
Understanding Self-supervised Learning with Dual Deep Networks

Abstract

1 We propose a novel theoretical framework to understand self-supervised learning
2 methods that employ dual pairs of deep ReLU networks (e.g., SimCLR, BYOL).
3 First, we prove that in each SGD update of SimCLR, the weights at each layer are
4 updated by a *covariance operator* that specifically amplifies initial random selec-
5 tivities that vary across data samples but survive averages over data augmentations,
6 which we show leads to the emergence of hierarchical features, if the input data
7 are generated from a hierarchical latent tree model. With the same framework,
8 we also show analytically that BYOL works due to an implicit contrastive term,
9 acting as an approximate covariance operator. The term is formed by the inter-
10 play between the zero-mean operation of BatchNorm and the extra predictor in
11 the online network. Extensive ablation studies justify our theoretical findings.

12 1 Introduction

13 While self-supervised learning (SSL) has achieved great empirical success across multiple domains,
14 including computer vision [17, 15, 8, 16, 25, 7], natural language processing [12], and speech recog-
15 nition [34, 4, 5], its theoretical understanding remains elusive, especially when multi-layer nonlinear
16 deep networks are involved. Unlike supervised learning (SL) that deals with labeled data, SSL learns
17 meaningful structures from randomly initialized networks without human-provided labels.

18 In this paper, we propose a systematic theoretical analysis of SSL with deep ReLU networks. Our
19 analysis imposes no parametric assumptions on the input data distribution and is applicable to state-
20 of-the-art SSL methods that typically involve two parallel (or *dual*) deep ReLU networks during
21 training (e.g., SimCLR [8], BYOL [16], etc). We do so by developing an analogy between SSL
22 and a theoretical framework for analyzing supervised learning, namely the student-teacher setting
23 [31, 1, 21, 28], which also employs a pair of dual networks. Our results indicate that SimCLR
24 weight updates at every layer are amplified by a fundamental positive semi definite (PSD) *covari-*
25 *ance operator* that only captures feature variability across data points that *survive* averages over
26 data augmentation procedures designed in practice to scramble semantically unimportant features
27 (e.g. random image crops, blurring or color distortions [13, 20, 25, 27]). This covariance operator
28 provides a principled framework to study how SimCLR amplifies initial random selectivity to obtain
29 distinctive features that vary *across* samples after surviving averages over data-augmentations.

30 Based on the covariance operator, we further show that (1) in a two-layer setting, a top-level covari-
31 ance operator helps accelerate the learning of low-level features, and (2) when the data are generated
32 by a hierarchical latent tree model, training deep ReLU networks leads to an emergence of the la-
33 tent variables in its intermediate layers. We also explain why BYOL works *without negative pairs*,
34 by showing analytically that an interplay between the zero-mean operation in BatchNorm and the
35 extra predictor in the online network creates an implicit contrastive term, thereby explaining em-
36 pirical observations in the recent blog [14]. We also discover that reinitializing the predictor every
37 a few epochs doesn’t hurt the performance measured by linear evaluation protocol, questioning the
38 hypothesis of “optimal predictor” in a recent version (v3) of BYOL [16].

39 To the best of our knowledge, we are the first to provide a systematic theoretical analysis of modern
40 SSL methods with deep ReLU networks that also incorporates data augmentation, and show how
41 the gradient-based SSL work *inside* the neural network. See Appendix I for more related works.

42 2 Overall framework

43 **Notation.** Consider an L -layer ReLU network obeying $\mathbf{f}_l = \psi(\tilde{\mathbf{f}}_l)$ and $\tilde{\mathbf{f}}_l = W_l \mathbf{f}_{l-1}$ for $l =$
44 $1, \dots, L$. Here $\tilde{\mathbf{f}}_l$ and \mathbf{f}_l are n_l dimensional pre-activation and activation vectors in layer l , with

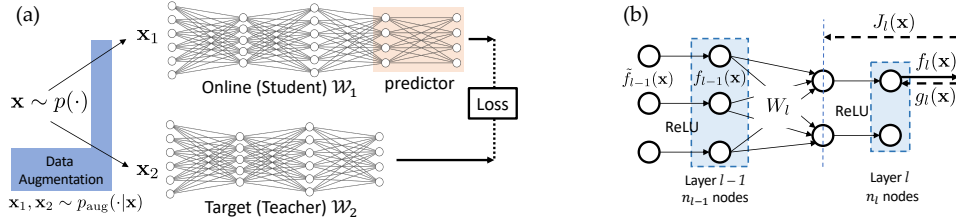


Figure 1: **(a)** Overview of the two SSL algorithms we study in this paper: SimCLR ($\mathcal{W}_1 = \mathcal{W}_2 = \mathcal{W}$, no predictor, NCE Loss) and BYOL (\mathcal{W}_1 has an extra predictor, \mathcal{W}_2 is a moving average), **(b)** Detailed Notation.

45 $f_0 = x$ being the input and $f_L = \tilde{f}_L$ the output (no ReLU at the top layer). $W_l \in \mathbb{R}^{n_l \times n_{l-1}}$
 46 are the weight matrices, and $\psi(u) := \max(u, 0)$ is the element-wise ReLU nonlinearity. We let
 47 $\mathcal{W} := \{W_l\}_{l=1}^L$ be all network weights. We also denote the gradient of any loss function with respect
 48 to f_l by $g_l \in \mathbb{R}^{n_l}$, and the derivative of the output f_L with respect to an earlier pre-activation f_l by
 49 the Jacobian matrix $J_l(x; \mathcal{W}) \in \mathbb{R}^{n_L \times n_l}$, as both play key roles in backpropagation (Fig. 1(b)).

50 **An analogy between self-supervised and supervised learning: the dual network scenario.**
 51 Many recent successful approaches to self-supervised learning (SSL), including SimCLR [8],
 52 BYOL [16] and MoCo [17], employ a dual “Siamese-like” pair [19] of such networks (Fig. 1(b)).
 53 Each network has its own set of weights \mathcal{W}_1 and \mathcal{W}_2 , receives respective inputs x_1 and x_2 and
 54 generates outputs $f_{1,L}(x_1; \mathcal{W}_1)$ and $f_{2,L}(x_2; \mathcal{W}_2)$. The pair of inputs $\{x_1, x_2\}$ can be either posi-
 55 tive or negative, depending on how they are sampled. For a positive pair, a *single* data point x is
 56 drawn from the data distribution $p(\cdot)$, and then two augmented views x_1 and x_2 are drawn from
 57 a conditional augmentation distribution $p_{\text{aug}}(\cdot|x)$. Possible image augmentations include random
 58 crops, blurs or color distortions, that ideally preserve semantic content useful for downstream tasks.
 59 In contrast, for a negative pair, two *different* data points $x, x' \sim p(\cdot)$ are sampled, and then each
 60 are augmented independently to generate $x_1 \sim p_{\text{aug}}(\cdot|x)$ and $x_2 \sim p_{\text{aug}}(\cdot|x')$. For SimCLR, the
 61 dual networks have tied weights with $\mathcal{W}_1 = \mathcal{W}_2$, and a loss function is chosen to encourage the
 62 representation of positive (negative) pairs to become similar (dissimilar). In BYOL, only positive
 63 pairs are used, and the first network \mathcal{W}_1 , called the online network, is trained to match the output of
 64 the second network \mathcal{W}_2 (the target), using an additional layer named *predictor*. The target network
 65 ideally provides training targets that can improve the online network’s representation and does not
 66 contribute gradient. The improved online network is gradually incorporated into the target network,
 67 yielding a bootstrapping procedure.

68 Our fundamental goal is to analyze the mechanisms governing how SSL methods like SimCLR
 69 and BYOL lead to the emergence of meaningful intermediate features, starting from random initial-
 70 izations, and how these features depend on the data distribution $p(x)$ and augmentation procedure
 71 $p_{\text{aug}}(\cdot|x)$. Interestingly, the analysis of *supervised* learning (SL) often employs a similar dual net-
 72 work scenario, called *teacher-student setting* [31, 1, 21, 28], where \mathcal{W}_2 are the ground truth weights
 73 of a *fixed* teacher network, which generates outputs in response to random inputs. These input-output
 74 pairs constitute training data for the first network, which is a student network. Only the student net-
 75 work’s weights \mathcal{W}_1 are trained to match the target outputs provided by the teacher. This yields an
 76 interesting mathematical parallel between SL, in which the teacher is fixed and only the student
 77 evolves, and SSL, in which both the teacher and student evolve with potentially different dynamics.
 78 This mathematical parallel opens the door to using techniques from SL (e.g., [31]) to analyze SSL.

79 **Gradient of ℓ_2 loss for dual deep ReLU networks.** As seen above, the (dis)similarity of repre-
 80 sentations between a pair of dual networks plays a key role in both SSL and SL. We thus consider
 81 minimizing a simple measure of dissimilarity, the squared ℓ_2 distance $r := \frac{1}{2} \|f_{1,L} - f_{2,L}\|^2$ be-
 82 tween the final outputs $f_{1,L}$ and $f_{2,L}$ of two multi-layer ReLU networks with weights \mathcal{W}_1 and \mathcal{W}_2
 83 and inputs x_1 and x_2 . Without loss of generality, we only analyze the gradient w.r.t \mathcal{W}_1 . For each
 84 layer l , we first define *connection* $K_l(x)$, a quantity that connects the bottom-up feature vector f_{l-1}
 85 with the top-down Jacobian J_l , which both contribute to the gradient at weight layer l .

86 **Definition 1** (The connection $K_l(x)$). The connection $K_l(x; \mathcal{W}) := f_{l-1}(x; \mathcal{W}) \otimes J_l^T(x; \mathcal{W}) \in$
 87 $\mathbb{R}^{n_l n_{l-1} \times n_L}$. Here \otimes is the Kronecker product.

88 **Theorem 1** (Squared ℓ_2 Gradient for dual deep ReLU networks). The gradient g_{W_l} of r w.r.t. $W_l \in$
 89 $\mathbb{R}^{n_l \times n_{l-1}}$ for a single input pair $\{x_1, x_2\}$ is (here $K_{1,l} := K_l(x_1; \mathcal{W}_1)$ and $K_{2,l} := K_l(x_2; \mathcal{W}_2)$):

$$g_{W_l} = \text{vec}(\partial r / \partial W_{1,l}) = K_{1,l} \left[K_{1,l}^T \text{vec}(W_{1,l}) - K_{2,l}^T \text{vec}(W_{2,l}) \right]. \quad (1)$$

90 See Appendix for proofs of all theorems in main text.

3 Analysis of SimCLR

As discussed above, SimCLR [8] employs both positive and negative input pairs, and a symmetric network structure with $\mathcal{W}_1 = \mathcal{W}_2 = \mathcal{W}$. Let $\{\mathbf{x}_1, \mathbf{x}_+\}$ be a positive input pair, and let $\{\mathbf{x}_1, \mathbf{x}_{k-}\}$ for $k = 1, \dots, H$ be H negative pairs. These input pairs induce corresponding squared ℓ_2 distances in output space, $r_+ := \frac{1}{2} \|\mathbf{f}_{1,L} - \mathbf{f}_{+,L}\|_2^2$, and $r_{k-} := \frac{1}{2} \|\mathbf{f}_{1,L} - \mathbf{f}_{k-,L}\|_2^2$. The InfoNCE loss [26] with temperature τ then minimizes (maximizes) the positive (negative) pair distances:

$$L(r_+, r_{1-}, r_{2-}, \dots, r_{H-}) := -\log \frac{e^{-r_+/\tau}}{e^{-r_+/\tau} + \sum_{k=1}^H e^{-r_{k-}/\tau}} \quad (2)$$

When $\|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1$, we have $-\frac{1}{2}\|\mathbf{u} - \mathbf{v}\|_2^2 = \text{sim}(\mathbf{u}, \mathbf{v}) - 1$ where $\text{sim}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{u}\|_2 \|\mathbf{v}\|_2}$, and Eqn. 2 reduces to what the original SimCLR uses (the term $e^{-1/\tau}$ cancels out). We leave the topmost ℓ_2 -normalization to future work. One property of InfoNCE is important for our analysis:

Theorem 2 (Property of Contrastive Loss). *For InfoNCE loss (Eqn. 2), we have $\frac{\partial L}{\partial r_+} > 0$ and*

$$\frac{\partial L}{\partial r_{k-}} < 0 \text{ for } k = 1, \dots, H \text{ and } \frac{\partial L}{\partial r_+} + \sum_{k=1}^H \frac{\partial L}{\partial r_{k-}} = 0.$$

Now, under an approximation in which we neglect variations in $\frac{\partial L}{\partial r_{k-}}$ across different k in a large minibatch, Theorem 1 and 2 imply the SimCLR gradient is governed by positive semi-definite (PSD) covariance operator at any layer (some simple contrastive loss (e.g., $r_+ - r_-$) satisfies it trivially):

Theorem 3 (Covariance Operator from Contrastive Loss). *Assume that for any $\{\mathbf{x}_1, \mathbf{x}_+, \mathbf{x}_{k-}\}$, $\frac{\partial L}{\partial r_{k-}} = -\beta/H$ is a constant, then in a large batch limit, the update for W_l is*

$$W_l(t+1) = W_l(t) + \alpha \Delta W_l(t), \quad \text{where } \text{vec}(\Delta W_l(t)) = \beta \mathbb{V}_{\mathbf{x}}[\bar{K}_l(\mathbf{x})] \text{vec}(W_l(t)). \quad (3)$$

where α is the learning rate and $\bar{K}_l(\mathbf{x}; \mathcal{W}) := \mathbb{E}_{\mathbf{x}' \sim p_{\text{aug}}(\cdot|\mathbf{x})} [K_l(\mathbf{x}'; \mathcal{W})]$ is the expected connection under the augmentation distribution, conditional on the datapoint \mathbf{x} .

Above, we use the notation $\text{Cov}[X, Y] := \mathbb{E}[XY^\top] - \mathbb{E}[X]\mathbb{E}[Y]^\top$ and $\mathbb{V}[X] := \text{Cov}[X, X]$. The covariance operator $\mathbb{V}_{\mathbf{x}}[\bar{K}_l(\mathbf{x})] \in \mathbb{R}^{n_l m_{l-1} \times n_l m_{l-1}}$ is a time-varying PSD matrix over the entire training procedure. Therefore, all its eigenvalues are non-negative and at any time t , W_l is most amplified along its largest eigenmodes. Intuitively, this covariance operator ignores different views of the same sample \mathbf{x} by averaging over the augmentation distribution to compute $\bar{K}_l(\mathbf{x})$, and then computes the expected covariance of this *augmentation averaged* connection with respect to the data distribution $p(\mathbf{x})$. Thus, at all layers, any variability in the connection across different data points, that survives augmentation averages, leads to weight amplification. This amplification of weights by the PSD data covariance of an augmentation averaged connection constitutes a fundamental description of SimCLR learning dynamics for *arbitrary* data and augmentation distributions.

To understand the intuitions of covariance operator, please check Appendix A for details.

4 Analysis of ingredients underlying the success of BYOL

In BYOL, the two networks are no-longer identical and, interestingly, only positive pairs are used for training. The first network with weights $\mathcal{W}_1 = \mathcal{W} := \{\mathcal{W}_{\text{base}}, \mathcal{W}_{\text{pred}}\}$ is an *online* network that is trained to predict the output of the second *target* network with weights $\mathcal{W}_2 = \mathcal{W}'$, using a learnable predictor $\mathcal{W}_{\text{pred}}$ to map the online to target outputs (Fig. 1(a) and Fig. 1 in [16]). In contrast, the target network has $\mathcal{W}' := \{\mathcal{W}'_{\text{base}}, \mathcal{W}'_{\text{pred}}\}$, where $\mathcal{W}'_{\text{pred}}$ is an identity and $\mathcal{W}'_{\text{base}}$ is exponential moving average (EMA) of $\mathcal{W}_{\text{base}}$: $\mathcal{W}'_{\text{base}}(t+1) = \gamma_{\text{ema}} \mathcal{W}'_{\text{base}}(t) + (1 - \gamma_{\text{ema}}) \mathcal{W}_{\text{base}}(t)$.

We now theoretically analyze BYOL to explain why the predictor and an additional ingredient BatchNorm (BN) [18] are simultaneously required for BYOL's success. Without EMA ($\gamma_{\text{ema}} = 0$) and the predictor, we obtain $\mathcal{W}' = \mathcal{W}$ and the BYOL update without BN for \mathcal{W} is (derivation in Appendix F.1):

$$\text{vec}(\Delta W_l)_{\text{sym}} = -\mathbb{E}_{\mathbf{x}} [\mathbb{V}_{\mathbf{x}' \sim p_{\text{aug}}(\cdot|\mathbf{x})} [K_l(\mathbf{x}')]] \text{vec}(W_l). \quad (4)$$

This update only promotes variance minimization in the representations of different augmented views of the same data samples and therefore would yield model collapse. Motivated by a recent blogpost [14], we now consider BN in addition to the predictor. In particular, we analyze a simplified version of BN in which the mini-batch mean only is subtracted. When adding predictor, Theorem 1 can still be applied by adding identity layers on top of the target network \mathcal{W}' so that the

online and the target networks have the same depth. Theorem 5 in [30] demonstrates this version of BN shifts the downward gradients so their mini-batch mean is 0:

$$\tilde{g}_l^i := g_l^i - \frac{1}{|B|} \sum_{i \in B} g_l^i = g_l^i - \bar{g}_l \quad (5)$$

Here g_l^i is the i -th sample in a batch and \bar{g}_l is the batch average (same for \bar{f}_l). Backpropagating through this BN (vs. just subtracting the mean only in the forward pass¹), leads to a correction term:

Theorem 4. *If (1) the network is linear from layer l to the topmost and (2) the downward gradient g_l undergoes Eqn. 5, then with large batch limits, the correction of the update is² (for brevity, dependency on \mathcal{W} is omitted, while dependency on \mathcal{W}' is made explicit):*

$$\text{vec}(\delta W_l^{\text{BN}}) = \mathbb{E}_{\mathbf{x}} [\bar{K}_l(\mathbf{x})] \{ \mathbb{E}_{\mathbf{x}} [\bar{K}_l^T(\mathbf{x})] \text{vec}(W_l) - \mathbb{E}_{\mathbf{x}} [\bar{K}_l^T(\mathbf{x}; \mathcal{W}')] \text{vec}(W_l') \} \quad (6)$$

and the corrected weight update is $\widehat{\Delta W}_l := \Delta W_l + \delta W_l^{\text{BN}}$. Using Eqn. 4, we have:

$$\text{vec}(\widehat{\Delta W}_l) = \text{vec}(\Delta W_l)_{\text{sym}} - \mathbb{V}_{\mathbf{x}} [\bar{K}_l(\mathbf{x})] \text{vec}(W_l) + \text{Cov}_{\mathbf{x}} [\bar{K}_l(\mathbf{x}), \bar{K}_l(\mathbf{x}; \mathcal{W}')] \text{vec}(W_l') \quad (7)$$

Theorem 4 makes several predictions.

SimCLR. In this case, both networks use the same weight and there is no predictor. This means $\mathcal{W}' = \mathcal{W}$. From the analysis above, we have $\delta W_l^{\text{BN}} = 0$ and BN should not matter. This is justified in the recent blogpost [14].

BYOL. When the predictor is present, $\mathcal{W}' \neq \mathcal{W}$ and BN is present, from the analysis above we know that $\delta W_l^{\text{BN}} \neq 0$, which provides an implicit contrastive term. Note that $\mathcal{W}' \neq \mathcal{W}$ means there is a predictor, the online network uses EMA, or both.

The Predictor. We first discuss the predictor without EMA. To see why $\mathcal{W}_{\text{pred}}$ plays a critical role, consider the last two terms (denoted as $\widehat{\Delta W}_l$) in Eqn. 7:

$$\widehat{\Delta W}_l = -\mathbb{V}_{\mathbf{x}} [\bar{K}_l(\mathbf{x}; \mathcal{W})] \text{vec}(W_l) + \text{Cov}_{\mathbf{x}} [\bar{K}_l(\mathbf{x}; \mathcal{W}), \bar{K}_l(\mathbf{x}; \mathcal{W}')] \text{vec}(W_l') \quad (8)$$

If there is no EMA (i.e., $\mathcal{W}'_{\text{base}} = \mathcal{W}_{\text{base}}$) and the weights of the predictor are all small positive numbers (e.g., $\approx \beta > 0$), then $+\text{Cov}_{\mathbf{x}} [\bar{K}_l(\mathbf{x}), \bar{K}_l(\mathbf{x}; \mathcal{W}')] \approx \beta \mathbb{V}_{\mathbf{x}} [\bar{K}_l(\mathbf{x}; \mathcal{W}_{\text{base}})]$ and for all layer l in $\mathcal{W}_{\text{base}}$, Eqn. 8 becomes:

$$\widehat{\Delta W}_l \approx \beta(1 - \beta) \mathbb{V}_{\mathbf{x}} [\bar{K}_l(\mathbf{x}; \mathcal{W}_{\text{base}})] \text{vec}(W_l) \quad (9)$$

Intuitively, in Eqn. 8, the first term $-\mathbb{V}_{\mathbf{x}} [\bar{K}_l(\mathbf{x}; \mathcal{W})]$ is second order in the Jacobian of $\mathcal{W}_{\text{pred}}$, while the second term $+\text{Cov}_{\mathbf{x}} [\bar{K}_l(\mathbf{x}; \mathcal{W}), \bar{K}_l(\mathbf{x}; \mathcal{W}')] \text{vec}(W_l')$ is first-order with respect to $\mathcal{W}_{\text{pred}}$. So if the predictor weight $\mathcal{W}_{\text{pred}}$ is “small” in magnitude (e.g., $\beta \ll 1$), then the latter term dominates and $\widehat{\Delta W}_l$ becomes the covariance operator. In this regime, BYOL with predictor+BN sensibly amplifies variance across data samples through the covariance operator, and minimizes variance across different augmented views of the same data sample through $-\mathbb{E}_{\mathbf{x}} [\mathbb{V}_{\mathbf{x}' \sim p_{\text{aug}}(\cdot|\mathbf{x})} [K_l(\mathbf{x}'; \mathcal{W})]]$, which is the first term in Eqn. 7. Interestingly, SimCLR doesn’t have such a term.

The Exponential Moving Average (EMA). On the other hand, the EMA part might play a different role. Consider the following linear dynamic system, which is a simplified version of Eqn. 7:

$$\mathbf{w}(t+1) - \mathbf{w}(t) = \Delta \mathbf{w}(t) = \alpha [-\mathbf{w}(t) + (1 - \lambda) \mathbf{w}_{\text{ema}}(t)] \quad (10)$$

Using z -transform, we could compute its two roots (See Appendix F.3 for the derivation):

$$z_{\min, \max} = 1 - \frac{1}{2} \left[(1 - \gamma_{\text{ema}} + \alpha) \pm \sqrt{(1 - \gamma_{\text{ema}} + \alpha)^2 - 4\alpha(1 - \gamma_{\text{ema}})\lambda} \right] \quad (11)$$

As analyzed above, the magnification factor of $\mathbf{w}_{\text{ema}}(t)$ is larger than that of $\mathbf{w}(t)$, and thus $\lambda < 0$, and $z_{\max} > 1 > z_{\min}$. As a result, $\mathbf{w}(t) \propto z_{\max}^t$ will have exponential growth and learning happens. Compared to no EMA case (i.e., $\gamma_{\text{ema}} = 0$), with a $\gamma_{\text{ema}} < 1$ but close to 1, z_{\max} becomes smaller (but still > 1) and the exponential growth is less aggressive, which stabilizes the training.

Indeed, empirical findings in a recent blogpost [14] as well as our own experiments (Tbl. 4) suggests that standard BYOL without BN fails. In addition, we also initialize the predictor with small positive weights (See Appendix H), as well as reinitialize the predictor weight once in a while (Tbl. 6), and BYOL still works well. These empirical evidences are all consistent with the theoretical prediction.

¹In PyTorch, the former is `x-x.mean()` and the latter is `x-x.mean().detach()`.

²A formal treatment requires Jacobian J to incorporate BatchNorm’s contribution and is left for future work.

References

- [1] Zeyuan Allen-Zhu and Yuanzhi Li. Backward feature correction: How deep learning performs deep learning. *arXiv preprint arXiv:2001.04413*, 2020.
- [2] Humam Alwassel, Dhruv Mahajan, Lorenzo Torresani, Bernard Ghanem, and Du Tran. Self-supervised learning by cross-modal audio-video clustering. *arXiv preprint arXiv:1911.12667*, 2019.
- [3] Sanjeev Arora, Hrishikesh Khandeparkar, Mikhail Khodak, Orestis Plevrakis, and Nikunj Saunshi. A theoretical analysis of contrastive unsupervised representation learning. February 2019.
- [4] Alexei Baevski and Abdelrahman Mohamed. Effectiveness of self-supervised pre-training for asr. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7694–7698. IEEE, 2020.
- [5] Alexei Baevski, Steffen Schneider, and Michael Auli. vq-wav2vec: Self-supervised learning of discrete speech representations. *arXiv preprint arXiv:1910.05453*, 2019.
- [6] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 132–149, 2018.
- [7] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *arXiv preprint arXiv:2006.09882*, 2020.
- [8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020.
- [9] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.
- [10] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 215–223, 2011.
- [11] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [13] William Falcon and Kyunghyun Cho. A framework for contrastive self-supervised learning and designing a new approach, 2020.
- [14] Abe Fetterman and Josh Albrecht. Understanding self-supervised and contrastive learning with "bootstrap your own latent" (byol), 2020.
- [15] Priya Goyal, Dhruv Mahajan, Abhinav Gupta, and Ishan Misra. Scaling and benchmarking self-supervised visual representation learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6391–6400, 2019.
- [16] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*, 2020.
- [17] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020.
- [18] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *ICML*, 2015.

- 222 [19] Gregory Koch, Richard Zemel, and Ruslan Salakhutdinov. Siamese neural networks for one-
223 shot image recognition. In *ICML deep learning workshop*, volume 2. Lille, 2015.
- 224 [20] Alexander Kolesnikov, Xiaohua Zhai, and Lucas Beyer. Revisiting self-supervised visual rep-
225 resentation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern*
226 *Recognition*, pages 1920–1929, 2019.
- 227 [21] Andrew K Lampinen and Surya Ganguli. An analytic theory of generalization dynamics and
228 transfer learning in deep linear networks. In *International Conference on Learning Representations (ICLR)*, 2018.
- 230 [22] Jason D Lee, Qi Lei, Nikunj Saunshi, and Jiacheng Zhuo. Predicting what you already know
231 helps: Provable self-supervised learning. *arXiv preprint arXiv:2008.01064*, 2020.
- 232 [23] Junnan Li, Pan Zhou, Caiming Xiong, Richard Socher, and Steven CH Hoi. Prototypical
233 contrastive learning of unsupervised representations. *arXiv preprint arXiv:2005.04966*, 2020.
- 234 [24] JG Liao and Arthur Berg. Sharpening jensen’s inequality. *The American Statistician*, 2018.
- 235 [25] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant rep-
236 resentations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
237 *Recognition*, pages 6707–6717, 2020.
- 238 [26] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive
239 predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- 240 [27] Senthil Purushwalkam and Abhinav Gupta. Demystifying contrastive self-supervised learning:
241 Invariances, augmentations and dataset biases. *arXiv preprint arXiv:2007.13916*, 2020.
- 242 [28] David Saad and Sara A Solla. Dynamics of on-line gradient descent learning for multilayer
243 neural networks. In *Advances in neural information processing systems*, pages 302–308, 1996.
- 244 [29] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *arXiv preprint*
245 *arXiv:1906.05849*, 2019.
- 246 [30] Yuandong Tian. A theoretical framework for deep locally connected relu network. *arXiv*
247 *preprint arXiv:1809.10829*, 2018.
- 248 [31] Yuandong Tian. Student specialization in deep relu networks with finite width and input di-
249 mension. *ICML*, 2020.
- 250 [32] Christopher Tosh, Akshay Krishnamurthy, and Daniel Hsu. Contrastive learning, multi-view
251 redundancy, and linear models. *arXiv preprint arXiv:2008.10150*, 2020.
- 252 [33] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through
253 alignment and uniformity on the hypersphere. *arXiv preprint arXiv:2005.10242*, 2020.
- 254 [34] Anne Wu, Changhan Wang, Juan Pino, and Jiatao Gu. Self-supervised representations improve
255 end-to-end speech translation. *arXiv preprint arXiv:2006.12124*, 2020.

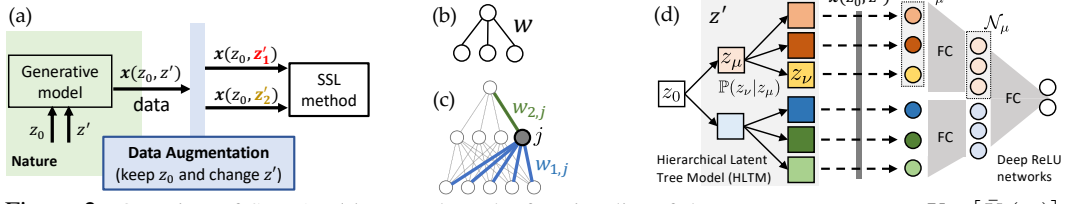


Figure 2: Overview of Sec. A. (a) To analyze the functionality of the covariance operator $\mathbb{V}_{z_0} [\bar{K}_l(z_0)]$ (Eqn. 3), we assume that Nature generates the data from a certain generative model with latent variable z_0 and z' , while data augmentation takes $x(z_0, z')$, changes z' but keeps z_0 intact. (b) Sec. A.1: one layer one neuron example. (c) Sec. A.2: two-layer case where $\mathbb{V}[\bar{K}_1]$ and $\mathbb{V}[\bar{K}_2]$ interplay. (d) Sec. A.3: Hierarchical Latent Tree Models and deep ReLU networks trained with SimCLR. A latent variable z_μ , and its corresponding nodes \mathcal{N}_μ in multi-layer ReLU side, covers a subset of input x , resembling local receptive fields in ConvNet.

A How the covariance operator drives the emergence of features

To concretely illustrate how the fundamental covariance operator derived in Theorem 3 drives feature emergence in SimCLR, we setup the following paradigm for analysis. We assume the input $x = x(z_0, z')$ is generated by two groups of latent variables, *class/sample-specific* latent z_0 and *nuisance* latent z' . After data augmentation, z_0 remains the same while z' changes (Fig. 2(a)). In this case, the covariance operator is simply $\mathbb{V}_{z_0} [\bar{K}_l(z_0)]$ since all nuisance latents z' are integrated out in $\bar{K}_l(z_0)$.

In this setting, we first show that a linear neuron performs PCA within an augmentation preserved subspace. We then consider how nonlinear neurons with local receptive fields (RFs) can learn to detect simple objects. Finally, we extend our analysis to deep ReLU networks exposed to data generated by a hierarchical latent tree model (HLT), proving that, with sufficient over-parameterization, there exist lucky nodes at initialization whose activation is correlated with latent variables underlying the data, and that SimCLR amplifies these initial lucky patterns during learning.

A.1 Self-supervised learning and the single neuron: illustrative examples

A single linear neuron performs PCA in a preserved subspace. For a single linear neuron ($L = 1, n_L = 1$), the connection in definition 1 is simply $K_1(x) = x$. Now imagine the input space x can be decomposed into the direct sum of a semantically relevant subspace, and its orthogonal complement, which corresponds to a subspace of nuisance features. Furthermore, suppose the augmentation distribution $p_{\text{aug}}(\cdot|x)$ is obtained by multiplying x by a random Gaussian matrix that acts *only* in the nuisance subspace, thereby identically preserving the semantic subspace. Then the augmentation averaged connection $\bar{K}_1(x) = Q^s x$ where Q^s is a projection operator onto the semantic subspace. In essence, only the projection of data onto the semantic subspace survives augmentation averaging, as the nuisance subspace is scrambled. Then the covariance operator in Theorem 3 is $\mathbb{V}_x [\bar{K}_1(x)] = Q^s \mathbb{V}_x [x] Q^{s\top}$. Thus the covariance of the data distribution, projected onto the semantic subspace, governs the growth of the weight vector W_1 , demonstrating SimCLR on a single linear neuron performs PCA within a semantic subspace preserved by data augmentation.

A single linear neuron cannot detect localized objects. We now consider a generative model in which data vectors can be thought of as images of objects of the form $x(z_0, z')$ where z_0 is an important latent semantic variable denoting object identity, while z' is an unimportant latent variable denoting nuisance features, like object pose or location. The augmentation procedure scrambles pose/position while preserving object identity. Consider a simple concrete example (Fig. 3(a)):

$$x(z_0, z') = \begin{cases} e_{z'} + e_{(z'+1) \bmod d} & z_0 = 1 \\ e_{z'} + e_{(z'+2) \bmod d} & z_0 = 2, \end{cases} \quad (12)$$

Here $0 \leq z' \leq d-1$ denotes d discrete translational object positions on a periodic ring and $z_0 \in \{1, 2\}$ denotes two possible objects 11 and 101. The distribution is uniform both over objects and positions: $p(z_0, z') = \frac{1}{2d}$. Augmentation shifts the object to a uniformly random position via $p_{\text{aug}}(z'|z_0) = 1/d$. For a single linear neuron $K_1(x) = x$, and the augmentation-averaged connection is $\bar{K}_1(z_0) = \frac{2}{d} \mathbf{1}$, and is actually independent of object identity z_0 (both objects activate two pixels at any location). Thus $\mathbb{V}_{z_0} [\bar{K}_1(z_0)] = 0$ and no learning happens.

A local receptive field (RF) does not help. In the same generative model, now consider a linear neuron with a local RF of width 2. Within the RF only four patterns can arise: 00, 01, 10, 11. Taking the expectation over z' given z_0 (Fig. 3(a2)) yields $\bar{K}_1(z_0=1) = \frac{1}{d} [x_{11} + x_{01} + x_{10} + (d-3)x_{00}]$

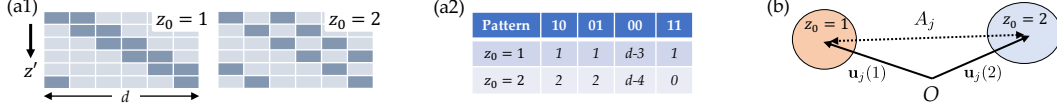


Figure 3: **(a)** Two 1D objects under translation: (a1) Two different objects 11 ($z_0 = 1$) and 101 ($z_0 = 2$) located at different locations specified by z' . (a2) The frequency table for a neuron with local receptive field of size 2. **(b)** In two-layer case (Fig. 2(c)), $\mathbb{V}[\bar{K}_1]$ and $\mathbb{V}[\bar{K}_2]$ interplay in two-cluster data distribution.

and $\bar{K}_1(z_0=2) = \frac{1}{d} [2x_{01} + 2x_{10} + (d-4)x_{00}]$. Here, $x_{11} \in \mathbb{R}^2$ denotes pattern 11. This yields

$$\mathbb{V}_{z_0} [\bar{K}_1(z_0)] = \frac{1}{4d^2} \mathbf{u} \mathbf{u}^\top \quad \text{where } \mathbf{u} := \mathbf{x}_{11} + \mathbf{x}_{00} - \mathbf{x}_{01} - \mathbf{x}_{10}. \quad (13)$$

and $\mathbb{V}_{z_0} [\bar{K}_1(z_0)] \in \mathbb{R}^{2 \times 2}$ since the RF has width 2. Note that the signed sum of the four pattern vectors in \mathbf{u} actually cancel, so that $\mathbf{u} = \mathbf{0}$, implying $\mathbb{V}_{z_0} [\bar{K}_1(z_0)] = 0$ and no learning happens. Interestingly, although the conditional distribution of the 4 input patterns depends on the object identity z_0 (Fig. 3(a2)), a linear neuron cannot learn to discriminate the objects.

A nonlinear neuron with local RF can learn to detect object selective features. With a ReLU neuron with weight vector \mathbf{w} , from Def. 1, the connection is now $K_1(\mathbf{x}, \mathbf{w}) = \psi'(\mathbf{w}^\top \mathbf{x})$. Suppose $\mathbf{w}(t)$ happens to be selective for a *single* pattern \mathbf{x}_p (where $p \in \{00, 01, 10, 11\}$), i.e., $\mathbf{w}(t)^\top \mathbf{x}_p > 0$ and $\mathbf{w}(t)^\top \mathbf{x}_{p'} < 0$ for $p' \neq p$. The augmentation averaged connection is then $\bar{K}_1(z_0) \propto \mathbf{x}_p$ where the proportionality constant depends on object identity z_0 and can be read off (Fig. 3(a2)). Since this averaged connection varies with object identity z_0 for all p , the covariance operator does not vanish and is given by $\mathbb{V}_{z_0} [\bar{K}_1(z_0)] = c_p \mathbf{x}_p \mathbf{x}_p^\top$ where the constant $c_p > 0$ depends on the selective pattern p and can be computed from Fig. 3(a2). By Theorem 3, the dot product $\mathbf{x}_p^\top \mathbf{w}(t)$ grows over time:

$$\mathbf{x}_p^\top \mathbf{w}(t+1) = \mathbf{x}_p^\top (I_{2 \times 2} + \alpha \beta c_p \mathbf{x}_p \mathbf{x}_p^\top) \mathbf{w}(t) = (1 + \alpha \beta c_p \|\mathbf{x}_p\|^2) \mathbf{x}_p^\top \mathbf{w}_j(t) > \mathbf{x}_p^\top \mathbf{w}_j(t) > 0. \quad (14)$$

Thus the learning dynamics amplifies the initial selectivity to the object selective feature vector \mathbf{x}_p in a way that cannot be done with a linear neuron. Note this argument also holds with bias terms and initial selectivity for more than one patterns. Moreover, with a local RF, the probability of weak initial selectivity to local object sensitive features is high, and we may expect amplification of such weak selectivity in real neural network training.

313 A.2 Two-layer case with multiple hidden neurons

Now consider a two-layer network ($L = 2$) with 1 output ($n_L = 1$) and 1-hidden layer with n_1 ReLU neurons (Fig. 2(c)). In this case, the augmentation-averaged connection $\bar{K}_1(z_0)$ at the lowest layer $l = 1$ can be written as

$$\bar{K}_1(z_0) = [w_{2,1} \mathbf{u}_1^\top(z_0), w_{2,2} \mathbf{u}_2^\top(z_0), \dots, w_{2,n_1} \mathbf{u}_{n_1}^\top(z_0)]^\top \in \mathbb{R}^{n_1 d} \quad (15)$$

where $\mathbf{u}_j(z_0) := \mathbb{E}_{z'|z_0} [\mathbf{x}(z_0, z') \mathbb{I}(\mathbf{w}_{1,j}^\top \mathbf{x}(z_0, z') \geq 0)]$ is the expected activation of the hidden layer. Note that each $\mathbf{u}_j(z_0) \in \mathbb{R}^d$ and $w_{2,j}$ is the scalar weight of the second weight matrix that connects node j to the output. Then we have the following theorem:

Theorem 5. *If $\text{Cov}_{z_0} [\mathbf{u}_j, \mathbf{u}_k] = 0$ for $j \neq k$, then the time derivative of $w_{2,j}$ and $\mathbf{w}_{1,j}$ satisfies:*

$$\dot{w}_{2,j} = w_{2,j} \mathbf{w}_{1,j}^\top A_j \mathbf{w}_{1,j}, \quad \dot{\mathbf{w}}_{1,j} = w_{2,j}^2 A_j \mathbf{w}_{1,j}, \quad \text{where } A_j := \mathbb{V}_{z_0} [\mathbf{u}_j(z_0)]. \quad (16)$$

It is easy to see that $d\|w_{2,j}\|^2/dt = d\|\mathbf{w}_{1,j}\|^2/dt$ and thus $w_{2,j}^2 = \|\mathbf{w}_{1,j}\|^2 + c$ where c is some constant which doesn't change over time. Since A_j is PSD, $\mathbf{w}_{1,j}^\top A_j \mathbf{w}_{1,j} \geq 0$ and $w_{2,j}$ always grows, which in turn accelerates the convergence of $\mathbf{w}_{1,j}$ to the largest eigenvector of A_j . This shows *top-down modulation*: the existence of top-layer weights accelerate the training of the lower layer. Previous works [1] also mention a similar concept in supervised learning, called “backward feature corrections”, for deep polynomial networks with quadratic activations. Here we provide a rigorous formulation showing similar behaviors happen in SSL with SGD training.

Note that if we consider ReLU neurons, A_j changes with $\mathbf{w}_{1,j}$; while for linear nodes, A_j is a constant, since the gating $\mathbb{I}(\mathbf{w}_{1,j}^\top \mathbf{x} > 0)$ is always 1. If the data \mathbf{x} forms a mixture of two Gaussians:

Symbol	Definition	Size	Description
$\mathcal{N}_l, \mathcal{Z}_l$			The set of all nodes and all latent variables at layer l .
$\mathcal{N}_\mu, \mathcal{N}_\mu^{\text{ch}}$			Nodes corresponding to latent variable z_μ . $\mathcal{N}_\mu^{\text{ch}}$ are children under \mathcal{N}_μ .
$P_{\mu\nu}$	$[\mathbb{P}(z_\nu z_\mu)]$	2×2	The top-down transition probability from z_μ to z_ν .
$v_j(z_\mu), \mathbf{v}_j$	$\mathbb{E}_z[f_j z_\mu], [v_j(z_\mu)]$	scalar, 2	Expected activation f_j given z_μ (z_μ 's descendants are marginalized).
$\mathbf{f}_\mu, \mathbf{f}_{\mathcal{N}_\mu^{\text{ch}}}$	$[f_j]_{j \in \mathcal{N}_\mu}, [f_k]_{k \in \mathcal{N}_\mu^{\text{ch}}}$	$ \mathcal{N}_\mu , \mathcal{N}_\mu^{\text{ch}} $	Activations for all nodes $j \in \mathcal{N}_\mu$ and for the children of \mathcal{N}_μ .
$\rho_{\mu\nu}$	$2\mathbb{P}(z_\nu=1 z_\mu=1) - 1$	scalar in $[-1, 1]$	Polarity of the transitional probability.
ρ_0	$\mathbb{P}(z_0=1) - \mathbb{P}(z_0=0)$	scalar	Polarity of probability of root latent z_0 .
s_k	$\frac{1}{2}(v_k(1) - v_k(0))$	scalar	Discrepancy of node k w.r.t its latent variable $z_{\nu(k)}$.
$\mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}}$	$[\rho_{\mu\nu(k)} s_k]_{k \in \mathcal{N}_\mu^{\text{ch}}}$	$ \mathcal{N}_\mu^{\text{ch}} $	Child selectivity vector.

Table 1: Notation for Sec. A.3 (binary symmetric HLTM).

330 $\mathbf{x} \sim \frac{1}{2}\mathbb{I}(z_0=1)N(\mathbf{w}_1^*, \sigma^2) + \frac{1}{2}\mathbb{I}(z_0=2)N(\mathbf{w}_2^*, \sigma^2)$ and let $\Delta\mathbf{w}^* := \mathbf{w}_1^* - \mathbf{w}_2^*$, then in the linear case,
331 $A_j \sim \Delta\mathbf{w}^* \Delta\mathbf{w}^{*\top}$ and $\mathbf{w}_{1,j}$ converges to $\pm\Delta\mathbf{w}^*$ (Fig. 3(b)). In the nonlinear case with multiple
332 Gaussians, if one of the Gaussians sits at the origin (e.g., background noise), then dependent on
333 initialization, A_j evolves into $\mathbf{w}_k^* \mathbf{w}_k^{*\top}$ for some center k , and $\mathbf{w}_{1,j} \rightarrow \mathbf{w}_k^*$. Note this dynamics is
334 insensitive to specific parametric forms of the input data.

335 A.3 Deep ReLU SSL training with Hierarchical Latent Tree Models (HLTMs)

336 For multi-layer case, we mainly study a hierarchical latent tree model, where each leaf variable is
337 sampled via a Markov Chain from the root latent variable z_0 . The model is motivated by the structure
338 of ConvNet: each latent variable at the intermediate layer generates a *subset* of the observable data \mathbf{x}
339 that is sent to the deep ReLU networks, in which the nodes also have local receptive fields (Fig. 2(d)).

340 We define symbols in Tbl. 1. At layer l , we have categorical latent variables $\{z_\mu\}$, where $\mu \in \mathcal{Z}_l$
341 is the Greek letter index for latent variables. Each z_μ can take discrete values. The topmost latent
342 variable is z_0 . Following the tree structure, for $\mu \in \mathcal{Z}_l$ and $\nu_1, \nu_2 \in \mathcal{Z}_{l-1}$, conditional independence
343 holds: $\mathbb{P}(z_{\nu_1}, z_{\nu_2}|z_\mu) = \mathbb{P}(z_{\nu_1}|z_\mu)\mathbb{P}(z_{\nu_2}|z_\mu)$. The final sample \mathbf{x} is the instantiation of the leaf
344 latents (Fig. 2(d)), and thus depends on all latent variables. Corresponding to the hierarchical tree
345 model, each neural network node $j \in \mathcal{N}_l$ maps to a unique $\mu = \mu(j) \in \mathcal{Z}_l$. Let \mathcal{N}_μ be all nodes that
346 map to μ . For $j \in \mathcal{N}_\mu$, its activation f_j only depends on the value of z_μ and its descendant latent
347 variables. We also define $v_j(z_\mu) := \mathbb{E}_z[f_j|z_\mu]$ as the expected activation w.r.t z_μ . Given a sample
348 \mathbf{x} , the data augmentation is done by *resampling* all z_μ (which are z' in Fig. 2), given the root z_0 .

349 **Symmetric Binary HLTMs.** Here we consider a symmetric binary case: each $z_\mu \in \{0, 1\}$ and for
350 $\mu \in \mathcal{Z}_l, \nu \in \mathcal{Z}_{l-1}$, $\mathbb{P}(z_\nu=1|z_\mu=1) = \mathbb{P}(z_\nu=0|z_\mu=0) = (1 + \rho_{\mu\nu})/2$, where the *polarity*
351 $\rho_{\mu\nu} \in [-1, 1]$ measures how informative $z_\mu = 1$ is. If $\rho_{\mu\nu} = \pm 1$ then there is no stochasticity in
352 the top-down generation process; $\rho_{\mu\nu} = 0$ means no information in the downstream latents and the
353 posterior of z_0 given the observation \mathbf{x} can only be uniform. See Appendix for more general cases.

354 Now we compute covariance operator $\mathbb{V}_{z_0}[\bar{K}_\mu(z_0)]$ at different layers, where $\bar{K}_\mu(z_0) =$
355 $\mathbb{E}_{z'}[\mathbf{f}_{\mathcal{N}_\mu^{\text{ch}}} \otimes J_\mu^\top | z_0]$, we first check the term $\mathbb{E}_{z'}[\mathbf{f}_{\mathcal{N}_\mu^{\text{ch}}} | z_0]$ and incorporate the Jacobian later.

356 **Theorem 6** (Covariance operator in binary HLTMs). $\mathbb{V}_{z_0}[\mathbb{E}_{z'}[\mathbf{f}_{\mathcal{N}_\mu^{\text{ch}}} | z_0]] = \rho_{0\mu}^2(1 - \rho_{0\mu}^2)\mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}}\mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}}^\top$.
357 Here $\mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}} := [\rho_{\mu\nu(k)} s_k]_{k \in \mathcal{N}_\mu^{\text{ch}}}$. If $\max_{\alpha\beta} |\rho_{\alpha\beta}| \leq \gamma < 1$, then $\lim_{L \rightarrow +\infty} \rho_{0\mu} \rightarrow 0$ for $\mu \in \mathcal{N}_l$.

358 Theorem 6 suggests when $\rho_{0\mu}$ and $\|\mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}}\|$ is large, the covariance operator is large and training
359 is faster. For deep HLTMs and deep networks, at lower layers, $\rho_{0\mu} \rightarrow 0$ and $P_{0\mu}$ becomes uniform
360 due to mixing from the progression of the Markov Chain, making $\mathbb{V}_{z_0}[\bar{K}_l(z_0)]$ small. Therefore,
361 training in SSL is faster at the top layers where the covariance operators have large magnitude. On
362 the other hand, large $\|\mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}}\|$ means that $s_k := (v_k(1) - v_k(0))/2$ is large, or the expected activation
363 $v_k(z_\nu)$ is *selective* among different values of z_ν for $\nu \in \text{ch}(\mu)$. Interestingly, this can be achieved
364 by *over-parameterization* (i.e., $|\mathcal{N}_\mu| > 1$):

365 **Theorem 7** (Lucky nodes in deep ReLU networks regarding to binary HLTMs). Suppose each element of the weights W_l between layer $l+1$ and l are initialized with
366 Uniform $[-\sigma_w \sqrt{3/|\mathcal{N}_\mu^{\text{ch}}|}, \sigma_w \sqrt{3/|\mathcal{N}_\mu^{\text{ch}}|}]$. There exists σ_l^2 so that $\mathbb{V}[f_k|z_\nu] \leq \sigma_l^2$ for any $k \in \mathcal{N}_l$.
367 For any $\mu \in \mathcal{Z}_{l+1}$, if $|\mathcal{N}_\mu| = O(\exp(c))$, then with high probability, there exists at least one node
368 $j \in \mathcal{N}_\mu$ so that their pre-activation gap $|\tilde{v}_j(1) - \tilde{v}_j(0)| = 2\mathbf{w}_j^\top \mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}} > 0$ and the activations satisfy:

$$\left| v_j^2(1) - v_j^2(0) \right| \geq 3\sigma_w^2 \left[\frac{1}{4|\mathcal{N}_\mu^{\text{ch}}|} \sum_{k \in \mathcal{N}_\mu^{\text{ch}}} |v_k(1) - v_k(0)|^2 \left(\frac{c+6}{6} \rho_{\mu\nu}^2 - 1 \right) - \sigma_l^2 \right]. \quad (17)$$

Intuitively, this means that with large polarity $\rho_{\mu\nu}$ (or strong signals sending top-down), over-parameterized ReLU networks with random initialization yields selective neurons, if the lower layer also contains selective ones. For example, when $\sigma_l = 0$, $c = 9$, if $\rho_{\mu\nu} \geq 63.3\%$ then there is a gap between expected activations $v_j(1)$ and $v_j(0)$, and the gap is larger when the selectivity in the lower layer l is higher. Note that at the lowest layer, $\{v_k\}$ are themselves observable leaf latent variables and are selective by definition. So a bottom-up mathematical induction will happen.

After the gradient update, if we assume the top-down Jacobian is just 1, then for the “lucky” node j we will have (for brevity, let $c := \alpha\beta\rho_{0\mu}^2(1 - \rho_0^2)$):

$$\mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}}^\top \mathbf{w}_j(t+1) = \mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}}^\top \left[I + c \mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}} \mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}}^\top \right] \mathbf{w}_j(t) = (1 + c \|\mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}}\|_2^2) \mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}}^\top \mathbf{w}_j(t) > \mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}}^\top \mathbf{w}_j(t) > 0$$

which means that the pre-activation gap $\tilde{v}_j(1) - \tilde{v}_j(0) = 2\mathbf{w}_j^\top \mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}}$ grows over time and the latent variable z_μ is *learned* (instantiated as f_j) during training, even if the network is never supervised with its true value. Similar to Sec. A.2, once the top layer starts to have large weights, $\tilde{K}_l(z_0)$ for lower layer becomes larger due to large Jacobian and training is accelerated.

We implement the HLTm and confirm, as predicted by our theory, that the intermediate layers of deep ReLU networks do indeed learn the latent variables of the HLTm (see Tbl. 2 below).

B Proofs: Background and Basic Setting (Section 2)

B.1 Theorem 1

Definition 2 (reversibility). *A layer l is reversible if there is a $G_l(\mathbf{x}; \mathcal{W}) \in \mathbb{R}^{n_l \times n_{l-1}}$ so that $\mathbf{f}_l(\mathbf{x}; \mathcal{W}) = G_l(\mathbf{x}; \mathcal{W}) \mathbf{f}_{l-1}(\mathbf{x}; \mathcal{W})$ and $\mathbf{g}_{l-1} = G_l^\top(\mathbf{x}; \mathcal{W}) Q_l \mathbf{g}_l$ for some constant PSD matrix $\mathbb{R}^{n_l \times n_l} \ni Q_l \succeq 0$. A network is reversible if all layers are.*

Note that many different kinds of layers have this reversible property, including linear layers (MLP and Conv) and (leaky) ReLU nonlinearity. For multi-layer ReLU network, for each layer l , we have:

$$G_l(\mathbf{x}; \mathcal{W}) = D_l(\mathbf{x}; \mathcal{W}) W_l, \quad Q_l \equiv I_{n_l \times n_l} \quad (18)$$

where $D_l \in \mathbb{R}^{n_l \times n_l}$ is a binary diagonal matrix that encodes the gating of each neuron at layer l . The gating $D_l(\mathbf{x}; \mathcal{W})$ depends on the current input \mathbf{x} and current weight \mathcal{W} .

In addition to ReLU, other activation function also satisfies this condition, including linear, LeakyReLU and monomial activations. For example, for power activation $\psi(x) = x^p$ where $p > 1$, we have (where $\tilde{\mathbf{f}}_l$ is the pre-activation at layer l):

$$G_l(\mathbf{x}; \mathcal{W}) = \text{diag}^{p-1}(\tilde{\mathbf{f}}_l) W_l, \quad Q_l \equiv p I_{n_l \times n_l} \quad (19)$$

Remark. Note that the reversibility is not the same as invertible. Specifically, reversibility only requires the transfer function of a backpropagation gradient is a transpose of the forward function.

Lemma 1 (Recursive Gradient Update (Extension to Lemma 1 in [31])). *Let (pseudo)-Jacobian matrix $J_L(\mathbf{x}) = I_{n_L \times n_L}$, and recursively define $J_{l-1}(\mathbf{x}) := J_l(\mathbf{x}) \sqrt{Q_l} G_l(\mathbf{x}) \in \mathbb{R}^{n_l \times n_{l-1}}$. Here $\sqrt{Q_l}$ is the constant PSD matrix so that $\sqrt{Q_l} \sqrt{Q_l}^\top = Q_l \succeq 0$.*

If (1) the network is reversible (Def. 2) and (2) $\sqrt{Q_l}$ commutes with $J_l(\mathbf{x}_1)^\top J_l(\mathbf{x}_1)$ and $J_l(\mathbf{x}_1)^\top J_l(\mathbf{x}_2)$, then minimizing the ℓ_2 objective

$$r(\mathcal{W}_1) := \frac{1}{2} \|\mathbf{f}_L(\mathbf{x}_1; \mathcal{W}_1) - \mathbf{f}_L(\mathbf{x}_2; \mathcal{W}_2)\|_2^2 \quad (20)$$

with respect to weight matrix W_l at layer l yields the following gradient at layer l :

$$\mathbf{g}_l = J_l^\top(\mathbf{x}_1; \mathcal{W}_1) [J_l(\mathbf{x}_1; \mathcal{W}_1) \mathbf{f}_l(\mathbf{x}_1; \mathcal{W}_1) - J_l(\mathbf{x}_2; \mathcal{W}_2) \mathbf{f}_l(\mathbf{x}_2; \mathcal{W}_2)] \quad (21)$$

406 *Proof.* We prove by induction. Note that our definition of W_l is the transpose of W_l defined in [31].
 407 Also our $\mathbf{g}_l(\mathbf{x})$ is the gradient *before* nonlinearity, while [31] uses the same symbol for the gradient
 408 after nonlinearity.

409 For notation brevity, we let $\mathbf{f}_l(\mathbf{x}_1) := \mathbf{f}_l(\mathbf{x}_1; \mathcal{W}_l)$ and $G_l(\mathbf{x}_1) := G_l(\mathbf{x}_1; \mathcal{W}_l)$. Similar for \mathbf{x}_2 and
 410 W_2 .

411 When $l = L$, by the property of ℓ_2 -loss, we know that $\mathbf{g}_L = \mathbf{f}_L(\mathbf{x}_1; \mathcal{W}_1) - \mathbf{f}_L(\mathbf{x}_2; \mathcal{W}_2)$, by setting
 412 $J_L(\mathbf{x}_1) = J_L(\mathbf{x}_2) = I$, the condition holds. Now suppose for layer l , we have:

$$\mathbf{g}_l = J_l^\top(\mathbf{x}_1) [J_l(\mathbf{x}_1) \mathbf{f}_l(\mathbf{x}_1) - J_l(\mathbf{x}_2) \mathbf{f}_l(\mathbf{x}_2)] \quad (22)$$

413 Then:

$$\mathbf{g}_{l-1} = G_l^\top(\mathbf{x}_1) Q_l \mathbf{g}_l \quad (23)$$

$$= G_l^\top(\mathbf{x}_1) Q_l J_l^\top(\mathbf{x}_1) [J_l(\mathbf{x}_1) \mathbf{f}_l(\mathbf{x}_1) - J_l(\mathbf{x}_2) \mathbf{f}_l(\mathbf{x}_2)] \quad (24)$$

$$= \underbrace{G_l^\top(\mathbf{x}_1) \sqrt{Q_l} J_l^\top(\mathbf{x}_1)}_{J_{l-1}^\top(\mathbf{x}_1)} \cdot [J_l(\mathbf{x}_1) \sqrt{Q_l} \mathbf{f}_l(\mathbf{x}_1) - J_l(\mathbf{x}_2) \sqrt{Q_l} \mathbf{f}_l(\mathbf{x}_2)] \quad (25)$$

$$= J_{l-1}^\top(\mathbf{x}_1) \left[\underbrace{J_l(\mathbf{x}_1) \sqrt{Q_l} G_l(\mathbf{x}_1)}_{J_{l-1}(\mathbf{x}_1)} \mathbf{f}_{l-1}(\mathbf{x}_1) - \underbrace{J_l(\mathbf{x}_2) \sqrt{Q_l} G_l(\mathbf{x}_2)}_{J_{l-1}(\mathbf{x}_2)} \mathbf{f}_{l-1}(\mathbf{x}_2) \right] \quad (26)$$

$$= J_{l-1}^\top(\mathbf{x}_1) [J_{l-1}(\mathbf{x}_1) \mathbf{f}_{l-1}(\mathbf{x}_1) - J_{l-1}(\mathbf{x}_2) \mathbf{f}_{l-1}(\mathbf{x}_2)] \quad (27)$$

414 Note that for multi-layered ReLU network, $G_l(\mathbf{x}) = D_l(\mathbf{x}) W_l$, $Q_l = I$ for each ReLU+Linear
 415 layer, if we set $\mathbf{x}_1 = \mathbf{x}_2 = \mathbf{x}$, $\mathcal{W}_1 = \mathcal{W}$, $\mathcal{W}_2 = \mathcal{W}^*$ (teacher weights), then we go back to the
 416 original Lemma 1 in [31]. \square

417 **Remark on the top-most ℓ_2 -normalization layer.** For ℓ_2 -normalized layer $\mathbf{f}_l := \mathbf{f}_{l-1} / \|\mathbf{f}_{l-1}\|_2$,
 418 we have $G_l := 1 / \|\mathbf{f}_{l-1}\|_2 \cdot I_{n_l \times n_l}$ and due to the following identity (here $\tilde{\mathbf{y}} := \mathbf{y} / \|\mathbf{y}\|_2$):

$$\frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{y}} = \frac{1}{\|\mathbf{y}\|_2} (I - \tilde{\mathbf{y}} \tilde{\mathbf{y}}^\top) \quad (28)$$

419 Therefore we have $\partial \mathbf{f}_l / \partial \mathbf{f}_{l-1} = (I_{n_l \times n_l} - \mathbf{f}_l \mathbf{f}_l^\top) G_l$ and we could set $Q_l := I - \mathbf{f}_l \mathbf{f}_l^\top$, which is a
 420 projection matrix and thus PSD. Furthermore, since the normalization layer is at the topmost, $J_l = I$
 421 and Q_l trivially commutes with $J_l^\top J_l$.

422 The *only issue* is that Q_l is not a constant PSD matrix and can change over training. Therefore
 423 Lemma 1 doesn't apply exactly to such a layer but can be regarded as an approximate way to model.

424 **Remark on ResNet.** Note that the same structure holds for blocks of ResNet with ReLU activation.

425 Now we prove Theorem 1 in a more general setting where the network is reversible (note that deep
 426 ReLU networks are included and its Q_l is a simple identity matrix):

427 **Lemma 2** (Squared ℓ_2 Gradient for dual deep reversible networks). *The gradient g_{W_l} of the squared*
 428 *loss r with respect to $W_l \in \mathbb{R}^{n_l \times n_{l-1}}$ for a single input pair $\{\mathbf{x}_1, \mathbf{x}_2\}$ is:*

$$g_{W_l} = \text{vec}(\partial r / \partial W_{1,l}) = K_{1,l} [K_{1,l}^\top \text{vec}(W_{1,l}) - K_{2,l}^\top \text{vec}(W_{2,l})]. \quad (29)$$

429 Here $K_l(\mathbf{x}; \mathcal{W}) := \mathbf{f}_{l-1}(\mathbf{x}; \mathcal{W}) \otimes J_l^\top(\mathbf{x}; \mathcal{W})$, $K_{1,l} := K_l(\mathbf{x}_1; \mathcal{W}_1)$ and $K_{2,l} := K_l(\mathbf{x}_2; \mathcal{W}_2)$.

430 *Proof.* We consider more general case where the two towers have different parameters, namely \mathcal{W}_1
 431 and \mathcal{W}_2 . Applying Lemma 1 for the branch with input \mathbf{x}_1 at the linear layer l , and we have:

$$\mathbf{g}_{1,l} = J_{1,l}^\top [J_{1,l} W_{1,l} \mathbf{f}_{1,l-1} - J_{2,l} W_{2,l} \mathbf{f}_{2,l-1}] \quad (30)$$

432 where $\mathbf{f}_{1,l-1} := \mathbf{f}_{l-1}(\mathbf{x}_1; \mathcal{W}_1)$ is the activation of layer $l-1$ just below the linear layer at tower 1
 433 (similar for other symbols). The gradient (and the weight update, according to gradient descent) of
 434 the weight W_l between layer l and layer $l-1$ is:

$$\frac{\partial r}{\partial W_{1,l}} = \mathbf{g}_{1,l} \mathbf{f}_{1,l-1}^\top \quad (31)$$

$$= J_{1,l}^\top J_{1,l} W_{1,l} \mathbf{f}_{1,l-1} \mathbf{f}_{1,l-1}^\top - J_{1,l}^\top J_{2,l} W_{2,l} \mathbf{f}_{2,l-1} \mathbf{f}_{1,l-1}^\top \quad (32)$$

435 Using $\text{vec}(AXB) = (B^\top \otimes A)\text{vec}(X)$ (where \otimes is the Kronecker product), we have:

$$\text{vec}\left(\frac{\partial r}{\partial W_{1,l}}\right) = \left(\mathbf{f}_{1,l-1}\mathbf{f}_{1,l-1}^\top \otimes J_{1,l}^\top J_{1,l}\right) \text{vec}(W_{1,l}) - \left(\mathbf{f}_{1,l-1}\mathbf{f}_{2,l-1}^\top \otimes J_{1,l}^\top J_{2,l}\right) \text{vec}(W_{2,l}) \quad (33)$$

436 Let

$$K_l(\mathbf{x}; W) := \mathbf{f}_{l-1}(\mathbf{x}; \mathcal{W}) \otimes J_l^\top(\mathbf{x}; \mathcal{W}) \in \mathbb{R}^{n_l n_{l-1} \times n_L} \quad (34)$$

437 Note that $K_l(\mathbf{x}; W)$ is a function of the current weight W , which includes weights at all layers. By
438 the mixed-product property of Kronecker product $(A \otimes B)(C \otimes D) = AC \otimes BD$, we have:

$$\text{vec}\left(\frac{\partial r}{\partial W_{1,l}}\right) = K_l(\mathbf{x}_1)K_l(\mathbf{x}_1)^\top \text{vec}(W_{1,l}) - K_l(\mathbf{x}_1)K_l(\mathbf{x}_2)^\top \text{vec}(W_{2,l}) \quad (35)$$

$$= K_l(\mathbf{x}_1) [K_l(\mathbf{x}_1)^\top \text{vec}(W_{1,l}) - K_l(\mathbf{x}_2)^\top \text{vec}(W_{2,l})] \quad (36)$$

439 where $K_l(\mathbf{x}_1) = K_l(\mathbf{x}_1; \mathcal{W}_1)$ and $K_l(\mathbf{x}_2) = K_l(\mathbf{x}_2; \mathcal{W}_2)$.

440 In SimCLR case, we have $\mathcal{W}_1 = \mathcal{W}_2 = \mathcal{W}$ so

$$\text{vec}\left(\frac{\partial r}{\partial W_l}\right) = K_l(\mathbf{x}_1) [K_l(\mathbf{x}_1) - K_l(\mathbf{x}_2)]^\top \text{vec}(W_l) \quad (37)$$

441

□

442 C Proofs: Analysis of SimCLR using Teacher-Student Setting (Section 3)

443 C.1 Theorem 2

444 *Proof.* we have:

$$\frac{\partial L}{\partial r_+} = \frac{1}{\tau} \left(1 - \frac{e^{-r_+/\tau}}{e^{-r_+/\tau} + \sum_{k'=1}^H e^{-r_{k'}/\tau}} \right) > 0 \quad (38)$$

$$\frac{\partial L}{\partial r_{k-}} = -\frac{1}{\tau} \left(\frac{e^{-r_{k-}/\tau}}{e^{-r_+/\tau} + \sum_{k'=1}^H e^{-r_{k'}/\tau}} \right) < 0, \quad k = 1, \dots, H \quad (39)$$

445 and obviously we have:

$$\frac{\partial L}{\partial r_+} + \sum_{k=1}^H \frac{\partial L}{\partial r_{k-}} = 0 \quad (40)$$

446

□

447 C.2 Theorem 3

448 *Proof.* First we have:

$$\text{vec}(g_{W_l}) = \frac{\partial L}{\partial W_l} = \frac{\partial L}{\partial r_+} \frac{\partial r_+}{\partial W_l} + \sum_{k=1}^H \frac{\partial L}{\partial r_{k-}} \frac{\partial r_{k-}}{\partial W_l} \quad (41)$$

449 Then we compute each terms. Using Theorem 1, we know that:

$$\frac{\partial r_+}{\partial W_l} = K_l(\mathbf{x}_1)(K_l(\mathbf{x}_1) - K_l(\mathbf{x}_+))^\top \text{vec}(W_l) \quad (42)$$

$$\frac{\partial r_{k-}}{\partial W_l} = K_l(\mathbf{x}_1)(K_l(\mathbf{x}_1) - K_l(\mathbf{x}_{k-}))^\top \text{vec}(W_l), \quad k = 1, \dots, n \quad (43)$$

450 Since Eqn. 40 holds, $K_l(\mathbf{x}_1)K_l^\top(\mathbf{x}_1)$ will be cancelled out and we have:

$$\text{vec}(g_{W_l}) = K_l(\mathbf{x}_1) \sum_{k=1}^H \left[\frac{\partial L}{\partial r_{k-}} (K_l(\mathbf{x}_+) - K_l(\mathbf{x}_{k-}))^\top \right] \text{vec}(W_l) \quad (44)$$

451 Then we use the assumption that $\frac{\partial L}{\partial r_{k-}} = -\beta/H$ with $\beta > 0$:

$$\text{vec}(g_{W_l}) = -\beta K_l(\mathbf{x}_1) \left[K_l^\top(\mathbf{x}_+) - \frac{1}{n} \sum_{k=1}^H K_l^\top(\mathbf{x}_{k-}) \right] \text{vec}(W_l) \quad (45)$$

452 Taking large batch limits, we know that $\mathbb{E}[K_l(\mathbf{x}_1)K_l^\top(\mathbf{x}_+)] = \mathbb{E}_{\mathbf{x}}[\bar{K}_l(\mathbf{x})\bar{K}_l^\top(\mathbf{x})]$ since $\mathbf{x}_1, \mathbf{x}_+ \sim$
 453 $p_{\text{aug}}(\cdot|\mathbf{x})$ are all augmented data points from a common sample \mathbf{x} . On the other hand,
 454 $\mathbb{E}[K_l(\mathbf{x}_1)K_l^\top(\mathbf{x}_{k-})] = \mathbb{E}_{\mathbf{x}}[\bar{K}_l(\mathbf{x})] \mathbb{E}_{\mathbf{x}}[\bar{K}_l^\top(\mathbf{x})]$ since \mathbf{x}_1 and \mathbf{x}_{k-} are generated from indepen-
 455 dent samples and data augmentation. Therefore,

$$\text{vec}(g_{W_l}) = -\beta \left\{ \mathbb{E}_{\mathbf{x}}[\bar{K}_l(\mathbf{x})\bar{K}_l^\top(\mathbf{x})] - \mathbb{E}_{\mathbf{x}}[\bar{K}_l(\mathbf{x})] \mathbb{E}_{\mathbf{x}}[\bar{K}_l^\top(\mathbf{x})] \right\} \text{vec}(W_l) \quad (46)$$

$$= -\beta \mathbb{V}_{\mathbf{x}}[\bar{K}_l(\mathbf{x})] \text{vec}(W_l) \quad (47)$$

456 \square

457 **D Proofs: The dynamics of two-layer ReLU network and the interplays of** 458 **covariance operators between nearby layers (Section A.2)**

459 **D.1 Theorem 5**

460 *Proof.* For convenience, we define the centralized version of $\mathbf{u}_j(z_0)$: $\hat{\mathbf{u}}_j(z_0) = \mathbf{u}_j(z_0) -$
 461 $\mathbb{E}_{z_0}[\mathbf{u}_j(z_0)]$ and the matrices $A_{jk} := \text{Cov}_{z_0}[\mathbf{u}_j(z_0), \mathbf{u}_k(z_0)] = \mathbb{E}_{z_0}[\hat{\mathbf{u}}_j(z_0)\hat{\mathbf{u}}_k^\top(z_0)]$.

462 At layer $l = 1$ the covariance operator is $\mathbb{V}_{z_0}[\bar{K}_1(z_0)] = [w_{2,j}w_{2,k}A_{jk}] \in \mathbb{R}^{n_1 \times n_1}$.

463 On the other hand, if we check the second layer $l = 2$, we could compute $\bar{K}_2(z_0) \in \mathbb{R}^{n_1}$:

$$\bar{K}_2(z_0) = \mathbb{E}_{z'|z_0}[\mathbf{f}_1|z_0] = \mathbb{E}_{z'|z_0}[D_1 W_1 \mathbf{x}|z_0] = \mathbb{E}_{z'|z_0}[\mathbf{x}^\top \otimes D_1(\mathbf{x})|z_0] \text{vec}(W_1) \quad (48)$$

464 For each component j , we have $[\bar{K}_2(z_0)]_j = \mathbf{w}_{1,j}^\top \mathbf{u}_j(z_0)$, or:

$$\bar{K}_2(z_0) = [\mathbf{w}_{1,1}^\top \mathbf{u}_1(z_0), \mathbf{w}_{1,2}^\top \mathbf{u}_2(z_0), \dots, \mathbf{w}_{1,n_1}^\top \mathbf{u}_{n_1}(z_0)] \quad (49)$$

465 So at layer $l = 2$ we can compute the covariance operator $\mathbb{V}_{z_0}[\bar{K}_2(z_0)] = [\mathbf{w}_{1,j}^\top A_{jk} \mathbf{w}_{1,k}] \in$
 466 $\mathbb{R}^{n_1 \times n_1}$.

467 Using the assumption that $A_{jk} := \mathbb{E}_{z_0}[\hat{\mathbf{u}}_j(z_0)\hat{\mathbf{u}}_k^\top(z_0)] = \mathbb{E}_{z_0}[\hat{\mathbf{u}}_j(z_0)] \mathbb{E}_{z_0}[\hat{\mathbf{u}}_k^\top(z_0)] = 0$ for $j \neq k$.
 468 We arrive at the conclusion that \mathbf{w}_j is only coupled with v_j , yielding the following dynamical
 469 system:

$$\dot{w}_{2,j} = w_{2,j} \mathbf{w}_{1,j}^\top A_{jj} \mathbf{w}_{1,j}, \quad \dot{\mathbf{w}}_{1,j} = w_{2,j}^2 A_{jj} \mathbf{w}_{1,j} \quad (50)$$

470 \square

471 **E Proofs: Hierarchical Latent Tree Models (Section A.3)**

472 **E.1 Lemmas**

473 **Lemma 3** (Variance Squashing). *Suppose a function $\phi : \mathbb{R} \mapsto \mathbb{R}$ is L -Lipschitz continuous: $|\phi(x) -$
 474 $\phi(y)| \leq L|x - y|$, then for $x \sim p(\cdot)$, we have:*

$$\mathbb{V}_p[\phi(x)] \leq L^2 \mathbb{V}_p[x] \quad (51)$$

475 *Proof.* Suppose $x, y \sim p(\cdot)$ are independent samples and $\mu_\phi := \mathbb{E}[\phi(x)]$. Note that $\mathbb{V}[\phi(x)]$ can be
 476 written as the following:

$$\begin{aligned} \mathbb{E}[|\phi(x) - \phi(y)|^2] &= \frac{1}{2} \mathbb{E}[|(\phi(x) - \mu_\phi) - (\phi(y) - \mu_\phi)|^2] \\ &= \mathbb{E}[|\phi(x) - \mu_\phi|^2] + \mathbb{E}[|\phi(y) - \mu_\phi|^2] - 2\mathbb{E}[(\phi(x) - \mu_\phi)(\phi(y) - \mu_\phi)] \\ &= 2\mathbb{V}_p[\phi(x)] \end{aligned} \quad (52)$$

Symbol	Definition	Size	Description
$\mathcal{N}_l, \mathcal{Z}_l$			The set of all nodes and all latent variables at layer l .
$\mathcal{N}_\mu, \mathcal{N}_\mu^{\text{ch}}$			Nodes corresponding to latent variable z_μ . $\mathcal{N}_\mu^{\text{ch}}$ are children under \mathcal{N}_μ .
m_μ			Number of possible categorical values taken by z_μ . $0 \leq z_\mu < m_\mu$.
$\mathbf{0}_\mu, \mathbf{1}_\mu$		m_μ	All-one and all-zero vectors.
$P_{\mu\nu}$	$[\mathbb{P}(z_\nu z_\mu)]$	$m_\mu \times m_\nu$	The top-down transition probability from z_μ to z_ν .
$\rho_{\mu\nu}$	$2\mathbb{P}(z_\nu=1 z_\mu=1) - 1$	scalar in $[-1, 1]$	Polarity of the transitional probability in the binary case.
P_0	$\text{diag}[\mathbb{P}(z_0)]$	$m_0 \times m_0$	The diagonal matrix of probability of z_0 taking different values.
$v_j(z_\mu)$	$\mathbb{E}_z[f_j z_\mu]$	scalar	Expectation of activation f_j given z_μ (z_μ 's descendants are marginalized).
\mathbf{v}_j	$[v_j(z_\mu)]$	m_μ	Vector form of $v_j(z_\mu)$.
$\mathbf{f}_\mu, \mathbf{f}_{\mathcal{N}_\mu^{\text{ch}}}$	$[f_j]_{j \in \mathcal{N}_\mu}, [f_k]_{k \in \mathcal{N}_\mu^{\text{ch}}}$	$ \mathcal{N}_\mu , \mathcal{N}_\mu^{\text{ch}} $	Activations for all nodes $j \in \mathcal{N}_\mu$ and for the children of \mathcal{N}_μ
$\mathbf{v}_{0k}, \mathbf{V}_{0, \mathcal{N}_\mu^{\text{ch}}}$	$[\mathbb{E}_z[f_k z_0]], [\mathbf{v}_{0k}]_{k \in \mathcal{N}_\mu^{\text{ch}}}$	$m_0, m_0 \times \mathcal{N}_\mu^{\text{ch}} $	Expected activation conditioned on z_0
s_k	$\frac{1}{2}(v_k(1) - v_k(0))$	scalar	Discrepancy of node k w.r.t its latent variable $z_{\nu(k)}$.
$\mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}}$	$[\rho_{\mu\nu(k)} s_k]_{k \in \mathcal{N}_\mu^{\text{ch}}}$	$ \mathcal{N}_\mu^{\text{ch}} $	Child selectivity vector in the binary case.

Table 2: Extended notation in HLTM.

Therefore we have:

$$\mathbb{V}_p[\phi(x)] = \frac{1}{2} \mathbb{E} [|\phi(x) - \phi(y)|^2] \leq \frac{L^2}{2} \mathbb{E} [|x - y|^2] = L^2 \mathbb{V}_p[x] \quad (53)$$

□

Lemma 4 (Sharpened Jensen's inequality [24]). *If function ϕ is twice differentiable, and $x \sim p(\cdot)$, then we have:*

$$\frac{1}{2} \mathbb{V}[x] \inf \phi'' \leq \mathbb{E}[\phi(x)] - \phi(\mathbb{E}[x]) \leq \frac{1}{2} \mathbb{V}[x] \sup \phi'' \quad (54)$$

Lemma 5 (Sharpened Jensen's inequality for ReLU activation). *For ReLU activation $\psi(x) := \max(x, 0)$ and $x \sim p(\cdot)$, we have:*

$$0 \leq \mathbb{E}[\psi(x)] - \psi(\mathbb{E}[x]) \leq \sqrt{\mathbb{V}_p[x]} \quad (55)$$

Proof. Since ψ is a convex function, by Jensen's inequality we have $\mathbb{E}[\psi(x)] - \psi(\mathbb{E}[x]) \geq 0$. For the other side, let $\mu := \mathbb{E}_p[x]$ and we have (note that for ReLU, $\psi(x) - \psi(\mu) \leq |x - \mu|$):

$$\mathbb{E}[\psi(x)] - \psi(\mathbb{E}[x]) = \int (\psi(x) - \psi(\mu)) p(x) dx \quad (56)$$

$$\leq \int |x - \mu| p(x) dx \quad (57)$$

$$\leq \left(\int |x - \mu|^2 p(x) dx \right)^{1/2} \left(\int p(x) dx \right)^{1/2} \quad (58)$$

$$= \sqrt{\mathbb{V}_p[x]} \quad (59)$$

where the last inequality is due to Cauchy-Schwarz. □

E.2 More general assumption of conditional distribution in HLTM

We consider a more general case of HLTM where each z_μ is categorical: $0 \leq z_\mu < m_\mu$. For convenience, we define the following symbols for $k \in \mathcal{N}_\mu^{\text{ch}}$ (note that $|\mathcal{N}_\mu^{\text{ch}}| = \mathcal{N}_\mu^{\text{ch}}$ is the number of the children of the node set \mathcal{N}_μ):

$$\mathbf{v}_{\mu k} := \mathbb{E}_z[f_k|z_\mu] = P_{\mu\nu(k)} \mathbf{v}_k \in \mathbb{R}^{m_\mu} \quad (60)$$

$$\mathbf{V}_{\mu, \mathcal{N}_\mu^{\text{ch}}} := [\mathbf{v}_{\mu k}]_{k \in \mathcal{N}_\mu^{\text{ch}}} \quad (61)$$

$$\tilde{\mathbf{v}}_j := \left[\mathbb{E}_z[\tilde{f}_j|z_\mu] \right] = \mathbf{V}_{\mu, \mathcal{N}_\mu^{\text{ch}}} \mathbf{w}_j \in \mathbb{R}^{m_\mu} \quad (62)$$

As an extension of binary symmetric HLTM, we make an assumption for the transitional probability:

491 **Assumption 1.** For $\mu \in \mathcal{Z}_l$ and $\nu \in \mathcal{Z}_{l-1}$, the transitional probability matrix $P_{\mu\nu} := [\mathbb{P}(z_\nu|z_\mu)]$
 492 has decomposition $P_{\mu\nu} = \frac{1}{m_\nu} \mathbf{1}_\mu \mathbf{1}_\nu^\top + C_{\mu\nu}$ where $C_{\mu\nu} \mathbf{1}_\nu = \mathbf{0}_\mu$ and $\mathbf{1}_\mu^\top C_{\mu\nu} = \mathbf{0}_\nu$.

493 Note that $C_{\mu\nu} \mathbf{1} = \mathbf{0}$ is obvious due to the property of conditional probability. The real condition is
 494 $\mathbf{1}_\mu^\top C_{\mu\nu} = \mathbf{0}_\nu$. If $m_\mu = m_\nu$, then $P_{\mu\nu}$ is a square matrix and Assumption 1 is equivalent to $P_{\mu\nu}$ is
 495 double-stochastic. Assumption 1 makes computation of $P_{\mu\nu}$ easy for any z_μ and z_ν .

496 **Lemma 6** (Transition Probability). *If Assumption 1 holds, then for $\mu \in \mathcal{Z}_l$, $\nu \in \mathcal{Z}_{l-1}$ and $\alpha \in$
 497 \mathcal{Z}_{l-2} , we have:*

$$P_{\mu\alpha} = P_{\mu\nu} P_{\nu\alpha} = \frac{1}{m_\alpha} \mathbf{1}_\mu \mathbf{1}_\alpha^\top + C_{\mu\nu} C_{\nu\alpha} \quad (63)$$

498 In general, for any $\mu \in \mathcal{N}_{l_1}$ and $\alpha \in \mathcal{N}_{l_2}$ with $l_1 > l_2$, we have:

$$P_{\mu\alpha} = \frac{1}{m_\alpha} \mathbf{1}_\mu \mathbf{1}_\alpha^\top + \prod_{\mu, \dots, \xi, \zeta, \dots, \alpha} C_{\xi\zeta} \quad (64)$$

499 *Proof.* Using Assumption 1, we have

$$P_{\mu\alpha} = P_{\mu\nu} P_{\nu\alpha} \quad (65)$$

$$= \left(\frac{1}{m_\nu} \mathbf{1}_\mu \mathbf{1}_\nu^\top + C_{\mu\nu} \right) \left(\frac{1}{m_\alpha} \mathbf{1}_\nu \mathbf{1}_\alpha^\top + C_{\nu\alpha} \right) \quad (66)$$

500 since $\mathbf{1}_\nu^\top \mathbf{1}_\nu = m_\nu$, $C_{\mu\nu} \mathbf{1}_\nu = \mathbf{0}_\mu$ and $\mathbf{1}_\nu^\top C_{\nu\alpha} = \mathbf{0}_\alpha$, the conclusion follows. \square

501 **Remark.** In the symmetric binary HLTM mentioned in the main text, all $C_{\mu\nu}$ can be parameterized
 502 as (here $\mathbf{q} := [-1, 1]^\top$):

$$C_{\mu\nu} = C_{\mu\nu}(\rho_{\mu\nu}) = \frac{1}{2} \begin{bmatrix} \rho_{\mu\nu} & -\rho_{\mu\nu} \\ -\rho_{\mu\nu} & \rho_{\mu\nu} \end{bmatrix} = \frac{1}{2} \rho_{\mu\nu} \mathbf{q} \mathbf{q}^\top \quad (67)$$

503 This is because $\mathbf{1}_2^\top C_{\mu\nu} = \mathbf{0}_2$ and $C_{\mu\nu} \mathbf{1}_2 = \mathbf{0}_2$ provides 4 linear constraints (1 redundant), leaving
 504 1 free parameter, which is the polarity $\rho_{\mu\nu} \in [-1, 1]$ of latent variable z_ν given its parent z_μ .
 505 Moreover, since $\mathbf{q}^\top \mathbf{q} = 2$, the parameterization is close under multiplication:

$$C(\rho_{\mu\nu}) C(\rho_{\nu\alpha}) = \frac{1}{4} \mathbf{q} \mathbf{q}^\top \mathbf{q} \mathbf{q}^\top \rho_{\mu\nu} \rho_{\nu\alpha} = \frac{1}{2} \mathbf{q} \mathbf{q}^\top \rho_{\mu\nu} \rho_{\nu\alpha} = C(\rho_{\mu\nu} \rho_{\nu\alpha}) \quad (68)$$

506 E.3 Theorem 6

507 *Proof.* First note that for each node $k \in \mathcal{N}_\mu^{\text{ch}}$:

$$\mathbf{v}_{0k} := [\mathbb{E}_z [f_k | z_0]] = \sum_{z_\nu} \mathbb{E}_z [f_k | z_\nu] \mathbb{P}(z_\nu | z_0) = P_{0\nu} \mathbf{v}_k \quad (69)$$

$$= \left(\frac{1}{m_\nu} \mathbf{1}_0 \mathbf{1}_\nu^\top + C_{0\nu} \right) \mathbf{v}_k \quad (70)$$

$$= \frac{1}{m_\nu} \mathbf{1}_0 \mathbf{1}_\nu^\top \mathbf{v}_k + C_{0\nu} \mathbf{v}_k \quad (71)$$

508 Note that $\mathbf{1}_0 \mathbf{1}_\nu^\top \mathbf{v}_k$ is a constant regarding to change of z_0 . So we could remove it when computing
 509 covariance operator. On the other hand, for a categorical distribution

$$(\mathbb{P}(z_0 = 0), u(0)), \quad (\mathbb{P}(z_0 = 1), u(1)), \quad \dots, \quad (\mathbb{P}(z_0 = m_0 - 1), u(m_0 - 1))$$

510 With $P_0 := \text{diag}[\mathbb{P}(z_0)]$, the mean is $\mathbb{E}_{z_0} [u] = \mathbf{1}^\top P_0 \mathbf{u}$ and its covariance can be written as (here
 511 $\mathbf{1} = \mathbf{1}_0$):

$$\mathbb{V}_{z_0} [u] = (\mathbf{u} - \mathbf{1} \mathbf{1}^\top P_0 \mathbf{u})^\top P_0 (\mathbf{u} - \mathbf{1} \mathbf{1}^\top P_0 \mathbf{u}) = \mathbf{u}^\top (P_0 - P_0 \mathbf{1} \mathbf{1}^\top P_0) \mathbf{u} \quad (72)$$

512 Note that each column of $V_{0, \mathcal{N}_\mu^{\text{ch}}}$ is \mathbf{v}_{0k} . Setting $\mathbf{u} = \mathbf{v}_{0k}$ and we have:

$$\mathbb{V}_{z_0} [\mathbb{E}_z [\mathbf{f}_{\mathcal{N}_\mu^{\text{ch}}} | z_0]] = V_{0, \mathcal{N}_\mu^{\text{ch}}}^\top (P_0 - P_0 \mathbf{1} \mathbf{1}^\top P_0) V_{0, \mathcal{N}_\mu^{\text{ch}}} \quad (73)$$

513 Note that $\mathbb{E}_{z_0} [\mathbf{v}_{0k}] = \frac{1}{m_\nu} \mathbf{1}_\nu^\top \mathbf{v}_k + \mathbb{E}_{z_0} [C_{0\nu} \mathbf{v}_k]$, since $\mathbf{1}^\top P_0 \mathbf{1} = 1$. With some computation, we could
 514 see $\text{Cov}_{z_0} [\mathbf{v}_{0k}, \mathbf{v}_{0k'}] = \text{Cov}_{z_0} [C_{0\nu(k)} \mathbf{v}_k, C_{0\nu(k')} \mathbf{v}_{k'}]$.

515 The equation above can be applied for any cardinality of latent variables. In the binary symmetric
 516 case, we have (note here we define $\rho_0 := \mathbb{P}(z_0 = 1) - \mathbb{P}(z_0 = 0)$, and $\mathbf{q} := [-1, 1]^\top$):

$$P_0 - P_0 \mathbf{1} \mathbf{1}^\top P_0 = \frac{1}{4} (1 - \rho_0^2) \mathbf{q} \mathbf{q}^\top \quad (74)$$

517 Note that in the binary symmetric case, according to remarks in Lemma 6, all $C_{\mu\nu} = \frac{1}{2} \rho_{\mu\nu} \mathbf{q} \mathbf{q}^\top$ and
 518 we could compute $C_{0\nu} \mathbf{v}_k$:

$$C_{0\nu} \mathbf{v}_k = \frac{1}{2} \rho_{0\nu} \mathbf{q} \mathbf{q}^\top \mathbf{v}_k = \rho_{0\nu} \frac{1}{2} (v_k(1) - v_k(0)) \mathbf{q} = \rho_{0\nu} s_k \mathbf{q} \quad (75)$$

519 where according to Eqn. 68, we have:

$$\rho_{0\nu} := \prod_{0, \dots, \alpha, \beta, \dots, \nu} \rho_{\alpha\beta} \quad (76)$$

520 and the covariance between node k and k' can be computed as:

$$\text{Cov}_{z_0}[\mathbf{v}_{0k}, \mathbf{v}_{0k'}] = \text{Cov}_{z_0}[C_{0\nu(k)} \mathbf{v}_k, C_{0\nu(k')} \mathbf{v}_{k'}] \quad (77)$$

$$= \rho_{0\nu(k)} \rho_{0\nu(k')} s_k s_{k'} \frac{1}{4} \mathbf{q}^\top \mathbf{q} \mathbf{q}^\top \mathbf{q} (1 - \rho_0^2) \quad (78)$$

$$= \rho_{0\nu(k)} \rho_{0\nu(k')} s_k s_{k'} (1 - \rho_0^2) \quad (79)$$

$$= \rho_{0\mu}^2 \rho_{\mu\nu(k)} \rho_{\mu\nu(k')} s_k s_{k'} (1 - \rho_0^2) \quad (80)$$

521 The last equality is due to the fact that due to tree structure, the path from z_0 to all child nodes in
 522 $\mathcal{N}_\mu^{\text{ch}}$ must pass z_μ .

523 Therefore we can compute the covariance operator:

$$\mathbb{V}_{z_0}[\mathbb{E}_z[\mathbf{f}_{\mathcal{N}_\mu^{\text{ch}} | z_0}]] = \rho_{0\mu}^2 (1 - \rho_0^2) \mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}} \mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}}^\top \quad (81)$$

524 When $L \rightarrow +\infty$, we have:

$$\rho_{0\nu} := \prod_{0, \dots, \alpha, \beta, \dots, \nu} \rho_{\alpha\beta} \rightarrow 0 \quad (82)$$

525 and thus the covariance becomes zero as well. \square

526 E.4 Theorem 7

527 *Proof.* According to our setting, for each node $k \in \mathcal{N}_\mu$, there exists a unique latent variable z_ν with
 528 $\nu = \nu(k)$ that corresponds to it. In the following we omit its dependency on k for brevity.

529 Since we are dealing with binary case, we define the following for convenience:

$$v_k^+ := v_k(1) \quad (83)$$

$$v_k^- := v_k(0) \quad (84)$$

$$\bar{v}_k := \frac{1}{2} (v_k^+ + v_k^-) = \frac{1}{2} (\mathbb{E}_z[f_k | z_\nu = 1] + \mathbb{E}_z[f_k | z_\nu = 0]) \quad (85)$$

$$\bar{\mathbf{v}}_{\mathcal{N}_\mu^{\text{ch}}} := [\bar{v}_k]_{k \in \mathcal{N}_\mu^{\text{ch}}} \in \mathbb{R}^{|\mathcal{N}_\mu^{\text{ch}}|} \quad (86)$$

$$s_k := \frac{1}{2} (v_k^+ - v_k^-) = \frac{1}{2} (\mathbb{E}_z[f_k | z_\nu = 1] - \mathbb{E}_z[f_k | z_\nu = 0]) \quad (87)$$

$$\mathbf{s}_{\mathcal{N}_\mu^{\text{ch}}} := [s_k]_{k \in \mathcal{N}_\mu^{\text{ch}}} \in \mathbb{R}^{|\mathcal{N}_\mu^{\text{ch}}|} \quad (88)$$

530 We also define the *sensitivity* of node k to be $\lambda_k := |(v_k^+)^2 - (v_k^-)^2|$. Intuitively, a large λ_k means
 531 that the node k is sensitive for changes of latent variable z_ν . If $\lambda_k = 0$, then the node k is invariant
 532 to latent variable z_ν .

533 We first consider pre-activation $\tilde{f}_j := \sum_k w_{jk} f_k$ and its expectation with respect to latent variable
 534 z :

$$\tilde{v}_j^+ := \mathbb{E}_z[\tilde{f}_j | z_\mu = 1], \quad \tilde{v}_j^- := \mathbb{E}_z[\tilde{f}_j | z_\mu = 0] \quad (89)$$

535 Note that for each node $k \in \mathcal{N}_\mu^{\text{ch}}$ we have:

$$v_{\mu k}^+ = \bar{v}_k + \rho_{\mu\nu} s_k, \quad v_{\mu k}^- = \bar{v}_k - \rho_{\mu\nu} s_k \quad (90)$$

536 Let $\mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}} := [a_k]_{k \in \mathcal{N}_\mu^{\text{ch}}} := [\rho_{\mu\nu} s_k]_{k \in \mathcal{N}_\mu^{\text{ch}}}$ and

$$\mathbf{u}_{\mathcal{N}_\mu^{\text{ch}}}^+ := V_{\mu, \mathcal{N}_\mu^{\text{ch}}}^+ := \left[\mathbb{E}[f_k | z_\mu = 1] \right] = \bar{\mathbf{v}}_{\mathcal{N}_\mu^{\text{ch}}} + \mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}} \quad (91)$$

$$\mathbf{u}_{\mathcal{N}_\mu^{\text{ch}}}^- := V_{\mu, \mathcal{N}_\mu^{\text{ch}}}^- := \left[\mathbb{E}[f_k | z_\mu = 0] \right] = \bar{\mathbf{v}}_{\mathcal{N}_\mu^{\text{ch}}} - \mathbf{a}_{\mathcal{N}_\mu^{\text{ch}}} \quad (92)$$

537 Then we have $\tilde{v}_j^+ = \mathbf{w}_j^\top \mathbf{u}_{\mathcal{N}_\mu^{\text{ch}}}^+$ and $\tilde{v}_j^- = \mathbf{w}_j^\top \mathbf{u}_{\mathcal{N}_\mu^{\text{ch}}}^-$.

538 Note that \mathbf{w}_j is a random variable with each entry $w_{jk} \sim \text{Uniform} \left[-\sigma_w \sqrt{\frac{3}{|\mathcal{N}_\mu^{\text{ch}}|}}, \sigma_w \sqrt{\frac{3}{|\mathcal{N}_\mu^{\text{ch}}|}} \right]$.

539 It is easy to verify that $\mathbb{E}[w_{jk}] = 0$ and $\mathbb{V}[w_{jk}] = \sigma_w^2 / |\mathcal{N}_\mu^{\text{ch}}|$. Therefore, for two dimensional
540 vector $\tilde{\mathbf{v}}_j = [\tilde{v}_j^+, \tilde{v}_j^-]^\top$, we can compute its first and second order moments: $\mathbb{E}_w[\tilde{\mathbf{v}}_j] = \mathbf{0}$ and

541 $\mathbb{V}_w[\tilde{\mathbf{v}}_j] = \frac{\sigma_w^2}{|\mathcal{N}_\mu^{\text{ch}}|} V_{\mu, \mathcal{N}_\mu^{\text{ch}}} V_{\mu, \mathcal{N}_\mu^{\text{ch}}}^\top = \frac{\sigma_w^2}{|\mathcal{N}_\mu^{\text{ch}}|} [\mathbf{u}_{\mathcal{N}_\mu^{\text{ch}}}^+, \mathbf{u}_{\mathcal{N}_\mu^{\text{ch}}}^-]^\top [\mathbf{u}_{\mathcal{N}_\mu^{\text{ch}}}^+, \mathbf{u}_{\mathcal{N}_\mu^{\text{ch}}}^-].$

542 Define the positive and negative set (note that $a_k := \rho_{\mu\nu} s_k$):

$$A_+ = \{k : a_k \geq 0\}, \quad A_- = \{k : a_k < 0\} \quad (93)$$

543 Without loss of generality, assume that $\sum_{k \in A_+} a_k^2 \geq \sum_{k \in A_-} a_k^2$. In the following, we show there
544 exists j with λ_j is greater than some positive threshold. Otherwise the proof is symmetric and we
545 can show λ_j is lower than some negative threshold.

546 When $|\mathcal{N}_\mu^{\text{ch}}|$ is large, by Central Limit Theorem, \mathbf{v} can be regarded as zero-mean 2D Gaussian
547 distribution and we have for some $c > 0$:

$$\mathbb{P} \left(\tilde{v}_j^+ \geq \frac{\sqrt{c} \sigma_w}{\sqrt{|\mathcal{N}_\mu^{\text{ch}}|}} \|\mathbf{u}_{\mathcal{N}_\mu^{\text{ch}}}^+\| \right) = \frac{1 - \text{erf}(\sqrt{c}/2)}{2} \quad (94)$$

548 Moreover, if $\mathbf{a}_l \neq \mathbf{0}$, then the following probability is also not small :

$$\mathbb{P} \left(\tilde{v}_j^+ \geq \frac{\sqrt{c} \sigma_w}{\sqrt{|\mathcal{N}_\mu^{\text{ch}}|}} \|\mathbf{u}_{\mathcal{N}_\mu^{\text{ch}}}^+\| \quad \text{and} \quad \tilde{v}_j^- < 0 \right) \quad (95)$$

549 Therefore, when $|\mathcal{N}_\mu| = O(\exp(c))$, with high probability, there exists \mathbf{w}_j so that

$$\tilde{v}_j^+ = \mathbf{w}_j^\top \mathbf{u}_{\mathcal{N}_\mu^{\text{ch}}}^+ \geq \frac{\sqrt{c} \sigma_w}{\sqrt{|\mathcal{N}_\mu^{\text{ch}}|}} \|\mathbf{u}_{\mathcal{N}_\mu^{\text{ch}}}^+\|, \quad \tilde{v}_j^- = \mathbf{w}_j^\top \mathbf{u}_{\mathcal{N}_\mu^{\text{ch}}}^- < 0 \quad (96)$$

550 Since $\bar{\mathbf{v}}_{\mathcal{N}_\mu^{\text{ch}}} \geq 0$ (all f_k are after ReLU and non-negative), this leads to:

$$\tilde{v}_j^+ \geq \frac{\sqrt{c} \sigma_w}{\sqrt{|\mathcal{N}_\mu^{\text{ch}}|}} \|\mathbf{u}_{\mathcal{N}_\mu^{\text{ch}}}^+\| \geq \frac{\sqrt{c} \sigma_w}{\sqrt{|\mathcal{N}_\mu^{\text{ch}}|}} \sqrt{\sum_{k \in A_+} a_k^2} \geq \sigma_w \sqrt{\frac{c}{2|\mathcal{N}_\mu^{\text{ch}}|} \sum_{k \in \mathcal{N}_\mu^{\text{ch}}} \rho_{\mu\nu}^2 s_k^2} \quad (97)$$

551 By Jensen's inequality, we have (note that $\psi(x) := \max(x, 0)$ is the ReLU activation):

$$v_j^+ = \mathbb{E}_z[f_j | z_\mu = 1] = \mathbb{E}_z[\psi(\tilde{f}_j) | z_\mu = 1] \quad (98)$$

$$\geq \psi \left(\mathbb{E}_z[\tilde{f}_j | z_\mu = 1] \right) = \psi(\tilde{v}_j^+) \geq \sigma_w \sqrt{\frac{c}{2|\mathcal{N}_\mu^{\text{ch}}|} \sum_{k \in \mathcal{N}_\mu^{\text{ch}}} \rho_{\mu\nu}^2 s_k^2} \quad (99)$$

On the other hand, we also want to compute $v_j^- := \mathbb{E}_z [f_j | z_\mu = 0]$ using sharpened Jensen's inequality (Lemma 5). For this we need to compute the conditional covariance $\mathbb{V}_z[\tilde{f}_j | z_\mu]$:

$$\mathbb{V}_z[\tilde{f}_j | z_\mu] \stackrel{\textcircled{2}}{=} \sum_k w_{jk}^2 \mathbb{V}_z[f_k | z_\mu] \stackrel{\textcircled{3}}{\leq} \frac{3\sigma_w^2}{|\mathcal{N}_\mu^{\text{ch}}|} \sum_k \mathbb{V}_z[f_k | z_\mu] \quad (100)$$

$$= \frac{3\sigma_w^2}{|\mathcal{N}_\mu^{\text{ch}}|} \sum_k (\mathbb{E}_{z_\nu | z_\mu} [\mathbb{V}[f_k | z_\nu]] + \mathbb{V}_{z_\nu | z_\mu} [\mathbb{E}_z [f_k | z_\nu]]) \quad (101)$$

$$\leq 3\sigma_w^2 \left(\sigma_l^2 + \frac{1}{|\mathcal{N}_\mu^{\text{ch}}|} \sum_k \mathbb{V}_{z_\nu | z_\mu} [\mathbb{E}_z [f_k | z_\nu]] \right) \quad (102)$$

Note that $\textcircled{2}$ is due to conditional independence: f_k as the computed activation, only depends on latent variable z_ν and its descendants. Given z_μ , all z_ν and their respective descendants are independent of each other and so does f_k . $\textcircled{3}$ is due to the fact that each w_{jk} are sampled from uniform distribution and $|w_{jk}| \leq \sigma_w \sqrt{\frac{3}{|\mathcal{N}_\mu^{\text{ch}}|}}$.

Here $\mathbb{V}_{z_\nu | z_\mu} [\mathbb{E}_z [f_k | z_\nu]] = s_k^2(1 - \rho_{\mu\nu}^2)$ can be computed analytically. It is the variance of a binomial distribution: with probability $\frac{1}{2}(1 + \rho_{\mu\nu})$ we get v_k^+ otherwise get v_k^- . Therefore, we finally have:

$$\mathbb{V}_z[\tilde{f}_j | z_\mu] \leq 3\sigma_w^2 \left(\sigma_l^2 + \frac{1}{|\mathcal{N}_\mu^{\text{ch}}|} \sum_k s_k^2(1 - \rho_{\mu\nu}^2) \right) \quad (103)$$

As a side note, using Lemma 3, since ReLU function ψ has Lipschitz constant ≤ 1 (empirically it is smaller), we know that:

$$\mathbb{V}_z[f_j | z_\mu] \leq 3\sigma_w^2 \left(\sigma_l^2 + \frac{1}{|\mathcal{N}_\mu^{\text{ch}}|} \sum_k s_k^2(1 - \rho_{\mu\nu}^2) \right) \quad (104)$$

Finally using Lemma 5 and $\tilde{v}_j^- < 0$, we have:

$$v_j^- = \mathbb{E}_z [f_j | z_\mu = 0] = \mathbb{E}_z [\psi(\tilde{f}_j) | z_\mu = 0] \quad (105)$$

$$\leq \psi \left(\mathbb{E}_z [\tilde{f}_j | z_\mu = 0] \right) + \sqrt{\mathbb{V}_z[\tilde{f}_j | z_\mu = 0]} \quad (106)$$

$$= \sqrt{\mathbb{V}_z[\tilde{f}_j | z_\mu = 0]} \quad (107)$$

$$\leq \sigma_w \sqrt{3\sigma_l^2 + \frac{3}{|\mathcal{N}_\mu^{\text{ch}}|} \sum_k s_k^2(1 - \rho_{\mu\nu}^2)} \quad (108)$$

Combining Eqn. 99 and Eqn. 108, we have a bound for λ_j :

$$\lambda_j = (v_j^+)^2 - (v_j^-)^2 \geq 3\sigma_w^2 \left[\frac{1}{|\mathcal{N}_\mu^{\text{ch}}|} \sum_k s_k^2 \left(\frac{c+6}{6} \rho_{\mu\nu}^2 - 1 \right) - \sigma_l^2 \right] \quad (109)$$

□

F Proofs: The Analysis of BYOL 4

F.1 Derivation of BYOL gradient

Note that for BYOL, we have:

$$\text{vec} \left(\frac{\partial r}{\partial W_l} \right) = K_l(\mathbf{x}_1; \mathcal{W}) [K_l^\top(\mathbf{x}_1; \mathcal{W}) \text{vec}(W_l) - K_l^\top(\mathbf{x}_2; \mathcal{W}') \text{vec}(W'_l)] \quad (110)$$

under large batchsize, we have (note that we omit \mathcal{W} for any term that depends on \mathcal{W} , but make dependence of \mathcal{W}' explicit in the math expression):

$$\text{vec} \left(\frac{\partial r}{\partial W_l} \right) = \mathbb{E}_{\mathbf{x} \sim p(\cdot)} [\mathbb{E}_{\mathbf{x}' \sim p_{\text{aug}}(\cdot | \mathbf{x})} [K_l(\mathbf{x}') K_l^\top(\mathbf{x}')] \text{vec}(W_l) - \bar{K}_l(\mathbf{x}) \bar{K}_l^\top(\mathbf{x}; \mathcal{W}') \text{vec}(W'_l)]$$

For brevity, we write $\mathbb{E}_{\mathbf{x}}[\cdot] := \mathbb{E}_{\mathbf{x} \sim p(\cdot)}[\cdot]$ and $\mathbb{E}_{\mathbf{x}'}[\cdot] := \mathbb{E}_{\mathbf{x}' \sim p_{\text{aug}}(\cdot|\mathbf{x})}[\cdot]$. Similar for \mathbb{V} . And the equation above can be written as:

$$\text{vec}\left(\frac{\partial r}{\partial W_l}\right) = \mathbb{E}_{\mathbf{x}}\{\mathbb{V}_{\mathbf{x}'}[K_l(\mathbf{x}')]\}\text{vec}(W_l) \quad (111)$$

$$+ \mathbb{E}_{\mathbf{x}}\{\bar{K}_l(\mathbf{x})[\bar{K}_l^\top(\mathbf{x})\text{vec}(W_l) - \bar{K}_l^\top(\mathbf{x}; \mathcal{W}')\text{vec}(W'_l)]\} \quad (112)$$

In terms of weight update by gradient descent, since $\Delta W_l = -\frac{\partial r}{\partial W_l}$, we have:

$$\text{vec}(\Delta W_l) = -\mathbb{E}_{\mathbf{x}}\{\mathbb{V}_{\mathbf{x}'}[K_l(\mathbf{x}')]\}\text{vec}(W_l) \quad (113)$$

$$- \mathbb{E}_{\mathbf{x}}\{\bar{K}_l(\mathbf{x})[\bar{K}_l^\top(\mathbf{x})\text{vec}(W_l) - \bar{K}_l^\top(\mathbf{x}; \mathcal{W}')\text{vec}(W'_l)]\} \quad (114)$$

If we consider the special case $\mathcal{W} = \mathcal{W}'$, then the last two terms cancelled out, yielding:

$$\text{vec}(\Delta W_l)_{\text{sym}} = -\mathbb{E}_{\mathbf{x}}\{\mathbb{V}_{\mathbf{x}'}[K_l(\mathbf{x}')]\}\text{vec}(W_l) \quad (115)$$

And the general update (Eqn. 117) can be written as:

$$\text{vec}(\Delta W_l) = \text{vec}(\Delta W_l)_{\text{sym}} \quad (116)$$

$$- \mathbb{E}_{\mathbf{x}}\{\bar{K}_l(\mathbf{x})[\bar{K}_l^\top(\mathbf{x})\text{vec}(W_l) - \bar{K}_l^\top(\mathbf{x}; \mathcal{W}')\text{vec}(W'_l)]\} \quad (117)$$

F.2 Theorem 4

Proof. When BN is present, Eqn. 110 needs to be corrected with an additional term, $\widetilde{\frac{\partial r}{\partial W_l}} := \frac{\partial r}{\partial W_l} - \delta W_l^{\text{BN}}$, where δW_l^{BN} is defined as follows:

$$\delta W_l^{\text{BN}} := \frac{1}{|B|} \sum_{i \in B} D_l^i \bar{\mathbf{g}}_l \mathbf{f}_{l-1}^{i\top} \quad (118)$$

From the proof of Theorem 1, we know that for each sample $i \in B$:

$$D_l^i \mathbf{g}_l^i = J_l^{i\top} [J_l^i W_l \mathbf{f}_{l-1}^i - J_l^i(\mathcal{W}') W'_l \mathbf{f}_{l-1}^i(\mathcal{W}')] \quad (119)$$

Since the network is linear from layer l to the topmost layer L , we have $D_l^i = \bar{D}_l$. Since the only input dependent part in J_l^i is the gating function between the current layer l and the topmost layer L , for linear network the gating is always 1 and thus $\bar{J}_l = J_l^i$ and is independent of input data. We now have (note that we omit \mathcal{W} for any terms that are dependent on \mathcal{W} , but will write \mathcal{W}' explicitly for terms that are depend on \mathcal{W}'):

$$\delta W_l^{\text{BN}} := \frac{1}{|B|} \sum_{i \in B} D_l^i \bar{\mathbf{g}}_l \mathbf{f}_{l-1}^{i\top} = -\bar{D}_l \bar{\mathbf{g}}_l \bar{\mathbf{f}}_{l-1}^\top \quad (120)$$

$$= \bar{J}_l^\top [\bar{J}_l W_l \bar{\mathbf{f}}_{l-1} - \bar{J}_l(\mathcal{W}') W'_l \bar{\mathbf{f}}_{l-1}(\mathcal{W}')] \bar{\mathbf{f}}_{l-1}^\top \quad (121)$$

Therefore we have:

$$\text{vec}(\delta W_l^{\text{BN}}) = (\bar{\mathbf{f}}_{l-1} \otimes \bar{J}_l^\top) [(\bar{\mathbf{f}}_{l-1} \otimes \bar{J}_l^\top)\text{vec}(W_l) - (\bar{\mathbf{f}}_{l-1}(\mathcal{W}') \otimes \bar{J}_l^\top(\mathcal{W}'))\text{vec}(W'_l)] \quad (122)$$

Note that by assumption, since \bar{J}_l doesn't depend on the input data, we have

$$\bar{\mathbf{f}}_{l-1} \otimes \bar{J}_l^\top = \mathbb{E}_B[\mathbf{f}_{l-1}] \otimes \bar{J}_l^\top = \mathbb{E}_B[\mathbf{f}_{l-1} \otimes \bar{J}_l^\top] \quad (123)$$

Taking large batchsize limits and notice that the batch B could contain any augmented data generated from independent samples from $p(\cdot)$, we have:

$$\text{vec}(\delta W_l^{\text{BN}}) = \mathbb{E}_{\mathbf{x}, \mathbf{x}'}[K_l(\mathbf{x}')] \mathbb{E}_{\mathbf{x}, \mathbf{x}'}[K_l^\top(\mathbf{x}')] \text{vec}(W_l) \quad (124)$$

$$- \mathbb{E}_{\mathbf{x}, \mathbf{x}'}[K_l(\mathbf{x}')] \mathbb{E}_{\mathbf{x}, \mathbf{x}'}[K_l^\top(\mathbf{x}'; \mathcal{W}')] \text{vec}(W'_l) \quad (125)$$

An important thing is that the expectation is taking over $\mathbf{x} \sim p(\mathbf{x})$ and $\mathbf{x}' \sim p_{\text{aug}}(\cdot|\mathbf{x})$. Intuitively, this is because $\bar{\mathbf{f}}_{l-1}$ and $\bar{\mathbf{g}}_l$ are averages over the entire batch, which has both intra-sample and inter-sample variation.

591 With augment-mean connection $\bar{K}_l(\mathbf{x})$ we could write:

$$\begin{aligned} \text{vec}(\delta W_l^{\text{BN}}) &= \mathbb{E}_{\mathbf{x}} [\bar{K}_l(\mathbf{x})] \mathbb{E}_{\mathbf{x}} [\bar{K}_l^\top(\mathbf{x})] \text{vec}(W_l) - \mathbb{E}_{\mathbf{x}} [\bar{K}_l(\mathbf{x})] \mathbb{E}_{\mathbf{x}} [\bar{K}_l^\top(\mathbf{x}; \mathcal{W}')] \text{vec}(W_l') \\ &= \mathbb{E}_{\mathbf{x}} [\bar{K}_l(\mathbf{x})] \{ \mathbb{E}_{\mathbf{x}} [\bar{K}_l^\top(\mathbf{x})] \text{vec}(W_l) - \mathbb{E}_{\mathbf{x}} [\bar{K}_l^\top(\mathbf{x}; \mathcal{W}')] \text{vec}(W_l') \} \end{aligned} \quad (126)$$

592 Plug in $\delta W_{l, \text{BN}}$ into Eqn. 117 and we have corrected gradient for BYOL:

$$\text{vec}\left(\frac{\partial r}{\partial W_l}\right) = \text{vec}\left(\frac{\partial r}{\partial W_l}\right) - \text{vec}(\delta W_l^{\text{BN}}) \quad (127)$$

$$= \mathbb{E}_{\mathbf{x}} [\mathbb{V}_{\mathbf{x}' \sim p_{\text{aug}}(\cdot|\mathbf{x})} [K_l(\mathbf{x}')]] \text{vec}(W_l) + \mathbb{V}_{\mathbf{x}} [\bar{K}_l(\mathbf{x})] \text{vec}(W_l) \quad (128)$$

$$- \text{Cov}_{\mathbf{x}} [\bar{K}_l(\mathbf{x}), \bar{K}_l(\mathbf{x}; \mathcal{W}')] \text{vec}(W_l') \quad (129)$$

593 And the weight update $\widetilde{\Delta W}_l = \Delta W_l + \delta W_l^{\text{BN}}$ is:

$$\text{vec}\left(\widetilde{\Delta W}_l\right) = -\mathbb{E}_{\mathbf{x}} [\mathbb{V}_{\mathbf{x}' \sim p_{\text{aug}}(\cdot|\mathbf{x})} [K_l(\mathbf{x}')]] \text{vec}(W_l) - \mathbb{V}_{\mathbf{x}} [\bar{K}_l(\mathbf{x})] \text{vec}(W_l) \quad (130)$$

$$+ \text{Cov}_{\mathbf{x}} [\bar{K}_l(\mathbf{x}), \bar{K}_l(\mathbf{x}; \mathcal{W}')] \text{vec}(W_l') \quad (131)$$

594 Using Eqn. 115, we have:

$$\text{vec}\left(\widetilde{\Delta W}_l\right) = \text{vec}(\Delta W_l)_{\text{sym}} \quad (132)$$

$$- \mathbb{V}_{\mathbf{x}} [\bar{K}_l(\mathbf{x})] \text{vec}(W_l) + \text{Cov}_{\mathbf{x}} [\bar{K}_l(\mathbf{x}), \bar{K}_l(\mathbf{x}; \mathcal{W}')] \text{vec}(W_l') \quad (133)$$

595

□

596 F.3 Analysis on EMA

597 Consider the following discrete dynamics of a weight vector $\mathbf{w}(t)$:

$$\mathbf{w}(t+1) - \mathbf{w}(t) = \alpha [-\mathbf{w}(t) + (1-\lambda)\mathbf{w}_{\text{ema}}(t)] \quad (134)$$

598 where α is the learning rate, $\mathbf{w}_{\text{ema}}(t+1) = \gamma_{\text{ema}}\mathbf{w}_{\text{ema}}(t) + (1-\gamma_{\text{ema}})\mathbf{w}(t)$ is the exponential
599 moving average of $\mathbf{w}(t)$. For convenience, we use $\eta := 1 - \gamma_{\text{ema}}$.

600 Since it is a recurrence equation, we apply z -transform on the temporal domain, where $\mathbf{w}(z) :=$
601 $\mathbb{Z}[\mathbf{w}(t)] = \sum_{t=0}^{+\infty} \mathbf{w}(t)z^{-t}$. This leads to:

$$z(\mathbf{w}(z) - \mathbf{w}(0)) = \mathbf{w}(z) - \alpha(\mathbf{w}(z) - \mathbb{Z}[\mathbf{w}_{\text{ema}}(t)](1-\lambda)) \quad (135)$$

602 Note that for $\mathbf{w}_{\text{ema}}(t)$ we have:

$$z(\mathbf{w}_{\text{ema}}(z) - \mathbf{w}_{\text{ema}}(0)) = (1-\eta)\mathbf{w}_{\text{ema}}(z) + \eta\mathbf{w}(z) \quad (136)$$

603 If we set $\mathbf{w}_{\text{ema}}(0) = 0$, i.e., the target network is all zero at the beginning, then it gives $\mathbf{w}_{\text{ema}}(z) =$
604 $\frac{\eta}{z-1+\eta}\mathbf{w}(z)$. Plugging it back to Eqn. 135 and we have:

$$z(\mathbf{w}(z) - \mathbf{w}(0)) = \mathbf{w}(z) - \alpha\mathbf{w}(z) \left(1 - \frac{\eta}{z-1+\eta}(1-\lambda)\right) \quad (137)$$

605 And then we could solve $\mathbf{w}(z)$:

$$\mathbf{w}(z) = \frac{z(z-1+\eta)}{(z-1)^2 + (\eta+\alpha)(z-1) + \alpha\eta\lambda} \mathbf{w}(0) \quad (138)$$

606 Note that the denominator has two roots z_1 and z_2 :

$$z_{1,2} = 1 - \frac{1}{2} \left(\eta + \alpha \pm \sqrt{(\eta+\alpha)^2 - 4\alpha\eta\lambda} \right) \quad (139)$$

607 and $\mathbf{w}(z)$ can be written as

$$\mathbf{w}(z) = \frac{z(z-1+\eta)}{(z-z_1)(z-z_2)} \mathbf{w}(0) \quad (140)$$

Without loss of generality, let $z_1 < z_2$. The larger root $z_2 > 1$ when $\lambda < 0$, so the zero ($z = 1 - \eta = \gamma_{\text{ema}}$) in the nominator won't cancel out the pole at z_2 . And we have:

$$\frac{z}{(z - z_1)(z - z_2)} = \frac{z}{z_2 - z_1} \frac{(z - z_1) - (z - z_2)}{(z - z_1)(z - z_2)} \quad (141)$$

$$= \frac{z}{z_2 - z_1} \left(\frac{1}{z - z_2} - \frac{1}{z - z_1} \right) \quad (142)$$

$$= \frac{1}{z_2 - z_1} \left(\frac{1}{1 - z_2 z^{-1}} - \frac{1}{1 - z_1 z^{-1}} \right) \quad (143)$$

where $1/(1 - z_2 z^{-1})$ corresponds to a power series z_2^t in the temporal domain. Therefore, we could see $w(t)$ has exponential growth due to $z_2 > 1$.

Table 3: Normalized Correlation between the topmost latent variables in binary HLTM and topmost nodes in deep ReLU networks ($L = 5$) go up when training with SimCLR with NCE loss. We see higher correlations at both initialization and end of training, with more over-parameterization (**Left**: $|\mathcal{N}_\mu| = 2$, **Right**: $|\mathcal{N}_\mu| = 5$).

$\rho_{\mu\nu}$	Initial	1 epoch	20 epochs	$\rho_{\mu\nu}$	Initial	1 epoch	20 epochs
$\sim \text{Uniform}[0.7, 1]$	0.51	0.69	0.76	$\sim \text{Uniform}[0.7, 1]$	0.60	0.72	0.88
$\sim \text{Uniform}[0.8, 1]$	0.65	0.76	0.79	$\sim \text{Uniform}[0.8, 1]$	0.73	0.80	0.87
$\sim \text{Uniform}[0.9, 1]$	0.81	0.85	0.86	$\sim \text{Uniform}[0.9, 1]$	0.87	0.90	0.95

Table 4: Top-1 STL performance with different combination of predictor (P), EMA and BatchNorm using BYOL. When EMA is on, we use $\gamma_{\text{ema}} = 0.996$. Batch size is 128 and all experiments run on 5 seeds.

-	EMA	BN	EMA, BN	P	P, EMA	P, BN	P, EMA, BN
38.7 ± 0.6	39.3 ± 0.9	33.0 ± 0.3	32.8 ± 0.5	39.5 ± 3.1	44.4 ± 3.2	63.6 ± 1.06	78.1 ± 0.3

Table 5: Top-1 STL performance using different BatchNorm components in the predictor and the projector of BYOL ($\gamma_{\text{ema}} = 0.996$, 100 epochs). There is no affine part. “ μ ” = zero-mean normalization only, “ μ, σ ” = BN without affine, “ μ, σ^\dagger ” = normalization with mean and std but only backpropagating through mean. All variants with detached zero-mean normalization (in red) yield similar poor performance as no normalization.

-	μ	σ	μ, σ	μ^\dagger	σ^\dagger	μ^\dagger, σ	μ, σ^\dagger	$\mu^\dagger, \sigma^\dagger$
43.9 ± 4.2	64.8 ± 0.6	72.2 ± 0.9	78.1 ± 0.3	44.2 ± 7.0	54.2 ± 0.6	48.3 ± 2.7	76.3 ± 0.4	47.0 ± 8.1

Table 6: Top-1 performance of BYOL using reinitialization of the predictor every T epochs.

	Original BYOL	ReInit $T = 5$	ReInit $T = 10$	ReInit $T = 20$
STL-10 (100 epochs)	78.1	78.6	79.1	79.0
ImageNet (60 epochs)	60.9	61.9	62.4	62.4

G Experiments

We test our theoretical findings through experiments on STL-10 [10] and ImageNet [11]. We use a simplified linear evaluation protocol: the linear classifier is trained on frozen representations computed *without* data augmentation. This reuses pre-computed representations and accelerates evaluation by 10x.

Hierarchical Latent Tree Model (HLTM). We implement the HLTM and check whether the intermediate layers of deep ReLU networks learn the corresponding latent variables at the same layer. The degree of learning is measured by the normalized correlations between the ground truth latent variable z_μ and its best corresponding node $j \in \mathcal{N}_\mu$ at each layer. Tbl. 2 indicates this measure increases with over-parameterization and learning, consistent with the analysis in Sec. A.3.

Factors underlying BYOL performance. To test our theory, we perform an ablation study of BYOL on STL-10 by modifying three key components: predictor, EMA and BN. Tbl. 4 shows that BN and predictor are critical and EMA further improves the performance. First, without a predictor, neither BN nor EMA give good performance. A predictor without BN still doesn't work. A predictor with BN starts to show good performance (63.6%) and further adding EMA leads to the best performance (78.1%). This is consistent with our theoretical findings in Sec. 4, in which we show that using a predictor with BN yields $\delta W_t^{\text{BN}} \neq 0$ and leads to an implicit contrastive term.

To further test our understanding of the role played by BN, we fractionate BN into several sub-components: subtract by batch mean (`mean-norm`), divide by batch standard deviation (`std-norm`) and `affine`, and do ablation studies (Tbl. 5). Surprisingly, removing `affine` yields slightly better performance on STL-10 (from 78.1% to 78.7%). We also find that variants of `mean-norm` performs reasonably, while variants of *detached* `mean-norm` has similar poor performance as no normalization, supporting that centralizing backpropagated gradient leads to implicit contrastive terms (Sec. 4). Note that `std-norm` also helps, which we leave for future analysis.

We also check whether the online network requires an “optimal predictor” as suggested by recent version (v3) of BYOL. For this, we reinitialize the predictor (`ReInit`) every T epochs and compare the final performance under linear evaluation protocol. Interestingly, as shown in Tbl. 6, `ReInit` actually improves the performance a bit, compared to the original BYOL that keeps training the same predictor, which should therefore be closer to optimal. Moreover, if we shrink the initial weight range of the predictor to make $\text{Cov}_{\mathbf{x}}[\bar{K}_l(\mathbf{x}), \bar{K}_l(\mathbf{x}; \mathcal{W}')] (third term in Eqn. 7) more dominant, and reduce the learning rate, the performance further improves (See Tbl. 10 in Appendix H), thereby corroborating our analysis.$

H Further Experiments on Predictor

Based on our theoretical analysis, we try training the predictor in different ways and check whether it still works.

From the analysis in Sec. 4, we know that the reason why BYOL works is due to the dominance of $\text{Cov}_{\mathbf{x}}[\bar{K}_l(\mathbf{x}), \bar{K}_l(\mathbf{x}; \mathcal{W})]$ and its resemblance of the covariance operator $\mathbb{V}_{\mathbf{x}}[\bar{K}_l(\mathbf{x})]$, which is a PSD matrix.

The dominance should be stronger if the predictor has smaller weights than normally initialized using Xavier/Kaiming initialization. Also, $\text{Cov}_{\mathbf{x}}[\bar{K}_l(\mathbf{x}), \bar{K}_l(\mathbf{x}; \mathcal{W}')] should behave more like a PSD matrix, if the predictor’s weights are all small positive numbers and no BN is used.$

The following table justifies our theoretical findings. In particular, Tbl. 10 shows better performance in STL-10 with smaller learning rate and smaller sample range of the predictor weights.

Table 7: Training one-layer predictor with positive initial weights and no EMA ($\gamma_{\text{ema}} = 0$). All experiments run for 3 seeds.

Sample range of predictor weight	[0, 0.01]	[0, 0.02]	[0, 0.05]
With BN in predictor	62.78 ± 1.40	62.94 ± 1.03	62.31 ± 1.80
Without BN in predictor	71.95 ± 0.27	72.06 ± 0.44	71.91 ± 0.59

Table 8: Training one-layer predictor with positive initial weights with EMA ($\gamma_{\text{ema}} = 0.996$) and predictor resetting every $T = 10$ epochs. All experiments run for 3 seeds. Note that Xavier range is $\text{Uniform}[-0.15, 0.15]$ and our initialization range is much smaller than that.

Sample range of predictor weight	[0, 0.003]	[0, 0.005]	[0, 0.007]
With BN in predictor	65.61 ± 1.34	70.56 ± 0.57	70.87 ± 1.51
Without BN in predictor	74.39 ± 0.67	74.52 ± 0.63	74.80 ± 0.57

Table 9: Same as Tbl. 8 but with different weight range. All experiments run for 3 seeds.

Sample range of predictor weight	[0, 0.01]	[0, 0.02]	[0, 0.05]
With BN in predictor	68.98 ± 2.34	66.56 ± 1.70	68.41 ± 1.19
Without BN in predictor	74.66 ± 0.81	73.60 ± 0.32	74.34 ± 0.77

I Related Works

In addition to SimCLR and BYOL, many concurrent SSL frameworks exist to learn good representations for computer vision tasks. MoCo [17, 9] keeps a large bank of past representations in a queue as the slow-progressing target to train from. DeepCluster [6] and SwAV [7] learn the representations by iteratively or implicitly clustering on the current representations and improving representations using the cluster label. [2] applies similar ideas to multi-modality tasks. Contrastive Predictive

Table 10: Top-1 Performance on STL-10 with a two-layer predictor with BN and EMA ($\gamma_{\text{ema}} = 0.996$). Learning rate is smaller (0.02) and predictor weight sampled from $\text{Uniform}[-\text{range}, \text{range}]$. Note that for this, Xavier range is $\text{Uniform}[-0.097, 0.097]$ and our range is smaller.

Weight range	0.01	0.02	0.03	0.05
$T = 3$	79.48 ± 0.40	79.70 ± 0.47	79.66 ± 0.37	78.63 ± 0.10
$T = 5$	78.97 ± 0.62	79.63 ± 0.23	79.65 ± 0.37	79.01 ± 0.27
$T = 10$	79.25 ± 0.20	79.63 ± 0.22	79.58 ± 0.25	79.18 ± 0.22
$T = 20$	79.15 ± 0.66	79.91 ± 0.10	79.78 ± 0.05	79.54 ± 0.25

661 Coding [26] learns the representation by predicting the future of a sequence in the latent space with
 662 autoregressive models and InfoNCE loss. Contrastive MultiView Coding [29] uses multiple sensory
 663 channels (called “view”) of the same scene as the positive pairs and independently sampled views as
 664 the negative pairs to train the model. Recently, [23] moves beyond instance-wise pairs and proposes
 665 to use prototypes to construct training pairs that are more semantically meaningful.

666 In contrast, the literature on the (theoretical) analysis of SSL is sparse. [33] shows directly optimiz-
 667 ing the alignment/uniformity of the positive/negative pairs leads to comparable performance against
 668 contrastive loss. [3] proposes an interesting analysis of how contrastive learning aids downstream
 669 classification tasks, given assumptions about data generation. [22] analyzes how learning pretext
 670 tasks could reduce the sample complexity of the downstream task and [32] analyzes contrastive loss
 671 with multi-view data in the semi-supervised setting, with different generative models. However,
 672 they either work on linear model or treat deep models as a black-box function approximator with
 673 sufficient capacity. In comparison, we incorporate self-supervision, deep models, contrastive loss,
 674 data augmentation and generative models together into the same theoretical framework, and makes
 675 an attempt to understand how the intermediate features are emerged in modern SSL architectures
 676 with deep models that achieve SoTA.