

Automatically and Efficiently Illustrating Polynomial Equalities in Agda

Donnacha Oisín Kidney

Supervisor: Professor Gregory Provan

Final-Year Project—BSc in Computer Science

Department of Computer Science
University College Cork

April 8, 2019

Abstract

We present a new library which automates the construction of equivalence proofs between polynomials over commutative rings and semirings in the programming language Agda [20]. It is significantly faster than Agda's existing solver. We use reflection to provide a simple interface to the solver, and demonstrate how to use the constructed proofs to provide step-by-step solutions.

1 Declaration of Originality

In signing this declaration, you are conforming, in writing, that the submitted work is entirely your own original work, except where clearly attributed otherwise, and that it has not been submitted partly or wholly for any other educational award.

I hereby declare that:

- this is all my own work, unless clearly indicated otherwise, with full and proper accreditation;
- with respect to my own work: none of it has been submitted at any educational institution contributing in any way to an educational award;
- with respect to another's work: all text, diagrams, code, or ideas, whether verbatim, paraphrased or otherwise modified or adapted, have been duly attributed to the source in a scholarly manner, whether from books, papers, lecture notes or any other student's work, whether published or unpublished, electronically or in print.

Signed: _____

Date: _____

Contents

1	Declaration of Originality	3
2	Introduction	6
2.1	Background	6
2.2	Our Contributions	8
2.3	Scope	9
3	Overview of the Proof Technique	11
4	The Interface	12
4.1	Implementation	13
4.2	Maintaining Invariants	16
5	Performance	17
5.1	Normalisation	17
5.1.1	Horner Normal Form	17
5.1.2	Sparse Encodings	18
5.1.3	Hanging Indices	19
5.2	Unification	21
5.2.1	Avoid Progress at all Costs	21
5.2.2	Avoid Identities	22
5.3	Benchmarks	24
6	Verification	25
7	Pedagogical Solutions	25
8	Related Work	27
9	Conclusion	28

2 Introduction

2.1 Background

The Foundational Crisis In the early 20th century, the foundations of mathematics began to crumble. The first crack was Russell’s paradox, discovered in 1901. Similar paradoxes soon followed: each represented core errors in the foundational theories that had underpinned all of mathematics up until that point.

The work of papering over these cracks began in earnest: Russell himself (along with Whitehead) published the *Principia Mathematica* (PM) [26], a formalisation of contemporary mathematics based on a theory that accounted for the discovered paradoxes. The task proved to be far more difficult than was anticipated: the PM was infamously verbose (350 pages of preparation precede the proof that $1 + 1 = 2$), and both men eventually gave up before reaching their goal of a full formalisation.

These repeated attempts and failures to shore up the foundations of mathematics became known as the “Foundational Crisis”. There were three competing schools of thought on how to address the crisis: Russell and Whitehead belonged to the logicism camp, although this came to be overshadowed by the other two philosophies. David Hilbert championed the “formalists”, the leading school of thought at the time. These mathematicians envisioned a solution to the crisis in the form of a finite set of axioms, satisfying a set of requirements (here simplified):

Completeness Any true statement can be derived from the axioms.

Consistency No untrue statement can be derived from the axioms.

Decidability Any statement can be decided as true or false using an algorithm.

These requirements became known as Hilbert’s program.

While intended to be a scaffolding on which to rebuild, these requirements acted more as targets for later attacks. Firstly, in 1931, Gödel published his incompleteness theorems [10], which proved that the first two requirements (completeness and consistency) were *impossible* to achieve. And finally, in 1936, Alonzo Church and Alan Turing independently showed that the Entscheidungsproblem was unsolvable [5, 6], showing that the third requirement was also impossible to satisfy.

Intuitionism Among the rubble stood intuitionism: the third school of thought, once maligned, now was more attractive in light of the failure of Hilbert’s program. Formalists and intuitionists had deep philosophical disagreements: the formalist

held that mathematics was the pursuit of external truth. Hilbert, in particular, insisted that the axioms and rules of a mathematical formalism were not arbitrary:

We are not speaking here of arbitrariness in any sense. Mathematics is not like a game whose tasks are determined by arbitrarily stipulated rules. Rather, it is a conceptual system possessing internal necessity that can only be so and by no means otherwise. (David Hilbert [12])

Intuitionists, on the other hand, believed that mathematics was a pure construct of the human mind, and had little or nothing to do with objective reality. Its main proponent was L. E. J. Brouwer, who became engaged in a bitter dispute with Hilbert, eventually culminating in Brouwer’s removal from the *Mathematische Annalen*, on the basis of Hilbert’s claim that his theories represented a “danger to mathematics” [23].

In terms of nuts-and-bolts, the defining feature separating intuitionism from classical logic is the absence of the law of the excluded middle, or “choice”:

$$\forall P. P \vee \neg P \tag{1}$$

This axiom has a deep relation with constructiveness and computability. Intuitionism is a constructive theory: to say that something is true is to say that we have a proof of it. Allowing this axiom, then, would imply that once we could state a theorem intuitionistically, we would be able to pluck—either from thin air or from an algorithm—a proof of its truth or falsehood. But of course, via the Entscheidungsproblem, we know that this is impossible.

Intuitionistic Type Theory and Agda Intuitionism has had many incarnations since Brouwer (it gained widespread popularity after Bishop), but the particular version of interest to us comes from a strange source: the type systems of programming languages. In formal terms, we are talking about the “Curry-Howard isomorphism” (Fig. 1): a way to bridge the world of programming languages and formalised mathematics.

This brings us, at long last, to Agda, the subject of this work. Agda is a programming language and intuitionistic theory based on Per Martin-Löf’s intuitionistic type theory [15]. Its syntax and evaluation strategy is similar to Haskell, but as a formalism it can be used to prove mathematical statements.

$$Type \iff Proposition$$

$$Program \iff Proof$$

Figure 1: The Curry-Howard Isomorphism

lemma : $\forall x y \rightarrow x + y * 1 + 3 \approx 2 + 1 + y + x$

```
lemma x y = begin
  x + y * 1 + 3 ≈⟨ refl ⟨ +-cong ⟩ *-identityr y ⟨ +-cong ⟩ refl {3} ⟩
  x + y + 3    ≈⟨ +-comm x y ⟨ +-cong ⟩ refl ⟩
  y + x + 3    ≈⟨ +-comm (y + x) 3 ⟩
  3 + (y + x)  ≈⟨ sym (+-assoc 3 y x) ⟩
  2 + 1 + y + x ■ lemma = solve NatRing
```

(a) A Tedious Proof

(b) Our Solver

Figure 2: Comparison Between A Manual Proof and The Automated Solver

2.2 Our Contributions

It is over a hundred years since the publication Principia Mathematica, and we are still a long way away from formalising all of mathematics. There is a popular website which tracks the progress of this formalisation against the “100 greatest theorems” [27]: at time of writing, the number stands at 93.

There are many reasons for why we haven’t managed to formalise all 100 yet: chief among them is that, while we have come far from the days of the PM, proofs are still verbose and tedious.

This work intends to alleviate some of the tedium of constructing a particular kind of proof: identities over commutative rings. We write a verified and proven library for automation in Agda, to automate the construction of proofs like the one in Fig. 2a, making them as simple as Fig. 2b.

The main contributions of our library are as follows:

Ease of Use Proofs like the one in Fig. 2a are long, difficult to write, and uninteresting. Our solver, in contrast, is extremely simple to use: the single line in Fig. 2b solves the lemma.

This interface (section 4) is implemented using a lightweight reflection system (section 4.1), which does not require the user to write *any* reflection code, even if they use the solver with their own custom type.

Performance Our solver is significantly faster than Agda’s current ring solver, cutting type-checking time down from minutes to seconds in several use cases (section 5.3).

As described in [11], we use a sparse internal representation of polynomials (section 5.1.2). However, because of differences between Agda and Coq’s type checker, we found that this optimisation—on its own—did not deliver a significant speedup, and actually damaged performance in a number of cases. Achieving the performance we did required an entirely separate kind of optimisation, described in section 5.2.

Pedagogical Solutions Computer algebra systems (CASs) outside the rigorous world of dependently-typed languages do far more than just check proofs for mistakes: they have a wealth of other features which can help with learning mathematics as well as verifying it.

We hope that similar systems developed in Agda can do the same: as a demonstration, we implement “step-by-step solutions”, one of the most popular features of modern CASs. Far from being ill-suited to Agda, we show that the constructive nature of our proofs allows for a natural implementation (section 7).

2.3 Scope

In this section, we will explain the intended uses and necessary limitations of the solver.

Equivalences First, an inflexibility: the solver deals very specifically with the domain of *equivalence* proofs, like the one in Fig. 2. While it may be of use in other settings (finding roots, etc.), that is not explored here.

Setoids On the other hand, we are very flexible about what kind of “equivalence” we prove. In fact, the solver will work with any equivalence relation, as long as it comes with proofs of the relevant ring axioms. This is useful for all of the usual things (approximating quotients and so on), but it also provides the basis for our “step-by-step solutions” implementation in section 7.

Almost Rings As in [11, section 5], we use a peculiar algebraic structure which lies somewhere between a semiring and a ring. These “almost-rings” have all of the usual laws of a commutative ring, but instead of demanding additive inverses, they require the comparatively permissive “pseudo-inverse” operation, which obeys the following equations:

$$-(x * y) = -x * y \tag{2}$$

$$-(x + y) = -x + -y \tag{3}$$

This allows the solver to work on types which don't have an additive inverse (like \mathbb{N}): such types just supply the identity function instead of negation, and the two laws above are satisfied.

A potential worry is that because we don't require $x + -x = 0$ axiomatically, it won't be provable in our system. Happily, this is not the case: as long as $1 + -1$ reduces to 0 in the coefficient set, the solver will verify the identity.

Weak Decidability A core optimisation in our solver (section 5.1.2) relies on the ability to test arbitrary coefficients for zero. Instead of requiring decidable equality (which would greatly diminish the number of types the solver can work with), we instead ask for weakly decidable equivalence with zero:

$$\text{is-zero} : \forall x \rightarrow \text{Maybe } (0 \# \approx x)$$

Just as in Agda's current solver, this allows users to avail of the optimisation if their type supports it, or skip it (`is-zero = const nothing`) if not.

Correctness The nature of the solver means it is intrinsically sound (i.e. it cannot prove an equivalence unless there is one): since all it does is rearrange and join together the ring axioms, it cannot prove anything that does not derive from them. We have not, however, proven completeness (that every equivalence will be found by our solver).

In the internal representation of the solver, we prove several data structure invariants (like sparsity) intrinsically.

The reflection-based interface is unproven, but since the output is type checked our claim of soundness still stands: a bug in our reflection code can only cause the solver to miss a solution, never to prove something it shouldn't.

3 Overview of the Proof Technique

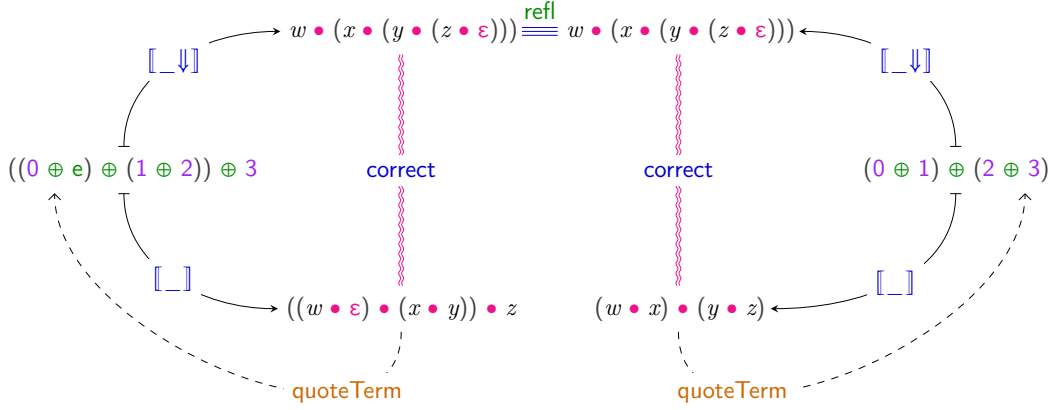


Figure 3: The Reflexive Proof Process

Before diving into specifics, we'll first give a quick overview of how the solver works, so it's clear how the bits of implementation described later in the paper fit together.

The technique we use for automating equivalence proofs comes from [1]: the general idea is that we prove two expressions equivalent by proving that they're both equivalent to the same canonical form.

The diagram in Fig. 3 demonstrates this for the identity from Fig. 2: on the bottom of the diagram you can see the left and right hand side of the identity we want to prove, and on the top we can see their normal forms. The actual proof the solver provides is represented by the \approx path.

To prove that each expression is equivalent to the canonical form we first represent the expressions using the type in Fig. 4. This type represents the Abstract Syntax Tree (AST) for expressions in the almost-ring algebra: it has constructors corresponding to each ring operation ($x + y = x \oplus y$, $x * y = x \otimes y$), and it can refer to variables via their de Bruijn index (so x becomes $\text{I } 0$).

There are two ways to evaluate the AST: the $[_]$ function converts the AST to the expression we want to prove, whereas $[_ \Downarrow]$ converts it to a canonical form. The implementation of the $[_ \Downarrow]$ function is described in section 5.1.

```
data Expr {ℓ} (A : Set ℓ) (n : ℕ) : Set ℓ where
  K   : A → Expr A n
  I   : Fin n → Expr A n
  _⊕_ : Expr A n → Expr A n → Expr A n
  _⊗_ : Expr A n → Expr A n → Expr A n
  _⊗_ : Expr A n → ℕ → Expr A n
  ⊖_  : Expr A n → Expr A n
```

Figure 4: A Type for Ring Expressions

Equivalence of the canonical forms is proven via `correct`: some of the details of this are explained in section 6.

Finally, instead of asking the user to construct the AST themselves, we use reflection to automate it. This is described in the following section.

4 The Interface

We felt an easy-to-use interface was one of the most important components of the library as a whole. Since we wanted to minimise the amount a user would have to learn to use the solver, we kept the surface area of the library quite small: aside from the almost-ring type, the rest of the interface consists of just two macros (`solve` and `solveOver`). We tried to make their usage as obvious as possible: just stick one of them (with the required arguments) in the place you need a proof, and the solver will do the rest for you.

`solve` is demonstrated in Fig. 2b. It takes a single argument: an implementation of the algebra. `solveOver` is designed to be used in conjunction with manual proofs, so that a programmer can automate a “boring” section of a larger more complex proof (Fig. 6). As well as the algebra implementation, this macro takes a list of free variables to use to compute the solution.

Because this interface is quite small, it’s worth pointing out what’s missing, or rather, what we *don’t* require from the user:

- We don’t ask the user to construct the `Expr` AST which represents their proof

```
lemma : ∀ x y → x + y * 1 + 3 ≡ 2 + 1 + y + x
lemma = +*-Solver.solve 2
(λ x y → x :+ y :* con 1 :+ con 3 := con 2 :+ con 1 :+ y :+ x)
refl
```

Figure 5: The Old Solver

```
lemma : ∀ x y → x + y * 1 + 3 ≈ 2 + 1 + y + x
lemma x y =
begin
  x + y * 1 + 3 ≈ { +-comm (x + y * 1) 3 }
  3 + (x + y * 1) ≈ { solveOver (x :: y :: []) Nat.ring }
  3 + y + x      ≡ { }
  2 + 1 + y + x   ■
```

Figure 6: The `solveOver` Macro

obligation. Compare this to Fig. 5: we had to write the type of the proof twice (once in the signature and again in the AST), and we had to learn the syntax for the solver’s AST.

As well as being more verbose, this approach is less composable: every change to the proof type has to be accompanied by a corresponding change in the call to the solver. In contrast, the call to `solveOver` above effectively amounts to a demand for the compiler to “figure it out!” Any change to the expressions on either side will result in an *automatic* change to the proof constructed.

- We don’t ask the user to write any kind of “reflection logic” for their type. In other words, we don’t require a function which (for instance) recognises and parses the user’s type in the reflected AST, or a function which does the opposite, converting a concrete value into the AST that (when unquoted) would produce an expression equivalent to the quoted value.

This kind of logic is complex, and very difficult to get right. While some libraries can assist with the task [19, 25] it is still not fully automatic.

4.1 Implementation

Agda has powerful metaprogramming facilities, which allow programs to manipulate their own code. Here, we’ll use reflection to implement the interface to our solver.

Agda’s reflection API is mostly encapsulated by the following three types:

Term The representation of Agda’s AST, retrievable via `quoteTerm`.

Name The representation of identifiers, retrievable via `quote`.

TC The type-checker monad, which includes scoping and environment information, can raise type errors, unify variables, or provide fresh names. Computations in the **TC** monad can be run with `unquote`.

While `quote`, `quoteTerm`, and `unquote` provide all the functionality we need, they’re somewhat low-level and noisy (syntactically speaking). Agda also provides a mechanism (which it calls macros) to package metaprogramming code so it looks like a normal function call (as in `solve`).

Reflection is obviously a powerful tool, but it has a reputation for being unsafe and error-prone. Agda’s reflection system does not break type safety, but we *are* able to construct **Terms** which are ill-typed, which often result in confusing error-messages on the user’s end. Unfortunately, constructing ill-typed terms is quite

easy to do: the `Term` type itself does not contain a whole lot of type information, and it's quite fragile and sensitive to context. Variables, for instance, are referred to by their de Bruijn indices, meaning that the same `Term` can break if it's simply moved under a lambda.

Building a robust interface using reflection required a great deal of care. To demonstrate some of the techniques we used, we'll look at two functions from the core of the interface. First, `toExpr`:

```

toExpr : Term → Term
toExpr (def (quote AlmostCommutativeRing._+_ ) xs) = getBinOp (quote _⊕_) xs
toExpr (def (quote AlmostCommutativeRing._*_ ) xs) = getBinOp (quote _⊗_) xs
toExpr (def (quote AlmostCommutativeRing._^_ ) xs) = getExp xs
toExpr (def (quote AlmostCommutativeRing._-_ ) xs) = getUnOp (quote _⊖_) xs
toExpr v@(var x _) with x ℕ.<? numVars
... | yes p = v
... | no ¬p = constExpr v
toExpr t = constExpr t

```

This function is called on the `Term` representing one side of the target equivalence. It converts it to the corresponding `Expr`. In other words, it performs the following transformation:

$$((w \bullet \varepsilon) \bullet (x \bullet y)) \bullet z \mapsto ((0 \oplus e) \oplus (1 \oplus 2)) \oplus 3$$

When it encounters one of the ring operators, it calls the corresponding helper function (`getBinOp`, `getExp`, or `getUnOp`) which finds the important subterms from the operator's argument list.

If it *does not* manage to match an operator or a variable, it assumes that what it has must be a constant, and wraps it up in the `K` constructor. This is the key trick which allows us to avoid ever asking the user to quote their own type. While it may seem unsafe at first glance, we actually found it to be more robust (for our use case) than the alternative:

Principle 1 (Don't reimplement the typechecker) *While it may seem good and fastidious to rigorously check the structure and types of arguments given to a macro, we found better results by avoiding validity-checking in metaprogramming code. Instead, we preferred to proceed as if there were no errors (if possible), but arrange the output so that the user would still see a type error where the input was incorrect.*

Taking this case as an example, if the user indeed manages to supply something other than the correct type, Agda will catch the error, as an incorrect argument to `K`.

If, on the other hand, we had asked the user to quote their own type, we would have trouble handling (for instance) closed applications of functions, references to names outside the lambda, etc. This approach, on the other hand, has no such difficulty.

Next, we'll look at one of the helper functions: `getExp`, which deals with exponentiation.

```
getExp : List (Arg Term) → Term
getExp (x <::> y <::> []) = quote _⊗_ { con } 3 ...<::> toExpr x <::> y <::> []
getExp (x :: xs) = getExp xs
getExp _ = unknown
```

It extracts the last two arguments to the exponentiation operator, and wraps them up with the `⊗`. Before the two visible arguments to the exponentiation operator, we first apply `3 ...<::>`. This applies three hidden arguments as “unknown”, i.e. asks Agda to infer them. We could guess them ourselves: the first is the universe level of the carrier type, the second is the carrier type, and the third is the number of variables in the expression. We decided against it, though, instead being intentionally unspecific:

Principle 2 (Supply the minimal amount of information) *There were several instances where, in constructing a term, we were tempted to supply explicitly some argument that Agda usually infers. Universe levels were a common example. We found this approach to be error-prone, however: as it turns out, the compiler is better at guessing implicits than we are. Instead, we preferred to leverage the compiler, relying on inference over direct metaprogramming as much as possible.*

The final point to make is that the entire interface implementation is itself quite small (fewer than 100 lines). This isn't because our code was terse: rather, we intentionally minimised the amount of metaprogramming we did.

Principle 3 (Keep Metaprogramming to the Edges) *With great power comes poor error messages, fragility, and a loss of first-class status. Therefore, If something can be done without reflection, do it, and use reflection as the glue to get from one standard representation to another.*

4.2 Maintaining Invariants

One obvious benefit of reflection is a terse interface. However, we feel that another benefit—resilience to change—is just as important. This section illustrates that resilience with an example.

Agda allows us to encode program correctness in types, so we can *prove* properties we would have otherwise only been able to test. Unfortunately, these kinds of proofs tend to be very tightly coupled to the implementation of the algorithms they verify. This can make iteration difficult, where small optimisations or bug fixes can invalidate proofs for other invariants.

To demonstrate the problem, and how our solver can reduce some of the burden, we'll look at size-indexed binary trees:

```
data Tree : ℕ → Set a where
  leaf : Tree 0
  node : ∀ {n m} → A → Tree n → Tree m → Tree (n + m + 1)
```

We've deliberately chosen an awkward type here: in contrast to the more common size-indexed lists, the index (the size) does not match the shape of the data structure. As a result, almost every function which manipulates the tree in some way will have to come accompanied by a verbose, complex proof. Take this line, for instance, which performs a left-rotation on the tree:

```
rot! (node {a} x xl (node {b} {c} y yr yl)) = node y (node x xl yl) yr
```

A sensible invariant to encode here is that the function does not change the size of the tree. Unfortunately, to *prove* that invariant, we have to prove the following:

$$1 + (1 + a + c) + b = 1 + a + (1 + b + c)$$

Though simple, this is precisely the kind of proof which requires many fussy applications of the ring axioms. Here, our solver can help:

```
rot! (node {a} x xl (node {b} {c} y yr yl)) = node y (node x xl yl) yr
⇒ ∀ { a :: b :: c :: [] }
```

While cutting down on the amount of code we need to write is always a good thing, the real strength of this method is that it automatically infers the input type. This makes it resilient to small changes in the code. So, when we notice the bug in the code above (*yl* and *yr* are swapped in the pattern-match), we can simply *fix it*, without having to touch any of the proof code.

$$\begin{aligned} \text{rot}^! (\text{node } \{a\} \ x \ xl (\text{node } \{b\} \ \{c\} \ y \ yl \ yr)) &= \text{node } y (\text{node } x \ xl \ yl) \ yr \\ &\Rightarrow \forall \langle a :: b :: c :: [] \rangle \end{aligned}$$

If we hadn't used the solver, this fix would have necessitated a totally new proof. By automating the proof, we allow the compiler to automatically check what we *mean* ("does the size of the tree stay the same?"), while we worry about other details.

5 Performance

Our solver is significantly faster than the current solver: the following sections will detail how we achieved that speedup. We will start by describing the naive implementation; in section 5.1.2 we demonstrate how we added the optimisations from [11]; and in section 5.2 we will describe the Agda-specific optimisations which account for the bulk of our speedup. Finally, section 5.3 contains some benchmarks against the current solver.

5.1 Normalisation

Most of code written for the solver is concerned with normalisation: the $\llbracket _ \rrbracket$ function in Fig. 3. This converts from the expression AST (Fig. 4) to a canonical form.

5.1.1 Horner Normal Form

The particular "canonical form" we'll start with is the same as in Agda's current ring solver: Horner normal form. A polynomial (more specifically, a monomial) in x is represented as a list of coefficients of increasing powers of x . As an example, the following polynomial:

$$3 + 2x^2 + 4x^5 + 2x^7 \tag{4}$$

Is represented by this list:

$$3 :: 0 :: 2 :: 0 :: 0 :: 4 :: 0 :: 2 :: []$$

Operations on these polynomials are similar to operations in positional number systems.


```

_⊞_ : Poly → Poly → Poly
[] ⊞ ys = ys
(x :: xs) ⊞ [] = x :: xs
(x :: xs) ⊞ (y :: ys) = x + y :: xs ⊞ ys

_⊗_ : Poly → Poly → Poly
_⊗_ [] _ = []
_⊗_ (x :: xs) =
  foldr (λ y ys → x * y :: map (_ * y) xs ⊞ ys) []

```

So to get from **Expr** to **Poly** we map each constructor to the relevant polynomial operation. Then, to get from **Poly** to an expression in the underlying ring, we use Horner’s rule: a classic example of the **foldr** function.

```

[[_]] : Poly → Carrier → Carrier
[ xs ] ρ = foldr (λ y ys → ρ * ys + y) 0# xs

```

5.1.2 Sparse Encodings

Our first avenue for optimisation comes from [11]. Notice that the encoding above is quite wasteful: it always stores an entry for each coefficient, even if it’s zero. In practice, we’re likely to often find long strings of zeroes (in expressions like x^{10}), meaning that our representation will contain long “gaps” between the coefficients we’re actually interested in (non-zero ones).

To fix the problem we’ll switch to a *sparse* encoding, by storing a “power index” with every coefficient. This will represent the size of the gap from the previous non-zero coefficient. Taking 4 again as an example, we would now represent it as follows:

```

(3, 0) :: (2, 1) :: (4, 2) :: (2, 1) :: []

```

Next, we turn our attention to the task of adding multiple variables. Luckily, there’s an easy way to do it: nesting. Multivariate polynomials will be represented as “polynomials of polynomials”, where each level of nesting corresponds to one variable. It’s perhaps more clearly expressible in types:

```

Poly : ℕ → Set c
Poly zero = Carrier
Poly (suc n) = List (Poly n × ℕ)

```

Inductively speaking, a “polynomial” in 0 variables is simply a constant, whereas a polynomial in n variables is a list of coefficients, which are themselves polynomials in $n - 1$ variables.

Before running off to use this representation, though, we should notice that we have created another kind of “gap” which we should avoid with a sparse encoding. For a polynomial with n variables, we will always have n levels of nesting, even if the polynomial does not actually refer to all n variables. In the extreme case, representing the constant 6 in a polynomial of 3 variables looks like the following:

((((6, 0) :: []), 0) :: [], 0) :: [])

The solution is another index: this time an “injection” index. This represents “how many variables to skip over before you get to the interesting stuff”. In contrast to the previous index, though, this one is type-relevant: we can’t just store a \mathbb{N} next to the nested polynomial to represent the gap. Because the polynomial is indexed by the number of variables it contains, any encoding of the gap will have provide the proper type information to respect that index.

5.1.3 Hanging Indices

The problem is a common one: we have a piece of code that works efficiently, and we now want to make it “more typed”, by adding more information to it, *without* changing the complexity class.

We found the following strategy to be useful: first, write the untyped version of the code, forgetting about the desired invariants as much as possible. Then, to add the extra type information, look for an inductive type which participates in the algorithm, and see if you can “hang” some new type indices off of it.

In our case, the injection index (distance to the next “interesting” polynomial) was simply stored as an \mathbb{N} , and the information we needed was the number of variables in the inner polynomial, and the number of variables in the outer. All of that is stored in the following proof of \leq :

```
data _≤_ (m : ℕ) : ℕ → Set where
  m≤m : m ≤ m
  ≤-s  : ∀ {n} → m ≤ n → m ≤ suc n
```

A value of type $n \leq m$ mimics the inductive structure of the \mathbb{N} we were storing to represent the distance between n and m . We were able to take this analogy quite far: in a few functions, for instance, we needed to compare these gaps. By mimicking the inductive structure of \mathbb{N} , we were able to directly translate `Ordering` and `compare` on \mathbb{N} :

```

data Ordering : ℕ → ℕ → Set where
  less      : ∀ m k → Ordering m (suc (m + k))
  equal     : ∀ m   → Ordering m m
  greater   : ∀ m k → Ordering (suc (m + k)) m

```

into equivalent functions on \leq :

```

data ≤-Ordering {n : ℕ} : ∀ {i j}
  → (i ≤ n : i ≤ n)
  → (j ≤ n : j ≤ n)
  → Set

where
  ≤-lt : ∀ {i j-1}
    → (i ≤ j-1 : i ≤ j-1)
    → (j ≤ n : suc j-1 ≤ n)
    → ≤-Ordering (≤-trans (≤-s i ≤ j-1) j ≤ n)
               j ≤ n

  ≤-gt : ∀ {i-1 j}
    → (i ≤ n : suc i-1 ≤ n)
    → (j ≤ i-1 : j ≤ i-1)
    → ≤-Ordering i ≤ n
               (≤-trans (≤-s j ≤ i-1) i ≤ n)

  ≤-eq : ∀ {i}
    → (i ≤ n : i ≤ n)
    → ≤-Ordering i ≤ n
               i ≤ n

  ≤-compare : ∀ {i j n}
    → (x : i ≤ n)
    → (y : j ≤ n)
    → ≤-Ordering x y

  ≤-compare m ≤ m m ≤ m = ≤-eq m ≤ m
  ≤-compare m ≤ m (≤-s y) = ≤-gt m ≤ m y
  ≤-compare (≤-s x) m ≤ m = ≤-lt x m ≤ m
  ≤-compare (≤-s x) (≤-s y)
  with ≤-compare x y
  ... | ≤-lt i ≤ j-1 _ = ≤-lt i ≤ j-1 (≤-s y)
  ... | ≤-gt _ j ≤ i-1 = ≤-gt (≤-s x) j ≤ i-1
  ... | ≤-eq _       = ≤-eq (≤-s x)

```

5.2 Unification

After applying the previous optimisations, we might expect an immediate speedup in the solver: unfortunately, this isn't the case. Without some careful adjustments, the optimisations in the previous section can actually *slow down* the solver. In this section, we'll try and explain the problem and how we fixed it, and give general guidelines on how to write Agda code which typechecks quickly.

Up until now, we have focused on the *operations* performed on the polynomial. Remember, though, the reflexive proof process has several steps: only one of them containing the operations ($\llbracket _ \Downarrow \rrbracket$ in Fig. 3). Despite having the most complex implementation, this isn't the most expensive step: surprisingly, the innocuous-looking `refl` takes the bulk of the time! Typechecking this step involves unifying the two normalised expressions, a task which is quite expensive, with counterintuitive performance characteristics.

First, the good news. In the general case, unifying two expressions takes time proportional to the size of those expressions, so our hard-won optimisations do indeed help us.

Unfortunately, though, the “general case” isn't really that general: Agda's unification algorithm has a very important shortcut which we *must* make use of if we want our code to typecheck quickly: *syntactic equality*.

Before the full unification algorithm, Agda runs a quick check to see if the two expressions it's testing for equality are *syntactically* equal. This can make a big difference in unification problems like the following:

$$\text{sum } [1..100] \stackrel{?}{=} \text{sum } [1..100]$$

By noticing that these expressions are syntactically equal, we can avoid actually computing the `sum` function. Taking advantage of that shortcut is key to achieving decent performance. With that in mind, there are two main strategies we'll use to encourage syntactic equality:

5.2.1 Avoid Progress at all Costs

First, we will consider something which may seem inconsequential: the order of arguments to the evaluation functions.

$$\llbracket xs \rrbracket_l \rho = \text{foldr } (\lambda y \, ys \rightarrow \rho * ys + y) \, 0 \, xs$$

$$\llbracket xs \rrbracket_r \rho = \text{foldr } (\lambda y \, ys \rightarrow y + ys * \rho) \, 0 \, xs$$

$\llbracket _ \rrbracket_r$ is the definition we've been working with so far. Some readers might find $\llbracket _ \rrbracket_r$ more natural, however. The reason is that it's more productive: in lazy languages,

the usual convention is that functions which take multiple arguments should scrutinise those arguments from left to right. The $*$ and $+$ functions (on \mathbb{N} , at any rate) follow that convention, meaning that $\llbracket _ \rrbracket_r$ is able to make more progress without a concrete x . Taking the polynomial $x^2 + 2$ as an example:

$$\begin{array}{ll}
 \llbracket 2 :: 0 :: 1 :: _ \rrbracket_l x & \equiv \langle \rangle \\
 x * (x * (x * 0 + 1) + 0) + 2 & \equiv \langle \rangle \\
 x * (x * (x * 0 + 1) + 0) + 2 & \blacksquare \\
 \llbracket 2 :: 0 :: 1 :: _ \rrbracket_r x & \equiv \langle \rangle \\
 2 + (0 + (1 + 0 * x) * x) * x & \equiv \langle \rangle \\
 \text{succ}(\text{succ}((x + 0) * x)) & \blacksquare
 \end{array}$$

In $\llbracket _ \rrbracket_l$, we're blocked pretty much straight away, as x is the first thing we try to scrutinise. In $\llbracket _ \rrbracket_r$, since all of the constants are kept to the left, they're scrutinised first, allowing us to perform much more normalisation before being blocked.

This is exactly what you *don't* want! Since both expressions will be coming out of the same evaluation function, they should have the same structure, meaning that we don't *need* the reduction of outer terms that $\llbracket _ \rrbracket_r$ gives us. We only need to perform normalisation on the coefficients: these are computed during the manipulations of the polynomial, and so may contain unevaluated expressions. If we used $\llbracket _ \rrbracket_r$ as our definition, then the type checker will likely hit an inequality on the *first* term, and as a result we lose all opportunity for syntactic equality. $\llbracket _ \rrbracket_l$, on the other hand, front-loads all of the variables, maintaining syntactic equality for as long as possible.

As well as that, we don't have any control of the structure we get from the ring operators. This means that any reduction, as well as being unnecessary, can destroy the structural similarity between the two expressions, and as a result their syntactic equality.

The (counterintuitive) lesson learned is as follows: to speed up unification, keep things which are likely to be syntactically equal to the left, and *don't* structure your functions to encourage progress. Simply swapping the arguments (as we do above) resulted in a performance improvement of several orders of magnitude.

5.2.2 Avoid Identities

It's a good idea to avoid identities (expressions like $0 + x$ or $1 * x$) in the normalised expression. This will reduce the size of your expression, which is helpful in general, but more importantly it increases the likelihood of finding syntactic equality in the argument to the identity (x in the examples above).

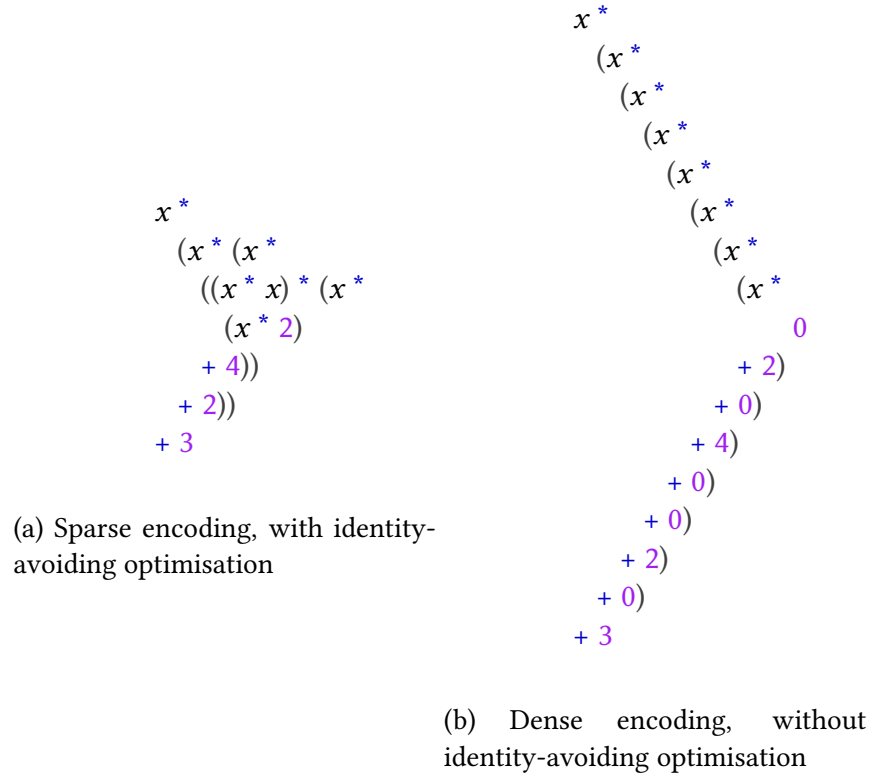


Figure 7: Comparison of the normal forms of equation 4

Our sparse representation helps significantly in this case by entirely removing 0 from the generated expression. Another place we can make improvements is in the base cases for recursive functions. Take exponentiation, for example:

$$\begin{aligned}
 _ \wedge _ &: \text{Carrier} \rightarrow \mathbb{N} \rightarrow \text{Carrier} \\
 x \wedge \text{zero} &= 1 \\
 x \wedge \text{suc } i &= x^* (x \wedge i)
 \end{aligned}$$

We can avoid that 1 in the majority of cases by rewriting the function to have an extra base case:

$ \begin{aligned} _ \wedge _ + 1 &: \text{Carrier} \rightarrow \mathbb{N} \rightarrow \text{Carrier} \\ x \wedge \text{zero} + 1 &= x \\ x \wedge \text{suc } i + 1 &= (x \wedge i + 1)^* x \end{aligned} $	$ \begin{aligned} _ \wedge _ &: \text{Carrier} \rightarrow \mathbb{N} \rightarrow \text{Carrier} \\ x \wedge \text{zero} &= 1 \\ x \wedge \text{suc } i &= x \wedge i + 1 \end{aligned} $
---	---

In the library, we employ this idea extensively, avoiding unnecessary identities as much as we could. This has a significant effect on the size of the resulting normal form, but also ensures that normalisation stops exactly where we want it to, preserving the structure of the expressions as much as is possible. This makes a significant difference to both size and syntactic similarity as can be seen in Fig. 7.

5.3 Benchmarks

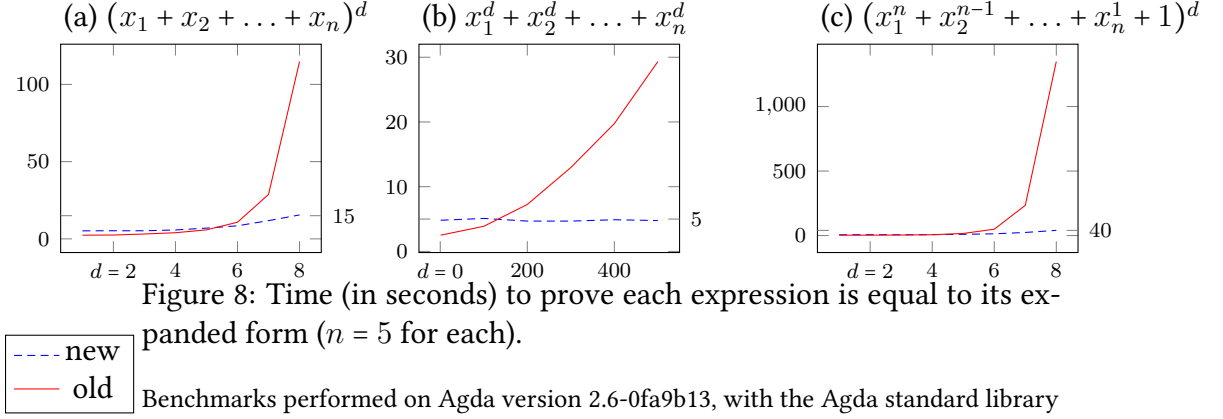


Figure 8: Time (in seconds) to prove each expression is equal to its expanded form ($n = 5$ for each).

Benchmarks performed on Agda version 2.6-0fa9b13, with the Agda standard library at commit-3bd3334a9552490e396f73f96812105a27e5917b, on a 2016 MacBook Pro, with a 2.9 GHz Intel Core i7 and 16 GB of RAM.

The main performance benefit of the new solver is that it reduces the size of the polynomials, both during manipulation and for unification. This reduction is dramatic, even for small expressions: multiplication, for example, generates polynomials with sizes proportional to the product of the sizes of its arguments.

We expect that this will yield 5- to 10-fold speedups in many common use cases. Fig. 8a shows time taken to type check a proof that $(x_1 + x_2 + x_3 + x_4 + x_5)^d$ is equal to its expanded form. The new representation is clearly faster overall, with a factor of 7.5 speedup for $d = 8$. However, the speedup is even more dramatic when the powers of the terms are mixed: Fig. 8c demonstrates this, with a 30-factor speedup at $d = 8$. We feel that this represents a much more common use-case.

These benchmarks cover a broad range of polynomials, with a mix of the three operations provided (addition, multiplication, and exponentiation). In each, the new solver exhibits an order of magnitude speedup at higher powers. We have not been able to find a case where the old solver is significantly faster than the new; however, the old solver does exhibit a small lead (roughly 2-3 seconds, which narrows to about 1 second without reflection) on very simple expressions, possibly caused by the overhead of the new solver's more complex implementation. 1 second is quite small in the context of Agda type checking (the standard library, for instance, takes several minutes to type check), so we feel this slight loss is more than made up for by the gains. Nonetheless, if a user really wants to use the old solver, the other components described here are entirely modular, and can work with any underlying solver which uses the reflexive technique.

6 Verification

The output of the solver is a constructive proof of equivalence: this is *derived* from a generic proof that the operations on the solver are a ring homomorphism from the carrier type. Put another way, for the solver to work properly, we would need to prove that addition (and multiplication, and negation, etc.) on Horner normal forms corresponds with addition on the carrier type.

These proofs are long (about 1000 lines) and complex. Without careful structuring of the proofs, every new optimisation would require a whole new round of proof code, with very little reuse.

To avoid this problem, we took inspiration from [18], and relied heavily on abstraction and folds to improve the reuse in proof code. In particular, we defined many operations as *metamorphisms* [9]. So, instead of defining (say) negation over the polynomial type itself, we will define a metamorphism to express negation, and then call some higher-order function to run that metamorphism over a polynomial.

$$\begin{aligned} \text{Meta} &: \mathbb{N} \rightarrow \text{Set } c \\ \text{Meta } n &= \text{Poly } n \times \text{Coeff } n \star \rightarrow \text{Poly } n \times \text{Coeff } n \star \end{aligned}$$

From here, we can define the *semantics* of a metamorphism. As an example, Fig. 9 shows the semantics of `poly-map`, a simple morphism which behaves something like `map` on lists.

Now, each operation only has to be proven up to the semantics defined above. Crucially, optimisations like the sparse encoding *respect* these semantics, so we only have to change our proof in one place: the definition of `poly-mapR`.

$$\begin{aligned} \text{poly-mapR} &: \forall \{n\} \rho \rho s \\ &\rightarrow ([f] : \text{Poly } n \rightarrow \text{Poly } n) \\ &\rightarrow (f : \text{Carrier} \rightarrow \text{Carrier}) \\ &\rightarrow (\forall x y \rightarrow x * f y \approx f(x * y)) \\ &\rightarrow (\forall x y \rightarrow f(x + y) \approx f x + f y) \\ &\rightarrow (\forall y \rightarrow \llbracket [f] y \rrbracket \rho s \approx f(\llbracket y \rrbracket \rho s)) \\ &\rightarrow (f 0\# \approx 0\#) \\ &\rightarrow \forall xs \\ &\rightarrow \Sigma? \llbracket \text{poly-map } [f] xs \rrbracket (\rho, \rho s) \\ &\quad \approx f(\Sigma \llbracket xs \rrbracket (\rho, \rho s)) \end{aligned}$$

7 Pedagogical Solutions

One of the core aims of this work is to take a step towards making Agda a useful tool for doing mathematics. The rest of this paper has described our efforts to compensate for Agda's disadvantage in this area: namely, a pedantic typechecker. This section will attempt to show the other side of the coin, and demonstrate some of the unique benefits that come from using a programming language to do your proofs.

Figure 9: The Semantics of `poly-map`

Outside of computer scientists and mathematicians, most people’s experience of computer algebra probably amounts to the step-by-step solutions from Wolfram|Alpha [22] or some similar system.

Something so high-level and user-facing hardly seems like it’s in a dependently-typed language’s wheelhouse; on the other hand, the very nature of a proof in Agda is that it has computational content: why not make some of that content an explanation for the equality?

Prior work in this area includes [14]: there, the problem is reformulated reformulates the problem as one of *path-finding*. The left-hand-side and right-hand-side of the equation are vertices in a graph, where the edges are single steps to rewrite an expression to an equivalent form. A^* is used to search.

Unfortunately, this approach has to deal with a huge search space: every vertex will have an edge for almost every one of the ring axioms, and as such a good heuristic is essential. Furthermore, what this should be is not clear: [14] uses a measure of the “simplicity” of an expression.

Notice, however, that paths in undirected graphs form a perfectly reasonable equivalence relation: transitivity is the concatenation of paths, reflexivity is the empty path, and symmetry is *reversing* a path. Equivalence classes, in this analogy, are connected components of the graph.

More practically speaking, we implement these “paths” as lists, where the elements of the list are elementary ring axioms. When we want to display a step-by-step solution, we simply print out each element of the list in turn, interspersed with the states of the expression (the vertices in the graph).

If we stopped there, however, the solver would output incredibly verbose “solutions”: far too verbose to be human-readable. Instead, we must apply a number of path-compression heuristics to cut down on the solution length:

1. First, we remove loops from the graph. Fig 10 shows an example solution without this heuristic applied: it crosses the same point multiple times, creating useless steps in the output. In contrast to using just A^* on its own, the search space is minimal (with only one outward edge for each vertex).
2. Then, we filter out “uninteresting” steps. These are steps which are obvious to a human, like associativity, or evaluation of closed terms. When a step is divided over two sides of an operator, it is deemed “interesting” if either side is interesting.

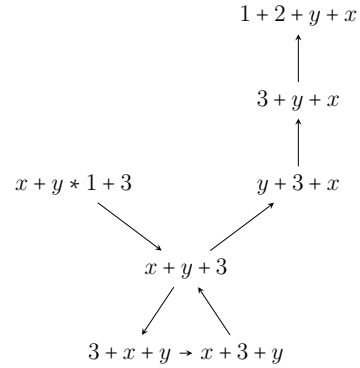


Figure 10: Graph Containing Loops

After applying those heuristics, our solver outputs the explanation in Fig. 11 for the lemma in Fig. 2

$$\begin{aligned}
 & x + y + 3 \\
 &= \{ \text{+-comm}(x, y + 3) \} \\
 & y + 3 + x \\
 &= \{ \text{+-comm}(y, 3) \} \\
 & 3 + y + x
 \end{aligned}$$

Figure 11: Step-by-Step Output From Our Solver

8 Related Work

In independently-typed programming languages, the state-of-the-art solver for polynomial equalities (over commutative rings) was originally presented in [11], and is used in Coq’s `ring` solver. This work improved on the already existing solver [7] in both efficiency and flexibility. In both the old and improved solvers, a reflexive technique is used to automate the construction of the proof obligation (as described in [1]).

Agda [20] is a dependently-typed programming language based on Martin-Löf’s Intuitionistic Type Theory [15]. Its standard library [8] currently contains a ring solver which is similar in flexibility to Coq’s `ring`, but does not support the reflection-based interface, and is less efficient to the one presented here.

In [21], an implementation of an automated solver for the dependently-typed language Idris [2] is described. The solver is implemented with a “correct-by-construction” approach, in contrast to [11]. The solver is defined over *noncommutative* rings, meaning that it is more general (can work with more types) but less powerful (meaning it can prove fewer identities). It provides a reflection-based interface, but internally uses a dense representation.

Reflection and metaprogramming are relatively recent additions to Agda, but form an important part of the interfaces to automated proof procedures. Reflection in dependent types in general is explored in [4], and specific to Agda in [24].

Formalisation of mathematics in general is an ongoing project. [27] tracks how much of “The 100 Greatest Theorems” [13] have so far been formalised (at time of writing, the number stands at 93). DoCon [17] is a notable Agda library in this regard: it contains many tools for basic maths, and implementations of several CAS algorithms. Its implementation is described in [16]. [3] describes the manipulation of polynomials in both Haskell and Agda.

Finally, the study of *pedagogical* CASs which provide step-by-step solutions is explored in [14]. One of the most well-known such system is Wolfram Alpha [28],

which has step-by-step solutions [22].

9 Conclusion

We have presented a ring solver for the programming language and proof assistant Agda. It is faster than the existing solver: common use-cases can expect to see 5-fold increases in speed, and some pathological cases can see improvements of a factor of 30 or more. The interface is easy to use: the solver can be accessed with a single macro call, and requires no knowledge of its inner workings. The solver is flexible: it works out-of-the-box with any commutative ring or semiring defined over a setoid, with no special additions required. Finally, we demonstrated this flexibility by implementing, with no modification to the solver’s code, automatic generation of step-by-step solutions.

We think the future of formalised mathematics looks bright: as well as helping keep us honest, we believe computers and proof assistants can make mathematics easier and more fun. We hope this project gets us closer to that future.

References

- [1] Samuel Boutin. Using reflection to build efficient and certified decision procedures. In Martín Abadi and Takayasu Ito, editors, *Theoretical Aspects of Computer Software*, Lecture Notes in Computer Science, pages 515–529. Springer Berlin Heidelberg, 1997.
- [2] Edwin Brady. Idris, a general-purpose dependently typed programming language: Design and implementation. *Journal of Functional Programming*, 23(05):552–593, September 2013.
- [3] Chen-Mou Cheng, Ruey-Lin Hsu, and Shin-Cheng Mu. Functional Pearl: Folding Polynomials of Polynomials. In *Functional and Logic Programming*, Lecture Notes in Computer Science, pages 68–83. Springer, Cham, May 2018.
- [4] David Raymond Christiansen. *Practical Reflection and Metaprogramming for Dependent Types*. PhD thesis, IT University of Copenhagen, November 2015.
- [5] Alonzo Church. An Unsolvability Problem of Elementary Number Theory. *American Journal of Mathematics*, 58(2):345–363, 1936.
- [6] Alonzo Church. A. M. Turing. On computable numbers, with an application to the Entscheidungsproblem. Proceedings of the London Mathematical Society.

- ciety, 2 s. vol. 42 (1936–1937), pp. 230–265. *The Journal of Symbolic Logic*, 2(1):42–43, March 1937.
- [7] The Coq Development Team. *The Coq Proof Assistant Reference Manual, Version 7.2*. 2002.
 - [8] Nils Anders Danielsson. The Agda standard library, June 2018.
 - [9] Jeremy Gibbons. Metamorphisms: Streaming Representation-Changers. *Science of Computer Programming*, 65(2):108–139, 2007.
 - [10] Kurt Gödel. Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme I. *Monatshefte für Mathematik und Physik*, 38(1):173–198, December 1931.
 - [11] Benjamin Grégoire and Assia Mahboubi. Proving Equalities in a Commutative Ring Done Right in Coq. In *Theorem Proving in Higher Order Logics*, volume 3603 of *Lecture Notes in Computer Science*, pages 98–113, Berlin, Heidelberg, 2005. Springer Berlin Heidelberg.
 - [12] David Hilbert. *Natur und mathematisches Erkennen: Vorlesungen, gehalten 1919–1920 in Göttingen*. Birkhäuser, 1992.
 - [13] Nathan W. Kahl. *The Hundred Greatest Theorems*, 2004.
 - [14] Dmitrij Lioubartsev. *Constructing a Computer Algebra System Capable of Generating Pedagogical Step-by-Step Solutions*. PhD thesis, KTH Royal Institute of Technology, Stockholm, Sweden, 2016.
 - [15] Per Martin-Löf. *Intuitionistic Type Theory*. Padua, June 1980.
 - [16] Sergei D Meshveliani. Dependent Types for an Adequate Programming of Algebra. Technical report, Program Systems Institute of Russian Academy of sciences, Pereslavl-Zalessky, Russia, 2013.
 - [17] Sergei D. Meshveliani. DoCon-A a Provable Algebraic Domain Constructor, April 2018.
 - [18] Shin-Cheng Mu, Hsiang-Shang Ko, and Patrik Jansson. Algebra of programming in Agda: Dependent types for relational program derivation. *Journal of Functional Programming*, 19(5):545–579, September 2009.
 - [19] Ulf Norell. Agda-prelude: Programming library for Agda, August 2018.

- [20] Ulf Norell and James Chapman. Dependently Typed Programming in Agda. Technical report, 2008.
- [21] Franck Slama and Edwin Brady. Automatically Proving Equivalence by Type-Safe Reflection. In Herman Geuvers, Matthew England, Osman Hasan, Florian Rabe, and Olaf Teschke, editors, *Intelligent Computer Mathematics*, volume 10383, pages 40–55. Springer International Publishing, Cham, 2017.
- [22] The Development Team. Step-by-Step Math, December 2009.
- [23] D. van Dalen. The War of the Frogs and the Mice, or the Crisis of the Mathematische Annalen. *The Mathematical Intelligencer*, 12(4):17–31, September 1990.
- [24] P. D. van der Walt. *Reflection in Agda*. Master’s Thesis, Universiteit of Utrecht, October 2012.
- [25] Paul van der Walt and Wouter Swierstra. Engineering Proof by Reflection in Agda. In Ralf Hinze, editor, *Implementation and Application of Functional Languages*, volume 8241, pages 157–173. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013.
- [26] A. N. Whitehead and B. Russell. *Principia Mathematica. Vol. I*. 1910.
- [27] Freek Wiedijk. Formalizing 100 Theorems, October 2018.
- [28] Wolfram Research, Inc. Wolfram|Alpha. Wolfram Research, Inc., 2019.