

---

# What Architectural Inductive Bias Makes Diffusion Models Succeed? A Perspective from the Implicit Regularization of Gradient Descent

---

Anonymous Authors<sup>1</sup>

## Abstract

Diffusion and flow-based models succeed by training a neural network to predict noise or velocity from corrupted inputs. But why this training succeeds is not fully explained by the denoising objective alone, because the same objective can fail completely when the architecture of the denoiser changes. We study the role of architecture through the lens of gradient dynamics. The key property we identify is sparse connectivity: each neuron receives input from only a small subset of coordinates, a design shared across convolutional and transformer denoisers. We prove that sparse connectivity makes memorization strictly harder than in fully connected networks by shifting the implicit regularization of gradient descent away from the ambient input geometry and onto a collection of low-dimensional patches. Controlled denoising experiments corroborate this theory, and an extension to deep denoisers shows that clean-data prediction keeps internal representations lower-dimensional across layers. Our results point to a concrete mechanism of architectural inductive bias: the architecture determines the geometry on which gradient descent operates, and through this geometry it shapes which solutions training can find.

## 1. Introduction

Diffusion and flow-based generative models learn to sample from a data distribution by training a neural network to reverse a noise process (Ho et al., 2020; Song et al., 2021; Lipman et al., 2023; Liu et al., 2023). At each training step, the network receives a corrupted input—a clean sample with Gaussian noise added—and is asked to recover information about the original data. The standard objective trains the

network to predict the noise component directly, which is equivalent to learning the score function of the noisy distribution (Song et al., 2021). The denoising neural network is thus the central learnable component of these systems.

Substantial effort has gone into understanding what makes the denoising objective effective. But the objective is only part of the story. A recent experiment by Li & He (2025) makes this visible in a simple setting. They trained a fully connected network (FCN) to predict noise on a two-dimensional spiral distribution embedded in a high-dimensional space. The objective, noise schedule, and optimizer were all standard. Yet the network failed: it could not recover the spiral structure, producing diffuse samples that spread into irrelevant ambient directions. This is the exact same objective that powers state-of-the-art image generators. The variable that changes between failure and success is not the objective. It is the architecture. The denoisers used in practice are never fully connected. They are either convolutional U-Nets (Ronneberger et al., 2015; Song & Ermon, 2020; Dhariwal & Nichol, 2021) or patch-based transformers (Peebles & Xie, 2023; Bao et al., 2023; Ma et al., 2024), and when the architecture is right, the same objective can learn complex distributions in far higher dimensions.

Why does the architecture matter so much, and what property of these architectures is responsible? In this paper, we approach this question through the lens of how architectural design changes the gradient dynamics of training. We argue that the architecture is not merely a means of restricting the function class. It determines the data geometry on which gradient descent operates, and through this geometry, it shapes the implicit regularization of gradient descent.

We isolate one property shared across these architectures: each neuron receives input from only a small subset of coordinates. We formalize this as *sparse connectivity* and analyze how it reshapes the implicit regularization of gradient descent. Under this framework, we prove that sparse connectivity fundamentally changes how the implicit regularization acts: instead of operating on the full high-dimensional input, gradient descent operates on a collection of low-dimensional patches, and the geometry of this patch collection governs the strength of the resulting regularization (Theorem 4.2). We prove that when the patches are small and concentrated,

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the FoGen Workshop at ICML 2026. Do not distribute.

055 sparsely connected networks are harder to memorize on and  
 056 admit non-vacuous generalization guarantees on spherical  
 057 data where fully connected networks provably fail (Theorem  
 058 4.3), and we construct a matching failure mode when the  
 059 patch geometry is itself unstructured (Theorem 4.4).

060 We complement this theory with empirical evidence on three  
 061 fronts. First, we reproduce the fully connected failure in  
 062 a controlled denoising setting, and show that performance  
 063 degrades systematically as the local input dimension per neuron  
 064 increases toward the fully connected limit (Section 3).  
 065 Second, we verify the quantitative predictions of the stability  
 066 framework on synthetic data (Section 4). Third, we extend  
 067 this geometric perspective to deep denoisers, and show that  
 068 the output parameterization plays a role analogous to the  
 069 input interface: clean-data prediction provides an analytic  
 070 long skip that keeps the internal representations of the network  
 071 lower-dimensional across layers (Section 5). This  
 072 provides a geometric interpretation, from the perspective of  
 073 gradient dynamics, for the empirical success of clean-data  
 074 prediction reported by Li & He (2025).  
 075

076 Taken together, our results identify a concrete mechanism  
 077 through which architectural inductive bias operates in denoising  
 078 generative models. The architecture does not merely  
 079 restrict the function class. It determines the data geometry  
 080 on which gradient descent acts, and through this geometry,  
 081 it shapes which solutions training can find.  
 082

## 083 2. Related Work

084  
 085 **Diffusion and Flow generative modeling.** Diffusion-based  
 086 generative models (Ho et al., 2020; Song et al., 2021) cast  
 087 generation as the iterative denoising of a Gaussian-corrupted  
 088 input by a learned denoiser. Subsequent work has refined  
 089 individual components: noise schedules and sampler design  
 090 (Karras et al., 2022; Lin et al., 2024), timestep-weighted  
 091 optimization (Hang et al., 2023), latent-space modeling with  
 092 cross-attention conditioning (Rombach et al., 2022), few-  
 093 step distillation (Song et al., 2023), and continuous-time  
 094 transport through flow matching and rectified flow (Lipman  
 095 et al., 2023; Liu et al., 2023). Comparative studies that hold  
 096 the denoising objective fixed and vary the remaining components  
 097 (Karras et al., 2022; Ma et al., 2024) establish that  
 098 the denoiser network itself is a first-order determinant of  
 099 generative performance, leaving open the question of *which*  
 100 architectural properties are responsible.

101 **Architectures for denoising generative models.** The U-  
 102 Net (Ronneberger et al., 2015) has served as the canonical  
 103 denoiser since the earliest score-based generative models  
 104 (Song & Ermon, 2020); Dhariwal & Nichol (2021) subsequently  
 105 characterized the contributions of its residual blocks,  
 106 skip connections, and attention placements through systematic  
 107 ablations. Beyond U-Net variants, transformer-based  
 108  
 109

denoisers operating on tokenized image patches have been  
 introduced as alternatives: U-ViT preserves the U-Net’s long  
 skip connections in token space (Bao et al., 2023); DiT establishes  
 power-law scaling for latent-patch denoisers (Peebles & Xie,  
 2023); DiffiT introduces time-conditional self-attention  
 tailored to the denoising trajectory (Hatamizadeh et al.,  
 2024); and SiT decouples backbone, objective, and noise  
 transport to isolate the contribution of each (Ma et al.,  
 2024). A common structural feature across these architectures  
 is that each hidden unit receives input from a spatially  
 localized region of the noisy input, a property our analysis  
 formalizes as *sparse connectivity*.

### Generalization and memorization in diffusion models.

A recent line of diffusion-specific theory attributes generalization  
 to locality in the denoiser. Kamb & Ganguli (2025) derive an  
 analytic patch-mosaic optimum under locality and equivariance  
 constraints that closely predicts trained U-Nets; Niedoba et al.  
 (2025) approximate pretrained denoisers across architectures  
 by patch-level empirical aggregation; and An et al. (2025) report  
 that early-layer attention locality is predictive of generalization  
 in DiTs. From a training-dynamics perspective, Bonnaire et al.  
 (2025) identify a timescale separation between generation and  
 memorization. Memorization itself has been localized to cross-  
 attention magnitudes, specific neurons, and local image regions  
 (Wen et al., 2024; Ren et al., 2024; Hintersdorf et al.,  
 2024; Chen & collaborators, 2025; Zhang et al., 2026). We  
 provide optimization-theoretic counterpart of these observations.

## 3. Successful Denoisers are Patch-Based

Diffusion models are typically defined by their training objective.  
 A neural network receives a corrupted input and predicts a quantity  
 related to either the clean data or the injected noise. This view  
 explains why denoising is a statistically meaningful learning  
 problem, but it does not explain why the denoising network used  
 in practice can learn the structure of the target distribution  
 through gradient descent. In this section, we take architecture as  
 the variable and show that patch-based processing can substantially  
 change the outcome of the same noise-prediction task.

### 3.1. Background: diffusion and flows

Diffusion models define a path from clean data to noise and  
 train a neural network to reverse this path (Ho et al., 2020).  
 Let  $\mathbf{x}_0 \sim p_{\text{data}}$  be a clean sample and let  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$  be  
 Gaussian noise. A common forward process produces a noisy  
 input

$$\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \epsilon,$$

where  $t$  indexes the noise level. In the noise-prediction  
 parameterization, the denoising network  $\epsilon_\theta(\mathbf{x}_t, t)$  is trained

by

$$\mathcal{L}_{\text{noise}}(\theta) = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [\|\epsilon_{\theta}(\mathbf{x}_t, t) - \epsilon\|_2^2].$$

This is the objective used in our experiments.

Noise prediction is directly connected to score estimation. For the additive noising model  $\mathbf{x}_t = \mathbf{x}_0 + \sigma_t \epsilon$ , Tweedie’s formula gives

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) = \frac{\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t] - \mathbf{x}_t}{\sigma_t^2}.$$

Since  $\mathbf{x}_t = \mathbf{x}_0 + \sigma_t \epsilon$ , learning the conditional mean of  $\epsilon$  is equivalent to learning the score up to a scalar factor. For general affine schedules  $\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \epsilon$ , the same relationship holds after rescaling. This is why noise prediction became a standard way to train score-based diffusion models (Ho et al., 2020; Song et al., 2021).

The same family of models also has an ODE or flow view. If  $p_t$  denotes the marginal distribution of  $\mathbf{x}_t$ , sampling can be written through a probability-flow ODE whose vector field is determined by the score (Song et al., 2021). Flow matching and rectified flow formulate this idea more directly by choosing an interpolation  $\mathbf{z}_t = a_t \mathbf{x}_0 + b_t \epsilon$  and training a neural network to predict the velocity  $\mathbf{v}_t = a'_t \mathbf{x}_0 + b'_t \epsilon$  along the path (Lipman et al., 2023; Liu et al., 2023). These formulations differ in parameterization and weighting, but they share the same basic structure. The model is trained on corrupted high-dimensional inputs and must learn a vector field that separates data structure from injected noise.

Modern diffusion models typically parameterize this denoising network with architectures built around spatial structure. Early score-based and diffusion models largely used convolutional U-Net backbones (Ronneberger et al., 2015; Song & Ermon, 2020; Song et al., 2021; Dhariwal & Nichol, 2021). More recent transformer-based denoisers process noisy images or latents as patch tokens, as in U-ViT, DiT, DiffT, and SiT (Bao et al., 2023; Peebles & Xie, 2023; Hatamizadeh et al., 2024; Ma et al., 2024). Thus, across convolutional and transformer families, patch-based processing has become a common design pattern in denoising generative models.

### 3.2. A toy experiment with architecture as the variable

Inspired by the toy setup of Li & He (2025), we construct a simple high-dimensional denoising problem whose intrinsic structure is easy to visualize. Let  $\hat{\mathbf{x}} \in \mathbb{R}^2$  be a sample from a two-dimensional spiral distribution and let  $\mathbf{P} \in \mathbb{R}^{D \times 2}$  be a fixed column-orthogonal matrix. The observed clean sample is

$$\mathbf{x}_0 = \mathbf{P} \hat{\mathbf{x}} \in \mathbb{R}^D.$$

The model never observes  $\hat{\mathbf{x}}$  or  $\mathbf{P}$ . It only sees corrupted samples  $\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \epsilon$  in the ambient space and is trained with the noise-prediction loss  $\mathcal{L}_{\text{noise}}$ . We set  $D = 256$ , which makes the observed space much larger than the

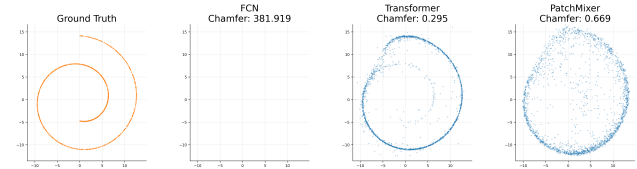


Figure 1. **Architecture affects noise prediction on a randomly embedded spiral.** The ground-truth spiral is embedded into  $\mathbb{R}^{256}$ , models are trained in the ambient space, and samples are visualized after projection back to the original two-dimensional coordinates. Chamfer distance compares generated and ground-truth point clouds, where lower is better. The FCN fails, while the Transformer and PatchMixer recover the spiral much more accurately under the same noise-prediction objective.

intrinsic dimension while keeping the generated samples easy to visualize through projection by  $\mathbf{P}^\top$ .

All models predict noise, so the prediction target, loss, noise schedule, optimizer, and sampling procedure are fixed. We vary only the architecture of the score estimator. We compare an FCN with two patch-based architectures, a vision Transformer (Dosovitskiy et al., 2021) and an MLP-Mixer-style model (Tolstikhin et al., 2021), which we call PatchMixer. Both patch-based models partition the 256 ambient coordinates into 16 patches of dimension 16. The Transformer linearly embeds each patch into a 96-dimensional token and processes the token sequence with single-head self-attention layers. PatchMixer uses the same patch tokens and alternates token-mixing and channel-mixing ReLU layers. The comparison therefore keeps the denoising task fixed while changing how the high-dimensional noisy vector is first presented to the network.

Figure 1 shows a clear separation between the architectures. The FCN does not recover the spiral and produces samples that spread away from the intrinsic support, while the Transformer and PatchMixer produce much more faithful samples under the same objective. The training curves and PCA spectra in Figure 2 give a more detailed view of the failure.

### 3.3. Patch size controls the difficulty of noise prediction

The previous comparison separates full-vector processing from patch-based processing. We now vary the patch size within the patch-based models to test whether the local input dimension itself matters. We embed the same two-dimensional spiral into a  $32 \times 32$  dimensional ambient space and train Transformer and PatchMixer denoisers with the same noise-prediction objective. Only the patch size is changed.

Figure 3 shows that the quality of noise prediction depends strongly on patch size. For both Transformer and PatchMixer, smaller patches generally produce samples that bet-

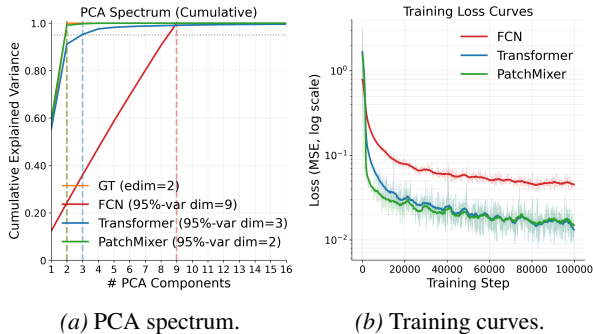


Figure 2. **Optimization and geometry of generated samples.** (a) The ground-truth distribution is intrinsically two-dimensional. PatchMixer reaches 95% explained variance with 2 principal components and the Transformer with 3, while the FCN requires 9, indicating that its generated samples spread into extra ambient directions. (b) The models have comparable parameter counts, with the FCN slightly larger in our implementation, but the FCN decreases the noise-prediction loss more slowly and plateaus at a higher value.

ter concentrate around the ground-truth spiral, while larger patches make each token process a higher-dimensional noisy input and lead to more diffuse generations. The benefit of patch-based architectures therefore comes not from patching itself, but from the patches being sufficiently small. It comes from presenting the denoising network with sufficiently small local views of the corrupted input.

The takeaway of this section is that architecture changes the learnability of the same noise-prediction task. Patch-based models can recover low-dimensional structure from high-dimensional noisy observations in settings where an FCN fails, and reducing the patch size improves this effect. The next section formalizes this observation through a stability-based analysis of sparsely connected networks and shows how patch size enters the implicit regularization induced by gradient descent.

#### 4. Theoretical Analysis: Sparsely Connected Networks Below the Edge of Stability

The experiments in Section 3 isolate a clean empirical pattern: under the same noise-prediction objective, optimizer, and noise schedule, a fully connected denoiser fails to learn the low-dimensional signal buried in high-dimensional noise, while a patch-based denoiser succeeds. This section explains why. The answer lies in how the data geometry interacts with the implicit regularization of gradient descent via the design of network architecture.

For FCNs on high-dimensional isotropic inputs, Liang et al. (2025; 2026) identified a failure mode called *neural shattering*: if the inputs of training data sets concentrate near a hypersphere (e.g. high-dimensional Gaussian distribution),

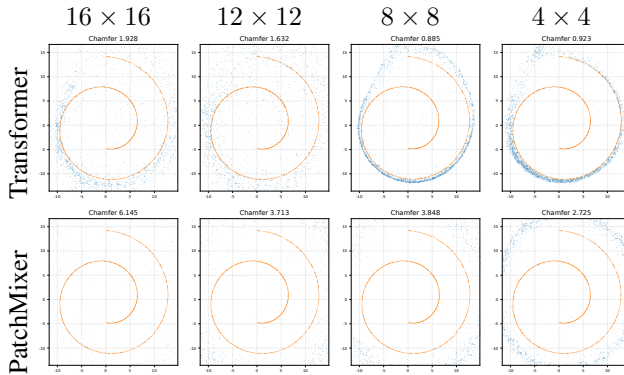


Figure 3. **Patch-size ablation for noise prediction.** A two-dimensional spiral is embedded into a  $32 \times 32$  dimensional ambient space and generated samples are visualized after projection to the intrinsic two-dimensional coordinates. Columns vary the patch size and rows compare Transformer and PatchMixer denoisers. Smaller patches generally improve sample quality, while larger patches produce more diffuse generations and higher Chamfer distances.

one can separate any single training point from all others by a hyperplane that cuts off a tiny cap. A ReLU neuron oriented along that direction and thresholded at the base of the cap activates only on very few samples. Because the neuron fires on a negligible fraction of the data, its gradient updates are controlled entirely by that single sample. The neuron learns a dedicated per-sample correction rather than a shared feature. The resulting network consists of many rare-activation, high-magnitude units that collectively interpolate the training labels.

Crucially, Liang et al. (2025; 2026) prove that when the input distribution is spherical, these memorization ReLU FCNs exist in a regime where gradient descent on neural networks typically enters, the *below the edge of stability* (BEoS) regime: the largest eigenvalue of the training loss Hessian increases, and then hovers around a stability threshold controlled by the step size, where (Cohen et al., 2021). They also empirically verified that GD indeed finds these memorization solutions below the edge of stability. Moreover, they proved that the generalization performance smoothly degrades as data concentrates near a hypersphere. However, when it comes to diffusion, at moderate to high noise levels, each training sample  $\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \epsilon$  is dominated by the Gaussian term, which yields this bad data geometry making FCNs vulnerable to memorization. We propose that this could be a potential reason to explain why gradient-descent trained FCNs cannot be a good denoiser/score estimator in practice.

**Definition 4.1** (Below Edge of Stability (Qiao et al., 2024, Definition 2.3)). Let  $\{\theta_t\}_{t \geq 1}$  be a GD trajectory on  $\mathcal{L}$  with step size  $\eta$ . Any parameter  $\theta_t$  with  $\lambda_{\max}(\nabla^2 \mathcal{L}(\theta_t)) \leq \frac{2}{\eta}$  is called a solution *Below the Edge of Stability* (BEoS).

The question then becomes: what changes when the denoiser is patch-based? A patch-based architecture departs from the fully connected case in one essential way. Instead of presenting each neuron with the full  $d$ -dimensional input, it partitions the input into patches. To capture this property, we study *sparingly connected networks* (SCNs).

#### 4.1. Problem Setup

**Notations.** Throughout the paper, we use  $O(\cdot)$  and  $\Omega(\cdot)$  to absorb constants, while  $\tilde{O}(\cdot)$  absorbs logarithmic factors.  $\mathbb{S}^{d-1}$  denotes the unit hypersphere in  $\mathbb{R}^d$  and  $\mathbb{B}_R^d$  denotes a  $d$ -dimensional ball of radius  $R$ . The activation function is the ReLU denoted by  $\phi(z) = \max\{0, z\}$ . We write  $[n] = \{1, 2, \dots, n\}$ .

**Sparse connectivity pattern.** Fix integers  $d \geq 1$  and  $1 \leq m \leq d$ . A *sparse connectivity pattern* is specified by a collection of coordinate subsets (receptive fields)  $\mathcal{S} = \{S_j\}_{j=1}^J$ , where each  $S_j \subset [d]$  has size  $m$ . Each subset induces a coordinate projection  $\pi_j : \mathbb{R}^d \rightarrow \mathbb{R}^m$ , and we call  $\mathbf{x}^{(S_j)} := \pi_j(\mathbf{x})$  a *patch*. This formalizes the property that each hidden neuron only receives input from a restricted subset of the ambient coordinates. Patch-based architectures, such as the patch based denoisers in Section 3, are instances of this design with  $m \ll d$ .

**Sparingly connected network with weight sharing (SCN).** This is the minimal architecture that captures patch-based processing. Given receptive fields  $\mathcal{S}$  and width  $K \in \mathbb{N}$ ,

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{k=1}^K v_k \left( \frac{1}{J} \sum_{j=1}^J \phi(\mathbf{w}_k^T \pi_j(\mathbf{x}) - b_k) \right) + \beta. \quad (1)$$

We refer to this architecture as a *sparingly connected network* (SCN) or *SCN* if the sparse connectivity pattern is prescribed. A unit  $\phi(\mathbf{w}_k^T \cdot - b_k)$  is called a *filter* or *neuron*. It is said to be *activated* on a patch  $\mathbf{p}_j$  if  $\mathbf{w}_k^T \mathbf{p}_j > b_k$ . Let  $\Theta^{\text{SCN}}$  be the parameter set of all SCNs of *any finite* width. When  $m = d$  and  $J = 1$ , the network reduces to a fully connected architecture, and our analysis recovers the FCN results of (Liang et al., 2026).

**Data and loss.** A dataset is  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  with  $\mathbf{x}_i \in \mathbb{B}_R^d$  and  $y_i \in [-D, D]$ . We use the squared loss and define the training objective  $\mathcal{L}(\boldsymbol{\theta}) := \frac{1}{2n} \sum_{i=1}^n (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i)^2$ . Throughout, GD refers to vanilla gradient descent with learning rate  $0 < \eta < 2$ .

Now we define the stability regularized parameter subset:

$$\Theta_{\text{BEoS}}^{\text{SCN}}(\eta, \mathcal{D}) := \left\{ \boldsymbol{\theta} \in \Theta^{\text{SCN}} \mid \lambda_{\max}(\nabla^2 \mathcal{L}(\boldsymbol{\theta})) \leq \frac{2}{\eta} \right\}. \quad (2)$$

Our goal is to understand how the architecture, encoded by the receptive field system  $\mathcal{S}$ , changes the effective capacity of this BEoS class.

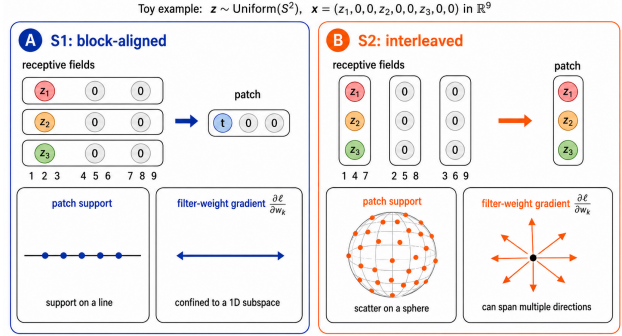


Figure 4. **Data  $\times$  architecture determines the geometry that gradient descent interacts with.** The same spherical ambient data produces a one-dimensional patch cloud under  $\mathcal{S}_1$  and a three-dimensional spherical patch cloud under  $\mathcal{S}_2$ .

**Statistical learning framework.** We assume the data points  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  are i.i.d. samples from a distribution  $\mathcal{P}$  on  $\mathbb{B}_R^d \times [-D, D]$ . The *population risk* of a predictor  $f$  is  $\mathcal{R}(f) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} [(f(\mathbf{x}) - y)^2]$ . The *empirical risk* on dataset  $\mathcal{D}$  is  $\hat{\mathcal{R}}_{\mathcal{D}}(f) = \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2$ . The *generalization gap* is  $\text{GenGap}(f, \mathcal{D}) := |\mathcal{R}(f) - \hat{\mathcal{R}}_{\mathcal{D}}(f)|$ .

#### 4.2. A Toy Example: How Sparse Connectivity Shapes the Gradient Dynamics

We use a toy example to preview how sparse connectivity in the network shapes the gradient dynamics itself. Consider an ambient input  $z \sim \text{Uniform}(\mathbb{S}^2)$  embedded into  $\mathbb{R}^9$  by padding each coordinate with two zeros. The ambient data is spherical, so under the analysis of (Liang et al., 2026), a fully connected network trained on this data receives no meaningful constraint from the BEoS condition. Now equip a two-layer SCN with two different receptive field systems:

- $\mathcal{S}_1$  (block-aligned): each receptive field groups 3 consecutive coordinates. Every patch has the form  $(t, 0, 0)$ , so the patch geometry is one-dimensional.
- $\mathcal{S}_2$  (interleaved): each receptive field groups coordinates modulo the stride. A patch recovers  $(z_1, z_2, z_3)$ , so the patch geometry is spherical.

The gradient of the training loss with respect to a shared filter  $\mathbf{w}_k$  exposes the difference. For the model (1),

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \mathbf{w}_k} = \frac{v_k}{nJ} \sum_{i,j} r_i \underbrace{\mathbb{1}\{\mathbf{w}_k^T \pi_j(\mathbf{x}_i) > b_k\}}_{\text{residue and activation scalars}} \underbrace{\pi_j(\mathbf{x}_i)}_{\text{patch vectors}} \quad (3)$$

where  $r_i = f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i$  is the residual. The update to  $\mathbf{w}_k$  is a linear combination of the patches that activate it. Under  $\mathcal{S}_1$ , every activated patch lies in a one-dimensional subspace, so each filter gradient is confined to that subspace.

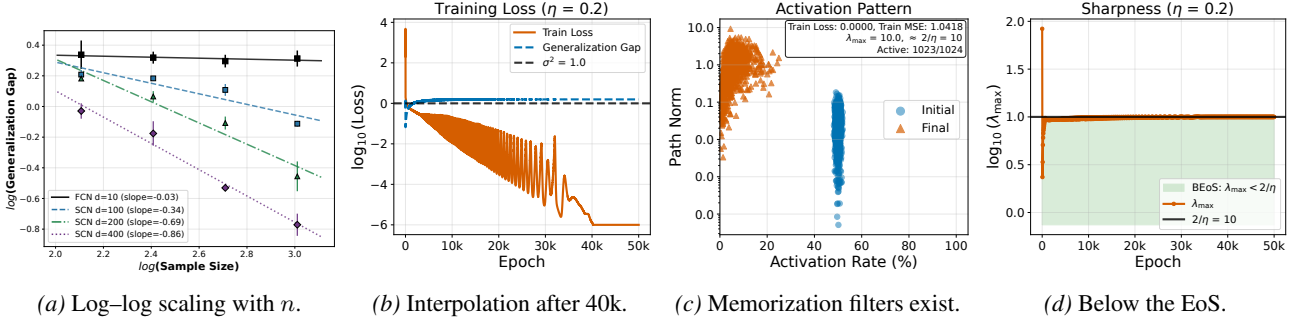


Figure 5. **Numerical validations via synthetic experiments.** (a)  $\text{GenGap}(f_\theta, \mathcal{D})$  versus the sample size  $n$  on a log–log scale. The fitted slope summarizes the empirical rate: if  $\text{GenGap} \lesssim n^{-c}$ , then  $\log(\text{GenGap}) \leq -c \log n + b$ , so a more negative slope indicates faster decay (better generalization). In our experiments, the FCN slope is nearly flat (slope =  $-0.03$  at  $d = 10$ ), whereas SCN exhibits increasingly negative slopes as  $d$  grows (slope =  $-0.34$  at  $d = 100$ ,  $-0.69$  at  $d = 200$ , and  $-0.86$  at  $d = 400$ ), indicating faster decay with  $n$ . (b) Training loss and generalization gap, showing that the SCN interpolates noisy labels in the random spherical-patch setting. (c) Scatter plot of per-neuron path norm  $|v_k| \|\mathbf{w}_k\|$  against activation rate  $\frac{1}{nJ} \sum_{i,j} \mathbb{1}\{\mathbf{w}_k^\top \pi_j(\mathbf{x}_i) > b_k\}$ . Many neurons have large path norms while activating on only a small fraction of the patch multiset, consistent with the failure mode in Theorem 4.4: rare activations can fit labels while keeping the sharpness small. (d) The largest Hessian eigenvalue hovers near  $2/\eta \approx 10$ , indicating the EoS regime.

The patch multiset is highly concentrated, and any filter that activates on a meaningful fraction of patches must align with the shared direction. Under  $\mathcal{S}_2$ , activated patches span the entire three-dimensional patch space.

The two receptive field systems thus produce fundamentally different gradient dynamics on the same ambient data, and the sets of solutions that GD can stably reach differ dramatically. The rest of this section makes this intuition precise through the lens of the BEOs condition. Theorem 4.2 shows that the BEOs condition bounds a weighted path norm whose weight is governed by the patch geometry. Theorem 4.3 and Theorem 4.4 then develop the two contrasting regimes previewed by  $\mathcal{S}_1$  and  $\mathcal{S}_2$ : one enables generalization, and the other fails.

### 4.3. Characterization of the Stability Regularization

Given a dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  and receptive fields  $\mathcal{S} = \{S_j\}_{j=1}^J$ , let  $\mathbf{X}_D^S$  be a random vector drawn uniformly from the patch multiset  $\{\pi_j(\mathbf{x}_i)\}_{(i,j) \in [n] \times [J]} \subset \mathbb{R}^m$ . Define the weight function  $g_{\mathcal{D}, \mathcal{S}} : \mathbb{S}^{m-1} \times \mathbb{R} \rightarrow \mathbb{R}$  by

$$g_{\mathcal{D}, \mathcal{S}}(\mathbf{u}, t) := \mathbb{E} [\phi(\mathbf{u}^\top \mathbf{X}_D^S - t)] \times \sqrt{\mathbb{P}(\mathbf{u}^\top \mathbf{X}_D^S > t)^2 + \left\| \mathbb{E} [\mathbf{X}_D^S \mathbb{1}\{\mathbf{u}^\top \mathbf{X}_D^S > t\}] \right\|_2^2}. \quad (4)$$

This function evaluates a neuron’s activation boundary  $(\mathbf{u}, t)$  against the patch multiset. The term  $\mathbb{P}(\mathbf{u}^\top \mathbf{X}_D^S > t)^2$  is the squared activation mass: a neuron that fires on a large fraction of patches incurs a large penalty on its path norm, while one that fires on few patches incurs a weak penalty.

**Theorem 4.2** (BEOs implies patch-geometry-weighted path

norm bound). *Fix  $\mathcal{D}$  and  $\mathcal{S}$ . For any  $\theta \in \Theta^{\text{SCN}}$ ,*

$$\sum_{k=1}^K |v_k| \|\mathbf{w}_k\| g_{\mathcal{D}, \mathcal{S}} \left( \frac{\mathbf{w}_k}{\|\mathbf{w}_k\|}, \frac{b_k}{\|\mathbf{w}_k\|} \right) \leq \frac{1}{2} \left( \lambda_{\max}(\nabla^2 \mathcal{L}(\theta)) + 2(R+1)\sqrt{2\mathcal{L}(\theta)} - 1 \right). \quad (5)$$

*In particular, for any  $\theta \in \Theta_{\text{BEOs}}^{\text{SCN}}(\eta, \mathcal{D})$ ,*

$$\sum_{k=1}^K |v_k| \|\mathbf{w}_k\| g_{\mathcal{D}, \mathcal{S}} \left( \frac{\mathbf{w}_k}{\|\mathbf{w}_k\|}, \frac{b_k}{\|\mathbf{w}_k\|} \right) \leq \frac{1}{\eta} - \frac{1}{2} + (R+1)\sqrt{2\mathcal{L}(\theta)}. \quad (6)$$

The proof is in Appendix D. The core connection between the Hessian and  $g_{\mathcal{D}, \mathcal{S}}$  is detailed in Proposition D.10. The architecture  $\mathcal{S}$  enters solely through this weight function: different receptive field systems produce different patch multisets, hence different penalty structures. When  $m = d$  and  $J = 1$ , the network reduces to a fully connected architecture and (6) recovers the result of (Liang et al., 2026).

Returning to our toy example makes this concrete. Under  $\mathcal{S}_1$ , the patch multiset is essentially one-dimensional and highly concentrated. Any neuron that activates on a non-negligible fraction of patches must cut through the dense interior, so the activation probability  $\mathbb{P}(\mathbf{u}^\top \mathbf{X}_D^S > t)$  is bounded away from zero and  $g_{\mathcal{D}, \mathcal{S}}$  has a positive lower bound. The BEOs condition therefore enforces a meaningful path-norm constraint. Under  $\mathcal{S}_2$ , the patch multiset is a three-dimensional sphere. A neuron can isolate an individual patch with a carefully placed hyperplane, driving the activation probability arbitrarily close to zero and making  $g_{\mathcal{D}, \mathcal{S}}$  arbitrarily small. The BEOs constraint becomes vacuous. The next two subsections develop the quantitative consequences of this dichotomy.

#### 4.4. Two Regimes: Generalization versus Memorization

The next two theorems formalize the contrasting behaviors previewed by  $\mathcal{S}_1$  and  $\mathcal{S}_2$ .

**Theorem 4.3** (Generalization under structured patch geometry). *Let the marginal distribution of  $\mathbf{x}$  be  $\text{Uniform}(\mathbb{S}^{d-1})$ , and let  $\mathcal{D}$  be a dataset of  $n$  i.i.d. samples. Assume  $d > 3$  and  $1 \leq m < \frac{d(d-3)}{d+3}$ . For any  $\theta \in \Theta_{\text{BEoS}}^{\text{SCN}}(\eta, \mathcal{D})$  with  $\|f_\theta\|_\infty \leq M$ , with probability  $\geq 1 - 2\delta$ ,*

$$\text{GenGap}(f_\theta, \mathcal{D}) \lesssim_d \text{poly}\left(\frac{1}{\eta}, J, M\right) n^{-\frac{(d-m)(d+3)}{6d^2-2md+6d-6m}}. \quad (7)$$

The crucial feature of this bound is its scaling with ambient dimension. When the patch size  $m$  is fixed and  $d \rightarrow \infty$ , the exponent approaches  $-\frac{1}{6} + O(m/d)$ , so there is *no curse of dimensionality*. In the same spherical setting, FCNs admit no non-trivial guarantee under the BEoS condition (Liang et al., 2026). The geometric reason, formalized in Proposition D.11, is that when  $m \ll d$ , the patch projections concentrate near the origin of  $\mathbb{R}^m$ . The patch multiset becomes low-dimensional and concentrated, resisting the kind of pointwise isolation that weakens the stability constraint. The path-norm penalty in Theorem 4.2 therefore becomes effective. The formal proof appears in Appendix E, and Figure 5 numerically verifies the predicted scaling.

**Theorem 4.4** (Memorization below the edge of stability). *Assume all training patches have norm at most one, and every nonzero-labeled training example contains a unit-norm patch that appears nowhere else in the full training patch multiset. There exists a network of the form (1) with width  $K \leq n$  that interpolates  $\mathcal{D}$  and satisfies*

$$\lambda_{\max}(\nabla_{\theta}^2 \mathcal{L}) \leq 1 + \frac{D^2 + 2/J^2}{n}. \quad (8)$$

Theorem 4.4 shows that the BEoS condition alone does not prevent memorization when the patch multiset admits single-patch isolation. From the weighted path-norm perspective of Theorem 4.2, the same rare activations correspond to small values of  $g_{\mathcal{D}, \mathcal{S}}$ , making the stability-induced constraint weak. The detailed proof can be found in Appendix F. An empirical demonstration is shown in Figure 5. The experiment uses random spherical patches (case B in Figure 4). Since these patches are uniformly sampled from  $\mathbb{S}^{m-1}$ , the uniqueness condition in Theorem 4.4 holds generically.

Taken together, Theorems 4.3 and 4.4 delineate two regimes governed by the patch geometry induced by the receptive-field system  $\mathcal{S}$ . When  $\mathcal{S}$  and the data produce a concentrated patch multiset that resists single-patch isolation, the BEoS condition yields an effective regularity constraint and generalization is possible. When the patch multiset contains

patches that can be isolated by ReLU half-spaces, the constraint collapses and memorization becomes compatible with stability. For a fixed training dataset, the architectural design  $\mathcal{S}$  determines the induced patch multiset and hence which regime the network operates in.

## 5. Towards Deep Denoisers: Output Skips Shape Representation Geometry

The previous sections focus on the geometry of the input interface. In a two-layer sparsely connected model, each filter acts on a patch of the corrupted input, so patch size controls the dimension of the object processed by gradient descent. Deep denoisers introduce another effect. A block at depth  $k$  does not process the raw noisy input directly; it processes the representation produced by all previous blocks. Thus the relevant geometry is no longer fixed at the input layer, but evolves along the residual stream.

This perspective is useful for interpreting the clean-prediction result of Li & He (2025). They show that large-patch pixel Transformers can fail catastrophically when directly trained to predict noise or velocity, while the same architecture becomes effective when it predicts the clean image. Their explanation is based on the manifold assumption: clean images are structured and low-dimensional, whereas noise and velocity are high-dimensional quantities. We agree with this target-space intuition, but we use it to highlight an architectural mechanism. Clean prediction changes the output parameterization of the denoiser. It gives the model an analytic long skip from the noisy input to the final noise or velocity predictor, so that the learned branch can focus on producing clean structure.

We illustrate this effect in the same  $D = 256$  randomly embedded spiral setup used above. Both models use the same deep fully connected backbone and the same training procedure in Section 3. The only change is the output parameterization. Direct  $\epsilon$ -prediction produces diffuse samples, while  $x$ -prediction recovers the spiral much more accurately in Figure 6.

To see the architectural content of this parameterization, consider the interpolation

$$\mathbf{z}_t = (1-t)\mathbf{x} + t\epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I}),$$

where  $t = 0$  is clean data and  $t = 1$  is pure noise. If the learned branch predicts the clean sample,  $f_\theta(\mathbf{z}_t, t) \approx \mathbf{x}$ , then the corresponding noise predictor is

$$\hat{\epsilon} = \frac{\mathbf{z}_t - (1-t)f_\theta(\mathbf{z}_t, t)}{t} = \frac{1}{t}\mathbf{z}_t - \frac{1-t}{t}f_\theta(\mathbf{z}_t, t).$$

Thus, as a noise predictor,  $x$ -prediction is a fixed output-side skip from  $\mathbf{z}_t$  plus a learned clean branch. The same

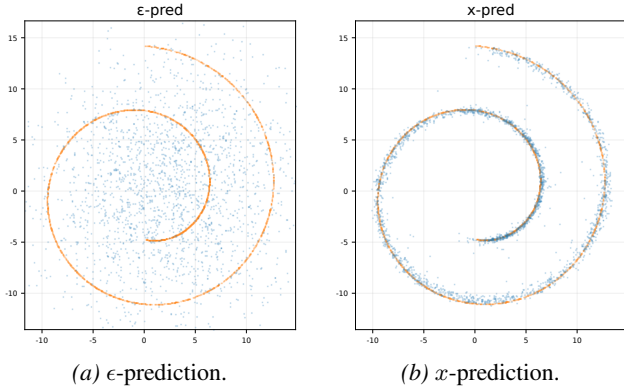


Figure 6. **Clean prediction rescues the same deep denoising backbone.** Both models are trained on a two-dimensional spiral embedded into  $\mathbb{R}^{256}$  and visualized after projection to the intrinsic coordinates. Direct  $\epsilon$ -prediction produces diffuse samples, while  $x$ -prediction recovers the spiral support much more accurately.

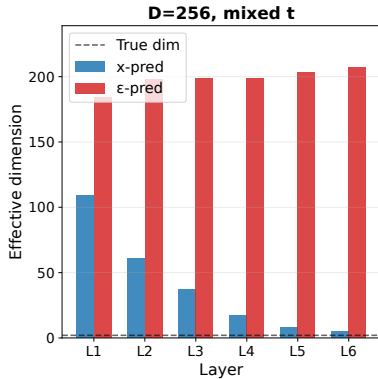


Figure 7. **Effective dimension under mixed noise levels.** Hidden representations are aggregated under the training distribution of  $t$ . The  $x$ -prediction model forms lower-dimensional representations than the direct  $\epsilon$ -prediction model across layers.

observation holds for velocity prediction after the corresponding change of coefficients. The coefficients are determined by the schedule, so this is a specific parameterization rather than a generic residual connection. Under the noise-prediction loss, the learned branch is trained by

$$\|\hat{\epsilon} - \epsilon\|_2^2 = \frac{(1-t)^2}{t^2} \|f_\theta(\mathbf{z}_t, t) - \mathbf{x}\|_2^2.$$

The high-dimensional noisy input is carried to the output analytically, while the network branch receives a clean-image regression objective.

This changes the forward pressure on the residual stream. Write a depth- $L$  backbone schematically as

$$\mathbf{h}^0 = E\mathbf{z}_t, \mathbf{h}^k = \mathbf{h}^0 + \sum_{\ell=1}^k B_\ell(\mathbf{h}^{\ell-1}, t), f_\theta(\mathbf{z}_t, t) = R(\mathbf{h}^L).$$

The input to block  $k$  is  $\mathbf{h}^{k-1} = \mathbf{h}^0 + \sum_{\ell=1}^{k-1} B_\ell(\mathbf{h}^{\ell-1}, t)$ . Under direct  $\epsilon$ -prediction, the final readout must recover

high-dimensional noise. Even if early blocks begin to extract the clean low-dimensional signal, later blocks are pressured to transform the residual stream back into a representation from which noise can be read out. The same stream has to both build useful denoising features and carry high-dimensional noisy information.

Under  $x$ -prediction, the learned branch has a consistent target across depth. Each block improves the representation only insofar as it helps the final readout approximate  $\mathbf{x}$ . The fixed skip already supplies the explicit  $\mathbf{z}_t$  term needed to form  $\hat{\epsilon}$ , so the residual stream does not need to preserve every high-dimensional noisy direction. Forward propagation is therefore biased toward a progressively cleaner representation: as depth increases, the input to the next block becomes more structured and lower-dimensional. This is the sense in which clean prediction can be viewed as an architectural design. It shapes the geometry seen by later layers.

We measure this effect by computing the effective dimension of hidden states, defined as the number of principal components needed to explain 95% of the variance. Figure 7 reports the layerwise effective dimension. For  $x$ -prediction, the dimension decreases across layers, while for direct  $\epsilon$ -prediction, the representation remains high-dimensional through most of the network. This matches the residual-stream interpretation: the clean branch learns a progressively denoised geometry, while the noise-prediction stream stays tied to noisy sources. In this way, the  $x$ -prediction model forms lower-dimensional internal representations across the network, while direct  $\epsilon$ -prediction stays far from the intrinsic dimension.

This lower-dimensional representation geometry is also relevant to memorization. A model whose intermediate features align with the shared low-dimensional structure has less incentive to allocate capacity to isolated high-dimensional residuals. In contrast, a high-dimensional residual stream gives the later layers more room to fit idiosyncratic noise directions. Thus output skips, like sparse connectivity, can be understood as architectural choices that alter the geometry on which gradient descent operates.

## 6. Discussion

This work shows that the architecture of a denoiser matters because it determines the geometry on which gradient descent operates: sparse connectivity at the input restricts each neuron to patches, making memorization harder; output skips such as clean-data prediction keep the internal representations lower-dimensional. Together, these findings suggest a unified geometric principle for understanding architectural inductive bias in generative models: the crucial variable is not just what functions can be represented, but what geometry the optimizer sees during training.

## References

- 440  
441  
442 NIST digital library of mathematical functions.  
443 <http://dlmf.nist.gov/>. See §5.6(i), Eq. 5.6.E4 for  
444 Gautschi’s inequality (Gamma ratio bounds).
- 445 Aharon, M., Elad, M., and Bruckstein, A. K-svd: An algo-  
446 rithm for designing overcomplete dictionaries for sparse  
447 representation. *IEEE Transactions on signal processing*,  
448 54(11):4311–4322, 2006.
- 449  
450 An, J., Wang, D., Guo, P., Luo, J., and Schwing, A. On  
451 inductive biases that enable generalization in diffusion  
452 transformers. In *Advances in Neural Information Pro-*  
453 *cessing Systems*, 2025. URL [https://openreview.](https://openreview.net/forum?id=lE2cD7C9fk)  
454 [net/forum?id=lE2cD7C9fk](https://openreview.net/forum?id=lE2cD7C9fk).
- 455 Anthony, M. and Bartlett, P. L. *Neural Network Learning:*  
456 *Theoretical Foundations*. Cambridge University Press,  
457 1999.
- 458  
459 Arora, S., Du, S. S., Hu, W., Li, Z., Salakhutdinov, R. R.,  
460 and Wang, R. On exact computation with an infinitely  
461 wide neural net. *Advances in neural information process-*  
462 *ing systems*, 32, 2019.
- 463  
464 Bao, F., Nie, S., Xue, K., Cao, Y., Li, C., Su, H., and Zhu, J.  
465 All are worth words: A vit backbone for diffusion models.  
466 In *Proceedings of the IEEE/CVF Conference on Com-*  
467 *puter Vision and Pattern Recognition*, pp. 22669–22679,  
468 2023. URL [https://openaccess.thecvf.](https://openaccess.thecvf.com/content/CVPR2023/html/Bao_All_Are_Worth_Words_A_ViT_Backbone_for_Diffusion_Models_CVPR_2023_paper.html)  
469 [com/content/CVPR2023/html/Bao\\_All\\_](https://openaccess.thecvf.com/content/CVPR2023/html/Bao_All_Are_Worth_Words_A_ViT_Backbone_for_Diffusion_Models_CVPR_2023_paper.html)  
470 [Are\\_Worth\\_Words\\_A\\_ViT\\_Backbone\\_for\\_](https://openaccess.thecvf.com/content/CVPR2023/html/Bao_All_Are_Worth_Words_A_ViT_Backbone_for_Diffusion_Models_CVPR_2023_paper.html)  
471 [Diffusion\\_Models\\_CVPR\\_2023\\_paper.html](https://openaccess.thecvf.com/content/CVPR2023/html/Bao_All_Are_Worth_Words_A_ViT_Backbone_for_Diffusion_Models_CVPR_2023_paper.html).
- 472  
473 Beyer, L., Izmailov, P., Kolesnikov, A., Caron, M., Korn-  
474 blith, S., Zhai, X., Minderer, M., Tschannen, M., Alab-  
475 dulumohsin, I., and Pavetic, F. Flexivit: One model for all  
476 patch sizes. In *Proceedings of the IEEE/CVF Conference*  
477 *on Computer Vision and Pattern Recognition (CVPR)*, pp.  
478 14496–14506, 2023.
- 479  
480 Bonnaire, T., Urfin, R., Biroli, G., and Mezard, M. Why  
481 diffusion models don’t memorize: The role of implicit  
482 dynamical regularization in training. In *Advances in Neu-*  
483 *ral Information Processing Systems*, 2025. URL [https://openreview.](https://openreview.net/forum?id=BSZqpqqgM0)  
484 [net/forum?id=BSZqpqqgM0](https://openreview.net/forum?id=BSZqpqqgM0).
- 485  
486 Brutzkus, A., Globerson, A., Malach, E., Netser, A. R.,  
487 and Shalev-Schwartz, S. Efficient learning of cnns using  
488 patch based features. In *International Conference on*  
489 *Machine Learning*, pp. 2336–2356. PMLR, 2022.
- 490  
491 Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J.,  
492 Bojanowski, P., and Joulin, A. Emerging properties in  
493 self-supervised vision transformers. In *Proceedings of*  
494 *the IEEE/CVF International Conference on Computer*  
*Vision (ICCV)*, pp. 9650–9660, 2021.
- Chen, C. and collaborators. Exploring local memo-  
rization in diffusion models via bright ending atten-  
tion. In *International Conference on Learning Repre-*  
*sentations*, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2410.21665)  
[2410.21665](https://arxiv.org/abs/2410.21665). Accepted spotlight; author list in source  
should be checked against final proceedings if needed.
- Coates, A., Ng, A., and Lee, H. An analysis of single-  
layer networks in unsupervised feature learning. In *Pro-*  
*ceedings of the fourteenth international conference on*  
*artificial intelligence and statistics*, pp. 215–223. JMLR  
Workshop and Conference Proceedings, 2011.
- Cohen, J., Kaur, S., Li, Y., Kolter, J. Z., and Talwalkar,  
A. Gradient descent on neural networks typically occurs  
at the edge of stability. In *International Conference on*  
*Learning Representations*, 2021.
- Dhariwal, P. and Nichol, A. Diffusion models beat  
gans on image synthesis. In *Advances in Neural*  
*Information Processing Systems*, volume 34, pp. 8780–  
8794, 2021. URL [https://proceedings.](https://proceedings.neurips.cc/paper/2021/hash/49ad23d1ec9fa4bd8d77d02681df5cfa-Abstract.html)  
[neurips.cc/paper/2021/hash/](https://proceedings.neurips.cc/paper/2021/hash/49ad23d1ec9fa4bd8d77d02681df5cfa-Abstract.html)  
[49ad23d1ec9fa4bd8d77d02681df5cfa-Abstract.](https://proceedings.neurips.cc/paper/2021/hash/49ad23d1ec9fa4bd8d77d02681df5cfa-Abstract.html)  
html.
- Ding, L., Drusvyatskiy, D., Fazel, M., and Harchaoui, Z.  
Flat minima generalize for low-rank matrix recovery. *In-*  
*formation and Inference: A Journal of the IMA*, 13(2):  
iaae009, 2024.
- Ding, X., Zhang, X., Zhou, Y., Han, J., Ding, G., and Sun,  
J. Scaling up your kernels to 31x31: Revisiting large  
kernel design in cnns. In *Proceedings of the IEEE/CVF*  
*Conference on Computer Vision and Pattern Recognition*  
*(CVPR)*, pp. 11963–11975, 2022.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn,  
D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer,  
M., Heigold, G., Gelly, S., Uszkoreit, J., and Houslyby,  
N. An image is worth 16x16 words: Transformers for  
image recognition at scale. In *International Conference*  
*on Learning Representations*, 2021.
- Gunasekar, S., Lee, J., Soudry, D., and Srebro, N. Charac-  
terizing implicit bias in terms of optimization geometry.  
In *International Conference on Machine Learning*, pp.  
1832–1841. PMLR, 2018a.
- Gunasekar, S., Lee, J. D., Soudry, D., and Srebro, N. Im-  
plicit bias of gradient descent on linear convolutional  
networks. *Advances in neural information processing*  
*systems*, 31, 2018b.
- Hang, T., Gu, S., Li, C., Bao, J., Chen, D., Hu, H.,  
Geng, X., and Guo, B. Efficient diffusion training  
via min-snr weighting strategy. In *Proceedings of the*

- 495 *IEEE/CVF International Conference on Computer Vision*,  
496 pp. 7441–7451, 2023. URL [https://arxiv.org/](https://arxiv.org/abs/2303.09556)  
497 [abs/2303.09556](https://arxiv.org/abs/2303.09556).
- 498 Hatamizadeh, A., Song, J., Liu, G., Kautz, J., and Vahdat,  
499 A. Diffit: Diffusion vision transformers for image gener-  
500 ation. In *European Conference on Computer Vision*,  
501 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=uAKk0I3xxm)  
502 [id=uAKk0I3xxm](https://openreview.net/forum?id=uAKk0I3xxm).
- 503 Haussler, D. Decision-theoretic generalizations of the pac  
504 model for neural net and other learning applications. *In-*  
505 *formation and Computation*, 100(1):78–150, 1992.
- 506 He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learn-  
507 ing for image recognition. In *Proceedings of the IEEE*  
508 *conference on computer vision and pattern recognition*,  
509 pp. 770–778, 2016.
- 510 Hintersdorf, D., Struppek, L., Kersting, K., Dziedzic, A.,  
511 and Boenisch, F. Finding nemo: Localizing neurons  
512 responsible for memorization in diffusion models. In  
513 *Advances in Neural Information Processing Systems*,  
514 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=YAEKMFZyJm)  
515 [id=YAEKMFZyJm](https://openreview.net/forum?id=YAEKMFZyJm).
- 516 Ho, J., Jain, A., and Abbeel, P. Denoising diffu-  
517 sion probabilistic models. In *Advances in Neural*  
518 *Information Processing Systems*, volume 33, pp. 6840–  
519 6851, 2020. URL [https://proceedings.](https://proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html)  
520 [neurips.cc/paper/2020/hash/](https://proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html)  
521 [4c5bcfec8584af0d967f1ab10179ca4b-Abstract.](https://proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html)  
522 [html](https://proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html).
- 523 Jacot, A., Gabriel, F., and Hongler, C. Neural tangent ker-  
524 nel: Convergence and generalization in neural networks.  
525 *Advances in neural information processing systems*, 31,  
526 2018.
- 527 Kamb, M. and Ganguli, S. An analytic theory of  
528 creativity in convolutional diffusion models. In  
529 *Proceedings of the 42nd International Conference*  
530 *on Machine Learning*, volume 267 of *Proceedings*  
531 *of Machine Learning Research*, pp. 28795–28831,  
532 2025. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v267/kamb25a.html)  
533 [v267/kamb25a.html](https://proceedings.mlr.press/v267/kamb25a.html).
- 534 Karras, T., Aittala, M., Aila, T., and Laine, S. Elucidating  
535 the design space of diffusion-based generative models.  
536 In *Advances in Neural Information Processing Systems*,  
537 volume 35, pp. 26565–26577, 2022. URL [https://](https://arxiv.org/abs/2206.00364)  
538 [arxiv.org/abs/2206.00364](https://arxiv.org/abs/2206.00364).
- 539 Lahoti, A., Karp, S., Winston, E., Singh, A., and Li, Y. Role  
540 of locality and weight sharing in image-based tasks: A  
541 sample complexity separation between cnns, lcns, and  
542 fens. In *The Twelfth International Conference on Learn-*  
543 *ing Representations*, 2024.
- 544 Li, T. and He, K. Back to basics: Let denoising generative  
545 models denoise, 2025. arXiv preprint arXiv:2511.13720.
- 546 Li, Z., Zhang, Y., and Arora, S. Why are convolutional  
547 nets more sample-efficient than fully-connected nets? In  
548 *International Conference on Learning Representations*,  
549 2021.
- Liang, T., Qiao, D., Wang, Y.-X., and Parhi, R. Stable  
minima of ReLU neural networks suffer from the curse  
of dimensionality: The neural shattering phenomenon.  
In *Advances in Neural Information Processing Systems*  
(*NeurIPS*), 2025.
- Liang, T., Cloninger, A., Parhi, R., and Wang, Y.-X. Gen-  
eralization below the edge of stability: The role of data  
geometry. In *International Conference on Learning Rep-*  
*resentations (ICLR)*, 2026.
- Lin, S., Liu, B., Li, J., and Yang, X. Common diffusion noise  
schedules and sample steps are flawed. In *Proceedings*  
*of the IEEE/CVF Winter Conference on Applications of*  
*Computer Vision*, pp. 5404–5411, 2024. URL [https://](https://arxiv.org/abs/2305.08891)  
[arxiv.org/abs/2305.08891](https://arxiv.org/abs/2305.08891).
- Lipman, Y., Chen, R. T. Q., Ben-Hamu, H., Nickel, M.,  
and Le, M. Flow matching for generative modeling. In  
*International Conference on Learning Representations*,  
2023. URL [https://openreview.net/forum?](https://openreview.net/forum?id=PqvMRDCJT9t)  
[id=PqvMRDCJT9t](https://openreview.net/forum?id=PqvMRDCJT9t).
- Liu, H., Chen, M., Zhao, T., and Liao, W. Besov function ap-  
proximation and binary classification on low-dimensional  
manifolds using convolutional residual networks. In *In-*  
*ternational Conference on Machine Learning*, pp. 6770–  
6780. PMLR, 2021a.
- Liu, X., Gong, C., and Liu, Q. Flow straight and fast:  
Learning to generate and transfer data with rectified flow.  
In *International Conference on Learning Representations*,  
2023. URL [https://openreview.net/forum?](https://openreview.net/forum?id=XVjTT1nw5z)  
[id=XVjTT1nw5z](https://openreview.net/forum?id=XVjTT1nw5z).
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin,  
S., and Guo, B. Swin transformer: Hierarchical vision  
transformer using shifted windows. In *IEEE/CVF In-*  
*ternational Conference on Computer Vision (ICCV)*, pp.  
9992–10002, 2021b.
- Liu, Z., Mao, H., Wu, C.-Y., Feichtenhofer, C., Darrell, T.,  
and Xie, S. A ConvNet for the 2020s. In *IEEE/CVF*  
*Conference on Computer Vision and Pattern Recognition*  
(*CVPR*), pp. 11966–11976, 2022.
- Ma, N., Goldstein, M., Albergo, M. S., Boffi, N. M., Vanden-  
Eijnden, E., and Xie, S. Sit: Exploring flow and diffusion-  
based generative models with scalable interpolant trans-  
formers. In *European Conference on Computer Vi-*

- 550 *sion*, 2024. URL [https://arxiv.org/abs/2401.](https://arxiv.org/abs/2401.08740)  
551 08740.
- 552
- 553 Mao, T., Shi, Z., and Zhou, D.-X. Theory of deep convo-  
554 lutional neural networks iii: Approximating radial func-  
555 tions. *Neural Networks*, 144:778–790, 2021.
- 556
- 557 Mhaskar, H., Liao, Q., and Poggio, T. When and why are  
558 deep networks better than shallow ones? In *Proceed-*  
559 *ings of the AAAI conference on artificial intelligence*,  
560 volume 31, 2017.
- 561
- 562 Mohri, M., Rostamizadeh, A., and Talwalkar, A. *Founda-*  
563 *tions of Machine Learning*. MIT Press, second edition,  
564 2018.
- 565
- 566 Mulayoff, R., Michaeli, T., and Soudry, D. The implicit  
567 bias of minima stability: A view from function space.  
568 *Advances in Neural Information Processing Systems*, 34:  
569 17749–17761, 2021.
- 570
- 571 Nacson, M. S., Mulayoff, R., Ongie, G., Michaeli, T., and  
572 Soudry, D. The implicit bias of minima stability in multi-  
573 variate shallow ReLU networks. In *International Confer-*  
574 *ence on Learning Representations*, 2023.
- 575
- 576 Nguyen, D.-K., Assran, M., Jain, U., Oswald, M. R., Snoek,  
577 C. G. M., and Chen, X. An image is worth more than  
578 16x16 patches: Exploring transformers on individual pix-  
579 els. *arXiv preprint arXiv:2406.09415*, 2024.
- 580
- 581 Niedoba, M., Zwartsenberg, B., Murphy, K. P., and Wood,  
582 F. Towards a mechanistic explanation of diffusion model  
583 generalization. In *International Conference on Machine*  
584 *Learning*, 2025. URL [https://openreview.net/](https://openreview.net/forum?id=Hrp6jRIKdX)  
585 [forum?id=Hrp6jRIKdX](https://openreview.net/forum?id=Hrp6jRIKdX).
- 586
- 587 Oono, K. and Suzuki, T. Approximation and non-parametric  
588 estimation of resnet-type convolutional neural networks.  
589 In *International conference on machine learning*, pp.  
590 4922–4931. PMLR, 2019.
- 591
- 592 Parhi, R. and Nowak, R. D. Near-minimax optimal estima-  
593 tion with shallow ReLU neural networks. *IEEE Transac-*  
594 *tions on Information Theory*, 69(2):1125–1139, 2023.
- 595
- 596 Paulin, M., Mairal, J., Douze, M., Harchaoui, Z., Perronnin,  
597 F., and Schmid, C. Convolutional patch representations  
598 for image retrieval: an unsupervised approach. *Inter-*  
599 *national Journal of Computer Vision*, 121(1):149–168,  
600 2017.
- 601
- 602 Peebles, W. and Xie, S. Scalable diffusion models with trans-  
603 formers. In *Proceedings of the IEEE/CVF International*  
604 *Conference on Computer Vision*, pp. 4195–4205, 2023.  
URL <https://arxiv.org/abs/2212.09748>.
- Peyré, G. Manifold models for signals and images. *Com-*  
*puter vision and image understanding*, 113(2):249–260,  
2009.
- Poggio, T. Foundations of deep learning: Compositional  
sparsity of computable functions. Technical report,  
CBMM memo 138, 2022.
- Poggio, T. and Fraser, M. Compositional sparsity of learn-  
able functions. *Bulletin of the American Mathematical*  
*Society*, 61(3):438–456, 2024.
- Qiao, D., Zhang, K., Singh, E., Soudry, D., and Wang,  
Y.-X. Stable minima cannot overfit in univariate ReLU  
networks: Generalization by large step sizes. In *Advances*  
*in Neural Information Processing Systems*, volume 37,  
pp. 94163–94208, 2024.
- Ren, J., Li, Y., Zeng, S., Xu, H., Lyu, L., Xing, Y., and  
Tang, J. Unveiling and mitigating memorization in text-  
to-image diffusion models through cross attention. In  
*European Conference on Computer Vision*, 2024. URL  
<https://arxiv.org/abs/2403.11052>.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and  
Ommer, B. High-resolution image synthesis with la-  
tent diffusion models. In *Proceedings of the IEEE/CVF*  
*Conference on Computer Vision and Pattern Recogni-*  
*tion*, pp. 10684–10695, 2022. URL [https://arxiv.](https://arxiv.org/abs/2112.10752)  
[org/abs/2112.10752](https://arxiv.org/abs/2112.10752).
- Ronneberger, O., Fischer, P., and Brox, T. U-net: Convolu-  
tional networks for biomedical image segmentation. In *In-*  
*ternational Conference on Medical image computing and*  
*computer-assisted intervention*, pp. 234–241. Springer,  
2015.
- Shi, Z., Fang, Z., and Cao, Y. Approximation and estimation  
capability of vision transformers for hierarchical compo-  
sitional models. *Applied and Computational Harmonic*  
*Analysis*, pp. 101849, 2025.
- Siegel, J. W. and Xu, J. Characterization of the variation  
spaces corresponding to shallow neural networks. *Con-*  
*structive Approximation*, pp. 1–24, 2023.
- Song, Y. and Ermon, S. Improved techniques for training  
score-based generative models. In *Advances in Neural*  
*Information Processing Systems*, volume 33, pp. 12438–  
12448, 2020. URL [https://proceedings.](https://proceedings.neurips.cc/paper/2020/hash/92c3b916311a5517d9290576e3ea37ad-Abstract.html)  
[neurips.cc/paper/2020/hash/](https://proceedings.neurips.cc/paper/2020/hash/92c3b916311a5517d9290576e3ea37ad-Abstract.html)  
[92c3b916311a5517d9290576e3ea37ad-Abstract.](https://proceedings.neurips.cc/paper/2020/hash/92c3b916311a5517d9290576e3ea37ad-Abstract.html)  
html.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A.,  
Ermon, S., and Poole, B. Score-based generative mod-  
eling through stochastic differential equations. In *In-*  
*ternational Conference on Learning Representations*,

2021. URL <https://openreview.net/forum?id=PxtTIG12RRHS>.
- Song, Y., Dhariwal, P., Chen, M., and Sutskever, I. Consistency models. In *International Conference on Machine Learning*, 2023. URL <https://openreview.net/forum?id=FmqFfMTNnv>.
- Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., and Srebro, N. The implicit bias of gradient descent on separable data. *Journal of Machine Learning Research*, 19(70):1–57, 2018.
- Thiry, L., Arbel, M., Belilovsky, E., and Oyallon, E. The unreasonable effectiveness of patches in deep convolutional kernels methods. In *International Conference on Learning Representations*, 2021.
- Tolstikhin, I. O., Houlsby, N., Kolesnikov, A., Beyer, L., Zhai, X., Unterthiner, T., Yung, J., Steiner, A., Keysers, D., Uszkoreit, J., et al. Mlp-mixer: An all-mlp architecture for vision. *Advances in neural information processing systems*, 34:24261–24272, 2021.
- Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., and Jégou, H. Training data-efficient image transformers & distillation through attention. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 10347–10357. PMLR, 2021.
- Trockman, A. and Kolter, J. Z. Patches are all you need? *arXiv preprint arXiv:2201.09792*, 2022.
- van der Vaart, A. W. and Wellner, J. A. *Weak Convergence and Empirical Processes*. Springer, 1996.
- Wainwright, M. J. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press, 2019.
- Wang, F., Yang, T., Yu, Y., Ren, S., Wei, G., Wang, A., Shao, W., Zhou, Y., Yuille, A., and Xie, C. Adventurer: Optimizing vision mamba architecture designs for efficiency. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 30157–30166, June 2025a.
- Wang, F., Yu, Y., Wei, G., Shao, W., Zhou, Y., Yuille, A., and Xie, C. Scaling laws in patchification: An image is worth 50,176 tokens and more. In *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pp. 65278–65290. PMLR, 2025b.
- Wang, W., Xie, E., Li, X., Fan, D.-P., Song, K., Liang, D., Lu, T., Luo, P., and Shao, L. Pvt v2: Improved baselines with pyramid vision transformer. *Computational Visual Media*, 8(3):415–424, 2022.
- Wang, Z. and Wu, L. Theoretical analysis of the inductive biases in deep convolutional networks. *Advances in Neural Information Processing Systems*, 36:74289–74338, 2023.
- Wen, Y., Liu, Y., Chen, C., and Lyu, L. Detecting, explaining, and mitigating memorization in diffusion models. In *International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=84n3UwkH7b>.
- Wu, H., Xiao, B., Codella, N., Liu, M., Dai, X., Yuan, L., and Zhang, L. Cvt: Introducing convolutions to vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 22–31, 2021.
- Wu, L. and Su, W. J. The implicit regularization of dynamical stability in stochastic gradient descent. In *International Conference on Machine Learning*, pp. 37656–37684. PMLR, 2023.
- Xiao, T., Singh, M., Mintun, E., Darrell, T., Dollár, P., and Girshick, R. Early convolutions help transformers see better. In *Advances in Neural Information Processing Systems*, volume 34, pp. 30392–30400, 2021.
- Xie, S., Girshick, R. B., Dollár, P., Tu, Z., and He, K. Aggregated residual transformations for deep neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5987–5995, 2017.
- Yao, Z., Gholami, A., Keutzer, K., and Mahoney, M. W. Pyhessian: Neural networks through the lens of the hessian. In *2020 IEEE international conference on big data (Big data)*, 2020.
- Yarotsky, D. Error bounds for approximations with deep relu networks. *Neural networks*, 94:103–114, 2017.
- Yuan, L., Chen, Y., Wang, T., Yu, W., Shi, Y., Jiang, Z., Tay, F. E. H., Feng, J., and Yan, S. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 558–567, 2021.
- Zhang, Z., Zhang, K., Chen, M., Takeda, Y., Wang, M., Zhao, T., and Wang, Y.-X. Nonparametric classification on low dimensional manifolds using overparameterized convolutional residual networks. *Advances in Neural Information Processing Systems*, 37:65738–65764, 2024.
- Zhang, Z., Li, X., Li, X., Shi, L., Wu, M., Tao, M., and Qu, Q. Generalization of diffusion models arises with a balanced representation space. In *International Conference on Learning Representations*, 2026. URL <https://openreview.net/forum?id=57TheGgNAN>.

660 Zhou, D.-X. Universality of deep convolutional neural net-  
661 works. *Applied and computational harmonic analysis*, 48  
662 (2):787–794, 2020.

663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714

## A. More Related Work

**Patch-based representation learning.** Early unsupervised image representation methods decomposed images into local patches and learned sparse dictionaries or manifolds from these primitives (Aharon et al., 2006; Peyré, 2009). Patch based representations have also shown success in unsupervised learning (Paulin et al., 2017) and as preprocessing for image classification (Coates et al., 2011; Thiry et al., 2021; Brutzkus et al., 2022). These results highlight that patch space is highly structured, but it remains open how neural network training exploits such structure. Our work helps bridge this gap by analyzing generalization below the edge of stability.

**Geometric awareness in deep learning.** Recent theoretical work has explored how architectural design incorporates geometric awareness to improve neural network training and approximation. Approximation theory research has demonstrated that CNN architectures overcome the curse of dimensionality by exploiting compositional structure in natural images (Mhaskar et al., 2017; Zhou, 2020; Poggio, 2022; Poggio & Fraser, 2024), with Mao et al. (2021) formally proving the superiority of CNNs over FCNs for learning certain composite functions. More recently, Vision Transformers (Dosovitskiy et al., 2021) have been analyzed through the same lens: Shi et al. (2025) extend these approximation guarantees to ViTs, showing they outperform FCNs for compositional functions, and Trockman & Kolter (2022) identify patch extraction, rather than attention, as the critical component, further reinforcing the strength of sparsely connected architectures. These results, however, either assume the labeling function satisfies a compositional structure or operate in settings with explicit regularization. Our work requires no assumptions on the labeling function and operates in the overparameterized regime.

**Theoretical analysis of CNNs and separation from FCNs.** The advantages of CNNs over FCNs have been established from both approximation-theoretic and statistical perspectives. Approximation theory shows that CNNs with appropriate sparse weights can achieve near-minimax optimal sample complexity (Oono & Suzuki, 2019) and adapt to intrinsic dimension when data lie on a low-dimensional manifold (Liu et al., 2021a), with Zhang et al. (2024) extending these results to overparameterized regimes under weight decay. However, these analyses rely on reducing CNNs to fully connected networks (Yarotsky, 2017) and therefore do not reveal the architectural insights specific to CNNs that we describe in this paper. Specifically, we do not require the data to be supported on a low-dimensional manifold to avoid the curse of dimensionality, nor do we require explicit regularization such as weight decay or a sparsity constraint. On the statistical side, a sample complexity separating CNNs, locally connected networks without weight sharing, and FCNs has been established (Li et al., 2021; Wang & Wu, 2023; Lahoti et al., 2024). Our results are for the same model family, but differ in that we consider an overparameterized regime without explicit regularization, which allows us to prove a stronger separation between FCNs and CNNs. Lastly, the significance of the input data distribution, rather than the labeling function, was not discovered in prior work.

**Implicit bias of gradient descent.** Many existing work on the implicit bias of gradient descent relies on strong assumptions on the data distributions (e.g., linearly separable data), simplified model architectures (e.g., linear activation) to make a gradient dynamics analysis tractable (Soudry et al., 2018; Gunasekar et al., 2018a;b), or weight initialization schemes that keep the model in the kernel regimes (Jacot et al., 2018; Arora et al., 2019). While CNNs were studied (Gunasekar et al., 2018b; Arora et al., 2019), the nature of the results are different from ours.

**Edge of stability, minima stability, and generalization by large stepsizes.** Our approach builds upon a recent line of work that studies the set of solutions that gradient descent can visit (or converge to) via either Edge-of-stability observation or the minima stability theory (Ding et al., 2024; Mulayoff et al., 2021; Nacson et al., 2023; Wu & Su, 2023; Qiao et al., 2024; Liang et al., 2025; 2026). These approaches enable formal analysis of the generalization properties without having to analyze gradient dynamics. To the best of our knowledge, they all focused on feedforward neural networks, and we are the first to study CNNs and the impact of model architecture choices in the implicit regularization of large stepsizes.

## B. Broader (Informal) Discussion on Network Architecture Design: Bridge our Theoretical Insights to Empirical Observations

This appendix expands on the architectural implications discussed in Section 6. The purpose is not to claim that the stability analysis in the main text provides a complete theory of modern deep vision architectures. Rather, the point is that stability gives a tractable lens through which one can see a more general object: the geometry of the local vectors exposed to gradient descent. This object appears explicitly in the shallow model analyzed in this paper, and it has close analogues in the patchification and local-operator design choices used in modern CNNs and Vision Transformers.

**The patch matrix as the geometry seen by backpropagation.** In the setting of this paper, consider the patch matrix

$$\mathbf{X} \in \mathbb{R}^{(nJ) \times m},$$

whose row indexed by  $(i, j)$  equals  $\pi_j(\mathbf{x}_i)^\top$ . Let  $r_i := f_\theta(\mathbf{x}_i) - y_i$ , and define the lifted residual vector  $\bar{\mathbf{r}} \in \mathbb{R}^{nJ}$  by

$$(\bar{\mathbf{r}})_{(i,j)} := r_i.$$

Let  $\mathbf{S}(\theta) \in \{0, 1\}^{(nJ) \times K}$  be the gating matrix

$$\mathbf{S}(\theta)_{(i,j),k} := \mathbb{1}\{\mathbf{w}_k^\top \pi_j(\mathbf{x}_i) > b_k\},$$

and let  $\mathbf{v} := (v_1, \dots, v_K)^\top$ . Stack first-layer weights as

$$\mathbf{W} := [\mathbf{w}_1^\top; \dots; \mathbf{w}_K^\top] \in \mathbb{R}^{K \times m}.$$

A direct calculation gives a backprop-aligned factorization of the gradient. Define the patch-level backprop signal

$$\mathbf{R}(\theta) := \frac{1}{nJ} \bar{\mathbf{r}} \mathbf{v}^\top \in \mathbb{R}^{(nJ) \times K}, \quad \mathbf{G}(\theta) := \mathbf{R}(\theta) \odot \mathbf{S}(\theta),$$

where  $\odot$  denotes entrywise multiplication. Then

$$\nabla_{\mathbf{W}} \mathcal{L}(\theta) = \mathbf{G}(\theta)^\top \mathbf{X}. \tag{9}$$

Equivalently, writing  $\mathbf{V} := \text{diag}(v_1, \dots, v_K)$ , one can express the same signal as

$$\mathbf{G}(\theta) = \frac{1}{nJ} \text{diag}(\bar{\mathbf{r}}) \mathbf{S}(\theta) \mathbf{V}.$$

The factorization (9) separates two roles. The matrix  $\mathbf{G}(\theta)$  aggregates the learning signal generated by residuals, readout weights, and gates. The patch matrix  $\mathbf{X}$  then maps this signal into parameter updates. Thus the architecture does not only restrict the number of parameters, it determines the geometry through which backpropagated signals are converted into motion in parameter space. In this sense,  $\mathbf{X}$  acts as a signal rectifier in the sense that data priors and receptive-field structure jointly shape the directions along which gradient descent can effectively move. For example, Figure 9a shows that, for CIFAR-10, 90% of the energy of the convolutional patch matrix is concentrated in only a few principal directions, suggesting a strong geometric constraint on the filter dynamics.

This viewpoint also clarifies why the stability condition in the main text is informative. The terms appearing in the definition of  $g_{\mathcal{D}, \mathcal{S}}$  are drawn from the same ingredients as (9): gate statistics induced by  $\mathbf{S}(\theta)$  and geometric moments of the patch cloud encoded by  $\mathbf{X}$ . Theorem 4.3 can therefore be viewed as a static stability proxy for the backpropagation geometry in (9). It does not say that stability alone explains generalization. Indeed, Theorem 4.4 shows that stable solutions may still memorize when the patch cloud admits isolating half-spaces. If corresponding backpropagation direction may enable a filter move to a location such that it only activates on a rare, nearly isolated patch, then the network may fit labels while paying little stability cost. If the patch cloud is concentrated and resistant to such isolation, the same stability lens yields an effective regularity constraint.

**Patch size as local-operator dimension in Vision Transformers.** Vision Transformers make the role of patch geometry especially explicit. The original ViT converts an image into a sequence of non-overlapping patches and applies a shared patch embedding before global token mixing (Dosovitskiy et al., 2021). This patchification step is often introduced as an efficiency device, since larger patches shorten the token sequence and reduce the cost of attention. The factorization (9) suggests a second interpretation. Patchification also chooses the first local geometry seen by training. A larger patch presents each shared local operator with a higher-dimensional and more compressed local vector. A smaller patch reduces this local dimension and preserves finer spatial information, while increasing the number of tokens that later layers must mix.

Several empirical lines of work support the idea that patch size is a genuine training and generalization variable, not only an implementation detail. FlexiViT trains a single ViT over a range of patch sizes and demonstrates that patch size controls an accuracy–compute tradeoff at deployment time (Beyer et al., 2023). Self-supervised ViTs also reveal the importance of this axis: DINO identifies small patches as one ingredient behind the strong emergent properties of self-supervised ViTs (Caron

Table 1. Selected patchification-scaling results adapted from (Wang et al., 2025b). Patch size denotes the side length of the square image patch used to form one token. Panel A shows that reducing patch size improves ImageNet accuracy. Panel B separates true patch-size reduction from only increasing the sequence length by interpolating coarse tokens. Here we only list their results on DeiT-B (Touvron et al., 2021) and Adventurer-B (Wang et al., 2025a).

Panel A: ImageNet classification top-1 accuracy					
Model and input	$16 \times 16$	$8 \times 8$	$4 \times 4$	$2 \times 2$	$1 \times 1$
DeiT-B, $64 \times 64$ input	68.2	76.9	80.1	80.8	81.3
DeiT-B, $128 \times 128$ input	78.1	81.0	82.3	82.9	–
Adventurer-B, $224 \times 224$ input	82.6	83.9	84.3	84.5	84.6
Panel B: Sequence-length ablation on ImageNet top-1 accuracy					
Model and input	Seq. length	Interpolated coarse tokens	True patch scaling		
DeiT-B, $128 \times 128$ input	256	78.2 (+0.1)	81.0 (+2.9)		
DeiT-B, $128 \times 128$ input	1,024	78.2 (+0.1)	82.3 (+4.2)		
Adventurer-B, $224 \times 224$ input	784	82.7 (+0.1)	83.9 (+1.3)		
Adventurer-B, $224 \times 224$ input	3,136	82.8 (+0.2)	84.3 (+1.7)		
Adventurer-B, $224 \times 224$ input	12,544	82.8 (+0.2)	84.5 (+1.9)		

et al., 2021). Early-convolution studies make the same point from the optimization side. The standard ViT patchify stem is a large-kernel, large-stride convolution, and replacing it by a short stack of small-stride convolutions significantly changes optimization stability and improves performance under common training recipes (Xiao et al., 2021). These results are consistent with the idea that the first patch matrix is not an innocuous preprocessing choice: it shapes the geometry of the signals that subsequent layers and backpropagation operate on.

Recent patchification scaling studies push this observation further. Wang et al. (Wang et al., 2025b) systematically reduce the patch size from the standard coarse-token regime toward pixel-level tokenization and report smooth gains across classification, semantic segmentation, object detection, and instance segmentation. Table 1 summarizes the results most relevant to our discussion. The most informative part is their sequence-length ablation. When coarse patch tokens are merely interpolated into a longer sequence, the gains remain small. When genuinely smaller patches are extracted from the image, the gains are much larger. This suggests that patch-size scaling changes the local visual primitives exposed to the model, not only the number of tokens. In the notation of this paper, reducing patch size changes the initial patch matrix itself: it lowers spatial compression in each local vector and changes the multiset of local rows through which the network forms its updates. This aligns with the gradient factorization  $\nabla_{\mathbf{W}} \mathcal{L} = \mathbf{G}^T \mathbf{X}$ : the relevant object is the geometry of the rows in the patch or token matrix, not only the nominal sequence length. Related pixel-token studies push the same question to the extreme by asking how far visual models can go when patchification is removed almost entirely (Nguyen et al., 2024).

**Why small patches require hierarchy and mixing.** Smaller patches improve the local geometry exposed to the model, but they also increase the burden of mixing information across locations. This tradeoff helps explain a major trend in vision-transformer architecture design. Hierarchical Transformers such as Swin do not keep a fixed high-dimensional token space throughout the network; they begin with finer local tokens, restrict early mixing to local windows, and then gradually merge tokens while increasing representation dimension (Liu et al., 2021b). PVTv2 follows a related logic through overlapping patch embeddings, linear-complexity attention, and convolutional feed-forward layers (Wang et al., 2022). CvT introduces convolutional token embeddings and convolutional projections inside Transformer blocks (Wu et al., 2021). Tokens-to-Token ViT replaces one-shot patchification by a progressive tokenization module that recursively aggregates local neighboring tokens (Yuan et al., 2021). These designs are not only remedies for a defective large-patch ViT. They can be read as different ways of managing the same geometric tradeoff: expose early layers to controlled local vectors, then recover expressivity through hierarchy, overlap, convolutional projection, or stage-wise mixing.

In the notation of (9), a deep network replaces the fixed raw patch matrix  $\mathbf{X}$  by a sequence of representation matrices. At layer  $\ell$ , one should imagine a matrix

$$\mathbf{X}_{\ell} = [\mathbf{z}_{\ell,j}(\mathbf{x}_i)^T]_{(i,j)}$$

whose rows are local tokens or local feature vectors produced by the preceding layers. Backpropagation through a local

Table 2. The design trend of sparse connectivity in modern vision backbones.

Representative backbone	Core sparse connectivity pattern	Dimension of the processed unit
ResNet-50 (He et al., 2016)	$3 \times 3$ bottleneck convolution	$3 \times 3 \times C_{\text{in}}$ (up to $C_{\text{in}} = 512$ )
ResNeXt-50 (Xie et al., 2017)	$3 \times 3$ grouped convolution	$3 \times 3 \times C_{\text{in}}/g$ (up to $g = 32$ )
ConvNeXt-T (Liu et al., 2022)	$7 \times 7$ depthwise convolution	$7 \times 7 = 49$
ViT-B/16 (Dosovitskiy et al., 2021)	uniform Transformer	768 for all blocks
Swin-T (Liu et al., 2021b)	hierarchical Transformer	$96 \rightarrow 192 \rightarrow 384 \rightarrow 768$

operator at layer  $\ell$  again has the schematic form

$$\nabla_{\mathbf{w}_\ell} \mathcal{L} \approx \mathbf{G}_\ell(\boldsymbol{\theta})^\top \mathbf{X}_\ell,$$

up to the additional Jacobian factors introduced by normalization, residual connections, attention, and nonlinear mixing. Thus the relevant geometry is no longer fixed by the raw data distribution. It evolves during training. Patch size, windowing, overlap, token merging, and convolutional stems can all be viewed as mechanisms for shaping the sequence

$$\mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_L$$

of optimization-facing geometries. **Patch geometry in CNNs: from spatial-channel mixing to depthwise convolutional kernel.** CNN architecture design provides a sequence of examples in which the vector seen by each local operator is progressively controlled. In a residual bottleneck block, the  $3 \times 3$  spatial convolution acts after a channel-reducing  $1 \times 1$  projection, so the spatial operator does not see the full channel dimension of the block (He et al., 2016). ResNeXt makes this control explicit through grouped convolutions: each group applies a  $3 \times 3$  operator to only a fraction of the channels, reducing the local vector from  $3 \times 3 \times C_{\text{in}}$  to  $3 \times 3 \times C_{\text{in}}/g$  (Xie et al., 2017). ConvNeXt pushes the same principle further by using depthwise spatial convolutions inside an inverted-bottleneck block, with expressive channel mixing supplied by pointwise layers (Liu et al., 2022). Thus, from ResNet to ResNeXt to ConvNeXt, the spatial operator is increasingly restricted to a structured local geometry, while expressivity is recovered through channel mixing and depth.

The kernel-size ablation in ConvNeXt gives a useful empirical instance of the trade-off suggested by our theory. Increasing the depthwise kernel from  $3 \times 3$  to  $7 \times 7$  improves ImageNet accuracy, while increasing it further to  $9 \times 9$  or  $11 \times 11$  yields little additional gain (Liu et al., 2022). In the language of this paper, enlarging the spatial kernel increases the context available to each local operator, reducing the bias induced by overly local processing. At the same time, a larger local vector exposes gradient descent to a richer patch geometry, which can weaken the regularizing effect of locality. The saturation around  $7 \times 7$  is consistent with a bias–variance trade-off for the geometry seen by local operators: enough spatial support improves representation, while excessive support brings less benefit once the local computation is already sufficiently expressive.

RepLKNet develops this direction further by treating very large kernels as a separate architectural scaling dimension (Ding et al., 2022). It shows that replacing ConvNeXt’s  $7 \times 7$  kernels with much larger stage-wise kernels can further improve performance, especially on dense prediction tasks. The key point for our purposes is not only the larger spatial support, but the training-time parameterization used to make such kernels effective. RepLKNet does not train a bare large kernel. It introduces *small-kernel structural re-parameterization*: during training, a large kernel branch is paired with a small kernel branch, such as a  $3 \times 3$  or  $5 \times 5$  branch. After training, the small kernel is padded to the center of the large kernel and algebraically merged, so inference uses a single large kernel.

This mechanism has a direct interpretation in our gradient-geometry framework. Let  $\mathbf{x}_L \in \mathbb{R}^{K^2}$  denote the large local patch and let  $\mathbf{x}_s = \mathbf{C}\mathbf{x}_L \in \mathbb{R}^{s^2}$  be the centered small patch selected from it. Ignoring normalization for clarity, the training-time operator has the form

$$h(\mathbf{x}) = \mathbf{w}_L^\top \mathbf{x}_L + \mathbf{w}_s^\top \mathbf{x}_s = (\mathbf{w}_L + \mathbf{C}^\top \mathbf{w}_s)^\top \mathbf{x}_L, \quad \mathbf{w}_{\text{eff}} = \mathbf{w}_L + \mathbf{C}^\top \mathbf{w}_s.$$

At inference time this is just a large kernel with weight  $\mathbf{w}_{\text{eff}}$ . During training, however, the two branches expose different local geometries to gradient descent. If  $r$  denotes the backpropagated scalar signal, then

$$\nabla_{\mathbf{w}_L} \mathcal{L} = r \mathbf{x}_L, \quad \nabla_{\mathbf{w}_s} \mathcal{L} = r \mathbf{C} \mathbf{x}_L, \quad \Delta \mathbf{w}_{\text{eff}} = -\eta r (\mathbf{I} + \mathbf{C}^\top \mathbf{C}) \mathbf{x}_L,$$

up to the branch-dependent scaling introduced by normalization. The small branch therefore gives the centered local subspace an additional optimization path. It biases training toward small-scale local structure while the large branch captures broad spatial context.

This example sharpens the message of the patch-matrix factorization  $\nabla_{\mathbf{W}} \mathcal{L} = \mathbf{G}^\top \mathbf{X}$ . The geometry relevant to learning is not determined solely by the inference-time operator. RepLKNet uses an inference-time large kernel, but its training-time architecture exposes gradient descent to both a large-patch geometry and a centered small-patch geometry. Thus structural re-parameterization is a practical mechanism for modifying the optimization-facing patch geometry without changing the final operator class. Together with the ConvNeXt ablation, this supports the broader view that CNN design balances spatial context against implicit regularization induced by controlled local geometry.

## C. Functional Analysis of Shallow ReLU Networks

### C.1. Path-norm and Variation Semi-norm of ReLU Networks

This subsection collects several basic facts from (Parhi & Nowak, 2023) and (Siegel & Xu, 2023) that we will use later.

**Definition C.1.** Let  $f_\theta(\mathbf{x}) = \sum_{k=1}^K v_k \phi(\mathbf{w}_k^\top \mathbf{x} - b_k) + \beta$  be a fully-connected two-layer neural network. The (unweighted) path-norm of  $f_\theta$  is defined to be

$$\|f_\theta\|_{\text{path}} := \sum_{k=1}^K |v_k| \|\mathbf{w}_k\|_2. \quad (10)$$

**Dictionary representation of ReLU networks.** Using the positive 1-homogeneity of ReLU, one may rescale each hidden unit while leaving the realized function unchanged:

$$v_k \phi(\mathbf{w}_k^\top \mathbf{x} - b_k) = a_k \phi(\mathbf{u}_k^\top \mathbf{x} - t_k), \quad \mathbf{u}_k := \frac{\mathbf{w}_k}{\|\mathbf{w}_k\|_2} \in \mathbb{S}^{d-1}, \quad t_k := \frac{b_k}{\|\mathbf{w}_k\|_2}, \quad a_k := v_k \|\mathbf{w}_k\|_2.$$

Consequently,  $f_\theta$  can be written in the normalized finite-sum form

$$f(\mathbf{x}) = \sum_{k=1}^{K'} a_k \phi(\mathbf{u}_k^\top \mathbf{x} - t_k) + \mathbf{c}^\top \mathbf{x} + c_0. \quad (11)$$

Define the (ReLU) ridge dictionary as  $\mathcal{D}_\phi := \{\phi(\mathbf{u}^\top \cdot -t) : \mathbf{u} \in \mathbb{S}^{d-1}, t \in \mathbb{R}\}$ . We focus on the *overparameterized, width-agnostic* collection obtained by taking the *union over all finite widths*

$$\mathcal{F}_{\text{fin}} := \bigcup_{K \geq 1} \left\{ \sum_{k=1}^K a_k \phi(\mathbf{u}_k^\top \cdot -t_k) + \mathbf{c}^\top(\cdot) + c_0 \right\}, \quad (12)$$

and quantify complexity via the smallest path-norm among all realizations of  $f$ :

$$\|f\|_{\text{path}, \min} := \inf \{ \|f_\theta\|_{\text{path}} : f_\theta \equiv f \text{ of the form (11)} \}.$$

**From finite sums to a width-agnostic integral representation.** Rather than fixing a particular width  $K$ , it is convenient to work with a convex, measure-theoretic formulation that *captures the closure/convex hull of (12)*. Concretely, let  $\nu$  be a finite signed Radon measure on  $\mathbb{S}^{d-1} \times [-R, R]$  and consider

$$f(\mathbf{x}) = \int_{\mathbb{S}^{d-1} \times [-R, R]} \phi(\mathbf{u}^\top \mathbf{x} - t) d\nu(\mathbf{u}, t) + \mathbf{c}^\top \mathbf{x} + c_0. \quad (13)$$

Every finite network of the form (11) corresponds to an *atomic* (hence sparse) measure  $\nu = \sum_{k=1}^K a_k \delta_{(\mathbf{u}_k, t_k)}$ , and conversely any such atomic  $\nu$  yields a finite network. Therefore, (13) should be viewed as a *width-agnostic relaxation* aligned with (10), rather than as an assumption of an infinite-width limit.

**Definition C.2.** The (unweighted) variation (semi)norm

$$|f|_V := \inf \{ \|\nu\|_{\mathcal{M}} : f \text{ admits (11) for some } (\nu, c, c_0) \}, \quad (14)$$

where  $\|\nu\|_{\mathcal{M}}$  is the total variation of  $\nu$ .

For the compact region  $\Omega = \mathbb{B}_R^d$ , we define the bounded variation function class as

$$V_C(\Omega) := \left\{ f: \Omega \rightarrow \mathbb{R} \mid f = \int_{\mathbb{S}^{d-1} \times [-R, R]} \phi(\mathbf{u}^\top \mathbf{x} - t) d\nu(\mathbf{u}, t) + \mathbf{c}^\top \mathbf{x} + b, |f|_V \leq C \right\}. \quad (15)$$

In particular, identifying (11) with the atomic measure  $\nu = \sum_k a_k \delta_{(\mathbf{u}_k, t_k)}$  yields

$$|f|_V \leq \sum_k |a_k| = \|f\theta\|_{\text{path}}, \quad \text{hence} \quad |f|_V \leq \|f\|_{\text{path}, \min}.$$

Moreover, the minimal total variation required to represent  $f$  coincides with the minimal path-norm over all finite decompositions:

$$\|f\|_{\text{path}, \min} = |f|_V. \quad (16)$$

Thus, (14) provides a *nonparametric* analogue of the path-norm: it encodes the same complexity notion while *not committing* to a fixed width  $K$ .

**Remark C.3** (“Arbitrary width”  $\neq$  “infinite width”). All statements here are about  $\mathcal{F}_{\text{fin}}$  in (12), namely networks of *finite* (but unconstrained) width. The integral representation (13) is introduced as a convenient convexification/closure of this union for analysis and regularization; it *does not* posit an infinite-width limit. In particular, when training in the variational form with a total-variation penalty on  $\nu$ , first-order optimality implies that optimal measures are sparse (i.e., have finite support), which corresponds exactly to *finite-width* networks. Therefore, our results hold for *arbitrary (yet finite) width*; the continuum measure serves only as a tool to characterize and control  $\|f\|_{\text{path}, \min}$ .

## C.2. The Metric Entropy of Variation Spaces

Metric entropy is a standard way to describe how “compact” a subset  $A$  is inside a metric space  $(X, \rho_X)$ . We recall the notions of covering numbers and metric entropy.

**Definition C.4** (Covering Number and Entropy). Let  $A$  be a compact subset of a metric space  $(X, \rho_X)$ . For  $t > 0$ , the *covering number*  $N(A, t, \rho_X)$  is the smallest number of closed balls of radius  $t$  whose union contains  $A$ :

$$N(t, A, \rho_X) := \min \left\{ N \in \mathbb{N} : \exists x_1, \dots, x_N \in X \text{ s.t. } A \subset \bigcup_{i=1}^N \mathbb{B}(x_i, t) \right\}, \quad (17)$$

where  $\mathbb{B}(x_i, t) = \{y \in X : \rho_X(y, x_i) \leq t\}$ . The *metric entropy* of  $A$  at scale  $t$  is then

$$H_t(A)_X := \log N(t, A, \rho_X). \quad (18)$$

Covering/entropy bounds for bounded-variation-type classes have been established in the literature. In what follows, we will use the estimate stated below.

**Proposition C.5** (Parhi & Nowak 2023, Appendix D). *The metric entropy of  $V_C(\mathbb{B}_R^d)$  (see Definition C.2) with respect to the  $L^\infty(\mathbb{B}_R^d)$ -distance  $\|\cdot\|_\infty$  satisfies*

$$\log N(t, V_C(\mathbb{B}_R^d), \|\cdot\|_\infty) \lesssim_d \left( \frac{C}{t} \right)^{\frac{2d}{d+3}}. \quad (19)$$

where  $\lesssim_d$  hides constants (which could depend on  $d$ ) and logarithmic factors.

### C.3. Generalization Gap of Unweighted Variation Function Class

As a middle step towards bounding the generalization gap of the weighted variation function class, we bound the generalization gap of the unweighted variation function class using chaining and Gaussian complexity, together with the  $L^\infty$  metric entropy bound in Proposition C.5.

**Proposition C.6** (Wainwright 2019, Chapter 13). *Fix design points  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and denote the empirical norm*

$$\|f\|_n^2 := \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i)^2.$$

Let  $\mathcal{F}$  be a class of real-valued functions on  $\{\mathbf{x}_i\}_{i=1}^n$ , and define

$$\widehat{\mathcal{G}}_n(\mathcal{F}) := \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(\mathbf{x}_i), \quad \varepsilon_1, \dots, \varepsilon_n \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1), \quad \mathcal{G}_n(\mathcal{F}) := \mathbb{E} \widehat{\mathcal{G}}_n(\mathcal{F}).$$

Then

$$\mathcal{G}_n(\mathcal{F}) \lesssim \frac{16}{\sqrt{n}} \int_0^{\text{diam}(\mathcal{F}, \|\cdot\|_n)} \sqrt{\log N(t, \mathcal{F}, \|\cdot\|_n)} dt, \quad (20)$$

where  $\text{diam}(\mathcal{F}, \|\cdot\|_n) := \sup_{f_1, f_2 \in \mathcal{F}} \|f_1 - f_2\|_n$ . Moreover, with probability at least  $1 - \delta$ ,

$$\widehat{\mathcal{G}}_n(\mathcal{F}) \leq \mathcal{G}_n(\mathcal{F}) + \text{diam}(\mathcal{F}, \|\cdot\|_n) \frac{\sqrt{\log(1/\delta)}}{\sqrt{n}} \quad (\delta > 0). \quad (21)$$

**Lemma C.7.** *Let  $\mathcal{F}_{M,C} = \{f \in V_C(\mathbb{B}_R^d) \mid \|f\|_\infty \leq M\}$  with  $M \geq D$ , and let  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathcal{P}^{\otimes n}$  where  $|Y| \leq D$  a.s. Then with probability at least  $1 - \delta$ ,*

$$\sup_{f \in \mathcal{F}_{M,C}} |R(f) - \widehat{R}_{\mathcal{D}}(f)| \lesssim_d C^{\frac{d}{d+3}} M^{\frac{d+6}{d+3}} n^{-\frac{1}{2}} + M^2 \left( \frac{\log(1/\delta)}{n} \right)^{\frac{1}{2}}. \quad (22)$$

*Proof.* Let  $\ell_f(\mathbf{x}, y) := (y - f(\mathbf{x}))^2$  and  $\mathcal{L}_{M,C} := \{\ell_f : f \in \mathcal{F}_{M,C}\}$ . Since  $|y| \leq D \leq M$  and  $\|f\|_\infty \leq M$ , for any  $f, g \in \mathcal{F}_{M,C}$  and any  $(\mathbf{x}, y)$ ,

$$|\ell_f(\mathbf{x}, y) - \ell_g(\mathbf{x}, y)| = |f(\mathbf{x}) - g(\mathbf{x})| |f(\mathbf{x}) + g(\mathbf{x}) - 2y| \leq 4M |f(\mathbf{x}) - g(\mathbf{x})|.$$

Therefore, for the empirical norm on  $\{\mathbf{x}_i\}_{i=1}^n$ ,

$$\|\ell_f - \ell_g\|_n \leq 4M \|f - g\|_n \leq 4M \|f - g\|_\infty, \quad (23)$$

and  $\text{diam}(\mathcal{L}_{M,C}, \|\cdot\|_n) \leq 8M^2$ .

A standard Gaussian symmetrization/concentration argument (see, e.g., Wainwright 2019, Chapter 5) yields that with probability at least  $1 - \delta$ ,

$$\sup_{f \in \mathcal{F}_{M,C}} |R(f) - \widehat{R}_{\mathcal{D}}(f)| \lesssim \widehat{\mathcal{G}}_n(\mathcal{L}_{M,C}) + M^2 \left( \frac{\log(1/\delta)}{n} \right)^{\frac{1}{2}}. \quad (24)$$

It remains to bound  $\widehat{\mathcal{G}}_n(\mathcal{L}_{M,C})$  by chaining. By (23), any  $t$ -cover of  $\mathcal{F}_{M,C}$  in  $\|\cdot\|_n$  induces a  $(4Mt)$ -cover of  $\mathcal{L}_{M,C}$  in  $\|\cdot\|_n$ , hence

$$\log N(t, \mathcal{L}_{M,C}, \|\cdot\|_n) \leq \log N\left(\frac{t}{4M}, \mathcal{F}_{M,C}, \|\cdot\|_n\right) \leq \log N\left(\frac{t}{4M}, V_C(\mathbb{B}_R^d), \|\cdot\|_\infty\right),$$

where we used  $\|h\|_n \leq \|h\|_\infty$ . Proposition C.5 then gives, up to logarithmic factors,

$$\log N(t, \mathcal{L}_{M,C}, \|\cdot\|_n) \lesssim_d \left( \frac{MC}{t} \right)^{\frac{2d}{d+3}}.$$

Applying Proposition C.6 with  $\mathcal{G} = \mathcal{L}_{M,C}$  and  $\text{diam}(\mathcal{L}_{M,C}, \|\cdot\|_n) \leq 8M^2$  yields

$$\begin{aligned} \widehat{\mathcal{G}}_n(\mathcal{L}_{M,C}) &\lesssim_d \frac{1}{\sqrt{n}} \int_0^{8M^2} \left(\frac{MC}{t}\right)^{\frac{d}{d+3}} dt + M^2 \left(\frac{\log(1/\delta)}{n}\right)^{\frac{1}{2}} \\ &\asymp_d C^{\frac{d}{d+3}} M^{\frac{d+6}{d+3}} n^{-\frac{1}{2}} + M^2 \left(\frac{\log(1/\delta)}{n}\right)^{\frac{1}{2}}. \end{aligned}$$

Combining with (24) proves (22).  $\square$

## D. Proof of Theorem 4.2

Let  $\iota_j : \mathbb{R}^m \rightarrow \mathbb{R}^d$  be the dual embedding such that  $\pi_j \circ \iota_j = \text{id}_{\mathbb{R}^m}$ . Then (1) is equivalent to a fully connected neural network in a form of

$$f_{\theta}(\mathbf{x}) = \sum_{k=1}^K \frac{1}{J} \sum_{j=1}^J v_k \phi(\iota_j(\mathbf{w}_k)^{\top} \mathbf{x} - b_k) + \beta. \quad (25)$$

Note that for any fixed  $k$ ,  $\|\iota_j(\mathbf{w}_k)\|_2 = \|\mathbf{w}_k\|_2$  for all  $j = 1, \dots, J$ . Therefore, the notion of path norm and variation norm (together with their weighted version) still make sense for the CNN model

$$\|f_{\theta}\|_{\text{path}} = \frac{1}{J} \sum_{k=1}^K \sum_{j=1}^J |v_k| \|\iota_j(\mathbf{w}_k)\|_2 = \sum_{k=1}^K |v_k| \|\mathbf{w}_k\|_2. \quad (26)$$

By direct computation, the Hessian matrix of the loss function is expressed as

$$\nabla_{\theta}^2 \mathcal{L} = \underbrace{\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} f(\mathbf{x}_i) \nabla_{\theta} f(\mathbf{x}_i)^{\top}}_{\mathbf{T}_{\mathcal{D}}} + \underbrace{\frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i) \nabla_{\theta}^2 f(\mathbf{x}_i)}_{\mathbf{R}_{\mathcal{D}}}. \quad (27)$$

**Definition D.1.** Given a dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  and a set of local receptive fields  $\mathcal{S} = \{S_j\}_{j=1}^J$ , we define a random vector  $\mathbf{X}_{\mathcal{D}}^{\mathcal{S}}$  uniformly draw from  $\{\pi_j(\mathbf{x}_i)\}_{(i,j)}^{n \times J} \subset \mathbb{R}^m$ . For any  $\mathbf{u} \in \mathbb{S}^{m-1}$ ,  $t \in \mathbb{R}$ , we define the weight function

$$g_{\mathcal{D},\mathcal{S}}(\mathbf{u}, t) = \min \left\{ \tilde{g}_{\mathcal{D},\mathcal{S}}(\mathbf{u}, t), \tilde{g}_{\mathcal{D},\mathcal{S}}(-\mathbf{u}, -t) \right\}. \quad (28)$$

where

$$g_{\mathcal{D},\mathcal{S}}(\mathbf{u}, t) := \mathbb{E} [\phi(\mathbf{u}^{\top} \mathbf{X}_{\mathcal{D}}^{\mathcal{S}} - t)] \sqrt{\mathbb{P}(\mathbf{u}^{\top} \mathbf{X}_{\mathcal{D}}^{\mathcal{S}} > t)^2 + \left\| \mathbb{E} \left[ \mathbf{X}_{\mathcal{D}}^{\mathcal{S}} \mathbb{1} \left\{ \mathbf{u}^{\top} \mathbf{X}_{\mathcal{D}}^{\mathcal{S}} > t \right\} \right] \right\|_2^2}. \quad (29)$$

**Proposition D.2.** Given a data set  $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^n$  and a network model (1) with local receptive fields  $\mathcal{S}$ . Then we have

$$\lambda_{\max}(\mathbf{T}_{\mathcal{D}}) \geq 1 + 2 \sum_{k=1}^K |v_k| \|\mathbf{w}_k\| g_{\mathcal{D},\mathcal{S}} \left( \frac{\mathbf{w}_k}{\|\mathbf{w}_k\|}, \frac{b_k}{\|\mathbf{w}_k\|} \right). \quad (30)$$

*Proof.* We write  $\mathbf{T}_{\mathcal{D}}$  in terms of tangent features

$$\mathbf{T}_{\mathcal{D}} = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} f(\mathbf{x}_i) \nabla_{\theta} f(\mathbf{x}_i)^{\top} = \frac{1}{n} \Phi \Phi^{\top}.$$

Consequently,

$$\lambda_{\max}(\mathbf{T}_{\mathcal{D}}) = \max_{\mathbf{u} \in \mathbb{S}^{n-1}} \frac{1}{n} \|\Phi \mathbf{u}\|^2 \geq \frac{1}{n^2} \|\Phi \mathbf{1}\|^2. \quad (31)$$

For any point  $\mathbf{x}$  and abbreviate its patch extraction on  $S_j$  by  $\mathbf{x}^{(S_j)} := \pi_j(\mathbf{x})$ . Define the ReLU gate

$$m_k^{(S_j)}(\mathbf{x}) := \mathbb{1} \left\{ \mathbf{w}_k^\top \mathbf{x}^{(S_j)} > b_k \right\}.$$

For any sample  $\mathbf{x}_i$ , denote  $m_{k,i}^{(S_j)} := m_k^{(S_j)}(\mathbf{x}_i)$ . Then the partial derivatives are

$$\frac{\partial f(\mathbf{x}_i)}{\partial v_k} = \frac{1}{J} \sum_{j=1}^J m_{k,i}^{(S_j)} \cdot (\mathbf{w}_k^\top \mathbf{x}_i^{(S_j)} - b_k), \quad \frac{\partial f(\mathbf{x}_i)}{\partial \mathbf{w}_k} = \frac{1}{J} \sum_{j=1}^J m_{k,i}^{(S_j)} \cdot v_k \cdot \mathbf{x}_i^{(S_j)},$$

$$\frac{\partial f(\mathbf{x}_i)}{\partial b_k} = -\frac{1}{J} \sum_{j=1}^J m_{k,i}^{(S_j)} \cdot v_k, \quad \frac{\partial f(\mathbf{x}_i)}{\partial \beta} = 1.$$

Stacking these over  $i$  and plugging  $u = \mathbf{1}/\sqrt{n}$  in (31), we get

$$\begin{aligned} \frac{1}{n^2} \|\Phi \mathbf{1}\|^2 &= 1 + \frac{1}{n^2} \sum_{k=1}^K \left[ (v_k)^2 \left( \left\| \sum_{i=1}^n \frac{1}{J} \sum_{j=1}^J m_{k,i}^{(S_j)} \mathbf{x}_i^{(S_j)} \right\|^2 + \left( \sum_{i=1}^n \frac{1}{J} \sum_{j=1}^J m_{k,i}^{(S_j)} \right)^2 \right) \right. \\ &\quad \left. + \left( \sum_{i=1}^n \frac{1}{J} \sum_{j=1}^J \phi \left( \mathbf{w}_k^\top \mathbf{x}_i^{(S_j)} - b_k \right) \right)^2 \right] \\ &= 1 + \frac{1}{n^2} \sum_{k=1}^K \left[ (v_k)^2 \left( \left\| \sum_{i=1}^n \frac{1}{J} \sum_{j=1}^J m_{k,i}^{(S_j)} \mathbf{x}_i^{(S_j)} \right\|^2 + \left( \sum_{i=1}^n \frac{1}{J} \sum_{j=1}^J m_{k,i}^{(S_j)} \right)^2 \right) \right. \\ &\quad \left. + \|\mathbf{w}_k\|_2^2 \left( \sum_{i=1}^n \frac{1}{J} \sum_{j=1}^J \phi \left( \mathbf{u}_k^\top \mathbf{x}_i^{(S_j)} - t_k \right) \right)^2 \right] \tag{32} \\ &\geq 1 + \frac{2}{n^2} \sum_{k=1}^K |v_k| \|\mathbf{w}_k\|_2 \left( \sum_{i=1}^n \frac{1}{J} \sum_{j=1}^J \phi \left( \mathbf{u}_k^\top \mathbf{x}_i^{(S_j)} - t_k \right) \right) \\ &\quad \cdot \sqrt{\left\| \sum_{i=1}^n \frac{1}{J} \sum_{j=1}^J m_{k,i}^{(S_j)} \mathbf{x}_i^{(S_j)} \right\|^2 + \left( \sum_{i=1}^n \frac{1}{J} \sum_{j=1}^J m_{k,i}^{(S_j)} \right)^2} \quad (\text{since } a^2 + b^2 \geq 2ab) \\ &= 1 + 2 \sum_{k=1}^K |v_k| \|\mathbf{w}_k\|_2 \left( \frac{1}{nJ} \sum_{i=1}^n \sum_{j=1}^J \phi \left( \mathbf{u}_k^\top \mathbf{x}_i^{(S_j)} - t_k \right) \right) \\ &\quad \cdot \sqrt{\left\| \frac{1}{nJ} \sum_{i=1}^n \sum_{j=1}^J m_{k,i}^{(S_j)} \mathbf{x}_i^{(S_j)} \right\|^2 + \left( \frac{1}{nJ} \sum_{i=1}^n \sum_{j=1}^J m_{k,i}^{(S_j)} \right)^2} \\ &= 1 + 2 \sum_{k=1}^K |v_k| \|\mathbf{w}_k\|_2 \mathbb{E} [\phi(\mathbf{u}^\top \mathbf{X}_D^S - t)] \sqrt{\mathbb{P}(\mathbf{u}^\top \mathbf{X}_D^S > t)^2 + \left\| \mathbb{E}[\mathbf{X}_D^S \mathbb{1}\{\mathbf{u}^\top \mathbf{X}_D^S > t\}] \right\|^2} \\ &= 1 + 2 \sum_{k=1}^K |v_k| \|\mathbf{w}_k\|_2 g_{\mathcal{D},S}(\mathbf{u}_k, t_k), \quad \mathbf{u}_k = \frac{\mathbf{w}_k}{\|\mathbf{w}_k\|_2}, \quad t_k = \frac{b_k}{\|\mathbf{w}_k\|_2}. \end{aligned}$$

Then combining (31) and the above inequality yields the claim.  $\square$

**Lemma D.3.** Consider the model (1)

$$f(\mathbf{x}) = \sum_{k=1}^K \frac{v_k}{J} \sum_{j=1}^J \phi(\mathbf{w}_k^\top \pi_j(\mathbf{x}) - b_k) + \beta.$$

where the input satisfies  $\|\mathbf{x}\|_2 \leq R$ , each patch extractor  $\pi_j: \mathbb{R}^d \rightarrow \mathbb{R}^m$  is a coordinate projection, and  $\phi(t) = \max\{0, t\}$ .

Let  $\boldsymbol{\theta} = (\mathbf{w}_1^\top, \dots, \mathbf{w}_K^\top, b_1, \dots, b_K, v_1, \dots, v_K, \beta)^\top$  collect all parameters. Assume  $f_{\boldsymbol{\theta}}(\mathbf{x})$  is twice differentiable with respect to  $\boldsymbol{\theta}$  at  $\mathbf{x}$ , i.e., for all  $k$  and  $S_j \in \mathcal{S}$  we have  $\mathbf{w}_k^\top \mathbf{x}^{(S_j)} \neq b_k$ . Then for any perturbation vector  $\boldsymbol{\omega}$  with  $\|\boldsymbol{\omega}\|_2 = 1$ , it holds that

$$|\boldsymbol{\omega}^\top \nabla_{\boldsymbol{\theta}}^2 f_{\boldsymbol{\theta}}(\mathbf{x}) \boldsymbol{\omega}| \leq 2(R+1).$$

*Proof.* Write  $\boldsymbol{\theta} = (\mathbf{w}_1^\top, \dots, \mathbf{w}_K^\top, b_1, \dots, b_K, v_1, \dots, v_K, \beta)^\top$ . The total number of parameters is  $N = K \cdot m + K + K + 1 = K(m+2) + 1$ .

Let the corresponding perturbation vector be

$$\boldsymbol{\omega} = (\boldsymbol{\alpha}_1^\top, \dots, \boldsymbol{\alpha}_K^\top, \delta_1, \dots, \delta_K, \gamma_1, \dots, \gamma_K, \iota)^\top \in \mathbb{R}^N,$$

where  $\boldsymbol{\alpha}_k \in \mathbb{R}^m$  corresponds to  $\mathbf{w}_k$ ,  $\delta_k \in \mathbb{R}$  to  $b_k$ ,  $\gamma_k \in \mathbb{R}$  to  $v_k$ , and  $\iota \in \mathbb{R}$  to  $\beta$ . The normalization constraint is

$$\|\boldsymbol{\omega}\|_2^2 = \sum_{k=1}^K \|\boldsymbol{\alpha}_k\|_2^2 + \sum_{k=1}^K \delta_k^2 + \sum_{k=1}^K \gamma_k^2 + \iota^2 = 1.$$

For the fixed input  $\mathbf{x}$ , set  $\mathbf{x}^{(S_j)} := \pi_j(\mathbf{x}) \in \mathbb{R}^m$  and define the ReLU gate

$$m_k^{(S_j)} := \mathbf{1}\{\mathbf{w}_k^\top \mathbf{x}^{(S_j)} > b_k\} \in \{0, 1\}.$$

By the twice-differentiability assumption, all gates are constant in a neighborhood of  $\boldsymbol{\theta}$ .

Within this gate-fixed region,  $f_{\boldsymbol{\theta}}$  is affine in  $(\mathbf{w}_k, b_k)$  once  $v_k$  is held fixed, and affine in  $v_k$  once  $(\mathbf{w}_k, b_k)$  are held fixed. Therefore the only nonzero second partial derivatives inside the  $k$ -th neuron block are the mixed ones with  $v_k$

$$\frac{\partial^2 f_{\boldsymbol{\theta}}}{\partial \mathbf{w}_k \partial v_k} = \frac{\partial}{\partial v_k} \left( \frac{1}{J} \sum_{S_j \in \mathcal{S}} v_k m_k^{(S_j)} \mathbf{x}^{(S_j)} \right) = \frac{1}{J} \sum_{j=1}^J m_k^{(S_j)} \mathbf{x}^{(S_j)} =: \mathbf{s}_k \in \mathbb{R}^m, \quad (33)$$

$$\frac{\partial^2 f_{\boldsymbol{\theta}}}{\partial b_k \partial v_k} = \frac{\partial}{\partial v_k} \left( \frac{1}{J} \sum_{j=1}^J v_k m_k^{(S_j)} \right) = -\frac{1}{J} \sum_{j=1}^J m_k^{(S_j)} =: t_k \in \mathbb{R}. \quad (34)$$

All other second derivatives inside the block vanish (as do any cross-neuron second derivatives and all those involving  $\beta$ ). Hence, with the block variable  $\boldsymbol{\theta}_k := (\mathbf{w}_k^\top, b_k, v_k)^\top$ , we have

$$\nabla_{(\boldsymbol{\theta}_k)}^2 f_{\boldsymbol{\theta}}(\mathbf{x}) = \begin{pmatrix} \mathbf{0}_{m \times m} & \mathbf{0}_m & \mathbf{s}_k \\ \mathbf{0}_m^\top & 0 & t_k \\ \mathbf{s}_k^\top & t_k & 0 \end{pmatrix}. \quad (35)$$

The full quadratic form splits over neuron blocks:

$$\begin{aligned} \boldsymbol{\omega}^\top \nabla_{\boldsymbol{\theta}}^2 f_{\boldsymbol{\theta}}(\mathbf{x}) \boldsymbol{\omega} &= \sum_{k=1}^K (\boldsymbol{\alpha}_k^\top \quad \delta_k \quad \gamma_k) \begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{s}_k \\ \mathbf{0}^\top & 0 & t_k \\ \mathbf{s}_k^\top & t_k & 0 \end{pmatrix} \begin{pmatrix} \boldsymbol{\alpha}_k \\ \delta_k \\ \gamma_k \end{pmatrix} \\ &= \sum_{k=1}^K 2\gamma_k (\boldsymbol{\alpha}_k^\top \mathbf{s}_k + \delta_k t_k). \end{aligned} \quad (36)$$

Recall the definition of the notations (33) and (34),

$$|t_k| = \frac{1}{J} \sum_{j=1}^J m_k^{(S_j)} \leq 1, \quad (37)$$

$$\|\mathbf{s}_k\|_2 = \left\| \frac{1}{J} \sum_{j=1}^J m_k^{(S_j)} \mathbf{x}^{(S_j)} \right\|_2 \leq \frac{1}{J} \sum_{j=1}^J m_k^{(S_j)} \|\mathbf{x}^{(S_j)}\|_2 \leq \frac{1}{J} \sum_{j=1}^J m_k^{(S_j)} R = R|t_k|, \quad (38)$$

because  $\mathbf{x}^{(S_j)} = \pi_j(\mathbf{x})$  is a coordinate projection of  $\mathbf{x}$ ,  $\|\mathbf{x}^{(S_j)}\|_2 \leq \|\mathbf{x}\|_2 \leq R$ .

Then for each  $k$ ,

$$\begin{aligned} |2\gamma_k(\boldsymbol{\alpha}_k^\top \mathbf{s}_k + \delta_k t_k)| &\leq 2|\gamma_k| \left( \|\boldsymbol{\alpha}_k\|_2 \|\mathbf{s}_k\|_2 + |\delta_k| t_k \right) \quad (\text{Cauchy-Schwarz}) \\ (38) \implies &\leq 2|\gamma_k| \left( R t_k \|\boldsymbol{\alpha}_k\|_2 + |\delta_k| t_k \right) \\ (37) \implies &\leq 2|\gamma_k| \left( R \|\boldsymbol{\alpha}_k\|_2 + |\delta_k| \right), \end{aligned} \quad (39)$$

Summing over  $k$  for (39) and plug in (36), we have

$$\begin{aligned} |\boldsymbol{\omega}^\top \nabla_{\boldsymbol{\theta}}^2 f_{\boldsymbol{\theta}}(\mathbf{x}) \boldsymbol{\omega}| &\leq 2 \left( R \sum_k |\gamma_k| \|\boldsymbol{\alpha}_k\|_2 + \sum_k |\gamma_k| |\delta_k| \right) \\ (\text{Cauchy-Schwarz}) &\leq 2 \left( R \sqrt{\sum_k \gamma_k^2} \sqrt{\sum_k \|\boldsymbol{\alpha}_k\|_2^2} + \sqrt{\sum_k \gamma_k^2} \sqrt{\sum_k \delta_k^2} \right). \end{aligned} \quad (40)$$

Using the normalization  $\sum_k \|\boldsymbol{\alpha}_k\|_2^2 + \sum_k \delta_k^2 + \sum_k \gamma_k^2 + \iota^2 = 1$ , we have  $\sqrt{\sum_k \gamma_k^2} \leq 1$ ,  $\sqrt{\sum_k \|\boldsymbol{\alpha}_k\|_2^2} \leq 1$ , and  $\sqrt{\sum_k \delta_k^2} \leq 1$ . Thus  $|\boldsymbol{\omega}^\top \nabla_{\boldsymbol{\theta}}^2 f_{\boldsymbol{\theta}}(\mathbf{x}) \boldsymbol{\omega}| \leq 2(R+1)$ .  $\square$

**Theorem D.4** (Restate Theorem 4.2). *Given a data set  $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^n$  and a network model (1) with the set of local receptive fields  $\mathcal{S}$ , we have*

$$\sum_{k=1}^K |v_k| \|\mathbf{w}_k\|_{g_{\mathcal{D}, \mathcal{S}}} \left( \frac{\mathbf{w}_k}{\|\mathbf{w}_k\|}, \frac{b_k}{\|\mathbf{w}_k\|} \right) \leq \frac{1}{2} \left( \lambda_{\max}(\nabla_{\boldsymbol{\theta}}^2 \mathcal{L}(\boldsymbol{\theta})) + 2(R+1)\sqrt{2\mathcal{L}(\boldsymbol{\theta})} - 1 \right).$$

*Proof.* It suffices to prove the first assertion. Recall that by direct computation, the Hessian matrix of the loss function is expressed as

$$\nabla_{\boldsymbol{\theta}}^2 \mathcal{L} = \underbrace{\frac{1}{n} \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_i) \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_i)^\top}_{\mathbf{T}_{\mathcal{D}}} + \underbrace{\frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i) \nabla_{\boldsymbol{\theta}}^2 f(\mathbf{x}_i)}_{\mathbf{R}_{\mathcal{D}}}. \quad (41)$$

Let  $\boldsymbol{\omega}$  be the unit eigenvector (i.e.,  $\|\boldsymbol{\omega}\|_2 = 1$ ) corresponding to the largest eigenvalue of the matrix  $\mathbf{T}_{\mathcal{D}}$ , the maximum eigenvalue of the Hessian matrix of the loss can be lower-bounded as follows:

$$\begin{aligned} \lambda_{\max}(\nabla_{\boldsymbol{\theta}}^2 \mathcal{L}(\boldsymbol{\theta})) &\geq \boldsymbol{\omega}^\top \nabla_{\boldsymbol{\theta}}^2 \mathcal{L}(\boldsymbol{\theta}) \boldsymbol{\omega} \\ &\geq \underbrace{\lambda_{\max}(\mathbf{T}_{\mathcal{D}})}_{(\text{Term A})} + \underbrace{\frac{1}{n} \sum_{i=1}^n (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i) \boldsymbol{\omega}^\top \nabla_{\boldsymbol{\theta}}^2 f_{\boldsymbol{\theta}}(\mathbf{x}_i) \boldsymbol{\omega}}_{(\text{Term B})}. \end{aligned} \quad (42)$$

According to Proposition D.2, Term A is lower bounded by

$$\lambda_{\max}(\mathbf{T}_{\mathcal{D}}) \geq 1 + 2 \sum_{k=1}^K |v_k| \|\mathbf{w}_k\|_{g_{\mathcal{D}, \mathcal{S}}} \left( \frac{\mathbf{w}_k}{\|\mathbf{w}_k\|}, \frac{b_k}{\|\mathbf{w}_k\|} \right), \quad (43)$$

For (Term B), an upper bound can be established using the training loss  $\mathcal{L}(\theta)$  via the Cauchy-Schwarz inequality. This also employs a notable uniform upper bound for  $|\omega^\top \nabla_{\theta}^2 f_{\theta}(\mathbf{x}_n) \omega|$ , as detailed in Lemma D.3:

$$|(\text{Term B})| \leq \sqrt{\frac{1}{n} \sum_{i=1}^n (f_{\theta}(\mathbf{x}_i) - y_i)^2} \cdot \sqrt{\frac{1}{n} \sum_{i=1}^n (\omega^\top \nabla_{\theta}^2 f_{\theta}(\mathbf{x}_i) \omega)^2} \leq 2(R+1) \sqrt{2\mathcal{L}(\theta)}. \quad (44)$$

Thus, we have

$$\lambda_{\max}(\mathbf{T}_{\mathcal{D}}) \leq \lambda_{\max}(\nabla_{\theta}^2 \mathcal{L}(\theta)) + |(\text{Term B})| \leq \lambda_{\max}(\nabla_{\theta}^2 \mathcal{L}(\theta)) + 2(R+1) \sqrt{2\mathcal{L}(\theta)}. \quad (45)$$

Finally, we plug (43) into (45)

$$\sum_{k=1}^K |v_k| \|\mathbf{w}_k\| g_{\mathcal{D},S} \left( \frac{\mathbf{w}_k}{\|\mathbf{w}_k\|}, \frac{b_k}{\|\mathbf{w}_k\|} \right) \leq \frac{1}{2} \left( \lambda_{\max}(\nabla_{\theta}^2 \mathcal{L}(\theta)) + 2(R+1) \sqrt{2\mathcal{L}(\theta)} - 1 \right).$$

□

### D.1. Empirical Process for the Weight Function

We study uniform deviations of the empirical weight function  $g_{\mathcal{D},S}$  from its population counterpart  $g_{\mathcal{P},S}$  under i.i.d. sampling  $\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathcal{P}_{\mathcal{X}}$ . Although the collection  $\{\pi_j(\mathbf{x}_i)\}_{i,j}$  is not i.i.d. across  $(i, j)$  in general, for each fixed  $j$  the patches  $\{\pi_j(\mathbf{x}_i)\}_{i=1}^n$  are i.i.d. in  $\mathbb{R}^m$ .

**Definition D.5.** For any  $S_j \in \mathcal{S}$ , let  $\mathbf{X}_{\mathcal{D}}^{(S_j)}$  be a random vector drawn uniformly at random from the training examples  $\{\pi_j(\mathbf{x}_i)\}_{i=1}^n$  with patch extraction to the local receptive field  $S_j$ . For any  $\mathbf{u} \in \mathbb{S}^{m-1}$ ,  $t \in \mathbb{R}$ , define the aggregated gate probability, ReLU margin, and gated first moment

$$\begin{aligned} \bar{p}_{\mathcal{D},S}(\mathbf{u}, t) &:= \frac{1}{J} \sum_{j=1}^J \mathbb{P} \left( \mathbf{u}^\top \mathbf{X}_{\mathcal{D}}^{(S_j)} > t \right), \\ \bar{r}_{\mathcal{D},S}(\mathbf{u}, t) &:= \frac{1}{J} \sum_{j=1}^J \mathbb{E} \left[ \phi \left( \mathbf{u}^\top \mathbf{X}_{\mathcal{D}}^{(S_j)} - t \right) \right], \\ \bar{\mathbf{a}}_{\mathcal{D},S}(\mathbf{u}, t) &:= \frac{1}{J} \sum_{j=1}^J \mathbb{E} \left[ \mathbf{X}_{\mathcal{D}}^{(S_j)} \mathbb{1} \left\{ \mathbf{u}^\top \mathbf{X}_{\mathcal{D}}^{(S_j)} > t \right\} \right] \in \mathbb{R}^m. \end{aligned} \quad (46)$$

Based on these quantities, define the empirical weight function

$$g_{\mathcal{D},S}(\mathbf{u}, t) := \bar{r}_{\mathcal{D},S}(\mathbf{u}, t) \sqrt{\bar{p}_{\mathcal{D},S}(\mathbf{u}, t)^2 + \|\bar{\mathbf{a}}_{\mathcal{D},S}(\mathbf{u}, t)\|_2^2}. \quad (47)$$

For each  $j \in [J]$  and  $(\mathbf{u}, t) \in \mathbb{S}^{m-1} \times [-1, 1]$ , define the population components

$$p_j(\mathbf{u}, t) := \mathbb{P}(\mathbf{u}^\top \pi_j(\mathbf{X}) > t), \quad r_j(\mathbf{u}, t) := \mathbb{E}[\phi(\mathbf{u}^\top \pi_j(\mathbf{X}) - t)], \quad \mathbf{a}_j(\mathbf{u}, t) := \mathbb{E}[\pi_j(\mathbf{X}) \mathbb{1}\{\mathbf{u}^\top \pi_j(\mathbf{X}) > t\}].$$

Let  $\bar{p}_{\mathcal{P},S}$ ,  $\bar{r}_{\mathcal{P},S}$ , and  $\bar{\mathbf{a}}_{\mathcal{P},S}$  denote their averages over  $j \in [J]$ . Define

$$g_{\mathcal{P},S}(\mathbf{u}, t) := \bar{r}_{\mathcal{P},S}(\mathbf{u}, t) \sqrt{\bar{p}_{\mathcal{P},S}(\mathbf{u}, t)^2 + \|\bar{\mathbf{a}}_{\mathcal{P},S}(\mathbf{u}, t)\|_2^2}.$$

**Lemma D.6** (Complexity of the component classes). *Under the constraint  $\|\mathbf{z}\|_2 \leq 1$ , the following classes are uniformly bounded VC-subgraph classes with VC-subgraph dimension  $O(m)$ :*

$$\mathcal{H} := \{ \mathbf{z} \mapsto \mathbb{1}\{\mathbf{u}^\top \mathbf{z} > t\} : \mathbf{u} \in \mathbb{S}^{m-1}, t \in [-1, 1] \},$$

$$\mathcal{F}_{\text{ReLU}} := \{ \mathbf{z} \mapsto \phi(\mathbf{u}^\top \mathbf{z} - t) : \mathbf{u} \in \mathbb{S}^{m-1}, t \in [-1, 1] \},$$

and

$$\mathcal{A} := \{ \mathbf{z} \mapsto \mathbf{v}^\top \mathbf{z} \mathbb{1}\{\mathbf{u}^\top \mathbf{z} > t\} : \mathbf{u}, \mathbf{v} \in \mathbb{S}^{m-1}, t \in [-1, 1] \}.$$

Moreover,  $\mathcal{H}$  and  $\mathcal{A}$  are bounded by 1, and  $\mathcal{F}_{\text{ReLU}}$  is bounded by 2.

*Proof.* The class  $\mathcal{H}$  is the class of halfspace indicators in  $\mathbb{R}^m$ , hence has VC-dimension at most  $m + 1$ .

For  $\mathcal{F}_{\text{ReLU}}$ , the subgraph of  $z \mapsto \phi(\mathbf{u}^\top z - t)$  is

$$\{(z, s) : s < \phi(\mathbf{u}^\top z - t)\} = (\{\mathbf{u}^\top z > t\} \cap \{s < \mathbf{u}^\top z - t\}) \cup (\{\mathbf{u}^\top z \leq t\} \cap \{s < 0\}).$$

For  $\mathcal{A}$ , the subgraph of  $z \mapsto \mathbf{v}^\top z \mathbb{1}\{\mathbf{u}^\top z > t\}$  is

$$\{(z, s) : s < \mathbf{v}^\top z \mathbb{1}\{\mathbf{u}^\top z > t\}\} = (\{\mathbf{u}^\top z > t\} \cap \{s < \mathbf{v}^\top z\}) \cup (\{\mathbf{u}^\top z \leq t\} \cap \{s < 0\}).$$

Both subgraph classes are constant-size Boolean combinations of halfspaces in  $\mathbb{R}^{m+1}$ . By the standard growth-function closure argument for finite Boolean combinations of VC classes, their VC dimensions are  $O(m)$ . Equivalently, their underlying real-valued function classes have pseudo-dimension  $O(m)$ , see (Anthony & Bartlett, 1999, Ch. 3 and Definition 11.2). The boundedness claims follow from  $\|z\|_2 \leq 1$ ,  $t \in [-1, 1]$ , and  $\|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1$ .  $\square$

**Theorem D.7** (Uniform deviation for  $g_{\mathcal{D}, S}$ ). *Assume  $\|\pi_j(\mathbf{X})\|_2 \leq 1$  almost surely for all  $j \in [J]$ . There exists a universal constant  $C_{\text{ep}} > 0$  such that for every  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,*

$$\sup_{\mathbf{u} \in \mathbb{S}^{m-1}, t \in [-1, 1]} |g_{\mathcal{D}, S}(\mathbf{u}, t) - g_{\mathcal{P}, S}(\mathbf{u}, t)| \leq C_{\text{ep}} \sqrt{\frac{m + \log(2J/\delta)}{n}}.$$

*Proof.* For each fixed  $j \in [J]$ , let  $P_j$  denote the law of  $\mathbf{Z}_j := \pi_j(\mathbf{X}) \in \mathbb{R}^m$ , and let  $P_{n,j} := n^{-1} \sum_{i=1}^n \delta_{\pi_j(\mathbf{x}_i)}$ . Define

$$\hat{p}_j(\mathbf{u}, t) := P_{n,j} \mathbb{1}\{\mathbf{u}^\top z > t\}, \quad \hat{r}_j(\mathbf{u}, t) := P_{n,j} \phi(\mathbf{u}^\top z - t), \quad \hat{\mathbf{a}}_j(\mathbf{u}, t) := P_{n,j} [z \mathbb{1}\{\mathbf{u}^\top z > t\}].$$

By the VC-subgraph entropy bound and maximal inequality (van der Vaart & Wellner, 1996, Theorems 2.6.7 and 2.14.1), together with McDiarmid's inequality, any class  $\mathcal{F}$  uniformly bounded by  $B$  with VC-subgraph dimension  $V$  satisfies

$$\sup_{f \in \mathcal{F}} |(P_n - P)f| \lesssim B \sqrt{\frac{V + \log(1/\eta)}{n}}$$

with probability at least  $1 - \eta$ .

Applying this bound to  $\mathcal{H}$  and  $\mathcal{F}_{\text{ReLU}}$  from Lemma D.6, for every fixed  $j$  and  $\eta \in (0, 1)$ , with probability at least  $1 - \eta$ ,

$$\sup_{\mathbf{u}, t} |\hat{p}_j(\mathbf{u}, t) - p_j(\mathbf{u}, t)| \lesssim \sqrt{\frac{m + \log(1/\eta)}{n}}, \quad \sup_{\mathbf{u}, t} |\hat{r}_j(\mathbf{u}, t) - r_j(\mathbf{u}, t)| \lesssim \sqrt{\frac{m + \log(1/\eta)}{n}},$$

where all suprema are over  $\mathbf{u} \in \mathbb{S}^{m-1}$  and  $t \in [-1, 1]$ .

For the vector moment, use duality:

$$\begin{aligned} \sup_{\mathbf{u}, t} \|\hat{\mathbf{a}}_j(\mathbf{u}, t) - \mathbf{a}_j(\mathbf{u}, t)\|_2 &= \sup_{\mathbf{u}, t} \sup_{\mathbf{v} \in \mathbb{S}^{m-1}} |\mathbf{v}^\top (\hat{\mathbf{a}}_j(\mathbf{u}, t) - \mathbf{a}_j(\mathbf{u}, t))| \\ &= \sup_{\mathbf{u}, \mathbf{v}, t} |(P_{n,j} - P_j) [\mathbf{v}^\top z \mathbb{1}\{\mathbf{u}^\top z > t\}]|. \end{aligned}$$

The last supremum is over the class  $\mathcal{A}$  from Lemma D.6. Hence, with probability at least  $1 - \eta$ ,

$$\sup_{\mathbf{u}, t} \|\hat{\mathbf{a}}_j(\mathbf{u}, t) - \mathbf{a}_j(\mathbf{u}, t)\|_2 \lesssim \sqrt{\frac{m + \log(1/\eta)}{n}}.$$

Taking  $\eta = \delta/(3J)$  and applying a union bound over the three component bounds and over  $j \in [J]$ , we obtain that, with probability at least  $1 - \delta$ , simultaneously for all  $j \in [J]$ ,

$$\sup_{\mathbf{u}, t} |\hat{p}_j(\mathbf{u}, t) - p_j(\mathbf{u}, t)| \leq \varepsilon, \quad \sup_{\mathbf{u}, t} |\hat{r}_j(\mathbf{u}, t) - r_j(\mathbf{u}, t)| \leq \varepsilon, \quad \sup_{\mathbf{u}, t} \|\hat{\mathbf{a}}_j(\mathbf{u}, t) - \mathbf{a}_j(\mathbf{u}, t)\|_2 \leq \varepsilon,$$

where all suprema are over  $\mathbf{u} \in \mathbb{S}^{m-1}$  and  $t \in [-1, 1]$ , and  $\varepsilon := C_0 \sqrt{(m + \log(3J/\delta))/n}$  for a universal constant  $C_0 > 0$ . Averaging over  $j$  does not increase the deviations. Therefore, uniformly over  $(\mathbf{u}, t) \in \mathbb{S}^{m-1} \times [-1, 1]$ ,

$$|\bar{p}_{\mathcal{D},S}(\mathbf{u}, t) - \bar{p}_{\mathcal{P},S}(\mathbf{u}, t)| \leq \varepsilon, \quad |\bar{r}_{\mathcal{D},S}(\mathbf{u}, t) - \bar{r}_{\mathcal{P},S}(\mathbf{u}, t)| \leq \varepsilon, \quad \|\bar{\mathbf{a}}_{\mathcal{D},S}(\mathbf{u}, t) - \bar{\mathbf{a}}_{\mathcal{P},S}(\mathbf{u}, t)\|_2 \leq \varepsilon.$$

It remains to pass from the components to  $g$ . For both empirical and population quantities,

$$0 \leq \bar{p} \leq 1, \quad 0 \leq \bar{r} \leq 2, \quad \|\bar{\mathbf{a}}\|_2 \leq \bar{p} \leq 1.$$

Define  $F(r, p, \mathbf{a}) := r\sqrt{p^2 + \|\mathbf{a}\|_2^2}$ . On this domain,

$$|F(r, p, \mathbf{a}) - F(r', p', \mathbf{a}')| \leq \sqrt{2}|r - r'| + 2\sqrt{(p - p')^2 + \|\mathbf{a} - \mathbf{a}'\|_2^2}.$$

Thus, uniformly over  $(\mathbf{u}, t) \in \mathbb{S}^{m-1} \times [-1, 1]$ ,

$$|g_{\mathcal{D},S}(\mathbf{u}, t) - g_{\mathcal{P},S}(\mathbf{u}, t)| \lesssim \varepsilon \lesssim \sqrt{\frac{m + \log(2J/\delta)}{n}},$$

after adjusting the universal constant. This proves the theorem.  $\square$

## D.2. The Computation of the Population Weight Functions for Uniform Distributions on Spheres

Assume that  $\mathcal{P}_{\mathbf{X}}$  is a uniform distribution on  $\mathbb{S}^{d-1} \subset \mathbb{R}^d$ . Here we compute the population version of the components in the weighted function  $g_{\mathcal{D},S}$  defined in Definition D.5.

Let  $d \geq 2$ ,  $m < d$ . Draw  $\mathbf{X} \sim \text{Uniform}(\mathbb{S}^{d-1}) \subset \mathbb{R}^d$ , set  $\mathbf{Z} = \text{proj}(\mathbf{X}) \in \mathbb{R}^m$  by projecting the first  $m$  coordinates, and fix any unit  $\mathbf{u} \in \mathbb{R}^m$ . Define

$$Z := \mathbf{u}^\top \mathbf{Z} \in [-1, 1]. \quad (48)$$

By rotational invariance in  $\mathbb{R}^d$ ,  $Z$  has the one-coordinate marginal of  $\text{Uniform}(\mathbb{S}^{d-1})$ , hence its density is

$$f_d(z) = c_d (1 - z^2)^{\frac{d-3}{2}} \mathbb{1}\{-1 \leq z \leq 1\}, \quad c_d := \frac{\Gamma(\frac{d}{2})}{\sqrt{\pi} \Gamma(\frac{d-1}{2})}. \quad (49)$$

All bounds below are independent of  $m$  and  $\mathbf{u}$ , depending only on  $d$ .

For shorthand, write  $\alpha := \frac{d-3}{2} \geq -\frac{1}{2}$  and observe

$$(1 - z^2)^\alpha = [(1 - z)(1 + z)]^\alpha, \quad z \in [-1, 1], \quad (50)$$

together with the elementary bounds  $1 \leq 1 + z \leq 2$  for  $z \in [0, 1]$ .

**Lemma D.8** ((Tail of  $Z$ )). *For all  $t \in [0, 1]$ ,*

$$\frac{c_d}{\alpha + 1} (1 - t)^{\alpha+1} \leq \mathbb{P}(Z > t) \leq \frac{c_d 2^\alpha}{\alpha + 1} (1 - t)^{\alpha+1}. \quad (51)$$

Consequently  $\mathbb{P}(Z > t) \asymp (1 - t)^{\frac{d-1}{2}}$  as  $t \rightarrow 1^-$ .

*Proof.* Using (49) and (50),

$$\mathbb{P}(Z > t) = \int_t^1 c_d (1 - z^2)^\alpha dz = c_d \int_t^1 [(1 - z)(1 + z)]^\alpha dz. \quad (52)$$

Since  $1 \leq 1 + z \leq 2$  for  $z \in [t, 1]$ , we have

$$c_d \int_t^1 (1 - z)^\alpha dz \leq \mathbb{P}(Z > t) \leq c_d 2^\alpha \int_t^1 (1 - z)^\alpha dz. \quad (53)$$

Evaluating  $\int_t^1 (1 - z)^\alpha dz = \frac{(1-t)^{\alpha+1}}{\alpha+1}$  (valid for  $\alpha > -1$ ; here  $\alpha \geq -\frac{1}{2}$ ), we obtain (51).  $\square$

1485 **Lemma D.9.** For all  $t \in [0, 1)$ ,

$$1486 \frac{c_d}{(\alpha + 1)(\alpha + 2)} (1 - t)^{\alpha + 2} \leq \mathbb{E}[\phi(Z - t)] \leq \frac{c_d 2^\alpha}{(\alpha + 1)(\alpha + 2)} (1 - t)^{\alpha + 2}. \quad (54)$$

1487  
1488  
1489 Consequently  $\mathbb{E}[\phi(Z - t)] \asymp (1 - t)^{\frac{d+1}{2}}$  as  $t \uparrow 1$ .

1490  
1491 *Proof.* By definition and (49),

$$1492 \mathbb{E}[\phi(Z - t)] = \int_t^1 (z - t) f_d(z) dz = c_d \int_t^1 (z - t) [(1 - z)(1 + z)]^\alpha dz. \quad (55)$$

1493  
1494 Bounding  $1 \leq 1 + z \leq 2$  for  $z \in [t, 1]$  yields

$$1495 c_d \int_t^1 (z - t)(1 - z)^\alpha dz \leq \mathbb{E}[\phi(Z - t)] \leq c_d 2^\alpha \int_t^1 (z - t)(1 - z)^\alpha dz. \quad (56)$$

1496  
1497 Substitute  $y = 1 - z$  to compute the shared integral:

$$1498 \int_t^1 (z - t)(1 - z)^\alpha dz = \int_0^{1-t} [(1 - t) - y] y^\alpha dy = \frac{(1 - t)^{\alpha + 2}}{(\alpha + 1)(\alpha + 2)}, \quad (57)$$

1499 valid for  $\alpha > -1$ . This gives (54). □

1500  
1501 **Proposition D.10.** There exist absolute constants (depending only on  $d$ )

$$1502 c_L(d) := \frac{c_d^2}{(\alpha + 1)^2 (\alpha + 2)}, \quad c_U(d) := \frac{c_d^2 2^{2\alpha}}{(\alpha + 1)^2 (\alpha + 2)}, \quad (58)$$

1503  
1504 such that for all  $t \in [0, 1)$ ,

$$1505 c_L(d) (1 - t)^{2\alpha + 3} \leq \mathbb{P}(Z > t) \cdot \mathbb{E}[\phi(Z - t)] \leq c_U(d) (1 - t)^{2\alpha + 3}. \quad (59)$$

1506  
1507 Equivalently, since  $\alpha = \frac{d-3}{2}$ ,

$$1508 \mathbb{P}(Z > t) \cdot \mathbb{E}[\phi(Z - t)] \asymp (1 - t)^d \quad \text{as } t \uparrow 1, \quad (60)$$

1509  
1510 with explicit

$$1511 c_L(d) = \frac{c_d^2}{\left(\frac{d-1}{2}\right)^2 \left(\frac{d+1}{2}\right)}, \quad c_U(d) = \frac{c_d^2 2^{d-3}}{\left(\frac{d-1}{2}\right)^2 \left(\frac{d+1}{2}\right)}, \quad c_d = \frac{\Gamma\left(\frac{d}{2}\right)}{\sqrt{\pi} \Gamma\left(\frac{d-1}{2}\right)}. \quad (61)$$

1512  
1513 *Proof.* The inequality can be deduced by combining (51) and (54). □

1514  
1515 **Proposition D.11** (Boundary tail for projected radius). Fix integers  $d \geq 2$  and  $1 \leq m < d$ . Let  $\mathbf{X} \sim \text{Uniform}(\mathbb{S}^{d-1}) \subset \mathbb{R}^d$   
1516 and  $\mathbf{Z} = \text{proj}(\mathbf{X}) \in \mathbb{R}^m$ . Then there exist universal constants

$$1517 c_L(d, m), c_U(d, m) \in (0, \infty)$$

1518 depending only on  $d, m$  (and independent of the choice of projection and of  $t$ ) such that for all  $t \in (0, \frac{1}{4}]$ ,

$$1519 c_L(d, m) t^{\frac{d-m}{2}} \leq \mathbb{P}(\|\mathbf{Z}\| > 1 - t) \leq c_U(d, m) t^{\frac{d-m}{2}}. \quad (62)$$

1520  
1521 In particular,

$$1522 \mathbb{P}(\|\mathbf{Z}\| > 1 - t) \asymp t^{\frac{d-m}{2}} \quad \text{as } t \downarrow 0, \quad (63)$$

1523  
1524 with one admissible choice

$$1525 c_L(d, m) = \frac{2^{-|\frac{m}{2}-1|}}{\frac{d-m}{2} B\left(\frac{m}{2}, \frac{d-m}{2}\right)}, \quad c_U(d, m) = \frac{2^{|\frac{m}{2}-1|+\frac{d-m}{2}}}{\frac{d-m}{2} B\left(\frac{m}{2}, \frac{d-m}{2}\right)}, \quad (64)$$

1526  
1527 where  $B(a, b) = \Gamma(a)\Gamma(b)/\Gamma(a + b)$  denotes the Beta function.

1540 *Proof.* Write  $R := \|\mathbf{Z}\| \in [0, 1]$  and  $S := R^2$ . A standard Gaussian-ratio representation gives

$$1541 \quad S \sim \text{Beta}(a, b), \quad a := \frac{m}{2}, \quad b := \frac{d-m}{2}, \quad (65)$$

1542 with density

$$1543 \quad f_S(x) = \frac{x^{a-1}(1-x)^{b-1}}{B(a, b)} \mathbf{1}_{\{0 < x < 1\}}. \quad (66)$$

1544 For  $t \in (0, 1)$ ,

$$1545 \quad \mathbb{P}(R > 1-t) = \mathbb{P}(S > (1-t)^2) = \frac{1}{B(a, b)} \int_{(1-t)^2}^1 x^{a-1}(1-x)^{b-1} dx. \quad (67)$$

1546 Set  $y := 1 - (1-t)^2 = 2t - t^2$ . Then  $y \in (0, 2t)$  and for  $t \in (0, \frac{1}{4}]$  we have  $y \leq \frac{1}{2}$  and  $(1-y) \geq \frac{1}{2}$ . Since  $x \in [1-y, 1]$ , the elementary bound

$$1547 \quad 2^{-|a-1|} \leq x^{a-1} \leq 2^{|a-1|} \quad \text{for all } x \in [\frac{1}{2}, 1] \quad (68)$$

1548 holds deterministically (because  $\ln x \in [-\ln 2, 0]$  and  $x^{a-1} = e^{(a-1)\ln x}$ ). Using (68) in (67) and integrating the  $(1-x)^{b-1}$  part exactly gives

$$1549 \quad \frac{2^{-|a-1|}}{B(a, b)} \int_{1-y}^1 (1-x)^{b-1} dx \leq \mathbb{P}(R > 1-t) \leq \frac{2^{|a-1|}}{B(a, b)} \int_{1-y}^1 (1-x)^{b-1} dx. \quad (69)$$

1550 Since  $\int_{1-y}^1 (1-x)^{b-1} dx = \frac{y^b}{b}$ , we obtain

$$1551 \quad \frac{2^{-|a-1|}}{b B(a, b)} y^b \leq \mathbb{P}(R > 1-t) \leq \frac{2^{|a-1|}}{b B(a, b)} y^b. \quad (70)$$

1552 Finally, because  $t \leq y \leq 2t$ , we have  $t^b \leq y^b \leq (2t)^b$ , and (70) yields

$$1553 \quad \frac{2^{-|a-1|}}{b B(a, b)} t^b \leq \mathbb{P}(R > 1-t) \leq \frac{2^{|a-1|+b}}{b B(a, b)} t^b, \quad (71)$$

1554 which is exactly (62) with the explicit constants in (64). This proves  $\mathbb{P}(\|\mathbf{Z}\| > 1-t) \asymp t^b = t^{(d-m)/2}$  as  $t \downarrow 0$ .  $\square$

## 1555 E. Proof of Theorem 4.3

1556 **Theorem E.1** (Detailed version of Theorem 4.3). *Suppose  $\mathcal{P}$  is a joint distribution of  $(\mathbf{x}, y)$ . Assume that the marginal distribution of  $\mathbf{x}$  is  $\text{Uniform}(\mathbb{S}^{d-1})$  and the marginal distribution of  $y$  is supported on  $[-D, D]$  for some  $D > 0$ . Fix a dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , where each data point is drawn i.i.d. from  $\mathcal{P}$ . Let  $M \geq D$ . If  $\theta \in \Theta_{\text{BEoS}}^{\text{SCN}}(\eta, \mathcal{D})$  and  $\|f_\theta\|_\infty \leq M$ , then, with probability  $\geq 1 - 2\delta$ , we have that for the plug-in risk estimator  $\widehat{\mathcal{R}}_{\mathcal{D}}(f) := \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2$ ,*

$$1557 \quad \left| \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} [(f_\theta(\mathbf{x}) - y)^2] - \widehat{\mathcal{R}}(f_\theta) \right| \quad (72)$$

$$1558 \quad \lesssim_d J M^2 \varepsilon^{\frac{d-m}{2}} + \left( A \varepsilon^{-d} \right)^{\frac{d}{d+3}} M^{\frac{d+6}{d+3}} n^{-\frac{1}{2}},$$

1559 for any  $\varepsilon \in (0, 1)$  such that  $\varepsilon^d \gtrsim d^2 \sqrt{\frac{m+\log(J/\delta)}{n}}$ <sup>1</sup>. Here  $A = \frac{1}{\eta} - \frac{1}{2} + 4M$ , and  $\lesssim_d$  hides constants depending only on  $d$  and logarithmic factors in  $n$  and  $(J/\delta)$ .

1560 *Proof.* For  $\varepsilon \in (0, 1)$ , we define

$$1561 \quad \mathcal{I}_\varepsilon^{\text{all}} := \left\{ \mathbf{x} \in \mathbb{S}^{d-1} : \max_{S_j \in \mathcal{S}} \|\mathbf{x}^{(S_j)}\| \leq 1 - \varepsilon \right\}, \quad (73)$$

$$1562 \quad \mathcal{O}_\varepsilon^{\text{any}} := \left\{ \mathbf{x} \in \mathbb{S}^{d-1} : \exists S_j \in \mathcal{S}, \|\mathbf{x}^{(S_j)}\| > 1 - \varepsilon \right\}. \quad (74)$$

1563 <sup>1</sup>We only need  $\text{poly}(d)$  samples to make the feasible choice of  $\varepsilon$  non-vacuous. Here we hide the universal constant.

1595 Then  $\mathbb{S}^{d-1} = \mathbb{I}_\varepsilon^{\text{all}} \cup \mathbb{O}_\varepsilon^{\text{any}}$  (disjoint). According to the law of total expectation, the population risk is decomposed into

$$1596$$

$$1597 \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} \left[ (f(\mathbf{x}) - y)^2 \right] = \mathbb{P}(\mathbf{x} \in \mathbb{O}_\varepsilon^{\text{any}}) \cdot \mathbb{E}_{\mathbb{O}} \left[ (f(\mathbf{x}) - y)^2 \right] + \mathbb{P}(\mathbf{x} \in \mathbb{I}_\varepsilon^{\text{all}}) \cdot \mathbb{E}_{\mathbb{I}} \left[ (f(\mathbf{x}) - y)^2 \right], \quad (75)$$

1599 where  $\mathbb{E}_{\mathbb{O}}$  means that  $\{\mathbf{x}, y\}$  is a new sample from the data distribution conditioned on  $\mathbf{x} \in \mathbb{O}_\varepsilon^{\text{any}}$  and  $\mathbb{E}_{\mathbb{I}}$  means that  $(\mathbf{x}, y)$   
1600 is a new sample from the data distribution conditioned on  $\mathbf{x} \in \mathbb{I}_\varepsilon^{\text{all}}$ .

1602 Similarly, we also have this decomposition for empirical risk

$$1603$$

$$1604 \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2 = \frac{1}{n} \left( \sum_{i \in I} (f(\mathbf{x}_i) - y_i)^2 + \sum_{i \in O} (f(\mathbf{x}_i) - y_i)^2 \right)$$

$$1605 = \frac{n_I}{n} \frac{1}{n_I} \sum_{i \in I} (f(\mathbf{x}_i) - y_i)^2 + \frac{n_O}{n} \frac{1}{n_O} \sum_{i \in O} (f(\mathbf{x}_i) - y_i)^2, \quad (76)$$

1607 where  $I$  is the set of data points with  $\mathbf{x}_i \in \mathbb{I}_\varepsilon^{\text{all}}$  and  $O$  is the set of data points with  $\mathbf{x}_i \in \mathbb{O}_\varepsilon^{\text{any}}$ . Then the generalization gap  
1610 can be decomposed into

$$1611$$

$$1612 |R(f) - \hat{R}_{\mathcal{D}}(f)| \leq \mathbb{P}(\mathbf{x} \in \mathbb{O}_\varepsilon^{\text{any}}) \cdot \mathbb{E}_{\mathbb{O}} \left[ (f(\mathbf{x}) - y)^2 \right] + \frac{n_O}{n} \frac{1}{n_O} \sum_{i \in O} (f(\mathbf{x}_i) - y_i)^2 \quad (77)$$

$$1613 + \left| \mathbb{P}(\mathbf{x} \in \mathbb{I}_\varepsilon^{\text{all}}) - \frac{n_I}{n} \right| \frac{1}{n_I} \sum_{i \in I} (f(\mathbf{x}_i) - y_i)^2 \quad (78)$$

$$1614 + \mathbb{P}(\mathbf{x} \in \mathbb{I}_\varepsilon^{\text{all}}) \cdot \left| \mathbb{E}_{\mathbb{I}} \left[ (f(\mathbf{x}) - y)^2 \right] - \frac{1}{n_I} \sum_{i \in I} (f(\mathbf{x}_i) - y_i)^2 \right|. \quad (79)$$

1620 Using the property that the marginal distribution of  $\mathbf{x}$  is Uniform( $\mathbb{S}^{d-1}$ ) and its concentration property (Proposition D.11 +  
1621 union bound),

$$1622 \mathbb{P}(\mathbb{O}_\varepsilon^{\text{any}}) = \mathbb{P} \left( \bigcup_{S_j \in \mathcal{S}} \{\|\mathbf{x}^{(S_j)}\| > 1 - \varepsilon\} \right) \lesssim_d J \cdot \varepsilon^{\frac{d-m}{2}}. \quad (80)$$

1625 The Hoeffding inequality guarantees that, with probability at least  $1 - \delta/3$ ,

$$1626 \frac{n_O}{n} \leq J \cdot \varepsilon^{\frac{d-m}{2}} + \sqrt{\frac{\log(6/\delta)}{n}} \quad (81)$$

1629 Therefore, we may conclude that

$$1630 (77) \lesssim_d JM^2 \varepsilon^{\frac{d-m}{2}}, \quad (82)$$

1631 where  $\lesssim_d$  hides the constants that could depend on  $d$  and logarithmic factors of  $1/\delta$ .

1633 For the term (78), with probability  $1 - \delta/3$

$$1634 \begin{cases} \left| \mathbb{P}(\mathbf{x} \in \mathbb{I}_\varepsilon^{\text{all}}) - \frac{n_I}{n} \right| & \lesssim \sqrt{\frac{\varepsilon^{(d-m)/2} \log(6J/\delta)}{n}}, \\ \frac{1}{n_I} \sum_{i \in I} (f(\mathbf{x}_i) - y_i)^2 & \lesssim M^2 \end{cases} \quad (83)$$

1638 so we may also conclude that

$$1639 (78) \lesssim M^2 \sqrt{\frac{\varepsilon^{(d-m)/2} \log(6J/\delta)}{n}} \leq M^2 \sqrt{\frac{\log(6J/\delta)}{n}} \quad (84)$$

1642 For the part of the interior (79), the scalar  $\mathbb{P}(\mathbf{x} \in \mathbb{I}_\varepsilon^{\text{all}})$  is less than 1 with high-probability. Therefore, we just need to deal  
1643 with the term

$$1644 \mathbb{E}_{\mathbb{I}} \left[ (f(\mathbf{x}) - y)^2 \right] - \frac{1}{n_I} \sum_{i \in I} (f(\mathbf{x}_i) - y_i)^2. \quad (85)$$

1647 Define the projected core index set  $\mathcal{C}_\varepsilon := \{(\mathbf{u}, t) \in \mathbb{S}^{m-1} \times [-1, 1] : |t| \leq 1 - \varepsilon\}$ . For any  $S_j \in \mathcal{S}$ ,  $\mathbf{X} \sim \text{Uniform}(\mathbb{S}^{d-1})$   
1648 and  $\mathbf{X}^{(S_j)} = \pi_j(\mathbf{X})$ ,

1649

- (51) shows that  $\mathbb{P}(\mathbf{u}^\top \mathbf{X}^{(S_j)} > t) \asymp (1-t)^{\frac{d-1}{2}}$ , and
- (54) shows that  $\mathbb{E}[\phi(\mathbf{u}^\top \mathbf{X}^{(S_j)} - t)] \asymp (1-t)^{\frac{d+1}{2}}$ .

Therefore we have

$$\left( \frac{1}{J} \sum_{j=1}^J \mathbb{P}(\mathbf{u}^\top \mathbf{X}^{(S_j)} > t) \right) \left( \frac{1}{J} \sum_{j=1}^J \mathbb{E}[\phi(\mathbf{u}^\top \mathbf{X}^{(S_j)} - t)] \right) \asymp (1-t)^d, \quad t \uparrow 1. \quad (86)$$

According to (47) and the definition  $g = \min\{\tilde{g}(\mathbf{u}, t), \tilde{g}(-\mathbf{u}, -t)\}$ , we have the pointwise lower bound

$$g_{\mathcal{P}, S}(\mathbf{u}, t) \geq \min \left\{ \bar{r}_{\mathcal{P}, S}(\mathbf{u}, t) \bar{p}_{\mathcal{P}, S}(\mathbf{u}, t), \bar{r}_{\mathcal{P}, S}(-\mathbf{u}, -t) \bar{p}_{\mathcal{P}, S}(-\mathbf{u}, -t) \right\},$$

since  $\sqrt{\bar{p}^2 + \|\bar{\mathbf{a}}\|_2^2} \geq \bar{p}$ .

Under  $\mathbf{X} \sim \text{Uniform}(\mathbb{S}^{d-1})$ , by rotational invariance, for each fixed  $j$  and unit  $\mathbf{u}$  the quantities  $\mathbb{P}(\mathbf{u}^\top \mathbf{X}^{(S_j)} > t)$  and  $\mathbb{E}[\phi(\mathbf{u}^\top \mathbf{X}^{(S_j)} - t)]$  depend on  $t$  only through the scalar marginal  $Z$  in (49). Moreover both are non-increasing in  $t$ . Hence the product

$$h(t) := \mathbb{P}(Z > t) \cdot \mathbb{E}[\phi(Z - t)]$$

is non-increasing for  $t \in [0, 1]$ , and for any  $t \in [-1, 1]$  we have

$$g_{\mathcal{P}, S}(\mathbf{u}, t) \gtrsim_d h(|t|).$$

Therefore, for the core  $\mathcal{C}_\varepsilon = \{(\mathbf{u}, t) : |t| \leq 1 - \varepsilon\}$ ,

$$g_{\mathcal{P}, S, \min}(\varepsilon) := \inf_{(\mathbf{u}, t) \in \mathcal{C}_\varepsilon} g_{\mathcal{P}, S}(\mathbf{u}, t) \gtrsim_d h(1 - \varepsilon).$$

By Proposition D.10,  $h(1 - \varepsilon) \asymp \varepsilon^d$ , hence there exists  $c_g(d) > 0$  such that

$$g_{\mathcal{P}, S, \min}(\varepsilon) \geq c_g(d) \varepsilon^d. \quad (87)$$

On the simultaneous core  $\mathbb{I}_\varepsilon^{\text{all}}$ , for every local receptive field  $S_j$  and any unit  $\mathbf{u} \in \mathbb{R}^m$ ,

$$\|\mathbf{x}^{(S_j)}\| \leq 1 - \varepsilon \Rightarrow \begin{cases} t \geq 1 - \varepsilon \Rightarrow \phi(\mathbf{u}^\top \mathbf{x}^{(S_j)} - t) = 0, \\ t \leq -1 + \varepsilon \Rightarrow \phi(\mathbf{u}^\top \mathbf{x}^{(S_j)} - t) = \mathbf{u}^\top \mathbf{x}^{(S_j)} - t \text{ (affine)}. \end{cases} \quad (88)$$

Therefore all *large-offset* units ( $|t| \geq 1 - \varepsilon$ ) are affine on  $\mathbb{I}_\varepsilon^{\text{all}}$  *simultaneously across all views* and can be absorbed into a global affine term. The remaining *core-offset* units with  $|t| \leq 1 - \varepsilon$  are controlled by the  $g_S$ -weighted variation. Using  $|f_\theta|_{V_g} \leq A$  and (87),

$$\left| f_\theta \Big|_{V(\mathbb{I}_\varepsilon^{\text{all}})} \leq \frac{|f_\theta|_{V_g}}{g_{\mathcal{P}, S, \min}(\varepsilon)} \lesssim A \varepsilon^{-d}. \quad (89)$$

Therefore, we may leverage the generalization bounds for the unweighted path-norm constraint (see Lemma C.7) to deduce that with probability at least  $1 - \delta/3$ ,

$$(79) \lesssim_d (A \varepsilon^{-d})^{\frac{d}{d+3}} M^{\frac{d+6}{d+3}} n^{-\frac{1}{2}}, \quad (90)$$

Combining (82), (84) and (90), we obtain

$$\sup_{\theta \in \Theta_{\text{BEoS}}(\eta, \mathcal{D})} \text{GenGap}(f_\theta; \mathcal{D}) \lesssim_d JM^2 \varepsilon^{\frac{d-m}{2}} + (A \varepsilon^{-d})^{\frac{d}{d+3}} M^{\frac{d+6}{d+3}} n^{-\frac{1}{2}}. \quad (91)$$

According to standard empirical process theory (see Theorem D.7), for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$\sup_{\substack{\mathbf{u} \in \mathbb{S}^{m-1} \\ t \in [-1, 1]}} |g_{\mathcal{D}, S}(\mathbf{u}, t) - g_{\mathcal{P}, S}(\mathbf{u}, t)| \leq C_{\text{ep}} \sqrt{\frac{m + \log(2J/\delta)}{n}} =: \zeta_n. \quad (92)$$

Consequently, for any  $\varepsilon \in (0, 1)$ ,

$$g_{\mathcal{D}, \mathcal{S}, \min}(\varepsilon) := \inf_{\substack{\mathbf{u} \in \mathbb{S}^{m-1} \\ |t| \leq 1-\varepsilon}} g_{\mathcal{D}, \mathcal{S}}(\mathbf{u}, t) \geq g_{\mathcal{P}, \mathcal{S}, \min}(\varepsilon) - \zeta_n, \quad (93)$$

where  $g_{\mathcal{P}, \mathcal{S}, \min}(\varepsilon)$  is the corresponding population quantity.

For  $\mathcal{P}_{\mathbf{X}} = \text{Uniform}(\mathbb{S}^{d-1})$ , Proposition D.10 (see (62)) gives the explicit lower bound

$$g_{\mathcal{P}, \mathcal{S}, \min}(\varepsilon) \geq c_L(d) \varepsilon^d. \quad (94)$$

Therefore, a sufficient **validity condition** ensuring  $g_{\mathcal{D}, \mathcal{S}, \min}(\varepsilon) \geq 2\zeta_n$  is

$$\varepsilon^d \geq \frac{2C_{\text{ep}}}{c_L(d)} \sqrt{\frac{m + \log(2J/\delta)}{n}}. \quad (95)$$

The constant  $C_{\text{ep}} > 0$  is universal (it does not depend on  $d, m, n, J$ ), and can be taken explicitly from any standard VC-/pseudo-dimension uniform convergence inequality; see, e.g., Mohri et al. (2018, Chapter 3) or Haussler (1992). Moreover, the constant  $c_L(d)$  defined in Proposition D.10 satisfies

$$c_L(d) \gtrsim \frac{1}{d^2} \quad (d \geq 3), \quad (96)$$

which follows from standard bounds on Gamma-function ratios (e.g. Gautschi's inequality; see dlm, §5.6(i)). Hence the admissible range for  $\varepsilon$  is nonempty (e.g.  $\varepsilon \in [\varepsilon_{\min}, 1)$ ) as soon as  $n \gtrsim \text{poly}(d) (m + \log(2J/\delta))$ .  $\square$

**Corollary E.2.** *Under the same conditions as Theorem 4.3, assume  $d > 3$  and  $1 \leq m < \frac{d(d-3)}{d+3}$ . Let*

$$Q := 3d^2 + 3d - md - 3m \quad \text{and} \quad \varepsilon_{\min} := \left( c_d d^2 \sqrt{\frac{m + \log(J/\delta)}{n}} \right)^{1/d}.$$

Define the optimal choice

$$\varepsilon^* \asymp A^{\frac{2d}{Q}} J^{-\frac{2(d+3)}{Q}} M^{-\frac{2d}{Q}} n^{-\frac{d+3}{D}}, \quad (97)$$

and choose the feasible/truncated value

$$\varepsilon^\dagger := \max\{\varepsilon_{\min}, \varepsilon^*\}.$$

Then, with probability at least  $1 - \delta$ ,

$$\sup_{\boldsymbol{\theta} \in \Theta_{\text{BEoS}}^{\text{SCN}}(\eta, \mathcal{D})} \left| \mathbb{E}[(f_{\boldsymbol{\theta}}(\mathbf{x}) - y)^2] - \widehat{\mathcal{R}}_{\mathcal{D}}(f_{\boldsymbol{\theta}}) \right| \lesssim_d JM^2(\varepsilon^\dagger)^{\frac{d-m}{2}} + \left( A(\varepsilon^\dagger)^{-d} \right)^{\frac{d}{d+3}} M^{\frac{d+6}{d+3}} n^{-1/2}. \quad (98)$$

Moreover, in the regime where  $\varepsilon^* \geq \varepsilon_{\min}$  (e.g., for  $n$  sufficiently large, treating  $A, J, M$  as constants), plugging  $\varepsilon = \varepsilon^*$  into (7) yields the optimized rate

$$\sup_{\boldsymbol{\theta} \in \Theta_{\text{BEoS}}^{\text{SCN}}(\eta, \mathcal{D})} \left| \mathbb{E}[(f_{\boldsymbol{\theta}}(\mathbf{x}) - y)^2] - \widehat{\mathcal{R}}_{\mathcal{D}}(f_{\boldsymbol{\theta}}) \right| \lesssim_d A^{\alpha_A} J^{\alpha_J} M^{\alpha_M} n^{-\alpha_n}, \quad (99)$$

where

$$\alpha_A = \frac{d(d-m)}{Q}, \quad \alpha_J = \frac{2d^2}{Q}, \quad \alpha_M = \frac{4d^2}{Q} + \frac{(d+6)(d-m)(d+3)}{(d+3)Q}, \quad \alpha_n = \frac{(d-m)(d+3)}{2Q}.$$

In particular, if  $m$  is fixed and  $d \rightarrow \infty$ , then

$$\alpha_n \rightarrow \frac{1}{6}, \quad \alpha_A \rightarrow \frac{1}{3}, \quad \alpha_J \rightarrow \frac{2}{3}, \quad \alpha_M \rightarrow \frac{5}{3}.$$

1760 *Proof.* The optimal choice of  $\varepsilon$  minimizing the RHS of (91) is

$$1761 \varepsilon^* \asymp A^{\frac{2d}{Q}} J^{-\frac{2(2d+3)}{Q}} M^{-\frac{2d(d+3)}{(2d+3)Q}} n^{-\frac{d+3}{Q}}, \quad (100)$$

1762 where

$$1763 Q = (d-m)(d+3) + 2d^2 = 3d^2 + 3d - md - 3m.$$

1764 Plugging  $\varepsilon^*$  into (91) yields

$$1765 \sup_{\theta \in \Theta_{\text{BEOs}}(\eta, \mathcal{D})} \text{GenGap}(f_{\theta}; \mathcal{D}) \lesssim_d A^{\frac{2d}{Q}} J^{-\frac{2(d+3)}{Q}} M^{-\frac{2d}{Q}} n^{-\frac{d+3}{Q}} \quad (101)$$

$$1766 \lesssim_{A, M, J, d} O\left(n^{-\frac{(d-m)(d+3)}{2((d-m)(d+3)+2d^2)}}\right).$$

1767 In particular, when  $m$  is fixed and  $d \rightarrow \infty$ ,

$$1768 \frac{(d-m)(d+3)}{2((d-m)(d+3)+2d^2)} \rightarrow \frac{1}{6}.$$

1769 Therefore, this rate does not suffer from the curse of dimensionality when  $m$  is fixed.

1770 Finally, we verify the validity of the plug-in choice  $\varepsilon^*$ . The validity condition (95) requires  $\varepsilon^d \gtrsim_d n^{-1/2}$  (up to logarithmic factors). Since  $(\varepsilon^*)^d \asymp n^{-d(d+3)/Q}$  (treating  $A, J, M$  as constants), plugging (100) into (95) gives

$$1771 n^{-\frac{d(d+3)}{Q}} \gtrsim_d n^{-1/2}.$$

1772 This requires

$$1773 \frac{1}{2} > \frac{d(d+3)}{Q} = \frac{d(d+3)}{(d-m)(d+3)+2d^2},$$

1774 which is equivalent to

$$1775 (d-m)(d+3) + 2d^2 > 2d(d+3) \iff d^2 - 3d > m(d+3) \iff m < \frac{d(d-3)}{d+3}.$$

1776  $\square$

## 1777 F. Proof of Theorem 4.4

1778 Let  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  be a dataset with  $\mathbf{x}_i \in \mathbb{R}^d$ , and let  $\mathcal{S} = \{\pi_j : \mathbb{R}^d \rightarrow \mathbb{R}^m\}_{j=1}^J$  be a collection of coordinate projections (patch extractors). For each  $j \in [J]$ , abbreviate the extracted patch by  $\mathbf{x}^{(S_j)} := \pi_j(\mathbf{x})$ , and for each sample write  $\mathbf{x}_i^{(S_j)} := \pi_j(\mathbf{x}_i)$ .

1779 Assume labels are uniformly bounded:  $|y_i| \leq D$  for all  $i \in [n]$ . Define  $I_{\neq 0} := \{i \in [n] : y_i \neq 0\}$ .

1780 Consider width- $K$  two-layer sparsely connected ReLU models with Global Average Pooling (GAP),

$$1781 f_{\theta}(\mathbf{x}) = \sum_{k=1}^K \frac{v_k}{J} \sum_{j=1}^J \phi(\mathbf{w}_k^{\top} \mathbf{x}^{(S_j)} - b_k) + \beta, \quad \phi(t) = \max\{t, 0\}, \quad (102)$$

1782 where  $\mathbf{w}_k \in \mathbb{R}^m$ ,  $b_k \in \mathbb{R}$  are the shared filter weights and bias,  $v_k \in \mathbb{R}$  is the output weight, and  $\beta \in \mathbb{R}$ . We write  $\theta = \{(\mathbf{w}_k, b_k, v_k)\}_{k=1}^K \cup \{\beta\}$ .

1783 **Theorem F.1** (Flat interpolation with width  $\leq n$ ). *Assume that  $\|\mathbf{x}_i^{(S_j)}\|_2 \leq 1$  for all  $i \in [n]$ ,  $j \in [J]$ , and there exists a map  $\tau : I_{\neq 0} \rightarrow [J]$  that assigns  $\mathbf{p}_i := \mathbf{x}_i^{(S_{\tau(i)})}$  such that  $\|\mathbf{p}_i\|_2 = 1$  and  $\mathbf{p}_i \neq \mathbf{x}_\ell^{(S_j)}$  for all  $(\ell, j) \neq (i, \tau(i))$ . There exists a width  $K \leq n$  network of the form (1) that interpolates the dataset and whose Hessian operator norm satisfies*

$$1784 \lambda_{\max}(\nabla_{\theta}^2 \mathcal{L}) \leq 1 + \frac{D^2 + 2/J^2}{n}. \quad (103)$$

**Construction F.2.** Let  $K := |I_{\neq 0}| \leq n$  and index the hidden units by  $k \in I_{\neq 0}$ . For each  $k \in I_{\neq 0}$  define the anchor patch  $\mathbf{p}_k := \mathbf{x}_k^{(S_{\tau(k)})}$  and set

$$\rho_k := \max_{(\ell, j) \neq (k, \tau(k))} (\mathbf{x}_\ell^{(S_j)})^\top \mathbf{p}_k < 1, \quad b_k \in (\rho_k, 1), \quad \mathbf{w}_k := \mathbf{p}_k. \quad (104)$$

Set the output bias and output weights

$$\beta := 0, \quad v_k := \frac{J y_k}{1 - b_k}, \quad k \in I_{\neq 0}. \quad (105)$$

(Justification of  $\rho_k < 1$ .) According to the assumption, for any  $(\ell, j) \neq (k, \tau(k))$ ,

$$(\mathbf{x}_\ell^{(S_j)})^\top \mathbf{p}_k \leq \|\mathbf{x}_\ell^{(S_j)}\|_2 \|\mathbf{p}_k\|_2 \leq 1.$$

If equality held then necessarily  $\|\mathbf{x}_\ell^{(S_j)}\|_2 = 1$  and  $\mathbf{x}_\ell^{(S_j)} = \mathbf{p}_k$ , contradicting to the assumption. Hence  $\rho_k < 1$ .

By (104), for any sample index  $i$  and any patch index  $j$ ,

$$\mathbf{w}_k^\top \mathbf{x}_i^{(S_j)} - b_k = \begin{cases} 1 - b_k > 0, & (i, j) = (k, \tau(k)), \\ \leq \rho_k - b_k < 0, & (i, j) \neq (k, \tau(k)). \end{cases} \quad (106)$$

Thus neuron  $k$  is activated on exactly one patch in the entire collection  $\{\mathbf{x}_i^{(S_j)}\}_{i,j}$ , namely  $\mathbf{x}_k^{(S_{\tau(k)})}$ , and inactivated on all other patches. Using (106) and (105), at  $\mathbf{x}_k$  we have

$$f_\theta(\mathbf{x}_k) = \frac{v_k}{J} \phi(1 - b_k) = \frac{v_k}{J} (1 - b_k) = y_k, \quad k \in I_{\neq 0},$$

and for  $i \notin I_{\neq 0}$  all constructed units are inactivated on all patches of  $\mathbf{x}_i$ , hence  $f_\theta(\mathbf{x}_i) = \beta = 0 = y_i$ . Therefore,  $f_\theta(\mathbf{x}_i) = y_i$  for all  $i \in [n]$ .

For each constructed unit, define

$$\tilde{v}_k := \text{sign}(v_k) \in \{\pm 1\}, \quad \tilde{\mathbf{w}}_k := |v_k| \mathbf{w}_k, \quad \tilde{b}_k := |v_k| b_k. \quad (107)$$

Then for any input  $\mathbf{x}$ ,

$$\tilde{v}_k \sum_{j=1}^J \phi(\tilde{\mathbf{w}}_k^\top \mathbf{x}^{(S_j)} - \tilde{b}_k) = \frac{v_k}{J} \sum_{j=1}^J \phi(\mathbf{w}_k^\top \mathbf{x}^{(S_j)} - b_k), \quad (108)$$

so interpolation is preserved. Moreover, the activation pattern on the full patch collection is unchanged because (106) has strict inequalities and  $|v_k| > 0$ . At the unique activated patch  $(k, \tau(k))$ , the (post-rescaling) pre-activation is

$$\tilde{z}_k := \tilde{\mathbf{w}}_k^\top \mathbf{x}_k^{(S_{\tau(k)})} - \tilde{b}_k = |v_k| (1 - b_k) = J |y_k| > 0, \quad |\tilde{v}_k| = 1. \quad (109)$$

In what follows we work with the reparameterized network and drop tildes for readability, implicitly assuming  $|v_k| = 1$  for all  $k \in I_{\neq 0}$  and

$$z_k := \mathbf{w}_k^\top \mathbf{x}_k^{(S_{\tau(k)})} - b_k = J |y_k|. \quad (110)$$

**Proposition F.3.** Let  $\theta$  be the model in Construction F.2 after the reparameterization (107). Then

$$\lambda_{\max}(\nabla_{\theta}^2 \mathcal{L}) \leq 1 + \frac{D^2 + 2/J^2}{n}.$$

*Proof.* By direct computation, the Hessian  $\nabla_{\theta}^2 \mathcal{L}$  is given by

$$\nabla_{\theta}^2 \mathcal{L} = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} f_{\theta}(\mathbf{x}_i) \nabla_{\theta} f_{\theta}(\mathbf{x}_i)^\top + \frac{1}{n} \sum_{i=1}^n (f_{\theta}(\mathbf{x}_i) - y_i) \nabla_{\theta}^2 f_{\theta}(\mathbf{x}_i). \quad (111)$$

Since the model interpolates  $f_{\theta}(\mathbf{x}_i) = y_i$  for all  $i$ , we have

$$\nabla_{\theta}^2 \mathcal{L} = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} f_{\theta}(\mathbf{x}_i) \nabla_{\theta} f_{\theta}(\mathbf{x}_i)^{\top}. \quad (112)$$

Denote the tangent features matrix by

$$\Phi = [\nabla_{\theta} f_{\theta}(\mathbf{x}_1), \nabla_{\theta} f_{\theta}(\mathbf{x}_2), \dots, \nabla_{\theta} f_{\theta}(\mathbf{x}_n)]. \quad (113)$$

Then (112) can be written as  $\nabla_{\theta}^2 \mathcal{L} = \Phi \Phi^{\top} / n$ , and thus

$$\lambda_{\max}(\nabla_{\theta}^2 \mathcal{L}) = \max_{\gamma \in \mathbb{S}^{(m+2)K}} \frac{1}{n} \|\Phi^{\top} \gamma\|^2 = \max_{\mathbf{u} \in \mathbb{S}^{n-1}} \frac{1}{n} \|\Phi \mathbf{u}\|^2. \quad (114)$$

For the gate  $m_k^{(S_j)}(\mathbf{x}) := \mathbb{1}\{\mathbf{w}_k^{\top} \mathbf{x}^{(S_j)} > b_k\}$ , let  $m_{k,i}^{(S_j)} := m_k^{(S_j)}(\mathbf{x}_i)$ . From direct computation,

$$\begin{aligned} \frac{\partial f_{\theta}(\mathbf{x}_i)}{\partial v_k} &= \frac{1}{J} \sum_{j=1}^J m_{k,i}^{(S_j)} \cdot (\mathbf{w}_k^{\top} \mathbf{x}_i^{(S_j)} - b_k), & \frac{\partial f_{\theta}(\mathbf{x}_i)}{\partial \mathbf{w}_k} &= \frac{1}{J} \sum_{j=1}^J m_{k,i}^{(S_j)} \cdot v_k \cdot \mathbf{x}_i^{(S_j)}, \\ \frac{\partial f_{\theta}(\mathbf{x}_i)}{\partial b_k} &= -\frac{1}{J} \sum_{j=1}^J m_{k,i}^{(S_j)} \cdot v_k, & \frac{\partial f_{\theta}(\mathbf{x}_i)}{\partial \beta} &= 1. \end{aligned} \quad (115)$$

By the one-to-one activation property (106), each sample  $\mathbf{x}_i$  with  $i \in I_{\neq 0}$  activates exactly one unit (the unit  $k = i$ ) on exactly one patch  $S_{\tau(i)}$ , and samples with  $i \notin I_{\neq 0}$  activate none. Hence the sample-wise gradient  $\nabla_{\theta} f_{\theta}(\mathbf{x}_i)$  has support only on the parameter triplet  $(\mathbf{w}_i, b_i, v_i, \beta)$  when  $i \in I_{\neq 0}$ , and is zero on  $(\mathbf{w}_k, b_k, v_k)$  for all  $k$  when  $i \notin I_{\neq 0}$ . Writing the nonzero gradient block explicitly (recall  $|v_i| = 1$  after reparameterization and (110)),

$$\begin{aligned} \nabla_{(\mathbf{w}_i, b_i, v_i, \beta)} f_{\theta}(\mathbf{x}_i) &= \begin{pmatrix} \nabla_{(\mathbf{w}_i, b_i, v_i)} f_{\theta}(\mathbf{x}_i) \\ 1 \end{pmatrix}, \\ \nabla_{(\mathbf{w}_i, b_i, v_i)} f_{\theta}(\mathbf{x}_i) &= \begin{cases} \begin{pmatrix} \frac{v_i}{J} \mathbf{x}_i^{(S_{\tau(i)})} \\ -\frac{v_i}{J} \\ \frac{1}{J} (\mathbf{w}_i^{\top} \mathbf{x}_i^{(S_{\tau(i)})} - b_i) \end{pmatrix} = \begin{pmatrix} \frac{v_i}{J} \mathbf{x}_i^{(S_{\tau(i)})} \\ -\frac{v_i}{J} \\ |y_i| \end{pmatrix}, & (i \in I_{\neq 0}), \\ \mathbf{0}, & (i \notin I_{\neq 0}). \end{cases} \end{aligned} \quad (116)$$

Let  $\mathbf{u} = (u_1, \dots, u_n) \in \mathbb{S}^{n-1}$  and plug (116) into (114). As in the fully connected case, after a row permutation (grouping neuron parameters) we obtain

$$\begin{aligned} \lambda_{\max}(\nabla_{\theta}^2 \mathcal{L}) &= \max_{\mathbf{u} \in \mathbb{S}^{n-1}} \frac{1}{n} \|\Phi \mathbf{u}\|^2 \\ &= \frac{1}{n} \max_{\mathbf{u} \in \mathbb{S}^{n-1}} \sum_{i=1}^n u_i^2 \|\nabla_{(\mathbf{w}_i, b_i, v_i)} f_{\theta}(\mathbf{x}_i)\|_2^2 + \left( \sum_{i=1}^n u_i \right)^2 \end{aligned} \quad (117)$$

$$= \frac{1}{n} \max_{\mathbf{u} \in \mathbb{S}^{n-1}} \sum_{i \in I_{\neq 0}} u_i^2 \left( \frac{1}{J^2} \|\mathbf{x}_i^{(S_{\tau(i)})}\|_2^2 + \frac{1}{J^2} + y_i^2 \right) + \left( \sum_{i=1}^n u_i \right)^2. \quad (118)$$

According to the assumption,  $\|\mathbf{x}_i^{(S_{\tau(i)})}\|_2 = 1$  for all  $i \in I_{\neq 0}$ , and by bounded labels  $y_i^2 \leq D^2$  for all  $i$ . Thus

$$\begin{aligned} \lambda_{\max}(\nabla_{\theta}^2 \mathcal{L}) &\leq \frac{1}{n} \left( \max_{i \in [n]} \left( \frac{1}{J^2} \|\mathbf{x}_i^{(S_{\tau(i)})}\|_2^2 + \frac{1}{J^2} + y_i^2 \right) + \max_{\mathbf{u} \in \mathbb{S}^{n-1}} \left( \sum_{i=1}^n u_i \right)^2 \right) \\ &\leq \frac{1}{n} \left( \frac{2}{J^2} + D^2 + n \right) = 1 + \frac{D^2 + 2/J^2}{n}. \end{aligned}$$

If we remove the output bias term  $\beta$  from the parameters, then the last term  $(\sum_i u_i)^2$  in (117) is removed.  $\square$

## G. Extending the Framework beyond Weight Sharing and Global Average Pooling

The main text analyzes the minimal convolutional model (1) with weight sharing and global average pooling (GAP). These two design choices simplify the exposition but are not prerequisites for the stability-induced regularization mechanism. To demonstrate that *sparse connectivity alone* is the essential ingredient, this appendix analyzes a variant that retains the sparse receptive-field pattern of the main-text SCN but removes both weight sharing and global pooling.

### G.1. Sparsely connected network without weight sharing

**Setup and notation.** Fix the collection of local receptive fields  $\mathcal{S} = \{S_j\}_{j=1}^J$  (each  $S_j \subset [d]$  with  $|S_j| = m$ ) and the corresponding coordinate projections  $\pi_j : \mathbb{R}^d \rightarrow \mathbb{R}^m$ .

**Definition G.1** (SCN without weight sharing). A *sparsely connected network without weight sharing*, with receptive fields  $\mathcal{S}$  and width  $K$ , is a function of the form

$$f_{\theta}(\mathbf{x}) = \sum_{k=1}^K \sum_{j=1}^J v_{k,j} \phi(\mathbf{w}_{k,j}^{\top} \pi_j(\mathbf{x}) - b_{k,j}) + \beta, \quad (119)$$

where  $\mathbf{w}_{k,j} \in \mathbb{R}^m$ ,  $b_{k,j} \in \mathbb{R}$ ,  $v_{k,j} \in \mathbb{R}$  are *independent* parameters for each neuron index  $k \in [K]$  and location index  $j \in [J]$ , and  $\beta \in \mathbb{R}$  is a scalar bias. We write  $\theta = \{(\mathbf{w}_{k,j}, b_{k,j}, v_{k,j})\}_{k \in [K], j \in [J]} \cup \{\beta\}$  for the full parameter vector.

The key structural property shared with the main-text SCN (1) is *sparse connectivity*: neuron  $(k, j)$  receives input only from the  $m$ -dimensional patch  $\pi_j(\mathbf{x})$ . The architecture (119) differs from (1) in two respects: (i) filter weights  $\mathbf{w}_{k,j}$  and biases  $b_{k,j}$  are *not* shared across locations  $j$ , and (ii) output weights  $v_{k,j}$  may vary freely with  $j$ , so there is no global average pooling.

**Data and loss.** As in the main text, the dataset is  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  with  $\mathbf{x}_i \in B_R^d$  and  $y_i \in [-D, D]$ , and the training objective is  $\mathcal{L}(\theta) := \frac{1}{2n} \sum_{i=1}^n (f_{\theta}(\mathbf{x}_i) - y_i)^2$ .

**Location-specific weight function.** Because neurons at different locations  $j$  now operate with independent parameters on different patch distributions, the relevant weight function is defined *per location*. For each  $j \in [J]$ , let  $\mathbf{X}_{\mathcal{D}}^{(j)}$  be a random vector drawn uniformly from the training patches at position  $j$ , i.e. from  $\{\pi_j(\mathbf{x}_i)\}_{i=1}^n \subset \mathbb{R}^m$ . Define

$$g_{\mathcal{D},S}^{(j)}(\mathbf{u}, t) := \mathbb{E}[\phi(\mathbf{u}^{\top} \mathbf{X}_{\mathcal{D}}^{(j)} - t)] \sqrt{\mathbb{P}(\mathbf{u}^{\top} \mathbf{X}_{\mathcal{D}}^{(j)} > t)^2 + \left\| \mathbb{E}[\mathbf{X}_{\mathcal{D}}^{(j)} \mathbf{1}\{\mathbf{u}^{\top} \mathbf{X}_{\mathcal{D}}^{(j)} > t\}] \right\|^2}. \quad (120)$$

This is the direct analogue of the global weight function  $g_{\mathcal{D},S}$  in Definition D.1, specialized to a single location.

### G.2. From the BEoS condition to a location-wise weighted path norm

We establish the counterpart of Theorem 4.2 for the unshared architecture (119). The proof follows the same strategy—lower-bounding  $\lambda_{\max}(\mathbf{T}_{\mathcal{D}})$  via the all-ones direction—but now each neuron  $(k, j)$  contributes independently, which allows the bound to decompose over locations.

**Lemma G.2** (Hessian residual bound for the unshared SCN). *For the architecture (119) and any  $\mathbf{x}$  with  $\|\mathbf{x}\|_2 \leq R$ ,*

$$\|\nabla_{\theta}^2 f_{\theta}(\mathbf{x})\|_{\text{op}} \leq 2(R + 1).$$

*Proof.* The proof of Lemma D.3 bounds  $|\omega^{\top} \nabla_{\theta}^2 f(\mathbf{x}) \omega|$  for a unit perturbation  $\omega$  by examining each neuron independently: the only nonzero second derivatives within neuron  $(k, j)$  are the mixed partials between  $(\mathbf{w}_{k,j}, b_{k,j})$  and  $v_{k,j}$ , exactly as in (33)–(34). Since different neurons occupy orthogonal parameter blocks, their contributions to  $\omega^{\top} \nabla_{\theta}^2 f(\mathbf{x}) \omega$  sum without interaction. Each block contributes at most  $2|\gamma_{k,j}|(R\|\alpha_{k,j}\|_2 + |\delta_{k,j}|)$  (cf. (39)), where  $(\alpha_{k,j}, \delta_{k,j}, \gamma_{k,j})$  is the perturbation restricted to neuron  $(k, j)$ . Summing over  $(k, j)$  and applying Cauchy–Schwarz with the normalization  $\|\omega\|_2 = 1$  yields the same bound  $2(R + 1)$ .  $\square$

**Proposition G.3** (Tangent-feature lower bound for the unshared SCN). *Let  $\mathbf{T}_{\mathcal{D}} = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} f_{\theta}(\mathbf{x}_i) \nabla_{\theta} f_{\theta}(\mathbf{x}_i)^{\top}$ . Then*

$$\lambda_{\max}(\mathbf{T}_{\mathcal{D}}) \geq 1 + 2 \sum_{k=1}^K \sum_{j=1}^J |v_{k,j}| \|\mathbf{w}_{k,j}\| g_{\mathcal{D},S}^{(j)}\left(\frac{\mathbf{w}_{k,j}}{\|\mathbf{w}_{k,j}\|}, \frac{b_{k,j}}{\|\mathbf{w}_{k,j}\|}\right). \quad (121)$$

1980 *Proof.* Write the gate indicators  $m_{k,j,i} := \mathbb{1} \left\{ \mathbf{w}_{k,j}^\top \pi_j(\mathbf{x}_i) > b_{k,j} \right\}$ , the normalized direction  $\mathbf{u}_{k,j} := \mathbf{w}_{k,j} / \|\mathbf{w}_{k,j}\|_2$ , and  
 1981 the normalized threshold  $t_{k,j} := b_{k,j} / \|\mathbf{w}_{k,j}\|_2$ . The partial derivatives of  $f_\theta$  at sample  $\mathbf{x}_i$  are

$$1983 \frac{\partial f}{\partial \mathbf{w}_{k,j}}(\mathbf{x}_i) = v_{k,j} m_{k,j,i} \pi_j(\mathbf{x}_i), \quad \frac{\partial f}{\partial b_{k,j}}(\mathbf{x}_i) = -v_{k,j} m_{k,j,i}, \quad (122)$$

$$1984 \frac{\partial f}{\partial v_{k,j}}(\mathbf{x}_i) = \phi(\mathbf{u}_{k,j}^\top \pi_j(\mathbf{x}_i) - b_{k,j}), \quad \frac{\partial f}{\partial \beta}(\mathbf{x}_i) = 1. \quad (123)$$

1985 Since each neuron  $(k, j)$  operates on the patch  $\pi_j(\mathbf{x}_i)$  with *independent* parameters  $(\mathbf{w}_{k,j}, b_{k,j}, v_{k,j})$ , the gradient blocks  
 1986 are orthogonal across different  $(k, j)$  pairs.

1987 Denoting the tangent-feature matrix  $\Phi = [\nabla_\theta f(\mathbf{x}_1), \dots, \nabla_\theta f(\mathbf{x}_n)] \in \mathbb{R}^{p \times n}$ , we use the standard lower bound (cf. (31))

$$1988 \lambda_{\max}(\mathbf{T}_D) = \frac{1}{n} \max_{\|\mathbf{u}\|=1} \|\Phi \mathbf{u}\|^2 \geq \frac{1}{n^2} \|\Phi \mathbf{1}\|^2 = \frac{1}{n^2} \left\| \sum_{i=1}^n \nabla_\theta f(\mathbf{x}_i) \right\|^2. \quad (124)$$

1989 Expanding the squared norm and using the orthogonality of neuron blocks gives

$$1990 \frac{1}{n^2} \|\Phi \mathbf{1}\|^2 = 1 + \sum_{k=1}^K \sum_{j=1}^J \underbrace{\left[ v_{k,j}^2 \left( \left\| \frac{1}{n} \sum_i m_{k,j,i} \pi_j(\mathbf{x}_i) \right\|^2 + \left( \frac{1}{n} \sum_i m_{k,j,i} \right)^2 \right) \right]}_{=: P_{k,j}} + \underbrace{\|\mathbf{w}_{k,j}\|_2^2 \left( \frac{1}{n} \sum_i \phi(\mathbf{u}_{k,j}^\top \pi_j(\mathbf{x}_i) - t_{k,j}) \right)^2}_{=: Q_{k,j}}. \quad (125)$$

1991 Here the “1” comes from  $(\frac{1}{n} \sum_i \partial f / \partial \beta)^2 = 1$ . The term  $P_{k,j}$  collects the squared sums of the  $\mathbf{w}_{k,j}$ - and  $b_{k,j}$ -gradients  
 1992 (both proportional to  $v_{k,j}$ ), while  $Q_{k,j}$  comes from the  $v_{k,j}$ -gradient (proportional to  $\|\mathbf{w}_{k,j}\|_2$ ). Crucially, the double sum  
 1993 decomposes over  $(k, j)$  because the parameters are not shared across locations.

1994 We apply  $a^2 + b^2 \geq 2ab$  to each  $(k, j)$  independently (cf. the analogous step in the proof of Proposition D.2):

$$1995 P_{k,j} + Q_{k,j} \geq 2 |v_{k,j}| \|\mathbf{w}_{k,j}\|_2 \underbrace{\left( \frac{1}{n} \sum_i \phi(\mathbf{u}_{k,j}^\top \pi_j(\mathbf{x}_i) - t_{k,j}) \right) \sqrt{\left( \frac{1}{n} \sum_i m_{k,j,i} \right)^2 + \left\| \frac{1}{n} \sum_i m_{k,j,i} \pi_j(\mathbf{x}_i) \right\|^2}}_{=: g_{D,S}^{(j)}(\mathbf{u}_{k,j}, t_{k,j})}. \quad (126)$$

1996 To verify the identification with (120), note that the empirical averages over  $i \in [n]$  are precisely the expectation, probability,  
 1997 and conditional first moment under the uniform draw  $\mathbf{X}_D^{(j)} \sim \text{Uniform}(\{\pi_j(\mathbf{x}_i)\}_{i=1}^n)$ :

$$1998 \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{u}^\top \pi_j(\mathbf{x}_i) - t) = \mathbb{E}[\phi(\mathbf{u}^\top \mathbf{X}_D^{(j)} - t)], \quad \frac{1}{n} \sum_{i=1}^n m_{k,j,i} = \mathbb{P}(\mathbf{u}^\top \mathbf{X}_D^{(j)} > t),$$

$$1999 \frac{1}{n} \sum_{i=1}^n m_{k,j,i} \pi_j(\mathbf{x}_i) = \mathbb{E}[\mathbf{X}_D^{(j)} \mathbb{1} \{ \mathbf{u}^\top \mathbf{X}_D^{(j)} > t \}].$$

2000 Summing (126) over  $(k, j)$  and combining with (124)–(125) yields (121).  $\square$

2001 **Remark G.4** (Why the per-neuron decomposition succeeds here). In the GAP model (1), the output weight  $v_k$  is shared  
 2002 uniformly across all locations via the  $1/J$  average. This forces  $\nabla_{\mathbf{w}_k} f = \frac{v_k}{J} \sum_j m_{k,i}^{(S_j)} \mathbf{x}_i^{(S_j)}$ , in which contributions from  
 2003 different locations are combined with a *common* coefficient  $v_k/J$ . The common coefficient allows  $v_k^2$  to be factored out,  
 2004 and the resulting  $\frac{1}{n^2} \|\Phi \mathbf{1}\|^2$  factors through the *global* patch distribution  $\mathbf{X}_D^S$ , producing the global weight function  $g_{D,S}$  of  
 2005 Proposition D.2.

2006 In the unshared architecture (119), each neuron  $(k, j)$  has its own output weight  $v_{k,j}$ , but since  $\mathbf{w}_{k,j}$  is also independent  
 2007 across  $j$ , the gradient blocks are orthogonal and the squared norm decomposes cleanly into per- $(k, j)$  terms (125). This  
 2008 per-neuron decomposition makes the AM–GM step (126) straightforward.

An intermediate architecture—shared filters  $\mathbf{w}_k$  with location-dependent output weights  $v_{k,j}$ —would couple locations through  $\nabla_{\mathbf{w}_k} f = \sum_j v_{k,j} m_{k,j,i} \pi_j(\mathbf{x}_i)$ , in which contributions from different locations can cancel when  $v_{k,j}$  have varying signs. This coupling prevents both the per-location decomposition used above and the common-coefficient factorization available under GAP. Analyzing such architectures would require different proof techniques.

**Theorem G.5** (BEoS constraint for the unshared SCN). *For the architecture (119), if  $\boldsymbol{\theta} \in \Theta_{\text{BEoS}}^{\text{SCN}}(\eta, \mathcal{D})$ , i.e.,  $\lambda_{\max}(\nabla_{\boldsymbol{\theta}}^2 \mathcal{L}) \leq 2/\eta$ , then*

$$\sum_{k=1}^K \sum_{j=1}^J |v_{k,j}| \|\mathbf{w}_{k,j}\| g_{\mathcal{D},\mathcal{S}}^{(j)}\left(\frac{\mathbf{w}_{k,j}}{\|\mathbf{w}_{k,j}\|}, \frac{b_{k,j}}{\|\mathbf{w}_{k,j}\|}\right) \leq \frac{1}{\eta} - \frac{1}{2} + (R+1)\sqrt{2\mathcal{L}(\boldsymbol{\theta})}. \quad (127)$$

*Proof.* Decompose the Hessian as  $\nabla_{\boldsymbol{\theta}}^2 \mathcal{L} = \mathbf{T}_{\mathcal{D}} + \mathbf{R}_{\mathcal{D}}$  (cf. (27)) with

$$\mathbf{R}_{\mathcal{D}} = \frac{1}{n} \sum_{i=1}^n (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i) \nabla_{\boldsymbol{\theta}}^2 f_{\boldsymbol{\theta}}(\mathbf{x}_i).$$

Let  $\boldsymbol{\omega}$  be the unit eigenvector of  $\mathbf{T}_{\mathcal{D}}$  corresponding to  $\lambda_{\max}(\mathbf{T}_{\mathcal{D}})$ . Then (cf. (42))

$$\lambda_{\max}(\nabla_{\boldsymbol{\theta}}^2 \mathcal{L}) \geq \boldsymbol{\omega}^{\top} \nabla_{\boldsymbol{\theta}}^2 \mathcal{L} \boldsymbol{\omega} = \lambda_{\max}(\mathbf{T}_{\mathcal{D}}) + \underbrace{\frac{1}{n} \sum_{i=1}^n (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i) \boldsymbol{\omega}^{\top} \nabla_{\boldsymbol{\theta}}^2 f_{\boldsymbol{\theta}}(\mathbf{x}_i) \boldsymbol{\omega}}_{=: \Delta}.$$

By the Cauchy–Schwarz inequality and Lemma G.2,

$$|\Delta| \leq \sqrt{\frac{1}{n} \sum_i (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i)^2} \cdot \sqrt{\frac{1}{n} \sum_i (\boldsymbol{\omega}^{\top} \nabla_{\boldsymbol{\theta}}^2 f_{\boldsymbol{\theta}}(\mathbf{x}_i) \boldsymbol{\omega})^2} \leq 2(R+1)\sqrt{2\mathcal{L}(\boldsymbol{\theta})}.$$

Therefore,

$$\lambda_{\max}(\mathbf{T}_{\mathcal{D}}) \leq \lambda_{\max}(\nabla_{\boldsymbol{\theta}}^2 \mathcal{L}) + |\Delta| \leq \frac{2}{\eta} + 2(R+1)\sqrt{2\mathcal{L}(\boldsymbol{\theta})}.$$

Inserting the lower bound from Proposition G.3 and rearranging yields (127).  $\square$

### G.3. Generalization analysis

We now derive the generalization guarantee for the unshared architecture (119), following the same interior/exterior decomposition as in Appendix E.

**Interior/exterior decomposition.** For  $\varepsilon \in (0, 1)$ , define

$$\mathbb{I}_{\varepsilon}^{\text{all}} := \left\{ \mathbf{x} \in \mathbb{S}^{d-1} : \max_{j \in [J]} \|\pi_j(\mathbf{x})\| \leq 1 - \varepsilon \right\}, \quad \mathbb{O}_{\varepsilon}^{\text{any}} := \left\{ \mathbf{x} \in \mathbb{S}^{d-1} : \exists j \in [J], \|\pi_j(\mathbf{x})\| > 1 - \varepsilon \right\}.$$

By a union bound and Proposition D.11,  $\mathbb{P}(\mathbb{O}_{\varepsilon}^{\text{any}}) \lesssim J \varepsilon^{(d-m)/2}$ .

**Lower bound on  $g_{\mathcal{D},\mathcal{S}}^{(j)}$  over the interior.** On  $\mathbb{I}_{\varepsilon}^{\text{all}}$ , every patch satisfies  $\|\pi_j(\mathbf{x})\| \leq 1 - \varepsilon$  for all  $j$ . When the marginal of  $\mathbf{x}$  is Uniform( $\mathbb{S}^{d-1}$ ), each projected patch  $\pi_j(\mathbf{x}) \in \mathbb{R}^m$  inherits the same rotational-invariance structure analyzed in Appendix D. By the same computation as in Proposition D.10 (see (59)), the population weight function at each location satisfies

$$g_{\mathcal{P}}^{(j)}(\mathbf{u}, t) \geq c_g(d) \varepsilon^d \quad \text{for all } \mathbf{u} \in \mathbb{S}^{m-1}, |t| \leq 1 - \varepsilon,$$

with a constant  $c_g(d) > 0$  depending only on  $d$ . The uniform deviation bound of Theorem D.7, applied separately at each location  $j$  (the patches  $\{\pi_j(\mathbf{x}_i)\}_{i=1}^n$  are i.i.d. for fixed  $j$ ), gives

$$\sup_{\mathbf{u} \in \mathbb{S}^{m-1}, |t| \leq 1} |g_{\mathcal{D},\mathcal{S}}^{(j)}(\mathbf{u}, t) - g_{\mathcal{P}}^{(j)}(\mathbf{u}, t)| \leq C_{\text{ep}} \sqrt{\frac{m + \log(2J/\delta)}{n}} =: \zeta_n,$$

with probability at least  $1 - \delta$ , simultaneously for all  $j \in [J]$  (via a union bound absorbing  $J$  into the logarithm). Hence, for any  $\varepsilon$  satisfying  $\varepsilon^d \geq 2\zeta_n/c_g(d)$ , we have

$$g_{\mathcal{D},\mathcal{S}}^{(j)}(\mathbf{u}, t) \geq \frac{1}{2} c_g(d) \varepsilon^d \quad \text{for all } j \in [J], \mathbf{u} \in \mathbb{S}^{m-1}, |t| \leq 1 - \varepsilon. \quad (128)$$

**From the BEOs constraint to an unweighted path norm.** On the interior  $\mathbb{I}_\varepsilon^{\text{all}}$ , every neuron with  $|t_{k,j}| \geq 1 - \varepsilon$  is either identically zero or fully affine on all training patches (cf. (88)), and can be absorbed into the affine component. For the remaining neurons with  $|t_{k,j}| \leq 1 - \varepsilon$ , inserting the lower bound (128) into the BEOs constraint (127) and setting  $A := \frac{1}{\eta} - \frac{1}{2} + 4M$  yields

$$\sum_{\substack{k,j: \\ |t_{k,j}| \leq 1 - \varepsilon}} |v_{k,j}| \|\mathbf{w}_{k,j}\| \leq \frac{A}{\frac{1}{2} c_g(d) \varepsilon^d} \asymp A \varepsilon^{-d}. \quad (129)$$

**Controlling the generalization gap.** The function class on the interior set  $\mathbb{I}_\varepsilon^{\text{all}}$  has variation norm bounded by  $A\varepsilon^{-d}$  (equation (129)). The generalization gap decomposes as in (77)–(79):

- *Exterior contribution.* Since  $\|f_\theta\|_\infty \leq M$  and  $|y| \leq D \leq M$ , the exterior contributes at most  $O(JM^2\varepsilon^{(d-m)/2})$  (cf. (82)).
- *Interior contribution.* By Lemma C.7 (metric entropy of variation spaces), the interior generalization gap satisfies

$$\left| \mathbb{E}_I [(f_\theta - y)^2] - \frac{1}{n_I} \sum_{i \in I} (f_\theta(\mathbf{x}_i) - y_i)^2 \right| \lesssim_d (A\varepsilon^{-d})^{\frac{d}{d+3}} M^{\frac{d+6}{d+3}} n^{-1/2}$$

(cf. (90)).

- *Concentration of the interior fraction.* By Hoeffding’s inequality,  $|\mathbb{P}(\mathbb{I}_\varepsilon^{\text{all}}) - n_I/n| \lesssim \sqrt{\log(1/\delta)/n}$  with high probability, contributing a lower-order term  $O(M^2 \sqrt{\log(1/\delta)/n})$  (cf. (84)).

Combining these three terms gives the following result.

**Theorem G.6** (Generalization for the unshared SCN). *Assume  $\mathbf{x} \sim \text{Uniform}(\mathbb{S}^{d-1})$  and  $|y| \leq D$ . Let  $M \geq D$ ,  $\|f_\theta\|_\infty \leq M$ , and suppose  $\theta$  satisfies the BEOs condition  $\lambda_{\max}(\nabla_\theta^2 \mathcal{L}) \leq 2/\eta$ . Assume  $d > 3$  and  $1 \leq m < \frac{d(d-3)}{d+3}$ . Then, with probability at least  $1 - 2\delta$ ,*

$$\text{GenGap}(f_\theta; \mathcal{D}) \lesssim_d JM^2 \varepsilon^{\frac{d-m}{2}} + (A\varepsilon^{-d})^{\frac{d}{d+3}} M^{\frac{d+6}{d+3}} n^{-1/2}, \quad (130)$$

for any  $\varepsilon \in (0, 1)$  satisfying  $\varepsilon^d \gtrsim d^2 \sqrt{(m + \log(J/\delta))/n}$ , where  $A = \frac{1}{\eta} - \frac{1}{2} + 4M$ .

*Proof.* The proof follows the decomposition and estimates described above. The exterior bound uses Proposition D.11 and Hoeffding’s inequality (see (82)). The interior bound uses the path-norm estimate (129) together with the metric entropy of variation spaces (Lemma C.7; see (90)). The concentration term  $|\mathbb{P}(\mathbb{I}_\varepsilon^{\text{all}}) - n_I/n|$  is controlled as in (84) and is of smaller order.  $\square$

Optimizing over  $\varepsilon$  as in Corollary E.2 yields the rate exponent

$$-\frac{(d-m)(d+3)}{2((d-m)(d+3) + 2d^2)},$$

which tends to  $-\frac{1}{6}$  for fixed  $m$  and  $d \rightarrow \infty$ , identical to the rate obtained for the weight-shared SCN in Theorem 4.3.

**Conclusion.** The only architectural prerequisite for the stability-induced implicit regularization is the *sparse connectivity pattern*: each hidden neuron receives input from a small subset of the ambient coordinates. Weight sharing and GAP change how the patch geometry is aggregated. In the GAP model, the shared filter and common output coefficient couple all

locations and yield a single global patch-multiset weight  $g_{\mathcal{D},\mathcal{S}}$ . In the unshared model, the same mechanism appears as a sum of location-wise weights  $g_{\mathcal{D},\mathcal{S}}^{(j)}$ . Thus weight sharing and GAP are not necessary for stability-induced regularization. They determine whether the patch-geometry penalty is global and pooled or location-wise. This conclusion substantiates the design principle articulated in Section 6: the effective vector dimension seen by each local operator governs the strength of implicit regularization, regardless of whether parameters are shared across locations.

## H. The Role of Weight Sharing

The analysis in Section 4 shows that sparse connectivity strengthens implicit regularization by shifting the geometric input to gradient descent. But the SCN architecture bundles two distinct inductive biases together: sparse fan-in and weight sharing. Does sparse fan-in alone produce the effect, or is weight sharing essential?

To test this, we introduce a third architecture identical to the SCN except that the filters are no longer shared across locations. This *sparsely connected network without weight sharing* (SCN-noWS) has independent parameters per patch position:

$$f_{\theta}(\mathbf{x}) = \sum_{k=1}^K \sum_{j=1}^J v_{k,j} \phi(\mathbf{w}_{k,j}^{\top} \pi_j(\mathbf{x}) - b_{k,j}) + \beta.$$

Under weight sharing, a filter is driven by patches from all locations. To memorize a sample, the filter must isolate a patch that is distinguishable from every other patch in the entire multiset. Without weight sharing, a location-specific neuron only needs to find a patch that stands out within its own location. In typical data such as images, recurring local structures make cross-location isolation far harder than within-location isolation. Thus, weight sharing raises the threshold for memorization, pushing the network toward features that are common across the patch multiset.

The following experiment tests this hypothesis in a controlled denoising setting. **Data.** We construct  $D$ -dimensional inputs composed of  $J$  equal-sized disjoint patches. Within the patch multiset, most patches are pure Gaussian noise, but a small fraction carry signal. Concretely, fix two signal centers  $\mathbf{v}_+, \mathbf{v}_- \in \mathbb{S}^{m-1}$ . For each sample, a label  $y \in \{\mathbf{v}_+, \mathbf{v}_-\}$  is drawn uniformly, and a random location  $j^*$  is selected. The patch at  $j^*$  is set to  $y$  plus a small Gaussian perturbation; all other patches are independent Gaussian noise with matched total energy. This creates a sharp cluster structure: the vast majority of patches form a dense noise cluster around the origin, while two small signal clusters lie near  $\mathbf{v}_+$  and  $\mathbf{v}_-$ . Spatially, the signal is extremely sparse (one informative patch per sample).

**Setup.** We train three two-layer ReLU networks of identical width and learning rate: a fully connected network (FCN), an SCN-noWS, and an SCN. The task is regression: predict the label from the noisy input. The training objective is standard denoising MSE, without explicit regularization.

**Results.** Figure 8 shows training and test loss. All three architectures fit the training data to near-zero loss. However, only the SCN generalizes to held-out samples; the FCN and SCN-noWS both overfit. This is exactly what the memorization-threshold argument predicts. In the SCN-noWS, a neuron at location  $j$  needs only to find a signal patch that is distinguishable from other patches at the *same* location; with high probability, the one signal patch in a training sample appears as a local outlier. The SCN, by contrast, forces each filter to respond to patches across all locations. A filter that fires on a signal patch at one location will also fire on the many noise patches elsewhere. This drives up its activation mass and the associated stability penalty, pushing the network away from sample-specific memorization and toward filters tuned to the recurring signal clusters.

This experiment reveals a deeper property of weight sharing in denoising settings. Because the added noise is independent across coordinates, patches at different locations are perturbed by independent noise realizations. A filter restricted to a single location cannot separate signal from noise using cross-location statistics; it is forced to rely on the local noise realization. Weight sharing lifts this restriction: a shared filter sees many independent noise realizations of the same underlying patch structure, and its aggregated gradient suppresses the noise while amplifying the consistent signal directions. Thus weight sharing acts as an implicit denoising mechanism at the level of gradient dynamics.

**Remark H.1** (Weight sharing as a geometric form of translation invariance). Translation invariance is usually invoked to explain parameter efficiency and equivariance. Our analysis suggests a complementary geometric perspective. By coupling each filter to all spatial locations, weight sharing effectively redefines the “unit of isolation” for memorization. A training sample can be memorized without sharing if it has *any* distinguishable patch at *any* location. With sharing, merely local distinctiveness is not enough; the patch must be globally distinctive across the entire multiset. For natural images, where

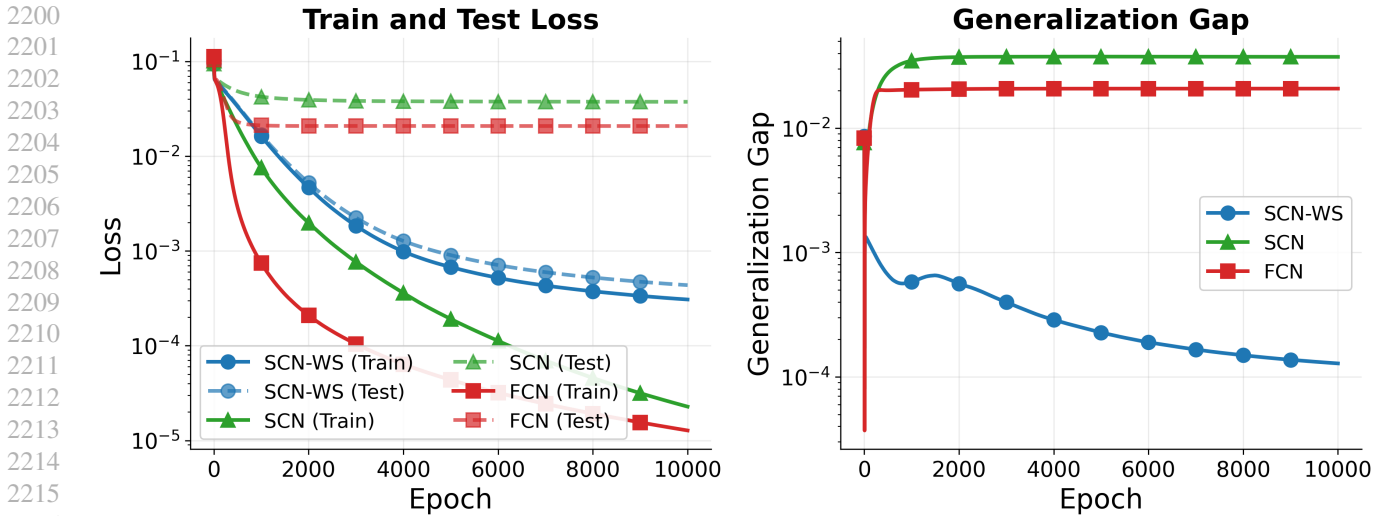


Figure 8. Training loss (solid) and test loss (dashed) for the three architectures on the clustered-patch denoising task. Only the SCN resists overfitting.

patches repeat across space, this condition is much harder to satisfy. Weight sharing thus funnels gradient descent away from sample-specific memorization and toward generalizable features, providing a form of implicit regularization rooted in the geometry of the patch multiset rather than in parameter counts.

### I. Patch Geometry of Natural Images

We now examine which side of the dichotomy established in Figure 4 real image data falls on. We extract  $3 \times 3$  convolutional patches (stride 1, no padding) from the CIFAR-10 training set, sample  $10^7$  patches. Figure 9a compares the patch cloud with the cloud of full images. The patch cloud is strongly concentrated: 90% of its variance lives in just 4 directions of  $\mathbb{R}^{27}$ , while the full images in  $\mathbb{R}^{3072}$  need over 100 directions. Recalling the gradient view (3), such a low-dimensional patch cloud implies that the gradients themselves are constrained to a low-rank subspace—much like the one-dimensional  $\mathcal{S}_1$  example in Figure 4. Moreover, Figure 9b indicates that SCN indeed benefit from the convolutional patch geometry.

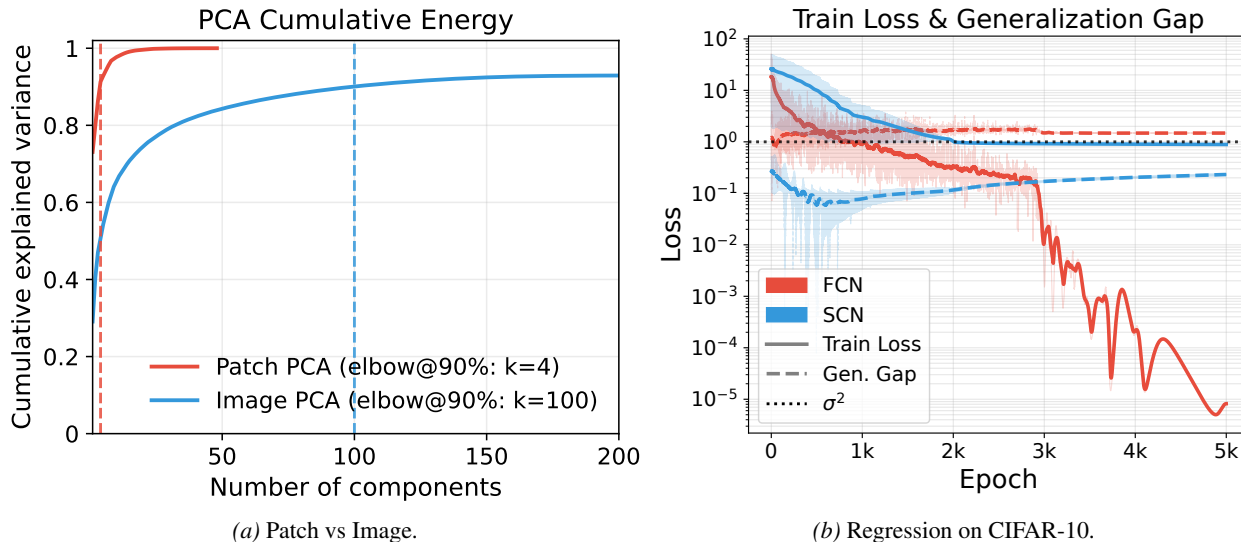


Figure 9. (a) PCA explained-variance ratio for the patch cloud ( $\mathbb{R}^{27}$ , solid) and the ambient image cloud ( $\mathbb{R}^{3072}$ , dashed). Patches concentrate 90% variance in 3 principal directions; images need over 100. (b) The SCN stabilizes near the noise floor; the FCN interpolates.

The takeaway message is that local receptive-field design can affect implicit regularization, not only feature extraction. Even under the same ambient data distribution, different receptive-field systems can induce very different patch geometries, and hence can change the balance between expressivity and regularization. In particular, the size of patch/local receptive field size in the hidden layers is a concrete perspective.

## J. Experimental Details

We adopt the random-design nonparametric regression setting.

**The setting of nonparametric regression.** The non-parametric regression with noisy labels on the random design is one typical setting in this program. Suppose  $\{\mathbf{x}_i\}_{i=1}^n$  are i.i.d. sampled from a distribution  $\mathcal{P}_{\mathbf{X}}$  supported on  $\mathbb{B}_R^d$  and  $y_i = f_{true}(\mathbf{x}_i) + \xi_i$  for  $i \in [n]$ , where  $f_{true}: \mathbb{R}^d \rightarrow \mathbb{R}$  is the ground-true function and  $\{\xi_i\}_{i=1}^n$  are i.i.d. Gaussian noises  $\mathcal{N}(0, \sigma^2)$ . In this setting, the goal of the regression task is to find a predictor  $f$  to minimize the mean squared error (MSE):

$$\text{MSE}(f) = \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - f_{true}(\mathbf{x}_i))^2. \quad (131)$$

The population level of (131) is known as *excess risk*,

$$\text{Excess}(f) := \mathbb{E}_{\mathbf{x} \sim \mathcal{P}_{\mathbf{X}}} [(f(\mathbf{x}) - f_{true}(\mathbf{x}))^2], \quad (132)$$

which is also called the *estimation error* under  $L^2(\mathcal{P}_{\mathbf{X}})$ .

In this setting, the population risk (under squared loss) of a predictor  $f$  decomposes as

$$\mathcal{R}(f) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} [(f(\mathbf{x}) - y)^2] = \text{Excess}(f) + \sigma^2, \quad (133)$$

The additive term  $\sigma^2$  is the irreducible error contributed by the label noise, and it is achieved by the Bayes predictor  $f^*(\mathbf{x}) = \mathbb{E}[y | \mathbf{x}] = f_{true}(\mathbf{x})$ , namely  $\mathcal{R}(f^*) = \sigma^2$ . Consequently,  $\mathcal{R}(f) - \mathcal{R}(f^*) = \mathbb{E}_{\mathbf{x} \sim \mathcal{P}_{\mathbf{X}}} [(f(\mathbf{x}) - f_{true}(\mathbf{x}))^2]$ , so controlling the excess risk is equivalent to controlling the population regret in squared-loss regression.

Moreover, the generalization gap in this random-design regression model admits the following equivalent form

$$\text{GenGap}(f, \mathcal{D}) = \left| \text{Excess}(f) + \sigma^2 - \widehat{\mathcal{R}}_{\mathcal{D}}(f) \right|. \quad (134)$$

This identity makes explicit how the generalization gap compares the *population* performance (excess risk plus irreducible noise) against the *training* performance measured by the empirical squared loss.

### J.1. Experimental details for the interpolation below the edge of stability

We use synthetic data to empirically validate our claim that SCN generalizes well on spherical data when  $m \ll d$ , and the rate does not deteriorate as the ambient dimension increases.

**Data generation:** Fix a collection of receptive fields  $\mathcal{S} = \{S_j\}_{j=1}^J$  with patch size  $m$ . In all experiments we use disjoint coordinate patches  $S_j = \{(j-1)m+1, \dots, jm\}$  so that  $J = \lfloor d/m \rfloor$ . We first sample a ground-truth predictor  $f_{true} \in \Theta^{\text{SCN}}$  from the architecture in (1) with a moderate width  $K_{true} = 20$  (fixed across all  $n$ ), and generate training inputs  $\{\mathbf{x}_i\}_{i=1}^n$  i.i.d. from  $\text{Uniform}(\mathbb{S}^{d-1})$ . The labels are generated by  $y_i = f_{true}(\mathbf{x}_i) + \xi_i$ , where  $\xi_i \sim \mathcal{N}(0, \sigma^2)$  are i.i.d. Gaussian noise. We also generate an independent test set  $\{\tilde{\mathbf{x}}_r\}_{r=1}^N \sim \text{Uniform}(\mathbb{S}^{d-1})$  (with  $N \gg n$ ) for Monte Carlo evaluation of generalization gap.

**Architecture:** We compare two overparameterized two-layer ReLU models trained by full-batch GD: (i) **SCN** (two-layer sparsely connected ReLU network in (1)), and (ii) **FCN** obtained by taking  $m = d$  and  $J = 1$  in (1). Unless otherwise stated, both models use width  $K = 1024$ , which places the training in an overparameterized regime for the sample sizes we consider. Both use Kaiming-normal initialization for hidden weights and zero initialization for all biases.

**Metrics:** For each trained predictor  $f$ , we report (i) **train loss**  $\widehat{\mathcal{R}}_{\mathcal{D}}(f)$  This quantity measures how well the predictor fit the training set; (ii) **estimated generalization gap**,  $\widehat{\text{GenGap}}(f, \mathcal{D})$ , which is the main object predicted by our theory (Theorem 4.3). (iii) **estimated excess risk**  $\widehat{\text{Excess}}(f) = \frac{1}{N} \sum_{r=1}^N (f(\tilde{\mathbf{x}}_r) - f_{true}(\tilde{\mathbf{x}}_r))^2$ , the Monte Carlo estimator of

Excess( $f$ ) defined in Eq (132); (iv) **Hessian Sharpness**  $\lambda_{\max}(\nabla^2 \mathcal{L}(\theta_t))$ , computed every 500 epochs via PyHessian (Yao et al., 2020) to verify the trajectory satisfies BEoS;

**Optimization** All models are trained with full-batch gradient descent (no momentum, no weight decay) on the squared loss  $\mathcal{L}(\theta)$ , using learning rate  $\eta = 0.2$  for 30000 epochs; the corresponding BEoS threshold is  $2/\eta = 10$ .

**Sweep and Scaling** For SCN, we sweep (Figure 5) ambient dimension  $d \in \{100, 200, 400\}$ ; for FCN, we fix  $d = 10$  as a baseline where fully connected networks are expected to perform well. Sample sizes range over  $n \in \{128, 256, 512, 1024\}$ . Results averaged over 5 seeds. To estimate the sample-size scaling, we sweep  $n \in \{128, 256, 512, 1024\}$  and repeat the full pipeline over multiple random seeds. We plot  $\log \widehat{\text{GenGap}}(f_{\hat{\theta}}, \mathcal{D})$  versus  $\log n$  and report the least-squares fitted slope (SCN:  $d = 100, 200, 400$ ; FCN:  $d = 10$ ).

## J.2. Experimental details for Figure 5

Similar setup and metrics as the spherical-data setup above in J.1, with two changes. First, each patch is sampled independently from the unit sphere and concatenated:  $\mathbf{x} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(J)}]$  with  $\mathbf{x}^{(j)} \sim \text{Uniform}(\mathbb{S}^{m-1})$ , instead of the whole-vector sampling  $\mathbf{x} \sim \text{Uniform}(\mathbb{S}^{d-1})$ . Second, we fix  $m = 10$ ,  $J = 8$  (so  $d = 80$ ),  $K_{\text{true}} = 50$ ,  $n = 512$ ,  $\sigma = 1$ , and train at width  $K = 1024$  for 50,000 epochs at  $\eta = 0.2$  with gradient clipping. Additional experiments with similar setup in Fig 10, 11 and Table 3.

## J.3. Experimental details for Figure 9a

We evaluate four patchification schemes on CIFAR-10: (i) *Conv Patches*,  $4 \times 4$  patches with stride 2 (baseline); (ii) *Random Patches*,  $4 \times 4$  contiguous crops at random spatial locations (same patch content as Conv but irregular grid); (iii) *ShufflePixel*, conv-style patches taken after a global pixel permutation (negative control); and (iv) *Full Image Space*, the flattened 3072-D image. For each scheme we draw 2,000,000 reference and 50,000 query patches and report two geometry metrics: (a) the PCA *effective rank*, taken as the elbow of the cumulative explained variance at 90% (lower values indicate more structured distributions); and (b) the area under the half-space (Tukey) depth concentration curve  $\Psi(T) = \Pr[\hat{d}(\mathbf{z}) \geq T]$  for  $T \in [0, 0.5]$  (higher values indicate distributions that are harder to shatter), computed via the random-projection approximation of Liang et al. (2026) with 512 uniformly random unit directions and 4096 histogram bins per direction. Full Image Space uses randomized PCA for the spectrum and 10000 query points.

## J.4. Experimental details for Figure 9b

We repeat the above setup in J.1 with two changes. First, inputs are CIFAR-10 images  $\mathbf{x} \in \mathbb{R}^{3072}$ , and SCN uses patches ( $3 \times 3 \times 3$ , stride 1;  $L = 900$ ,  $m = 27$ ). Second, targets are generated by a *frozen* ground-truth SCN  $f_{\text{true}}$  of width  $K_{\text{true}} = 20$  with the same patchification,  $y_i = f_{\text{true}}(\mathbf{x}_i) + \xi_i$ ,  $\xi_i \sim \mathcal{N}(0, 1)$ ; excess risk is reported as the MSE against  $f_{\text{true}}$  on the standard CIFAR-10 test split ( $N = 10000$ ). We train SCN and FCN, both at width  $K = 1024$ , on  $n_{\text{train}} = 1024$  images for 5000 epochs with  $\eta = 0.2$  and gradient clipping. Both networks are overparameterized in width ( $K \gg n_{\text{train}}$ ), and SCN’s training loss reaches the noise floor  $\sigma^2$ , ruling out underfitting as the explanation for its small excess risk.

Table 3. Generalization gap and excess risk for SCN on spherical data (corresponding to Figure 5a). As ambient dimension  $d$  increases with fixed patch size  $m$ , both metrics decrease—confirming the “blessing of dimensionality”.

$d$	$J$	Gen Gap	Excess Risk
100	10	$1.283 \pm 0.114$	$0.613 \pm 0.101$
200	20	$0.780 \pm 0.076$	$0.352 \pm 0.047$
400	40	$0.295 \pm 0.004$	$0.147 \pm 0.010$

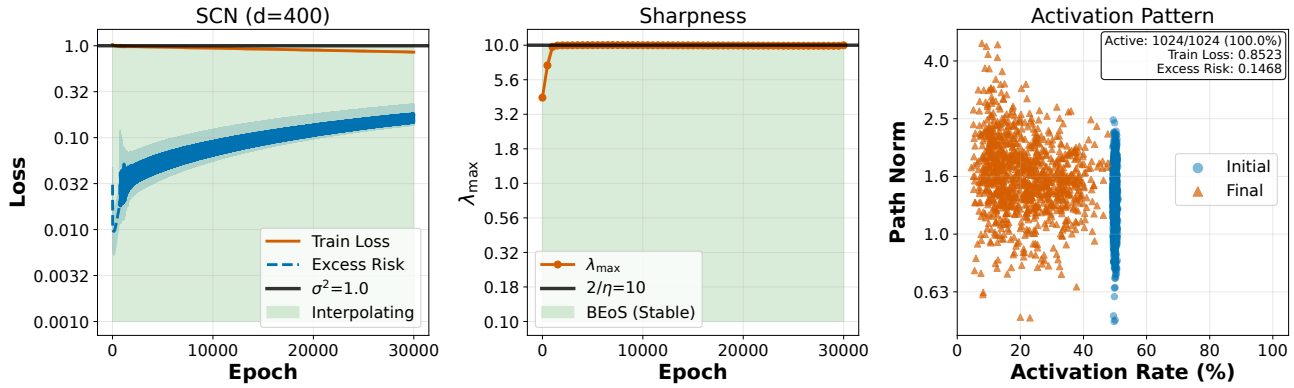


Figure 10. SCN generalizes on high-dimensional spherical data. (Left) With  $d = 400$  and fixed patch size  $m \ll d$ , train loss plateaus near  $\sigma^2$  while excess risk decreases to 0.15, confirming generalization rather than memorization. (Middle) Sharpness saturates at BEoS ( $\lambda_{\max} \approx 2/\eta$ ). (Right) Neurons spread across moderate activation rates, unlike the sparse isolation in flat interpolation (Figure 5). This validates Theorem 4.4: when  $m \ll d$ , stability-induced regularization prevents overfitting. Results averaged over 5 seeds with  $\eta = 0.2$ , trained for 30k epochs.

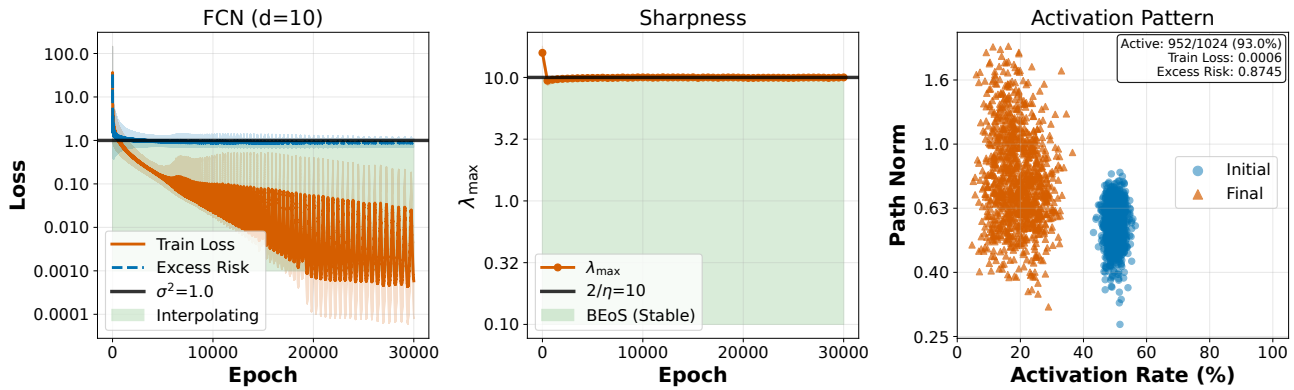


Figure 11. FCN satisfies BEoS but still memorizes (Left) FCN ( $d = 10$ ) interpolates noisy labels (train loss  $\rightarrow 0$ ) while excess risk remains  $\approx \sigma^2$ . (Middle) Sharpness saturates at BEoS. (Right) Activation pattern after training. Despite satisfying BEoS, FCN fails to generalize—confirming that on spherical data, stability constraints alone are insufficient without convolutional structure. Averaged over 5 seeds; with similar setting as Figure 10