

---

# Computerising Mathematical Text with MathLang

FAIROUZ KAMAREDDINE, *School of Mathematical and Computer Sciences, Heriot-Watt University, Edinburgh, Scotland,*  
*E-Mail: fairouz@macs.hw.ac.uk*

JOE B. WELLS, *School of Mathematical and Computer Sciences, Heriot-Watt University, Edinburgh, Scotland,*  
*E-Mail: jbw@macs.hw.ac.uk*

CHRISTOPH ZENGLER, *School of Mathematical and Computer Sciences, Heriot-Watt University, Edinburgh, Scotland,*  
*E-Mail: zengler@macs.hw.ac.uk*

## Abstract

Mathematical texts can be computerised in many ways that capture differing amounts of the mathematical meaning. At one end, there is document imaging, which captures the arrangement of black marks on paper, while at the other end there are proof assistants (e.g., Mizar, Isabelle, Coq, etc.), which capture the full mathematical meaning and have proofs expressed in a formal foundation of mathematics. In between, there are computer typesetting systems (e.g., L<sup>A</sup>T<sub>E</sub>X and Presentation MathML) and semantically oriented systems (e.g., Content MathML, OpenMath, OMDoc, etc.).

The MathLang project was initiated in 2000 by Fairouz Kamareddine and Joe Wells with the aim of developing an approach for computerising mathematical texts which is flexible enough to connect the different approaches to computerisation, which allows various degrees of formalisation, and which is compatible with different logical frameworks (e.g., set theory, category theory, type theory, etc.) and proof systems. The approach is embodied in a computer representation, which we call MathLang, and associated software tools, which are being developed by ongoing work. Four Ph.D. students (Manuel Maarek (2002/2007), Krzysztof Retel (since 2004), Robert Lamar (since 2006)), and Christoph Zengler (since 2008) and over a dozen master's degree and undergraduate students have worked on MathLang. The project's progress and design choices are driven by the needs for computerising representative mathematical texts chosen from various branches of mathematics.

Currently, MathLang supports entry of mathematical text either in an XML format or using the T<sub>E</sub>X<sub>MACS</sub> editor. Methods are provided for adding, checking, and displaying various information *aspects*. One aspect is a kind of weak type system that assigns categories (term, statement, noun (class), adjective (class modifier), etc.) to parts of the text, deals with binding names to meanings, and checks that a kind of grammatical sense is maintained. Another aspect allows weaving together mathematical meaning and visual presentation and can associate natural language text with its mathematical meaning. Another aspect allows identifying chunks of text, marking their roles (theorem, definition, explanation, example, section, etc.), and indicating relationships between the chunks (A uses B, A contradicts B, A follows from B, etc.). Software tool support can use this aspect to check and explain the overall logical structure of a text.

Further aspects are being designed to allow adding additional formality to a text such as proof structure and details of how a human-readable proof is encoded into a fully formalised version (previously we used Mizar and Isabelle but here, for the first time we develop the MathLang formalisation into Coq). [24] surveyed the status of the MathLang project up to November 2007. This paper

## 2 Computerising Mathematical Text with MathLang

picks on from that survey, fills in a number of formalisation and implementation gaps and creates a formalisation path via MathLang into Coq. We show for the first time how the DRa information can be used to automatically generate proof skeletons for different theorem provers, we formalise and implement the textual order of a text and explain how it can be derived from the original text. Our proposed generic algorithm (for generating the proof skeleton which depends on the original mathematical text and the desired theorem prover), is highly configurable and caters for arbitrary theorem provers. This generic algorithm as well as all the new algorithms and concepts we present here, are implemented in our software tool. Furthermore, we give hints for the development of an algorithm which is able to convert parts of a CGa annotated text automatically into the syntax of a special theorem prover.

To test our approach we create the complete path of encoding in and formalising through MathLang, for the first chapter of Landau's book "Grundlagen der Analysis". We started with the plain text document, annotated categories and mathematical roles to chunks of text and automatically generated in MathLang, a Coq and a Mizar proof skeleton for the chapter. We used our hints to convert parts of the CGa annotated text of Landau's first chapter into Coq. Both the Coq proof skeleton and the converted CGa parts into Coq, simplified the process of the full formalisation of the first chapter of Landau's book in Coq.

This full formalisation in Coq illustrates the use of MathLang in different theorem proving settings. Previously we have only used Mizar and Isabelle. With the work in this paper, we show that MathLang is able to handle different formal foundations. Since November 2007 [24], significant progress has been made on MathLang which is reported here for the first time.

*Keywords:* mathematical knowledge management, mathematical vernacular, mathematical type-setting, logical foundations of mathematics, proof assistants, proof checkers, theorem provers

## 1 Background and motivation

The MathLang project is based on considering these two questions:

1. What is the relationship between the logical foundations of mathematical reasoning and the actual practice of mathematicians?
2. In what ways can computers support the development and communication of mathematical knowledge?

### 1.1 Logical Foundations

Our first question, of the relationship between the practice of mathematics and its logical foundations, has been an issue for at least two millennia. Logic was already influential in the study and development of mathematics since the time of the ancient Greeks. One of the main issues was already known by Aristotle, namely that for a logical/mathematical proposition  $\Phi$ ,

- given a purported proof of  $\Phi$ , it is not hard to check whether the argument really proves  $\Phi$ , but
- in contrast, if one is asked to find a proof of  $\Phi$ , the search may take a very long time (or even go forever without success) even if  $\Phi$  is true.

Aristotle used logic to reason about everything (mathematics, law, farming, medicine, etc.). A formal logical style of deductive reasoning about mathematics was introduced in Euclid's geometry [14].

The 1600s saw an increase in the importance of logic. Researchers like Leibniz wanted to use logic to address not just mathematical questions but also more esoteric

questions like the existence of God. In the 1800s, the need for a more precise style in mathematics arose, because controversial results had appeared in analysis [15]. Some controversies were solved by Cauchy's precise definition of convergence in his *Cours d'Analyse* [4], others benefitted from the more exact definition of real numbers given by Dedekind [9], while at the same time Cantor was making a tremendous contribution to the formalisation of set theory and number theory [2, 3] and Peano was making influential steps in formalised arithmetic [32] (albeit without an extensive treatment of logic or quantification).

In the last decades of the 1800s, the contributions of Frege made the move toward formalisation much more serious. Frege found

*“... the inadequacy of language to be an obstacle; no matter how unwieldy the expressions I was ready to accept, I was less and less able, as the relations became more and more complex, to attain precision”*

Based on this understanding of a need for greater preciseness, Frege presented *Begriffsschrift* [10], the first formalisation of logic giving logical concepts via symbols rather than natural language. “Begriffsschrift” is the name both of the book and of the formal system the book presents. Frege wrote this about the Begriffsschrift:

*“Its first purpose, therefore, is to provide us with the most reliable test of the validity of a chain of inferences and to point out every presupposition that tries to sneak in unnoticed, so that its origin can be investigated.”*

Later, Frege wrote the *Die Grundlagen der Arithmetik* and *Grundgesetze der Arithmetik* [11, 12, 38] where he argued that mathematics is a branch of logic and described arithmetic in the Begriffsschrift. Frege's Grundgesetze was the culmination of his work on building a formal foundation for mathematics.

One of the major issues in the logical foundations of mathematics is that the naive approach of Frege's Grundgesetze (and Cantor's earlier set theory) is inconsistent. Russell discovered a paradox in Frege's system (and also Russell's own system) that allows proving a contradiction, from which everything can be proven, including all the false statements [38, 15]. The need to build logical foundations for mathematics that do not suffer from such paradoxes has led to many diverging approaches. Russell invented a form of *type theory* which he used in the famous *Principia Mathematica* [39]. Others have subsequently introduced many kinds of type theories and modern type theories are quite different from Russell's. Brouwer introduced a different direction, that of *intuitionism*. Later, ideas from intuitionism and type theory were combined, and even extended to cover the power of classical logic (which Brouwer's intuitionism rejects). Zermelo followed a different direction in introducing an axiomatisation of set theory [40], later extended by Fraenkel and Skolem to form the well known Zermelo/Fraenkel (ZF) system. In yet another direction, it is possible to use category theory as a foundation. And there are other proposed foundations, too many to discuss here.

Despite the variety of possible foundations for mathematics, in practice real mathematicians do not express their work in terms of a foundation. It seems that most modern mathematicians tend to *think* in terms that are compatible with ZFC (which is ZF extended with the Axiom of Choice), but in practice they almost never write the full formal details. And it is quite rare for mathematicians to do their thinking

#### 4 Computerising Mathematical Text with MathLang

while regarding a type theory as the foundation, even though type theories are among the most thoroughly developed logical foundations (in particular with well developed computer proof software systems). Instead, mathematicians write in a kind of *common mathematical language* (CML) (sometimes called a *mathematical vernacular*), for a number of reasons:

- Mathematicians have developed conventional ways of using nouns, adjectives, verbs, sentences, and larger chunks of text to express mathematical meaning. However, the existing logical foundations do not address the convenient use of natural language text to express mathematical meanings.
- Using a foundation requires picking one specific foundation, and any foundation commits to some number of fixed choices. Such choices include what kinds of mathematical objects to take as the primitives (e.g., sets, functions, types, categories, etc.), what kinds of logical rules to use (e.g., “natural deduction” vs. “logical deduction”, whether to allow the full power of classical logic, etc.), what kinds of syntax and semantics to allow for logical propositions (first-order vs. higher-order), etc. Having made some initial choices, further choices follow, e.g., for a set theory one must then choose the axioms (Zermelo/Fraenkel, Tarski/Grothendieck, etc.), or for a type theory the kinds of types and the typing rules (Calculus of Constructions, Martin-Löf, etc.). Fixed choices make logical foundations undesirable to use for three reasons:
  - Much of mathematics can be built on top of all of the different foundations. Hence, committing to a particular foundation would seem to unnecessarily limit the applicability of mathematical results.
  - The details of how to build some mathematical concepts can vary quite a bit from foundation to foundation. Issues that cause difficulty include how to handle “partial functions”, induction, reasoning modulo equations, etc. Because these issues *can* be handled in all foundations, practising mathematicians tend to see the low-level details of these issues as inessential, irrelevant, and uninteresting, and are not willing to write the low-level details.
  - Some mathematics only works for some foundations. Hence, for a mathematician to develop the specialised expertise needed to express mathematics in terms of one particular foundation would seem to unnecessarily limit the scope of mathematics the mathematician could address. An ordinary mathematician is happy to be reassured by a mathematical logician that what they are doing *can* be expressed in some foundation, but the ordinary mathematician usually does not care to work out precisely how.

Moreover there is no universal agreement as to which is the best logical foundation.

- In practice, formalising a mathematical text in any of the existing foundations is an extremely time-consuming, costly, and mentally painful activity. Formalisation also requires special expertise in the particular foundation used that goes far beyond the ordinary expertise of even extremely good mathematicians. Furthermore, mathematical texts formalised in any of the existing foundations are generally structured in a way which is radically different from what is optimal for the human reader’s understanding, and which is difficult for ordinary mathematicians to use. (Some proof software systems like Mizar (which is based on Tarski/Grothendieck set theory) attempt to reduce this problem, and partially

succeed.) What is a single step in a usual human-readable mathematical text may turn into a multitude of smaller steps in a formalised version. New details completely missing from the human-readable version may need to be woven throughout the entire text. The original text may need to be reorganised and reordered so radically that it seems like it is almost turned inside out in the formal version.

So, although mathematics was a driving force for the research in logic in the 19th or 20th century, mathematics and logic have kept a distance from each other. Practising mathematicians do not want to use formal mathematical logic and have for centuries done most mathematical work outside of the strict boundaries of formal logic.

## 1.2 Computerisation of Mathematical Knowledge

Our second question, of how to use mechanical computers to support mathematical knowledge, is more recent but is unavoidable since automation and computation can provide tremendous services to mathematics. (There are also extensive opportunities for combining progress in logic and computerisation not only in mathematics but also in other areas: program verification, bio-informatics, chemistry, music, etc.)

Mechanical computers have been used from their beginning for mathematical purposes. Starting in the 1960s, computers began to play a role in handling not just computations, but abstract mathematical knowledge. Nowadays, computers can represent mathematical knowledge in various ways:

- Pixel map images of pages of mathematical articles may be stored on the computer. While useful, it is extremely difficult for computer programs to access the semantics of mathematical knowledge represented this way. Even keyword searching is difficult, because first OCR (Optical Character Recognition) must be performed and high quality OCR for mathematical material is still an area with significant research challenges rather than a proven technology (e.g., there is great difficulty with matrices [25]).
- Typesetting systems like  $\text{\LaTeX}$  or  $\text{\TeX}_{\text{MACS}}$  [37], can be used with mathematical texts for editing them and formatting them for viewing or printing. The document formats of these systems can also be used for storage and archiving. Such systems provide good defaults for visual appearance and allow fine control when needed. They support commonly needed document structures and allow custom structures to be created, at least to the extent of being able to produce the correct visual appearance.

Unfortunately, unless the mathematician is amazingly disciplined, the logical structure of symbolic formulas is not directly represented. Furthermore, the logical structure of mathematics as embedded in natural language text is not represented at all. This makes it difficult for computer programs to access document semantics because fully automated discovery of the semantics of natural language text still performs too poorly to use in practical systems. Even human-assisted semi-automated semantic analysis of natural language is primitive, and we are aware of no such systems with special support for mathematical text. As a consequence, there is generally no computer support for checking the correctness of mathematics represented this way or for doing searching based on semantics (as opposed to keywords).

## 6 Computerising Mathematical Text with MathLang

- Mathematical texts can be written in more semantically oriented document representations like OpenMath [1] and OMDoc [26], Content MathML [5], etc. There is generally support for converting from these representations to typesetting systems like  $\text{\LaTeX}$  or Presentation MathML in order to produce readable and printable versions of the mathematical text. These systems are better than the typesetting systems at representing the knowledge in a computer-accessible way. Some aspects of the semantics of symbolic formulas can be represented in some of these systems. Unfortunately, in practice it is still difficult to have enough control over visual presentation with representations like OMDoc, so practising mathematicians still prefer to use the typesetting systems.

Systems like OMDoc share the same difficulties with accessing the logical structure of mathematics embedded in natural language text as mentioned above for typesetting systems. Although systems like OMDoc have ways to associate symbolic formulas with uninterpreted chunks of natural language text, these chunks are opaque to the computer and there is no method for checking that these associations are correct.

A separate weakness is that although there is support for semantics, unlike proof systems (see below), OMDoc and similar systems do not have good support for expressing the semantics in terms of a logical foundation of mathematics. Also, type checking symbolic formulas (beyond mere arity checking), which is an important tool for ensuring that symbolic formulas even have meaningful semantics, is not generally handled by these systems.

- There are software systems like *proof assistants* (sometimes called *proof checkers*, these include Coq, Isabelle, Mizar, Isar, etc.) and automated *theorem provers* (Boyer-Moore, Otter, etc.), which we collectively call *proof systems*. Each proof system provides a formal languages for writing mathematics based on some foundation of logic and mathematics. Work on computer support for formal foundations began in the late 1960s with work by de Bruijn on Automath (AUTOMating MATHEmatics) [31]. Automath supported automated checking of the full correctness of a mathematical text written in Automath’s formal language. Since then, many proof systems have been built to mechanically check logic, mathematics, and computer software (e.g., Boyer-Moore, Isabelle, HOL, Coq, etc.). Generally, these systems support checking of full correctness, and it is possible in theory (although not necessarily easy) for computer programs to access and manipulate the semantics of the mathematical statements.

Unfortunately, there are great disadvantages in using these systems. First, all of the problems mentioned for logical foundations in section 1.1 are incurred, e.g., the enormous expense of formalisation. Furthermore, one must choose a specific proof system (Isabelle, Coq, Mizar, PVS, etc.) and each software system has its own advantages and pitfalls and takes quite some time to learn. In practice, some of these systems are only ever learnt from a “master” in an “apprenticeship” setting. Most proof systems have no meaningful support for the mathematical use of natural language text. A notable exception is Mizar, which however requires the use of natural language in a rigid and somewhat inflexible way. Most proof systems suffer from the use of proof tactics, which make it easier to construct proofs and make proofs smaller, but obscure the reasoning for readers because the meaning of each tactic is often *ad hoc* and implementation-dependent. As a result

of these and other disadvantages, ordinary mathematicians do not generally read mathematics written in the language of a proof system, and are usually not willing to spend the effort to formalise their own work in a proof system.

- Computer algebra systems (e.g., Maxima, Maple, Mathematica, etc.) are widely used software environments designed for carrying out computations, primarily symbolic but sometimes also numeric. Each CAS has a language for writing mathematical expressions and statements and for describing computations. The languages can also be used for representing mathematical knowledge. The main advantage for such a language is integration with a CAS.

Typically, a CAS language is not tied to any specific foundation and has little or no support for guaranteeing correctness of mathematical statements. A CAS language also typically has little or no support for embedded natural language text, or for precise control over typesetting. So a CAS is often used for calculating results, but these results are usually converted into some other language or format for dissemination or verification. Nonetheless, there are useful possibilities for using a CAS for archiving and communicating mathematical knowledge.

We are gradually developing a system named MathLang which we hope will be usable as a bridge between more than one of the above categories of ways of representing mathematical knowledge. We also aim for MathLang to make easier (without requiring) the partial or full formalisation of mathematical texts in some foundation.

## 2 An overview of MathLang

### 2.1 The goals of MathLang

Sections 1.1 and 1.2 described issues with the practice of mathematics: the difficulty for the normal mathematician in directly using a formal foundation, and the disadvantages of the various computer representations of mathematics. To learn how to address these issues, in 2000 we (Kamareddine and Wells) began the MathLang project to develop a new mathematical language called MathLang<sup>1</sup>, so that texts written in CML (the common mathematical language, expressed either with pen and paper, or L<sup>A</sup>T<sub>E</sub>X) is written instead in MathLang in a way that satisfies these goals:

1. A MathLang text should support the usual features of CML: natural language text, symbolic formulas, images, document structures, control over visual presentation, etc. And the usual computer support for editing such texts should be available.
2. It should be possible to write a MathLang text in a way that is significantly less ambiguous than the corresponding CML text. A MathLang text should *somehow* support representing the text's mathematical semantics and structure. The support for semantics should cover not just individual pieces of text and symbolic formulas but also the entire document and the document's relationship to other documents (to allow building connected libraries). The degree of formality in representing the mathematical semantics should be flexible, and at least one choice of degree of formality should be both inexpensive *and* useful. There should be some automated checking of the well-formedness of the mathematical semantics.

---

<sup>1</sup>We always named the project MathLang but initially named the proposed language NML (New Mathematical Language).

## 8 Computerising Mathematical Text with MathLang

3. The structure of a MathLang text should follow the structure of the corresponding CML, so that the experience of reading and writing MathLang should be close to that of reading and writing CML. This should make it easier for an author to see and have confidence that a MathLang text correctly represents their intentions. Thus, if any foundational formal systems are used in MathLang, then MathLang should somehow adapt the formal systems to the needs of the authors and readers, rather than requiring the authors and readers to adapt their thinking to fit the rigid confines of any existing foundations.
4. The structure of a MathLang text should make it easier to support further post-authorship computer manipulations that respect its mathematical structure and meaning. Examples include semantics-based searches, computations via computer algebra systems, extraction of proof sketches (to be completed into a full formalisation in a proof system), etc.
5. A particular important case of the previous point is that MathLang should support (but not require) interfacing with proof systems so that a MathLang text can contain full formal details in some foundation and the formalisation can be automatically verified.
6. Authoring of a MathLang text should not be significantly harder for the ordinary mathematician than authoring  $\text{\LaTeX}$ . Features of MathLang that the author does not want (such as formalisation in a proof system) should not require any extra effort from an author.
7. The design of MathLang should be compatible with (as yet undetermined) future extensions to support additional uses of mathematical knowledge. Also, the design of MathLang should make it easy to combine with existing languages (e.g., OM-`Doc`, `\TeX_{MACS}`). This way, MathLang might end up being a method for extending an existing language in addition to (or instead of) a language on its own.

None of the previously existing representations for mathematical texts satisfies our goals, so we have been developing new techniques.

MathLang is intended to support different degrees of formalisation. Furthermore, for those documents where full formalisation is a goal, MathLang is intended to allow this to be accomplished in gradual steps. Some of the motivations for varying degrees of formalisation have already been discussed in sections 1.1 and 1.2. Full formalisation is sometimes desirable, but also is often undesirable due to its expense and the requirement to commit to many inessential foundational details. Partial formalisation can sometimes be desirable for various reasons; as examples, it has the potential to be helpful with automated checking, semantics-based searching and querying, and interfacing with computer algebra systems (and other mathematical computation environments). Partial formalisation can be carried out to different degrees:

- The abstract syntax trees of symbolic formulas can be represented accurately. This is usually missing when using systems like  $\text{\LaTeX}$  or Presentation MathML, while more semantically oriented systems usually provide this to some degree. This can be used to provide editing support for algebraic rearrangements and simplifications, and can help with interfacing with computer algebra systems.
- The mathematical structure of natural language text can be represented in a way similar to how symbolic formulas are handled. Furthermore, mixed text and

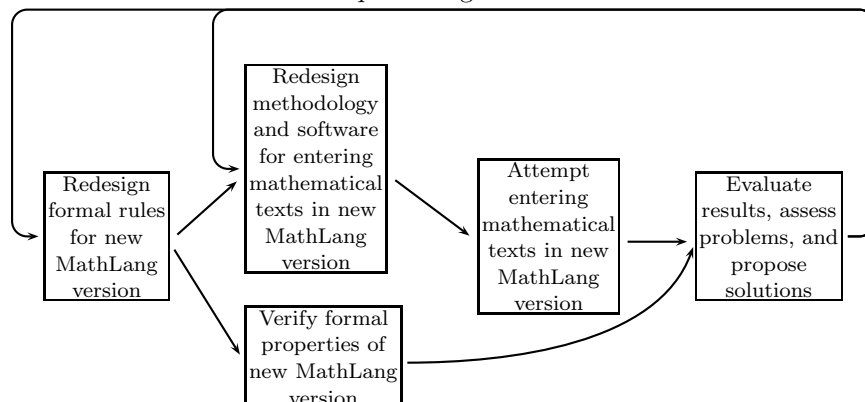


FIG. 1. Iterative MathLang development process

symbols can be handled. This can help in the same way as capturing the structure of symbolic formulas can help. Nearly all previous systems do not support handling natural language text in this way.

- A weak type system can be used to check simple grammatical conditions without checking full semantic sensibility.
- Justifications (inside proofs and between formal statements) can be linked (without necessarily always indicating precisely how they are used). Some examples of potential uses of this feature include the following:
  - Extracting only those parts of a document that are relevant to specific results. (This could be useful in educational systems.)
  - Checking that each instance of apparently circular reasoning is actually handled via induction.
  - Calculating proof gaps as a first step toward fuller formalisation.
- If one commits to a foundation (or in some cases, to a family of foundations), one can start to use more sophisticated type systems in formulas and statements for checking more aspects of well-formedness.
- And there are further possibilities.

## 2.2 The process of designing MathLang

We are gradually refining the design of MathLang based on experience testing the use of MathLang for representative mathematical texts. This iterative process is depicted in figure 1. We started from an initial design based on ideas from WTT (see below). Each iteration of the overall procedure is to test the design by evaluating encodings of real mathematical texts, during which issues and difficulties are encountered, which lead to new needs being discovered and corresponding design adjustments. The design includes formal rules for the representation of mathematical texts, as well as patterns and methodology for entering texts in this representation, and supporting software.

Our choice of mathematical texts for testing is primarily oriented toward texts that represent the variety of mathematical writing by ordinary mathematicians rather than texts that represent the interests of formalists and mathematical logicians. Much of

our testing has been with pre-existing texts. In some cases, we have chosen texts that have previously been formalised by others so that we can compare the representations, e.g., *A Compendium of Continuous Lattices* [13] of which at least 60% has been formalised in Mizar [33], and Landau’s *Foundations of Analysis* [28] which was fully formalised in Automath [35]. In other cases, we have chosen texts of historical importance which are known to have errors to ensure that MathLang’s design will not exclude them, e.g., Euclid’s *Elements* [14]. We have chosen other texts to exercise other aspects of MathLang. We have also been testing the authoring of new texts.

### 2.3 *The original starting point: WTT*

For purely historical reasons it is interesting to point out that the initial design of MathLang was based heavily on ideas from the Weak Type Theory (WTT) of Nederpelt and Kamareddine [23], which in turn was heavily inspired by the Mathematical Vernacular (MV) of de Bruijn [8].

In the terminology of WTT, a document is a *book* which is a sequence of *lines*, each of which is a pair of a *sentence* (a *statement* or a *definition*) and a *context* of facts (*declarations* or *statements*) assumed in the sentence. WTT has four ways of introducing names. A *definition* introduces a name whose scope is the rest of the book and associates the name with its meaning. A name introduced by a definition can have *parameters* whose scope is the body of the definition. A *declaration* in a context introduces a name (with no parameters) whose scope is only the current line. Finally, a *preface* gives names whose scope is the document; names introduced by prefaces have parameters but unlike definitions their meanings are not provided (and thus presumed to be given externally to the document). Declarations, definitions, and statements can contain *phrases* which are built from *terms*, *sets*, *nouns*, and *adjectives*. Using the terminology of object-oriented programming languages, nouns act like *classes* and adjectives act like *mixins* (a special kind of function from classes to classes). WTT uses a weak type system with types like `noun`, `set`, `term`, `adjective`, `statement`, `definition`, `context`, and `book` to check some basic well-formedness conditions. Sets are used when something is definitely known to be a set and the richer structure of a noun is not needed, and terms are used for things that are not sets (and sometimes for sets in cases where the type system is too weak).

Although WTT provides many useful ideas, the definition of WTT has many limitations. The many different ways of introducing names are too complicated and awkward. WTT provides no way to indicate which statements are used to justify other statements and in general does not deal with *proofs* and logical correctness. WTT provides no ways to present the structure of a text to human readers; there is no way of grouping statements and identifying their mathematical/discourse roles such as *theorem*, *lemma*, *conjecture*, *proof*, *section*, *chapter*. WTT provides no way to give human names to statements (e.g., “Newman’s Lemma”). WTT provides no way to use in one document concepts defined in another document.

### 2.4 *The current MathLang design*

The current MathLang design has developed through the experience of a large number of students, including both shorter projects (over a dozen projects by either 4th year

undergraduate students or M.Sc. students) and Ph.D. studies (by 4 students: Maarek, Retel, Lamar and Zengler). Every student has done work to write in MathLang some piece of mathematical text. The experience gained from this has led to the current design of MathLang which is (currently) divided into three *aspects*:

- The Core Grammatical aspect (CGa) [21, 22, 29, 17] takes the best features of WTT [23] and MV [8], simplifies difficult aspects of WTT, and enhances the nouns and adjectives of WTT with ideas from object-oriented programming so that nouns are more like classes and adjectives are more like mixins. In CGa, the different kinds of name-introducing forms of WTT are unified; all definitions by default have indefinite forward scope and a local scope operator is used to allow local definitions. The basic unit becomes the *step*, which can be either a definition, a statement (a phrase that asserts something), or a *block* which is merely a grouping of steps. CGa keeps WTT’s notions of nouns, adjectives, terms, sets, definitions, and statements. CGa provides a kind of *grammar* for well-formed mathematics with grammatical categories and allows checking for some basic well-formedness conditions (e.g., the origin of all names and symbols can be tracked).
- The Text and Symbol aspect (TSa) [20, 16, 29, 17] allows integrating normal typesetting and authoring software with the mathematical structure represented with CGa. TSa allows weaving together usual mathematical authoring representations such as L<sup>A</sup>T<sub>E</sub>X, XML, or T<sub>E</sub>X<sub>MACS</sub> with CGa data. Thanks to a notion of *souring* rules (called “souring” because it does the opposite of what is usually called *syntactic sugar*), TSa allows the structure of the mathematical text to follow the structure of the CML as conceived by the mathematician.
- The Document Rhetorical aspect (DRa) [19] supports identifying portions of a text and expressing the relationships between them. Any portion of text (e.g., a phrase, a step, a block, etc.) can be given an identity. Many kinds of relationships can be expressed between identified pieces of text. For example, a chunk of text can be identified as a “theorem”, and another can be identified as the “proof” of that theorem. Similarly, one chunk of text can be a “subsection” or “chapter” of another. Given these identified relationships, it becomes possible to do computations to check whether all dependencies are identified, to check whether the relationships are sensible or possibly problematic (and whether therefore the author should be warned), and to extract and explain the logical structure of a text. Dependencies identified this way have been used in generating formal proof sketches and identifying the proof holes that remain to be filled. This paper presents further formalisation and implementation of notions related to DRa.

In addition to the design of MathLang itself, there has been work on relating a MathLang text to a fully formalised version of the text. Using the information in the CGa and DRa aspects of a MathLang text, we have developed in [18] a procedure for producing a corresponding Mizar document, first as a proof sketch with holes and then as a fully completed proof. We have recently begun to work also on doing this with Isabelle. In this paper, we make further progress in completing the path in MathLang in order to reach full formalisation and we introduce a third theorem prover (Coq) as a test bed for MathLang (in addition to Mizar and Isabelle). We develop the proof skeleton idea presented earlier in [18] specifically for Mizar, into an automatically generated proof skeleton in a choice of theorem provers (including

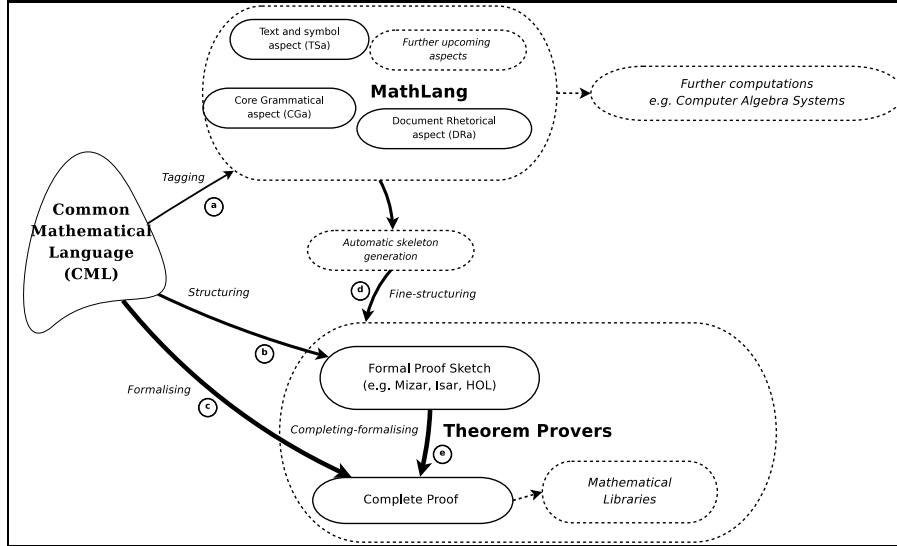


FIG. 2. Overall situation of work in MathLang

Mizar, Isar and Coq). To achieve this, we give a generic algorithm for proof skeleton generation which takes the required prover as one of its arguments. We also give hints for the development of a generic algorithm which automatically converts parts of a CGa annotated text into the syntax of the theorem prover it is given as an argument.

Figure 2 (adapted from [18]) diagrams the overall current situation of work on MathLang. In the rest of this paper, we discuss the aspects CGa, TSa, and DRa in more detail, we introduce the generic automatic proof skeleton generator and how parts of the CGa annotated text can be formalised into a theorem prover. We also discuss the work on interfacing MathLang with the Coq system. CGa and TSa have not changed since November 2007 [24], but we have since made progress on the formalisation and implementation of notions related to the DRa and on the full formalisation path into an arbitrary theorem prover. We have also introduced a third theorem prover (Coq) as a test bed for Mathlang where we carry out the most extensive example to date of how a text can pass through all the path of MathLang from the version written by the mathematician into a fully formalised text. Since November 2007 [24], significant progress has been made on MathLang which is reported here for the first time.

### 3 The aspects of MathLang

#### 3.1 The Core Grammatical aspect (CGa)

CGa [21, 22, 29, 17] is a formal language inspired initially by WTT [23] and MV [8], and then later shaped by repeated experiences of entering mathematical texts in early versions of CGa.

The basic constructs of CGa are the *step* and the *expression*. The tasks handled in WTT by books, prefaces, lines, declarations, definitions, and statements are all represented as steps in CGa. A step can be a *block*  $\{s_1, \dots, s_n\}$ , which is merely a

sequence of steps. A step can be a *local scoping*  $s_1 \triangleright s_2$ , which is a pair of steps  $s_1$  and  $s_2$  where the definitions and declarations of  $s_1$  are restricted in scope to  $s_2$  and the assertions of  $s_1$  are assumptions of  $s_2$ . A step can also be a *definition*, a *declaration*, or an expression (which asserts a truth). Expressions are also used for the bodies of definitions and inside the types in declarations. The possibilities for expressions include uses of defined identifiers, identifier declarations, and *noun descriptions*. A noun description allows specifying *characteristics* of a class of entities.

Here is an example. Consider this (silly) CML text:

“Given that  $\mathfrak{M}$  is a set,  $y$  and  $x$  are natural numbers, and  $x$  belongs to  $\mathfrak{M}$ , it holds that  $x + y = y + x$ .”

A straightforward encoding of the above text in CGa would be the following:

```
{M : set; y : natural_number; x : natural_number; ∈(x,M)}
▷ =(+(x,y), +(y,x))
```

This example assumes that somewhere earlier in the document there will be declarations like these:

```
...; ∈(term, set) : stat; =(term, term) : stat; natural_number : noun;
+ (natural_number, natural_number) : natural_number; ...
```

Here,  $M$ ,  $y$ ,  $x$ ,  $\in$ ,  $=$ , and  $+$  are *identifiers*<sup>2</sup> while *term*, *set*, *stat*, and *noun* are *keywords* of CGa. The semicolon, colon, comma, parentheses, braces, and right triangle ( $\triangleright$ ) symbols are part of the syntax of CGa. The statements like  $\in(\mathbf{term}, \mathbf{set}) : \mathbf{stat}$  are declarations; this example declares  $\in$  to be an operator that takes two arguments, one of type *term* and one of type *set*, and yields a result of type *stat* (statement). The statement  $M : \mathbf{set}$  is an abbreviation for  $M() : \mathbf{set}$  which declares the identifier  $M$  to have zero parameters.

CGa uses grammatical/linguistic/syntactic *categories* (also called *types*) to make explicit the grammatical role played by the elements of a mathematical text. In the above example, we see the category expressions *term*, *set*, *stat*, *noun*, and *natural\_number*. In fact, the category expression *natural\_number* acts as an abbreviation for *term(natural\_number)*, and *term*, *set*, and *noun* are abbreviations for *term(Noun {})*, *set(Noun {})*, and *noun(Noun {})*, which all use the *uncharacterised* noun description *Noun {}*. A noun description is of the form *Noun s* and describes a class of entities with *characteristics* (declared operations and true facts) defined by the step  $s$ . The arguments of the category constructors *term*, *set*, and *noun* are expressions which evaluate to noun descriptions. The category *term(e)* describes individual entities belonging to the class described by the noun expression  $e$ , and the category *set(e)* describes any set of such entities. The category *noun(e)* describes any noun which defines all the operations described by  $e$  with the same types. So in the above example, the abbreviation *term* is the type of all mathematical entities, the abbreviation *set* is the type of any set, *noun* is the type of any noun (and specifies no characteristics for it), and *natural\_number* is the type of any mathematical entity having the characteristics described by the noun *natural\_number*.<sup>3</sup> The behaviour of

<sup>2</sup>Our current implementation only allows ASCII characters in identifiers, but we plan to support any graphic Unicode characters.

<sup>3</sup>CGa has other mechanisms that allow specifying additional characteristics of the noun *natural\_number* separate from its declaration, and we assume in this example that this is done.

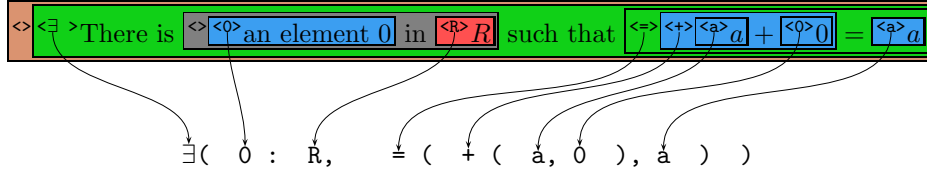


FIG. 3. Example of CGa encoding of CML text

nouns in CGa is similar to that of *classes* in object-oriented programming languages. CGa also has *adjectives* which are like object-oriented *mixins* and act as functions from nouns to nouns.

Here are some further examples of categories (see also [23]), where we put parts that do not determine the indicated category inside boxes:

Terms: the triangle  $ABC$ ; the centre of  $\boxed{ABC}$ ;  $d(\boxed{x}, \boxed{y})$ .

Nouns: a triangle; an edge of  $\boxed{ABC}$ ; a group.

Adjectives: equilateral  $\boxed{\text{triangle}}$ ; prime  $\boxed{\text{number}}$ ; Abelian  $\boxed{\text{group}}$ .

Statements:  $\boxed{P}$  lies between  $\boxed{Q}$  and  $\boxed{R}$ ;  $\boxed{5} \geq \boxed{3}$ ;  $\boxed{AB}$  is  $\boxed{\text{an edge of } ABC}$ .

Definition: a number  $p$  is prime whenever  $\boxed{\dots}$ .

These categories are all inspired from WTT. Full details of the rules of CGa are in other papers [22, 29]. For this paper it is crucial to mention the following colour codings for these CGa categories: **term** **set** **noun** **adjective** **statement** **definition** **declaration** **step** **context**.

The types of CGa are more sophisticated than the *weak types* of WTT and allow tracking which operations are meaningful in some additional cases. Although CGa's types are more powerful than WTT's, there are still significant limitations. One limitation is that higher-order types are not allowed. For example, although CGa allows the type  $(\text{term}, \text{term}) \rightarrow \text{term}$ , which is the type of an operator that takes two arguments of type **term** and returns a result of type **term**, CGa does not allow using the type  $((\text{term}) \rightarrow \text{term}, \text{term}) \rightarrow \text{term}$ , which would be the type of an operator that takes another operator as its first argument. Higher-order types can be awkwardly and crudely emulated in CGa by encapsulation with noun types, but this emulation does not work well due to the fact that CGa's type polymorphism is shallow, which is another significant limitation. To work around the weakness of CGa's type polymorphism, in practice we find ourselves often giving entities the type **term** instead of a more precise type. We continue to work on making the type system more flexible without making it too complex. It is important to understand that the goal of CGa's type system is *not* to ensure full correctness, but merely to check whether the reasoning parts of a document are coherently built in a sensible way.

The design of CGa is due to Kamareddine, Maarek and Wells [22]. The implementation of CGa is due to Maarek [29].

### 3.2 The Text and Symbol aspect (TSa)

TSa [20, 16, 29, 17] is a representation that allows interleaving pieces of CGa with pieces of CML in the form of mixtures of natural language, symbolic formulas, and formatting instructions for visual presentation. The interleaving can be at any level of granularity: meanings can be associated at a coarse grain with entire paragraphs

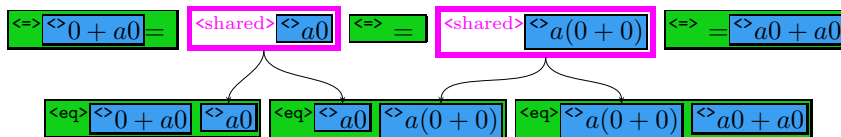


FIG. 4. Example of using sourcing in TSA to support sharing

or sections, or at a fine grain with individual words, phrases, and symbols. Arbitrary amounts of mathematically uninterpreted text can be included. The TSA representation is inspired by the XQuery/XPath Data Model (XDM) [7] used for representing the information content of XML documents. In TSA, a document  $d$  is built from the empty document ( $[]$ ) by sequencing ( $d_1, d_2$ ) and labelling ( $\ell\langle d \rangle$ ).

As an example of TSA, consider the piece of CML text and its CGa representation given in figure 3.<sup>4</sup> The example could be represented in TSA by the following fine-grained interleaving of CGa<sup>5</sup> and  $\text{\LaTeX}$ :

```

    “There is #1 such that #2.”
    <math>\exists\langle\langle\text{\#1 in \#2}\rangle\langle\langle\text{\text{“an element \#0\#”}\langle 0 \rangle, \text{\text{“\#R\#”}\langle R \rangle}\rangle\rangle\rangle\rangle,
    \text{\text{“\#1 = \#2\#”}}\langle\langle\text{\text{“\#1 + \#2”}}\langle+\langle\text{\text{“a”}}\langle a \rangle, \text{\text{“0”}}\langle 0 \rangle\rangle\rangle, \text{\text{“a”}}\langle a \rangle\rangle\rangle
    
```

This example uses the abbreviation that  $\ell$  stands for  $\ell\langle [] \rangle$ . For example, “a” $\langle a \rangle$  actually stands for “a” $\langle a\langle [] \rangle \rangle$ .

Associated with TSA are methods for extracting separately the CGa and the typesetting instructions or other visual representation. For the example, from the TSA above can be extracted the following TSA representation of just the CGa portion:

$$\exists\langle\langle 0, R \rangle, =\langle+\langle a, 0 \rangle, a \rangle\rangle$$

The CGa portion of this text can be type checked and used for processing that needs to know the mathematical meaning of the text. Similarly, the following pieces of  $\text{\LaTeX}$  can also be extracted:

```

    “There is #1 such that #2.”
    <math>\langle\langle\text{\#1 in \#2}\rangle\langle\langle\text{\text{“an element \#0\#”}, \text{\text{“\#R\#”}}\rangle\rangle\rangle\rangle,
    \text{\text{“\#1 = \#2\#”}}\langle\langle\text{\text{“\#1 + \#2”}}\langle\langle\text{\text{“a”}, \text{\text{“0”}}, \text{\text{“a”}}\rangle\rangle\rangle\rangle
    
```

This tree of  $\text{\LaTeX}$  typesetting instructions can be further flattened for actual processing by  $\text{\LaTeX}$  into a string such as:

“There is an element  $\$0\$$  in  $\$R\$$  such that  $\$a + 0 = a\$$ .”

The idea of the TSA representation is independent of the visual formatting language used. Although we use  $\text{\LaTeX}$  in our example here, in our implementations so far we have used the  $\text{\TeX}_{\text{MACS}}$  internal representation and also XML.

As part of the task of using TSA to interleave CGa and more traditional natural language and typesetting information, we needed to develop techniques for handling

<sup>4</sup>This example comes from [16].

<sup>5</sup>The representation shown here omits type/category annotations that we usually include with the CGa identifiers used in the TSA representation.

certain challenging CML formations where the mathematical structure and the CML representation do not nicely match. For example, in the text  $0 + a0 = a0 = a(0 + 0) = a0 + a0$ , the terms  $a0$  and  $a(0 + 0)$  are each shared between two equations. Most formal representations would require either duplicating these shared terms, like for example  $0 + a0 = a0 \wedge a0 = a(0 + 0) \wedge a(0 + 0) = a0 + a0$ , or explicitly abstracting the shared terms. To allow the TSa representation to be as close to CML as possible, we instead solve this issue by using “*souring*” annotations in the TSa representation [16]. These annotations are a third kind of node label used in TSa, in addition to the CGa and formatting labels. We have developed methods using souring annotations for extracting the correct mathematical meaning and the nice visual presentation in the CML style. For the above example this is given in figure 4.

We have developed more sophisticated annotations that can handle more complicated cases of sharing of terms between equations. Souring annotations have also been developed to support several other common CML formulations. Support for folding and mapping over lists allows using forms like  $\forall a, b, c \in S.P$  as shorthand for  $\forall a \in S. \forall b \in S. \forall c \in S. P$  and  $\{a, b, c\}$  as shorthand for  $\{a\} \cup (\{b\} \cup (\{c\} \cup \emptyset))$ . We have not yet developed folding that is sophisticated enough to handle ellipsis (...) as in CML formulations like the following example (from [34]):

$$f[\underbrace{x, \dots, x}_{n+1 \text{ arguments}}] = \frac{f^{(n)}(x)}{n!}$$

We have implemented a user interface as an extension of the  $\text{T}_{\text{E}}\text{X}_{\text{MACS}}$  editor for entering the TSa MathLang representation. The author can use mouse and keyboard commands to annotate CML text entered in  $\text{T}_{\text{E}}\text{X}_{\text{MACS}}$  with boxes representing the CGa grammatical categories in order to assign CGa identifiers and thereby explicitly indicate mathematical meanings. The user interface allows displaying either a pure CML view which hides the TSa and CGa information, a pure CGa view, or various combined views including a view like the one depicted in figure 3. The same interface allows adding souring annotations like those depicted in figure 4.

In future work, we would like to develop techniques for not just pairing a single CML presentation with its CGa meaning, but also allowing multiple parallel visual presentations such as multiple natural languages (not just English), both natural language and symbolic formula presentations, and presentations in different symbolic notations. We would like to develop better software support to aid in semi-automatically converting existing CML texts into MathLang via TSa and CGa.

The design of TSa is due to Kamareddine, Maarek, and Wells with contributions by Lamar to the souring rules [16, 29]. The implementation is primarily by Maarek [29].

### 3.3 *The Document Rhetorical aspect (DRa)*

DRa [19, 18] is a system for attaching annotations to mathematical documents that indicate the roles played by different parts of a document. DRa assumes the underlying mathematical representation (which can be the MathLang aspects CGa or TSa) has some mechanism for identifying document parts.

Some DRa annotations can be unary predicates on parts; these include annotations indicating ordinary document sectioning roles such as *part*, *chapter*, *section*, etc. (like

the sectioning supported by L<sup>A</sup>T<sub>E</sub>X, OMDoc, DocBook, etc.) and others indicating special mathematical roles such as *theorem*, *lemma*, *proof*, etc.

Other DRa annotations can be binary predicates on parts; these include such relationships between parts as “*justifies*”, “*uses*”, “*inconsistent with*”, and “*example of*”. Regarding the annotation of justifications, remember that a CML text is usually incomplete: a mathematical thought process makes jumps from one interesting point to the next, skipping over details. This does not mean that many mistakes can occur; these details are usually so obvious for the mathematician that a couple of words are enough (e.g., “*apply theorem 35*”). The mathematician knows that too many details hinder concentration. To allow MathLang text to be close to CML text, DRa allows informal justifications, which can be seen as *hints* about which statements would be used in the proof of another statement.

Figure 5 gives an example (taken from [18] and implemented by Retel) where the mathematician has identified parts of the text (indicated by letters A through I in the figure). Figure 6, shows the underlying mathematical representation of some example DRa annotations for the example in figure 5. Here, the mathematician has given each identified part a structural (e.g., chapter, section, etc.) and/or mathematical (e.g., lemma, corollary, proof, etc.) rhetorical role, and has indicated the relation between wrapped chunks of texts (e.g., justifies, uses, etc.). Note that all the DRa annotations are represented as triples; this allows using the machinery of RDF [6] (a W3C standard that is aimed at the “semantic web”) to represent and manipulate them.

The DRa structure of a text can be represented as a tree (which is exactly the tree of the XML representation of the DRa annotated MathLang document). Due to the tree structure of a DRa annotated document, we refer to an annotated part of a text as a DRa node. We see an example of such a DRa node in figure 9. The role of this node is **declaration** and its name is **decA**. Note that the content of a DRa node is the user’s CGa and TSa annotation.

In the DRa annotation of a document, there is a dedicated root node (the Document node) where each top-level DRa node is a child of this root node. For example in figure 7, we see a tree consisting of 10 nodes. The root node (labelled Document) has four children nodes and five grandchildren nodes (which are all children of B).

We distinguish between proved nodes (*theorem*, *lemma*, etc.) which have a solid line in the picture and unproved nodes (*axiom*, *definition*, etc.) which have a broken line. We introduce this distinction in this paper because with the current implementation of DRa the user has the possibility to create its own mathematical and structural roles. Because we want to check a DRa annotated document for validity, the information whether a node is to be proved or not is important. For example such information would result in an error if someone tries to prove an unproved node e.g. by proving a definition or an axiom. When document  $D_2$  references document  $D_1$  it can reference the root node of document  $D_1$  to include all of its mathematical text.

In figure 5 one can see, that there are four top-level nodes in the example: A, B, C and D, representing respectively lemma 1, a proof of lemma 1, corollary 2 and a proof of corollary 2. The proof of lemma 1 has five children: E, F, G, H, I representing respectively the definition of the predicate, a claim, the proof of the claim, case 1 and case 2. The visual representation of this tree can be seen in figure 7.

By traversing the tree in pre-order we derive the original linear order of the DRa nodes of the text. Pre-order means that the traversal starts with the root node and



(A, hasMathematicalRhetoricalRole, lemma)	(B, justifies, A)
(E, hasMathematicalRhetoricalRole, definition)	(D, justifies, C)
(F, hasMathematicalRhetoricalRole, claim)	(D, uses, A)
(G, hasMathematicalRhetoricalRole, proof)	(G, uses, E)
(B, hasMathematicalRhetoricalRole, proof)	(F, uses, E)
(H, hasMathematicalRhetoricalRole, case)	(H, uses, E)
(I, hasMathematicalRhetoricalRole, case)	(H, caseOf, B)
(C, hasMathematicalRhetoricalRole, corollary)	(H, caseOf, I)
(D, hasMathematicalRhetoricalRole, proof)	

FIG. 6. Example of DRa relationships between chunks of text in figure 5

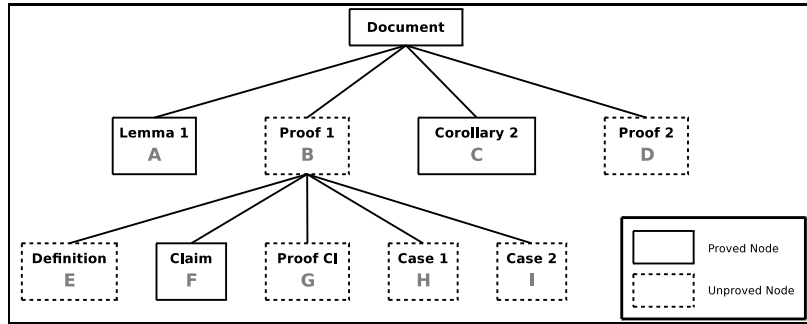


FIG. 7. Example of a tree of the DRa nodes of a document

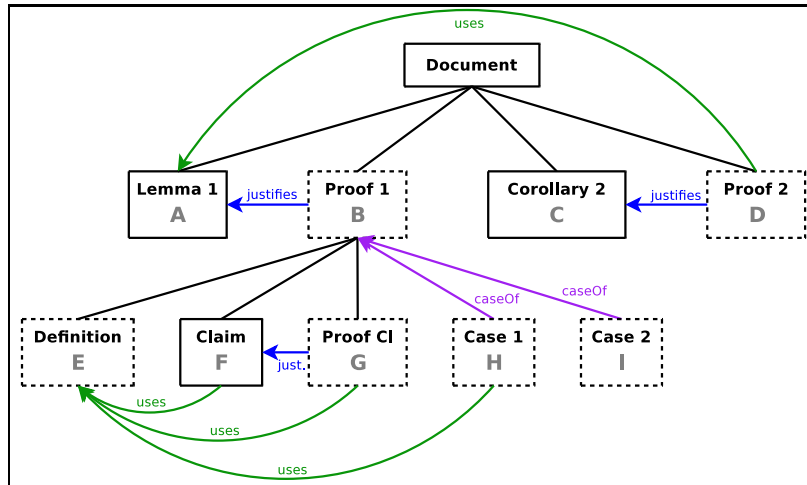


FIG. 8. Dependency graph for an example DRa tree

### 3.3.1 Textual Order

To be able to examine the proper structure of a DRa tree we introduce the concept of textual order between two nodes in the tree. The concept of textual order is a modification of the logical precedence presented in [19]. In this article we formalise this concept of order for the first time and show how we can process this information to automatically generate a proof skeleton. The textual order expresses the dependencies

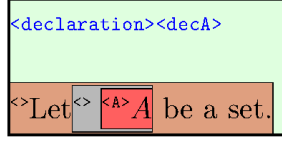


FIG. 9. An example for a single DRa node

between parts of the text. For example if a node  $A$  uses a part of a node  $B$ , then in a sequence of reasoning steps,  $B$  has to be before  $A$ . In order to have a formal definition of textual order, we introduce some notions for DRa nodes.

Recall that the content of a DRa node is its CGa and TSa part and that we have nine different kinds of CGa annotations: **term** **set** **noun** **adjective** **statement** **declaration** **definition** **step** **context**. A DRa node can also have further DRa nodes as children (e.g.  $B$  in figure 7 has the children  $E, F, G, H$  and  $I$ ). We introduce different sets for a DRa node  $n$ . All these sets can be automatically generated from the user's CGa and TSa annotations of the document. Table 1 defines these sets and gives examples for the CGa annotated text in figure 10 which is the definition of the subset relation.



FIG. 10. CGa annotations for the definition of the subset relation

Set	Description	Example for figure 10
$T(n)$	$\{x \mid x \text{ is part of } n \text{ and } x \text{ is annotated as } \text{term}\}$	$\{x\}$
$S(n)$	$\{x \mid x \text{ is part of } n \text{ and } x \text{ is annotated as } \text{set}\}$	$\{A, B\}$
$N(n)$	$\{x \mid x \text{ is part of } n \text{ and } x \text{ is annotated as } \text{noun}\}$	$\{\}$
$A(n)$	$\{x \mid x \text{ is part of } n \text{ and } x \text{ is annotated as } \text{adjective}\}$	$\{\}$
$ST(n)$	$\{x \mid x \text{ is part of } n \text{ and } x \text{ is annotated as } \text{statement}\}$	$\{A \subset B, x \in A, x \in B, x \in A \implies x \in B, \forall x(x \in A \implies x \in B)\}$
$DC(n)$	$\{x \mid \exists q \text{ part of } n, q \text{ is annotated as } \text{declaration} \text{ and } x \text{ is the declared symbol of } q\}$	$\{x\}$
$DF(n)$	$\{x \mid \exists q \text{ part of } n, q \text{ is annotated as } \text{definition} \text{ and } x \text{ is the defined symbol of } q\}$	$\{\subset\}$
$SP(n)$	$\{x \mid x \text{ is part of } n \text{ and } x \text{ is annotated as } \text{step}\}$	$\{A \subset B \iff \forall x(x \in A \implies x \in B)\}$
$C(n)$	the set of all parts of $n$ annotated as <b>context</b>	$\{\}$
$\mathcal{ENV}(n)$	$\{x \mid \exists m \neq n, m \text{ is a node in the pre-order path from the root node to the node } n, x \text{ is a part of } m, \text{ and } x \text{ is annotated as } \text{statement}\}$	

TABLE 1. Sets for a DRa node  $n$  and examples

Let us give further examples of  $DC(n)$  and  $DF(n)$ . In section 3 we had the following example of a list of declarations (call it  $ex$ ):

```
...; ∈(term, set) : stat; =(term, term) : stat; natural_number : noun;
+ (natural_number, natural_number) : natural_number; ...
```

For  $ex$ , we have that  $\mathcal{DC}(ex) = \{\in, =, \text{natural\_number}, +\}$ .

Now take the example of figure 11 (and call it  $ex'$ ). This example introduces the definition of  $\neg$  (Definition 1). We have that  $\mathcal{DF}(ex') = \{\neg\}$ .

The syntax of a definition in the internal representation of CGa (which is not necessarily the same as the syntax introduced by the reader), is an identifier with a (possibly empty list of arguments) on the left-hand side followed by the symbol  $:=$  and an expression on the right-hand side. The introduced symbol is again the identifier of the left-hand side. For the example of figure 10 (call it  $ex''$ ), the internal CGa representation is:

$$\subset(A, B) := \text{forall}(a, \text{impl}(\text{in}(a, A), \text{in}(a, B)));$$

the only introduced symbol is  $\subset$  and hence  $\mathcal{DF}(ex'') = \{\subset\}$ .

Note that  $\mathcal{ENV}(n)$  is the environment of all mathematical statements that occur before the statements of  $n$  (from the root node). Note furthermore that in the CGa syntax of MathLang, a definition or a declaration can only introduce a **term**, **set**, **noun**, **adjective**, or **statement**. Furthermore, recall that mathematical symbols or notions can only be introduced by **definition** or **declaration** and that mathematical facts can only be introduced by a **statement**. We define the set  $\mathcal{IN}(n)$  of *introduced symbols and facts* of a DRa node  $n$  as follows:

$$\mathcal{IN}(n) := \mathcal{DF}(n) \cup \mathcal{DC}(n) \cup \{s \mid s \in \mathcal{ST}(n) \wedge s \notin \mathcal{ENV}(n)\} \cup \bigcup_{c \text{ childOf } n} \mathcal{IN}(c)$$

At the heart of a **context**, **step**, **definition**, or **declaration**, one finds a set of **term**, **set**, **noun**, **adjective**, and **statement**. A DRa node  $n$  *uses* the set  $\mathcal{USE}(n)$  defined by:

$$\mathcal{USE}(n) := \mathcal{T}(n) \cup \mathcal{S}(n) \cup \mathcal{N}(n) \cup \mathcal{A}(n) \cup \mathcal{ST}(n) \cup \bigcup_{c \text{ childOf } n} \mathcal{USE}(c)$$

LEMMA 3.3

$\mathcal{DF}(n) \cup \mathcal{DC}(n) \subseteq \mathcal{T}(n) \cup \mathcal{S}(n) \cup \mathcal{N}(n) \cup \mathcal{A}(n) \cup \mathcal{ST}(n)$  for every DRa node  $n$ .

LEMMA 3.4

$\mathcal{IN}(n) \subseteq \mathcal{USE}(n)$  for every DRa node  $n$ .

PROOF. By induction on the depth of parenthood of  $n$ . If  $n$  has no children then use lemma 3.3. Assume the property holds for all children  $c$  of  $n$ . By lemma 3.3 and the induction hypothesis, we have  $\mathcal{IN}(n) \subseteq \mathcal{USE}(n)$ . ■

We demonstrate these notions with an example. Consider a part of a mathematical text and its corresponding DRa tree with relations as in figure 11.

We assume the document starts with an environment which contains two statements: **<True>True** and **<False>False**. Hence  $\mathcal{ENV}(\text{def1}) = \{True, False\}$ . When traversing the tree we start with the given environment for the node **def1**:

$$\mathcal{ENV}(\text{def1}) = \{True, False\}$$

The environment for **case1** consists of the environment of **def1** and all new statements of **def1**. In **def1** there is only the new statement  $\neg$  which is added to the environment:

$$\mathcal{ENV}(\text{case1}) = \{\neg\} \cup \mathcal{ENV}(\text{def1})$$



- **Weak textual order**  $\preceq$ : This order describes a subpart relation between two nodes ( $A$  is a subpart of  $B$ , written as  $A \preceq B$ ). More formally:

$$A \preceq B := \mathcal{IN}(A) \subseteq \mathcal{IN}(B) \wedge \mathcal{USE}(A) \subseteq \mathcal{USE}(B)$$

When  $A \preceq B$ , we also write  $B \succeq A$ .

- **Common textual order**  $\leftrightarrow$ : This order describes the relation that two nodes use at least one common symbol or statement. More formally:

$$A \leftrightarrow B := \exists x(x \in \mathcal{USE}(A) \wedge x \in \mathcal{USE}(B))$$

A DRa relation between two nodes induces a textual order. Table 3 gives some standard relations and their textual order.

Relation	Meaning	Order
$A$ uses $B$	$A$ uses a statement or a symbol of $B$	$B \prec A$
$A$ inconsistentWith $B$	some statement in $A$ contradicts a statement in $B$	$B \prec A$
$A$ justifies $B$	$A$ is the proof for $B$	$A \leftrightarrow B$
$A$ relatesTo $B$	There is a connection between $A$ and $B$ but no dependence	$A \leftrightarrow B$
$A$ caseOf $B$	$A$ is a case of $B$	$A \preceq B$

TABLE 3. Example of DRa relations and their textual order

We can now verify the relations of the example of figure 11 and their textual orders (Table 4). It is obvious that all five conditions hold and hence the relations

Relation	Condition	Order
( <a href="#">case1</a> , <a href="#">caseOf</a> , <a href="#">def1</a> )	$\mathcal{IN}(\text{case1}) \subseteq \mathcal{IN}(\text{def1}) \wedge \mathcal{USE}(\text{case1}) \subseteq \mathcal{USE}(\text{def1})$	<a href="#">case1</a> $\preceq$ <a href="#">def1</a>
( <a href="#">case2</a> , <a href="#">caseOf</a> , <a href="#">def1</a> )	$\mathcal{IN}(\text{case2}) \subseteq \mathcal{IN}(\text{def1}) \wedge \mathcal{USE}(\text{case2}) \subseteq \mathcal{USE}(\text{def1})$	<a href="#">case2</a> $\preceq$ <a href="#">def1</a>
( <a href="#">pr1</a> , <a href="#">justifies</a> , <a href="#">lem1</a> )	$\exists x(x \in \mathcal{USE}(\text{pr1}) \wedge x \in \mathcal{USE}(\text{lem1}))$	<a href="#">pr1</a> $\leftrightarrow$ <a href="#">lem1</a>
( <a href="#">lem1</a> , <a href="#">uses</a> , <a href="#">def1</a> )	$\exists x(x \in \mathcal{USE}(\text{lem1}) \wedge x \in \mathcal{IN}(\text{def1}))$	<a href="#">def1</a> $\prec$ <a href="#">lem1</a>
( <a href="#">pr1</a> , <a href="#">uses</a> , <a href="#">def1</a> )	$\exists x(x \in \mathcal{USE}(\text{pr1}) \wedge x \in \mathcal{IN}(\text{def1}))$	<a href="#">def1</a> $\prec$ <a href="#">pr1</a>

TABLE 4. Conditions for the relations of the example

are valid. For example the relation ([case2](#), [uses](#), [lem1](#)) would not be valid, because  $\neg \exists x(x \in \mathcal{USE}(\text{case1}) \wedge x \in \mathcal{IN}(\text{lem1}))$ .

Note that these conditions are only of a syntactical form. There is no semantical checking if e.g. a “justifies” relation really connects a proved node and its proof.

### 3.3.2 The GoTO

The GoTO is the Graph of textual order. For each kind of relation in the dependency graph (DG) of a DRa tree we can provide a corresponding textual order  $\prec$ ,  $\preceq$  or  $\leftrightarrow$ . These different kinds of order can be interpreted as edges in a directed graph. So we can transform the dependency graph into a GoTO by transforming each edge of the DG. So far there are two reasons why the GoTO is produced:

1. Automatic Checking of the GoTO can reveal errors in the document (e.g. loops in the structure of the document).

2. The GoTO is used to automatically produce a proof skeleton for a certain prover.

To transform an edge of the DG we need to know which textual order it induces. Each relation has a specific order  $\prec, \succ, \preceq, \succeq, \leftrightarrow$ . Table 5 shows the graphical representation of such edges and an example relation we have already seen in our examples. There




$(A, \text{uses}, B)$	$A \succ B$	
$(A, \text{caseOf}, B)$	$A \preceq B$	
$(A, \text{justifies}, B)$	$A \leftrightarrow B$	

TABLE 5. Graphical representation of edges in the GoTO

is also a relation between a DRa node and its children: For each child  $c$  of a node  $n$  we have the edge  $c \preceq n$  in the GoTO. This “childOf” relation is added automatically when producing the GoTO. But it can be added manually by the user. This can be useful e.g. in papers with a page restriction, where some parts of the text are relocated in the appendix but would be originally within the main text.

The algorithm for producing the GoTO out of the DG works in two steps:

1. transform each relation of the DG into its corresponding edge in the GoTO
2. for each child  $c$  of a node  $n$  add the edge  $c \preceq n$  to the GoTO

When performing this algorithm on the example of figure 8 we get the GoTO as demonstrated in figure 12. Each relation of the DG which induces a  $\leftrightarrow$  textual order is replaced by the corresponding edge in the GoTO. We can see these edges between a proved node and its proof where the “justifies” relation induces a  $\leftrightarrow$  order (e.g. between A and B, C and D, and F and G). The children of the node B are connected to B via  $\preceq$  edges in the GoTO. For the “caseOf” relation, the user has manually specified the relation, the other edges were added automatically by the algorithm generating the GoTO. The relations which induce the order  $\prec$  are transformed into the corresponding directed edges in the GoTO. We see that the direction of the nodes has changed with respect to the DG. This is because we only have “uses” relations, and for a relation  $(A, \text{uses}, B)$  we have the textual order  $B \prec A$  which means, that the direction of the edge changes.

### 3.3.3 Automatic checking of DG and GoTO

We implemented two kinds of failures: warnings and errors. At the current development of DRa we check for four different kinds of failures:

1. Loops in the GoTO (error)
2. Proof of an unproved node (error)
3. More than one proof for a proved node (warning)
4. Missing proof for a proved node (warning)

The checks for ii) - iv) are performed in the DG. For ii) we check for every node of type “unproved” if there is an incoming edge of type “justifies”. If so, an error

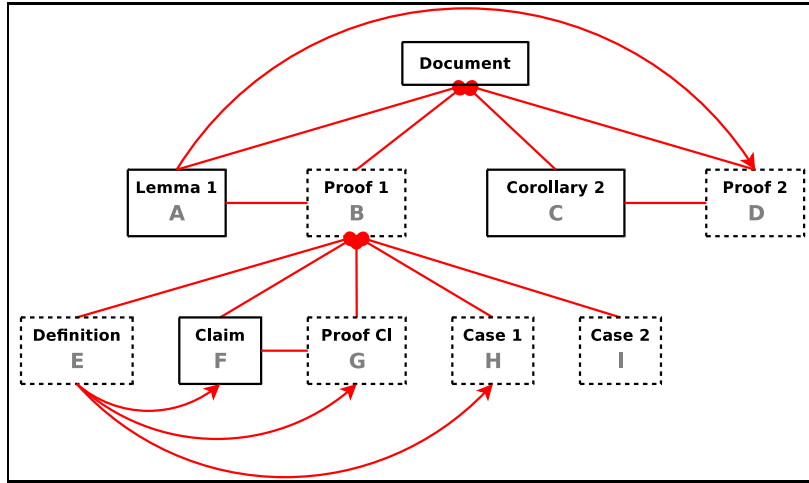


FIG. 12. Graph of Textual Order for an example DRa tree

is returned (e.g. when someone tries to prove an axiom or a definition). For iii) and iv) we check for each node of type “proved” if there is an incoming edge of type “justifies”. If not, we return a warning (this can be a deliberate omission of the proof or just a mistake). If there is more than one proof for one node we return also a warning (most formal systems cannot handle multiple proofs).

For i) we search for cycles in the GoTO. Therefore we have to define how we treat the three different kinds of edges. Edges of type  $\prec$  and  $\preceq$  are treated as directed edges. Edges of type  $\leftrightarrow$  are in principal undirected edges, which means for an edge  $A \leftrightarrow B$ , one can get from  $A$  to  $B$  and from  $B$  to  $A$  in the GoTO. But it is important, that within one cycle such an edge is only used in one direction. Otherwise we would have a trivial cycle between two nodes connected by a  $\leftrightarrow$  edge.

As we will see in the next section, a single node in the DRa tree can first be translated when all its children nodes are ready to be translated. To reflect this circumstance we have to add certain nodes in the GoTO for the cycle check. Let us demonstrate this with an example. Consider a DG and GoTO as in figure 13

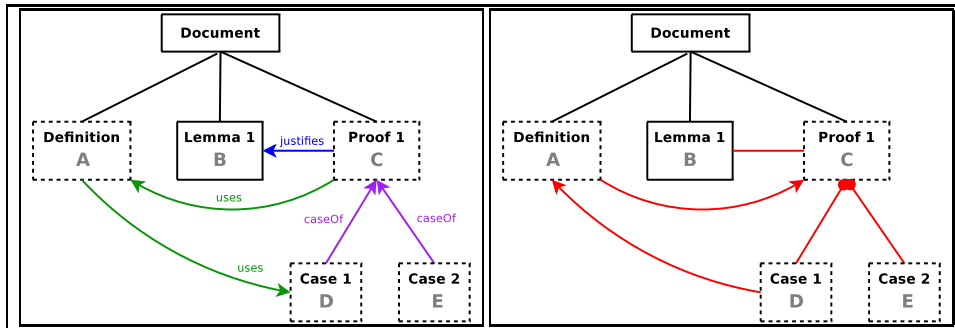


FIG. 13. Example of a not recognised loop in a DRa (left DG, right GoTO)

Apparently there is a cycle in this tree, because to be able to translate C we need to translate its children D and E. But before we can translate C we must have translated A because C uses A. But the child D of C is used by A. There we have a deadlock situation. Neither A nor C can be processed. To recognise such kinds of cycles we add certain edges to the GoTO when checking for cycles. Therefore we have to look at the children of a node  $n$ : hidden cycles can only evolve, when there are edges  $e_i$  from a child node  $c_i$  to a target node  $t_i$  which is not a sibling of  $c_i$ . Hence we add an edge  $c_i \succ n$  for each such node  $e_i$  to the GoTO. This can be done via algorithm 1.

```

foreach node  $n$  of the tree do
  foreach child  $c$  of  $n$  do
    foreach outgoing edge  $e$  of  $c$  do
      if target node  $t$  of  $e$  is no sibling of  $c$  then
        add a Strong textual precedence edge from  $n$  to  $t$ ;
      end
    end
  end
end

```

**Algorithm 1:** Adding additional edges to the GoTO

We could also add new edges for all incoming edges of the children  $c_i$  but this is not necessary because the textual order of the “childOf” relation is a directed edge from each child  $c_i$  to its parent node  $n$ . There we can use the transitivity of the edges and find a cycle anyway.

In the example from figure 13, algorithm 1 would add one edge to the GoTO: The child node D of C has an outgoing node to the non-sibling node A. So a new directed edge from C to A is added which yields the result of figure 14. There you can clearly see a cycle between the nodes A, C and A with the edges A-C and C-A.

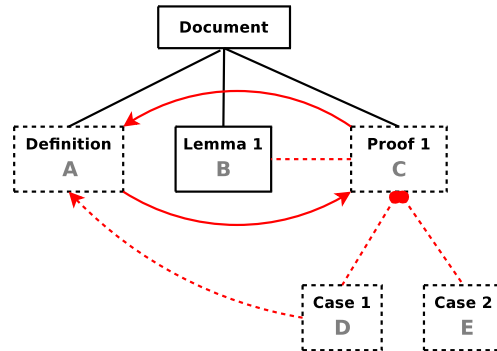


FIG. 14. GoTO graph of the example of figure 13 with added edges

Figure 15 demonstrates another situation of a cycle in a DRa annotated text. The problem is mainly, that lemma 1 uses lemma 2 but the proof of lemma 2 uses a part of the proof of lemma 1. This situation would end up in a deadlock when processing the GoTO e.g. when producing the proof skeleton. We see a cycle between the nodes A, C, D, F, B and A with the edges A-C, C-D, D-F, F-B, and B-A. Here we also

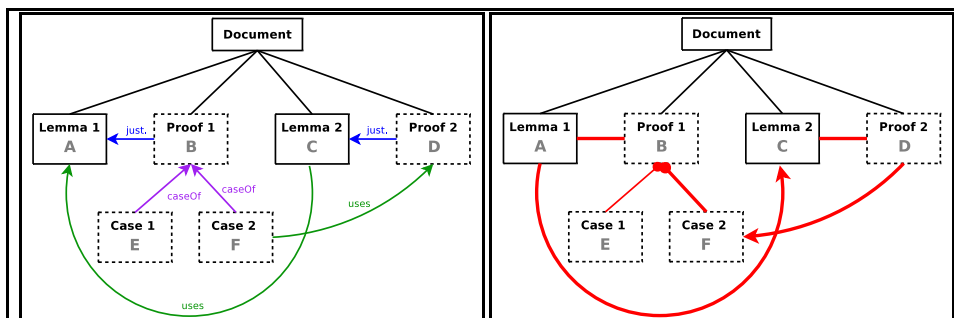


FIG. 15. Example of a loop in the GoTO (DG left, GoTO right)

see why we do not need to add incoming edges to the parent nodes. For node F we have an incoming edge but due to the direction of the “childOf” edge from F to B, we can use the transitivity. In both examples, an error would be returned with the corresponding nodes and edges.

## 4 Connecting MathLang to formal foundations

### 4.1 Goals for formalisation

Current approaches to formalising CML texts generally involve rewriting the text from scratch; there is no clear methodology in which the text can gradually change in small steps into its formal version. One of MathLang’s goals is to support formalising a text in small steps that do not require radically reorganising the text. Also, a text with fully formal content should continue to be able to be presented in the same way as a less formal version originally developed by a mathematician. We envision formalisation as working by adding additional layers of information to a MathLang document to support embedding formal proofs. Ideally, there should be flexible control over how much of the additional information is presented to the reader; the additional information could form part of the visual presentation, or could exist “behind the scenes” to provide assurance of correctness.

As part of the goal of supporting formalisation in MathLang, we desire to keep MathLang independent of any particular formal foundation. However, as proofs embedded in a MathLang document become more formal, it will be necessary to tie them more closely to a particular proof system. It might be possible that fully formal documents could be kept independent of any particular foundation by allowing the most formal parts of a document to be expressed redundantly in multiple proof systems. (This is similar in spirit to the way the natural language portion of a document might be expressed simultaneously in multiple natural languages.)

Following the general MathLang development/design strategy in figure 1, we have been developing methodology and software for connecting a MathLang document with formal versions of its content.

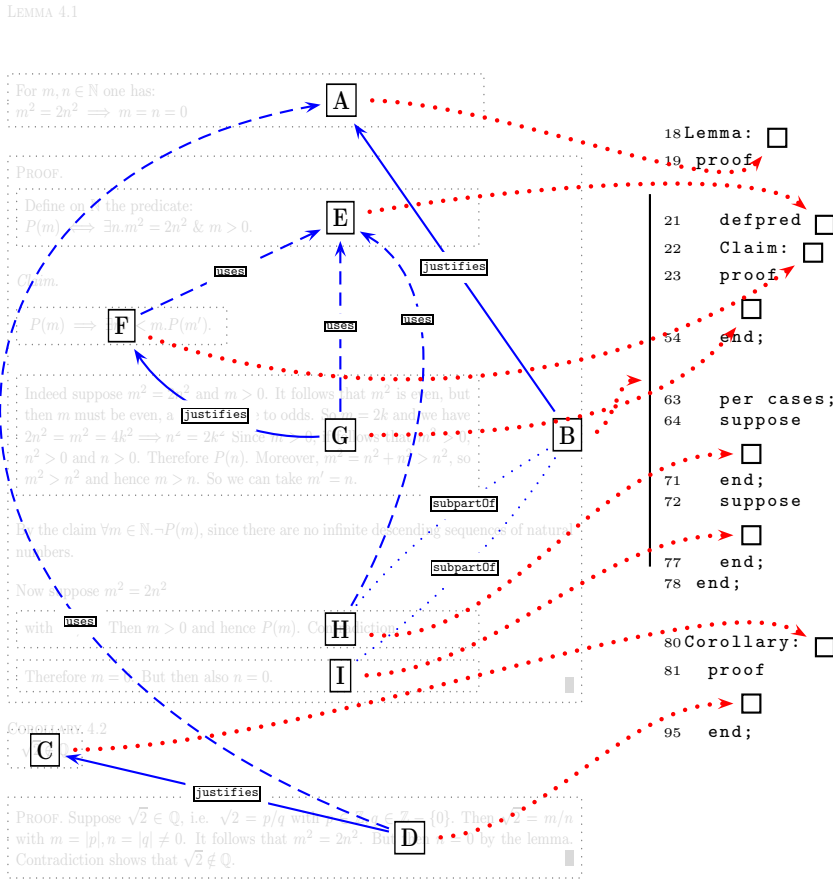


FIG. 16. Generating a Mizar Text-Propser skeleton from MathLang DRa and CGa

## 4.2 Current developments

When the MathLang project started, we planned to start with Coq for our initial development of support for formalisation, but because Krzysztof Retel joined us with his previous Mizar experience, we started with Mizar instead. Our work so far with Retel on extending MathLang documents into Mizar formalisations [18] involves constructing a skeleton of a Mizar document (e.g. figure 16) from a MathLang document, and then completing the Mizar skeleton separately. A Mizar document consists of an Environment-Declaration and a Text-Propser. In Mizar, the Environment-Declaration is used to generate the Environment which has the needed knowledge from MML (Mizar’s Mathematical Library). The Text-Propser is checked for correctness using the knowledge in the Environment.

The new PhD student Christoph Zengler currently works on an automation of the skeleton generation for arbitrary theorem provers. Therefore he designed a generic algorithm for transforming the DRa tree into a proof skeleton. Since at this stage of formalisation we do not want to tie to any particular foundation, the algorithm is

highly configurable which means it takes the desired theorem prover as an argument and generates the proof skeleton within this theorem prover.

The aim of this skeleton generation is once again to stay as close as possible to the mathematician’s original CML text. But due to certain restrictions for different theorem provers the original order cannot always be respected. We give some classical examples when this can happen:

- Nested lemmas/theorems: Sometimes mathematicians define new lemmas or theorems inside proofs. Not every theorem prover can handle such an approach (e.g. Coq). In the case of such theorem provers, it is necessary to “de-nest” the theorems/lemmas.
- Forward references: Sometimes a paper first gives an example for a theorem before it states the theorem. Some theorem provers (e.g. Mizar) do not support such forward references. The text has to be rewritten so that it only has backward references (i.e. references to already stated mathematical constructs).
- Outsourced proofs: A usual practise in mathematical writing is to outsource complex proofs which are not mandatory for the central results, in the appendix. When formalising such documents, these proofs need to be put in the right place.

The algorithm for re-arranging the parts of the text and generating the proof skeleton performs reordering only when it is necessary for the theorem prover at hand.

#### 4.2.1 The generic automated Skeleton Generation Algorithm (gSGA)

The algorithm for the generation of the proof skeleton has two parameters. 1) The input MathLang XML file with DRa annotations and 2) a configuration file (in XML format) for the theorem prover (cf. table 6).

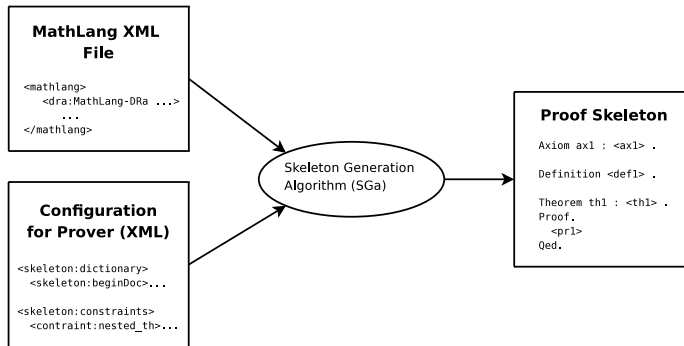


TABLE 6. The skeleton generation algorithm

Before we look at the configurations for a particular theorem prover, we will look at the algorithm in general. The algorithm works on the DRa tree as seen in the last section. A DRa node can have one of three states: processed (black), in-process (grey) and unprocessed (white). A processed node has already been translated into a part of the proof skeleton, a node in-process is one that is being checked, while an unprocessed node is still awaiting translation. We need this information to know which nodes we have already translated and which nodes are still to be translated.

This information is completely independent of the proved/unproved property of a node. The method for generating the output of a single node is shown in algorithm 2.

```

while foundwhite do
  foreach child c of the node do
    if c is unprocessed & isReady(c) then
      processNode(c);
      generateOutput(c);
      foundwhite := true;
      break;
    end
  end
end

```

**Algorithm 2:** generateOutput(Node node)

The algorithm searches a node that can be processed and processes it recursively. Process in this context means that the node at hand is translated and added to the proof skeleton. The generateOutput function has to be performed for the Document root node. Then it is recursively performed for all the nodes in the DRa tree. The important function in this algorithm is the one which decides whether a node is ready to be processed or not. This decision is only dependent on the GoTO of the DRa tree. A node is ready to be processed if it fulfils the following three criteria:

1. it has no incoming  $\prec$  edges (in the GoTO) of unprocessed (white) nodes
2. all its children are ready to be processed
3. if the node is a proved node: its proof is ready to be processed

Algorithm 3 tests these three properties of a node and returns the result.

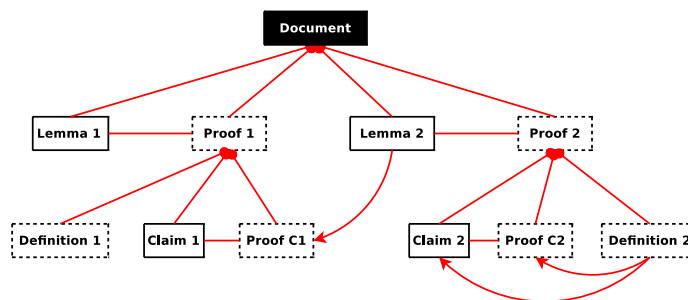
It is important when checking if each child of the  $n$  children is ready, to perform the test  $n$  times because a rearrangement can also be required for the children. If there are still white children after  $n$  steps, then the children cannot be yet processed and so the node cannot be processed.

To illustrate the algorithm we will now look at a detailed example. Consider the DRa tree of a (typical and not well structured) mathematical text with its DG and GoTO edges as in figure 17.

```

foreach incoming edge  $e$  of the node do
  if type of  $e$  is  $\prec$   $\mathcal{E}\mathcal{E}$  source of  $e$  is unprocessed (white) then
    return false
  end
end
mark node  $n$  as grey;
 $n$  = number of children of the node;
for 1.. $n$  do
  foreach child  $c$  of the node do
    if  $c$  is not processed  $\mathcal{E}\mathcal{E}$  isReady( $c$ ) then
      mark  $c$  as grey;
      break;
    end
  end
end
if still a white node is among the children of the node then
  reset all grey nodes back to white;
  return false
end
if node is a proved node then
  proof = proof of the node;
  if not isReady(proof) then
    reset all grey nodes back to white;
    return false
  end
end
reset all grey nodes back to white;
return true
    
```

Algorithm 3: isReady(Node node)



The root node of the document can be marked as processed and the algorithm starts at this node. The first child is Lemma 1. Criterion 1) is fulfilled, since the node has no incoming  $\prec$  edges in the GoTO. Criterion 2) is fulfilled because the node has no children. For criterion 3) Proof 1 has to be ready to be processed before we can mark Lemma 1 as ready to be processed. Proof 1 has no incoming  $\prec$  edges. So criterion 1) is fulfilled. For criterion 2) the children of the proof have to be ready to be processed. Definition 1 is ready, but the proof of Claim 1, Proof C1 has an incoming node of an unprocessed node (Lemma 2). So Claim 1 is not ready and hence, neither are Proof 1 and Lemma 1.

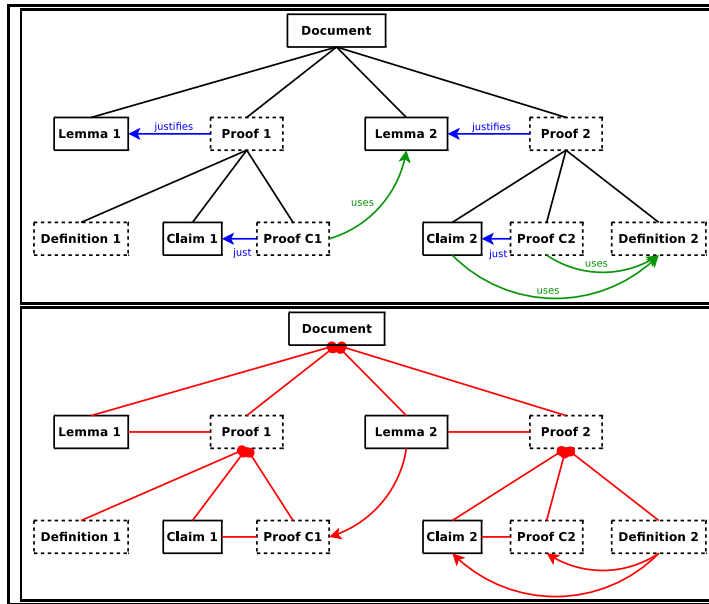
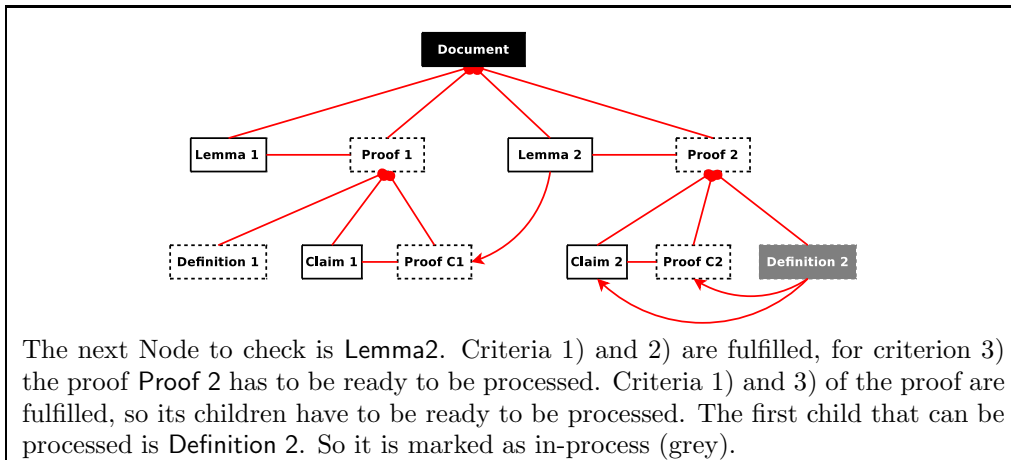
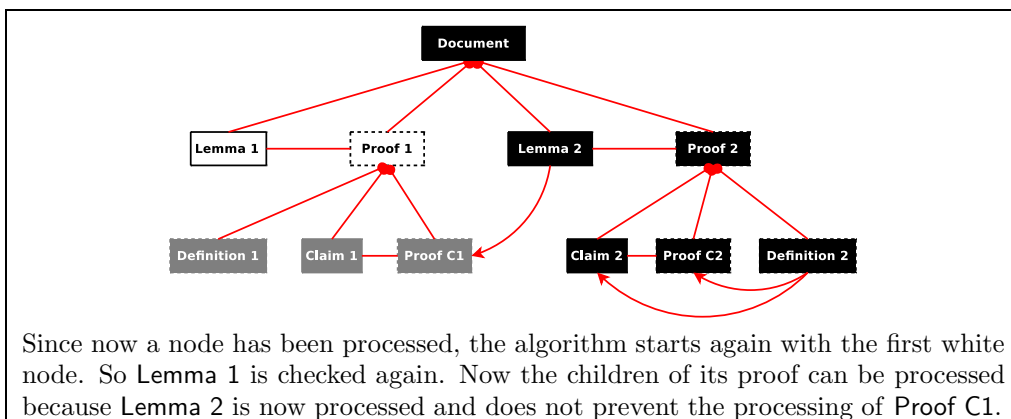
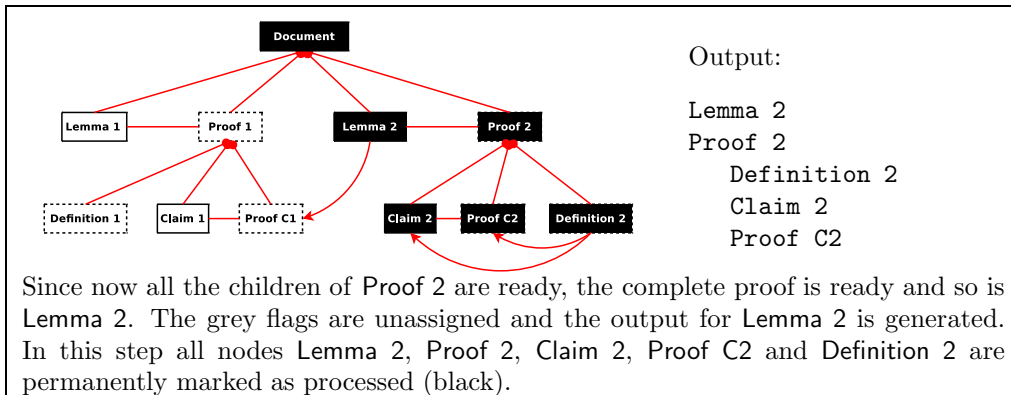
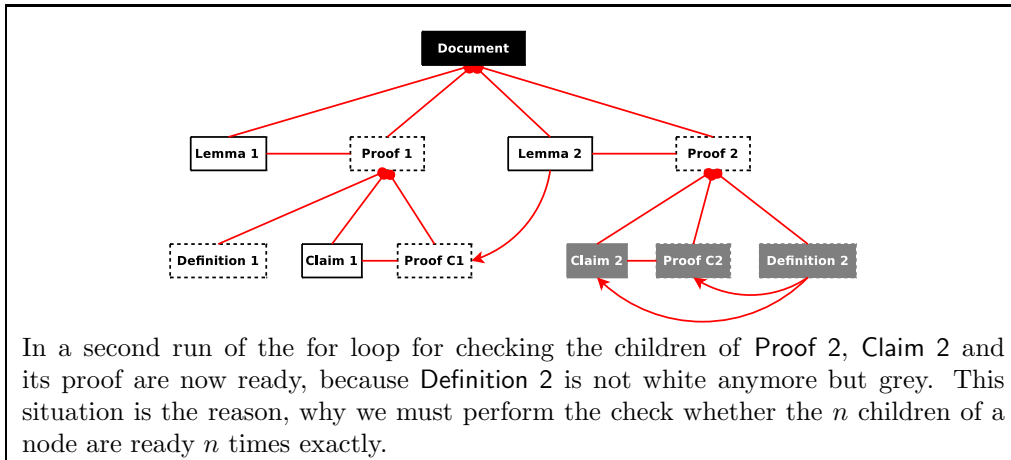
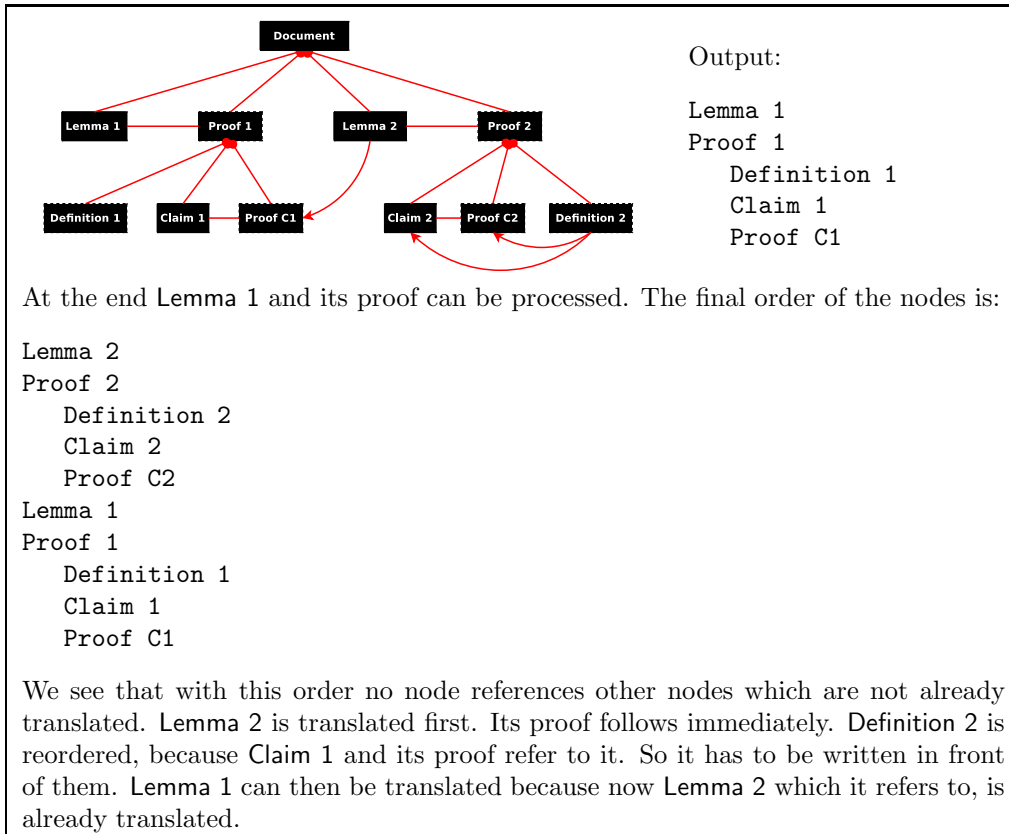


FIG. 17. Example to illustrate Skeleton generation (DG at top, GoTO at bottom)







#### 4.2.2 The configuration of gSGA

The transformation of the DRa annotated text into a proof skeleton has two steps.

- reorder the text to satisfy the constraints of the particular theorem prover
- translate each DRa annotation to its corresponding syntax in the language of the theorem prover

The configuration file for a particular theorem prover for the gSGA reflects these two steps: there is a dictionary part and a constraints part. The dictionary contains a rule for each mathematical or structural role of DRa. A single DRa node has two important properties: a name and a content. This information is used in the translation. Within the configuration file we can refer to the name of a node with `%name` and to the body with `%body`. A new line (for better readability) can be inserted with `%nl`. Consider the example of a DRa node which we saw in figure 9. The role of this node is declaration and its name is `decA`. The body of this node is the sentence *Let A be a set* or its CGa annotation. A translation into Mizar could be:

```
reserve <body of decA> ;
```

The rule for this translation would be:

```
reserve %body ;
```

Such a kind of declaration in Coq would be:

```
Variable <body of decA> .
```

And the for this translation would be:

```
Variable %body .
```

Here, we let a single rule be embedded in an XML tag whose attribute "name" is the corresponding keyword:

```
<skeleton:keyword name="declaration">
  reserve %body ;
</skeleton:keyword>
```

The constraints section of the configuration file for a theorem prover configures two main properties: the allowance of forward properties and of nested mathematical constructs. Forward references can be allowed via the tag:

```
<skeleton:forwardrefs>true</skeleton:forwardrefs>
```

Changing the content of the tag to "false" forbids forward references. If there is no such tag, the default value is "false".

For a configuration of nested constructs there are two possibilities:

- Either allow in general the nesting of constructs defining those exceptions for which nesting is not allowed;
- Or forbid in general the nesting of constructs defining those exceptions for which nesting is allowed.

The next configuration allows nesting in general but not for definitions and axioms:

```
<skeleton:nesting>true</skeleton:nesting>

<skeleton:nest role="definition">false</skeleton:nest>
<skeleton:nest role="axiom">false</skeleton:nest>
```

### 4.2.3 The flattening of the DRa graph

The next question we have to deal with, is how to perform changes to the tree when certain nestings are not allowed. We call this a flattening of the graph, because certain nodes are removed from their original position and inserted as direct children of the DRa top-level node. Algorithm 4 achieves this effect.

We refer to every child of the DRa top-level node as a node at level 1. Every child of such a node is at level 2 and so on. If a mathematical role must not be nested, it can only appear at level 1. So we check for each node at a level greater than level 1, if its corresponding mathematical role can be nested. If not, then the node and all its required siblings are removed from this level and put in front of their parent node. Since there is no "childOf" relation between this no-longer-child and its parent node, the relation between child and parent changes form  $\preceq$  to  $\prec$ .

The required sibling nodes are determined again in the GoTO. When a node is moved in front of its parent node, then there is a  $\prec$  edge between this node and its

```

foreach (child c of the node) do
  flattenNode(c);
  if c cannot be nested then
    nodelist := transitive closure of incoming nodes of c;
    foreach node n of nodelist do
      remove n of list of children of node;
      add n in front of node as a sibling;
    end
  end
end

```

Algorithm 4: flattenNode(Node node)

former parent. So each sibling of the removed node from whom there is a incoming node must be moved with the node. This includes its children or - for a proved node - its proof. Since for these children we have to move the related nodes too, we can build the transitive closure over the incoming nodes of the node which has to be moved. All nodes in this closure have to be relocated in front of the parent node.

We demonstrate this algorithm again on the example from figure 17. For a first demonstration we assume that the nesting of definitions is not allowed. So Definition 1 and Definition 2 have to be removed from level 2 and be relocated in front of their parent nodes. The transitive closure over incoming edges in the GoTO yields no new nodes for removing (because the definitions have no incoming edges in the GoTO). The resulting new flattened graph can be seen in figure 18. We see that the two definition are now at level 1 and their edges to their former parent nodes have changed from  $\preceq$  to  $\prec$ . The output for this graph according to the algorithm from the last section is given on the lefthandside of table 7.

Definition 1	Definition 2
Definition 2	Claim 2
Lemma 2	Proof C2
Proof 2	Lemma 2
Claim 2	Proof 2
Proof C2	Claim 1
Lemma 1	Proof C1
Proof 1	Lemma 1
Claim 1	Proof 1
Proof C1	Definition 1

TABLE 7. Outputs of the graphs of figures 18 (lefthandside) and 19 (righthandside)

On the other hand, if we allow definitions to be nested but forbid nested claims, we get the graph of figure 19. The first claim which is found in the graph is Claim 1. The transitive closure yields that Proof C1 needs also to be removed since there is a  $\leftrightarrow$  edge to the claim. The second claim which is found is Claim 2. The transitive closure yields again that its proof as well as Definition 2 have to be removed.

The output for this graph is given on the righthandside of figure 7.

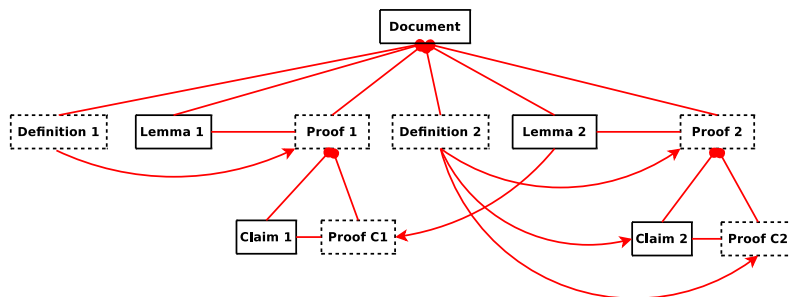


FIG. 18. A flattened graph of the GoTO of figure 17 without nested definitions

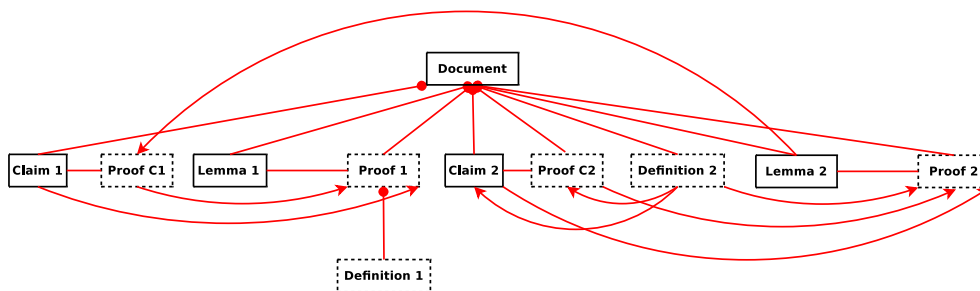


FIG. 19. A flattened graph of the GoTO of figure 17 without nested claims

## 5 A full formalisation in Coq via MathLang: chapter 1 of Landau’s “Grundlagen der Analysis”

Landau’s “Grundlagen der Analysis” [28] remains the only book which has been fully formalised in a theorem prover [36]. We intend to fully formalise and encode this book in MathLang. Throughout the years, a number of undergraduate/MSc students (most notably, Mona Petrie, Amalia Retzepe, Anastasios Tsaousis and Anastasios Asimakopoulos) contributed to encoding parts/chapters of this book at the CGa level of MathLang (even when CGa was still an under-developed variation of WTT).

Throughout we continued to use Landau’s text as a test bed for encoding at the CGa or TSa level of MathLang [29, 20]. In this section we will take for the first time, the first chapter of Landau’s book into all encoding levels in MathLang up to a full formalisation in Coq.

We will deal with the first chapter of the book - natural numbers - to show how an encoding of an existing document works. We have given a complete CGa, TSa and DRa annotation for the chapter. Furthermore, we have generated a proof skeleton automatically with the gSGA for Mizar and Coq and a complete formalised version of the chapter in Coq was created. We used the MathLang  $\text{T}_{\text{EX}}^{\text{MACS}}$  plugin to annotate the existing plaintext of the book.

To clarify the path we took, we look once again at the overall diagram of the different paths in MathLang (figure 2). We first used path ③ and annotated the complete text with CGa, TSa and DRa annotations with the help of the MathLang  $\text{T}_{\text{EX}}^{\text{MACS}}$  plugin. The second step was to automatically generate a proof skeleton of the annotated text. With the help of the proof skeleton and the CGa annotations we

Group	Meaning	Encoding	
Quantifiers	$\forall$	for all	<code>&lt;forall&gt;</code> $\forall$ <code>&lt;#&gt;</code> <code>&lt;#&gt;</code> . <code>&lt;#&gt;</code>
	$\exists$	exists	<code>&lt;exists&gt;</code> $\exists$ <code>&lt;#&gt;</code> <code>&lt;#&gt;</code> . <code>&lt;#&gt;</code>
	$\exists!$	exists exactly one	<code>&lt;exists_one&gt;</code> $\exists!$ <code>&lt;#&gt;</code> <code>&lt;#&gt;</code> . <code>&lt;#&gt;</code>
Boolean connectives	$\wedge$	and	<code>&lt;and&gt;</code> <code>&lt;#&gt;</code> $\wedge$ <code>&lt;#&gt;</code>
	$\vee$	or	<code>&lt;or&gt;</code> <code>&lt;#&gt;</code> $\vee$ <code>&lt;#&gt;</code>
	$\implies$	implication	<code>&lt;impl&gt;</code> <code>&lt;#&gt;</code> $\implies$ <code>&lt;#&gt;</code>
	$\oplus$	exclusive or	<code>&lt;xor&gt;</code> <code>&lt;#&gt;</code> $\oplus$ <code>&lt;#&gt;</code>
Set theory	$\in$	element of	<code>&lt;in&gt;</code> <code>&lt;#&gt;</code> $\in$ <code>&lt;#&gt;</code>
	$\subset$	subset of	<code>&lt;subset&gt;</code> <code>&lt;#&gt;</code> $\subset$ <code>&lt;#&gt;</code>
	$\{\}$	constructor for a set	<code>&lt;&lt;Set&gt;&gt;</code> { <code>&lt;#&gt;</code> <code>&lt;#&gt;</code> }
	$\emptyset$	empty set	<code>&lt;emptyset&gt;</code> $\emptyset$
	$=$	equality of sets	<code>&lt;seteq&gt;</code> <code>&lt;#&gt;</code> $=$ <code>&lt;#&gt;</code>
	$\neq$	inequality of sets	<code>&lt;setneq&gt;</code> <code>&lt;#&gt;</code> $\neq$ <code>&lt;#&gt;</code>
Special functions	$:=$	is a	<code>&lt;isa&gt;</code> <code>&lt;#&gt;</code> <code>isa</code> <code>&lt;#&gt;</code>
	$1$	one	<code>&lt;1&gt;</code> $1$
	$S(x)$	the successor function	<code>&lt;succ&gt;</code> <code>&lt;#&gt;</code>
	$-$	function for indexing	<code>&lt;index&gt;</code> <code>&lt;#&gt;</code> $-$ <code>&lt;#&gt;</code>

TABLE 8. The preface for the first chapter of Landau’s book

fully formalised the proofs in Coq completing the paths  $\textcircled{1}$  and  $\textcircled{2}$ . The final result was a complete formalised version of the first chapter of Landau’s book in Coq.

## 5.1 CGa and TSa annotations

### 5.1.1 The Preface

In the preface of a MathLang document we introduce symbols that are not defined in the text but are used throughout it. These are often quantifiers or Boolean connectives like  $\wedge$  or  $\vee$ . These symbols are often pre-encoded in theorem provers (e.g. Coq has special symbols for the logical and, or, implication, etc.). The preface of the first chapter of Landau’s book consists of 17 different symbols as given in table 5.1.1.

Two functions deserve further explanation:

1. The “is a” function is used to express that a particular term is an instance of a noun. E.g. the first axiom of the book is that *1 is a natural number*, so the encoding of this axiom is

```
<isa><1>1 is a <natural_number>natural number
```

2. The “index” function is used to express a notion in the style of  $a_b = c$  which can be defined as a function  $index(a, b) = c$ . So the index function has two terms as argument and yields a term as a result.

### 5.1.2 The first section

The first section of the first chapter introduces the natural numbers, equality on natural numbers and five axioms (an extension of the Peano axioms). We introduce a noun `<natural_numbers>natural numbers` and the set `<N>N` of natural numbers. Equality `<eq><#> = <#>` and inequality `<neq><#> = <#>` between natural numbers are declared rather than defined. There are three properties of equality, which we encoded. We will show one encoding of these to recapitulate TSa annotations with sharing and to see how to use the symbols given in the preface. The original statement is

$$x = x \text{ for every } x$$

We see that this is a universal quantification of  $x$  and a well formed equivalent statement would be:  $\forall x(x = x)$ . Since the positions are swapped in Landau's text we use the position sourcing. The sourcing annotation of this statement is:

$$\text{<forall><2> } x = x \text{ for every } \text{<1> } x$$

This yields the final statement

$$\text{<> <forall> <2> <eq> <x> } x = \text{<x> } x \text{ for every } \text{<1> <> <x> } x$$

Next we show how to encode axiom 2 which illustrates that “wordy” parts of the text can also be annotated, not only mathematical statements. The original statement is:

*For each  $x$  there exists exactly one natural number, called the successor of  $x$ , which will be denoted by  $x'$*

The “for each” can be translated with a universal quantifier, the “exactly one” with the  $\exists!$  quantifier. So we get the general structure:

`<forall>For each  $x$  <exists_one>there exists exactly one natural number, called the successor of  $x$ , which will be denoted by  $x'$`

The complete statement can be e.g. encoded corresponding to the following formal statement  $\forall x(\exists!x'(succ(x) = x'))$ :

`<forall>For each <x> $x$  <exists_one>there exists exactly <x'>one natural number, called the successor of <eq><x> $x$  <x'>, which will be denoted by  $x'$`

### 5.1.3 Sections 2 - 4

Within the next sections, addition (section 2), ordering (section 3) and multiplication (section 4) are introduced. There are 36 theorems with proofs and 6 definitions: addition, greater than, less than, greater or equal than, less or equal than and multiplication. There are many simple structured theorems like that of figure 20. We want to examine our way of annotating these theorems. The main theorem  $x + y = y + x$  is annotated in a straightforward manner.  $x$  and  $y$  are annotated as terms, plus as a

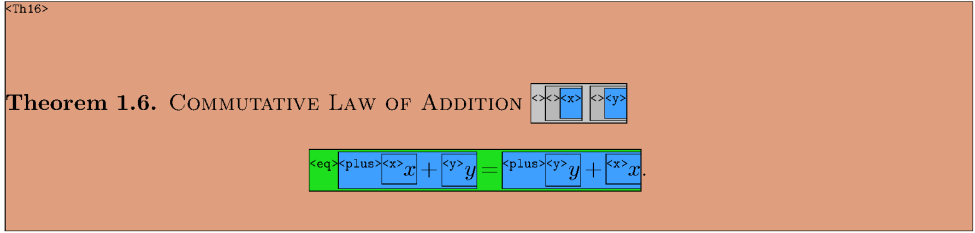


FIG. 20. Simple Theorem of the second section

function, taking two terms as arguments and yielding a term as a result. The equality between these terms is a statement. Since we did not declare  $x$  and  $y$  in the preface or in a global context we do this with a local scoping. This information is added in the first annotated line where we declare  $x$  and  $y$  as terms and put these two annotations into a context which means that this binding holds within the whole step.

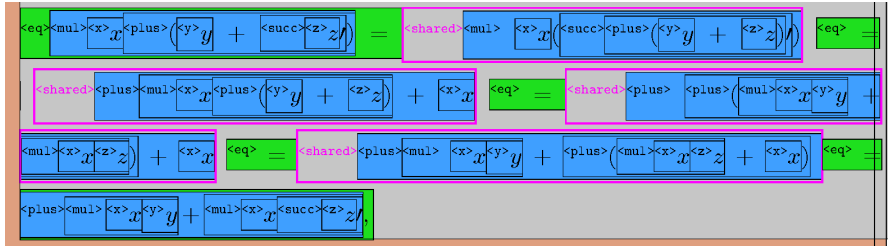


FIG. 21. Souring in chains of equations

Landau often used chains of equations for proofs as in this proof for the equality of  $x(y + z')$  and  $xy + xz'$  in the proof of Theorem 30 of the first chapter of the book:

$$x(y + z') = x((y + z)') = x(y + z) + x = (xy + xz) = x = xy = (xz + x) = xy + xz'$$

Here we benefit from our souring methods - in this case especially the sharing of variables (See figure 21). There are also often hidden quantification like the one in the example  $x = x$  for every  $x$ , where we need the souring for swapping positions. These two functionalities of T<sub>S</sub>a save a lot of time in annotating mathematical documents.

For some theorems we use the Boolean connectives although they are not mentioned explicitly in the text. E.g. Theorem 16 states:

$$\text{If } x \leq y, y < z \text{ or } x < y, y \leq z \\ \text{then } x < z$$

We annotate the premise of the theorem as a disjunction of two conjunctions as seen in figure 22. Another use of Boolean connectives is when we have formulations like “*exactly one of the following must be the case...*”. There we use the exclusive or  $\oplus$  to annotate the fact that exactly one of the cases must hold. We defined the exclusive or in the preface and therefore have to take care that we find a corresponding construct in the used theorem prover (see table 5.1.1).



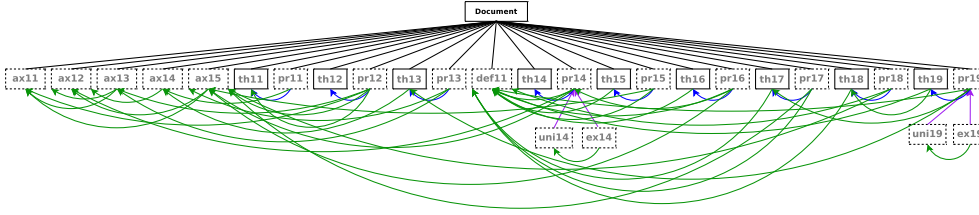


FIG. 24. The DG of sections 1 and 2 of chapter 1 of Landau's book

The GoTO of the example is also produced automatically.<sup>7</sup> There are no errors or warning in the document which means we have no loops in the GoTO, no proofs for unproved nodes, no double proofs for a node and no missing proofs for proved nodes.

### 5.3 Generation of the proof skeleton

Since there are no errors in the GoTO, the proof skeleton can be produced without any warnings. We have 8 different mathematical roles in the document: axioms, definitions, theorems, proofs, cases, case, existenceParts and uniquenessParts. We make a distinction between cases and case because e.g. in Mizar we have a special keyword introducing cases (`per cases;`) and then keywords for each case (`suppose ...`). So we annotated the cases as child nodes of the case node. Table 9 gives an overview of the rules that were used to generate the Mizar and the Coq proof skeleton. Since in Coq there are no special keywords for uniqueness, existence or cases, these rules translate only the body of these nodes and add no keywords.

Role	Mizar rule	Coq rule
axiom	<code>%name : %body ;</code>	<code>Axiom %name : %body .</code>
definition	<code>definition %name : %nl %body %nl end;</code>	<code>Definition : %body .</code>
theorem	<code>theorem %name: %nl %body</code>	<code>Theorem %name : %body .</code>
proof	<code>proof %nl %body %nl end;</code>	<code>Proof %name : %body .</code>
cases	<code>per cases; %nl</code>	<code>%body</code>
case	<code>suppose %nl %body %nl end;</code>	<code>%body</code>
existencePart	<code>existence %nl %body</code>	<code>%body</code>
uniquenessPart	<code>uniqueness %nl %body</code>	<code>%body</code>

TABLE 9. The Mizar and Coq rules for the dictionary

In table 10 we give a part of the skeleton for section 4 (Mizar on the left, Coq on the right).<sup>8</sup>

### 5.4 Completing the proofs in Coq

As we already explained, MathLang aims to remain as independent as possible of a particular foundation, while in addition facilitating the process of formalising mathematics in different theorem provers. Two PhD students of the MathLang project

<sup>7</sup>The GoTO of the whole chapter can be found in the extended article on the web pages of the authors.

<sup>8</sup>The complete output of the skeleton for Mizar and Coq for the whole chapter can be found in the extended article on the web pages of the authors.

theorem th131: <th131>	Theorem th131: <th131> .
proof <pr131> end;	Proof. <pr131> Qed.
theorem th132: <th132>	Theorem th132: <th132> .
proof	Proof.
per cases;	
suppose <pr132case1> end;	<pr132case1>
suppose <pr132case2> end;	<pr132case2>
suppose <pr132case3> end;	<pr132case3>
end;	Qed.

TABLE 10: Part of the Mizar (on the left) and Coq (on the right) output from gSGA

(Retel, respectively Lamar) are concerned with the MathLang paths into Mizar, respectively Isar. In this paper we study for the first time the MathLang path into Coq. We show how the CGa, TSa and DRa encoding of chapter one of Landau’s book is taken into a fully formalised Coq code.

Currently we use the proof skeleton produced in the last section and fill all the %body parts by hand. We intend to investigate in the future how parts of the CGa and DRa annotations can be transformed automatically to Coq. In this section we explain why the process of formalising a mathematical text into Coq through MathLang is simpler than the formalisation of the text directly into Coq.

To begin with, we code the preface of the document (see table 5.1.1). The most complicated section to code in Coq was the first one, because we had to translate the axioms in a way we can use them productively in Coq. We defined the natural numbers as an inductive set - just as Landau does in his book.

```
Inductive nats : Set :=
| I : nats
| succ : nats -> nats
```

Then we translate axioms 2 - 4 almost literally from our CGa annotations. For example the annotation of Axiom 3 (“ax13”) in our document is:



By just viewing the interpretations of the annotations we get:

$$\text{forall } x \text{ (neq (succ}(x), 1)) \tag{a}$$

The automatically generated Coq proof skeleton for this axiom is:

Axiom ax13 : <ax13> . (b)

Now, we simply replace the <ax13> placeholder of (b) with the literal translation of the interpretations in (a) to get the valid Coq axiom (this literal translation could also be done by an algorithm that we plan to implement soon):

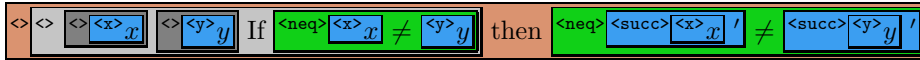
Axiom ax13 : forall x:nats, neq (succ x) I .

The other axioms could be completed in a similar way and as seen, this is a very simple process that can be carried out using automated tools that reduce the burden on the user (the proof skeleton is automated, the interpretations are obtained automatically from the CGa annotations which are simple to do, and for many parts of the text, the combination of the proof skeleton with the interpretations can also be automated).

Similarly for the theorems of chapter 1 of Landau’s book, the work needed to get the full formalisation is straightforward: E.g. Theorem 1 is written by Landau as:

If  $x \neq y$  then  $x' \neq y'$

Its annotation in MathLang CGa is:



The CGa annotation of the context (also called local scoping) can also be seen as the premise of an implication. So the upper statement can be translated to (again we see here that this is only a simple rewriting of the interpretations of the annotations):

decl(x), decl(y) : neq x y -> neq (succ x) (succ y)

And when we compare this line with its Coq translation we see again, it is just a literal transcription of the interpretation parts of CGa and therefore could be easily performed by an algorithm.

Theorem th11 (x y:nats) : neq x y -> neq (succ x) (succ y) .

From the 36 theorems of the chapter 28 could be translated literally into their corresponding Coq theorems.

Now, we want also to look at a simple proof and how it can be translated into Coq. The encoding of Theorem 2 of the first chapter in Coq is

theorem th12 (x:nats) : neq (succ x) x .

Landau proves this theorem with induction. He first shows, that  $1' \neq 1$  and then that with the assumption of  $x' \neq x$  it also holds that  $(x')' \neq x'$ .

Since we defined the natural numbers as an inductive set, we can also do our proof in the Landau style. We introduce the variable  $x$  and eliminate it, which yields two subgoals that we need to prove. These subgoals are exactly the induction basis and the induction step.

Proof.

intro x. elim x.

2 subgoals

```
x : nats
----- (1/2)
neq (succ I) I

----- (2/2)
forall n : nats, neq (succ n) n -> neq (succ (succ n)) (succ n)
```

Landau proved the first case with the help of Axiom 3 which states, that for all  $x$  it holds that  $x' \neq 1$ . We can just apply this axiom in Coq to prove the first case:

```
apply ax13.
```

```
1 subgoal
x : nats
----- (1/1)
forall n : nats, neq (succ n) n -> neq (succ (succ n)) (succ n)
```

The next step is to introduce  $n$  as natural number and to introduce the induction hypothesis:

```
intros N H.
```

```
1 subgoal
x : nats
n : nats
H : neq (succ n) n
----- (1/1)
neq (succ (succ n)) (succ n)
```

We see that this is exactly the second case of Landau's proof. He proved this case with Theorem 1 - we do the same:

```
apply th11.
```

```
1 subgoal
x : nats
n : nats
H : neq (succ n) n
----- (1/1)
neq (succ n) n
```

And of course this is exactly the induction hypotheses which we already have as an assumption and we can finish the proof:

```
assumption.
```

```
Proof completed.
```

The complete theorem and its proof in Coq finally look like this:

```
Theorem th12 (x:nats) : neq (succ x) x .
```

```
Proof.
intro x. elim x.
apply ax13.
intros n H.
apply th11.
assumption.
Qed.
```

We also used another hint for translating from the CGa part to the Coq formalisation. When we have a Theorem of the following kind:

```
Theorem th11 (x y:nats) : neq x y -> neq (succ x) (succ y) .
```

This is equivalent to:

```
Theorem th11 : forall x y:nats, neq x y -> neq (succ x) (succ y) .
```

A proof of such a theorem always starts with the introduction of the universal quantified variables, so in this case  $x$  and  $y$ . In terms of Coq this means:

```
intros x y.
```

We can do this for every proof. If it is a proof by induction we can also choose the induction variable in the next step. For example if we have an induction variable  $x$  we would write:

```
elim x.
```

We took the proof skeleton for Coq and extended it with these hints and the straightforward encoding of the 28 theorems. The result can be found in the extended article on the authors' web pages. With the help of these hints we were able to produce 234 lines of correct Coq lines. The completed proof has 957 lines. In other words, we could automatically generate one fourth of the complete formalised text. This is a large simplification of the formalisation process, even for an expert in Coq who can then better devote his attention to the important issues of formalisation: the proofs.

Of course there are some proofs within this chapter whose translation is not as easy and straightforward as the proof of Theorem 2 given above. But with the help of the CGa annotations and the automatically generated proof skeleton, we have completed the Coq proofs of the whole of chapter one in a couple of hours. As we said above, the combination of interpretations and proof skeletons can be implemented so that it leads for parts of the text, into automatically generated Coq proofs. This will speed further the formalisation and again will remove more burdens from the user. The complete Coq proof of chapter 1 of Landau's book can again be found in the extended article on the authors' web pages.

## 6 Conclusion

### 6.1 *Work and accomplishments so far*

The work on MathLang has been led by Kamareddine and Wells, has involved the hard work of 4 Ph.D. students (Maarek, Retel, Lamar, and Zengler), and has benefited from

implementation and evaluation work by numerous undergraduate and master's degree students. This work has led to a number of publications [21, 20, 22, 16, 17, 18, 19] and Maarek's Ph.D. thesis [29]. Retel's Ph.D. thesis is expected imminently, while Lamar's Ph.D. studies are nearing the halfway point. Zengler has just started with his studies in 2008.

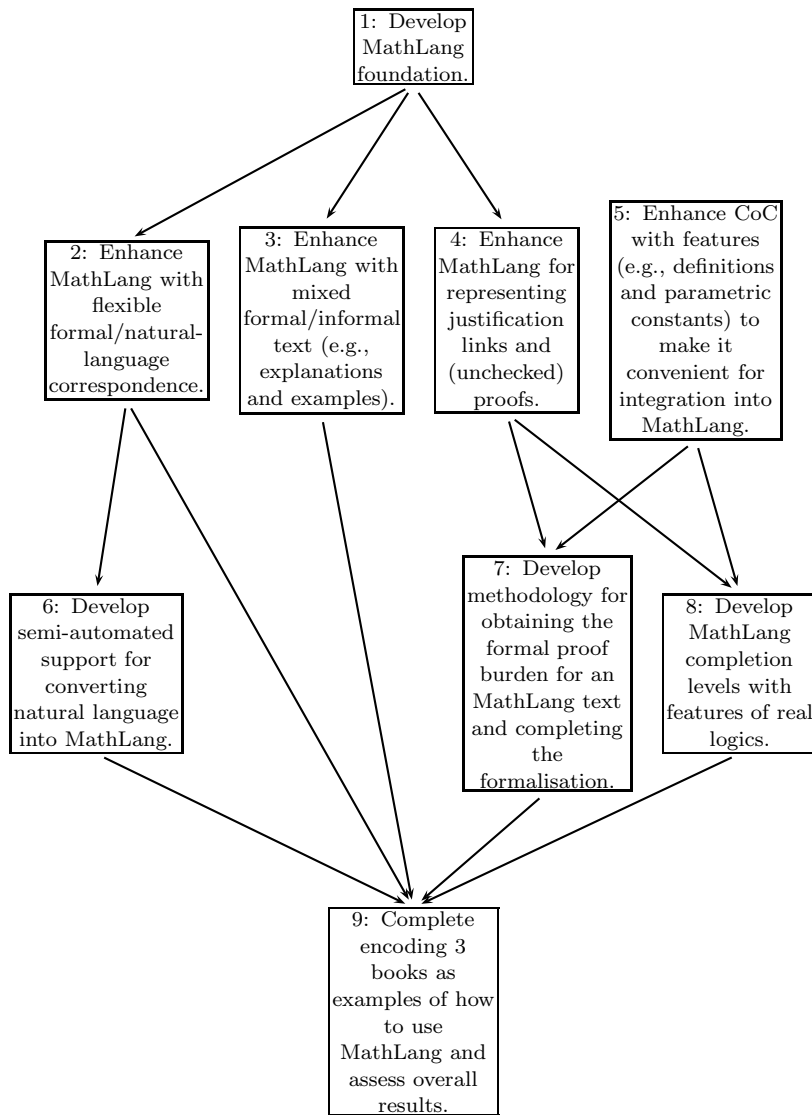


FIG. 25. MathLang project planning diagram from the year 2001

To compare our initial plans and our actual achievements, figure 25 contains a project planning dependency diagram from an early MathLang project plan that

shows the tasks and the dependencies between them as we imagined them in 2001.<sup>9</sup> Seven years later, we can see that the development of CGa [21, 22, 29, 17], TSa [20, 16, 29, 17], and DRa [19, 18] represent substantial progress on tasks 1, 2, 3, and 4, although more remains to be done on these tasks.

[24] summarised this progress on tasks 1..4 up to November 2007, stated that little progress was made on tasks 5 and 6, and summarised the partial progress on tasks 7 and 8. Task 9 has been ongoing since the start of the MathLang project and will continue for years to come.

Since [24], we have achieved further progress on these tasks especially on tasks 1, 7 and 9. In this paper we reported on this progress. For task 1, we have given the foundations behind the DRa annotation of a text and implemented the automatic checking of the DRa annotated text. However, the biggest progress since [24] is on tasks 7 and 9. For task 7, we developed the proof skeleton idea presented earlier in [18] specifically for Mizar, into an automatically generated proof skeleton in a choice of theorem provers (including Mizar, Isar and Coq). To achieve this, we gave a generic algorithm for proof skeleton generation which takes the required prover as one of its arguments. We also gave hints for the development of a generic algorithm which is able to convert parts of a CGa annotated text automatically into the syntax of the theorem prover it is given as an argument.

For task 9, we have given the complete MathLang encoding and full formalisation into Coq via the MathLang path of the first chapter of Landau’s book [28]. Landau’s book is one of the 3 books proposed for task 9 right at the start of the MathLang project. As such, parts of this book have been the subject of numerous encodings within the MathLang project (as part of undergraduate, MSc and PhD projects within MathLang). However, these encodings have only taken place at the lower aspects of MathLang. Here, we encode the first chapter through all the aspects and follow the whole MathLang path all the way to develop a full formalisation in the Coq prover. This is the most extensive example to date of how a text can pass through all the path of MathLang from the version written by the mathematician into a fully formalised text. We give examples to illustrate that this formalisation into Coq using the MathLang path is indeed easier than a direct formalisation into Coq.

We have not yet worked on task 5, which envisioned altering the structure of a specific mathematical foundation to make it more suitable for embedding in ordinary mathematical documents. We might still do this, but we are now more likely to work on techniques that can be used with more than one foundation and that can avoid altering the foundation.

Our initial plans for task 8 have replaced a concept we called “completion levels” by our notion of “aspects”. The semi-automatic extraction of mathematical structure from CML proposed in task 6 remains for future work.

## 6.2 *Future work and planned results*

MathLang is a long-term project and we expect there will be years of design, implementation, and evaluation, followed by repeated redesign, reimplementing, and

---

<sup>9</sup>The only change from our original diagram is that in addition to calling the project “MathLang” (as we always have), we now also call the language “MathLang” instead of the older name of “NML” (New Mathematical Language).

re-evaluation. There are many areas which we have identified as needing more work and investigation. One area is improvements to the MathLang software (currently based on the  $\text{\TeX}_{\text{MACS}}$  editor) to make it easier to enter information for the core MathLang aspects (currently CGa and DRa). This is likely to include work on semi-automatically recognising the mathematical meaning of natural language text. A second area is further designing and developing the portions of MathLang needed for better support of formalisation. An issue here is how much expertise in any particular target proof system will be needed for authoring. It may be possible to arrange things in MathLang to make it easy for an expert in a proof system to collaborate with an ordinary mathematician in completing a formalisation. A third area where work is needed is in the overall evaluation process needed to ensure MathLang meets actual needs. This will require testing MathLang with ordinary mathematicians, mathematics students, and other users. And there are additional areas where work will be needed, including areas we have not yet anticipated.

The MathLang project aims for a number of outcomes. MathLang aims to support mathematics as practised by the ordinary mathematician, which is generally not formalised, as well as work toward full formalisation. We expect that after further improvements on the MathLang design and software, writing MathLang documents (without formalising them) will be easy for ordinary mathematicians. MathLang will support various kinds of consistency checking even for non-formalised mathematics. MathLang will be independent of any particular logical foundation of mathematics; individual documents will be able to be formal in one or more particular foundations, or not formalised.

MathLang hopes to open a new useful era of collaboration between ordinary mathematicians, logicians (who ordinarily stay apart from other mathematicians), and computer science researchers working in such areas as theorem proving and mathematical knowledge management who can develop tools to link them together. The MathLang project's outputs will include a document representation, software suitable for manipulating this representation, and documentation and guidance for how to use the representation and the software. MathLang's document representation is intended to help with various kinds of automated computerised processing of mathematical knowledge. It should be possible to link MathLang documents together to form a public library of reusable mathematics. MathLang aims to better support translation between natural languages of mathematical texts and multi-lingual texts. MathLang aims to better support the differing uses of mathematical knowledge by different kinds of people, including ordinary practising mathematicians, students, computer scientists, logicians, linguists, etc.

## References

- [1] J. Abbott, A. van Leeuwen, and A. Strotmann. Objectives of openmath. Technical Report 12, RIACA (Research Institute for Applications of Computer Algebra), 1996. The TR archives of RIACA are incomplete. Earlier versions of this paper can be found at the "old OpenMath Home Pages" archived at the Uni. Köln.
- [2] Georg Cantor. Beiträge zur Begründung der transfiniten Mengenlehre (part 1). *Mathematische Annalen*, 46:481–512, 1895.
- [3] Georg Cantor. Beiträge zur Begründung der transfiniten Mengenlehre (part 2). *Mathematische Annalen*, 49:207–246, 1897.

- [4] Augustin-Louis Cauchy. *Cours d'Analyse de l'École Royale Polytechnique*. Debure, Paris, 1821. Also in *Œuvres Complètes* (2), volume III, Gauthier-Villars, Paris, 1897.
- [5] W3C (World Wide Web Consortium). Mathematical markup language (MathML) version 2.0. W3C Recommendation, October 2003.
- [6] W3C (World Wide Web Consortium). RDF Primer. W3C Recommendation, February 2004.
- [7] W3C (World Wide Web Consortium). XQuery 1.0 and XPath 2.0 data model (XDM). W3C Recommendation, 2007.
- [8] N.G. de Bruijn. The mathematical vernacular, a language for mathematics with typed sets. In *Workshop on Programming Logic*, 1987. Reprinted in [31, F.3].
- [9] Richard Dedekind. *Stetigkeit und irrationale Zahlen*. Vieweg & Sohn, Braunschweig, 1872. Fourth edition published in 1912.
- [10] Gottlob Frege. *Begriffsschrift: eine der arithmetischen nachgebildete Formelsprache des reinen Denkens*. Nebert, Halle, 1879. Can be found on pp. 1–82 in [38].
- [11] Gottlob Frege. *Grundgesetze der Arithmetik*, volume 1. Hermann Pohle, Jena, 1893. Republished 1962 (Olms, Hildesheim).
- [12] Gottlob Frege. *Grundgesetze der Arithmetik*, volume 2. Hermann Pohle, Jena, 1903. Republished 1962 (Olms, Hildesheim).
- [13] G. Gierz, K. H. Hofmann, K. Keimel, J. D. Lawson, M. W. Mislove, and D. S. Scott. *A Compendium of Continuous Lattices*. Springer-Verlag, 1980.
- [14] Thomas L. Heath. *The 13 Books of Euclid's Elements*. Dover, 1956. In 3 volumes. Sir Thomas Heath originally published this in 1908.
- [15] Fairouz Kamareddine, Twan Laan, and Rob Nederpelt. *A Modern Perspective on Type Theory from Its Origins Until Today*, volume 29 of *Kluwer Applied Logic Series*. Kluwer Academic Publishers, May 2004.
- [16] Fairouz Kamareddine, Robert Lamar, Manuel Maarek, and J. B. Wells. Restoring natural language as a computerised mathematics input method. In MKM '07 [30], pages 280–295.
- [17] Fairouz Kamareddine, Manuel Maarek, Krzysztof Retel, and J. B. Wells. Digitised mathematics: Computerisation vs. formalisation. In *Review of the National Center for Digitization*, volume 10, pages 1–8, Faculty of Mathematics, Belgrade, Serbia, 2007.
- [18] Fairouz Kamareddine, Manuel Maarek, Krzysztof Retel, and J. B. Wells. Gradual computerisation/formalisation of mathematical texts into Mizar. In Roman Matuszewski and Anna Zalewska, editors, *From Insight to Proof: Festschrift in Honour of Andrzej Trybulec*, volume 10(23) of *Studies in Logic, Grammar and Rhetoric*, pages 95–120. University of Białystok, 2007. Under the auspices of the Polish Association for Logic and Philosophy of Science.
- [19] Fairouz Kamareddine, Manuel Maarek, Krzysztof Retel, and J. B. Wells. Narrative structure of mathematical texts. In MKM '07 [30], pages 296–311.
- [20] Fairouz Kamareddine, Manuel Maarek, and J. B. Wells. Flexible encoding of mathematics on the computer. In *Mathematical Knowledge Management, 3rd Int'l Conf., Proceedings*, volume 3119 of *Lecture Notes in Computer Science*, pages 160–174. Springer, 2004.
- [21] Fairouz Kamareddine, Manuel Maarek, and J. B. Wells. Mathlang: Experience-driven development of a new mathematical language. In *Proc. [MKMNET] Mathematical Knowledge Management Symposium*, volume 93 of *ENTCS*, pages 138–160, Edinburgh, UK (2003-11-25/--29), February 2004. Elsevier Science.
- [22] Fairouz Kamareddine, Manuel Maarek, and J. B. Wells. Toward an object-oriented structure for mathematical text. In *Mathematical Knowledge Management, 4th Int'l Conf., Proceedings*, volume 3863 of *Lecture Notes in Artificial Intelligence*, pages 217–233. Springer, 2006.
- [23] Fairouz Kamareddine and Rob Nederpelt. A refinement of de Bruijn's formal language of mathematics. *J. Logic Lang. Inform.*, 13(3):287–340, 2004.
- [24] Fairouz Kamareddine and J. B. Wells. Computerizing mathematical text with mathlang. *Electron. Notes Theor. Comput. Sci.*, 205:5–30, 2008.
- [25] Toshihiro Kanahori, Alan Sexton, Volker Sorge, and Masakazu Suzuki. Capturing abstract matrices from paper. In *Mathematical Knowledge Management, 5th Int'l Conf., Proceedings*, volume 4108 of *Lecture Notes in Computer Science*, pages 124–138. Springer, 2006.
- [26] Michael Kohlhase. *An Open Markup Format for Mathematical Documents, OMDoc (Version 1.2)*, volume 4180 of *Lecture Notes in Artificial Intelligence*. Springer-Verlag, 2006.

- [27] Edmund Landau. *Grundlagen der Analysis*. Chelsea, 1930.
- [28] Edmund Landau. *Foundations of Analysis*. Chelsea, 1951. Translation of [27] by F. Steinhardt.
- [29] Manuel Maarek. *Mathematical Documents Faithfully Computerised: the Grammatical and Text & Symbol Aspects of the MathLang Framework*. PhD thesis, Heriot-Watt University, Edinburgh, Scotland, June 2007.
- [30] *Towards Mechanized Mathematical Assistants (Calcuemus 2007 and MKM 2007 Joint Proceedings)*, volume 4573 of *Lecture Notes in Artificial Intelligence*. Springer, 2007.
- [31] Rob Nederpelt, J. H. Geuvers, and Roel C. de Vrijer. *Selected Papers on Automath*, volume 133 of *Studies in Logic and the Foundations of Mathematics*. North-Holland, Amsterdam, 1994.
- [32] Giuseppe Peano. *Árithmetices Principia, Nova Methodo Exposita*. Bocca, Turin, 1889. An English translation can be found on pp. 83–97 in [38].
- [33] P. Rudnicki. An overview of the Mizar project. In *Proceedings of the 1992 Workshop on Types for Proofs and Programs*, 1992.
- [34] Alan Sexton and Volker Sorge. The ellipsis in mathematical documents. Talk overhead images presented at the IMA (Institute for Mathematics and its Applications, University of Minnesota) “Hot Topic” Workshop The Evolution of Mathematical Communication in the Age of Digital Libraries held on 2006-12-08/---09.
- [35] Lambert S. van Benthem Jutting. *Checking Landau’s “Grundlagen” in the AUTOMATH System*. PhD thesis, Eindhoven, 1977. Partially reprinted in [31, B.5,D.2,D.3,D.5,E.2].
- [36] Lambert S. van Benthem Jutting. *Checking Landau’s “Grundlagen” in the AUTOMATH system*. PhD thesis, Eindhoven, 1977.
- [37] Joris van der Hoeven. GNU TeXmacs. *SIGSAM Bulletin*, 38(1):24–25, 2004.
- [38] J. van Heijenoort. *From Frege to Gödel: A Source Book in Mathematical Logic, 1879–1931*. Harvard University Press, 1967.
- [39] Alfred North Whitehead and Bertrand Russel. *Principia Mathematica*. Cambridge University Press, 1910–1913. In three volumes published from 1910 through 1913. Second edition published from 1925 through 1927. Abridged edition published in 1962.
- [40] Ernst Zermelo. Untersuchungen über die Grundlagen der Mengenlehre (part 1). *Mathematische Annalen*, 65:261–281, 1908. An English translation can be found on pp. 199–215 in [38].

Received 05 May 2008