

Utrecht University

MASTER THESIS

FORMALIZED PROOF OF AUTOMATIC DIFFERENTIATION IN COQ

Student:
Curtis Chin Jen Sem

Supervisors:
Mathijs Vákár
Wouter Swierstra

Department of Information and Computing Science
Last updated: June 13, 2020

Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aenean quis dapibus lorem. Praesent volutpat feugiat erat. Quisque quis hendrerit lectus, et malesuada nisi. Quisque id elementum lectus. Phasellus sit amet ante ornare, hendrerit orci non, consectetur erat. Phasellus pulvinar orci urna. Donec fringilla fringilla ornare. Sed ut tempus arcu, eget ornare ligula.

Ut varius pretium pellentesque. Nam sit amet sapien lobortis nisl faucibus tempor. Curabitur non enim venenatis, euismod elit convallis, auctor sapien. Praesent eget urna sed justo luctus malesuada. In scelerisque metus nibh, ullamcorper efficitur eros fermentum non. Phasellus accumsan congue diam, non fringilla lorem fringilla sit amet. Ut molestie feugiat sagittis. Integer in lobortis justo, et euismod augue. Suspendisse nec euismod lectus, at condimentum magna. Pellentesque eu elementum dui.

Contents

1	Introduction	4
2	Background	5
2.1	Automatic differentiation	5
2.2	Denotational semantics	6
2.3	Coq	8
2.3.1	Language representation	8
2.3.2	Dependently-typed programming in Coq	10
2.4	Logical relations	12
2.5	Related work	13
3	Formalizing Forward-Mode AD	13
3.1	Simply Typed Lambda Calculus	14
3.2	Adding Sums and Primitive Recursion	21
3.3	Arrays	22
4	Optimization	22
4.1	Program Transformations	22
5	Reverse-Mode AD	22
6	Discussion	22
6.1	Problems	22
6.2	Future Work	22
7	Conclusion	22
A	Language Definitions	22
B	Forward-Mode Macro	22
C	Denotations	22

1 Introduction

AI and machine learning research has sparked a lot of new interest in recent times. It has been used in fields such as computer vision, natural language processing, and as opponents in various games such as chess and Go. In machine learning and more specifically neural network research, researchers set up functions referred to as layers between the input and output data and through an algorithm called back propagation, try to optimize the network such that it learns how to solve the problem implied by the data. Back propagation makes heavy use of automatic differentiation, but programming in an environment which allows for automatic differentiation can be limited.

Frameworks such as Tangent¹ or autograd² make use of source code transformations and operator overloading, which can restrict which high-level optimizations one is able to apply to generated code. Support for higher-order derivatives is also limited.

Programming language research has a rich history with many well-known both high- and low-level optimization techniques such as partial evaluation and deforestation. If instead of a framework, we were to have a programming language that is able to facilitate automatic differentiation, we would be able to apply many of these techniques. Through the use of higher-order functions and type systems, we would also get additional benefits such as code-reusability and correctness.

In this thesis, we will aim to formalize an extendable correctness proof of an implementation of automatic differentiation on a simply-typed lambda calculus in the **Coq** proof assistant, opening up further possibilities for formally proving the correctness of more complex language features in the future. Our formalization is based on a recent proof by Huot, Staton, and Vákár [26]. They proved, using a denotational model of diffeological spaces, that their forward mode emulating macro is correct when applied to a simply-typed lambda calculus with products, co-products and inductive types.

With this thesis we will aim for the following goals:

- Formalize the proofs of both the forward mode and continuation-based automatic differentiation algorithms specified by Huot, Staton, and Vákár [26] in **Coq**.
- Prove that well-known compile-time optimizations such as the partial evaluation, are correct with respect to the semantics of automatic differentiation.

¹ <https://github.com/google/tangent>

² <https://github.com/HIPS/autograd>

- Extend the proof with the array types and compile-time optimization rules by Shaikhha, et. al.[20].

Chapter 2 includes a background section explaining many of the topics and techniques used in this thesis. The formalization of forward-mode automatic differentiation is given in Chapter 3, starting from a base simply-typed lambda calculus extended with product types and incrementally adding new types and language constructs. Chapters 4 and 5 give formalizations of optimization avenues using respectively program transformations and continuation-based reverse-mode automatic differentiation.

As a notational convention, we will use specialized notation in the definitions themselves. Coq normally requires that pretty printed notations be defined separately from the definitions they reference. The letter Γ is used for typing contexts while lowercase Greek letters are usually used for types.

2 Background

2.1 Automatic differentiation

One of the principal techniques used in machine learning is back propagation, which calculates the gradient of a function. This gradient is used in the gradient descent algorithm to optimize objective functions by determining the direction of steepest descent[16]. Automatic differentiation has a long and rich history, where its driving motivation is to efficiently calculate the derivatives of functions in a manner that is both correct and fast. There are several different methods of implementing automatic differentiation algorithms, such as source-code transformations or operator overloading. These algorithms usually transform any program which implements some function to one that calculates its derivative. In addition to the standard semantics present in most programming languages, concepts relevant to differentiation such the chain rule are needed.

There are two main variants of automatic differentiation, namely forward mode and reverse mode automatic differentiation. In forward mode automatic differentiation every term in the function trace is annotated with the corresponding derivative of that term. These are also known as the respectively the primal and tangent traces. So every partial derivative of every sub-function is calculated parallel to its counterpart. We will take the function $f(x, y) = x^2 + (x - y)$ as an example. The dependencies between the terms and operations of the function is visible in the computational graph in Figure 1. The corresponding traces are filled in Table 1 for the input val-

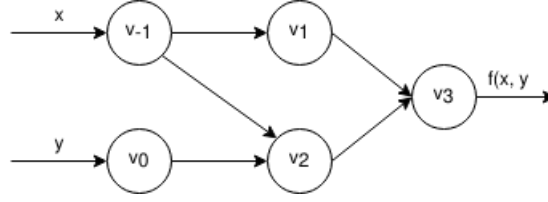


Figure 1: Computational graph of $f(x, y) = x^2 + (x - y)$

ues $x = 2, y = 1$. We can calculate the partial derivative $\frac{\delta f}{\delta x}$ at this point by setting $x' = 1$ and $y' = 0$. In this paper we will prove the correctness of a simple forward mode automatic differentiation algorithm with respect to the semantics of a simply-typed lambda calculus.

Reverse mode automatic differentiation takes a different approach. It works backwards from the output by annotating each intermediate variable v_i with an adjoint $v'_i = \frac{\delta y_i}{\delta v_i}$. To do this, two passes are necessary. Like the forward mode variant, a primal trace is needed to determine the intermediate variables and subfunction dependencies. The second pass calculates the derivatives by working backwards from the output using the adjoints, also called the adjoint trace.

The optimal choice between automatic differentiation variant is heavily dependent on the specific function being differentiated. The number of applications of the forward mode algorithm is dependent on the number of input variables, as it has to be rerun for each partial derivative of the function. On the other hand, as reverse mode AD works backwards, the reverse-pass needs to be redone for each output variable. In machine learning research, reverse mode AD is generally preferred as the objective functions regularly contain a small number of output variables. How one does reverse mode automatic differentiation on a functional language is still an active area of research.

2.2 Denotational semantics

The notion of denotational semantics tries to find underlying mathematical models able to underpin the concepts known in programming languages. The most well-known example is the solution given by Dana Scott and Christopher Strachey[1] for lambda calculi, also called domain theory. To be able to formalize non-termination and partiality, they thought to use concepts such as partial orderings and least fixed points[24]. In this model, programs are interpreted as partial functions, and recursive computations

Primal trace			Tangent trace		
v_{-1}	$= x$	$= 2$	v'_{-1}	$= x'$	$= 1$
v_0	$= y$	$= 1$	v'_0	$= y'$	$= 0$
v_1	$= v_{-1}^2$	$= 4$	v'_1	$= 2 * v_{-1}$	$= 4$
v_2	$= v_{-1} - v_0$	$= 1$	v'_2	$= v'_{-1} - v'_0$	$= 1$
v_3	$= v_1 + v_2$	$= 5$	v'_3	$= v'_1 + v'_2$	$= 5$
f	$= v_3$	$= 5$	f'	$= v'_3$	$= 5$

Table 1: Primal and tangent traces of $f(x, y) = x^2 + (x - y)$

as taking the fixpoint of such functions. Non-termination, on the other hand, is formalized as a value bottom that is lower in the ordering relation than any other element.

Automatic differentiation introduces a complication with picking a proper denotational model as the notion of derivability needs to be established. If the language under consideration were to be restricted to real-typed terms, cartesian spaces would have been sufficient as any well-typed term $x_1 : R, \dots, x_n : R \vdash t : R$ could be interpreted as the corresponding smooth function $\llbracket t \rrbracket : R^n \rightarrow R$. This, however, does not work when function types are added as their denotational equivalent, function spaces, are not supported by cartesian spaces[26]. In the original pen and paper proof of automatic differentiation this thesis is based on by Huot, Staton and Vákár[26], the mathematical models used were diffeological spaces.

For the purpose of this thesis, however, we were able to avoid using diffeological spaces by directly encoding the property of differentiability in the logical relation itself. We were also able to avoid the domain theoretical models such as ω -cpo by excluding language constructs such as recursion and iteration where non-termination and partiality come into play. As a part of its type system, **Coq** contains a set-theoretical model available under the sort `Set`, which is satisfactory as the denotational semantics for our language.

Because we use the real numbers as the ground type in our language, we also needed an encoding of the real numbers in **Coq**. The library for real numbers in **Coq** has improved in recent times from one based on a completely axiomatic definition to one involving Cauchy sequences³. For the purposes of this thesis, however, we also needed differentiability as the denotational result of applying the macro operation. Instead of encoding this by hand, we opted for the more comprehensive library `Coquelicot`[17], which contains many useful definitions for differentiating functions.

³ <https://coq.inria.fr/library/Coq.Reals.ConstructiveCauchyReals.html>

$$\begin{array}{c}
\frac{elem\ n\ \Gamma = \tau}{\Gamma \vdash var\ n : \tau} \text{TVAR} \qquad \frac{(\sigma, \Gamma) \vdash t : \tau}{\Gamma \vdash t : \sigma \rightarrow \tau} \text{TAbs} \\
\\
\frac{\Gamma \vdash t1 : \sigma \rightarrow \tau \quad \Gamma \vdash t2 : \sigma}{\Gamma \vdash t1\ t2 : \tau} \text{TApp}
\end{array}$$

Figure 2: Type-inference rules for a simply-typed lambda calculus using De-Brujin indices

2.3 Coq

Coq is a proof assistant based on the calculus of constructions type theory created by Thierry Coquand and Gérard Huet[2]. In the past 30 years since it has been released, research has contributed to extending the proof assistant with additional features such as inductive and co-inductive data types[3], dependent pattern matching[12] and advanced modular constructions for organizing large mathematical proofs[11][15].

The core of this type theory is based on constructive logic and so many of the laws known in classical logic are not provable. An example includes the law of the excluded middle, $\forall A, A \vee \neg A$. In some cases they can, however, be safely added to **Coq** without making its logic inconsistent. These are readily available in the standard library. Due to its usefulness in proving propositions over functions, we will make use of the functional extensionality axiom in **Coq**, which states that functions are equal if they are extensionally equivalent.

2.3.1 Language representation

When defining a simply-typed lambda calculus, there are two main possibilities[23]. The arguably simpler variant, known as an extrinsic representation, is traditionally the one introduced to new students learning **Coq**. In the extrinsic representation, the terms themselves are untyped and typing judgments are defined separately as relations between the types and terms. A basic example of working with this is given in Software Foundations[18]. This, however, required many additional lemmas and machinery to be proved to be able to work with both substitutions and contexts as these are defined separate from the terms. As an example, the preservation property which states

that reduction does not change the type of a term, needs to be proven explicitly. The other approach, also called an intrinsic representation, makes use of just a single well-typed definition. Ill-typed terms are made impossible by the type checker. This representation, while beneficial in the proof load, however complicates much of the normal machinery involved in programming language theory. One example is how one would define operations such as substitutions or weakening.

But even when choosing an intrinsic representation, the problem of variable binding persists. Much meta-theoretical research has been done on possible approaches to this problem each with their own advantages and disadvantages. The POPLmark challenge gives a comprehensive overview of each of the possibilities in various proof assistants[7]. An example of an approach is the nominal representation where every variable is named. While this does follow the standard format used in regular mathematics, problems such as alpha-conversion and capture-avoidance appears.

```

Inductive ty : Type :=
  | unit : ty
  |  $\Rightarrow$  : ty  $\rightarrow$  ty  $\rightarrow$  ty.

Inductive tm : Type :=
  | var : string  $\rightarrow$  tm
  | abs : string  $\rightarrow$  ty  $\rightarrow$  tm  $\rightarrow$  tm
  | app : tm  $\rightarrow$  tm  $\rightarrow$  tm.

```

Code snippet 1: Simply typed λ -calculus using an extrinsic nominal representation.

The approach used in the rest of this thesis is an extension of the De-Brujin representation which numbers variables relative to the binding lambda term. In this representation the variables are referred to as well-typed De-Brujin indices. A significant benefit of this representation is that the problems of capture avoidance and alpha equivalence are avoided. As an alternative, instead of using numbers to represent the distance, indices within the typing context can be used to ensure that a variable is always well-typed and well-scoped. While the idea of using type indexed terms has been both described and used by many authors[5][6][8], the specific formulation used in this thesis using separate substitutions and rename operations was fleshed out in Coq by Nick Benton, et. al.[13], and was also used as one of the examples

in the second POPLmark challenge which deals with logical relations[19]. While this does avoid the problems present in the nominal representation, it unfortunately does have some problems of its own. Variable substitutions have to be defined using two separate renaming and substitution operations. Renaming is formulated as extending the typing context of variables, while substitution actually swaps the variables for terms. Due to using indices from the context as variables, some lifting boilerplate is also needed to manipulate contexts.

```

Inductive  $\tau \in \Gamma$  : Type :=
  | Top :  $\forall \Gamma \tau, \tau \in (\tau :: \Gamma)$ 
  | Pop :  $\forall \Gamma \tau \sigma, \tau \in \Gamma \rightarrow \tau \in (\sigma :: \Gamma)$ .

Inductive  $\text{tm } \Gamma \tau$  : Type :=
  | var :  $\forall \Gamma \tau, \tau \in \Gamma \rightarrow \text{tm } \Gamma \tau$ 
  | abs :  $\forall \Gamma \tau \sigma, \text{tm } (\sigma :: \Gamma) \tau \rightarrow \text{tm } \Gamma (\sigma \Rightarrow \tau)$ 
  | app :  $\forall \Gamma \tau \sigma, \text{tm } \Gamma (\sigma \Rightarrow \tau) \rightarrow \text{tm } \Gamma \sigma \rightarrow \text{tm } \Gamma \tau$ .

```

Code snippet 2: Basis of a simply-typed λ -calculus using a strongly typed intrinsic formulation.

2.3.2 Dependently-typed programming in Coq

In **Coq**, one can normally write function definitions using either case-analysis as is done in other functional languages, or using **Coq**'s tactics language. Using the standard case-analysis functionality can cause the code to be complicated and verbose, even more so when proof terms are present in the function signature. This has been caused by the previously poor support in Coq for dependent pattern matching. Using the return keyword, one is able to vary the result type of a match expression. But due to requirement Coq used to have that case expressions be syntactically total, this could be very annoying to work with. One other possibility would be to write the function as a relation between its input and output. This also has its limitations as you then lose computability as Coq treats these definitions opaquely. In this case the standard `simpl` tactic which invokes **Coq**'s reduction mechanism is not able to reduce instances of the term. This often requires the user to write many more proofs to be able to work with the definitions.

```

Inductive ilist : Type → nat → Type :=
  | nil : ∀ A, ilist A 0
  | cons : ∀ A n, A → ilist A n → ilist A (S n)

```

```

Definition hd' {A} n (ls : ilist A n) :=
  match ls in (ilist A n) return
    (match n with
      | 0 => unit
      | S _ => A end) with
  | nil => tt
  | cons h _ => h
end.

```

```

Definition hd {A} n (ls : ilist A (S n)) : A := hd' n ls.

```

Code snippet 3: Definition of a length indexed list and hd using the return keyword, adapted from Certified Programming with Dependent Types[14].

As an example, we will work through defining a length indexed list and a corresponding head function limited to lists of length at least one in Snippet 3. Using the **Coq** keyword **return**, it is possible to let the return type of a match expressions depend on the result of one of the type arguments. This makes it possible to define an auxiliary function which, while total on the length of the list, has an incorrect return type. It namely returns the type `unit` if the input list had the length zero. We can then use this auxiliary function in the actual head function by specifying that the list has length at least one. It should be noted that more recent versions of Coq do not require that case expressions be syntactically total, so specifying that the input list has a length of at least zero is enough to eliminate the requirement for the zero-case.

Mathieu Sozeau introduces an extension to **Coq** via a new keyword **Program** which allows the use of case-analysis in more complex definitions[27][10]. To be more specific, it allows definitions to be specified separately from their accompanying proofs, possibly filling them in automatically if possible. While this does improve on the previous situation, using the definitions in proofs can often be unwieldy due to the amount of boilerplate introduced. This makes debugging error messages even harder than it already is in a proof assistant. This approach was used by Benton in

```

Equations hd {A n} (ls : ilist A n) (pf : n <> 0) : A :=
hd nil pf with pf eq_refl := {};
hd (cons h n) _ := h.

```

Code snippet 4: Definition of hd using Equations

his formulation of strongly typed terms.

Sozeau further improves on this introducing a method for user-friendlier dependently-typed pattern matching in **Coq** in the form of the Equations library[12][22]. This introduces **Agda**-like dependent pattern matching with with-clauses. It does this by using a notion called coverings, where a covering is a set of equations such that the pattern matchings of the type signature are exhaustive. There are two main ways to integrate this in a dependently typed environment, externally where it is integrated as high-level constructs in the pattern matching core as **Agda** does it, or internally by using the existing type theory and finding witnesses of the covering to prove the definition correct, which is the approach used by Sozeau. Due to the intrinsic typeful representation this paper uses, much of this was invaluable when defining the substitution operators as the regular type checker in Coq often had difficulty unifying dependently typed terms in certain cases.

2.4 Logical relations

Logical relations is a technique often employed when proving programming language properties of statically typed languages[21]. There are two main ways they are used, namely as unary and binary relations. Unary logical relations, also known as logical predicates, are predicates over single terms and are typically used to prove language characteristics such as type safety or strong normalization. Binary logical relations on the other hand are used to prove program equivalences, usually in the context of denotational semantics as we will do. There have been many variations on the versatile technique from syntactic step-indexed relations which have been used to solve recursive types[9], to open relations which enable working with terms of non-ground type[25][26]. Logical relations in essence are relations between terms defined by induction on their types. A logical relations proof consists of 2 main steps. The first states the terms for which the property is expected to hold are in the relation, while the second states that the property of interest follows from the relation. The second step is easier to prove as it usually

```

Equations SN {Γ} τ (t : tm Γ τ): Prop :=
  SN unit t := halts t;
  SN (τ ⇒ σ) t := halts t ∧
    (∀ (s : tm Γ τ), SN τ s → SN σ (app Γ σ τ t s));

```

Code snippet 5: Example of a logical predicate used in a strong normalizations proof in the intrinsic strongly-typed formulation

follows from the definition of the relation. The first on the other hand, will often require proving a generalized variant, called the fundamental property of the logical relation. In most cases this requires that the relation is correct with respect to applying substitutions.

A well-known logical relations proof is the proof of strong normalization of well-typed terms, which states that all terms eventually terminate. An example of a logical relation used in such a proof using the intrinsic strongly-typed formulation is given in Snippet 5. Noteworthy is the case for function types, which indicates that an application should maintain the strongly normalization property. If one were to attempt the proof of strong normalization without using logical relations, they would get stuck in the cases dealing with function types. More specifically when reducing an application, the induction hypothesis is not strong enough to prove that substituting the argument into the body of the abstraction also results in a terminating term. The proof given in the paper this thesis is based on, is a logical relations proof on the denotation semantics using diffeological spaces as its domains[26]. A similar, independent proof of correctness was given by Barthe, et. al.[25] using a syntactic relation on the operational semantics.

2.5 Related work

3 Formalizing Forward-Mode AD

The formalization of the forward-mode automatic differentiation macro will be explained in the following sections. The formal proof will start from a base simply-typed lambda calculus enriched with product types and incrementally add both sum and array types. Also included in the final language is an implementation of primitive recursion using integer types. Many of the theorems and lemmas will stay consistent between sections as

the overarching correctness statement does not change. Later section only add their respective cases to the corresponding proofs of the lemmas.

3.1 Simply Typed Lambda Calculus

As mentioned in background section 2.3.1, we will make use of De-Brujin indices in a intrinsic representation to formulate our language. Both addition and multiplication are included as operations on the real numbers. Our base language consists of the simply-typed lambda calculus with product types and real numbers, hereby referenced as $\Lambda_{\delta}^{\times, \rightarrow, R}$, where δ consists of the addition and multiplication operators.

Both the language constructs and the typing rules for this language are standard for a simply typed lambda calculus, as shown in 3. As expected we include variables, applications and abstractions in the `var`, `app` and `abs` constructors. Product types are added to the language in the form of binary projections, `first` and `second` to represent respectively the left and right projections of tuple terms. For real numbers, `rval` is used to introduce real numbered constants and `add` and `mul` will be used to respectively encode addition and multiplication.

These can be translated into Coq definitions in a reasonably straightforward manner, with each case keeping track of both how the typing context and types change. In the `var` case we need some way to determine what type the variable is referencing. Like many others previously [13][4], instead of using indices into the list accompanied with a proof that the index does not exceed the length of the list, we make use of an inductively defined type evidence to type our variables as shown in 2. The cases for `app` and `abs` are as expected, where variables in the body of abstraction are able to reference their respective arguments.

Notably, in the original proof by Huot, Staton, and Vákár [26], they make use of n-ary products accompanied with pattern matching expressions. We opted to implement binary projection products, as they are conceptually simpler while still retaining much of the same functionality expected of product types.

$$\begin{array}{c}
\frac{elem\ n\ \Gamma = \tau}{\Gamma \vdash var\ n : \tau} \text{TVAR} \qquad \frac{(\sigma, \Gamma) \vdash t : \tau}{\Gamma \vdash abs\ t : \sigma \rightarrow \tau} \text{TAbs} \\
\\
\frac{\Gamma \vdash t1 : \sigma \rightarrow \tau \quad \Gamma \vdash t2 : \sigma}{\Gamma \vdash app\ t1\ t2 : \tau} \text{TApp} \\
\\
\frac{\Gamma \vdash t1 : \tau \quad \Gamma \vdash t2 : \sigma}{\Gamma \vdash tuple\ t1\ t2 : \tau \times \sigma} \text{TTuple} \\
\\
\frac{\Gamma \vdash t : \tau \times \sigma}{\Gamma \vdash first\ t : \tau} \text{TFST} \qquad \frac{\Gamma \vdash t : \tau \times \sigma}{\Gamma \vdash second\ t : \sigma} \text{TSND} \\
\\
\frac{r \in \mathcal{R}}{\Gamma \vdash rval\ r : R} \text{TRVAL} \\
\\
\frac{\Gamma \vdash r1 : R \quad \Gamma \vdash r2 : R}{\Gamma \vdash add\ r1\ r2 : R} \text{TADD} \qquad \frac{\Gamma \vdash r1 : R \quad \Gamma \vdash r2 : R}{\Gamma \vdash mul\ r1\ r2 : R} \text{TMULL}
\end{array}$$

Figure 3: Type-inference rules for the base simply-typed lambda calculus

```

Inductive tm (Γ : Ctx) : ty → Type :=
  ...
  (* Binary projection products *)
  | tuple : forall τ σ,
    tm Γ τ →
    tm Γ σ →
    tm Γ (τ × σ)
  | first : forall τ σ, tm Γ (τ × σ) → tm Γ τ
  | second : forall τ σ, tm Γ (τ × σ) → tm Γ σ
  ...

```

We use the same inductively defined macro on types and terms used by many previous authors to implement the forward-mode automatic differentiation macro[26][25][20]. The forward-mode macro, $\vec{\mathcal{D}}$, keeps track of both primal and tangent traces using tuples as respectively its first and second elements. In most cases, the macro simply preserves the structure of the language. The cases for real numbers such as addition and multiplication are the exception. Here, the element encoding the tangent trace needs to contain the proper syntactic translation of the derivative of the operation.

$$\begin{aligned}
 \vec{\mathcal{D}}(\mathbb{R}) &= \mathbb{R} \times \mathbb{R} & \vec{\mathcal{D}}(\text{rval } n) &= \text{tuple } (\text{rval } n) (\text{rval } 0) \\
 \vec{\mathcal{D}}(\tau \times \sigma) &= \vec{\mathcal{D}}(\tau) \times \vec{\mathcal{D}}(\sigma) & \vec{\mathcal{D}}(\text{add } n \ m) &= \text{tuple } (\text{add } n \ m) (\text{add } n' \ m') \\
 \vec{\mathcal{D}}(\tau \rightarrow \sigma) &= \vec{\mathcal{D}}(\tau) \rightarrow \vec{\mathcal{D}}(\sigma) & \vec{\mathcal{D}}(\text{mul } n \ m) &= \text{tuple } (\text{mul } n \ m) \\
 & & & (\text{add } (\text{mul } n' \ m) (\text{mul } m' \ n)))
 \end{aligned}$$

Figure 4: Macro on base simply-typed lambda calculus

This is implemented by destructing the recursive calls to \mathcal{D} to access the syntactic counterparts of the primal and tangent denotations. Note that applying the macro to variables do nothing exceptional as the macro is also applied to the typing context, so variables implicitly already reference macro-applied terms.

Due to restricting our language to be total and excluding constructs related to partiality such as general recursion and iteration, it suffices to give our language a set-theoretic denotational semantics. Like the type evidences, well-typed terms will denote to functions $\llbracket \Gamma \rrbracket \rightarrow \llbracket \tau \rrbracket$.

For any type τ , simply swap out the syntactic type to its corresponding

$$\begin{array}{ll}
\llbracket \mathbf{R} \rrbracket = & \llbracket \text{var } v \rrbracket = \lambda x. \text{lookup } \llbracket v \rrbracket x \\
\llbracket \tau \times \sigma \rrbracket = \llbracket \tau \rrbracket * \llbracket \sigma \rrbracket & \llbracket \text{app } t_1 t_2 \rrbracket = \lambda x. (\llbracket t_1 \rrbracket(x))(\llbracket t_2 \rrbracket(x)) \\
\llbracket \tau \rightarrow \sigma \rrbracket = \llbracket \tau \rrbracket \rightarrow \llbracket \sigma \rrbracket & \llbracket \text{abs } t \rrbracket = \lambda x y. \llbracket t \rrbracket(y :: x) \\
\llbracket \text{Top} \rrbracket = \text{hd} & \llbracket \text{add } t_1 t_2 \rrbracket = \lambda x. \llbracket t_1 \rrbracket(x) + \llbracket t_2 \rrbracket(x) \\
\llbracket \text{Pop } v \rrbracket = \text{tl} \circ \llbracket v \rrbracket & \llbracket \text{mul } t_1 t_2 \rrbracket = \lambda x. \llbracket t_1 \rrbracket(x) * \llbracket t_2 \rrbracket(x) \\
& \llbracket \text{tuple } t_1 t_2 \rrbracket = \lambda x. (\llbracket t_1 \rrbracket(x), \llbracket t_2 \rrbracket(x)) \\
& \llbracket \text{first } t \rrbracket = \lambda x. \text{let } (x, y) = \llbracket t \rrbracket(x) \text{ in } x \\
& \llbracket \text{second } t \rrbracket = \lambda x. \text{let } (x, y) = \llbracket t \rrbracket(x) \text{ in } y
\end{array}$$

Figure 5: Denotations of the base simply-typed lambda calculus

Coq variant in **Type**. So the syntactic **R** will denote to the **type** in **Coq**. Denoting the terms in our language now corresponds to finding the appropriate inhabitants in the denotated types. We denote our typing contexts Γ , lists of types, as heterogeneous lists containing their corresponding denotations. The specific implementation of heterogeneous lists used, correspond to the one given by Adam Chlipala[14].

When working through giving the constructs in our language the proper denotations, most of the cases are straightforward. Notable is the case for variables, where we made use of an inductively defined evidence to type our terms. As denotations, these evidences will correspond to lookups into our heterogeneous lists to their appropriate types.

As mentioned by Barthe, et. al.[25], this small calculus, $\Lambda_\delta^{\times, \rightarrow, \mathbb{R}}$, accompanied with the arguably simple set-theoretic denotational semantics is expressive enough to encode the higher-order polynomials containing the addition and multiplication operators.

Example 1 (Square). *abs (mul (var Top) (var Top)) denotes to the square function $\lambda x. x * x$.*

Proof. This follows from the definition of our denotation functions.

$$\begin{aligned}
& \llbracket \text{abs (mul (var Top) (var Top))} \rrbracket [] \\
& \equiv \lambda x. \llbracket \text{mul (var Top) (var Top)} \rrbracket [x] \\
& \equiv \lambda x. \llbracket \text{var Top} \rrbracket [x] * \llbracket \text{var Top} \rrbracket [x] \\
& \equiv \lambda x. x * x
\end{aligned}$$

□

As we work with denotations, smooth functions $f : R^n \rightarrow R$ can be interpreted as the denotations of a corresponding syntactic term $x_1, \dots, x_n \vdash t : R$. Intuitively, the free variables in the term t denote the usages of the parameters of the function and as such are restricted to terms of type R , same as each of the arguments in the function f . Note that while both the arguments and result type of t are restricted to R , t itself can consist of higher order types.

Although Barthe, et. al.[25] gave a syntactic proof of correctness of the macro, our proof follows the more denotational style of proof given by Huot, Staton and Vákár[26]. Correctness of the forward-mode macro consists of the assertion that the denotation of any macro-applied term of type $x_1 : R, \dots, x_n : R \vdash t : R$ will result in a pair of both the denotation of the original term and the derivative of that denotation.

Likewise, the proof of correctness will follow a logical relations argument, first generalizing the statement for both the typing contexts and the result type. The relation will ensure that both the derivability property and the derivatives are preserved over higher-order types. We define the logical relation as a type-indexed relation between denotations of both terms and their macro-applied variants, so for any type τ , S_τ is the relation between functions $R \rightarrow \llbracket \tau \rrbracket$ and $R \rightarrow \llbracket \vec{\mathcal{D}}(\tau) \rrbracket$.

$$S_\tau(f, g) = \begin{cases} \text{smooth } f \wedge g = \lambda x. (f(x), \frac{\partial f}{\partial x}(x)) & : \tau = R \\ \exists f_1, f_2, g_1, g_2, & : \tau = \sigma \times \rho \\ \quad S_\sigma(f_1, f_2), S_\rho(g_1, g_2). \\ \quad f = \lambda x. (f_1(x), g_1(x)) \wedge \\ \quad g = \lambda x. (f_2(x), g_2(x)) & \\ \forall f_1, f_2. & : \tau = \sigma \rightarrow \rho \\ \quad S_\sigma(f_1, f_2) \Rightarrow \\ \quad S_\rho(\lambda x. f(x)(f_1(x)), \lambda x. f(x)(f_2(x))) & \end{cases} \quad (1)$$

The next step involves proving that syntactically well-typed terms are semantically correct. In other words, the relation needs to be proven valid for any term $x_1 : R, \dots, x_n : R \vdash t : \tau$ and argument function $f : R \rightarrow R^n$ such that $S_\tau(\llbracket t \rrbracket \circ f, \llbracket \vec{\mathcal{D}}(t) \rrbracket \circ \vec{\mathcal{D}}_n \circ f)$.

To properly instantiate the arguments to the denotation of the macro-applied term, an auxiliary function is needed that pairs each constant with their derivative 0. So it transforms $f : R \rightarrow \llbracket R^n \rrbracket$ into $\vec{\mathcal{D}}_n(f, x) : R \rightarrow \llbracket \vec{\mathcal{D}}(R^n) \rrbracket$. The full type signature of the function becomes $\vec{\mathcal{D}}_n : (R \rightarrow \llbracket R^n \rrbracket) \rightarrow R \rightarrow \llbracket \vec{\mathcal{D}}(R^n) \rrbracket$, which essentially accompanies each argument

supplied by f with its accompanying derivative.

$$\vec{D}_n(f, x) = \begin{cases} f(x) & : n = 0 \\ ((hd \circ f)(x)), \frac{\partial(hd \circ f)}{\partial x}(x) :: \vec{D}_{n'}(tl \circ f, x) & : n = S(n') \end{cases} \quad (2)$$

Proving this statement directly by induction on the typing derivation, however, does not work. As expected in a logical relations proof, the inductive issue lies in both the case for applications and abstractions. To make this work, the correctness statement needs to be generalized to arbitrary contexts and implicitly, substitutions.

If this were a syntactic proof, one would need to show that relation is preserved when applying substitutions consisting of arbitrary terms, possibly containing higher-order constructs. In this style of proof, the same concept needs to be incorporated in the argument function f , which intuitively speaking, supplies the terms referenced by variables through the typing context.

To prove this statement, it first needs to be generalized to arbitrary substitutions. The key in formulating these denotationally lies in what was previously the argument function $f : R \rightarrow R^n$. Previously the function was used to indicate the open variables or function arguments. If generalized to $\Gamma = x_1 : \tau_1, \dots, x_n : \tau_n$, this same function could be seen as a function which supplies terms for each open variable x_1, \dots, x_n with their appropriate types. So the argument function now becomes the pair of functions $s : R \rightarrow \llbracket \Gamma \rrbracket$ and $s_D : R \rightarrow \llbracket \vec{D}(\Gamma) \rrbracket$. Note that the functions s and s_D are built out of the denotations of terms such that these same denotations follow the logical relation (3) for our language. We phrase this requirement as a definition.

Definition 1. (Instantiation) Substitutional functions $s : R \rightarrow \llbracket \Gamma \rrbracket$ and $s_D : R \rightarrow \llbracket \vec{D}(\Gamma) \rrbracket$ are inductively instantiated such that they follow

$$inst_\Gamma(f, g) = \begin{cases} f = const(\Box) \wedge g = const(\Box) & : \Gamma = \Box \\ \forall f_1, f_2, g_1, g_2. & : \Gamma = (\tau :: \Gamma') \\ inst_{\Gamma'}(f_1, g_1) \wedge S_\tau(f_2, g_2) & \\ \implies f = \lambda x. (f_2(x) :: f_1(x)) \wedge & \\ g = \lambda x. (g_2(x) :: g_1(x)) & \end{cases} \quad (3)$$

Using this notion of substitution instantiations we can now formulate our substitution lemma.

Lemma 1 (Substitution). For any well-typed term $\Gamma \vdash t : \tau$, and instantiation functions $s : R \rightarrow \llbracket \Gamma \rrbracket$ and $s_D : R \rightarrow \llbracket \vec{D}(\Gamma) \rrbracket$ such that they follow $inst_\Gamma(s, s_D)$, then $S_\tau(\llbracket t \rrbracket \circ s, \llbracket \vec{D}(t) \rrbracket \circ s_D)$.

The proof proceeds by induction on the typing derivation of the well-typed term t . The majority of cases follow from the induction hypothesis. The case for `var` follows from *inst* which ensures that any term referenced is semantically well-typed with respect to the relation. Proving the cases used to encode the operators on reals such as `add` and `mul` involve proving both smoothness and giving the witness of the derivative.

We can derive the fundamental property of the base logical relation directly from the substitution lemma. This involves proving the prerequisite *inst* we used previously. Note that the correctness of both the macro and the fundamental property is dependent on the requirement that the denotations supplied by the argument function are smooth.

Lemma 2 (Fundamental property). *For any term $x_1 : R, \dots, x_n : R \vdash t : R$, $\vec{\mathbb{D}}(t)$ gives the dual number representation of $\llbracket t \rrbracket$, such that for any argument function $f : R \rightarrow R^n$, then $S_\tau(\llbracket t \rrbracket \circ f, \vec{\mathbb{D}}(t) \circ \vec{\mathcal{D}}_n \circ f)$.*

This is proven by the substitution lemma. The remaining goal of *inst_{Rⁿ}* is proven by induction on n , the number of arguments. With the case where $n = 0$, the goal is trivial due to the argument function f being extensionally equal to *const* \square , which directly corresponds to *inst_{\square}*. Similarly the induction step is proven by both the induction hypothesis and the assumption that the denotations of the arguments supplied are smooth.

Theorem 1 (Macro correctness). *For any term $x_1 : R, \dots, x_n : R \vdash t : R$, $\vec{\mathbb{D}}(t)$ gives the dual number representation of $\llbracket t \rrbracket$, such that for any argument function $f : R \rightarrow R^n$, then $\vec{\mathbb{D}}(t) \circ \vec{\mathcal{D}}_n \circ f = \lambda x. (\llbracket t \rrbracket \circ f, \partial(\llbracket t \rrbracket \circ f) / \partial x)$.*

Proof. This is proven by showing that the goal follows from the logical relation which itself implied by the fundamental property (2) for well-typed terms.

$$\begin{aligned}
& \vec{\mathbb{D}}(t) \circ \vec{\mathcal{D}}_n \circ f = \lambda x. (\llbracket t \rrbracket \circ f, \partial(\llbracket t \rrbracket \circ f) / \partial x) \\
& \Vdash \text{(By definition of } S_R \text{ with } f := \llbracket t \rrbracket \circ f \text{ and } g := \vec{\mathbb{D}}(t) \circ \vec{\mathcal{D}}_n \circ f) \\
& S_R(\llbracket t \rrbracket \circ f, \vec{\mathbb{D}}(t) \circ \vec{\mathcal{D}}_n \circ f) \\
& \Vdash \text{(Fundamental property (2))}
\end{aligned}$$

□

3.2 Adding Sums and Primitive Recursion

Now that correctness has been verified for the base simply-typed lambda calculus, the next goal will be to add in both sum and integer types. In the interest of testing the flexibility of both the representation and the proofing technique, integer types and primitive recursion were also added. The inference rules for the new language constructs added for sum and number types are given in figure 6

$$\begin{array}{c}
 \frac{\Gamma \vdash e : \tau + \sigma \quad \Gamma \vdash t1 : \tau \rightarrow \rho \quad \Gamma \vdash t2 : \sigma \rightarrow \rho}{\Gamma \vdash \text{case } e \text{ } t1 \text{ } t2 : \rho} \text{TCASE} \\
 \\
 \frac{\Gamma \vdash t : \tau}{\Gamma \vdash \text{inl } t : \tau + \sigma} \text{TIINL} \qquad \frac{\Gamma \vdash t : \sigma}{\Gamma \vdash \text{inr } t : \tau + \sigma} \text{TIINR} \\
 \\
 \frac{n \in \mathcal{N}}{\Gamma \vdash \text{nval } n : \mathbb{N}} \text{TNVAL} \qquad \frac{\Gamma \vdash t : \mathbb{N}}{\Gamma \vdash \text{nsucc } t : \mathbb{N}} \text{TNSUCC} \\
 \\
 \frac{\Gamma \vdash f : \tau \rightarrow \tau \quad \Gamma \vdash n : \mathbb{N} \quad \Gamma \vdash t : \tau}{\Gamma \vdash \text{nrec } f \text{ } n \text{ } t : \tau} \text{TPRIM}
 \end{array}$$

Figure 6: Type-inference rules for language constructs for sum types and primitive recursion

3.3 Arrays

4 Optimization

4.1 Program Transformations

5 Reverse-Mode AD

6 Discussion

6.1 Problems

6.2 Future Work

7 Conclusion

A Language Definitions

B Forward-Mode Macro

C Denotations

References

- [1] D. Scott, “Outline of a mathematical theory of computation”, *Kiberneticheskij Sbornik. Novaya Seriya*, vol. 14, Jan. 1977 (cit. on p. 6).
- [2] T. Coquand and G. Huet, “The calculus of constructions”, *Information and Computation*, vol. 76, no. 2-3, pp. 95–120, Feb. 1988. doi: 10.1016/0890-5401(88)90005-3. [Online]. Available: [https://doi.org/10.1016/0890-5401\(88\)90005-3](https://doi.org/10.1016/0890-5401(88)90005-3) (cit. on p. 8).
- [3] T. Coquand and C. Paulin, “Inductively defined types”, in *COLOG-88*, Springer Berlin Heidelberg, 1990, pp. 50–66. doi: 10.1007/3-540-52335-9_47. [Online]. Available: https://doi.org/10.1007/3-540-52335-9_47 (cit. on p. 8).
- [4] T. Coquand and P. Dybjer, “Inductive definitions and type theory an introduction (preliminary version)”, in *Foundation of Software Technology and Theoretical Computer Science*, P. S. Thiagarajan, Ed., Berlin, Heidelberg: Springer Berlin Heidelberg, 1994, pp. 60–76, ISBN: 978-3-540-49054-8 (cit. on p. 14).

- [5] T. Altenkirch and B. Reus, “Monadic presentations of lambda terms using generalized inductive types”, in *Proceedings of the 13th International Workshop and 8th Annual Conference of the EACSL on Computer Science Logic*, ser. CSL ’99, Berlin, Heidelberg: Springer-Verlag, 1999, pp. 453–468, ISBN: 3540665366 (cit. on p. 9).
- [6] C. McBride and J. McKinna, “The view from the left”, *J. Funct. Program.*, vol. 14, no. 1, pp. 69–111, Jan. 2004, ISSN: 0956-7968. DOI: 10.1017/S0956796803004829. [Online]. Available: <https://doi.org/10.1017/S0956796803004829> (cit. on p. 9).
- [7] B. E. Aydemir, A. Bohannon, M. Fairbairn, J. N. Foster, B. C. Pierce, P. Sewell, D. Vytiniotis, G. Washburn, S. Weirich, and S. Zdancewic, “Mechanized metatheory for the masses: The PoplMark challenge”, in *Lecture Notes in Computer Science*, Springer Berlin Heidelberg, 2005, pp. 50–65. DOI: 10.1007/11541868_4. [Online]. Available: https://doi.org/10.1007/11541868_4 (cit. on p. 9).
- [8] R. Adams, “Formalized metatheory with terms represented by an indexed family of types”, in *Types for Proofs and Programs*, J.-C. Filliâtre, C. Paulin-Mohring, and B. Werner, Eds., Berlin, Heidelberg: Springer Berlin Heidelberg, 2006, pp. 1–16, ISBN: 978-3-540-31429-5 (cit. on p. 9).
- [9] A. Ahmed, “Step-indexed syntactic logical relations for recursive and quantified types”, in *Programming Languages and Systems*, P. Sestoft, Ed., Berlin, Heidelberg: Springer Berlin Heidelberg, 2006, pp. 69–83, ISBN: 978-3-540-33096-7 (cit. on p. 12).
- [10] M. Sozeau, “Program-ing finger trees in coq”, *SIGPLAN Not.*, vol. 42, no. 9, pp. 13–24, Oct. 2007, ISSN: 0362-1340. DOI: 10.1145/1291220.1291156. [Online]. Available: <https://doi.org/10.1145/1291220.1291156> (cit. on p. 11).
- [11] M. Sozeau and N. Oury, “First-class type classes”, in *Lecture Notes in Computer Science*, Springer Berlin Heidelberg, 2008, pp. 278–293. DOI: 10.1007/978-3-540-71067-7_23. [Online]. Available: https://doi.org/10.1007/978-3-540-71067-7_23 (cit. on p. 8).
- [12] M. Sozeau, “Equations: A dependent pattern-matching compiler”, in *Interactive Theorem Proving*, Springer Berlin Heidelberg, 2010, pp. 419–434. DOI: 10.1007/978-3-642-14052-5_29. [Online]. Available: https://doi.org/10.1007/978-3-642-14052-5_29 (cit. on pp. 8, 12).

- [13] N. Benton, C.-K. Hur, A. J. Kennedy, and C. McBride, “Strongly typed term representations in coq”, *Journal of Automated Reasoning*, vol. 49, no. 2, pp. 141–159, Mar. 2011. doi: 10.1007/s10817-011-9219-0. [Online]. Available: <https://doi.org/10.1007/s10817-011-9219-0> (cit. on pp. 9, 14).
- [14] A. Chlipala, *Certified Programming with Dependent Types: A Pragmatic Introduction to the Coq Proof Assistant*. The MIT Press, 2013, ISBN: 0262026651 (cit. on pp. 11, 17).
- [15] A. Mahboubi and E. Tassi, “Canonical structures for the working coq user”, in *Interactive Theorem Proving*, Springer Berlin Heidelberg, 2013, pp. 19–34. doi: 10.1007/978-3-642-39634-2_5. [Online]. Available: https://doi.org/10.1007/978-3-642-39634-2_5 (cit. on p. 8).
- [16] A. G. Baydin, B. A. Pearlmutter, A. Radul, and J. M. Siskind, “Automatic differentiation in machine learning: A survey”, *J. Mach. Learn. Res.*, vol. 18, 153:1–153:43, 2015 (cit. on p. 5).
- [17] S. Boldo, C. Lelay, and G. Melquiond, “Coquelicot: A user-friendly library of real analysis for coq”, *Mathematics in Computer Science*, vol. 9, pp. 41–62, 2015 (cit. on p. 7).
- [18] B. C. Pierce, A. A. de Amorim, C. Casinghino, M. Gaboardi, M. Greenberg, C. Hrițcu, V. Sjöberg, A. Tolmach, and B. Yorgey, *Programming Language Foundations*, ser. Software Foundations series, volume 2. Electronic textbook, May 2018 (cit. on p. 8).
- [19] A. Abel, G. Allais, A. Hameer, B. Pientka, A. Momigliano, S. Schäfer, and K. Stark, “Poplmark reloaded: Mechanizing proofs by logical relations”, *Journal of Functional Programming*, vol. 29, e19, 2019. doi: 10.1017/S0956796819000170 (cit. on p. 10).
- [20] A. Shaikhha, A. Fitzgibbon, D. Vytiniotis, and S. Jones, “Efficient differentiable programming in a functional array-processing language”, *Proceedings of the ACM on Programming Languages*, vol. 3, pp. 1–30, Jul. 2019. doi: 10.1145/3341701 (cit. on pp. 5, 16).
- [21] L. Skorstengaard, *An introduction to logical relations*, 2019. arXiv: 1907.11133 [cs.PL] (cit. on p. 12).
- [22] M. Sozeau and C. Mangin, “Equations reloaded: High-level dependently-typed functional programming and proving in coq”, *Proceedings of the ACM on Programming Languages*, vol. 3, no. ICFP, pp. 1–29, Jul. 2019. doi: 10.1145/3341690. [Online]. Available: <https://doi.org/10.1145/3341690> (cit. on p. 12).

- [23] P. Wadler and W. Kokke, *Programming Language Foundations in Agda*, 2019, Available at <http://plfa.inf.ed.ac.uk/> (cit. on p. 8).
- [24] A. Aaby, “Introduction to programming languages”, *Syntax Imperative Programming Concurrent Programming Object-Oriented Programming Evaluation*, vol. 3, Apr. 2020 (cit. on p. 6).
- [25] G. Barthe, R. Crubillé, U. D. Lago, and F. Gavazzo, *On the versatility of open logical relations: Continuity, automatic differentiation, and a containment theorem*, 2020. arXiv: 2002.08489 [cs.PL] (cit. on pp. 12, 13, 16–18).
- [26] M. Huot, S. Staton, and M. Vákár, *Correctness of automatic differentiation via diffeologies and categorical gluing*, 2020. arXiv: 2001.02209 [cs.PL] (cit. on pp. 4, 7, 12–14, 16, 18).
- [27] M. Sozeau, “Subset coercions in coq”, in *Lecture Notes in Computer Science*, Springer Berlin Heidelberg, pp. 237–252. doi: 10.1007/978-3-540-74464-1_16. [Online]. Available: https://doi.org/10.1007/978-3-540-74464-1_16 (cit. on p. 11).