Correctness of Automatic Differentiation via Diffeologies and Categorical Gluing

Mathieu Huot ⊠^{1*}, Sam Staton^{1*}, and Matthijs Vákár^{2*}

¹ University of Oxford, UK
² Utrecht University, The Netherlands
*Equal contribution mathieu.huot@stx.ox.ac.uk

Abstract. We present semantic correctness proofs of Automatic Differentiation (AD). We consider a forward-mode AD method on a higher order language with algebraic data types, and we characterise it as the unique structure preserving macro given a choice of derivatives for basic operations. We describe a rich semantics for differentiable programming, based on diffeological spaces. We show that it interprets our language, and we phrase what it means for the AD method to be correct with respect to this semantics. We show that our characterisation of AD gives rise to an elegant semantic proof of its correctness based on a gluing construction on diffeological spaces. We explain how this is, in essence, a logical relations argument. Finally, we sketch how the analysis extends to other AD methods by considering a continuation-based method.

1 Introduction

Automatic differentiation (AD), loosely speaking, is the process of taking a program describing a function, and building the derivative of that function by applying the chain rule across the program code. As gradients play a central role in many aspects of machine learning, so too do automatic differentiation systems such as TensorFlow [1] or Stan [6].

Differentiation has a well developed mathematical theory in terms of differential geometry. The aim of this paper is to formalize this connection between differential geometry and the syntactic operations of AD. In this way we achieve two things: (1) a com-

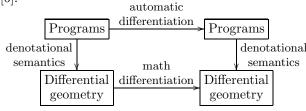


Fig. 1. Overview of semantics/correctness of AD.

positional, denotational understanding of differentiable programming and AD; (2) an explanation of the correctness of AD.

This intuitive correspondence (summarized in Fig. 1) is in fact rather complicated. In this paper we focus on resolving the following problem: higher order functions play a key role in programming, and yet they have no counterpart in traditional differential geometry. Moreover, we resolve this problem while retaining the compositionality of denotational semantics.

Higher order functions and differentiation. A major application of higher order functions is to support disciplined code reuse. Code reuse is particularly acute in machine learning. For example, a multi-layer neural network might be built of millions of near-identical neurons, as follows.

neuron_n: (realⁿ*(realⁿ*real))
$$\rightarrow$$
 real

neuron_n $\stackrel{\text{def}}{=} \lambda \langle x, \langle w, b \rangle \rangle$. $\varsigma(w \cdot x + b)$

layer_n: $((\tau_1 * P) \to \tau_2) \to (\tau_1 * P^n) \to \tau_2^n$

layer_n $\stackrel{\text{def}}{=} \lambda f$. $\lambda \langle x, \langle p_1, \dots, p_n \rangle \rangle$. $\langle f \langle x, p_1 \rangle, \dots, f \langle x, p_n \rangle \rangle$

comp: $(((\tau_1 * P) \to \tau_2) * ((\tau_2 * Q) \to \tau_3)) \to (\tau_1 * (P * Q)) \to \tau_3$
 τ

comp $\stackrel{\text{def}}{=} \lambda \langle f, g \rangle$. $\lambda \langle x, \langle p, g \rangle \rangle$. $\langle g \langle f \langle x, p \rangle, g \rangle$

(Here $\zeta(x) \stackrel{\text{def}}{=} \frac{1}{1+e^{-x}}$ is the sigmoid function, as illustrated.) We can use these functions to build a network as follows (see also Fig. 2):

$$\operatorname{comp}\langle \operatorname{layer}_m(\operatorname{neuron}_k), \operatorname{comp}\langle \operatorname{layer}_n(\operatorname{neuron}_m), \operatorname{neuron}_n \rangle \rangle : (\operatorname{real}^k * P) \to \operatorname{real}^k$$

Here $P \cong \mathbf{real}^p$ with p = (m(k+1) + n(m+1) + n + 1). This program (1) describes a smooth (infinitely differentiable) function. The goal of automatic differentiation is to find its derivative.

If we β -reduce all the λ 's, we end up with a very long function expression just built from the sigmoid function and linear algebra. We can then find a program for calculating its derivative by applying the chain rule. However, automatic differentiation can also be expressed without first β -reducing, in a compositional way, by explaining how higher order functions like (layer) and (comp) propagate derivatives.

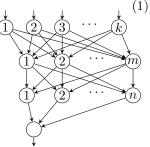


Fig. 2. The network in (1) with k inputs and two hidden layers.

This paper is a semantic analysis of this compositional approach.

The general idea of denotational semantics is to interpret types as spaces and programs as functions between the spaces. In this paper, we propose to use diffeological spaces and smooth functions [31,15] to this end. These satisfy the following three desiderata:

- $-\mathbb{R}$ is a space, and the smooth functions $\mathbb{R} \to \mathbb{R}$ are exactly the functions that are infinitely differentiable;
- The set of smooth functions $X \to Y$ between spaces again forms a space, so we can interpret function types.
- The disjoint union of a sequence of spaces again forms a space, so we can interpret variant types and inductive types.

We emphasise that the most standard formulation of differential geometry, using manifolds, does not support spaces of functions. Diffeological spaces seem to us the simplest notion of space that satisfies these conditions, but there are other

candidates [3,32]. A diffeological space is in particular a set X equipped with a chosen set of curves $C_X \subseteq X^{\mathbb{R}}$ and a smooth map $f: X \to Y$ must be such that if $\gamma \in C_X$ then $\gamma; f \in C_Y$. This is remiscent of the method of logical relations.

From smoothness to automatic derivatives at higher types. Our denotational semantics in diffeological spaces guarantees that all definable functions are smooth. But we need more than just to know that a definable function happens to have a mathematical derivative: we need to be able to find that derivative.

In this paper we focus on a simple, forward mode automatic differentiation method, which is a macro translation on syntax (called \vec{D} in §2). We are able to show that it is correct, using our denotational semantics.

Here there is one subtle point that is central to our development. Although differential geometry provides established derivatives for first order functions (such as neuron above), there is no canonical notion of derivative for higher order functions (such as layer and comp) in the theory of diffeological spaces (e.g. [7]). We propose a new way to resolve this, by interpreting types as triples (X, X', S) where, intuitively, X is a space of inhabitants of the type, X' is a space serving as a chosen bundle of tangents over X, and $S \subseteq X^{\mathbb{R}} \times X'^{\mathbb{R}}$ is a binary relation between curves, informally relating curves in X with their tangent curves in X'. This new model gives a denotational semantics for automatic differentiation.

In §3 we boil this new approach down to a straightforward and elementary logical relations argument for the correctness of automatic differentiation. The approach is explained in detail in §5.

Related work and context. AD has a long history and has many implementations. AD was perhaps first phrased in a functional setting in [25], and there are now a number of teams working on AD in the functional setting (e.g. [33, 30, 12]), some providing efficient implementations. Although that work does not involve formal semantics, it is inspired by intuitions from differential geometry and category theory.

This paper adds to a very recent body of work on verified automatic differentiation. Much of this is concurrent with and independent from the work in this article. In the first order setting, there are recent accounts based on denotational semantics in manifolds [13] and based on synthetic differential geometry [9], as well as work making a categorical abstraction [8] and work connecting operational semantics with denotational semantics [2,27]. Recently there has also been significant progress at higher types. The work of Brunel et al. gives formal correctness proofs for reverse-mode derivatives on computation graphs [5]. The work of Barthe et al. [4] provides a general discussion of some new syntactic logical relations arguments including one very similar to our syntactic proof of Theorem 1. We understand that the authors of [9] are working on higher types.

The differential λ -calculus [11] is related to AD, and explicit connections are made in [21, 22]. One difference is that the differential λ -calculus allows addition of terms at all types, and hence vector space models are suitable, but this appears peculiar with the variant and inductive types that we consider here.

Finally we emphasise that we have chosen the neural network (1) as our running example mainly for its simplicity. There are many other examples of AD

outside the neural networks literature: AD is useful whenever derivatives need to be calculated on high dimensional spaces. This includes optimization problems more generally, where the derivative is passed to a gradient descent method (e.g. [29, 17, 28, 18, 10, 20]). Other applications of AD are in advanced integration methods, since derivatives play a role in Hamiltonian Monte Carlo [24, 14] and variational inference [19].

Summary of contributions. We have provided a semantic analysis of automatic differentiation. Our syntactic starting point is a well-known forward-mode AD macro on a typed higher order language (e.g. [30, 33]). We recall this in §2 for function types, and in §4 we extend it to inductive types and variants. The main contributions of this paper are as follows.

- We give a denotational semantics for the language in diffeological spaces, showing that every definable expression is smooth (§3).
- We show correctness of the AD macro by a logical relations argument (Th. 1).
- We give a categorical analysis of this correctness argument with two parts: canonicity of the macro in terms of syntactic categories, and a new notion of glued space that abstracts the logical relation (§5).
- We then use this analysis to state and prove a correctness argument at all first order types (Th. 2).
- We show that our method is not specific to one particular AD macro, by also considering a continuation-based AD method (§6).

$\mathbf{2}$ A simple forward-mode AD translation

Rudiments of differentiation and dual numbers. Recall that the derivative of a function $f: \mathbb{R} \to \mathbb{R}$, if it exists, is a function $\nabla f: \mathbb{R} \to \mathbb{R}$ such that $\nabla f(x_0) = \frac{\mathrm{d}f(x)}{\mathrm{d}x}(x_0)$ is the gradient of f at x_0 . To find ∇f in a compositional way, two generalizations are reasonable:

- We need both f and ∇f when calculating $\nabla (f;g)$ of a composition f; g, using the chain rule, so we are really interested in the pair $(f, \nabla f) : \mathbb{R} \to \mathbb{R} \times \mathbb{R}$;
- In building f we will need to consider functions of multiple arguments, such as $+: \mathbb{R}^2 \to \mathbb{R}$, and these functions should propagate derivatives.

Thus we are more generally interested in transforming a function $g: \mathbb{R}^n \to \mathbb{R}$ into a function $h: (\mathbb{R} \times \mathbb{R})^n \to \mathbb{R} \times \mathbb{R}$ in such a way that for any $f_1 \dots f_n : \mathbb{R} \to \mathbb{R}$,

$$(f_1, \nabla f_1, \dots, f_n, \nabla f_n); h = ((f_1, \dots, f_n); g, \nabla ((f_1, \dots, f_n); g)).$$
 (2)

An intuition for h is often given in terms of dual numbers. The transformed function operates on pairs of numbers, (x, x'), and it is common to think of such a pair as $x + x'\epsilon$ for an 'infinitesimal' ϵ . But while this is a helpful intuition, the formalization of infinitesimals can be intricate, and the development in this paper is focussed on the elementary formulation in (2).

The reader may also notice that h encodes all the partial derivatives of g. For example, if $g: \mathbb{R}^2 \to \mathbb{R}$, then with $f_1(x) \stackrel{\text{def}}{=} x$ and $f_2(x) \stackrel{\text{def}}{=} x_2$, by applying (2) to x_1 we obtain $h(x_1, 1, x_2, 0) = (g(x_1, x_2), \frac{\partial g(x, x_2)}{\partial x}(x_1))$ and similarly $h(x_1, 0, x_2, 1) = (g(x_1, x_2), \frac{\partial g(x_1, x)}{\partial x}(x_2))$. And conversely, if g is differentiable in each argument, then a unique h satisfying (2) can be found by taking linear combinations of partial derivatives:

$$h(x_1, x_1', x_2, x_2') = (g(x_1, x_2), x_1' \cdot \frac{\partial g(x, x_2)}{\partial x}(x_1) + x_2' \cdot \frac{\partial g(x_1, x)}{\partial x}(x_2)).$$

In summary, the idea of differentiation with dual numbers is to transform a differentiable function $g: \mathbb{R}^n \to \mathbb{R}$ to a function $h: \mathbb{R}^{2n} \to \mathbb{R}^2$ which captures g and all its partial derivatives. We packaged this up in (2) as a sort-of invariant which is useful for building derivatives of compound functions $\mathbb{R} \to \mathbb{R}$ in a compositional way. The idea of forward mode automatic differentiation is to perform this transformation at the source code level.

A simple language of smooth functions. We consider a standard higher order typed language with a first order type real of real numbers. The types (τ, σ) and terms (t, s) are as follows.

$$(au,\sigma)$$
 and terms (t,s) are as follows.
 $au,\sigma,
ho$::= types | $(au_1*\dots* au_n)$ finite product | real real numbers | $au\to\sigma$ function t,s,r ::= terms | x variable | $\underline{c} \mid t+s \mid t*s \mid \varsigma(t)$ operations/constants | $\langle t_1,\dots,t_n \rangle$ | case t of $\langle x_1,\dots,x_n \rangle \to s$ tuples/pattern matching | $\lambda x.t \mid ts$ function abstraction/app.

The typing rules are in Figure 3. We have included a minimal set of operations for the sake of illustration, but it is not difficult to add further operations. We add some simple syntactic sugar $t-u \stackrel{\text{def}}{=} t + \underline{(-1)} * u$. We intend ς to stand for the sigmoid function, $\varsigma(x) \stackrel{\text{def}}{=} \frac{1}{1+e^{-x}}$. We further include syntactic sugar $\det x = t \ln s$ for $(\lambda x.s) t$ and $\lambda \langle x_1, \ldots, x_n \rangle$. t for $\lambda x.\mathbf{case} x$ of $\langle x_1, \ldots, x_n \rangle \to t$.

Syntactic automatic differentiation: a functorial macro. The aim of forward mode AD is to find the dual numbers representation of a function by syntactic manipulations. For our simple language, we implement this as the following inductively defined macro $\vec{\mathcal{D}}$ on both types and terms (see also [33, 30]):

lowing inductively defined macro
$$\overrightarrow{\mathcal{D}}$$
 on both types and terms (see also [33, 30]) $\overrightarrow{\mathcal{D}}(\mathbf{real}) \stackrel{\text{def}}{=} (\mathbf{real*real})$ $\overrightarrow{\mathcal{D}}(\tau \to \sigma) \stackrel{\text{def}}{=} \overrightarrow{\mathcal{D}}(\tau) \to \overrightarrow{\mathcal{D}}(\sigma)$ $\overrightarrow{\mathcal{D}}((\tau_1 * \cdots * \tau_n)) \stackrel{\text{def}}{=} (\overrightarrow{\mathcal{D}}(\tau_1) * \cdots * \overrightarrow{\mathcal{D}}(\tau_n))$

$$\frac{\Gamma \vdash \underline{c} : \mathbf{real}}{\Gamma \vdash \underline{c} : \mathbf{real}} (c \in \mathbb{R}) \quad \frac{\Gamma \vdash \underline{t} : \mathbf{real}}{\Gamma \vdash \underline{t} + s : \mathbf{real}} \quad \frac{\Gamma \vdash \underline{t} : \mathbf{real}}{\Gamma \vdash \underline{t} * s : \mathbf{real}} \quad \frac{\Gamma \vdash \underline{t} : \mathbf{real}}{\Gamma \vdash \underline{t} * s : \mathbf{real}} \quad \frac{\Gamma \vdash \underline{t} : \mathbf{real}}{\Gamma \vdash \zeta(\underline{t}) : \mathbf{real}} \quad \frac{\Gamma \vdash \underline{t} : \mathbf{real}}{\Gamma \vdash \zeta(\underline{t}) : \mathbf{real}} \quad \frac{\Gamma \vdash \underline{t} : \mathbf{real}}{\Gamma \vdash \zeta(\underline{t}) : \mathbf{real}} \quad \frac{\Gamma \vdash \underline{t} : \mathbf{real}}{\Gamma \vdash \zeta(\underline{t}) : \mathbf{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \mathbf{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \mathbf{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \mathbf{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \mathbf{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \mathbf{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \mathbf{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} : \underline{real}}{\Gamma \vdash \zeta(\underline{t}) : \underline{real}} \quad \frac{\Gamma \vdash \underline{t} :$$

Fig. 3. Typing rules for the simple language.

$$\overrightarrow{\mathcal{D}}(x) \stackrel{\mathrm{def}}{=} x \qquad \overrightarrow{\mathcal{D}}(\underline{c}) \stackrel{\mathrm{def}}{=} \langle \underline{c}, 0 \rangle$$

$$\overrightarrow{\mathcal{D}}(t+s) \stackrel{\mathrm{def}}{=} \mathbf{case} \, \overrightarrow{\mathcal{D}}(t) \, \mathbf{of} \, \langle x, x' \rangle \to \mathbf{case} \, \overrightarrow{\mathcal{D}}(s) \, \mathbf{of} \, \langle y, y' \rangle \to \langle x+y, x'+y' \rangle$$

$$\overrightarrow{\mathcal{D}}(t*s) \stackrel{\mathrm{def}}{=} \mathbf{case} \, \overrightarrow{\mathcal{D}}(t) \, \mathbf{of} \, \langle x, x' \rangle \to \mathbf{case} \, \overrightarrow{\mathcal{D}}(s) \, \mathbf{of} \, \langle y, y' \rangle \to \langle x*y, x*y'+x'*y \rangle$$

$$\overrightarrow{\mathcal{D}}(\varsigma(t)) \stackrel{\mathrm{def}}{=} \mathbf{case} \, \overrightarrow{\mathcal{D}}(t) \, \mathbf{of} \, \langle x, x' \rangle \to \mathbf{let} \, y = \varsigma(x) \, \mathbf{in} \, \langle y, x'*y*(1-y) \rangle$$

$$\overrightarrow{\mathcal{D}}(\lambda x.t) \stackrel{\mathrm{def}}{=} \lambda x. \overrightarrow{\mathcal{D}}(t) \quad \overrightarrow{\mathcal{D}}(ts) \stackrel{\mathrm{def}}{=} \overrightarrow{\mathcal{D}}(t) \, \overrightarrow{\mathcal{D}}(s) \quad \overrightarrow{\mathcal{D}}(\langle t_1, \dots, t_n \rangle) \stackrel{\mathrm{def}}{=} \langle \overrightarrow{\mathcal{D}}(t_1), \dots, \overrightarrow{\mathcal{D}}(t_n) \rangle$$

$$\overrightarrow{\mathcal{D}}(\mathbf{case} \, t \, \mathbf{of} \, \langle x_1, \dots, x_n \rangle \to s) \stackrel{\mathrm{def}}{=} \mathbf{case} \, \overrightarrow{\mathcal{D}}(t) \, \mathbf{of} \, \langle x_1, \dots, x_n \rangle \to \overrightarrow{\mathcal{D}}(s)$$

We extend $\overrightarrow{\mathcal{D}}$ to contexts: $\overrightarrow{\mathcal{D}}(\{x_1:\tau_1,...,x_n:\tau_n\}) \stackrel{\text{def}}{=} \{x_1:\overrightarrow{\mathcal{D}}(\tau_1),...,x_n:\overrightarrow{\mathcal{D}}(\tau_n)\}.$ This turns $\overrightarrow{\mathcal{D}}$ into a well-typed, functorial macro in the following sense.

Lemma 1 (Functorial macro). If
$$\Gamma \vdash t : \tau$$
 then $\overrightarrow{\mathcal{D}}(\Gamma) \vdash \overrightarrow{\mathcal{D}}(t) : \overrightarrow{\mathcal{D}}(\tau)$. If $\Gamma, x : \sigma \vdash t : \tau$ and $\Gamma \vdash s : \sigma$ then $\overrightarrow{\mathcal{D}}(\Gamma) \vdash \overrightarrow{\mathcal{D}}(t[s/x]) = \overrightarrow{\mathcal{D}}(t)[\overrightarrow{\mathcal{D}}(s)/x]$.

Example 1 (Inner products). Let us write τ^n for the *n*-fold product $(\tau * \dots * \tau)$. Then, given $\Gamma \vdash t, s : \mathbf{real}^n$ we can define their inner product

$$\Gamma \vdash t \cdot_n s \stackrel{\text{def}}{=} \mathbf{case} \, t \, \mathbf{of} \, \langle z_1, \dots, z_n \rangle \to \\ \mathbf{case} \, s \, \mathbf{of} \, \langle y_1, \dots, y_n \rangle \to z_1 * y_1 + \dots + z_n * y_n : \mathbf{real} \\ \text{To illustrate the calculation of } \overrightarrow{\mathcal{D}}, \text{ let us expand (and } \beta\text{-reduce) } \overrightarrow{\mathcal{D}} \, (t \cdot_2 s) : \\ \mathbf{case} \, \overrightarrow{\mathcal{D}} \, (t) \, \mathbf{of} \, \langle z_1, z_2 \rangle \to \mathbf{case} \, \overrightarrow{\mathcal{D}} \, (s) \, \mathbf{of} \, \langle y_1, y_2 \rangle \to \mathbf{case} \, z_1 \, \mathbf{of} \, \langle z_{1,1}, z_{1,2} \rangle \to \\ \mathbf{case} \, y_1 \, \mathbf{of} \, \langle y_{1,1}, y_{1,2} \rangle \to \mathbf{case} \, z_2 \, \mathbf{of} \, \langle z_{2,1}, z_{2,2} \rangle \to \mathbf{case} \, y_2 \, \mathbf{of} \, \langle y_{2,1}, y_{2,2} \rangle \to \\ \langle z_{1,1} * y_{1,1} + z_{2,1} * y_{2,1}, \ z_{1,1} * y_{1,2} + z_{1,2} * y_{1,1} + z_{2,1} * y_{2,2} + z_{2,2} * y_{2,1} \rangle$$

Example 2 (Neural networks). In our introduction (1), we provided a program in our language to build a neural network out of expressions neuron, layer, comp; this program makes use of the inner product of Ex. 1. We can similarly calculate $\overrightarrow{\mathcal{D}}$ of such deep neural nets by mechanically applying the macro.

3 Semantics of differentiation

Consider for a moment the first order fragment of the language in § 2, with only one type, **real**, and no λ 's or pairs. This has a simple semantics in the category of cartesian spaces and smooth maps. Indeed, a term $x_1 \dots x_n : \mathbf{real} \vdash t : \mathbf{real}$ has a natural reading as a function $[\![t]\!] : \mathbb{R}^n \to \mathbb{R}$ by interpreting our operation symbols by the well-known operations on $\mathbb{R}^n \to \mathbb{R}$ with the corresponding name. In fact, the functions that are definable in this first order fragment are smooth, which means that they are continuous, differentiable, and their derivatives are continuous, differentiable, and so on. Let us write **CartSp** for this category of cartesian spaces (\mathbb{R}^n for some n) and smooth functions.

The category **CartSp** has cartesian products, and so we can also interpret product types, tupling and pattern matching, giving us a useful syntax for constructing functions into and out of products of \mathbb{R} . For example, the interpretation of (neuron_n) in (1) becomes

$$\mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R} \xrightarrow{[\![\cdot_n]\!] \times \mathrm{id}_{\mathbb{R}}} \mathbb{R} \times \mathbb{R} \xrightarrow{[\![+]\!]} \mathbb{R} \xrightarrow{[\![\varsigma]\!]} \mathbb{R}.$$

where $[\![\cdot_n]\!]$, $[\![+]\!]$ and $[\![\varsigma]\!]$ are the usual inner product, addition and the sigmoid function on \mathbb{R} , respectively.

Inside this category, we can straightforwardly study the first order language without λ 's, and automatic differentiation. In fact, we can prove the following by plain induction on the syntax:

The interpretation of the (syntactic) forward AD $\vec{\mathcal{D}}(t)$ of a first-order term t equals the usual (semantic) derivative of the interpretation of t as a smooth function.

However, as is well known, the category CartSp does not support function spaces. To see this, notice that we have polynomial terms

$$x_1, \dots, x_d : \mathbf{real} \vdash \lambda y. \sum_{n=1}^d x_n y^n : \mathbf{real} \to \mathbf{real}$$

for each d, and so if we could interpret (real \to real) as a Euclidean space \mathbb{R}^p then, by interpreting these polynomial expressions, we would be able to find continuous injections $\mathbb{R}^d \to \mathbb{R}^p$ for every d, which is topologically impossible for any p, for example as a consequence of the Borsuk-Ulam theorem (see Appx. A).

This means that we cannot interpret the functions (layer) and (comp) from (1) in CartSp, as they are higher order functions, even though they are very useful and innocent building blocks for differential programming! Clearly, we could define neural nets such as (1) directly as smooth functions without any higher order subcomponents, though that would quickly become cumbersome for deep networks. A problematic consequence of the lack of a semantics for higher order differential programs is that we have no obvious way of establishing compositional semantic correctness of $\overrightarrow{\mathcal{D}}$ for the given implementation of (1).

Diffeological spaces. This motivates us to turn to a more general notion of differential geometry for our semantics, based on diffeological spaces [15]. The key idea will be that a higher order function is called smooth if it sends smooth functions to smooth functions, meaning that we can never use it to build first order functions that are not smooth. For example, (comp) in (1) has this property.

Definition 1. A diffeological space (X, \mathcal{P}_X) consists of a set X together with, for each n and each open subset U of \mathbb{R}^n , a set $\mathcal{P}_X^U \subseteq [U \to X]$ of functions, called plots, such that

- all constant functions are plots;
- if $f: V \to U$ is a smooth function and $p \in \mathcal{P}_X^U$, then $f; p \in \mathcal{P}_X^V$; if $\left(p_i \in \mathcal{P}_X^{U_i}\right)_{i \in I}$ is a compatible family of plots $(x \in U_i \cap U_j \Rightarrow p_i(x) = p_j(x))$ and $(U_i)_{i \in I}$ covers U, then the gluing $p: U \to X: x \in U_i \mapsto p_i(x)$ is a plot.

We call a function $f: X \to Y$ between diffeological spaces smooth if, for all plots $p \in \mathcal{P}_X^U$, we have that $p; f \in \mathcal{P}_Y^U$. We write $\mathbf{Diff}(X,Y)$ for the set of smooth maps from X to Y. Smooth functions compose, and so we have a category **Diff** of diffeological spaces and smooth functions.

A diffeological space is thus a set equipped with structure. Many constructions of sets carry over straightforwardly to diffeological spaces.

Example 3 (Cartesian diffeologies). Each open subset U of \mathbb{R}^n can be given the structure of a diffeological space by taking all the smooth functions $V \to U$

as \mathcal{P}_U^V . It is easily seen that smooth functions from $V \to U$ in the traditional sense coincide with smooth functions in the sense of diffeological spaces. Thus diffeological spaces have a profound relationship with ordinary calculus.

In categorical terms, this gives a full embedding of CartSp in Diff.

Example 4 (Product diffeologies). Given a family $(X_i)_{i\in I}$ of diffeological spaces, we can equip the product $\prod_{i\in I} X_i$ of sets with the product diffeology in which U-plots are precisely the functions of the form $(p_i)_{i\in I}$ for $p_i \in \mathcal{P}^U_{X_i}$.

This gives us the categorical product in **Diff**.

Example 5 (Functional diffeology). We can equip the set $\mathbf{Diff}(X,Y)$ of smooth functions between diffeological spaces with the functional diffeology in which U-plots consist of functions $f:U\to\mathbf{Diff}(X,Y)$ such that $(u,x)\mapsto f(u)(x)$ is an element of $\mathbf{Diff}(U\times X,Y)$.

This specifies the categorical function object in **Diff**.

Semantics and correctness of AD. We can now give a denotational semantics to our language from § 2. We interpret each type τ as a set $\llbracket \tau \rrbracket$ equipped with the relevant diffeology, by induction on the structure of types:

$$\llbracket \mathbf{real} \rrbracket \stackrel{\text{def}}{=} \mathbb{R} \qquad \llbracket (\tau_1 * \dots * \tau_n) \rrbracket \stackrel{\text{def}}{=} \prod_{i=1}^n \llbracket \tau_i \rrbracket \qquad \llbracket \tau \to \sigma \rrbracket \stackrel{\text{def}}{=} \mathbf{Diff}(\llbracket \tau \rrbracket, \llbracket \sigma \rrbracket)$$

A context $\Gamma = (x_1 \colon \tau_1 \dots x_n \colon \tau_n)$ is interpreted as a diffeological space $\llbracket \Gamma \rrbracket \stackrel{\text{def}}{=} \prod_{i=1}^n \llbracket \tau_i \rrbracket$. Now well typed terms $\Gamma \vdash t \colon \tau$ are interpreted as smooth functions $\llbracket t \rrbracket \colon \llbracket \Gamma \rrbracket \to \llbracket \tau \rrbracket$, giving a meaning for t for every valuation of the context. This is routinely defined by induction on the structure of typing derivations. Constants $\underline{c} \colon \mathbf{real}$ are interpreted as constant functions; and the first order operations $(+,*,\varsigma)$ are interpreted by composing with the corresponding functions, which are smooth. For example, $\llbracket \varsigma(t) \rrbracket (\rho) \stackrel{\text{def}}{=} \varsigma(\llbracket t \rrbracket (\rho))$, where $\rho \in \llbracket \Gamma \rrbracket$. Variables are interpreted as $\llbracket x_i \rrbracket (\rho) \stackrel{\text{def}}{=} \rho_i$. The remaining constructs are interpreted as follows, and it is straightforward to show that smoothness is preserved.

$$[\![\langle t_1, \dots, t_n \rangle]\!](\rho) \stackrel{\text{def}}{=} ([\![t_1]\!](\rho), \dots, [\![t_n]\!](\rho)) \qquad [\![\lambda x : \tau.t]\!](\rho)(a) \stackrel{\text{def}}{=} [\![t]\!](\rho, a) \ (a \in [\![\tau]\!])$$
$$[\![\mathbf{case} \ t \ \mathbf{of} \ \langle \dots \rangle \to s]\!](\rho) \stackrel{\text{def}}{=} [\![s]\!](\rho, [\![t]\!](\rho)) \qquad [\![t \ s]\!](\rho) \stackrel{\text{def}}{=} [\![t]\!](\rho)([\![s]\!](\rho))$$

Notice that a term x_1 : **real**, ..., x_n : **real** $\vdash t$: **real** is interpreted as a smooth function $\llbracket t \rrbracket : \mathbb{R}^n \to \mathbb{R}$, even if t involves higher order functions (like (1)). Moreover the macro differentiation $\overrightarrow{\mathcal{D}}(t)$ is a function $\llbracket \overrightarrow{\mathcal{D}}(t) \rrbracket : (\mathbb{R} \times \mathbb{R})^n \to (\mathbb{R} \times \mathbb{R})$. This enables us to state a limited version of our main correctness theorem:

Theorem 1 (Semantic correctness of $\overrightarrow{\mathcal{D}}$ (limited)). For any term $x_1 : \mathbf{real}, \ldots, x_n : \mathbf{real} \vdash t : \mathbf{real}$, the function $[\![\overrightarrow{\mathcal{D}}(t)]\!]$ is the dual numbers representation (2) of $[\![t]\!]$. In detail: for any smooth functions $f_1 \ldots f_n : \mathbb{R} \to \mathbb{R}$,

$$(f_1, \nabla f_1, \dots, f_n, \nabla f_n); \llbracket \overrightarrow{\mathcal{D}}(t) \rrbracket = ((f_1 \dots f_n); \llbracket t \rrbracket, \nabla ((f_1 \dots f_n); \llbracket t \rrbracket)).$$

(For instance, if
$$n=2$$
, then $[\![\overrightarrow{\mathcal{D}}(t)]\!](x_1,1,x_2,0)=([\![t]\!](x_1,x_2),\frac{\partial [\![t]\!](x,x_2)}{\partial x}(x_1)).)$

Proof. We prove this by logical relations. Although the following proof is elementary, we found it by using the categorical methods in § 5.

For each type τ , we define a binary relation S_{τ} between curves in $\llbracket \tau \rrbracket$ and curves in $\llbracket \overrightarrow{\mathcal{D}}(\tau) \rrbracket$, i.e. $S_{\tau} \subseteq \mathcal{P}_{\llbracket \tau \rrbracket}^{\mathbb{R}} \times \mathcal{P}_{\llbracket \overrightarrow{\mathcal{D}}(\tau) \rrbracket}^{\mathbb{R}}$, by induction on τ :

```
 -S_{\mathbf{real}} \stackrel{\text{def}}{=} \{ (f, (f, \nabla f)) \mid f : \mathbb{R} \to \mathbb{R} \text{ smooth} \}; 
 -S_{(\tau * \sigma)} \stackrel{\text{def}}{=} \{ ((f_1, g_1), (f_2, g_2)) \mid (f_1, f_2) \in S_{\tau}, (g_1, g_2) \in S_{\sigma} \}; 
 -S_{\tau \to \sigma} \stackrel{\text{def}}{=} \{ (f_1, f_2) \mid \forall (g_1, g_2) \in S_{\tau}. (x \mapsto f_1(x)(g_1(x)), x \mapsto f_2(x)(g_2(x))) \in S_{\sigma} \}. 
Then, we establish the following 'fundamental lemma':
```

If
$$x_1:\tau_1,...,x_n:\tau_n \vdash t : \sigma$$
 and, for all $1 \le i \le n$, for all smooth $f_i : \mathbb{R} \to \llbracket \tau_i \rrbracket$ and $g_i : \mathbb{R} \to \llbracket \overrightarrow{\mathcal{D}}(\tau_i) \rrbracket$ such that (f_i,g_i) is in S_{τ_i} , we have that $((f_1,...,f_n);\llbracket t \rrbracket,(g_1,...,g_n);\llbracket \overrightarrow{\mathcal{D}}(t) \rrbracket)$ is in S_{σ} .

This is proved routinely by induction on the typing derivation of t. The case for * relies on the precise definition of $\overrightarrow{\mathcal{D}}(t*s)$, and similarly for $+, \varsigma$.

We conclude the theorem from the fundamental lemma by considering the case where $\tau_i = \sigma = \mathbf{real}$, m = n and $s_i = y_i$.

4 Extending the language: variant and inductive types

In this section, we show that the definition of forward AD and the semantics generalize if we extend the language of §2 with variants and inductive types. As an example of inductive types, we consider lists. This specific choice is only for expository purposes and the whole development works at the level of generality of arbitrary algebraic data types generated as initial algebras of (polynomial) type constructors formed by finite products and variants.

Similarly, our choice of operations is for expository purposes. More generally, assume given a family of operations $(\mathsf{Op}_n)_{n\in\mathbb{N}}$ indexed by their arity n. Further assume that each $\mathsf{op}\in\mathsf{Op}_n$ has type $\mathsf{real}^n\to\mathsf{real}$. We then ask for a certain closure of these operations under differentiation, that is we define

$$\overrightarrow{\mathcal{D}}(\mathsf{op}(t_1,\ldots,t_n)) \stackrel{\mathrm{def}}{=} \mathbf{case} \, \overrightarrow{\mathcal{D}}(t_1) \, \mathbf{of} \, \langle x_1, x_1' \rangle \to \ldots \to \mathbf{case} \, \overrightarrow{\mathcal{D}}(t_n) \, \mathbf{of} \, \langle x_n, x_n' \rangle \to \langle \mathsf{op}(x_1,\ldots,x_n), \sum_{i=1}^n x_i' * \partial_i \mathsf{op}(x_1,\ldots,x_n) \rangle$$
 where $\partial_i \mathsf{op}(x_1,\ldots,x_n)$ is some chosen term in the language, involving free vari-

where $\partial_i \operatorname{op}(x_1, \dots, x_n)$ is some chosen term in the language, involving free variables from x_1, \dots, x_n , which we think of as implementing the partial derivative of op with respect to its *i*-th argument. For constructing the semantics, every op must be interpreted by some smooth function, and, to establish correctness, the semantics of $\partial_i \operatorname{op}(x_1, \dots, x_n)$ must be the semantic *i*-th partial derivative of the semantics of $\operatorname{op}(x_1, \dots, x_n)$.

Language. We additionally consider the following types and terms:

$$\tau, \sigma, \rho ::= types | list(\tau) list
| {\ell_1 \tau_1 | \dots | \ell_n \tau_n} variant$$

$$\begin{array}{lll} t, s, r & ::= & \text{terms} \\ & | & \tau.\ell\,t & \text{variant constructor} \\ & | & [] & | & t :: s & \text{empty list and cons} \\ & | & \textbf{case}\,t\,\textbf{of}\,\{\ell_1\,x_1 \to s_1 \ \big| \, \cdots \ \big|\,\ell_{\mathsf{n}}\,x_n \to s_n\} & \text{pattern matching: variants} \\ & | & \textbf{fold}\,(x_1,x_2).t\,\textbf{over}\,s\,\textbf{from}\,r & \text{list fold} \end{array}$$

We extend the type system according to:

$$\frac{\Gamma \vdash t : \tau_i}{\Gamma \vdash \tau . \ell_i t : \tau} ((\ell_i \tau_i) \in \tau) \quad \frac{\Gamma \vdash t : \tau}{\Gamma \vdash [] : \mathbf{list}(\tau)} \quad \frac{\Gamma \vdash t : \tau}{\Gamma \vdash t : s : \mathbf{list}(\tau)} \frac{\Gamma \vdash t : s : \mathbf{list}(\tau)}{\Gamma \vdash t : s : \mathbf{list}(\tau)}$$

$$\frac{\Gamma \vdash t : \{\ell_1 \tau_1 \mid \dots \mid \ell_n \tau_n\} \quad \text{for each } 1 \le i \le n : \Gamma, x_i : \tau_i \vdash s_i : \tau}{\Gamma \vdash \mathbf{case} \, t \, \mathbf{of} \, \{\ell_1 x_1 \to s_1 \mid \dots \mid \ell_n x_n \to s_n\} : \tau}$$

$$\frac{\Gamma \vdash \mathbf{case} \, t \, \mathbf{of} \, \{\ell_1 x_1 \to s_1 \mid \dots \mid \ell_n x_n \to s_n\} : \tau}{\Gamma \vdash s : \mathbf{list}(\tau) \quad \Gamma \vdash r : \sigma \quad \Gamma, x_1 : \tau, x_2 : \sigma \vdash t : \sigma}$$

$$\frac{\Gamma \vdash \mathbf{fold} \, (x_1, x_2) . \mathbf{t} \, \mathbf{over} \, s \, \mathbf{from} \, r : \sigma}{\Gamma \vdash \mathbf{fold} \, (x_1, x_2) . \mathbf{t} \, \mathbf{over} \, s \, \mathbf{from} \, r : \sigma}$$

We can then extend $\overrightarrow{\mathcal{D}}$ to our new types and terms by

To demonstrate the practical use of expressive type systems for differential programming, we consider the following two examples.

Example 6 (Lists of inputs for neural nets). Usually, we run a neural network on a large data set, the size of which might be determined at runtime. To evaluate a neural network on multiple inputs, in practice, one often sums the outcomes. This can be coded in our extended language as follows. Suppose that we have a network $f: (\mathbf{real}^n * P) \to \mathbf{real}$ that operates on single input vectors. We can construct one that operates on lists of inputs as follows:

$$g \stackrel{\text{def}}{=} \lambda \langle l, w \rangle$$
.fold $(x_1, x_2).f \langle x_1, w \rangle + x_2$ over l from $0: (\text{list}(\text{real}^n) * P) \to \text{real}$

Example 7 (Missing data). In practically every application of statistics and machine learning, we face the problem of missing data: for some observations, only partial information is available. In an expressive typed programming language like we consider, we can model missing data conveniently using the data type $\mathbf{maybe}(\tau) = \{ \text{Nothing ()} \mid \text{Just } \tau \}$. In the context of a neural network, one might use it as follows. First, define some helper functions

from Maybe
$$\stackrel{\text{def}}{=} \lambda x. \lambda m.$$
 case m of $\{\text{Nothing } \underline{\ } \to x \mid \text{Just } x' \to x' \}$
from Maybe $\stackrel{\text{def}}{=} \lambda \langle x_1, ..., x_n \rangle. \lambda \langle m_1, ..., m_n \rangle. \langle \text{from Maybe } x_1 m_1, ..., \text{from Maybe } x_n m_n \rangle$
 $: (\text{maybe}(\text{real}))^n \to \text{real}^n \to \text{real}^n$
 $\max \stackrel{\text{def}}{=} \lambda f. \lambda l.$ fold $(x_1, x_2). f(x_1 :: x_2)$ over l from $[]: (\tau \to \sigma) \to \text{list}(\tau) \to \text{list}(\sigma)$

Given a neural network $f: (\mathbf{list}(\mathbf{real}^k) * P) \to \mathbf{real}$, we can build a new one that operates on on a data set for which some covariates (features) are missing, by passing in default values to replace the missing covariates:

$$\lambda \langle l, \langle m, w \rangle \rangle . f \langle \text{map (fromMaybe}^k m) \, l, w \rangle$$

: (list((maybe(real))^k)*(real^k*P)) \rightarrow real

Then, given a data set l with missing covariates, we can perform automatic differentiation on this network to optimize, simultaneously, the ordinary network parameters w and the default values for missing covariates m.

Semantics. In \S 3 we gave a denotational semantics for the simple language in diffeological spaces. This extends to the language in this section, as follows. As before, each type τ is interpreted as a diffeological space, which is a set equipped with a family of plots:

– A variant type $\{\ell_1 \tau_1 \mid \dots \mid \ell_n \tau_n\}$ is inductively interpreted as the disjoint union $[\![\{\ell_1 \tau_1 \mid \dots \mid \ell_n \tau_n\}]\!] \stackrel{\text{def}}{=} \biguplus_{i=1}^n [\![\tau_i]\!]$ with U-plots

$$\mathcal{P}^{U}_{\llbracket\{\ell_1\,\tau_1\big|\dots\big|\ell_n\,\tau_n\}\rrbracket} \stackrel{\text{def}}{=} \left\{ \left[U_j \stackrel{f_j}{\longrightarrow} \llbracket\tau_j\rrbracket \to \biguplus_{i=1}^n \llbracket\tau_i\rrbracket \right]_{j=1}^n \middle| U = \biguplus_{j=1}^n U_j, \ f_j \in \mathcal{P}^{U_j}_{\llbracket\tau_j\rrbracket} \right\}.$$

- A list type $\mathbf{list}(\tau)$ is interpreted as the set of lists, $[\![\mathbf{list}(\tau)]\!] \stackrel{\text{def}}{=} \bigcup_{i=1}^{\infty} [\![\tau]\!]^i$ with U-plots

$$\mathcal{P}_{\llbracket \mathbf{list}(\tau) \rrbracket}^U \stackrel{\mathrm{def}}{=} \left\{ \left[U_j \xrightarrow{f_j} \llbracket \tau \rrbracket^j \to \biguplus_{i=1}^{\infty} \llbracket \tau \rrbracket^i \right]_{j=1}^{\infty} \mid U = \biguplus_{j=1}^{\infty} U_j, \ f_j \in \mathcal{P}_{\llbracket \tau \rrbracket^j}^{U_j} \right\}.$$

The constructors and destructors for variants and lists are interpreted as in the usual set theoretic semantics. It is routine to show inductively that these interpretations are smooth. Thus every term $\Gamma \vdash t : \tau$ in the extended language is interpreted as a smooth function $\llbracket t \rrbracket : \llbracket \Gamma \rrbracket \to \llbracket \tau \rrbracket$ between diffeological spaces.

(In this section we focused on a language with lists, but other inductive types are easily interpreted in the category of diffeological spaces in much the same way; the categorically minded reader may regard this as a consequence of **Diff** being a concrete Grothendieck quasitopos, e.g. [3].)

5 Categorical analysis of forward AD and its correctness

This section has three parts. First, we give a categorical account of the functoriality of AD (Ex. 8). Then we introduce our gluing construction, and relate it to the correctness of AD (dgm. (3)). Finally, we state and prove a correctness theorem for all first order types by considering a category of manifolds (Th. 2).

Syntactic categories. Our language induces a syntactic category as follows.

Definition 2. Let Syn be the category whose objects are types, and where a morphism $\tau \to \sigma$ is a term in context $x : \tau \vdash t : \sigma$ modulo the $\beta\eta$ -laws (Fig. 4). Composition is by substitution.

For simplicity, we do not impose arithmetic identities such as x + y = y + x in **Syn**. As is standard, this category has the following universal property.

Lemma 2 (e.g. [26]). For every bicartesian closed category C with list objects, and every object $F(\mathbf{real}) \in C$ and morphisms $F(\underline{c}) \in C(1, F(\mathbf{real}))$, $F(+), F(*) \in C(F(\mathbf{real}) \times F(\mathbf{real}))$, $F(\mathbf{real})$), $F(\varsigma) \in \mathbf{Syn}(F(\mathbf{real}), F(\mathbf{real}))$ in C, there is a unique functor $F: \mathbf{Syn} \to C$ respecting the interpretation and preserving the bicartesian closed structure as well as list objects.

Proof (notes). The functor $F: \mathbf{Syn} \to \mathcal{C}$ is a canonical denotational semantics for the language, interpreting types as objects of \mathcal{C} and terms as morphisms. For instance, $F(\tau \to \sigma) \stackrel{\text{def}}{=} (F\tau \to F\sigma)$, the function space in the category \mathcal{C} , and $F(ts) \stackrel{\text{def}}{=}$ is the composite (Ft, Fs); eval. When $\mathcal{C} = \mathbf{Diff}$, the denotational semantics of the language in diffeological spaces (§3,4) can be understood as the unique structure preserving functor $[\![-]\!]: \mathbf{Syn} \to \mathbf{Diff}$ satisfying $[\![\mathbf{real}]\!] = \mathbb{R}$, $[\![\varsigma]\!] = \varsigma$ and so on.

Example 8 (Canonical definition forward AD). The forward AD macro $\overrightarrow{\mathcal{D}}$ (§2,4) arises as a canonical cartesian closed functor on **Syn**. Consider the unique cartesian closed functor $F: \mathbf{Syn} \to \mathbf{Syn}$ such that $F(\mathbf{real}) = \mathbf{real} * \mathbf{real}, F(\underline{c}) = \overrightarrow{\mathcal{D}}(\underline{c}), F(\varsigma) = \overrightarrow{\mathcal{D}}(\varsigma(x)),$ and $F(+) = \gamma : F(\mathbf{real}) * F(\mathbf{real}) + Case \circ of / x, y \to \overrightarrow{\mathcal{D}}(x+y) : F(\mathbf{real})$

```
F(+) = z : F(\mathbf{real}) * F(\mathbf{real}) \vdash \mathbf{case} \ z \ \mathbf{of} \ \langle x,y \rangle \to \overrightarrow{\mathcal{D}}(x+y) : F(\mathbf{real})

F(*) = z : F(\mathbf{real}) * F(\mathbf{real}) \vdash \mathbf{case} \ z \ \mathbf{of} \ \langle x,y \rangle \to \overrightarrow{\mathcal{D}}(x*y) : F(\mathbf{real})

Then for any type \tau, F(\tau) = \overrightarrow{\mathcal{D}}(\tau), and for any term x : \tau \vdash t : \sigma, F(t) = \overrightarrow{\mathcal{D}}(t) as morphisms F(\tau) \to F(\sigma) in the syntactic category.
```

Categorical gluing and logical relations. Gluing is a method for building new categorical models which has been used for many purposes, including logical relations and realizability [23]. Our logical relations argument in the proof of Th. 1 can be understood in this setting. In this subsection, for the categorically minded, we explain this, and in doing so we quickly recover a correctness result for the more general language in \S 4 and for arbitrary first order types.

We define a category \mathbf{Gl}_U whose objects are triples (X, X', S) where X and X' are diffeological spaces and $S \subseteq \mathcal{P}_X^U \times \mathcal{P}_{X'}^U$ is a relation between their U-plots. A morphism $(X, X', S) \to (Y, Y', T)$ is a pair of smooth functions

```
\begin{array}{c} \mathbf{case} \ \langle t_1, \dots, t_n \rangle \ \mathbf{of} \ \langle x_1, \dots, x_n \rangle \rightarrow s = s[^{t_1}/_{x_1}, \dots, ^{t_n}/_{x_n}] \\ s[^t/_y] \ ^\#x_1 = \dots x_n \ \mathbf{case} \ t \ \mathbf{of} \ \langle x_1, \dots, x_n \rangle \rightarrow s[^{(x_1, \dots, x_n)}/_y] \\ \mathbf{case} \ \ell_i \ t \ \mathbf{of} \ \{\ell_1 \ x_1 \rightarrow s_1 \ | \ \cdots \ | \ \ell_n \ x_n \rightarrow s_n\} = s_i [^t/_{x_i}] \\ s[^t/_y] \ ^\#x_1 = \dots x_n \ \mathbf{case} \ t \ \mathbf{of} \ \{\ell_1 \ x_1 \rightarrow s[^{\ell_1 \ x_1}/_y] \ | \ \cdots \ | \ \ell_n \ x_n \rightarrow s[^{\ell_n \ x_n}/_y]\} \\ \mathbf{fold} \ (x_1, x_2) . t \ \mathbf{over} \ [] \ \mathbf{from} \ r = r \\ \mathbf{fold} \ (x_1, x_2) . t \ \mathbf{over} \ s_1 :: s_2 \ \mathbf{from} \ r = t[^{s_1}/_{x_1}, \ \mathbf{fold} \ (x_1, x_2) . t \ \mathbf{over} \ s_2 \ \mathbf{from} \ r \\ u = s[^{l}/_y], r[^s/_{x_2}] = s[^{x_1 :: y}/_y] \Rightarrow s[^t/_y] \ ^\#x_1 = 2 \\ \mathbf{fold} \ (x_1, x_2) . r \ \mathbf{over} \ t \ \mathbf{from} \ u \end{array}
```

Fig. 4. Standard $\beta\eta$ -laws (e.g. [26]) for products, functions, variants and lists.

 $f: X \to Y$, $f': X' \to Y'$, such that if $(g, g') \in S$ then $(g; f, g'; f') \in T$. The idea is that this is a semantic domain where we can simultaneously interpret the language and its automatic derivatives.

Proposition 1. The category Gl_U is bicartesian closed, has list objects, and the projection functor $proj : Gl_U \to Diff \times Diff$ preserves this structure.

Proof (notes). The category \mathbf{Gl}_U is a full subcategory of the comma category $\mathrm{id}_{\mathbf{Set}} \downarrow \mathbf{Diff}(U, -) \times \mathbf{Diff}(U, -)$. The result thus follows by the general theory of categorical gluing (e.g. [16, Lemma 15]).

We give a semantics $(-) = ((-)_0, (-)_1, S_-)$ for the language in $\mathbf{Gl}_{\mathbb{R}}$, interpreting types τ as objects $((\tau)_0, (\tau)_1, S_\tau)$, and terms as morphisms. We let $(\mathbf{real})_0 \stackrel{\text{def}}{=} \mathbb{R}$ and $(\mathbf{real})_1 \stackrel{\text{def}}{=} \mathbb{R}^2$, with the relation $S_{\mathbf{real}} \stackrel{\text{def}}{=} \{(f, (f, \nabla f)) \mid f : \mathbb{R} \to \mathbb{R} \text{ smooth}\}$. We interpret the constants \underline{c} as pairs $(\underline{c})_0 \stackrel{\text{def}}{=} \underline{c}$ and $(\underline{c})_1 \stackrel{\text{def}}{=} (\underline{c}, 0)$, and we interpret $+, \times, \varsigma$ in the standard way (meaning, like [-]) in $(-)_0$, but according to the derivatives in $(-)_1$, for instance, $(*)_1 : \mathbb{R}^2 \times \mathbb{R}^2 \to \mathbb{R}^2$ is

$$(*)_1((x,x'),(y,y')) \stackrel{\text{def}}{=} (xy,xy'+x'y).$$

At this point one checks that these interpretations are indeed morphisms in $\mathbf{Gl}_{\mathbb{R}}$. This amounts to checking that these interpretations are dual numbers representations in the sense of (2). The remaining constructions of the language are interpreted using the categorical structure of $\mathbf{Gl}_{\mathbb{R}}$, following Lem. 2.

Notice that the diagram below commutes. One can check this by hand or note that it follows from the initiality of \mathbf{Syn} (Lem. 2): all the functors preserve all the structure.

$$\mathbf{Syn} \xrightarrow{(\mathrm{id}, \overrightarrow{\mathcal{D}}(-))} \mathbf{Syn} \times \mathbf{Syn}
\downarrow \mathbb{I} - \mathbb{I} \times \mathbb{I} - \mathbb{I}$$

$$\mathbf{Gl}_{\mathbb{R}} \xrightarrow{\mathrm{proj}} \mathbf{Diff} \times \mathbf{Diff}$$
(3)

We thus arrive at a restatement of the correctness theorem (Th. 1), which holds even for the extended language with variants and lists, because for any $x_1...x_n$: $\mathbf{real} \vdash t : \mathbf{real}$, the interpretations ($\llbracket t \rrbracket$, $\llbracket \overrightarrow{\mathcal{D}}(t) \rrbracket$) are in the image of the projection $\mathbf{Gl}_{\mathbb{R}} \to \mathbf{Diff} \times \mathbf{Diff}$, and hence $\llbracket \overrightarrow{\mathcal{D}}(t) \rrbracket$ is a dual numbers encoding of $\llbracket t \rrbracket$.

Correctness at all first order types, via manifolds. We now generalize Theorem 1 to hold at all first order types, not just the reals. To do this, we need to define the derivative of a smooth map between the interpretations of first order types. We do this by recalling the well known theory of manifolds and tangent bundles.

For our purposes, a smooth manifold M is a second-countable Hausdorff topological space together with a smooth atlas: an open cover \mathcal{U} together with homeomorphisms $(\phi_U: U \to \mathbb{R}^{n(U)})_{U \in \mathcal{U}}$ (called charts) such that $\phi_U^{-1}; \phi_V$ is smooth

on its domain of definition for all $U, V \in \mathcal{U}$. A function $f: M \to N$ between manifolds is smooth if $\phi_U^{-1}; f; \psi_V$ is smooth for all charts ϕ_U and ψ_V of M and N, respectively. Let us write **Man** for this category.

Our manifolds are slightly unusual because different charts in an atlas may have different finite dimension n(U). Thus we consider manifolds with dimensions that are potentially unbounded, albeit locally finite. This does not affect the theory of differential geometry as far as we need it here.

Each open subset of \mathbb{R}^n can be regarded as a manifold. This lets us regard the category of manifolds \mathbf{Man} as a full subcategory of the category of diffeological spaces. We consider a manifold $(X, \{\phi_U\}_U)$ as a diffeological space with the same carrier set X and where the plots \mathcal{P}_X^U are the smooth functions in $\mathbf{Man}(U, X)$. A function $X \to Y$ is smooth in the sense of manifolds if and only if it is smooth in the sense of diffeological spaces [15]. For the categorically minded reader, this means that we have a full embedding of \mathbf{Man} into \mathbf{Diff} . Moreover, the natural interpretation of the first order fragment of our language in \mathbf{Man} coincides with that in \mathbf{Diff} . That is, the embedding of \mathbf{Man} into \mathbf{Diff} preserves finite products and countable coproducts (hence initial algebras of polynomial endofunctors).

Proposition 2. Suppose that a type τ is first order, i.e. it is just built from reals, products, variants, and lists (or, again, arbitrary inductive types), and not function types. Then the diffeological space $\llbracket \tau \rrbracket$ is a manifold.

Proof (notes). This is proved by induction on the structure of types. In fact, one may show that every such $\llbracket \tau \rrbracket$ is isomorphic to a manifold of the form $\biguplus_{i=1}^n \mathbb{R}^{d_i}$ where the bound n is either finite or ∞ , but this isomorphism is typically not an identity function.

The constraint to first order types is necessary because, e.g. the space $[real \rightarrow real]$ is not a manifold, because of a Borsuk-Ulam argument (see Appx. A).

We recall that the derivative of any morphism $f: M \to N$ of manifolds is given as follows. For each point x in a manifold M, define the tangent space $\mathcal{T}_x M$ to be the set $\{\gamma \in \mathbf{Man}(\mathbb{R}, M) \mid \gamma(0) = x\} / \sim$ of equivalence classes $[\gamma]$ of smooth curves γ in M based at x, where we identify $\gamma_1 \sim \gamma_2$ iff $\nabla(\gamma_1; f)(0) = \nabla(\gamma_2; f)(0)$ for all smooth $f: M \to \mathbb{R}$. The tangent bundle of M is the set $\mathcal{T}(M) \stackrel{\text{def}}{=} \biguplus_{x \in M} \mathcal{T}_x(M)$. The charts of M equip $\mathcal{T}(M)$ with a canonical manifold structure. Then for smooth $f: M \to N$, the derivative $\mathcal{T}(f): \mathcal{T}(M) \to \mathcal{T}(N)$ is defined as $\mathcal{T}(f)(x, [\gamma]) \stackrel{\text{def}}{=} (f(x), [\gamma; f])$. All told, the derivative is a functor $\mathcal{T}: \mathbf{Man} \to \mathbf{Man}$.

As is standard, we can understand the tangent bundle of a composite space in terms of that of its parts.

Lemma 3. There are canonical isomorphisms $\mathcal{T}(\biguplus_{i=1}^{\infty} M_i) \cong \biguplus_{i=1}^{\infty} \mathcal{T}(M_i)$ and $\mathcal{T}(M_1 \times \ldots \times M_n) \cong \mathcal{T}(M_1) \times \ldots \times \mathcal{T}(M_n)$.

We define a canonical isomorphism $\phi_{\tau}^{\overrightarrow{\mathcal{D}}\mathcal{T}}: [\![\overrightarrow{\mathcal{D}}(\tau)]\!] \to \mathcal{T}([\![\tau]\!])$ for every type τ , by induction on the structure of types. We let $\phi_{\mathbf{real}}^{\overrightarrow{\mathcal{D}}\mathcal{T}}: [\![\overrightarrow{\mathcal{D}}(\mathbf{real})]\!] \to \mathcal{T}([\![\mathbf{real}]\!])$ be

given by $\phi_{\mathbf{real}}^{\overrightarrow{\mathcal{D}}\mathcal{T}}(x,x') \stackrel{\text{def}}{=} (x,[t\mapsto x+x't])$. For the other types, we use Lemma 3. We can now phrase correctness at all first order types.

Theorem 2 (Semantic correctness of $\overrightarrow{\mathcal{D}}$ (full)). For any ground τ , any first order context Γ and any term $\Gamma \vdash t : \tau$, the syntactic translation $\overrightarrow{\mathcal{D}}$ coincides with the tangent bundle functor, modulo these canonical isomorphisms:

$$\begin{bmatrix} \overrightarrow{\mathcal{D}}(\Gamma) \end{bmatrix} \xrightarrow{\overrightarrow{\mathcal{D}}(t)} \begin{bmatrix} \overrightarrow{\mathcal{D}}(\tau) \end{bmatrix} \\
\phi_{\Gamma}^{\overrightarrow{\mathcal{D}}\tau} \downarrow \cong \qquad \cong \downarrow_{\phi_{\tau}^{\overrightarrow{\mathcal{D}}\tau}} \\
\mathcal{T}(\llbracket \Gamma \rrbracket) \xrightarrow{\mathcal{T}(\llbracket t \rrbracket)} \mathcal{T}(\llbracket \tau \rrbracket)$$

Proof (notes). For any curve $\gamma \in \mathbf{Man}(\mathbb{R}, M)$, let $\bar{\gamma} \in \mathbf{Man}(\mathbb{R}, \mathcal{T}(M))$ be the tangent curve, given by $\bar{\gamma}(x) = (\gamma(x), [t \mapsto \gamma(x+t)])$. First, we note that a smooth map $h: \mathcal{T}(M) \to \mathcal{T}(N)$ is of the form $\mathcal{T}(g)$ for some $g: M \to N$ if for all smooth curves $\gamma: \mathbb{R} \to M$ we have $\bar{\gamma}; h = (\gamma; g): \mathbb{R} \to \mathcal{T}(N)$. This generalizes (2). Second, for any first order type τ , $S_{\llbracket\tau\rrbracket} = \{(f, \tilde{f}) \mid \tilde{f}; \phi_{\tau}^{\vec{\mathcal{D}}\mathcal{T}} = \bar{f}\}$. This is shown by induction on the structure of types. We conclude the theorem from diagram (3), by putting these two observations together.

6 A continuation-based AD algorithm

We now illustrate the flexibility of our framework by briefly describing an alternative syntactic translation $\overleftarrow{\mathcal{D}}_{\rho}$. This alternative translation uses aspects of continuation passing style, inspired by recent developments in reverse mode AD [33, 5]. In brief, $\overleftarrow{\mathcal{D}}_{\rho}$ works by $\overleftarrow{\mathcal{D}}_{\rho}(\mathbf{real}) = (\mathbf{real}*(\mathbf{real} \to \rho))$. Thus instead of using a pair of a number and its tangent, we use a pair of a number and a continuation. The answer type $\rho = \mathbf{real}^k$ needs to have the structure of a vector space, and the continuations that we consider will turn out to be linear maps. Because we work in continuation passing style, the chain rule is applied contravariantly. If the reader is familiar with reverse-mode AD algorithms, they may think of the dimension k as the number of memory cells used to store the result.

Computing the whole gradient of a term $x_1 : \mathbf{real}, ..., x_k : \mathbf{real} \vdash t : \mathbf{real}$ at once is then achieved by running $\overleftarrow{\mathcal{D}}_k(t)$ on a k-tuple of basis vectors for \mathbf{real}^k .

We define the continuation-based AD macro $\overleftarrow{\mathcal{D}}_k$ on types and terms as the unique structure preserving functor $\mathbf{Syn} \to \mathbf{Syn}$ with $\overleftarrow{\mathcal{D}}_k(\mathbf{real}) = (\mathbf{real*}(\mathbf{real} \to \mathbf{real}^k))$ and

$$\overleftarrow{\mathcal{D}}_{k}(\underline{c}) \stackrel{\text{def}}{=} \langle \underline{c}, \lambda z. \langle \underline{0}, \dots, \underline{0} \rangle \rangle \\
\overleftarrow{\mathcal{D}}_{k}(t+s) \stackrel{\text{def}}{=} \mathbf{case} \overleftarrow{\mathcal{D}}_{k}(t) \mathbf{of} \langle x, x' \rangle \to \mathbf{case} \overleftarrow{\mathcal{D}}_{k}(s) \mathbf{of} \langle y, y' \rangle \to \langle x+y, \lambda z. x' z + y' z \rangle \\
\overleftarrow{\mathcal{D}}_{k}(t*s) \stackrel{\text{def}}{=} \mathbf{case} \overleftarrow{\mathcal{D}}_{k}(t) \mathbf{of} \langle x, x' \rangle \to \mathbf{case} \overleftarrow{\mathcal{D}}_{k}(s) \mathbf{of} \langle y, y' \rangle \to \\
\langle x*y, \lambda z. x' (y*z) + y' (x*z) \rangle \\
\overleftarrow{\mathcal{D}}_{k}(\varsigma(t)) \stackrel{\text{def}}{=} \mathbf{case} \overleftarrow{\mathcal{D}}_{k}(t) \mathbf{of} \langle x, x' \rangle \to \mathbf{let} y = \varsigma(x) \mathbf{in} \langle y, \lambda z. x' (y*(1-y)*z) \rangle.$$
Here, we use sugar $x: \mathbf{real}^{k}, y: \mathbf{real}^{k} \vdash x + y \stackrel{\text{def}}{=} \mathbf{case} x \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle \to \mathbf{case} (x) = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{case} (x) \mathbf{of} \langle x_{1}, \dots, x_{k} \rangle = \mathbf{o$

case y of $\langle y_1, \ldots, y_k \rangle \to \langle x_1 + y_1, \ldots, x_k + y_k \rangle$. (We could easily expand this definition by making $\overleftarrow{\mathcal{D}}_k$ preserve all other term and type formers, as we did for $\overrightarrow{\mathcal{D}}$.) Note that the corresponding scheme for an arbitrary n-ary operation op would be (c.f. the scheme for forward AD in §4)

$$\overleftarrow{\mathcal{D}}_k(\mathsf{op}(t_1,\ldots,t_n)) \stackrel{\mathrm{def}}{=} \mathbf{case} \overleftarrow{\mathcal{D}}_k(t_1) \, \mathbf{of} \, \langle x_1, x_1' \rangle \to \ldots \to \mathbf{case} \, \overleftarrow{\mathcal{D}}_k(t_n) \, \mathbf{of} \, \langle x_n, x_n' \rangle \to \langle \mathsf{op}(x_1,\ldots,x_n), \lambda z. \sum_{i=1}^n x_i' (\partial_i \mathsf{op}(x_1,\ldots,x_n) * z) \rangle.$$

The idea is that $\overleftarrow{\mathcal{D}}_k(t)$ is a higher order function that simultaneously computes t (the forward pass) and defines as a continuation the reverse pass which computes the gradient. In order to actually run the algorithm, we need two auxiliary definitions

$$\operatorname{lamR}_{\mathbf{real}}^{k} \stackrel{\text{def}}{=} \lambda z. \operatorname{\mathbf{case}} z \operatorname{\mathbf{of}} \langle x, x' \rangle \to \operatorname{\mathbf{case}} x' \operatorname{\mathbf{of}} \langle x'_{1}, \dots, x'_{k} \rangle \to \\
\langle x, \lambda y. \langle x'_{1} * y, \dots, x'_{k} * y \rangle : \overrightarrow{\mathcal{D}}_{k}(\mathbf{real}) \to \overleftarrow{\mathcal{D}}_{k}(\mathbf{real})$$

$$\operatorname{evR}^{k} \stackrel{\text{def}}{=} \lambda z. \operatorname{\mathbf{case}} z \operatorname{\mathbf{of}} \langle x, x' \rangle \to \langle x, x' | \lambda : \overleftarrow{\mathcal{D}}_{k}(\mathbf{real}) \to \overrightarrow{\mathcal{D}}_{k}(\mathbf{real})$$

 $\operatorname{evR}_{\mathbf{real}}^k \stackrel{\text{def}}{=} \lambda z. \operatorname{case} z \operatorname{of} \langle x, x' \rangle \to \langle x, x' \, \underline{1} \rangle : \overleftarrow{\mathcal{D}}_k(\mathbf{real}) \to \overrightarrow{\mathcal{D}}_k(\mathbf{real}).$ Here, $\overrightarrow{\mathcal{D}}_k$ is a macro on types (and terms) with exactly the same inductive definition as $\overrightarrow{\mathcal{D}}$ except for the base case $\overrightarrow{\mathcal{D}}_k(\mathbf{real}) = (\mathbf{real*real}^k)$. By noting that both $\overrightarrow{\mathcal{D}}_k$ and $\overleftarrow{\mathcal{D}}_k$ preserve all type formers, we can extend these definitions to all first order types $\tau \colon z : \overrightarrow{\mathcal{D}}_k(\tau) \vdash \operatorname{lamR}_{\tau}^k(z) : \overleftarrow{\mathcal{D}}_k(\tau), \ z : \overleftarrow{\mathcal{D}}_k(\tau) \vdash \operatorname{evR}_{\tau}^k(z) : \overrightarrow{\mathcal{D}}_k(\tau).$ We can think of $\operatorname{lamR}_{\tau}^k(z)$ as encoding k tangent vectors $z : \overrightarrow{\mathcal{D}}_k(\tau)$ as a closure, so it is suitable for running $\overleftarrow{\mathcal{D}}_k(t)$ on, and $\operatorname{evR}_{\tau}^k(z)$ as actually evaluating the reverse pass defined by $z : \overleftarrow{\mathcal{D}}_k(\tau)$ and returning the result as k tangent vectors. The idea is that given some $x : \tau \vdash t : \sigma$ between first order types τ, σ , we run our continuation-based AD by running $\operatorname{evR}_{\sigma}^k(\overleftarrow{\mathcal{D}}_k(t)[^{\operatorname{lamR}_{\tau}^k(z)}/_x])$.

The correctness proof closely follows that for forward AD. In particular, one defines a binary logical relation $(\mathbf{real})^{r,k} = (\mathbb{R}, \mathbb{R} \times (\mathbb{R}^k)^{\mathbb{R}}, S^{r,k}_{\mathbf{real}})$, where $S^{r,k}_{\mathbf{real}} = \left\{ (f, x \mapsto (f(x), y \mapsto (\partial_1 f(x) * y, \dots, \partial_k f(x) * y))) \mid f \in \mathcal{P}^{\mathbb{R}^k}_{\mathbb{R}} \right\}$, on the plots $\mathcal{P}^{\mathbb{R}^k}_{\mathbb{R}} \times \mathcal{P}^{\mathbb{R}^k}_{\mathbb{R} \times ((\mathbb{R}^k)^{\mathbb{R}})}$ and verifies that $[\underline{c}] \times [\overleftarrow{\mathcal{D}}_k(\underline{c})]$, $[x + y] \times [\overleftarrow{\mathcal{D}}_k(x + y)]$, $[x*y] \times [\overleftarrow{\mathcal{D}}_k(x*y)]$ and $[\varsigma(x)] \times [\overleftarrow{\mathcal{D}}_k(\varsigma(x))]$ respect this logical relation. It follows that this relation extends to a functor $(-)^{r,k} : \mathbf{Syn} \to \mathbf{Gl}_{\mathbb{R}^k}$ such that id $\times \overleftarrow{\mathcal{D}}_k$ factors over $(-)^{r,k}$, implying the correctness of the continuation-based AD by the following lemma.

Lemma 4. For all first order types τ (i.e. types not involving function types), we have that $[evR_{\tau}^{k}(lamR_{\tau}^{k}(t))] = [t]$.

Proof (notes). This follows by an induction on the structure of τ . The idea is that $\operatorname{lam} R_{\tau}^k$ embeds reals into function spaces as linear maps, which is undone by $\operatorname{evR}_{\tau}^k$ by evaluating the linear maps at $\underline{1}$.

To phrase correctness, in this setting, however, we need a few definitions. Keeping in mind the canonical projection $\mathcal{T}(M) \to M$, we define $\mathcal{T}^k(M)$ as the k-fold categorical pullback (fibre product) $\mathcal{T}(M) \times_M \ldots \times_M \mathcal{T}(M)$. To be explicit, $\mathcal{T}^k_x M$ consists of k-tuples of tangent vectors at the base point x. Again, \mathcal{T}^k extends to a functor $\mathbf{Man} \to \mathbf{Man}$ by defining $\mathcal{T}^k(f)(x,(v_1,\ldots,v_k)) \stackrel{\mathrm{def}}{=} (f(x),(\mathcal{T}_x(f)(v_1),\ldots,\mathcal{T}_x(f)(v_k)))$. As \mathcal{T}^k preserves countable coproducts and

finite products (like \mathcal{T}), it follows that the isomorphisms $\phi_{\tau}^{\overrightarrow{\mathcal{D}}\mathcal{T}}$ generalize to canonical isomorphisms $\phi_{\tau,k}^{\overrightarrow{\mathcal{D}}\mathcal{T}}: \llbracket \overrightarrow{\mathcal{D}}_k(\tau) \rrbracket \to \mathcal{T}^k(\llbracket \tau \rrbracket)$ for first order types τ . This leads to the following correctness statement for continuation-based AD.

Theorem 3 (Semantic correctness of $\overleftarrow{\mathcal{D}}_k$). For any ground τ , any first order context Γ and any term $\Gamma \vdash t : \tau$, syntactic translation $t \mapsto \operatorname{evR}_{\tau}^k(\overleftarrow{\mathcal{D}}_k(t)[^{\operatorname{lamR}_{\Gamma}^k(z)}/...])$ coincides with the tangent bundle functor, modulo these canonical isomorphisms:

$$\begin{bmatrix} \overrightarrow{\mathcal{D}}_{k}(\Gamma) \end{bmatrix} \xrightarrow{\mathbb{I} \operatorname{lamR}_{\Gamma}^{k}; \overleftarrow{\mathcal{D}}_{k}(t); \operatorname{evR}_{\tau}^{k} \end{bmatrix}}
\begin{bmatrix} \overrightarrow{\mathcal{D}}_{k}(\tau) \end{bmatrix} \\
\phi_{\Gamma, k}^{\overrightarrow{\mathcal{D}}, \tau} \swarrow \cong \qquad \qquad \cong \bigvee_{\tau} \phi_{\tau, k}^{\overrightarrow{\mathcal{D}}, \tau} \\
\mathcal{T}^{k}(\llbracket \Gamma \rrbracket) \xrightarrow{\mathcal{T}^{k}(\llbracket t \rrbracket)} \mathcal{T}^{k}(\llbracket \tau \rrbracket)$$

For example, when $\tau = \mathbf{real}$ and $\Gamma = x, y : \mathbf{real}$, we can run our continuation-based AD to compute the gradient of a program $x, y : \mathbf{real} \vdash t : \mathbf{real}$ at values x = V, y = W by evaluating

$$\text{evR}_{\mathbf{real}}^2(\overleftarrow{\mathcal{D}}_2(t)[^{(\text{lamR}_{x:\mathbf{real}}^2v)}/_x,^{(\text{lamR}_{y:\mathbf{real}}^2w)}/_y])[^{\langle V,\langle\underline{1},\underline{0}\rangle\rangle}/_v,^{\langle W,\langle\underline{0},\underline{1}\rangle\rangle}/_w].$$

Indeed,

$$\begin{split} & [[\operatorname{evR}^2_{\mathbf{real}} \left(\overleftarrow{\mathcal{D}}_2(t) [^{(\operatorname{lamR}^2_{x:\mathbf{real}} \, v)} /_x, (^{\operatorname{lamR}^2_{y:\mathbf{real}} \, w)} /_y] \right) [\langle V, \langle \underline{1}, \underline{0} \rangle \rangle /_v, \langle W, \langle \underline{0}, \underline{1} \rangle \rangle /_w]]] = \\ & ([\![t]\!] ([\![V]\!], [\![W]\!]), \partial_1 [\![t]\!] ([\![V]\!], [\![W]\!]), \partial_2 [\![t]\!] ([\![V]\!], [\![W]\!])). \end{split}$$

7 Discussion and future work

Summary. We have shown that diffeological spaces provide a denotational semantics for a higher order language with variants and inductive types (§3,4). We have used this to show correctness of a simple AD translation (Thm. 1, Thm. 2). But the method is not tied to this specific translation, as we illustrated in Section 6.

The structure of our elementary correctness argument for Theorem 1 is a typical logical relations proof. As explained in Section 5, this can equivalently be understood as a denotational semantics in a new kind of space obtained by categorical gluing.

Overall, then, there are two logical relations at play. One is in diffeological spaces, which ensures that all definable functions are smooth. The other is in the correctness proof (equivalently in the categorical gluing), which explicitly tracks the derivative of each function, and tracks the syntactic AD even at higher types.

Connection to the state of the art in AD implementation. As is common in denotational semantics research, we have here focused on an idealized language and simple translations to illustrate the main aspects of the method. There are a number of points where our approach is simplistic compared to the advanced current practice, as we now explain.

Representation of vectors. In our examples we have treated n-vectors as tuples of length n. This style of programming does not scale to large n. A better solution would be to use array types, following [30]. Our categorical semantics and correctness proofs straightforwardly extend to cover them, in a similar way to our treatment of lists.

Efficient forward-mode AD. For AD to be useful, it must be fast. The syntactic translation \overrightarrow{D} that we use is the basis of an efficient AD library [30]. However, numerous optimizations are needed, ranging from algebraic manipulations, to partial evaluations, to the use of an optimizing C compiler. A topic for future work would be to validate some of these manipulations using our semantics. The resulting implementation is performant in experiments [30].

Efficient reverse-mode AD. Our sketch of continuation-based AD is primarily intended to emphasise that our denotational approach is not tied to any specific translation \overrightarrow{D} . Nonetheless, it is worth noting that this algorithm shares similarities with advanced reverse-mode implementations: (1) it calculates derivatives in a (contravariant) "reverse pass" in which derivatives of operations are evaluated in the reverse order compared to their order in calculating the function value; (2) it can be used to calculate the full gradient of a function $\mathbb{R}^n \to \mathbb{R}$ in a single reverse pass (while n passes of fwd AD would be necessary). However, it lacks important optimizations and the continuation scales with the size of the input n where it should scale with the size of the output. This adds an important overhead, as pointed out in [25]. Speed being the main attraction of reverse-mode AD, its implementations tend to rely on mutable state, control operators and/or staging [25, 6, 33, 5], which we have not considered here.

Other language features. The idealized languages that we considered so far do not touch on several useful language constructs. For example: the use of functions that are partial (such as division) or partly-smooth (such as RelU); phenomena such as iteration, recursion; and probabilities. There are suggestions that the denotational approach using diffeological spaces can be adapted to these features using standard categorical methods. We leave this for future work.

Acknowledgements. We have benefited from discussing this work with many people, including B. Pearlmutter, O. Kammar, C. Mak, L. Ong, G. Plotkin, A. Shaikhha, J. Sigal, and others. Our work is supported by the Royal Society and by a Facebook Research Award. In the course of this work, MV has also been employed at Oxford (EPSRC Project EP/M023974/1) and at Columbia in the Stan development team. This project has received funding from the European Unions Horizon 2020 research and innovation programme under the Marie Skodowska-Curie grant agreement No. 895827.

References

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al.: Tensorflow: A system for large-scale machine learning. In: 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16). pp. 265–283 (2016)
- 2. Abadi, M., Plotkin, G.D.: A simple differentiable programming language. In: Proc. POPL 2020. ACM (2020)
- 3. Baez, J., Hoffnung, A.: Convenient categories of smooth spaces. Transactions of the American Mathematical Society **363**(11), 5789–5825 (2011)
- 4. Barthe, G., Crubillé, R., Lago, U.D., Gavazzo, F.: On the versatility of open logical relations: Continuity, automatic differentiation, and a containment theorem. In: Proc. ESOP 2020. Springer (2020), to appear
- Brunel, A., Mazza, D., Pagani, M.: Backpropagation in the simply typed lambdacalculus with linear negation. In: Proc. POPL 2020 (2020)
- 6. Carpenter, B., Hoffman, M.D., Brubaker, M., Lee, D., Li, P., Betancourt, M.: The Stan math library: Reverse-mode automatic differentiation in C++. arXiv preprint arXiv:1509.07164 (2015)
- Christensen, J.D., Wu, E.: Tangent spaces and tangent bundles for diffeological spaces. arXiv preprint arXiv:1411.5425 (2014)
- 8. Cockett, J.R.B., Cruttwell, G.S.H., Gallagher, J., Lemay, J.S.P., MacAdam, B., Plotkin, G.D., Pronk, D.: Reverse derivative categories. In: Proc. CSL 2020 (2020)
- Cruttwell, G., Gallagher, J., MacAdam, B.: Towards formalizing and extending differential programming using tangent categories. In: Proc. ACT 2019 (2019)
- Duchi, J., Hazan, E., Singer, Y.: Adaptive subgradient methods for online learning and stochastic optimization. Journal of Machine Learning Research 12(Jul), 2121– 2159 (2011)
- 11. Ehrhard, T., Regnier, L.: The differential lambda-calculus. Theoretical Computer Science **309**(1-3), 1–41 (2003)
- 12. Elliott, C.: The simple essence of automatic differentiation. Proceedings of the ACM on Programming Languages 2(ICFP), 70 (2018)
- 13. Fong, B., Spivak, D., Tuyéras, R.: Backprop as functor: A compositional perspective on supervised learning. In: 2019 34th Annual ACM/IEEE Symposium on Logic in Computer Science (LICS). pp. 1–13. IEEE (2019)
- 14. Hoffman, M.D., Gelman, A.: The No-U-Turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo. Journal of Machine Learning Research 15(1), 1593–1623 (2014)
- 15. Iglesias-Zemmour, P.: Diffeology. American Mathematical Soc. (2013)
- Johnstone, P.T., Lack, S., Sobocinski, P.: Quasitoposes, quasiadhesive categories and Artin glueing. In: Proc. CALCO 2007 (2007)
- 17. Kiefer, J., Wolfowitz, J., et al.: Stochastic estimation of the maximum of a regression function. The Annals of Mathematical Statistics **23**(3), 462–466 (1952)
- 18. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
- 19. Kucukelbir, A., Tran, D., Ranganath, R., Gelman, A., Blei, D.M.: Automatic differentiation variational inference. The Journal of Machine Learning Research 18(1), 430–474 (2017)
- 20. Liu, D.C., Nocedal, J.: On the limited memory BFGS method for large scale optimization. Mathematical programming 45(1-3), 503–528 (1989)

- 21. Mak, C., Ong, L.: A differential-form pullback programming language for higher-order reverse-mode automatic differentiation (2020), arxiv:2002.08241
- 22. Manzyuk, O.: A simply typed λ -calculus of forward automatic differentiation. In: Proc. MFPS 2012 (2012)
- Mitchell, J.C., Scedrov, A.: Notes on sconing and relators. In: International Workshop on Computer Science Logic. pp. 352–378. Springer (1992)
- 24. Neal, R.M., et al.: MCMC using Hamiltonian dynamics. Handbook of Markov Chain Monte Carlo 2(11), 2 (2011)
- 25. Pearlmutter, B.A., Siskind, J.M.: Reverse-mode AD in a functional framework: Lambda the ultimate backpropagator. ACM Transactions on Programming Languages and Systems (TOPLAS) **30**(2), 7 (2008)
- Pitts, A.M.: Categorical logic. Tech. rep., University of Cambridge, Computer Laboratory (1995)
- Plotkin, G.D.: Some principles of differential programming languages (2018), invited talk, POPL 2018
- 28. Qian, N.: On the momentum term in gradient descent learning algorithms. Neural networks **12**(1), 145–151 (1999)
- 29. Robbins, H., Monro, S.: A stochastic approximation method. The annals of mathematical statistics pp. 400–407 (1951)
- 30. Shaikhha, A., Fitzgibbon, A., Vytiniotis, D., Peyton Jones, S.: Efficient differentiable programming in a functional array-processing language. Proceedings of the ACM on Programming Languages 3(ICFP), 97 (2019)
- 31. Souriau, J.M.: Groupes différentiels. In: Differential geometrical methods in mathematical physics, pp. 91–128. Springer (1980)
- 32. Stacey, A.: Comparative smootheology. Theory Appl. Categ. 25(4), 64–117 (2011)
- 33. Wang, F., Wu, X., Essertel, G., Decker, J., Rompf, T.: Demystifying differentiable programming: Shift/reset the penultimate backpropagator. Proceedings of the ACM on Programming Languages 3(ICFP) (2019)

Open Access This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (http://creativecommons.org/licenses/by/4.0/), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.



A CartSp and Man are not cartesian closed categories

Lemma 5. There is no continuous injection $\mathbb{R}^{d+1} \to \mathbb{R}^d$.

Proof. If there were, it would restrict to a continuous injection $S^d \to \mathbb{R}^d$. The Borsuk-Ulam theorem, however, tells us that every continuous $f: S^d \to \mathbb{R}^d$ has some $x \in S^d$ such that f(x) = f(-x), which is a contradiction.

Let us define the terms:

$$x_0 : \mathbf{real}, \dots, x_n : \mathbf{real} \vdash t_n = \lambda y. x_0 + x_1 * y + \dots + x_n * y^n : \mathbf{real} \rightarrow \mathbf{real}$$

Assuming that $\mathbf{CartSp}/\mathbf{Man}$ is cartesian closed, observe that these get interpreted as injective continuous (because smooth) functions $\mathbb{R}^n \to \llbracket \mathbf{real} \to \mathbf{real} \rrbracket$ in \mathbf{CartSp} and \mathbf{Man} .

Theorem 4. CartSp is not cartesian closed.

Proof. In case **CartSp** were cartesian closed, we would have $[\![\mathbf{real} \to \mathbf{real}]\!] = \mathbf{real}^n$ for some n. Then, we would get, in particular a continuous injection $[\![t_{n+1}]\!] : \mathbb{R}^{n+1} \to \mathbb{R}^n$, which contradicts Lemma 5.

Theorem 5. Man is not cartesian closed.

Proof. Observe that we have $\iota_n: \mathbb{R}^n \to \mathbb{R}^{n+1}$; $\langle a_0, \dots, a_n \rangle \mapsto \langle a_0, \dots, a_n, 0 \rangle$ and that ι_n ; $\llbracket t_{n+1} \rrbracket = \llbracket t_n \rrbracket$. Let us write A_n for the image of $\llbracket t_n \rrbracket$ and $A = \cup_{n \in \mathbb{N}} A_n$. Then, A_n is connected because it is the continuous image of a connected set. Similarly, A is connected because it is the non-disjoint union of connected sets. This means that A lies in a single connected component of $\llbracket \mathbf{real} \to \mathbf{real} \rrbracket$, which is a manifold with some finite dimension, say d.

Take some $x \in \mathbb{R}^{d+1}$ (say, 0), take some open d-ball U around $\llbracket t_{d+1} \rrbracket(x)$, and take some open d+1-ball V around x in $\llbracket t_{d+1} \rrbracket^{-1}(U)$. Then, $\llbracket t_{d+1} \rrbracket$ restricts to a continuous injection from V to U, or equivalently, \mathbb{R}^{d+1} to \mathbb{R}^d , which contradicts Lemma 5.