

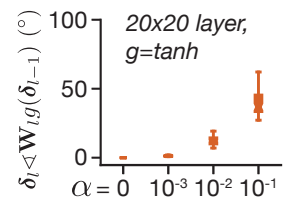
1 We thank the reviewers for their constructive and fair criticism of our submission. Overall we noted four main areas of
2 critique, addressed below with reference to each reviewer’s specific comments (reviewers referred to as R1-4).

3 **1. Biological plausibility (R1):** *R1 questioned the biological plausibility of Fig 2, the mechanisms of switching or*
4 *gating of feedback in the brain, and the concurrent use of feedback weights for other purposes.*

5 We insist that our claims of plausibility are reasonable. First, solving the weight transport problem alone should make
6 any algorithm more plausible than backprop. Second, we hypothesize that the remaining issues with Fig 2 (signed
7 errors, separate passes, connectivity gating) could be mitigated in combination with insights from other studies (e.g.,
8 propagation of targets, signal multiplexing with dendritic compartments), but are not central to dynamic inversion as a
9 novel solution to weight transport. Furthermore, our hypotheses of feedback gating are not far from neuroscientific
10 theories about inhibition and neuromodulation — e.g., acetylcholine may control the relative strengths of feedforward
11 and feedback connectivity (Hasselmo 2006, “The role of acetylcholine in learning and memory.”). We are willing
12 to soften our claims of plausibility in the text, though we note that the term “biologically-plausible” is often loosely
13 defined, and it is common in the literature to tackle one problem at a time. Lastly, in our statements of using feedback
14 for other purposes, we simply wished to point out that this may lead to conflicts in weight values, and the flexibility
15 of dynamic inversion could help (we can soften this claim too). We did not mean to claim that dynamic inversion is
16 always better than learning feedback (e.g., Kolen & Pollack 1994), which we can elaborate on in the discussion.

17 **2. Stability and leak (R2,R3):** *R2 argued that stability is crucial and under-stressed; R3 suggested that the leak*
18 *parameter α may have a large effect on the steady state.*

19 We find the questions raised about stability and leak to be of less concern than the reviewers
20 fear. Stability was rarely (if ever) a problem in our experiments — only for linear regression
21 was it necessary to perform stability optimization during training (and only during the
22 first ~ 100 iterations of 2000); for nonlinear regression, MNIST classification, and MNIST
23 autoencoding, the dynamics remained stable once the initialization was set (also note that
24 we keep track of and optimize stability for non-dynamic inversion too). Furthermore, it
25 is both common in models and well within the hypothesized capabilities of biology to
26 specify precise initializations for learning algorithms (e.g., see Zador 2019, “A critique of pure learning...”). Thus, we
27 contend that R2’s suggestion of testing a plausible stability-enforcing mechanism is not critical here, and testing stability
28 optimization without dynamics would be equivalent to feedback alignment with a particular weight initialization.
29 Overall, we argue that stability enforcement (through dynamics and/or plasticity) seems reasonable given current
30 neuroscientific theories (e.g., Zenke, Ganguli & Gerstner 2017, “The temporal paradox of Hebbian learning and
31 homeostatic plasticity.”). As for the leak and scaling of $\alpha \mathbf{u}_i$ — we can include simulations showing that the leak
32 produces gradual increases in error (example on the right shows angle between δ_i and output reconstruction $\tilde{\delta}_i$).



33 **3. MNIST results and scalability (R2,R3,R4):** *R2-4 were concerned with performance on MNIST classification; R2*
34 *and R4 pointed out that only non-dynamic inversion (NDI) was implemented on the MNIST tasks; R2 and R3 wanted*
35 *more explanation of single-loop inversion (SDI); R4 noted concerns of scalability in general.*

36 We stress that our aim in this submission was to introduce a proof-of-concept idea, and so we find the deficiencies in
37 performance to be concerning, but not critical. As the reviewers suggest, fine-tuning the architecture, hyperparameters,
38 and optimization procedure may change these results (as well as further study into the conditioning of the inversion).
39 Furthermore, the results on MNIST autoencoding show potential benefits over feedback alignment. We justify the fact
40 that only non-dynamic inversion (NDI) was used in the MNIST examples with the observation that NDI behaves nearly
41 identically to DI (Fig 3d,i; when $\alpha > 0$), though we would agree to verify this. We also agree with the reviewers that
42 not enough detail was given for single-loop inversion (SDI), and we would be happy to elaborate on it.

43 **4. Theoretical justification and link to Gauss-Newton (GN) optimization (R3,R4):** *R3 argued that the paper lacks*
44 *theoretical justification, and questioned the link to GN optimization; R4 pointed out some omitted literature.*

45 We again reiterate that the main aim of the submission was to introduce a novel idea. While it is true that a theoretical
46 understanding of our algorithm is lacking (and we thank R3 for the interesting suggestions), it is common for such
47 rigorous theoretical work to follow publication of the original idea. For example, target propagation has recently
48 received more rigorous analysis which may be applicable to dynamic inversion of targets, which we can cite (R4
49 mentioned Meulemans et al. 2020 “A theoretical framework for target propagation”, and see Bengio 2020, “Deriving
50 differential target propagation from iterating approximate inverses.”). In any case, we do agree to soften our claims
51 that dynamic inversion approximates GN optimization, and instead point it out as an interesting link for future study.
52 Lastly, R3’s mention that the pre-activation variable \mathbf{a}_i doesn’t appear in the update can be readily explained. In
53 backpropagation, this variable is only needed to estimate the slope of the nonlinearity, $g'(\mathbf{a}_i)$. In dynamic inversion, the
54 nonlinearity is included implicitly in the dynamics (which is arguably more biologically-plausible), though this may be
55 problematic for approximate inversions and thus merits further study.