1 We thank all reviewers for their comments. We will correct all typos and address all minor comments in the final paper.

2 **Reviewer #1:** Biological plausibility of C1-C3 and $\phi$? **C1:** Given an input-target pair, C1 means that a model
3 makes a prediction from the input, before it learns from the target, which is quite natural for biological neural systems.
4 **C2 and $\phi$:** C2 is realized by $\phi$, and it is very plausible that $\phi$ can be performed by biological neurons, because some
5 types of neurons are well-known to respond predominantly to changes in their input [71]. **C3:** Though it is unnatural
6 for biological neural systems to have a specific integration step, we show in the supplementary material that relaxing C3
7 results in BP with a different learning rate for different layers, which does not violate the core of BP.
8 Weight transport problem? In the predictive coding model, the errors are back-propagated by correct weights, because
9 the model includes feedback connections that also learn. The weight modification rules for corresponding feedforward
10 and feedback weights are the same, which ensures that they remain equal if initialized to equal values (see [48]).
11 Other criticisms for BP (spiking and recurrent neural networks) remain? As discussed in the related work section, they
12 are left as future work, while this work addresses two of the most crucial ones: local plasticity and autonomy.
13 Related work of https://arxiv.org/abs/2006.04182? This work was published on 7 Jun 2020 (after the NeurIPS'20
14 deadline). It includes another approximation to BP but no equivalence; we will add it to the related work.

15 **Reviewer #2:** Weight transport problem? Please see line 8 in the response to Reviewer #1.
16 Locality of learning rule? (the update of $\theta_{i,j}^{l+1}$ depends on $x_i^l$)? According to Eq. 9, the update of $\theta_{i,j}^{l+1}$ actually depends
17 on $x_j^{l+1}$ rather than $x_i^l$ (we suspect the reviewer might have misread Eq. 9), so it is local.
18 Confusion about inference? Alg. 2 explicitly states C1 in the second "Require". In line 158, we make a note that C1 is
19 omitted in Fig. 2 for simplicity. This note will be included in the caption of Fig. 2.
20 Full autonomy during the initial convergence of prediction? During the prediction phase, the error nodes change due
21 to feedforward input, while during learning, the error nodes change due to feedback input. We will clarify that, in order
22 to prevent learning during prediction, $\phi$ is equal to 1 only if the change in error node is caused by feedback input.

23 **Reviewer #3:** Limitations of PCNs. They are well-discussed in [10,48]; we will add a summary of them. Note that
24 some limitations have been addressed (e.g., 1-to-1 connections are addressed in [10]). We will also review key studies
25 illustrating that PCNs are widely used and informative models of information processing in the brain.
26 Model in [48] already autonomous, because plasticity triggered by network convergence? The plasticity trigger in [48]
27 requires global information (i.e., total error in all error nodes), while our trigger $\phi$ needs only local information, thus, is
28 more plausible for biological implementation.
29 Experiment of Fa-Z-IL? $\phi$ with $t_d > 4$ succeeds in all detections (Table 3), i.e., Fa-Z-IL with such $\phi$ coincides with
30 Z-IL (BP). We will include the classification results of such Fa-Z-IL, which produces the same results as Z-IL (BP).
31 Clarification of full autonomy? Fa-Z-IL does require the input to be presented before the teacher to satisfy C1. We
32 consider this to be a requirement of the learning setup and will make it explicit in the paper, moderating our claim of
33 autonomy. However, we consider such requirement to be much weaker, compared to switching computational rules
34 (BP) and detecting convergence of global variables (IL). We leave the study of removing this requirement or putting it
35 inside an autonomous neural system as future research.
36 Moderate the title? We will modify the title to the more specific "Can the Brain Do Backpropagation? — Exact
37 Implementation of Backpropagation in Predictive Coding Networks".
38 From Z-IL to Fa-Z-IL? We will add the statement that Fa-Z-IL loses formal equivalence to BP, but with $t_d > 4$,
39 empirical equivalence always remains.
40 C1 in [48] should be acknowledged? We will acknowledge this.
41 Clarification of [48]? In line 42, we have stated that some previous works are equivalent to BP when feedback is
42 sufficiently weak (i.e., teacher weakly perturbs the network); we will add such clarification when introducing [48].
43 Discuss about [46]? We will add a section outlining the differences of learning rules between [46] and Z-IL, along
44 which we point out that some steps may serve similar general purposes. A deeper study on the connections of the two
45 substantially different learning rules is left as future research.

46 **Reviewer #4:** Moderate the title? Will be changed to "Can the Brain Do Backpropagation? — Exact Implementation
47 of Backpropagation in Predictive Coding Networks".
48 Strong claim? We will moderate the sentence pointed out by the reviewer.
49 Classification accuracy? The classification accuracy of all situations in Fig. 3 will be added to the final paper: the
50 averaged accuracies of BP, IL, Z-IL, and Fa-Z-IL are 94.16%, 93.78%, 94.16%, and 94.16%, respectively.
51 Details of the two criteria? The divergence of the test error is the L1 distance between the corresponding test errors,
52 averaged over 64 training iterations (the test error is evaluated after each training iteration). The divergence of the final
53 weights is the sum of the L2 distance between the corresponding weights, after the last training iteration.
54 Extra description of Fig. 2? Lines 133–141 are the description of Fig. 2; this will be included in the caption of Fig. 2.