

1 We thank the reviewers for their constructive and inspiring feedback. We will improve the paper by incorporating the
 2 following responses. As we cannot see R2 (i.e., Reviewer #2), we respond to the reviews by R1, R3, and R4 only.

3 **[R1/R4]** The correctness of autodiff systems defined in the paper could be misleading to practitioners. Any relationship
 4 between the correctness of autodiff systems and that of applications (e.g., gradient descent) built upon autodiff systems?

5 We agree with the reviewers’ points that (i) the correctness of the applications built upon autodiff systems is as important
 6 as that of the underlying autodiff systems, but (ii) the former does not necessarily follow from the latter (especially
 7 as defined in the paper). These two notions of correctness address separate issues, and our work is mainly about the
 8 second notion (i.e., the correctness of autodiff systems). Also, we do not claim that our correctness condition is “the”
 9 gold standard. Rather we are just suggesting “a” correctness condition that can serve as a reasonable (possibly minimal)
 10 requirement for existing and future autodiff systems. We will clarify this limitation in the revised version of the paper.

11 Here are detailed responses to the point (ii) on the applications mentioned in the reviews.

- 12 • Gradient descent: As illustrated by R1’s example of $f(x)$, our correctness condition for autodiff systems does not
 13 necessarily imply the correctness of the gradient descent based on those systems (i.e., that the gradient descent
 14 converges to Clarke critical points). This gives a partial answer to R3’s question on possible drawbacks of using
 15 intensional derivatives. On the other hand, we conjecture that if gradient descent is based on intensional derivatives
 16 and starts at randomly chosen initial points, it would be correct almost surely. This is an open problem.
- 17 • Hamiltonian Monte Carlo (HMC) and variational inference (VI), possibly for probabilistic programming: We
 18 reiterate that PAP functions enjoy a nice property that they are analytic on each piece of domain, whose boundary is
 19 measure-zero. The property has been crucially used to design various methods of HMC and VI for non-differentiable
 20 densities and prove their correctness (e.g., [1, 2]). This signifies the importance of studying PAP functions. Whether
 21 the correctness claims in those works would still hold if intensional derivatives are used in place of standard ones, is
 22 another open problem. We will discuss these interesting open problems in the revised version of the paper.

23 Although there remain a few open problems, we strongly believe that our work would serve as an important first step
 24 towards understanding and resolving those problems. Above all, as far as we know, this is the first work that (i) raises
 25 subtleties in the well-known chain rule when applied to almost-everywhere differentiable functions, (ii) gives concrete
 26 counterexamples illuminating those subtleties, and (iii) proves some reasonable (possibly minimal) correctness of
 27 existing autodiff systems that permits non-differentiable functions, using only elementary mathematics.

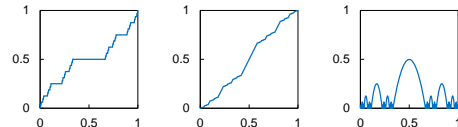
28 [1] Discontinuous Hamiltonian Monte Carlo for discrete parameters and discontinuous likelihoods. *Biometrika*, 2020.
 29 [2] Reparameterization Gradient for Non-Differentiable Models. In *NeurIPS*, 2018.

30 **[R1]** Is there any better correctness condition for autodiff systems that permits some non-differentiable functions but, at
 31 the same time, ensures nice behaviors of autodiff systems when applied to differentiable functions?

32 This is a good question that would lead to interesting future work. We do not have concrete results, but a possible
 33 approach would be to divide PAP representations and intensional derivatives into “good” and “bad” ones, and consider
 34 a new correctness condition that involves only those “good” ones. A desired property of “good” intensional derivatives
 35 might be that they must be identical to the standard derivative if a given function is differentiable everywhere. Under the
 36 property, R1’s $f(x)$ is considered a “bad” PAP representation. A promising idea to construct “good” PAP representations
 37 that induce “good” intensional derivatives is to put additional requirements to their domain partitions, such as that each
 38 piece of a domain partition should be a half-open interval if the entire domain is \mathbb{R} ; R1’s $f(x)$ violates this requirement.

39 **[R3]** Plots of functions used in the proof of Proposition 1?

40 Shown right are draft plots of the 1-Cantor function, and f and g in
 41 the proof of Proposition 1. We will add them to the paper.



42 **[R4]** Implementation of `relu`, `reciprocal_no_nan`, and `sqrt`? Any connections to observed values of their “derivatives”?

43 We checked that TensorFlow and PyTorch compute the “derivatives” of the three functions f , not by applying autodiff
 44 (or symbolic differentiation) to the implementation of f , but by evaluating a separately written implementation of $\tilde{D}f$.
 45 The implementations of f and $\tilde{D}f$, denoted by \mathbf{f} and $\tilde{D}\mathbf{f}$, are as follows.

	<code>relu(x)</code>	$\tilde{D}\text{relu}(x)$	<code>recip(x)</code>	$\tilde{D}\text{recip}(x)$	<code>sqrt(x)</code>	$\tilde{D}\text{sqrt}(x)$
46 TensorFlow	<code>_max(x, 0)</code>	<code>if (x > 0) 1 0</code>	<code>_div(1, x)</code>	<code>_div(-1, x × x)</code>	<code>_sqrt(x)</code>	<code>0.5/_sqrt(x)</code>
PyTorch	<code>if (x ≤ 0) 0 x</code>	<code>if (x ≤ 0) 0 1</code>	N/A	N/A	<code>_sqrt(x)</code>	<code>1/(2 × _sqrt(x))</code>

47 Here `_max` and `_sqrt` are functions provided by a standard math library, and `_div(x1, x2)` is implemented as `if (x2 =`
 48 `0) 0 (x1/x2)`. If we interpret `relu` and $\tilde{D}\text{relu}$ as PAP representations, one can see that $(\tilde{D}\text{relu}) = D(\text{relu})$ may
 49 fail depending on the implementation of `_max` (e.g., consider `_max(x1, x2) = if (x1 ≥ x2) x1 x2)`. This suggests that
 50 “derivatives” of primitive functions f in autodiff systems could have nothing to do with the implementation of f .