

ROBUST LEARNING IN BAYESIAN PARALLEL BRANCHING GRAPH NEURAL NETWORKS: THE NARROW WIDTH LIMIT

Anonymous authors

Paper under double-blind review

ABSTRACT

The infinite width limit of random neural networks is known to result in Neural Networks as Gaussian Process (NNGP) (Lee et al. (2018)), characterized by task-independent kernels. It is widely accepted that larger network widths contribute to improved generalization (Park et al. (2019)). However, this work challenges this notion by investigating the narrow width limit of the Bayesian Parallel Branching Graph Neural Network (BPB-GNN), an architecture that resembles residual GCN. We demonstrate that when the width of a BPB-GNN is significantly smaller compared to the number of training examples, each branch exhibits more robust learning due to a symmetry breaking of branches in kernel renormalization. Surprisingly, the performance of a BPB-GNN in the narrow width limit is generally superior to or comparable to that achieved in the wide width limit in bias-limited scenarios. Furthermore, the readout norms of each branch in the narrow width limit are mostly independent of the architectural hyperparameters but generally reflective of the nature of the data. We also extend the results to other more general architectures such as the residual-MLP and demonstrate that the narrow width effect is a general feature of the branching networks. Our results characterize a newly defined narrow-width regime for parallel branching networks in general.

1 INTRODUCTION

The study of neural network architectures has seen substantial growth, particularly in understanding how network width impacts learning and generalization. In general, wider networks are believed to perform better (Allen-Zhu et al. (2019); Jacot et al. (2018); Gao et al. (2024)). However, this work challenges the prevailing assumption by exploring the narrow width limit of Bayesian Parallel Branching Graph Neural Networks (BPB-GNNs), an architecture inspired by residual GCN networks (Chen et al. (2020a; 2022)). We show theoretically and empirically that narrow-width networks can perform better than their wider counterparts due to a symmetry-breaking effect in kernel renormalization, in bias-limited scenarios. This paper presents a detailed theoretical analysis of BPB-GNNs in the narrow-width regime, highlighting realistic conditions under which these networks demonstrate robust learning and comparable generalization.

Contributions :

1. We introduce a novel yet simple GCN architecture with parallel independent branches, and derive the exact generalization error for node regression in the statistical limit as the sample size $P \rightarrow \infty$ and network width $N \rightarrow \infty$, with their ratio a finite number $\alpha = P/N$, in the over-parameterized regime.
2. We show that in the Bayesian setting the bias will decrease and saturate at a narrow hidden layer width, a surprising phenomenon due to kernel renormalization. We demonstrate that this can be understood as a robust learning effect of each branch in the student-teacher task, where each student's branch is learning the teacher's branch.
3. We demonstrate this narrow-width limit in the real-world dataset Cora and understand each branch's importance as a nature of the dataset.

- 054 4. We further show that this narrow-width effect is a general feature of branching networks,
 055 with the residual-MLP architecture as an example.
 056

057 2 RELATED WORKS

059 **Infinitely wide neural networks:** Our work follows a long tradition of mathematical analysis
 060 of infinitely-wide neural networks (Neal (2012); Jacot et al. (2018); Lee et al. (2018); Bahri et al.
 061 (2024)), resulting in NTK or NNGP kernels. Recently, such analysis has been extended to structured
 062 neural networks, including GNNs (Du et al. (2019); Walker & Glocker (2019); Huang et al. (2021)).
 063 However, they do not provide an analysis of feature learning in which the kernel depends on the
 064 tasks.
 065

066 **Kernel renormalization and feature learning:** There has been progress in understanding simple
 067 MLPs in the feature-learning regime as the shape of the kernel changes with task or time (Li &
 068 Sompolinsky (2021); Atanasov et al. (2021); Avidan et al. (2023); Wang & Jacot (2023)). We
 069 develop such understanding in graph-based networks.
 070

071 **Theoretical analysis of GCN:** There is a long line of works that theoretically analyze the expres-
 072 siveness (Xu et al. (2018); Geerts & Reutter (2022)) and generalization performance (Tang & Liu
 073 (2023); Garg et al. (2020); Aminian et al. (2024)) of GNN. However, it is challenging to calcu-
 074 late the dependence of generalization errors on tasks. In particular, the PAC-Bayes approach Liao
 075 et al. (2020); Ju et al. (2023) results in generalization bounds that are too large and that can be only
 076 computed with norms of learned weights. To our knowledge, our work is first to decompose the
 077 generalization error into bias and variance *a priori* (not dependent on learned weights) for linear
 078 GCNs with residual-like structures. The architecture closest to our linear BPB-GCN is the linearly
 079 decoupled GCN proposed by Cong et al. (2021); however, the overall readout vector is shared for
 080 all branches, which will not result in kernel renormalization for different branches.

081 3 BPB-GNN

082 3.1 PARALLEL BRANCHING GNN ARCHITECTURE

085 We are motivated to study the parallel branching networks as they resemble residual blocks in com-
 086 monly used architectures (Kipf & Welling (2016); He et al. (2016); Chen et al. (2020b)) and are
 087 easy to study analytically with our Bayesian framework. Given graph $G = (A, X)$, where A is the
 088 adjacency matrix and X the node feature matrix, the final readout for node μ is a scalar $f^\mu(G; \Theta)$
 089 that depends on the graph and network parameters Θ . The parallel branching GNN is an ensemble of
 090 GNN branches, where each branch operates independently with no weight sharing. In this work, we
 091 focus on the simple setup of branches made of linear GCN with one hidden layer, but with different
 092 number of convolutions A^l on the input node features (Figure 1(a)). In this way, parallel branching
 093 GNN is analogous to GCN with residual connections, for which the final node readout can also be
 094 thought of as an ensemble of convolution layers (Veit et al. (2016)). Concretely, the overall readout
 095 $f^\mu(G; \Theta)$ for node μ is a sum of L branch readouts f_i^μ :

$$096 f^\mu(G; \Theta) = \sum_{l=0}^{L-1} f_i^\mu(G; \Theta_l = \{W^{(l)}, a^{(l)}\}), \quad (1)$$

098 where

$$099 f_i^\mu(G, \Theta_l) = \frac{1}{\sqrt{L}} \sum_{i=1}^N \frac{1}{\sqrt{N}} a_i^{(l)} \sum_{j=1}^{N_0} \frac{1}{\sqrt{N_0}} W_{ij}^{(l)} \sum_{\nu=1}^n (A^l)_{\mu\nu} x_j^\nu \quad (2)$$

102 In matrix notation,

$$103 F = \sum_l \frac{1}{\sqrt{LNN_0}} A^l X W^{(l)} a^{(l)} \quad (3)$$

106 ¹In this paper, the convolution operation is normalized as $A = D^{-1/2}(\hat{A} + I)D^{-1/2}$, where \hat{A} is the
 107 original Adjacency matrix and D is the degree matrix. We also use feature standardization after convolution
 for each branch to normalize the input

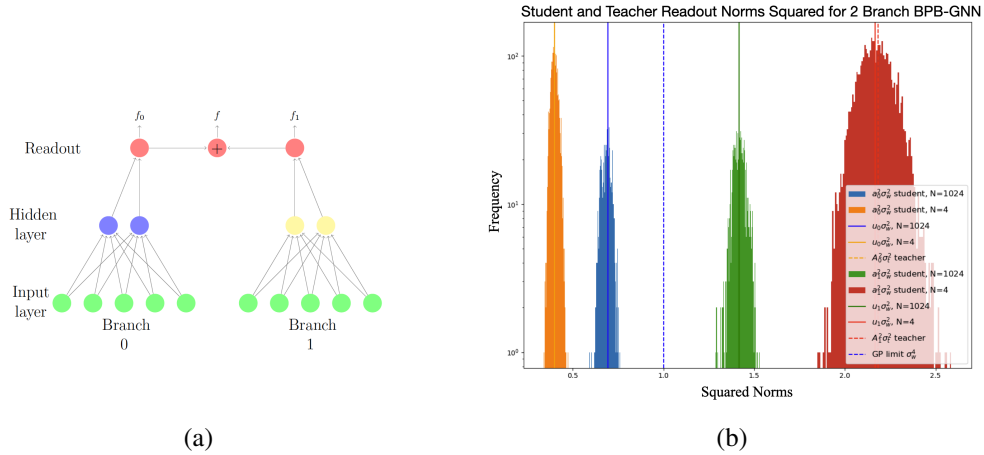


Figure 1: Overview of the main takeaway: BPB-GNN learns robust representations for each branch at narrow width. (a) The parallel branching GNN architecture, with 2 branches. The independent branches have non-sharing weights and produce the final output f as a sum of branch-level readouts f_l . (b) Student and teacher readout norms squared for wide and narrow student BPB-GNN networks. The student network with width N is trained with the teacher network’s output. Histograms correspond to the samples from Hamiltonian Monte Carlo simulations and solid lines correspond to the order parameters calculated theoretically. $\sigma_t = \sigma_w = 1$. At $N = 4$, the HMC samples of branch readout norms squared (orange and red histograms) for the student network $\frac{\|A_l\|^2}{N} \sigma_w^2$ concentrate at their respective theoretical values $u_l \sigma_w^2$ and overlap with the teacher’s readout norms squared $\frac{\|A_l\|^2}{N} \sigma_t^2$ (orange and red dashed lines) for corresponding branches. At $N = 1024$ the samples for the student network (blue and green histograms) concentrate at their respective theoretical values but remain far from the teacher’s values, instead approaching the GP limit σ_w^4 (blue dashed line).

Note that when $L = 2$, the network reduces exactly to a 2-layer residual GCN (Chen et al. (2020c)). Here N_0 is the input dimension, N the width of the hidden layer, $W^{(l)}$ the weight of the hidden layer for branch l , A^l the l th power of the adjacency matrix and $a^{(l)}$ the final readout weight for branch l . We will consider only the linear activation function for this paper and provide a summary of notations in Appendix A.1.

3.2 BAYESIAN NODE REGRESSION

We consider a Bayesian semi-supervised node regression problem, for which the posterior probability distribution for the weight parameters is given by

$$P(\Theta) = \frac{1}{Z} e^{-E(\Theta; G, Y)/T} = \frac{1}{Z} \exp\left(-\frac{1}{2T} \sum_{\mu=1}^P (f^\mu(G, \Theta) - y^\mu)^2 - \frac{1}{2\sigma_w^2} \Theta^T \Theta\right), \quad (4)$$

where the first term in the exponent corresponds to the likelihood term induced by learning P node labels y_μ with squared loss and the second term corresponds to the Gaussian prior with variance σ_w^2 . $Z = \int e^{-E(\Theta)/T} d\Theta$ is the normalization constant. This Bayesian setup is well motivated, as the Langevin dynamics trained with energy potential E and temperature T that results in this equilibrium posterior distribution shares a lot in common with the Gradient Descent (Avidan et al. (2023); Naveh et al. (2021)) and Stochastic GD optimizers (Mignacco & Urbani (2022); Wu et al. (2020)). In fact, Li & Sompolinsky (2021) shows empirically that the Bayesian equilibrium is a statistical distribution of the usual gradient descent with early stopping optimization and with random initializations at 0 temperature in DNNs, where the L_2 regularization strength σ_w^2 corresponds to the Gaussian initialization variance.

We are interested in understanding the weight and predictor statistics of each branch and how they contribute to the overall generalization performance of the network. In the following theoretical derivations, our working regime is in the overparameterizing high-dimensional limit (Li & Som-

polinsky (2021); Montanari & Subag (2023); Bordelon & Pehlevan (2022); Howard et al. (2024)): $P, N, N_0 \rightarrow \infty$, $\frac{P}{N} = \alpha$ finite and the capacity $\alpha_0 = \frac{P}{LN_0} < 1$. As we will show later, this limit is practically true even with P, N not so large (our smallest N is 4). We will also use near-0 temperature, in which case the training error will be near 0.

3.3 KERNEL RENORMALIZATION AND ORDER PARAMETERS

The normalization factor, or the partition function, $Z = \int e^{-E(\Theta)/T} d\Theta$ carries all the information to calculate the predictor statistics and the generalization dependence on network hyperparameters N, L, σ_w^2 . Using Eq. 2,4, we can integrate out the readout weights a_l 's first, resulting in

$$Z = \int dW e^{-H(W)}, \quad (5)$$

with effective Hamiltonian $H(W)$ in terms of the hidden layer weights for all branches

$$H(W) = \frac{1}{2\sigma_w^2} \sum_{l=0}^{L-1} \text{Tr} W_l^T W_l + \frac{1}{2} Y^T (K(W) + TI)^{-1} Y + \frac{1}{2} \log \det(K(W) + TI), \quad (6)$$

where

$$K(W) = \frac{1}{L} \sum_l \frac{\sigma_w^2}{N} (H_l(W_l) H_l(W_l)^T)|_P \quad (7)$$

is the $P \times P$ kernel matrix dependent on the observed P nodes with node features $H_l = A^l X W_l$ and we denote $|_P$ as the matrix restricting to the elements generated by the training nodes.

As shown in Appendix A.2, we can further integrate out the W_l 's and get the partition function $Z = \int Du e^{-H(u)}$ described by a final effective Hamiltonian independent of weights

$$H(u) = S(u) + E(u), \quad (8)$$

where we call $S(u)$ the entropic term

$$S(u) = - \sum_l \frac{N}{2} \log u_l + \sum_l \frac{N}{2\sigma_w^2} u_l \quad (9)$$

and $E(u)$ the energetic term

$$E(u) = \frac{N\alpha}{2P} Y^T \left(\sum_l \frac{1}{L} u_l K_l + TI \right)^{-1} Y + \frac{N\alpha}{2P} \log \det \left(\sum_l \frac{1}{L} u_l K_l + TI \right), \quad (10)$$

where $K_l = \frac{\sigma_w^2}{N_0} [A^l X X^T A^l]|_P$ is the $(P \times P)$ input node feature kernel for branch l .

Therefore, the final effective Hamiltonian has the overall kernel

$$K = \sum_l \frac{1}{L} u_l K_l, \quad (11)$$

where u_l 's are order parameters which are the minimum of the effective Hamiltonian Eq. 8 that satisfy the saddle point equations:

$$N \left(1 - \frac{u_l}{\sigma_w^2} \right) = -r_l + \text{Tr}_l, \quad (12)$$

where $r_l = Y^T (K + TI)^{-1} \frac{u_l K_l}{L} (K + TI)^{-1} Y$ and $\text{Tr}_l = \text{Tr} [K^{-1} \frac{u_l K_l}{L}]$.

GP kernel vs. renormalized kernel: Observe that as $\alpha = P/N \rightarrow 0$, the entropic term dominates, and thus $u_l = \sigma_w^2$ for all branches from Eq. 12. This is the usual Gaussian Process (GP) limit, or the infinite-width limit. In this case, the kernel K_l for each branch is not changed by the training data, and each branch has the same contribution in terms of its strength.

However, when $\alpha = P/N$ is large, there is a correction to the GP prediction, and u_l 's in general depend on the training data; therefore, we have feature learning in each branch with kernel renormalization. It turns out that u_l is exactly the statistical average of readout norm squared over the posterior distribution (Appendix A.4), ie.

$$u_l = \langle \|a_l\|^2 \rangle / N \quad (13)$$

for each branch l , where we use the bracket notation to represent the expectation value over the posterior distribution. Therefore, branches in general become more and more different as the width of the hidden layer N decreases. We call this phenomenon *symmetry breaking*, which is discussed in Section 4.2.

3.4 PREDICTOR STATISTICS AND GENERALIZATION

Under the theoretical framework, we obtain analytically (Appendix A.3) the mean predictor $\langle f^\nu(G) \rangle$ and variance $\langle \delta f_\nu(G)^2 \rangle$ for a single test node ν as Eq.54 and Eq.55, respectively. We use this to calculate the generalization performance, defined by the MSE on t test nodes

$$\epsilon_g = \left\langle \frac{1}{t} \sum_{\nu=1}^t (f^\nu(G) - y^\nu)^2 \right\rangle_\Theta = \text{Bias} + \text{Variance}, \quad (14)$$

where

$$\text{Bias} = \frac{1}{t} \sum_{\nu=1}^t (\langle f^\nu(G) \rangle_\Theta - y^\nu)^2, \text{Variance} = \frac{1}{t} \sum_{\nu=1}^t \langle \delta f_\nu^2 \rangle. \quad (15)$$

Note that our definition of bias and variance is a statistical average over the posterior weight distribution, which is slightly different from the usual definition from GD.

4 THE NARROW WIDTH LIMIT

As we discussed briefly in Section 3.3, the kernel becomes highly renormalized at narrow width. In fact, in the extreme scenario as $\alpha = P/N \rightarrow \infty$, the energetic term in the Hamiltonian completely dominates, and we would expect that the generalization performance saturates as the order parameters in the energetic terms become independent of width N . Therefore, just as infinitely wide networks correspond to the GP limit, we propose that there exists a *narrow width limit* when the network width is extremely small compared to the number of training samples.

4.1 ROBUST LEARNING OF BRANCHES: THE EQUIPARTITION CONJECTURE

What happens in the narrow width limit? In the following, we demonstrate that each branch will learn robustly at narrow width.

The equipartition theorem Consider a student-teacher network setup, where the teacher network is given by

$$f^*(G; \Theta^*) = \sum_l f_l^*(G; W_l^*) = \sum_l \frac{1}{\sqrt{N_t L}} \sum_i a_{i,l}^* h_i^l(G; W_l^*). \quad (16)$$

$W_{ij,l}^* \sim \mathcal{N}(0, \sigma_t^2)$ and $a_{i,l}^* \sim \mathcal{N}(0, \beta_l^2)$, where β_l^2 is the variance assigned to the readout weight for the teacher branch l and N_t is the width of the hidden layer. Similarly, the student network is given by the same architecture, with layer width N and learns from P node labels from the teacher $Y_\mu = f_\mu^*(G, \Theta^*)$ in the Bayesian regression setup of Eq. 4 with prior variance σ_w^2 . We state that at $T = 0$ temperature, in the limit $N_t \rightarrow \infty$ and $\alpha = P/N \rightarrow \infty$, there exists a solution to the saddle point equations 12

$$u_l = \sigma_t^2 \beta_l^2 / \sigma_w^2. \quad (17)$$

Proof:

At the narrow width limit $\alpha \rightarrow \infty$, the saddle equation 12 becomes $r_l = \text{Tr}_l$. Now we calculate

$$Y_\mu^* Y_\mu^{*T} = \sum_{i,j=1}^{N_t} \sum_{l_1, l_2} \frac{1}{N_t L} a_{i,l_1}^* a_{j,l_2}^* h_i^{l_1, \mu}(G; W_{l_1}) h_j^{l_2, \nu}(G; W_{l_2}) \quad (18)$$

by law of large numbers, this quantity concentrates at its expectation value

$$Y^* Y^{*T} = \lim_{N_t \rightarrow \infty} \mathbb{E}_{a^*, W^*} (Y Y^T) = \sum_l \frac{\beta_l^2 \sigma_t^2}{N_0 L} (A^l X X^T A^l)|_P = \frac{\beta_l^2}{L} K_l(\sigma_t^2), \quad (19)$$

where $K_l(\sigma^2)$ represents the GP input kernel with prior variance σ^2 . Therefore, r_l becomes

$$r_l = Y^{*T} K^{-1} \frac{u_l K_l}{L} K^{-1} Y^* = \text{Tr}(K^{-1} \frac{u_l K_l}{L} Y^* Y^{*T}) = \text{Tr}_l((\sum_l u_l K_l(\sigma_w^2)/L)^{-1} \sum_l \beta_l^2 \frac{K_l}{L}) \quad (20)$$

Thus there exists a solution that satisfies the saddle point equations

$$u_l \sigma_w^2 = \beta_l^2 \sigma_l^2. \quad (21)$$

Furthermore, by Eq.13, we show that the student branch norms learn exactly the teacher branch norms, ie. $\langle a_l^2 \rangle \sigma_w^2 = a_l^{*2} \beta_l^2$. We call this equipartition theorem, as the mean-squared readout and the variance (A.4) have to exactly balance each other, which contribute to the energy term in the Hamiltonian. Although the theorem is proven for the BPB-GNN architecture, we show that the result is general in Appendix B.1.

4.2 STUDENT-TEACHER EXPERIMENT ON ROBUST BRANCH LEARNING

We demonstrate this robust learning phenomenon and provide a first evidence of the equipartition conjecture with the student-teacher experiment setup introduced in the previous section. We use the contextual stochastic block model (CSBM) (Deshpande et al. (2018)) to generate the graph, where the adjacency matrix is given by a stochastic block model with two blocks, and the node feature is generated with latent vectors corresponding to the two blocks (Appendix C.1). Both student and teacher network has $L = 2$ branches. We calculate the order parameters u_l 's for the student network with the saddle point equations, as well as perform Hamiltonian Monte Carlo (HMC) sampling (Appendix C.5) from the posterior distribution of the student's network weights to generate a distribution of student readout norms squared a_l^2 . We compare the values of the student branch readout norms to their corresponding teacher branches for the student network with either narrow or wide width. As shown in Figure 1(b), an extremely narrow student network ($N = 4$) learns the teacher's branch readout norms very robustly, while a much wider network ($N = 1024$) fails to learn the branches' norms and approaches the GP ($N \rightarrow \infty$) limit where the two branches are indistinguishable.

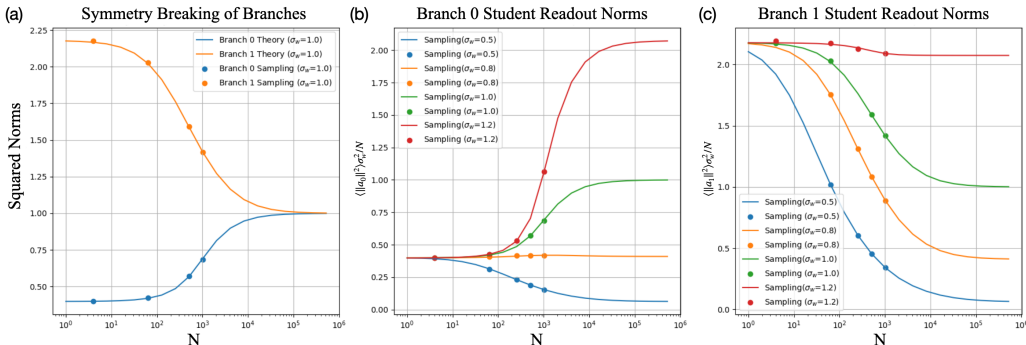


Figure 2: Statistical average of student readout norms squared as a function of network width from theory and HMC sampling, for student-teacher tasks described in Section 4.2. (a): $\langle \|a_l\|^2 \rangle \sigma_w^2 / N$ as a function of network width N for a fixed σ_w . The branch norms break the GP symmetry as it goes to the narrow width limit. (b)(c): Branch 0 and branch 1 readout norm squared respectively for a range of σ_w regularization values. The student branch norms with different regularization strengths all converge to the same teacher readout norm values at narrow width.

Symmetry breaking and convergence of branches: We perform theoretical calculations and HMC sampling for the student branch squared readout norms as we vary the student network width N and the prior regularization strength σ_w . Figure 2(a) shows the statistical average of the student branch squared readout norms, ie. $\langle \|a_l\|^2 \rangle \sigma_w^2 / N$ as a function of the network width N , where the branch norms split as the network width gets smaller, which we call symmetry breaking. The symmetry breaking of branch norms from the GP limit to the narrow width limit accompanies the convergence to learning teacher's norms at narrow width for different σ_w 's as shown in Figure 2(b)(c), supporting Eq. 62.

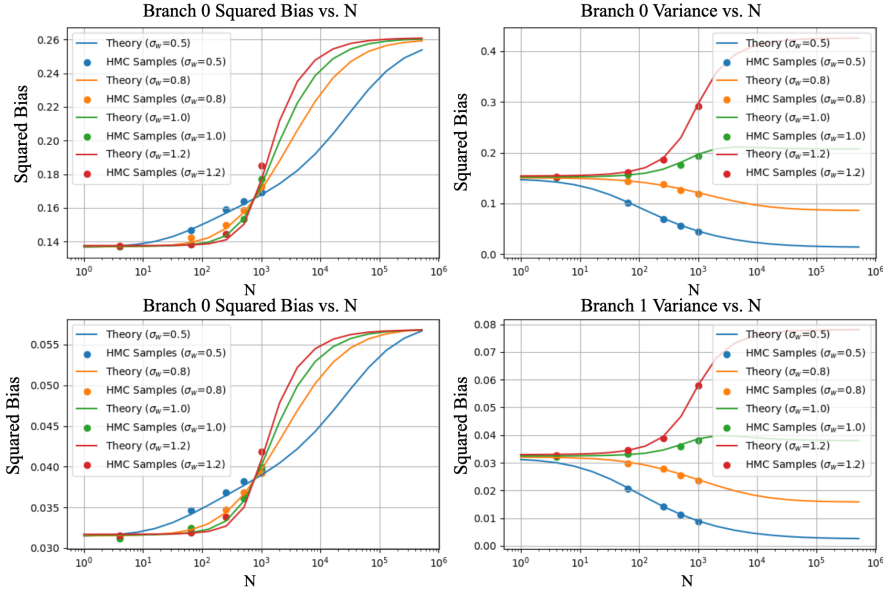


Figure 3: Student network squared bias and variance for individual branches as a function of network width N and regularization strength σ_w . The mean and variance of branch l readout f_l^μ for node μ is calculated in A.3 and the bias and variance for branch l can be inferred similarly as Eq. 15. Generalization values are normalized over the average true readout labels.

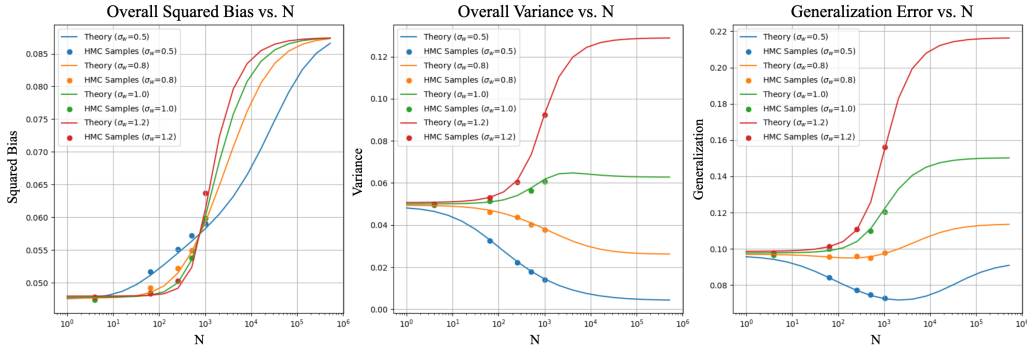


Figure 4: Student network generalization performance as a function of network width N and regularization strength σ_w . Generalization is normalized over the average true readout labels.

Narrow-to-wide width transition: We can also determine the generalization properties of the student readout labels at both the branch level f_l (Appendix A.3) and overall level f using Eq. 1415. We perform theoretical calculations and HMC sampling of generalization errors as a function of the network width N and regularization strength σ_w , with results shown in Figure 3 and 4. At narrow width, we expect individual branch f_l to learn the corresponding teacher’s branch output f_l^* independently, causing the bias to increase with network width. This is observed for both branches, with a transition from the narrow-width regime to the GP regime. The regularization strength σ_w controls the transition window, with larger σ_w ’s leading to sharper transitions. This aligns with our analysis of the entropic and energetic contributions, where the larger σ_w amplifies the distinction between the two terms. In contrast, the variance decreases with network width for small σ_w ’s, resulting in a trade-off between the contributions of bias and variance to overall generalization performance, as shown clearly with the graph of network generalization vs. N with $\sigma_w = 0.5$ in Figure 4.

378
 379
 380
 381
 382
 383
 384
 385
 386
 387
 388
 389
 390
 391
 392
 393
 394
 395
 396
 397
 398
 399
 400
 401
 402
 403
 404
 405
 406
 407
 408
 409
 410
 411
 412
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431

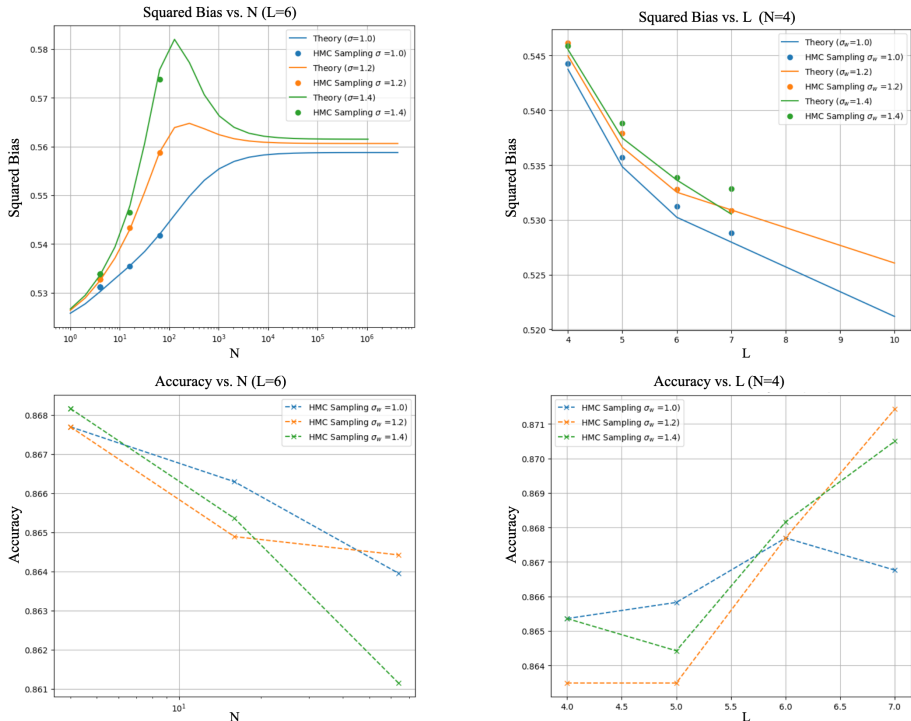


Figure 5: Cora generalization performance vs. network width N and branch number L , for various regularization strength σ_w 's. The accuracy is computed by turning the mean predictor from HMC samples into a class label using its sign.

5 BPB-GNN ON CORA

We also perform experiments on the Cora benchmark dataset (McCallum et al. (2000)) by training the BPB-GNN with binary node regression, for a range of L, N, σ_w values (Appendix C.2 for details). We observe a similar narrow-to-wide width transition for the bias term. As shown in Figure 5, the bias increases with network width, transitioning to the GP regime, and we observe the trend extending to a potential narrow width limit.² Additionally, it is demonstrated that using more branches that involve higher-order convolutions improves performance.

Convergence of branch readout norms at narrow width An interesting aspect of the BPB-GNN network is that the branch readout norms converge at the narrow width for different hyperparameters σ_w and L , reflecting the natural branch importance for the task.

As shown in Figure 6, the BPB-GNN with branches $L = 6$ robustly learns the readout norms at narrow width independently of σ_w 's, consistent with the student-teacher results. This suggests that we can recast the data as generated from a ground-truth teacher network even for real-world datasets. The last branch of the BPB-GNN network has a larger contribution, reflecting the presence of higher-order convolutions in the Cora dataset. From a kernel perspective, increasing branches better distinguish the nodes, as shown in Figure 7. This could explain the selective turn-off of intermediate branches and the increased contribution of the last branch. Our results also indicate that there is no oversmoothing problem at narrow width: as L increases, individual branches can still learn robustly.

Furthermore, the first two branches are learned most robustly at narrow width, as shown in Figure 8, where the branch norms converge for the first two branches even for BPB-GNNs with different

²In this case, the narrow width limit is hard to demonstrate as the transition window is below realistic minimum of network width.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

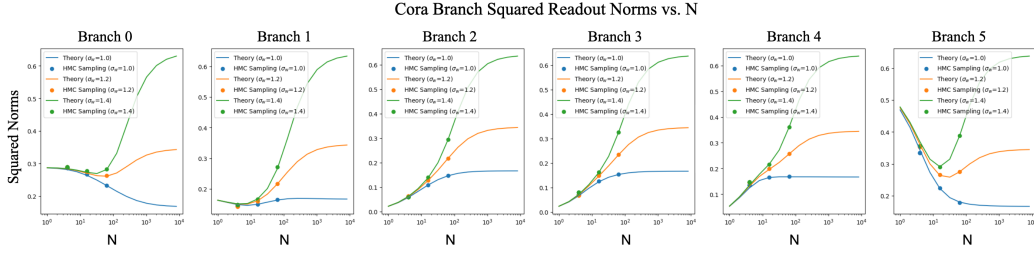


Figure 6: Cora experiment: statistical average of squared readout norms $\langle \|a_l\|^2 \rangle \sigma_w^2 / N$ for each branch l as a function of the network width N , and regularization strength σ_w . The BPB-GNN has $L = 6$ branches.

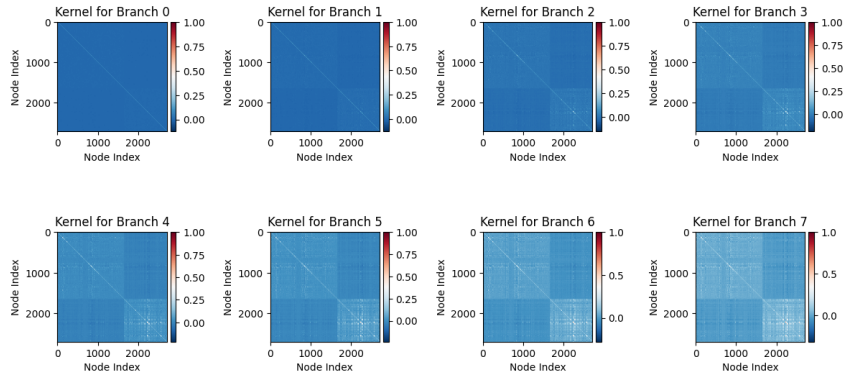


Figure 7: Kernel $K_l = \frac{\sigma_w^2}{N_0} A^l X_0 X_0^T A^l$ for each branch l shown for the first 8 branches on Cora dataset, sorted by node labels. A is the normalized adjacency matrix of the Cora graph, X_0 the node feature matrix and N_0 the node feature dimension. Initialization variance $\sigma_w^2 = 1$ and total node number $n = 2708$.

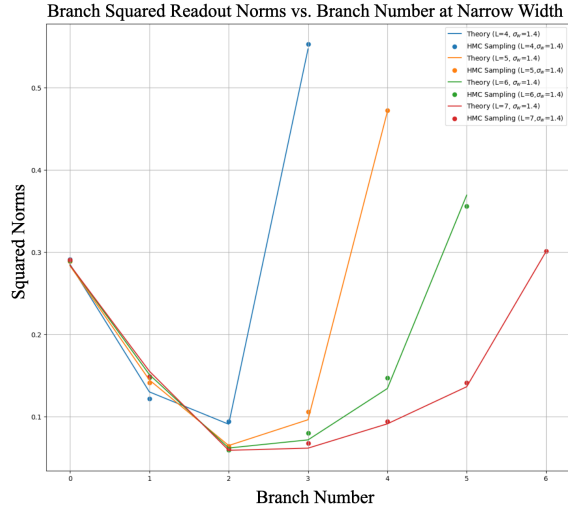


Figure 8: Branch Importance vs. branch number l on Cora. Legends represent different total number of branches L . The branch importance is defined as the statistical average of branch squared readout norms $\langle \|a_l\|^2 \rangle \sigma_w^2 / N$ at the narrow width limit; here we take the empirical branch norm values at $N = 4$ and fixed $\sigma_w = 1.4$.

L. This suggests that the branch importance, as reflected by the norms learned at narrow width, indicates the contribution of the bare data and the first convolution layer.

6 DISCUSSION

The findings presented in this paper reveal that BPB-GNNs exhibit unique characteristics in the narrow width limit. Unlike the infinite-width limit, where neural networks behave as Gaussian Process (GP) with task-independent kernels, narrow-width BPB-GNNs undergo significant kernel renormalization. This renormalization leads to breaking of the symmetry between the branches, resulting in more robust and differentiated learning. Our experiments demonstrate that narrow-width BPB-GNNs can retain and, in some cases, improve generalization performance compared to their wider counterparts, particularly in bias-limited scenarios where the regularization effects dominate. Therefore, there is a trade-off between the expressive power of GNN Xu et al. (2018); Loukas (2019) and stability Stogin et al. (2020). Additionally, the observed independence of readout norms from architectural hyperparameters suggests that narrow-width BPB-GNNs can capture the intrinsic properties of the data more effectively. These insights suggest potential new strategies for optimizing neural network architectures in practical applications, challenging the traditional emphasis on increasing network width for better performance. Although our work is focused on GNN architectures, our findings are more general for a family of branching networks as shown in Appendix B, and in particular the Bayesian theory with kernel renormalization is shown to be successful for CNN Aiudi et al. (2023) as well as transformers Tiberi et al. (2024). Intriguingly, Vyas et al. (2024) finds that an ensemble of ResNets averaged over random initializations accidentally shows this narrow width effect, for which our findings provide an explanation.

Limitations: We demonstrate that the bias robustly decreases at narrow width for the branching networks due to kernel renormalization, which defies the common belief that wider networks are better (Allen-Zhu et al. (2019); Lee et al. (2018)). However, we do find that for real-world scenarios such as BPB-GNN on Cora, and residual-MLP on Cifar10, the generalization error is usually variance-limited. In principle, training an ensemble of branching networks averaged over random initialization helps to reduce the variance significantly such that the narrow-width effect is more pronounced as in Vyas et al. (2024). Furthermore, since our Bayesian setup corresponds to an ensemble of fully trained over-parametrized networks in an offline fashion, the learning is still in the so-called lazy regime with minimal feature learning (Atanasov et al. (2022); Karkada (2024)). We hypothesize that narrower width serves as a regularizer in this regime which might break down as the network is trained with online learning and approaches the rich regime (Yang et al. (2022)) with more feature learning. In that case, it might be possible that even the bias does not decrease at the narrow width.

7 CONCLUSION

In conclusion, this paper introduces and investigates the concept of narrow width limits in Bayesian branching networks. Our results indicate that branching networks with significantly narrower widths can achieve better or competitive performance, contrary to common beliefs. This is attributed to effective symmetry breaking and kernel renormalization in the narrow-width limit, which leads to robust learning. Our theoretical analysis, supported by empirical evidence, establishes a new understanding of how network architecture influences learning outcomes. This work provides a novel perspective on the infinite width limit of neural networks and suggests further research directions in understanding narrow width as a regularization effect.

REFERENCES

- R Aiudi, R Pacelli, Alessandro Vezzani, Raffaella Burioni, and Pietro Rotondo. Local kernel renormalization as a mechanism for feature learning in overparametrized convolutional neural networks. *arXiv preprint arXiv:2307.11807*, 2023.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparametrized neural networks, going beyond two layers. *Advances in neural information processing systems*, 32, 2019.

- 540 Gholamali Aminian, Yixuan He, Gesine Reinert, Łukasz Szpruch, and Samuel N Cohen. General-
541 ization error of graph neural networks in the mean-field regime. *arXiv preprint arXiv:2402.07025*,
542 2024.
- 543 Alexander Atanasov, Blake Bordelon, and Cengiz Pehlevan. Neural networks as kernel learners:
544 The silent alignment effect. *arXiv preprint arXiv:2111.00034*, 2021.
- 545 Alexander Atanasov, Blake Bordelon, Sabarish Sainathan, and Cengiz Pehlevan. The on-
546 set of variance-limited behavior for networks in the lazy and rich regimes. *arXiv preprint*
547 *arXiv:2212.12147*, 2022.
- 548 Yehonatan Avidan, Qianyi Li, and Haim Sompolinsky. Connecting ntk and nngp: A unified the-
549 oretical framework for neural network learning dynamics in the kernel regime. *arXiv preprint*
550 *arXiv:2309.04522*, 2023.
- 551 Yasaman Bahri, Boris Hanin, Antonin Brossollet, Vittorio Erba, Christian Keup, Rosalba Pacelli,
552 and James B. Simon. Les houches lectures on deep learning at large and infinite width, 2024.
- 553 Michael Betancourt. A conceptual introduction to hamiltonian monte carlo. *arXiv preprint*
554 *arXiv:1701.02434*, 2017.
- 555 Blake Bordelon and Cengiz Pehlevan. Self-consistent dynamical field theory of kernel evolution
556 in wide neural networks. *Advances in Neural Information Processing Systems*, 35:32240–32256,
557 2022.
- 558 Rohitash Chandra, Ayush Bhagat, Manavendra Maharana, and Pavel N Krivitsky. Bayesian graph
559 convolutional neural networks via tempered mcmc. *IEEE Access*, 9:130353–130365, 2021.
- 560 Lei Chen, Le Wu, Richang Hong, Kun Zhang, and Meng Wang. Revisiting graph based collaborative
561 filtering: A linear residual graph convolutional network approach. In *Proceedings of the AAAI*
562 *conference on artificial intelligence*, volume 34, pp. 27–34, 2020a.
- 563 Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph
564 convolutional networks. In *International conference on machine learning*, pp. 1725–1735. PMLR,
565 2020b.
- 566 Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph
567 convolutional networks. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th In-*
568 *ternational Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning*
569 *Research*, pp. 1725–1735. PMLR, 13–18 Jul 2020c. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v119/chen20v.html)
570 [press/v119/chen20v.html](https://proceedings.mlr.press/v119/chen20v.html).
- 571 Rong Chen, Li Guanghui, and Chenglong Dai. Feature fusion via deep residual graph convolutional
572 network for hyperspectral image classification. *IEEE Geoscience and Remote Sensing Letters*,
573 19:1–5, 2022.
- 574 Youngmin Cho and Lawrence Saul. Kernel methods for deep learning. *Advances in neural informa-*
575 *tion processing systems*, 22, 2009.
- 576 Weilin Cong, Morteza Ramezani, and Mehrdad Mahdavi. On provable benefits of depth in training
577 graph convolutional networks. *Advances in Neural Information Processing Systems*, 34:9936–
578 9949, 2021.
- 579 Yash Deshpande, Subhabrata Sen, Andrea Montanari, and Elchanan Mossel. Contextual stochastic
580 block models. *Advances in Neural Information Processing Systems*, 31, 2018.
- 581 Simon S Du, Kangcheng Hou, Russ R Salakhutdinov, Barnabas Poczos, Ruosong Wang, and Keyulu
582 Xu. Graph neural tangent kernel: Fusing graph neural networks with graph kernels. *Advances in*
583 *neural information processing systems*, 32, 2019.
- 584 Tianxiang Gao, Xiaokai Huo, Hailiang Liu, and Hongyang Gao. Wide neural networks as gaussian
585 processes: Lessons from deep equilibrium models. *Advances in Neural Information Processing*
586 *Systems*, 36, 2024.

- 594 Vikas Garg, Stefanie Jegelka, and Tommi Jaakkola. Generalization and representational limits
595 of graph neural networks. In *International Conference on Machine Learning*, pp. 3419–3430.
596 PMLR, 2020.
- 597 Floris Geerts and Juan L Reutter. Expressiveness and approximation properties of graph neural
598 networks. *arXiv preprint arXiv:2204.04661*, 2022.
- 600 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
601 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
602 770–778, 2016.
- 603 Jessica N Howard, Ro Jefferson, Anindita Maiti, and Zohar Ringel. Wilsonian renormalization of
604 neural network gaussian processes. *arXiv preprint arXiv:2405.06008*, 2024.
- 606 Wei Huang, Yayong Li, Weitao Du, Jie Yin, Richard Yi Da Xu, Ling Chen, and Miao Zhang.
607 Towards deepening graph neural networks: A gntk-based optimization perspective. *arXiv preprint*
608 *arXiv:2103.03113*, 2021.
- 609 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and gen-
610 eralization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- 612 Haotian Ju, Dongyue Li, Aneesh Sharma, and Hongyang R Zhang. Generalization in graph neural
613 networks: Improved pac-bayesian bounds on graph diffusion. In *International Conference on*
614 *Artificial Intelligence and Statistics*, pp. 6314–6341. PMLR, 2023.
- 615 Dhruva Karkada. The lazy (ntk) and rich (μ) regimes : *agentletutorial*. *arXiv preprint arXiv:2404.19719*, 2024.
- 619 Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.
620 *arXiv preprint arXiv:1609.02907*, 2016.
- 621 Alex Krizhevsky et al. Learning multiple layers of features from tiny images. 2009.
- 623 Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S. Schoenholz, Jeffrey Pennington, and Jascha
624 Sohl-Dickstein. Deep neural networks as gaussian processes. 2018.
- 625 Qianyi Li and Haim Sompolsky. Statistical mechanics of deep linear neural networks: The back-
626 propagating kernel renormalization. *Phys. Rev. X*, 11:031059, Sep 2021. doi: 10.1103/PhysRevX.
627 11.031059. URL <https://link.aps.org/doi/10.1103/PhysRevX.11.031059>.
- 629 Renjie Liao, Raquel Urtasun, and Richard Zemel. A pac-bayesian approach to generalization bounds
630 for graph neural networks. *arXiv preprint arXiv:2012.07690*, 2020.
- 631 Andreas Loukas. What graph neural networks cannot learn: depth vs width. *arXiv preprint*
632 *arXiv:1907.03199*, 2019.
- 634 Andrew K. McCallum, Kamal Nigam, Jason Rennie, et al. Automating the construction of inter-
635 net portals with machine learning. *Information Retrieval*, 3(2):127–163, 2000. doi: 10.1023/A:
636 1009953814988.
- 637 Francesca Mignacco and Pierfrancesco Urbani. The effective noise of stochastic gradient descent.
638 *Journal of Statistical Mechanics: Theory and Experiment*, 2022(8):083405, 2022.
- 640 Andrea Montanari and Eliran Subag. Solving overparametrized systems of random equations: I. model
641 and algorithms for approximate solutions. 2023.
- 642 Gadi Naveh, Oded Ben David, Haim Sompolsky, and Zohar Ringel. Predicting the outputs of finite
643 deep neural networks trained with noisy gradients. *Phys. Rev. E*, 104:064301, Dec 2021. doi: 10.
644 1103/PhysRevE.104.064301. URL [https://link.aps.org/doi/10.1103/PhysRevE.
645 104.064301](https://link.aps.org/doi/10.1103/PhysRevE.104.064301).
- 646 Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business
647 Media, 2012.

- 648 Daniel Park, Jascha Sohl-Dickstein, Quoc Le, and Samuel Smith. The effect of network width on
 649 stochastic gradient descent and generalization: an empirical study. In *International Conference on*
 650 *Machine Learning*, pp. 5042–5051. PMLR, 2019.
- 651 John Stogin, Ankur Mali, and C Lee Giles. A provably stable neural network turing machine. *arXiv*
 652 *preprint arXiv:2006.03651*, 2020.
- 654 Huayi Tang and Yong Liu. Towards understanding generalization of graph neural networks. In An-
 655 dreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan
 656 Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202
 657 of *Proceedings of Machine Learning Research*, pp. 33674–33719. PMLR, 23–29 Jul 2023. URL
 658 <https://proceedings.mlr.press/v202/tang23f.html>.
- 659 Lorenzo Tiberi, Francesca Mignacco, Kazuki Irie, and Haim Sompolinsky. Dissecting the interplay of
 660 attention paths in a statistical mechanics theory of transformers. *arXiv preprint arXiv:2405.15926*,
 661 2024.
- 663 Andreas Veit, Michael J Wilber, and Serge Belongie. Residual networks behave like ensembles of
 664 relatively shallow networks. *Advances in neural information processing systems*, 29, 2016.
- 665 Nikhil Vyas, Alexander Atanasov, Blake Bordelon, Depen Morwani, Sabarish Sainathan, and Cengiz
 666 Pehlevan. Feature-learning networks are consistent across widths at realistic scales. *Advances in*
 667 *Neural Information Processing Systems*, 36, 2024.
- 669 Ian Walker and Ben Glocker. Graph convolutional gaussian processes. In *International Conference on*
 670 *Machine Learning*, pp. 6495–6504. PMLR, 2019.
- 671 Zihan Wang and Arthur Jacot. Implicit bias of sgd in $l_{-}\{2\}$ -regularized linear dnns: One-way jumps
 672 from high to low rank. *arXiv preprint arXiv:2305.16038*, 2023.
- 674 Jingfeng Wu, Wenqing Hu, Haoyi Xiong, Jun Huan, Vladimir Braverman, and Zhanxing Zhu. On the
 675 noisy gradient descent that generalizes as sgd. In *International Conference on Machine Learning*,
 676 pp. 10367–10376. PMLR, 2020.
- 677 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural net-
 678 works? *arXiv preprint arXiv:1810.00826*, 2018.
- 680 Greg Yang, Edward J Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder,
 681 Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large neural networks
 682 via zero-shot hyperparameter transfer. *arXiv preprint arXiv:2203.03466*, 2022.

685 A DETAILS ON THEORY OF BPB-GNN

687 A.1 SUMMARY OF NOTATIONS

688 **Hyperparameters and Dimensions**

689	P	Number of training nodes
690	n	Total number of nodes for a graph
691	N_0	Input node feature dimension
692	N	BPB-GNN hidden layer width
693	L	Total number of branches
694	σ_w	L_2 prior regularization strength
695	T	Temperature
696	$\alpha_0 = \frac{P}{LN_0}$	Network capacity
697	$\alpha = \frac{P}{N}$	Width ratio
698		
699		
700		
701		

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

Network Architecture and Input/Output

$\hat{A} \in \mathbb{R}^{n \times n}$	Adjacency matrix
$A \in \mathbb{R}^{n \times n}$	Normalized adjacency matrix by its degree matrix D
$X \in \mathbb{R}^{n \times N_0}$	Input node feature matrix
$G = (X, A)$	Graph
$W^{(l)} \in \mathbb{R}^{N_0 \times N}$	Hidden layer weight for branch l
$a^{(l)}, a_l \in \mathbb{R}^N$	Readout vector for branch l
$\Theta_l = (W_l, a_l)$	Collection of parameters for branch l
$h_l^\mu \in \mathbb{R}^N$	Activation vector for branch l and node μ
$f_l^\mu \in \mathbb{R}$	Readout prediction for branch l and node μ
$f^\mu \in \mathbb{R}$	Overall readout prediction for node μ
$y^\mu \in \mathbb{R}$	Overall node label for node μ
$H_l(W_l) \in \mathbb{R}^{P \times N}$	Activation feature matrix for branch l
$Y \in \mathbb{R}^P$	Training node labels
$F \in \mathbb{R}^P$	Readout predictions

Statistical Theory

$E(\Theta; G, Y)$	Energy loss function
Z	Partition function
$I \in \mathbb{R}^{P \times P}$	Identity matrix
$H(W)$	Hamiltonian after integrating out readout a_l 's
$u \in \mathbb{R}^L$	Order parameters as saddle point solution
$u_l \in \mathbb{R}$	Order parameter for branch l
$H(u)$	Hamiltonian as a function of order parameters
$S(u)$	Entropy as a function of order parameters
$E(u)$	Energy as a function of order parameters
r_l	Mean squared readout
Tr_l	Variance-related of readout
$\langle x \rangle$	Statistical average of any quantity x over the posterior distribution

Kernels

756	$K(W) \in \mathbb{R}^{P \times P}$	Hidden layer weight dependent overall kernel
757		
758	K_l	Branch l kernel
759	$K \in \mathbb{R}^{P \times P}$	Overall kernel averaged over W 's, $K = \sum_l \frac{w_l K_l}{L}$
760	$k_l^\nu \in \mathbb{R}^{P \times 1}$	Branch l kernel column for the P training nodes against test node ν
761		
762	$k^\nu \in \mathbb{R}^{P \times 1}$	Overall kernel column for the P training nodes against test node ν
763	$ _P$	Kernel restricted to the P training nodes against P training nodes
764		
765	$ _{(P,\nu)}$	Kernel restricted to P training nodes against the test node ν
766		
767	$ _{(\nu,\nu)}$	Kernel restricted to the test node ν against test node ν

Student-teacher Setup

770	$Y_l^* \in \mathbb{R}^P$	Teacher network readout prediction for branch l for P
771		nodes
772	$Y^* \in \mathbb{R}^P$	Teacher network overall readout for P nodes
773		
774	$W_l^* \in \mathbb{R}^{N_0 \times N}$	Teacher hidden layer weight for branch l
775	$a_l^* \in \mathbb{R}^N$	Teacher readout vector for branch l
776		
777	β_l^2	Teacher readout variance
778	σ_t^2	Teacher hidden layer weight variance
779		

A.2 KERNEL RENORMALIZATION

782 Following similar derivations as the first kernel renormalization work Li & Sompolinsky (2021), we
 783 will integrate out the weights in the partition function $Z = \int d\theta \exp(-E(\Theta)/T)$, from the readout
 784 layer weights a_l 's to the hidden layer weights W_l 's and arrive at an effective Hamiltonian shown in
 785 the main text.

786 First, we linearize the energy in terms of a_l 's and transform the integral by introducing the auxiliary
 787 variables $t^\mu, \mu = 1, \dots, P$.

$$788 \quad Z = \int d\Theta \int \prod_{\mu=1}^P dt_\mu \exp \left[-\frac{1}{2\sigma_w^2} \Theta^\top \Theta - \sum_{\mu=1}^P it_\mu \left(\frac{1}{\sqrt{LN}} \sum_{i=1}^N \sum_{l=0}^{L-1} a_i^{(l)} h_i^\mu(G) - Y^\mu \right) - \frac{T}{2} t^\top t \right] \quad (22)$$

792 Now we can integrate out a_l 's as they are linearized and the partition function becomes

$$793 \quad Z = \int DW e^{-H(W)}, \quad (23)$$

794 with effective Hamiltonian

$$795 \quad H(W) = \frac{1}{2\sigma_w^2} \sum_{l=0}^{L-1} \text{Tr} W_l^T W_l + \frac{1}{2} Y^T (K(W) + TI)^{-1} Y + \frac{1}{2} \log \det(K(W) + TI), \quad (24)$$

796 where

$$800 \quad K(W) = \frac{1}{L} \sum_l \frac{\sigma_w^2}{N_0} (H_l(W_l) H_l(W_l)^T) |_P \quad (25)$$

801 is the $P \times P$ kernel matrix dependent on the observed P nodes with node features $H_l = A^l X W_l$
 802 and denote $|_P$ as the matrix restricting to the elements generated by the training nodes.

803 Now we perform the integration on W_l 's and introduce the auxiliary variables t^μ 's again using the
 804 Gaussian trick and get

$$805 \quad Z = \int \prod_{l=0}^{L-1} dW_l \int dt \exp \left[-\frac{1}{2} t^T (K(W) + TI) t - \frac{1}{2\sigma_w^2} \sum_l \text{Tr}(W_l^T W_l) + it^T Y \right] \quad (26)$$

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

$$= \int dt \exp \left[it^T Y + G(t) - \frac{T}{2} t^T t \right], \quad (27)$$

where $G(t)$ is in terms of the kernel averaged over the Gaussian measure

$$G(t) = \log \left\langle \exp \left(-\frac{1}{2N} t^T K(W) t \right) \right\rangle_W \quad (28)$$

Writing out the integral explicitly, we have

$$G(t) = \log \int \prod_{l=0}^{L-1} DW_l \exp \left(-\frac{1}{2N} \sum_{\mu, \nu} t^\mu t^\nu \sum_{j,l} \frac{\sigma_w^2}{N_0 L} \left(\sum_{\mu', i} A_{\mu, \mu'}^l W_{ij}^l x_i^{\mu'} \right) \left(\sum_{\nu', i'} A_{\nu, \nu'}^l W_{i'j}^l x_{i'}^{\nu'} \right) \right) \quad (29)$$

where DW_l is the Gaussian measure

$$DW_l = \prod_{i,j} e^{-W_{i,j}^2 / 2\sigma_w^2} dW_{l,i,j} \quad (30)$$

By recognizing that this is just a product of the same N integrals on N_0 dimensional weight vector \vec{W}_l^j for a fixed j neuron, we get

$$G(t) = N \log \int \prod_{l=0}^{L-1} d\vec{W}_l^j \exp \left(\sum_{l,i,i'} \left(-\frac{1}{2N} \sum_{\mu, \nu} t^\mu t^\nu \frac{\sigma_w^2}{N_0 L} \xi_{l,i}^\mu \xi_{l,i'}^\nu + \frac{1}{2\sigma_w^2} \delta_{i,i'} \right) W_{i,l}^j W_{i',l}^j \right) \quad (31)$$

where $\xi_{l,i}^\mu = \sum_{\mu'} A_{\mu, \mu'}^l x_i^{\mu'}$ is the l th branch convolution on node μ . Since the network is linear, we are able to perform the Gaussian integration exactly and get

$$G(t) = -\frac{N}{2} \sum_l \log(\det(I_{i,i'} + \frac{\sigma_w^2}{N_0 L} t^\mu \xi_{l,i}^\mu \xi_{l,i'}^\nu)) \quad (32)$$

where the determinant is over the i, i' matrix, and thus performing the determinant we arrive at

$$G(t) = -\frac{N}{2} \sum_l \log(1 + t^T \frac{K_l}{L} t), \quad (33)$$

where $K_{\mu, \nu}^l = \frac{\sigma_w^2}{N} \sum_i t^\mu \xi_{l,i}^\mu \xi_{l,i}^\nu$ is the Gaussian process kernel for branch l .

Next, we insert L delta functions $\int \prod_l dh_l \delta(h_l - t^T K_l t)$ to enforce the identity $G(t) = \sum_l -\frac{N}{2} \log(1 + h_l)$ and use the Fourier representation of it with auxiliary variables u_l to get:

$$\begin{aligned} Z &= \int \prod_{l=0}^{L-1} dh_l du_l dt \exp \left(it^T Y - \sum_l \frac{N}{2} \log(1 + h_l) + \sum_l \frac{N}{2\sigma_w^2} u_l h_l - \frac{1}{2} t^T \left(\sum_l \frac{1}{L} u_l K_l + TI \right) t \right) \\ &= \int \prod_{l=0}^{L-1} dh_l du_l \exp \left(-\sum_l \frac{N}{2} \log(1 + h_l) + \sum_l \frac{N}{2\sigma_w^2} u_l h_l \frac{1}{2} Y^T \left(\sum_l \frac{1}{L} u_l K_l + TI \right)^{-1} Y \right) \end{aligned} \quad (34)$$

where

$$K_l = \frac{\sigma_w^2}{N_0} [A^l X X^T A^{lT}]_P \quad (35)$$

is the input kernel for branch l . Now as $N \rightarrow \infty$ and $\alpha = \frac{P}{N}$ fixed, we can perform the saddle point approximation and get the saddle points for h_l as

$$1 + h_l = \frac{\sigma_w^2}{u_l} \quad (36)$$

864 Plugging this back to the equation, we get

$$865 \quad Z = \int \prod_l du_l e^{-H_{eff}(u)}, \quad (37)$$

866 with the effective Hamiltonian

$$867 \quad H_{eff}(u) = S(u) + E(u), \quad (38)$$

868 where we call $S(u)$ the entropic term

$$869 \quad S(u) = - \sum_l \frac{N}{2} \log u_l + \sum_l \frac{N}{2\sigma_w^2} u_l \quad (39)$$

870 and $E(u)$ the energetic term

$$871 \quad E(u) = \frac{1}{2} Y^T \left(\sum_l \frac{1}{L} u_l K_l + TI \right)^{-1} Y + \frac{1}{2} \log \det \left(\sum_l \frac{1}{L} u_l K_l + TI \right) \quad (40)$$

872 Therefore, after integrating out W_l , the effective kernel is given by

$$873 \quad K = \sum_l \frac{1}{L} u_l K_l, \quad (41)$$

874 and the saddle point equations for u_l 's are determined by

$$875 \quad N \left(1 - \frac{u_l}{\sigma_w^2} \right) = -Y^T (K + TI)^{-1} \frac{u_l K_l}{L} (K + TI)^{-1} Y + \text{Tr} \left[K^{-1} \frac{u_l K_l}{L} \right], \quad (42)$$

876 where we call

$$877 \quad r_l = Y^T (K + TI)^{-1} \frac{u_l K_l}{L} (K + TI)^{-1} Y \quad (43)$$

878 and

$$879 \quad \text{Tr}_l = \text{Tr} \left[K^{-1} \frac{u_l K_l}{L} \right] \quad (44)$$

880 As we will show later, these represent the mean and variance of the readout norm squared respectively. In the $T = 0$ case, the saddle point equation becomes

$$881 \quad N \left(1 - \frac{u_l}{\sigma_w^2} \right) = -Y^T K^{-1} \frac{u_l K_l}{L} K^{-1} Y + \text{Tr} \left[K^{-1} \frac{u_l K_l}{L} \right] \quad (45)$$

882 A.3 PREDICTOR STATISTICS AND GENERALIZATION

883 We can get the predictor statistics of each branch readout $y_l^\nu(G)$ on a new test node ν by considering the generating function:

$$884 \quad Z(\eta_1, \dots, \eta_L) = \int D\Theta \exp \left\{ -\frac{\beta}{2} \sum_\mu (f^\mu(G; \Theta) - y^\mu)^2 \right. \\ 885 \quad \left. + \sum_l i\eta_l \frac{1}{\sqrt{NL}} \sum_i a_i^{(l)} h_i^{(l),\nu}(G, W_l) - \frac{T}{2\sigma_w^2} \Theta^T \Theta \right\} \quad (46)$$

886 Therefore, by taking the derivative with respect to each η_l , we arrive at the statistics for $y_l(x)$ as:

$$887 \quad \langle f_l^\nu(G) \rangle = \partial_{i\eta_l} \log Z \Big|_{\vec{\eta}=0} \quad (47)$$

$$888 \quad \langle \delta f_{l,\nu}^2(G) \rangle = \partial_{i\eta_l}^2 \log Z \Big|_{\vec{\eta}=0} \quad (48)$$

889 After integrating out the weights Θ layer by layer, we have:

918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934
 935
 936
 937
 938
 939
 940
 941
 942
 943
 944
 945
 946
 947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971

$$\begin{aligned}
 Z(\eta_1, \dots, \eta_L) = & \int \Pi_i du_i \exp \left\{ \sum_l \left(\frac{N}{2} \log u_l - \frac{N}{2\sigma_w^2} u_l \right) \right. \\
 & + \frac{1}{2} (iY + \sum_l \frac{1}{L} \eta_l u_l k_l^\nu)^T (\sum_l \frac{1}{L} u_l K_l + TI)^{-1} (iY + \sum_l \eta_l \frac{1}{L} u_l k_l^\nu \\
 & \left. - \frac{1}{2} \log \det (\sum_l \frac{1}{L} u_l K_l + TI) - \frac{1}{2} \sum_l \eta_l^2 \frac{1}{L} K_l^{\nu, \nu} \right\}. \quad (49)
 \end{aligned}$$

Here

$$k_l^\nu = \frac{\sigma_w^2}{N_0} [A^l X X^T A^l]_{(P, \nu)} \quad (50)$$

is the $P \times 1$ column kernel matrix for test node ν and all training nodes, and

$$K_l^{\nu, \nu} = \frac{\sigma_w^2}{N_0} [A^l X X^T A^l]_{(\nu, \nu)} \quad (51)$$

is the single matrix element for the test node. Therefore, eventually, we have:

$$\langle f_l^\nu \rangle = \frac{u_l k_{l, \nu}^T}{L} (K + TI)^{-1} Y \quad (52)$$

and

$$\langle \delta f_{l, \nu}^2 \rangle = \frac{u_l K_l^{\nu, \nu}}{L} - \frac{u_l k_{l, \nu}^T}{L} (K + TI)^{-1} \frac{u_l k_{l, \nu}}{L} \quad (53)$$

The predictor statistics of the overall readout $f = \sum_l f_l$ is given by:

$$\langle f^\nu(G) \rangle = \sum_l \frac{u_l k_{l, \nu}^T}{L} (K + TI)^{-1} Y = k_\nu^T (K + TI)^{-1} Y \quad (54)$$

$$\langle \delta f(G)_\nu^2 \rangle = \sum_l u_l K_l^{\nu, \nu} - \sum_{l, l'} u_l k_{l, \nu}^T (K + TI)^{-1} u_{l'} k_{l', \nu} = K_{\nu, \nu} - k_\nu^T (K + TI)^{-1} k_\nu \quad (55)$$

A.4 STATISTICS OF BRANCH READOUT NORMS

From the partition function Eq.22, we can relate the mean of readout weights a_l to the auxiliary variable t by

$$\langle a_l \rangle_W = -i \frac{\sigma_w^2}{\sqrt{N}} \Phi_l^T \langle t \rangle = -\frac{\sigma_w^2}{\sqrt{NL}} \Phi_l^T (K + TI)^{-1} Y, \quad (56)$$

where Φ_l is the node feature matrix for the hidden layer nodes. We have

$$\langle a_l^T \rangle \langle a_l \rangle = \sigma_w^2 Y^T (K + TI)^{-1} \frac{u_l K_l}{L} (K + TI)^{-1} Y = r_l \sigma_w^2 \quad (57)$$

We can calculate the second-order statistics of a_l : the variance is

$$\langle \delta a_l^T \delta a_l \rangle = \sigma_w^2 \text{Tr} (I + \frac{\sigma_w^2 \beta}{NL} \Phi_l \Phi_l^T)^{-1} = \sigma_w^2 (N - \text{Tr} (K + TI)^{-1} \frac{u_l K_l}{L}) = \sigma_w^2 (N - \text{Tr}_l) \quad (58)$$

Therefore,

$$\langle a_l^2 \rangle = \langle \delta a_l^T \delta a_l \rangle + \langle \delta a_l^T \delta a_l \rangle = N \sigma_w^2 + \sigma_w^2 r_l + 1 - \sigma_w^2 \text{Tr}_l = N u_l \quad (59)$$

Therefore, we have proved the main text claim that the order parameter u_l 's are really the mean squared readout norms of the branches.

B THE NARROW WIDTH LIMIT FOR GENERAL ARCHITECTURES

In this section, we demonstrate that our results on the narrow-width limit can be generalized to any network with branching structures, supported by both theoretical and empirical evidence.

Consider a general branching network, with L independent branches each with hidden layer weight W_l and readout weight a_l , ie.

$$f(x; \Theta) = \sum_{l=0}^{L-1} \frac{1}{\sqrt{L}} f_l(x; \Theta_l) = \sum_l \frac{1}{\sqrt{NL}} a_l \phi_l(x; W_l), \quad (60)$$

where $\phi_l(x; W_l)$ represents any feed-forward function on the dataset X , ie. in matrix form the branch l readout is $F_l = \phi_l(\frac{1}{\sqrt{N_0}} X W_l)$. As an example, $\phi_l(X W_l) = \frac{1}{\sqrt{N_0}} A^l X W_l$ represents l convolutions on the node features in the BPB-GNN architecture. $\phi_l(X W_l) = \frac{1}{\sqrt{N_0}} A_l X W_l$ can also represent an attention head l for attention-networks with frozen attention A_l , as well as a kernel convolution for patch l in the CNN architecture. In the residual-MLP example below, ϕ_l represents either the linear or the ReLU branch.

With the general branching network, we can similarly consider the posterior distribution in the Bayesian setup as Eq. 4 with squared error as the likelihood term and L_2 regularization as the prior term. Provided the network is *linear*, the partition function Eq.5 can be integrated in a similar fashion and arrive at the saddle point equation

$$N(1 - \frac{u_l}{\sigma_w^2}) = -Y^T (K + TI)^{-1} \frac{u_l K_l}{L} (K + TI)^{-1} Y + \text{Tr}[K^{-1} \frac{u_l K_l}{L}], \quad (61)$$

with $K = \sum_l u_l K_l / L$ and K_l corresponding to the NNGP kernel (Lee et al. (2018)). In fact, even with non-linear activations, we conjecture that the kernel renormalization still holds with the order parameters u_l that satisfies the saddle point equation, and the posterior weight statistics calculated with the renormalized kernel agrees with the experiments as empirically shown in the residual-MLP section.

B.1 THE GENERALIZED EQUIPARTITION THEOREM

With the general branching network, we show the narrow width limit result for the student-teacher setup. Specifically we prove the following

The generalized equipartition theorem : For the general branching network Eq. 60, where the teacher network width is N_t with $W_{l,ij}^* \sim \mathcal{N}(0, \sigma_t^2)$ and $a_{i,l}^* \sim \mathcal{N}(0, \beta_l^2)$ and the student network with layer width N in the Bayesian regression setup of Eq. 4, at $T = 0$ temperature, in the limit $N_t \rightarrow \infty$ and $\alpha = P/N \rightarrow \infty$, there exists a solution to the saddle point equations 12

$$u_l = \sigma_t^2 \beta_l^2 / \sigma_w^2. \quad (62)$$

Proof: The proof is exactly the same as before, except for changing K_l as a general NNGP kernel that corresponds to the feature map ϕ_l for branch l .

B.2 RESIDUAL-MLP

In the following, we show a simple 2-layer residual network as an example for the general branching networks, with a linear branch $\phi_0(x) = \frac{1}{\sqrt{N_0}} W_0 x$ and a ReLU branch $\phi_1(x) = \frac{1}{\sqrt{N_0}} \text{ReLU}(W_1 x)$. Therefore, the overall kernel is $K = \frac{1}{2}(u_0 K_0 + u_1 K_1)$, with the NNGP kernels K_l given by

$$K_0 = K_{linear} = \frac{\sigma^2}{N_0} X_0 X_0^T \quad (63)$$

and the relu kernel (Cho & Saul (2009))

$$K_1(x, x') = K_{relu}(x, x') = \frac{1}{2\pi} \sqrt{K_0(x, x) K_0(x', x')} J(\theta), \quad (64)$$

where

$$J(\theta) = \sin(\theta) + (\pi - \theta) \cos(\theta) \quad (65)$$

and

$$\theta = \cos^{-1}\left(\frac{K_0(x, x')}{\sqrt{K_0(x, x)K_0(x', x')}}\right) \quad (66)$$

For the student-teacher task (details in Appendix C.3), we observe the same narrow-width limit phenomenon as in the BPB-GNN case, where each student branch robustly learns the teacher’s corresponding branch but approaches the same GP limit as the network width gets larger.

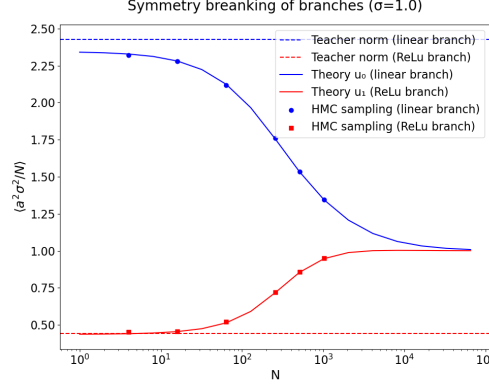


Figure 9: Statistical average of student readout norms squared as a function of network width from theory and HMC sampling, for the residual-MLP student-teacher task, with fixed $\sigma_w = 1$.

As shown in Figure 9 and 10, there is a symmetry breaking of branches from the GP-limit to the narrow-width limit, with the ReLU branch and the linear branch learning their corresponding teacher branch norms as N decreases. Furthermore, narrower network performs better in terms of the bias

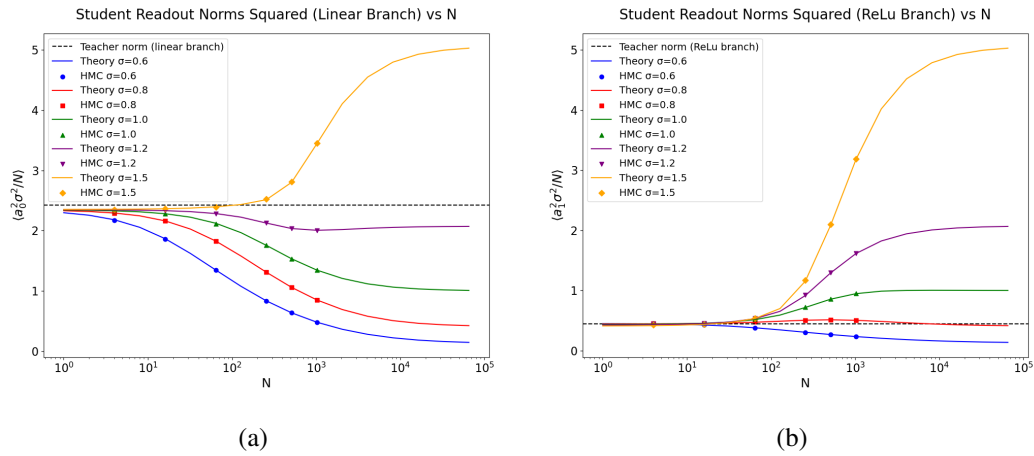


Figure 10: $\langle \|a_l\|^2 \rangle \sigma_w^2 / N$ as a function of network width N for a range of σ_w values. (a) Linear branch; (b) Relu branch. The student branch norms with different regularization strengths all converge to the same teacher readout norm values at narrow width.

and larger σ_w corresponds to a shorter transition window from the GP-limit to the narrow-width limit, similar to the BPB-GNN case. As shown in Figure 11 and 12, the trade-off between the bias and variance contribution is more pronounced for small σ_w ’s, usually with an optimal width for the generalization error. These results support our generalized equipartition theorem and demonstrates that the narrow width limit is a general phenomenon for branching networks.

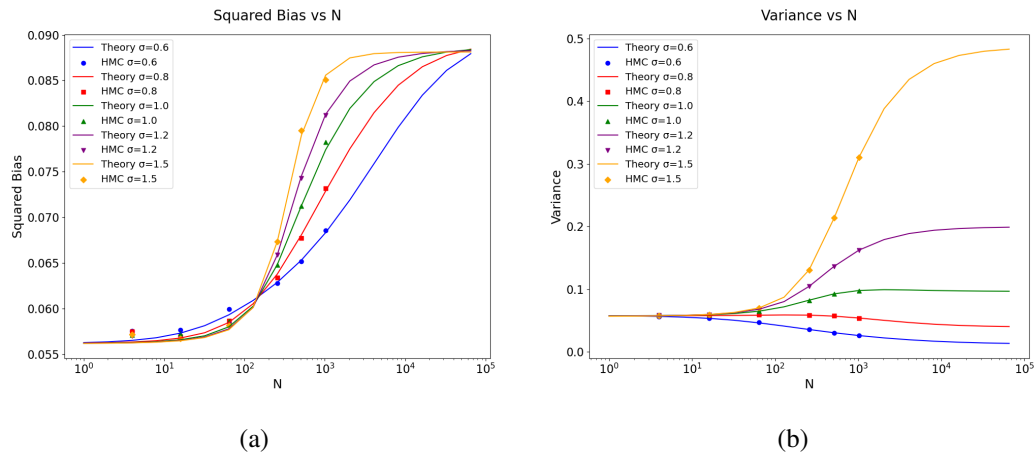


Figure 11: Student residual-MLP network generalization performance as a function of network width N and regularization strength σ_w . (a) Bias; (b) Variance.

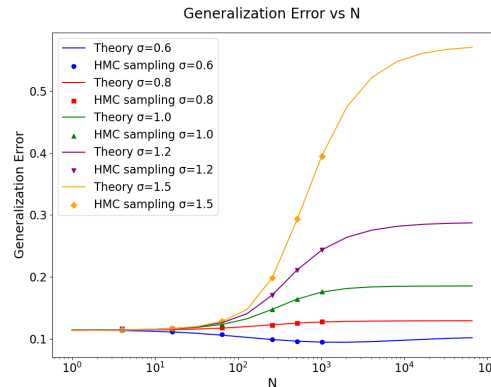


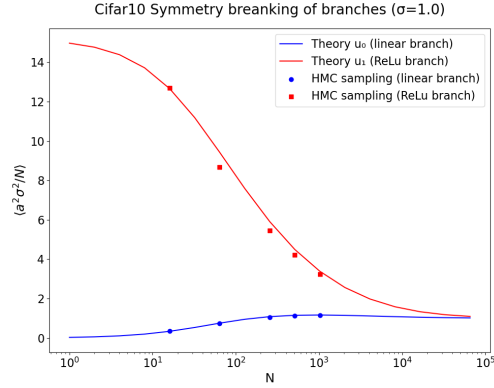
Figure 12: Overall generalization as a sum of bias and variance in Figure 11.

B.3 RESIDUAL-MLP ON CIFAR10

We also perform theory and experiments on the Cifar10 (Krizhevsky et al. (2009)) dataset for a binary classification task (translated to Bayesian regression), with details elaborated in Appendix C.4. Similarly, there exists symmetry breaking of branches and the branch norms converge at the narrow width for different σ_w 's, as shown in Figure 13,14. Interestingly, the ReLu branch completely dominates at the narrow width, which means the ReLu branch is much more important for the residual-MLP when the training sample size is much bigger than the network width.

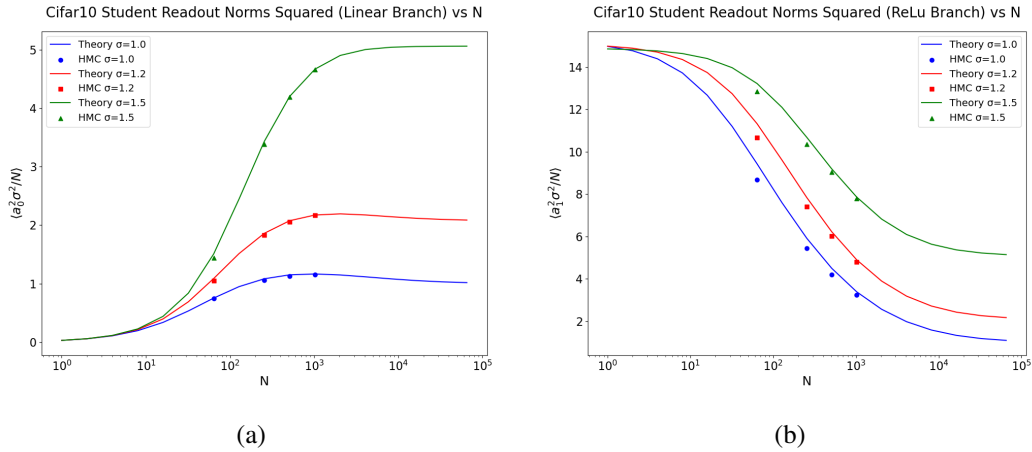
Furthermore, we also show in Figure 15,16(a) that the dominance of ReLu branch at narrow width drives the bias to decrease, albeit overall generalization is dominated by the variance. Using the average of the label predictors over different experimental samples, we show that the test accuracy indeed decreases with larger network width N in Figure 16(b).

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147



1148 Figure 13: Statistical average of student readout norms squared as a function of network width from
1149 theory and HMC sampling, for the Cifar10 dataset with residual-MLP, with fixed $\sigma_w = 1$.

1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165



1166 Figure 14: Cifar10 $\langle \|a_l\|^2 \rangle \sigma_w^2 / N$ as a function of network width N for a range of σ_w values. (a)
1167 Linear branch; (b) Relu branch. The student branch norms with different regularization strengths all
1168 converge to the same teacher readout norm values at narrow width.

1169

1170 C EXPERIMENTAL DETAILS

1171

1172 C.1 STUDENT-TEACHER CSBM

1173

1174 For the student-teacher task, we use the contextual stochastic block model introduced by Deshpande
1175 et al. (2018) to generate the graph G . The adjacency matrix is given by

1176
1177
1178
1179

$$A_{ij} = \begin{cases} 1 & \text{with probability } p = c_{in}/n, \text{ if } i, j \leq n/2 \\ 1 & \text{with probability } p = c_{in}/n, \text{ if } i, j \geq n/2 \\ 1 & \text{with probability } q = c_{out}/n, \text{ otherwise} \end{cases} \quad (67)$$

1180 where

1181

$$c_{in,out} = d \pm \sqrt{d}\lambda \quad (68)$$

1182

d is the average degree and λ the homophily factor.

1183

The feature vector \vec{x}^μ for a particular node μ is given by

1184

$$\vec{x}_\mu = \sqrt{\frac{\mu}{n}} y^\mu \vec{u} + \vec{\xi}_\mu, \quad (69)$$

1185

1186

where

1187

$$\vec{u} \sim \mathcal{N}(0, I_{N_0}), \vec{\xi}_\mu \sim \mathcal{N}(0, I_{N_0}) \quad (70)$$

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

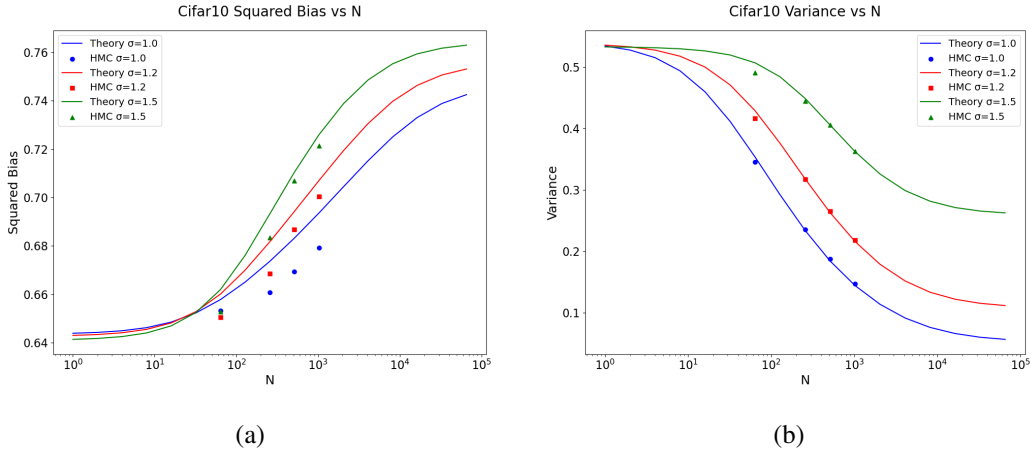


Figure 15: Residual-MLP generalization performance on Cifar10 as a function of network width N and regularization strength σ_w . (a) Bias; (b) Variance.

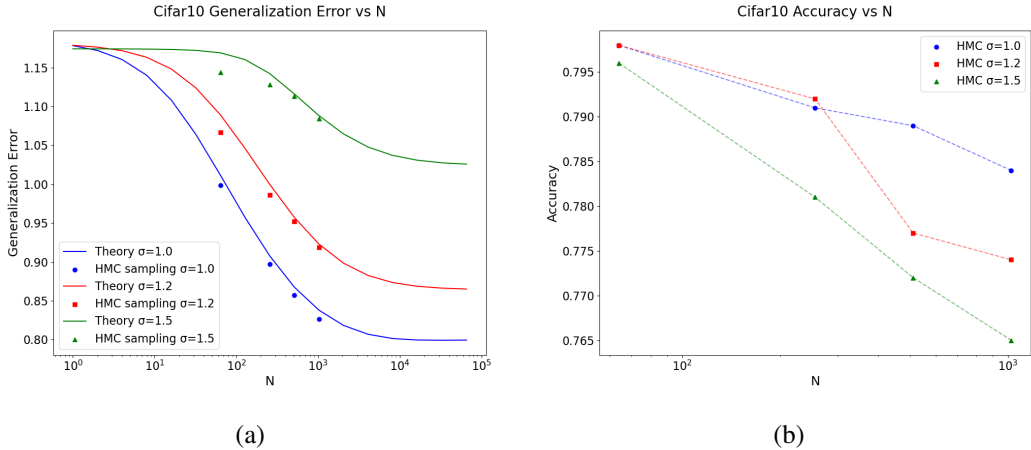


Figure 16: (a) Overall generalization error for residual-MLP on Cifar10. (b) Test accuracy on Cifar10, calculated using the class label by averaging the predictions from HMC samples.

In the experiment, we use $N_0 = 950, d = 20, \lambda = 4$ and $\mu = 4$. The teacher network parameters are variance $\sigma_t^2 = 1$, width $N_t = 1024$, branch norms variance $\beta_0^2 = 0.4, \beta_1^2 = 2$ for individual element of the readout vector a_l . We used total number of nodes $n = 2600$ with train ratio 0.65 for semi-supervised node regression. Temperature $T = 0.0005\sigma_w^2$ for each σ_w value in $\sigma_w = \{0.5, 0.8, 1.0, 1.2\}$.

C.2 CORA

For the Cora dataset, we use a random split of the data (total nodes $n = 2708$) into 21% as training set and 79% as test set. We group the classes (1, 2, 4) into one group and the rest for the other group for binary node regression, with labels as ± 1 's. The Bayesian theory and HMC sampling follows the same design as in the student-teacher setup. We use temperature $T = 0.01$ for both theory and sampling as the sampling becomes more difficult for smaller temperature. This explains the discrepancy of the GP limit bias for different σ_w values in $\sigma_w = \{1.0, 1.2, 1.4\}$.

1242 C.3 STUDENT-TEACHER RESIDUAL-MLP

1243
1244 For the student-teacher task of the residual MLP, we use the data matrix X generated from i.i.d
1245 unit Gaussian distribution, with input dimension $N_0 = 1024$, training sample size $P = 1280$ and
1246 test sample size $P_{test} = 320$. The teacher network weights are also generated with i.i.d Gaussian,
1247 with the hidden layer width $N_t = 1024$, hidden layer weight variance $\sigma_w^2 = 1$, and readout branch
1248 variance $\beta_l^2 = (2.4, 0.4)$ for the two branches. The temperature is set as $T = 0.001\sigma_w^2$ for a range
1249 of σ_w values in $\sigma_w = \{0.6, 0.8, 1.0, 1.2, 1.5\}$.

1250 C.4 CIFAR10

1251
1252 For the Cifar10 dataset, we filter the category of airplane and that of cat for the regression task.
1253 We use a random split of training sample size $P = 1000$ and test sample size $P_{test} = 1000$. The
1254 input data is only from the red channel with input dimension $N_0 = 1024$. The temperature is set
1255 at $T = 0.005\sigma_w^2$ for a range of σ_w values in $\sigma_w = \{1.0, 1.2, 1.5\}$. Within the compute budget, we
1256 observe significant sample correlations, which has a large impact on measuring the bias term. This
1257 is why we observe relatively large deviation of the bias results from theory in Figure 15(a).

1258 C.5 HAMILTONIAN MONTE CARLO

1259
1260 The sampling experiments in the paper are all done with Hamiltonian Monte Carlo (Betancourt
1261 (2017); Neal (2012); Chandra et al. (2021)) simulations, a popular method for sampling a proba-
1262 bility distribution. HMC has faster convergence to the posterior distribution compared to Langevin
1263 dynamics. We used Numpyrho to set up chains and run the simulations on the GPU cluster. All
1264 experiments are within 20 GPU hours budget on Nvidia A100-40G. Due to memory constraint, we
1265 only sampled up to $N = 1024$ hidden layer width for the student-teacher CSBM experiment and
1266 $N = 64$ for the Cora experiment. Since we mainly aim to demonstrate the narrow width effect in
1267 this paper, this suffices the purpose.

1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295