specifies which systems she wants to use). On the other hand, for workflows where we do not need or want service-discovery, alignments can be “compiled” into  $n^2$  transport-efficient direct translations that may even eliminate the need for serialization and parsing.

For local applications, where all involved systems and the mediator run on the same machine, we have started investigating programming language bridges that would allow direct communication between systems and the MMT mediator. This would allow by-passing the overhead of serialization and parsing and network communication. Concretely, MMT can now communicate bidirectionally with Python, the underlying language of SAGEMATH.

Another path to increased efficiency is to statically compile alignments into translation modules that would allow MitM-aware communication without the need for a dynamic mediator. This would speed up communication as only one recursive translation function would be needed instead of two.

**Acknowledgments.** The authors gratefully acknowledge the fruitful discussions with other participants of work package WP6, in particular Alexander Konovalov on SCSCP, Paul Dehayé on the SAGEMATH export and the organization of the MitM ontology, Luca de Feo on OpenMath phrasebooks and the SCSCP library in python.

## Bibliography

- [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>.
- [Bus+04] Stephen Buswell et al. *The Open Math Standard, Version 2.0*. Tech. rep. The OpenMath Society, 2004. URL: <http://www.openmath.org/standard/om20>.
- [Cre16] John Cremona. “The L-Functions and Modular Forms Database Project”. In: *Foundations of Computational Mathematics* 16.6 (2016), pp. 1541–1553. ISSN: 1615-3383. DOI: [10.1007/s10208-016-9306-z](https://doi.org/10.1007/s10208-016-9306-z).
- [D6.818] John Cremona et al. *Report on OpenDreamKit deliverable D6.8: GCurated Math-in-the-Middle Ontology and Alignments for GAP/SAGE/LMFDB*. Deliverable D6.8. OpenDreamKit, 2018. URL: <https://github.com/OpenDreamKit/OpenDreamKit/raw/master/WP6/D6.8/report-final.pdf>.
- [Deh+16] Paul-Olivier Dehaye et al. “Interoperability in the OpenDreamKit Project: The Math-in-the-Middle Approach”. In: *Intelligent Computer Mathematics 2016*. Conferences on Intelligent Computer Mathematics. (Bialystok, Poland, July 25–29, 2016). Ed. by Michael Kohlhasse et al. LNAI 9791. Springer, 2016. ISBN: 978-3-319-08434-3. URL: <https://github.com/OpenDreamKit/OpenDreamKit/blob/master/WP6/CICM2016/published.pdf>.
- [FGT92] William M. Farmer, Josuah Guttman, and Javier Thayer. “Little Theories”. In: *Proceedings of the 11<sup>th</sup> Conference on Automated Deduction*. Ed. by D. Kapur. LNCS 607. Saratoga Springs, NY, USA: Springer Verlag, 1992, pp. 467–581.
- [Fre+] Sebastian Freundt et al. *Symbolic Computation Software Composability Protocol (SCSCP)*. Version 1.3. URL: [https://github.com/OpenMath/scscp/blob/master/revisions/SCSCP\\_1\\_3.pdf](https://github.com/OpenMath/scscp/blob/master/revisions/SCSCP_1_3.pdf) (visited on 08/27/2017).
- [GAP] The GAP Group. *GAP – Groups, Algorithms, and Programming*. URL: <http://www.gap-system.org> (visited on 08/30/2016).
- [Gon+13] Georges Gonthier et al. “A Machine-Checked Proof of the Odd Order Theorem”. In: *Interactive Theorem Proving*. Ed. by Sandrine Blazy, Christine Paulin-Mohring, and David Pichardie. Vol. 7998. LNCS. Springer, 2013, pp. 163–179. ISBN: 978-3-642-39633-5. DOI: [10.1007/978-3-642-39634-2\\_14](https://doi.org/10.1007/978-3-642-39634-2_14).
- [Gro] The Object Management Group. *The Common Object Request Broker Architecture*. <http://www.corba.org/>.
- [Jup] *Project Jupyter*. URL: <http://www.jupyter.org> (visited on 08/22/2017).
- [Koh+11] Michael Kohlhasse et al. “The Planetary System: Web 3.0 & Active Documents for STEM”. In: *Procedia Computer Science* 4 (2011): *Special issue: Proceedings of the International Conference on Computational Science (ICCS)*. Ed. by Mitsuhiro Sato et al. Finalist at the Executable Paper Grand Challenge, pp. 598–607. DOI: [10.1016/j.procs.2011.04.063](https://doi.org/10.1016/j.procs.2011.04.063). URL: <http://kwarc.info/kohlhasse/papers/epc11.pdf>.

- [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<sup>A</sup>T<sub>E</sub>X as a Semantic Markup Format”. In: *Mathematics in Computer Science 2.2* (2008), pp. 279–304. URL: <https://kwarc.info/kohlhase/papers/mcs08-stex.pdf>.
- [LK16] Enxhell Luzhnica and Michael Kohlhase. “Formula Semantification and Automated Relation Finding in the OEIS”. In: *Mathematical Software - ICMS 2016 - 5th International Congress*. Ed. by Gert-Martin Greuel et al. Vol. 9725. LNCS. Springer, 2016. DOI: [10.1007/978-3-319-42432-3](https://doi.org/10.1007/978-3-319-42432-3). URL: <http://kwarc.info/kohlhase/papers/icms16-oeis.pdf>.
- [Lmf] *LMFDB - API*. <http://www.lmfdb.org/api/>. visited on: 09/17/2017.
- [LMFDB] The LMFDB Collaboration. *The L-functions and Modular Forms Database*. URL: <http://www.lmfdb.org> (visited on 02/01/2016).
- [Mit] *MitM/Foundation*. URL: <https://gl.mathhub.info/MitM/Foundation> (visited on 09/01/2017).
- [Mit03] Nilo Mitra. *SOAP 1.2 Part 0: Primer*. W3C Recommendation. 2003. URL: <http://www.w3.org/TR/2003/REC-soap12-part0-20030624>.
- [MMT] *MMT – Language and System for the Uniform Representation of Knowledge*. project web site. URL: <https://uniformal.github.io/> (visited on 08/30/2016).
- [MRK18] Dennis Müller, Florian Rabe, and Michael Kohlhase. “Theories as Types”. In: ed. by Didier Galmiche, Stephan Schulz, and Roberto Sebastiani. Springer Verlag, 2018. URL: <http://kwarc.info/kohlhase/papers/ijcar18-records.pdf>.
- [Mül+17a] Dennis Müller et al. “Alignment-based Translations Across Formal Systems Using Interface Theories”. In: *Fifth Workshop on Proof eXchange for Theorem Proving - PxTP 2017*. 2017. URL: <http://jazzpirate.com/Math/AlignmentTranslation.pdf>.
- [Mül+17b] Dennis Müller et al. “Classification of Alignments between Concepts of Formal Mathematical Systems”. In: *Intelligent Computer Mathematics (CICM) 2017*. Conferences on Intelligent Computer Mathematics. Ed. by Herman Geuvers et al. LNAI 10383. Springer, 2017. ISBN: 978-3-319-62074-9. DOI: [10.1007/978-3-319-62075-6](https://doi.org/10.1007/978-3-319-62075-6). URL: <http://kwarc.info/kohlhase/papers/cicm17-alignments.pdf>.
- [OC] *OSCAR Computer Algebra System*. URL: <https://oscar.computeralgebra.de/> (visited on 09/04/2018).
- [OEIS] OEIS Foundation Inc., ed. *The On-Line Encyclopedia of Integer Sequences*. URL: <http://oeis.org> (visited on 05/28/2017).
- [OMCP] *OpenMath CD Group: polygrp*. URL: <http://www.openmath.org/cdgroups/polygrp.html> (visited on 09/01/2017).
- [PB] *Protocol Buffers*. URL: <https://developers.google.com/protocol-buffers/> (visited on 09/04/2018).
- [POMa] *An OpenMath 2.0 implementation in Python*. URL: <https://github.com/OpenMath/py-openmath> (visited on 09/04/2016).
- [POMb] *An SCSCP module for Python*. URL: <https://github.com/OpenMath/py-scscp> (visited on 09/04/2016).
- [Rab12] Florian Rabe. “A Query Language for Formal Mathematical Libraries”. In: *Intelligent Computer Mathematics*. Conferences on Intelligent Computer Mathematics (CICM). (Bremen, Germany, July 9–14, 2012). Ed. by Johan Jeuring et al. LNAI

7362. Berlin and Heidelberg: Springer Verlag, 2012, pp. 142–157. ISBN: 978-3-642-31373-8. arXiv: [1204.4685 \[cs.LO\]](#).
- [Rab13] Florian Rabe. “The MMT API: A Generic MKM System”. In: *Intelligent Computer Mathematics*. Conferences on Intelligent Computer Mathematics. (Bath, UK, July 8–12, 2013). Ed. by Jacques Carette et al. Lecture Notes in Computer Science 7961. Springer, 2013, pp. 339–343. ISBN: 978-3-642-39319-8. DOI: [10.1007/978-3-642-39320-4](#).
- [RK13] Florian Rabe and Michael Kohlhase. “A Scalable Module System”. In: *Information & Computation* 0.230 (2013), pp. 1–54. URL: <http://kwarc.info/frabe/Research/mmt.pdf>.
- [Sage] The Sage Developers. *SageMath, the Sage Mathematics Software System*. URL: <http://www.sagemath.org> (visited on 09/30/2016).
- [Slo03] Neil J. A. Sloane. “The On-Line Encyclopedia of Integer Sequences”. In: *Notices of the AMS* 50.8 (2003), p. 912.
- [SNG] *Singular*. URL: <https://www.singular.uni-kl.de/> (visited on 08/22/2017).
- [sTeX] *KWARC/sTeX*. URL: <https://github.com/KWARC/sTeX> (visited on 05/15/2015).
- [Wie17] Tom Wiesing. “Enabling Cross-System Communication Using Virtual Theories and QMT”. Master’s Thesis. Bremen, Germany: Jacobs University Bremen, Aug. 2017. URL: <https://github.com/tkw1536/MasterThesis/raw/master/thesis.pdf>.
- [WSDL07] David Booth and Canyang Kevin Liu. *Web Services Description Language (WSDL) Version 2.0 Part 0: Primer*. W3C Recommendation. 2007. URL: <http://www.w3.org/TR/wsdl20-primer>.
- [LMF] The LMFDB Collaboration. *The L-functions and Modular Forms Database*.

Disclaimer: this report, together with its annexes and the reports for the earlier deliverables, is self contained for auditing and reviewing purposes. Hyperlinks to external resources are meant as a convenience for casual readers wishing to follow our progress; such links have been checked for correctness at the time of submission of the deliverable, but there is no guarantee implied that they will remain valid.