

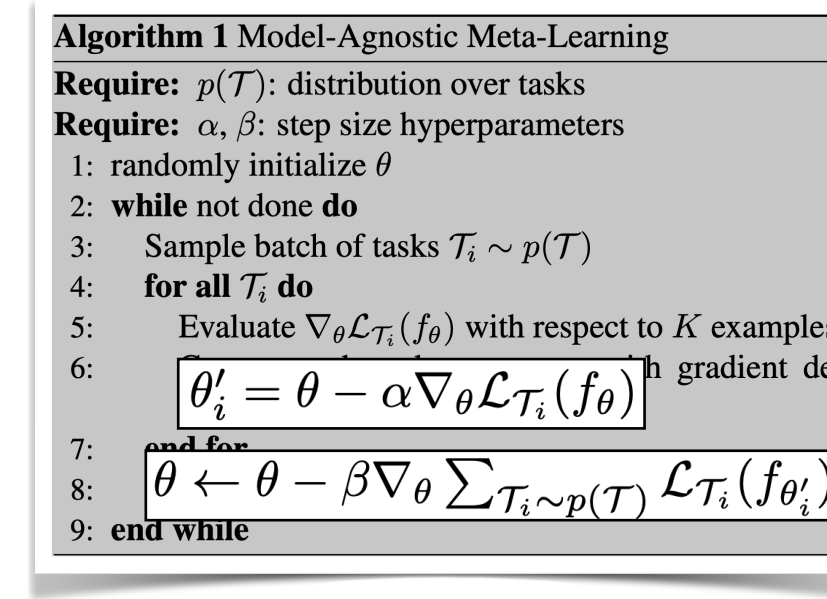
Nominal Semantics of First-Class Automatic Differentiation

Jack Czenszak and Alexander K. Lew

Yale

Motivation

- First-class automatic differentiation (AD) is a desirable language feature used throughout the machine learning literature (see right);
- Many modern machine learning frameworks (JAX, PyTorch, and JuliaDiff [4]) implement first-class AD using the *tagged forward-mode algorithm* to avoid *perturbation confusion* (see below) [6];
- However, despite effort, this algorithm has *not* been formally proven correct, making these heavily-used systems untrustworthy.



Finn et al. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. 2017.



Chandra et al. Designing Perceptual Puzzles by Differentiating Probabilistic Programs. 2022.

First-Class AD with Standard Dual Numbers (Bad)

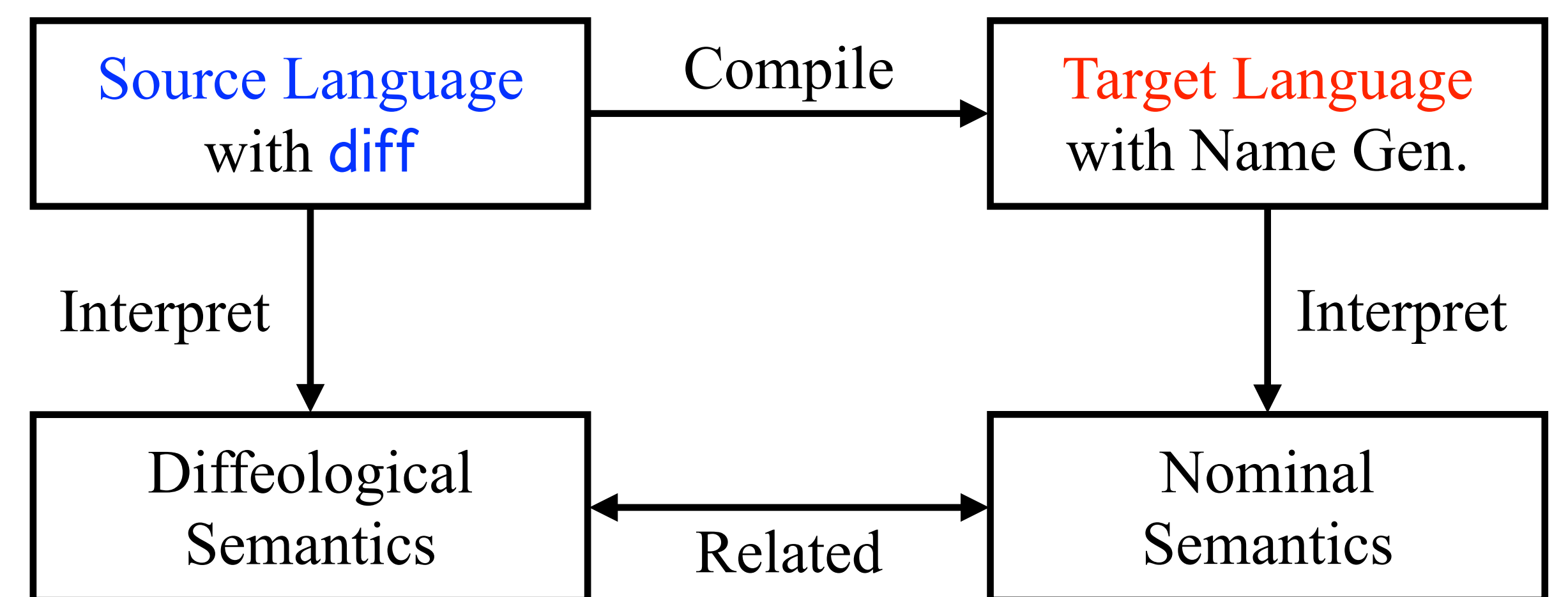
```
diff (λx.x · (diff (λy.x + y) 1)) 1
= coeff ((λx.x · (diff (λy.x + y) 1)) (1 + ε))
= coeff ((1 + ε) · (diff (λy.(1 + ε) + y) 1))
= coeff ((1 + ε) · (coeff ((λy.(1 + ε) + y) (1 + ε))))
= coeff ((1 + ε) · (coeff ((1 + ε) + (1 + ε))))
= coeff ((1 + ε) · (coeff (2 + 2ε))) Combined!
= coeff ((1 + ε) · 2) = coeff (2 + 2ε) = 2 ❌
```

First-Class AD with Tagged Forward-Mode (Correct)

```
diff (λx.x · (diff (λy.x + y) 1)) 1
= coeff1 ((λx.x · (diff (λy.x + y) 1)) (1 + ε1))
= coeff1 ((1 + ε1) · (diff (λy.(1 + ε1) + y) 1))
= coeff1 ((1 + ε1) · (coeff2 ((λy.(1 + ε1) + y) (1 + ε2))))
= coeff1 ((1 + ε1) · (coeff2 ((1 + ε1) + (1 + ε2))))
= coeff1 ((1 + ε1) · (coeff2 (2 + ε1 + ε2))) Distinct
= coeff1 ((1 + ε1) · 1) = coeff1 (1 + ε1) = 1 ✅
```

Formalizing the Tagged Forward-Mode Algorithm

- Compile a high-level *source language* with first-class *diff* to a *target language* (with an *obvious* implementation strategy) that can express the tagged forward-mode algorithm using name generation;
- Source language* is given semantics $\llbracket - \rrbracket_1$ in the category of diffeological spaces **Diff** [1];
- Target language* is given semantics $\llbracket - \rrbracket_2$ in the Schanuel topos **Sch** = **[I, Set]_{pb}**, a well-studied semantic model for name generation [7];
- Compiler is proven correct using categorical logical relations.



Source and Target Languages

Source Language Syntax

Type $\ni A, B ::= \text{real} \mid \text{unit} \mid A \times B \mid A \Rightarrow B$
Expr $\ni M, N ::= x \mid r \in \mathbb{R} \mid \text{op} \mid \text{diff} \mid \langle \rangle \mid \langle M, N \rangle$
 $\mid \text{fst } M \mid \text{snd } M \mid \lambda x : A. M \mid MN$

Target Language Syntax

Type $\ni A, B ::= \text{dual} \mid \text{name} \mid \text{unit} \mid A \times B \mid A \Rightarrow B \mid \mathcal{T}A$
Expr $\ni M, N ::= x \mid r \in \mathbb{R} \mid \text{op} \mid \text{coeff} \mid \text{epsil} \mid \text{new} \mid \langle \rangle$
 $\mid \langle M, N \rangle \mid \text{fst } M \mid \text{snd } M \mid \lambda x : A. M$
 $\mid MN \mid x \leftarrow M; N \mid \text{ret } M$

Source Language Semantics in **Diff**

$\llbracket \text{real} \rrbracket_1 = \mathbb{R} \quad \llbracket A \Rightarrow B \rrbracket_1 = \mathbf{Diff}(\llbracket A \rrbracket_1, \llbracket B \rrbracket_1)$
 $\llbracket \text{unit} \rrbracket_1 = \{ \star \} \quad \llbracket \text{diff} \rrbracket_1(\gamma) = d$

- $d : \mathbf{Diff}(\mathbb{R}, \mathbb{R}) \rightarrow \mathbf{Diff}(\mathbb{R}, \mathbb{R})$ is first-order differentiation.

Target Language Semantics in **Sch**

$\llbracket \text{dual} \rrbracket_2 = \mathbb{D} \quad \llbracket \text{coeff } M \ N \rrbracket_2 = \text{coeff} \circ \langle \llbracket M \rrbracket_2, \llbracket N \rrbracket_2 \rangle$
 $\llbracket \text{name} \rrbracket_2 = \mathbb{A} \quad \llbracket \text{epsil } M \rrbracket_2 = \text{epsil} \circ \llbracket M \rrbracket_2$
 $\llbracket \mathcal{T}A \rrbracket_2 = \mathcal{T}\llbracket A \rrbracket_2 \quad \llbracket \text{new} \rrbracket_2 = \text{new} \circ !_{\llbracket \Gamma \rrbracket_2}$

- $\mathcal{T}F = \text{colim}_{X \in \mathbf{I}} F(- + X)$ is the *name generation monad* [7];
- $\mathbb{A}(X) = X$ is the *names object* [7];
- $\text{new}_X(\star) = [1, \text{inr}_{X,1}(\star)]$ generates a *fresh name* [7];
- Dual numbers* \mathbb{D} , **coeff** : $\mathbb{A} \times \mathbb{D} \rightarrow \mathbb{D}$, and **epsil** : $\mathbb{A} \rightarrow \mathbb{D}$ are:

$$\mathbb{D}(X) = \left\{ \sum_{S \subseteq X} p_S \prod_{x \in S} x \in \mathbb{R}[X] \mid p_S \in \mathbb{R} \right\}$$

$$\text{coeff}_X(x, p) = p_1 \text{ where } p = p_1 x + p_2 \quad \text{epsil}_X(x) = x$$

Acknowledgements. We would like to thank the anonymous LAFI 2026 reviewers for their comments and Pedro H. Azevedo de Amorim, Mathieu Huot, John M. Li, Cameron Moy, and Sam Staton for helpful discussions. **References.** [1] Huot et al. Correctness of automatic differentiation via diffeologies and categorical gluing. 2020. [2] Katsumata. A semantic formulation of \top -lifting and logical predicates for computational metalanguage. 2005. [3] Krawiec et al. Provably correct, asymptotically efficient, higher-order reverse-mode automatic differentiation. 2022. [4] Manzyuk et al. Perturbation confusion in forward automatic differentiation of higher-order functions. 2019. [5] Mitchell and Scedrov. Notes on scoring and relators. 1992. [6] Siskind and Pearlmutter. Perturbation confusion and referential transparency: Correct functional implementation of forward-mode AD. 2005. [7] Stark. Names and Higher-Order Functions. 1994.

Compiler Correctness

Source-to-Target Compiler

- Well-typed *source* programs $\Gamma \vdash M : A$ are compiled to well-typed *target* programs $\mathcal{C}(\Gamma) \vdash \mathcal{C}(M) : \mathcal{T}\mathcal{C}(A)$ according to:

$\mathcal{C}(\text{real}) = \text{dual} \quad \mathcal{C}(\langle M, N \rangle) = x \leftarrow \mathcal{C}(M);$
 $\mathcal{C}(A \times B) = \mathcal{C}(A) \times \mathcal{C}(B) \quad y \leftarrow \mathcal{C}(N);$
 $\mathcal{C}(A \Rightarrow B) = \mathcal{C}(A) \Rightarrow \mathcal{T}\mathcal{C}(B) \quad \text{ret } \langle x, y \rangle$
 $\mathcal{C}(r) = \text{ret } r \quad \mathcal{C}(MN) = f \leftarrow \mathcal{C}(M);$
 $\mathcal{C}(\lambda x. M) = \text{ret } \lambda x. \mathcal{C}(M) \quad x \leftarrow \mathcal{C}(N);$
 $\mathcal{C}(\text{fst } M) = x \leftarrow \mathcal{C}(M); \text{ret } x \quad fx$
 $\mathcal{C}(\text{diff}) = \text{ret } \lambda f : \text{dual} \Rightarrow \mathcal{T}\text{dual}. \text{ret } \lambda x : \text{dual}.$
 $n \leftarrow \text{new}; d \leftarrow f(x + \text{epsil } n); \text{ret } \text{coeff } n \ d$

Correctness Proof by Logical Relations

Theorem (Correctness). Suppose $\cdot \vdash M : \text{real}$ has $\llbracket M \rrbracket_1 = \llbracket r \rrbracket_1$ for some $r \in \mathbb{R}$, then $\llbracket \mathcal{C}(M) \rrbracket_2 = \llbracket \text{ret } r \rrbracket_2$.

Lemma (Fundamental Property). Consider $\Gamma \vdash M : A$. For any pair $(\gamma_1, \gamma_2) \in \mathcal{V}(\Gamma)$, we get $(\llbracket M \rrbracket_1 \circ \gamma_1, \llbracket \mathcal{C}(M) \rrbracket_2 \circ \gamma_2) \in \dot{\mathcal{V}}\mathcal{V}(A)$.

- Proven by categorical *logical relations* [5] indexed by Kripke worlds (similar to [3]) of finite sets of tags:

$$\mathcal{V}(\text{real})(X) = \left\{ (f, g) \in K(\mathbb{R}, \mathbb{D}) \mid g(i, y) = \sum_{S \subseteq X} \frac{\partial |S|}{\partial S} f(y_{i(x)})_{x \in X} \prod_{x \in S} i(x) \right\}$$

$$\mathcal{E}(\text{real})(X) = \left\{ (f, g) \in K(\mathbb{R}, \mathcal{T}\mathbb{D}) \mid g = \eta_{\mathbb{D}} \circ h, (f, h) \in \mathcal{V}(\text{real})(X) \right\}$$

with $K(X, F) = \mathbf{Diff}(\mathbb{R}^{A(-)}, X) \times (\mathbb{P} \Rightarrow F)$ and $\mathbb{P}(X) = \mathbb{R}^{A(X)}$ s.t. $\mathbb{P}(f)$ pads with zeroes.

- Construct a *gluing category* **Gl** and *fibration for logical relations* [2] $p : \mathbf{Gl} \rightarrow \mathbf{Diff} \times \mathbf{Sch}$ by pulling back the evident sub-object fibration along $K : \mathbf{Diff} \times \mathbf{Sch} \rightarrow [\mathbf{I}, \mathbf{Set}]$, and lift $\text{id}_{\mathbf{Diff}} \times \mathcal{T}$ using \top \top -lifting [3] with param. $(\langle \llbracket - \rrbracket_1, \llbracket \mathcal{C}(-) \rrbracket_2 \rangle, \mathcal{E})$ to get $\dot{\mathcal{T}}$ on **Gl**;
- Proof is then by induction, with *diff* being the main difficulty.