Semi-Supervised Learning for Molecular Graphs via Ensemble Consensus

Rasmus H. Tirsgaard¹

Laurits Fredsgaard¹

Marisa Wodrich¹

Mikkel Jordahn¹

Mikkel N. Schmidt¹

¹Technical University of Denmark

Section for Cognitive Systems, Department of Applied Mathematics and Computer Science {rhti, laula, mawod, mikkjo, mnsc}@dtu.dk

Abstract

Machine learning is transforming molecular sciences by accelerating property prediction, simulation, and the discovery of new molecules and materials. Acquiring labeled data in these domains is often costly and time-consuming, whereas large collections of unlabeled molecular data are readily available. Standard semi-supervised learning methods often rely on label-preserving augmentations, which are challenging to design in the molecular domain, where minor changes can drastically alter properties. In this work, we show that semi-supervised methods that rely on an ensemble consensus can boost predictive accuracy across a diverse range of molecular datasets, task types, and graph neural network architectures. Notably, we show that training with an ensemble consensus objective results in an effect similar to knowledge distillation; an individual member of an ensemble trained this way outperforms a full ensemble trained in a traditional supervised fashion in almost all cases. In addition, this type of semi-supervised training reduces calibration error and is robust over different datasets.

1 Introduction

In recent years, machine learning has emerged as a transformative tool in the molecular sciences, accelerating discovery in areas ranging from predicting quantum mechanical properties [Schütt et al., 2021, 2017, Musaelian et al., 2023, Wood et al., 2025] to discovering novel drugs [Wong et al., 2024, Kellenberger et al., 2007, Vidler et al., 2013, Zhuang et al., 2014, Ren et al., 2023] and catalysts [Pillai et al., 2023, Sun et al., 2024, Bai et al., 2025]. However, despite recent efforts to curate large labeled datasets [Merchant et al., 2023, Levine et al., 2025], the scarcity of labeled data remains a fundamental bottleneck.

In materials and drug discovery, labels often come from computationally expensive simulations, such as density functional theory (DFT), or resource-intensive laboratory measurements. Consequently, datasets with specialized high-quality labels are typically small, while large databases of unlabeled molecules (e.g., ZINC, PubChem [Irwin et al., 2012, Kim et al., 2024]) are not fully exploited. This scenario—abundant unlabeled data coupled with scarce labeled data—is an ideal setting for semi-supervised learning (SSL).

Yet, many state-of-the-art methods are poorly suited for the molecular domain. Dominant techniques such as consistency training [Berthelot et al., 2019, Sohn et al., 2020] critically depend on data augmentation strategies that create perturbed copies of an input while preserving its label. Such

39th Conference on Neural Information Processing Systems (NeurIPS 2025) Workshop: AI4Mat Workshop (Neurips 2025).

augmentations are notoriously difficult to design for molecules, where minor structural changes can drastically alter the chemical properties we aim to predict. Meanwhile, approaches such as iterative pseudo-labeling [Scudder, 1965, Riloff and Wiebe, 2003, Huang et al., 2022] hinges on the ability to reliably rank predictions by confidence in order to select the best candidates for pseudo-labeling and to avoid reinforcing model errors. This highlights a critical gap where standard SSL benchmarks and algorithms do not translate well to the practical challenges of molecular science.

In this work, we build upon a class of SSL methods that does not require the explicit design of data augmentations, but rather relies on an *ensemble consistency loss*. Specifically, we train a model ensemble where each member learns from labeled data using a standard supervised loss and from unlabeled data using a loss that promotes agreement among the ensemble members. While ensemble coupling in self-supervised learning has been explored in previous work [Sajjadi et al., 2016, Tarvainen and Valpola, 2017, Platanios, 2018], our formulation is theoretically grounded in an ensemble loss ambiguity decomposition, trains in a single run, and exhibits a knowledge distillation-like effect that has not previously been discussed. Surprisingly, we find that a single model from the coupled ensemble often achieves greater accuracy than an entire decoupled ensemble. We demonstrate the effectiveness of this approach on a wide array of molecular graph datasets.

2 Methods

We address a standard semi-supervised learning problem with a small set of labeled data, $\mathcal{D}_L = \{(x_i,y_i)\}_{i=1}^{N_L}$, and a large set of unlabeled data, $\mathcal{D}_U = \{u_j\}_{j=1}^{N_U}$. We assume that both datasets are drawn from the same underlying distribution. Our method utilizes a deep ensemble of M models, $\bar{f} = \{f_{\theta_m}\}_{m=1}^{M}$, initialized with different random weights.

The training objective is defined on each model f_{θ_m} within the ensemble. At each training step, its parameters θ_m are updated to minimize a composite loss, \mathcal{L}_m , which combines a standard supervised signal \mathcal{L}_{sup} with an ensemble-driven consistency signal $\mathcal{L}_{\text{consistency}}$:

$$\mathcal{L}_m = \mathcal{L}_{\text{sup}}(f_{\theta_m}, B_L) + \gamma \mathcal{L}_{\text{consistency}}(f_{\theta_m}, \bar{f}, B_U), \tag{1}$$

where γ is the coupling weight, and B_L and B_U are mini-batches of labeled and unlabeled data. During training, all models are updated simultaneously by minimizing the sum of their individual losses i.e. $\mathcal{L} = \sum_{m=1}^{M} \mathcal{L}_m$. The first term, \mathcal{L}_{sup} , is the standard task-specific loss for model f_{θ_m} on the labeled batch, such as mean squared error (MSE) for regression or cross-entropy (CE) for classification. The second term, $\mathcal{L}_{\text{consistency}}$, provides the semi-supervised signal. It is calculated for the model f_{θ_m} but depends on the outputs of the entire ensemble. For each unlabeled sample $u \in B_U$, a consensus prediction, $\bar{f}(u)$, is computed by averaging the predictions of all M models:

$$\bar{f}(u) = \frac{1}{M} \sum_{m=1}^{M} f_{\theta_m}(u).$$
 (2)

The Consensus prediction serves as the augmentation-free consistency target for model f_{θ_m} . We penalize the discrepancy between model prediction and the ensemble consensus as

$$\mathcal{L}_{\text{consistency}}(f_{\theta_m}, \bar{f}, B_U) = \frac{1}{|B_U|} \sum_{u \in B_U} D\left(f_{\theta_m}(u), \bar{f}(u)\right). \tag{3}$$

Here, D is a suitable distance metric, for example, the task-specific supervised loss (e.g., L2 or KL-divergence). In practice, when minimizing the loss we detach the gradient through $\bar{f}(u)$, as the consensus prediction is at least as accurate as the individual members' predictions on average (see Appendix B), ensuring that the ensemble is not encouraged to match the less accurate individual predictions. Note, detaching the gradient has been observed to result in failure cases such as *learner collusion* [Jeffares et al., 2023], but in our experiments the results are not affected negatively.

2.1 Consensus-Diversity Dynamics

Our proposed SSL training scheme directly manipulates the trade-off between accurate individual models and high diversity among them. The unsupervised loss term, $\mathcal{L}_u(x_u) = \mathcal{L}(f_{\theta_i}(x_u), f_e(x_u))$, creates a pull towards consensus by guiding each model f_{θ_i} to agree with the more stable ensemble

prediction f_e . This directly reduces the average individual error by providing a high-quality supervisory signal for unlabeled data. Simultaneously, this pull is counteracted by forces that preserve diversity. Each model begins from a unique random initialization and follows a distinct optimization path due to the stochastic nature of mini-batch SGD. This dynamic allows the models to converge to different solutions in parameter space while still agreeing in function space. Therefore, our method does not eliminate diversity but rather regulates it.

We speculate that the continuous learning between models means they should be less likely get stuck on early bad predictions, as could happen with many forms of pseudo-labeling. This is because the ensemble targets are "moving" with the ensemble. This suggests the ensemble prediction does not need a warm-startup strategy, as other works have observed by Tarvainen and Valpola [2017] and used in Filipiak et al. [2021], Platanios [2018].

3 Experimental setup

We evaluate our method in two settings: First, on a quantum chemistry benchmark to demonstrate its relevance for 3D-geometry-based molecular property prediction, and then across a diverse suite of graph-level tasks to assess its broader applicability. All ensemble members were trained on identical mini-batches of supervised data to simplify implementation. While this strategy reduces ensemble diversity, potentially limiting the ensemble's predictive power, it allows for a fair direct comparison with single models.

Semi-supervised protocol To simulate the common scenario of data scarcity, we restrict the supervised portion of our training to a small fraction for each task (10%). The remaining training data (90%) is treated as unlabeled and is used exclusively for our ensemble consistency loss. Our primary baseline is a standard deep ensemble of the same architecture, trained only on this small labeled data subset. This setup allows us to directly measure the performance gain from leveraging unlabeled data.

Datasets We evaluate our method across several domains. First, we predict 12 molecular properties on the QM9 dataset [Wu et al., 2018] using a four-member (M=4) PaiNN [Schütt et al., 2021] ensemble, where we also study performance scaling by varying the ensemble size of a single target. For broader validation, we use GCN [Kipf and Welling, 2016], GIN [Xu et al., 2019] and GatedGCN [Bresson and Laurent, 2017] architectures on a suite of graph-level tasks, adapting the code from Rampásek et al. [2022] and Luo et al. [2024], again with M=4. To demonstrate general applicability, further experiments on non-molecular and non-graph datasets are included in the Appendix. All datasets are split 80/10/10 for train/validation/test.

Hyperparameter tuning To ensure well-tuned models for datasets, the training hyperparameters (learning rate and weight decay) were optimized for each target and model based on the validation performance of a single model in the supervised setting on the reduced labeled data. These hyperparameters were kept fixed across different SSL methods tested to ensure fair comparison. The parameters associated with each specific SSL method (coupling weight, mean-teacher decay, etc.) were optimized based on validation accuracy for each target on QM9, and selected for the GNN+datasets based on the best value of ZINC. Details about the tuning procedures and selected hyperparameters can be found in Appendix D.

Evaluation We evaluate the predictive performance for a standard ensemble, an ensemble using SSL via ensemble consensus (ours), the individual members from the ensembles, mean-teacher [Tarvainen and Valpola, 2017], and PSEUDO σ [Huang et al., 2022]. All results are reported as the mean along with 1.96 times the standard error of the mean across different seeds.

4 Results and Discussion

4.1 Molecular Property Prediction on QM9

The performance of our method on the 12 regression targets of the QM9 dataset is presented in Table 1. The results indicate that training with the ensemble consistency loss ("Supervised + SSL")

Table 1: PaiNN performance (MAE) on QM9 targets. Results are reported as mean ± 1.96 standard error of the mean over 5 seeds.

		Individ	Individual Member		nble (M=4)		
Target	Unit	Supervised	Supervised + SSL	Supervised	Supervised + SSL	Mean-teacher	$PSEUD\sigma$
μ	D	$.07390 \pm .00077$	$.06191 \pm .00024$.06808±.00056	$.06136 \pm .00024$.07211±.00106	$.06487 \pm .00088$
α	a_0^3	$.1622 \pm .0011$	$.1322 \pm .0011$	$.1419 \pm .0009$	$.1303 \pm .0011$	$.1570 \pm .0009$	$.1454 \pm .0004$
ϵ_{HOMO}	meV	$80.61 \pm .5062$	$73.98 \pm .4368$	$76.47 \pm .5361$	$73.08 \pm .4472$	$80.61 \pm .9079$	78.72 ± 1.1196
ϵ_{LUMO}	meV	$62.04 \pm .4253$	$57.72 \pm .2247$	$59.31 \pm .4609$	$57.24 \pm .2159$	$61.97 \pm .3648$	59.74 ± 0.6248
$\Delta\epsilon$	meV	$125.2 \pm .4734$	$117.0 \pm .4988$	$119.4 \pm .4468$	$115.7 \pm .5100$	125.0 ± 1.155	$122.5 \pm .7501$
$\langle R^2 \rangle$	a_0^2	$.7922 \pm .0284$	$.6100 \pm .0206$	$.6246 \pm .0205$	$.5605 \pm .0206$	$.7987 \pm .0197$	$.9099 \pm .0157$
ZPVE	meV	$2.220 \pm .0055$	$2.014 \pm .0054$	$2.074 \pm .0054$	$1.991 \pm .0055$	$2.182 \pm .0170$	$2.141 \pm .0123$
U_0	meV	$24.88 \pm .1477$	$19.96 \pm .1291$	$20.91 \pm .1557$	$19.38 \scriptstyle{\pm .1278}$	$24.73 \pm .3803$	$24.00 \pm .2929$
U	meV	$25.19 {\scriptstyle\pm .2025}$	$20.17 \pm .1577$	$21.10 {\pm}.1914$	$19.59 {\scriptstyle \pm .1574}$	$24.96 {\scriptstyle\pm} .4452$	$24.28 \pm .4701$
H	meV	$25.12 \pm .1981$	$20.14 \pm .1268$	$21.09 \pm .1946$	$19.55 {\scriptstyle \pm .1328}$	$24.85 \pm .4295$	$24.33 {\scriptstyle \pm .6181}$
G	meV	$25.38 {\pm}.1856$	$20.31 \pm .1571$	$21.41 {\scriptstyle \pm .1811}$	$19.75 \pm .1634$	$25.16 {\pm} .3359$	$24.33 {\scriptstyle \pm .4412}$
C_v	cal mol K	$.05668 {\pm}.00038$	$.04498 {\pm}.00019$	$.04884 {\pm}.00035$	$.04392 {\scriptstyle \pm .00020}$	$.05570 \pm .00021$	$.05408 {\pm} .00027$

reduces the MAE across all evaluated targets when compared to the supervised-only baseline. This is observed for both the individual PaiNN models and the four-member ensembles. Furthermore, the individual model from the coupled ensemble consistently outperforms the traditional supervised ensemble on all targets. This is also the case across different ensemble sizes, as explored in Table 3.

4.2 Generalization across graph benchmarks

Our experiments on QM9 and the more varied GNN+ benchmark (see Table 2) show that our ensemble-based SSL framework consistently improves model performance in low-data regimes. The most significant finding is the substantial boost in accuracy for individual models, a direct result of the knowledge transferred from the ensemble's consensus on unlabeled data. This finding is alike that of ensemble distilling [Hinton et al., 2015], where the knowledge of an ensemble is transferred to a single, smaller model, except that our method inherently produces knowledgeable single models. This is explained through the semi-supervised effect on the entire ensemble, resulting in even better ensemble consensus targets for individual models to learn from. This has a key practical benefit: while the method requires an ensemble during training, a single, improved model can be deployed for inference. This offers a valuable trade-off, where an increased one-time training cost yields a final model that is both highly accurate and computationally efficient at inference time. Chemical property screening or MD simulations are compelling use case, where models are called many times and can introduce computational bottlenecks if very expensive. It is noteworthy that for datasets where the parameter (γ or the mean-teacher decay) related to SSL was not directly tuned, the improvement in predictive accuracy was noticeably smaller. This indicates the SSL parameter is highly dependent on the specific dataset.

In the table 3 and C.2 in the appendix, the predictive performance scales with the number of members in the coupled ensemble. Individual models from the ensemble trained with our method consistently perform at a similar level to an entire traditional ensemble across all ensemble sizes.

Limitations The primary limitation of our approach is the computational overhead associated with training an ensemble consensus model. Transfer learning is another method that is often used in sparsely-labeled settings, which we have not compared against.

Future work Our findings suggest several promising avenues for future research. While this work created a semi-supervised split from a fully labeled dataset, a compelling next step would be to test out of domain generalization, by using all available labeled data for supervision while introducing a separate, fully unlabeled dataset. This would more directly quantify the benefit of leveraging vast, external chemical libraries and be of interest in a practical setting. Using ensembles for semi-supervised learning also opens the direction for improving accuracy in a principled manner by diversifying the ensemble members through existing techniques. Furthermore, different strategies for how to couple our ensemble can be investigated. Only including the unsupervised data in the later part during training could potentially result in similar predictive performance, while reducing computational cost.

Acknowledgments and Disclosure of Funding

The authors acknowledge support from the Novo Nordisk Foundation under grant no NNF22OC0076658 (Bayesian neural networks for molecular discovery).

We acknowledge the Danish e-infrastructure Consortium (DeiC) for awarding this project access to the LUMI supercomputer, owned by the EuroHPC Joint Undertaking, hosted by CSC (Finland) and the LUMI consortium.

Finally, we also thank the reviewers for their time and effort in reviewing this work.

References

- E. Arazo, D. Ortego, P. Albert, N. E. O'Connor, and K. McGuinness. Pseudo-labeling and confirmation bias in deep semi-supervised learning. In 2020 International Joint Conference on Neural Networks, IJCNN 2020, Glasgow, United Kingdom, July 19-24, 2020, pages 1–8. IEEE, 2020. doi: 10.1109/IJCNN48605.2020.9207304. URL https://doi.org/10.1109/IJCNN48605.2020.9207304.
- R. Bai, Y. Yao, Q. Lin, L. Wu, Z. Li, H. Wang, M. Ma, D. Mu, L. Hu, H. Yang, W. Li, S. Zhu, X. Wu, X. Rui, and Y. Yu. Preferable single-atom catalysts enabled by natural language processing for high energy density Na-S batteries. *Nature Communications*, 16(1):5827, 7 2025. ISSN 2041-1723. doi: 10.1038/s41467-025-60931-x.
- D. Berthelot, N. Carlini, I. J. Goodfellow, N. Papernot, A. Oliver, and C. Raffel. Mixmatch: A holistic approach to semi-supervised learning. In H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. B. Fox, and R. Garnett, editors, Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pages 5050-5060, 2019. URL https://proceedings.neurips.cc/paper/2019/hash/1cd138d0499a68f4bb72bee04bbec2d7-Abstract.html.
- X. Bresson and T. Laurent. Residual gated graph convnets. *CoRR*, abs/1711.07553, 2017. URL http://arxiv.org/abs/1711.07553.
- C. Bucilu, R. Caruana, and A. Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 535–541, 2006.
- X. Chen, Y. Yuan, G. Zeng, and J. Wang. Semi-supervised semantic segmentation with cross pseudo supervision. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages 2613–2622. Computer Vision Foundation / IEEE, 2021. doi: 10.1109/CVPR46437.2021.00264. URL https://openaccess.thecvf.com/content/CVPR2021/html/Chen_Semi-Supervised_Semantic_Segmentation_With_Cross_Pseudo_Supervision_CVPR_2021_paper.html.
- D. Filipiak, P. Tempczyk, and M. Cygan. n-cps: Generalising cross pseudo supervision to n networks for semi-supervised semantic segmentation. *CoRR*, abs/2112.07528, 2021. URL https://arxiv.org/abs/2112.07528.
- T. Fukuda, M. Suzuki, G. Kurata, S. Thomas, J. Cui, and B. Ramabhadran. Efficient knowledge distillation from an ensemble of teachers. In *Interspeech*, pages 3697–3701, 2017.
- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016, pages 770–778. IEEE Computer Society, 2016. doi: 10.1109/CVPR.2016.90. URL https://doi.org/10.1109/CVPR.2016.90.
- B. Heo, M. Lee, S. Yun, and J. Y. Choi. Knowledge transfer via distillation of activation boundaries formed by hidden neurons. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01): 3779–3787, Jul. 2019. doi: 10.1609/aaai.v33i01.33013779. URL https://ojs.aaai.org/index.php/AAAI/article/view/4264.

- G. E. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *CoRR*, abs/1503.02531, 2015. URL http://arxiv.org/abs/1503.02531.
- K. Huang, V. Sresht, B. Rai, and M. Bordyuh. Uncertainty-aware pseudo-labeling for quantum calculations. In J. Cussens and K. Zhang, editors, *Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*, volume 180 of *Proceedings of Machine Learning Research*, pages 853–862. PMLR, 01–05 Aug 2022. URL https://proceedings.mlr.press/v180/huang22a.html.
- J. J. Irwin, T. Sterling, M. M. Mysinger, E. S. Bolstad, and R. G. Coleman. ZINC: A Free Tool to Discover Chemistry for Biology. *Journal of Chemical Information and Modeling*, 52(7):1757– 1768, 7 2012. ISSN 1549-9596. doi: 10.1021/ci3001277.
- A. Jeffares, T. Liu, J. Crabbé, and M. van der Schaar. Joint training of deep ensembles fails due to learner collusion. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/2bde8fef08f7ebe42b584266cbcfc909-Abstract-Conference.html.
- E. Kellenberger, J.-Y. Springael, M. Parmentier, M. Hachet-Haas, J.-L. Galzi, and D. Rognan. Identification of Nonpeptide CCR5 Receptor Agonists by Structure-based Virtual Screening. *Journal of Medicinal Chemistry*, 50(6):1294–1303, 3 2007. ISSN 0022-2623. doi: 10.1021/jm061389p.
- S. Kim, J. Chen, T. Cheng, A. Gindulyte, J. He, S. He, Q. Li, B. A. Shoemaker, P. A. Thiessen, B. Yu, L. Zaslavsky, J. Zhang, and E. E. Bolton. Pubchem 2025 update. *Nucleic Acids Research*, 53(D1):D1516-D1525, 11 2024. ISSN 1362-4962. doi: 10.1093/nar/gkae1059. URL https://doi.org/10.1093/nar/gkae1059.
- T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. *CoRR*, abs/1609.02907, 2016. URL http://arxiv.org/abs/1609.02907.
- A. Krogh and J. Vedelsby. Neural Network Ensembles, Cross Validation, and Active Learning. In G. Tesauro, D. Touretzky, and T. Leen, editors, *Advances in Neural Information Processing Systems*, volume 7. MIT Press, 1994. URL https://proceedings.neurips.cc/paper_files/paper/1994/file/b8c37e33defde51cf91e1e03e51657da-Paper.pdf.
- S. Laine and T. Aila. Temporal ensembling for semi-supervised learning. *CoRR*, abs/1610.02242, 2016. URL http://arxiv.org/abs/1610.02242.
- D. S. Levine, M. Shuaibi, E. W. C. Spotte-Smith, M. G. Taylor, M. R. Hasyim, K. Michel, I. Batatia, G. Csányi, M. Dzamba, P. Eastman, N. C. Frey, X. Fu, V. Gharakhanyan, A. S. Krishnapriyan, J. A. Rackers, S. Raja, A. Rizvi, A. S. Rosen, Z. Ulissi, S. Vargas, C. L. Zitnick, S. M. Blau, and B. M. Wood. The open molecules 2025 (omol25) dataset, evaluations, and models, 2025. URL https://arxiv.org/abs/2505.08762.
- Y. Luo, L. Shi, and X. Wu. Classic gnns are strong baselines: Reassessing gnns for node classification. In A. Globersons, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. M. Tomczak, and C. Zhang, editors, Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 15, 2024, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/b10ed15ff1aa864f1be3a75f1ffc021b-Abstract-Datasets_and_Benchmarks_Track.html.
- A. Malinin, B. Mlodozeniec, and M. J. F. Gales. Ensemble distribution distillation. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net, 2020. URL https://openreview.net/forum?id=BygSP6Vtvr.
- A. Merchant, S. Batzner, S. S. Schoenholz, M. Aykol, G. Cheon, and E. D. Cubuk. Scaling deep learning for materials discovery. *Nature*, 2023. doi: 10.1038/s41586-023-06735-9.
- S. Mishra, B. Murugesan, I. B. Ayed, M. Pedersoli, and J. Dolz. Do not trust what you trust: Miscalibration in semi-supervised learning. *CoRR*, abs/2403.15567, 2024. doi: 10.48550/ARXIV. 2403.15567. URL https://doi.org/10.48550/arXiv.2403.15567.

- A. Musaelian, S. Batzner, A. Johansson, L. Sun, C. J. Owen, M. Kornbluth, and B. Kozinsky. Learning local equivariant representations for large-scale atomistic dynamics. *Nature Communications*, 14(1):579, 2 2023. ISSN 2041-1723. doi: 10.1038/s41467-023-36329-y.
- H. S. Pillai, Y. Li, S.-H. Wang, N. Omidvar, Q. Mu, L. E. K. Achenie, F. Abild-Pedersen, J. Yang, G. Wu, and H. Xin. Interpretable design of Ir-free trimetallic electrocatalysts for ammonia oxidation with graph neural networks. *Nature Communications*, 14(1):792, 2023. ISSN 2041-1723. doi: 10.1038/s41467-023-36322-5. URL https://doi.org/10.1038/s41467-023-36322-5.
- E. A. Platanios. Agreement-based learning. CoRR, abs/1806.01258, 2018. URL http://arxiv.org/abs/1806.01258.
- L. Rampásek, M. Galkin, V. P. Dwivedi, A. T. Luu, G. Wolf, and D. Beaini. Recipe for a general, powerful, scalable graph transformer. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 December 9, 2022, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/5d4834a159f1547b267a05a4e2b7cf5e-Abstract-Conference.html.
- F. Ren, X. Ding, M. Zheng, M. Korzinkin, X. Cai, W. Zhu, A. Mantsyzov, A. Aliper, V. Aladinskiy, Z. Cao, S. Kong, X. Long, B. H. Man Liu, Y. Liu, V. Naumov, A. Shneyderman, I. V. Ozerov, J. Wang, F. W. Pun, D. A. Polykovskiy, C. Sun, M. Levitt, A. Aspuru-Guzik, and A. Zhavoronkov. AlphaFold accelerates artificial intelligence powered drug discovery: efficient discovery of a novel CDK20 small molecule inhibitor. *Chemical Science*, 14(6):1443–1452, 2023. ISSN 2041-6520. doi: 10.1039/D2SC05709C.
- E. Riloff and J. Wiebe. Learning extraction patterns for subjective expressions. In *Proceedings* of the 2003 Conference on Empirical Methods in Natural Language Processing, pages 105–112, 2003. URL https://aclanthology.org/W03-1014/.
- M. Sajjadi, M. Javanmardi, and T. Tasdizen. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. In D. D. Lee, M. Sugiyama, U. von Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 1163–1171, 2016. URL https://proceedings.neurips.cc/paper/2016/hash/30ef30b64204a3088a26bc2e6ecf7602-Abstract.html.
- K. Schütt, P. Kindermans, H. E. S. Felix, S. Chmiela, A. Tkatchenko, and K. Müller. Schnet: A continuous-filter convolutional neural network for modeling quantum interactions. In I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, editors, Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 991–1001, 2017. URL https://proceedings.neurips.cc/paper/2017/hash/303ed4c69846ab36c2904d3ba8573050-Abstract.html.
- K. Schütt, O. T. Unke, and M. Gastegger. Equivariant message passing for the prediction of tensorial properties and molecular spectra. In M. Meila and T. Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 9377–9388. PMLR, 2021. URL http://proceedings.mlr.press/v139/schutt21a.html.
- H. Scudder. Probability of error of some adaptive pattern-recognition machines. *IEEE Transactions on Information Theory*, 11(3):363–371, 1965. doi: 10.1109/TIT.1965.1053799.
- K. Sohn, D. Berthelot, N. Carlini, Z. Zhang, H. Zhang, C. Raffel, E. D. Cubuk, A. Kurakin, and C. Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/06964dce9addb1c5cb5d6e3d9838f733-Abstract.html.

- J. Sun, R. Tu, Y. Xu, H. Yang, T. Yu, D. Zhai, X. Ci, and W. Deng. Machine learning aided design of single-atom alloy catalysts for methane cracking. *Nature Communications*, 15(1):6036, 7 2024. ISSN 2041-1723. doi: 10.1038/s41467-024-50417-7.
- A. Tarvainen and H. Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, editors, Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 1195–1204, 2017. URL https://proceedings.neurips.cc/paper/2017/hash/68053af2923e00204c3ca7c6a3150cf7-Abstract.html.
- L. R. Vidler, P. Filippakopoulos, O. Fedorov, S. Picaud, S. Martin, M. Tomsett, H. Woodward, N. Brown, S. Knapp, and S. Hoelder. Discovery of Novel Small-Molecule Inhibitors of BRD4 Using Structure-Based Virtual Screening. *Journal of Medicinal Chemistry*, 56(20):8073–8088, 10 2013. ISSN 0022-2623. doi: 10.1021/jm4011302.
- F. Wong, E. J. Zheng, J. A. Valeri, N. M. Donghia, M. N. Anahtar, S. Omori, A. Li, A. Cubillos-Ruiz, A. Krishnan, W. Jin, A. L. Manson, J. Friedrichs, R. Helbig, B. Hajian, D. K. Fiejtek, F. F. Wagner, H. H. Soutter, A. M. Earl, J. M. Stokes, L. D. Renner, and J. J. Collins. Discovery of a structural class of antibiotics with explainable deep learning. *Nature*, 626(7997):177–185, 2 2024. ISSN 0028-0836. doi: 10.1038/s41586-023-06887-8.
- B. M. Wood, M. Dzamba, X. Fu, M. Gao, M. Shuaibi, L. Barroso-Luque, K. Abdelmaqsoud, V. Gharakhanyan, J. R. Kitchin, D. S. Levine, K. Michel, A. Sriram, T. Cohen, A. Das, A. Rizvi, S. J. Sahoo, Z. W. Ulissi, and C. L. Zitnick. UMA: A Family of Universal Models for Atoms. 6 2025.
- D. Wood, T. Mu, A. M. Webb, H. W. J. Reeve, M. Luján, and G. Brown. A unified theory of diversity in ensemble learning. *J. Mach. Learn. Res.*, 24:359:1–359:49, 2023. URL http://jmlr.org/papers/v24/23-0041.html.
- Z. Wu, B. Ramsundar, E. N. Feinberg, J. Gomes, C. Geniesse, A. S. Pappu, K. Leswing, and V. Pande. MoleculeNet: A Benchmark for Molecular Machine Learning. 10 2018.
- Q. Xie, Z. Dai, E. H. Hovy, T. Luong, and Q. Le. Unsupervised data augmentation for consistency training. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/44feb0096faa8326192570788b38c1d1-Abstract.html.
- K. Xu, W. Hu, J. Leskovec, and S. Jegelka. How powerful are graph neural networks? In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019. URL https://openreview.net/forum?id=ryGs6iA5Km.
- D. Yarowsky. Unsupervised word sense disambiguation rivaling supervised methods. In *33rd Annual Meeting of the Association for Computational Linguistics*, pages 189–196, Cambridge, Massachusetts, USA, June 1995. Association for Computational Linguistics. doi: 10.3115/981658. 981684. URL https://aclanthology.org/P95-1026/.
- S. Zagoruyko and N. Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. *CoRR*, abs/1612.03928, 2016. URL http://arxiv.org/abs/1612.03928.
- C. Zhuang, S. Narayanapillai, W. Zhang, Y. Y. Sham, and C. Xing. Rapid Identification of Keap1Nrf2 Small-Molecule Inhibitors through Structure-Based Virtual Screening and Hit-Based Substructure Search. *Journal of Medicinal Chemistry*, 57(3):1121–1126, 2 2014. ISSN 0022-2623. doi: 10.1021/jm4017174.

A Background

A.1 Background on Semi-Supervised Learning

Semi-Supervised Learning Semi-supervised learning (SSL) is a machine learning paradigm designed for settings with a small amount of labeled data and a much larger amount of unlabeled data. The idea is to leverage the unlabeled data to learn about the underlying structure of the data distribution p(x), which in turn improves the model's ability to learn the mapping from inputs to outputs, p(y|x). Effective SSL methods are typically built upon one or more of the following assumptions:

- Smoothness Assumption: If two points x_1, x_2 are close in a high-density region of the underlying data manifold, their corresponding labels y_1, y_2 should also be close or identical.
- Cluster Assumption: The data tends to form distinct clusters, and points within the same cluster are likely to share the same label. This implies that a good decision boundary should lie in the low-density region between clusters.

Consistency Loss Consistency regularization is currently the most dominant family of SSL methods. The core idea is that the model's prediction for an unlabeled data point should remain consistent under small perturbations. This directly enforces the smoothness assumption. A successful perturbation or data augmentation is one that explores the local neighborhood of a data point on the manifold without changing its label. The objective is typically formulated as minimizing a distance measure (e.g., Mean Squared Error or KL-Divergence) between the model's predictions for two different augmentations of the same input:

$$\mathcal{L}_{\text{consistency}} = \mathbb{E}_{x_u \sim X_u} [D(f_{\theta}(\text{aug}_1(x_u)) || f_{\theta}(\text{aug}_2(x_u))].$$

Different choices of the perturbations give rise to a wide range of methods. Π -models [Sajjadi et al., 2016] enforce that two predictions should be the same under transformations to the data, the use of dropout and random pooling for perturbations to the model. Each unlabeled datapoint is passed through the network twice and penalized for the difference in the predictions between the passes. The benefit of consistency loss is highly linked to the quality of the data augmentation techniques, as shown in [Xie et al., 2020]. Temporal ensembling [Laine and Aila, 2016] builds upon this by maintaining an exponential moving average of predictions for each unlabeled example to create a more stable consistency target. Instead of applying a temporal averaging over the predictions, the mean-teacher method [Tarvainen and Valpola, 2017] averages the model weights and uses the predictions of that model as the consistency target. In the above works, the predictions can be seen as coming from a sort of pseudo-ensemble. As the members of this pseudo-ensemble are based on the trajectory or perturbation of a single network, the diversity of the predictions is reduced and biased, which reduces the prediction accuracy as we later highlight.

This problem can be mitigated by introducing multiple different initial weightings of the same architecture and training them in parallel to use as consistency targets. Chen et al. [2021] (cross pseudo supervision) proposes to do this for pixel-wise segmentation, where the prediction of each of the two ensemble members is hard labeled and used as the consistency target. Filipiak et al. [2021] further extends this for pixel-wise segmentation by using n ensemble models and taking all combinations of hard labeled predictions as the consistency targets. Another paper that explores different ensemble predictions is Platanios [2018]. Here the ensemble members are restarted multiple times during training, and the consensus target is computed from a trainable majority vote or Restricted Boltzmann Machine. All the above methods can be seen as stemming from a broad class of SSL methods that rely on the prediction of an ensemble to guide the training of the individual models to improve predictive accuracy.

In many applications, there exist few or no data augmentations that preserve the label of a data point. Examples include molecules, where the chemical properties can be changed significantly under small changes to the molecule. This restricts the consistency loss methods to only rely on perturbations to the model and not the data. This makes the class of ensemble-based SSL methods well-suited for the problem.

Pseudo-labeling Pseudo-labeling [Yarowsky, 1995, Scudder, 1965, Riloff and Wiebe, 2003], also known as self-training or entropy minimization, is a process where an initial model is trained on

the labeled data points and then used to predict labels for a large unlabeled dataset. The primary risk of this method is confirmation bias: if the model generates an incorrect pseudo-label with high confidence, it will reinforce its own mistake during retraining, leading to error propagation. To mitigate this risk, modern SSL methods often integrate more sophisticated frameworks. For example, one uncertainty-aware approach uses a model's evidential uncertainty to estimate the quality of each pseudo-label. This enables an adaptive weighting scheme where high-uncertainty (low-quality) pseudo-labels are given a smaller weight in the loss function, reducing their biasing effect. While this can be effective, such a strategy requires an initial, full training phase on the labeled data before the episodic pseudo-labeling can begin. It also introduces several additional tunable hyperparameters related to its episodic schedule, which require careful tuning [Huang et al., 2022].

Knowledge Distillation Knowledge distillation [Bucilu et al., 2006, Hinton et al., 2015] was proposed as a way of using a complex "teacher" model to transfer its knowledge to a simpler "student" model. Usually, the teacher model is either a model with more parameters or the same model with multiple predictions averaged over multiple augmentations of the input, but the use of an ensemble as the teacher has also been explored [Hinton et al., 2015, Fukuda et al., 2017, Malinin et al., 2020]. The transfer of knowledge can be enforced at different levels, such as feature representations [Heo et al., 2019] or intermediate layers [Zagoruyko and Komodakis, 2016]. Approaches that match predictions are most closely related to our work. Aligning student and teacher predictions resembles the use of consistency targets in semi-supervised learning, with the key distinction that distillation is typically applied post-hoc, and thus lacks a bootstrapping effect where the teacher also benefits from the students progress. Furthermore, knowledge distillation is often focused on preserving the uncertainty calibration of the teacher or achieving computational efficiency by deploying the smaller student model instead of the larger one.

B Theoretical Motivation

The theoretical motivation for our method is grounded in the formal relationship between an ensemble's performance and that of its individual members. Ensemble performance is governed by a fundamental trade-off between the accuracy of the individual models and the diversity of their predictions. This relationship can be expressed through a loss decomposition, which shows that for any convex loss function, the ensemble's loss is guaranteed to be less than or equal to the average of the individual losses [Wood et al., 2023]. This stems from Jensen's inequality and takes the general form:

Ensemble Loss = Average Individual Loss
$$-$$
 Ambiguity (4)

The ambiguity (or diversity) term is a non-negative quantity measuring disagreement among the members. This decomposition reveals that optimal ensemble performance requires not only accurate individual models but also beneficial diversity.

Mean Squared Error This principle is most clearly illustrated in regression with Mean Squared Error (MSE), where the decomposition is exact and well-established [Krogh and Vedelsby, 1994]. For an ensemble of M models $\{f_{\theta_m}\}_{m=1}^M$ with a mean prediction $\bar{f}(x)$, the decomposition is:

$$\underbrace{(y - \bar{f}(x))^2}_{\text{Ensemble MSE}} = \underbrace{\frac{1}{M} \sum_{m=1}^{M} (y - f_m(x))^2}_{\text{Average Individual MSE}} - \underbrace{\frac{1}{M} \sum_{m=1}^{M} (\bar{f}(x) - f_m(x))^2}_{\text{Ambiguity (Prediction Variance)}}.$$
 (5)

Here, the ambiguity is simply the variance of the predictions around the ensemble mean, providing a clear, label-independent measure of diversity.

Cross-Entropy The same principle extends to classification, though the decomposition for Cross-Entropy (CE) loss is more nuanced. Using the geometric mean to average probabilities across the ensemble yields a clean, label-independent decomposition, as in regression [Wood et al., 2023]. An

exact decomposition is also available for the arithmetic mean:

$$\underbrace{-\mathbf{y} \cdot \ln \bar{\mathbf{f}}}_{\text{Ensemble CE Loss}} = \underbrace{-\frac{1}{M} \sum_{m=1}^{M} \mathbf{y} \cdot \ln \mathbf{f}_{m}}_{\text{Avg. Individual CE Loss}} - \underbrace{\sum_{c=1}^{C} y_{c} \ln \frac{\frac{1}{M} \sum_{m=1}^{M} f_{m,c}}{(\prod_{m=1}^{M} f_{m,c})^{1/M}}}_{\text{Ambiguity (Label-Dependent)}},$$
(6)

although here the ambiguity term is explicitly a function of the true label vector \mathbf{y} (where y_c is the true probability of class c), making it label-dependent [Wood et al., 2023]. Crucially, this ambiguity term is still guaranteed to be non-negative, ensuring that the ensemble loss is always less than or equal to the average individual loss.

Because the ensemble consensus is provably superior to the average individual model, using it as a consistency target for unlabeled data is both effective and theoretically well-justified. In addition, the ensemble prediction will be a useful signal as long as the models are better than random. This suggests the ensemble prediction does not need to incorporate a warm-startup to provide a useful predictive signal, as other works have observed [Tarvainen and Valpola, 2017] and used [Filipiak et al., 2021, Platanios, 2018].

C Extended Studies

C.1 GNN+ benchmark

To assess the broader applicability of our method, we evaluate it on several molecule-related benchmarks using three different GNN architectures. The results are summarized in Table 2, and are consistent with the performance on QM9. Looking at a single model, the addition of the SSL task consistently improves performance over the supervised-only baseline across all datasets and architectures. This performance gain also translates to the full ensembles, which show improvement when trained with the consistency loss. The performance of a single model trained with our SSL method often exceeds that of an entire ensemble trained only on labeled data.

Table 2: Performance on molecule-related benchmarks using different GNN architectures averaged across 5 seeds.

			GC	CN	G	IN	Gated	IGCN
Dataset	Training	Metric	Individual	Ensemble	Individual	Ensemble	Individual	Ensemble
ZINC	Supervised Consensus Pairwise Mean teacher	MAE↓	$\begin{array}{c} .3163 \pm .0121 \\ .2406 \pm .0150 \\ .2462 \pm .0108 \\ .2884 \pm .0128 \end{array}$	$.2934 \pm .0094 \\ .2367 \pm .0148 \\ .2390 \pm .0102 \\ -$	$.2765 \pm .0247 \\ .2519 \pm .0246 \\ .2500 \pm .0083 \\ .2791 \pm .0117$	$.2516 \pm .0136 \\ .2485 \pm .0232 \\ .2462 \pm .0092 \\ -$	$.2920 \pm .0113 \\ .2717 \pm .0230 \\ .2653 \pm .0158 \\ .2830 \pm .0159$	$.2646 \pm .0235 \\ .2658 \pm .0177 \\ .2597 \pm .0171 \\ -$
Peptides-struct	Supervised Consensus Pairwise Mean teacher	MAE ↓	$.3047 \pm .0098 \\ .2868 \pm .0062 \\ .2933 \pm .0031 \\ .2985 \pm .0029$	$.2932 \pm .0084 \\ .2866 \pm .0061 \\ .2892 \pm .0029 \\ -$	$.2966 \pm .0067 \\ .2944 \pm .0072 \\ .2916 \pm .0030 \\ .2948 \pm .0023$	$.2918 \pm .0058 \\ .2938 \pm .0068 \\ .2901 \pm .0029 \\ -$	$.2994 \pm .0105 \\ .2854 \pm .0061 \\ .2898 \pm .0042 \\ .2953 \pm .0034$	$.2908 \pm .0101 \\ .2848 \pm .0068 \\ .2870 \pm .0041 \\ -$
Peptides-func	Supervised Consensus Pairwise Mean teacher	AP↑	$.4931 \pm .0346 \\ .5070 \pm .0141 \\ .5055 \pm .0151 \\ .4893 \pm .0169$	$.5105 \pm .0342 \\ .5160 \pm .0141 \\ .5163 \pm .0150 \\ -$	$.4566 \pm .0224 \\ .4756 \pm .0180 \\ .4739 \pm .0110 \\ .4611 \pm .0130$	$.4765 \pm .0327 \\ .4815 \pm .0179 \\ .4811 \pm .0117 \\ -$	$.4289 \pm .0051 \\ .4509 \pm .0144 \\ .4463 \pm .0067 \\ .4352 \pm .0058$	$.4444 {\pm}.0200 \\ .4580 {\pm}.0062 \\ .4548 {\pm}.0069 \\ -$
ogbg-molhiv	Supervised Consensus Pairwise Mean teacher	AUROC↑	$.7216 \pm .0193 \\ .7308 \pm .0218 \\ .7247 \pm .0160 \\ .7213 \pm .0161$	$.7357 \pm .0212 \\ .7357 \pm .0212 \\ .7336 \pm .0146 \\ -$	$.7329 \pm .0166 \\ .7339 \pm .0149 \\ .7273 \pm .0128 \\ .6996 \pm .0207$	$.7346 \pm .0165 \\ .7347 \pm .0153 \\ .7294 \pm .0128 \\ -$	$.7312 \pm .0081 \\ .7361 \pm .0069 \\ .7375 \pm .0052 \\ .7295 \pm .0165$	$.7341 \scriptstyle{\pm .0107} \\ .7383 \scriptstyle{\pm .0073} \\ .7403 \scriptstyle{\pm .0050} \\ -$
ogbg-molpcba	Supervised Consensus Pairwise Mean teacher	AP↑	$.1368 \pm .0025 \\ .1476 \pm .0023 \\ .1471 \pm .0027 \\ .1435 \pm .0016$	$.1578 \pm .0030 \\ .1585 \pm .0026 \\ .1597 \pm .0028 \\ -$	$.1421 \pm .0026 \\ .1496 \pm .0033 \\ .1498 \pm .0021 \\ .1479 \pm .0037$	$.1567 \pm .0029 \\ .1567 \pm .0039 \\ .1574 \pm .0024 \\ -$	$.1615 \pm .0034 \\ .1701 \pm .0036 \\ .1674 \pm .0032 \\ .1669 \pm .0028$	$.1779 \pm .0043 \\ .1781 \pm .0034 \\ .1765 \pm .0034 \\ -$

C.2 Scaling with Number of Ensemble Members

Using the same setup as section 4.1, we investigate how the predictive accuracy scale with the number of models in the ensemble. The results are detailed in Table C.2.

Table 3: PaiNN performance (MAE) on QM9 internal energy at 0K in eV (U_0) for different ensemble sizes averaged across 5 seeds, with mean ± 1.96 standard error of the mean.

	Individu	ıal member	Ensemble		
Size (M)	Supervised	Supervised + SSL	Supervised	Supervised + SSL	
1	24.8752 ± 0.1477	_	_	_	
2	_	$20.8847 {\pm} 0.2947$	$22.1755 {\pm} 0.4394$	$20.4256 {\pm} 0.2876$	
3	_	$20.4414 {\pm} 0.1625$	$21.3140{\pm} 0.2289$	$19.9342 {\pm} 0.1661$	
4	_	$19.9642 {\pm} 0.1291$	$20.9101 {\pm} 0.1557$	$19.3816 {\pm} 0.1278$	

Scaling past 4 members We also investigate the predictive accuracy scaling with the number of ensemble members to larger than 4 sizes. Ensembles of these sizes were not feasible to do on any of the graph datasets, so we instead use the original computer vision version of CIFAR-10. This also validates that our method works for other domains than graphs. We use ResNet-18 [He et al., 2016] with 5,000 labeled and 40,000 unlabeled data-points without any data augmentations. We performed an exhaustive hyperparameter sweep using a single seed over learning rate (0.1, 0.075, 0.05, 0.025, 0.01, 0.0075, 0.005, 0.0025, 0.001), and weight decay (0.01, 0.025, 0.05, 0.075, 0.1, 0.25, 0.5) for the purely supervised model. The number of epochs and learning rate annealing was fixed at a number informally found to work. The parameters of best performing model on validation accuracy at the last epoch was selected. The optimal values can be found in 14. The coupling weight was fixed kept at $\gamma=1$.

The hyper-parameters can be found in Appendix D.3. From the accuracy results in Table 4 and calibration scores in section C.3, we see a significant increase in accuracy and calibration scores going from a single model to a coupled ensemble with just two models. Interestingly, the individual prediction accuracy of a model trained in a coupled ensemble of two models outperforms the ensemble prediction from all decoupled ensemble sizes tested. This highlights the semi-supervised effect from using unlabeled data for training. Looking at the calibration metrics in Appendix C.3, we see that the calibration results for the coupled ensemble are worse than the uncoupled one. This is often seen in self-supervised learning, as the "self-validating" training can result in worse calibration from confirmation bias [Arazo et al., 2020, Mishra et al., 2024]. Surprisingly, we see the individual calibration improving over the decoupled model (i.e., a single model), and also improving as the number of ensemble members increases.

Table 4: Predictive accuracy (%) on CIFAR-10 validation, comparing Decoupled and Coupled models. The values represent mean \pm 1.96 standard error of the mean.

	Individual Accuracy %		Ensemble Accuracy (%)	
Ensemble size	Decoupled	Coupled	Decoupled	Coupled
1	$59.08{\scriptstyle\pm1.35}$			
2	÷	$66.36{\scriptstyle\pm0.45}$	$62.51{\scriptstyle\pm0.40}$	$66.96{\scriptstyle\pm0.47}$
4	÷	$67.24{\scriptstyle\pm0.40}$	$64.65{\scriptstyle\pm0.46}$	$67.92{\scriptstyle\pm0.49}$
8	:	$67.64{\scriptstyle\pm0.35}$	$65.73{\scriptstyle\pm0.51}$	$68.34{\scriptstyle\pm0.34}$
16	÷	$67.75{\scriptstyle\pm0.32}$	$66.41{\scriptstyle\pm0.57}$	$68.54{\scriptstyle\pm0.45}$
32	:	$67.75{\scriptstyle\pm0.30}$	$66.64{\scriptstyle\pm0.37}$	$68.52{\scriptstyle\pm0.35}$

C.3 Calibration Metrics on CIFAR-10

Table 5: NLL on CIFAR-10, comparing decoupled and coupled models. The values represent mean \pm 1.96 standard error of the mean.

	NLL				
	Individua	l member	Ense	mble	
Ensemble size	Decoupled	Coupled	Decoupled	Coupled	
1	1.543 ± 0.109	• • •	• • •		
2	÷	$1.217{\scriptstyle\pm0.021}$	$1.267 \scriptstyle{\pm 0.020}$	$1.161{\scriptstyle\pm0.021}$	
4	÷	$1.169{\scriptstyle\pm0.019}$	$1.121{\scriptstyle\pm0.012}$	$1.096 {\pm} 0.019$	
8	÷	$1.142{\scriptstyle\pm0.017}$	$1.048 \scriptstyle{\pm 0.015}$	$1.064 \scriptstyle{\pm 0.016}$	
16	÷	$1.126{\scriptstyle\pm0.015}$	$1.007 \scriptstyle{\pm 0.015}$	$1.047{\scriptstyle\pm0.014}$	
32	:	1.123 ± 0.019	$0.990{\scriptstyle\pm0.011}$	$1.042{\scriptstyle\pm0.018}$	

Table 6: AUC-ROC on CIFAR-10, comparing decoupled and coupled models. The values represent mean \pm 1.96 standard error of the mean.

	AUC-ROC				
	Individua	l member	Ense	mble	
Ensemble size	Decoupled	Coupled	Decoupled	Coupled	
1	$.8885 {\pm} .0075$				
2	÷	$.9250 {\scriptstyle \pm .0020}$	$.9125 {\scriptstyle \pm .0021}$	$.9292 {\pm} .0020$	
4	:	$.9295 {\pm} .0019$	$.9266 {\pm} .0016$	$.9349 {\scriptstyle \pm .0019}$	
8	:	$.9316 {\pm}.0019$	$.9336 {\pm} .0021$	$.9377 {\pm}.0018$	
16	:	$.9323 {\scriptstyle \pm .0016}$	$.9384 {\pm} .0019$	$.9386 {\pm}.0015$	
32	:	$.9329 \scriptstyle{\pm .0019}$	$.9409 {\scriptstyle \pm .0017}$	$.9394 {\pm}.0019$	

Table 7: ECE on CIFAR-10, comparing decoupled and coupled models. The values represent mean \pm 1.96 standard error of the mean.

	ECE				
	Individua	l member	Ensemble		
Ensemble size	Decoupled	Decoupled Coupled		Coupled	
1	$.2210 \scriptstyle{\pm .0357}$	• • •	• • •	• • •	
2	:	$.1713 \scriptstyle{\pm .0041}$	$.1128 \pm .0057$	$.1512 {\pm} .0049$	
4	÷	$.1609 {\pm} .0034$	$.0591 {\pm} .0043$	$.1369 \scriptstyle{\pm .0040}$	
8	÷	$.1548 {\pm}.0031$	$.0320 {\pm} .0043$	$.1301 {\pm} .0035$	
16	÷	$.1494 \scriptstyle{\pm .0028}$	$.0243 \pm .0033$	$.1235 {\pm}.0030$	
32	÷	$.1485 {\pm} .0044$	$.0207 {\pm} .0043$	$.1226 {\pm} .0043$	

Table 8: Brier score on CIFAR-10, comparing decoupled and coupled models. The values represent mean \pm 1.96 standard error of the mean.

	Brier				
	Individua	l member	Ense	mble	
Ensemble size	Decoupled	Coupled	Decoupled	Coupled	
1	$.4854 {\pm} .0271$		• • •		
2	:	$.5594 \scriptstyle{\pm .0042}$	$.4585 {\pm} .0041$	$.5530 {\pm} .0040$	
4	:	$.5654 {\pm} .0040$	$.4422 {\pm} .0047$	$.5572 {\pm} .0040$	
8	:	$.5676 {\pm} .0048$	$.4316 \scriptstyle{\pm .0050}$	$.5588 {\pm} .0047$	
16	÷	$.5652 {\pm} .0044$	$.4283 {\pm} .0049$	$.5563 {\pm} .0043$	
32	:	$.5649 \scriptstyle{\pm .0030}$	$.4263 {\scriptstyle \pm .0033}$	$.5558 {\pm} .0030$	

C.4 Non-Chemical GNN+ datasets

Results for non-chemical GNN+ datasets are shown in Table 9. Note the consensus and mean-teacher run for the GatedGCN models were not computed, as the models were too large to fit in memory.

Table 9: Performance on non-molecule-related benchmarks, comparing supervised models with those using additional self-supervised learning (SSL). Results are shown for individual models (Individual) and the full ensemble (Ensemble). Results are the mean ± 1.96 standard error of the mean over 5 different seeds.

			GCN		GIN		GatedGCN	
Dataset	Training	Metric	Individual	Ensemble	Individual	Ensemble	Individual	Ensemble
CIFAR-10	Supervised Consensus Mean teacher	Acc (%)†	$\begin{array}{c} 50.44{\pm}0.33 \\ 55.33{\pm}0.31 \\ 50.64{\pm}0.28 \end{array}$	$55.38{\scriptstyle \pm 0.49}\atop 57.11{\scriptstyle \pm 0.42}$			57.69±0.34	61.23±0.45
MNIST	Supervised Consensus Mean teacher	Acc (%)†	$\begin{array}{c} 96.61{\scriptstyle \pm 0.07} \\ 96.82{\scriptstyle \pm 0.08} \\ 96.55{\scriptstyle \pm 0.06} \end{array}$	$96.97 \pm 0.04 \ 96.93 \pm 0.11$			$\begin{array}{c} 96.96 {\pm} 0.05 \\ 97.48 {\pm} 0.06 \\ 96.84 {\pm} 0.13 \end{array}$	

D Hyperparameters

D.1 QM9

Our hyperparameter search for QM9 followed a two-step process. First, we started with baseline hyperparameters from a fully supervised setting and tuned the learning rate and weight decay for a single model on the 10% labeled data subset. Second, using these optimized parameters, we then tuned the coupling weight (γ) for the size-4 ensemble by searching over $\{1.0, 0.1, 0.01, 0.001, 0.0001\}$. The coupling weight swept for the mean-teacher was $\{0.9, 0.95, 0.99, 0.995, 0.999\}$. Final a architectural and training configurations are detailed in Table 10 and Table 11.

D.1.1 PSEUD σ **Baseline Implementation**

To provide a robust comparison, we re-implemented the Uncertainty-Aware Pseudo-labeling (PSEUD σ) method [Huang et al., 2022] for the PaiNN architecture, as the original work did not test this model. We used our exact 10% labeled / 90% unlabeled training data split to ensure a direct and fair comparison. Our implementation used a PaiNN backbone with an added evidential head to output the required prior parameters $(\gamma, v, \alpha, \beta)$. We trained using AdamW (1e-4 LR, 1e-4 WD) with a batch size of 32. The training schedule consisted of a 1000-epoch initial training phase on the labeled data, followed by 15 outer-loop episodes (M) of 100 inner-loop epochs (K) each. We

adopted the original paper's recommended low-data hyperparameters, including an evidential regularization coefficient (λ) of 0.5 and epistemic uncertainty for adaptive weighting. Consistent with the PSEUD σ strategy, our cosine annealing learning rate scheduler was re-initialized at the start of each of the 15 episodes.

Table 10: Hyperparameter Configuration for QM9. These are fixed across all targets.

Hyperparameter	Value
Training	
Batch size	32
Epochs	1000
Optimizer	AdamW
Scheduler	Cosine annealing
Coupling	
Unsupervised loss criterion	L2

Table 11: Additional hyperparameter Configuration for QM9 for different targets.

Target	Learning rate	Weight decay	Coupling weight	Mean teacher decay
μ	1e-3	1e-3	0.1	0.995
α	1e-4	1e-3	0.1	0.99
ϵ_{HOMO}	1e-3	0	0.01	0.95
$\epsilon_{ m LUMO}$	5e-4	1e-6	0.01	0.9
$\Delta\epsilon$	1e-3	0	0.01	0.99
$\langle R^2 \rangle$	5e-4	1e-4	0.1	0.99
ZPVE	5e-4	1e-5	0.001	0.99
U_0	1e-4	1e-4	0.01	0.99
U	1e-4	0	0.01	0.9
H	1e-4	1e-4	0.01	0.9
G	1e-4	1e-5	0.01	0.995
C_v	1e-4	1e-5	0.01	0.995

D.2 GNN+ Datasets

We keep the hyperparameters for the different datasets and models the same as in the original paper, except for the number of epochs, weight decay, and learning rate. As we are training with 10% of the original data, we double the number of epochs to mitigate the fewer parameter updates. We then made a two-step hyper-parameter sweep; initially the learning rate using original weight decay values, and afterwards the weight decay using the found best learning rates. The learning rates investigated were (0.25, 0.5, 1.0, 2.0, 4.0) times the original learning rate value for that model and dataset. The weight decays investigated was $(10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 0)$. We could not simply multiply the weight decay values by a fixed factor, as some of the original weight decay values were 0. These sweeps were performed for a single uncoupled model following the same tuning procedure as in the original paper. Notably, this means that the predictive accuracy report from each run is the best validation performance seen during any of the epochs. The found learning rates are listed in Table 12, and weight decays Table 13 below. The train, validation, and test splits follow the same procedure as Luo et al. [2024]. Each seed shuffles the labeled and unlabeled part of the training data.

The SSL parameters were selected based on the best performing values on the validation score on ZINC. The mean-teacher values investigated was (0.9, 0.99, 0.995, 0.999), and the coupling weight for the consensus and pair-wise methods were (0.25, 0.5, 0.75, 1, 1.25, 1.5, 1.75, 2.0). The optimal value of mean-teacher was found to be 0.999, and coupling weight for the consensus learning was 1.0, and the pairwise loss was tied between 0.5 and 0.75, so we went with 0.5 based on the recommendations in Filipiak et al. [2021].

Table 12: Tuned learning rates for GNN models across datasets.

Dataset	GCN	GINE	GATEDGCN
CIFAR-10	0.002	0.0005	0.001
CLUSTER	0.0005	0.0005	0.002
ogbg-molhiv	0.0001	0.00005	0.0004
MalNet-Tiny	0.00025	0.002	0.002
MNIST	0.001	0.002	0.001
PATTERN	0.004	0.001	0.000125
ogbg-molpcba	0.000125	0.000125	0.00025
peptides-func	0.0005	0.002	0.002
peptides-struct	0.002	0.0005	0.002
ogbg-ppa	0.0006	0.0012	0.0003
PascalVOC-SP	0.004	0.002	0.0005
ZINC	0.004	0.001	0.004

Table 13: Tuned weight decays for GNN models across datasets.

Dataset	GCN	GINE	GATEDGCN
CIFAR-10	10^{-2}	10^{-1}	10^{-2}
CLUSTER	0	10^{-1}	10^{-6}
ogbg-molhiv	10^{-3}	10^{-1}	10^{-5}
MalNet-Tiny	10^{-4}	10^{-2}	10^{-4}
MNIST	10^{-1}	10^{-2}	10^{-5}
PATTERN	10^{-3}	10^{-2}	10^{-1}
ogbg-molpcba	10^{-1}	10^{-2}	10^{-5}
peptides-func	0	10^{-1}	10^{-3}
peptides-struct	10^{-3}	10^{-5}	10^{-1}
ogbg-ppa	10^{-1}	10^{-1}	10^{-2}
PascalVOC-SP	10^{-1}	10^{-4}	10^{-2}
ZINC	10^{-1}	10^{-5}	10^{-3}

D.3 CIFAR-10

The hyperparameter configurations for CIFAR-10 are shown in Table 14.

Table 14: Hyperparameter Configuration for CIFAR-10.

Hyperparameter	Value				
Learning Rate					
Learning rate Annealing method Step size Learning rate reduction	0.005 Step 1 0.975				
Regularization					
L2 Weight Decay	0.075				
Optimizer					
Optimizer Momentum	SGD 0.9				
Training					
Epochs	250				
Loss Function					
Coupled loss weighting Ensemble coupled loss Supervised loss	1.0 KL-divergence Cross-entropy				

E Calibration Scores for the ogbg-molhiv

We also investigate the calibration on the ogbg-molhiv benchmark. We do not investigate the datasets ogbg-pcba and peptides functional due to the to the large skewing of classes and missing values. The results are included in Table 15 and Table 16. We see across different architectures that the coupling of the ensemble improves the calibration scores, especially NLL. One notable exception is the MCE score for the GIN ensemble model, where the coupled ensemble becomes significantly worse.

Table 15: Individual Performance on the ogbg-molhiv dataset

	GCN		GIN		GatedGCN	
Metric	Decoupled	Coupled	Decoupled	Coupled	Decoupled	Coupled
Accuracy	$95.78{\scriptstyle\pm0.38}$	96.18 ± 0.48	$95.97{\scriptstyle\pm0.68}$	96.30 ± 0.35	$95.66{\scriptstyle\pm0.72}$	96.01 ± 0.57
ROC-AUC	$.721 \scriptstyle{\pm .0193}$	$.731 \scriptstyle{\pm .0218}$	$.733 \scriptstyle{\pm .017}$	$.734 \scriptstyle{\pm .015}$	$.731 \pm .008$	$.736 \pm .007$
NLL	$.375 {\scriptstyle \pm .185}$	$.230 \scriptstyle{\pm .0662}$	$.147 \scriptstyle{\pm .015}$	$.140 \pm .012$	$.200 \pm .033$	$.180 {\pm} .023$
ECE	$.0312 \scriptstyle{\pm .0092}$	$.0246 \scriptstyle{\pm .0039}$	$.0113 \scriptstyle{\pm .0045}$	$.0105 {\pm} .0048$	$.0232 {\scriptstyle\pm .0069}$	$.0201 {\scriptstyle\pm .0049}$
MCE	$.2041 {\scriptstyle\pm .0994}$	$.2058 {\pm}.0763$	$.1113 \scriptstyle{\pm .0620}$	$.1058 {\scriptstyle \pm .0246}$	$.1154 \scriptstyle{\pm .0399}$	$.0985 {\scriptstyle\pm .0287}$

Table 16: Ensemble Performance on the ogbg-molhiv dataset

	GCN		GIN		GatedGCN	
Metric	Decoupled	Coupled	Decoupled	Coupled	Decoupled	Coupled
Accuracy	96.66 ± 0.33	96.60 ± 0.20	96.11 ± 0.66	96.39 ± 0.33	96.03 ± 0.62	96.12 ± 0.56
ROC-AUC	$.7350 {\scriptstyle \pm .0228}$	$.7357 {\scriptstyle\pm .0212}$	$.7346 {\pm} .0165$	$.7347 {\scriptstyle\pm .0153}$	$.7341 \scriptstyle{\pm .0107}$	$.7383 {\pm} .0073$
NLL	$.2437 {\scriptstyle \pm .1051}$	$.1760 \scriptstyle{\pm .0275}$	$.1432 {\pm} .0130$	$.1383 {\pm} .0108$	$.1821 \pm .0249$	$.1729 \scriptstyle{\pm .0208}$
ECE	$.0261 \pm .0057$	$.0224 \pm .0046$	$.0121 \pm .0039$	$.0109 \pm .0037$	$.0201 \pm .0051$	$.0193 \pm .0045$
MCE	$.2587 {\pm} .0793$	$.2617 {\pm}.0564$	$.1585 {\pm} .0760$	$.1933 {\pm} .0852$	$.1566 {\pm} .0576$	$.1533 {\pm} .0251$

F Ablation Studies

F.1 Soft or Hard labels for Classification

Often semi-supervised methods use some form of "hard-labeling" as the consistency target. Usually, this is implemented as setting the ensemble target for an unlabeled datapoint to be the most likely label, as predicted by the individual model [Filipiak et al., 2021, Tarvainen and Valpola, 2017] or the ensemble [Platanios, 2018]. This removes the underlying uncertainty information of the estimates, and risking drastically reducing the calibration of the model by making it overconfident. The motivation for using hard-labeling is the assumption of label smoothness, as it forces the model to pick the same label for data points close together. We investigate this assumption in table 17. The results on accuracy show that hard-labelling slightly benefits the accuracy, it comes at the cost of worse calibration metrics such as ECE and MCE for the individual models. The reason for such a small increase in accuracy can be explained by the label-smoothens assumption can be violated for graphs and especially molecules.

	Non-Ensemble		Ensemble		
Metric	Mean	Hard Label	Mean	Hard Label	
Accuracy (%)↑	56.0220 ± 0.2233	56.2020 ± 0.5595	56.7640 ± 0.2742	57.1920 ± 0.4124	
ROC ↑	$.9040 \pm .0017$	$.8936 {\pm}.0025$	$.7598 {\pm}.0015$	$.7621 {\pm} .0022$	
F1 ↑	$.5586 {\pm} .0021$	$.5607 {\pm} .0051$	$.5661 {\pm} .0023$	$.5706 \pm .0034$	
ECE ↓	$.1514 {\pm}.0030$	$.3034 {\pm} .0052$	$.4324 {\pm} .0027$	$.4281 {\pm} .0041$	
$MCE \downarrow$	$.2307 {\pm}.0030$	$.4252 {\scriptstyle \pm .0141}$	$.4324 {\scriptstyle\pm .0027}$	$.4281 {\scriptstyle \pm .0041}$	

Table 17: Calibration metrics on graph CIFAR-10.

F.2 Pairwise or Mean Ensemble Loss?

There is a strong theoretical connection between the pairwise loss between ensemble members used in n-CPS and the coupled ensemble loss presented in this work. For a convex loss \mathcal{L} that can be written on the form $\mathcal{L}(x-y)$, then Jensen's inequality yields

$$L(f_{\theta_i}(x) - \mathbb{E}_m[f_{\theta_m}(x)]) = L(\mathbb{E}_m[f_{\theta_i}(x) - f_{\theta_m}(x)])$$

$$\leq \mathbb{E}_m[L(f_{\theta_i}(x) - f_{\theta_m}(x))]$$

$$= \frac{1}{M} \sum_{m=1}^M L(f_{\theta_i}(x) - f_{\theta_m}(x))$$

$$\leq \frac{1}{M-1} \sum_{m=1}^M L(f_{\theta_i}(x) - f_{\theta_m}(x)).$$

As $f_{\theta_i}(x) - f_{\theta_m}(x) = 0$ if i = m this upper bound is exactly the n-CPS loss. In general this upper bound is not tight, but if M = 2 and \mathcal{L} is of the form $(x - y)^l$, e.g. the l_1 or l_2 -loss we get

$$\mathcal{L}(f_{\theta_1} - \mathbb{E}_m[f_{\theta_m}(x)]) = \left(f_{\theta_1} - \frac{f_{\theta_1} + f_{\theta_2}}{2}\right)^l$$
$$= \frac{1}{2^l}(f_{\theta_1} - f_{\theta_2})^l.$$

We see that the two losses are equal up to a scaling factor that disappears if we tune the learning rate.

F.3 Robustness of Coupled Weighting

To investigate the robustness of the coupled weighting γ , we followed the same experimental setup on CIFAR-10 with a Resnet18 model. The results can be seen in Figure 1. From the figure, we see that the validation accuracy is somewhat flat as soon as $\gamma > 1$, but there is a small optimum around $\gamma = 6$. This illustrates that at least for CIFAR-10, the choice of γ is robust.

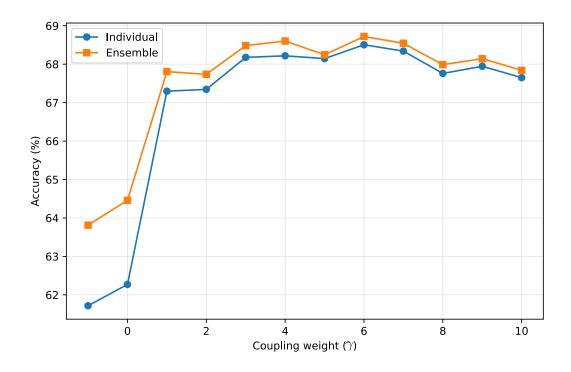


Figure 1: Validation accuracy as a function of the weighting of the ensemble consistency loss.

F.4 How to Schedule the coupled loss

Initially, during training, the members of the ensemble models only have weak prediction strength. This results in the ensemble prediction serving only as a weak signal guiding the models. Intuitively, this suggests that the weighting of the coupled loss should be added or increased as training progresses. We investigate if this is the case in the same CIFAR-10 setting. We let the ensemble coupling weighting be a linear function of the number of epochs, and vary the starting value and slope of the ensemble coupling weighting. The results can be seen in Figure 2, where negative coupling weights are clipped to 0, while Figure 3 shows the un-clipped results (in the relevant area). From Figure 2, we see that for CIFAR-10, there is no large benefit to begin coupling later compared to selecting a good constant coupling value. Note that a delayed start corresponds to a negative start value and a positive increase pr. epoch, as an initial coupling of -1 and a pr. epoch increase of 0.1 means it starts at epoch 10, due to clipping.

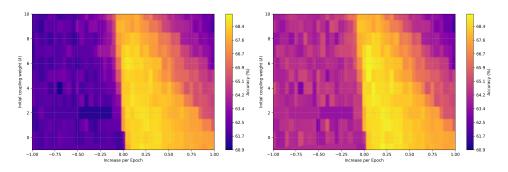


Figure 2: Validation accuracy as a function of the initial coupling weight and the increase in coupling weights per epoch for an individual model (left) and a coupled ensemble with two members (right). The results are averaged over 3 seeds.

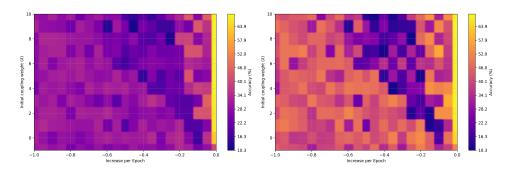


Figure 3: Validation accuracy as a function of the weighting of the ensemble consistency loss.

F.5 Different Losses

We also investigated the sensitivity to different formulations of the ensemble consistency loss. The results are shown in Table 18. We ran with the same setup for the computer vision CIFAR-10 and two ensemble members. While the best performing loss function was KL-divergence (the same form as the supervised loss), the "regression" functions (L_1, L_2, L_∞) performed about the same. Only the reversed KL-divergence, $D_{KL}(E||I)$, resulted in lower accuracy, at around the same level as a decoupled model (see Table 4).

Table 18: Validation accuracy with different ensemble consistency loss functions. Results averaged over 10 seeds. Here, I is the individual prediction and E is the ensemble consensus.

Ensemble Loss	Individual Accuracy
$\begin{array}{c} L_{\infty} \\ D_{KL}(I E) \\ D_{KL}(E I) \\ L_{1} \\ L_{2} \end{array}$	$66.23 \pm 0.29 \\ 66.62 \pm 0.51 \\ 59.37 \pm 0.78 \\ 66.01 \pm 0.51 \\ 66.12 \pm 0.45$

F.6 Different coupling strategies

We investigated different strategies for coupling the unsupervised loss on QM9. This includes various combinations of three parameters: the *coupling weight*, the *coupling start* and the *coupling schedule*.

Coupling weight The coupling weight parameter defines how much the unsupervised loss should contribute to the total loss. When set to 0, only the supervised loss will be taken into account.

Coupling start The coupling start refers to when the unsupervised loss in included during training, i.e. for the first x% of epochs, the model is only trained on the labeled data and only afterwards, the unsupervised loss with be included via coupling. Depending on the dataset and task, it intuitively can make sense to first let the model learn a little bit before evaluating the loss on unlabeled data. Specifically, in regression tasks this can be the case, since the model output is not bounded, as opposed to classification tasks. When set to 0, coupling will be used through the whole training. This parameter is given in percentage, i.e. percentage of total training epochs after which the coupling should start.

Coupling schedule Three different coupling schedules were tested: *constant*, *increase* and *bell*. *Constant* refers to the the coupling weight being constant from onset until the end of training. *Increase* means that the there will be a smooth ramp up until the coupling weight reaches its maximum (i.e. the coupling weight parameter). *Bell* means that there is a smooth bell curve over the coupling weight, i.e. first in increases, then decreases. Here, it will start and end at 0, and peak at a maximum which is set via the coupling weight parameter.

Figure 4 and Figure 5 shows the impact of different coupling strategies on the model performance, here for target 4 and 7 of QM9 respectively. We can see that a good choice of the coupling weight is crucial for our method to result in a significant improvement in MAE compared to the fully supervised baseline. The optimal coupling weight seems to differ per task, as both targets have a different optimum (0.1 for target 4 and 0.01 for target 7). A good value for the coupling start seems to depend on the choice of coupling weight, however a trend can be observed that for the best coupling weight options for each target, the optimal coupling start is 0, i.e. using coupling from the start of training. The optimal choice of coupling schedule seems to depend on both of the other choices, but in the specific case of target 4, the *increase* schedule led to the best performance. For target 7, the *bell* schedule resulted in the best ensemble performance, while the *constant* schedule led to the best individual performance.

One interesting finding here is that if we couple too strongly, meaning we are weighing the unsupervised loss to high, the ensemble performance gets worse than the baseline, while at the same time the individual members from the ensemble are outperforming the baseline. This is due to the models collapsing, so while each individual model is better than an individual model that was not coupled, ensembling has no significant benefit anymore.

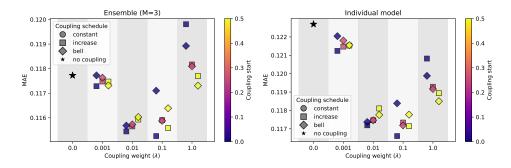


Figure 4: Performance (MAE) of coupled ensembles (left) and individual models from coupled ensembles (right) for different coupling strategies, for QM9 target 4.

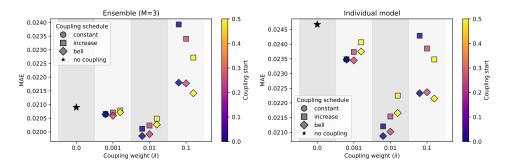


Figure 5: Performance (MAE) of coupled ensembles (left) and individual models from coupled ensembles (right) for different coupling strategies, for QM9 target 7.

F.7 Evaluating Overfitting on Unlabeled Data

To evaluate potential overfitting to the unlabeled data, we compare the final model's performance on the unlabeled training set against its performance on the unseen test set. For this analysis, we leverage our access to the ground-truth labels of the unlabeled set to compute its MAE. As presented in Table 19, the performance is nearly identical across both datasets for all 12 QM9 targets. This strong correspondence indicates that our method avoids overfitting to the unlabeled data used during training. This has a significant practical benefit, as it means the model's predictions on the entire unlabeled set can be reliably used for downstream tasks.

Table 19: PaiNN performance (MAE) on QM9 targets, comparing the held-out test set with the unlabeled dataset used during training. Results are reported for 5 seeds.

Target	Unit	Data	Individual Member	Ensemble (M=4)
μ	D	Test Unlabeled	$.0619 \pm .0003$ $.0596 \pm .0003$	$.0613 \pm .0003 \\ .0596 \pm .0003$
α	a_0^3	Test Unlabeled	$.1322 \pm .0011 \\ .1268 \pm .0008$	$.1303 \pm .0011 \\ .1261 \pm .0008$
$\epsilon_{ ext{HOMO}}$	meV	Test Unlabeled	$73.9789 {\pm}.4368 \\ 71.7113 {\pm}.4012$	$73.0755 {\pm} .4472 \\ 71.6826 {\pm} .4018$
$\epsilon_{ m LUMO}$	meV	Test Unlabeled	$\begin{array}{c} 57.7186 {\pm}.2247 \\ 56.8810 {\pm}.1844 \end{array}$	$\begin{array}{c} 57.2369 {\pm}.2159 \\ 56.8676 {\pm}.1839 \end{array}$
$\Delta\epsilon$	meV	Test Unlabeled	$117.0365 {\pm}.4988 \\ 114.1592 {\pm}.3078$	$115.7195 {\scriptstyle \pm .5100} \\ 114.1303 {\scriptstyle \pm .3091}$
$\langle R^2 \rangle$	a_0^2	Test Unlabeled	$.6100 \pm .0206 \\ .5918 \pm .0205$	$.5605 \pm .0206 \\ .5552 \pm .0202$
ZPVE	meV	Test Unlabeled	$\begin{array}{c} 2.0138 \pm .0054 \\ 1.9925 \pm .0066 \end{array}$	$1.9907 {\pm}.0055 \\ 1.9883 {\pm}.0066$
U_0	meV	Test Unlabeled	$19.9642 {\pm}.1291 \\ 19.3096 {\pm}.1434$	$19.3816 {\scriptstyle \pm .1278} \\ 18.9715 {\scriptstyle \pm .1416}$
\overline{U}	meV	Test Unlabeled	$20.1731 {\scriptstyle\pm .1577} \\ 19.5288 {\scriptstyle\pm .1248}$	$19.5886 {\scriptstyle \pm .1574} \\ 19.1908 {\scriptstyle \pm .1234}$
Н	meV	Test Unlabeled	$20.1407 {\pm}.1268 \\ 19.5028 {\pm}.1370$	$19.5509 {\pm}.1328 \\ 19.1620 {\pm}.1355$
\overline{G}	meV	Test Unlabeled	$20.3142 {\pm}.1571 \\ 19.7490 {\pm}.1384$	$19.7479 {\scriptstyle \pm .1634} \atop 19.4296 {\scriptstyle \pm .1400}$
C_v	cal mol K	Test Unlabeled	$.0449 \pm .0002$ $.0443 \pm .0001$	$.0439 {\pm} .0002 \\ .0439 {\pm} .0001$