

---

000 WHY AND WHEN DEEP IS BETTER THAN SHALLOW:  
001 AN IMPLEMENTATION-AGNOSTIC STATE-TRANSITION  
002 VIEW OF DEPTH SUPREMACY  
003  
004  
005

006 **Anonymous authors**

007 Paper under double-blind review  
008  
009

010  
011 ABSTRACT  
012

013 Why and when is deep better than shallow? We answer this question in a frame-  
014 work that is agnostic to network implementation. We formulate a deep model as an  
015 abstract state-transition semigroup acting on a general metric space, and separate  
016 the implementation (e.g., ReLU nets, transformers, and chain-of-thought) from  
017 the abstract state transition. We prove a bias-variance decomposition in which the  
018 variance depends only on the abstract depth- $k$  network and not on the implemen-  
019 tation (Theorem 1). We further split the bounds into output and hidden parts to tie  
020 the depth dependence of the variance to the metric entropy of the state-transition  
021 semigroup (Theorem 2). We then investigate implementation-free conditions under  
022 which the variance grow polynomially or logarithmically with depth (Section  
023 4). Combining these with exponential or polynomial bias decay identifies four  
024 canonical bias-variance trade-off regimes (EL/EP/PL/PP) and produces explicit  
025 optimal depths  $k^*$ . Across regimes,  $k^* > 1$  typically holds, giving a rigorous  
026 form of depth supremacy. The lowest generalization error bound is achieved un-  
027 der the EL regime (exp-decay bias + log-growth variance), explaining why and  
028 when deep is better, especially for iterative or hierarchical concept classes such as  
029 neural ODEs, diffusion/score-matching models, and chain-of-thought reasoning.

030  
031 1 INTRODUCTION  
032

033 Deep learning has achieved remarkable empirical success, yet its theoretical foundations remain in-  
034 complete. Empirically, deeper architectures excel across modalities, yet classical complexity mea-  
035 sures predict that increasing depth should worsen generalization. Prior studies using VC-dimension  
036 bounds (Bartlett et al., 2019), norm-based complexity measures (Neyshabur et al., 2015; Bartlett  
037 et al., 2017; Golowich et al., 2020), compression arguments (Arora et al., 2018; Suzuki et al.,  
038 2020; Lotfi et al., 2022), and nonparametric regression theory (Schmidt-Hieber, 2020; Suzuki, 2019;  
039 Nakada & Imaizumi, 2020; Imaizumi & Fukumizu, 2022), have derived estimation error bounds that  
040 grows exponentially, polynomially, or logarithmically with depth, but almost always within restric-  
041 tive architectural assumptions (e.g., ReLU networks under strong norm controls). This leaves open  
042 a central question: why and when is deep better than shallow, in a way that is not tied to a particular  
043 implementation. We address this gap by recasting deep learning as abstract state-transition models,  
and by analyzing depth through the geometry and algebra of the transition semigroups.

044 We start from a bias-variance decomposition in which the variance term is independent of imple-  
045 mentation: it depends only on the ideal depth- $k$  hypothesis class  $\mathcal{H}_k = H \circ B(k, F)$ , where  $B(k, F)$   
046 is the word ball of length  $k$  in the generator set  $F$  (Theorem 1). We then split the Rademacher bound  
047 into output and hidden contributions and investigate the hidden-layer effect via the entropy integral  
048 (Theorem 2). This reduction shows that the growth of covering numbers of word balls—a property  
049 of the state-transition semigroup—governs the depth dependence of variance. We identify general  
050 conditions that suppress growth to polynomial or logarithmic order, giving a structural explanation  
051 of when deeper models remain statistically benign (Section 4).

052 Why and when is deep better? We pair these variance results with bias decay that is either exponen-  
053 tial or polynomial in depth, yielding four canonical trade-off regimes (EL/EP/PL/PP). We compute  
the optimal depth  $k^*$  in each regime and discuss representative examples (Sections 5-6). Two mes-

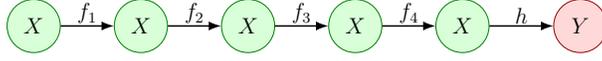


Figure 1: Example of a neural network (depth  $k = 4$ ) in consideration. The input layer is formulated as state space  $\mathcal{X}$ , the hidden layers as state transition functions  $f_i : \mathcal{X} \rightarrow \mathcal{X}$ , and the output layer as readout function  $h : \mathcal{X} \rightarrow \mathbb{R}$ . The entire network is formulated as a state transition model.

sages follow. First, depth supremacy: across regimes the optimal depth typically satisfies  $k^* > 1$ , so shallow ( $k = 1$ ) is not optimal in general. Second, the EL regime (exponential bias decay with only logarithmic variance growth) is where depth helps the most, a pattern naturally realized in hierarchical or iterative concept classes—neural ODEs, diffusion/score models, chain-of-thought reasoning—where composition is intrinsic to the data-generation or inference process.

## 2 SETTING

The most important feature of deep neural networks is function composition. Especially in modern AI, there are cases where operations include repeatedly calling an AI model itself, such as autoregression, in-context learning, test-time computation, and diffusion models. Although such operations may not have explicit network implementations, they can also be regarded as a kind of deep structure in the sense of sequential information processing. From this perspective, we formulate neural networks as *state transition models*. Precisely, to cover as wide a range of deep learning models as possible, we formulate the data domain, denoted  $\mathcal{X}$ , as an arbitrary metric space, and the deep network as a state-transition model by identifying the data domain  $\mathcal{X}$  with state space, hidden layer  $f : \mathcal{X} \rightarrow \mathcal{X}$  with transition, and output layer  $h : \mathcal{X} \rightarrow \mathbb{R}$  with observation/readout. (see Fig. 1)

### 2.1 COMMON DATA DOMAINS AND FUNCTION SPACES

We employ *one-based numbering* in this study, such as  $\mathbb{N} = \{1, 2, \dots\}$  and  $[n] := \{1, 2, \dots, n\}$ .

Let  $(\mathcal{X}, d)$  be a metric space serving as the state space, and let  $\mathcal{Y} = \mathbb{R}$  be the observation space. Let  $(\mathcal{X} \times \mathcal{Y}, P)$  be a probability space, and let  $S_n := \{(X_i, Y_i)\}_{i=1}^n \sim P$  denote an i.i.d. sample.

Let  $C(\mathcal{X})$  be the normed space of all continuous real-valued functions,  $h : \mathcal{X} \rightarrow \mathbb{R}$ , equipped with uniform norm  $\|h\|_\infty := \sup_{x \in \mathcal{X}} |h(x)|$ . Let  $C(\mathcal{X}, \mathcal{X})$  be the metric space of all continuous self-maps,  $f : \mathcal{X} \rightarrow \mathcal{X}$ , equipped with uniform metric  $d_\infty(f, g) := \sup_{x \in \mathcal{X}} d(f(x), g(x))$ . In particular, let  $\text{id} : \mathcal{X} \rightarrow \mathcal{X}$  denote the identity map. We note that both  $C(\mathcal{X})$  and  $C(\mathcal{X}, \mathcal{X})$  need not be complete in this study. We summarized (pre)compactness of function sets in Appendix I.

### 2.2 NEURAL NETWORKS (HYPOTHESIS CLASS)

We denote by  $F \subset C(\mathcal{X}, \mathcal{X})$  the set of functions corresponding to hidden layers, and by  $H \subset C(\mathcal{X})$  the set of functions corresponding to output layers. For example, in the case of a ReLU network, we may take  $\mathcal{X}$  to be the Euclidean space  $\mathbb{R}^m$  or the cube  $[0, 1]^m$ ,  $F$  to be the set of affine transformations with ReLU activations  $\{f(x) = \text{ReLU}(Wx) \mid W \in \mathbb{R}^{m \times m}\}$ , and  $H$  to be the set of linear functions  $\{h(x) = w \cdot x \mid w \in \mathbb{R}^m\}$ . Since our setting is quite general, we can also take  $\mathcal{X}$  to be a graph or a manifold, and  $F$  to be CNNs or LLMs.

We then define the depth- $k$  neural network (hypothesis class) as

$$\mathcal{H}_k := H \circ B(k, F). \quad (1)$$

Here,  $B(k, F)$  denotes the *word ball*, i.e. the set of all continuous self-maps on  $\mathcal{X}$  obtained by composing elements of  $F$  at most  $k$  times. For example,  $B(2, F) = \{\text{id}, f_1, f_2 \circ f_1 : f_1, f_2 \in F\}$ , and thus  $\mathcal{H}_2 = \{h, h \circ f_1, h \circ f_2 \circ f_1 : h \in H, f_1, f_2 \in F\}$ . By construction,  $\mathcal{H}_0 = H$ ,  $\mathcal{H}_k \subset \mathcal{H}_{k+1}$ , and  $\mathcal{H}_k \subset C(\mathcal{X}, \mathcal{X})$ . We note when the generator  $F$  contains the identity element  $\text{id}$ , then  $F^k = B(k, F)$ . So, the following identity always holds:  $(F \cup \{\text{id}\})^k = B(k, F)$ .

Besides the abstract network model  $\mathcal{H}_k$ , we introduce another class  $\mathcal{H}_{\text{imp}} \subset C(\mathcal{X}, \mathcal{X})$  called the *implementation model*. While  $\mathcal{H}_k = H \circ B(k, F)$  is not supposed to have implementation,  $\mathcal{H}_{\text{imp}}$  is supposed to have specific implementations such as ReLU networks, ConvNets, and LLMs.

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161

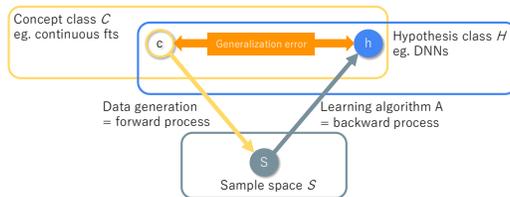


Figure 2: Framework of machine learning in consideration. The *concept class*  $\mathcal{C}$  is the class of (unseen) data generators called concept, the *hypothesis class*  $\mathcal{H}$  is the class of learning models, the *sample space*  $\mathcal{S}$  is the space of datasets. The learning (algorithm)  $A : \mathcal{S} \rightarrow \mathcal{H}$  is a mapping that assigns hypotheses to data. Generalization error is the discrepancy between the concept  $c \in \mathcal{C}$  and the outcome  $h = A(S_n) \in \mathcal{H}$ . *Complexity (of learning model)* refers to an absolute size of hypothesis class  $\mathcal{H}$ , while *expressive power* refers to a size of  $\mathcal{H}$  relative to concept class  $\mathcal{C}$ .

### 2.3 REGULARIZED EMPIRICAL RISK MINIMIZATION (LEARNING ALGORITHM)

As the learning algorithm, denoted  $A$ , we assume (*restricted*) *empirical risk minimization (RERM)* over abstract hypothesis class  $\mathcal{H}_k$ . Precisely, we consider a two-stage learning: First minimize the empirical risk, denoted  $\hat{L}_n$ , over abstract depth- $k$  networks  $\mathcal{H}_k$ , then approximate the minimizer in  $\mathcal{H}_k$  by specific networks  $\mathcal{H}_{\text{imp}}$ . In practice, trained networks obtained by optimization are not identical to exact empirical risk minimizers, but we assume they are identical for simplicity.

For the loss function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$ , we assume boundedness and 1-Lipschitz continuity at a.e. second argument (as well as measurability):  $0 \leq \ell(y', y) \leq b$ , and  $|\ell(y, y_o) - \ell(y', y_o)| \leq \beta_\ell |y - y'|$  a.e.  $y_o \in \mathcal{Y}$ . Examples include smoothed 0-1 loss in classification, or truncated and normalized squared loss in regression.

Given an i.i.d. sample  $S_n := \{(X_i, Y_i)\}_{i=1}^n$  drawn from  $P$ , put the population risk as  $L[h] := \mathbb{E}_{(X, Y) \sim P}[\ell(h(X), Y)]$ , and the empirical risk as  $\hat{L}(S_n)[h] := \frac{1}{n} \sum_{i=1}^n \ell(h(X_i), Y_i)$ , respectively.

By embed :  $\mathcal{H}_k \rightarrow \mathcal{H}_{\text{imp}}$ , we denote an embedding operator from abstract network  $\mathcal{H}_k$  to specific network  $\mathcal{H}_{\text{imp}}$ . With this notation, the two-stage learning algorithm is formulated as

$$A(S_n) := \text{embed} \left( \arg \min_{f \in \mathcal{H}_k} \hat{L}(S_n)[f] \right), \quad S_n \in (\mathcal{X} \times \mathcal{Y})^n. \quad (2)$$

### 2.4 GENERALIZATION ERROR

To formulate the generalization error, we introduce a standard machine learning framework as summarized in Fig. 2, which is composed of three classes and two processes: concept class  $\mathcal{C}$ , sample space  $\mathcal{S}$ , and hypothesis class  $\mathcal{H}$ . A concept  $c$  generates a sample (or a dataset)  $S_n \in \mathcal{S}$ , and learning algorithm  $A : \mathcal{S} \rightarrow \mathcal{H}$  mapping the sample  $S_n$  to a hypothesis  $\hat{h} \in \mathcal{H}$ . The concept  $c$  itself is assumed to be unknown, and so the learning algorithm  $A$  is an estimator of  $c$  under the constraint that it can only access  $S_n$ . Generalization error is the discrepancy between the concept  $c \in \mathcal{C}$  and outcome  $\hat{h} = A(S_n) \in \mathcal{H}$  of the learning algorithm.

In this study, the concept class  $\mathcal{C}$  is a certain collection of measurable functions  $c : \mathcal{X} \rightarrow \mathbb{R}$ , the sample space  $\mathcal{S}$  is the product space  $(\mathcal{X} \times \mathcal{Y})^n$ , and the hypothesis class is composed of two sub-classes: abstract depth- $k$  network  $\mathcal{H}_k$  and specific network  $\mathcal{H}_{\text{imp}}$ , and our two-stage learning algorithm  $A$  is a composition of ERM  $\mathcal{S} \rightarrow \mathcal{H}_k$  followed by embedding  $\mathcal{H}_k \rightarrow \mathcal{H}_{\text{imp}}$ .

The generalization error refers to three related quantities: Population risk  $L[\hat{h}]$ , Excess risk  $L[\hat{h}] - L_C$ , and Generalization gap  $L[\hat{h}] - \hat{L}[\hat{h}]$ . Here,  $L_C := \inf_{c \in \mathcal{C}} L[c]$  is the infimum risk attainable over all the data generators in  $\mathcal{C}$ . In all the quantities, the focus is on the population risk  $L[\hat{h}]$ , but it is intractable because the data distribution  $P$  is assumed to be unknown. Thus in the theoretical analysis we estimate the discrepancy from either the Bayes risk  $L_C$  or the training loss  $\hat{L}[\hat{h}]$ .

### 3 MAIN RESULTS

To state the generalization error bounds, the main theorem of this study, we additionally introduce two quantities: Rademacher complexity, and covering number. Further details are summarized in brief notes: Appendices G and H.

**Rademacher Complexity.** Let  $\mathcal{H}$  be a separable set of real-valued measurable functions on a probability space  $(\mathcal{X}, P)$  and let  $S := \{X_1, \dots, X_n\}$  be a sample drawn from distribution  $P^n$ . Let  $\sigma := \{\sigma_1, \dots, \sigma_n\}$  be independent Rademacher variables (i.e.  $\Pr\{\sigma_i = \pm 1\} = \frac{1}{2}$ ). The *empirical Rademacher complexity*  $\hat{\mathfrak{R}}_S(\mathcal{H})$  of  $\mathcal{H}$  on  $S$  and the (*population*) *Rademacher complexity*  $\mathfrak{R}_n(\mathcal{H})$  are respectively given by

$$\hat{\mathfrak{R}}_S(\mathcal{H}) = \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(X_i) \right], \quad \text{and} \quad \mathfrak{R}_n(\mathcal{H}) := \mathbb{E}_{S \sim P^n} \left[ \hat{\mathfrak{R}}_S(\mathcal{H}) \right].$$

**Covering Number.** Given a (pseudo-)metric space  $(\mathcal{M}, d)$  and positive number  $\varepsilon > 0$ , an  $\varepsilon$ -*cover* is a finite set  $C \subset \mathcal{M}$  such that every  $m \in \mathcal{M}$  lies within distance  $\varepsilon$  of some  $c \in C$ . The covering number  $N(\mathcal{M}, d, \varepsilon)$  is the size of the smallest such cover.

#### 3.1 BIAS-VARIANCE DECOMPOSITION

Let  $\varepsilon_{\text{imp}}$  be the uniform approximation error of embedding from  $\mathcal{H}_k \rightarrow \mathcal{H}_{\text{imp}}$ , i.e.  $\varepsilon_{\text{imp}} := \sup_{h \in \mathcal{H}_k} \|h - \text{embed}[h]\|_\infty$ , and let  $\varepsilon_{\text{model}}$  be the model bias of  $\mathcal{H}_k$  relative to  $\mathcal{C}$ , i.e.  $\varepsilon_{\text{model}} := \inf_{h \in \mathcal{H}_k} L[h] - \inf_{c \in \mathcal{C}} L[c]$ .

Then, we have the following bias-variance decompositions of generalization error bound.

**Theorem 1** (Bias-Variance Decomposition for ERM over Depth- $k$  Networks). *With probability at least  $1 - \delta$  over the draw of i.i.d. sample  $S \sim P^n$ , the excess risk and generalization gap are respectively decomposed as follows:*

$$L[\hat{h}] - L_C \leq \varepsilon_{\text{imp}} + \varepsilon_{\text{model}} + 4\hat{\mathfrak{R}}_S(\mathcal{H}_k) + 2b\sqrt{\frac{2 \log 1/\delta}{n}}, \quad (3)$$

$$L[\hat{h}] - \hat{L}[\hat{h}] \leq 2\varepsilon_{\text{imp}} + 2\hat{\mathfrak{R}}_S(\mathcal{H}_k) + b\sqrt{\frac{2 \log 1/\delta}{n}}. \quad (4)$$

See Appendix B for the proof. Here, *bias* refers to approximation errors  $\varepsilon_{\text{imp}}$  and  $\varepsilon_{\text{model}}$ , and *variance* refers to estimation error  $\hat{\mathfrak{R}}_S(\mathcal{H}_k)$ . The bias-variance decomposition is carefully designed so that the variance term depends only on the abstract network  $\mathcal{H}_k$  and does not depend on specific network  $\mathcal{H}_{\text{imp}}$ . Thus for variance analysis, we can concentrate on abstract network  $\mathcal{H}_k$ , independent of network architectures/implementations.

#### 3.2 HIDDEN-OUTPUT DECOMPOSITION OF RADEMACHER COMPLEXITY

We can further decompose the variance term into hidden and output layers as follows. Given a finite sample  $S := \{X_i\}_{i=1}^n$ , let  $d_S$  denote the *empirical (pseudo-)metric* of maps  $f, g \in C(X, X)$  defined by  $d_S(f, g) := \left(\frac{1}{n} \sum_{i=1}^n d(f(X_i), g(X_i))^2\right)^{1/2}$ . Similarly, let  $\|\cdot\|_S$  denote the *empirical (pseudo-)norm* of function  $h \in C(X)$  defined by  $\|h\|_S := \left(\frac{1}{n} \sum_{i=1}^n |h(X_i)|^2\right)^{1/2}$ .

**Theorem 2** (Hidden-Output Decomposition of Rademacher Complexity for Depth- $k$  Network). *Assume that  $H \subset C(\mathcal{X})$  is uniformly  $L$ -Lipschitz, and  $F \subset C(\mathcal{X}, \mathcal{X})$  is totally-bounded with identity element  $\text{id}_{\mathcal{X}}$ . Write  $B_k := B(k, F)$  for short. The Rademacher complexity of composition class  $\mathcal{H}_k := H \circ B_k$  is decomposed into the sum of the complexities of output class  $H$  and hidden class  $F$  as follows:*

$$\hat{\mathfrak{R}}_S(\mathcal{H}_k) \leq \hat{\mathfrak{R}}_S(H) + \frac{12L}{\sqrt{n}} \int_0^{\text{diam}(B_k)} \sqrt{\log N(B_k, d_S, \varepsilon)} d\varepsilon. \quad (5)$$

Here,  $\text{diam}(B_k) := \sup_{f \in B_k} d_S(f, \text{id}_{\mathcal{X}})/2$  denotes the diameter of  $B_k$  in  $d_S$ .

See Appendix C for the proof. Based on this decomposition, we can analyze the effects of the output and hidden layers on estimation error separately. Particularly, the effect of hidden layer  $F$  is much clearer than the composite form  $H \circ F$ .

*Remark 1.* While upper-bounding Rademacher complexity via covering numbers is classically known as Dudley’s entropy integral (Theorem 8), to the best of our knowledge our bound is novel for two reasons. First, the second term is the entropy integral of  $\mathcal{X}$ -valued maps  $F \subset C(\mathcal{X}, \mathcal{X})$ , for which Rademacher complexity cannot be defined. By the construction, the Rademacher complexity can be defined only for real-valued functions such as  $H \subset C(\mathcal{X})$ . Thus prior deep learning studies based on the Rademacher complexity analysis have largely avoided such a formulation as “ $\hat{\mathfrak{R}}_n(F)$ ”.

Second, the first term retains the form of a Rademacher complexity rather than an entropy integral, which is convenient because it leaves open computational routes other than covering-number bounds. The proof of this mixed bound is somewhat technical. A simpler (but looser) derivation replaces the Rademacher complexity of  $\mathcal{H}_k$  by the entropy integral over the whole class and then factors the covering number of  $\mathcal{H}_k$  into that of the output class  $H$  and of the depth- $k$  hidden maps  $B_k$ , yielding a two-entropy-integral estimate:

$$\hat{\mathfrak{R}}_S(\mathcal{H}_k) \leq \frac{12}{\sqrt{n}} \int_0^{\text{diam}(\mathcal{H}_k)} \left[ \sqrt{\log N(H, \|\cdot\|_S, \varepsilon/2)} + \sqrt{\log N(B_k, d_S, \varepsilon/2L)} \right] d\varepsilon. \quad (6)$$

In Appendix D, we have supplemented the simpler proof.

## 4 GROWTH RATES ANALYSIS

In this section we quantify how the covering number of the depth- $k$  network,  $N(B(k, F), d_\infty, \varepsilon)$ , grows with  $k$ . Through Dudley’s entropy integral (Theorem 2), this growth governs the depth dependence of the variance term in our generalization bounds (Theorem 1): exponential growth of  $N(B(k, F), d_\infty, \varepsilon)$  yields polynomial growth of the variance, whereas polynomial growth of  $N$  yields only logarithmic variance growth. We therefore seek implementation-agnostic conditions that force saturation (no growth), polynomial growth, or (super/double-)exponential growth in  $k$ .

Technically, we focus on the covering number  $N$  in  $d_\infty$ , rather than in  $d_S$ . Here  $N(\bullet, d_\infty, \varepsilon)$  upper-bounds  $N(\bullet, d_S, \varepsilon)$  because  $d_S \leq d_\infty$ . Besides, we focus on the case where generator  $F$  is infinite and compact, namely *compactly generated* semigroup  $\langle F \rangle$ , since typical neural network parameters form finite-dimensional manifolds, thus the induced semigroups are infinitely generated in general.

Below, we only summarize conditions. The proofs and examples are provided in Appendix E. Arzelà–Ascoli theorems are summarized in Appendix I.

### 4.1 OVERVIEW OF CONDITIONS

Conditions P1 and P2 explain variance at most logarithmic in  $k$ ; Conditions E1, E2 and E3 identify polynomial or worse variance growth. The relationships between each condition can be organized as shown in Table 1

Table 1: Simplified look-up table mapping from conditions on space  $\mathcal{X}$  and generator  $F$  to growth rates (saturate, polynomial, or (at least) exponential). Note that this is simplified and incomplete; for example, the trichotomy  $\text{Lip } F \lesseqgtr 1$  is easy to check but cannot completely classify the conditions.

	$\text{Lip } F < 1$ (contractive)	$\text{Lip } F = 1$ (isometric)	$\text{Lip } F > 1$ (expansive)
$\mathcal{X}$ compact	saturate (P1)	saturate (P1)	(poly, P2) or (exp, E2)
$\mathcal{X}$ non-compact	(saturate, P1’)	(poly, P2) or (exp, E1)	(super-exp, E3)

When the underlying space  $\mathcal{X}$  is *compact* and the generating system  $F$  is *non-expansive*, the Arzelà–Ascoli theorem implies that the generated semigroup is (relatively) compact; hence the covering numbers saturate at a depth-independent constant (Condition P1). Therefore, for the covering numbers to grow with depth, one needs either *non-compactness* or *non-contracting* behavior.

Indeed, when at least one of these two compact/non-expansive requirements fails—namely, in the cases (*non-compact + isometric*) and (*compact + expanding*)—one can show *exponential growth* under suitable additional assumptions (Conditions E1 and E2).

For *finitely generated groups*, the *polynomial growth* is attained when and only when the group is *virtually nilpotent* (Gromov, 1981); and *exponential growth* is attained when the group is *free*.

Similarly, for *compactly generated semigroups*, at most *polynomial growth* is attained when the semigroup embeds into a *nilpotent Lie group* such as an abelian group and the Heisenberg group (Condition P2); and at least *exponential growth* is attained when the semigroup contains a *free group* with satisfying appropriate *separation conditions* (Conditions E1 and E2). Moreover, one can even attain *super- or double-exponential growth* when the graph exhibits *uniformly-expanding* (Condition E3).

On the other hand, unlike in the case of finitely generated groups, the growth rate cannot be determined solely from the *independent complexity* (or simplicity) of the base space  $\mathcal{X}$  or of the generating system  $F$ . For example, although both Examples 3 and 11 treat the Cantor set (a complex space) and a free action (a complex action), the former saturates to a depth-independent constant (i.e., grows polynomially with degree 0), whereas the latter grows exponentially. The reason is that merely having a free-group or tree-like structure in the generators is insufficient: if the *separation condition* fails and distances between deep nodes *contract*, then by Arzelà–Ascoli (Condition P1) the dependence collapses to a constant rather than becoming exponential.

#### 4.2 SUFFICIENT CONDITIONS FOR POLYNOMIAL UPPERBOUNDS OF COVERING NUMBERS (LOGARITHMIC GROWTH OF ENTROPIES)

**Condition P1 (Equicontinuous on Compact Saturates).** If the state space  $\mathcal{X}$  is compact and the generated semigroup  $\langle F \rangle$  is precompact/equicontinuous—e.g., if maps are uniformly Lipschitz or non-expanding—then Arzelà–Ascoli implies  $\langle F \rangle$  is compact and the covering numbers of  $B(k, F)$  saturate: they do not depend on  $k$  at any fixed  $\varepsilon > 0$ .

*Examples.* (Example 1) Rotations on a circle,  $\mathcal{X} = \mathbb{S}^1$  and  $F \subset \text{Iso}(\mathbb{S}^1)$ , yield saturation. (Example 2) Finite isometry groups on finite set  $\mathcal{X}$ : growth is trivially bounded. (Example 3) A contractive iterated-function system on  $[0, 1]$  producing the Cantor set: once depth exceeds an  $\varepsilon$ -dependent threshold,  $N(B(k, F), \varepsilon)$  is controlled by the Cantor attractor and no longer grows with  $k$ .

**Condition P2 (Nilpotent Grows Polynomially).** Suppose  $\langle F \rangle$  embeds bi-Lipschitzly into a connected, simply connected nilpotent Lie group  $H$  that acts on  $\mathcal{X}$ . Then there exists a constant  $D$  (the Guivarc’h–Bass homogeneous dimension of  $H$ ) such that

$$N(B(k, F), d_\infty, \varepsilon) \lesssim \left(1 + \frac{k}{\varepsilon}\right)^D.$$

Thus the hidden-layer contribution to the variance grows at most logarithmically with  $k$ .

*Examples.* (Example 4) Translations on Euclidean space  $\mathcal{X} = \mathbb{R}^d$  (non-compact, abelian, isometric):  $O((k/\varepsilon)^D)$ ,  $D \leq d$ . (Example 5) Shear on torus  $\mathcal{X} = \mathbb{T}^2$  (compact, abelian, expanding):  $= 1 + k$  (independent of  $\varepsilon < 1$ ). (Example 6) Discrete Heisenberg group actions on torus  $\mathcal{X} = \mathbb{T}^3$  (compact, nilpotent, expanding):  $\Theta(1 + k)^D$ ,  $D = 4$  (independent of  $\varepsilon < \varepsilon_0$ ). (Example 7) Upper-Triangular group bi-Lipschitz actions on  $\mathcal{X} = \mathbb{R}^d$  (non-compact, nilpotent, expanding):  $\Theta(1 + k/\varepsilon)^D$ ,  $D = \frac{d(d-1)(d+1)}{6}$ .

#### 4.3 SUFFICIENT CONDITIONS FOR EXPONENTIAL LOWERBOUNDS OF COVERING NUMBERS (POLYNOMIAL GROWTH FOR ENTROPIES)

**Condition E1 (Free Semigroup with Point Separation Grows Exponentially)** If the semigroup generated by  $F = \{f_1, \dots, f_r\}$  is free at every length and there exist  $x_* \in \mathcal{X}$  and  $\delta > 0$  such that for all distinct words  $u, v$  of the same length  $d(f_u(x_*), f_v(x_*)) \geq \delta$ , then for every  $\varepsilon < \delta/2$  and every depth  $k$ ,

$$N(B(k, F), d_\infty, \varepsilon) \geq r^k.$$

Separation at a single probe transfers the exponential word count to the covering number.

324 *Examples.* (Example 8) (Free group, standard word metric; non-compact, isometric) Let  $\mathcal{X} = F_r$   
325 (non-compact) with standard word metric and  $F = \{L_{a_i}\}$  (isometric); at the basepoint  $e$ , distinct  
326 words of the same length are  $\geq 2$  apart, so for  $\varepsilon < 1$ ,  $N(B(k, F), d_\infty, \varepsilon) \geq r^k$ , even though  
327  $d_\infty(L_u, L_v) = \infty$  for  $u \neq v$ . Example 9 (Free group, bi-invariant metric; non-compact, isometric)  
328 Let  $\mathcal{X} = F_r$  with a conjugacy-invariant (bi-invariant) word metric and  $F = \{L_{a_i}\}$  (isometric); then  
329  $d_\infty(L_g, L_h) = d(e, g^{-1}h) < \infty$  and at the basepoint  $e$ , distinct words satisfy  $d \geq 1$ , so for  $\varepsilon < 1/2$ ,  
330  $N(B(k, F), d_\infty, \varepsilon) \geq r^k$ .

331  
332 **Condition E2 (Ping-Pong Grows Exponentially)** Suppose there are pairwise separated ‘‘cham-  
333 bers’’  $U_i$  and anchors  $a_i$  with  $f_i$  mapping  $U_i$  surjectively and expansively across chambers while  
334 collapsing  $\mathcal{X} \setminus U_i$  near  $a_i$ , and  $\min_{i \neq j} d(a_i, a_j) > 0$ ; then for some  $\delta > 0$ , every  $\varepsilon < \delta/2$  and depth  
335  $k$  satisfy

$$336 \quad N(B(k, F), d_\infty, \varepsilon) \geq r^k.$$

337 The rightmost mismatch in two words always yields a uniform output gap.

338 *Examples.* (Example 10) (interval, piecewise-linear expand/reset; compact, expansive) Let  $\mathcal{X} =$   
339  $[0, 1]$  (compact) and  $F = \{f_0, f_1\}$  be PL maps that expand (slope  $> 1$ ) on their own subintervals  
340 and reset to constants elsewhere; the ping-pong conditions give a uniform  $\delta = 1/3 - 2\eta > 0$ , hence  
341 for  $\varepsilon < \delta/2$ ,  $N(B(k, F), d_\infty, \varepsilon) \geq 2^k$ . (Example 11) (Subshift ping-pong; compact, expansive).  
342 Let  $\mathcal{X} = \mathcal{A}^{\mathbb{N}}$  (compact Cantor space with ultrametric  $d_\theta$ ) and define  $f_a$  to act as the shift on the  
343 cylinder  $\llbracket a \rrbracket$  (expansive) and as the constant  $\bar{a}$  off  $\llbracket a \rrbracket$ ; chambers and anchors are uniformly separated  
344 ( $\delta = 1$ ), so for  $\varepsilon < 1/2$ ,  $N(B(k, F), d_\infty, \varepsilon) \geq r^k$ .

345  
346 **Condition E3 (Uniform Expansion for Super/Double-Exponential).** Assume there is a *reset*  
347 map  $r$  sending all of  $\mathcal{X}$  into a compact  $k$ , a uniformly expanding map  $A$  on  $k$  with co-Lipschitz  
348 constant  $\lambda > 1$ , and a compact family of 1-Lipschitz *writers*  $G \subset C(K, K)$ . Then, for words of the  
349 form  $A \circ g_k \circ A \circ \dots \circ A \circ g_1 \circ r$ ,

$$350 \quad N(B(Ck + O(1), F), \varepsilon) \gtrsim \prod_{j=1}^k N(G, \varepsilon / \lambda^{k-j+1}),$$

353 which amplifies layerwise entropy multiplicatively. Consequences include super-exponential growth  
354 when  $G$  sits on a  $D$ -dimensional  $C^0$ -submanifold (giving  $\log N \gtrsim k^2$ ) and double-exponential  
355 growth when  $G$  is a Hölder ball on a  $d$ -dimensional manifold (giving  $\log N \gtrsim \lambda^{(d/\alpha)k}$ ).

356 *Example.* (Example 12) On  $\mathcal{X} = \mathbb{R}^d$  (with a bounded base metric), one can take a radial reset  $r$ , a  
357 near-linear expansion  $A$  on  $K = B(0, 1)$ , and writers  $G = \{\text{id} + u \mid \|u\|_\infty \leq 1, [u]_{C^\alpha} \leq 1\}$ , which  
358 realizes the double-exponential case.

## 360 5 OPTIMAL DEPTH

361  
362 As an application of our general results, we estimate the *optimal depth* by balancing the approxi-  
363 mation and estimation errors. As the bias-variance decomposition suggested, the estimation term  
364 (variance), denoted  $\text{var}(k, n)$ , is an intrinsic quantity depending solely on the hypothesis class  $\mathcal{H}_k$   
365 itself, whereas the approximation term (bias), denoted  $\text{bias}(k)$ , is an extrinsic quantity depending not  
366 only on  $\mathcal{H}_k$  but also on the concept class  $\mathcal{C}$ . Thus the same hypothesis class  $\mathcal{H}_k$  can yield different  
367 approximation error rate depending on the choice of concept class  $\mathcal{C}$ .

368 Here we analyze regimes where the approximation error decays with depth  $k$  either exponentially or  
369 polynomially:

$$370 \quad \text{bias}(k) = \exp(-\alpha k) \text{ (Exponential decay),} \quad \text{or} \quad k^{-\beta} \text{ (Polynomial decay)}$$

371 with parameters  $\alpha, \beta > 0$ , and where the covering number (governing estimation error) grows poly-  
372 nomially or exponentially, yielding estimation error bounds of root-log or root-polynomial growths:

$$373 \quad \text{var}(k, n) = \sqrt{\log(k)/n} \text{ (root-Logarithmic growth),} \quad \text{or} \quad \sqrt{k^\gamma/n} \text{ (root-Polynomial growth)}$$

374 with parameter  $\gamma > 0$ . Thus by combining the *two* bias decay rates with the *two* variance growth  
375 rates, we obtain the  $(2 \times 2 =)$  *four regimes: EL, EP, PL, and PP.*

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

Table 2: Example of the optimal depth balancing bias and variance

Regime	bias( $k$ )	var( $k, n$ )	Optimal depth $k^*$	Optimal gen( $k^*, n$ )
EL	$\exp(-\alpha k)$	$\sqrt{\log(k)/n}$	$\log(n)/(2\alpha) + \tilde{o}(1)$	$\asymp \sqrt{\log(\log(n))/n}$
EP	$\exp(-\alpha k)$	$\sqrt{k^\gamma/n}$	$\log(n)/(2\alpha) + \tilde{O}(1)$	$\asymp \sqrt{(\log n)^\gamma/(2\alpha)^\gamma n}$
PL	$k^{-\beta}$	$\sqrt{\log(k)/n}$	$\sim (2\beta n/\log(2\beta n))^{1/(2\beta)}$	$\asymp \sqrt{\log(n)/(2\beta n)}$
PP	$k^{-\beta}$	$\sqrt{k^\gamma/n}$	$\asymp n^{1/(2\beta+\gamma)}$	$\asymp n^{-\beta/(2\beta+\gamma)}$

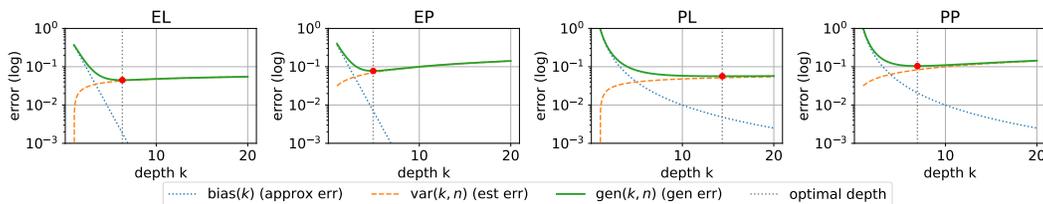


Figure 3: Typical examples of approximation error and estimation error ( $\alpha = 1.0, \beta = 2.0, \gamma = 1.0$ )

By equating leading terms, the optimal depths are estimated as in Table 2. Typical examples are also visualized in Fig. 3. The details of calculations are described in Appendix F.

In every regime, the optimal depth  $k^*$  grows with training sample size  $n$ : larger datasets require deeper networks. Precisely, the optimized generalization error increases (= gets worse) in the order:  $\text{EL} \lesssim \text{EP} \lesssim \text{PL} \lesssim \text{PP}$ , and the exp-decay bias yields a shallower  $O(\log n)$ -depth, while poly-decay bias yields a deeper  $O(\text{poly } n)$ -depth. Therefore, in general, the approximation efficiency has a stronger influence on the generalization performance than the model complexity.

## 6 EXAMPLES

Here, we discuss specific examples that fall under the four regimes.

### 6.1 CONTRACTIVE TEACHER-STUDENT SETTING (EL)

A teacher-student setting refers to the scenario where the hypothesis class (student) and the concept class (teacher) share the same compositional structure. Let the depth- $k$  student class be  $\mathcal{H}_k := H \circ B(k, F)$  and take the concept class to be the infinitely deep limit  $\mathcal{C} := \mathcal{H}_\infty := H \circ B(\infty, F)$ . Assume the input space  $\mathcal{X}$  is bounded with  $\text{diam}(\mathcal{X}) = D < \infty$ ; the output layer  $H$  is 1-Lipschitz; and the intermediate-layer semigroup  $F \subset C(\mathcal{X}, \mathcal{X})$  is *contractive*, i.e., there exists  $\lambda \in (0, 1)$  such that  $\text{Lip}(f) \leq \lambda$  for all  $f \in F$ . Then any teacher  $c \in \mathcal{H}_\infty$  can be written as  $c = h \circ u \circ v$  with  $h \in H, u \in B(k, F)$ , and  $v \in B(\ell, F)$  for some  $\ell$ . Truncating the tail yields a depth- $k$  student  $h \circ u \in \mathcal{H}_k$  that approximates  $c$  with error

$$\|h \circ u \circ v - h \circ u\|_\infty \leq \text{Lip}(h) \text{Lip}(u) d_\infty(v, \text{id}) \leq 1 \cdot \lambda^k \cdot D = D\lambda^k,$$

since  $\text{Lip}(u) \leq \lambda^k$  and  $d_\infty(v, \text{id}) \leq D$ . Hence the approximation error decays exponentially:

$$\sup_{c \in \mathcal{H}_\infty} \inf_{h \in \mathcal{H}_k} \|c - h\|_\infty \leq D\lambda^k.$$

Moreover, the contraction assumption implies the logarithmic variance growth due to Condition P1. Therefore, a teacher-student setting with contractive assumption falls under the EL regime.

### 6.2 NEURAL OPERATOR (EL)

Furuya et al. (2025) investigated the approximation error of Neural Operators (NOs) that learn the solution operator of nonlinear parabolic PDEs on a bounded domain  $\mathcal{X} \subset \mathbb{R}^d$ . By aligning a single hidden layer of the NO with one step of the Picard iteration for the PDE, a textbook iterative arguments for DE solving, they show that the approximation error decays at a *sub-exponential* rate

432  $\exp O(-\sqrt{k})$  in the network depth  $k$ . Since Picard iteration is *contractive* and, moreover, generated  
 433 by a *single* operator, the resulting NO forms a *commutative contractive semigroup*, implying that the  
 434 covering numbers grow only polynomially with depth (Conditions P1 and P2). The function space  
 435 for approximation is a *mixed Lebesgue space*  $L_t^r L_x^s$ , which we can regard as the concept class  $\mathcal{C}$ .

436 This NO setting can be viewed as a concrete instance of the teacher-student framework discussed  
 437 earlier: both the hypothesis class and the concept class share a hierarchical (iterative) structure,  
 438 enabling efficient depth-driven approximation. It thus exemplifies the *PL regime*, in which deep  
 439 learning is particularly advantageous. Together with the previous example, it is suggested that deep  
 440 learning is most effective when the underlying data space, namely the concept class  $\mathcal{C}$ , possesses a  
 441 hierarchical structure such as that induced by deep compositions or differential equations.

### 443 6.3 PARTIALLY-FREE TEACHER-STUDENT SETTING (EP)

444 As mentioned in the EL example (and Condition P1), networks generated by a contractive semigroup  
 445 on a compact space are precompact, so their covering numbers do not grow with depth. While slow  
 446 growth of estimation error is preferable in practice, we present one artificial construction where  
 447 estimation complexity grows exponentially in depth, even though approximation error still decays  
 448 exponentially. The basic idea is again in a teacher-student setting to make the hypothesis class  $\mathcal{H}_k$   
 449 deliberately *larger* (redundant) than the concept class  $\mathcal{C}$ .

450 Let the input space  $\mathcal{X}$  be bounded with  $\text{diam}(\mathcal{X}) = D < \infty$ . Assume the output layer  $H$  is  
 451 1-Lipschitz. For hidden layer maps, take the union  $F := S \cup T$  of a *contractive* semigroup  
 452  $S \subset C(\mathcal{X}, \mathcal{X})$  with  $\sup_{s \in S} \text{Lip}(s) =: \lambda < 1$  and a *free* part  $T$  with finite generators whose  
 453 generated semigroup  $\langle T \rangle := B(\infty, T)$  is  $\rho$ -separated in the uniform metric: there exists  $\rho > 0$   
 454 such that  $d_\infty(f, g) \geq \rho$  for all distinct  $f, g \in \langle T \rangle$ . This implies a covering-number lower bound  
 455  $N(B(k, T), \varepsilon) \gtrsim |T|^k$  for any  $\varepsilon < \rho/2$ .

456 Define the concept class by the contractive part only,  $\mathcal{C} := H \circ B(\infty, S)$ , and the depth- $k$  hypothesis  
 457 class by the full union,  $\mathcal{H}_k := H \circ B(k, S \cup T)$ . Then we are in the *EP regime*: the approximation  
 458 error decays exponentially by the same truncation argument as before:

$$459 \sup_{c \in H \circ B(\infty, S)} \inf_{h \in H \circ B(k, S \cup T)} \|c - h\|_\infty \leq D\lambda^k.$$

460 On the other hand, the covering numbers of the hidden maps grow at least exponentially with depth,  
 461 because (Condition E1 holds, or)

$$462 N(B(k, S \cup T), \varepsilon) \geq N(B(k, T), \varepsilon) \gtrsim |T|^k \quad (\varepsilon < \rho/2),$$

463 so the estimation term exhibits exponential dependence on  $k$ . Such an example can be made fully  
 464 concrete by taking  $\mathcal{X}$  to be the Hilbert cube  $[0, 1]^{\mathbb{N}}$  equipped with  $\ell^\infty$ -norm and constructing explicit  
 465 choices of  $S$  (contractive) and  $T$  (finite,  $\rho$ -separated free generators) on this space.

### 470 6.4 RELU NETWORKS IN HÖLDER-SMOOTH SPACE (PP/PL)

471 A canonical setting where approximation error decays only at a *polynomial* rate is given by *Jackson-*  
 472 *type* bounds for *Hölder*  $C^s([0, 1]^d)$ , *Sobolev*  $W^{s,p}([0, 1]^d)$ , and *Besov*  $B_q^{s,p}([0, 1]^d)$  spaces: For  
 473 such a function  $f$  in these spaces, the best  $m$ -parameter approximation achieves order  $O(m^{-s/d})$ ,  
 474 with matching lower bounds  $\Omega(m^{-s/d})$  under very mild assumptions ( DeVore et al., 1989). Thus,  
 475 exponentially fast approximation cannot be expected in these spaces.

476 A line of expressive power analysis of ReLU networks initiated by Yarotsky (Yarotsky, 2017; 2018;  
 477 Yarotsky & Zhevnerchuk, 2020; Siegel, 2023; Yang & He, 2024) shows that deep ReLU networks  
 478 can attain the so-called *super-convergence*, or surpass the Jackson's rates, by combining piecewise-  
 479 polynomial approximation with *bit-extraction*, a highly compressed, discontinuous encoding tech-  
 480 nique (from function to parameter); the speedup hinges on violating the regularity assumptions  
 481 underlying the Jackson-type lower bounds, yet the decay remains polynomial rather than exponen-  
 482 tial. Although a rigorous estimation error analysis was not provided by the authors, bit-extraction  
 483 is a Cantor attractor, which may yield both a polynomial growth (Example 3) and an exponential  
 484 growth (Example 11) of covering numbers. Thus, we expect the settings where the hypothesis class  
 485 is ReLU network and concept class is Hölder/Sobolev/Besov spaces fall under *PP/PL regime*.

---

486 Apart from the Jackson’s regime, estimation error for ReLU networks has been shown to grow  
487 *polynomially* in depth  $k$  via VC-dimension arguments (Bartlett et al., 2019) and compression-based  
488 generalization bounds (Arora et al., 2018; Suzuki et al., 2020; Lotfi et al., 2022). Putting these  
489 observations together: taking the concept class  $\mathcal{C}$  as Hölder/Sobolev/Besov, and the hypothesis class  
490  $\mathcal{H}_k$  as depth- $k$  ReLU networks yields the *PP regime*.

491 We remark that Suzuki (2019) investigated both approximation and estimation error rates for ReLU  
492 networks in both Besov and mixed-smooth Besov spaces, and obtained exactly the *PP regime*: poly-  
493 nomial bias decay with polynomial variance growth.

## 495 6.5 RELU NETWORKS IN HIERARCHICAL CLASS (PL)

497 Schmidt-Hieber (2020) developed a hierarchical class, named *composite function class*  $G$ , ob-  
498 tained by compositions of Hölder-smooth maps and showed that deep ReLU networks achieve the  
499 minimax-optimal rate. Their argument bounds the covering numbers of deep ReLU classes, yield-  
500 ing estimation terms that increase only *logarithmically* in depth  $k$ . On the approximation side, they  
501 obtain a bound with two terms: an *exponentially* decaying term in depth  $k$  (from compositional  
502 structure) plus a *polynomial* Jackson’s rate term in the number of parameters  $m$ . While this is not  
503 purely a polynomial decay in depth, the overall picture fits within a *PL regime*.

## 505 7 CONCLUSION

507 We developed a unified framework for analyzing the depth dependence of generalization error en-  
508 compassing a variety of network architectures by splitting deep learning models into abstract state  
509 transitions composed with concrete implementations. We analyzed generalization through a bias-  
510 variance decomposition in which the variance term depends only on the state-transition semigroup.  
511 We characterized when the covering numbers of semigroups grow polynomially or exponentially  
512 with depth, yielding the root-logarithmic or root-polynomial variance growth. We combined these  
513 results with exponential or polynomial bias decay to obtain four bias-variance trade-off regimes  
514 (EL, EP, PL, PP) and their optimal depths. The optimal depth  $k^*$  is typically greater than one, which  
515 constitutes a rigorous form of *depth supremacy*. Among the four regimes, depth performs the best  
516 in EL, where bias decays Exponentially while variance grows only Logarithmically. Hierarchical or  
517 iterative concept classes  $\mathcal{C}$  such as Neural ODEs, diffusion/score models, and chain-of-thought rea-  
518 soning models naturally realize the EL regime. On the other hand, classical Hölder/Sobolev/Besov  
519 spaces tend to yield PL/PP regimes and make depth supremacy harder to establish. Our analysis  
520 highlights *compactly-generated semigroups* as a unifying mathematical object and connects depth-  
521 generalization to ideas from coarse geometry and dynamical systems. We hope this bridge will  
522 motivate further interaction between modern mathematics and the theory of deep learning.

## 523 REFERENCES

- 525 Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger Generalization Bounds for  
526 Deep Nets via a Compression Approach. In *Proceedings of the 35th International Conference on*  
527 *Machine Learning*, volume 80, pp. 254–263, Stockholm, 2018. PMLR.
- 528 Peter Bartlett, Dylan J. Foster, and Matus Telgarsky. Spectrally-normalized margin bounds for  
529 neural networks. In *Advances in Neural Information Processing Systems 31*, pp. 6240–6249,  
530 Long Beach, 2017.
- 532 Peter L Bartlett and Shahar Mendelson. Rademacher and Gaussian Complexities: Risk Bounds and  
533 Structural Results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- 535 Peter L Bartlett, Nick Harvey, Christopher Liaw, and Abbas Mehrabian. Nearly-tight VC-dimension  
536 and Pseudodimension Bounds for Piecewise Linear Neural Networks. *Journal of Machine Learn-*  
537 *ing Research*, 20(63):1–17, 2019.
- 538 Emmanuel Breuillard. Geometry of locally compact groups of polynomial growth and shape of  
539 large balls. *Group, Geometry, and Dynamics*, 8(3):669–732, 2014.

- 
- 540 Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural Ordinary Dif-  
541 ferential Equations. In *Advances in Neural Information Processing Systems*, volume 31, pp.  
542 6572–6583, Montréal, Canada, 2018.
- 543 Ronald A DeVore, Ralph Howard, and Charles Micchelli. Optimal nonlinear approximation.  
544 *manuscripta mathematica*, 63(4):469–478, 1989.
- 545 Ronen Eldan and Ohad Shamir. The Power of Depth for Feedforward Neural Networks. In *29th*  
546 *Annual Conference on Learning Theory*, volume 49, pp. 1–34, 2016.
- 547 Takashi Furuya, Koichi Taniguchi, and Satoshi Okuda. Quantitative Approximation for Neural Op-  
548 erators in Nonlinear Parabolic Equations. In *The Thirteenth International Conference on Learning*  
549 *Representations*, 2025.
- 550 Noah Golowich, Alexander Rakhlin, and Ohad Shamir. Size-Independent Sample Complexity of  
551 Neural Networks. In *Proceedings of the 31st Conference On Learning Theory*, volume 75, pp.  
552 297–299. PMLR, 2018.
- 553 Noah Golowich, Alexander Rakhlin, and Ohad Shamir. Size-independent sample complexity of  
554 neural networks. *Information and Inference: A Journal of the IMA*, 9(2):473–504, 2020.
- 555 Mikhael Gromov. Groups of polynomial growth and expanding maps. *Publications Mathématiques*  
556 *de l’Institut des Hautes Études Scientifiques*, 53(1):53–78, 1981.
- 557 Masaaki Imaizumi and Kenji Fukumizu. Deep Neural Networks Learn Non-Smooth Functions Ef-  
558 fectively. In *Proceedings of the Twenty-Second International Conference on Artificial Intelligence*  
559 *and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pp. 869–878. PMLR,  
560 2019.
- 561 Masaaki Imaizumi and Kenji Fukumizu. Advantage of Deep Neural Networks for Estimating Func-  
562 tions with Singularity on Hypersurfaces. *Journal of Machine Learning Research*, 23(111):1–54,  
563 2022.
- 564 Nikola Kovachki, Zongyi Li, Burigede Liu, Kamyar Aizzadenesheli, Kaushik Bhattacharya, An-  
565 drew Stuart, and Anima Anandkumar. Neural Operator: Learning Maps Between Function  
566 Spaces. *arXiv preprint: 2108.08481*, 2021.
- 567 Viktor Losert. On the structure of groups with polynomial growth. *Mathematische Zeitschrift*, 195  
568 (1):109–117, 1987. ISSN 1432-1823. doi: 10.1007/BF01161604.
- 569 Sanae Lotfi, Marc Finzi, Sanyam Kapoor, Andres Potapczynski, Micah Goldblum, and Andrew G  
570 Wilson. PAC-Bayes Compression Bounds So Tight That They Can Explain Generalization. In  
571 *Advances in Neural Information Processing Systems*, volume 35, pp. 31459–31473, 2022.
- 572 Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Uni-  
573 versity Press, 2019.
- 574 Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*.  
575 Adaptive Computation and Machine Learning series. MIT Press, second edition, 2018.
- 576 Ryumei Nakada and Masaaki Imaizumi. Adaptive Approximation and Generalization of Deep Neu-  
577 ral Network with Intrinsic Dimensionality. *Journal of Machine Learning Research*, 21(174):1–38,  
578 2020.
- 579 Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. Norm-Based Capacity Control in Neural  
580 Networks. In *Proceedings of The 28th Conference on Learning Theory*, volume 40, pp. 1–26,  
581 Paris, France, 2015. JMLR W&CP.
- 582 Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with ReLU acti-  
583 vation function. *The Annals of Statistics*, 48(4):1875–1897, 2020.
- 584 Hui Shen, Jingxuan Zhang, Boning Xiong, Rui Hu, Shoufa Chen, Zhongwei Wan, Xin Wang,  
585 Yu Zhang, Zixuan Gong, Guangyin Bao, Chaofan Tao, Yongfeng Huang, Ye Yuan, and Mi Zhang.  
586 Efficient Diffusion Models: A Survey. *Transactions on Machine Learning Research*, 2025. ISSN  
587 2835-8856.

---

594 Jonathan W Siegel. Optimal Approximation Rates for Deep ReLU Neural Networks on Sobolev and  
595 Besov Spaces. *Journal of Machine Learning Research*, 24(357):1–52, 2023.  
596

597 Jonathan W Siegel and Jinchao Xu. Sharp Bounds on the Approximation Rates, Metric Entropy,  
598 and n-Widths of Shallow Neural Networks. *Foundations of Computational Mathematics*, 2022.

599 Taiji Suzuki. Adaptivity of deep ReLU network for learning in Besov and mixed smooth Besov  
600 spaces: optimal rate and curse of dimensionality. In *International Conference on Learning Rep-*  
601 *resentations*, 2019.  
602

603 Taiji Suzuki, Hiroshi Abe, and Tomoaki Nishimura. Compression based bound for non-compressed  
604 network: unified generalization error analysis of large compressible deep neural network. In  
605 *International Conference on Learning Representations*, 2020.

606 Matus Telgarsky. Benefits of depth in neural networks. In *29th Annual Conference on Learning*  
607 *Theory*, pp. 1–23, 2016.  
608

609 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc V  
610 Le, and Denny Zhou. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models.  
611 In *Advances in Neural Information Processing Systems*, volume 35, pp. 24824–24837, 2022.

612 Yahong Yang and Juncai He. Deeper or Wider: A Perspective from Optimal Generalization Error  
613 with Sobolev Loss. In *Proceedings of the 41st International Conference on Machine Learning*,  
614 volume 235 of *Proceedings of Machine Learning Research*, pp. 56109–56138, 2024.

615 Dmitry Yarotsky. Error bounds for approximations with deep ReLU networks. *Neural Networks*,  
616 94:103–114, 2017.  
617

618 Dmitry Yarotsky. Optimal approximation of continuous functions by very deep ReLU networks. In  
619 *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine*  
620 *Learning Research*, pp. 639–649. PMLR, 2018.

621 Dmitry Yarotsky and Anton Zhevnerchuk. The phase diagram of approximation rates for deep neural  
622 networks. In *Advances in Neural Information Processing Systems*, volume 33, pp. 13005–13015,  
623 2020.  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

---

## A LITERATURE OVERVIEW

**Depth separation and expressivity.** A classical line of work shows that modest increases in depth can yield exponential representational advantages. Eldan & Shamir (2016) proved a three-vs-two-layer separation for a simple radial function, requiring exponential width for any depth-2 approximant, while Telgarsky (2016) established families exhibiting exponential gaps between networks of depth  $O(k^3)$  and  $O(k)$  with semi-algebraic gates (including ReLU) and provided constructive hard instances; subsequent work extended separations beyond radial constructions. These results clarify when *expressivity* favors depth, but do not by themselves pin down estimation behavior.

**Generalization via capacity control.** Combinatorial analyses give nearly tight VC/pseudodimension bounds for piecewise-linear networks, scaling roughly linearly in depth for fixed width/weights, providing a baseline picture of depth in classical uniform-convergence frameworks (Bartlett et al., 2019). Norm- and margin-based approaches bound estimation error through products of layer norms (path/spectral) (Neyshabur et al., 2015; Bartlett et al., 2017), sometimes yielding size-independent or gently depth-dependent bounds under additional structure (e.g., margin normalization). These strands highlight multiple possible depth dependencies—linear, polynomial, or even milder—depending on how complexity is measured.

**Rademacher/covering and size-independent bounds.** A complementary thread controls depth via data-dependent complexities (Rademacher, covering). Golowich et al. (2018; 2020) obtained bounds that (under norm constraints) improve the depth dependence and can be independent of width and depth in certain regimes; later refinements further reduced explicit depth factors. These works show how the estimation side may be decoupled from naively counted parameters and instead tied to geometric quantities of the hypothesis class.

**Compression and PAC-Bayes.** A productive viewpoint explains generalization via compressibility: if a trained network admits a succinct reparametrization, one can transfer that compression into generalization guarantees. Arora et al. (2018) formalized this link and demonstrated strong bounds in practice; follow-ups convert compression bounds to the original (non-compressed) networks and sharpen them via PAC-Bayes with subspace quantization, yielding state-of-the-art nonvacuous estimates (Suzuki et al., 2020; Lotfi et al., 2022). While depth typically enters these bounds through compressibility or margin quantities, the methodology is agnostic to architecture details.

**Nonparametric regression with deep networks.** Another large body of work analyzes approximation–estimation trade-offs of ReLU networks on smoothness classes. Schmidt-Hieber (2020) showed near-optimal rates in nonparametric regression, with depth playing an essential role; Suzuki (2019) established optimal adaptivity on (mixed) Besov spaces and improvements over linear/kernel baselines; Nakada & Imaizumi (2020) tied generalization to intrinsic (Minkowski) dimension; and Imaizumi & Fukumizu (2019; 2022) identified regimes with singularities where DNNs are minimax-superior to traditional estimators. Recent approximation results (e.g., optimal Sobolev/Besov rates) further sharpen the expressivity side (Yarotsky, 2017; 2018; Yarotsky & Zhevnerchuk, 2020; Siegel & Xu, 2022; Siegel, 2023; Yang & He, 2024). These analyses, however, are typically architecture-specific (ReLU feedforward) and hinge on smoothness assumptions.

**Iterative/hierarchical models and continuous depth.** Many modern systems are naturally modeled as *compositions* or *flows*—precisely the setting of our state-transition abstraction. Neural ODEs (Chen et al., 2018) as well as Neural Operators (Kovachki et al., 2021) treat depth as continuous time evolution; diffusion/score-based models (Shen et al., 2025) implement long iterative refinement; and chain-of-thought (Wei et al., 2022) prompting in LLMs explicitly unfolds multi-step reasoning. These families motivate studying depth-dependent generalization at the level of abstract state transitions rather than fixed architectures.

**This Study.** Relative to these threads, this study is *implementation-agnostic*: instead of parameterizing a specific architecture, we analyze *state-transition semigroups on metric spaces*, derive depth dependence of the *variance* via covering/Rademacher complexity of word balls, give conditions for *polynomial/logarithmic* growth, and couple them with *exponential/polynomial* bias decay to compute *optimal depth* across four regimes. This yields a unified lens for when and why *depth supremacy* (optimal  $k^* > 1$ ) emerges—particularly in iterative/hierarchical settings suggested above.

---

## 702 B PROOF OF THEOREM 1

703  
704 Here we show a slightly generalized version of Theorem 1 in the main text. Since the abstract model  
705 class, denoted  $\mathcal{H}_k$  in the main text, need not be a state-transition model  $H \circ B(k, F)$ , we simply  
706 write  $\mathcal{H}_{\text{abs}}$  instead. The main claim of this theorem is that in the setting where a hypothesis class  
707  $\mathcal{H}_{\text{imp}}$  with a specific implementation, such as ReLU, approximates a hypothesis class  $\mathcal{H}_{\text{abs}}$  without  
708 an implementation, like a state-transition model, the generalization error bound can be decomposed  
709 into bias and variance so that the bias depends on both  $\mathcal{H}_{\text{imp}}$  and  $\mathcal{H}_{\text{abs}}$ , while the variance depends  
710 only on  $\mathcal{H}_{\text{abs}}$  (and not on  $\mathcal{H}_{\text{imp}}$ ).

### 711 Assumptions for the high probability bounds:

- 712 • Let  $\mathcal{X}$  be a measurable space, and let  $T(\mathcal{X})$  be a topological space of real-valued measur-
- 713 able functions  $h : \mathcal{X} \rightarrow \mathbb{R}$
- 714 • Let  $\mathcal{Y}$  be a metric space
- 715 • Let  $\mathcal{Z} := \mathcal{X} \times \mathcal{Y}$ , and let  $(\mathcal{Z}, P)$  be a probability space
- 716 • Let  $S_n := (Z_i)_{i \in [n]} \sim P$  be an i.i.d. sample
- 717 • Let  $\mathcal{H}_{\text{abs}}$  be a separable subspace of  $T(\mathcal{X})$
- 718 • Let  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, b]$  be a bounded non-negative measurable function, and assume that
- 719 the first argument is  $\beta_\ell$ -Lipschitz uniformly in the second argument
- 720 • Write  $\ell \cdot h(x, y) := \ell(h(x), y)$ , and  $\ell \cdot \mathcal{H}_{\text{abs}} := \{\ell \cdot h(x, y) \mid h \in \mathcal{H}_{\text{abs}}\}$
- 721 • Let  $L[h] := \mathbb{E}_{(X, Y) \sim P}[\ell \cdot h(X, Y)]$  and  $\hat{L}(S_n)[h] := \frac{1}{n} \sum_{i=1}^n \ell \cdot h(X_i, Y_i)$

722  $\implies$  (the Lipschitz assumption yields  $\mathfrak{R}(\ell \cdot \mathcal{H}_{\text{abs}}) \leq \beta_\ell \mathfrak{R}(\mathcal{H}_{\text{abs}})$  and) Theorem 6 holds at Eqs. (13)  
723 and (15).

### 730 Assumptions for minimization:

- 731 • Let  $\mathcal{C}$  be a subspace of  $T(\mathcal{X})$
- 732 • Assume  $L$  is bounded below over  $\mathcal{C}$
- 733 • Assume  $\mathcal{H}_{\text{abs}}$  is compact in  $T(\mathcal{X})$ , and both  $L$  and  $\hat{L}(S_n)$  are (lower)semicontinuous on  
734  $\mathcal{H}_{\text{abs}}$
- 735 • Write  $L_{\mathcal{C}} := \inf_{c \in \mathcal{C}} L[c]$ ,  $L_{\mathcal{H}_{\text{abs}}} := \inf_{h \in \mathcal{H}_{\text{abs}}} L[h]$ , and  $\varepsilon_{\text{model}} := L_{\mathcal{H}_{\text{abs}}} - L_{\mathcal{C}}$
- 736 • Let  $f^* := \arg \min_{f \in \mathcal{H}_{\text{abs}}} L[f]$  and  $\hat{f} := \arg \min_{f \in \mathcal{H}_{\text{abs}}} \hat{L}(S_n)[f]$

737  $\implies$  Exist the infima  $L_{\mathcal{C}}$  and  $L_{\mathcal{H}_{\text{abs}}}$ , error  $\varepsilon_{\text{model}}$ , and the minimizers  $f^*$  and  $\hat{f}$ .

### 741 Assumptions for implementation error bound:

- 742 • Assume  $(T(\mathcal{X}), d_T)$  is a (pseudo-)metric space of real-valued measurable functions  $h :$   
743  $\mathcal{X} \rightarrow \mathbb{R}$
- 744 • Let  $\mathcal{H}_{\text{imp}}$  be a compact subspace in  $T(\mathcal{X})$
- 745 • Assume  $\varepsilon_{\text{imp}} := \sup_{f \in \mathcal{H}_{\text{abs}}} \inf_{h \in \mathcal{H}_{\text{imp}}} d_T(h, f) < \infty$
- 746 • Define embed :  $\mathcal{H}_{\text{abs}} \rightarrow \mathcal{H}_{\text{imp}}$  by embed $[f] := \arg \min_{h \in \mathcal{H}_{\text{imp}}} d_T(h, f)$
- 747 • Assume both  $L$  and  $\hat{L}(S_n)$  are  $\beta_L$ - and  $\beta_{\hat{L}(S_n)}$  Lipschitz over  $\mathcal{H}_{\text{abs}}$ , respectively

748  $\implies$  embed is well-defined, and for any  $f \in \mathcal{H}_{\text{abs}}$  and  $h := \text{embed}(f)$ ,  $|L[h] - L[f]| \leq$   
749  $\beta_L d_T(h, f) \leq \beta_L \varepsilon_{\text{imp}}$  and similarly  $|\hat{L}(S_n)[h] - \hat{L}(S_n)[f]| \leq \beta_{\hat{L}(S_n)} d_T(h, f) \leq \beta_{\hat{L}(S_n)} \varepsilon_{\text{imp}}$ .

**Conclusion: (Bias-Variance Decomposition)** Then, the excess risk and generalization gap are bounded as follows: With probability at least  $1 - \delta$  over an i.i.d. draw of sample  $S_n \sim P$ ,

$$L[\hat{h}] - L_C \leq \beta_L \varepsilon_{\text{imp}} + \varepsilon_{\text{model}} + 4\beta_\ell \hat{\mathfrak{R}}_{S_n}(\mathcal{H}_{\text{abs}}) + 2b\sqrt{\frac{2 \log 1/\delta}{n}} \quad (7)$$

$$L[\hat{h}] - \hat{L}(S_n)[\hat{h}] \leq (\beta_L + \beta_{\hat{L}(S_n)})\varepsilon_{\text{imp}} + 2\beta_\ell \hat{\mathfrak{R}}_{S_n}(\mathcal{H}_{\text{abs}}) + b\sqrt{\frac{2 \log 1/\delta}{n}}. \quad (8)$$

**Example of sufficient conditions:** For example, the following setting satisfy the above assumptions:

- $\mathcal{Y} := \mathbb{R}$
- $\ell$  is 1-Lipschitz in the first argument a.e., i.e.  $|\ell(y, y_0) - \ell(y', y_0)| \leq |y - y'|$  a.e.  $y_0 \in \mathcal{Y}$
- $T(\mathcal{X}) := L^\infty(\mathcal{X})$ ,  $\mathcal{C} := T(\mathcal{X})$ , and  $\mathcal{H}_{\text{abs}}, \mathcal{H}_{\text{imp}}$  are a separable and compact subspaces in  $L^\infty(\mathcal{X})$

so that  $|L[h] - L[f]| \leq \mathbb{E}_X |h(X) - f(X)| \leq \|h - f\|_\infty$  and  $|\hat{L}(S_n)[h] - \hat{L}(S_n)[f]| \leq \frac{1}{n} \sum_{i=1}^n |h(X_i) - f(X_i)| \leq \|h - f\|_\infty$  for every  $f, h \in L^\infty(\mathcal{X})$ ,  $\beta_L = \beta_{\hat{L}(S_n)} = 1$ ,  $L_C = \inf_{c \in \mathcal{C}} L[c] = 0$ ,  $\varepsilon_{\text{model}} = L_{\mathcal{H}_{\text{abs}}} - L_C = \inf_{h \in \mathcal{H}_{\text{abs}}} L[h]$ , and  $\varepsilon_{\text{imp}} = \sup_{f \in \mathcal{H}_{\text{abs}}} \inf_{h \in \mathcal{H}_{\text{imp}}} \|h - f\|_\infty$ .

*Proof.* Write  $\hat{h} := A(S_n) = \text{embed}(\hat{f})$ .

$$L[\hat{h}] - L_C = \underbrace{L[\hat{h}] - L[\hat{f}]}_{\text{implementation error} \leq \varepsilon_{\text{imp}}} + \underbrace{L[\hat{f}] - L_C}_{\text{(i) excess risk wrt } \mathcal{H}_{\text{abs}}} \quad (9)$$

$$\text{(i) } L[\hat{f}] - L_C = \underbrace{L[\hat{f}] - L_{\mathcal{H}_{\text{abs}}}}_{\text{(ii) variance}} + \underbrace{L_{\mathcal{H}_{\text{abs}}} - L_C}_{\text{model bias } \varepsilon_{\text{model}}} \quad (10)$$

$$\text{(ii) } L[\hat{f}] - L_{\mathcal{H}_{\text{abs}}} = \underbrace{L[\hat{f}] - \hat{L}(S_n)[\hat{f}]}_{\text{generalization gap wrt } \mathcal{H}_{\text{abs}}} + \underbrace{\hat{L}(S_n)[\hat{f}] - \hat{L}(S_n)[f^*]}_{\leq 0 \text{ by the minimality}} + \underbrace{\hat{L}(S_n)[f^*] - L[f^*]}_{\text{(estimation error)}} \quad (11)$$

$$\leq 2 \sup_{f \in \mathcal{H}_{\text{abs}}} |\hat{L}(S_n)[f] - L[f]| \quad (12)$$

$$\leq 4\beta_\ell \hat{\mathfrak{R}}_{S_n}(\mathcal{H}_{\text{abs}}) + 2b\sqrt{\frac{2 \log 1/\delta}{n}} \quad \text{with prob. } \geq 1 - \delta \text{ by Theorem 6} \quad (13)$$

So, we have the following high probability bound:

$$L[\hat{h}] - L_C \leq \beta_L \varepsilon_{\text{imp}} + \varepsilon_{\text{model}} + 4\hat{\mathfrak{R}}_{S_n}(\mathcal{H}_{\text{abs}}) + 2b\sqrt{\frac{2 \log 1/\delta}{n}}. \quad (14)$$

Similarly, the generalization gap (= test error - training error) is decomposed as follows:

$$\begin{aligned} L[\hat{h}] - \hat{L}(S_n)[\hat{h}] &= \underbrace{L[\hat{h}] - L[\hat{f}]}_{\text{implementation error}} + \underbrace{L[\hat{f}] - \hat{L}(S_n)[\hat{f}]}_{\text{generalization gap wrt } \mathcal{H}_{\text{abs}}} + \underbrace{\hat{L}(S_n)[\hat{f}] - \hat{L}(S_n)[\hat{h}]}_{\text{implementation error}} \\ &\leq \beta_L \varepsilon_{\text{imp}} + \sup_{f \in \mathcal{H}_{\text{abs}}} |L[f] - \hat{L}(S_n)[f]| + \beta_{\hat{L}(S_n)} \varepsilon_{\text{imp}} \\ &\leq (\beta_L + \beta_{\hat{L}(S_n)})\varepsilon_{\text{imp}} + 2\beta_\ell \hat{\mathfrak{R}}_{S_n}(\mathcal{H}_{\text{abs}}) + b\sqrt{\frac{2 \log 1/\delta}{n}} \quad \text{w. p. } \geq 1 - \delta \quad (15) \end{aligned}$$

□

## C PROOF OF THEOREM 2

Let  $(\mathcal{X}, d)$  be a pseudo-metric space, let  $C(\mathcal{X})$  be the normed space of real-valued continuous functions  $h : \mathcal{X} \rightarrow \mathbb{R}$  equipped with uniform norm  $\|\bullet\|_\infty$ , and let  $C(\mathcal{X}, \mathcal{X})$  be the pseudo-metric space of continuous maps  $f : \mathcal{X} \rightarrow \mathcal{X}$  equipped with the induced uniform pseudo-metric  $d_\infty(f, f') := \sup_{x \in \mathcal{X}} d(f(x), f'(x))$ . Additionally, given a finite sample  $S = \{X_1, \dots, X_n\} \subset \mathcal{X}$ , let  $d_S$  denote an induced empirical pseudo-metric  $d_S(f, f') := (\frac{1}{n} \sum_{x \in S} d(f(x), f'(x))^2)^{1/2}$ .

**Theorem 3** (Theorem 2, restated). *For any  $L$ -Lipschitz class  $H \subset C(\mathcal{X})$  and any totally-bounded continuous class  $F \subset C(\mathcal{X}, \mathcal{X})$  with identity element  $\text{id}_{\mathcal{X}}$  equipped with pseudo-metric  $d_S$ , the Rademacher complexity of composition class  $H \circ F := \{h \circ f : \mathcal{X} \rightarrow \mathbb{R} \mid h \in H, f \in F\}$  is decomposed into the sum of the complexities of output class  $H$  and hidden class  $F$  as follows:*

$$\hat{\mathfrak{R}}_S(H \circ F) \leq \hat{\mathfrak{R}}_S(H) + \frac{12L}{\sqrt{n}} \int_0^{\text{diam}(F)} \sqrt{\log N(F, d_S, \varepsilon)} d\varepsilon, \quad (16)$$

$$\mathfrak{R}_n(H \circ F) \leq \mathfrak{R}_n(H) + \frac{12L}{\sqrt{n}} \mathbb{E}_S \left[ \int_0^{\text{diam}(F)} \sqrt{\log N(F, d_S, \varepsilon)} d\varepsilon \right]. \quad (17)$$

Here  $\mathfrak{R}_n(H)$  denotes the Rademacher complexity of  $H$ ,  $\log N(F, d_{\infty}, \varepsilon)$  denotes the metric entropy of  $F$  in pseudo-metric  $d_{\infty}$ , and  $\text{diam}(F) := \sup_{f \in F} d_S(f, \text{id}_{\mathcal{X}})/2$ .

*Proof.* For the fixed sample  $S = \{X_1, \dots, X_n\}$  and the Rademacher vector  $\sigma = (\sigma_1, \dots, \sigma_n) \in \{\pm 1\}^n$  define the process

$$Z_f(S, H, \sigma) := \frac{1}{n} \sup_{h \in H} \sum_{i=1}^n \sigma_i h(f(X_i)), \quad f \in F. \quad (18)$$

Take an arbitrary anchor element  $f_0 \in F$  and write the centered version

$$\tilde{Z}_f := Z_f - Z_{f_0}, \quad \text{so that} \quad \tilde{Z}_{f_0} = 0. \quad (19)$$

Because every  $h \in H$  is  $L$ -Lipschitz w.r.t.  $d$ ,

$$|h(f(x)) - h(f'(x))| \leq Ld(f(x), f'(x)), \quad \forall x \in \mathcal{X} \quad (20)$$

hence for the difference of the processes

$$\begin{aligned} |\tilde{Z}_f - \tilde{Z}_{f'}| &= |Z_f - Z_{f'}| \\ &= \left| \frac{1}{n} \sup_{h \in H} \sum_{i=1}^n \sigma_i h(f(X_i)) - \frac{1}{n} \sup_{h \in H} \sum_{i=1}^n \sigma_i h(f'(X_i)) \right| \\ &\leq \frac{1}{n} \sup_{h \in H} \sum_{i=1}^n |h(f(X_i)) - h(f'(X_i))| \\ &\leq \frac{L}{n} \sum_{i=1}^n d(f(X_i), f'(X_i)) \leq Ld_S(f, f'). \end{aligned} \quad (21)$$

Conditioned on the sample  $S$ , the random variable  $\tilde{Z}_f - \tilde{Z}_{f'}$  is a Rademacher (hence sub-Gaussian) sum with variance proxy at most  $L^2 d_S(f, f')^2$ . Formally,

$$\mathbb{E}_{\sigma} \exp \left( \lambda [\tilde{Z}_f - \tilde{Z}_{f'}] \right) \leq \exp \left( \frac{\lambda^2}{2} L^2 d_S(f, f')^2 \right), \quad \forall \lambda \in \mathbb{R}, \quad (22)$$

so  $(\tilde{Z}_f)_{f \in F}$  is a centered sub-Gaussian process indexed by the pseudo-metric space  $(F, d_S)$  with variance parameter  $L$ .

The Dudley's integral bound (Martin J. Wainwright, 2019, Theorem 5.22) states that for such a centered sub-Gaussian process

$$\mathbb{E}_{\sigma} \sup_{f \in F} \tilde{Z}_f \leq \frac{CL}{\sqrt{n}} \int_0^{\text{diam}(F)} \sqrt{\log N(F, d_S, u)} du, \quad (23)$$

where  $C \leq 12$  is a universal constant. (The integral starts at 0 but the integrand is clipped at the bound of  $\tilde{Z}$ , hence the upper limit  $\sup_{f \in F} d_S(f, f_0) \leq \text{diam}(F)$ .)

864 Un-centering  $\tilde{Z}_f$ , we get

$$865 \mathbb{E}_\sigma \sup_{f \in F} Z_f \leq \mathbb{E}_\sigma \left[ \sup_{f \in F} [Z_f - Z_{f_0}] + \sup_{f \in F} Z_{f_0} \right] \\ 866 \\ 867 \\ 868 \\ 869 \leq \mathbb{E}_\sigma Z_{f_0} + \frac{12L}{\sqrt{n}} \int_0^{\text{diam}(F)} \sqrt{\log N(F, d_S, u)} du, \quad (24) \\ 870$$

871 Because the expectation is only over  $\sigma$ , Eq. (24) already gives the empirical Rademacher bound  
872  $\hat{\mathfrak{R}}_S(H \circ F)$  stated previously; taking the outer expectation over samples produces the population  
873 bound.  
874

875 If the identity map lies in the class  $F$ , i.e.

$$876 \text{id}_X \in F, \quad (25)$$

877 then we can simply take  $f_0 = \text{id}_X$ . Two consequences follow: First, the leading term becomes the  
878 ordinary empirical complexity of  $H$ .  
879

$$880 \mathbb{E}_\sigma Z_{f_0} = \frac{1}{n} \sup_{h \in H} \sum_{i=1}^n \sigma_i h(X_i) = \hat{\mathfrak{R}}_S(H). \quad (26) \\ 881 \\ 882$$

883 No extra optimization or data-dependent choice of  $f_0$  is needed. Second, the right-hand integral is  
884 unchanged. The second term still measures only how flexibly the maps in  $F$  can transport points  
885 (through  $N(F, d_S, u)$ ) and is totally independent of the special anchor we used. Hence, with  $\text{id}_X \in$   
886  $F$  the empirical bound specializes to

$$887 \hat{\mathfrak{R}}_S(H \circ F) \leq \hat{\mathfrak{R}}_S(H) + \frac{12L}{\sqrt{n}} \int_0^{\text{diam}(F)} \sqrt{\log N(F, d_S, u)} du. \quad (27) \\ 888$$

889 The algebra and all probabilistic ingredients are exactly the same; the anchor is just chosen once and  
890 for all, which simplifies both notation and interpretation.  $\square$   
891

## 892 D A SIMPLER VERSION OF THEOREM 2

893 We assume the composition structure  $\mathcal{H}_{\text{abs}} = H \circ F$  for the abstract model class, where  $H$  and  
894  $F$  respectively are collections of real-valued continuous functions  $h_o : \mathcal{X} \rightarrow \mathbb{R}$  and continuous  
895 self-maps  $f : \mathcal{X} \rightarrow \mathcal{X}$ , i.e.  $H \subset C(\mathcal{X})$  and  $F \subset C(\mathcal{X}, \mathcal{X})$ . In the following, we assume all  
896  $\mathcal{X}, C(\mathcal{X}), C(\mathcal{X}, \mathcal{X})$  are (pseudo)metric spaces to use covering numbers, every  $h \in H$  is  $B$ -bounded  
897 and  $L$ -Lipschitz, and  $F$  is equicontinuous and its *range*  $F(\mathcal{X})$  is compact, so that  $F$  and  $H$  are  
898 totally bounded thus their covering numbers are finite.  
899

900 Let  $S = (X_1, \dots, X_n)$  be a sample from  $\mathcal{X}$  and define the empirical  $L_2(S)$  metric  $\|h -$   
901  $h'\|_S := \sqrt{\frac{1}{n} \sum_{i=1}^n (h(X_i) - h'(X_i))^2}$ . According to *Dudley's entropy integral* (Theorem 8), the  
902 Rademacher complexity  $\hat{\mathfrak{R}}_S(H \circ F)$  is bounded by the logarithm of the covering number (*metric*  
903 *entropy*) of hypothesis class. Since  $\|h - h'\|_S \leq \|h - h'\|_\infty$ , the covering number in  $L_2(S)$  is no  
904 larger than that in  $\|\cdot\|_\infty$ . So Dudley's bound yields  
905

$$906 \hat{\mathfrak{R}}_S(H \circ F) \leq \frac{12}{\sqrt{n}} \int_0^{\text{diam}(H \circ F, \|\cdot\|_S)} \sqrt{\log N(H \circ F, \|\cdot\|_S, \varepsilon)} d\varepsilon \quad (28) \\ 907$$

$$908 \leq \frac{12}{\sqrt{n}} \int_0^B \sqrt{\log N(H \circ F, \|\cdot\|_\infty, \varepsilon)} d\varepsilon. \quad (29) \\ 909 \\ 910$$

911 Then, according to Lemma 8, the covering number in the uniform norm of the composite class  $H \circ F$   
912 can be decomposed into a product of covering numbers; or in terms of the Rademacher complexity  
913 bound, the sum of metric entropies: For every positive number  $\varepsilon > 0$ ,  
914

$$915 \log N(H \circ F, \|\cdot\|_\infty, \varepsilon) \leq \log N(H, \|\cdot\|_\infty, \varepsilon/2) + \log N(F, d_\infty, \varepsilon/(2L)). \quad (30) \\ 916$$

917 Insert Eq. (30) and use  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ , we have the following decomposition of the Rademacher  
complexity for composition class:

**Theorem 4** (Theorem 2, simpler parallel). *For any sample  $S$ ,*

$$\hat{\mathfrak{R}}_S(H \circ F) \leq \frac{12}{\sqrt{n}} \int_0^B \left\{ \sqrt{\log N(H, \|\cdot\|_\infty, \varepsilon/2)} + \sqrt{\log N(F, d_\infty, \varepsilon/(2L))} \right\} d\varepsilon. \quad (31)$$

*Averaging over samples gives the same bound for  $\mathfrak{R}_n(H \circ F)$ . Changing variables ( $u = \varepsilon/2, v = \varepsilon/(2L)$ ), yields a completely explicit decomposition:*

$$\mathfrak{R}_n(H \circ F) \leq \frac{24}{\sqrt{n}} \left[ \int_0^{B/2} \sqrt{\log N(H, \|\cdot\|_\infty, u)} du + L \int_0^{B/(2L)} \sqrt{\log N(F, d_\infty, v)} dv \right]. \quad (32)$$

These inequalities decompose the statistical complexity of the composite class into separate contributions from the feature layer  $F$  and the output layer  $H$ .

## E DETAILS ON GROWTH RATE ANALYSIS

We provide details on the conditions and examples overviewed in Section 4.

### E.1 ASSUMPTIONS AND NOTATION

Throughout,  $(\mathcal{X}, d)$  is a (not necessarily compact) metric space, and  $C(\mathcal{X}, \mathcal{X})$  denotes the set of continuous self-maps of  $\mathcal{X}$ .

**Uniform metric** On  $C(\mathcal{X}, \mathcal{X})$  we use the uniform metric

$$d_\infty(f, g) := \sup_{x \in \mathcal{X}} d(f(x), g(x)).$$

When  $\mathcal{X}$  is non-compact,  $d_\infty$  is not well-defined as the metric may diverge. We will always guarantee finiteness in one of the following two ways (either assumption suffices):

- ( $d_b$ ) replace  $d$  by its bounded version  $d_b(x, y) := \min\{1, d(x, y)\}$ ;
- ( $C_b$ ) fix a compact set  $K_0 \subset \mathcal{X}$  and assume all maps under discussion (and all words built from them) take values in  $K_0$ .

All the three metrics— $C(\mathcal{X}, \mathcal{X})$  with compact  $\mathcal{X}$  and  $d_\infty$ ,  $C(\mathcal{X}, \mathcal{X})$  with non-compact  $\mathcal{X}$  and  $d_b$ , and  $C_b(\mathcal{X}, \mathcal{X})$  with non-compact  $\mathcal{X}$  and  $d_\infty$ —induce the uniform topology.

**Arzelà–Ascoli and compact-open topology** Arzelà–Ascoli theorems assert relative compactness in the *compact-open topology* (COT), which is *weaker* than uniform topology (UT) in general. COT is equivalent to the compact-convergence topology, or the topology induced from the uniform convergence on every compact sets. COT *coincides with* UT when  $\mathcal{X}$  is compact. See Appendix I for more details on Arzelà–Ascoli theorems.

**Homeo and isometry groups**  $\text{Homeo}(\mathcal{X}) \subset C(\mathcal{X}, \mathcal{X})$  denotes the homeomorphism group of  $\mathcal{X}$ , that is, the set of all *bijective* bi-continuous self-maps with the function composition as group operation.  $\text{Isom}(\mathcal{X}) \subset C(\mathcal{X}, \mathcal{X})$  denotes the isometry group of  $(\mathcal{X}, d)$ , that is, the set of all *bijective* distance-preserving self-maps with the function composition as group operation. The above group operations are continuous in compact-open topology, or the topology induced by the following metric  $d_\infty$  (restricted to an arbitrary compact subset  $K \subset \mathcal{X}$ ). Hence  $\text{Homeo}(\mathcal{X})$  and  $\text{Iso}(\mathcal{X})$  are topological groups.

**Semigroup and word balls** For  $F \subset C(\mathcal{X}, \mathcal{X})$ , write  $\langle F \rangle$  for the semigroup generated by  $F$  under composition. For  $k \in \mathbb{N}$ , let  $B(k, F)$  be the set of maps obtainable as a composition of at most  $k$  maps from  $F$ .

**Covering and packing numbers** For a metric space  $(M, \rho)$ , write

$$N(A, \rho, \varepsilon) := \min \left\{ |C| \mid C \subset M \text{ s.t. } A \subset \bigcup_{y \in C} B_\rho(y, \varepsilon) \right\}$$

for the  $\varepsilon$ -covering number of  $A \subset M$ , and

$$M(A, \rho, \varepsilon) := \max \{|S| \mid S \subset A, \rho(x, y) \geq \varepsilon \forall x \neq y \in S\}$$

for the  $\varepsilon$ -packing number. See also Appendix H for more details on covering and packing numbers.

**Lipschitz and co-Lipschitz constants** For  $f \in C(\mathcal{X}, \mathcal{X})$  and  $S \subset \mathcal{X}$  nonempty,

$$\text{Lip}_S(f) := \sup_{x \neq y \in S} \frac{d(f(x), f(y))}{d(x, y)}, \quad \text{coLip}_S(f) := \inf_{x \neq y \in S} \frac{d(f(x), f(y))}{d(x, y)}.$$

We say “ $f$  is uniformly expanding on  $S$ ” if  $\text{coLip}_S(f) > 1$ .

**Word length in  $\langle F \rangle$ .** For  $h \in \langle F \rangle$ , denote by  $\ell_F(h)$  the least  $\ell$  such that  $h$  is a composition of  $\ell$  maps from  $F$ . In all results below, when we assert that certain auxiliary maps  $h$  belong to  $\langle F \rangle$ , we also assume that  $\ell_F(h)$  is bounded by a constant independent of  $k$ .

**Subshift with ultrametric** On the *one-sided full shift*  $\Sigma_m^+ = [m]^{\mathbb{N}}$  (with  $[m] = \{1, \dots, m\}$  discrete), fix  $\theta \in (0, 1)$  and define the ultrametric

$$d_\theta(x, y) = \begin{cases} 0, & x = y, \\ \theta^{n(x, y)-1}, & x \neq y, \end{cases} \quad \text{where } n(x, y) := \min\{i \geq 1 \mid x_i \neq y_i\}.$$

This is an ultrametric (satisfies *strong triangle inequality*  $d(x, z) \leq \max\{d(x, y), d(y, z)\}$ ), it induces the product topology, and it makes  $\Sigma_m^+$  a compact, totally disconnected, perfect (Cantor-type) space. Distances take only the values  $1, \theta, \theta^2, \dots$ ; in particular, sequences that differ in the first symbol are at distance 1.

For a finite word  $u = u_1 \cdots u_\ell$ , the (*prefix*) *cylinder* of depth  $\ell$  is given by

$$\llbracket u \rrbracket := \{w \in \Sigma_m^+ \mid w_1 = u_1, \dots, w_\ell = u_\ell\},$$

a clopen set; cylinders form a basis and the depth- $\ell$  cylinders partition  $\Sigma_m^+$  into  $m^\ell$  pieces. In the metric  $d_\theta$ ,  $\text{diam}(\llbracket u \rrbracket) = \theta^\ell$ , and any two distinct depth- $\ell$  cylinders are separated by at least  $\theta^{\ell-1}$ . Moreover, balls coincide with cylinders: for  $w \in \Sigma_m^+$  and  $\ell \geq 0$ ,  $\overline{B}_{d_\theta}(w, \theta^\ell) = \llbracket w_1 \cdots w_\ell \rrbracket$ . The left shift  $\sigma(w)_i = w_{i+1}$  satisfies  $\sigma(\llbracket au \rrbracket) = \llbracket u \rrbracket$  and  $\sigma^{-1}(\llbracket u \rrbracket) = \bigsqcup_{a \in [m]} \llbracket au \rrbracket$ .

## E.2 BASIC FACTS ON COVERING AND PACKING NUMBERS

**Lemma 1** (Packing-Covering). *For any totally bounded metric space  $A$  and  $\varepsilon > 0$ ,*

$$M(A, 2\varepsilon) \leq N(A, \varepsilon) \leq M(A, \varepsilon).$$

**Lemma 2** (Lipschitz Embedding). *Let  $(X, d_X), (Y, d_Y)$  be metric spaces. Suppose  $\phi : X \rightarrow Y$  is  $L$ -Lipschitz, then for every subset  $S \subset X$ ,*

$$N(\phi(S), d_Y, L\varepsilon) \leq N(S, d_X, \varepsilon), \quad M(\phi(S), d_Y, L\varepsilon) \leq M(S, d_X, \varepsilon).$$

**Lemma 3** (Subadditivity). *For any metric space  $(M, d)$ , subsets  $F, G \subset M$ , and  $\varepsilon > 0$ , we have*

$$N(F \cup G, \varepsilon) \leq N(F, \varepsilon) + N(G, \varepsilon).$$

*Proof.* Let  $A, B$  be  $\varepsilon$ -coverings of  $F, G$  respectively. Then,  $A \cup B$  is an  $\varepsilon$ -covering of  $F \cup G$  because for any  $h \in F \cup G$  there exists  $c \in A \cup B$  satisfying  $d(h, c) \leq \varepsilon$ . Thus,  $N(F \cup G, \varepsilon) \leq |A \cup B| \leq |A| + |B| = N(F, \varepsilon) + N(G, \varepsilon)$ .  $\square$

**Lemma 4** (Sub-multiplicativity). *For any bi-Lipschitz metric semigroup  $(M, d)$  with  $\sup_{f \in M} d(fx, fy) \leq \lambda d(x, y)$  (left Lipschitz) and  $\sup_{f \in M} d(xf, yf) \leq \rho d(x, y)$  (right Lipschitz), for any subsets  $F, G \subset M$ , and  $\varepsilon, \delta > 0$ , we have*

$$N(FG, \varepsilon + \delta) \leq N(F, \varepsilon/\rho)N(G, \delta/\lambda).$$

1026 *Proof.* Let  $A, B$  be  $\alpha, \beta$ -coverings of  $F, G$  respectively. Then,  $AB$  is an  $(\rho\alpha + \lambda\beta)$ -covering of  $FG$   
1027 because for any  $fg \in FG$  there exists  $ab \in AB$  satisfying  $d(fg, ab) \leq d(fg, ag) + d(ag, ab) \leq$   
1028  $\rho d(f, a) + \lambda d(g, b) = \rho\alpha + \lambda\beta$ . Thus,  $N(FG, \rho\alpha + \lambda\beta) \leq |AB| \leq |A||B| = N(F, \alpha)N(G, \beta)$ .  
1029 Letting  $\alpha = \varepsilon/\rho, \beta = \delta/\lambda$  yields the assertion.  $\square$

1030

1031 Also recall:

1032 **Lemma 5** (Right-Composition is 1-Lipschitz).  $d_\infty(a \circ f, b \circ f) \leq d_\infty(a, b)$   
1033

1034 *Proof.*  $d_\infty(a \circ f, b \circ f) = \sup_{x \in \mathcal{X}} d(a(f(x)), b(f(x))) \leq \sup_{y \in \mathcal{X}} d(a(y), b(y)) = d_\infty(a, b)$   $\square$   
1035

1036 **Lemma 6** (Left-Composition is Lipschitz).  $d_\infty(f \circ a, f \circ b) \leq \text{Lip}(f)d_\infty(a, b)$   
1037

1038 *Proof.*  $d_\infty(f \circ a, f \circ b) = \sup_{x \in \mathcal{X}} d(f(a(x)), f(b(x))) \leq \sup_{x \in \mathcal{X}} \text{Lip}(f)d(a(x), b(x)) =$   
1039  $\text{Lip}(f)d_\infty(a, b)$   $\square$

1040

1041 **Lemma 7** (Packing–Covering via Finite Probes). Let  $P = \{x_1, \dots, x_m\} \subset \mathcal{X}$  be finite and de-  
1042 fine  $d_P((y_j), (z_j)) := \max_{1 \leq j \leq m} d(y_j, z_j)$  on  $\mathcal{X}^m$ , and  $\text{ev}_P : C(\mathcal{X}, \mathcal{X}) \rightarrow \mathcal{X}^m, \text{ev}_P(f) =$   
1043  $(f(x_1), \dots, f(x_m))$ . Then  $\text{ev}_P$  is 1-Lipschitz:  $d_P(\text{ev}_P(f), \text{ev}_P(g)) \leq d_\infty(f, g)$ . If  $\text{ev}_P(S) \subset \mathcal{X}^m$   
1044 contains  $M$  points that are pairwise  $\delta$ -separated (in  $d_P$ ), then for every  $\varepsilon < \delta/2$ ,

1045

$$N(S, d_\infty, \varepsilon) \geq M.$$

1046

1047 *Proof.* The 1-Lipschitz claim is immediate:

1048

$$d(f(x_j), g(x_j)) \leq \sup_{x \in \mathcal{X}} d(f(x), g(x)) = d_\infty(f, g).$$

1049

1051 If  $\text{ev}_P(s_1), \dots, \text{ev}_P(s_M)$  are  $\delta$ -separated, then  $d_\infty(s_i, s_j) \geq d_P(\text{ev}_P(s_i), \text{ev}_P(s_j)) \geq \delta$ . Any ball  
1052 of radius  $\varepsilon < \delta/2$  in the  $d_\infty$ -metric can contain at most one of the  $s_i$ 's, so at least  $M$  balls are  
1053 needed.  $\square$

1054

### 1055 E.3 CONDITIONS FOR SATURATION AND POLYNOMIAL GROWTH

1056

#### 1057 E.3.1 P1. EQUICONTINUOUS SEMIGROUP ON COMPACT DOMAIN $\implies$ SATURATION IN $k$

1058

1059 **Condition P1** (Equicontinuous Semigroup on Compact Domain Saturates). Assume  $\mathcal{X}$  is compact.  
1060 Suppose (at least) one of the following assumptions is satisfied:

1061

1. Semigroup  $\langle F \rangle$  is (pre)compact,

1062

2a. Semigroup  $\langle F \rangle$  is equicontinuous,

1063

2b. Semigroup  $\langle F \rangle$  is uniformly Lipschitz:  $\text{Lip}\langle F \rangle < \infty$ , or

1064

2c. Generators  $F$  are non-expanding:  $\text{Lip} F \leq 1$ .

1065

1066 Then, closure semigroup  $G := \overline{\langle F \rangle}^{d_\infty} \subset C(\mathcal{X}, \mathcal{X})$  is compact and equicontinuous, and for all  $\varepsilon > 0$   
1067 and all  $k$

1068

$$N(B(k, F), \varepsilon) \leq N(G, \varepsilon) \quad (< \infty),$$

1069

1070 hence no dependence on  $k$ .

1071

1072 *Proof.* If  $\mathcal{X}$  is compact, then  $C(\mathcal{X}, \mathcal{X})$  is complete in  $d_\infty$ . So the closure of any precompact subset  
1073 in  $C(\mathcal{X}, \mathcal{X})$  is compact. Besides, Arzelà–Ascoli yields if  $\mathcal{X}$  is compact, then (1) any compact  
1074 subset of  $C(\mathcal{X}, \mathcal{X})$  is equicontinuous, and (2) any equicontinuous subset is relatively compact (thus  
1075 compact). (Namely, any subset in  $C(\mathcal{X}, \mathcal{X})$  with compact  $\mathcal{X}$  is precompact  $\iff$  compact  $\iff$   
1076 equicontinuous.) Both Assumptions 1 and 2a are straightforward. Since uniform Lipschitz implies  
1077 equicontinuous, Assumptions 2b and 2c are reduced to Assumption 2a. Hence each assumption  
1078 implies  $G$  is compact (and equicontinuous). Trivial inclusion  $B(k, F) \subseteq G$  yields the bound, and  
1079 compactness of  $G$  gives finiteness of  $N(G, \varepsilon)$ .  $\square$

1080 *Example 1* (Rotations on Circle). Let  $\mathcal{X} = \mathbb{S}^1$ ,  $A \subset \mathbb{S}^1$  compact, and put  $F := \{R_\alpha \mid \alpha \in A\}$   
 1081 (rotations). All are isometries, and the closure of the generated subgroup  $\langle F \rangle$  is a compact torus  
 1082 (either a finite set or the full circle group depending on rational relations) contained in the rotation  
 1083 group  $G = O(1)$ , so

$$1084 \quad N(B(k, F), d_\infty, \varepsilon) \leq N(O(1), d_\infty, \varepsilon) = N(\mathbb{S}^1, d, \varepsilon) \quad (k\text{-independent}).$$

1086 *Example 2* (Finite Isometry Group on Finite Set). If  $\mathcal{X}$  is finite and  $F \subset \text{Iso}(\mathcal{X})$ , then  $G$  is finite;  
 1087  $N(B(k, F), \varepsilon)$  is bounded by  $|G|$  for all  $k$ .

1088 **Condition P1'** (Contraction to Compact Invariant Set). *Assume*

- 1089
- 1090 1. (uniform contraction)  $\sup_{f \in F} \text{Lip}(f) \leq c < 1$ ,
  - 1091 2. (compact attractor) there exists a nonempty compact  $F$ -invariant set  $A \subset \mathcal{X}$  (i.e.,  $f(A) \subset$   
 1092  $A$  for all  $f \in F$ ), and
  - 1093 3. (compact absorbing set) there exist  $L \in \mathbb{N}$  and bounded set  $K \subset \mathcal{X}$  such that  $f(\mathcal{X}) \subset K$   
 1094 for all depth- $L$  map  $f \in F^L$ .

1095  
 1096 Then for every  $\varepsilon > 0$  there exists

$$1097 \quad m(\varepsilon) := L + \left\lceil \log_{1/c} \left( \frac{2 \text{diam}(K)}{\varepsilon} \right) \right\rceil$$

1098 such that for all  $k \geq m(\varepsilon)$ ,

$$1099 \quad N(B(k, F), d_\infty, \varepsilon) \leq N(A, d, \frac{\varepsilon}{2}) + \sum_{\ell=0}^{m(\varepsilon)-1} [N(F, d_\infty, \frac{\varepsilon}{2})]^\ell.$$

1100  
 1101 In particular, the right-hand side is independent of  $k$ , so growth in  $k$  saturates at any fixed  $\varepsilon$  to the  
 1102 entropy of the attractor  $A$ .

1103 *Remark 2.* When  $\mathcal{X}$  is compact, we can simply set  $A = K = \mathcal{X}$  and  $L = 0$ .

1104 *Remark 3.* A single saturation map  $\sigma$  such as  $\tanh$ , logistic map, and clipping yields bounded  
 1105 absorbing set  $K = \sigma(\mathcal{X})$  with  $L = 1$ .

1106  
 1107 *Proof.* Pick any  $a_0 \in A$ . For any word  $w \in F^\ell$ , we have  $\text{Lip}(w) \leq c^\ell$  and  $w(a_0) \in A$  (invariance).  
 1108 If  $\ell \geq m(\varepsilon)$ , then

$$1109 \quad \sup_{x \in \mathcal{X}} d(w(x), w(a_0)) \leq c^{\ell-L} \text{diam}(K), \quad (\text{bounded absorbing set})$$

1110 and the right-hand side is  $\leq \varepsilon/2$ , i.e.

$$1111 \quad d_\infty(w, \text{const}_{w(a_0)}) \leq \varepsilon/2.$$

1112 Hence all sufficiently deep words, i.e.  $\bigcup_{\ell \geq m(\varepsilon)} F^\ell$ , are covered by the  $\varepsilon/2$ -thickening of the set of  
 1113 constant maps landing in  $A$ :  $\text{const}_q : q \in A$ . These constants at resolution  $\varepsilon/2$  are parameterized  
 1114 by an  $\varepsilon/2$ -net of  $A$ , giving the bound  $N(A, \varepsilon/2)$ . For the finitely many shallow layers  $\ell < m(\varepsilon)$ ,  
 1115 use the submultiplicative covering inequality for composition to get  $N(F^\ell, \varepsilon) \leq [N(F, \varepsilon/2)]^\ell$ , and  
 1116 sum over  $\ell$ .  $\square$

1117 Intuitively, deep words are *almost constant* (their images have diameter  $\leq c^{\ell-L} \text{diam}(K)$ ), so at  
 1118 scale  $\varepsilon$  only the landing point  $w(a_0) \in A$  matters; hence the dependence on  $k$  disappears once  $\ell$   
 1119 exceeds an  $\varepsilon$ -dependent “memory length”  $m(\varepsilon)$ .

1120 *Example 3* (Cantor Attractor (Compact, Free, Contraction)). Let  $\mathcal{X} = [0, 1]$  and  $F =$   
 1121  $f_0(x) = x/3, f_1(x) = (x+2)/3$ . Then  $\sup \text{Lip}(f_i) = 1/3 = c < 1$ . The attractor  $A$  is the  
 1122 middle-third Cantor set  $C$  (compact,  $F$ -invariant). By Condition P1', for every  $\varepsilon > 0$  and  
 1123  $k \geq m(\varepsilon) = \lceil \log_3(2/\varepsilon) \rceil$ ,

$$1124 \quad N(B(k, F), \varepsilon) \lesssim N(C, \varepsilon/2),$$

1125 independent of  $k$ . (Indeed, every  $w \in F^\ell$  has the form  $w(x) = 3^{-\ell}x + b_w$  and hence is  $\varepsilon/2$ -close  
 1126 to the constant map  $x \mapsto b_w \in C$ .)

1134 E.3.2 P2. BI-LIPSCHITZ NILPOTENT LIE GROUP ACTION  $\implies$  POLYNOMIAL  
 1135 UPPERBOUNDS IN  $k$   
 1136

1137 **Condition P2** (Bi-Lipschitz Nilpotent Lie Group Action Grows Polynomially). *Let  $H$  be a con-*  
 1138 *ected, simply-connected nilpotent Lie group, let  $d_H$  be a left-invariant Riemannian metric on  $H$ ,*  
 1139 *and let  $D$  be the Guivarc'h–Bass homogeneous dimension of  $H$ . Assume that there is a faithful*  
 1140 *action  $\alpha : H \rightarrow \text{Homeo}(\mathcal{X})$  such that  $\langle F \rangle \subset \alpha(H)$ , and suppose bi-Lipschitz control: There exist*  
 1141 *constants  $0 < c \leq C < \infty$  such that for all  $g, h \in H$ ,*

$$1142 \quad cd_H(g, h) \leq d_\infty(\alpha(g), \alpha(h)) \leq Cd_H(g, h).$$

1143 *Then, there exist a constant  $C$  such that for every  $\varepsilon > 0$  and  $k \geq 1$ ,*

$$1144 \quad N(B(k, F), \varepsilon) \leq C \left(1 + \frac{k}{\varepsilon}\right)^D.$$

1145  
 1146  
 1147 *Remark 4.* The Guivarc'h–Bass homogeneous dimension bridges the algebraic property such as the  
 1148 nilpotency of Lie group  $H$  with the geometric property, namely the polynomial volume growth rate  
 1149 of  $H$ . We briefly reviewed the theory in Appendix J based on Breuillard (2014).  
 1150

1151 *Proof.* Let  $S := \alpha^{-1}(F) \subset H$ , a compact generating set for the subgroup  $H_0 := \langle S \rangle$  (closed in  
 1152  $H$ ). Let  $\ell_S$  be the word length with respect to  $S$ , let  $B[\ell_S](k) := \{h \in H \mid \ell_S(h) \leq k\}$  be the  
 1153 word ball of length at most  $k$ , and let  $B[d_H](k) := \{h \in H \mid d_H(e, h) \leq k\}$  be the geodesic ball of  
 1154 radius  $k$ . Any word of length  $\leq k$  in  $F$  equals  $\alpha(g)$  with  $g \in B[\ell_S](k)$ . By standard comparability  
 1155 of word and geodesic metrics on nilpotent Lie groups, there is  $A \geq 1$  so that

$$1156 \quad B[\ell_S](k) = \{h \in H \mid \ell_S(h) \leq k\} \subset B[d_H](Ak) = \{h \in H \mid d_H(e, h) \leq Ak\}.$$

1157 By the bi-Lipschitz control,  
 1158

$$1159 \quad N(\alpha(B[d_H](Ak)), d_\infty, \varepsilon) \leq N(B[d_H](Ak), d_H, \varepsilon/C).$$

1160 Balls in  $H$  have polynomial metric entropy: since  $H$  has homogeneous dimension  $D$ , one has  
 1161  $N(B[d_H](R), \delta) \leq C'(1 + \frac{R}{\delta})^D$  for every  $R >$  and  $\delta > 0$  (cover by  $\delta$ -lattices in exponential  
 1162 coordinates). Therefore

$$1163 \quad N(B(k, F), d_\infty, \varepsilon) \leq N(\alpha(B[d_H](Ak)), d_\infty, \varepsilon) \leq C''(1 + k/\varepsilon)^D.$$

1165  $\square$

1166 *Example 4* (Translations on Euclidean space (Non-compact, Abelian, Isometric)). Let  $\mathcal{X} = \mathbb{R}^d$ ,  
 1167  $H = \mathbb{R}^d$  acting by translation  $\alpha(v) := x \mapsto x + v, v \in H$ . For any compact  $A \subset H$ , put  $F :=$   
 1168  $\alpha(A) \subset C(\mathcal{X}, \mathcal{X})$ . The homogeneous dimension is  $D = \dim \text{span } A (\leq d)$ ; and  $d_\infty(\alpha(v), \alpha(w)) =$   
 1169  $\|v - w\|$  (independent of  $x$ ), so the bi-Lipschitz constants are  $c = C = 1$ . Therefore,  
 1170

$$1171 \quad N(B(k, F), d_\infty, \varepsilon) = N(B(k, A), \|\cdot\|, \varepsilon) \lesssim (1 + k/\varepsilon)^D.$$

1172 We note that if  $\mathcal{X} = \mathbb{T}^d = \mathbb{R}^d/\mathbb{Z}^d$  (torus, compact) instead of  $\mathbb{R}^d$  (non-compact), then it grows at  
 1173  $O((k/\varepsilon)^D)$  for early  $k$ , and saturates to  $O(1/\varepsilon^D)$  by Condition P1.

1174 *Example 5* (Shear on Torus (Compact, Abelian, Expanding)). Let  $\mathcal{X} = \mathbb{T}^2$ ,  $H = \mathbb{Z}$  acting by shear  
 1175  $\alpha(n) := (x, y) \mapsto (x, x + ny), n \in H$ . Put  $F := \alpha(1) \subset C(\mathcal{X}, \mathcal{X})$ . The homogeneous dimension  
 1176 is  $D = 1$ ; and  $d_\infty(\alpha(n), \alpha(m)) = \delta_{nm}$ , so the bi-Lipschitz constants are  $c = C = 1$ . Therefore for  
 1177  $\varepsilon < 1$ ,

$$1178 \quad N(B(k, F), d_\infty, \varepsilon) = N(\{0\} \cup [k], \delta_{\bullet\bullet}, \varepsilon) = 1 + k.$$

1179 *Example 6* (Discrete Heisenberg Group on Torus (Compact, Nilpotent, Expanding)). Take  $H$  the  
 1180 Heisenberg group acting affine-linearly on  $\mathbb{T}^3$  (or  $\mathbb{R}^3/\mathbb{Z}^3$ ) by translations along orbits. The homo-  
 1181 geneous dimension is  $D = 4$ ; with a compact generating set  $F \subset \alpha(H)$ , one gets  $N(B(k, F), \varepsilon) \lesssim$   
 1182  $(1 + k/\varepsilon)^4$ .

1183 *Example 7* (Upper-Triangular Unipotent Group on Euclidean Space (Non-compact, Nilpotent, Ex-  
 1184 panding)). Let  $\mathcal{X} = \mathbb{R}^d$ ,  $H = UT_d(\mathbb{R}) = \{I_d + N_d \mid N_d \in GL_d(\mathbb{R}) \text{ strictly upper triangular}\}$   
 1185 acting by a bi-Lipschitz  $\alpha$ . Suppose a compact generating set  $F \subset \alpha(H)$ . The homogeneous  
 1186 dimension of  $H$  is  $D = \frac{d(d-1)(d+1)}{6}$ . So,  
 1187

$$1188 \quad N(B(k, F), d_\infty, \varepsilon) \asymp (1 + k/\varepsilon)^D.$$

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

#### E.4 CONDITIONS FOR (SUPER-/DOUBLE-)EXPONENTIAL GROWTHS

##### E.4.1 E1: FREE SEMIGROUP + ONE-POINT UNIFORM SEPARATION $\implies$ EXPONENTIAL LOWERBOUNDS IN $k$

**Condition E1** (Free Semigroup + Point Uniform Separation). *Let  $F = \{f_1, \dots, f_r\} \subset C(\mathcal{X}, \mathcal{X})$  with  $r \geq 2$ . Assume:*

1. (Freeness at each length) *For every  $k$ , the map  $u \mapsto f_u := f_{i_k} \circ \dots \circ f_{i_1}$  is injective on words  $u \in [r]^k$  (i.e., the semigroup generated by  $F$  is free on these generators).*
2. (Uniform separation at a base point) *There exist  $x_* \in \mathcal{X}$  and  $\delta > 0$  such that for all  $k$  and all distinct words  $u, v$  of length  $k$ ,  $d(f_u(x_*), f_v(x_*)) \geq \delta$ .*

*Then, for every  $k$  and every  $\varepsilon < \delta/2$ ,*

$$N(B(k, F), d_\infty, \varepsilon) \geq r^k.$$

*Proof.* Apply Lemma 7 with  $P = \{x_*\}$ . For each fixed  $k$ , the  $r^k$  words of length  $k$  give  $r^k$  maps whose images at  $x_*$  are  $\delta$ -separated; hence for  $\varepsilon < \delta/2$  at least  $r^k$  balls are needed to cover  $B(k, F)$ .  $\square$

*Example 8* (Free Group with Word Metric (Non-Compact, Free, Isometric)). Let  $\mathcal{X} = F_r$  ( $r \geq 2$ ) with the standard word metric  $d_G(u, v) = |u^{-1}v|$  (reduced word length), and let  $F = \{L_{a_1}, \dots, L_{a_r}\}$  where  $L_{a_i}(x) = a_i x$ . Each  $L_{a_i}$  is an isometry of  $(\mathcal{X}, d_G)$ ; the semigroup is free. Take  $x_* = e$  (identity). If  $u \neq v$  are words of the same length  $k$ , then  $u^{-1}v$  is a nontrivial reduced word of length  $\geq 2$  (recall “first rightmost mismatch gives  $a^{-1}b$ ”). Hence

$$d_G(L_u(x_*), L_v(x_*)) = d_G(u, v) = |u^{-1}v| \geq 2,$$

so Condition E1 applies with  $\delta = 2$ : for all  $\varepsilon < 1$ ,

$$N(B(k, F), d_\infty, \varepsilon) \geq r^k.$$

*Remark 5.* For distinct  $g, h \in F_r$ ,

$$d_\infty(L_g, L_h) = \sup_x d_G(gx, hx) = \sup_x d_G(e, x^{-1}g^{-1}hx) = \infty,$$

since  $|x^{-1}g^{-1}hx| \rightarrow \infty$  along  $x = s^n$  with  $s$  avoiding the boundary letters of  $g^{-1}h$ . This does not harm the lower bound: we only need that  $d_\infty(L_u, L_v) \geq d_G(L_u(e), L_v(e)) \geq 2$ , then apply Lemma 7 with  $\varepsilon < 1$ .

**Condition E1'** (Isometry with Coarse Coding). *Assume  $F \subset \text{Iso}(\mathcal{X})$  and that  $d_\infty(g, h) < \infty$  for all  $g, h \in \langle F \rangle$  (this holds, for example, when  $d$  is bi-invariant:  $d(axb, ayb) = d(x, y)$ ). Suppose furthermore:*

1. (Freeness) *The semigroup generated by  $F$  is free (no relations in positive words).*
2. (Coarse embedding at a base point) *There exist  $x_* \in \mathcal{X}$  and  $c > 0$  such that for all  $u, v \in \langle F \rangle$ ,*

$$d(u(x_*), v(x_*)) \geq c \cdot \text{dist}_{\text{word}}(u, v),$$

*where  $\text{dist}_{\text{word}}$  is the usual combinatorial distance on the free semigroup.*

*Then, for every  $k$  and every  $\varepsilon < c/2$ ,*

$$N(B(k, F), d_\infty, \varepsilon) \geq r^k.$$

*Proof.* Since  $u \mapsto u(x_*)$  is  $c$ -Lipschitz from below with respect to word distance, distinct words of the same length are  $\geq c$ -separated at  $x_*$ . Apply Lemma 7 with  $P = \{x_*\}$ . The finiteness of  $d_\infty$  guarantees that covering numbers are meaningful.  $\square$

1242 *Example 9* (Free Group with Bi-Invariant Metric (Non-compact, Free, Isometric)). Let  $\mathcal{X} = G = F_r$   
 1243 ( $r \geq 2$ ). Let

$$1244 \quad S := \{xa_i x^{-1}, xa_i^{-1} x^{-1} \mid x \in F_r, i = 1, \dots, r\}$$

1245 and define the *conjugacy-invariant* word metric

$$1247 \quad d_S(g, h) := |g^{-1}h|_S,$$

1248 the shortest length in the alphabet  $S$ . This metric is *bi-invariant*, hence all left translations  $L_g$  are  
 1249 isometries and, crucially,

$$1251 \quad d_\infty(L_g, L_h) = \sup_x d_S(gx, hx) = \sup_x d_S(e, x^{-1}g^{-1}hx) = d_S(e, g^{-1}h) < \infty$$

1252 (the supremum is independent of  $x$ ).

1253 Let  $F = \{L_{a_1}, \dots, L_{a_r}\}$ . The positive semigroup is free. With  $x_* = e$ , for words  $u \neq v$  of the  
 1254 same length,

$$1255 \quad d_S(L_u(x_*), L_v(x_*)) = d_S(u, v) = |u^{-1}v|_S \geq 1,$$

1256 so Condition E1' applies with  $c = 1$ : for every  $\varepsilon < 1/2$ ,

$$1260 \quad N(B(k, F), d_\infty, \varepsilon) \geq r^k.$$

1261 Unlike Example 8,  $d_\infty$  is **finite** for all pairs.

#### 1264 E.4.2 E2: PING-PONG WITH NON-CONTRACTION $\implies$ EXPONENTIAL LOWERBOUNDS IN $k$

1265 **Condition E2** (Ping-Pong). Let  $F = \{f_1, \dots, f_r\} \subset C(\mathcal{X}, \mathcal{X})$  with  $r \geq 2$ . Suppose there are  
 1266 pairwise separated non-empty sets (chambers)  $U_1, \dots, U_r \subset \mathcal{X}$ , points (anchors)  $a_1, \dots, a_r \in \mathcal{X}$ ,  
 1267 and constants  $\Delta, \delta_0 > 0$  such that:

- 1269 • (PP1: Pairwise separation)  $\text{dist}(U_i, U_j) \geq \Delta$  for all  $i \neq j$ .
- 1271 • (PP2: Expansion on own domain) For each  $i$ ,  $f_i$  maps  $U_i$  onto a set that meets every  $U_j$ .
- 1272 • (PP3: Reset off domain) On  $\mathcal{X} \setminus U_i$ ,  $f_i$  stays within  $\delta_0$  of the constant  $a_i$ .
- 1273 • (PP4: Anchors separated)  $\min_{i \neq j} d(a_i, a_j) > 0$ .

1274 Then, there exists  $\delta > 0$  (depending only on  $\Delta, \delta_0, \{a_i\}$ ) such that for every  $k$  and every  $\varepsilon < \delta/2$ ,

$$1277 \quad N(B(k, F), d_\infty, \varepsilon) \geq r^k.$$

1278 Moreover, the map  $u \mapsto f_u$  is injective (freeness).

1281 *Proof.* For a word  $u = i_k \cdots i_1$  construct  $x_u \in U_{i_1}$  by backwards induction so that

$$1283 \quad f_{i_j}(x_{i_j \cdots i_1}) \in U_{i_{j+1}} \quad (1 \leq j \leq k-1),$$

1284 which is possible by (PP2). Now fix distinct words  $u = i_k \cdots i_1$  and  $v = j_k \cdots j_1$  of the same length  
 1285 and let  $t$  be the rightmost index with  $i_t \neq j_t$ . Then  $f_{i_{t-1}} \cdots f_{i_1}(x_u) \in U_{i_t}$ , but  $f_{j_{t-1}} \cdots f_{j_1}(x_u) \notin$   
 1286  $U_{i_t}$ . Applying  $f_{i_t}$  vs.  $f_{j_t}$  at that step, one output is close to  $a_{j_t}$  (by reset), the other remains in a set  
 1287 well separated by  $\Delta$ . Hence

$$1288 \quad d(f_u(x_u), f_v(x_u)) \geq \delta$$

1289 for some  $\delta > 0$  determined by  $\Delta, \delta_0, \{a_i\}$ . Distinct words are therefore distinct maps (freeness),  
 1290 and Lemma 7 with  $P = \{x_u \mid |u| = k\}$  (or even a single carefully chosen probe) gives the claimed  
 1291 lower bound, uniform in  $k$ .  $\square$

1292 *Remark 6.* The assumptions formalize the common “expand inside the chamber; collapse to an  
 1293 anchor outside” motif. Uniform expansion  $\text{coLip} > 1$  is needed; only local expansion and a reset  
 1294 suffice.

1296 *Example 10 (Piecewise-Linear Ping-Pong (Compact, Expansive)).* Let  $\mathcal{X} = [0, 1]$  with Euclidean  
 1297 distance. Fix  $\eta \in (0, 1/12)$  and set

$$1299 \quad I_0 = \left[0, \frac{1}{3} - \eta\right], \quad I_1 = \left[\frac{2}{3} + \eta, 1\right], \quad J_0 = \left[\frac{1}{3} - \eta, \frac{1}{3} + \eta\right], \quad J_1 = \left[\frac{2}{3} - \eta, \frac{2}{3} + \eta\right].$$

1301 Let  $U_0 := I_0 \cup J_0$  and  $U_1 := I_1 \cup J_1$ . Define anchors  $a_0 = \frac{1}{6}$ ,  $a_1 = \frac{5}{6}$ . Define  $f_0, f_1 \in$   
 1302  $C([0, 1], [0, 1])$  by

$$1303 \quad f_0(x) = \begin{cases} 3x & x \in I_0 \\ \text{linear interpolation} & x \in J_0 \\ a_0 & x \in [\frac{1}{3} + \eta, 1] \end{cases}, \quad f_1(x) = \begin{cases} 3x - 2 & x \in I_1 \\ \text{linear interpolation} & x \in J_1 \\ a_1 & x \in [0, \frac{2}{3} - \eta]. \end{cases}$$

1307 Then:

- 1309 • On its own chamber  $U_i$ ,  $f_i$  is (piecewise) expansive with slope  $3 > 1$  on  $I_i$ .
- 1310 • Off  $U_i$ ,  $f_i$  is constant (reset) up to a thin collar.
- 1311 • Because the images  $f_i(U_i)$  cover long subintervals  $[0, 1 - 3\eta]$  and  $[3\eta, 1]$ , each meets both  
 1312  $U_0$  and  $U_1$  when  $\eta < 1/6$ .
- 1313 • The chambers are separated by  $\text{dist}(U_0, U_1) = \Delta = \frac{1}{3} - 2\eta > 0$ , and  $|a_0 - a_1| = 2/3$ .

1316 The ping-pong assumptions (PP1)–(PP4) hold, so Condition E2 yields: with  $\delta := \min\{\Delta, |a_0 -$   
 1317  $a_1|\} = \frac{1}{3} - 2\eta$  and any  $\varepsilon < \delta/2$ ,

$$1319 \quad N(B(k, \{f_0, f_1\}), d_\infty, \varepsilon) \geq 2^k.$$

1321 *Example 11 (Ping-Pong over Subshift (Compact, Expansive)).* Fix an alphabet  $\mathcal{A} = [r] \cup \{\bullet\}$  with  
 1322  $r \geq 2$  active symbols  $[r]$  and one extra padding symbol  $\bullet$ . Let

$$1323 \quad \mathcal{X} = \mathcal{A}^{\mathbb{N}} = \{x = (x_0, x_1, x_2, \dots) \mid x_j \in \mathcal{A}\}$$

1325 be the (one-sided) full shift. Equip  $\mathcal{X}$  with the standard ultrametric  $d_\theta$  for some fixed  $\theta \in (0, 1)$ :

$$1326 \quad d_\theta(x, y) = \begin{cases} 0 & x = y \\ \theta^{\min\{n \geq 0 \mid x_n \neq y_n\}} & x \neq y \end{cases}$$

1329 Then  $(\mathcal{X}, d_\theta)$  is compact, totally disconnected, and the left shift  $\sigma(x)_n = x_{n+1}$  is  $L$ -Lipschitz with  
 1330  $L = \theta^{-1} > 1$  (hence expansive).

1331 For each  $a \in \mathcal{A}$ , write the clopen 1-cylinder  $\llbracket a \rrbracket = \{x \in \mathcal{X} \mid x_0 = a\}$ , and the anchor sequence  
 1332  $\bar{a} = (a, a, a, \dots)$ .

1334 Define  $F = \{f_1, \dots, f_r\} \subset C(\mathcal{X}, \mathcal{X})$  by

$$1335 \quad f_a(x) = \begin{cases} \sigma(x) & x \in \llbracket a \rrbracket \\ \bar{a} & x \notin \llbracket a \rrbracket \end{cases} \quad a \in [r].$$

1339 Because  $\llbracket a \rrbracket$  is clopen, each  $f_a$  is continuous; on  $\llbracket a \rrbracket$  it is  $\sigma$  (Lipschitz constant  $\theta^{-1} > 1$ ), and on  
 1340  $\mathcal{X} \setminus \llbracket a \rrbracket$  it is constant. Thus each  $f_a$  is expansive (uniform Lipschitz with constant  $\theta^{-1} > 1$ ).

1341 Let  $U_a := \llbracket a \rrbracket$  (clopen chambers),  $a \in [r]$ . Then:

- 1343 • Pairwise separation. If  $a \neq b$ , then for any  $x \in U_a, y \in U_b$ ,  $\min d_\theta(x, y) = \theta^0 = 1$ . Hence  
 1344  $\text{dist}(U_a, U_b) = 1$ .
- 1345 • Expansion/coverage on own domain. For each  $a$ ,  $f_a|_{U_a} = \sigma$  maps  $U_a$  bijectively onto  $\mathcal{X}$   
 1346 (surjective and expanding).
- 1347 • Reset off domain. On  $\mathcal{X} \setminus U_a$ ,  $f_a = \text{const}_{\bar{a}}$  (exact reset; the “reset diameter” is 0).
- 1348 • Separated anchors.  $d_\theta(\bar{a}, \bar{b}) = 1$  for  $a \neq b$ .

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

These are exactly the ping-pong assumptions (PP1)–(PP4) in Condition E2 with

$$\Delta = \min_{a \neq b} \text{dist}(U_a, U_b) = 1, \quad \delta_0 = 0, \quad \min_{a \neq b} d_\theta(\bar{a}, \bar{b}) = 1.$$

Hence our separation scale can be taken as  $\delta = 1$ .

Let a word  $u = a_k \cdots a_1 \in [r]^k$ . Define the *tail-padded probe*

$$x_u := (a_1, a_2, \dots, a_k, \underbrace{\bullet, \bullet, \bullet, \dots}_{\text{all } \bullet}) \in \mathcal{X}.$$

Then, by construction,

$$f_u(x_u) = \sigma^k(x_u) = \bar{\bullet},$$

because each of the first  $k$  steps sees the correct chamber and applies  $\sigma$ , peeling off the  $k$ -letter prefix and revealing the all- $\bullet$  tail.

If  $v = b_k \cdots b_1 \neq u$ , let  $t$  be the *rightmost* index with  $a_t \neq b_t$ . When applying  $f_v$  to  $x_u$ , the first  $t-1$  letters match and act as shifts; at step  $t$  we apply  $f_{b_t}$  to a sequence whose 0-coordinate is  $a_t \neq b_t$ , so  $f_{b_t}$  resets to the anchor  $\bar{b}_t$ . From then on, the state never contains  $\bullet$  at the 0-coordinate (subsequent resets only use anchors  $\bar{b}_s$  with  $b_s \in [r]$ , and  $\sigma$  preserves 0-coordinate  $b_s$  on the constant sequence  $\bar{b}_s$ ). Consequently

$$f_v(x_u) \neq \bar{\bullet}.$$

Therefore

$$d_\theta(f_u(x_u), f_v(x_u)) = d_\theta(\bar{\bullet}, f_v(x_u)) = 1.$$

This shows *freeness*: distinct words  $u \neq v$  define distinct maps  $f_u \neq f_v$ .

Fix  $k \in \mathbb{N}$  and consider the finite probe set

$$P_k := \{x_u \mid u \in [r]^k\} \subset \mathcal{X},$$

of cardinality  $r^k$ . For any distinct  $u, v$  we have just seen that at the coordinate  $x_u$ ,

$$d_\theta(f_u(x_u), f_v(x_u)) = 1 \geq \delta.$$

Thus the  $r^k$  vectors  $E_{P_k}(f_u)|_{|u|=k} \subset \mathcal{X}^{P_k}$  are pairwise 1-separated in the max metric. By Lemma 7, for every  $\varepsilon < \frac{1}{2}$ ,

$$N(B(k, F), d_\infty, \varepsilon) \geq r^k.$$

#### E.4.3 E3. UNIFORM EXPANSION $\implies$ SUPER-/DOUBLE-EXPONENTIAL LOWERBOUNDS IN $k$

**Condition E3** (Uniformly Expansive Semigroup Grows Super-/Double Exponentially). *Assume:*

1. (*reset map*) there exist a compact  $K \subset \mathcal{X}$  and  $r \in B(C_r, F)$  and  $r(\mathcal{X}) \subset K$ .
2. (*uniform expanding maps*) there exists  $A \in B(C_A, F)$  and

$$\text{coLip}_K(A) \geq \lambda > 1, \quad \text{Lip}_K(A) \leq \Lambda < \infty.$$

3. (*writers*) there exists a compact set  $G \subset B(C_G, F) \cap C(K, K)$  with  $\sup_{g \in G} \text{Lip}_K(g) \leq 1$ .

Fix  $\varepsilon > 0$  and  $k \in \mathbb{N}$ . Consider words of the form

$$W_k = \{w = A \circ g_k \circ A \circ g_{k-1} \circ \cdots \circ A \circ g_1 \circ r \mid g_j \in G\}.$$

Then  $W_k$  is a subset of  $B(Ck + C_r, F)$  with  $C = C_A + C_G$ , and

$$N(W_k, d_\infty, \varepsilon) \geq N(B(Ck + O(1), F), d_\infty, \varepsilon) \geq \prod_{j=1}^k N\left(G, \frac{\varepsilon}{\lambda^{k-j+1}}\right).$$

1404 *Proof.* Let  $\delta_j := \varepsilon/\lambda^{k-j+1}$ . Choose for each  $j$  a  $\delta_j$ -net  $\mathcal{N}_j \subset G$  of maximal cardinality  $N(G, \delta_j)$ .  
 1405 Consider the set of words  $w$  obtained by picking  $g_j \in \mathcal{N}_j$  independently. If two such words  $w, w'$   
 1406 first differ at position  $j$  (counting from the right), then for all  $x \in \mathcal{X}$ ,

$$1407 \quad d((g_j \circ r)(x), (g'_j \circ r)(x)) \geq \delta_j$$

1408  
 1409 after evaluating inside  $K$ , and subsequent composition with  $A$  at least multiplies distances by  
 1410  $\lambda^{k-j+1}$ . The 1-Lipschitz  $g_\ell$  does not increase distances, so overall

$$1411 \quad d_\infty(w, w') \geq \lambda^{k-j+1} \delta_j = \varepsilon.$$

1412  
 1413 Thus we obtain a packing of size  $\prod_j N(G, \delta_j)$  at scale  $\varepsilon$ , and the claim follows from  $N \geq M$ .  $\square$

1414  
 1415 Condition E3 converts information about the layer-wise covering numbers  $N(G, \cdot)$  into a covering  
 1416 lower bound for deep compositions. Two regimes are of special interest.

1417 **Condition E3a** (Super-exponential Growth for Finite-dimensional Families). *Suppose  $G$  is con-*  
 1418 *tained in a  $D$ -dimensional  $C^0$  submanifold of  $C(K, K)$  and  $\sup_{g \in G} \text{Lip}_K(g) \leq 1$ . Then there is a*  
 1419 *constant  $c_1 > 0$  such that for all sufficiently small  $\delta > 0$ ,*

$$1420 \quad N(G, d_\infty, \delta) \geq c_1 \delta^{-D}.$$

1421  
 1422 *Consequently, for every fixed  $\varepsilon > 0$ ,*

$$1423 \quad \log N(B(Ck+O(1), F), d_\infty, \varepsilon)$$

$$1424 \quad \geq \sum_{j=1}^k (D(k-j+1) \log \lambda + D \log(1/\varepsilon) - c_2)$$

$$1425 \quad = \frac{D}{2} (\log \lambda) k(k+1) + Dk \log(1/\varepsilon) + O(k).$$

1426  
 1427  
 1428  
 1429 *In particular,*

$$1430 \quad N(B(Ck+O(1), F), d_\infty, \varepsilon) \geq \lambda^{cDk^2} \varepsilon^{-Dk} = \exp(\Theta(k^2 + k \log(1/\varepsilon)))$$

1431  
 1432 *for some  $c > 0$ : growth is super-exponential in  $k$ .*

1433  
 1434  
 1435 *Proof.* The entropy estimate  $N(G, \delta) \geq c_1 \delta^{-D}$  is standard for compact  $D$ -dimensional submani-  
 1436 folds of a Banach space under the  $C^0$ -norm. Apply Condition E3 and sum the geometric-arithmetic  
 1437 progression.  $\square$

1438 **Condition E3b** (Double-exponential Growth for Hölder Balls). *Suppose  $\mathcal{X}$  is a  $d$ -dimensional man-*  
 1439 *ifold with bi-Lipschitz coordinates, and  $K \subset \mathcal{X}$  be compact. Let*

$$1440 \quad G = \{g \in C(K, K) \mid g = \text{id}_K + u, \|u\|_\infty \leq 1, [u]_{C^\alpha} \leq 1\}$$

1441  
 1442 *with  $\alpha \in (0, 1]$ . Then there exists  $c > 0$  such that for  $\delta \in (0, \delta_0]$ ,*

$$1443 \quad \log N(G, d_\infty, \delta) \geq c \delta^{-d/\alpha}.$$

1444  
 1445 *Hence for every fixed  $\varepsilon > 0$ ,*

$$1446 \quad \log N(B(Ck+O(1), F), d_\infty, \varepsilon) \geq \sum_{j=1}^k c \left( \frac{\lambda^{k-j+1}}{\varepsilon} \right)^{d/\alpha} \geq c' \varepsilon^{-\frac{d}{\alpha}} \lambda^{\frac{d}{\alpha} k},$$

1447  
 1448  
 1449 *and therefore*

$$1450 \quad N(B(Ck+O(1), F), d_\infty, \varepsilon) \geq \exp \left( C \lambda^{\frac{d}{\alpha} k} \varepsilon^{-\lambda^{\frac{d}{\alpha}}} \right) = \exp(\exp(\Theta(k + \log(1/\varepsilon)))).$$

1451  
 1452  
 1453 *This is double-exponential growth in  $k$ . (the exponent itself grows exponentially in  $k$ ).*

1454  
 1455  
 1456 *Proof.* The entropy bound for Hölder balls in the sup norm is classical and follows by dyadic par-  
 1457 titioning of  $K$  at mesh size  $\delta^{1/\alpha}$ , prescribing values on the grid with resolution  $\delta$ , and using the  
 Hölder constraint to extend; Then apply Condition E3 and sum a geometric series.  $\square$

1458 *Example 12* (Hölder Balls on Euclidean Space). Let  $\mathcal{X} = \mathbb{R}^d$  with the bounded metric  $d_b$ . Fix  
 1459  $K = B(0, 1)$ .

- 1461 • (Reset)  $r(x) := x / \max\{1, |x|\}$  maps  $\mathbb{R}^d$  into  $K$ , continuously.
- 1462
- 1463 • (Expansion) pick  $\lambda > 1$  and define in polar coordinates an  $A$  that equals  $x \mapsto \lambda x$  on  
 1464  $B(0, 1/2)$ , smoothly interpolates on the annulus to the identity, and equals the identity  
 1465 outside  $B(0, 1)$ . Then  $\text{coLip}_K(A) \geq \lambda' > 1$  for some  $\lambda'$  close to  $\lambda$ .
- 1466
- 1467 • (Writers) take  $G$  to be the  $C^\alpha$ -Hölder ball of small perturbations supported in  $B(0, 1/2)$ ,  
 1468 with  $\|u\|_\infty \leq 1$  and  $[u]_{C^\alpha} \leq 1$ , and define  $g = \text{id} + u$ . Each  $g$  is 1-Lipschitz on  $K$  after  
 1469 shrinking the ball if needed.

1470 All three ingredients lie in a compact  $F := \overline{\{r, A\} \cup G}$ . Condition E3b applies and yields a double-  
 1471 exponential lower bound.

#### 1473 E.4.4 EXPONENTIAL UPPERBOUNDS

1474 Let  $L := \sup_{f \in F} \text{Lip}(f)$ . A standard chaining argument gives

$$1475 \quad N(B(k, F), d_\infty, \varepsilon) \leq \prod_{j=0}^{k-1} N\left(F, \frac{\varepsilon}{L^j}\right).$$

1476 Indeed, approximate the rightmost factor in a word within  $\varepsilon/L^{k-1}$ , the next within  $\varepsilon/L^{k-2}$ , and  
 1477 so on; the Lipschitz constants propagate the errors to at most  $\varepsilon$  at the output. This complements  
 1478 Condition E3, which provides a matching-flavor lower bound under expansion.

## 1484 F PROOFS FOR BALANCING BIAS-VARIANCE TRADE-OFF IN DEPTH

1485 Throughout, we minimize the upper bound

$$1486 \quad \text{gen}(k, n) \lesssim \text{bias}(k) + \text{var}(k, n),$$

1487 treat  $k$  as a positive real (round to the nearest integer at the end), and use the standard heuristic  
 1488 that—because  $\text{bias}(k)$  is decreasing in  $k$  while  $\text{var}(k, n)$  is increasing—the minimizer occurs where  
 1489 the two terms are of the same order:

$$1490 \quad \text{bias}(k^*) \asymp \text{var}(k^*, n).$$

1491 Solving that equation gives  $k^*$ ; plugging back yields the minimized rate. (If a term does not cross,  
 1492 the optimum is at a boundary, but in all four regimes below they do cross for large  $n$ .)

### 1498 F.1 EP (EXP-DECAY BIAS, POLY-GROWTH VARIANCE)

$$1499 \quad \text{bias}(k) = e^{-\alpha k}, \quad \text{var}(k, n) = n^{-1/2} k^{\gamma/2}.$$

1500 Balance:

$$1501 \quad e^{-\alpha k} \asymp n^{-1/2} k^{\gamma/2} \iff \alpha k = \frac{1}{2} \log n - \frac{\gamma}{2} \log k.$$

1502 As  $n \rightarrow \infty$ ,  $\log k \ll \log n$ , so an asymptotic solution is

$$1503 \quad k^* = \frac{1}{2\alpha} \left( \log n - \gamma \log \log n + O(1) \right).$$

1504 Plugging back (either term) gives

$$1505 \quad \text{gen}(k^*, n) \asymp n^{-1/2} (\log n)^{\gamma/2} \quad (\text{more precisely } \approx (2\alpha)^{-\gamma/2} n^{-1/2} (\log n)^{\gamma/2} \text{ up to a factor } \asymp 1).$$

1512 F.2 EL (EXP-DECAY BIAS, LOG-GROWTH VARIANCE)

1513  
1514 
$$\text{bias}(k) = e^{-\alpha k}, \quad \text{var}(k, n) = \sqrt{\log k/n}.$$

1515  
1516 **Balance:**  
1517 
$$e^{-\alpha k} \asymp \sqrt{\log k/n} \iff \alpha k = \frac{1}{2} \log n - \frac{1}{2} \log \log k.$$

1518  
1519 **Hence**  
1520 
$$k^* = \frac{1}{2\alpha} \left( \log n - \log \log \log n + o(1) \right),$$

1521  
1522 **and**  
1523 
$$\text{gen}(k^*, n) \asymp \sqrt{\log \log n/n}.$$

1524  
1525 F.3 PP (POLY-DECAY BIAS, POLY-GROWTH VARIANCE)

1526  
1527 
$$\text{bias}(k) = k^{-\beta}, \quad \text{var}(k, n) = n^{-1/2} k^{\gamma/2}.$$

1528  
1529 **Balance:**  
1530 
$$k^{-\beta} \asymp n^{-1/2} k^{\gamma/2} \iff k^{\beta+\gamma/2} \asymp n^{1/2}.$$

1531  
1532 **Thus**  
1533 
$$k^* \asymp n^{1/(2\beta+\gamma)}, \quad \text{gen}(k^*, n) \asymp n^{-\beta/(2\beta+\gamma)}.$$

1534 F.4 PL (POLY-DECAY BIAS, LOG-GROWTH VARIANCE)

1535  
1536 
$$\text{bias}(k) = k^{-\beta}, \quad \text{var}(k, n) = \sqrt{\log k/n}.$$

1537  
1538 **Balance (square both sides):**

1539  
1540 
$$k^{-2\beta} \asymp \frac{\log k}{n} \iff k^{2\beta} \log k \asymp n.$$

1541  
1542 Let  $k = e^t$ . Then  $te^{2\beta t} \asymp n$ , so

1543  
1544 
$$2\beta t = W(2\beta n) \implies k^* = \exp\left(\frac{1}{2\beta} W(2\beta n)\right) = \left(\frac{2\beta n}{W(2\beta n)}\right)^{1/(2\beta)},$$

1545  
1546 where  $W$  is the Lambert  $W$  function. Consequently,

1547  
1548 
$$\text{gen}(k^*, n) \asymp \sqrt{\frac{W(2\beta n)}{2\beta n}} \sim \sqrt{\frac{\log n}{2\beta n}} \quad (\text{since } W(x) \sim \log x).$$

1551  
1552 G NOTES ON RADEMACHER COMPLEXITY AND GENERALIZATION ERROR  
1553 BOUNDS

1554  
1555 We refer to (Martin J. Wainwright, 2019; Mohri et al., 2018) for more details on Rademacher com-  
1556 plexity and its application to generalization error bounds.

1557  
1558 G.1 RADEMACHER COMPLEXITY

1559  
1560 **Sample space** Let  $(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$  be a measurable space that serves as the sample space, and let  
1561  $P$  be an unknown Borel probability measure on  $\mathcal{B}(\mathcal{Z})$ . We observe an i.i.d. sample  $S =$   
1562  $(Z_1, \dots, Z_n) \stackrel{\text{iid}}{\sim} P.$

1563  
1564 **Hypothesis class** A hypothesis is a measurable function  $h : \mathcal{Z} \rightarrow \mathbb{R}$ . The hypothesis class is a  
1565 separable set  $\mathcal{H} \subset \{h : \mathcal{Z} \rightarrow \mathbb{R}\}$ . Here, separability is assumed to prevent pathologies where the  
supremum over uncountable  $\mathcal{H}$  is non-measurable.

1566 **Rademacher complexities** Introduce an independent sequence of Rademacher variables  $\sigma =$   
 1567  $(\sigma_1, \dots, \sigma_n)$  with  $\Pr(\sigma_i = 1) = \Pr(\sigma_i = -1) = \frac{1}{2}$ . For the fixed sample  $S$  the empirical  
 1568 Rademacher complexity of  $\mathcal{H}$  is

$$1569 \hat{\mathfrak{R}}_S(\mathcal{H}) = \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(Z_i) \right].$$

1572 Averaging this quantity over all samples of size  $n$  drawn from  $P$  gives the (distribution-dependent)  
 1573 Rademacher complexity

$$1574 \mathfrak{R}_n(\mathcal{H}) = \mathbb{E}_{S \stackrel{\text{iid}}{\sim} P} [\hat{\mathfrak{R}}_S(\mathcal{H})].$$

### 1576 Properties

- 1578 • **Symmetrization:**  $\mathbb{E}_S \left[ \sup_{h \in \mathcal{H}} (P - \hat{P})h \right] \leq 2\mathfrak{R}_n(\mathcal{H})$  (Bartlett & Mendelson, 2002)
- 1581 • **Ledoux–Talagrand contraction lemma:** if  $\psi : \mathbb{R} \rightarrow \mathbb{R}$  is  $\beta$ -Lipschitz and  $\psi(0) = 0$  then  
 1582  $\hat{\mathfrak{R}}_S(\psi \circ \mathcal{H}) \leq \beta \hat{\mathfrak{R}}_S(\mathcal{H})$  (Mohri et al., 2018, Lemma 5.7)
- 1583 • **Monotonicity and convexity preservation:** if  $\mathcal{H} \subset \mathcal{H}'$  or  $\text{conv}(\mathcal{H})$  is taken, complexity  
 1584 cannot increase (Bartlett & Mendelson, 2002)

## 1585 G.2 UNIFORM DEVIATION FOR LIPSCHITZ, BOUNDED LOSS

1587 **Sample space** Let  $\mathcal{X}$  be a measurable space. Take an input-output sample space  $\mathcal{X} \times \mathcal{Y}$  with  
 1588  $\mathcal{Y} = \mathbb{R}$ .

1590 **Hypothesis class** Let  $\mathcal{H} \subset \{h : \mathcal{X} \rightarrow \mathbb{R}\}$  be any class of measurable functions.

1592 **Loss function** Fix a measurable loss  $\ell : \mathbb{R} \rightarrow [0, b]$  that is  $\beta_\ell$ -Lipschitz in its *first* argument:

$$1593 |\ell(a, y) - \ell(a', y)| \leq \beta_\ell |a - a'| \quad (a, a' \in \mathbb{R}, y \in \mathbb{R}).$$

1594 For brevity write  $\ell \cdot h(x, y) = \ell(h(x), y)$  and  $\ell \cdot \mathcal{H} = \{\ell \cdot h : h \in \mathcal{H}\} \subset [0, b]^{\mathcal{X} \times \mathcal{Y}}$ .

1596 **Risks** Given a sample  $S = ((X_i, Y_i))_{i=1}^n$ , define empirical and population risks

$$1598 \hat{L}_n[h] = \frac{1}{n} \sum_{i=1}^n \ell(h(X_i), Y_i), \quad L[h] = \mathbb{E}_{(X, Y) \sim P} [\ell(h(X), Y)].$$

1601 **Theorem 5** (uniform deviation theorem (in expectation)).

$$1602 \mathbb{E}_S \left[ \sup_{h \in \mathcal{H}} |L[h] - \hat{L}_n[h]| \right] \leq 2\beta_\ell \mathfrak{R}_n(\mathcal{H}). \quad (33)$$

1605 *Proof.* The symmetrization identity applied to  $\ell \cdot \mathcal{H}$  gives

$$1606 \mathbb{E}_S \left[ \sup_{h \in \mathcal{H}} (L[h] - \hat{L}_n[h]) \right] \leq 2\mathfrak{R}_n(\ell \cdot \mathcal{H}). \quad (34)$$

1609 Using the contraction lemma with scale  $\beta_\ell$  and the fact  $\ell(a, y) - \ell(0, y) = \psi_y(a)$  has Lipschitz  
 1610 constant  $\beta_\ell$ , one gets

$$1611 \mathfrak{R}_n(\ell \cdot \mathcal{H}) \leq \beta_\ell \mathfrak{R}_n(\mathcal{H}). \quad (35)$$

1612 Combining Eqs. (34) and (35) yields the assertion.  $\square$

1614 **Theorem 6** (uniform deviation theorem (high probability)). *With probability at least  $1 - \delta$  over the*  
 1615 *i.i.d. draw  $S \sim P^n$ ,*

$$1616 \sup_{h \in \mathcal{H}} |L[h] - \hat{L}_n[h]| \leq 2\beta_\ell \hat{\mathfrak{R}}_S(\mathcal{H}) + b \sqrt{\frac{2 \log(1/\delta)}{n}} \quad (36)$$

$$1619 \sup_{h \in \mathcal{H}} |L[h] - \hat{L}_n[h]| \leq 2\beta_\ell \mathfrak{R}_n(\mathcal{H}) + b \sqrt{\frac{2 \log(1/\delta)}{n}} \quad (37)$$

1620 *Proof.* Define  $G(S) = \sup_{h \in \mathcal{H}} |L[h] - \hat{L}_n[h]|$ . Replacing one coordinate  $(X_j, Y_j)$  in  $S$  can change  
 1621 each summand by at most  $b/n$ , so  
 1622

$$1623 \quad |G(S) - G(S')| \leq \frac{b}{n} \quad \text{whenever } S, S' \text{ differ in exactly one entry.}$$

1624  
 1625 Hence  $G$  satisfies the bounded-difference condition and McDiarmid's inequality gives  
 1626

$$1627 \quad \Pr \left\{ G(S) \geq \mathbb{E}[G(S)] + \sqrt{\frac{2b^2 \log(1/\delta)}{n}} \right\} \leq \delta.$$

1628  
 1629 Inserting Eq. (33) and simplifying constants we obtain, with probability  $1 - \delta$ ,  
 1630

$$1631 \quad \sup_{h \in \mathcal{H}} |L[h] - \hat{L}_n[h]| \leq 2\beta_\ell \hat{\mathfrak{R}}_S(\mathcal{H}) + b\sqrt{\frac{2 \log(1/\delta)}{n}}. \quad (38)$$

1632  
 1633 Taking expectation and applying Jensen's inequality yields  
 1634

$$1635 \quad \sup_{h \in \mathcal{H}} |L[h] - \hat{L}_n[h]| \leq 2\beta_\ell \mathfrak{R}_n(\mathcal{H}) + b\sqrt{\frac{2 \log(1/\delta)}{n}}, \quad (39)$$

1636  
 1637 with the same probability. □  
 1638  
 1639

1640 The bound controls the worst-case *generalization gap* uniformly over  $\mathcal{H}$ . If  $\hat{\mathfrak{R}}_S(\mathcal{H})$  (data-dependent)  
 1641 or  $\mathfrak{R}_n(\mathcal{H})$  (distribution-dependent) is small—scaling, say, like  $O(1/\sqrt{n})$ —then empirical risk min-  
 1642 imization on  $S$  provably finds a hypothesis whose population risk is close to optimal.  
 1643

1644  $\hat{\mathfrak{R}}_S$  quantifies how well the function class can *correlate with random noise* on the sample. A low  
 1645 value implies strong regularization, yielding smaller sample complexity. The second term is a uni-  
 1646 versal concentration penalty that decays at the Monte-Carlo rate  $O(\sqrt{\log(1/\delta)/n})$ .  
 1647

### 1648 G.3 BOUNDING RADEMACHER COMPLEXITY

1649 Let  $S = (X_1, \dots, X_n)$  be an i.i.d. sample from  $P$ . Throughout this section every  $h \in \mathcal{H}$  is assumed  
 1650 to satisfy  $\|h\|_\infty \leq b$  for some constant  $b > 0$ .  
 1651

#### 1652 **Massart's finite class lemma**

1653 **Theorem 7** (Massart's finite class lemma). *Suppose  $\mathcal{H}$  is a finite family,  $|\mathcal{H}| < \infty$ . Then,*  
 1654

$$1655 \quad \hat{\mathfrak{R}}_S(\mathcal{H}) \leq b\sqrt{\frac{2 \log |\mathcal{H}|}{n}}. \quad (40)$$

1656  
 1657 We refer to (Mohri et al., 2018, Theorem 3.7) for the proof.  
 1658

1659 The bound depends logarithmically on the cardinality of the class, illustrating how modeling with a  
 1660 finite but exponentially large dictionary can still be statistically benign.  
 1661  
 1662

1663 **Covering numbers and Dudley's entropy integral** When  $\mathcal{H}$  is infinite a combinatorial bound  
 1664 like Eq. (40) is no longer adequate. A refined control is obtained by chaining the increments of the  
 1665 empirical process in the empirical  $L^2$ -metric  $\|h - h'\|_S = \left(\frac{1}{n} \sum_{i=1}^n (h(X_i) - h'(X_i))^2\right)^{1/2}$ .  
 1666

1667 For any  $\varepsilon > 0$  let  $N(\varepsilon, \mathcal{H}, \|\cdot\|_S)$  denote the minimal number of  $\|\cdot\|_S$ -balls of radius  $\varepsilon$  needed to  
 1668 cover  $\mathcal{H}$ .

1669 **Theorem 8** (Dudley's entropy integral inequality).

$$1670 \quad \hat{\mathfrak{R}}_S(\mathcal{H}) \leq \frac{12}{\sqrt{n}} \int_0^{\text{diam}(\mathcal{H}, \|\cdot\|_S)} \sqrt{\log N(\varepsilon, \mathcal{H}, \|\cdot\|_S)} d\varepsilon. \quad (41)$$

1671  
 1672 We refer to (Martin J. Wainwright, 2019, Theorem 5.22) for the proof.  
 1673

1674 **Consequences and examples** For Hölder-smooth functions. If  $\mathcal{H}$  is a unit ball of a Hölder class  
 1675 of order  $s > d/2$  on  $[0, 1]^d$ , classical approximation theory gives  $\log N(\varepsilon, \mathcal{H}, \|\cdot\|_S) \lesssim \varepsilon^{-d/s}$ .  
 1676 Inserting this into Eq. (41) yields

$$1677 \hat{\mathfrak{R}}_S(\mathcal{H}) \lesssim \frac{1}{\sqrt{n}} \int_0^1 \varepsilon^{-d/(2s)} d\varepsilon = \frac{C_{d,s}}{n^{s/(2s+d)}}.$$

1680 Thus Rademacher complexity reproduces the minimax rate  $n^{-2s/(2s+d)}$  for non-parametric regres-  
 1681 sion.

1682 For linear prediction in  $\mathbb{R}^p$ . With  $\mathcal{H} = \{x \mapsto \langle w, x \rangle : \|w\|_2 \leq 1\}$  and  $\|X_i\|_2 \leq 1$ , the covering  
 1683 number in  $\|\cdot\|_S$  is bounded by  $N(\varepsilon) \leq (3/\varepsilon)^p$ . Dudley's integral then recovers  $\hat{\mathfrak{R}}_S(\mathcal{H}) \leq \sqrt{p/n}$ ,  
 1684 i.e. the familiar parametric rate.

1685 Dudley's entropy integral Eq. (41) and Massart's finite-class lemma Eq. (40) together form a bridge  
 1686 between combinatorial and geometric measures of capacity: the first quantifies richness by counting  
 1687 distinguishable hypotheses at each scale, the second by raw cardinality. In practice one often proves  
 1688 (data-dependent) covering-number bounds and plugs them into Eq. (41); the resulting Rademacher  
 1689 complexity then feeds directly into high-probability risk bounds through the symmetrization and  
 1690 concentration arguments.

## 1692 H NOTES ON COVERING AND PACKING NUMBERS

### 1693 H.1 PSEUDO-METRIC SPACE

1694 **Definition 1** (pseudo-metric space). A pseudo-metric space  $(X, \rho)$  is a pair of set  $X$  and non-  
 1695 negative real-valued function  $\rho : X \times X \rightarrow \mathbb{R}_{\geq 0}$  (called pseudo-metric) satisfying the following  
 1696 three axioms:

- 1700 •  $\forall x, y \in X : \rho(x, y) = \rho(y, x)$
- 1701 •  $\forall x, y, z \in X : \rho(x, z) \leq \rho(x, y) + \rho(y, z)$
- 1702 •  $\forall x \in X : \rho(x, x) = 0$ .

1703 Note that it is a metric space if it further satisfies

- 1704 •  $\forall x, y \in X : \rho(x, y) = 0 \implies x = y$ .

1705 In other words, in a pseudo-metric space, two distinct points  $x, y \in X$  may achieve zero distance  
 1706  $\rho(x, y) = 0$ .

### 1707 H.2 COVERING AND PACKING NUMBERS

1708 Let  $(X, \rho)$  be a pseudo-metric space. Let  $B(x; \varepsilon) := \{y \in X \mid \rho(x, y) \leq \varepsilon\}$  denote the closed ball  
 1709 in  $X$  of radius  $\varepsilon$  centered at a point  $x \in X$ .

1710 **Definition 2.** A subset  $C \subset X$  is a  $\varepsilon$ -covering if for every  $x \in X$  there exists  $y \in C$  satisfying  
 1711  $\rho(x, y) \leq \varepsilon$ . In other words,

$$1712 X \subset \bigcup_{x \in C} B(x; \varepsilon).$$

1713 The minimal cardinality of  $\varepsilon$ -covering, i.e.

$$1714 N(X, \rho, \varepsilon) := \inf \{|C| \mid C \subset X \text{ is } \varepsilon\text{-covering}\}$$

1715 is called the covering number of  $X$ .

1716 **Definition 3.** A subset  $P \subset X$  is an  $\varepsilon$ -packing if for every  $x, y \in P$   $\rho(x, y) > \varepsilon$ . In other words,

$$1717 \forall x, y \in P, x \neq y \implies B(x; \varepsilon) \cap B(y; \varepsilon) = \emptyset$$

1718 The maximal cardinality of  $\varepsilon$ -packing, i.e.

$$1719 M(X, \rho, \varepsilon) := \sup \{|P| \mid P \subset X \text{ is } \varepsilon\text{-packing}\}$$

1720 is called the packing number of  $X$ .

1728 While the covering number involves a minimization problem, the packing number involves a maxi-  
 1729 mization problem. The following theorem relates the duality of two quantities.

1730 **Theorem 9.** *Let  $(X, \rho)$  be a pseudo-metric space. For any positive number  $\varepsilon > 0$ , we have*

$$1731 N(X, \rho, \varepsilon) \leq M(X, \rho, \varepsilon) \leq N(X, \rho, \varepsilon/2).$$

### 1733 H.3 COVERING NUMBER OF COMPOSITIONS

1734 Let  $(X, d_X), (Y, d_Y), (Z, d_Z)$  be pseudo-metric spaces. Let  $C(X, Y), C(Y, Z), C(X, Z)$  denote the  
 1735 complete metric spaces of continuous maps equipped with uniform metrics:

$$1736 d_{C(X,Y)}(f, f') := \sup_{x \in X} d_Y(f(x), f'(x)), \quad f, f' \in C(X, Y)$$

$$1737 d_{C(Y,Z)}(g, g') := \sup_{y \in Y} d_Z(g(y), g'(y)), \quad g, g' \in C(Y, Z)$$

$$1738 d_{C(X,Z)}(h, h') := \sup_{x \in X} d_Z(h(x), h'(x)), \quad h, h' \in C(X, Z)$$

1739 respectively. Given a subclasses of continuous maps

$$1740 F \subset C(X, Y), \quad G \subset C(Y, Z),$$

1741 define their composition class

$$1742 G \circ F := \{g \circ f \mid g \in G, f \in F\} \subset C(X, Z).$$

1743 **Lemma 8** (composition lemma). *Assume every  $g \in G$  is at most  $L$ -Lipschitz: There exists  $L > 0$   
 1744 for every  $g \in G$  such that*

$$1745 d_Z(g(y), g(y')) \leq L d_Y(y, y') \quad (\forall y, y' \in Y).$$

1746 Then for any positive numbers  $\varepsilon, \delta_G, \delta_F > 0$  satisfying  $\varepsilon = \delta_G + L\delta_F$ ,

$$1747 N(\varepsilon, G \circ F, d_{C(X,Z)}) \leq N(\delta_G, G, d_{C(Y,Z)}) N(\delta_F, F, d_{C(X,Y)}) \quad (42)$$

1748 Taking logs gives a chain rule of metric entropies:

$$1749 \log N(\varepsilon, G \circ F) \leq \log N(\delta_G, G) + \log N(\delta_F, F). \quad (43)$$

1750 *Proof.* Build nets for  $G$  and  $F$ . Choose

$$1751 G_{\text{net}} := \{g_1, \dots, g_M\}, \quad F_{\text{net}} := \{f_1, \dots, f_N\}$$

1752 such that

$$1753 d_{C(Y,Z)}(g, g_m) \leq \delta_G, \quad d_{C(X,Y)}(f, f_n) \leq \delta_F.$$

1754 Bound the composition error. For any  $g \in G, f \in F$  pick the nearest net points  $g_m, f_n$  and use  
 1755 triangle inequalities:

$$1756 d_{C(X,Z)}(g \circ f, g_m \circ f_n) \leq d_{C(X,Z)}(g \circ f, g \circ f_n) + d_{C(X,Z)}(g \circ f_n, g_m \circ f_n)$$

$$1757 \leq L d_{C(X,Y)}(f, f_n) + d_{C(Y,Z)}(g, g_m)$$

$$1758 \leq L\delta_F + \delta_G = \varepsilon.$$

1759 Thus the composite error is at most  $\varepsilon$ .

1760 Counting. The set  $G_{\text{net}} \circ F_{\text{net}}$  is an  $\varepsilon$ -net of size  $MN$ , yielding the stated bound.  $\square$

## 1761 I NOTES ON TOTALLY BOUNDEDNESS OF FUNCTION SETS AND 1762 ARZELÀ–ASCOLI PRINCIPLE

1763 Here we list sufficient conditions for totally boundedness of a set  $H \subset C(X)$  in the uniform norm  
 1764  $\|f\|_\infty = \sup_{x \in X} |f(x)|$ . Throughout,  $X$  is a topological space and

$$1765 C(X) := \{f : X \rightarrow \mathbb{R} \mid f \text{ is continuous and bounded}\}, \quad \|f\|_\infty = \sup_{x \in X} |f(x)|.$$

1782 **Total boundedness.** A subset  $H \subset C(X)$  is totally bounded if for every  $\varepsilon > 0$  there exist finitely  
1783 many points  $f_1, \dots, f_N \in C(X)$  such that  $H \subset \bigcup_{i=1}^N B_\infty(f_i, \varepsilon)$ , where  $B_\infty(f, \varepsilon) = \{g \mid \|g -$   
1784  $f\|_\infty < \varepsilon\}$ .  
1785

1786 **Equicontinuity.** A family  $H \subset C(X)$  is *equicontinuous at*  $x \in X$  if for every  $\varepsilon > 0$  there is a  
1787 neighbourhood  $U$  of  $x$  with  $|f(y) - f(x)| < \varepsilon$  for all  $y \in U$  and  $f \in H$ . It is *equicontinuous on*  $X$   
1788 if this holds at every point.  
1789

1790 **Uniform boundedness.**  $H$  is uniformly bounded if  $\sup_{f \in H} \|f\|_\infty < \infty$ .  
1791

1792 **Compact exhaustion.** For a locally compact Hausdorff space we write  $C_0(X) \subset C(X)$  for the  
1793 subspace of functions *vanishing at infinity*, i.e.  $\forall \varepsilon > 0 \exists$  compact  $K \subset X : |f(x)| < \varepsilon$  for all  
1794  $x \notin K$ .  
1795

## 1796 I.1 A GENERIC ARZELÀ–ASCOLI PRINCIPLE

1797 **Theorem 10** (Abstract Arzelà–Ascoli). *Let  $X$  be a Tychonoff space (completely regular Hausdorff).*  
1798 *A set  $H \subset C(X)$  is totally bounded in  $\|\cdot\|_\infty$  provided that*  
1799

- 1800 1.  $H$  is uniformly bounded, and
- 1801 2. for every compact  $K \subset X$  the restriction  $H|_K := \{f|_K \mid f \in H\}$  is equicontinuous on  
1802  $K$ .  
1803

1804 The proof combines (i) uniform boundedness to control the range and (ii) the classical Arzelà–Ascoli  
1805 theorem on each compact  $K$ ; a diagonal argument then yields a finite  $\varepsilon$ -net on  $X$ .  
1806

## 1807 I.2 CASE-BY-CASE CONDITIONS

### 1809 I.2.1 $X$ COMPACT HAUSDORFF

1810 Because  $X$  itself is compact, condition (2) above is just *global* equicontinuity. Hence

1811 **Corollary 1** (Classical Arzelà–Ascoli). *If  $X$  is compact Hausdorff and  $H \subset C(X)$  is uniformly*  
1812 *bounded and equicontinuous on  $X$ , then  $H$  is totally bounded (and its closure is compact in  $C(X)$ ).*

1813 (For real-valued functions “pointwise relative compactness” in the usual statements reduces to uni-  
1814 form boundedness.)  
1815

### 1817 I.2.2 $X$ LOCALLY COMPACT HAUSDORFF

1818 Now  $X$  need not be compact. Total boundedness in  $C(X)$  fails unless one controls behavior at  
1819 infinity. A convenient framework is  $C_0(X)$ .  
1820

1821 **Theorem 11** (Arzelà–Ascoli for  $C_0(X)$ ). *Let  $X$  be locally compact Hausdorff and let  $H \subset C_0(X)$ .*  
1822 *Suppose*

- 1823 1. *Uniform boundedness:*  $\sup_{f \in H} \|f\|_\infty < \infty$ .
- 1824 2. *Local equicontinuity:* for every compact  $K \subset X$  the family  $H|_K$  is equicontinuous.
- 1825 3. *Uniform vanishing at infinity:*  $\forall \varepsilon > 0 \exists$  compact  $K \subset X : \sup_{f \in H} \sup_{x \notin K} |f(x)| < \varepsilon$ .  
1826  
1827

1828 Then  $H$  is totally bounded in  $C_0(X)$ .  
1829

1830 Condition (3) ensures that all functions become uniformly small outside a common compact set,  
1831 enabling one to restrict attention to a compact domain where Corollary 1 applies.  
1832

### 1833 I.2.3 $X$ METRIC

1834 Let  $(X, d)$  be a metric space. Because metric spaces are paracompact and first-countable, local  
1835 equicontinuity simplifies to the existence of a *uniform modulus of continuity*.

1836 **Theorem 12** (Metric Arzelà–Ascoli). *If there exists a modulus of continuity  $\omega : [0, \infty) \rightarrow [0, \infty)$*   
 1837 *with  $\omega(r) \rightarrow 0(r \rightarrow 0)$  such that*

$$1838 |f(x) - f(y)| \leq \omega(d(x, y)) \quad \forall f \in H, \forall x, y \in X,$$

1840 *and  $H$  is uniformly bounded, then  $H$  is totally bounded in  $C(X)$ .*

1842 A common—and often sufficient—specialization is a *uniform Lipschitz bound*:

$$1844 \sup_{f \in H} \text{Lip}(f) < \infty, \quad \text{Lip}(f) = \sup_{x \neq y} \frac{|f(x) - f(y)|}{d(x, y)}.$$

#### 1847 I.2.4 $X$ PSEUDO-METRIC

1849 A pseudo-metric  $d$  allows zeros for  $x \neq y$ . Let  $\tilde{X} = X / \sim$  with  $x \sim y \iff d(x, y) = 0$ ; then  $d$   
 1850 induces a genuine metric  $\tilde{d}$  on  $\tilde{X}$ , and composition with the quotient map identifies  $C(\tilde{X}) \cong C(X)$ .  
 1851 Apply Theorem 12 to  $(\tilde{X}, \tilde{d})$ . Explicitly:

1852 **Corollary 2.** *If a modulus of continuity  $\omega$  as in Theorem 12 exists with respect to the pseudo-metric*  
 1853  *$d$ , and  $H$  is uniformly bounded, then  $H$  is totally bounded in  $C(X)$ .*

#### 1856 I.2.5 WEAKER COMPACTNESS HYPOTHESES ON $H$

1857 Frequently  $H$  is known to satisfy a property *weaker* than full relative compactness—e.g. *uniform*  
 1858 *boundedness plus equicontinuity*, but *without* an a priori modulus common to all  $f \in H$ . The follow-  
 1859 ing lemma fills the gap when  $X$  is metric (or after reducing to a metric space as in Appendix I.2.4).

1860 **Lemma 9** (Modulus extraction). *Let  $(X, d)$  be a totally bounded metric space and let  $H \subset C(X)$*   
 1861 *be uniformly bounded and equicontinuous. Then there exists a single modulus  $\omega$  satisfying the*  
 1862 *assumption of Theorem 12. Consequently  $H$  is totally bounded.*

1864 *Sketch.* Fix a dense sequence  $(x_n)$  in  $X$ . For each  $k \in \mathbb{N}$  use equicontinuity at  $x_k$  and uniform  
 1865 boundedness to find a local modulus; take a maximum over finitely many  $x_k$ 's to build a global  $\omega$ .  
 1866 Details follow standard proofs of the classical Ascoli theorem.  $\square$

### 1869 I.3 POSITIVE EXAMPLES

1870 *Example 13* (compact domain). Let  $X = [0, 1]$  (compact metric space), and

$$1872 H := \{\text{polynomials } p \text{ of degree } \leq m \text{ such that } |p(x)| \leq 1 \text{ for all } x \in [0, 1]\}.$$

1874 Then,  $H$  is totally bounded because all the criteria of Theorem 12 are satisfied as follows:

- 1876 • Uniform boundedness:  $\|p\|_\infty \leq 1$ .
- 1877
- 1878 • Common Lipschitz constant exists on the compact interval (classical Markov-type esti-  
 1879 mates).

1880 *Example 14* (vanishing at infinity). Let  $X = \mathbb{R}$ , and

$$1882 H = \{x \mapsto \sin(x/n)\}_{n \in \mathbb{N}}.$$

1883 Then  $H$  is totally bounded because all conditions of Theorem 11 hold as follows:

- 1885 • Uniform boundedness:  $\|f\|_\infty \leq 1$ .
- 1886
- 1887 • Common Lipschitz bound:  $\text{Lip}(f) \leq 1$ .
- 1888
- 1889 • Uniformly vanishing at infinity: for any  $\varepsilon > 0$  choose compact  $K = [-R, R]$  with  $R$  large;  
 then  $|f(x)| < \varepsilon$  for  $|x| > R$  uniformly in  $n$ .

1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

### I.3.1 FAILURE EXAMPLES

*Example 15* (due to behavior at infinity). Let  $X = \mathbb{R}$  (locally compact but not compact), and

$$H = \{x \mapsto \sin(nx)\}_{n \in \mathbb{N}}.$$

Then,  $H$  is *not* totally bounded in  $C(\mathbb{R})$  because The functions do *not* vanish at infinity uniformly; Condition (3) of Theorem 11 fails.

*Example 16* (due to lack of equicontinuity). Let  $X = [0, 1]$ , and

$$H = \{x \mapsto \sqrt{x} + 1/n\}_{n \in \mathbb{N}}.$$

Then  $H$  is *not* totally bounded because uniform boundedness holds, but the functions are *not* equicontinuous at  $x = 0$  (the derivative blows up); Criterion (2) of Corollary 1 fails.

## J GUIVARC’H–BASS FORMULA AND HOMOGENEOUS DIMENSION

Following Breuillard (2014), we quickly overview homogeneous dimension and Guivarc’h–Bass formula for locally compact groups with polynomial growth.

### J.1 HOMOGENEOUS DIMENSION

**Definition 4** (Homogeneous dimension of nilpotent Lie groups). Let  $N$  be a simply connected nilpotent Lie group with Lie algebra  $\mathfrak{n}$  and central descending series  $C_i(\mathfrak{n})$  (so  $C_1(\mathfrak{n}) = \mathfrak{n}$ ,  $C_{i+1}(\mathfrak{n}) = [\mathfrak{n}, C_i(\mathfrak{n})]$ ). The *homogeneous dimension* of  $N$  is

$$d(N) = \sum_{i \geq 1} \dim C_i(\mathfrak{n}).$$

Equivalently, for a grading  $\mathfrak{n} = \bigoplus_{i \geq 1} \mathfrak{m}_i$  adapted to the lower central series, one has

$$d(N) = \sum_{i \geq 1} i \cdot \dim \mathfrak{m}_i.$$

The second formula is the Guivarc’h–Bass form; see below.

### J.2 THE GUIVARC’H–BASS FORMULA

**Theorem 13** (Corollary 2.9 and Equation 17). *For a (simply connected) nilpotent Lie group  $N$  with a compact neighborhood  $\mathcal{U}$  of the identity and any  $n \geq 1$ , volume of the product sets satisfies*

$$C_1 n^d \leq \text{vol}_N(\mathcal{U}^n) \leq C_2 n^d$$

for some constants  $C_1, C_2 > 0$ , where

$$d = \sum_{i \geq 1} i \dim \mathfrak{m}_i \quad (\text{Guivarc’h–Bass})$$

with  $\mathfrak{n} = \bigoplus_{i \geq 1} \mathfrak{m}_i$  a grading compatible with the lower central series; equivalently  $d = \sum_{i \geq 1} \dim C_i(\mathfrak{n})$ .

*Remark 7.* In the general case (Theorem 1.1), the exponent  $d(G)$  for any locally compact  $G$  of polynomial growth is exactly the Guivarc’h–Bass homogeneous dimension of the graded nilpotent Lie group that arises as the asymptotic cone (graded nilshadow) of  $G$ .

### J.3 HOW HOMOGENEOUS DIMENSION CONTROLS VOLUME GROWTH

**Theorem 14** (General  $G$ , Theorem 1.1).  *$G$  locally compact of polynomial growth,  $\Omega$  compact symmetric generating set. Then,  $\text{vol}_G(\Omega^n) \sim c(\Omega) n^{d(G)}$  with  $d(G)$  equal to the Guivarc’h–Bass homogeneous dimension of the graded nilshadow (asymptotic cone).*

1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997

#### J.4 SUFFICIENT CONDITIONS FOR POLYNOMIAL GROWTH

**Connected Lie groups of type  $(R)$ .** A connected Lie group  $S$  has polynomial growth iff it is of type  $(R)$ , i.e.  $\text{ad}(x)$  has only purely imaginary eigenvalues for all  $x \in \mathfrak{s}$ . In particular, such  $S$  is solvable-by-compact, and every connected nilpotent Lie group is of type  $(R)$ , hence of polynomial growth. (Guivarc’h–Jenkins characterization)

**Passing to cocompact subgroups / compact quotients.** Polynomial growth is preserved in both directions when passing to a *cocompact subgroup* and when taking quotients by *compact normal subgroups*. (Lemma 7.7)

**Discrete subgroups inside solvable type  $(R)$ .** If  $\Gamma$  is a discrete subgroup of a connected solvable Lie group of type  $(R)$ , then  $\Gamma$  is *virtually nilpotent* (hence of polynomial growth). (Remark 7.8)

**Virtually nilpotent (Gromov, 1981; Losert, 1987)** Finitely generated groups of polynomial growth are *virtually nilpotent*. Losert extends the structural picture to locally compact groups.

These give practical sufficient hypotheses: *nilpotent* (or graded nilpotent/Carnot), *connected type  $(R)$* , *cocompact embedding in a polynomial-growth group*, and *virtually nilpotent discrete subgroups* of solvable type  $(R)$ .

#### J.5 CONCRETE EXAMPLES OF HOMOGENEOUS DIMENSION

The homogeneous dimension  $d(G)$  coincides with the algebraic/geometric dimension  $d$  if the group is abelian (e.g.  $\mathbb{Z}^d$  and  $\mathbb{R}^d$ ); while  $d(G)$  is larger if the group is non-abelian. Thus the volume growth faster when  $G$  is non-abelian.

*Example 17* ( $\mathbb{Z}^d$  and  $\mathbb{R}^d$  (abelian)). Grading has only layer  $\mathfrak{m}_1$  of dimension  $d$ ; thus  $d(N) = 1 \cdot d = d$ . Hence volume grows like  $t^d$ . general formula above.)

*Example 18* (Heisenberg group  $H_3$ , Section 9).  $\dim \mathfrak{m}_1 = 2$  and  $\dim \mathfrak{m}_2 = 1$ , so  $d(H_3) = 1 \cdot 2 + 2 \cdot 1 = 4$ .

*Example 19* (Heisenberg group  $H_5$ , Section 9.2).  $\dim \mathfrak{m}_1 = 4$ ,  $\dim \mathfrak{m}_2 = 1$ , so  $d(H_5) = 4 + 2 = 6$ .

*Example 20* (Unitriangular group  $UT(n)$  — strictly upper triangular  $n \times n$  matrices, step  $n - 1$ ). Using Guivarc’h–Bass, the  $i$ -th layer has  $\dim \mathfrak{m}_i = n - i$  (entries on the  $i$ -th superdiagonal). Hence

$$d(UT(n)) = \sum_{i=1}^{n-1} i(n-i) = \frac{n(n-1)(n+1)}{6}.$$

*Example 21* (A worked solvable example, Example 3.3). For  $G = \mathbb{R} \ltimes_{\varphi} \mathbb{R}^n$  where the unipotent part of  $\varphi_t$  has  $n_k$  Jordan blocks of size  $k$ , the paper computes

$$d(G) = 1 + \sum_{k \geq 1} \frac{k(k+1)}{2} n_k.$$