



Automatic Differentiation in PCF

Damiano Mazza, Michele Pagani

► To cite this version:

Damiano Mazza, Michele Pagani. Automatic Differentiation in PCF. Proceedings of the ACM on Programming Languages, 2021, 5, pp.1 - 27. 10.1145/3434309 . hal-03379031

HAL Id: hal-03379031

<https://hal.science/hal-03379031>

Submitted on 14 Oct 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Automatic Differentiation in PCF

DAMIANO MAZZA, CNRS, France

MICHELE PAGANI, Université de Paris, France

We study the correctness of automatic differentiation (AD) in the context of a higher-order, Turing-complete language (PCF with real numbers), both in forward and reverse mode. Our main result is that, under mild hypotheses on the primitive functions included in the language, AD is almost everywhere correct, that is, it computes the derivative or gradient of the program under consideration *except* for a set of Lebesgue measure zero. Stated otherwise, there are inputs on which AD is incorrect, but the probability of randomly choosing one such input is zero. Our result is in fact more precise, in that the set of failure points admits a more explicit description: for example, in case the primitive functions are just constants, addition and multiplication, the set of points where AD fails is contained in a countable union of zero sets of polynomials.

CCS Concepts: • **Theory of computation** → **Program semantics**; *Theory and algorithms for application domains*.

Additional Key Words and Phrases: Differentiable Programming, Lambda-Calculus, Linear Logic

ACM Reference Format:

Damiano Mazza and Michele Pagani. 2021. Automatic Differentiation in PCF. *Proc. ACM Program. Lang.* 5, POPL, Article 28 (January 2021), 27 pages. <https://doi.org/10.1145/3434309>

1 INTRODUCTION

Automatic differentiation (AD) provides efficient methods for computing the derivative (or, more generally, the gradient or Jacobian) of a function specified by a computer program. Since computing derivatives is a key ingredient in the resolution of all sorts of optimization problems, it is not surprising that AD grew into a large field with applications to a host of scientific domains, most notably machine learning [Baydin et al. 2018].

Traditionally, AD focused on first-order imperative programs and, although its techniques allow the presence of flow control instructions and loops [Beck and Fischer 1994; Joss 1976; Speelpenning 1980], its scope was often limited to straight-line programs (also known as *computational graphs* [Goodfellow et al. 2016]), which were enough for most practical purposes, such as expressing neural networks. After the advances in deep learning of the last years, this is no longer the case: neural network architectures are now “dynamic”, in the sense that the input may influence the shape of the net, and expressing such architectures requires resorting a priori to the full power of a modern programming language, yielding what some have called *differentiable programming* [LeCun 2018]. This evolution of deep learning spurred the rapid development of differentiable programming frameworks [Abadi et al. 2016; Paszke et al. 2017] and, at the same time, received much attention in programming languages (PL) research for establishing its theoretical foundations [Abadi and Plotkin 2020; Brunel et al. 2020; Elliott 2018; Huot et al. 2020; Shaikhha et al. 2019; Wang et al. 2019].

Authors’ addresses: Damiano Mazza, CNRS, UMR 7030, LIPN, Université Sorbonne Paris Nord, France, Damiano.Mazza@lipn.univ-paris13.fr; Michele Pagani, IRIF UMR CNRS 8243, Université de Paris, France, pagani@irif.fr.



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2021 Copyright held by the owner/author(s).

2475-1421/2021/1-ART28

<https://doi.org/10.1145/3434309>

From the viewpoint of PL theory, AD methods boil down to *program transformations*: writing \mathbb{R} and R for the set and type of real numbers, respectively, we receive as input a program $M : \mathbb{R}^n \rightarrow \mathbb{R}$ computing a (possibly partial) function $\llbracket M \rrbracket : \mathbb{R}^n \rightarrow \mathbb{R}$ whose gradient $\nabla \llbracket M \rrbracket$ exists in a set $d(M) \subseteq \mathbb{R}^n$ (called the *domain of differentiability*), and we must output another program $grad(M)$ computing $\nabla \llbracket M \rrbracket$. The crucial features that one typically asks of such transformations are:

- (i) **efficiency**: asymptotically, evaluating $grad(M)$ is not more costly than evaluating M ;
- (ii) **soundness**: $grad(M)(r)$ evaluates to $\nabla \llbracket M \rrbracket(r)$ for all $r \in d(M)$.

Notice that there is a tension between efficiency and soundness: implementing the definition of derivative as a limit gives a trivially sound (to an arbitrary degree of precision) but unacceptably inefficient method. Conversely, more efficient transformations tend to be more complex (e.g., reverse mode is more complex than forward mode, see below) and their soundness more difficult to prove. Also observe that we are only interested in the correctness of the result when $\nabla \llbracket M \rrbracket$ is defined; in case $r \notin d(M)$, the evaluation of $grad(M)(r)$ may give anything, including (but not necessarily!) divergence.

Another highly desirable feature of $grad$ is *modularity*: if P is a subprogram of M then $grad(P)$ is a subprogram of $grad(M)$ or, if this is not literally the case, the computation of the former may be reused in computing the latter. Indeed, it has been known from the early days of AD that modularity offers a path to attaining both efficiency and soundness: the program M is decomposed into elementary blocks whose gradients are immediately computable, and $grad(M)$ is obtained by assembling these transformed blocks following the structure of M . In this way, the execution of $grad(M)$ mimics that of M , yielding efficiency, and soundness relies on the so-called *chain rule* of calculus, which assures us that the derivative of a compound function may be expressed in terms of the derivative of its components. Furthermore, one sees that there are two “dual” ways of assembling the transformed blocks to form $grad(M)$: a covariant way, yielding *forward mode* AD, and a contravariant way, yielding *reverse mode* AD (this will be explained in Sect. 2.2).

The theory of AD transformations has by now been developed to considerable depth by several authors: Pearlmutter and Siskind first pointed out that reverse mode AD, commonly known as *backpropagation*, may be naturally expressed in terms of higher-order programs, and used this idea to develop a differentiable variant of Scheme [Pearlmutter and Siskind 2008]; more recently, Elliott emphasized functoriality as a systematic way of understanding the modular nature of AD transformations [Elliott 2018]; the work [Wang et al. 2019] introduced Lantern, a fully general differentiable programming framework in which the notion of delimited continuation is used to correctly handle memory updates during backpropagation; finally, Brunel, Mazza and Pagani showed that the continuation-passing machinery at work in [Wang et al. 2019] (and, implicitly, in [Pearlmutter and Siskind 2008]) may be understood in terms of linear negation (in the sense of Girard’s linear logic), giving a purely functional transformation for reverse mode AD and a conceptually clean analysis of its efficiency in terms of a “linear factoring” evaluation rule [Brunel et al. 2020]. On the semantics side, Abadi and Plotkin studied denotational semantics for a first order differentiable language [Abadi and Plotkin 2020] and Huot, Staton and Vákár gave a uniform approach to proving soundness of AD (forward and reverse) for simply-typed programs based on a diffeology semantics [Huot et al. 2020].

Nevertheless, an analysis of soundness of AD transformations for a fully general programming language is currently missing: the above-mentioned work is either fully general but lacks soundness proofs [Pearlmutter and Siskind 2008; Wang et al. 2019] or proves soundness in a restricted setting (first order [Abadi and Plotkin 2020] or simply-typed λ -calculi [Barthe et al. 2020; Brunel et al. 2020; Huot et al. 2020]). Filling this gap is precisely the contribution of the present paper: we study the soundness of AD transformations in the setting of the (idealized) functional programming language

PCF_R, a variant with real numbers of Plotkin's famous Turing-complete language [Plotkin 1977]. For forward mode, we use the standard transformation described for instance in [Wang et al. 2019]. For reverse mode, since PCF_R is purely functional, we consider an extension of the transformation introduced in [Brunel et al. 2020], albeit simplified in that we leave linearity aside, since that is only needed for efficiency and here we are merely interested in soundness. The extension, which is a contribution of this paper in its own right, concerns conditional statements and fixpoints, which are not dealt with in *loc. cit.*

The first relevant observation is that, in presence of conditionals, soundness in the sense of statement (ii) above actually fails. Consider the program

$$\text{SillyId} \quad := \quad \lambda x^{\mathbb{R}}. \text{if } x = 0 \text{ then } 0 \text{ else } x.$$

We clearly have that SillyId : $\mathbb{R} \rightarrow \mathbb{R}$ and that $\llbracket \text{SillyId} \rrbracket$ is the identity function. We therefore expect $\text{grad}(\text{SillyId})$ to compute the constant function 1. And yet, by modularity/functoriality, AD transformations will give something like

$$\text{grad}(\text{SillyId}) \quad = \quad \lambda x^{\mathbb{R}}. \text{if } x = 0 \text{ then } 0 \text{ else } 1,$$

which obviously gives the wrong result for $x = 0$. This phenomenon, which is well known in the AD community [Beck and Fischer 1994], is due to functoriality turning a syntactic discontinuity into a semantic one. Notice that, although the above example is indeed quite silly, similar situations may happen in a non-trivial neural network with rectified linear unit activation: if

$$\text{ReLU} \quad := \quad \lambda x^{\mathbb{R}}. \text{if } x \leq 0 \text{ then } 0 \text{ else } x,$$

then $\text{ReLU}(x) - \text{ReLU}(-x)$ behaves exactly as $\text{SillyId}(x)$. Also, using recursive definitions, it is easy to obtain programs on which AD fails on infinitely many inputs, even uncountably many in case of programs of type $\mathbb{R}^n \rightarrow \mathbb{R}$ with $n > 1$ (simply consider $\lambda x^{\mathbb{R}}. \lambda y^{\mathbb{R}}. \text{if } x \cdot y = 0 \text{ then } 0 \text{ else } x \cdot y$).

So, to each given PCF_R program $M : \mathbb{R}^n \rightarrow \mathbb{R}$, we may assign a set $\text{Fail}(M) \subseteq d(M)$ of points on which AD is unsound. Notice once again that we disregard what lies outside of $d(M)$, where $\text{grad}(M)$ is free to behave arbitrarily. For instance, $\text{grad}(\text{ReLU})$ is something like $\lambda x^{\mathbb{R}}. \text{if } x \leq 0 \text{ then } 0 \text{ else } 1$, which evaluates to 0 when $x = 0$, even though $\llbracket \text{ReLU} \rrbracket$ is not differentiable in 0. Morally, we cannot say that AD is “wrong” when there is no “right” value to compare it to.

After toying with more examples, one is led to conjecture that $\text{Fail}(M)$ is always of measure zero (in the sense of the standard Lebesgue measure on \mathbb{R}^n), so one may hope to establish an “almost-everywhere” relaxation of (ii):

- (ii') **ae-soundness:** $\text{grad}(M)(\mathbf{r})$ evaluates to $\nabla \llbracket M \rrbracket(\mathbf{r})$ for all $\mathbf{r} \in d(M)$ *except* on a set of measure zero.

This is exactly the main result of our paper. Let us stress that, considering that “full” soundness is impossible, such a result is quite meaningful in practice because of the link between the Lebesgue measure and the standard understanding of randomness on \mathbb{R}^n . In typical deep learning applications, weights are initialized “at random” and later updated via gradient descent in order to minimize a loss function. Technically, this means that the weights evolve following some standard probability distribution, which always arises from integrating a probability density function with respect to the Lebesgue measure. So, an informal way of stating (ii') is that *AD almost never fails*, in the sense that, according to the standard definition of probability, the likelihood of computing wrong derivatives during gradient descent is zero.

The main result itself is articulated in Theorem 33 and Theorem 42. The first result takes care of soundness proper: we define, for any given PCF_R program $M : \mathbb{R}^n \rightarrow \mathbb{R}$ whose domain (*i.e.*, the inputs on which it converges) is $\Downarrow M$, a set $S(M) \subseteq \Downarrow M$ of *stable points*, and Theorem 33 affirms

that statement (ii) holds on $d(M) \cap S(M)$. Then, we establish in Theorem 42 that the set $\Downarrow M \setminus S(M)$ of *unstable points* of M is of measure zero. Since $d(M) \subseteq \Downarrow M$, this proves statement (ii').

The intuition behind stable points is the following. Take $M : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{r} \in \mathbb{R}^n$ and trace the execution of $M(\mathbf{r})$. This means, in particular, unfolding the recursive definitions in M and choosing, for each instance of a conditional statement of M , the “then” or “else” branch. One obtains thus a program with no conditionals and no fixpoints, *i.e.*, a *simply-typed λ -term*, which is said to *trace* $M(\mathbf{r})$. Such a program depends of course on \mathbf{r} . If, however, there exists a simply-typed λ -term t and an open neighborhood $U \subseteq \mathbb{R}^n$ of \mathbf{r} (in the standard topology) such that $t(\mathbf{r}')$ traces $M(\mathbf{r}')$ for all $\mathbf{r}' \in U$, then \mathbf{r} is *stable* (Definition 26). For instance, any $r \neq 0$ is stable for the ReLU program given above: there always exists an open interval I around r such that either $\text{Zero} := \lambda x^{\mathbb{R}}.0$ or $\text{Id} := \lambda x^{\mathbb{R}}.x$ traces ReLU on I , depending on whether $r < 0$ or $r > 0$, respectively. On the other hand, 0 is an unstable point of ReLU: any open interval around 0 must contain negative points, on which ReLU is traced by Zero, and positive points, on which ReLU is traced by Id, and of course $\text{Zero} \neq \text{Id}$.

The proof of Theorem 33 uses stability to reduce the soundness of AD on $\text{PCF}_{\mathbb{R}}$ to the soundness of AD on simply-typed λ -terms, which may be established in various ways [Brunel et al. 2020; Huot et al. 2020]. The idea is simple: if $\mathbf{r} \in S(M)$, then M “behaves like” a simply-typed λ -term t in an open neighborhood of \mathbf{r} , and we know that AD works for t everywhere, so it “must” work for M on \mathbf{r} . Although intuitively clear, the actual argument is surprisingly subtle. First, the definition of trace (Definition 25) is not obvious, due to non-uniformity issues introduced by higher types: two copies of the same higher-order subterm may be traced in different ways, as explained in the example given at the beginning of Sect. 3.1. Second, knowledge of the correctness of $\text{grad}(t)(\mathbf{r})$ does not immediately imply the correctness of $\text{grad}(M)(\mathbf{r})$ and some non-trivial work is needed to show that they behave similarly.

The measure-zero bound on unstable points (Theorem 42), albeit obtained via a standard logical predicate argument, also requires non-trivial elements, most notably the notion of *complete quasi-continuity*, which is needed to account for the behavior of unstable points under composition, and the related notion of *quasivariety*. Although the exact meaning of these notions depends on the choice of primitives of $\text{PCF}_{\mathbb{R}}$ (*i.e.*, the basic real functions included in the language), under mild assumptions a quasivariety is always of measure zero, and Theorem 42 states precisely that $\text{Fail}(M)$ is a quasivariety. For example, when the primitives are just constants, addition and multiplication, quasivarieties are arbitrary subsets of countable unions of zero sets of polynomials. Furthermore, our Lemma 41 implies properties of $\text{PCF}_{\mathbb{R}}$ -definable functions which, as far as we can tell, were previously unknown. For example, if $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is definable in $\text{PCF}_{\mathbb{R}}$, and if $U \subseteq \mathbb{R}$ is open, then in general the border of $f^{-1}(U)$ (*i.e.*, $f^{-1}(U)$ minus its interior) is not empty because conditionals introduce discontinuities, but *it is always a quasivariety*. Similarly, *the set of zeros of f is always the disjoint union of an open set and a quasivariety*.

Related work. We already mentioned some of the relevant previous work at the interface between AD and PL theory and stressed that our contribution here is to study the soundness of AD in a fully general setting (higher-order, Turing-complete language), something which, as far as we know, was lacking.

We also mentioned that the unsoundness of AD in presence of conditional statements is well known [Beck and Fischer 1994]. Surprisingly, though, recent PL work on the subject acknowledges this problem only sporadically, *e.g.* [Abadi and Plotkin 2020]. The solution proposed therein is restricting to *continuous* Boolean conditions, meaning that the inverse images of the two Boolean values along such conditions are open. For example, testing for zero or for non-positivity are not continuous, because the inverse image of true is the singleton $\{0\}$ or $]-\infty, 0]$, respectively, which

are not open. With this limitation, the authors prove statement (ii) for a Turing-complete first-order language. The benefit of Abadi and Plotkin's approach is allowing a denotational semantics modeling the *grad* operator, but it has the drawback of representing standard total functions with programs diverging on singularities (for example, ReLU yields a program diverging in 0), which is somewhat unexpected. We discuss this further at the end of Sect. 2.2 and in Sect. 5.

Our approach, in the wake of a large part of the AD literature, is to stick to the standard semantics and provide a bound on the unsoundness of AD, as precise as possible. This approach dates back to the Seventies, as far as we know to Joss's Ph.D. thesis [Joss 1976], who proved statement (ii') for forward mode AD in the context of an imperative language with variable assignments, basic arithmetic functions (sum, multiplication, division), conditional statements and *gotos*. So our result may be seen as an extension Joss's theorem in several directions: to a higher-order language; to a wider set of primitive functions (as long as they form an *admissible clone*, Definition 11); and to reverse mode AD. Additionally, Theorem 42 is more precise than (ii'), because it characterizes the set of failure points $\text{Fail}(M)$ as a quasivariety, which is a rather special example of negligible set. Some discussion about this point is given in Sect. 5, in particular the proof of Theorem 42 hints to methods for automatically computing at least some overapproximation of $\text{Fail}(M)$ statically from the structure of M .

From a broader perspective, variants of PCF with real numbers similar to the one studied here have been considered in the literature, e.g. [Escardó 1996] and [Di Gianantonio and Edalat 2013]. The latter actually also considers AD, but it is not about correctness and is quite different in spirit, being more focused on denotational semantics. There is also a recent line of work whose goal is to understand the non-differentiable points of program-defined functions, such as [Mak et al. 2020; Zhou et al. 2019], including in the context of AD [Lee et al. 2020], where it is a natural and important question [Griewank and Walther 2008]. Non-differentiability and unsoundness of AD have an important point in common: they are both introduced by conditionals. This explains why some notions used in our paper also crop up in the study of non-differentiability, such as zero sets of analytic functions [Zhou et al. 2019]. However, let us underline that the two issues are orthogonal: from our perspective, PCF-definable functions might as well have been differentiable *everywhere* (as in the Sillyld example), what matters is that AD still makes mistakes and we wish to understand them. Whether our techniques also yield tools for describing the set of non-differentiable points of PCF-definable functions and, in that case, exactly how they relate to the above-mentioned work is an interesting question which we leave for the future.

Finally, let us mention that our notion of stable point (Definition 26), which is new as far as we know, is based on a concept of trace (Definition 25) belonging to the same circle of ideas as Ehrhard and Regnier's Taylor expansion [Ehrhard and Regnier 2006, 2008] and the modern understanding of intersection types in the spirit of Mazza, Pellissier and Vial's work [Mazza 2017; Mazza et al. 2018]. This extremely general perspective allows the definition of finitary approximations of programs at the level of the operational semantics, rather than denotational, as was the case traditionally. Our proof techniques are therefore not *ad hoc* for our current purposes and may be expected to have applications beyond the present paper.

Contents of the paper. Sect. 2 introduces the language PCF_R with its rewriting relation (Fig. 1) and the AD transformations (Fig. 2 and Equations (9), (10)). Our results are quite general and do not depend on a specific operational semantics but apply to a wide family of them (Proposition 10), including the standard ones (call-by-value, call-by-name, etc.). We also introduce admissibility of primitives (Definition 11) and the associated topological notions, among which quasivarieties (Definition 13). Sect. 3 and Sect. 4 are the heart of the paper, containing the proofs of Theorem 33 and Theorem 42, respectively, as described above. Sect. 5 concludes the paper, discussing the results.

Notations. We write $f : A \rightarrow B$ to say that f is a partial function from a set A to a set B . In that case, $\Downarrow f$ will denote the subset of A on which f is defined. Given two partial functions $f, g : A \rightarrow B$ and $a \in A$, the equality $f(a) = g(a)$ means that either both $f(a)$ and $g(a)$ are defined and equal, or that both $f(a)$ and $g(a)$ are undefined. The notation $f(a) \neq g(a)$ of course is understood as the logical negation of that. Given a subset $A' \subseteq A$, we write $f|_{A'}$ for the restriction of f to A' . If A is endowed with a complete measure λ (typically A is \mathbb{R} and λ is the Lebesgue measure), then we say that f and g are *almost everywhere equal*, and we write $f \sim g$, if $\lambda(\{a \in A \mid f(a) \neq g(a)\}) = 0$. This is equivalent to the existence of two subsets A', Z of A such that $\Downarrow f \cup \Downarrow g \subseteq A' \cup Z$, $\lambda(Z) = 0$ and $f|_{A'} = g|_{A'}$, a fact which is implicitly used in the proof of Proposition 10.¹

We write $B_\varepsilon(\mathbf{r})$ for the open ball of \mathbb{R}^n of radius ε centered at $\mathbf{r} \in \mathbb{R}^n$. Given $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, we denote by $d(f)$ the *domain of differentiability* of f , defined to be the set of all $\mathbf{r} \in \mathbb{R}^n$ where f is differentiable in the sense that the total derivative of f at \mathbf{r} exists, *i.e.*, f admits a best linear approximation at \mathbf{r} . We denote by $\mathcal{J}f$ the *Jacobian* of f . We recall that, if $\partial_i f_j$ denotes the partial derivative of the j -th component of f with respect to its i -th parameter, then within $d(f)$ the Jacobian is equal to the $m \times n$ matrix $(\partial_i f_j)$. If $m = 1$, the Jacobian is called *gradient* and denoted by ∇f . As a special case of the above, within $d(f)$ we have $\nabla f = (\partial_1 f, \dots, \partial_n f)$. Although it may happen in general that all partial derivatives exist without the Jacobian/gradient being defined, it will never be the case in what follows because *we will always work within $d(f)$* .

2 PCF WITH REAL NUMBERS

2.1 Terms and Semantics

The programming language we use, called PCF_R , is introduced in Fig. 1. There is only one base type, R , for real numbers. We use n -ary products for convenience. If one prefers, these may be seen as syntactic sugar defined from nullary and binary products. Throughout the paper, we stipulate that unary products are just identities: $\langle M \rangle$ and $\pi_1^1 M$ both stand for M . We often omit the index n in a projection π_i^n , when inessential. The empty product type is denoted by 1 . Given a type A , we write A^n for the n -fold product $A \times \dots \times A$.

The metavariables ϕ, χ, ψ range over a set of function symbols, each coming with a type of the form $R^k \rightarrow R$, where k is the *arity* of the symbol. We suppose that the set of function symbols contains at least all real numbers $r \in \mathbb{R}$ as nullary symbols, called *numerals*, as well as binary addition and multiplication, for which we use infix notation, *i.e.*, $M + N$ and $M \cdot N$ stand for $+(M, N)$ and $\cdot(M, N)$, respectively. We also write n -ary sums as syntactic sugar. Each function symbol $\phi : R^k \rightarrow R$ comes with a function $\llbracket \phi \rrbracket : R^k \rightarrow \mathbb{R}$ and we assume that $\llbracket r \rrbracket$, $\llbracket + \rrbracket$ and $\llbracket \cdot \rrbracket$ are the corresponding numbers and operations on real numbers. The functions $\llbracket \phi \rrbracket$ for ϕ ranging over function symbols will be referred to as the *primitive functions* of PCF_R .

The notation $\text{if}(P, M, N)$ is just a compact form of $\text{if } P \leq 0 \text{ then } M \text{ else } N$. We call P the *guard* of the conditional.

The primitive functions one considers are usually very regular, typically analytic (e.g. exponential, logarithm, trigonometric functions, sigmoid maps...). Sect. 2.3 details the precise conditions that primitives must enjoy in order for our results to hold. The expressive power of PCF_R considerably enlarges the set of definable functions, in particular introducing singularities. The following

¹Proof of the equivalence: let $D := \{a \in A \mid f(a) \neq g(a)\}$. If $\lambda(D) = 0$, then we may take $A' := (\Downarrow f \cup \Downarrow g) \setminus D$ and $Z := D$. Conversely, given A' and Z with the required properties, notice that $D \subseteq \Downarrow f \cup \Downarrow g$, hence $D \subseteq A' \cup Z$. But observe that $D \cap A' = \emptyset$, so $D \subseteq Z$, which gives us $\lambda(D) = 0$.

$$A, B ::= \mathbf{R} \mid A \rightarrow B \mid A_1 \times \cdots \times A_n$$

(a) Types.

$$\frac{}{\Gamma, x^A \vdash x : A} \quad \frac{\phi : \mathbf{R}^k \rightarrow \mathbf{R}, \quad \Gamma \vdash M_1 : \mathbf{R}, \quad \dots, \quad \Gamma \vdash M_k : \mathbf{R}}{\Gamma \vdash \phi(M_1, \dots, M_k) : \mathbf{R}}$$

$$\frac{\Gamma, x^A \vdash M : B}{\Gamma \vdash \lambda x^A. M : A \rightarrow B} \quad \frac{\Gamma \vdash M : A \rightarrow B, \quad \Gamma \vdash N : A}{\Gamma \vdash MN : B}$$

$$\frac{\Gamma \vdash M_1 : A_1, \quad \dots, \quad \Gamma \vdash M_k : A_k}{\Gamma \vdash \langle M_1, \dots, M_k \rangle : A_1 \times \cdots \times A_k} \quad \frac{\Gamma \vdash M : A_1 \times \cdots \times A_k}{\Gamma \vdash \pi_i^k M : A_i}$$

$$\frac{\Gamma \vdash P : \mathbf{R} \quad \Gamma \vdash M : A \quad \Gamma \vdash N : A}{\Gamma \vdash \text{if}(P, M, N) : A} \quad \frac{\Gamma, f : A \rightarrow B \vdash M : A \rightarrow B}{\Gamma \vdash \text{fix } f^{A \rightarrow B}. M : A \rightarrow B}$$

(b) Terms and typing rules.

$$(\lambda x. M)N \rightarrow M\{N/x\} \quad \pi_i^k \langle M_1, \dots, M_k \rangle \rightarrow M_i \quad \phi(r_1, \dots, r_k) \rightarrow \llbracket \phi \rrbracket(r_1, \dots, r_k)$$

$$\text{if}(r, M, N) \rightarrow \begin{cases} M & \text{if } r \leq 0 \\ N & \text{if } r > 0 \end{cases} \quad \text{fix } f. M \rightarrow M\{\lambda x. (\text{fix } f. M)x/f\}$$

(c) Rewriting steps.

Fig. 1. The language $\text{PCF}_{\mathbf{R}}$ over the ground type \mathbf{R} of real numbers.

examples will be useful in the sequel:

$$\begin{aligned} \text{ReLU} &:= \lambda x^{\mathbf{R}}. \text{if}(x, 0, x) & \text{Int} &:= \lambda x^{\mathbf{R}}. \lambda y^{\mathbf{R}}. \text{if}(y - x, \text{if}(y - x + 1, 1, 0), 1) \\ \text{SillyId} &:= \lambda x^{\mathbf{R}}. \text{if}(x, \text{if}(-x, 0, x), x) & \text{Floor} &:= \lambda x^{\mathbf{R}}. (\text{fix } f. \lambda n^{\mathbf{R}}. \text{if}(\text{Int } x \, n, n, f(\text{if}(x, n - 1, n + 1)))) 0 \end{aligned} \quad (1)$$

ReLU and SillyId are the $\text{PCF}_{\mathbf{R}}$ terms corresponding to the namesake examples discussed in the Introduction. ReLU is a typical example of a continuous non-differentiable function, having a corner in 0. We know that AD fails on SillyId ; Sect. 2.2 will elaborate on this point. The program Int takes two inputs x and y and gives 0 if $x \in [y, y + 1]$, or 1 otherwise. It is auxiliary to the definition of the Floor function, mapping a real number to the greatest integer less than or equal to it. Floor is an example of how recursive definitions may yield maps with an infinite number of non-differentiable points, being discontinuous on the integers.

We denote by $M\{N/x\}$ the capture-free substitution of a term N to the free occurrences of the variable x in M .

Note that the presence of the additive structure on \mathbf{R} turns the product \mathbf{R}^n into a biproduct. Indeed, the injection i_i^n , for i such that $1 \leq i \leq n$, may be defined as

$$i_i^n := \lambda x^{\mathbf{R}}. \langle 0, \dots, 0, x, 0, \dots, 0 \rangle, \quad (2)$$

where the variable x is in the i -th position of the n -tuple. We omit the index n when inessential or clear from the context.

For every type $A \rightarrow B$, we set $\Omega_{A \rightarrow B} := \text{fix } f^{A \rightarrow B}.f$. We write just Ω when the type is irrelevant. Given a term $\Gamma, f : A \rightarrow B \vdash M : A \rightarrow B$ and $n \in \mathbb{N} \cup \{\infty\}$, we define $\text{fix}_n f.M$ of type $A \rightarrow B$ as follows:

$$\text{fix}_0 f.M := \Omega, \quad \text{fix}_{n+1} f.M := (\lambda f.M)(\lambda x.(\text{fix}_n f.M)x), \quad \text{fix}_\infty f.M := \text{fix } f.M. \quad (3)$$

Throughout the paper, we use boldface metavariables to denote sequences of metavariables, i.e., $\mathbf{x} = x_1, \dots, x_n$ is a sequence of variables, $\mathbf{M} = M_1, \dots, M_n$ is a sequence of terms, etc. The length of the sequence is specified only when necessary.

Let us introduce two particularly important classes of terms:

DEFINITION 1 (PROGRAM, SIMPLE TERM). A typing environment Γ is ground whenever all of its variables have type R . A PCF_R term M is called a program of arity n and coarity m whenever $x_1^R, \dots, x_n^R \vdash M : R^m$.

A term is called simple if it does not contain conditionals or fixpoints. Small Latin letters t, u, v range over simple terms. Note that the subset of the simple terms of PCF_R corresponds to the simply typed λ -calculus on the ground type R enriched with function symbols.

A context is a term with a single occurrence of a special variable $\{\cdot\}$, called the *hole*. We use metavariables C, D to range over contexts. Given a context C and a term M , we write $C\{M\}$ for the term obtained by replacing the hole $\{\cdot\}$ of C with M , allowing the capture of the free variables in M by the binders of C .

The reduction relation \rightarrow is defined by context closure of the rewriting rules in Fig. 1c:

DEFINITION 2 (REDUCTION). Fig. 1c defines the set of rewriting rules of PCF_R , which are pairs $R \rightarrow P$ of terms, with R called the redex and P the contractum of the rule.

A reduction step σ is a triple (C, R, P) such that C is a context, and $R \rightarrow P$ is a valid reduction rule. We also write $\sigma : C\{R\} \rightarrow C\{P\}$ or, when the context C is irrelevant, simply $\sigma : M \rightarrow N$, for $M = C\{R\}$ and $N = C\{P\}$, and say that σ fires the redex R in M . A term M without redexes, i.e. such that $M \neq C\{R\}$ for any C and any redex R , is said to be normal, or a normal form.

Given another context D , we denote by $D\{\sigma\}$ the step $(D\{C\}, R, P)$. Similarly we write $\sigma\{N/x\}$ for the triple $(C\{N/x\}, R\{N/x\}, P\{N/x\})$, which is still a valid reduction step.

A reduction sequence ρ from a term M to a term N , in symbols $\rho : M \rightarrow^* N$, is either empty, in which case $N = M$, or a sequence of reduction steps $(C_i, R_i, P_i)_{1 \leq i \leq n}$ with $n \geq 1$ such that $M = C_1\{R_1\}$, $N = C_n\{P_n\}$ and for all $i < n$, $C_i\{P_i\} = C_{i+1}\{R_{i+1}\}$. We call M the source of ρ , N its target and n the length of ρ . We often identify a single reduction step and the corresponding reduction sequence of length 1. The notations $D\{\rho\}$ and $\rho\{N/x\}$ are extended to reduction sequences in the obvious way.

Two reductions sequences $\rho : M \rightarrow^* M'$ and $\rho' : M' \rightarrow^* M''$ compose in the obvious way to yield a reduction sequence $\rho\rho' : M \rightarrow^* M''$. Empty reduction sequences are the identities of such an operation.

A sequence ρ is normalizing if there is no reduction step σ such that $\rho\sigma$ is a valid reduction sequence, or, equivalently, if the target of ρ is a normal form. A term M is called normalizing if there exists a normalizing reduction from M . Otherwise M is said to be diverging.

PCF_R is Turing-complete: usual PCF [Plotkin 1977] is essentially the fragment obtained by restricting to integer (including negative) numerals and sum. We recall some results about reduction which are completely standard (see e.g. [Amadio and Curien 1998]).

PROPOSITION 3 (CONFLUENCE). Whenever $M \rightarrow^* N_1$ and $M \rightarrow^* N_2$, there exist N such that $N_1 \rightarrow^* N$ and $N_2 \rightarrow^* N$. In particular, if N_1 and N_2 are normal forms, then $N_1 = N_2$.

PROPOSITION 4 (SUBJECT REDUCTION). If $\Gamma \vdash M : A$ and $M \rightarrow M'$, then $\Gamma \vdash M' : A$.

PROPOSITION 5 (STRONG NORMALIZATION FOR SIMPLE TERMS). *For every simple term t there exists $n \in \mathbb{N}$ such that the length of every reduction sequence starting from t is bounded by n .*

A reduction strategy \mathcal{S} is a relation between terms M of PCF_R and occurrences of redexes in M . A strategy is called *deterministic* whenever it is a partial function. We denote by $\xrightarrow{\mathcal{S}}$ the reduction relation defined by reducing only the redexes fired by \mathcal{S} and we write $\mathcal{S}\text{-nf}$ for a normal form of $\xrightarrow{\mathcal{S}}$. We write β for the maximal strategy, giving the reduction relation \rightarrow . Of course this strategy is not deterministic, however Proposition 3 assures that the normal form associated with a term is unique if it exists.

Reduction strategies are often defined by fixing a subset of redexes in Fig. 1c and a set of *evaluation contexts*. An example which will be useful in the sequel is *head reduction*, which is defined by taking all reduction rules but restricting their application to *head contexts*, generated by the following grammar:

$$H ::= \{\cdot\} \mid \phi(M_1, \dots, H, \dots, M_k) \mid HN \mid \langle H, N \rangle \mid \langle M, H \rangle \mid \pi_i H \mid \text{if}(H, M, N).$$

A *head reduction step* is of the form $H\{R\} \rightarrow H\{P\}$ with $R \rightarrow P$ a rewriting step of Fig. 1c and H a head context. A *head reduction sequence* is a reduction whose steps are all head reduction steps.

Observe that head reduction is not deterministic. Apart from the freedom in the order of the evaluation of the arguments of a function symbol, we allow to reduce within a pair as well as to project it: for instance, the term $\pi_1 \langle (\lambda x.x)y, M \rangle$ may be decomposed either as $H\{\pi_1 \langle (\lambda x.x)y, M \rangle\}$ with the empty head context $H = \{\cdot\}$, or as $H'\{(\lambda x.x)y\}$ with the head context $H' = \pi_1 \langle \{\cdot\}, M \rangle$, and the two decompositions fire different redexes.

A classic result [Amadio and Curien 1998; Barendregt 1985] is that head reduction is a “winning strategy” for finding the β -normal form of a closed program:

PROPOSITION 6. *Let M be a normalizing closed program (i.e., of type R^n) whose β -normal form is N . Then, there is a head reduction sequence $M \xrightarrow{*} N$.*

Let \mathcal{S} be a reduction strategy, let $\Gamma = x_1^R, \dots, x_n^R$ and let $\Gamma \vdash M : R^m$. We define the partial function $\llbracket M \rrbracket_{\Gamma}^{\mathcal{S}} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ as follows:

$$\llbracket M \rrbracket_{\Gamma}^{\mathcal{S}}(r_1, \dots, r_n) = \begin{cases} \langle q_1, \dots, q_m \rangle & \text{if } M\{r_1/x_1\} \dots \{r_n/x_n\} \xrightarrow{\mathcal{S}*} \langle q_1, \dots, q_m \rangle, \\ \perp & \text{otherwise,} \end{cases}$$

where \perp means undefined. By Proposition 3, $\llbracket M \rrbracket_{\Gamma}^{\mathcal{S}}$ is a well-defined partial function, even in case \mathcal{S} is not deterministic, in fact if $M\{r_1/x_1\} \dots \{r_n/x_n\} \xrightarrow{\mathcal{S}*} \langle q_1, \dots, q_m \rangle$ and $M\{r_1/x_1\} \dots \{r_n/x_n\} \xrightarrow{\mathcal{S}*} \langle q'_1, \dots, q'_m \rangle$ we have that $q_i = q'_i$ for each i , as $\xrightarrow{\mathcal{S}} \subseteq \rightarrow$ and this latter is confluent.

We write $\Downarrow^{\mathcal{S}} M$ for $\Downarrow \llbracket M \rrbracket_{\Gamma}^{\mathcal{S}}$ and $d^{\mathcal{S}}(M)$ for $d(\llbracket M \rrbracket_{\Gamma}^{\mathcal{S}})$. In case $\mathcal{S} = \beta$, we omit the indices from the above notations and write simply $\llbracket M \rrbracket_{\Gamma}$, $\Downarrow M$ and $d(M)$. The first of the next two statements follows immediately from Proposition 6, the second from the definitions and the confluence property.

PROPOSITION 7. *Let M, N and P be programs (P of coarity 1), then the following relations hold:*

$$\llbracket \phi(M_1, \dots, M_k) \rrbracket_{\Gamma} = \llbracket \phi \rrbracket_{\Gamma} \circ \langle \llbracket M_1 \rrbracket_{\Gamma}, \dots, \llbracket M_k \rrbracket_{\Gamma} \rangle, \quad \llbracket \langle M_1, \dots, M_k \rangle \rrbracket_{\Gamma} = \langle \llbracket M_1 \rrbracket_{\Gamma}, \dots, \llbracket M_k \rrbracket_{\Gamma} \rangle,$$

$$\llbracket \text{if}(P, M, N) \rrbracket_{\Gamma} = \mathbf{r} \mapsto \begin{cases} \llbracket M \rrbracket_{\Gamma}(\mathbf{r}) & \text{if } \llbracket P \rrbracket_{\Gamma}(\mathbf{r}) \leq 0, \\ \llbracket N \rrbracket_{\Gamma}(\mathbf{r}) & \text{if } \llbracket P \rrbracket_{\Gamma}(\mathbf{r}) > 0, \\ \perp & \text{if } \llbracket P \rrbracket_{\Gamma}(\mathbf{r}) = \perp. \end{cases}$$

$$\begin{aligned}\overrightarrow{\mathbf{D}}_n(\mathbf{R}) &:= \mathbf{R} \times \mathbf{R}^n & \mathbf{D}_n(A \rightarrow B) &:= \mathbf{D}_n(A) \rightarrow \mathbf{D}_n(B) \\ \overleftarrow{\mathbf{D}}_n(\mathbf{R}) &:= \mathbf{R} \times \mathbf{R}^{1-n} & \mathbf{D}_n(A_1 \times \dots \times A_k) &:= \mathbf{D}_n(A_1) \times \dots \times \mathbf{D}_n(A_k)\end{aligned}$$

(a) The action over types

$$\begin{aligned}\overrightarrow{\mathbf{D}}_n(\phi(\mathbf{M})) &:= \left(\lambda \mathbf{z}^{\mathbf{R} \times \mathbf{R}^n} . \left\langle \phi(\pi_1 \mathbf{z}), \sum_{i=1}^k \partial_i \phi(\pi_1 \mathbf{z}) \cdot \pi_2 z_i, \dots, \sum_{i=1}^k \partial_i \phi(\pi_1 \mathbf{z}) \cdot \pi_{n+1} z_i \right\rangle \right) \overrightarrow{\mathbf{D}}_n(\mathbf{M}) \\ \overleftarrow{\mathbf{D}}_n(\phi(\mathbf{M})) &:= \left(\lambda \mathbf{z}^{\mathbf{R} \times \mathbf{R}^{1-n}} . \left\langle \phi(\pi_1 \mathbf{z}), \lambda a^{\mathbf{R}} . \sum_{i=1}^k \pi_2 z_i (\partial_i \phi(\pi_1 \mathbf{z}) \cdot a) \right\rangle \right) \overleftarrow{\mathbf{D}}_n(\mathbf{M})\end{aligned}$$

(b) The action over a function symbol ϕ of arity k . We suppose that, for every $1 \leq i \leq k$, there is an associated function symbol $\partial_i \phi$ of arity k such that $\partial_i \llbracket \phi \rrbracket(\mathbf{r}) = \llbracket \partial_i \phi \rrbracket(\mathbf{r})$ for every $\mathbf{r} \in \mathbb{R}^k$ on which $\partial_i \llbracket \phi \rrbracket$ is defined. The writing $\phi(\mathbf{M})$ is a shortcut for $\phi(M_1, \dots, M_k)$, and, similarly, λz stands for the sequence of abstractions $\lambda z_1 \dots \lambda z_k$ and $\mathbf{D}_n(\mathbf{M})$ for the sequence of applications to $\mathbf{D}_n(M_1) \dots \mathbf{D}_n(M_k)$. These sequences are supposed empty if $k = 0$.

$$\begin{aligned}\mathbf{D}_n(x^A) &:= x^{\mathbf{D}_n(A)} & \mathbf{D}_n(\lambda x^A.M) &:= \lambda x^{\mathbf{D}_n(A)}. \mathbf{D}_n(M) & \mathbf{D}_n(MN) &:= \mathbf{D}_n(M) \mathbf{D}_n(N) \\ \mathbf{D}_n(\langle M_1, \dots, M_k \rangle) &:= \langle \mathbf{D}_n(M_1), \dots, \mathbf{D}_n(M_k) \rangle & \mathbf{D}_n(\pi_i M) &:= \pi_i \mathbf{D}_n(M) \\ \mathbf{D}_n(\text{if}(P, M, N)) &:= \text{if}(\pi_1 \mathbf{D}_n(P), \mathbf{D}_n(M), \mathbf{D}_n(N)) & \mathbf{D}_n(\text{fix } f^A.M) &:= \text{fix } f^{\mathbf{D}_n(A)}. \mathbf{D}_n(M)\end{aligned}$$

(c) The action over the other programming primitives.

Fig. 2. The forward and reverse AD transformations. The symbol \mathbf{D} denotes either one of them. The index n refers to the gradient dimension and acts only over ground type annotations. We will omit the index when inessential or clear from the context.

PROPOSITION 8. *For every program M and strategy \mathcal{S} , we have, for all $\mathbf{r} \in \Downarrow^{\mathcal{S}} M$, $\llbracket M \rrbracket_{\Gamma}^{\mathcal{S}}(\mathbf{r}) = \llbracket M \rrbracket_{\Gamma}(\mathbf{r})$. So, in particular, $\Downarrow^{\mathcal{S}} M \subseteq \Downarrow M$, $\mathbf{d}^{\mathcal{S}}(M) \subseteq \mathbf{d}(M)$ and, for every $\mathbf{r} \in \mathbf{d}^{\mathcal{S}}(M)$, $\mathcal{J} \llbracket M \rrbracket_{\Gamma}^{\mathcal{S}}(\mathbf{r}) = \mathcal{J} \llbracket M \rrbracket_{\Gamma}(\mathbf{r})$.*

2.2 Automatic Differentiation

As mentioned in the Introduction, the two modes of AD are performed by means of the program transformations $\overrightarrow{\mathbf{D}}$ (forward) and $\overleftarrow{\mathbf{D}}$ (reverse), outlined in full detail in Fig. 2. The cornerstone of these transformations is the chain rule, which describes the derivative of a composition of functions:

$$(f \circ g)'(x) = f'(g(x)) \cdot g'(x). \quad (4)$$

This says in particular that $(-)'$ is not functorial (i.e., modular/compositional): the right-hand side of Equation (4) does not use only $g'(x)$ but also $g(x)$. Functoriality may be achieved by transforming a map $f : \mathbb{R} \rightarrow \mathbb{R}$ into a map $\overrightarrow{\mathbf{D}}(f) : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ acting as follows:

$$\overrightarrow{\mathbf{D}}(f) : \langle z, \dot{z} \rangle \mapsto \langle f(z), f'(z) \cdot \dot{z} \rangle. \quad (5)$$

Referring to Equation (4), the input z , which is often called *primal* in the AD literature, corresponds to the output of g , whereas \dot{z} , called *tangent*, corresponds to the output of g' . The reader may easily check functoriality: $\overrightarrow{\mathbf{D}}(f \circ g) = \overrightarrow{\mathbf{D}}(f) \circ \overrightarrow{\mathbf{D}}(g)$. This generalizes to n -ary maps in the definition of

\vec{D}_n in Fig. 2b² and is the essential part of forward mode AD. The attribute *forward* refers to the fact that the computation of both the primal and the tangent follows the input-to-output flow, in particular the derivative $f'(z)$ is computed after having accumulated in z the derivative of its input function.

Notice that, if we denote by $|M|$ the size of a term M , we have that $|\vec{D}_n(\phi(M))| = O(n) + |\vec{D}_n(M)|$. So, supposing that F consists only of function symbols and variables, evaluating both F and $\vec{D}_n(F)$ on a given input requires a number of operations roughly equal to their size. But $|\vec{D}_n(F)| = O(n|F|)$, therefore the evaluation of $\vec{D}_n(F)$ is blown up by a factor of n with respect to the evaluation of F . In applications to deep learning, F is a loss function and n is the number of learning parameters, which may be huge (hundreds of millions).

Luckily, AD offers a more efficient method for applying the chain rule in these cases, called reverse mode AD, or *backpropagation*, because the idea is to accumulate the tangents in the reverse order with respect to the primals. More precisely, by taking the notation of (4), the backpropagation $\overleftarrow{D}(f)$ of f first computes $f'(g(x))$ and then waits for the derivative $g'(x)$ in order to perform the multiplication. As first observed by [Pearlmutter and Siskind 2008], this mode may be naturally expressed in a functional programming language by replacing the tangent variables z with *backpropagators* z^* representing functions (in fact, special forms of continuations):

$$\overleftarrow{D}(f) : \langle z, z^* \rangle \mapsto \langle f(z), \lambda a^R. z^*(f'(z) \cdot a) \rangle \quad (6)$$

A backpropagator is a map $\mathbb{R} \rightarrow \mathbb{R}^n$ *waiting* for a real number (the derivative of the next function) in order to achieve the computation of the gradient of the whole function. In (6), the second component of the returned pair is the backpropagator associated with f , the variable z^* being the awaited backpropagator associated with the input function of f (called g in (4)). This transformation generalizes to the definition of \overleftarrow{D}_n in Fig. 2b for n -ary maps.

We encourage the reader to check that $|\overleftarrow{D}_n(\phi(M))| = O(1) + |\overleftarrow{D}_n(M)|$, so if F consists only of function symbols and variables, the evaluation of $\overleftarrow{D}_n(F)$ is asymptotically as costly as the evaluation of F . However, in more complex cases, the sole transformation (6) is not enough to guarantee efficiency: if F contains sharing, *e.g.* a subroutine g called several times, the evaluation of $\overleftarrow{D}(F)$ may duplicate uselessly the computation associated with the backpropagator of g , and this may result in an exponential blowup. This highlights a key difference between forward and reverse mode: if F is a first order program (*i.e.*, every abstraction λx^A in F is such that A has no arrows), then $\vec{D}(F)$ is also a first order program, whereas $\overleftarrow{D}(F)$ is a higher order term. When backpropagation is expressed in an imperative language, as is usually the case, duplication is not a problem because efficiency is automatically achieved by accumulating the tangents (in reverse order) in memory. But for functional languages, subtle techniques have been introduced to avoid this problem, *e.g.* closure conversions [Pearlmutter and Siskind 2008] or memory references and delimited continuations [Wang et al. 2019].

In this paper we follow the approach of [Brunel et al. 2020], giving a purely functional solution based on linear logic types: backpropagators have type $\mathbb{R}^{\perp n}$, which corresponds to the set of *linear* maps from \mathbb{R} to \mathbb{R}^n . The efficiency of the transformation is then guaranteed by a *factoring rule* added to the operational semantics, which allows sharing the evaluation of different occurrences of a backpropagator z^* in an expression:

$$z^*M + z^*N \quad \rightarrow \quad z^*(M + N). \quad (7)$$

²Modulo some syntactic bureaucracy, *i.e.*, the pair $\langle z, \dot{z} \rangle$ of (5) corresponds in Fig. 2b to a single variable z of type $\mathbb{R} \times \mathbb{R}^n$, whose components are obtained by using the projections π_i

This rewriting rule is sound because backpropagators are linear maps, so they commute with sums. In particular, the normal form of a term obtained using (7) is the same one would have obtained, perhaps in more steps, without using it. As mentioned in the Introduction, in our present setting we are concerned only with soundness, not efficiency. For this reason, we adopt the definition

$$R^{\perp n} := R \rightarrow R^n$$

and do not consider the linear factoring rule (7). In fact, by the above remark, rule (7) only speeds up computation without introducing any error, so our almost-everywhere soundness result is transparent to its use, and enforcing linearity would only lead to unnecessary complications induced by a more sophisticated type system. We do retain the notation $(-)^{\perp}$ as a reminder that this is supposed to be a *linear* arrow (i.e., $R^{\perp n}$ should really be $R \multimap R^n$), but the transformation of [Brunel et al. 2020] remains well typed with the above “non-linear” definition of negation. Indeed, apart from the addition of conditional and fixpoints, the transformation of Fig. 2 is exactly that of *loc. cit.* and it is efficient as long as it is executed according to the operational semantics enriched with (7), so nothing is lost with respect to our previous work.

So far we have explained AD transformations only in regard to primitive functions, which are the “elementary blocks” of straight-line programs mentioned in the Introduction. It is an observation first formalized in [Wang et al. 2019] that the transformations may be extended to arbitrary programs simply by applying the functoriality principle: \vec{D} and \overleftarrow{D} are defined to commute with the programming constructs of the language, resulting in Fig. 2c. As a result, the abstract syntax tree of an expression is basically preserved by the two transformations,³ only the types of the variables are lifted so as to accommodate primals and tangents or backpropagators at the ground level. This behavior is often described by saying that AD is implemented via “operator overloading”. We prefer the term “functoriality” because we find it technically more appropriate.

Notice that, in the definition of $D(\text{if}(P, M, N))$, the transformation is applied also to the guard P in order to preserve typability, because P may share free variables with M and N . However, the computation of the gradient of P is useless and therefore $D(P)$ is projected to the first component. A possible optimization would be to define $D(\text{if}(P, M, N)) := \text{if}(D'(P), D(M), D(N))$ where D' is an auxiliary transformation such that $D'(x^A) = \pi_1 x^{D(A)}$ and which behaves homomorphically on every other term. We avoid introducing D' because it is not crucial for our results and because such an optimization is not so relevant at our level of abstraction (indeed, the evaluation of $D(P)$ is linear in the evaluation of P , as proved in [Brunel et al. 2020], so asymptotically there is no gain).

Some examples of the two modes of D are given in Fig. 3. In the case of subtraction (Fig. 3b), notice how the size of the term resulting from the forward transformation \vec{D}_n increases with the gradient dimension n , while it is constant in \overleftarrow{D}_n , as expected from the above discussion.

Given a term $\Gamma \vdash M : A$ with $\Gamma = x_1^{A_1}, \dots, x_n^{A_n}$, one can check that $D(\Gamma) \vdash D(M) : D(A)$, where $D(\Gamma) = x_1^{D(A_1)}, \dots, x_n^{D(A_n)}$. In particular, if M is a program, then

$$x_1^{R \times R^n}, \dots, x_n^{R \times R^n} \vdash \vec{D}_n(M) : R \times R^n \quad x_1^{R \times R^{\perp n}}, \dots, x_n^{R \times R^{\perp n}} \vdash \overleftarrow{D}_n(M) : R \times R^{\perp n} \quad (8)$$

If, furthermore, M is simple, then the computational behavior of the transformations \vec{D} and \overleftarrow{D} is given by the following result, which was proved in [Barthe et al. 2020; Brunel et al. 2020; Huot et al. 2020],⁴ and in which i_i^n are the injections of R into R^n as defined in Equation (2).

³This is the case for all constructs except the conditional, where a projection is added to the transformation of the guard. This minor technicality is due to the fact that we consider only the ground type of real numbers and not that of Booleans.

⁴In a personal communication, Mitchell Wand showed us that soundness of reverse mode AD may also be proved by means of “open” logical relations of the kind discussed in [Barthe et al. 2020] and used here for the unsoundness bound (Sect. 4).

$$\overrightarrow{\mathbf{D}}_1(\text{ReLU}) = \lambda x^{\mathbb{R} \times \mathbb{R}}. \text{if}(\pi_1 x, \langle 0, 0 \rangle, x) \quad \overleftarrow{\mathbf{D}}_1(\text{ReLU}) = \lambda x^{\mathbb{R} \times \mathbb{R}^\perp}. \text{if}(\pi_1 x, \langle 0, \lambda a. 0 \rangle, x)$$

(a) \mathbf{D}_1 of the rectified linear unit ReLU, defined in (1).

$$\begin{aligned} \overrightarrow{\mathbf{D}}_1(x^{\mathbb{R}} - y^{\mathbb{R}}) &= \left(\lambda z_1^{\mathbb{R} \times \mathbb{R}} z_2^{\mathbb{R} \times \mathbb{R}}. \langle \pi_1^2(z_1) - \pi_1^2(z_2), \pi_2^2(z_1) - \pi_2^2(z_2) \rangle \right) x^{\mathbb{R} \times \mathbb{R}} y^{\mathbb{R} \times \mathbb{R}} \\ \overrightarrow{\mathbf{D}}_2(x^{\mathbb{R}} - y^{\mathbb{R}}) &= \left(\lambda z_1^{\mathbb{R} \times \mathbb{R}^2} z_2^{\mathbb{R} \times \mathbb{R}^2}. \langle \pi_1^3(z_1) - \pi_1^3(z_2), \pi_2^3(z_1) - \pi_2^3(z_2), \pi_3^3(z_1) - \pi_3^3(z_2) \rangle \right) x^{\mathbb{R} \times \mathbb{R}^2} y^{\mathbb{R} \times \mathbb{R}^2} \\ \overleftarrow{\mathbf{D}}_n(x^{\mathbb{R}} - y^{\mathbb{R}}) &= \left(\lambda z_1^{\mathbb{R} \times \mathbb{R}^{1+n}} z_2^{\mathbb{R} \times \mathbb{R}^{1+n}}. \langle \pi_1^2(z_1) - \pi_1^2(z_2), \lambda a^{\mathbb{R}}. (\pi_2^2(z_1) 1 \cdot a + \pi_2^2(z_2) (-1) \cdot a) \rangle \right) x^{\mathbb{R} \times \mathbb{R}^{1+n}} y^{\mathbb{R} \times \mathbb{R}^{1+n}} \end{aligned}$$

(b) \mathbf{D}_n of the subtraction $x - y$, with $n = 1, 2$, where we suppose $\partial_1(x - y) := 1$ and $\partial_2(x - y) := -1$.

Fig. 3. Some examples of the \mathbf{D} transformations. Notice that we take the liberty of using the same name for the ground variables $x^{\mathbb{R}}$ and $y^{\mathbb{R}}$ and their images $x^{\mathbf{D}(\mathbb{R})}$ and $y^{\mathbf{D}(\mathbb{R})}$ under \mathbf{D} .

PROPOSITION 9 (SOUNDNESS OF AD FOR SIMPLE TERMS). *Let $\Gamma = x_1^{\mathbb{R}}, \dots, x_n^{\mathbb{R}}$ and let $\Gamma \vdash t : \mathbb{R}$ be a simple program. Then, for all $\mathbf{r} = (r_1, \dots, r_n) \in \mathbf{d}(t)$, we have*

$$\begin{aligned} \overrightarrow{\mathbf{D}}_n(t) \{ \langle r_1, l_1^n 1 \rangle / x_1 \} \dots \{ \langle r_n, l_n^n 1 \rangle / x_n \} &\rightarrow^* \langle \llbracket t \rrbracket_{\Gamma}(\mathbf{r}), \nabla \llbracket t \rrbracket_{\Gamma}(\mathbf{r}) \rangle \\ \overleftarrow{\mathbf{D}}_n(t) \{ \langle r_1, l_1^n \rangle / x_1 \} \dots \{ \langle r_n, l_n^n \rangle / x_n \} &\rightarrow^* \langle \llbracket t \rrbracket_{\Gamma}(\mathbf{r}), u \rangle \end{aligned}$$

such that $u 1 \rightarrow^* \nabla \llbracket t \rrbracket_{\Gamma}(\mathbf{r})$.

Looking at Proposition 9, if M is an arbitrary program of arity n and coarity 1, it is reasonable to believe that the following programs compute $\nabla \llbracket M \rrbracket_{\Gamma}$ in $\mathbf{r} = (r_1, \dots, r_n) \in \mathbb{R}^n$ whenever this is defined:

$$\overrightarrow{\text{grad}}_n(M)(\mathbf{r}) := \pi_2^2 \overrightarrow{\mathbf{D}}_n(M) \{ \langle r_1, l_1^n 1 \rangle / x_1 \} \dots \{ \langle r_n, l_n^n 1 \rangle / x_n \}, \quad (9)$$

$$\overleftarrow{\text{grad}}_n(M)(\mathbf{r}) := (\pi_2^2 \overleftarrow{\mathbf{D}}_n(M) \{ \langle r_1, l_1^n \rangle / x_1 \} \dots \{ \langle r_n, l_n^n \rangle / x_n \}) 1. \quad (10)$$

In the sequel, we will omit the index n when inessential or clear from the context. We will also write $\text{grad}(M)$ for either one of the above terms.⁵

Referring to Fig. 3, it is immediate to check that, regardless of the mode, $\text{grad}_1(\text{ReLU } z^{\mathbb{R}})(r) \rightarrow^* 1$ if $r > 0$ and $\text{grad}_1(\text{ReLU } z^{\mathbb{R}})(r) \rightarrow^* 0$ if $r \leq 0$. This is the expected result except for $r = 0$, where the map $\llbracket \text{ReLU } z \rrbracket_{z^{\mathbb{R}}}$ is not differentiable. As discussed in the Introduction, our soundness result concerns only the domain of differentiability $\mathbf{d}(M)$ of a map $\llbracket M \rrbracket$ represented by a program and nothing is stated about the value of $\text{grad}(M)$ outside $\mathbf{d}(M)$. This situation is quite common in AD frameworks, where it may even be desirable to control the behavior of the “non-existent derivative” on singularities. For example, the following term cReLU implements the rectified linear unit (i.e., $\llbracket \text{ReLU } z \rrbracket_{z^{\mathbb{R}}} = \llbracket \text{ReLU } z \rrbracket_{z^{\mathbb{R}}}$) so that grad returns some arbitrarily chosen $q \in \mathbb{R}$ on 0:

$$\text{cReLU}_q := \lambda x^{\mathbb{R}}. \text{if}(x, \text{if}(-x, q \cdot x, 0), x), \quad \text{grad}_1(\text{cReLU}_q z^{\mathbb{R}})(r) \rightarrow^* \begin{cases} 1 & \text{if } r > 0, \\ q & \text{if } r = 0, \\ 0 & \text{if } r < 0. \end{cases} \quad (11)$$

We do not consider these computations as errors because $\nabla \llbracket \text{ReLU } z \rrbracket_{z^{\mathbb{R}}}$ is undefined at 0.

⁵Notice that the above definition of grad is slightly different from the informal one used in the Introduction: it applies to terms with free ground variables, whereas in the Introduction we abusively applied grad to closed terms.

We already discussed Sillyld (see (1)) as a first example of a mismatch between AD and the gradient in the domain of differentiability of a map. Let us consider here a refined example:

$$\text{EqProj} := \lambda x^R. \lambda y^R. \text{if}(x - y, \text{if}(y - x, x, y), y). \quad (12)$$

This term is extensionally equivalent to the binary projection $\lambda x^R. \lambda y^R. y$ and therefore its gradient should be $\langle 0, 1 \rangle$ on the whole domain \mathbb{R}^2 . By contrast, the reader may check that:

$$\begin{aligned} \overrightarrow{\text{grad}}_2(\text{EqProj } x_1^R x_2^R)(r_1, r_2) &= \pi_2^2 \left(\overrightarrow{\text{D}}_2(\text{EqProj}) \langle r_1, l_1^2 1 \rangle \langle r_2, l_2^2 1 \rangle \right) \\ &\rightarrow^* \pi_2^2 (\text{if}(r_1 - r_2, \text{if}(r_2 - r_1, \langle r_1, \langle 1, 0 \rangle \rangle, \langle r_2, \langle 0, 1 \rangle \rangle), \langle r_2, \langle 0, 1 \rangle \rangle)) \\ &\rightarrow^* \begin{cases} \langle 0, 1 \rangle & \text{if } r_1 \neq r_2, \\ \langle 1, 0 \rangle & \text{if } r_1 = r_2. \end{cases} \end{aligned}$$

The diagonal of \mathbb{R}^2 gives an uncountable set of errors (a similar computation yields the same result also for the reverse mode). However, this set is negligible, *i.e.*, of Lebesgue measure zero, in accordance with the claim (ii') stated in the Introduction.

Let us add a last comment on example (12). Consider the unary program $\text{EqProj } x_1^R x_1^R$, which is extensionally equivalent to the identity. One can check that $\overrightarrow{\text{grad}}_1(\text{EqProj } x_1 x_1)(r) \rightarrow^* 1$ for every $r \in \mathbb{R}$, so there is no error at all in this case. This is in sharp contrast with approaches based on partial conditionals, such as [Abadi and Plotkin 2020], in which conditionals diverge when the guard evaluates to 0: under such semantics, $\text{EqProj } x_1 x_1$ diverges everywhere.

Different reduction strategies change the convergence and differentiability domain of a program, so a priori the soundness of AD depends on the strategy. In the above examples, we considered the maximal reduction strategy β . If we wish to specialize to a more restrictive reduction strategy \mathcal{S} , we should prove that $\text{grad}(M)$ evaluates, in accordance with \mathcal{S} , to the gradient of $\llbracket M \rrbracket_{\Gamma}^{\mathcal{S}}$ at almost every point where this is defined. In fact, if we succeed in proving this with respect to β , then it follows also for any other “reasonable” strategy \mathcal{S} :

PROPOSITION 10. *Let $\Gamma \vdash M : R$ be a program. If $\llbracket \text{grad}(M) \rrbracket_{\Gamma} \mid_{d(M)} \sim \nabla(\llbracket M \rrbracket_{\Gamma})$, then for any reduction strategy \mathcal{S} such that $d^{\mathcal{S}}(M) \subseteq \Downarrow^{\mathcal{S}}(\text{grad}(M))$ we also have $\llbracket \text{grad}(M) \rrbracket_{\Gamma}^{\mathcal{S}} \mid_{d^{\mathcal{S}}(M)} \sim \nabla(\llbracket M \rrbracket_{\Gamma}^{\mathcal{S}})$.*

PROOF. By hypothesis we have that $d(M) = A \cup Z$ with Z of measure zero and $\llbracket \text{grad}(M) \rrbracket_{\Gamma} \mid_A = \nabla(\llbracket M \rrbracket_{\Gamma}) \mid_A$. Let \mathcal{S} be a reduction strategy satisfying the hypothesis and let $B := A \cap d^{\mathcal{S}}(M)$. We need to prove that $\llbracket \text{grad}(M) \rrbracket_{\Gamma}^{\mathcal{S}} \mid_B = \nabla(\llbracket M \rrbracket_{\Gamma}^{\mathcal{S}}) \mid_B$ and that $d^{\mathcal{S}}(M) \setminus B$ is negligible.

Let us start with this latter point. Notice that, by Proposition 8, $d^{\mathcal{S}}(M) \setminus B = d^{\mathcal{S}}(M) \setminus (A \cap d^{\mathcal{S}}(M)) = d^{\mathcal{S}}(M) \setminus A \subseteq d(M) \setminus A$, and the latter set is of measure zero by hypothesis.

Let us now take $\mathbf{r} \in B$. We have:

$$\begin{aligned} \llbracket \text{grad}(M) \rrbracket_{\Gamma}^{\mathcal{S}}(\mathbf{r}) &= \llbracket \text{grad}(M) \rrbracket_{\Gamma}(\mathbf{r}) && \text{by } B \subseteq d^{\mathcal{S}}(M) \subseteq \Downarrow^{\mathcal{S}}(\text{grad}(M)) \text{ and Proposition 8} \\ &= \nabla(\llbracket M \rrbracket_{\Gamma})(\mathbf{r}) && \text{because } B \subseteq A \\ &= \nabla(\llbracket M \rrbracket_{\Gamma}^{\mathcal{S}})(\mathbf{r}) && \text{by } B \subseteq d^{\mathcal{S}}(M) \text{ and Proposition 8.} \end{aligned}$$

□

The condition $d^{\mathcal{S}}(M) \subseteq \Downarrow^{\mathcal{S}}(\text{grad}(M))$ is reasonable, since $d^{\mathcal{S}}(M) \subseteq \Downarrow^{\mathcal{S}} M$ by definition and it is very likely that $\Downarrow^{\mathcal{S}} M \subseteq \Downarrow^{\mathcal{S}} \text{grad}(M)$, because the convergence of $\text{grad}(M)$ coincides with that of $D(M)$ and the latter essentially behaves like M , as we will prove in Sect. 3.2. Notice that, when t is simple, $\Downarrow^{\mathcal{S}} t = \Downarrow^{\mathcal{S}} \text{grad}(t)$ is trivially true because of strong normalization (Proposition 5), so

Proposition 9 in fact holds for any reduction strategy. Common strategies such as call-by-value, call-by-name and call-by-need are easily seen to enjoy the condition of Proposition 10. See Remark 30 below for a proof sketch in the case of call-by-value.

2.3 Primitive Functions and Complete Quasicontinuity

The only assumption we made so far about the function symbols of PCF_R is that for every ϕ of arity k and every $1 \leq i \leq k$, there is another k -ary function symbol $\partial_i \phi$ corresponding to the partial derivative of $\llbracket \phi \rrbracket$ with respect to its i -th argument (Fig. 2b). In order to prove one of our main results (Theorem 42), we will need to make some further topological and measure-theoretic assumptions, which we proceed to spell out.

In what follows, we always consider \mathbb{R}^n with its standard topology and the Lebesgue measure. We denote by $\text{int}(X)$ the interior of a set X (the largest open set contained in X) and by $\text{bor}(X)$ its *border*, defined as $\text{bor}(X) := X \setminus \text{int}(X)$. Equivalently, $\text{bor}(X) = \partial X \cap X$ where ∂X is the boundary of X (the closure of X minus $\text{int}(X)$). In case X is closed, $\text{bor}(X) = \partial X$.

We recall that a *clone* [Szendrei 1986] on a set A is a collection \mathbf{P} of functions $A^n \rightarrow A$ (for varying n) which is closed under composition⁶ and contains all projections (in particular, the identity on A). Notice that clones are stable under arbitrary intersections, hence every set \mathbf{F} of functions $A^n \rightarrow A$ (for possibly varying n) generates a clone $\langle \mathbf{F} \rangle$, the smallest clone containing \mathbf{F} .

DEFINITION 11 (ADMISSIBLE PRIMITIVE FUNCTIONS). *We say that a clone \mathbf{P} on \mathbb{R} is admissible if $f \in \mathbf{P}$ implies:*

- (1) *f is continuous on its domain;*
- (2) *if $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is not identically zero, then $f^{-1}(0)$ is of Lebesgue measure zero in \mathbb{R}^n .*

Fix a set \mathbf{F} of function symbols together with their semantics. We abusively denote by \mathbf{F} also the set of all $\llbracket \phi \rrbracket$ with ϕ ranging over the chosen function symbols. We say that \mathbf{F} forms an admissible set of primitive functions if $\langle \mathbf{F} \rangle$ is admissible.

It is well known [Mityagin 2015] that an example of admissible clone is provided by the collection of all real functions which are defined and analytic on some open set $U \subseteq \mathbb{R}^n$, for varying U and n . Notice that a subclone of an admissible clone is admissible. Therefore, a simple way of ensuring that the primitive functions of PCF_R are admissible is to ask that they are all analytic where they are defined. This is of course true of our “mandatory” primitive functions (constants, addition and multiplication), as well as all functions usually taken as primitive, such as division, square root, exponential, logarithm, the trigonometric functions and their inverses, Gaussian functions, many sigmoid functions (e.g. the error function), etc. Other desirable functions which are not analytic (the step function, the floor function, the rectified linear unit...) are usually programmable in PCF_R from these primitive functions.

The definitions that follow are parametric in a choice of admissible clone \mathbf{B} , so we should speak of \mathbf{B} -quasiopen set, \mathbf{B} -quasicontinuity, etc. However, for simplicity, we will omit the parameter \mathbf{B} , implicitly fixing once and for all an admissible set \mathbf{F} of primitive functions and letting $\mathbf{B} := \langle \mathbf{F} \rangle$. Functions in \mathbf{B} will be called *basic*.

DEFINITION 12 (QUASIOPEN SET). *We define the class of quasiopen sets of \mathbb{R}^n to be the smallest class of subsets of \mathbb{R}^n which:*

- (1) *contains every open set;*
- (2) *contains the zero set of every basic function $\mathbb{R}^n \rightarrow \mathbb{R}$;*
- (3) *is closed under countable unions and binary intersections.*

⁶We mean that if $f : A^k \rightarrow A$ and $g_1, \dots, g_k : A^n \rightarrow A$ are in \mathbf{P} , then so is the function $a \mapsto f(g_1(a), \dots, g_k(a))$.

Inductively, the quasiopen sets of \mathbb{R}^n may be defined as follows:

$$Q, Q' ::= U \mid h^{-1}(0) \mid \bigcup_{i \in I} Q_i \mid Q \cap Q',$$

where U ranges over the open sets of \mathbb{R}^n , $h : \mathbb{R}^n \rightarrow \mathbb{R}$ ranges over basic functions (which may further be supposed to be not identically zero) and I is countable.

DEFINITION 13 (QUASIVARIETY). A set $Z \subseteq \mathbb{R}^n$ is called a quasivariety if there exists a family $\{h_i\}_{i \in I}$ of basic functions $h_i : \mathbb{R}^n \rightarrow \mathbb{R}$ with I countable and such that $Z \subseteq \bigcup_{i \in I} h_i^{-1}(\{0\})$. In other words, a quasivariety is an arbitrary subset of a countable union of zero sets of basic functions.

The following result, which says that quasivarieties form a class of “negligible sets”, will be frequently used in the sequel, without explicit mention:

LEMMA 14. *Quasivarieties enjoy the following properties:*

- (1) **measure zero:** if $Z \subseteq \mathbb{R}^n$ is a quasivariety, then it has Lebesgue measure zero in \mathbb{R}^n ;
- (2) **stability under countable unions:** if $\{Z_i\}_{i \in I}$ is a countable family of quasivarieties, then $\bigcup_{i \in I} Z_i$ is a quasivariety;
- (3) **stability under subsets:** if Z is a quasivariety and $Z' \subseteq Z$, then Z' is a quasivariety.

PROOF. Immediate from the definition. □

LEMMA 15. *Let $Q \subseteq \mathbb{R}^n$ be quasiopen. Then:*

- (1) *there exists an open set U and a quasivariety Z such that $Q = U \cup Z$;*
- (2) *$\text{bor}(Q)$ is a quasivariety. Hence, in the above one may always take $U = \text{int}(Q)$ and $Z = \text{bor}(Q)$.*

The set of non-positive numbers $\mathbb{R}_{\leq 0}$ is an example of quasiopen subset of \mathbb{R} : to see why, simply notice that $\mathbb{R}_{\leq 0} = \mathbb{R}_{< 0} \cup \{0\}$, the first being open and the second being the zero set of the identity, which is always a basic function. In a sense, the key property of the class of quasiopen sets is that it includes both $\mathbb{R}_{\leq 0}$ and $\mathbb{R}_{> 0}$, a fact which will be used crucially in Lemma 38.

On the other hand, thick Cantor sets provide examples of non-quasiopen subsets of \mathbb{R} : such a set K is closed, of positive measure and has empty interior, so $K = \text{bor}(K)$, which would contradict Lemma 15.2 if K were quasiopen.

In what follows, if $f : A \rightarrow B$ and $g : C \rightarrow D$ are partial functions between sets, we write $f \times g$ for the function of type $A \times C \rightarrow B \times D$ such that $(f \times g)(a, c) = (f(a), g(c))$ whenever $f(a)$ and $g(c)$ are defined, and is undefined otherwise.

DEFINITION 16 ((COMPLETE) QUASICONTINUITY). A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is quasicontinuous⁷ if, for every quasiopen set $Q \subseteq \mathbb{R}^m$, $f^{-1}(Q)$ is quasiopen. We say that f is completely quasicontinuous (cqc) if $\text{id}_{\mathbb{R}^k} \times f$ is quasicontinuous, for all $k \in \mathbb{N}$.

Complete quasicontinuity is needed in order to have Lemma 17.4 below. It is worth pointing out, however, that we have not been able to find an example of a quasicontinuous function which is not completely quasicontinuous. So, while we conjecture that complete quasicontinuity is strictly stronger than quasicontinuity, the two notions might coincide in reality.

LEMMA 17. *We have the following properties:*

- (1) *a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is quasicontinuous iff for every Q which is either open or the zero set of a basic function, $f^{-1}(Q)$ is quasiopen.*

⁷The terminology “quasicontinuous” already has a standard meaning, unrelated to the one defined here. On the other hand, “quasiopen” and “completely quasicontinuous”, which are the fundamental notions used in this work, do not seem to have been used in the literature.

$$\frac{}{R \sqsubset R} \quad \frac{A' \sqsubset A \quad B' \sqsubset B}{A' \rightarrow B' \sqsubset A \rightarrow B} \quad \frac{A'_i \sqsubset A_i, \forall 1 \leq i \leq k}{A'_1 \times \cdots \times A'_k \sqsubset A_1 \times \cdots \times A_k} \quad \frac{A_i \sqsubset A, \forall 1 \leq i \leq n}{A_1 \times \cdots \times A_n \sqsubset A}$$

(a) The pre-trace relation on types.

$$\frac{}{\Xi \vdash \{\cdot\} \sqsubset \{\cdot\}} \quad \frac{A_1 \times \cdots \times A_n \sqsubset A}{\Xi, p^{A_1 \times \cdots \times A_n} \sqsubset x^A \vdash \pi_i^n p \sqsubset x} \quad i \in \{1, \dots, n\}$$

$$\frac{\Xi, p^{A'} \sqsubset x^A \vdash t \sqsubset M}{\Xi \vdash \lambda p^{A'}.t \sqsubset \lambda x.M} \quad \frac{\Xi \vdash t \sqsubset M, \quad \Xi \vdash u_1 \sqsubset N \quad \dots \quad \Xi \vdash u_n \sqsubset N}{\Xi \vdash t\langle u_1, \dots, u_n \rangle \sqsubset MN}$$

$$\frac{\Xi \vdash t_1 \sqsubset M_1, \quad \dots, \quad \Xi \vdash t_k \sqsubset M_k}{\Xi \vdash \langle t_1, \dots, t_k \rangle \sqsubset \langle M_1, \dots, M_k \rangle} \quad \frac{\Xi \vdash t \sqsubset M}{\Xi \vdash \pi_i^k t \sqsubset \pi_i^k M} \quad \frac{\Xi \vdash t_1 \sqsubset M_1 \quad \dots \quad \Xi \vdash t_k \sqsubset M_k}{\Xi \vdash \phi(t_1, \dots, t_k) \sqsubset \phi(M_1, \dots, M_k)}$$

$$\frac{\Xi \vdash t_i \sqsubset M_i}{\Xi \vdash \pi_i \langle t_i, t_i \rangle \sqsubset \text{if}(P, M_1, M_2)} \quad i \in \{1, 2\} \quad \frac{\Xi \vdash t \sqsubset \text{fix}_n f.M}{\Xi \vdash t \sqsubset \text{fix} f.M} \quad n > 0$$

(b) The pre-trace relation on terms. In the variable rule, if $n = 1$, then π_i is omitted. Recall the definition of $\text{fix}_n f.M$ in (3).

Fig. 4. The pre-trace relation between simple and arbitrary PCF_R terms.

- (2) Identities are cqc and cqc functions are stable under composition.
- (3) Basic functions are cqc. In particular, projections are cqc.
- (4) If $f : \mathbb{R}^k \rightarrow \mathbb{R}^m$ and $g : \mathbb{R}^k \rightarrow \mathbb{R}^n$ are cqc, then the function $\langle f, g \rangle : \mathbb{R}^k \rightarrow \mathbb{R}^{m+n}$ defined by $\langle f, g \rangle(z) := (f(z), g(z))$ if $z \in \Downarrow f \cap \Downarrow g$ and undefined otherwise, is also cqc.

Contrarily to what the name might suggest, a continuous function is *not* in general quasicontinuous. In fact, it is well known that any closed subset of \mathbb{R} may be the zero set of a map which is smooth everywhere (in particular, continuous). So let $\phi : \mathbb{R} \rightarrow \mathbb{R}$ be a smooth function whose zero set is a thick Cantor set K , which, as observed above, is not quasiopen. The set $\{0\}$ is quasiopen (it is the zero set of the identity, which is a basic map), and yet $\phi^{-1}(\{0\}) = K$, so ϕ is not quasicontinuous.

3 SOUNDNESS OF AD

We want to prove that $\text{grad}(M)$ computes the gradient of a program M almost everywhere in $d(M)$ (Theorem 42). The proof splits in two parts: Theorem 33 states that $\text{grad}(M)$ is sound for the set $S(M)$ of stable points of $d(M)$ (Definition 26) and Sect. 4 shows that $S(M)$ is actually almost all of $d(M) \subseteq \Downarrow M$, in the sense that $\Downarrow M \setminus S(M)$ is of measure zero.

Intuitively, a point $r \in \Downarrow M$ is stable whenever there exists a simple term t that “traces” the evaluation of M over an open ball $B_\varepsilon(r)$ of r . Such a t allows us to lift the soundness theorem for simple terms (Proposition 9) to $\text{grad}(M)$. This reasoning is based on the extrusion lemma (Lemma 31), which needs a notion of “trace” not only at level of terms (Definition 25), but also at the level of the reduction sequences (Definition 24).

3.1 Traces

The *pre-trace relation* is defined in Fig. 4. The judgments used in the definition are of the form $\Xi \vdash t \sqsubset M$, where t is a simple term or simple context, M is an arbitrary term or context of type, say, $\Gamma \vdash M : A$ and Ξ is a function mapping any variable x^A of Γ to a fresh variable $p^{A'}$ with $A' \sqsubset A$.

We usually denote this map as a list $p_1^{A'_1} \sqsubset x_1^{A_1}, \dots, p_n^{A'_n} \sqsubset x_n^{A_n}$, supposing the p_i 's and x_i 's to be pairwise different.

For brevity, we omit to specify the types of the terms in the subjects of the judgments in Fig. 4b, but we encourage the reader to verify that if $p_1 \sqsubset x_1, \dots, p_n \sqsubset x_n \vdash t \sqsubset M$ and $x_1^{C_1}, \dots, x_n^{C_n} \vdash M : A$, then $p_1^{C'_1}, \dots, p_n^{C'_n} \vdash t : A'$ with $A' \sqsubset A$ and, for all $1 \leq i \leq n$, $C'_i \sqsubset C_i$. In particular, t is a simply-typed λ -term or context. We write $t \sqsubset M$ when the typing environment Ξ is irrelevant.

Conditionals and fixpoints are the only additional features of PCFR with respect to the simply-typed λ -calculus, and the purpose of \sqsubset is to “trace” them with simply-typed terms themselves. The last two rules of Fig. 4b “slice out” a conditional with the traces of its two branches and unfold a fixpoint into its finite approximations. In fact, the conditional rule is a bit more convoluted as it uses a dummy projection in order to encode the index of the chosen branch, a crucial information for the extrusion property (Lemma 31). For example, $u_1 := \lambda x^R. \pi_1 \langle 0, 0 \rangle$ and $u_2 := \lambda x^R. \pi_2 \langle x, x \rangle$ are traces corresponding to the “then” and “else” branch, respectively, of ReLU defined in (1).

The variable rule of Fig. 4 also deserves an explanation. Its non-trivial shape, which is due to higher order types, may be understood as follows. Let $T := \lambda f^{R \rightarrow R}. f(f0 + 1)$ be a term using its (higher order) argument twice and consider the program $T \text{ReLU}$. Recall that ReLU contains a conditional controlled by its argument, and has two different traces u_1 and u_2 discussed above. However, the execution of $T \text{ReLU}$, as sketched in Fig. 5, explores both branches of ReLU , so we allow to trace T with $t := \lambda p. \pi_2 p((\pi_1 p)0 + 1)$ and $T \text{ReLU}$ with $t \langle u_1, u_2 \rangle$. That is, we allow different instances of the same variable f to be traced by different components of a tuple variable p , because, even if in the original program all occurrences of f are replaced by copies of the same term ReLU , different copies of ReLU might be traced by different simple terms, so the occurrences of f must be “separated” accordingly.

LEMMA 18. *Let M and t be normal forms of type $\mathbf{D}_n(\mathbf{R})$ whose free variables have type belonging to $\{\mathbf{D}_n(\mathbf{R}), \mathbf{R}, \mathbf{R}^n, \mathbf{R}^{\perp n}\}$. If $t \sqsubset M$ then $t = M$.*

LEMMA 19. *If $\Xi \vdash t \sqsubset M$, then:*

- (1) $\mathbf{D}(\Xi) \vdash \mathbf{D}(t) \sqsubset \mathbf{D}(M)$, where \mathbf{D} turns any assignment $p^{A'} \sqsubset x^A$ of Ξ into $p^{\mathbf{D}(A')} \sqsubset x^{\mathbf{D}(A)}$.
- (2) Let $\Xi = \Xi', x^A \sqsubset x^A$. For every closed simple term u of type A , we have $\Xi' \vdash t\{u/x\} \sqsubset M\{u/x\}$.

LEMMA 20. *We have that $\Xi \vdash w \sqsubset M\{N/x\}$ is equivalent to*

- $w = t\{u_1/x_1\} \dots \{u_n/x_n\}$, for some $n \in \mathbb{N}$ and terms t, u_1, \dots, u_n ,
- such that $\Xi, p^{A_1 \times \dots \times A_n} \sqsubset x^A \vdash t\{\pi_1 p/x_1\} \dots \{\pi_n p/x_n\} \sqsubset M$, p not free in t ,
- and $\Xi \vdash u_i \sqsubset N$ for all $1 \leq i \leq n$.

In particular, $\Xi, p^{A_1 \times \dots \times A_n} \sqsubset x^A \vdash w \sqsubset M$ implies $w = t\{\pi_1 p/x_1\} \dots \{\pi_n p/x_n\}$ for some t not containing p free.

DEFINITION 21 (TRACING REWRITING STEPS). *Let $\sigma : R \rightarrow P$ be a rewriting step of Fig. 1c, and v be a reduction sequence between simple terms. We define $v \sqsubset \sigma$ depending on R .*

- If $R = (\lambda x.M)N$: we ask that v is any reduction of the form

$$(\lambda p. t\{\pi p/x\})\langle \mathbf{u} \rangle \rightarrow t\{\pi \langle \mathbf{u} \rangle / x\} \rightarrow^* t\{\mathbf{u}/x\}$$

where $p \sqsubset x \vdash t\{\pi p/x\} \sqsubset M$, with $\{\pi p/x\}$ denoting the substitutions $\{\pi_1 p/x_1\} \dots \{\pi_n p/x_n\}$ for some $n \in \mathbb{N}$, and where \mathbf{u} is a sequence u_1, \dots, u_n of simple terms such that $u_i \sqsubset N$ for all i , and in which the redexes $\pi_i \langle \mathbf{u} \rangle \rightarrow u_i$ are reduced in any order.

- If $R = \pi_i \langle M_1, \dots, M_k \rangle$: we ask that v is any reduction of the form $\pi_i \langle t_1, \dots, t_k \rangle \rightarrow t_i$ for $t_1 \sqsubset M_1, \dots, t_k \sqsubset M_k$.

$$\begin{array}{c}
T \text{ReLU} \longrightarrow \text{ReLU}(\text{ReLU } 0 + 1) \longrightarrow \text{if}(\text{ReLU } 0 + 1, 0, \text{ReLU } 0 + 1) \\
\sqcup \qquad \qquad \qquad \sqcup \qquad \qquad \qquad \sqcup \\
t\langle u_1, u_2 \rangle \longrightarrow \pi_2\langle u_1, u_2 \rangle (\pi_1\langle u_1, u_2 \rangle 0 + 1) \xrightarrow{*} u_2(u_1 0 + 1) \longrightarrow \pi_2\langle u_1 0 + 1, u_1 0 + 1 \rangle \\
\\
\text{if}(\text{ReLU } 0 + 1, 0, \text{ReLU } 0 + 1) \xrightarrow{*} \text{if}(1, 0, \text{ReLU } 0 + 1) \longrightarrow \text{ReLU } 0 + 1 \longrightarrow \text{if}(0, 0, 0) + 1 \rightarrow 0 + 1 \rightarrow 1 \\
\sqcup \qquad \qquad \qquad \sqcup \qquad \qquad \qquad \sqcup \qquad \qquad \qquad \sqcup \qquad \qquad \qquad \sqcup \qquad \qquad \sqcup \\
\pi_2\langle u_1 0 + 1, u_1 0 + 1 \rangle = \pi_2\langle u_1 0 + 1, u_1 0 + 1 \rangle \longrightarrow u_1 0 + 1 \longrightarrow \pi_1\langle 0, 0 \rangle + 1 \rightarrow 0 + 1 \rightarrow 1
\end{array}$$

Fig. 5. Tracing the head reduction $T \text{ReLU} \xrightarrow{*} 1$. The term T is $\lambda f.f(f0 + 1)$, of which $t := \lambda p.\pi_2 p((\pi_1 p)0 + 1)$ is a trace; ReLU is given in (1), with traces $u_1 := \lambda x.\pi_1\langle 0, 0 \rangle$, $u_2 := \lambda x.\pi_2\langle x, x \rangle$.

- If $R = \phi(r)$: we ask that $v = \sigma$.
- If $R = \text{if}(r, M_1, M_2)$: we ask that v is any reduction of the form $\pi_i\langle t, t \rangle \rightarrow t$ with $t \sqsubset M_i$ and $i = 1$ if $r \leq 0$, otherwise $i = 2$.
- If $R = \text{fix } f.N$: we ask that $v \sqsubset \sigma'$ where σ' is any reduction of the form

$$(\lambda f.M)(\lambda x.(\text{fix}_n f.M)x) \rightarrow M\{\lambda x.(\text{fix}_n f.M)x / f\}$$

for some $n \in \mathbb{N}$, as defined in the first case.

LEMMA 22 (PULLBACK). Let $\sigma : R \rightarrow P$ be a rewriting step. For any $w \sqsubset P$, there exist $t \sqsubset R$ and $\xi : t \xrightarrow{*} w$ such that $\xi \sqsubset \sigma$.

We now extend Definition 21 to one-step head reductions (as defined in Sect. 2.1).

DEFINITION 23 (TRACING HEAD REDUCTION STEPS). Let $\sigma = (H, R, P)$ be a reduction step with H a head context, and let σ_0 denote the rewriting step $R \rightarrow P$. If $\xi : t \xrightarrow{*} u$ is a reduction sequence on simple terms, we write $\xi \sqsubset \sigma$ whenever one of the following holds:

- there exists a simple context $h \sqsubset H$ such that $\xi = h\{v\}$ with $v \sqsubset \sigma_0$ in the sense of Definition 21;
- the hole of H is in the guard of a conditional, ξ is the empty sequence and $t = u \sqsubset H\{R\}$.

Finally, we extend the relation \sqsubset to reduction sequences by reflexive-transitive closure.

DEFINITION 24 (TRACING HEAD REDUCTION SEQUENCES). Let ρ be a head reduction sequence starting from a term M and let ξ be a reduction sequence starting from a simple term t . We write $\xi \sqsubset \rho$ if either ρ and ξ are empty and $t \sqsubset M$, or if $\rho = \sigma_1 \cdots \sigma_n$ is of length $n > 0$ and $\xi = \xi_1 \cdots \xi_n$ such that $\xi_i \sqsubset \sigma_i$ for every $1 \leq i \leq n$, according to Definition 23.

Fig. 5 gives an example of tracing the head reduction of the term $T \text{ReLU}$ discussed above. Notice that the reduct of $t\langle u_1, u_2 \rangle$ is an intermediate term not corresponding to any trace. Notice also that all the reductions in the guard of a conditional (such as the first steps in the third line of Fig. 5) share the same traces. In general, $\xi \sqsubset \rho$ implies that ρ is a head reduction but not necessary ξ , because the first case of Definition 21 requires ξ to reduce projection redexes not in the “head position” of a simple term. In fact, the reduction of $t\langle u_1, u_2 \rangle$ in Fig. 5 is not head.

DEFINITION 25 (TRACE RELATION). Let M be a term and t a simple term. We say that t traces M , in symbols $t \sqsubseteq M$, whenever there exists a normalizing reduction ξ starting from t and a normalizing head reduction ρ starting from M such that $\xi \sqsubset \rho$.

As an example, let us consider the term SillyId of (1). Let us define the simple terms:

$$t_1 := \lambda x.\pi_1\langle \pi_1\langle 0, 0 \rangle, \pi_1\langle 0, 0 \rangle \rangle, \quad t_2 := \lambda x.\pi_1\langle \pi_2\langle x, x \rangle, \pi_2\langle x, x \rangle \rangle, \quad t_3 := \lambda x.\pi_2\langle x, x \rangle. \quad (13)$$

$$\begin{array}{c}
\frac{}{\{\cdot\} \triangleleft \{\cdot\}} \quad \frac{}{x^A \triangleleft x^{D(A)}} \quad \frac{M \triangleleft M'}{\lambda x.M \triangleleft \lambda x.M'} \quad \frac{M \triangleleft M' \quad N \triangleleft N'}{MN \triangleleft M'N'} \quad \frac{M \triangleleft M'}{\pi_i^k M \triangleleft \pi_i^k M'} \\
\\
\frac{M_1 \triangleleft M'_1 \quad \dots \quad M_k \triangleleft M'_k}{\langle M_1, \dots, M_k \rangle \triangleleft \langle M'_1, \dots, M'_k \rangle} \quad \frac{P \triangleleft P' \quad M \triangleleft M' \quad N \triangleleft N'}{\text{if}(P, M, N) \triangleleft \text{if}(\pi_1 P', M', N')} \quad \frac{M \triangleleft M'}{\text{fix} f.M \triangleleft \text{fix} f.M'} \\
\\
\frac{z_1^{D(R)} \dots z_k^{D(R)} \vdash t : R^{(\perp_n)} \text{ simple normal form}, \quad M_1 \triangleleft M'_1, \quad \dots, \quad M_k \triangleleft M'_k}{\phi(M_1, \dots, M_k) \triangleleft (\lambda z_1^{D(R)} \dots \lambda z_k^{D(R)}. \langle \phi(\pi_1 z_1, \dots, \pi_1 z_k), t \rangle) M'_1 \dots M'_k}
\end{array}$$

Fig. 6. The expansion relation. In the rule on the third line, if $k = 0$ then the right-hand term in the conclusion is just $\langle \phi, t \rangle$.

Notice that, for any $i \in \{1, 2, 3\}$, we have $t_i \sqsubset \text{Sillyld}$ (see Fig. 4), therefore also $t_i r \sqsubset \text{Sillyld } r$ for any real number r . By contrast, we have that $t_1 r \sqsubseteq \text{Sillyld } r$ iff $r = 0$, $t_2 r \sqsubseteq \text{Sillyld } r$ iff $r < 0$ and, symmetrically, $t_3 r \sqsubseteq \text{Sillyld } r$ iff $r > 0$. This highlights a sharp difference between the relations \sqsubset and \sqsubseteq : the former is static whereas the latter traces the execution of a term. Indeed, the projections of the t_i 's reflect the different choices in the conditionals of Sillyld .

3.2 AD Is Sound on Stable Points

Consider a PCF_R program $x_1^R, \dots, x_n^R \vdash M : R$. One can easily check that for every $\mathbf{r} \in \mathbb{R}^n$, $M\{\mathbf{r}/\mathbf{x}\}$ is normalizing if and only if there exists a simple term t tracing $M\{\mathbf{r}/\mathbf{x}\}$. However, this term t usually depends on the chosen \mathbf{r} . In the following definition of stable points, we consider a situation where t can be “uniformly” chosen in an open ball around \mathbf{r} .

DEFINITION 26. We define the set of stable points of a program $x_1^R, \dots, x_n^R \vdash M : R$ as follows:

$$S(M) := \left\{ \mathbf{r} \in \mathbb{R}^n \mid \exists \varepsilon > 0, \exists x_1^R, \dots, x_n^R \vdash t : R \text{ s.t. } t \sqsubset M \text{ and } \forall \mathbf{r}' \in B_\varepsilon(\mathbf{r}) \ t\{\mathbf{r}'/\mathbf{x}\} \sqsubseteq M\{\mathbf{r}'/\mathbf{x}\} \right\}.$$

Notice that we have restricted the tracing of M to terms t of the same type as M , in particular t does not split different occurrences of a free variable x_i of M . In fact, these variables are supposed to be replaced with numerals, which are not split by \sqsubset . Also observe that, by definition, we have $S(M) \subseteq \Downarrow M$, as $t\{\mathbf{r}/\mathbf{x}\} \sqsubseteq M\{\mathbf{r}/\mathbf{x}\}$ implies that $M\{\mathbf{r}/\mathbf{x}\}$ has a normal form. Moreover, $S(M)$ is open: it is easy to check that $S(M) = \bigcup_{\mathbf{r} \in S(M)} B_{\varepsilon_r}(\mathbf{r})$, where ε_r is the positive real whose existence is given by the very definition of stability of \mathbf{r} .

We already argued in the Introduction that $S(\text{ReLU } x^R) = \mathbb{R} \setminus \{0\}$. By recalling the above discussion about the simple terms t_1 , t_2 and t_3 in (13), we may infer that also $S(\text{Sillyld } x^R) = \mathbb{R} \setminus \{0\}$, in fact 0 is the border where one has to swap between $t_1 r$ and either $t_2 r$ or $t_3 r$ in tracing $\text{Sillyld } r$. Similarly, but with more involved simple terms, one can check that $S(\text{Floor } x^R) = \mathbb{R} \setminus \mathbb{Z}$ with Floor given in (1). As a last example, let us consider the term EqProj defined in (12) and the simple terms

$$t_1 := \lambda x. \lambda y. \pi_1 \langle \pi_1 \langle x, x \rangle, \pi_1 \langle x, x \rangle \rangle, \quad t_2 := \lambda x. \lambda y. \pi_1 \langle \pi_2 \langle y, y \rangle, \pi_2 \langle y, y \rangle \rangle, \quad t_3 := \lambda x. \lambda y. \pi_2 \langle y, y \rangle.$$

These are all such that $t_i \sqsubset \text{EqProj}$. However, $t_1 r q \sqsubseteq \text{EqProj } r q$ iff $r = q$, $t_2 r q \sqsubseteq \text{EqProj } r q$ iff $r < q$ and $t_3 r q \sqsubseteq \text{EqProj } r q$ iff $r > q$. So the diagonal splits the plane \mathbb{R}^2 in two open sets where either t_2 or t_3 uniformly traces the execution of EqProj , whereas the execution on the diagonal is traced by t_1 . But the diagonal contains no open set, therefore $S(\text{EqProj } x^R y^R) = \mathbb{R}^2 \setminus \{(r, r) \mid r \in \mathbb{R}\}$. Notice, finally, that $\text{EqProj } x^R x^R$ may be traced everywhere by $t_1 x^R x^R$, hence $S(\text{EqProj } x^R x^R) = \mathbb{R}$.

The tracing relation \sqsubseteq is defined over programs. However, in order to prove soundness (Theorem 33), we must move from a program M to its transformation $\mathbf{D}(M)$, this latter having a more complex type than M . The difficulty behind the proof of Theorem 33 is then to deduce from $t \sqsubseteq M$ a link between $\mathbf{D}(t)$ and $\mathbf{D}(M)$. In order to do that, we define a further relation \triangleleft (Fig. 6) catching an invariant between M and $\mathbf{D}(M)$ (as well as between t and $\mathbf{D}(t)$) stable under the evaluation of the two terms. Then the extrusion lemma (Lemma 31 and its iterated version Lemma 32) will use \triangleleft for deducing the needed link between $\mathbf{D}(t)$ and $\mathbf{D}(M)$ from the hypothesis $t \sqsubseteq M$.

DEFINITION 27 (EXPANSION). *The relation \triangleleft on terms and contexts of PCF_R is defined in Fig. 6.*

As mentioned above, expansion is used to establish that M and $\mathbf{D}(M)$ “behave similarly”. Such a link emerges from two of the main properties of \triangleleft , namely that $M \triangleleft \mathbf{D}(M)$ holds for every M (Lemma 28), and that it is a simulation (if $M \triangleleft M'$, then M' simulates M , Lemma 29). These justify the definition in the case $M = \phi(M_1, \dots, M_k)$: it mimics the definition of $\mathbf{D}(M)$ in the first component of the product, whereas the second component is chosen as an arbitrary closed simple normal form. The first gives us stability under reduction, in particular in the case of a conditional redex, while the second component cannot be asked to be linked with the partial derivatives of ϕ , because $\phi(M_1, \dots, M_k)$ eventually reduces to a numeral, which has zero derivative.

LEMMA 28. *For every program $x_1^R, \dots, x_n^R \vdash M : R, \mathbf{r} \in \mathbb{R}^n$ and sequence $\mathbf{u} = u_1, \dots, u_n$ of simple closed normal forms of suitable type, we have $M\{\mathbf{r}/\mathbf{x}\} \triangleleft \mathbf{D}(M)\{\langle \mathbf{r}, \mathbf{u} \rangle / \mathbf{x}\}$, where by $\{\langle \mathbf{r}, \mathbf{u} \rangle / \mathbf{x}\}$ we mean $\{\langle r_1, u_1 \rangle / x_1\} \cdots \{\langle r_n, u_n \rangle / x_n\}$.*

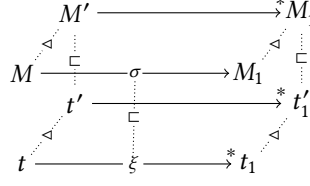
LEMMA 29. *Let $M \triangleleft M'$ and $M \rightarrow N$, then there exists N' such that $M' \rightarrow^* N'$ and $N \triangleleft N'$.*

PROOF SKETCH. Let (C, R, P) be the reduction step $M \rightarrow N$. The proof is an induction on C . The only non-trivial part is the base of the induction, i.e. $C = \{\cdot\}$, in which the reasoning splits following Fig. 1c: the case of a β -reduction is a consequence of a substitution lemma on \triangleleft . If $R = \phi(\mathbf{r})$, then $M' \rightarrow^* \langle \phi(\pi_1 L'), t\{L'/z\} \rangle$ for some L' of the same length as \mathbf{r} such that $r_i \triangleleft L'_i$ for all i . Notice that $r_i \triangleleft L'_i$ implies that $\pi_1 L'_i \rightarrow r_i$ as well as L'_i is a simple closed normal form, so in particular $t\{L'/z\}$ is normalizable by Proposition 5. We therefore have: $\langle \phi(\pi_1 L'), t\{L'/z\} \rangle \rightarrow^* \langle \llbracket \phi \rrbracket(\mathbf{r}), t' \rangle$ with t' a simple closed normal form, and we conclude by taking $N' := \langle \llbracket \phi \rrbracket(\mathbf{r}), t' \rangle$.

The cases of the other redexes (conditional, products and fixpoints) are immediate. \square

REMARK 30. *A variant of Lemma 29 can be used to prove that $\Downarrow^S M \subseteq \Downarrow^S \text{grad}(M)$, for some reduction strategy \mathcal{S} (see discussion after Proposition 10). For example, if \mathcal{S} is the call-by-value strategy, then one can prove the statement of Lemma 29 by replacing \rightarrow and \rightarrow^* with $\xrightarrow{\mathcal{S}}$ and $\xrightarrow{\mathcal{S}}^*$ (in the case of a $R = \phi(\mathbf{r})$ redex, one should notice that the terms L' are all values, and one should replace “normalizable” with “ \mathcal{S} -normalizable” and “simple closed normal form” with “simple closed \mathcal{S} -nf”). Then, consider a program M of arity 1 and suppose $\mathbf{r} \in \Downarrow^S M$, i.e. $M\{\mathbf{r}/\mathbf{x}\} \xrightarrow{\mathcal{S}}^* q$, for q a numeral. We have by Lemma 28 and this variant of Lemma 29 that $\mathbf{D}(M)\{\langle \mathbf{r}, \mathbf{u} \rangle / \mathbf{x}\} \xrightarrow{\mathcal{S}}^* N$ with $q \triangleleft N$ and each u_i either the normal form of $i_1^n 1$ or i_1^n , depending whether \mathbf{D} is $\overrightarrow{\mathbf{D}}$ or $\overleftarrow{\mathbf{D}}$ (recall Equations (9) and (10)). By inspecting Fig. 6, we can deduce $N = \langle q, t \rangle$ for some closed simple \mathcal{S} -nft. This means $\overleftarrow{\text{grad}}(M) \xrightarrow{\mathcal{S}}^* (\pi_2 \langle q, t \rangle) 1 \xrightarrow{\mathcal{S}} t 1$, this latter evaluating to a normal form because it is a simple term (Proposition 5). We conclude $\mathbf{r} \in \Downarrow^{\overleftarrow{\text{grad}}} (M)$. The case of $\overrightarrow{\text{grad}}(M)$ is simpler.*

LEMMA 31 (EXTRUSION). *Let $M \triangleleft M', t \triangleleft t', t' \sqsubset M'$. Let $\sigma : M \rightarrow M_1$ be a head reduction step and moreover let $\xi : t \rightarrow^* t_1$ be such that $\xi \sqsubset \sigma$ (so in particular $t \sqsubset M$ and $t_1 \sqsubset M_1$). Then there exist M'_1, t'_1 such that the following relations hold:*



PROOF SKETCH. By Definition 23, $\sigma = H\{\sigma_0\}$ for some head context H and reduction step $\sigma_0 : R \rightarrow P$. The proof is by induction on H .

The case $H = \{\cdot\}$ splits following Fig. 1c. For example, let σ be $M = (\lambda x.L)N \rightarrow L\{N/x\} = M_1$, so that ξ is the reduction $t = (\lambda p.w\{\pi p/x\})\langle u \rangle \rightarrow^* w\{u/x\} = t_1$. Then $M' = (\lambda x.L')N'$ with $L < L'$ and $N < N'$, and $t' = (\lambda p.\bar{w}')\langle u' \rangle$ with $w\{\pi p/x\} < \bar{w}'$ and $\langle u \rangle < \langle u' \rangle$. Moreover, since $t' \sqsubset M'$, by Lemma 20, we have $\bar{w}' = w'\{\pi p/x\}$, with $w'\{\pi p/x\} \sqsubset L'$, p not free in w' , and $u'_i \sqsubset N'$. Moreover, by induction on w , one can infer from $w\{\pi p/x\} < \bar{w}'$ that actually $w < w'$.

Of the induction cases, the only subtle one is $H = \text{if}(\bar{H}, N_1, N_2)$. Under this hypothesis, the reduction ξ is empty and $t_1 = t = \pi_i\langle u, u \rangle$ with $u \sqsubset N_i$ for some $i \in \{1, 2\}$, as well as $M' = \text{if}(\pi_1\bar{M}', N'_1, N'_2)$ with $\bar{H}\{R\} < \bar{M}'$, $N_1 < N'_1$, $N_2 < N'_2$ and $t' = \pi_i\langle u', u' \rangle$ with $u < u'$ and $u' \sqsubset N'_i$ (notice that the index i of the projection is the same in t and t' because $t < t'$). By Lemma 29, $\bar{M}' \rightarrow^* L$ such that $\bar{H}\{P\} < L$. We can then conclude by setting $M'_1 = \text{if}(\pi_1 L, N'_1, N'_2)$ and $t'_1 = t'$.

All of the remaining cases follow the same pattern. \square

LEMMA 32 (EXTRUSION TO NORMAL FORM). *Let $M < M'$, $t < t'$, $t \sqsubset M$ and $t' \sqsubset M'$, for t and M closed terms both of type R^n , for some $n \geq 0$. Then, $t' \rightarrow^* t''$ and $M' \rightarrow^* M''$ with t'' and M'' normal such that $t'' \sqsubset M''$.*

PROOF SKETCH. By Definition 25, there exist a normalizing reduction ξ starting from t and a normalizing reduction ρ starting from M , such that $\xi \sqsubset \rho$. The proof is by induction on ρ . \square

THEOREM 33 (SOUNDNESS). *For every program $\Gamma \vdash M : R$ and every $\mathbf{r} \in S(M) \cap d(M)$, we have:*

$$\text{grad}(M)(\mathbf{r}) \rightarrow^* \nabla(\llbracket M \rrbracket_\Gamma)(\mathbf{r}).$$

PROOF. The assumption $\mathbf{r} \in d(M)$ tells us that there is an open ball $B_{\varepsilon_0}(\mathbf{r})$ where $\nabla\llbracket M \rrbracket_\Gamma$ exists. By Definition 26, $\mathbf{r} \in S(M)$ means that there is a simple program $t \sqsubset M$ uniformly tracing M in an open ball of \mathbf{r} , i.e., there exists $\varepsilon > 0$ such that for all $\mathbf{r}' \in B_\varepsilon(\mathbf{r})$ we have $t\{\mathbf{r}'/x\} \sqsubset M\{\mathbf{r}'/x\}$, and of course we may take $\varepsilon \leq \varepsilon_0$. This implies in particular that $\llbracket t \rrbracket_\Gamma$ and $\llbracket M \rrbracket_\Gamma$ coincide on $B_\varepsilon(\mathbf{r})$ and, therefore, we have $\mathbf{r} \in d(t)$ as well.

By Lemma 28, $t\{\mathbf{r}/x\} < D(t)\{\langle \mathbf{r}, u \rangle/x\}$ and $M\{\mathbf{r}/x\} < D(M)\{\langle \mathbf{r}, u \rangle/x\}$, where u_i is either the normal form of $l_i^n 1$ or l_i^n , depending whether D is \overrightarrow{D} or \overleftarrow{D} (recall Equations (9) and (10)). Moreover, by Lemma 19 we have that $t \sqsubset M$ gives $D(t)\{\langle \mathbf{r}, u \rangle/x\} \sqsubset D(M)\{\langle \mathbf{r}, u \rangle/x\}$. We are then in position of applying Lemma 32 to $M\{\mathbf{r}/x\}$, $t\{\mathbf{r}/x\}$ (closed terms of ground type) and $D(t)\{\langle \mathbf{r}, u \rangle/x\}$, $D(M)\{\langle \mathbf{r}, u \rangle/x\}$. This gives us a normal form t' of $D(t)\{\langle \mathbf{r}, u \rangle/x\}$ and M' of $D(M)\{\langle \mathbf{r}, u \rangle/x\}$ such that $t' \sqsubset M'$. By subject reduction (Proposition 4), t', M' are closed normal forms of type $D(R)$, so Lemma 18 gives us $t' = M'$. Since $\mathbf{r} \in d(t)$, we conclude by Proposition 9. \square

4 UNSOUNDNESS OF AD

DEFINITION 34 (UNSTABLE POINT). *The set of unstable points of a program M , denoted by $U(M)$, is the complement of $S(M)$ (Definition 26) in $\llbracket M \rrbracket$, i.e., $U(M) := \llbracket M \rrbracket \setminus S(M)$.*

We know from the remark after Definition 26 that $U(M)$ is closed in $\llbracket M \rrbracket$ (with respect to the subspace topology). The goal of this section is to prove that it is a quasivariety, hence of

$$\begin{aligned}
P_\Gamma(R) &:= \{\Gamma \vdash M : R \mid \llbracket M \rrbracket_\Gamma \text{ is cqc and } U(M) \text{ is a quasivariety}\} \\
P_\Gamma(A \rightarrow B) &:= \{\Gamma \vdash M : A \rightarrow B \mid \forall N \in P_\Gamma(A), MN \in P_\Gamma(B)\} \\
P_\Gamma(A_1 \times \dots \times A_k) &:= \{\Gamma \vdash M : A_1 \times \dots \times A_k \mid \pi_i M \in P_\Gamma(A_i), \forall i \leq k\}
\end{aligned}$$

Fig. 7. The definition of the logical predicate $P_\Gamma(A)$, with Γ a ground context.

measure zero. The main tool is the logical predicate defined in Fig. 7 and its adequacy (Lemma 41). The structure of the proof is standard, but some new notions are needed. First, our programs are first-order functions, *i.e.*, terms with some free variables of ground type, so the logical predicate is indexed by a typing environment. Second, Lemma 35 states some properties of the notion of stability necessary to achieve the standard auxiliary lemmas of a logical predicate, such as closure under expansion (Lemma 37) or (a syntactic variant of) Scott-continuity (Lemma 39 and Lemma 40). Third, and more important, the standard lemma of logical predicates for the conditional (Lemma 38) is particularly subtle. This should not come as a surprise: as discussed above, the possibility of unsoundness of AD is due to conditionals. In particular, let us underline that Lemma 38 uses the notion of completely quasicontinuous map introduced in Sect. 2.3.

The logical predicate on which the proof is based is defined in Fig. 7. In the rest of the section, unless otherwise stated, $\Gamma := x_1^R, \dots, x_n^R$ is a ground context. Moreover, if $M : A_1 \rightarrow \dots \rightarrow A_p \rightarrow B$ and $\mathbf{L} = L_1, \dots, L_p$ such that $L_i : A_i$ for all $1 \leq i \leq p$, then the notation $M\mathbf{L}$ stands for $M L_1 \dots L_p$, which is of course a term of type B .

LEMMA 35. *We have the following inclusions, where the terms appearing in the statements are supposed to be typed under a ground context Γ .*

- (1) Let ϕ be a function symbol of arity k and let M_1, \dots, M_k be programs, then $\bigcap_i S(M_i) \subseteq S(\phi(M_1, \dots, M_k))$.
- (2) Let $R \rightarrow P$ be one of the rewriting rules in Fig. 1c, with R, P of type $B_1 \rightarrow \dots \rightarrow B_p \rightarrow R^m$. For all $1 \leq i \leq p$, let $\Gamma \vdash L_i : B_i$. Then, for all $1 \leq j \leq m$, $S(\pi_j(PL)) \subseteq S(\pi_j(RL))$.
- (3) Let $P : R$ and M_1, M_2 be of type $B_1 \rightarrow \dots \rightarrow B_p \rightarrow R^m$. For all $1 \leq i \leq p$, let $\Gamma \vdash L_i : B_i$. Let $X_1 := \llbracket P \rrbracket_\Gamma^{-1}(\mathbb{R}_{\leq 0})$ and $X_2 := \llbracket P \rrbracket_\Gamma^{-1}(\mathbb{R}_{> 0})$. Then, for all $1 \leq j \leq m$ and all $l \in \{1, 2\}$, we have that $S(\pi_j(M_l\mathbf{L})) \cap \text{int}(X_l) \subseteq S(\pi_j(\text{if}(P, M_1, M_2)\mathbf{L}))$.
- (4) Let $B = B_1 \rightarrow \dots \rightarrow B_p \rightarrow R^m$, let $\Gamma \vdash L_0 : A$ and $\Gamma \vdash L_i : B_i$ for all $1 \leq i \leq p$. For all $k \in \mathbb{N}$ and $1 \leq j \leq m$, $S(\pi_j((\text{fix}_k f^{A \rightarrow B}.M)\mathbf{L})) \subseteq S(\pi_j((\text{fix} f^{A \rightarrow B}.M)\mathbf{L}))$, where $\mathbf{L} := L_0, L_1, \dots, L_p$.

PROOF SKETCH. We only detail the proof of the branching case, the other cases are easy variants.

By taking the notations of point (3), let $l \in \{1, 2\}$ and let $\mathbf{r} \in S(\pi_j(M_l\mathbf{L})) \cap \text{int}(X_l)$. By definition, there exist $\varepsilon > 0$, $u \sqsubset \pi_j(M_l\mathbf{L})$ such that, for all $\mathbf{r}' \in B_\varepsilon(\mathbf{r})$, $u\{\mathbf{r}'/\mathbf{x}\} \sqsubset \pi_j(M_l\mathbf{L})\{\mathbf{r}'/\mathbf{x}\}$ and $\mathbf{r}' \in \text{int}(X_l)$. Notice that $u = \pi_j(u'\mathbf{u}'')$, with $u' \sqsubset M_l$ and $\mathbf{u}'' = \langle u'_1 \rangle \dots \langle u'_p \rangle$ such that, for every $1 \leq i \leq p$ and every element $u''_{i,h}$ of \mathbf{u}''_i , we have $u''_{i,h} \sqsubset L_i$. Let $t := \pi_j((\pi_\ell \langle u', u' \rangle)\mathbf{u}'')$ and notice that $t \sqsubset \pi_j(\text{if}(P, M_1, M_2)\mathbf{L})$. Let us prove that $t\{\mathbf{r}'/\mathbf{x}\} \sqsubset \pi_j(\text{if}(P, M_1, M_2)\mathbf{L})\{\mathbf{r}'/\mathbf{x}\}$.

Since $\mathbf{r}' \in \text{int}(X_l) \subseteq \llbracket P \rrbracket$, we have that $P\{\mathbf{r}'/\mathbf{x}\}$ is normalizing. Since P is ground, by Proposition 6 there is a head reduction sequence $\rho : P\{\mathbf{r}'/\mathbf{x}\} \rightarrow^* q$ such that q is a numeral. Let now $H = \pi_j(\text{if}(\{\cdot\}, M_1, M_2)\mathbf{u}'')\{\mathbf{r}'/\mathbf{x}\}$. Notice that $v \sqsubset H\{\rho\}$ for v the empty reduction sequence of $t\{\mathbf{r}'/\mathbf{x}\}$. Moreover, by hypothesis we have normalizing reduction sequences $v' \sqsubset \rho'$ from $u\{\mathbf{r}'/\mathbf{x}\} = \pi_j(u'\mathbf{u}'')\{\mathbf{r}'/\mathbf{x}\}$ and $\pi_j(M_l\mathbf{L})\{\mathbf{r}'/\mathbf{x}\}$, respectively. Furthermore, we have the head reduction steps:

$$v_0 : \pi_j(\pi_\ell \langle u', u' \rangle \mathbf{u}'')\{\mathbf{r}'/\mathbf{x}\} \rightarrow \pi_j(u'\mathbf{u}'')\{\mathbf{r}'/\mathbf{x}\}, \quad \rho_0 : \pi_j(\text{if}(q, M_1, M_2)\mathbf{L})\{\mathbf{r}'/\mathbf{x}\} \rightarrow \pi_j(M_l\mathbf{L})\{\mathbf{r}'/\mathbf{x}\}$$

such that $v_0 \sqsubset \rho_0$. We then have $v v_0 v' \sqsubset H\{\rho\} \rho_0 \rho'$, which allows us to conclude. \square

LEMMA 36 (FUNCTION SYMBOLS). *Let ϕ be a function symbol of arity k and, for each $1 \leq i \leq k$, $M_i \in P_\Gamma(R)$. If $\llbracket \phi \rrbracket$ is cqc, then $\phi(M_1, \dots, M_k) \in P_\Gamma(R)$.*

LEMMA 37 (CLOSURE UNDER EXPANSION). *Let $R \rightarrow P$ be one of the rewriting rules in Figure 1c. If $P \in P_\Gamma(A)$, then $R \in P_\Gamma(A)$.*

LEMMA 38 (CONDITIONAL). *If $P \in P_\Gamma(R)$ and $M_1, M_2 \in P_\Gamma(A)$, then $\text{if}(P, M_1, M_2) \in P_\Gamma(A)$.*

PROOF. Let $A = A_1 \rightarrow \dots \rightarrow A_p \rightarrow R^m$. It is enough to prove that, given $L = L_1, \dots, L_p$ such that $L_i \in P_\Gamma(A_i)$ for every $1 \leq i \leq p$, we have $\pi_j(\text{if}(P, M_1, M_2)L) \in P_\Gamma(R)$ for every $1 \leq j \leq m$.

Let $N := \pi_j(\text{if}(P, M_1, M_2)L)$, $Q_1 := \llbracket P \rrbracket_\Gamma^{-1}(\mathbb{R}_{\leq 0})$ and $Q_2 := \llbracket P \rrbracket_\Gamma^{-1}(\mathbb{R}_{> 0})$. Observe that both Q_1 and Q_2 are quasiopen, because $\mathbb{R}_{\leq 0}$ and $\mathbb{R}_{> 0}$ are quasiopen and $\llbracket P \rrbracket_\Gamma$ is cqc by hypothesis. Furthermore, $\mathbb{R}^l \times Q_i$ is quasiopen for every $l \geq 0$ and $i \in \{1, 2\}$, because it is the inverse image of Q_i via a projection $\mathbb{R}^l \times \mathbb{R}^p \rightarrow \mathbb{R}^p$, and these are cqc (Lemma 17.3).

To prove that $\llbracket N \rrbracket_\Gamma$ is cqc we need to show that, for every $l \geq 0$ and every quasiopen set $Q \subseteq \mathbb{R}^{l+1}$, the set $(\text{id}_{\mathbb{R}^l} \times \llbracket N \rrbracket_\Gamma)^{-1}(Q)$ is quasiopen. But Proposition 7 tells us that $(\text{id}_{\mathbb{R}^l} \times \llbracket N \rrbracket_\Gamma)^{-1}(Q)$ is equal to $((\mathbb{R}^l \times Q_1) \cap (\text{id}_{\mathbb{R}^l} \times \llbracket \pi_j(M_1 L) \rrbracket_\Gamma)^{-1}(Q)) \cup ((\mathbb{R}^l \times Q_2) \cap (\text{id}_{\mathbb{R}^l} \times \llbracket \pi_j(M_2 L) \rrbracket_\Gamma)^{-1}(Q))$, which is quasiopen because $\llbracket \pi_j(M_1 L) \rrbracket_\Gamma$ and $\llbracket \pi_j(M_2 L) \rrbracket_\Gamma$ are cqc by hypothesis.

For what concerns the unstable points, we first observe that, by Proposition 7, $\Downarrow N = \Downarrow P \cap ((Q_1 \cap \Downarrow M_1) \cup (Q_2 \cap \Downarrow M_2)) = (Q_1 \cap \Downarrow M_1) \cup (Q_2 \cap \Downarrow M_2)$, the second equality holding because $Q_1, Q_2 \subseteq \Downarrow P$. We may therefore write

$$U(N) = \left(\bigcup_{i \in \{1, 2\}} Q_i \cap \Downarrow M_i \right) \setminus S(N) = \bigcup_{i \in \{1, 2\}} (Q_i \cap \Downarrow M_i) \setminus S(N) \subseteq \bigcup_{i \in \{1, 2\}} (Q_i \cap \Downarrow M_i) \setminus (\text{int}(Q_i) \cap S(M_i))$$

where the inclusion is by Lemma 35.3. Now, for all $i \in \{1, 2\}$, let us write $A_i := (Q_i \cap \Downarrow M_i) \setminus (\text{int}(Q_i) \cap S(M_i))$. Notice that, by Lemma 15, $Q_i = \text{int}(Q_i) \cup Z_i$ with Z_i a quasivariety, hence

$$A_i = ((Q_i \cap \Downarrow M_i) \setminus \text{int}(Q_i)) \cup ((Q_i \cap \Downarrow M_i) \setminus S(M_i)) \subseteq (Q_i \setminus \text{int}(Q_i)) \cup (\Downarrow M_i \setminus S(M_i)) = Z_i \cup U(M_i),$$

so A_i is a quasivariety, because Z_i and $U(M_i)$ are. Since $U(N) \subseteq A_1 \cup A_2$, we are done. \square

LEMMA 39 (DIVERGENCE). *For every type $A \rightarrow B$, $\Omega_{A \rightarrow B} \in P_\Gamma(A \rightarrow B)$.*

LEMMA 40 (FIXPOINTS). *If $\forall k \in \mathbb{N}$, $\text{fix}_k f.M \in P_\Gamma(A \rightarrow B)$, then $\text{fix} f.M \in P_\Gamma(A \rightarrow B)$.*

LEMMA 41 (ADEQUACY). *Suppose that the primitive functions of PCF_R are admissible. Let $\Gamma := x_1^R, \dots, x_n^R, \Delta := y_1^{A_1}, \dots, y_m^{A_m}$, let $\Gamma, \Delta \vdash M : A$ and let $\Gamma \vdash N_i \in P_\Gamma(A_i)$ for all $1 \leq i \leq m$. Then,*

$$M\{N_1/y_1\} \cdots \{N_m/y_m\} \in P_\Gamma(A).$$

THEOREM 42. *Assuming that the primitive functions of PCF_R are admissible, for every program $\Gamma \vdash M : R$ the set*

$$\text{Fail}(M) := \{\mathbf{r} \in d(M) ; \text{grad}(M)(\mathbf{r}) \not\vdash^* \nabla(\llbracket M \rrbracket_\Gamma)(\mathbf{r})\}$$

is a quasivariety, hence of measure zero.

PROOF. By Theorem 33, we know that $\text{Fail}(M) \subseteq U(M)$, which is a quasivariety because $M \in P_\Gamma(R)$ by Lemma 41. \square

5 DISCUSSION AND PERSPECTIVES

On the significance of the measure zero bound. Since the set of real numbers representable on an actual computer is of measure zero in \mathbb{R} (it is finite!), one may feel skeptical about the significance of Theorem 42. This issue was already raised by Speelpenning [Speelpenning 1980] while commenting on Joss’s theorem [Joss 1976]. He defines a program similar to the following:

$$\text{SlowId} := \lambda x^{\mathbb{R}}. \left(\text{fix } f^{\mathbb{R} \rightarrow \mathbb{R}}. \lambda y^{\mathbb{R}}. \text{if}(x - y, \text{if}(y - x, y, f \text{ Next}), f \text{ Next}) \right) 0,$$

where $\vdash \text{Next} : \mathbb{R}$ is a primitive which cycles through machine-representable real numbers, based on some internal state (Speelpenning uses a random number generator in his example). So, given $r \in \mathbb{R}$, $\text{SlowId}(r)$ will eventually output r if this is machine-representable, and diverge otherwise. When executed on an actual computer, every r is necessarily representable and SlowId behaves like the identity. And yet, $\text{grad}(\text{SlowId } x)(r) \rightarrow^* 0$ for every machine-representable r , because Next is treated as a constant by AD transformations (there is no sensible alternative). Speelpenning concludes that, although SlowId is not a counterexample to Joss’s theorem (or to ours), from the practical viewpoint AD fails *everywhere* on it.

We believe that Speelpenning’s example is misleading. The reason why SlowId is not a counterexample to Theorem 42 is not that the set where AD fails on SlowId is of measure zero because it coincides with the set of machine-representable reals; it is because $\llbracket \text{SlowId}(x) \rrbracket_{x^{\mathbb{R}}}$ is nowhere differentiable! That is, in the notations of Theorem 42, we actually have $\text{Fail}(\text{SlowId}(x)) = \emptyset$ because $d(\text{SlowId}(x)) = \emptyset$, and this is because $\llbracket \text{SlowId}(x) \rrbracket_{x^{\mathbb{R}}}$ is defined only on a discrete set.

Anyway, if the set $R := \{r_1, \dots, r_c\}$ of machine-representable reals is finite, there are impractically large but straightforward programs achieving the intended behavior of Speelpenning’s example. For instance, with some syntactic sugar, define

$$\text{SlowId}' := \lambda x^{\mathbb{R}}. \text{if } x = r_1 \text{ then } r_1 \text{ else } (\dots \text{if } x = r_c \text{ then } r_c \text{ else } x \dots).$$

We have that $\llbracket \text{SlowId}'(x) \rrbracket_{x^{\mathbb{R}}}$ is actually the identity function, so $\text{Fail}(\text{SlowId}'(x)) = R$ and we may legitimately say that AD is wrong “everywhere”. But this is just a giant-sized version of the program SillyId of the Introduction, and speaks more of the contrivance of toying with $\text{PCF}_{\mathbb{R}}$ as a machine-executable language (which it is not) than of the value of our result. In general, questioning the significance of Theorem 42 on the grounds that computers are finite is like questioning Turing machines because of their infinite tape, or objecting to the whole idea of studying the asymptotic complexity of programs because in practice we only implement finite functions, whose asymptotic complexity is $O(1)$. In our opinion, there is little point in discussing this standpoint further.

More constructively, we may argue that the significance of Theorem 42 lies in the fact that it gives a finer bound than just measure zero. Let us call $\text{PCF}_{\mathbb{R}}$ with only the “mandatory” primitive functions (constants, addition, multiplication) *minimal* $\text{PCF}_{\mathbb{R}}$. This is already enough to express all differentiable programming architectures based on neural networks with rectified linear unit activation. Moreover, by using Taylor series, minimal $\text{PCF}_{\mathbb{R}}$ may also approximate every analytic function with arbitrary precision, so it has a wide range of potential applications. Theorem 42 tells us that, if $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a function definable in minimal $\text{PCF}_{\mathbb{R}}$, then the set of points on which ∇f exists but AD methods fail to compute it is contained in a countable union of algebraic varieties (i.e., zeros of polynomials). In particular, when $n = 1$, this set is countable.

Zero sets of polynomial equations have been studied for literally millennia as part of the vast field known as algebraic geometry. Albeit extremely complex in general, many results exist on their structure, which may be described or approximated very accurately in several cases. It is not excluded that, in the future, these results may be leveraged to develop static analysis techniques (e.g. type systems) for establishing the absence of errors in differentiable programs.

From Gradients to Jacobians. We limited our attention to programs implementing functions $\mathbb{R}^n \rightarrow \mathbb{R}^m$ with $m = 1$. The case $m > 1$, in which one would speak of Jacobians rather than gradients, is conceptually identical. First of all, observe that a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ may always be decomposed into m functions $f_i := \pi_i f : \mathbb{R}^n \rightarrow \mathbb{R}$, so restricting primitives to one output causes no loss of generality and Fig. 2 needs no modification. When $\Gamma \vdash M : \mathbb{R}^m$ with Γ containing n variables and $m > 1$, what needs to be modified are the Equations (9) and (10), which must yield an $m \times n$ matrix whose lines are built out of m expressions of the form $\text{grad}(\pi_i M)$. Theorem 33 and Theorem 42 lift to this setting because the Jacobian is just the collection of the m gradients.

Internalizing AD. The transformations of Fig. 2, and thus the definition of $\text{grad}(M)$ for a program M (Equations (9) and (10)) are external to PCF_R : a programmer may apply them for instance via a compiler, but the transformations are not accessible from within the program itself. For practical purposes, it would be interesting to have a programming language in which $\overrightarrow{\text{grad}}$ and $\overleftarrow{\text{grad}}$ are syntactic constructs, typed $x_1^R, \dots, x_n^R \vdash \text{grad}(M) : \mathbb{R}^n$ whenever $x_1^R, \dots, x_n^R \vdash M : \mathbb{R}$, and with $\text{grad}(M)$ being executed in such a way as to reflect the application of AD to M . A naive way of achieving this would be to turn the definition of Fig. 2 into rewriting rules; a more sophisticated approach was provided by Pearlmutter and Siskind for $\text{Stalin}\nabla$ [Pearlmutter and Siskind 2008].

The errors introduced by AD have the important consequence that *such an internalization is impossible without breaking the expected extensional semantics of programs*. This is because, as any denotational semantics, the standard semantics defined in Sect. 2.1 is contextual, in the sense that $\llbracket M \rrbracket = \llbracket N \rrbracket$ implies $\llbracket C\{M\} \rrbracket = \llbracket C\{N\} \rrbracket$ for any context C . Now, referring to (1), we have $\llbracket \text{SillyId}(x) \rrbracket_{x^R} = \llbracket x \rrbracket_{x^R}$ and yet we know that $\llbracket \text{grad}(\text{SillyId}(x)) \rrbracket_{x^R} \neq \llbracket \text{grad}(x) \rrbracket_{x^R}$, because the latter two functions differ at 0. So any denotational semantics of PCF_R with “internal AD” needs to interpret SillyId and $\text{Id} := \lambda x^R. x$ differently. “Resource-sensitive” semantics coming from linear logic do distinguish them, but it is easy to find other examples on which these semantics too fail. An example of denotational semantics which consistently works is the one introduced by Abadi and Plotkin [Abadi and Plotkin 2020], whose first order language does have internal AD. In that semantics, Id is the identity whereas SillyId is a “partial identity”, undefined at 0. It is not clear whether this extends to higher order (and thus to PCF_R , similarly to [Di Gianantonio and Edalat 2013]), but assuming it does, the meaning it gives to programs is somewhat unusual: for example, using the definition given in (1), $\text{Floor}(r)$ would diverge whenever r is an integer. This is a further drawback of partial conditional semantics, in addition to the one pointed out in Sect. 2.2 (concerning example (12)).

Another loosely related remark worth making at this point is that, seen as a functional on the Scott domain $\mathbb{R}_\perp \rightarrow \mathbb{R}_\perp$, where \mathbb{R}_\perp is the “flat” Scott domain typically used for interpreting \mathbb{R} as a ground type with total conditionals, the derivative operator ∂ is *not* Scott continuous.⁸ Although the technical consequences of this observation are not entirely clear, from an intuitive point of view it means that in the naive flat semantics the derivative operator is “not computable”, and therefore no recursive procedure (like those given by AD transformations) will be error-free.

ACKNOWLEDGMENTS

We would like to thank A. Brunel, T. Ehrhard and B.A. Pearlmutter for useful comments and discussions. This work was partially supported by ANR PRC project PPS (ANR-19-CE48-0014).

⁸Given $A \subseteq \mathbb{R}$, say that χ is the *indicator function* of A if $\chi(x) = 0$ when $x \in A$ and $\chi(x) = \perp$ otherwise. For $n > 0$, let $I_n :=] - \infty, 0] \cup] \frac{1}{n}, +\infty[$ and let φ_n be the indicator function of I_n . Notice that each φ_n is differentiable on $J_n := I_n \setminus \{0\}$ and $\partial\varphi_n$ is the indicator function of J_n . As elements of the Scott domain $\mathbb{R}_\perp \rightarrow \mathbb{R}_\perp$, the functions $(\varphi_n)_{n>0}$ and $(\partial\varphi_n)_{n>0}$ form two directed chains whose suprema are the identically zero function and the indicator function of $\mathbb{R} \setminus \{0\}$, respectively. In particular, $\partial(\sup_{n>0} \varphi_n) \neq \sup_{n>0} \partial\varphi_n$, so ∂ is not Scott continuous.

REFERENCES

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek Gordon Murray, Benoit Steiner, Paul A. Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2016. TensorFlow: A System for Large-Scale Machine Learning. In *Proceedings of OSDI*. USENIX Association, 265–283.
- Martín Abadi and Gordon D. Plotkin. 2020. A simple differentiable programming language. *Proc. ACM Program. Lang.* 4, POPL (2020), 38:1–38:28.
- Roberto Amadio and Pierre-Louis Curien. 1998. *Domains and Lambda-Calculi*. Vol. 46. Cambridge University Press.
- Henk P. Barendregt. 1985. *The Lambda Calculus, Its Syntax and Semantics*. North Holland.
- Gilles Barthe, Raphaëlle Crubillé, Ugo Dal Lago, and Francesco Gavazzo. 2020. On the Versatility of Open Logical Relations - Continuity, Automatic Differentiation, and a Containment Theorem. In *Proceedings of ESOP*. 56–83.
- Atilim Baydin, Barak Pearlmutter, Alexey Radul, and Jeffrey Siskind. 2018. Automatic differentiation in machine learning: A survey. *Journal of Machine Learning Research* 18 (2018), 1–43.
- Thomas Beck and Herbert Fischer. 1994. The if-problem in automatic differentiation. *J. Comput. Appl. Math.* 50, 1 (1994), 119–131.
- Aloïs Brunel, Damiano Mazza, and Michele Pagani. 2020. Backpropagation in the simply typed lambda-calculus with linear negation. *PACMPL* 4, POPL (2020), 64:1–64:27.
- Pietro Di Gianantonio and Abbas Edalat. 2013. A Language for Differentiable Functions. In *Proceedings of FOSSACS*. 337–352.
- Thomas Ehrhard and Laurent Regnier. 2006. Böhm Trees, Krivine’s Machine and the Taylor Expansion of Lambda-Terms. In *Proceedings of CIE*. 186–197.
- Thomas Ehrhard and Laurent Regnier. 2008. Uniformity and the Taylor expansion of ordinary lambda-terms. *Theor. Comput. Sci.* 403, 2-3 (2008), 347–372.
- Conal Elliott. 2018. The simple essence of automatic differentiation. *PACMPL* 2, ICFP (2018), 70:1–70:29.
- Martín Hötzel Escardó. 1996. PCF Extended with Real Numbers. *Theor. Comput. Sci.* 162, 1 (1996), 79–115.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press. <http://www.deeplearningbook.org>
- Andreas Griewank and Andrea Walther. 2008. *Evaluating derivatives - principles and techniques of algorithmic differentiation, Second Edition*. SIAM.
- Mathieu Huot, Sam Staton, and Matthijs Vákár. 2020. Correctness of Automatic Differentiation via Diffeologies and Categorical Gluing. In *Proceedings of FOSSACS*. 319–338.
- Johan Joss. 1976. *Algorithmisches Differenzieren*. Ph.D. Dissertation. ETH Zurich.
- Yann LeCun. 2018. Deep Learning est mort. Vive Differentiable Programming! (2018). <https://www.facebook.com/yann.lecun/posts/10155003011462143>
- Wonyeol Lee, Hangyeol Yu, Xavier Rival, and Hongseok Yang. 2020. On Correctness of Automatic Differentiation for Non-Differentiable Functions. *CoRR* abs/2006.06903 (2020).
- Carol Mak, C.-H. Luke Ong, Hugo Paquet, and Dominik Wagner. 2020. Densities of almost-surely terminating probabilistic programs are differentiable almost everywhere. *CoRR* abs/2004.03924 (2020).
- Damiano Mazza. 2017. *Polyadic Approximations in Logic and Computation*. Habilitation thesis. Université Paris 13.
- Damiano Mazza, Luc Pellissier, and Pierre Vial. 2018. Polyadic approximations, fibrations and intersection types. *Proc. ACM Program. Lang.* 2, POPL (2018), 6:1–6:28.
- Boris Mityagin. 2015. The Zero Set of a Real Analytic Function. *arXiv:1512.07276 [math.CA]* (2015).
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in PyTorch. (2017).
- Barak A. Pearlmutter and Jeffrey Mark Siskind. 2008. Reverse-mode AD in a Functional Framework: Lambda the Ultimate Backpropagator. *ACM Trans. Program. Lang. Syst.* 30, 2 (2008), 7:1–7:36.
- Gordon Plotkin. 1977. LCF Considered as a Programming Language. *Theoretical Computer Science* 5, 3 (1977), 223–255.
- Amir Shaikhha, Andrew Fitzgibbon, Dimitrios Vytiniotis, and Simon Peyton Jones. 2019. Efficient differentiable programming in a functional array-processing language. *PACMPL* 3, ICFP (2019), 97:1–97:30.
- Bert Speelpenning. 1980. *Compiling Fast Partial Derivatives of Functions Given by Algorithms*. Ph.D. Dissertation. University of Illinois at Urbana-Champaign.
- Agnes Szendrei. 1986. *Clones in Universal Algebra*. Presses de l’Université de Montréal.
- Fei Wang, Daniel Zheng, James M. Decker, Xilun Wu, Grégory M. Essertel, and Tiark Rompf. 2019. Demystifying differentiable programming: shift/reset the penultimate backpropagator. *PACMPL* 3, ICFP (2019), 96:1–96:31.
- Yuan Zhou, Bradley J. Gram-Hansen, Tobias Kohn, Tom Rainforth, Hongseok Yang, and Frank Wood. 2019. LF-PPL: A Low-Level First Order Probabilistic Programming Language for Non-Differentiable Models. In *Proceedings of AISTATS*. 148–157.