

Approximate Controllability of Continuity Equation of Transformers

Daniel Owusu Adu[✉], Member, IEEE, and Bahman Gharesifard[✉], Senior Member, IEEE

Abstract—Building on the recent work by Geshkovski et al. (2023) which provides an interacting particle system interpretation of Transformers with a continuous-time evolution, we study the controllability attributes of the corresponding continuity equation across the landscape of probability space curves. In particular, we consider the parameters of the Transformer’s continuous-time evolution as control inputs. We prove that given an absolutely continuous probability measure and a non-local Lipschitz velocity field that satisfy a continuity equation, there exist control inputs such that the measure and the non-local velocity field of the Transformer’s continuous-time evolution approximate them, respectively, in the p -Wasserstein and L^p -sense, where $1 \leq p < \infty$.

Index Terms—Distributed parameter systems, machine learning, neural networks.

I. INTRODUCTION

REMARKABLY, the recent work [1] shows that Transformers, which needless to say have revolutionized machine learning and beyond [2], [3], [4], [5], can be viewed as interacting particle systems. To be more precise, in [1], the authors considered the dynamics

$$\dot{x}_i(t) = \mathbf{P}_{x_i(t)} \left(\sum_{j=1}^n A_j(t, x_i(t), x_j(t)) V(t) x_j(t) \right). \quad (1)$$

Here:

- 1) $\mathbf{P}_{x_i} : \mathcal{S}^{d-1} \rightarrow T_{x_i} \mathcal{S}^{d-1}$, defined as

$$\mathbf{P}_{x_i} y = y - \langle x_i, y \rangle x_i \quad (2)$$

represents the projection of $y \in \mathcal{S}^{d-1} \subset \mathbb{R}^d$ the unit sphere onto the tangent space $T_{x_i} \mathcal{S}^{d-1}$. This projection map ensures that the relative positions of neighbouring states influence the dynamics of each state.

Manuscript received 8 March 2024; revised 4 May 2024; accepted 23 May 2024. Date of publication 27 May 2024; date of current version 13 June 2024. This work was supported in part by the U.S. Department of Commerce under Grant BS123456. Recommended by Senior Editor P. Tesi. (Corresponding author: Daniel Owusu Adu.)

Daniel Owusu Adu is with the Faculty of Mathematics, University of Georgia, Athens, GA 30605 USA (e-mail: daniel.adu@uga.edu).

Bahman Gharesifard is with the Department of Electrical and Computer Engineering, University of California at Los Angeles, Los Angeles, CA 90095 USA (e-mail: gharesifard@ucla.edu).

Digital Object Identifier 10.1109/LCSYS.2024.3406255

2475-1456 © 2024 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission.
See <https://www.ieee.org/publications/rights/index.html> for more information.

- 2) The self-attention $A_j(t, x_i(t), x_j(t))$ captures the importance or attention given by the i -th state to the j -th state relative to sequence of states $(x_i(t))_{i \in [n]} := (x_1(t), \dots, x_n(t)) \in (\mathbb{R}^d)^n$, where $[n] := \{1, 2, \dots, n\} \subset \mathbb{Z}$, at time t . Specifically,

$$A_j(t, x_i(t), x_j(t)) := \frac{e^{\langle Q(t)x_i(t), M(t)x_j(t) \rangle}}{\sum_{l=1}^n e^{\langle Q(t)x_i(t), M(t)x_l(t) \rangle}},$$

where $Q(t) \in \mathbb{R}^{p \times d}$ and $M(t) \in \mathbb{R}^{p \times d}$ determine the influence of the neighboring states.

- 3) The matrix $V(t) \in \mathbb{R}^{d \times d}$ scales the dot products in the self-attention mechanism, determining the strength of interactions among the states.

The analysis in [1] shows that due to the layer normalization property, the states evolve over the unit sphere $(x_i(t))_{i \in [n]} \in (\mathcal{S}^{d-1})^n$, for $t \geq 0$. Moreover, they showed that (1) inherits some clustering property almost surely, under some suitable condition.

Related to our work is the continuity equation

$$\partial_t \mu_t(x) + \nabla \cdot \left(\mathbf{P}_x \left(\int e^{\beta \langle x, y \rangle} y d\mu_t(y) \right) \mu_t(x) \right) = 0, \quad (3)$$

on $\mathbb{R}_{\geq 0} \times \mathcal{S}^{d-1}$, with initial distribution $\mu_0 \in \mathcal{P}(\mathcal{S}^{d-1})$, where $\beta > 0$ is fixed and determined by the magnitude of Q and M . In [1, Lemma 3.5] it is shown that (3) is a Wasserstein gradient flow of some energy function on $\mathcal{P}(\mathcal{S}^{d-1})$. The interacting particle systems interpretation of Transformers is also investigated in [6], where the clustering properties of the model is investigated, and it is shown that the absence of layer normalization leads to system instability.

This letter adopts a control-theoretic perspective. In particular, we consider

$$\partial_t \mu_{it}(x_i) + \nabla \cdot (F[\mu_i, W, A_i, V, b](t, x_i) \mu_{it}(x_i)) = 0, \quad (4)$$

with initial distributions $(\mu_{i0})_{i \in [n]} \in (\mathcal{P}(\mathcal{S}^{d-1}))^n$, where

$$F[\mu_i, W, A_i, V, b](t, x_i) := \mathbf{P}_{x_i} \left(W(t) \sigma \left(\sum_{k=1}^2 \int_{\mathcal{S}^{d-1}} A_{i,k}(t, x_i, y) V_k(t) y d\mu_{it}(y) + b(t) \right) \right) \quad (5)$$

is the trainable vector field and $A_i = (A_{i,1}, A_{i,2})$, where

$$A_{i,k}(t, x_i, y) = \frac{e^{\langle Q_k(t)x_i, M_k(t)y \rangle}}{\int_{\mathcal{S}^{d-1}} e^{\langle Q_k(t)x_i, M_k(t)y \rangle} d\mu_{it}(y)}, \quad (6)$$

with $k \in \{1, 2\}$ and $i \in [n]$. Here, $W(t) \in \mathbb{R}^{d \times d}$ and $b(t) \in \mathbb{R}^d$ are the weight matrix and bias vector for the neural network, respectively, and σ is the activation function. The specific σ will be made explicit later. Note that if μ_i has density f_i , then $F[\mu_i, W, A_i, V, b](t, x_i)$ in (5) is not uniquely determined by the value of $f_i(x_i)$ but rather the value of f_i on the whole sphere S^{d-1} . Therefore, F in (5) is non-local. If μ_i in (5) is an empirical measure, then

$$F[\mu_i, W, A_i, V, b](t, x_i) := P_{x_i} \left(W(t) \sigma \left(\sum_{k=1}^2 \sum_{j=1}^n A_{i,k}(t, x_i, x_j) V_k(t) x_j + b(t) \right) \right)$$

is the complete feed-forward layer (see [1, Sec. 2.3.2]). Note that since the trainable parameters $W(t), V_k(t) \in \mathbb{R}^{d \times d}, Q_k(t), M_k(t) \in \mathbb{R}^{p \times d}$ and $b(t) \in \mathbb{R}^d$ in (5) are independent of $i \in [n]$, the system in (4)–(5) is the i -th component of a sequence. From this point on, we sometimes use the shorthand notations

$$F := (F_i)_{i \in [n]} \quad \text{where} \quad F_i := F[\mu_i, W, A_i, V, b]. \quad (7)$$

The problem under study in this letter is the following:

Problem 1.1: Given $(\chi, v) := (\chi_i, v_i)_{i \in [n]}$, a sequence of Lipschitz velocity fields and absolutely continuous probability measures, where each pair (v_i, χ_i) satisfies the continuity equation of the form

$$\partial_t v_{it} + \nabla \cdot (\chi_i[v_{it}] v_{it}) = 0, \quad (8)$$

does there exist control inputs $W(t), V_k(t) \in \mathbb{R}^{d \times d}, Q_k(t), M_k(t) \in \mathbb{R}^{p \times d}$ and $b(t) \in \mathbb{R}^d$ such that the sequence $(F, \mu) := (F_i, \mu_i)_{i \in [n]}$, where each (F_i, μ_i) in (4), approximate (χ, v) in some proper sense?

The main objective of this letter is to provide an answer under appropriate regularity assumptions. Before we state this result, it is important to point out the technical difference between our work and that of Neural ODEs, for instance, in [7], [8], [9], [10], [11]. Beyond the fact that in the Neural ODE settings, one is often concerned with approximating one state, whereas here we deal with sequences, which somehow mimic ensemble control settings [12], complication arises from the fact that the velocity fields in (5) and (8) are *non-local* [13], [14], [15], [16]. Nevertheless, we show that $(F, \mu) = (F_i, \mu_i)_{i \in [n]}$, where each (F_i, μ_i) in (4), approximate $(\chi, v) := (\chi_i, v_i)_{i \in [n]}$, where each pair (χ_i, v_i) in (8) in some proper sense. As a result of the inherent non-local velocity fields F and χ , the approximation of the measures v is in p-Wasserstein sense and that of the velocity fields χ is in L^p -sense, where $1 \leq p < \infty$. For simplicity, we only deal with the case of 1-Wasserstein sense. We assert that akin to the significant impact Neural ODEs have had on the understanding of performance and training of neural networks, it is reasonable to anticipate that the novel transformer model will assume a comparable role.

This letter is organized as follows; in Section II, we state some Preliminary result on the existence of a general continuity equation with a non-local velocity field. We state the main result in Section III and follow with the proof in Section IV.

II. PRELIMINARIES ON TRANSPORT PDEs WITH NON-LOCAL VELOCITIES

This section provides some mathematical background that will be used throughout this letter. The following notations will be needed: let $L^\infty(\mathbb{R}^d; \mathbb{R}^d)$ denote the space of essentially bounded functions from \mathbb{R}^d to \mathbb{R}^d , $L_{loc}^\infty(\mathbb{R}; \mathbb{R})$ be the space of locally essentially bounded functions from \mathbb{R} to \mathbb{R} and $C(\mathbb{R}^d; \mathbb{R}^d)$ be the space of continuous functions from \mathbb{R}^d to \mathbb{R}^d . Let $\mathcal{P}(\mathbb{R}^d)$ and $\mathcal{P}_{ac}(\mathbb{R}^d)$ be the set of probability and absolutely continuous probability measures on \mathbb{R}^d , respectively. We define the metric \mathcal{W}_1 on $\mathcal{P}(\mathbb{R}^d)$ as

$$\mathcal{W}_1(\mu_0, \mu_f) := \sup \left\{ \int_{\mathbb{R}^d} f d(\mu_0 - \mu_f) : f \in C^\infty(\mathbb{R}^d) \cap \text{Lip}_1(\mathbb{R}^d) \right\}, \quad (9)$$

where $\text{Lip}_1(\mathbb{R}^d)$ is the set of Lipschitz functions on \mathbb{R}^d with Lipschitz constant less than 1. It is well-known, see for instance [17, Th. 7.12], that the topology induced by \mathcal{W}_1 on $\mathcal{P}(\mathbb{R}^d)$ coincides¹ with weak convergence on $\mathcal{P}(\mathbb{R}^d)$. Suppose that $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a measurable map, and $\mu_0, \mu_f \in \mathcal{P}(\mathbb{R}^d)$. Then the pushforward of μ_0 ala T is denoted by $T_\# \mu_0 = \mu_f$ and given by

$$\int_{\mathbb{R}^d} g(T(x)) d\mu_0(x) = \int_{\mathbb{R}^d} g(y) d\mu_0(T^{-1}(y)) = \int_{\mathbb{R}^d} g(y) d\mu_f \quad (10)$$

where $g \in L^1(\mathbb{R}^d, \mu_f)$ is any integrable function and as usual $L^1(\mathbb{R}^d, \mu_f)$ is the space of μ_f -integrable functions defined on \mathbb{R}^d .

We collect preliminary results of general non-local transport PDE [13], [14], [15], [16]. We say that a given pair (ϑ, ς) satisfies (8) if the function

$$t \mapsto \int_{\mathbb{R}^d} f(t, x) d\varsigma_t$$

is absolutely continuous for every $f \in C_c^\infty(\mathbb{R}_{\geq 0} \times \mathbb{R}^d; \mathbb{R})$ and

$$\begin{aligned} \int_{\mathbb{R}^d} f(t, x) d\varsigma_t(x) dx - \int_{\mathbb{R}^d} f(0, x) d\varsigma_0 = \\ \int_0^t \int_{\mathbb{R}^d} (\partial_t f(s, x) + \nabla f(s, x) \cdot \vartheta[\varsigma_s](s, x)) d\varsigma_s(x) ds \end{aligned}$$

holds, for almost every $t \in \mathbb{R}_{\geq 0}$. To guarantee the existence of a solution to (8), following from [13], [14], [15], [16], the following assumptions are needed:

Assumption 1: There exists positive functions $L_1, L_2, K \in L_{loc}^\infty(\mathbb{R}_{\geq 0}; \mathbb{R})$ such that

$$\vartheta : C(\mathbb{R}_{\geq 0}; \mathcal{P}(\mathbb{R}^d)) \rightarrow C(\mathbb{R}_{\geq 0} \times \mathbb{R}^d; \mathbb{R}^d) \cap L^\infty(\mathbb{R}_{\geq 0} \times \mathbb{R}^d; \mathbb{R}^d)$$

satisfies:

- 1) the inequality

$$\|\vartheta[\varsigma_t](t, x) - \vartheta[\varsigma_t](t, y)\|_{\mathbb{R}^d} \leq L_1(t) \|x - y\|_{\mathbb{R}^d}, \quad (11)$$

for all $\varsigma_t \in \mathcal{P}(\mathbb{R}^d)$ and $x, y \in \mathbb{R}^d$ and $t \in \mathbb{R}_{\geq 0}$.

¹The distance \mathcal{W}_1 metrizes the weak convergence of measures only if the measure has finite first moment. This condition is satisfied whenever the measures have compact support.

2) the inequality

$$\|\vartheta[\varsigma_t](t, x)\|_{\mathbb{R}^d} \leq L_2(t)(1 + \|x\|_{\mathbb{R}^d}), \quad (12)$$

for all $\varsigma_t \in \mathcal{P}(\mathbb{R}^d)$ and $x, y \in \mathbb{R}^d$ and $t \in \mathbb{R}_{\geq 0}$.

3) the inequality

$$\|\vartheta[\varsigma_0] - \vartheta[\varsigma_1]\|_{L^\infty(\mathbb{R}_{\geq 0}; C^0(\mathbb{R}^d))} \leq K(t)W_1(\varsigma_0, \varsigma_1) \quad (13)$$

for all $\varsigma_0, \varsigma_1 \in \mathcal{P}(\mathbb{R}^d)$ and $t \in \mathbb{R}_{\geq 0}$.

It is well known (see for instance [16, Th. 2.3] and reference therein) that if ϑ satisfies Assumption 1, there exists a unique solution $\varsigma \in C(\mathbb{R}_{\geq 0}; \mathcal{P}_{ac}(\mathbb{R}^d))$ to (8), whenever $\varsigma_0 \in \mathcal{P}_{ac}(\mathbb{R}^d)$.

III. MAIN RESULTS

This section presents the key findings of this letter. To this end, we fix the activation function $\sigma : \mathbb{R}^d \rightarrow \mathbb{R}^d$ to be the rectified linear function (ReLU) applied to each component of a vector $x_i = (x_{i1}, \dots, x_{id}) \in \mathbb{R}^d$:

$$\sigma(x_i) = \text{ReLU}(x_i) := (\text{ReLU}(x_{i1}), \dots, \text{ReLU}(x_{id})), \quad (14)$$

where $\text{ReLU}(x_{ij}) := \max\{x_{ij}, 0\}$, $j \in [d]$. Under the activation function, [18] has laid the groundwork by establishing that Transformers characterized as

$$\begin{aligned} \mathcal{T} &:= \{(x_i)_{i \in [n]} \mapsto (F[\mu_i, W, A_i, V, b](x_i))_{i \in [n]} \\ &:= (F[\mu_i, W, A_i, V, b](x_i + \omega_i))_{i \in [n]} : W, V_k \in \mathbb{R}^{d \times d}, \\ &\quad Q_k, M_k \in \mathbb{R}^{p \times d} \text{ and } b, \omega_i \in \mathbb{R}^d, \text{ where } k \in \{1, 2\}\}, \end{aligned} \quad (15)$$

where μ_i is a given empirical measure, serves as universal approximators of continuous sequence-to-sequence functions defined on a compact domain $f \in C((\mathcal{S}^{d-1})^n; (\mathbb{R}^d)^n)$, where

$$C((\mathcal{S}^{d-1})^n; (\mathbb{R}^d)^n) := \left\{ f : (\mathcal{S}^{d-1})^n \rightarrow (\mathbb{R}^d)^n : f \text{ is continuous functions} \right\}.$$

In particular, we later use the following universal approximation result, presented as [18, Th. 3].

Theorem 1: Let $1 \leq p < \infty$ and $\epsilon > 0$. Then for any $f \in C((\mathcal{S}^{d-1})^n; (\mathbb{R}^d)^n)$, there exists $g \in \mathcal{T}$ such that

$$d(f, g) := \left(\int_{(\mathcal{S}^{d-1})^n} \|f((x_i)_{i \in [n]}) - g((x_i)_{i \in [n]})\|^p dx_1 \dots dx_n \right)^{\frac{1}{p}} \leq \epsilon. \quad (16)$$

We state here that the universal approximation result also holds when $k \geq 3$ in (15), see [18]. We now state the main result of this letter.

Theorem 2: Consider $(\chi, v) := (\chi_i, v_i)_{i \in [n]}$, where for each $i \in [n]$, we have that χ_i satisfies Assumption 1 and $v_i \in C(\mathbb{R}_{\geq 0}; \mathcal{P}_{ac}(\mathcal{S}^{d-1}))$ is the corresponding unique solution to (8) with initial measure $v_{i0} \in \mathcal{P}_{ac}(\mathcal{S}^{d-1})$. Then, for any $\epsilon > 0$, there exist a final time $t_f > 0$ and piecewise constant in time control inputs

$$W, V_k \in L^\infty([0, t_f]; \mathbb{R}^{d \times d}) \quad Q_k, M_k \in L^\infty([0, t_f]; \mathbb{R}^{p \times d}),$$

where $k \in \{1, 2\}$ and $b, \omega \in L^\infty([0, t_f]; \mathbb{R}^d)$ such that the pair of sequence $(F, \mu) := (F_i, \mu_i)_{i \in [n]}$, where each (F_i, μ_i) is

characterized in (4)–(5), approximates (χ, v) in the following sense:

$$d(F(t, \cdot), \chi(t, \cdot)) \leq \epsilon,$$

where $\chi(t, (x_i)_{i \in [n]}) := (\chi_i[v_{it}](t, x_i))_{i \in [n]}$, for all $t \in [0, t_f]$ and

$$\sup_{t \in [0, t_f]} W_1(\mu_{it}, v_{it}) \leq \epsilon,$$

for all $i \in [n]$.

The theorem ensures that there exist appropriate control inputs such that (F, μ) characterized in (4)–(5) can approximate a given sequence of velocity fields and measures that solves (8).

IV. PROOFS OF MAIN RESULT

This section is devoted to presenting the proof of Theorem 2. The proof is organized into two subsections: firstly, we provide some properties on a generic sequence $\{(\chi_i^m, v_i^m)\}_{m \geq 1}$, where χ_i^m is piecewise constant in time vector field that approximate the solutions (χ_i, v_i) to the continuity equations (8) in Proposition 2. Next, we prove in Corollary 1 that the set of control inputs W, V_k, Q_k, M_k , where $k \in \{1, 2\}$ and b, ω can be selected in such a way as to ensure that the Transformer in (4) generates a sequence (F^N, μ^N) , where F^N is in (15) that have the required properties in the former subsection. This core idea is inspired from [8]; however, in our derivations, the technical details differ as we deal with non-local vector fields and approximating sequence-to-sequence functions. For that part, we utilize Theorem 1 and modified techniques in [13], [14], [15], [16], [19].

A. Approximating the Solutions to (8)

We start with a stepping stone result on the support of the solutions to (8).

Proposition 1: Let $i \in [n]$ and consider the continuity equation (8) where each χ_i satisfies Assumption 1 and $v_{i0} \in \mathcal{P}_{ac}(\mathcal{S}^{d-1})$. There exists a final time $t_f > 0$ such that the solution $v_i \in C(\mathbb{R}_{\geq 0}; \mathcal{P}_{ac}(\mathbb{R}^d))$ to (8), satisfies

$$\text{supp}(v_{it}) \subset \mathcal{S}^{d-1}, \quad (17)$$

for all $t \in [0, t_f]$.

Proof: Suppose that χ_i in (8) satisfies Assumption 1 and $v_{i0} \in \mathcal{P}_{ac}(\mathcal{S}^{d-1})$. Following from [16, Th. 2.3] the unique solution $v_i \in C(\mathbb{R}_{\geq 0}; \mathcal{P}_{ac}(\mathbb{R}^d))$ is characterized as

$$v_{it} = \Phi_i(t, \cdot)_\# v_{i0}, \quad (18)$$

where $\Phi_i(t, \cdot)$ is the diffeomorphic flow on \mathbb{R}^d that satisfies

$$\partial_t \Phi_i(t, x) = \chi_i[v_{it}](t, \Phi_i(t, x)) \quad \text{and} \quad \Phi_i(0, x) = x, \quad (19)$$

on $\mathbb{R}_+ \times \mathbb{R}^d$, where $i \in [n]$. We proceed to show that there exists t_f such that (19) admits a unique solution on $[0, t_f]$. To this end, let $T > 0$ be fixed. From (12), since $\|\chi_i[v_{it}](t, x)\|_{\mathbb{R}^d} \leq 2\bar{L}_{i2}$ holds for all $(t, x) \in [0, T] \times \mathcal{S}^{d-1}$, where $\bar{L}_{i2} = \text{essup}_{t \in [0, T]} L_{i2}(t)$. If $t_f \leq \min_{i \in [n]} \frac{1}{2\bar{L}_{i2}}$, then from [20, Ch. 1, Th. 1], for any $x \in \mathcal{S}^{d-1}$, we have that (19) admits a solution on $[0, t_f]$. For uniqueness, since the Lipschitz

condition (11) holds, for any $(t, x), (t, y) \in [0, T] \times \mathcal{S}^{d-1}$, from [20, Ch. 1, Th. 2] we have that the unique solution satisfies $\Phi_i(t, \text{supp}(v_{i0})) \subset \mathcal{S}^{d-1}$, where $t \in [0, t_f]$.

We now show (17). Let $z \in \text{supp}(v_{it})$. Then, from (18) and (10), since $v_{it}(\mathcal{B}_z(\epsilon)) = v_{i0}(\Phi_i^{-1}(t, \mathcal{B}_z(\epsilon))) > 0$ holds, for all $\epsilon > 0$, where $\mathcal{B}_z(\epsilon) \subset \mathbb{R}^d$ is the ball with radius ϵ centred at $z \in \text{supp}(v_{it})$, we have that $\Phi_i^{-1}(t, \mathcal{B}_z(\epsilon)) \subset \text{supp}(v_{i0})$ holds, for all $\epsilon > 0$. Since $\Phi(t, \cdot)$ is a diffeomorphic flow on \mathbb{R}^d and $t \mapsto \Phi(t, \cdot)$ is continuous map, we have that $\mathcal{B}_z(\epsilon) \subset \Phi_i(t, \text{supp}(v_{i0}))$, for all $\epsilon > 0$. Since $\Phi_i(t, \text{supp}(v_{i0})) \subset \mathcal{S}^{d-1}$, holds, for all $t \in [0, t_f]$, we have that $\mathcal{B}_z(\epsilon) \subset \mathcal{S}^{d-1}$, for all $\epsilon > 0$. This completes the proof. ■

Next, we provide an approximation result for (8).

Proposition 2: Suppose, for each $i \in [n]$, we have that χ_i satisfies Assumption 1 and $v_i \in C([0, t_f]; \mathcal{P}_{ac}(\mathcal{S}^{d-1}))$ is the corresponding solution in (8) with initial distribution $v_{i0} \in \mathcal{P}_{ac}(\mathcal{S}^{d-1})$. Then, for each $i \in [n]$, there exists $\{(\chi_i^m, v_i^m)\}_{m \geq 1}$ such that each χ_i^m is piecewise constant in time and (χ_i^m, v_i^m) solves (8) with initial distribution $v_{i0} \in \mathcal{P}_{ac}(\mathcal{S}^{d-1})$ and $\chi_i^m : C([0, t_f]; \mathcal{P}_{ac}(\mathcal{S}^{d-1})) \rightarrow C([0, t_f] \times \mathcal{S}^{d-1}; \mathcal{S}^{d-1}) \cap L^\infty([0, t_f] \times \mathcal{S}^{d-1}; \mathcal{S}^{d-1})$ uniformly converges to χ_i and $v_i^m \in C([0, t_f]; \mathcal{P}_{ac}(\mathcal{S}^{d-1}))$ weakly converges to v_i .

Proof: Let $0 := t_0^m < t_1^m < \dots < t_{2^m-1}^m := t_f$, where $t_\ell^m = \ell 2^{-m} t_f$ be a regular partition of $[0, t_f]$ into 2^m subintervals. Given (χ, v) , let

$$\chi_i^m[v_{it}] := \chi_i[v_{it}], \quad (20)$$

where $t \in [t_\ell^m, t_{\ell+1}^m]$, then χ_i^m is piecewise constant in time velocity field on $C([0, t_f]; \mathcal{P}_{ac}(\mathcal{S}^{d-1}))$. Since χ_i satisfies Assumption 1, we have that χ_i^m satisfies

$$\begin{aligned} \|\chi_i^m[v_{it}](t, x) - \chi_i^m[v_{it}](t, y)\|_{\mathbb{R}^d} &\leq L_{i2} \|x - y\|_{\mathbb{R}^d}, \quad \text{and} \\ \|\chi_i^m[v_{it}](t, x)\|_{\mathbb{R}^d} &\leq 2L_{i2}, \end{aligned} \quad (21)$$

where $L_{i\alpha} = \text{esssup}_{t \in [0, t_f]} L_{i\alpha}(t)$, with $\alpha \in \{1, 2\}$, hold for all $x, y \in \mathcal{S}^{d-1}$. Therefore, the unique solution v_i^m to

$$\partial_t v_i + \nabla \cdot (\chi_i^m[v_i] v_i) = 0 \quad \text{and} \quad v_{i0}^m = v_{i0}$$

is characterized by

$$\begin{aligned} v_{it}^m &= \Phi_i^m(t, \cdot)_\# v_{it}^m, \\ \partial_t \Phi_i^m(t, x) &= \chi_i[v_{it}^m](t_\ell^m, \Phi_i^m(t, x)), \end{aligned} \quad (22)$$

and $\Phi_i^m(t_\ell^m, x) := \Phi_i(t_\ell^m, x)$, for all $t \in [t_\ell^m, t_{\ell+1}^m]$, where Φ_i is as given in (18). We show that, for each $i \in [n]$, the pair (χ_i^m, v_i^m) , characterized in (20) and (22), respectively, converges to (χ_i, v_i) in some sense.

We show that for any $\epsilon > 0$, there exists $N_0 \in \mathbb{N}$ such that for any $m \geq N_0$, we have

$$\|\chi_i^m - \chi_i\| \leq \epsilon,$$

where

$$\|\chi_i^m - \chi_i\| := \sup_{v \in C([t_\ell^m, t_{\ell+1}^m]; \mathcal{P}(\mathcal{S}^{d-1}))} \|\chi_i^m[v_i] - \chi_i[v_i]\|_{L^\infty} \quad (23)$$

and $L^\infty := L^\infty([t_\ell^m, t_{\ell+1}^m] \times \mathcal{S}^{d-1}; \mathcal{S}^{d-1})$. To this end, using (20) and since

$$\|\chi_i^m[v_{it}] - \chi_i[v_{it}]\|_{L^\infty} = \|\chi_i[v_{it}^m] - \chi_i[v_{it}]\|_{L^\infty},$$

for all $t \in [t_\ell^m, t_{\ell+1}^m]$, using (13), we have that

$$\begin{aligned} \|\chi_i^m[v_{it}] - \chi_i[v_{it}]\|_{L^\infty} &\leq K_i \mathcal{W}_2(v_{it}^m, v_{it}) \\ &\leq K_i a_i \frac{t_f}{2^m}, \end{aligned} \quad (24)$$

for all $v_i \in C([t_\ell^m, t_{\ell+1}^m]; \mathcal{P}(\mathcal{S}^{d-1}))$, where the second inequality follows from the fact that the solution curve $t \mapsto v_{it}$ in (8) is a_i -Lipschitz continuous, for some positive constant $a_i \in \mathbb{R}$ (see for instance [15, Proposition 5]). Therefore, we conclude that χ_i^m uniformly converges to χ_i .

To show that the sequence $(v_i^m)_{m \geq 1}$ weakly converges to v_i , we first show that the family of functions $t \mapsto v_{it}^m$ is equi-continuous and equi-bounded. To this end, equi-bounded follows from the fact that from Proposition 1, we have that $\text{supp}(v_{it})$, $\text{supp}(v_{it}^m) \subset \mathcal{S}^{d-1}$, for all $t \in [0, t_f]$ and $m \geq 1$. For equi-continuous, from (20), since $\|\chi_i^m\| \leq 2L_{i2}$, following from [15, Proposition 1] we have that

$$\mathcal{W}_1(v_{it}^m, v_{it}^m) \leq 2L_{i2}|t - t_\ell^m|. \quad (25)$$

Therefore, the family of functions $t \mapsto v_{it}^m$ is equi-Lipschitz and hence equi-continuous. Hence, by Arzela-Ascoli Theorem, we have that the sequence $(v_i^m)_{m \geq 1}$ admits a subsequence $(v_i^{m_r})_{r \geq 1}$ that weakly converges to \tilde{v}_i . We proceed to show that the limit measure \tilde{v}_i satisfies

$$\int_0^{t_f} \int_{\mathcal{S}^{d-1}} (\partial_t f(t, x) + \nabla f(t, x) \cdot \chi_i[\tilde{v}_{it}](t, x)) d\tilde{v}_{it}(x) dt = 0, \quad (26)$$

for every $f \in C^\infty([0, t_f] \times \mathcal{S}^{d-1}; \mathbb{R})$, where $f(t_f, x) = 0$ for all $x \in \mathcal{S}^{d-1}$. We prove this by showing that the following statements hold:

1)

$$\lim_{r \rightarrow \infty} \int_0^{t_f} \int_{\mathcal{S}^{d-1}} \partial_t f(t, x) d(\tilde{v}_{it}(x) - v_{it}^{m_r}(x)) dt = 0.$$

2)

$$\lim_{r \rightarrow \infty} \int_0^{t_f} \int_{\mathcal{S}^{d-1}} \nabla f(t, x) \cdot (\chi_i[\tilde{v}_{it}](t, x) - \chi_i^{m_r}[v_{it}^{m_r}](t, x)) d\tilde{v}_{it}(x) dt = 0.$$

3)

$$\lim_{r \rightarrow \infty} \int_0^{t_f} \int_{\mathcal{S}^{d-1}} \nabla f(t, x) \cdot \chi_i^{m_r}[v_{it}^{m_r}](t, x) d(\tilde{v}_{it}(x) - v_{it}^{m_r}(x)) dt = 0.$$

For Statement 1: using (9) and since

$$\left| \int_0^{t_f} \int_{\mathcal{S}^{d-1}} \partial_t f(t, x) d(\tilde{v}_{it}(x) - v_{it}^{m_r}(x)) dt \right| \leq \max_{(t,x) \in [0, t_f] \times \mathcal{S}^{d-1}} \|\partial_t f(t, x)\|_{t_f} \sup_{t \in [0, t_f]} \mathcal{W}_1(\tilde{v}_{it}, v_{it}^{m_r})$$

and $v_{it}^{m_r}$ weakly converges to \tilde{v}_{it} , uniformly in $t \in [0, t_f]$, we conclude that Statement 1 holds.

For Statement 2, since

$$\left| \int_0^{t_f} \int_{\mathcal{S}^{d-1}} \nabla f(t, x) \cdot (\chi_i[\tilde{v}_{it}](t, x) - \chi_i^{m_r}[v_{it}^{m_r}](t, x)) d\tilde{v}_{it}(x) dt \right|$$

$$\left| \chi_i^{m_r} [v_{it}^{m_r}] (t, x) d\tilde{v}_{it}(x) dt \right| \leq \max_{(t,x) \in [0,t_f] \times \mathcal{S}^{d-1}} \|\partial_t f(t, x)\| \\ \sum_{l=0}^{2^m-1} \int_{t_l}^{t_{l+1}} \int_{\mathcal{S}^{d-1}} \|\chi_i [\tilde{v}_{it}] (t, x) - \chi_i^{m_r} [v_{it}^{m_r}] (t, x)\| d\tilde{v}_{it}(x) dt,$$

from (13), we have that

$$\left\| \chi_i [\tilde{v}_{it}] - \chi_i [v_{it}^{m_r}] \right\| \leq K \mathcal{W}_1 (\tilde{v}_{it}, v_{it}^{m_r}) \\ \leq K \left(\mathcal{W}_1 (\tilde{v}_{it}, v_{it}^{m_r}) + \mathcal{W}_1 (v_{it}^{m_r}, v_{it}^{m_r}) \right) \\ \leq K \left(\mathcal{W}_1 (\tilde{v}_{it}, v_{it}^{m_r}) + 2L_{i2} \frac{t_f}{2^{m_r}} \right),$$

where the last term in the later inequality follows from (25). Since $v_{it}^{m_r}$ weakly converges to \tilde{v}_{it} , we have that Statement 2 holds.

For Statement 3, from (12), we have that

$$\left| \int_0^{t_f} \int_{\mathcal{S}^{d-1}} \nabla f(t, x) \cdot \chi_i^{m_r} [v_{it}^{m_r}] (t, x) d(\tilde{v}_{it}(x) - v_{it}^{m_r}(x)) dt \right| \leq \max_{(t,x) \in [0,t_f] \times \mathcal{S}^{d-1}} \|\partial_t f(t, x)\| 2L_{i2} t_f \mathcal{W}_1 (\tilde{v}_{it}, v_{it}^{m_r}).$$

Using similar arguments in Statement 2, we conclude that Statement 3 holds. This completes the proof. ■

B. Using Transformers to Approximate the Solutions to a Continuity Equation

We are now ready to present our main proof. We first state a corollary of the universal approximation result which we later combine with the observation made in the previous section to prove Theorem 2.

Corollary 1: Given $(\chi, v) := (\chi_i, v_i)_{i \in [n]}$, where

$$\chi (t, (x_i)_{i \in [n]}) := (\chi_i [v_{it}] (t, x_i))_{i \in [n]}. \quad (27)$$

Suppose that χ_i satisfies Assumption 1 and $v_i \in C([0, t_f]; \mathcal{P}_{ac}(\mathcal{S}^{d-1}))$ is the solution in (8) with initial distribution $v_{i0} \in \mathcal{P}_{ac}(\mathcal{S}^{d-1})$. Furthermore, suppose that the time component of χ_i remain constant on $[t_\ell^N, t_{\ell+1}^N]$, where $\ell \in \{0, \dots, 2^N - 1\}$ with $t_\ell^N := \ell 2^{-N} t_f$. Then, there exist sequence of functions $F^N(t, \cdot) \in \mathcal{T}$, where $t \mapsto F^N(t, \cdot)$ is constant on the interval $[t_\ell^N, t_{\ell+1}^N]$, such that

$$\lim_{N \rightarrow \infty} d(F^N(t, \cdot), \chi(t, \cdot)) = 0$$

for all $t \in [0, t_f]$.

Proof: Given (χ, v) , consider (27). Then, by assumption and from (11), since $(t, (x_i)_{i \in [n]}) \mapsto \chi(t, (x_i)_{i \in [n]})$ is constant on $[t_\ell^N, t_{\ell+1}^N]$ and Lipschitz over $(\mathcal{S}^{d-1})^n$, we have that $\chi \in C((\mathcal{S}^{d-1})^n; (\mathbb{R}^d)^n)$ on $[t_\ell^N, t_{\ell+1}^N]$. Therefore, from Theorem 1 given an empirical measure μ_i , we have that, for every $\epsilon > 0$, there exists $F^N(t, \cdot) \in \mathcal{T}$, piecewise constant in time on $[0, t_f]$ such that $d(F^N(t, \cdot), \chi(t, \cdot)) \leq \epsilon$ holds, for all $t \in [t_\ell^N, t_{\ell+1}^N]$. This completes the proof. ■

We are now ready to provide the proof of Theorem 2.

Proof of Theorem 2: Given $(\chi, v) := (\chi_i, v_i)_{i \in [n]}$, where for each $i \in [n]$, we have that χ_i satisfies Assumption 1 and $v_i \in C(\mathbb{R}_{\geq 0}; \mathcal{P}_{ac}(\mathcal{S}^{d-1}))$ is the corresponding unique solution to (8) with initial measure $v_{i0} \in \mathcal{P}_{ac}(\mathcal{S}^{d-1})$. From Proposition 2,

we have that for each $i \in [n]$, there exists $(\chi_i^m, v_i^m)_{m \geq 1}$, where the pair (χ_i^m, v_i^m) is characterized in (20) and (22), respectively, such that

$$\lim_{m \rightarrow \infty} \|\chi_i^m - \chi_i\| = 0, \quad \text{and} \quad \lim_{m \rightarrow \infty} \mathcal{W}_1 (v_{it}^m, v_{it}) = 0, \quad (28)$$

for all $t \in [0, t_f]$. Here $\|\cdot\|$ is defined in (23) and, for any $m \geq 1$, we have that χ_i^m is constant in time on $[t_\ell^m, t_{\ell+1}^m]$ with $t_\ell^m = \ell 2^{-m} t_f$. Let $m \geq 1$ be fixed. Then, from Corollary 1, given an empirical measure μ_i , there exists a sequence $\{F^N\}_{N \geq 1} \subset \mathcal{T}$ in (15) such that

$$\lim_{N \rightarrow \infty} d(F^N(t, \cdot), \chi^m(t, \cdot)) = 0 \quad (29)$$

holds, for all $t \in [0, t_f]$. Furthermore, from (7), we have that

$$\left\{ F^N := (F_i^N)_{i \in [n]} \right\}_{N \geq 1}$$

admits the following characterization: for a given empirical measure μ_i , there exists piecewise time control inputs $W^N, V_k^N \in L^\infty([0, t_f]; \mathbb{R}^{d \times d})$, $Q_k^N, M_k^N \in L^\infty([0, t_f]; \mathbb{R}^{p \times d})$ and $b^N, \omega_i^N \in L^\infty([0, t_f]; \mathbb{R}^d)$, where $k \in \{1, 2\}$ such that for $N \geq 1$, we have that

$$F_i^N(t, x_i) := \mathbf{P}_{x_i} \left(W^N(t) \sigma \left(\sum_{k=1}^2 \sum_{j=1}^n A_{i,k}^N(t, x_i, x_j) V_k^N(t) x_j + b^N(t) \right) \right).$$

We proceed to generate the corresponding measures μ_i^N . To this end, we consider the piecewise constant in time vector field

$$F_i^N(t, x_i) := \mathbf{P}_{x_i} \left(W^N(t) \sigma \left(\sum_{k=1}^2 \int_{\mathcal{S}^{d-1}} A_{i,k}^N(t, x_i, y) V_k^N(t) y d\mu_{it}(y) + b^N(t) \right) \right), \quad (30)$$

and show that, for a fixed N and $i \in [n]$, we have that $F_i^N(t, x_i)$ above satisfies Assumption 1.

To this end, since the projection map \mathbf{P} in (2) is uniformly bounded and Lipschitz on \mathcal{S}^{d-1} , it is enough to show that

$$G_i^N(t, x_i) := W^N(t) \sigma \left(\sum_{k=1}^2 \int_{\mathcal{S}^{d-1}} A_{i,k}^N(t, x_i, y) V_k^N(t) y d\mu_{it}(y) + b^N(t) \right)$$

satisfies Assumption 1. For Statement 1 in Assumption 1, from (6) since

$$\nabla_{x_i} G_i^N(t, x_i) = W^N(t) \sigma \left(\sum_{k=1}^2 \int_{\mathcal{S}^{d-1}} \nabla_{x_i} A_{i,k}^N(t, x_i, y) V_k^N(t) y d\mu_{it}(y) \right)$$

for every $i \in [n]$ and $m \geq 1$, we have that

$$\|\nabla_{x_i} G_i^N\|_{L^\infty} \leq 2 \|W^N\|_{L^\infty} \left(\sum_{k=1}^2 \|M_k^N\|_{L^\infty} \|Q_k^N\|_{L^\infty} \|V_k^N\|_{L^\infty} \right).$$

Therefore, we conclude that $x_i \mapsto G_i^N(t, x_i)$ is Lipschitz.

For uniform boundedness, since

$$\left\| W^N(t) \sigma \left(\sum_{k=1}^2 \int_{S^{d-1}} A_{i,k}^N(t, x_i, y) V_k^N(t) y d\mu_{it}(y) + b^N(t) \right) \right\| \leq \| W^N \|_{L^\infty} \left(\sum_{k=1}^2 \| V_k^N \|_{L^\infty} + \| b^N \|_{L^\infty} \right)$$

we have that

$$\| G_i^N \|_{L^\infty} \leq \| W^N \|_{L^\infty} \left(\sum_{k=1}^2 \| V_k^N \|_{L^\infty} + \| b^N \|_{L^\infty} \right).$$

This concludes that $\mu_i \mapsto G_i^N$ is uniformly bounded.

Lastly, we show that $\mu_i \mapsto G_i^N$ is Lipschitz. From (14), since for $x_i = (x_{i1}, \dots, x_{id}) \in S^{d-1}$ we have that $\text{ReLU}(x_{ij})$ is Lipschitz and

$$\left\| \sum_{k=1}^2 \int_{S^{d-1}} A_{i,k}^N(t, x_i, y) V_k^N(t) y d(\mu_{it1}(y) - \mu_{it2}(y)) \right\| \leq \left(\sum_{k=1}^2 \| A_{i,k}^N(t) \|_{L^\infty} \| V_k^N \|_{L^\infty} \right) W_1(\mu_{it1}, \mu_{it2}).$$

we conclude that $\mu_i \mapsto G_i^N$ is Lipschitz. Therefore, $F_i^N(t, x_i)$ satisfies Assumption 1. From [16, Th. 2.3], we have that there exists a unique solution μ_i^N to

$$\partial_t \mu_{it}(x_i) + \nabla \cdot (F_i^N(t, x_i) \mu_{it}(x_i)) = 0, \quad (31)$$

with initial distribution $\mu_{i0} \in (\mathcal{P}_{ac}(S^{d-1}))^n$, where $i \in [n]$. Furthermore, from (29), for a fixed $m \geq 1$, by Arzela-Ascoli Theorem, we have that

$$\lim_{N \rightarrow \infty} W_1(\mu_{it}^N, v_{it}^m) = 0 \quad (32)$$

uniformly in $t \in [0, t_f]$.

Since $W_1(\mu_{it}^N, v_{it}) \leq W_1(\mu_{it}^N, v_{it}^m) + W_1(v_{it}^m, v_{it})$, for all $m \geq 1$, from (32), we have that $0 \leq \lim_{N \rightarrow \infty} W_1(\mu_{it}^N, v_{it}) \leq W_1(v_{it}^m, v_{it})$, for all $m \geq 1$, uniformly in $t \in [0, t_f]$. From (28), we have that $\lim_{N \rightarrow \infty} W_1(\mu_{it}^N, v_{it}) = 0$, uniformly in $t \in [0, t_f]$. Therefore, for every $\epsilon > 0$, there exists $N_0 \in \mathbb{N}$ such that for $N \geq N_0$, we have that $\sup_{t \in [0, t_f]} W_1(\mu_{it}^N, v_{it}) \leq \epsilon$, holds, for all $i \in [n]$. Similarly, since $d(F^N(t, \cdot), \chi(t, \cdot)) \leq d(F^N(t, \cdot), \chi^m(t, \cdot)) + d(\chi^m(t, \cdot), \chi(t, \cdot))$, from (29), we conclude that given (χ, v) , there exists sequence of piecewise constant control inputs $W^N, V_k^N \in L^\infty([0, t_f]; \mathbb{R}^{d \times d}), Q_k^N$,

$M_k^N \in L^\infty([0, t_f]; \mathbb{R}^{p \times d})$ and $b^N, \omega_i^N \in L^\infty([0, t_f]; \mathbb{R}^d)$ such that $\lim_{N \rightarrow \infty} d(F^N(t, \cdot), \chi(t, \cdot)) = 0$, for all $t \in [0, t_f]$. This finishes the proof. ■

REFERENCES

- [1] B. Geshkovski, C. Letrouit, Y. Polyanskiy, and P. Rigollet, "A mathematical perspective on transformers," 2023, *arXiv:2312.10794*.
- [2] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," 2018, *arXiv:1810.04805*.
- [3] T. Brown et al., "Language models are few-shot learners," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 33, 2020, pp. 1877–1901.
- [4] M. Chen et al., "Generative pretraining from pixels," in *Proc. Int. Conf. Mach. Learn.*, 2020, pp. 1691–1703.
- [5] A. Vaswani et al., "Attention is all you need," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 30, 2017, pp. 1–15.
- [6] B. Geshkovski, C. Letrouit, Y. Polyanskiy, and P. Rigollet, "The emergence of clusters in self-attention dynamics," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 36, 2024, pp. 1–55.
- [7] D. Ruiz-Balet and E. Zuazua, "Neural ODE control for classification, approximation, and transport," *SIAM Rev.*, vol. 65, no. 3, pp. 735–773, 2023.
- [8] K. Elamvazhuthi, B. Gharesifard, A. L. Bertozzi, and S. Osher, "Neural ODE control for trajectory approximation of continuity equation," *IEEE Control Syst. Lett.*, vol. 6, pp. 3152–3157, 2022.
- [9] P. Tabuada and B. Gharesifard, "Universal approximation power of deep residual neural networks via nonlinear control theory," 2020, *arXiv:2007.06007*.
- [10] A. Agrachev and A. Sarychev, "Control on the manifolds of mappings with a view to the deep learning," *J. Dyn. Control Syst.*, vol. 28, no. 4, pp. 989–1008, 2022.
- [11] C. Cuchiero, M. Larsson, and J. Teichmann, "Deep neural networks, generic universal interpolation, and controlled ODEs," *SIAM J. Math. Data Sci.*, vol. 2, no. 3, pp. 901–919, 2020.
- [12] J.-S. Li and N. Khaneja, "Control of inhomogeneous quantum ensembles," *Phys. Rev. A*, vol. 73, no. 3, 2006, Art. no. 30302.
- [13] L. Ambrosio and W. Gangbo, "Hamiltonian ODEs in the wasserstein space of probability measures," *Commun. Pure Appl. Math.*, vol. 61, no. 1, pp. 18–53, 2008.
- [14] H. Dong, "Well-posedness for a continuity equation with non-local velocity," *J. Funct. Anal.*, vol. 255, no. 11, pp. 3070–3097, 2008.
- [15] B. Piccoli and F. Rossi, "Continuity equation with non-local velocity in Wasserstein spaces: Convergence of numerical schemes," *Acta Applicandae Mathematicae*, vol. 124, pp. 73–105, Apr. 2013.
- [16] B. Piccoli, F. Rossi, and E. Trélat, "Control to flocking of the kinetic Cucker-Smale model," *SIAM J. Math. Anal.*, vol. 47, no. 6, pp. 4685–4719, 2015.
- [17] C. Villani, *Optimal Transport: Old and New*, vol. 338. Berlin, Germany: Springer, 2009.
- [18] C. Yun, S. Bhojanapalli, A. S. Rawat, S. J. Reddi, and S. Kumar, "Are transformers universal approximators of sequence-to-sequence functions?" 2019, *arXiv:1912.10077*.
- [19] H. O. Fattorini, *Infinite Dimensional Optimization and Control Theory*, vol. 54. Cambridge, U.K.: Cambridge Univ., 1999.
- [20] A. F. Filippov, *Differential Equations With Discontinuous Righthand Sides: Control Systems*, vol. 18. Springer Science & Business Media, 2013.