

Constructive interpolation and generalization rates for neural ODEs: a control perspective

Antonio Álvarez-López^{1,2,*}

Lorenzo Liverani¹

Enrique Zuazua^{1,2,3}

antonio.alvarezl@uam.es

lorenzo.liverani@fau.de

enrique.zuazua@fau.de

Abstract

We study supervised regression with neural ODEs (NODEs) from a control-theoretic perspective to derive explicit population-risk bounds. We focus on a widely used class of non-autonomous models with constant parameters and explicit time dependence, which we call semi-autonomous NODEs (SA-NODEs). We constructively prove that SA-NODEs are capable of *exact* interpolation of admissible finite datasets, and even satisfy a stronger property that we call *simultaneous cell controllability* (SCC): their flows can map prescribed disjoint cells into arbitrarily small target balls. This property is the mechanism that upgrades interpolation into quantitative generalization, by allowing SA-NODEs to emulate piecewise-constant nonparametric estimators. Consequently, our risk bounds recover the rates of histogram and nearest-neighbor estimators, provided the network width satisfies a conservative scaling with the sample size. Numerical experiments show that trained SA-NODEs achieve competitive—often lower—test errors than these baselines. Finally, we show that the explicit time dependence is essential. Although two-layer autonomous NODEs can interpolate geometrically nondegenerate datasets, structural obstructions prevent them from achieving SCC. These limitations, further confirmed numerically, support the view that SA-NODEs provide a minimal effective architecture for learning.

Keywords. Neural ODEs; controllability; interpolation; nonparametric regression; generalization bounds.

MSC. 68T07, 93B05, 62G08, 34H05, 68Q32.

1 Introduction

1.1 Expressivity and generalization

Supervised learning is one of the central paradigms in machine learning, with successful applications across a wide range of problems in science and engineering [20]. The standard framework relies on a dataset of size $N \geq 1$, consisting of pairwise distinct inputs $\mathbf{x}_i \in \mathcal{X}$ and corresponding labels $\mathbf{y}_i \in \mathcal{Y}$,

$$\mathcal{D}_N = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N, \quad \text{with } \mathbf{x}_i \neq \mathbf{x}_j \text{ for } i \neq j. \quad (1)$$

The goal is to approximate the unknown relationship between inputs and labels by fitting a predictor to the sample \mathcal{D}_N , typically via a parametric family $\{F_\theta\}_\theta$, where θ is the learnable parameter. The choice of \mathcal{Y} determines the task. In this work, we focus on regression in a deterministic, noiseless setting. Accordingly, we assume throughout that $\mathcal{X} \subseteq \mathbb{R}^d$ and $\mathcal{Y} = \mathbb{R}^d$ for $d \geq 2$, and that there exists a ground-truth map $y : \mathcal{X} \rightarrow \mathcal{Y}$ such that $y(\mathbf{x}_i) = \mathbf{y}_i$ for all i .

Of course, fitting \mathcal{D}_N alone does not constitute learning: the true goal is to approximate y accurately across \mathcal{X} . To assess performance beyond the training data, we fix a test probability measure $\mu \in \mathcal{P}(\mathcal{X})$, often viewed as the (unknown) data-generating distribution, and define the population (true) risk by

$$\mathcal{R}(\theta) := \int \ell(F_\theta(\mathbf{x}), y(\mathbf{x})) \, d\mu(\mathbf{x}), \quad (2)$$

where $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$ is a chosen loss function. Since μ is unknown, the population risk $\mathcal{R}(\theta)$ cannot be computed, let alone minimized, directly from \mathcal{D}_N . One therefore introduces the empirical risk

$$\mathcal{R}_N(\theta) := \frac{1}{N} \sum_{i=1}^N \ell(F_\theta(\mathbf{x}_i), \mathbf{y}_i) = \int \ell(F_\theta(\mathbf{x}), y(\mathbf{x})) \, d\mu_N(\mathbf{x}), \quad \mu_N := \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}_i} \in \mathcal{P}(\mathcal{X}). \quad (3)$$

¹Chair for Dynamics, Control, Machine Learning, and Numerics (Alexander von Humboldt Professorship), Department of Mathematics, Friedrich-Alexander-Universität Erlangen-Nürnberg, 91058 Erlangen, Germany.

²Departamento de Matemáticas, Universidad Autónoma de Madrid, 28049 Madrid, Spain.

³Chair of Computational Mathematics, Universidad de Deusto, Av. de las Universidades 24, 48007 Bilbao, Basque Country, Spain.

*Corresponding author.

16 Training then amounts to selecting parameters θ that minimize the empirical risk. This approach intro-
 17 duces a fundamental tension, which is revealed by the following error decomposition:

$$\mathcal{R}(\theta) \leq \mathcal{R}_N(\theta) + |\mathcal{R}(\theta) - \mathcal{R}_N(\theta)|. \quad (4)$$

18 The first term measures the fit on the dataset; driving it to zero (exact interpolation) requires enough model
 19 *expressivity*. The second term is the *generalization gap*, measuring the discrepancy between the expected
 20 risk under μ and its empirical proxy μ_N . Classical approaches usually control this gap by constraining
 21 the parameter space. This strategy is not designed for the interpolation regime, where models are in-
 22 herently overparameterized. Our purpose is to exhibit a different mechanism, based on controllability of
 23 macroscopic regions, by which exact interpolation and quantitative generalization can coexist.

24 Controlled flows and neural ODEs

25 We study expressivity and generalization within a class of predictors derived from dynamical systems.
 26 More precisely, we cast supervised learning under the guise of *controlled flows*, modeled by the ODE

$$\begin{cases} \dot{\mathbf{x}}(t) = v_{\theta(t)}(\mathbf{x}(t)), & t \in (0, T], \\ \mathbf{x}(0) = \mathbf{x}_0 \in \mathbb{R}^d, \end{cases} \quad (5)$$

27 with $T > 0$ fixed, and the vector field $v_{\theta(\cdot)}(\cdot) : [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ that we assume to be measurable in t
 28 and uniformly Lipschitz in \mathbf{x} . Under these conditions, (5) induces a unique flow

$$\Phi_t^\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad t \in [0, T],$$

29 and the predictor is $F_\theta = \Phi_T^\theta$. This formulation offers a control-theoretic interpretation [14, 28]. The
 30 parameters $\theta(\cdot)$ act as *controls* steering the dynamics, and empirical risk minimization becomes, in essence,
 31 a problem of *simultaneous controllability*: one seeks some $\theta(\cdot)$ such that Φ_T^θ drives each initial datum \mathbf{x}_i
 32 to its prescribed target \mathbf{y}_i . This perspective places the rich machinery of control theory at the service of
 33 machine learning and provides a unified framework to study both expressivity and generalization.

34 From a machine learning perspective, it is natural to parameterize v_θ as a neural network. With this
 35 choice, (5) is termed a neural ODE, with many variants differing in parameterizations, readouts or dis-
 36 cretizations [9, 14, 18, 23]. In the simple case of a single-layer ReLU network of width $p \in \mathbb{N}$, we obtain:

$$\dot{\mathbf{x}}(t) = \sum_{i=1}^p \mathbf{w}_i(t) (\mathbf{a}_i(t) \cdot \mathbf{x}(t) + b_i(t))_+, \quad t \in [0, T]. \quad (\text{NODE})$$

38 Here, $\theta = (\mathbf{w}_i, \mathbf{a}_i, b_i)_{i=1}^p : [0, T] \rightarrow (\mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R})^p$ are the controls, and $z_+ := \max\{z, 0\}$ is the ReLU.

39 While (NODE) is highly flexible, it is heavily parameterized, with its effective parameter count growing
 40 under finer temporal discretizations. This motivates the search for simpler architectures that preserve
 41 expressive power. We focus on one such model: the *semi-autonomous neural ODE*,

$$\dot{\mathbf{x}}(t) = \sum_{i=1}^p \mathbf{w}_i (\mathbf{a}_i \cdot \mathbf{x}(t) + b_i t + c_i)_+, \quad t \in [0, T], \quad (\text{SANODE})$$

42 derived from (NODE) by restricting \mathbf{w}_i and \mathbf{a}_i to be time-independent, and $b_i(t) = b_i t + c_i$. Equivalently,
 43 this corresponds to parameterizing v_θ via a time-invariant network applied to the extended state (\mathbf{x}, t) .
 44 This design makes (SANODE) an appealing middle ground [15, 27]: its relative simplicity eases implemen-
 45 tation with respect to (NODE) while retaining high expressivity [21].

46 1.2 Main results

47 Our goal is to establish expressivity and generalization guarantees for neural ODEs. We demonstrate that
 48 (SANODE) achieves both, avoiding structural limitations of purely autonomous models.

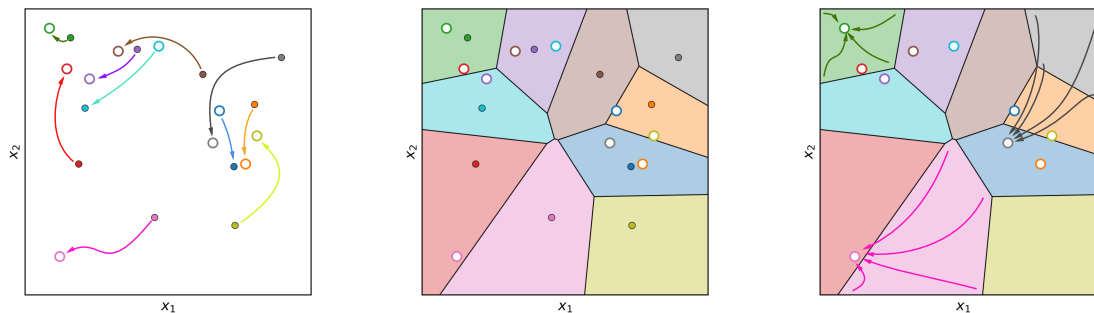


Figure 1: Overview of the proposed approach. **Left:** Interpolation of the training data from inputs (solid circles) to targets (hollow circles) via learned trajectories (arrows). **Center:** Voronoi partition of the input space, where each cell corresponds to a training sample. **Right:** The flow compresses entire cells toward their respective targets, illustrating the simultaneous cell controllability from Definition 2.

49 1.2.1 Expressivity.

50 The study of controllability for neural ODEs has garnered significant attention [1–3, 10, 28, 29]. A central
 51 concept in this literature is the *universal interpolation property* (UIP, [10]): given finitely many pairwise
 52 distinct inputs \mathbf{x}_i and targets \mathbf{y}_i , does there exist a control $\theta(\cdot)$ such that $\Phi_T^\theta(\mathbf{x}_i) \approx \mathbf{y}_i$ for all i ? While
 53 *approximate* interpolation (up to arbitrary precision) holds for various models, the exact UIP ($\Phi_T^\theta(\mathbf{x}_i) = \mathbf{y}_i$
 54 for all i) is much more rigid. For time-dependent controls as in (NODE) or its linear-in-control version (see
 55 [1]), it is typically achieved by introducing temporal discontinuities [3, 28]. However, this machinery
 56 breaks down for (semi-)autonomous systems, leaving their exact UIP as a challenging open question.

57 Beyond exact interpolation of points, we introduce the stronger notion of *simultaneous cell controllability*
 58 to control entire regions of the input space. As illustrated in Figure 1, this requires the flow to map
 59 disjoint convex cells (e.g., a Voronoi partition) toward their respective targets.

60 With these notions in place, our expressivity results for (SANODE) establish the following:

- 61 • We prove the exact UIP for (SANODE) in Theorem 2.3. Our approach is constructive and yields explicit
 62 controls rather than merely proving their existence.
- 63 • We formalize simultaneous cell controllability and prove that (SANODE) exhibits this property in Theo-
 64 rem 2.5. Via Barron constant estimation, we derive explicit bounds on the required network width.

65 1.2.2 Generalization.

66 Under the setting of (5) and the squared Euclidean loss, the population and empirical risks (2)–(3) become:

$$67 \quad \mathcal{R}(\theta) = \int \|\Phi_T^\theta(\mathbf{x}) - y(\mathbf{x})\|^2 d\mu(\mathbf{x}), \quad \mathcal{R}_N(\theta) = \frac{1}{N} \sum_{i=1}^N \|\Phi_T^\theta(\mathbf{x}_i) - \mathbf{y}_i\|^2. \quad (6)$$

68 The central challenge remains bounding the decomposition (4). Existing theoretical approaches typically
 69 bound the generalization gap uniformly over a restricted parameter space Θ . For instance, [22] proves
 70 that for the linear-in-control version of (NODE) under bounded Lipschitz controls, with probability $1 - \delta$:

$$\mathcal{R}(\theta) \leq \mathcal{R}_N(\theta) + \sup_{\theta \in \Theta} |\mathcal{R}(\theta) - \mathcal{R}_N(\theta)| \lesssim \mathcal{R}_N(\theta) + \sqrt{\frac{(m+1) \log(R_\Theta m N)}{N}} + \frac{m\sqrt{R_\Theta}}{N^{1/4}} + \sqrt{\frac{\log(1/\delta)}{N}}$$

71 where m is the number of scalar control functions and $R_\Theta \geq 0$ encapsulates the bounds on Θ . If the
 72 controls are time-independent, this bound improves by removing the $O(N^{-1/4})$ term.

73 While informative, such bounds highlight an inherent tension: keeping the complexity (m and R_Θ)
 74 fixed as $N \rightarrow \infty$ shrinks the generalization gap but may preclude exact interpolation, leaving the empirical
 75 risk $\mathcal{R}_N(\theta)$ large. Conversely, driving $\mathcal{R}_N(\theta)$ to zero requires increasing complexity when N grows large,
 76 which deteriorates the bound. This motivates a natural **question**:

77 *Can we design data-dependent controls θ_N , while appropriately scaling the network width, to*
 78 *guarantee that the population risk $\mathcal{R}(\theta_N)$ vanishes at a quantifiable rate as $N \rightarrow \infty$?*

79 We provide an affirmative answer by connecting (SANODE) to nonparametric statistics. Rather than assum-
 80 ing a fixed parametric form, nonparametric methods—such as kernel estimators [24, 33], histograms or
 81 k -nearest-neighbors [16, 31]—construct local approximations directly from the data. For Hölder spaces,
 82 these procedures yield explicit, minimax-optimal convergence rates [30]. This establishes a natural bench-
 83 mark: any learning method claiming to generalize should ideally match these rates. We show that (SANODE)
 84 achieves this benchmark by using simultaneous cell controllability to emulate these estimators.

85 **INFORMAL STATEMENT** (Propositions 4.5 and 4.6). Given a dataset \mathcal{D}_N , a compactly supported probability
 86 measure μ , and an α -Hölder continuous target y , there exists a control θ_N such that

$$\mathcal{R}(\theta_N) \lesssim \text{Err}_{\text{np}} + \text{Err}_{\text{node}},$$

87 where Err_{np} is the statistical error of a piecewise-constant nonparametric estimator built from \mathcal{D}_N , and
 88 Err_{node} measures how accurately the flow of (SANODE) realizes this estimator. Provided the network width
 89 p grows sufficiently fast with N ($p \geq p_N$), the realization error is absorbed. Taking the expectation over
 90 the random draw of \mathcal{D}_N (if the inputs \mathbf{x}_i are sampled i.i.d. from μ), we recover the classical rates:

$$\mathbb{E}_{\mathcal{D}_N} [\mathcal{R}(\theta_N)] \lesssim \begin{cases} N^{-\frac{2\alpha}{2\alpha+d}} & \text{for histogram estimators,} \\ \left(\frac{\log N}{N}\right)^{\frac{2\alpha}{d}} & \text{for nearest-neighbor estimators.} \end{cases}$$

91 Our width requirement p_N —which depends on N and the geometry of the data partition—acts as an
 92 achievability guarantee rather than a tight bound. It proves that the architecture inherently supports gen-
 93 eralization: given sufficient width, there exists a control whose true risk vanishes at an explicit rate. This
 94 circumvents the inherent tension of uniform bounds: exact interpolation and optimal statistical consis-
 95 tency can mathematically coexist.

96 **Related work.** Recent literature has explored generalization across various continuous-time models.
 97 Building on the uniform bounds of [22] for the linear-in-control version of (NODE), [32] extended this
 98 capacity-constrained approach to broader nonlinear neural ODEs. In parallel, other works have studied
 99 complementary mechanisms: [17] showed that incorporating feedback enhances robustness, while [8]
 100 derived generalization bounds for neural controlled differential equations driven by irregular time series.

101 1.2.3 The limits of autonomous models

102 Having established the capabilities of (SANODE), it is natural to ask whether explicit time dependence is
 103 truly necessary. We may remove it entirely, leading to the *autonomous neural ODE*:

$$\dot{\mathbf{x}}(t) = \sum_{i=1}^p \mathbf{w}_i (\mathbf{a}_i \cdot \mathbf{x}(t) + b_i)_+, \quad t \in [0, T]. \quad (\text{ANODE})$$

104 To compensate for the resulting structural constraints, we also consider increasing the compositional
 105 depth, yielding the *two-layer autonomous neural ODE* (2ANODE):

$$\dot{\mathbf{x}}(t) = \sum_{j=1}^{p_2} \mathbf{w}_j \left(\sum_{i=1}^{p_1} u_{ji} (\mathbf{a}_{ij} \cdot \mathbf{x}(t) + b_{ij})_+ + c_j \right)_+, \quad t \in [0, T], \quad (2\text{ANODE})$$

106 with $\mathbf{w}_j, \mathbf{a}_{ij} \in \mathbb{R}^d$ and $u_{ji}, b_{ij}, c_j \in \mathbb{R}$ for $p_1, p_2 \in \mathbb{N}$. Although not our primary focus, analyzing these
 107 autonomous models helps contextualize the advantages of explicit time dependence.

108 While we prove in Theorem 2.9 that (2ANODE) can achieve the exact UIP (provided the data satisfy the
 109 geometric condition (13)), we also discuss obstructions to simultaneous cell controllability in autonomous
 110 models like (ANODE) and (2ANODE). This reveals a gap between point interpolation and cell routing, a
 111 distinction that Section 5.2 illustrates experimentally. Thus, among the models we consider, (SANODE)
 112 provides the simplest setting in which both point interpolation and cell routing mechanisms coexist.

1.3 Organization

Section 2 introduces the core controllability notions and provides constructive proofs of the exact UIP for (SANODE) and (2ANODE). Section 3 reviews the nonparametric estimators that serve as our benchmarks. Section 4 then states and discusses our main generalization theorems. In Section 5, we numerically evaluate the capacity of neural ODEs to match or outperform these benchmarks on several regression tasks. Finally, Section 6 offers concluding remarks and future directions. Most proofs are deferred to Section 7.

Notation

For $n \in \mathbb{N}$, we use $[n] = \{1, \dots, n\}$. Vectors are bold (e.g. \mathbf{u}, \mathbf{v}) and we denote the Euclidean inner product of $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ by $\mathbf{u} \cdot \mathbf{v}$. The open ball in \mathbb{R}^d of radius $r > 0$ centered at \mathbf{x} is $B(\mathbf{x}, r)$, and the hypercube is $I_R := [-R, R]^d$ for $R > 0$. Let $\mathcal{P}_{\text{ac}}(\mathbb{R}^d)$ be the space of absolutely continuous probability measures on \mathbb{R}^d . We set

$$\log_+(r) := \max\{\log(r), 0\} \quad \text{for } r > 0, \quad \text{and} \quad \log_+(0) = 0. \quad (7)$$

Concerning inequalities, we write $a \lesssim b$ (and $a \gtrsim b$) to indicate that $a \leq Cb$ (and $a \geq Cb$) for some constant $C > 0$ whose independence from the parameters of interest (such as N, ε , or p) will be clear from the context. If $a \lesssim b$ and $b \lesssim a$, we write $a \asymp b$.

2 Controllability of neural ODEs

2.1 Simultaneous point controllability

We begin with our central notion of expressivity. Consider the general model (5) with flow map Φ_t^θ . Since $\Phi_T^\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a homeomorphism, and hence injective, exact interpolation is impossible if distinct inputs $\mathbf{x}_i \neq \mathbf{x}_j$ are assigned the same target. To avoid this, we restrict our attention to datasets satisfying

$$\mathbf{x}_i \neq \mathbf{x}_j \quad \text{and} \quad \mathbf{y}_i \neq \mathbf{y}_j \quad \text{for all } i \neq j. \quad (8)$$

Any dataset satisfying (8) will be called *admissible*.

Definition 1. We say that the control system associated with (5) possesses the approximate universal interpolation property (UIP) if for every $\varepsilon > 0$, every $N \in \mathbb{N}$, and every dataset $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \subset \mathbb{R}^d \times \mathbb{R}^d$ admissible in the sense of (8), there exist $T > 0$ and an admissible control θ such that the flow map Φ_T^θ generated by (5) satisfies

$$\|\Phi_T^\theta(\mathbf{x}_i) - \mathbf{y}_i\| \leq \varepsilon, \quad \text{for all } i \in [N]. \quad (9)$$

If (9) holds with $\varepsilon = 0$, we say that the system possesses the (exact) UIP.

When the vector field is a neural network, the admissible controls in Definition 1 are understood to range over all widths. Thus, the UIP is a property of the whole family, not of a fixed-width model.

Remark 2.1 (Time rescaling). In the definition above, the time horizon T is allowed to be chosen freely. However, this is equivalent to a fixed-horizon formulation whenever the family of vector fields under consideration is closed under positive scalar multiplication and time reparameterization. Indeed, for any $T_* > 0$, the rescaled curve $\mathbf{y}(s) := \mathbf{x}(sT/T_*)$ solves an analogous system on $[0, T_*]$ driven by the vector field $v_{\theta(sT/T_*)}T/T_*$, which remains in the same class. All models considered in this work satisfy this property, provided there are no uniform bounds imposed on the parameter space.

Before detailing our contributions, we briefly review known results for the general model (NODE). The UIP was proved in [28] using $p = 1$ neurons. We state its generalization to arbitrary $p \geq 1$ from [3].

Theorem 2.2 (Simultaneous point controllability for (NODE) [3, Theorem 1]). For every $d \geq 2, T > 0$ and $p \in \mathbb{N}$, the system (NODE) possesses the exact UIP. Furthermore, for any admissible dataset of size N (in the sense of (8)), the control can be chosen piecewise constant in time with $2 \lceil N/p \rceil - 1$ jumps.

152 The constructive proof of this theorem exploits the freedom afforded by time-dependent parameters.
 153 The key device is the orthogonality constraint $\mathbf{w}_i \cdot \mathbf{a}_i = 0$: geometrically, this forces the flow generated by
 154 each neuron to act as a shear transformation parallel to its own activation hyperplane. This decoupling
 155 makes the dynamics tractable and allows the control strategy to split naturally into two phases separated
 156 by a single switch. In the first interval, N parallel hyperplanes steer $d - 1$ coordinates of the data
 157 simultaneously; after the switch, the hyperplanes reorient to handle the remaining coordinate.

158 By contrast, for (SANODE) the vectors \mathbf{w}_i and \mathbf{a}_i are time-independent, so the reorientation mechanism
 159 described above is no longer available. The compensating source of flexibility is the explicit time dependence
 160 in the argument $\mathbf{a}_i \cdot \mathbf{x} + b_i t + c_i$, which makes the activation hyperplanes translate at a constant
 161 velocity given by b_i . Our first new result—whose proof is deferred to Section 7.1.1—shows that this single
 162 degree of freedom is sufficient to recover exact controllability without requiring any temporal switching.

163 **Theorem 2.3** (Simultaneous point controllability for (SANODE)). *For every $d \geq 2$ and $T > 0$, the system
 164 (SANODE) possesses the exact UIP. Moreover, for any admissible dataset of size N (in the sense of (8)), exact
 165 interpolation can be achieved with width $p = 2N$.*

166 2.2 Simultaneous cell controllability

167 The UIP captures the capacity to interpolate any admissible finite dataset. To control the population risk,
 168 however, one needs the flow to behave well across entire regions of the input space, rather than only at
 169 finitely many points. This leads to the following stronger controllability notion.

170 **Definition 2.** *We say that the control system associated with (5) possesses the property of simultaneous
 171 cell controllability (SCC) if, for every finite family of pairwise disjoint compact convex sets $\{A_k\}_{k=1}^K \subset \mathbb{R}^d$,
 172 every collection of target points $\{\mathbf{r}_k\}_{k=1}^K \subset \mathbb{R}^d$, and every $\eta > 0$, there exist $T > 0$ and an admissible control
 173 θ such that the flow map Φ_T^θ generated by (5) satisfies*

$$\Phi_T^\theta(A_k) \subset B(\mathbf{r}_k, \eta) \quad \text{for all } k \in [K]. \quad (10)$$

174 **Remark 2.4.** *Although framed here in the context of flows, simultaneous cell controllability is an architecture-
 175 independent notion. It applies to any parametric map F_θ , from neural networks to classical regression models.*

176 The property of SCC is stronger than the approximate UIP, as setting $A_k = \{\mathbf{x}_k\}$ in (10) directly
 177 recovers condition (9). The fundamental shift is that the flow must now simultaneously route whole
 178 macroscopic regions rather than isolated points. Despite this added complexity, the following theorem
 179 establishes that (SANODE) exhibits this property, and provides a lower bound on the required network
 180 width. This is derived in Section 7.1.3 by constructing a reference vector field and applying quantitative
 181 Barron-type estimates (with the geometric constants detailed in Remark 2.6).

182 **Theorem 2.5** (Simultaneous cell controllability for (SANODE)). *Let $d, K \geq 2$ and $R, T > 0$. Let $\{A_k\}_{k=1}^K$
 183 be a family of pairwise disjoint, compact, and convex subsets of $I_R := [-R, R]^d$, and let $\{\mathbf{r}_k\}_{k=1}^K \subset \mathbb{R}^d$ be
 184 pairwise distinct target vectors. Then there exist*

$$\mathfrak{C} = \mathfrak{C}(d, T, R, (\mathbf{r}_k)_k, (A_k)_k) > 0, \quad \mathfrak{M} = \mathfrak{M}(d, T, R, (\mathbf{r}_k)_k) > 0$$

185 and $\eta_0 \in (0, 1]$ such that for every $\eta \in (0, \eta_0]$ and every integer

$$p \geq \mathfrak{C} \eta^{-2-\mathfrak{C}} (1 - \log \eta)^2, \quad (11)$$

186 there exists a control θ of width p for (SANODE) such that

- 187 1. $\Phi_T^\theta(A_k) \subset B(\mathbf{r}_k, \eta)$ for all $k \in [K]$;
- 188 2. $\|\Phi_T^\theta\|_{L^\infty(I_R)} \leq \mathfrak{M}$.

189 *If the targets are not necessarily pairwise distinct, the same conclusions hold for any fixed $\eta > 0$; however,
 190 the constant \mathfrak{C} is no longer uniform and will depend on η .*

191 **Remark 2.6** (Explicit geometric constants). *The constants \mathfrak{C} , \mathfrak{M} , and η_0 in Theorem 2.5 may be chosen*
 192 *explicitly. A possible choice, obtained by combining Lemma 7.6 with Lemma 7.7, is governed by the geometry*
 193 *of the cells and their targets. Assume that the target vectors $\mathbf{r}_1, \dots, \mathbf{r}_K$ are pairwise distinct, let*

$$D_* := \max_{k \in [K]} \text{diam}(A_k), \quad s_* := \min_{i \neq j} \text{dist}(A_i, A_j),$$

194 *and let $\gamma_k : [0, T] \rightarrow \mathbb{R}^d$ be smooth pairwise disjoint curves such that $\gamma_k(0) \in A_k$, $\gamma_k(T) = \mathbf{r}_k$, γ_k is*
 195 *constant on $[0, 2T/3]$, and $\|\gamma_k(t)\| < 1 + \max_{k \in [K]} \|\mathbf{r}_k\| + \sqrt{d} R$ for all $t \in [0, T]$ —their existence is*
 196 *guaranteed for $d \geq 2$, see Lemma 7.5. Define*

$$m := \min_{t \in [0, T]} \min_{i \neq j} \|\gamma_i(t) - \gamma_j(t)\|, \quad G := 1 + \max_{k \in [K]} \max_{1 \leq j \leq d+5} \left\| \frac{d^j \gamma_k}{dt^j} \right\|_{L^\infty(0, T)},$$

197 *and let*

$$L := C_{d, T} \left[\frac{G}{m} + \frac{1 + D_*/s_*}{T} (1 + \log_+(4D_*)) \right],$$

198 *and*

$$B := C_{d, T} K \left[TG^{d+5}(m^d + m^{-4}) + \frac{(D_* + s_*)^{d+1}(1 + s_*^{-(d+4)})}{T} (1 + \log_+(4D_*)) \right].$$

Exploiting these quantities—see the proofs of Lemmas 7.6 and 7.7 for a derivation—one may take

$$\eta_0 := \min \left\{ 1, \frac{s_*}{2}, \frac{m}{2} \right\}, \quad \mathfrak{M} := 2 + \max \left\{ T, \max_{k \in [K]} \|\mathbf{r}_k\| + \frac{3}{2} \sqrt{d} R \right\}, \quad (12)$$

$$\mathfrak{C} := \max \{ 3, 2TL, 4C_d^2 \mathfrak{M}^4 T^2 B^2 e^{2TL} \}.$$

199 *Thus, while \mathfrak{M} depends solely on the domain size and target locations, \mathfrak{C} captures the geometric bottleneck of*
 200 *the partition: the term $1/s_*$ in L forces an exponential penalty $\exp(c/s_*)$ in the factor e^{2TL} as $s_* \rightarrow 0$.*

201 **Remark 2.7** (Width requirement and sharpness). *The lower bound (11) provides a sufficient condition rather*
 202 *than a sharp complexity estimate. As noted in Remark 2.6, the geometric constant \mathfrak{C} incurs an exponential*
 203 *penalty as the minimum cell separation s_* shrinks. Consequently, for the refined partitions used later to*
 204 *establish statistical consistency, this geometric bottleneck can force a super-polynomial threshold p_N as the*
 205 *dataset size N grows. This reflects the curse of dimensionality inherent in isolating many small regions.*

206 *We emphasize that this scaling acts purely as an achievability guarantee. The specific bound arises from*
 207 *our chosen proof strategy—specifically, applying the uniform approximation results of [19, Theorem 2]. Em-*
 208 *ploying alternative tools (e.g., [5, Proposition 6]) could alter the trade-off between the convergence rate in p*
 209 *and the magnitude of \mathfrak{C} . Determining the sharp intrinsic width required for simultaneous cell controllability*
 210 *remains an open quantitative question.*

211 2.3 Autonomous models

212 Having established the **UIP** and **SCC** for (**SANODE**), it is natural to ask whether explicit time dependence
 213 is truly necessary. To answer this, we remove it entirely and analyze the autonomous model (**ANODE**) and
 214 its two-layer variant (**2ANODE**).

215 While [3] establishes the approximate **UIP** for (**ANODE**)—including explicit error decay rates as the
 216 width grows—bridging the gap from approximate to exact interpolation remains unresolved:

217 **Open Problem.** *Does the exact **UIP** hold for the single-layer autonomous system (**ANODE**)?*

218 The difficulty is structural: unlike time-dependent models, which can concatenate flow stages by
 219 switching controls, autonomous single-layer fields are highly rigid. Specifically, they cannot be spatially
 220 localized. Because they are finite superpositions of ridge functions—whose Fourier transforms are sup-
 221 ported on one-dimensional sets—such nontrivial vector fields can never belong to $L^2(\mathbb{R}^d)$. Though ele-
 222 mentary, this fact seems absent from standard references, so we provide a proof.

223 **Lemma 2.8.** *Let $d \geq 2$ and let $\sigma \in \mathcal{C}^0(\mathbb{R})$ have at most polynomial growth. For any $p \geq 1$, define*

$$g(\mathbf{x}) = \sum_{j=1}^p w_j \sigma(\mathbf{a}_j \cdot \mathbf{x} + b_j), \quad (\mathbf{x} \in \mathbb{R}^d), \quad \text{with } \mathbf{a}_j \in \mathbb{R}^d, \quad w_j, b_j \in \mathbb{R}.$$

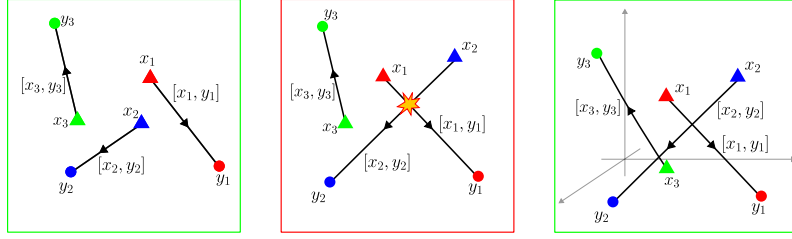


Figure 2: Illustration of the nondegeneracy condition (13). **Left & Right:** Admissible datasets for $d = 2$ and $d = 3$. **Center:** A degenerate data configuration. As discussed in Remark 2.10, condition (13) is generic for $d \geq 3$.

224 If $g \not\equiv 0$, then $g \notin L^2(\mathbb{R}^d)$. In particular, g cannot have compact support unless $g \equiv 0$.

225 *Proof.* For each j , the Fourier transform of the ridge function $\mathbf{x} \mapsto \sigma(\mathbf{a}_j \cdot \mathbf{x} + b_j)$ is supported, in the sense
 226 of tempered distributions, on the line $L_j = \{\lambda \mathbf{a}_j : \lambda \in \mathbb{R}\} \subset \mathbb{R}^d$. By linearity, \widehat{g} is then supported in the
 227 finite union $\bigcup_{j=1}^p L_j$, which has Lebesgue measure zero in \mathbb{R}^d since $d \geq 2$ by hypothesis.

228 Suppose $g \in L^2(\mathbb{R}^d)$. Then, by Plancherel, $\widehat{g} \in L^2(\mathbb{R}^d)$ as well. But $\widehat{g} = 0$ a.e. because an L^2 -function
 229 supported on a null set must vanish almost everywhere, and therefore $g \equiv 0$. \square

230 While Lemma 2.8 does not strictly rule out exact interpolation, it explains why “pointwise” steering
 231 strategies are exceptionally hard to implement with (ANODE) alone. A minimal way to restore localization
 232 without reintroducing time dependence is to add one compositional layer. By composing ReLU units, the
 233 two-layer field (2ANODE) can generate vector fields supported on convex polytopes (see Lemma 7.3 below),
 234 enabling the construction of localized controls.

235 For (2ANODE), we require the stronger assumption that input-output segments are mutually disjoint:

$$[\mathbf{x}_i, \mathbf{y}_i] \cap [\mathbf{x}_j, \mathbf{y}_j] = \emptyset \quad \text{for all } i \neq j, \quad (13)$$

236 where $[\mathbf{a}, \mathbf{b}]$ denotes the line segment between \mathbf{a} and \mathbf{b} . Note that this naturally implies $\mathbf{x}_i \neq \mathbf{y}_j$ for $i \neq j$.

237 **Theorem 2.9** (Simultaneous point controllability for (2ANODE)). *Let $N \geq 1$, $d \geq 2$, and $T > 0$ be fixed.*
 238 *Consider any dataset $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \subset \mathbb{R}^d \times \mathbb{R}^d$ admissible in the sense of (8), and additionally satisfying*
 239 *(13). Then there exists a (constant) control θ such that the flow map Φ_T^θ generated by (2ANODE), with $p_1 = 2d$*
 240 *and $p_2 = N$, satisfies*

$$\Phi_T^\theta(\mathbf{x}_i) = \mathbf{y}_i, \quad \text{for all } i \in [N].$$

241 The proof is deferred to Section 7.1.2. Although (2ANODE) achieves exact UIP without any form of time
 242 dependence, its mechanism is less efficient than that of (SANODE): it requires $p_1 = 2d$ and $p_2 = N$, hence
 243 $p_1 \cdot p_2 = 2dN$ first-layer units and $p_2 = N$ second-layer units. By contrast, (SANODE) uses $p = 2N$ units.

244 **Remark 2.10** (Dimensional constraints). *Condition (13) depends strongly on the dimension d . Because our*
 245 *current construction relies on vector fields localized strictly around straight segments $[\mathbf{x}_i, \mathbf{y}_i]$, these paths must*
 246 *not cross. In $d \geq 3$, linear segments generically do not intersect, so this condition is easily satisfied. In $d = 2$,*
 247 *however, straight-line collisions are frequent, making the assumption highly restrictive (see Figure 2).*

248 *Nevertheless, the UIP can theoretically be established for all $d \geq 2$ without requiring (13). Instead of*
 249 *moving along straight lines, particles could be routed through intermediate waypoints to avoid collisions,*
 250 *using localized two-layer “connectors” to steer the flow around obstacles. While mathematically feasible, the*
 251 *technical length of such a construction places it beyond the scope of this paper, and we defer it to future work.*

252 At this point, the distinction between point interpolation and cell routing becomes essential. While
 253 autonomous systems can map finite point clouds arbitrarily close to their targets by letting isolated tra-
 254 jectories bypass each other in $d \geq 2$ ([3, Theorem 7]), routing macroscopic cells A_k introduces structural
 255 obstructions. Since autonomous trajectories cannot cross, models like (ANODE) or (2ANODE) face a topo-
 256 logical barrier. For example, if the target of A_1 lies behind A_2 , the flow would need to push A_1 directly
 257 through the space initially occupied by A_2 , which violates uniqueness. The explicit time dependence in
 258 (SANODE) circumvents this by allowing trajectories to cross safely in space-time.

259 Furthermore, Definition 2 strictly requires $\Phi_T^\theta(A_k) \subset B(\mathbf{r}_k, \eta)$. A weaker L^2 -a.e. control might by-
 260 pass the obstruction by sacrificing small-measure subsets of the cells: failing to route these small regions to

261 their targets would free up physical space to create “corridors” for other cells to pass through, potentially
 262 accommodating autonomous systems.

263 In summary, while autonomous systems can interpolate points—especially with added depth as in
 264 (2ANODE)—their lack of a time coordinate to schedule motion makes them structurally ill-suited for cell
 265 routing. Thus, among the architectures considered here, (SANODE) provides the simplest setting where
 266 exact interpolation and generalization coexist.

267 3 Nonparametric estimation

268 In this section, we recall the construction and fundamental properties of piecewise-constant nonparamet-
 269 ric estimators. Our analysis focuses on two canonical approaches:

- 270 1. **Histogram estimators**, which partition the input domain using a fixed uniform grid.
- 271 2. **Nearest-neighbor estimators**, which rely on Voronoi tessellations [4, 7].

272 While the results in this section are standard, proofs are collected in Section 7.2, as the precise statements
 273 we need do not appear verbatim in the literature. We refer the interested reader to [16, 31] and references
 274 therein for a deeper dive into the theory of nonparametric estimators.

275 3.1 Histogram estimator

276 We first analyze the histogram approach. Fix $R > 0$, a resolution $h \in (0, 2R]$ and let

$$\mathcal{K} := \{0, \dots, \lceil 2R/h \rceil - 1\}^d \quad (14)$$

277 be the corresponding multi-index set. For each $k \in \mathcal{K}$, we define the cube

$$Q_k := I_R \cap \prod_{j=1}^d [-R + k_j h, -R + (k_j + 1)h), \quad (15)$$

278 ensuring that the collection $\{Q_k\}_{k \in \mathcal{K}}$ forms a partition of I_R up to a μ -null boundary set. The total
 279 number of cells is given by

$$K_h := |\mathcal{K}| = \left\lceil \frac{2R}{h} \right\rceil^d \lesssim h^{-d}, \quad (16)$$

280 and the diameter of each cell satisfies $\text{diam}(Q_k) \leq \sqrt{d}h$. We approximate $y(\cdot)$ by averaging its values
 281 within each cell. This leads to the definition of both a population-level and an empirical estimator:

- 282 • The population average $y_h: I_R \rightarrow \mathbb{R}^d$ is the best piecewise-constant approximation of y in $L^2(\mu)$:

$$y_h(\mathbf{x}) := \sum_{k \in \mathcal{K}} \mathbf{s}_k \mathbf{1}_{Q_k}(\mathbf{x}), \quad (17)$$

283 where the coefficients $\mathbf{s}_k \in \mathbb{R}^d$ are defined by

$$\mathbf{s}_k := \begin{cases} \frac{1}{\mu(Q_k)} \int_{Q_k} y(\mathbf{x}) \, d\mu(\mathbf{x}), & \text{if } \mu(Q_k) > 0, \\ 0, & \text{if } \mu(Q_k) = 0. \end{cases} \quad (18)$$

284 Note that $y_h(\cdot)$ is in fact the best piecewise-constant approximation of $y(\cdot)$ in $L^2(\mu)$, as each \mathbf{s}_k
 285 minimizes the local $L^2(\mu)$ -error over Q_k .

- 286 • The empirical average $y_{N,h}: I_R \rightarrow \mathbb{R}^d$, instead, is computed directly from the dataset:

$$y_{N,h}(\mathbf{x}) := \sum_{k \in \mathcal{K}} \mathbf{S}_k \mathbf{1}_{Q_k}(\mathbf{x}). \quad (19)$$

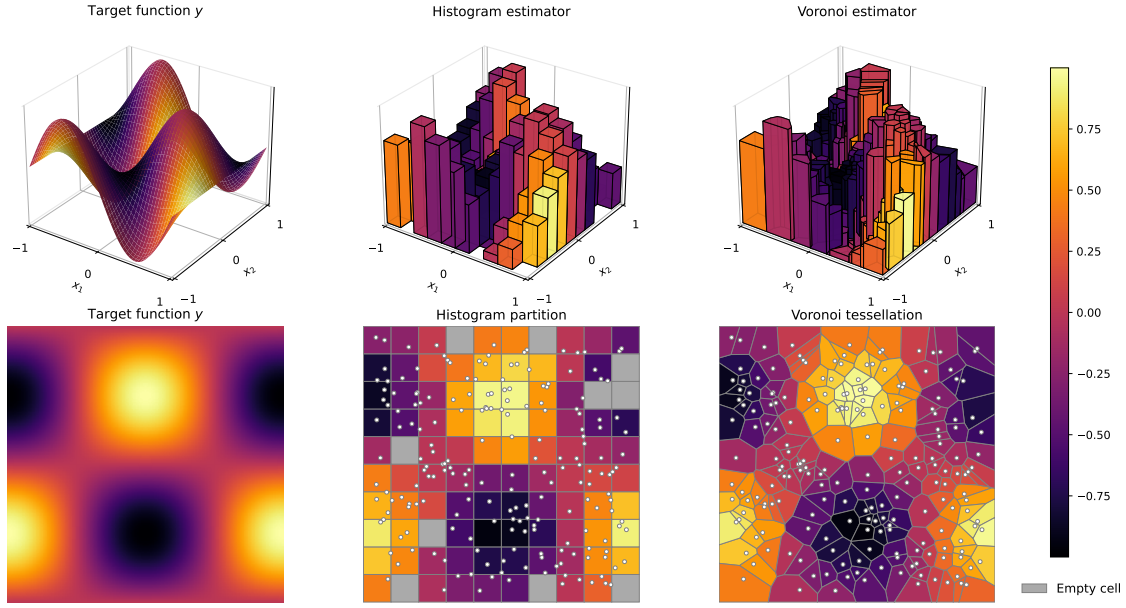


Figure 3: Histogram and Voronoi approximations of a function $y : \mathbb{R}^2 \rightarrow \mathbb{R}^2$, represented in both 3D (height and color) and 2D (color only). **Left:** The target function. **Center:** The histogram estimator, which assigns each cell the average value of its enclosed training points; gray cells indicate regions with no training data. **Right:** The Voronoi estimator, which assigns each cell the value of its nearest training point.

287 Letting $N_k := \sum_{i=1}^N \mathbf{1}_{Q_k}(\mathbf{x}_i)$ denote the number of samples falling into Q_k , the coefficients $\mathbf{S}_k \in \mathbb{R}^d$
 288 are now defined as

$$\mathbf{S}_k := \begin{cases} \frac{1}{N_k} \sum_{i: \mathbf{x}_i \in Q_k} y(\mathbf{x}_i), & \text{if } N_k \geq 1, \\ 0, & \text{if } N_k = 0. \end{cases} \quad (20)$$

289 A schematic representation of $y_{N,h}(\cdot)$ when $d = 2$ can be seen in Figure 3 (central column). We now
 290 establish the convergence rate for the empirical histogram estimator.

291 **Proposition 3.1.** Let $\mathcal{D}_N = \{(\mathbf{x}_i, y(\mathbf{x}_i))\}_{i=1}^N$ with $\mathbf{x}_1, \dots, \mathbf{x}_N$ drawn i.i.d. from $\mu \in \mathcal{P}_{ac}(I_R)$, and $y \in$
 292 $\mathcal{C}^0(I_R; \mathbb{R}^d)$. Then, for any $h \in (0, 2R]$, the estimator $y_{N,h}(\cdot)$ defined in (19) satisfies

$$\mathbb{E}_{\mathcal{D}_N} \left[\|y - y_{N,h}\|_{L^2(\mu)}^2 \right] \lesssim \omega_y(\sqrt{d}h)^2 + \frac{\|y\|_{L^\infty(I_R)}^2}{Nh^d}, \quad (21)$$

293 where $\omega_y(t) := \sup \{\|y(\mathbf{x}) - y(\mathbf{x}')\| : \mathbf{x}, \mathbf{x}' \in I_R, \|\mathbf{x} - \mathbf{x}'\| \leq t\}$ is the modulus of continuity of $y(\cdot)$.

294 The bound (21) captures the standard bias-variance trade-off in nonparametric estimation [16]. The
 295 bias term $\omega_y(\sqrt{d}h)^2$ measures the spatial approximation error and vanishes as $h \rightarrow 0$. Conversely, the
 296 statistical variance scales as $1/(Nh^d)$; as the grid becomes finer, fewer samples fall into each cell. Driving
 297 the total error to zero therefore requires the dataset size N to grow strictly faster than the number of cells.

298 The continuity assumption on y can be relaxed at the cost of losing an explicit algebraic decay rate
 299 for the bias. If $y \in L^\infty(I_R; \mathbb{R}^d)$, the bound (21) holds with the bias replaced by $\|y - y_h\|_{L^2(\mu)}^2$, which
 300 still vanishes as $h \rightarrow 0$ (e.g., via standard martingale convergence on nested grids [13, Theorem 4.4.6]).
 301 For $y \in L^2(\mu; \mathbb{R}^d)$, the bias similarly vanishes, though extending the $(Nh^d)^{-1}$ variance bound generally
 302 requires further assumptions. We do not pursue this level of generality here.

Conversely, by imposing stronger regularity on $y(\cdot)$ one can quantify the modulus of continuity. We
 focus on the Hölder case, characterized by

$$\omega_y(t) \leq L_y t^\alpha \quad \text{for some } \alpha \in (0, 1] \text{ and } L_y > 0.$$

303 Optimizing the grid resolution h yields the following minimal statistical rate.

304 **Corollary 3.2.** Suppose $y \in \mathcal{C}^{0,\alpha}(I_R; \mathbb{R}^d)$ for some $\alpha \in (0, 1]$. Then, the estimator y_h from (17) satisfies

$$\|y - y_h\|_{L^2(\mu)}^2 \lesssim h^{2\alpha}. \quad (22)$$

305 Moreover, $h = (2R) \wedge (N^{-\frac{1}{2\alpha+d}}) =: h_N$ yields the minimal convergence rate in (21),

$$\mathbb{E}_{\mathcal{D}_N} \left[\|y - y_{N,h_N}\|_{L^2(\mu)}^2 \right] \lesssim N^{-\frac{2\alpha}{2\alpha+d}}. \quad (23)$$

306 **Remark 3.3** (Sample complexity). The approximation by Riemann sums on a fixed grid of size h yields (22).
 307 To ensure a squared error below $\varepsilon > 0$, the number of cells must scale as $K_h \gtrsim \varepsilon^{-\frac{d}{2\alpha}}$. Consequently, the
 308 sample complexity to guarantee an expected risk below ε in (23) scales as:

$$N \gtrsim \varepsilon^{-\frac{2\alpha+d}{2\alpha}} = \varepsilon^{-1} \cdot \varepsilon^{-\frac{d}{2\alpha}}. \quad (24)$$

309 We observe that the sample complexity naturally decomposes into the product of the Monte Carlo rate ε^{-1}
 310 (associated with estimating a mean with fixed variance) and a geometric factor $\varepsilon^{-d/(2\alpha)}$. This reflects the
 311 curse of dimensionality: whereas numerical integration computes a single global average, function recovery
 312 requires estimating K_h distinct local averages simultaneously to resolve the spatial structure of $y(\cdot)$ in \mathbb{R}^d .

Remark 3.4 (Minimax optimality). In terms of the number of cells $K_h \asymp h^{-d}$, the bias estimate (22) yields

$$\|y - y_h\|_{L^2(\mu)} \lesssim K_h^{-\alpha/d}.$$

313 This matches the optimal rate among all linear approximation spaces of dimension K_h , as captured by the
 314 Kolmogorov n -width

$$d_n(\mathcal{F}) := \inf_{\dim(V)=n} \sup_{f \in \mathcal{F}} \inf_{g \in V} \|f - g\|_{L^2(\mu)}.$$

315 For the isotropic $\mathcal{C}^{0,\alpha}$ -Hölder ball \mathcal{F} , one has $d_n(\mathcal{F}) \asymp n^{-\alpha/d}$ [25], an order attained by uniform piecewise-
 316 constant grids such as y_h . Moreover, this exponent cannot be improved even by adaptive (nonlinear) methods
 317 with continuous parameter selection: the corresponding manifold n -widths exhibit the exact same decay rate
 318 on Besov (and hence Hölder) balls [12, Eq. (9.4)]. Adaptivity can only yield strictly faster rates on smaller,
 319 spatially inhomogeneous target classes [11].

320 The rate $N^{-\frac{2\alpha}{2\alpha+d}}$ matches the minimax lower bound for nonparametric regression under additive
 321 label noise [30]. When $\mathbf{y}_i = y(\mathbf{x}_i) + \text{noise}$, individual data points are unreliable, making local spatial
 322 averaging the optimal strategy to cancel out stochastic fluctuations.

323 In contrast, our setting is strictly *noiseless* because $\mathbf{y}_i = y(\mathbf{x}_i)$. Here, spatial averaging is overly
 324 conservative and wastes the precision of perfect data. Enforcing *strict interpolation* bypasses the statistical
 325 variance penalty, yielding faster rates. For instance, assuming $y \in \mathcal{C}^{k,\alpha}(I_R; \mathbb{R}^d)$ with $k \geq 0$ integer and
 326 $\alpha \in (0, 1]$, spline interpolants m_N achieve [6]

$$\|y - m_N\|_{L^2(\mu)}^2 \lesssim \|y - m_N\|_{L^\infty(I_R)}^2 \lesssim \left(\frac{\log N}{N} \right)^{\frac{2(k+\alpha)}{d}}. \quad (25)$$

327 To exploit the lack of noise, we turn to the simplest interpolation scheme: the nearest-neighbor estimator.

328 3.2 Nearest-neighbor estimator

329 We now partition I_R by assigning each point to its nearest neighbor within $\{\mathbf{x}_i\}_{i=1}^N$. Unlike the histogram
 330 approach, where the grid is fixed and independent of the data, this partition naturally adapts to the local
 331 density of \mathcal{D}_N : dense regions produce smaller cells, while sparse regions yield larger ones. This adaptivity
 332 is exactly what yields faster approximation rates.

333 Formally, let $\{V_i\}_{i=1}^N$ denote the covering of I_R by the closed Voronoi cells generated by $\{\mathbf{x}_i\}_{i=1}^N$:

$$V_i := \{\mathbf{x} \in I_R : \|\mathbf{x} - \mathbf{x}_i\| \leq \|\mathbf{x} - \mathbf{x}_j\| \text{ for all } j \neq i\}. \quad (26)$$

334 For any $\mathbf{x} \in I_R$, we define its nearest neighbor within the dataset by the function:

$$x_{\text{NN}}(\mathbf{x}) := \arg \min_{\mathbf{x}_i \in \{\mathbf{x}_1, \dots, \mathbf{x}_N\}} \|\mathbf{x} - \mathbf{x}_i\|. \quad (27)$$

335 where any ties on the cell boundaries are broken arbitrarily. Using this mapping, we define the nearest-
336 neighbor interpolant of y , denoted by y_N^V , as

$$y_N^V(\mathbf{x}) := y(x_{\text{NN}}(\mathbf{x})) = \sum_{i=1}^N y(\mathbf{x}_i) \mathbf{1}_{V_i}(\mathbf{x}). \quad (28)$$

337 A schematic representation of y_N^V for $d = 2$ can be seen in the rightmost column of Figure 3. The accuracy
338 of y_N^V is governed purely by the covering radius of \mathcal{D}_N , defined by

$$R_N := \sup_{\mathbf{x} \in I_R} \min_{i \in [N]} \|\mathbf{x} - \mathbf{x}_i\| = \sup_{\mathbf{x} \in I_R} \|\mathbf{x} - x_{\text{NN}}(\mathbf{x})\|. \quad (29)$$

339 The smaller R_N , the closer every point in I_R is to some input \mathbf{x}_i . For i.i.d. samples drawn from a density
340 bounded below on I_R , the expected scaling of R_N is characterized by [26, Theorem 2.1], which yields

$$\mathbb{E}_{\mathcal{D}_N}[R_N^q] \lesssim \left(\frac{\log N}{N}\right)^{q/d} \quad \text{for any } q > 0. \quad (30)$$

341 The following proposition uses this rate to quantify the L^2 -error for Hölder targets:

342 **Proposition 3.5.** *Let $N \geq 2$ and $\mathcal{D}_N = \{(\mathbf{x}_i, y(\mathbf{x}_i))\}_{i=1}^N$ with $\mathbf{x}_1, \dots, \mathbf{x}_N$ drawn i.i.d. from $\mu \in \mathcal{P}_{\text{ac}}(I_R)$.
343 Assume that μ admits a density ρ with $\inf_{\mathbf{x} \in I_R} \rho(\mathbf{x}) > 0$, and let $y \in \mathcal{C}^{0,\alpha}(I_R; \mathbb{R}^d)$ for some $\alpha \in (0, 1]$.
344 Then the nearest-neighbor interpolant y_N^V defined in (28) satisfies*

$$\mathbb{E}_{\mathcal{D}_N} \left[\|y - y_N^V\|_{L^2(\mu)}^2 \right] \lesssim \left(\frac{\log N}{N}\right)^{\frac{2\alpha}{d}}. \quad (31)$$

345 Note that the rate in (31) is strictly faster than the corresponding rate of $O\left(N^{-\frac{2\alpha}{2\alpha+d}}\right)$ derived in (23).

346 4 Generalization via simultaneous cell controllability

347 The previous two sections provide the ingredients needed for our subsequent analysis. On the one hand,
348 Section 2 shows that (SANODE) can map prescribed convex cells into small target balls. On the other
349 hand, Section 3 provides nonparametric estimators with explicit statistical risks. We now combine both
350 ingredients to derive generalization bounds.

351 4.1 General bounds

352 To derive quantitative rates, we assume that $\mathbf{x}_1, \dots, \mathbf{x}_N$ are drawn i.i.d. from μ .¹ This induces the expected
353 population risk over all realizations of the training sample:

$$\mathbb{E}_{\mathcal{D}_N} [\mathcal{R}(\theta_N)] := \int_{I_R} \cdots \int_{I_R} \mathcal{R}(\theta_N(\mathcal{D}_N)) \, d\mu(\mathbf{x}_1) \cdots d\mu(\mathbf{x}_N). \quad (32)$$

354 Our approach proceeds in two steps. First, we construct from \mathcal{D}_N a piecewise-constant proxy $y_N(\cdot)$ that
355 estimates $y(\cdot)$. Second, we design a control θ_N such that $\Phi_T^{\theta_N}$ approximates y_N on most of the domain.

Let $\{P_k\}_{k=1}^K$ be a partition of I_R into convex cells, let $\{\mathbf{r}_k\}_{k=1}^K \subset \mathbb{R}^d$ be target values computed from
 \mathcal{D}_N (e.g., via cell averages), and define the piecewise-constant estimator

$$y_N(\mathbf{x}) := \sum_{k=1}^K \mathbf{r}_k \mathbf{1}_{P_k}(\mathbf{x}).$$

356 Because Φ_T^θ is a homeomorphism, it cannot uniformly approximate the jump discontinuities of y_N . We
357 therefore restrict the accuracy requirements to regions safely away from the cell boundaries. For a margin
358 $\delta > 0$, we define the trimmed cores P_k^δ and the boundary layer Ω_δ (illustrated in Figure 4):

$$P_k^\delta := \{\mathbf{x} \in P_k : \text{dist}(\mathbf{x}, \partial P_k) \geq \delta\}, \quad \Omega_\delta := I_R \setminus \bigcup_{k=1}^K P_k^\delta. \quad (33)$$

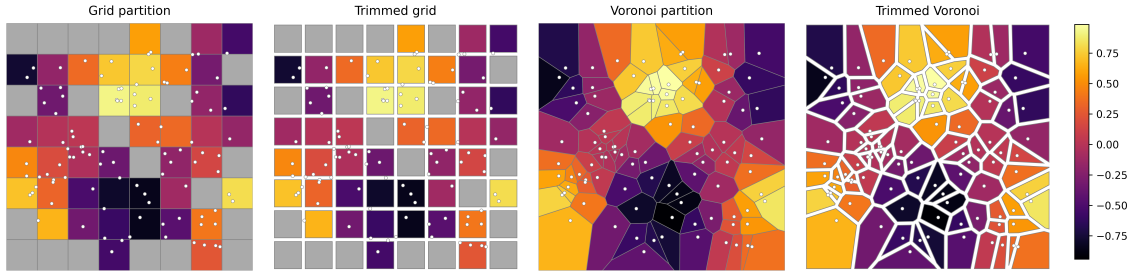


Figure 4: Comparison of grid and Voronoi partitions on the same dataset of $N = 100$ training points. **Left to right:** grid partition, trimmed grid partition, Voronoi partition, and trimmed Voronoi partition. In the trimmed cases, each cell is eroded inward by a margin $\delta > 0$, meaning points within distance δ of the boundary are removed to produce disjoint compact cores. The white regions indicate these removed boundary layers.

359 Since $y_N \equiv \mathbf{r}_k$ on each core P_k^δ , any control θ satisfying condition (10) up to tolerance η on these
 360 cores yields a natural decomposition of the error. This leads to the following model-agnostic bound.

361 **Theorem 4.1** (Template bound). Assume $d \geq 2$ and that the flow Φ_T^θ of (5) is cell-wise controllable in
 362 the sense of Definition 2. Let $\mathcal{D}_N = \{(\mathbf{x}_i, y(\mathbf{x}_i))\}_{i=1}^N$ with $\mathbf{x}_1, \dots, \mathbf{x}_N$ drawn i.i.d. from $\mu \in \mathcal{P}_{ac}(I_R)$, let
 363 $y \in \mathcal{C}^0(I_R; \mathbb{R}^d)$, and let y_N be a piecewise-constant estimator on $\{P_k\}_{k=1}^K$. Then, for any margin $\delta > 0$ and
 364 tolerance $\eta > 0$, there exists a data-dependent control $\theta_N = \theta_N(\mathcal{D}_N)$ such that

$$\mathbb{E}_{\mathcal{D}_N} [\mathcal{R}(\theta_N)] \lesssim \mathbb{E}_{\mathcal{D}_N} [\|y - y_N\|_{L^2(\mu)}^2] + \mathbb{E}_{\mathcal{D}_N} [\mu(\Omega_\delta)] + \eta^2, \quad (34)$$

365 provided $\Phi_T^{\theta_N}$ remains uniformly bounded on I_R .

Proof. By the triangle inequality, for every control θ we have

$$\|\Phi_T^\theta - y\|_{L^2(\mu)}^2 \leq 2\|y - y_N\|_{L^2(\mu)}^2 + 2\|y_N - \Phi_T^\theta\|_{L^2(\mu)}^2.$$

To bound the realization error, we split the domain into the trimmed cores and the boundary layer:

$$\begin{aligned} \|y_N - \Phi_T^\theta\|_{L^2(\mu)}^2 &\leq \sum_{k=1}^K \mu(P_k^\delta) \sup_{\mathbf{x} \in P_k^\delta} \|\Phi_T^\theta(\mathbf{x}) - \mathbf{r}_k\|^2 + \int_{\Omega_\delta} \|y_N(\mathbf{x}) - \Phi_T^\theta(\mathbf{x})\|^2 d\mu(\mathbf{x}) \\ &\leq \eta^2 + 2\mu(\Omega_\delta) \left(\|y_N\|_{L^\infty(I_R)}^2 + \|\Phi_T^\theta\|_{L^\infty(I_R)}^2 \right), \end{aligned}$$

366 where the first term exploits the condition of Definition 2. By the uniform boundedness of both the
 367 estimator and the flow, taking the expectation over \mathcal{D}_N yields the result. \square

368 The template bound (34) cleanly isolates the three sources of error. The statistical error depends enti-
 369 rely on the chosen proxy y_N , the control error is dictated by the network capacity ($\eta \rightarrow 0$ as width
 370 grows), and the geometric error depends on the volume of the boundary layer Ω_δ . We treat the two proxy
 371 partitions of interest separately.

- 372 • **Uniform grid.** The boundary layer Ω_δ is simply a union of δ -strips along the fixed cubic cell faces.
- 373 • **Voronoi.** The boundary layer Ω_δ is the δ -neighborhood of the Voronoi skeleton $\Sigma_N := \bigcup_{i=1}^N \partial V_i$.
 374 Because the cells V_i are intersections of half-spaces, they are convex polytopes, ensuring that the
 375 trimmed cores V_i^δ are pairwise disjoint compact convex sets. This allows Definition 2 to be applied
 376 directly, even though Σ_N depends on the random spatial configuration of the sample points.

377 The next two lemmas establish quantitative bounds on the expected measure of Ω_δ for each type.

378 **Lemma 4.2.** For $\mu \in \mathcal{P}_{ac}(I_R)$ with density bounded above, any partition $\{Q_k\}$ as in (15) with $h \in (0, 2R]$,
 379 and any margin $\delta \in (0, h/2)$, the boundary layer Ω_δ satisfies

$$\mu(\Omega_\delta) \lesssim \frac{\delta}{h}. \quad (35)$$

¹The deterministic core of the argument below, however, is purely geometric and applies to any fixed dataset.

380 **Lemma 4.3.** For any integers $d \geq 2$ and $N \geq 1$, any $\mu \in \mathcal{P}_{ac}(I_R)$ with density bounded above and below
 381 by positive constants, and any sufficiently small $\delta > 0$, the boundary layer Ω_δ generated by $\{\mathbf{x}_i\}_{i=1}^N$ satisfies

$$382 \quad \mathbb{E}_{\mathcal{D}_N} [\mu(\Omega_\delta)] \lesssim \delta N^{1/d}. \quad (36)$$

383 Substituting these geometric bounds and the statistical risks from Section 3 into the decomposition
 384 (34) yields explicit rates for any system possessing the SCC property.

385 **Corollary 4.4.** Under the hypotheses of Theorem 4.1, suppose that $y \in \mathcal{C}^{0,\alpha}(I_R; \mathbb{R}^d)$ for some $\alpha \in (0, 1]$.
 386 Furthermore, for the Voronoi estimator, assume that μ admits a strictly positive density on I_R . Then,

$$\mathbb{E}_{\mathcal{D}_N} [\mathcal{R}(\theta_N)] \lesssim \begin{cases} h^{2\alpha} + \frac{1}{Nh^d} + \frac{\delta}{h} + \eta^2 & \text{(Histogram),} \\ \left(\frac{\log N}{N}\right)^{2\alpha/d} + \delta N^{1/d} + \eta^2 & \text{(Voronoi).} \end{cases} \quad (37)$$

387 4.2 Nonparametric rates with semi-autonomous neural ODEs

388 We now specialize these bounds to (SANODE). By Theorem 2.5, this model satisfies the SCC and admits
 389 realizing flows that remain uniformly bounded on I_R . To match the baseline rates, it suffices to balance
 390 the statistical, geometric, and control errors in (34) by appropriately scaling the margin δ , the tolerance η ,
 391 and the width p_N . We first state the result for the histogram partition.

392 **Proposition 4.5** (Histogram rate). Let $d \geq 2$, let $y \in \mathcal{C}^{0,\alpha}(I_R; \mathbb{R}^d)$ for some $\alpha \in (0, 1]$, and let $\mathbf{x}_1, \dots, \mathbf{x}_N$
 393 be drawn i.i.d. from $\mu \in \mathcal{P}_{ac}(I_R)$ with bounded density. Define

$$h_N := N^{-\frac{1}{2\alpha+d}}, \quad \delta_N := h_N^{1+2\alpha}, \quad (38)$$

394 and let p_N satisfy condition (11) for the histogram partition with cells $\{Q_k\}$ of side h_N and margin δ_N . Then,
 395 for all sufficiently large N , there exists a data-dependent control θ_N of width p_N for (SANODE) such that

$$\mathbb{E}_{\mathcal{D}_N} [\mathcal{R}(\theta_N)] \lesssim N^{-\frac{2\alpha}{2\alpha+d}}.$$

396 Next, we establish the analogous result for the faster Voronoi partition.

397 **Proposition 4.6** (Voronoi rate). Let $d \geq 2$, let $y \in \mathcal{C}^{0,\alpha}(I_R; \mathbb{R}^d)$ for some $\alpha \in (0, 1]$, and let $\mathbf{x}_1, \dots, \mathbf{x}_N$
 398 be drawn i.i.d. from $\mu \in \mathcal{P}_{ac}(I_R)$ with density bounded above and below on I_R . Define

$$\delta_N := (\log N)^{\frac{2\alpha}{d}} N^{-\frac{2\alpha+1}{d}}, \quad (39)$$

399 and let p_N satisfy condition (11) for the Voronoi partition generated by $\{\mathbf{x}_i\}_{i=1}^N$ with margin δ_N . Then, for
 400 all sufficiently large N , there exists a data-dependent control θ_N of width p_N for (SANODE) such that

$$\mathbb{E}_{\mathcal{D}_N} [\mathcal{R}(\theta_N)] \lesssim \left(\frac{\log N}{N}\right)^{\frac{2\alpha}{d}}.$$

401 The proofs of both propositions are provided in Section 7.3. In both cases, the width p_N ensures the
 402 existence of the required flow for each N . This required width scaling is directly tied to the geometry of
 403 the chosen partition—specifically, the minimum cell separation and maximum cell diameter. As partitions
 404 refine, these worst-case geometric constraints drive the super-polynomial growth discussed in Remark 2.6.

405 5 Numerical experiments

406 We complement the theoretical analysis of Section 4 with numerical experiments on (SANODE). The first
 407 experiment investigates whether the trained model can match nonparametric estimator baselines at a
 408 network width well below the theoretical prescription. This tests the sharpness of the sufficient width
 409 scaling. The second demonstrates that time-dependence of the vector field is necessary for topologically

non-trivial tasks, providing empirical support for the autonomous/non-autonomous separation discussed in Section 2.3. The code to reproduce these experiments is available in the GitHub repository [ExGen](#).

Throughout, training points are sampled uniformly from $I_R = [-2, 2]^d$, and the population risk is estimated on a fixed test set using the squared ℓ^2 -loss. All models are trained with the Adam optimizer at a fixed learning rate of 10^{-4} . Training ends after a maximum number of epochs (specified per experiment), or earlier if the loss falls below 10^{-6} or fails to improve by more than 10^{-8} over 5×10^3 consecutive steps.

5.1 Width scaling

The width prescriptions derived in Section 4 scale explosively with the dimension d , much like the storage costs of the nonparametric estimators they are calibrated against (see also Remark 2.7). In practice, however, neural ODEs are inherently nonlinear function approximators and may exploit this structure to achieve the same risk with far fewer parameters. We investigate to what extent this is the case by fixing the dataset size N and sweeping the network width p to determine the minimal width at which (SANODE) matches classical nonparametric baselines.

We consider two target functions $f, g: \mathbb{R}^d \rightarrow \mathbb{R}^d$:

- **Smooth target** ($\alpha = 1$): $f(\mathbf{x})^{(2k-1)} = \sin(x^{(2k-1)}) \cos(x^{(2k)})$ and $f(\mathbf{x})^{(2k)} = \cos(x^{(2k-1)}) \sin(x^{(2k)})$ for $k = 1, \dots, \lfloor d/2 \rfloor$, with $f(\mathbf{x})^{(d)} = x^{(d)}$ if d is odd.

- **Hölder-1/2 target** ($\alpha = 1/2$): $g(\mathbf{x})^{(i)} = \text{sgn}(x^{(i)}) \sqrt{|x^{(i)}|}$ for $i \in [d]$.

We study two regimes: a low-dimensional setting with $d = 3$ and $N = 500$ training points, and a high-dimensional one where $d = 8$ with $N = 5000$ training points. The larger dataset for $d = 8$ compensates for the increased difficulty of the estimation problem in high dimension. The width is varied over a logarithmic grid $p \in \{5, 11, 23, 51, 109, 237, 512, 1024\}$. Training runs for at most 3×10^4 gradient steps, and the population risk is estimated on a fixed test set of 2×10^3 points. For reference, the histogram and Voronoi estimator errors are computed at the same (N, d) and displayed as horizontal baselines. The histogram estimator uses the optimal bandwidth $h_N = N^{-1/(2\alpha+d)}$, prescribed by Corollary 3.2.

In order to compare the complexity of (SANODE) with that of the nonparametric estimators, we report for each model the number of stored scalar values. For (SANODE) this coincides with the number of trainable parameters. For the histogram estimator this is at most Nd (at most N cells can be occupied, each storing a label in \mathbb{R}^d), and for the Voronoi estimator it is $2Nd$ (N input-label pairs, each in \mathbb{R}^d).

Figure 5 reports the test risk as a function of p for both targets and both dimensions. The quantitative results are collected in Table 1.

Case $d = 3$. For the smooth target, the (SANODE) already surpasses the histogram estimator at the smallest tested width ($p = 5$, 43 parameters) and converges to Voronoi-level error around $p = 512$ (4099 parameters), compared to the 3000 values stored by the Voronoi estimator. The crossover thus occurs at a parameter count comparable to the baseline, well below the theoretical prescription. The behavior on the Hölder-1/2 target is markedly different. Already at $p = 23$ (187 parameters), the (SANODE) achieves a test MSE of 1.2×10^{-4} , more than two orders of magnitude below the Voronoi baseline (2.8×10^{-2}), and the error plateaus around $p = 50$. The histogram estimator, with 500 stored values, achieves a test MSE far worse than the (SANODE) even at $p = 5$. This gap suggests that the (SANODE) exploits the component-wise, sign-symmetric structure of g , which the histogram partition cannot capture, and that the conservative worst-case width prescription is particularly loose for structured targets.

Case $d = 8$. The (SANODE) seems to do even better as the dimension increases, surpassing the histogram baseline (2.42×10^{-1}) around $p = 51$ (926 parameters) and approaching Voronoi-level error (8.95×10^{-2}) around $p = 237$ (4274 parameters), with almost an order of magnitude fewer parameters.

The picture for the Hölder-1/2 target is similar to the one in the low-dimensional setting. Despite the much higher dimension, the (SANODE) beats the histogram estimator (9.98×10^{-1}) already at $p = 5$ and drops below the Voronoi baseline (1.40×10^{-1}) between $p = 23$ and $p = 51$, reaching a plateau of around 10^{-5} from $p = 109$ onward, four orders of magnitude below the Voronoi baseline. This further confirms that the (SANODE) exploits the component-wise structure of g in a way that is largely insensitive to dimension, while the classical estimators degrade rapidly.

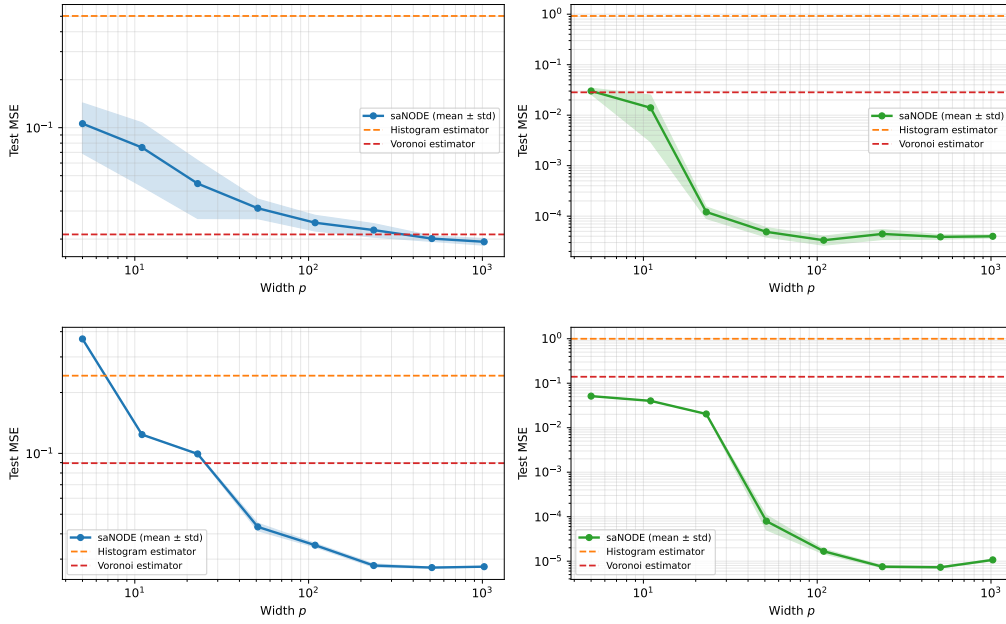


Figure 5: Test risk of (SANODE) as a function of network width p for the smooth target f (left) and the Hölder-1/2 target g (right), at $d = 3$, $N = 500$ (top) and $d = 8$, $N = 5000$ (bottom). Solid curves show the mean over 3 independent seeds and shaded bands indicate ± 1 standard deviation. Horizontal dashed lines mark the histogram and Voronoi estimator errors.

Table 1: Test MSE of (SANODE) versus p , averaged over 3 seeds. Complexity denotes the number of scalar degrees of freedom: trainable parameters for (SANODE), and Nd (histogram) or $2Nd$ (Voronoi) stored scalars for the baselines.

$d = 3, N = 500$					$d = 8, N = 5000$				
Method	p	Compl.	Smooth	Hölder-1/2	Method	p	Compl.	Smooth	Hölder-1/2
Histogram	–	1500	5.02×10^{-1}	9.22×10^{-1}	Histogram	–	40000	2.42×10^{-1}	9.98×10^{-1}
Voronoi	–	3000	2.14×10^{-2}	2.83×10^{-2}	Voronoi	–	80000	8.95×10^{-2}	1.40×10^{-1}
(SANODE)	5	43	1.06×10^{-1}	3.02×10^{-2}	(SANODE)	5	98	3.68×10^{-1}	5.13×10^{-2}
	11	91	7.51×10^{-2}	1.40×10^{-2}		11	206	1.24×10^{-1}	4.01×10^{-2}
	23	187	4.46×10^{-2}	1.21×10^{-4}		23	422	9.94×10^{-2}	2.04×10^{-2}
	51	411	3.12×10^{-2}	4.88×10^{-5}		51	926	4.33×10^{-2}	7.94×10^{-5}
	109	875	2.53×10^{-2}	3.33×10^{-5}		109	1970	3.52×10^{-2}	1.67×10^{-5}
	237	1899	2.28×10^{-2}	4.44×10^{-5}		237	4274	2.79×10^{-2}	7.51×10^{-6}
	512	4099	2.01×10^{-2}	3.88×10^{-5}		512	9224	2.72×10^{-2}	7.29×10^{-6}
	1024	8195	1.92×10^{-2}	3.99×10^{-5}		1024	18440	2.75×10^{-2}	1.07×10^{-5}

5.2 Necessity of time-dependence

As discussed in Section 2.3, because macroscopic trajectories cannot cross, autonomous flows face structural obstructions to cell routing whenever the target configuration is topologically incompatible with a continuous deformation of the input. We now provide direct empirical evidence for this limitation.

We consider a checkerboard sorting problem: given a $K \times K$ partition of $[-1, 1]^2$, the target assigns each point to $(+S, +S)$ if its cell indices (i, j) have an even sum, and to $(-S, -S)$ otherwise, setting $S = 0.7$. For $K \geq 2$, the alternating initial regions are heavily interleaved, forcing paths to cross to reach their respective targets. Time-dependent models like (SANODE) circumvent this spatial bottleneck by decoupling the motion in the extended space-time domain.

We compare (SANODE) against its autonomous counterpart (ANODE). To ensure a fair comparison, the width of (ANODE) is chosen so that the total parameter count matches the (SANODE):

$$\begin{cases} (\text{SANODE}) \text{ params} = (2d + 2)p + d \\ (\text{ANODE}) \text{ params} = (2d + 1)p + d \end{cases} \implies p_{\text{ANODE}} = \left\lceil \frac{2d + 2}{2d + 1} p_{\text{SANODE}} \right\rceil.$$

Both models are trained on $N = 1000$ points with identical optimizer settings for at most 5×10^4 gradient steps, and evaluated on a separate test set of 10^3 points. We sweep the resolution $K \in \{2, 3, 4\}$.

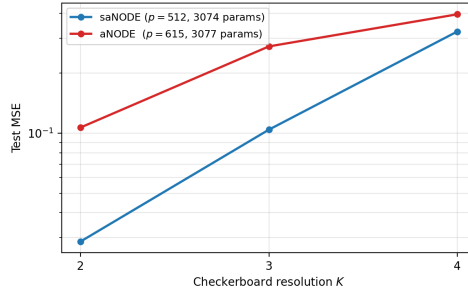


Figure 6: Test risk of (saNODE) and parameter-matched (aNODE) on the checkerboard sorting task with respect to K .

Figure 6 reports the test risk as a function of the grid resolution K . As expected, we observe a growing gap in test error between the two models. This behavior is consistent with the theoretical obstruction: the autonomous flow cannot cleanly resolve the increasingly fine interleaving of the cells.

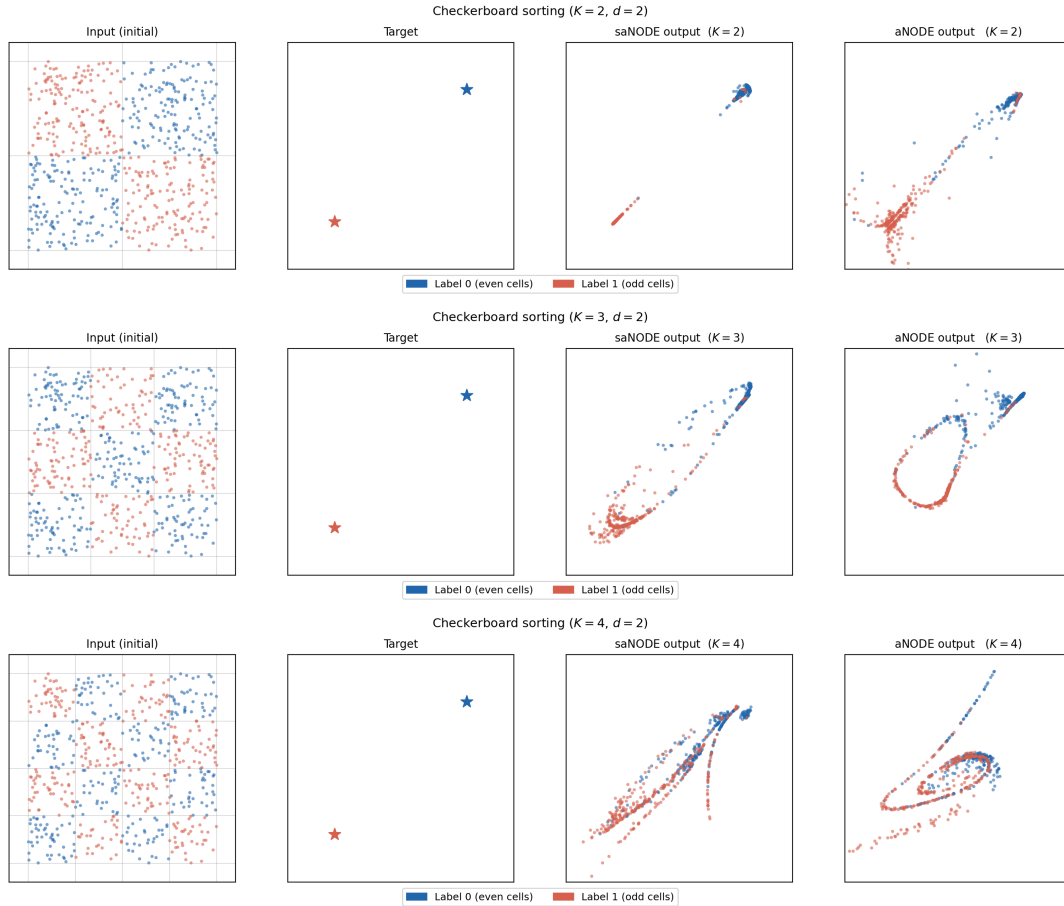


Figure 7: Checkerboard sorting for $d = 2$ at $K = 2$ (top), $K = 3$ (middle), and $K = 4$ (bottom). Each row shows, from left to right: input points, target assignment, (saNODE) output, and (aNODE) output. Points are colored by cell parity (even cells in blue, odd cells in red). The (saNODE) is able to separate the two classes at all resolutions, losing some precision only when $K = 4$. By contrast, the (aNODE) output becomes increasingly entangled as K grows.

Figure 7 provides a direct visualization of the learned mappings. The (saNODE) successfully routes the two classes into distinct clusters (showing only a minor loss of precision at $K = 4$). By contrast, the (aNODE) output remains entangled, with points from opposite classes overlapping. As the resolution in-

478 creases, this degradation becomes progressively more severe, confirming that the topological obstruction
479 fundamentally limits autonomous models in region-routing tasks.

480 6 Conclusions and perspectives

481 We have developed a controllability-based framework to study generalization in neural ODEs. A central
482 insight of our work is that point interpolation, while necessary for expressivity, cannot fully explain
483 generalization. To bridge this gap, we introduced *simultaneous cell controllability*—the ability of the flow
484 to compress entire input regions toward prescribed targets. We proved that the semi-autonomous model
485 ([SANODE](#)) satisfies this property, enabling it to approximate piecewise-constant nonparametric estimators.

486 This connection yields explicit population-risk bounds. By constructing a single, data-dependent flow
487 that realizes a statistically meaningful proxy, we demonstrate that exact interpolation and quantitative
488 generalization can coexist. Our results thus establish that overparameterized neural ODEs can emulate
489 nonparametric procedures, providing an achievability guarantee rather than an efficient scaling law. This
490 approach departs from standard uniform generalization bounds, which evaluate the worst-case error over
491 an entire hypothesis class and force a rigid trade-off between expressivity and global complexity.

492 Finally, we find that explicit time dependence is essential for our generalization mechanism. Without
493 it, topological obstructions break simultaneous cell controllability, preventing autonomous architectures
494 from routing macroscopic regions even when point interpolation is preserved.

495 We conclude by outlining natural directions opened for future work.

496 **Saturation of regularity and higher-order flows.** The bounds established here rely on piecewise-
497 constant (order-0) approximations, whose statistical bias saturates at Lipschitz continuity ($\alpha = 1$). We
498 deliberately favor this baseline because it naturally extends the exact [UIP](#). Nevertheless, if the ground truth
499 is smoother ($y \in \mathcal{C}^{m,\alpha}$), replacing the order-0 proxy with a degree- m cell-wise polynomial would yield
500 faster statistical rates. Transferring this higher-order rate to Φ_T^θ , however, would require a complex ana-
501 logue of [Definition 2](#) capable of realizing local higher-degree polynomials up to tolerance η . Formulating
502 and proving such higher-order controllability remains an interesting open problem.

503 **Quantitative sharpness.** While our proofs establish the vanishing of population risk, the required
504 network width grows super-polynomially with N due to worst-case geometric bottlenecks in the partition.
505 Closing the gap between these theoretical thresholds and the much smaller widths observed in practice
506 requires sharper quantitative controllability estimates.

507 **From static maps to trajectories.** Finally, our controllability viewpoint suggests a natural extension
508 from static regression to dynamical-system learning. While this work uses neural ODEs to approximate
509 fixed terminal maps, applications often involve continuous sequences. Extending simultaneous cell con-
510 trollability to steer entire trajectory tubes—rather than just terminal cells—would provide a framework
511 for trajectory-level generalization bounds in sequence modeling.

512 7 Proofs

513 7.1 Proofs of [Section 2](#)

514 7.1.1 Proof of [Theorem 2.3](#)

515 We show that ([SANODE](#)) is simultaneously controllable, and we give an explicit construction. In particular,
516 to steer N points to prescribed targets we use at most $p = 2N$ neurons. The proof requires several
517 intermediate results.

518 A first key observation, that we will use throughout this section, is that each triple of parameters
519 (\mathbf{a}_j, b_j, c_j) defines a moving hyperplane

$$H_j(t) := \{x \in \mathbb{R}^d : \mathbf{a}_j \cdot \mathbf{x} + b_j t + c_j = 0\}, \quad (40)$$

520 which translates with constant normal velocity $-b_j \mathbf{a}_j / \|\mathbf{a}_j\|^2$.

521 The following preliminary lemma provides exact controllability of $d - 1$ coordinates while keeping
522 one coordinate fixed.

523 **Lemma 7.1.** *Let $N \geq 1$, $d \geq 2$, and let $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \subset \mathbb{R}^d \times \mathbb{R}^d$ be any admissible finite dataset in the
524 sense of (8) and such that $x_i^{(1)} \neq x_j^{(1)}$ for all $i \neq j$. For any $T > 0$ and any $R > \max_i x_i^{(1)}$, there exists*

$$\theta_R = (\mathbf{w}_i, \mathbf{a}_i, b_i, c_i)_{i=1}^N \in (\mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R} \times \mathbb{R})^N$$

such that the flow map $\Phi_T^{\theta_R}$ of (SANODE) satisfies, for each $i \in [N]$,

$$\left[\Phi_T^{\theta_R}(\mathbf{x}_i) \right]^{(1)} = x_i^{(1)} \quad \text{and} \quad \left[\Phi_T^{\theta_R}(\mathbf{x}_i) \right]^{(k)} = y_i^{(k)} \quad \text{for } k = 2, \dots, d$$

525 and, for each $i \in [N]$, we have $(\mathbf{a}_i \cdot \mathbf{x} + b_i t + c_i)_+ = 0$ for all $t \geq T$ and any $\mathbf{x} \in \mathbb{R}^d$ such that $x^{(1)} \leq R$.

526 Furthermore, for every $R_0 > \max_i x_i^{(1)}$, the family of controls $\{\theta_R\}_{R \geq R_0}$ can be chosen so that

$$\sup_{R \geq R_0} \max_{i \in [N]} \max_{t \in [0, T]} \max_{k=2, \dots, d} \left| \left[\Phi_t^{\theta_R}(\mathbf{x}_i) \right]^{(k)} \right| < \infty. \quad (41)$$

527 *Proof.* Up to relabeling the points, we may assume

$$x_1^{(1)} < x_2^{(1)} < \dots < x_N^{(1)}. \quad (42)$$

528 We set $\mathbf{a}_i = \mathbf{e}_1$, and $w_i^{(1)} = 0$, for all $i \in [N]$. Then (SANODE) becomes

$$\dot{\mathbf{x}}(t) = \sum_{i=1}^N \mathbf{w}_i (x^{(1)}(t) + b_i t + c_i)_+, \quad t \in [0, T],$$

529 so, in components,

$$\dot{x}^{(1)}(t) = 0, \quad \dot{x}^{(k)}(t) = \sum_{i=1}^N w_i^{(k)} (x^{(1)}(t) + b_i t + c_i)_+, \quad k = 2, \dots, d. \quad (43)$$

530 In particular, $\left[\Phi_t^{\theta_R}(\mathbf{x}_i) \right]^{(1)} = x_i^{(1)}$ for all $t \in [0, T]$, so the first coordinate is frozen and (42) is preserved.

531 **Step 1.** Set

$$b_1 = \frac{x_1^{(1)} - 1 - R}{T}, \quad c_1 = 1 - x_1^{(1)}, \quad \text{and} \quad b_i = \frac{x_{i-1}^{(1)} - R}{T}, \quad c_i = -x_{i-1}^{(1)} \quad (i = 2, \dots, N).$$

532 Then each moving hyperplane $H_i(t) := \{x \in \mathbb{R}^d : x^{(1)} + b_i t + c_i = 0\}$ satisfies:

- 533 (i) $H_i(t)$ is orthogonal to \mathbf{e}_1 for all t .
- 534 (ii) For $i = 1$, $H_1(0) : x^{(1)} = x_1^{(1)} - 1$ and $H_1(T) : x^{(1)} = R$.
- 535 (iii) For $i \geq 2$, $H_i(0) : x^{(1)} = x_{i-1}^{(1)}$ and $H_i(T) : x^{(1)} = R$.

536 By the choice of (b_i, c_i) we have $b_i T + c_i = -R$ for all $i \in [N]$. Furthermore, the hyperplanes do not
537 intersect on $[0, T]$.

538 **Step 2.** Fix $i \in [N]$. Since $x^{(1)}(t) \equiv x_i^{(1)}$, for each $k \geq 2$ we have

$$\frac{d}{dt} \left[\Phi_t^{\theta_R}(\mathbf{x}_i) \right]^{(k)} = \sum_{j=1}^N w_j^{(k)} (x_i^{(1)} + b_j t + c_j)_+.$$

539 With the above choice of (b_j, c_j) and (42), if $j > i$ (hence $j \geq 2$), the argument $t \mapsto x_i^{(1)} + b_j t + c_j$ is
540 equal to $x_i^{(1)} - x_{j-1}^{(1)} \leq 0$ at $t = 0$, and equal to $x_i^{(1)} - R < 0$ at $t = T$. By linearity,

$$x_i^{(1)} + b_j t + c_j \leq 0 \quad \text{for all } t \in [0, T],$$

541 so neuron j is never active along the trajectory from \mathbf{x}_i . Therefore,

$$\left[\Phi_T^{\theta_R}(\mathbf{x}_i) \right]^{(k)} = x_i^{(k)} + \sum_{j=1}^i w_j^{(k)} \int_0^T (x_i^{(1)} + b_j t + c_j)_+ dt =: x_i^{(k)} + \sum_{j=1}^i w_j^{(k)} M_{i,j}.$$

542 Conversely, if $j = 1$, the argument at $t = 0$ is $x_i^{(1)} + c_1 = x_i^{(1)} + 1 - x_1^{(1)} > 0$, while for $2 \leq j \leq i$ it is

$$x_i^{(1)} + c_j = x_i^{(1)} - x_{j-1}^{(1)} \geq x_i^{(1)} - x_{i-1}^{(1)} > 0.$$

543 Thus, for all $1 \leq j \leq i$, the integrand is strictly positive at $t = 0$ and continuous, so $M_{i,j} > 0$.

544 **Step 3.** Imposing $\left[\Phi_T^{\theta_R}(\mathbf{x}_i) \right]^{(k)} = y_i^{(k)}$ for $k = 2, \dots, d$ yields, for each fixed k ,

$$y_i^{(k)} - x_i^{(k)} = \sum_{j=1}^i w_j^{(k)} M_{i,j}, \quad i \in [N],$$

545 which is a lower-triangular linear system in $(w_1^{(k)}, \dots, w_N^{(k)})$ with strictly positive diagonal. Hence it has
546 a unique solution for each $k = 2, \dots, d$.

547 Let $t \geq T$ and $\mathbf{x} \in \mathbb{R}^d$ with $x^{(1)} \leq R$. Since $R > \max_i x_i^{(1)}$, Step 1 yields $b_i < 0$ and $b_i T + c_i = -R$
548 for all $i \in [N]$. Thus, for $t \geq T$ it holds that $b_i t + c_i \leq -R$, which implies $(\mathbf{a}_i \cdot \mathbf{x} + b_i t + c_i)_+ = 0$ because

$$\mathbf{a}_i \cdot \mathbf{x} + b_i t + c_i \leq x^{(1)} - R \leq 0.$$

549 Finally, fix $R_0 > \max_i x_i^{(1)}$. The uniform boundedness follows from the explicit triangular construc-
550 tion. Indeed, the entries $M_{i,j}$ are of order R^{-1} as $R \rightarrow \infty$, while the corresponding weights grow at most
551 linearly in R . Hence

$$w_j^{(k)} \int_0^t (x_i^{(1)} + b_j s + c_j)_+ ds$$

552 remain uniformly bounded for $R \geq R_0$, uniformly in i, j, k and $t \in [0, T]$. Since the same quantities
553 depend continuously on R and no $M_{i,i}$ vanishes for $R \geq R_0$, the bound extends to all $R \geq R_0$. \square

554 *Proof of Theorem 2.3.* Since the dataset is finite and admissible, we can apply an arbitrarily small generic
555 rotation to ensure that

$$x_i^{(1)} \neq x_j^{(1)} \quad \text{and} \quad y_i^{(2)} \neq y_j^{(2)} \quad \text{for } i \neq j.$$

556 Up to relabeling the pairs, we assume

$$y_1^{(2)} < y_2^{(2)} < \dots < y_N^{(2)}. \quad (44)$$

557 We build $\theta \in (\mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R} \times \mathbb{R})^{2N}$ by concatenating two sets of N neurons which are active during the
558 two phases given by the time intervals $[0, T/2]$ and $[T/2, T]$.

559 **Step 1.** Fix $R_0 > \max_i x_i^{(1)}$. By the last assertion of Lemma 7.1, applied on the interval $[0, T/2]$, there
560 exists $C_0 > 0$ such that, for every $R_1 \geq R_0$,

$$\max_{i \in [N]} \max_{t \in [0, T/2]} \left| \left[\Phi_t^{\theta_{R_1}}(\mathbf{x}_i) \right]^{(2)} \right| \leq C_0.$$

561 Choose

$$R_2 < (-C_0) \wedge \min_{i \in [N]} y_i^{(2)}. \quad (45)$$

562 **Step 2.** Now, for $\ell = N + 1, \dots, 2N$, set $\mathbf{a}_\ell = -\mathbf{e}_2$ and $\mathbf{w}_\ell^{(k)} = 0$ for all $k = 2, \dots, d$. Choose (b_ℓ, c_ℓ)
563 so that the hyperplanes $\mathbf{H}_\ell(t) := \{x \in \mathbb{R}^d : -x^{(2)} + b_\ell t + c_\ell = 0\}$ all satisfy $\mathbf{H}_\ell(T/2) : x^{(2)} = R_2$. At
564 $t = T$, they must reach the following prescribed levels: for $\ell = N + 1, \dots, 2N - 1$, impose

$$b_\ell \frac{T}{2} + c_\ell = R_2 \quad \text{and} \quad b_\ell T + c_\ell = y_{\ell+1-N}^{(2)},$$

565 and for $\ell = 2N$, impose

$$b_{2N} \frac{T}{2} + c_{2N} = R_2 \quad \text{and} \quad b_{2N} T + c_{2N} = y_N^{(2)} + 1.$$

566 Now consider the auxiliary dynamics on $[T/2, T]$ generated only by the second block, starting from

$$\mathbf{z}_i = (x_i^{(1)}, y_i^{(2)}, \dots, y_i^{(d)}).$$

567 Let $\Psi_i(t)$ denote the corresponding trajectory. Since $w_\ell^{(2)} = 0$ for all $\ell \geq N+1$ we have $[\Psi_i(t)]^{(2)} \equiv y_i^{(2)}$
568 on this interval, and the first coordinate satisfies

$$\frac{d}{dt} [\Psi_i(t)]^{(1)} = \sum_{\ell=N+1}^{2N} w_\ell^{(1)} (-y_i^{(2)} + b_\ell t + c_\ell)_+.$$

569 Using (44) and the fact that $b_\ell t + c_\ell \leq b_\ell T + c_\ell$ for $t \in [T/2, T]$, we deduce that for $\ell = N+1, \dots, N+i-1$,

$$b_\ell t + c_\ell \leq b_\ell T + c_\ell = y_{\ell+1-N}^{(2)} \leq y_i^{(2)},$$

570 meaning $(-y_i^{(2)} + b_\ell t + c_\ell)_+ \equiv 0$ on $[T/2, T]$. Thus, only the neurons with $\ell \geq N+i$ contribute, yielding

$$[\Psi_i(T)]^{(1)} = x_i^{(1)} + \sum_{j=i}^N w_{N+j}^{(1)} \int_{T/2}^T (-y_i^{(2)} + b_{N+j} t + c_{N+j})_+ dt =: x_i^{(1)} + \sum_{j=i}^N w_{N+j}^{(1)} \widetilde{M}_{i,j}.$$

571 By construction, $\widetilde{M}_{i,j} > 0$ for all $1 \leq i \leq j \leq N$. Imposing $[\Psi_i(T)]^{(1)} = y_i^{(1)}$ yields the upper-triangular
572 linear system

$$y_i^{(1)} - x_i^{(1)} = \sum_{j=i}^N w_{N+j}^{(1)} \widetilde{M}_{i,j}, \quad i \in [N],$$

573 which has a unique solution since the diagonal entries $\widetilde{M}_{i,i}$ are strictly positive. Let

$$B := 1 + \max_{i \in [N]} \max_{t \in [T/2, T]} [\Psi_i(t)]^{(1)}.$$

574 **Step 3.** Choose $R_1 > \max\{R_0, B\}$ and apply Lemma 7.1 on $[0, T/2]$ with target level R_1 to construct the
575 first block. Since $R_1 \geq R_0$, the bound defining C_0 applies. Hence, along the first-stage trajectories,

$$\left[\Phi_t^{\theta_{R_1}}(\mathbf{x}_i) \right]^{(2)} \geq -C_0 \quad \text{for } t \in [0, T/2].$$

576 Moreover, since $b_\ell T/2 + c_\ell = R_2$, we have $b_\ell t + c_\ell \leq R_2$ for $t \in [0, T/2]$. Therefore, by $R_2 < -C_0$,

$$- \left[\Phi_t^{\theta_{R_1}}(\mathbf{x}_i) \right]^{(2)} + b_\ell t + c_\ell \leq C_0 + R_2 < 0,$$

577 so the second block is inactive on $[0, T/2]$. Since $R_1 > B$, the first block is inactive on $[T/2, T]$ along
578 the auxiliary trajectories Ψ_i . By uniqueness, the full trajectories coincide with the concatenation of the
579 first-stage trajectories and the auxiliary second-stage trajectories. Therefore

$$\Phi_T^\theta(\mathbf{x}_i) = \mathbf{y}_i, \quad i \in [N].$$

580 □

581 **Remark 7.2.** Theorem 2.3 extends to any neural ODE of the form

$$\dot{\mathbf{x}}(t) = \sum_{j=1}^p \mathbf{w}_j (\mathbf{a}_j \cdot \mathbf{x}(t) + f_j(t) + c_j)_+, \quad t \in [0, T], \quad (46)$$

582 provided each f_j is strictly monotone. Geometrically, the hyperplane velocities then vary over time rather
583 than remain constant. We emphasize, however, that the linear choice $f_j(t) = b_j t$ is the simplest.

584 For instance, one may take $f_j(t) = d_j \tanh(b_j t)$, where $b_j > 0$ rescales time and $d_j \in \mathbb{R}$ bounds the
585 hyperplane speed. In this case, the evolution splits into two phases:

586 **1. Control:** For small $b_j t$ one has $\tanh(b_j t) < 1$ and the exact-control objective is attained in this phase.

587 **2. Stationary:** As $\tanh(b_j t) \rightarrow 1$, the system smoothly transitions to an autonomous regime.

7.1.2 Proof of Theorem 2.9

The proof of Theorem 2.9 requires the following preliminary lemma.

Lemma 7.3. Fix $d \geq 2$. Let $f_1 : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be of the form (ANODE), let $\mathbf{u}_1, \dots, \mathbf{u}_d \in \mathbb{R}^d$ be linearly independent, and let $\alpha_i < \beta_i$ for each $i \in [d]$. For any $\delta_i > 0$ such that $\alpha_i + \delta_i < \beta_i - \delta_i$, define the sets

$$K := \{\mathbf{x} \in \mathbb{R}^d : \alpha_i \leq \mathbf{u}_i \cdot \mathbf{x} \leq \beta_i \forall i \in [d]\} \quad \text{and} \quad \tilde{K} := \{\mathbf{x} \in \mathbb{R}^d : \alpha_i + \delta_i \leq \mathbf{u}_i \cdot \mathbf{x} \leq \beta_i - \delta_i \forall i \in [d]\}.$$

Then, there exists $f_2 : \mathbb{R}^d \rightarrow \mathbb{R}^d$ of the form (2ANODE) such that $f_2 \equiv f_1$ on \tilde{K} and $f_2 \equiv 0$ on $\mathbb{R}^d \setminus K$.

Proof of Lemma 7.3. Let $\delta_{\min} := \min_{i \in [d]} \delta_i > 0$. We define the barrier function $B : \mathbb{R}^d \rightarrow \mathbb{R}$ by

$$B(\mathbf{x}) := \sum_{i=1}^d ((\alpha_i + \delta_i - \mathbf{u}_i \cdot \mathbf{x})_+ + (\mathbf{u}_i \cdot \mathbf{x} - (\beta_i - \delta_i))_+).$$

By definition, if $\mathbf{x} \in \tilde{K}$, all arguments inside the ReLUs are non-positive, yielding $B(\mathbf{x}) = 0$. Conversely, if $\mathbf{x} \notin \tilde{K}$, then for some index j , either $\mathbf{u}_j \cdot \mathbf{x} < \alpha_j$ or $\mathbf{u}_j \cdot \mathbf{x} > \beta_j$. In either case, the corresponding ReLU term is strictly greater than δ_j , ensuring that $B(\mathbf{x}) > \delta_{\min} > 0$ on $\mathbb{R}^d \setminus \tilde{K}$.

Let $f_1(\mathbf{x}) = \sum_{m=1}^{p_2} \mathbf{v}_m L_m(\mathbf{x})$, where $L_m(\mathbf{x}) := (\mathbf{c}_m \cdot \mathbf{x} + d_m)_+$ and $\mathbf{v}_m \in \mathbb{R}^d$. Since $\mathbf{u}_1, \dots, \mathbf{u}_d$ span \mathbb{R}^d , then $B(\mathbf{x})$ grows linearly as $\|\mathbf{x}\| \rightarrow \infty$. Because each L_m has at most linear growth, the ratio L_m/B is bounded on $\mathbb{R}^d \setminus \tilde{K}$, allowing us to choose $\kappa > 0$ such that

$$L_m(\mathbf{x}) \leq \kappa B(\mathbf{x}) \quad \text{for all } \mathbf{x} \notin \tilde{K} \text{ and } m \in [p_2].$$

Finally, define

$$f_2(\mathbf{x}) := \sum_{m=1}^{p_2} \mathbf{v}_m (L_m(\mathbf{x}) - \kappa B(\mathbf{x}))_+.$$

If $\mathbf{x} \in \tilde{K}$, then $B(\mathbf{x}) = 0$, which implies $f_2(\mathbf{x}) = \sum_m \mathbf{v}_m L_m(\mathbf{x}) = f_1(\mathbf{x})$. If $\mathbf{x} \notin \tilde{K}$, the bound $L_m(\mathbf{x}) \leq \kappa B(\mathbf{x})$ guarantees that every ReLU is inactive, yielding $f_2(\mathbf{x}) = 0$. Because $B(\mathbf{x})$ is a linear combination of ReLUs of affine functions, f_2 is representable in the two-layer form (2ANODE). \square

Remark 7.4. This result extends immediately to any pair of strictly nested convex polytopes: one replaces the defining inequalities of K in the construction of B with the affine constraints defining its facets.

We are now ready to prove Theorem 2.9.

Proof of Theorem 2.9. Fix $T > 0$ and an admissible dataset $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ satisfying (13). Without loss of generality, we may assume $\mathbf{x}_i \neq \mathbf{y}_i$ for all i . Let $\mathbf{u}_i := (\mathbf{y}_i - \mathbf{x}_i)/T$ for each $i \in [N]$. The curve $\gamma_i(t) := \mathbf{x}_i + t\mathbf{u}_i$ for $t \in [0, T]$ traces exactly the segment $S_i := [\mathbf{x}_i, \mathbf{y}_i]$. Since all S_i are compact and pairwise disjoint, the minimum distance between them is positive:

$$\delta := \min_{i \neq j} \text{dist}(S_i, S_j) > 0.$$

We construct disjoint neighborhoods around each S_i as follows. For each i , fix an orthonormal basis $(\mathbf{e}_{i,1}, \dots, \mathbf{e}_{i,d})$ with $\mathbf{e}_{i,1} = \frac{\mathbf{y}_i - \mathbf{x}_i}{\|\mathbf{y}_i - \mathbf{x}_i\|}$. We define for each $\rho > 0$,

$$K_{i,\rho} := \{\mathbf{x} \in \mathbb{R}^d : -\rho \leq (\mathbf{x} - \mathbf{x}_i) \cdot \mathbf{e}_{i,1} \leq \|\mathbf{y}_i - \mathbf{x}_i\| + \rho \quad \text{and} \quad |(\mathbf{x} - \mathbf{x}_i) \cdot \mathbf{e}_{i,k}| \leq \rho \quad \text{for all } k \geq 2\}.$$

Set $r := \delta/(4\sqrt{d})$. By definition, every $\mathbf{x} \in S_i$ satisfies

$$0 \leq (\mathbf{x} - \mathbf{x}_i) \cdot \mathbf{e}_{i,1} \leq \|\mathbf{y}_i - \mathbf{x}_i\|, \quad (\mathbf{x} - \mathbf{x}_i) \cdot \mathbf{e}_{i,k} = 0 \quad \text{for all } k \geq 2.$$

Hence $S_i \subset K_{i,r/2} \subset \text{int}(K_{i,r})$. Moreover, the maximum distance from any $\mathbf{x} \in K_{i,r}$ to S_i is $\leq \sqrt{r^2 + (d-1)r^2} = r\sqrt{d} = \delta/4$, so $K_{i,r}$ lies entirely within the $\delta/4$ -neighborhood of S_i .

By the triangle inequality, for any $i \neq j$,

$$\text{dist}(K_{i,r}, K_{j,r}) \geq \text{dist}(S_i, S_j) - 2(\delta/4) \geq \delta - \delta/2 = \delta/2 > 0.$$

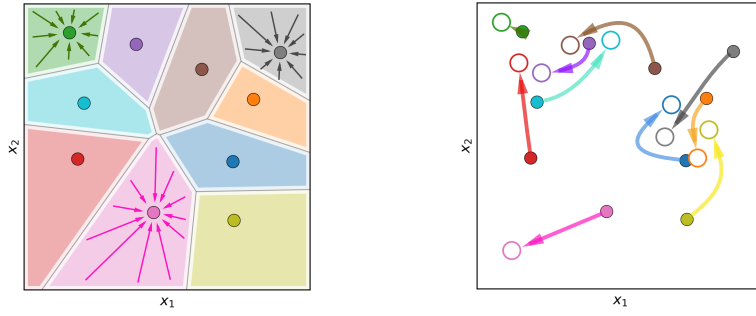


Figure 8: Two-step construction of the steering vector field. **Left:** Step 1 shows the initial pairwise disjoint convex sets A_k . The flow exponentially contracts these sets into small balls of radius δ centered around \mathbf{x}_k . **Right:** Step 2 depicts the rigid transport of these contracted balls along smooth, collision-free paths $\gamma_k(t)$ towards the target balls.

617 Hence, $K_{1,r}, \dots, K_{N,r}$ are pairwise disjoint.

618 We now construct the vector field. For each $i \in [N]$, we apply Lemma 7.3 to the constant field
 619 $g_i(\mathbf{x}) \equiv \mathbf{u}_i$ —which is trivially of the form (ANODE)—setting the inner region to $K_{i,r/2}$ and the outer
 620 region to $K_{i,r}$. This yields a localized field v_i of the form (2ANODE) such that

$$v_i(\mathbf{x}) = \mathbf{u}_i \quad \text{for all } \mathbf{x} \in K_{i,r/2}, \quad v_i(\mathbf{x}) = 0 \quad \text{for all } \mathbf{x} \in \mathbb{R}^d \setminus K_{i,r}.$$

621 Define the global vector field $v_\theta(\mathbf{x}) := \sum_{i=1}^N v_i(\mathbf{x})$. Since the neighborhoods $K_{1,r}, \dots, K_{N,r}$ are pairwise
 622 disjoint, for any fixed $i \in [N]$ we have

$$v_\theta(\mathbf{x}) = v_i(\mathbf{x}) + \sum_{j \neq i} \mathbf{0} = \mathbf{u}_i \quad \text{for all } \mathbf{x} \in K_{i,r/2}.$$

623 Recall that the trace of the trajectory is $\gamma_i([0, T]) = S_i \subset K_{i,r/2}$. Therefore, the ODE dynamics exactly
 624 follow the constant velocity: $\dot{\gamma}_i(t) = v_\theta(\gamma_i(t)) = \mathbf{u}_i$. Integrating this yields $\Phi_T^\theta(\mathbf{x}_i) = \gamma_i(T) = \mathbf{y}_i$. Since
 625 this holds for all $i \in [N]$, we deduce that Φ_T^θ perfectly interpolates the dataset.

626 Finally, we determine the architecture. By the explicit construction in Lemma 7.3, since $g_i(\mathbf{x}) \equiv \mathbf{u}_i$
 627 relies on a single constant term, the localized field takes the form

$$v_i(\mathbf{x}) = \mathbf{u}_i (1 - \kappa_i B_i(\mathbf{x}))_+,$$

628 where $B_i(\mathbf{x})$ is the barrier function. Because $K_{i,r}$ is defined by bounding the d projections onto the basis
 629 $\mathbf{e}_{i,k}$, the barrier B_i is a sum of exactly $2d$ ReLU terms (one for each upper and lower bound). Consequently,
 630 the global field v_θ exactly matches the two-layer architecture (2ANODE) with $p_1 = 2d$ and $p_2 = N$. \square

631 7.1.3 Proof of Theorem 2.5

632 We first state some auxiliary lemmas. The first establishes that we can construct a smooth, compactly
 633 supported, time-varying vector field that steers each A_k to its target ball. We then give an explicit estimate
 634 of its Barron norm as a function of η and of the geometry of the sets and connecting trajectories. To this
 635 end, for a compactly supported vector field $v \in \mathcal{C}_c^\infty(\mathbb{R} \times \mathbb{R}^d; \mathbb{R}^d)$, we define the Barron norm as

$$\|v\|_{\mathcal{B}_2^{d+1}} := \sum_{\ell=1}^d \int_{\mathbb{R} \times \mathbb{R}^d} \|(\tau, \omega)\|_1^2 |\widehat{v^{(\ell)}}(\tau, \omega)| \, d\tau \, d\omega, \quad (47)$$

636 where \widehat{v} is the Fourier transform of v . This definition and the subsequent lemma build upon the framework
 637 of [19, Theorem 2], which we adapt to the setting of time-dependent vector fields.

638 **Lemma 7.5** (Steering vector field). *Let $d \geq 2$, $K \geq 2$ and $T > 0$. Let $\{A_k\}_{k=1}^K$ be a family of pairwise
 639 disjoint compact convex subsets of \mathbb{R}^d , and let $\{\mathbf{r}_k\}_{k=1}^K \subset \mathbb{R}^d$ be pairwise distinct. Then, for every $\eta > 0$,
 640 there exists a vector field $v \in \mathcal{C}_c^\infty(\mathbb{R} \times \mathbb{R}^d; \mathbb{R}^d)$ whose flow map $\Phi_T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ from time 0 to time T
 641 satisfies*

$$\Phi_T(A_k) \subset B(\mathbf{r}_k, \eta) \quad \text{for all } k \in [K]. \quad (48)$$

642 *Proof of Lemma 7.5.* We first assume that the target points $\mathbf{r}_1, \dots, \mathbf{r}_K$ are pairwise distinct. Set

$$D_* := \max_{k \in [K]} \text{diam}(A_k), \quad s_* := \min_{i \neq j} \text{dist}(A_i, A_j) > 0. \quad (49)$$

643 By definition of s_* , the open neighborhoods

$$V_k := \left\{ \mathbf{x} \in \mathbb{R}^d : \text{dist}(\mathbf{x}, A_k) < \frac{s_*}{8} \right\} = A_k + \left(\frac{s_*}{8} \right) B(0, 1) \quad (50)$$

644 are pairwise disjoint. Since $d \geq 2$, we may choose arbitrary points $\mathbf{x}_k \in A_k$ and connect each \mathbf{x}_k to
 645 its target \mathbf{r}_k via pairwise disjoint polygonal arcs. This is a standard topological consequence of working
 646 in dimensions $d \geq 2$, where finitely many paths can avoid intersections by introducing small perturba-
 647 tions. By smoothing the corners of these arcs and parameterizing them appropriately over time, we obtain
 648 smooth, pairwise disjoint curves $\gamma_k : [0, T] \rightarrow \mathbb{R}^d$ such that

$$\gamma_k(0) = \mathbf{x}_k, \quad \gamma_k \text{ is constant on } [0, 2T/3], \quad \gamma_k(T) = \mathbf{r}_k \quad (51)$$

649 and, by ensuring the polygonal arcs are contained within a reasonably tight bounding ball before smooth-
 650 ing, we also satisfy

$$\max_{t \in [0, T]} \|\dot{\gamma}_k(t)\| < 1 + \max_{k \in [K]} \|\mathbf{r}_k\| + \max_{k \in [K]} \max_{\mathbf{x} \in A_k} \|\mathbf{x}\|. \quad (52)$$

651 Our proof is constructive. In particular, the final vector field will be the combination of a *compressing*
 652 vector field, collapsing each A_k to a small ball, and a *transporting* vector field, driving these balls to the
 653 targets \mathbf{r}_k . Figure 8 summarizes these steps visually. Accordingly, we divide this proof into three parts.

654 **Compression.** We fix a standard nonnegative radial mollifier $\psi \in \mathcal{C}_c^\infty(\mathbb{R}^d)$, with support in $B(0, 1)$
 655 and $\int \psi = 1$, and set

$$\psi_{s_*/32}(\mathbf{x}) := \left(\frac{32}{s_*} \right)^d \psi \left(\frac{32}{s_*} \mathbf{x} \right). \quad (53)$$

656 For each $k \in [K]$, define the smooth cutoff

$$\chi_k := \mathbf{1}_{A_k + (s_*/32)B(0,1)} * \psi_{s_*/32} \quad (54)$$

657 so that $\chi_k \in \mathcal{C}_c^\infty(\mathbb{R}^d; [0, 1])$ with $\chi_k \equiv 1$ on A_k and $\text{supp}(\chi_k) \subset A_k + (s_*/16)B(0, 1) \subset V_k$. Since the
 658 supports of χ_k are pairwise disjoint, the vector field

$$u(\mathbf{x}) := \sum_{k=1}^K \chi_k(\mathbf{x})(\mathbf{x}_k - \mathbf{x}) \quad (55)$$

659 is well-defined and satisfies $u \in \mathcal{C}_c^\infty(\mathbb{R}^d; \mathbb{R}^d)$. On each A_k , (55) simplifies exactly to $u(\mathbf{x}) = \mathbf{x}_k - \mathbf{x}$.

660 **Transport.** Let m be defined as

$$m := \min_{t \in [0, T]} \min_{i \neq j} \|\gamma_i(t) - \gamma_j(t)\|,$$

661 and denote

$$\delta := \min \left\{ \eta, \frac{s_*}{8}, \frac{m}{8} \right\}, \quad \lambda := \frac{3}{T} \log_+ \left(\frac{D_*}{\delta} \right) \geq 0. \quad (56)$$

662 Furthermore, let $\xi \in \mathcal{C}_c^\infty(\mathbb{R}^d; [0, 1])$ be a smooth radial cutoff with $\xi \equiv 1$ on $B(0, 1)$ and $\text{supp} \xi \subset B(0, 2)$,
 663 and define

$$w(t, \mathbf{x}) := \sum_{k=1}^K \dot{\gamma}_k(t) \xi \left(\frac{8(\mathbf{x} - \gamma_k(t))}{m} \right), \quad (t, \mathbf{x}) \in [0, T] \times \mathbb{R}^d. \quad (57)$$

664 For each $t \in [0, T]$, the spatial support of the k -th summand is contained in $B(\gamma_k(t), m/4)$, and these
 665 balls are pairwise disjoint. Moreover, $w(t, \mathbf{x}) = \dot{\gamma}_k(t)$ whenever $\|\mathbf{x} - \gamma_k(t)\| \leq \frac{m}{8}$.

666 **Concatenation.** To concatenate compression and transport, choose a nondecreasing $\varphi \in \mathcal{C}^\infty([0, T])$
 667 with $\varphi(t) = 0$ on $[0, T/3]$ and $\varphi(t) = 1$ on $[2T/3, T]$, and define the time-dependent vector field

$$v_{[0, T]}(t, \mathbf{x}) := (1 - \varphi(t)) \lambda u(\mathbf{x}) + \varphi(t) w(t, \mathbf{x}), \quad (t, \mathbf{x}) \in [0, T] \times \mathbb{R}^d. \quad (58)$$

668 Then $v_{[0, T]} \in \mathcal{C}_c^\infty([0, T] \times \mathbb{R}^d; \mathbb{R}^d)$, and thus (58) generates a unique flow Φ_t on $[0, T]$.

669 We verify that Φ_t satisfies the condition (48). Since γ_k is constant on $[0, 2T/3]$, we have $w \equiv 0$
 670 on $[0, 2T/3] \times \mathbb{R}^d$, and therefore (58) reduces there to $v_{[0, T]}(t, \mathbf{x}) = \lambda(1 - \varphi(t))u(\mathbf{x})$. Besides, solving
 671 $\dot{\mathbf{x}}(t) = \lambda(1 - \varphi(t))u(\mathbf{x})$, we see that the trajectory of any initial point $\mathbf{x}(0) \in A_k$ satisfies

$$\Phi_t(\mathbf{x}(0)) = \mathbf{x}_k + \exp\left(-\lambda \int_0^t (1 - \varphi(s)) ds\right) (\mathbf{x}(0) - \mathbf{x}_k) \in A_k \quad \text{for } t \in [0, 2T/3]. \quad (59)$$

Since $\varphi(s) = 0$ on $[0, T/3]$ and $\varphi \leq 1$ everywhere, the integral satisfies $\int_0^{2T/3} (1 - \varphi(s)) ds \geq T/3$.
 Recalling our choice of λ in (56), this guarantees that at $t = 2T/3$,

$$\|\Phi_{2T/3}(\mathbf{x}(0)) - \mathbf{x}_k\| \leq e^{-\lambda(T/3)} \|\mathbf{x}(0) - \mathbf{x}_k\| \leq e^{-\lambda(T/3)} D_* \leq \delta.$$

672 Consequently, we obtain

$$\Phi_{2T/3}(A_k) \subset B(\mathbf{x}_k, \delta). \quad (60)$$

673 On $[2T/3, T]$, we have $\varphi \equiv 1$, so $v_{[0, T]} = w$. Let $z_0 \in B(\mathbf{x}_k, \delta)$, and write $z_0 = \mathbf{x}_k + \mathbf{y} = \gamma_k(2T/3) + \mathbf{y}$
 674 with $\|\mathbf{y}\| \leq \delta \leq m/8$. The path $z(t) := \gamma_k(t) + \mathbf{y}$ stays in the region where $w(t, \cdot) = \dot{\gamma}_k(t)$, hence solves
 675 the ODE; by uniqueness, the flow starting from z_0 equals $z(t)$, and in particular reaches $\mathbf{r}_k + \mathbf{y}$ at $t = T$.
 676 Combined with (60), this gives

$$\Phi_T(A_k) \subset B(\mathbf{r}_k, \delta) \subset B(\mathbf{r}_k, \eta).$$

677 To conclude, by a standard smooth extension-and-cutoff argument in the time variable, there exists a
 678 global field

$$v \in \mathcal{C}_c^\infty(\mathbb{R} \times \mathbb{R}^d; \mathbb{R}^d) \quad \text{such that } v = v_{[0, T]} \quad \text{on } [0, T] \times \mathbb{R}^d, \quad (61)$$

679 Since $v = v_{[0, T]}$ on $[0, T] \times \mathbb{R}^d$, the flow of v on $[0, T]$ coincides with Φ_t as above. In particular, (48) holds.
 680 \square

681 **Lemma 7.6** (Quantitative estimates on the steering field). *Let $v = v_\eta$ be the vector field constructed in the*
 682 *proof of Lemma 7.5 with tolerance η when $\{\mathbf{r}_k\}_{k=1}^K \subset \mathbb{R}^d$ are pairwise distinct (see (61)) and let Φ_t be its*
 683 *flow on $[0, T]$. Define the minimum path separation m and the maximum path derivative G associated with*
 684 *the trajectories γ_k as*

$$m := \min_{t \in [0, T]} \min_{i \neq j} \|\gamma_i(t) - \gamma_j(t)\|, \quad G := 1 + \max_{1 \leq j \leq d+5} \max_{k \in [K]} \left\| \frac{d^j \gamma_k}{dt^j} \right\|_{L^\infty(0, T)}, \quad (62)$$

and let δ be as in (56). Then there exists a constant $C_{d, T} > 0$, independent of η and $(A_k, \gamma_k)_k$, such that

$$\sup_{t \in [0, T]} \text{Lip}_{\mathbf{x}}(v(t, \cdot)) \leq C_{d, T} \left[\frac{G}{m} + \frac{1 + D_*/s_*}{T} \log_+ \left(\frac{D_*}{\delta} \right) \right]. \quad (63)$$

Moreover, the Barron norm of v , as defined in (47), satisfies

$$\|v\|_{\mathcal{B}_2^{d+1}} \leq C_{d, T} K \left[TG^{d+5}(m^d + m^{-4}) + \frac{(D_* + s_*)^{d+1}(1 + s_*^{-(d+4)})}{T} \log_+ \left(\frac{D_*}{\delta} \right) \right]. \quad (64)$$

Finally, for $R > 0$ large enough so that $[-R, R]^d \supset \bigcup_{k \in [K]} A_k$, we also have:

$$\max_{(t, \mathbf{x}) \in [0, T] \times [-R, R]^d} \|\Phi_t(\mathbf{x})\| \leq 1 + \max_{k \in [K]} \|\mathbf{r}_k\| + \frac{3}{2} \sqrt{d} R. \quad (65)$$

685 *Proof.* Throughout the proof, C_d and $C_{d, T}$ denote positive constants that may change from line to line
 686 and depend only on d, T , and on the fixed smooth profiles ψ, ξ , and φ .

687 **Lipschitz estimate (63).** Using the decomposition (58) and the bound $0 \leq \varphi \leq 1$, we obtain

$$\sup_{t \in [0, T]} \text{Lip}_{\mathbf{x}}(v(t, \cdot)) \leq \lambda \|Du\|_{L^\infty(\mathbb{R}^d)} + \sup_{t \in [0, T]} \|D_{\mathbf{x}}w(t, \cdot)\|_{L^\infty(\mathbb{R}^d)}. \quad (66)$$

688 Recall the definition of u in (55). Writing $\chi_k = \mathbf{1}_{E_k} * \psi_{s_*/32}$ with $E_k := A_k + (s_*/32)B(0, 1)$, by the
689 classical Young's convolution inequality, we have

$$\|\nabla \chi_k\|_{L^\infty(\mathbb{R}^d)} = \|\mathbf{1}_{E_k} * \nabla \psi_{s_*/32}\|_{L^\infty(\mathbb{R}^d)} \leq \|\mathbf{1}_{E_k}\|_{L^\infty} \|\nabla \psi_{s_*/32}\|_{L^1} \leq C_d s_*^{-1}. \quad (67)$$

690 Since the functions χ_k have pairwise disjoint supports, and $\|\mathbf{x} - \mathbf{x}_k\| \leq D_* + s_*/16$ for $\mathbf{x} \in \text{supp}(\chi_k)$,
691 we obtain

$$\|Du\|_{L^\infty(\mathbb{R}^d)} \leq C_d \left(1 + \frac{D_*}{s_*}\right). \quad (68)$$

692 For the second term, differentiating (57) gives:

$$\partial_{x_j} w^{(\ell)}(t, \mathbf{x}) = \sum_{k=1}^K \dot{\gamma}_k^{(\ell)}(t) \frac{8}{m} \partial_j \xi \left(\frac{8(\mathbf{x} - \gamma_k(t))}{m} \right), \quad \text{for all } j, \ell \in [d].$$

693 For any fixed t , the support of the k -th summand is contained in $B(\gamma_k(t), m/4)$. These balls are pairwise
694 disjoint for all k , meaning at most one summand is nonzero at any given (t, \mathbf{x}) . Since $|\dot{\gamma}_k(t)| \leq G$ and
695 $\|\nabla \xi\|_{L^\infty} < \infty$, we obtain

$$\sup_{t \in [0, T]} \|D_{\mathbf{x}}w(t, \cdot)\|_{L^\infty(\mathbb{R}^d)} \leq C_d \frac{G}{m}. \quad (69)$$

696 Substituting (68), (69), and (56) into (66) yields the desired bound (63).

Barron estimate. Fix $n = d + 4$ and $g \in \mathcal{C}_c^\infty(\mathbb{R}^{d+1})$. Since $\frac{\|(\tau, \omega)\|_1^2}{(1 + \|(\tau, \omega)\|_1)^n}$ is integrable on \mathbb{R}^{d+1} ,

$$\int_{\mathbb{R}^{d+1}} \|(\tau, \omega)\|_1^2 |\widehat{g}(\tau, \omega)| \, d\tau \, d\omega \leq C_d \sup_{(\tau, \omega) \in \mathbb{R}^{d+1}} (1 + \|(\tau, \omega)\|_1)^n |\widehat{g}(\tau, \omega)|.$$

697 Expanding the weight as a sum of monomials and using the identity $(\tau, \omega)^\alpha \widehat{g} = c_\alpha \widehat{\partial^\alpha g}$ together with
698 $\|\widehat{h}\|_{L^\infty} \leq \|h\|_{L^1}$, we obtain

$$\int_{\mathbb{R}^{d+1}} \|(\tau, \omega)\|_1^2 |\widehat{g}(\tau, \omega)| \, d\tau \, d\omega \leq C_d \sum_{|\alpha| \leq n} \|\partial^\alpha g\|_{L^1(\mathbb{R}^{d+1})}. \quad (70)$$

699 Moreover, for every component $\ell \in [d]$, the extension operator and the cutoff yield

$$\sum_{a+|\beta| \leq n} \|\partial_t^a \partial_{\mathbf{x}}^\beta v^{(\ell)}\|_{L^1(\mathbb{R} \times \mathbb{R}^d)} \leq C_{d, T} \sum_{a+|\beta| \leq n} \|\partial_t^a \partial_{\mathbf{x}}^\beta v_{[0, T]}^{(\ell)}\|_{L^1((0, T) \times \mathbb{R}^d)}. \quad (71)$$

700 Applying (70) componentwise to v and using (71), we obtain

$$\|v\|_{\mathcal{B}_2^{d+1}} \leq C_{d, T} \sum_{\ell=1}^d \sum_{a+|\beta| \leq n} \|\partial_t^a \partial_{\mathbf{x}}^\beta v_{[0, T]}^{(\ell)}\|_{L^1((0, T) \times \mathbb{R}^d)}. \quad (72)$$

701 Thus, it suffices to bound the L^1 derivatives of $v_{[0, T]} = v^C + v^T$, where

$$v^C(t, \mathbf{x}) := (1 - \varphi(t)) \lambda u(\mathbf{x}), \quad v^T(t, \mathbf{x}) := \varphi(t) w(t, \mathbf{x}).$$

702 Since φ is a fixed smooth profile, its derivatives up to order n are bounded in $L^\infty(0, T)$ and in $L^1(0, T)$
703 by a constant $C_{d, T} > 0$; we will use this repeatedly below without further mention.

704 **Compression.** Recalling the definition of χ_k at (54) and of E_k above, Young's convolution inequality, com-
705 bined with the isodiametric bound $|E_k| \leq C_d (D_* + s_*)^d$ and the scaling $\|\partial^\gamma \psi_{s_*/32}\|_{L^1} \leq C_d s_*^{-|\gamma|}$ yields,
706 for every $|\gamma| \leq n$,

$$\|\partial^\gamma \chi_k\|_{L^1(\mathbb{R}^d)} \leq C_d (D_* + s_*)^d s_*^{-|\gamma|}. \quad (73)$$

707 Applying the product rule to $u^{(\ell)}(\mathbf{x}) = \sum_k \chi_k(\mathbf{x})(x_k^{(\ell)} - x^{(\ell)})$, and, using $\|\mathbf{x}_k - \mathbf{x}\| \leq D_* + s_*$ for
708 $\mathbf{x} \in \text{supp}(\chi_k)$, (73) gives

$$\sum_{|\beta| \leq n} \|\partial^\beta u^{(\ell)}\|_{L^1(\mathbb{R}^d)} \leq C_d K (D_* + s_*)^{d+1} (1 + s_*^{-n}). \quad (74)$$

709 Multiplying by the temporal factor $\lambda(1 - \varphi)$ and using the uniform bound on the L^1 -derivatives of φ ,

$$\sum_{a+|\beta| \leq n} \|\partial_t^a \partial_x^\beta v^{C,(\ell)}\|_{L^1((0,T) \times \mathbb{R}^d)} \leq C_{d,T} \lambda K (D_* + s_*)^{d+1} (1 + s_*^{-n}). \quad (75)$$

710 *Transport.* Differentiating $w^{(\ell)}(t, \mathbf{x}) = \sum_k \dot{\gamma}_k^{(\ell)}(t) \xi(8(\mathbf{x} - \gamma_k(t))/m)$ via the Leibniz and chain rules,
711 every spatial derivative extracts a factor m^{-1} , and every time derivative either hits ξ (extracting another
712 m^{-1}) or lands on the prefactor $\dot{\gamma}_k^{(\ell)}$. Furthermore, we note that

$$\text{supp } \xi(8(\cdot - \gamma_k(t))/m) \subset B(\gamma_k(t), m/4),$$

713 and that the latter are pairwise disjoint at each t . Thus, at most one summand contributes at any (t, \mathbf{x}) ,
714 so for all $a + |\beta| \leq n$,

$$\left| \partial_t^a \partial_x^\beta w^{(\ell)}(t, \mathbf{x}) \right| \leq C_d G^{n+1} \sum_{r=0}^a m^{-(|\beta|+r)} \sum_{k=1}^K \mathbf{1}_{B(\gamma_k(t), m/4)}(\mathbf{x}). \quad (76)$$

715 Integrating over $(0, T) \times \mathbb{R}^d$, where each ball has volume $\lesssim m^d$, yields

$$\|\partial_t^a \partial_x^\beta w^{(\ell)}\|_{L^1((0,T) \times \mathbb{R}^d)} \leq C_{d,T} K G^{n+1} \sum_{r=0}^a m^{d-|\beta|-r}. \quad (77)$$

716 Since $|\beta| + r \leq n = d + 4$, each exponent $d - |\beta| - r$ lies in $[-4, d]$, hence $m^{d-|\beta|-r} \leq m^d + m^{-4}$.
717 Summing over $a + |\beta| \leq n$ and absorbing the φ -factor from $v^T = \varphi w$, we obtain

$$\sum_{a+|\beta| \leq n} \|\partial_t^a \partial_x^\beta v^{T,(\ell)}\|_{L^1((0,T) \times \mathbb{R}^d)} \leq C_{d,T} K G^{n+1} (m^d + m^{-4}). \quad (78)$$

718 Combining (75), (78), and (72) with $n = d + 4$, and setting $\lambda \leq (3/T) \log_+(D_*/\delta)$ from (56), proves (64).

719 **Uniform bound.** Let $\mathbf{x}_0 \in [-R, R]^d$. For $t \in [0, 2T/3]$, we have $v(t, \mathbf{x}) = \lambda(1 - \varphi(t))u(\mathbf{x})$, since
720 $w \equiv 0$ on $[0, 2T/3] \times \mathbb{R}^d$. The flow either fixes \mathbf{x}_0 or moves it along a straight segment within the convex
721 set $[-R, R]^d$. Therefore, $\|\Phi_t(\mathbf{x}_0)\| \leq \sqrt{d} R$ for all $t \in [0, 2T/3]$.

For $t \in [2T/3, T]$ we have $v = w$, which vanishes outside $\bigcup_{k=1}^K B(\gamma_k(t), m/4)$. Therefore, whenever $\Phi_t(\mathbf{x}_0)$ is not stationary, it lies in $B(\gamma_k(t), m/4)$ for some k , yielding

$$\|\Phi_t(\mathbf{x}_0)\| \leq \|\gamma_k(t)\| + \frac{m}{4}.$$

By (52) we have $\|\gamma_k(t)\| < 1 + \max_{j \in [K]} \|\mathbf{r}_j\| + \sqrt{d} R$. Additionally, since $\gamma_k(0) \in A_k \subset [-R, R]^d$, we
trivially have $m \leq 2\sqrt{d} R$. Combining these estimates gives

$$\|\Phi_t(\mathbf{x}_0)\| \leq 1 + \max_{j \in [K]} \|\mathbf{r}_j\| + \frac{3}{2} \sqrt{d} R.$$

Since this bound dominates the compression phase and holds whenever the trajectory moves, we conclude

$$\max_{(t, \mathbf{x}) \in [0, T] \times [-R, R]^d} \|\Phi_t(\mathbf{x})\| \leq 1 + \max_{k \in [K]} \|\mathbf{r}_k\| + \frac{3}{2} \sqrt{d} R,$$

722 which proves (65). □

723 The next lemma establishes an $O(p^{-1/2})$ approximation bound for the flow generated by (SANODE).
 724 The result relies on [19, Theorem 2], which was subsequently adapted for dynamical systems in [21].
 725 Here, we extend [21, Theorem 2.3] by making the dependence on the Barron norm of the reference vector
 726 field explicit.

727 **Lemma 7.7.** Fix $d \geq 2$ and $T > 0$, let $v \in \mathcal{C}_c^\infty(\mathbb{R} \times \mathbb{R}^d; \mathbb{R}^d)$ with $\|v\|_{\mathcal{B}_2^{d+1}} < \infty$, and let $\Omega \subset \mathbb{R}^d$ be
 728 compact. Let $(\Psi_t)_{t \in [0, T]}$ be the flow induced on $[0, T]$ by $\dot{\mathbf{x}} = v(t, \mathbf{x})$, and assume that for some $R_* \geq 1$,

$$[0, T] \times \bigcup_{t \in [0, T]} \Psi_t(\Omega) \subset [-(R_* - 1), R_* - 1]^{d+1}. \quad (79)$$

729 Then there exists a constant $C_d > 0$, depending only on d , such that for every integer $p \geq 3$ satisfying

$$p > C_d^2 R_*^4 T^2 \|v\|_{\mathcal{B}_2^{d+1}}^2 \exp\left(2T \sup_{t \in [0, T]} \text{Lip}_{\mathbf{x}}(v(t, \cdot))\right), \quad (80)$$

730 there exists a control θ of width p for the field v_θ of the form (SANODE) such that the flow Φ_t^θ satisfies

$$\sup_{\mathbf{x} \in \Omega} \sup_{t \in [0, T]} \|\Psi_t(\mathbf{x}) - \Phi_t^\theta(\mathbf{x})\| \leq \frac{C_d R_*^2 T \|v\|_{\mathcal{B}_2^{d+1}}}{\sqrt{p}} \exp\left(T \sup_{t \in [0, T]} \text{Lip}_{\mathbf{x}}(v(t, \cdot))\right). \quad (81)$$

731 *Proof.* Denote $L := \sup_{t \in [0, T]} \text{Lip}_{\mathbf{x}}(v(t, \cdot))$. By [19, Theorem 2], there exists a control θ of width p and a
 732 constant $C_d > 0$ such that²

$$\sup_{(t, \mathbf{x}) \in [-R_*, R_*]^{d+1}} \|v(t, \mathbf{x}) - v_\theta(t, \mathbf{x})\| \leq \frac{C_d R_*^2 \|v\|_{\mathcal{B}_2^{d+1}}}{\sqrt{p}} =: \varepsilon.$$

733 For any $\mathbf{x} \in \Omega$, let $\tau \leq T$ be the maximal time such that $\Phi_t^\theta(\mathbf{x}) \in [-R_*, R_*]^d$ for all $t \in [0, \tau]$. On $[0, \tau]$,
 734 adding and subtracting $v(t, \Phi_t^\theta(\mathbf{x}))$ gives

$$\|\Psi_t(\mathbf{x}) - \Phi_t^\theta(\mathbf{x})\| \leq \int_0^t L \|\Psi_s(\mathbf{x}) - \Phi_s^\theta(\mathbf{x})\| ds + \varepsilon t,$$

735 so Grönwall's lemma yields

$$\sup_{t \in [0, \tau]} \|\Psi_t(\mathbf{x}) - \Phi_t^\theta(\mathbf{x})\| \leq \varepsilon T e^{LT} < 1,$$

736 where the last inequality follows from (80). Hence $\Phi_t^\theta(\mathbf{x})$ stays within distance 1 of $\Psi_t(\mathbf{x})$, which by (79)
 737 lies in $[-(R_* - 1), R_* - 1]^d$. Thus $\Phi_t^\theta(\mathbf{x}) \in (-R_*, R_*)^d$ on $[0, \tau]$, forcing $\tau = T$ by continuity. Taking
 738 the supremum over $\mathbf{x} \in \Omega$ gives (81). \square

739 We are now in a position to prove the theorem.

740 *Proof of Theorem 2.5.* Assume first that the targets \mathbf{r}_k are pairwise distinct. Since $d \geq 2$, we can choose
 741 smooth, pairwise disjoint $\gamma_k : [0, T] \rightarrow \mathbb{R}^d$ satisfying (51)–(52). We use the construction in the proof of
 742 Lemma 7.5 with these fixed paths and tolerance $\eta/4$. Let v be the resulting field and $(\Psi_t)_{t \in [0, T]}$ its flow.
 743 Let s_* be defined by (49) and (m, G) by (62), set $\eta_0 := \min\{1, s_*/2, m/2\}$, and fix $\eta \in (0, \eta_0]$. Therefore,

$$\Psi_T(A_k) \subset B(\mathbf{r}_k, \eta/4) \quad \text{for all } k \in [K]. \quad (82)$$

744 For $\mathfrak{M} > 0$ given by (12), it follows that:

$$[0, T] \times \bigcup_{t \in [0, T]} \Psi_t(I_R) \subset [-(\mathfrak{M} - 1), \mathfrak{M} - 1]^{d+1}. \quad (83)$$

²[19, Theorem 2] is scalar-valued, but applying it componentwise and distributing the hidden units among the d components only changes the dimensional constant C_d .

Moreover, if $\eta > 0$ is small enough that $2\eta < s_* \wedge m$, then $\delta = \min \left\{ \frac{\eta}{4}, \frac{s_*}{8}, \frac{m}{8} \right\} = \frac{\eta}{4}$. Lemma 7.6 then yields:

$$\sup_{t \in [0, T]} \text{Lip}_{\mathbf{x}}(v(t, \cdot)) \leq L(1 - \log \eta), \quad \|v\|_{\mathcal{B}_2^{d+1}} \leq B(1 - \log \eta),$$

Increasing \mathfrak{C} , if necessary, we also ensure that the lower bound (11) implies the size condition (80) for every $\eta \in (0, \eta_0]$, for constants $L, B \geq 0$ that depend only on $d, T, R, (\mathbf{r}_k)_k, (A_k)_k$. Applying Lemma 7.7 with $\Omega = I_R$ and $R_* = \mathfrak{M}$, we obtain that for any sufficiently large $p \geq 3$:

$$\begin{aligned} \sup_{\mathbf{x} \in I_R} \sup_{t \in [0, T]} \|\Psi_t(\mathbf{x}) - \Phi_t^\theta(\mathbf{x})\| &\leq \frac{C_d \mathfrak{M}^2 T B (1 - \log \eta)}{\sqrt{p}} \exp(LT(1 - \log \eta)) \\ &\leq \frac{C_d \mathfrak{M}^2 T B e^{LT} (1 - \log \eta)}{\sqrt{p}} \eta^{-LT}. \end{aligned}$$

Now, choose \mathfrak{C} sufficiently large so that $\mathfrak{C} > 2LT$ and $\sqrt{\mathfrak{C}} \geq 2C_d \mathfrak{M}^2 T B e^{LT}$. If p satisfies the bound in (11), we deduce:

$$\sup_{\mathbf{x} \in I_R} \sup_{t \in [0, T]} \|\Psi_t(\mathbf{x}) - \Phi_t^\theta(\mathbf{x})\| \leq \frac{C_d \mathfrak{M}^2 T B e^{LT}}{\sqrt{\mathfrak{C}}} \eta^{1+\mathfrak{C}/2-LT} \leq \frac{\eta}{2}.$$

For any $\mathbf{x} \in A_k$, combining this bound with (82) allows us to verify (i) via the triangle inequality:

$$\|\Phi_T^\theta(\mathbf{x}) - \mathbf{r}_k\| \leq \|\Phi_T^\theta(\mathbf{x}) - \Psi_T(\mathbf{x})\| + \|\Psi_T(\mathbf{x}) - \mathbf{r}_k\| \leq \frac{\eta}{2} + \frac{\eta}{4} < \eta.$$

Furthermore, using (83), we can easily verify (ii):

$$\|\Phi_T^\theta\|_{L^\infty(I_R)} \leq \|\Psi_T\|_{L^\infty(I_R)} + \|\Phi_T^\theta - \Psi_T\|_{L^\infty(I_R)} \leq (\mathfrak{M} - 1) + \frac{\eta}{2} < \mathfrak{M},$$

which holds because $\eta \leq 1$.

Finally, consider the case where the targets \mathbf{r}_k are not necessarily pairwise distinct. We can slightly perturb them to choose auxiliary targets $\tilde{\mathbf{r}}_k$ that are pairwise distinct and satisfy:

$$\|\tilde{\mathbf{r}}_k - \mathbf{r}_k\| < \eta/4 \quad \text{for all } k.$$

Applying the previous argument to these auxiliary targets $\tilde{\mathbf{r}}_1, \dots, \tilde{\mathbf{r}}_K$ with a tighter tolerance of $\eta/2$, we obtain, for sufficiently large p , a control θ such that $\Phi_T^\theta(A_k) \subset B(\tilde{\mathbf{r}}_k, \eta/2)$. By the triangle inequality, it immediately follows that $\Phi_T^\theta(A_k) \subset B(\mathbf{r}_k, \eta)$. Since the mutual separation of the auxiliary targets may be of order η , the path-separation parameter m , and hence the constants produced by the previous argument, need not remain uniform as $\eta \downarrow 0$. \square

7.2 Proofs of Section 3

Proof of Proposition 3.1. By the squared-sum inequality, the expected error can be decomposed into a bias and a variance component:

$$\mathbb{E}_{\mathcal{D}_N} \left[\|y - y_{N,h}\|_{L^2(\mu)}^2 \right] \leq 2 \|y - y_h\|_{L^2(\mu)}^2 + 2 \mathbb{E}_{\mathcal{D}_N} \left[\|y_h - y_{N,h}\|_{L^2(\mu)}^2 \right],$$

where y_h and $y_{N,h}$ are the population and empirical estimators defined in (17) and (19), respectively. We bound each term separately.

Bias. Fix a cell Q_k with $p_k := \mu(Q_k) > 0$. The population coefficient \mathbf{s}_k is the conditional expectation of $y(\cdot)$ on Q_k . By Jensen's inequality, for any $\mathbf{x} \in Q_k$:

$$\|y(\mathbf{x}) - \mathbf{s}_k\| = \left\| \frac{1}{p_k} \int_{Q_k} (y(\mathbf{x}) - y(\mathbf{x}')) \, d\mu(\mathbf{x}') \right\| \leq \frac{1}{p_k} \int_{Q_k} \|y(\mathbf{x}) - y(\mathbf{x}')\| \, d\mu(\mathbf{x}').$$

For any $\mathbf{x}, \mathbf{x}' \in Q_k$, the distance satisfies: $\|\mathbf{x} - \mathbf{x}'\| \leq \text{diam}(Q_k) \leq \sqrt{d}h$. Consequently, $\|y(\mathbf{x}) - y(\mathbf{x}')\| \leq \omega_y(\sqrt{d}h)$, which implies $\|y(\mathbf{x}) - \mathbf{s}_k\| \leq \omega_y(\sqrt{d}h)$. Since cells with $p_k = 0$ do not contribute to the $L^2(\mu)$ norm, integrating over the partition yields:

$$\|y - y_h\|_{L^2(\mu)}^2 = \sum_{k \in \mathcal{K}} \int_{Q_k} \|y(\mathbf{x}) - \mathbf{s}_k\|^2 \, d\mu(\mathbf{x}) \leq \omega_y(\sqrt{d}h)^2 \sum_{k \in \mathcal{K}} \mu(Q_k) = \omega_y(\sqrt{d}h)^2, \quad (84)$$

769 where \mathcal{K} is the index set defined in (14).

770 **Variance.** Since $y_h|_{Q_k} = \mathbf{s}_k$ and $y_{N,h}|_{Q_k} = \mathbf{S}_k$, the variance term becomes:

$$\mathbb{E}_{\mathcal{D}_N} \left[\|y_h - y_{N,h}\|_{L^2(\mu)}^2 \right] = \sum_{k \in \mathcal{K}} p_k \mathbb{E}_{\mathcal{D}_N} [\|\mathbf{s}_k - \mathbf{S}_k\|^2]. \quad (85)$$

771 Fix an index k and let $N_k := \sum_{i=1}^N \mathbf{1}_{Q_k}(\mathbf{x}_i)$ denote the number of samples falling into Q_k . Because the
772 inputs are drawn i.i.d. from μ , N_k follows a binomial distribution, $N_k \sim \text{Binom}(N, p_k)$.

773 Conditionally on $N_k = n \geq 1$, the coefficient \mathbf{S}_k is the empirical average of n i.i.d. random variables
774 $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ distributed as $y(\mathbf{X})$ for $\mathbf{X} \sim \mu(\cdot | Q_k)$. By definition, $\mathbb{E}[\mathbf{Y}_1] = \mathbf{s}_k$. The variance of this
775 empirical average is exactly:

$$\mathbb{E} [\|\mathbf{S}_k - \mathbf{s}_k\|^2 | N_k = n] = \mathbb{E} \left[\left\| \frac{1}{n} \sum_{j=1}^n (\mathbf{Y}_j - \mathbf{s}_k) \right\|^2 \right] = \frac{1}{n} \mathbb{E} [\|\mathbf{Y}_1 - \mathbf{s}_k\|^2] \leq \frac{\|y\|_{L^\infty(I_R)}^2}{n}.$$

776 If $N_k = 0$, we defined $\mathbf{S}_k = 0$, which gives $\|\mathbf{S}_k - \mathbf{s}_k\|^2 = \|\mathbf{s}_k\|^2 \leq \|y\|_{L^\infty(I_R)}^2$. Taking the unconditional
777 expectation over the binomial variable N_k , we obtain:

$$\mathbb{E}_{\mathcal{D}_N} [\|\mathbf{S}_k - \mathbf{s}_k\|^2] \leq \|y\|_{L^\infty(I_R)}^2 \left(\mathbb{E} \left[\frac{\mathbf{1}_{\{N_k > 0\}}}{N_k} \right] + \mathbb{P}(N_k = 0) \right).$$

778 We bound the two terms in the parenthesis. First, $\mathbb{P}(N_k = 0) = (1 - p_k)^N \leq e^{-Np_k}$. Second, utilizing
779 the elementary inequality $\frac{1}{n} \leq \frac{2}{n+1}$ for $n \geq 1$, we evaluate the expectation via an integral identity for the
780 binomial generating function:

$$\mathbb{E} \left[\frac{\mathbf{1}_{\{N_k > 0\}}}{N_k} \right] \leq 2 \mathbb{E} \left[\frac{1}{N_k + 1} \right] = 2 \int_0^1 \mathbb{E}[t^{N_k}] dt = 2 \int_0^1 (1 - p_k + p_k t)^N dt \leq \frac{2}{p_k(N+1)}.$$

781 Substituting these estimates into (85) yields:

$$\mathbb{E}_{\mathcal{D}_N} \left[\|y_h - y_{N,h}\|_{L^2(\mu)}^2 \right] \leq \|y\|_{L^\infty(I_R)}^2 \sum_{k \in \mathcal{K}} \left(\frac{2}{N+1} + p_k e^{-Np_k} \right).$$

782 Recall from (16) that the total number of cells is $K_h \lesssim h^{-d}$. Moreover, using $\sup_{p \geq 0} p e^{-Np} = (eN)^{-1}$,
783 we can bound the sum uniformly:

$$\mathbb{E}_{\mathcal{D}_N} \left[\|y_h - y_{N,h}\|_{L^2(\mu)}^2 \right] \leq \|y\|_{L^\infty(I_R)}^2 \left(\frac{2K_h}{N+1} + \frac{K_h}{eN} \right) \lesssim \frac{\|y\|_{L^\infty(I_R)}^2 K_h}{N} \lesssim \frac{\|y\|_{L^\infty(I_R)}^2}{Nh^d}.$$

784 Combining the bias and variance upper bounds directly gives (21). \square

785 *Proof of Corollary 3.2.* Hölder regularity gives $\omega_y(t) \leq L_y t^\alpha$ where $L_y > 0$ is the Hölder constant of $y(\cdot)$
786 on I_R . Hence (84) gives

$$\|y - y_h\|_{L^2(\mu)}^2 \leq \omega_y(\sqrt{d}h)^2 = L_y^2 (\sqrt{d}h)^{2\alpha} = L_y^2 d^\alpha h^{2\alpha} \lesssim h^{2\alpha}.$$

787 Substituting this into (21),

$$\mathbb{E}_{\mathcal{D}_N} \|y - y_{N,h}\|_{L^2(\mu)}^2 \lesssim h^{2\alpha} + (Nh^d)^{-1}.$$

788 The choice $h_N = (2R) \wedge N^{-1/(2\alpha+d)}$ balances both terms and gives

$$\mathbb{E}_{\mathcal{D}_N} \|y - y_{N,h_N}\|_{L^2(\mu)}^2 \lesssim N^{-2\alpha/(2\alpha+d)}. \quad \square$$

789 *Proof of Proposition 3.5.* If $y(\cdot)$ is α -Hölder with constant $L_y > 0$, then

$$\|y(\mathbf{x}) - y_N^V(\mathbf{x})\| = \|y(\mathbf{x}) - y(x_{\text{NN}}(\mathbf{x}))\| \leq L_y \|\mathbf{x} - x_{\text{NN}}(\mathbf{x})\|^\alpha \leq L_y R_N^\alpha, \quad (86)$$

790 where R_N is the covering radius defined in (29). Therefore,

$$\|y - y_N^V\|_{L^2(\mu)}^2 \leq \|y - y_N^V\|_{L^\infty(I_R)}^2 \leq L_y^2 R_N^{2\alpha}.$$

791 Since ρ is bounded below on the cube I_R , there exists $c > 0$ such that $\mu(B(\mathbf{x}, r) \cap I_R) \geq cr^d$ for all
792 $\mathbf{x} \in I_R$ and all sufficiently small $r > 0$. Combining this estimate with (30) directly proves (31). \square

7.3 Proofs of Section 4

Proof of Lemma 4.2. For a single cell Q_k of side length h , the volume of its trimmed boundary rim is $h^d - (h - 2\delta)^d \lesssim h^{d-1}\delta$. Summing over the $K_h \asymp h^{-d}$ cells comprising the partition, the total volume is bounded by $|\Omega_\delta| \lesssim (h^{-d})(h^{d-1}\delta) = \delta/h$. Because the probability density of μ is bounded from above by a constant, the measure directly satisfies $\mu(\Omega_\delta) \lesssim |\Omega_\delta| \lesssim \delta/h$. \square

To prove Lemma 4.3, we first isolate a preliminary estimate controlling the probability that the gap between the first two nearest-neighbor distances is small.

Lemma 7.8. *Assume $N \geq 2$, fix $\mathbf{x} \in I_R$, and let $D_1(\mathbf{x}) \leq D_2(\mathbf{x})$ be the distances from \mathbf{x} to its nearest and second-nearest points within $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$. For $r > 0$, denote $F_{\mathbf{x}}(r) := \mu(B(\mathbf{x}, r))$. Then, for every $\delta > 0$,*

$$\mathbb{P}_{\mathcal{D}_N}(D_2(\mathbf{x}) - D_1(\mathbf{x}) \leq 2\delta) \leq N(N-1) \int_0^\infty [F_{\mathbf{x}}(r+2\delta) - F_{\mathbf{x}}(r)] (1 - F_{\mathbf{x}}(r))^{N-2} dF_{\mathbf{x}}(r).$$

Proof. For each $i \in [N]$, set $Z_i := \|\mathbf{x} - \mathbf{x}_i\|$. Writing the order statistics as $Z_{[1]} \leq Z_{[2]} \leq \dots \leq Z_{[N]}$, we have $D_1(\mathbf{x}) = Z_{[1]}$ and $D_2(\mathbf{x}) = Z_{[2]}$. By symmetry,

$$\mathbb{P}(D_2(\mathbf{x}) - D_1(\mathbf{x}) \leq 2\delta) \leq N \mathbb{P}(Z_1 = Z_{[1]}, Z_{[2]} \leq Z_1 + 2\delta).$$

Conditioning with respect to Z_1 , we see that on the event $\{Z_1 = r, Z_1 = Z_{[1]}, Z_{[2]} \leq r + 2\delta\}$, at least one of the remaining $N-1$ distances must belong to $[r, r+2\delta]$, while all the others must be at least r . For a fixed index $j \in \{2, \dots, N\}$,

$$\mathbb{P}(r \leq Z_j \leq r + 2\delta) = F_{\mathbf{x}}(r+2\delta) - F_{\mathbf{x}}(r),$$

and for each $k \notin \{1, j\}$, $\mathbb{P}(Z_k \geq r) = 1 - F_{\mathbf{x}}(r)$. By independence, these probabilities multiply, and a union bound over the $N-1$ possible choices of j yields

$$\mathbb{P}(Z_1 = Z_{[1]}, Z_{[2]} \leq r + 2\delta \mid Z_1 = r) \leq (N-1) [F_{\mathbf{x}}(r+2\delta) - F_{\mathbf{x}}(r)] (1 - F_{\mathbf{x}}(r))^{N-2}.$$

Integrating with respect to the law of Z_1 , whose distribution function is $F_{\mathbf{x}}$, proves the claim. \square

Proof of Lemma 4.3. We first treat the case $N = 1$. Then $V_1 = I_R$, so $\Omega_\delta = \{\mathbf{x} \in I_R : \text{dist}(\mathbf{x}, \partial I_R) < \delta\}$, and since μ has density bounded above, $\mathbb{E}_{\mathcal{D}_N}[\mu(\Omega_\delta)] = \mu(\Omega_\delta) \lesssim \delta \leq \delta N^{1/d}$.

We may therefore assume $N \geq 2$. Let $\mathbf{x} \sim \mu$ be an independent test point; by Fubini's theorem,

$$\mathbb{E}_{\mathcal{D}_N}[\mu(\Omega_\delta)] = \mathbb{P}(\mathbf{x} \in \Omega_\delta(\mathcal{D}_N)),$$

where the probability is taken jointly over \mathcal{D}_N and \mathbf{x} . The bound is trivial whenever $\delta N^{1/d} \geq 1$ since probabilities are at most 1, so we restrict to

$$\delta N^{1/d} \leq 1. \tag{87}$$

Let $D_1(\mathbf{x})$ and $D_2(\mathbf{x})$ denote the distances from \mathbf{x} to its nearest and second-nearest sample points among $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, respectively. We claim that

$$\{\mathbf{x} \in \Omega_\delta, \text{dist}(\mathbf{x}, \partial I_R) > \delta\} \subset \{D_2(\mathbf{x}) - D_1(\mathbf{x}) \leq 2\delta\}. \tag{88}$$

Suppose that $\mathbf{x} \in \Omega_\delta$ and $\text{dist}(\mathbf{x}, \partial I_R) > \delta$, and let V_i be the Voronoi cell containing \mathbf{x} . Since $\mathbf{x} \notin V_i^\delta$, there exists $\mathbf{z} \in \partial V_i$ with $\|\mathbf{x} - \mathbf{z}\| \leq \delta$. The constraint $\text{dist}(\mathbf{x}, \partial I_R) > \delta$ excludes $\mathbf{z} \in \partial I_R$, so \mathbf{z} lies on an internal Voronoi boundary and is equidistant from at least two sample points $\mathbf{x}_i, \mathbf{x}_j$. The triangle inequality gives

$$\|\mathbf{x} - \mathbf{x}_j\| \leq \|\mathbf{x} - \mathbf{z}\| + \|\mathbf{z} - \mathbf{x}_i\| \leq \|\mathbf{x} - \mathbf{x}_i\| + 2\|\mathbf{x} - \mathbf{z}\|.$$

Taking \mathbf{x}_i as a nearest neighbor and \mathbf{x}_j as a second-nearest neighbor of \mathbf{x} , we obtain $D_2(\mathbf{x}) - D_1(\mathbf{x}) \leq 2\|\mathbf{x} - \mathbf{z}\| \leq 2\delta$, which proves (88). Consequently,

$$\mathbb{P}(\mathbf{x} \in \Omega_\delta(\mathcal{D}_N)) \leq \mathbb{P}(\text{dist}(\mathbf{x}, \partial I_R) \leq \delta) + \mathbb{P}(D_2(\mathbf{x}) - D_1(\mathbf{x}) \leq 2\delta). \tag{89}$$

823 Since μ has density bounded above and the Lebesgue measure of the δ -strip near ∂I_R is $O(\delta)$,

$$\mathbb{P}(\text{dist}(\mathbf{x}, \partial I_R) \leq \delta) \lesssim \delta. \quad (90)$$

824 For the second term, fix $\mathbf{x} \in I_R$ and define $F_{\mathbf{x}}(r) := \mu(B(\mathbf{x}, r))$. The density bound on μ controls the
825 μ -mass of a spherical annulus by its Euclidean volume:

$$F_{\mathbf{x}}(r + 2\delta) - F_{\mathbf{x}}(r) \lesssim \delta (r + 2\delta)^{d-1} \lesssim \delta r^{d-1} + \delta^d. \quad (91)$$

826 Substituting (91) into Lemma 7.8 and using the identity $(N-1) \int_0^\infty (1 - F_{\mathbf{x}})^{N-2} dF_{\mathbf{x}} = 1$ (immediate
827 from the change of variable $s = F_{\mathbf{x}}(r)$),

$$\mathbb{P}(D_2(\mathbf{x}) - D_1(\mathbf{x}) \leq 2\delta \mid \mathbf{x}) \lesssim N\delta I_{\mathbf{x}} + N\delta^d, \quad I_{\mathbf{x}} := (N-1) \int_0^\infty r^{d-1} (1 - F_{\mathbf{x}}(r))^{N-2} dF_{\mathbf{x}}(r). \quad (92)$$

828 It remains to bound $I_{\mathbf{x}}$ uniformly in $\mathbf{x} \in I_R$. Stieltjes integration by parts, whose boundary terms vanish
829 (at $r = 0$ since $d \geq 2$, at $r = \infty$ since $\text{supp } \mu \subset I_R$), gives

$$I_{\mathbf{x}} = (d-1) \int_0^\infty r^{d-2} (1 - F_{\mathbf{x}}(r))^{N-1} dr.$$

830 Let $\rho_{\min} > 0$ be a uniform lower bound for the density of μ . Since I_R is a cube, there exist constants
831 $c > 0$ and $r_0 > 0$ such that $|B(\mathbf{x}, r) \cap I_R| \geq c r^d$ for all $\mathbf{x} \in I_R$ and $0 < r \leq r_0$, hence $F_{\mathbf{x}}(r) \geq c r^d$ on
832 the same range. Using $1 - z \leq e^{-z}$, we get $(1 - F_{\mathbf{x}}(r))^{N-1} \leq e^{-c(N-1)r^d}$ for $r \leq r_0$, while for $r \geq r_0$,
833 $1 - F_{\mathbf{x}}(r) \leq 1 - F_{\mathbf{x}}(r_0)$ is uniformly bounded away from 1. Setting $D_R := \text{diam}(I_R)$ and splitting at r_0 ,

$$I_{\mathbf{x}} \lesssim \int_0^{r_0} r^{d-2} e^{-c(N-1)r^d} dr + \int_{r_0}^{D_R} r^{d-2} (1 - F_{\mathbf{x}}(r_0))^{N-1} dr.$$

834 The change of variables $u = (N-1)r^d$ in the first integral yields a bound $\lesssim N^{-(d-1)/d}$, and the second
835 is exponentially small in N and absorbed into the same bound. Hence

$$I_{\mathbf{x}} \lesssim N^{-(d-1)/d} \quad \text{uniformly in } \mathbf{x} \in I_R. \quad (93)$$

836 Substituting (93) into (92) yields

$$\mathbb{P}(D_2(\mathbf{x}) - D_1(\mathbf{x}) \leq 2\delta \mid \mathbf{x}) \lesssim \delta N^{1/d} + N\delta^d \lesssim \delta N^{1/d},$$

837 where the last step uses $N\delta^d = (\delta N^{1/d})^d \leq \delta N^{1/d}$ from (87) and $d \geq 2$. Taking expectations and
838 combining with (89) and (90), we conclude

$$\mathbb{E}_{\mathcal{D}_N} [\mu(\Omega_\delta)] \lesssim \delta + \delta N^{1/d} \lesssim \delta N^{1/d},$$

839 since $N^{1/d} \geq 1$. This proves (36). \square

840 *Proof of Corollary 4.4.* For the histogram partition, apply Theorem 4.1 with $y_N = y_{N,h}$. Then Proposition
841 3.1 and Lemma 4.2 give

$$\mathbb{E}_{\mathcal{D}_N} \left[\|y - y_{N,h}\|_{L^2(\mu)}^2 \right] \lesssim h^{2\alpha} + \frac{1}{Nh^d}, \quad \mathbb{E}_{\mathcal{D}_N} [\mu(\Omega_\delta)] \lesssim \frac{\delta}{h}.$$

842 Substituting these bounds into (34) gives the first estimate.

843 For the Voronoi partition, apply Theorem 4.1 with $y_N = y_N^V$. The expectation bound in Proposition 3.5
844 and the boundary-layer estimate in Lemma 4.3 give

$$\mathbb{E}_{\mathcal{D}_N} \left[\|y - y_N^V\|_{L^2(\mu)}^2 \right] \lesssim \left(\frac{\log N}{N} \right)^{2\alpha/d}, \quad \mathbb{E}_{\mathcal{D}_N} [\mu(\Omega_\delta)] \lesssim \delta N^{1/d}.$$

845 Substituting these bounds into (34) gives the Voronoi estimate. \square

846 **Proofs of Propositions 4.5 and 4.6**

847 *Proof of Proposition 4.5.* Define $\eta_N^2 := N^{-\frac{2\alpha}{2\alpha+d}}$, and $h_N := N^{-\frac{1}{2\alpha+d}}$, $\delta_N := h_N^{1+2\alpha}$, and apply the his-
848 togram bound in Corollary 4.4 with $h = h_N$ and $\delta = \delta_N$. By Theorem 2.5, there exists $p_N^* \in \mathbb{N}$ given by
849 (11), such that for every $p_N \geq p_N^*$ one can choose a control for (SANODE) of width p_N with $\eta^2 \asymp N^{-\frac{2\alpha}{2\alpha+d}}$,
850 whose flow is uniformly bounded on I_R . Moreover, Proposition 3.1 and Corollary 3.2 give

$$\mathbb{E}_{\mathcal{D}_N} \left[\|y - y_{N,h_N}\|_{L^2(\mu)}^2 \right] \lesssim h_N^{2\alpha} + \frac{1}{Nh_N^d} \asymp N^{-\frac{2\alpha}{2\alpha+d}},$$

851 and Lemma 4.2 yields $\mathbb{E}_{\mathcal{D}_N} [\mu(\Omega_{\delta_N})] \lesssim \delta_N h_N^{-1} = h_N^{2\alpha} = N^{-\frac{2\alpha}{2\alpha+d}}$. If $N > 2^{1+d/(2\alpha)}$ then $\delta_N < h_N/2$
852 and substituting these bounds into (37) yields the claim. \square

853 *Proof of Proposition 4.6.* Let p_N satisfy condition (11) for the Voronoi partition generated by $\{\mathbf{x}_i\}_{i=1}^N$, with
854 margin δ_N and tolerance η_N defined by

$$\eta_N^2 := \left(\frac{\log N}{N} \right)^{\frac{2\alpha}{d}}, \quad \delta_N := (\log N)^{\frac{2\alpha}{d}} N^{-\frac{2\alpha+1}{d}}.$$

855 Substituting the bounds (31) from Proposition 3.5 and (36) from Lemma 4.3 into (37) yields the conclusion.
856 \square

857 **References**

- 858 [1] A. Álvarez-López, B. Geshkovski, and D. Ruiz-Balet. Constructive approximate transport maps
859 with normalizing flows. *Applied Mathematics & Optimization*, 92(2):33, 2025. doi: 10.1007/
860 s00245-025-10299-7.
- 861 [2] A. Álvarez-López, R. Orive-Illera, and E. Zuazua. Cluster-based classification with neural ODEs via
862 control. *Journal of Machine Learning*, 4(2):128–156, 2025. doi: 10.4208/jml.241114.
- 863 [3] A. Álvarez-López, A. H. Slimane, and E. Zuazua. Interplay between depth and width for interpolation
864 in neural ODEs. *Neural Networks*, 180:106640, 2024. doi: 10.1016/j.neunet.2024.106640.
- 865 [4] F. Aurenhammer. Voronoi diagrams—a survey of a fundamental geometric data structure. *ACM*
866 *Computing Surveys*, 23(3):345–405, 1991. doi: 10.1145/116873.116880.
- 867 [5] F. Bach. Breaking the curse of dimensionality with convex neural networks. *Journal of Machine*
868 *Learning Research*, 18(19):1–53, 2017.
- 869 [6] B. Bauer, L. Devroye, M. Kohler, A. Krzyżak, and H. Walk. Nonparametric estimation of a function
870 from noiseless observations at random points. *Journal of Multivariate Analysis*, 160:93–104, 2017.
871 doi: 10.1016/j.jmva.2017.05.010.
- 872 [7] G. Biau and L. Devroye. *Lectures on the Nearest Neighbor Method*. Springer Series in the Data Sciences.
873 Springer, 2015. ISBN 9783319253862. doi: 10.1007/978-3-319-25388-6.
- 874 [8] L. Bleistein and A. Guilloux. On the generalization and approximation capacities of neural controlled
875 differential equations. In *International Conference on Learning Representations*, 2024.
- 876 [9] R. T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud. Neural ordinary differential equa-
877 tions. In *Advances in Neural Information Processing Systems*, volume 31, pages 6572–6583, 2018.
- 878 [10] J. Cheng, Q. Li, T. Lin, and Z. Shen. Interpolation, approximation, and controllability of deep neural
879 networks. *SIAM Journal on Control and Optimization*, 63(1):625–649, 2025. doi: 10.1137/23M1599744.
- 880 [11] A. Cohen, W. Dahmen, I. Daubechies, and R. DeVore. Tree approximation and optimal encoding.
881 *Applied and Computational Harmonic Analysis*, 11(2):192–226, 2001. doi: 10.1006/acha.2001.0336.
- 882 [12] R. A. DeVore. Nonlinear approximation. *Acta Numerica*, 7:51–150, 1998. doi: 10.1017/
883 S0962492900002816.
- 884 [13] R. Durrett. *Probability*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge
885 University Press, 5 edition, 2019. ISBN 9781108591034. doi: 10.1017/9781108591034.

- 886 [14] W. E. A proposal on machine learning via dynamical systems. *Communications in Mathematics and*
887 *Statistics*, 5(1):1–11, 2017. doi: 10.1007/s40304-017-0103-z.
- 888 [15] W. Grathwohl, R. T. Q. Chen, J. Bettencourt, I. Sutskever, and D. Duvenaud. FFJORD: Free-form con-
889 tinuous dynamics for scalable reversible generative models. In *International Conference on Learning*
890 *Representations*, 2019.
- 891 [16] L. Györfi, M. Kohler, A. Krzyżak, and H. Walk. *A Distribution-Free Theory of Nonparametric Regres-*
892 *sion*. Springer Series in Statistics. Springer, 2002. doi: 10.1007/b97848.
- 893 [17] J. Jia, Z. Yang, M. Wang, K. Guo, J. Yang, X. Yu, and L. Guo. Feedback favors the generalization of
894 Neural ODEs. In *International Conference on Learning Representations*, 2025.
- 895 [18] P. Kidger, J. Morrill, J. Foster, and T. Lyons. Neural controlled differential equations for irregular time
896 series. In *Advances in Neural Information Processing Systems*, volume 33, pages 6696–6707, 2020.
- 897 [19] J. Klusowski and A. Barron. Approximation by combinations of relu and squared relu ridge functions
898 with l1 and l0 controls. *IEEE Transactions on Information Theory*, 64(12):7649–7656, 2018. doi: 10.
899 1109/TIT.2018.2874447.
- 900 [20] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521:436–44, 2015.
- 901 [21] Z. Li, K. Liu, L. Liverani, and E. Zuazua. Universal approximation of dynamical systems by semiau-
902 tonomous neural odes and applications. *SIAM Journal on Numerical Analysis*, 64:193–223, 2026. doi:
903 10.1137/24M1679690.
- 904 [22] P. Marion. Generalization Bounds for Neural Ordinary Differential Equations and Deep Residual
905 Networks. In *Advances in Neural Information Processing Systems*, volume 36, pages 48918–48938,
906 2023.
- 907 [23] S. Massaroli, M. Poli, J. Park, A. Yamashita, and H. Asama. Dissecting neural ODEs. In *Advances in*
908 *Neural Information Processing Systems*, volume 33, pages 3952–3963, 2020.
- 909 [24] E. A. Nadaraya. On estimating regression. *Theory of Probability and its Applications*, 9(1):141–142,
910 1964. doi: 10.1137/1109020.
- 911 [25] A. Pinkus. *n-Widths in Approximation Theory*, volume 7 of *Ergebnisse der Mathematik und ihrer*
912 *Grenzgebiete (3)*. Springer, Berlin, Heidelberg, 1985. ISBN 978-3-540-13638-5. doi: 10.1007/
913 978-3-642-69894-1.
- 914 [26] A. Reznikov and E. B. Saff. The covering radius of randomly distributed points on a manifold. *Inter-*
915 *national Mathematics Research Notices*, 2016(19):6065–6094, 2016. doi: 10.1093/imrn/rnv342.
- 916 [27] Y. Rubanova, R. T. Q. Chen, and D. K. Duvenaud. Latent ordinary differential equations for
917 irregularly-sampled time series. In *Advances in Neural Information Processing Systems*, volume 32,
918 2019.
- 919 [28] D. Ruiz-Balet and E. Zuazua. Neural ODE Control for Classification, Approximation, and Transport.
920 *SIAM Review*, 65(3):735–773, 2023. doi: 10.1137/21M1411433.
- 921 [29] E. Sontag and H. Sussmann. Complete controllability of continuous-time recurrent neural networks.
922 *Systems & Control Letters*, 30(4):177–183, 1997. doi: [https://doi.org/10.1016/S0167-6911\(97\)00002-9](https://doi.org/10.1016/S0167-6911(97)00002-9).
- 923 [30] C. J. Stone. Optimal global rates of convergence for nonparametric regression. *The annals of statistics*,
924 pages 1040–1053, 1982.
- 925 [31] A. B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer Publishing Company, Incorporated,
926 1st edition, 2008. ISBN 978-0-387-79051-0.
- 927 [32] M. Verma and M. Kumar. Analysis of generalization capacities of Neural Ordinary Differential Equa-
928 tions. *Transactions on Machine Learning Research*, 2025.
- 929 [33] G. S. Watson. Smooth regression analysis. *Sankhyā: The Indian Journal of Statistics, Series A*, 26(4):
930 359–372, 1964.