

1 We thank all reviewers for their thoughtful comments and suggestions. We address each review separately.

2 **Reviewer #1.** Regarding the encoder architecture impact, in Table 1, we vary the classification model making it
3 increasingly more powerful, and demonstrate that our method produces improvements in all cases. For the agreement
4 model encoder, we found that the results are not as sensitive to the encoder choice (e.g., for CIFAR10 switching from an
5 MLP to a CNN did result in significant differences). We will include results for different agreement model architectures.

6 **Reviewer #5.** We thank R5 for the relevant references. We were not aware of them—especially the contemporaneous
7 ones from ICLR and ICML 2019. SNTG infers a similarity graph between samples, but it does so in a significantly
8 different way than GAM. Also, in contrast to SNTG, we propose an additional self-training component, and our method
9 is applicable when a graph is provided, whereas SNTG (as published) is not designed to use information from a provided
10 graph. We will include a thorough comparison to SNTG and Fast-SWA in the paper. We will also discuss the following:

- 11 – Parameters: Increasing the number of parameters of the baselines to match that of the respective GAMs results in
12 worse performance for the baselines (e.g., in Table 1, MLP₂₅₆ has the same number of parameters as MLP₁₂₈+GAM, as
13 we use an MLP₁₂₈ for the agreement model, but it performs worse than MLP₁₂₈+GAM and MLP₁₂₈). This is expected as
14 GAMs provide a robust form of regularization for training high capacity models that tend to overfit otherwise.
- 15 – Convergence: Prior work [e.g., Blum and Mitchell 1998, Balcan et al. 2005] proves that co-training converges if: (i)
16 the majority of the learners perform better than random guessing after the first iteration, and (ii) the mistakes they
17 make are weakly dependent. Our experiments indicate that (i) is true in our case. (ii) is harder to verify due to the
18 coupling between the models. However, our empirical evaluation shows that co-training converges successfully. Note
19 that in Fig. 5, even the worse iterations are well above chance, so it should not diverge under these assumptions.
- 20 – Experiments: The missing numbers for GCN₁₀₂₄+VAT are 83.4, 68.9, 79.5 on Cora, Citeseer, and Pubmed, while
21 GCN₁₀₂₄+VATENT obtains 32.5, 8.5, 18.0, which follow the same trend as our other results. For VATENT, we observed
22 that on the graph datasets the entropy term becomes large and dominates the loss. Decreasing its weight makes the
23 performance to converge to that of VAT. Our implementation works on CIFAR10 and SVHN, thus it seems unlikely
24 to be the reason behind the poor results. Interestingly, [2] reports only GCN+VAT results and not GCN+VATENT.

25 Regarding comparisons with other methods, we will add the results reported in [2] to Table 1. Their best numbers are
26 lower than our GCN+GAM. [4] tackles the same problem, but their evaluation is on random train/test splits rather than
27 the commonly used Planetoid splits. We observe that the GCN paper reports much better results on random splits than
28 [4], and we have demonstrated that GAM can be applied on top of GCN to improve it further. For completeness, we
29 will report results on random splits and compare with [4]. To compare with SNTG and Fast-SWA, we plan to run the
30 experiments with a 13-layer CNN suggested by R5 for the camera-ready. Note, however, that we do not necessarily see
31 these approaches as competitors to GAM, but rather as additional regularizers that, similar to VAT, can be applied in
32 conjunction with GAM to further improve generalization. To illustrate that GAM works for large networks too, here are
33 results (obtained after the submission deadline) using the WideResnet of Oliver et al. 2018 on CIFAR10-4000: baseline
34 79.69%, +II-Model 83.63%, +Mean teacher 84.13%, +VATENT 86.87%, +GAM* 87.42%.

35 **Reviewer #6.** R6 suggests a discussion on the challenges in simply replacing classification models in label propagation
36 with deep learning models. We address this through an example from our paper, and then explain how this example is
37 more broadly applicable. Replacing classification models in label propagation with deep learning models is exactly
38 what Neural Graph Machines (NGMs) do (described in Section 2): an NGM is a label propagation model complemented
39 by a deep learning classifier operating on the node features. Setting the regularization coefficients to 0 makes it a pure
40 deep learning model, while increasing their values brings it closer to label propagation. When the graph is noisy, the
41 regularization coefficients need to be small (otherwise the regularization forces connected nodes from different classes
42 to incorrectly have the same label), thereby reducing the effect of the graph on the model. However, with such minimal
43 regularization the model tends to overfit to the few available labeled examples. Our approach combines deep learning
44 with label propagation in a manner that allows us to handle noisy graphs in a robust fashion. Note that other methods
45 besides NGM also suffer from this problem (e.g., GCN, Planetoid)—see *Robustness* section. Our experiments show
46 how GAMs are able to learn in a much more robust manner.

47 Novelty: The novelty of our algorithm is the interaction between the agreement and classification models, which allows
48 it to benefit from both label propagation and deep learning even when dealing with noisy graphs (where most label
49 propagation algorithms fail), or no graphs at all. It is surprising and interesting that even though the two models learn
50 using the same features and same data, their interplay can produce such large increases in accuracy on a wide variety of
51 base networks (MLP, CNN, Resnet, GCN, and GAT), suggesting they learn complementary information.

52 Co-Training: We argue that our proposed training algorithm does indeed fit in the co-training framework. While
53 the original paper [Blum and Mitchell, 1998] proposed co-training in the setting described by R6, the same authors
54 subsequently proposed co-training settings where some classifiers predict label distributions and others predict coupling
55 constraints over these distributions (like in our setting). Perhaps the most notable and influential example of this is the
56 Never-Ending Language Learning (NELL) system [Mitchell 2015, 2018].