

# Non-Parametric and Regularized Dynamical Wasserstein Barycenters for Sequential Observations

Kevin C. Cheng\*, *IEEE Student Member*, Eric L. Miller\* *IEEE Fellow*,  
Michael C. Hughes†, Shuchin Aeron\* *IEEE Senior Member*

## Abstract

We consider probabilistic models for sequential observations which exhibit gradual transitions among a finite number of states. We are particularly motivated by applications such as human activity analysis where observed accelerometer time series contains segments representing distinct activities, which we call *pure states*, as well as periods characterized by continuous transition among these pure states. To capture this transitory behavior, the dynamical Wasserstein barycenter (DWB) model of [1] associates with each pure state a data-generating distribution and models the continuous transitions among these states as a Wasserstein barycenter of these distributions with dynamically evolving weights. Focusing on the univariate case where Wasserstein distances and barycenters can be computed in closed form, we extend [1] specifically relaxing the parameterization of the pure states as Gaussian distributions. We highlight issues related to the uniqueness in identifying the model parameters as well as uncertainties induced when estimating a dynamically evolving distribution from a limited number of samples. To ameliorate non-uniqueness, we introduce regularization that imposes temporal smoothness on the dynamics of the barycentric weights. A quantile-based approximation of the pure state distributions yields a finite dimensional estimation problem which we numerically solve using cyclic descent alternating between updates to the pure-state quantile functions and the barycentric weights. We demonstrate the utility of the proposed algorithm in segmenting both simulated and real world human activity time series.

## Index Terms

Wasserstein barycenter, displacement interpolation, dynamical model, sequential data, time series analysis, sliding window, non-parametric, quantile function, human activity analysis.

## I. INTRODUCTION

We consider a probabilistic model for sequentially observed data where the observation at each point in time depends on a dynamically evolving latent state. We are particularly motivated by systems that continuously move among a set of canonical behaviors, which we call *pure states*. Over some periods, the system may reside entirely in one of the pure states while over other periods, the system is transitioning among these pure states in a temporally smooth manner. There are many applications where such a model is appropriate including climate modeling [2], sleep analysis [3], simulating physical systems [4], as well as characterizing human activity from video [5] or wearable-derived accelerometry [6] data. Using the last case as an example, there will be periods when the individual will be engaged in a well-defined activity such as standing or running. During these intervals, the data can be modeled as drawn from a probability distribution specific to that canonical state. Given the high sampling rates of modern sensors, there also may be intervals where multiple consecutive observations reflect the gradual transition between or among pure states. Over these periods the distribution of the data is given by a suitable combination of the pure state distributions. Therefore, one possible model for these types of systems consists of three components: a set of distributions containing the data-generating distribution for each pure state, a continuously evolving latent state which captures the transition dynamics of the system as it moves among these pure states, and a means of interpolating among these pure state distributions to characterize the data distribution in the transition regions.

These types of systems pose some unique considerations that are not sufficiently addressed by prior work in time series modeling. The two most common methods for modeling latent state systems are continuous and discrete state-space models. Continuous state-space models [7], [8], [9] have no natural way to identify those pure states in which the system may persist for periods of time. In discrete state-space models such as hidden Markov models, [10], [11], [12], the dynamics are captured by a temporally varying state vector whose elements represent the probability that the system resides in each of a countable number of discrete (or in our terminology, pure) states. For these models, the data-generating distribution associated with this

\* Tufts University, Dept. of Electrical and Computer Engineering

† Tufts University, Dept. of Computer Science

This research was sponsored by the U.S. Army DEVCOM Soldier Center under the Measuring and Advancing Soldier Tactical Readiness and Effectiveness program and Cooperative Agreement Number W911QY-19-2-0003. We also acknowledge support from the U.S. National Science Foundation under award HDR-1934553 for the Tufts T-TRIPODS Institute. Shuchin Aeron is supported in part by NSF CCF:1553075, NSF RAISE 1931978, NSF ERC planning 1937057, and AFOSR FA9550-18-1-0465. Michael C. Hughes is supported in part by NSF IIS-1908617. Eric L. Miller is supported in part by NSF grants 1934553, 1935555, 1931978, and 1937057.

Code repository: [https://github.com/kevin-c-cheng/DWB\\_Nonparametric](https://github.com/kevin-c-cheng/DWB_Nonparametric)

©2023 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works. DOI: 10.1109/TSP.2023.3303616

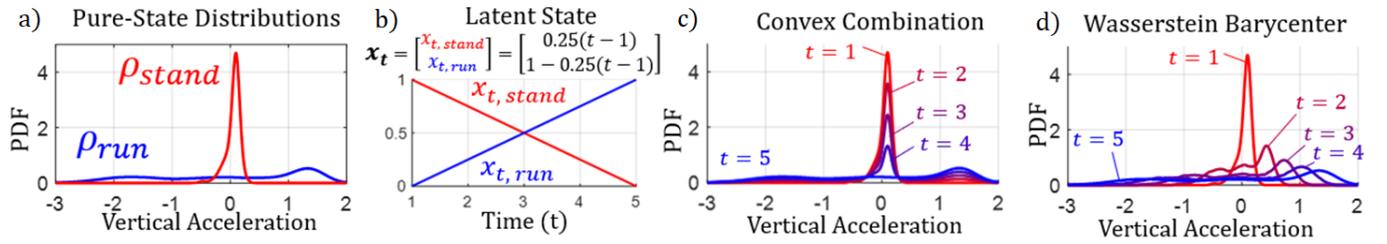


Fig. 1: **Comparison of convex combination vs Wasserstein barycenter for modeling human activity transitions.** The beep test (BT) dataset consists of a subject running back-and-forth between two points, stopping at each point (see Sec. V-B for more details). (a) The probability distribution functions (PDF) of the vertical acceleration of the system’s two pure states (stand, run) are estimated using a KDE with a Gaussian kernel whose mean corresponds to the observed data when the system resides in these pure states. Modeling a transition from stand to run via the time-varying weights (b) for  $t = 1, \dots, 5$ , we show the resulting data distributions during this transition region according to a convex combination (c) and Wasserstein barycenter (d) interpolation model.

latent state vector is given as a convex combination, i.e. a linear mixture, of the pure state distributions. As argued in [1] this is also an insufficient model for the problems which interest us. As an example, for modeling human activity, the data distribution illustrated in Fig. 1 produced by a convex combination of the underlying “standing” and “running” pure state distributions can be interpreted as “sometimes standing” and “sometimes running,” which is not a proper description of the gradual transition that actually occurs.

A better model for interpolating the data distribution during the transition period between standing and running would smoothly shift probability mass between the two pure state distributions. As illustrated in Fig. 1, such blending can be achieved through displacement interpolation [13] between two pure states or, for more than two pure states, as a *Wasserstein barycenter* of the associated distributions [14]. Using this perspective, Dynamical Wasserstein Barycenters (DWB) [1] were recently proposed to model a dynamically evolving distribution as a sequence of Wasserstein barycenters constructed as a time-varying convex combination of the pure state distributions. The dynamical weights, which lie on the probability simplex, are taken to be the latent state of the model. A Bayesian model is proposed in [1], whose parameters were determined via maximum *a posteriori* estimation.

Here we expand on [1] by highlighting two challenging characteristics of the DWB model and two improvements that address certain limitations of the original DWB approach. The first characteristic relates to uniqueness in DWB model identification, where multiple combinations of pure state distributions and barycentric weights can produce the same Wasserstein barycenter. Although this is true for multidimensional distributions, here we use a univariate formulation to more transparently demonstrate how this non-uniqueness is captured in an *inverse-scaling* relationship between the model’s latent state and pure state parameters. The second characteristic is related to a tradeoff in tracking and estimating an evolving data distribution from a single instance of a time series using a windowed approach to collect samples. Smaller windows lack the number of samples to ensure small statistical error in the estimation of the data distribution at a given point in time. On the other hand, larger windows span longer periods during which, under relatively faster dynamics, the data distribution can change significantly again increasing the estimation error. In a simulated example, where the dynamics consist of constant rate transitions between two Gaussian states, we show that there exists an optimal window size that balances these two effects and discuss the dependency of this window size tradeoff on the temporal dynamics of the latent state and pure states of the system.

Our first improvement addresses the limitations of the choice in [1] in using a probabilistic prior for the dynamics of the DWB weight vector. That approach may introduce additional unnecessary or potentially undesirable probabilistic properties on the latent state process such as a limiting distribution [15] which fails to adequately regularize the DWB estimation problem. Instead, we propose here a regularization scheme that imposes temporal smoothness by penalizing the difference between the simplex-constrained, latent state vectors at adjacent points in time. Drawing from the field of compositional analysis [16], the Bhattacharya-arccos distance [17] proves to be well-suited to our needs. As a consequence of the aforementioned inverse-scaling relationship, introducing this latent state regularizer impacts the model’s pure state distribution in a manner that causes them to diverge from the data. Therefore, we also introduce a regularizer to counteract this effect to ensure that the learned pure state distributions are representative of data while the system resides in each pure state.

Our second improvement removes the restriction in [1] where a parametric approach to model pure states with multivariate Gaussians was employed. Here we adopt a non-parametric approach and focus on the univariate case where the Wasserstein-2 distance between distributions is equivalent to the 2-norm between their respective quantile functions [18]. Using a discrete approximation to the pure state quantile functions leads to a convenient finite dimensional, regularized linear least squares problem for estimating the pure states.

Our numerical experiments empirically validate our analysis and improvements to the DWB model. Using simulated data, we demonstrate in a controlled setting how we effectively regularize our model parameters with proper consideration of the

inverse-scaling analysis and the impact of window size on the accuracy of the model parameters. Additionally, using real world human activity data, we show how our non-parametric approach leads to improved estimation of the system's pure state distributions as well as improved fit of the time-evolving distribution of the observed data compared to [1].

In summary, the primary contributions of this work consist of the following:

- 1) We highlight the non-uniqueness of the parameters corresponding to a Wasserstein barycenter by detailing the inverse-scaling relationship between the pure state distributions and the simplex-valued barycentric weights.
- 2) We explore the impact of the window size on the ability to accurately estimate a dynamically evolving data distribution by exploring the tradeoff between the errors associated with large and small windows and the dependency of this tradeoff on the dynamics and pure states of the system.
- 3) We propose regularizers for the model parameters that impose temporal smoothness in the latent states in a manner that addresses the non-uniqueness of the model.
- 4) We propose a flexible, non-parametric representation for univariate pure state distributions using a discrete approximation to the quantile function that results in a finite dimensional formulation for DWB learning.

The remainder of the paper is organized as follows: in Sec. II, we provide an overview of the Wasserstein distance and barycenter focusing on the univariate case. In Sec. III, we discuss the DWB model, highlighting the non-uniqueness and inverse-scaling property of the Wasserstein barycenter as well as the impact of the window size on the estimation of a dynamically evolving data distribution. In Sec. IV, we develop a variational problem for learning a DWB model, followed by a discussion of the regularization approach, and discretization of the pure state distributions required to obtain a finite dimensional estimation problem. We then formally state our non-parametric and regularized DWB variational problem and provide an algorithm to estimate the model parameters. In Sec. V we use simulated data to demonstrate the non-uniqueness, impact of window size and regularization terms discussed in this work and use real world human activity data to demonstrate the advantages of the non-parametric DWB approach relative to the Gaussian model.

## II. TECHNICAL BACKGROUND

The Wasserstein-2 distance is a metric on the space of probability distributions on  $\mathbb{R}^d$  with finite second moments [18], [19]. For two random variables  $q$  and  $s$  distributions  $\rho_q$  and  $\rho_s$ , the squared Wasserstein-2 distance is defined via,

$$\mathcal{W}_2^2(\rho_q, \rho_s) = \inf_{\pi \in \Pi(\rho_q, \rho_s)} \mathbb{E}_{q, s \sim \pi} \|q - s\|_2^2 \quad (1)$$

where  $\pi$  denotes the joint distribution of  $q$  and  $s$ , and  $\Pi(\rho_q, \rho_s)$  is the set of all joint distributions with marginals  $\rho_q, \rho_s$ . In this work, we refer to Eq. (1) as the squared Wasserstein distance.

Given a set of distributions  $\rho_{q_{1:K}} = \{\rho_{q_1}, \rho_{q_2}, \dots, \rho_{q_K}\}$  and a vector  $\mathbf{x} \in \Delta^K$ , where  $\Delta^K$  denotes the standard  $K$ -simplex, the *Wasserstein barycenter* is the distribution that minimizes the weighted (with respect to elements in  $\mathbf{x}$ ) squared Wasserstein distance to the set of distributions [14] and is given by,

$$\rho_B = B(\mathbf{x}, \rho_{q_{1:K}}) = \operatorname{argmin}_{\rho} \sum_{k=1}^K \mathbf{x}[k] \mathcal{W}_2^2(\rho, \rho_{q_k}), \quad (2)$$

where  $\mathbf{x}[k]$  denotes the  $k$ -th element of the vector  $\mathbf{x}$ . When  $\rho_q$  and  $\rho_s$  are univariate distributions with cumulative distribution functions  $P_q, P_s$ , the squared Wasserstein distance in Eq. (1) becomes [19], [20],

$$\mathcal{W}_2^2(\rho_q, \rho_s) = \int_0^1 (P_q^{-1}(\xi) - P_s^{-1}(\xi))^2 d\xi. \quad (3)$$

Here  $P_q^{-1}$  and  $P_s^{-1}$  are quantile functions, the generalized inverse [21] of the cumulative distribution function, given by,

$$P^{-1}(\xi) = \inf\{g \in \mathbb{R} : P(g) \geq \xi\}. \quad (4)$$

It follows from Eq. (3) and Eq. (2) that the Wasserstein barycenter of a set of univariate distributions with quantile functions  $P_{q_{1:K}}^{-1}$ , will have quantile function [20],

$$P_B^{-1} = \sum_{k=1}^K \mathbf{x}[k] P_{q_k}^{-1}. \quad (5)$$

## III. THE DYNAMICAL WASSERSTEIN BARYCENTER MODEL

Shown in Fig. 2, the DWB model [1] describes the distribution of a time series  $y_t$  at time  $t$  as,

$$y_t \sim \rho_{B_t} = B(\mathbf{x}_t, \rho_{q_{1:K}}) \quad (6)$$

where  $\rho_{q_k}$ ,  $k = 1, 2, \dots, K$  are the distributions of the pure states and the barycentric weight  $\mathbf{x}_t \in \Delta^K$  capture the dynamics of the transitions among these pure states.

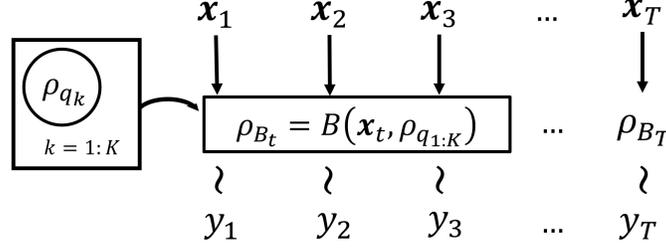


Fig. 2: **DWB model diagram.** The DWB models the distribution  $\rho_{B_t}$  from which the time series  $y_t$  is sampled as the Wasserstein barycenter of a set of pure state distributions  $\rho_{q_{1:K}}$  and barycentric weight  $\mathbf{x}_{1:T}$ , the latent state of the model.

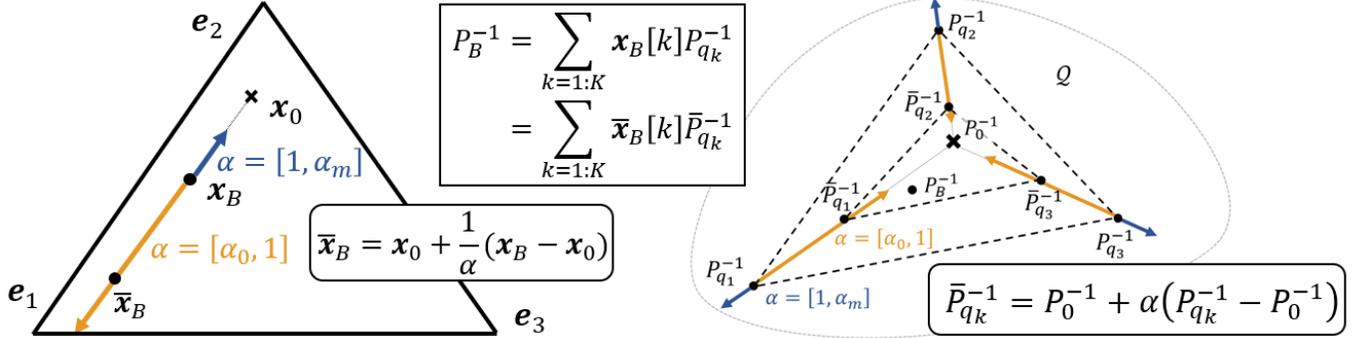


Fig. 3: **Diagram of the non-uniqueness and inverse-scaling effect of the parameters of a Wasserstein barycenter.** Consider a set of three pure states with quantile functions  $P_{q_{1:3}}^{-1}$  and simplex-valued weight  $\mathbf{x}_B \in \Delta^3$  where  $\rho_B = B(\mathbf{x}_B, \rho_{q_{1:3}})$  with quantile function  $P_B^{-1} = \sum_{k=1}^3 \mathbf{x}_B[k]P_{q_k}^{-1}$ . We construct a family of distinctly different pure state quantile functions  $\bar{P}_{q_{1:K}}^{-1}$  and barycentric weights  $\bar{\mathbf{x}}_B$  which produce the exact same barycenter  $P_B^{-1} = \sum_{k=1}^3 \bar{\mathbf{x}}_B[k]\bar{P}_{q_k}^{-1}$ . Let  $\mathbf{x}_0$  be another point on the simplex where  $\rho_0 = B(\mathbf{x}_0, \rho_{q_{1:3}})$  has quantile function  $P_0^{-1} = \sum_{k=1}^3 \mathbf{x}_0[k]P_{q_k}^{-1}$ . Given  $\mathbf{x}_0$  and  $\mathbf{x}_B$ , let  $\bar{\mathbf{x}}_B$  given by Eq. (7) be any point on the line connecting  $\mathbf{x}_0$  through  $\mathbf{x}_B$  to the edge of the simplex. Moving  $\bar{\mathbf{x}}_B$  away from  $\mathbf{x}_0$  along the line connecting  $\mathbf{x}_0$  and  $\mathbf{x}_B$  (orange segments), causes the pure state quantile functions  $\bar{P}_{q_k}^{-1}$  to move from  $P_{q_{1:3}}^{-1}$  towards  $P_0^{-1}$ . This corresponds to  $\alpha \in [\alpha_0, 1]$  where  $\alpha_0$  is the smallest value of  $\alpha$  such that  $\bar{\mathbf{x}}_B$  still lies on the simplex. Conversely moving  $\mathbf{x}_B$  towards  $\mathbf{x}_0$  (blue segments) results in the pure state quantile functions moving away from  $P_0^{-1}$ . This corresponds to  $\alpha \in [1, \alpha_m]$ , where  $\alpha_m$  is the largest value of  $\alpha$  such that all  $\bar{P}_{q_{1:3}}$  remain in the set of quantile functions  $\mathcal{Q}$ .

Given  $y_t, t = 1, 2, \dots, T$  modeled via equation (6), the problem is to estimate DWB model parameters which consist of the pure state distributions and the sequence of barycentric weights.

Below we discuss two key characteristics that pose challenges for estimating the parameters of the DWB model. The first is the non-uniqueness of the parameters (i.e., the pure state distributions and the barycentric weights) that yield a Wasserstein barycenter. The second relates to the complications that arise when we are provided only a single time series for learning a DWB model.

#### A. Non-uniqueness in the Parameters of a Wasserstein Barycenter

The issue of uniqueness refers to the fact that a Wasserstein barycenter is not described by a unique set of pure state distributions and barycentric weights. While the statement is true regardless of dimension (see Appendix A for an example), given the focus of this paper, we examine the univariate case in some detail. Specifically, we provide a construction that illustrates an inverse-scaling relation between the family of pure state distributions and barycentric weights that yields the same Wasserstein barycenter.

As shown in Fig. 3, assume we have a set of pure state distributions  $\rho_{q_{1:K}}$  indexed by  $k = 1, 2, \dots, K$  with quantile functions,  $P_{q_k}^{-1}$  and barycentric weights  $\mathbf{x}_B \in \Delta^K$  which give rise to the barycenter  $\rho_B = B(\mathbf{x}_B, \rho_{q_{1:K}})$  with quantile function  $P_B^{-1} = \sum_{k=1}^K \mathbf{x}_B[k]P_{q_k}^{-1}$  (Eq. (5)). For now,  $\mathbf{x}_B$  is assumed to lie in the interior of the simplex and we consider below the cases where  $\mathbf{x}_B$  is on a lower dimensional face or vertex. Let us choose another point  $\mathbf{x}_0 \neq \mathbf{x}_B$  corresponding to barycentric quantile function  $P_0^{-1} = \sum_{k=1}^K \mathbf{x}_0[k]P_{q_k}^{-1}$  (Eq. (5)). We construct a family of barycentric weights  $\bar{\mathbf{x}}_B$  and pure state quantile functions,  $\bar{P}_{q_{1:K}}^{-1}$  corresponding to distributions  $\bar{\rho}_{q_{1:K}}$  such that  $\rho_B = B(\bar{\mathbf{x}}_B, \bar{\rho}_{q_{1:K}})$ , or in other words,  $P_B^{-1} = \sum_{k=1}^K \bar{\mathbf{x}}_B[k]\bar{P}_{q_k}^{-1}$ .

Specifically, as seen in Fig. 3 we define  $\bar{x}_B$  to be a point on the line segment connecting  $x_0$  to the boundary of the simplex that passes through  $x_B$ . That is,

$$\bar{x}_B = x_0 + \frac{1}{\alpha}(x_B - x_0). \quad (7)$$

The parameter  $\alpha$  captures the scaling nature of this construction. When  $\alpha = \infty$  we have  $\bar{x}_B = x_0$ . As  $\alpha$  decreases,  $\bar{x}_B$  moves away from  $x_0$  along the blue segment connecting  $x_0$  and  $x_B$ , ultimately crossing  $x_B$  when  $\alpha = 1$ . Further reducing  $\alpha$  towards 0 moves  $\bar{x}_B$  along the orange component of the same line until some point  $\alpha \in (0, 1]$ , where  $\bar{x}_B$  reaches the boundary. We denote this point as  $\alpha_0$ .

With this definition of  $\bar{x}_B$  in Eq. (7) we construct quantile functions,  $\bar{P}_{q_1:K}^{-1}$  such that  $P_B^{-1} = \sum_{k=1}^K \bar{x}_B[k] \bar{P}_{q_k}^{-1}$ . Indeed,

$$\begin{aligned} P_B^{-1} &= \sum_{k=1}^K x_B[k] P_{q_k}^{-1} = \sum_{k=1}^K (\alpha \bar{x}_B[k] + (1 - \alpha) x_0[k]) P_{q_k}^{-1} \\ &= \alpha \sum_{k=1}^K \bar{x}_B[k] P_{q_k}^{-1} + (1 - \alpha) \sum_{k=1}^K x_0[k] P_{q_k}^{-1} \\ &= \alpha \sum_{k=1}^K \bar{x}_B[k] P_{q_k}^{-1} + (1 - \alpha) P_0^{-1} \quad \boxed{P_0^{-1} = \sum_{k=1}^K x[k] P_{q_k}^{-1}} \\ &= \sum_{k=1}^K (\alpha \bar{x}_B[k] P_{q_k}^{-1} + (1 - \alpha) \bar{x}_B[k] P_0^{-1}) \quad \boxed{\sum_{k=1}^K \bar{x}_B[k] = 1} \\ &= \sum_{k=1}^K \bar{x}_B[k] (\alpha P_{q_k}^{-1} + (1 - \alpha) P_0^{-1}) = \sum_{k=1}^K \bar{x}_B[k] \bar{P}_{q_k}^{-1} \end{aligned}$$

with for each  $k$ ,

$$\bar{P}_{q_k}^{-1} = P_0^{-1} + \alpha(P_{q_k}^{-1} - P_0^{-1}). \quad (8)$$

Eq. (8) bears a strong resemblance to Eq. (7), except now with reciprocal use of  $\alpha$ . For  $\alpha \in [\alpha_0, 1]$ ,  $\bar{P}_{q_k}^{-1}$  is a convex combination of  $P_0^{-1}$  and  $P_{q_k}^{-1}$  lying on the line segment connecting the two quantile functions. In this case, since the collection of monotone functions on  $[0, 1]$  is a convex set [22],  $\bar{P}_{q_1:K}^{-1}$  will be valid quantile functions. However, for  $\alpha > 1$ ,  $\bar{P}_{q_k}^{-1}$  extends beyond  $P_B^{-1}$  along the line that connects  $P_0^{-1}$  to  $P_B^{-1}$ . In this case,  $\bar{P}_{q_k}^{-1}$  is no longer a convex combination of  $P_0^{-1}$  and  $P_B^{-1}$  and is not guaranteed to be a quantile function. We denote  $\alpha_m$  as the maximum value of  $\alpha$  such that  $\bar{P}_{q_k}^{-1}$  is a valid quantile function.

Thus, the sets of  $\bar{x}_B$  and  $\bar{P}_{q_1:K}^{-1}$  corresponding to  $\alpha \in [\alpha_0, \alpha_m]$  according to Eqs. (7) and (8) describe the family of parameters that yield the same Wasserstein barycenter as  $x_B$  and  $P_{q_1:K}^{-1}$ . The reciprocal appearance of  $\alpha$  in these equations captures the inverse-scaling relationship for this family of parameters. As shown in Fig. 3, the case where  $\alpha \in [\alpha_0, 1]$  corresponds to the orange lines where *increasing* the distance of  $\bar{x}_B$  from  $x_0$  along the segment connecting  $x_B$  and  $x_0$ , results in the pure states *decreasing* their distances to  $P_0^{-1}$  each along linear trajectories connecting the  $P_{q_k}^{-1}$  to  $P_0^{-1}$ . Conversely, the case where  $\alpha \in [1, \alpha_m]$  corresponds to the blue lines in Fig. 3 where *decreasing* the distance between  $\bar{x}_B$  and  $x_0$  by moving along the line that connects  $x_B$  to  $x_0$ , results in pure state quantile functions  $\bar{P}_{q_k}^{-1}$  that are now *increasing* their distance from  $P_0^{-1}$  by extending linearly along the ray from  $P_0^{-1}$  to  $P_{q_k}^{-1}$ .

Should  $x_B$  lie on a face of the  $K$ -simplex of dimension greater than one but less than  $K$  (i.e., *not* a vertex), this construction may be repeated by placing  $x_0 \neq x_B$  in that same lower dimensional simplex. In such a case, we may also place  $x_0$  in the interior of the  $K$  dimensional simplex. More specifically, with  $x_B$  on a face and  $x_0$  located in the interior, it is clear that  $\alpha_0 = 1$  implying that the construction above holds only if  $\alpha_m > 1$ . This is also the case in the event that  $x_B$  is a vertex for any value of  $x_0 \neq x_B$ . With the constraint that  $\bar{P}_{q_1:K}^{-1}$  must remain valid quantile functions, it is possible to construct an example such that  $\alpha_m = 1$ , which when combined with the aforementioned case where  $\alpha_0 = 1$  means that we cannot employ this construction to show non-uniqueness. However, as discussed in Appendix A neither can we conclude that the barycenter is in fact unique.

This non-uniqueness and inverse-scaling relation implies that for models (such as the DWB) that require learning both the pure state quantile functions and the barycentric weights corresponding to one or a sequence of Wasserstein barycenters, introducing constraints on one set of parameters will have an impact on the other. We take this effect into account in Sec. IV-B when constructing regularizers to impose desirable properties on the latent state sequence and pure state distributions.

## B. Model Sampling

A second characteristic of learning a DWB model relates to how we incorporate observed data into the model. Ideally, we would directly observe the data distribution  $\rho_{B_t}$  as specified by Eq. (6) at each point in time or generate an estimate of these

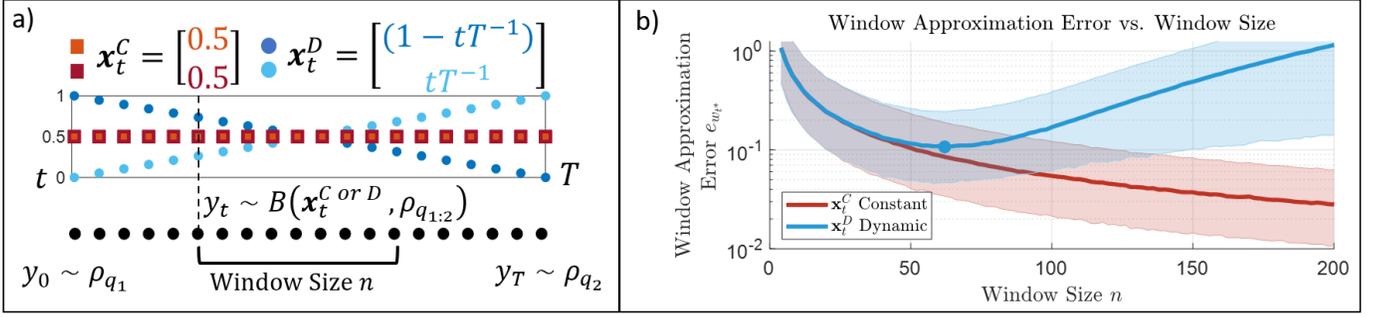


Fig. 4: **Approximating a distribution with a window of samples.** (a) Two simulated configurations, one where  $\mathbf{x}_t^C = [0.5, 0.5]^T$  is constant, and one where  $\mathbf{x}_t^D = [(1 - tT^{-1}), tT^{-1}]^T$  is dynamically evolving. We generate a time series  $y_t \sim B(\mathbf{x}_t^C \text{ or } D, \rho_{q_{1:2}})$  where  $\rho_{q_1} = \mathcal{N}(0, 5)$ , and  $\rho_{q_2} = \mathcal{N}(10, 0.2)$ <sup>1</sup>. (b) We estimate the distribution data distribution  $\rho_{B_{t^*}}$  where  $t^* = 0.5(T + 1)$ , with a distribution  $\rho_{y_{t^*}}$  comprised of a window of  $n$  samples centered at  $t^*$  according to Eq. (9). Varying the window size  $n$  we plot the average value and the 25/75-th quantile bands of  $e_{w_{t^*}} = \mathcal{W}_2^2(\rho_{y_{t^*}}, \rho_{B_{t^*}})$  from  $10^4$  simulations. In the constant state case where  $\rho_{y_{t^*}}$  consists of  $n$  IID samples from  $\rho_{B_{t^*}}$  the average error monotonically decreases. However, in the dynamic case, where the samples in the window are independent but not identically distributed, the U-shape curve highlights the window size tradeoff where  $n_0$  indicates the optimal window size.

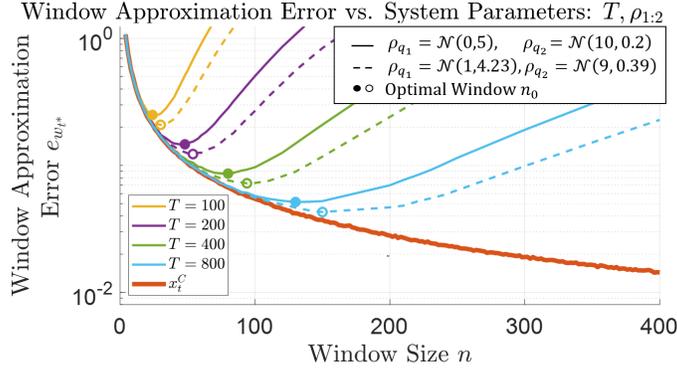


Fig. 5: **Impact of system dynamics and pure state distributions on WAE.** The average WAE is plotted for various configurations where solid lines refer to  $\rho_{q_1} = \mathcal{N}(0, 5)$ ,  $\rho_{q_2} = \mathcal{N}(10, 0.2)$  and dashed lines corresponds to a system where the Wasserstein distance between  $\rho_{q_{1:2}}$  is decreased,  $\rho_{q_1} = \mathcal{N}(1, 4.23)$   $\rho_{q_2} = \mathcal{N}(9, 0.39)$ . In both cases,  $\rho_{B_{t^*}} = \mathcal{N}(5, 1.8)$ . Decreasing the per-sample change in distribution by increasing  $T$  or decreasing  $\mathcal{W}_2^2(\rho_{q_1}, \rho_{q_2})$  (solid  $\rightarrow$  dashed) results in smaller  $e_{w_{t^*}}$ . The impact of these changes on  $e_{w_{t^*}}$  is greater for larger windows, which increases the optimal window  $n_0$  corresponding to the minimum of the U-curves.

distributions from multiple realizations of a time series. Unfortunately, in all the practical cases of interest to us, only a single instance of the time series is available for processing. Thus, we consider the problem of estimating the time-varying data distribution  $\rho_{B_t}$  from a single time series that is sampled from  $\rho_{B_t}$ . To do this, we consider a window of  $n$  samples centered at  $t$ , compiled into a vector  $\mathbf{y}_t = [y_{(t-\frac{n}{2})}, \dots, y_{(t+\frac{n}{2}-1)}]^T$ . For convenience of notation, we assume  $n$  to be even. We estimate the data distribution with the distribution  $\rho_{y_t}$  based on this sample window,

$$\rho_{y_t} = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{y}_t[i]}. \quad (9)$$

Here,  $\delta_{\mathbf{y}_t[i]}$  is the Dirac-delta measure located at  $\mathbf{y}_t[i]$ . We then define the *window approximation error* (WAE) as,

$$e_{w_t} = \mathcal{W}_2^2(\rho_{y_t}, \rho_{B_t}). \quad (10)$$

Since the samples that constitute  $\rho_{y_t}$  are random, the WAE is a random quantity.

Eq. (9) is an empirical measure when the samples in the window are drawn IID which is only possible when the barycentric weights  $\mathbf{x}_t$ , and thus  $\rho_{B_t}$ , is constant in time. In this case, as  $n \rightarrow \infty$ , it is known that the distribution function of  $\rho_{y_t}$  converges almost surely to the distribution function of  $\rho_{B_t}$ , and the expected value of  $e_{w_t}$  converges to zero [23]. The complication in

<sup>1</sup> $\mathcal{N}(\mu, \sigma^2)$ : Gaussian distribution with mean  $\mu$ , and variance  $\sigma^2$ .

our case comes from the fact that  $\mathbf{x}_t$  is changing with time, meaning that the samples in the window are no longer identically distributed as the data distribution changes from sample to sample.

This impact of a dynamically evolving distribution on the WAE is dependent on many factors including the size of the window  $n$  as well as properties of the system that impact the manner in which the data distribution changes between samples. To simplify matters, we consider a simple yet informative example in Fig. 4a of a system with two pure states and two latent state configurations, one where the latent state is *constant* with  $\mathbf{x}_t^C = [0.5, 0.5]^T$ , and one where the latent state is *dynamic*, changing at a constant rate where  $\mathbf{x}_t^D = [(1 - tT^{-1}), tT^{-1}]^T$ , for  $t = 0, 1, \dots, T$ . For two distributions  $\rho_{q_{1,2}}$ , we generate a time series by sampling independently  $y_t^{C \text{ or } D} \sim B(\mathbf{x}_t^{C \text{ or } D}, \rho_{q_{1,2}})$ . Thus, in the constant state case,  $y_t^C$  are IID. Let  $\rho_{y_{t^*}}$  be the empirically estimated distribution constructed according to Eq. (9) at  $t^* = 0.5(T + 1)$ , the half-way point of the transition. In both cases of  $\mathbf{x}_t^C$  and  $\mathbf{x}_t^D$  the distribution at this time is  $\rho_{B_{t^*}} = B([0.5, 0.5]^T, \rho_{q_{1,2}})$ .

When the window size is small, the distributional change over the sample window is also small. Therefore as seen in the left side of Fig. 4b, the dynamic case closely approximates the constant state where the average WAE decreases with respect to the window size [23]. When the window size is large, the samples being included in the window are further from  $t^*$  and thus the distributions of these samples increasingly diverge from  $\rho_{B_{t^*}}$ . Shown by the rising and expanded quantile bands on the right side of the plot, this effect causes the average and variability of  $e_{w_{t^*}}$  to increase for large windows. The resulting U-shape curve of the WAE highlights the window size tradeoff where the optimal window size  $n_0$  corresponding to the minimum of this curve balances the benefits of having more samples for robust estimation of a distribution with the effects of using samples farther from the point of interest.

One major factor that impacts this tradeoff is the magnitude of the distributional change from sample to sample. There are a number of factors that can decrease (resp. increase) the magnitude of this per-sample change in distribution including (1) decreasing (increasing) the rate of change of the system's continuous time dynamics; (2) increasing (decreasing) the sampling rate of the sensor which provides discrete measurements; or (3) decreasing (increasing) the Wasserstein distance between the pure state distributions which affects the total amount of distributional change during the transition among pure states. Continuing with the example, we demonstrate the impact of these factors on the window size tradeoff and optimal window size. Again, modeling the system as moving from  $\rho_{q_1}$  to  $\rho_{q_2}$  with dynamics specified by  $\mathbf{x}_t^D$ , factors (1) and (2) dictate the number of samples over which this transition occurs<sup>2</sup>, thus their combined impact can be understood by varying  $T$ . To understand the impact of (3), we simulate two different pure state configurations varying  $\mathcal{W}_2^2(\rho_{q_1}, \rho_{q_2})$ .

The results in Fig. 5 confirm that for a given window size, decreasing the per-sample change in distribution by increasing  $T$ , or by decreasing  $\mathcal{W}_2^2(\rho_{q_1}, \rho_{q_2})$  results in a decrease in the average WAE. Additionally, the increasing difference between the U-curves in Fig. 5 as we move towards larger windows confirms that the decreasing the per-sample change in distribution has an increasing benefit for larger windows where the dynamics of the system have a larger impact on the accuracy of the window estimate. This shifts the balance of the window size tradeoff as seen by the minimums of these U-curves moving to the right, implying that decreasing the per-sample change in distribution using any of the three methods mentioned increases the optimal window size. Indeed, in the limiting case where either  $T \rightarrow \infty$  or  $\mathcal{W}_2^2(\rho_{q_1}, \rho_{q_2}) \rightarrow 0$  in which case  $\rho_{q_1} = \rho_{q_2} = \rho_{B_{t^*}}$ , the samples will be drawn IID from a constant distribution. In this limit, the dynamic case converges to the constant case where the optimal window size  $n_0 \rightarrow \infty$  and  $e_{w_t} \rightarrow 0$ .

#### IV. NON-PARAMETRIC AND REGULARIZED DYNAMICAL WASSERSTEIN BARYCENTERS

In this section, we detail our proposed variational problem for estimating the parameters of the univariate DWB model. We discuss our regularization that ensures that the latent state evolves smoothly over time while taking into account its effect on the pure states through the inverse-scaling relationship. We also discuss our non-parametric representation for the pure state distributions using a discrete approximation to the pure state quantile function. Finally, we detail how this leads to a least-squares formulation of the variational DWB objective and propose an algorithm for learning the parameters of the model.

##### A. Variational Problem for DWB Model Estimation

Training a DWB model entails estimating the pure states distributions and latent states sequence to minimize a cost function that encourages both fidelity to the data as well as model parameters that conform to prior knowledge we may have concerning the general behavior of time series.

Building on the approach discussed in Sec. III-B, we create a sequence of  $N$  sample windows of length  $n$  to estimate the distribution of the time series at select points. For our simulations in Sec. V we use overlapping windows separated by a fixed stride length. Given  $t_i$  for  $i = 1, \dots, N$  as the starting index for these sample windows, let  $\mathbf{y}_i = [y_{t_i}, \dots, y_{t_i+n}]^T$  be the vector of samples and  $\rho_{y_i}$  the distribution according to Eq. (9) corresponding to this window of samples. Using the WAE in Eq. (10)

<sup>2</sup>The rate of change (1) may have units change in distribution per second, and the sampling rate (2) has units samples per second. Thus  $\frac{(1)}{(2)}$  will have units change in distribution per sample.

summed over  $i$  as a data fidelity term and encoding prior information in regularizers the details of which are discussed below, the variational problem we seek to solve is,

$$\hat{\rho}_{q_{1:K}}, \hat{\mathbf{x}}_{1:N} = \underset{\rho_{q_{1:K}}, \mathbf{x}_{1:N}}{\operatorname{argmin}} \underbrace{\sum_{i=1}^N \mathcal{W}_2^2(\rho_{y_i}, \rho_{B_i})}_{\text{Data Fit}} + \lambda_x \underbrace{R_x(\mathbf{x}_{1:N})}_{\text{x-regularization}} + \lambda_q N \underbrace{R_q(\rho_{q_{1:K}})}_{\text{q-regularization}}, \quad (11)$$

where  $\rho_{B_i} = B(\mathbf{x}_i, \rho_{q_{1:K}})$ . Here,  $\lambda_x, \lambda_q \geq 0$  are regularization weights. We multiply  $R_q$  by  $N$  so that it scales with the length of the time series along with the other terms in the cost function.

As highlighted by the inverse-scaling relation in Sec. III-A, any regularizer of either the latent state or pure state must consider its effect on the other parameter. Motivated by applications where the system evolves gradually over time, we propose a regularization scheme that imposes a gradually evolving latent state while ensuring that the learned pure state distributions accurately reflect the distribution of the data corresponding to when the system resides in a pure state.

### B. Parameter Regularization for DWB

As seen in our windowing simulations in Sec. III-B, even in a simple case, estimating a dynamically evolving distribution with a window of samples is a challenging problem. The ambiguities identified in our experiments in Sec. III-B may cause the learned latent state to vary greatly even if the system is constant or gradually evolving. Therefore, to limit its variability, we propose a regularizer that penalizes the sum of the squared distances,  $d^2(\mathbf{x}_i, \mathbf{x}_{i+1})$ , between successive latent states for a suitable distance  $d$  on the simplex [16]. While several choices are possible, for the purpose of this work, we choose the Bhattacharyya-arccos distance [17], one that is bounded and differentiable. We discuss alternative distances in Appendix B. Thus, this regularizer penalizes the total length of the latent state trajectory on the simplex according to

$$\begin{aligned} R_x(\mathbf{x}_{1:N}) &= \sum_{i=1}^{N-1} d^2(\mathbf{x}_i, \mathbf{x}_{i+1}) \\ &= \sum_{i=1}^{N-1} \left( \arccos \left( \sum_{k=1}^K \sqrt{\mathbf{x}_i[k] \mathbf{x}_{i+1}[k]} \right) \right). \end{aligned} \quad (12)$$

With the above choice for regularizing the latent state sequence, our choice for regularization of the pure state distributions is motivated by Eq. (7) and the inverse-scaling nature of barycentric non-uniqueness. Considering our non-uniqueness construction in Eq. (7), if we set the reference point to be the current value of the latent state  $\mathbf{x}_0 = \mathbf{x}_t$ , and constructed point as the next state  $\bar{\mathbf{x}}_B = \mathbf{x}_{t+1}$ , Eq. (7) takes the form  $\mathbf{x}_{t+1} = \mathbf{x}_t + \frac{1}{\alpha} \mathbf{d}_t$  where  $\mathbf{d}_t$  is a vector in the simplex along which we move the latent state. Since  $\mathbf{d}_t$  is essentially arbitrary, taking  $\alpha > 1$  will encourage small changes in  $\mathbf{x}_t$  as required by Eq. (12). Referring to Fig. 3, we see that in this  $\alpha > 1$  regime, barycentric non-uniqueness manifests in the divergence of the quantile functions; i.e., motion along the blue line segments. That this in fact can occur is verified in our experiments in Sec. V-A1. Now, for time series that reside in each pure state at some period in time, the observed data during those periods will be representative of each of the pure state distributions. Therefore, having estimated pure state distributions diverge is undesirable if we want to accurately learn these quantities. To counteract this diverging behavior, we propose a regularizer in the space of quantile functions that penalizes the sum of squared Wasserstein distances of the pure state quantile functions  $P_{q_k}^{-1}$  from a reference quantile function  $P_0^{-1}$ . Here, we choose  $P_0^{-1}$  be the quantile function of  $\rho_0 = B(\mathbf{x}_0, \rho_{q_{1:K}})$  where  $\mathbf{x}_0 = \frac{1}{K} \mathbb{1}_K$  (with  $\mathbb{1}_K$  the length  $K$  vector of all ones):

$$\begin{aligned} R_q(\rho_{q_{1:K}}) &= \sum_{k=1:K} \mathcal{W}_2^2 \left( \rho_{q_k}, B \left( \frac{1}{K} \mathbb{1}_K, \rho_{q_{1:K}} \right) \right) \\ &= \sum_{k=1:K} \int_0^1 \left( P_{q_k}^{-1}(\xi) - \frac{1}{K} \sum_{j=1:K} P_{q_j}^{-1}(\xi) \right)^2 d\xi, \end{aligned} \quad (13)$$

The two regularizers just discussed are designed to work in tandem. Here,  $R_x$  ensures that the latent state evolves gradually over the simplex while  $R_q$  ensures that the pure state quantile functions do not diverge in the ways predicted by the inverse-scaling analysis in Sec. III-A. Through our simulations in Sec. V-A, we demonstrate how by appropriately balancing the regularization weights one can reliably estimate the DWM model parameters.

### C. Discrete Quantile Approximation

As estimation of the infinite dimensional quantile functions  $P_{q_{1:K}}^{-1}$  in (11) requires a finite dimensional approximation we introduce the following two quantities:

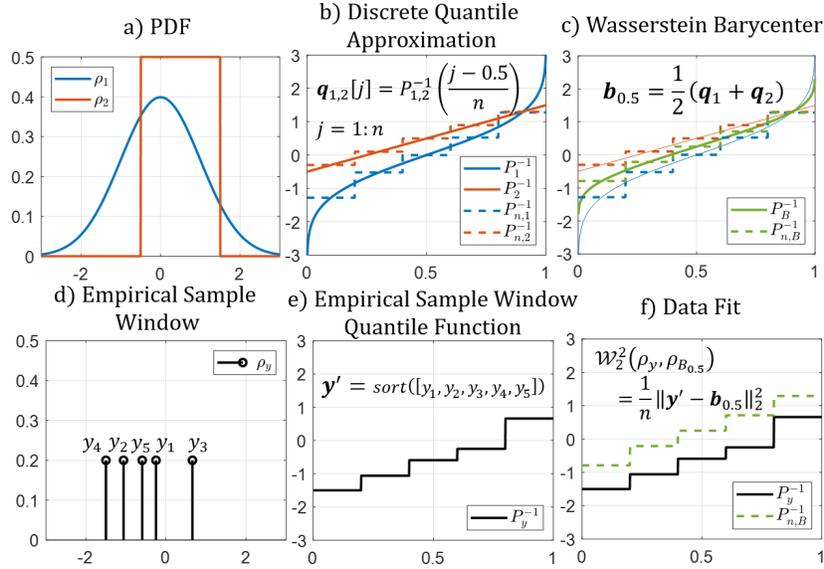


Fig. 6: **Discrete quantile representation for pure states, Wasserstein barycenter, and empirical distribution function.** (a) PDF of  $\rho_1, \rho_2$ . (b) Respective quantile function and  $n$ -DQA where  $n = 5$  with corresponding  $n$ -DQV  $\mathbf{q}_{1:2}$ . (c) Quantile function of the Wasserstein barycenter  $\rho_{B_{t^*}} = B([0.5, 0.5]^T, \rho_{1:2})$  and its  $n$ -DQV,  $\mathbf{b}_{t^*}$ , a weighted average of  $\mathbf{q}_1$  and  $\mathbf{q}_2$ . (d) Distribution according to Eq. (9) corresponding to sample window with  $n = 5$  points. (e), corresponding empirical quantile function  $P_y^{-1}$  where  $\mathbf{y}'$  is the sorted vector of samples. (f) Since  $\rho_{B_{t^*}}$  and  $P_y^{-1}$  are monotone step functions sharing the same set of discontinuities, from Eq. (3)  $\mathcal{W}_2^2(\rho_y, \rho_{n,B}) = \frac{1}{n} \|\mathbf{y}' - \mathbf{b}_{0.5}\|_2^2$

**Definition 1. [n-DQA]** Given a quantile function  $P^{-1} : [0, 1] \rightarrow \mathbb{R}$ , the  $n$ -point discrete quantile approximation ( $n$ -DQA) is a monotone step function  $P_n^{-1}(\xi) = P^{-1}\left(\frac{[\xi n] - 0.5}{n}\right)$  that is obtained by sampling the  $\{\frac{0.5}{n}, \frac{1.5}{n}, \dots, \frac{n-0.5}{n}\}$ -th quantiles of  $P^{-1}$ .

**Definition 2. [n-DQV]** The  $n$ -point discrete quantile vector ( $n$ -DQV)  $\mathbf{q} \in \mathbb{R}^n$  is comprised of the sampled quantiles from an  $n$ -point DQA:  $\mathbf{q}[j] = P^{-1}\left(\frac{j-0.5}{n}\right), j = 1, 2, \dots, n$ .

The first two plots in the top row of Fig. 6 illustrate these definitions. For a univariate distribution, the  $n$ -DQA approximates the quantile function with a monotone step function whose constant values are sampled from the quantile function on a uniform interval. Since the quantile function is monotone, these sampled quantiles and consequently the  $n$ -DQV is sorted in ascending order.

With this, we approximate the quantile functions of the pure state distributions  $P_{q_{1:K}}^{-1}$  with their respective  $n$ -DQA  $P_{n, q_{1:K}}^{-1}$ , parameterizing them according to their  $n$ -DQV, which are denoted as  $\mathbf{q}_{1:K}$ . Learning these  $n$ -DQVs amounts to estimating the  $\{\frac{0.5}{n}, \frac{1.5}{n}, \dots, \frac{n-0.5}{n}\}$ -quantiles of each pure state distribution. We also use this discrete quantile approach to estimate the quantile functions of the model Wasserstein barycenter  $\rho_{B_i} = B(\mathbf{x}_i, \rho_{q_{1:K}})$ . From Eq. (5) we see that the  $\xi$ -th quantile of the Wasserstein barycenter is a weighted combination of the  $\xi$ -th quantiles of  $\rho_{q_{1:K}}$ , with barycentric weight  $\mathbf{x}_i$ . Since the  $\mathbf{q}_{1:K}$  samples the quantile function of each pure state distribution at the same quantile values, the  $n$ -DQV of  $\rho_{B_i}$  is

$$\mathbf{b}_i = \sum_{k=1}^K \mathbf{x}_i[k] \mathbf{q}_k. \quad (14)$$

We denote corresponding  $n$ -DQA as  $P_{n, B_i}^{-1}$  and the corresponding distribution as  $\rho_{n, B_i}$ .

Using this  $n$ -DQV representation, by intentionally choosing  $n$ , the discretization level for the  $n$ -DQV, to be equal to the size of the sample window, we are able to pose a least-squares cost that approximates the data fit term in Eq. (11). To see this, we start by noting that all  $n$ -DQAs including  $P_{n, q_k}^{-1}, P_{n, B_i}^{-1}$  are monotone step functions on  $[0, 1]$  with discontinuities at  $\{\frac{1}{n}, \frac{2}{n}, \dots, \frac{n-1}{n}\}$ . Additionally, shown in the bottom row of Fig. 6, the distributions  $\rho_{y_i}$  corresponding to sample windows  $\mathbf{y}_i$  are discrete distributions comprised of Dirac-delta measures supported on a set of  $n$  points with uniform weights (Eq. (9)). By setting this window size  $n$  to be the same as the level of discretization used for the  $n$ -DQAs that approximate the pure state quantile functions, the quantile functions  $P_{y_i}^{-1}$  corresponding to  $\rho_{y_i}$  will have the same monotone piecewise constant structure with the same set of discontinuities as the  $n$ -DQA (bottom right of Fig. 6). With  $\mathbf{y}'_i$  being the vector obtained by sorting the elements of  $\mathbf{y}_i$  in increasing order, the Wasserstein distance between the  $\rho_{y_i}$  and  $\rho_{n, B_i}$  is simply

$$\mathcal{W}_2^2(\rho_{y_i}, \rho_{n, B_i}) = \frac{1}{n} \|\mathbf{y}'_i - \mathbf{b}_i\|_2^2, \quad (15)$$

where  $n$  is the size of the sample window corresponding to  $\rho_{y_i}$  and hence the length of  $\mathbf{y}'_i$ , as well as the level of discretization used for the  $n$ -DQA of the pure state distributions, and hence the length of  $\mathbf{b}_i$ .

#### D. Model Estimation and Algorithm

Using this discrete quantile parameterization for the pure state distributions allows us to pose an approximation to the variational objective function in Eq. (11) for the DWB model as a constrained nonlinear least squares problem. Let  $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_K] \in \mathbb{R}^{n \times K}$  denote the matrix whose columns correspond to the  $n$ -DQV of each of the pure states,  $\mathbf{Y} = [\mathbf{y}'_1, \dots, \mathbf{y}'_N] \in \mathbb{R}^{n \times N}$  the matrix whose columns correspond the sorted sample windows from the observed time series and  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{K \times N}$  the matrix whose columns correspond to the latent state vectors across time. Pulling together the regularizers discussed in Sec. IV-B, we pose the following constrained optimization problem,

$$\begin{aligned} \hat{\mathbf{Q}}, \hat{\mathbf{X}} &= \underset{\mathbf{Q}, \mathbf{X}}{\operatorname{argmin}} F(\mathbf{Q}, \mathbf{X}) \\ &= \underset{\mathbf{Q}, \mathbf{X}}{\operatorname{argmin}} \underbrace{\frac{1}{n} \|\mathbf{Y} - \mathbf{Q}\mathbf{X}\|_F^2}_{\text{Data fit}} + \lambda_x \underbrace{\sum_{i=1}^{N-1} d^2(\mathbf{X}[:, i+1], \mathbf{X}[:, i])}_{\text{Eq. (12)}} \\ &\quad + \lambda_q \underbrace{\frac{N}{n} \left\| \mathbf{Q} \left( \mathbf{I} - \frac{1}{K} \mathbb{1}_K \mathbb{1}_K^T \right) \right\|_F^2}_{\text{Eq. (13)}} \end{aligned} \quad (16)$$

subject to:

$$\mathbf{Q}[j+1, k] - \mathbf{Q}[j, k] \geq 0, \quad \begin{array}{l} k = 1, \dots, K, \\ j = 1, \dots, (n-1) \end{array} \quad (17)$$

$$\sum_{k=1}^K \mathbf{X}[k, i] = 1 \quad i = 1 : N \quad (18)$$

$$\mathbf{X}[k, i] \geq 0 \quad \begin{array}{l} k = 1, \dots, K \\ i = 1, \dots, N \end{array} \quad (19)$$

As detailed in Alg. 1, we minimize Eq. (16) using a block coordinate descent approach with two blocks, alternating between optimizing for  $\mathbf{Q}$  and  $\mathbf{X}$  while holding the other fixed. For both constrained optimization problems, we utilize the *sequential least squares programming* (SLSQP) optimizer [24] with python's `scipy` library [25].

The learned parameters of  $\hat{\mathbf{Q}}$  and  $\hat{\mathbf{X}}$  both have closed, non-empty, and convex constraints given by Eqs. (17)-(19). The  $n$ -DQV, and thus the columns of  $\mathbf{Q}$  must be sorted in ascending order (Eq. (17)). The barycentric weights, and thus the columns of  $\mathbf{X}$ , are constrained to the set of positive matrices with rows that sum to one (Eqs. (18), (19)). Thus, by [26], every limit point of an alternating block coordinate descent approach with two blocks is guaranteed to be a critical point. We show empirical convergence of this algorithm to a limit point in Fig. 10

Regarding the initialization of the parameters  $\hat{\mathbf{Q}}$  and  $\hat{\mathbf{X}}$ , we initialize the pure state distributions based on clustering the observation sample windows. Starting in line 14 in Alg. 1, we compute the similarity graph using the exponential of the negative square Wasserstein distance between the distributions estimated from any two sample windows and use spectral clustering from Python's `sklearn` [27] to learn  $K$  clusters. Denoting  $\mathcal{C}_k$  as the index set of sample windows that belong to cluster  $k$ , in line 22 we initialize the  $n$ -DQV of each pure state to be,  $\mathbf{q}_k = \frac{1}{|\mathcal{C}_k|} \sum_{i \in \mathcal{C}_k} \mathbf{y}'_i$ . We initialize the latent state  $\hat{\mathbf{X}}[t, k] = \frac{1}{K} \forall i = 1, \dots, N, k = 1, \dots, K$ , to be at the centroid of the simplex for all points in time.

## V. MODEL EVALUATION

Analyzing simulated and real world human activity data where the system gradually evolves among its pure states, we empirically demonstrate (1) how the proposed regularizers allow the DWB model to accurately recover the system parameters taking into account the inverse-scaling relationship between the model parameters, (2) the impact of the window size on the accuracy of estimating the system parameters of a dynamically evolving time series, and (3) how our non-parametric formulation of the DWB problem is able to accurately learn the pure state distributions.

### A. Simulated Experiments

Our approach proposes regularizers to the DWB model with the assumption that the latent state of the system transitions gradually among its pure states. Here, we consider a simulated system whose parameters reflect these properties, denoting the ground truth pure state distributions and latent states as  $\rho_{q_{1:K}}, \mathbf{x}_{1:T}$  and the model estimated parameters as  $\hat{\rho}_{q_{1:K}}, \hat{\mathbf{x}}_{1:T}$ .

---

**Algorithm 1: Non-parametric and Regularized DWB**


---

```

1 Input:
2  $y_1, \dots, y_t, \dots, y_T$ : Univariate time series
3  $t_1, \dots, t_i, \dots, t_N$ : Starting indices for sample windows
4  $K$ : Number of pure states
5 Hyperparameters:
6  $n$ : Sample window size    $\eta$ : Convergence threshold
7  $\lambda_x, \lambda_q$ : Regularization weights to define  $F$  (Eq. (16))
8 Output:
9  $\hat{\mathbf{Q}} \in \mathbb{R}^{n \times K}$ : Stacked pure state DQVs
10  $\hat{\mathbf{X}} \in \mathbb{R}^{K \times N}$ : Stacked barycentric weights
11 for  $i = 1 : N$  do
12    $\mathbf{y}'_i = \text{sort}(y_{t_i}, \dots, y_{t_i+n})$ ; // sorted windows
13 end
14 for  $i = 1 : N$  do // Window affinity matrix
15   for  $j = 1 : N$  do
16      $\mathbf{A}[i, j] = \frac{1}{n} \|\mathbf{y}'_i - \mathbf{y}'_j\|_2^2$ 
17   end
18 end
19  $c_{1:N} = \text{SpectralClustering}(K, \exp(-\mathbf{A}))$ ; // Cluster sample windows
20 for  $k = 1, \dots, K$  do
21    $\mathcal{C}_k = \{i : c_i = k\}$ 
22    $\mathbf{q}_k = \frac{1}{|\mathcal{C}_k|} \sum_{i \in \mathcal{C}_k} \mathbf{y}'_i$ ; // Init.  $n$ -DQV
23 end
24  $\mathbf{Y} = [\mathbf{y}'_1, \dots, \mathbf{y}'_N]$ ; // Stacked windows
25  $\hat{\mathbf{X}}^{(0)} = \frac{1}{K} \mathbf{1}_K \mathbf{1}_N^\top$ ; // Initialize  $\mathbf{X}$ 
26  $\hat{\mathbf{Q}}^{(0)} = [\mathbf{q}_1, \dots, \mathbf{q}_K]$ ; // Stacked  $n$ -DQVs
27 do
28    $\hat{\mathbf{X}}^{(i+1)} = \text{argmin}_{\mathbf{X}} F(\mathbf{Q} = \hat{\mathbf{Q}}^{(i)}, \mathbf{X} = \hat{\mathbf{X}}^{(i)})$ ; // SLSQP
29    $\hat{\mathbf{Q}}^{(i+1)} = \text{argmin}_{\mathbf{Q}} F(\mathbf{Q} = \hat{\mathbf{Q}}^{(i)}, \mathbf{X} = \hat{\mathbf{X}}^{(i+1)})$ ; // SLSQP
30 while  $F(\hat{\mathbf{Q}}^{(i)}, \hat{\mathbf{X}}^{(i)}) - F(\hat{\mathbf{Q}}^{(i+1)}, \hat{\mathbf{X}}^{(i+1)}) > \eta$ ;

```

---

Consider a system that consists of  $K = 3$  pure states with ground truth distributions,

$$\begin{aligned}
 \rho_{q_1} &= 0.5 \mathcal{N}(3, 0.25) + 0.25 \mathcal{N}(-3, 0.25) & (20) \\
 \rho_{q_2} &= \mathcal{U}[-4, 4] \\
 \rho_{q_3} &= \frac{1}{5} \sum_{i=1}^5 \mathcal{N}(\mathbf{q}[i], 1e^{-8}) \\
 \mathbf{q} &= [-2.88, -0.74, -0.64, -0.41, 1.82].
 \end{aligned}$$

Here,  $\mathcal{U}[a, b]$  denotes a uniform distribution on the interval  $[a, b]$ . We show the PDF in Fig. 7a, CDF in Fig. 7b, and quantile functions in Fig. 7c for each of these distributions

Now let us consider a time series that transitions among these three pure states following the trajectory outlined in Fig. 7d. The continuous-time latent state  $\mathbf{x}_\tau$  alternates between pausing at each pure state for 1 second, and then transitioning to another pure state for 2 seconds moving from  $\rho_1 \dots \rho_2 \dots \rho_3 \dots \rho_1$  over the course of continuous time  $\tau = [0, 9]$  seconds. To emulate the simulated setup in Sec. III-B of varying the per-sample change in distribution, we vary  $r$  (Hz), the rate at which we sample this continuous-time sequence to generate the ground truth latent state  $\mathbf{x}_t$  for  $t = 1, \dots, T$ . Time series are independently sampled  $y_t \sim \rho_{B_t} = B(\mathbf{x}_t, \rho_{q_{1:3}})$  by uniformly sampling a quantile  $\xi \in [0, 1]$  and evaluating the quantile function  $P_{B_t}^{-1}(\xi) = \sum_{k=1}^3 \mathbf{x}_t[k] P_{q_k}^{-1}(\xi)$ .

From this time series, we generate a sequence of sample windows using sliding window[28], spacing out the windows on a constant interval,  $t_{i+1} = t_i + \delta$ . For these experiments, we choose to set  $\delta = n$ , which partitions the time series into a sequence of  $N = \lfloor \frac{T}{n} \rfloor$  disjoint windows of size  $n$ .

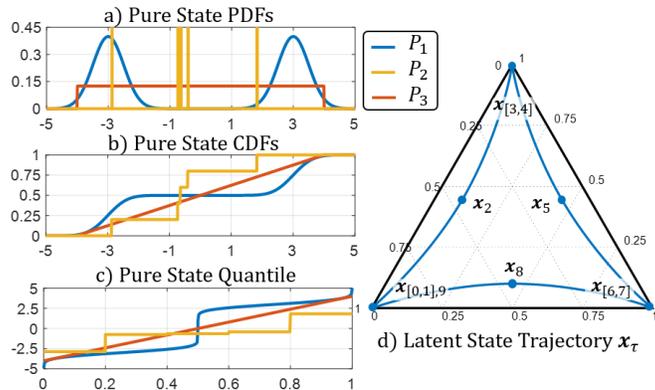


Fig. 7: **Simulated pure state distributions and latent state process.** (a) PDF, (b) CDF, and (c) quantile function for the pure state distributions used in the simulated experiments. (d) Simulated latent state trajectory as the system transitions among its three pure states.

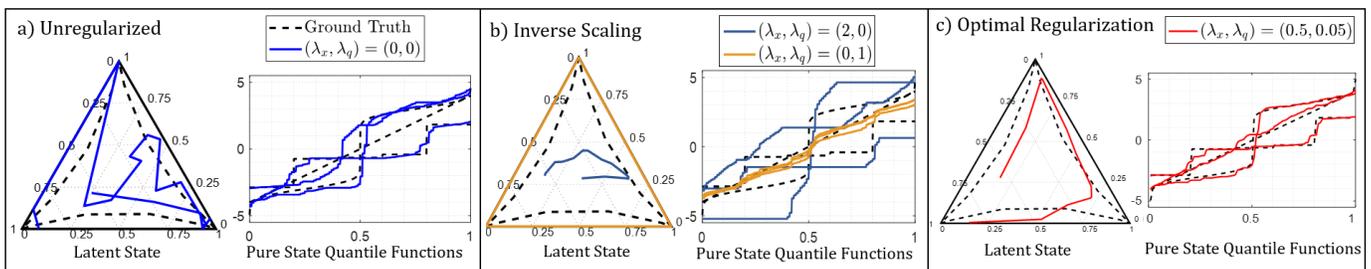


Fig. 8: **Impact of regularization and the inverse-scaling relationship.** (a) Without regularization, though the ground truth latent state evolves gradually, the learned latent state varies over the simplex. (b) Introducing just the latent state regularizer (blue) imposes some level of smoothness to the latent state, but consequently results in the pure state quantile functions diverging from their ground truth values. Conversely, including just the pure state regularizer (orange) causes the pure state quantile functions to move closer together and the latent states to move towards the boundary of the simplex. (c) Using regularization weights that minimize  $e_x + e_q$  balance the effects of these two regularizers to accurately recover the gradually evolving latent state sequence and the pure states quantile functions of the time series. Plots are shown for  $r = 200, n = 100$ .

With knowledge of ground truth, we can assess our model using the average distance of our learned parameters to these ground truth values according to,

$$e_q = \frac{1}{K} \sum_{k=1}^K \mathcal{W}_2^2(\rho_{q_k}, \hat{\rho}_{q_k}) \quad e_x = \frac{1}{N} \sum_{t=1}^N \|\mathbf{x}_t - \hat{\mathbf{x}}_t\|_2^2. \quad (21)$$

Since our model is unsupervised, there is no guarantee that the indexing of the learned pure states will match that of the ground truth. Therefore, when assessing our model, we assume that learned pure states (and subsequently the latent state vector) are reordered in a manner that minimizes Eq. (21).

1) *Regularization and Inverse-Scaling*: We demonstrate the inverse scaling relationship between the model parameters through the interaction between the latent state and pure state regularizers. As discussed in Sec. III-B, the difficulty in accurately estimating the model parameters stems from the problem of estimating dynamically evolving data distribution with a window of finite length. In Fig. 8, we see that in the absence of any regularization, the learned latent state can vary over the simplex, even when the ground truth latent state is stationary or gradually evolving. The blue lines of Fig. 8b show that introducing only the latent state regularizer that penalizes  $d^2(\mathbf{x}_t, \mathbf{x}_{t+1})$  imposes some level of smoothness on the latent state. However, as discussed in Sec. IV-B this regularizer has an additional effect that is similar to the  $\alpha > 1$  regime corresponding to the blue lines in Fig. 3 where the latent states trajectory contracts towards a point on the simplex. In this case, as specified by the inverse-scaling relationship, seen from the blue quantile plots in Fig. 8b, the pure state quantile functions diverge from the ground truth values. On the other hand, as seen in orange lines of Fig. 8b, having only the pure state regularizer causes the pure state quantile functions to be pulled closer together where by the inverse-scaling relationship, the latent states move away from the centroid of the simplex towards the boundary.

Only through the combination of these regularization terms can we recover the ground truth parameters of this simulated system that gradually evolves among its pure states. Using our knowledge of ground truth we perform a grid search by varying

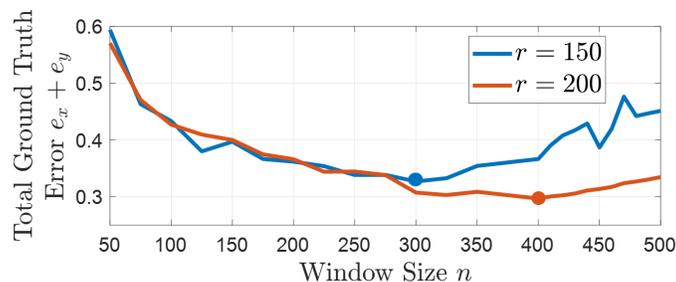


Fig. 9: **Impact of window size on simulated model accuracy.** Ground truth error averaged over 500 generated time series as a function of window size  $n$  for sampling rates of  $r = 150, 200$  Hz. The U-shape plots imply that the factors that impact the ability of small and large windows to estimate a dynamically evolving data distribution similarly impact the ground truth error of the learned DWB model parameters. Additionally, increasing the sampling rate  $r$  of the time series, which results in a smaller per-sample change in the data distribution, similarly improves the model accuracy for larger windows and shifts the minimum of the U-curve towards larger windows. The minimum error (solid dots) for  $r = 150$  was achieved with a window size of  $n = 300$ , while for  $r = 200$  the minimum window size was achieved at  $n = 400$ .

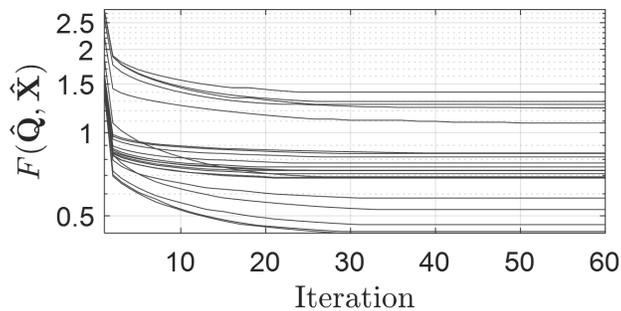


Fig. 10: **Convergence of simulated experiments:** Plot of the objective function shows convergence to a limit point. Shown for simulated data with  $n = 400, 410, \dots, 500$  and  $r = 150, 200$ .

the regularization weights  $\lambda_x, \lambda_q = \{1e^{-4}, 2e^{-4}, 5e^{-4}, 1e^{-3}, \dots, 1e^{-1}\}$  and picking the pair that minimizes the total ground truth error  $e_x + e_q$ . As seen from Fig. 8c, the DWB model corresponding to these optimal regularization weights more accurately recover the ground truth parameters of this simulated system. In Sec. V-B, we perform parameter selection in the absence of ground truth using L-surfaces [29].

2) *Impact of Window Size:* In Sec. III-B, we highlighted the effect of window size and the rate of change in the data distribution on the accuracy of estimating a dynamically evolving distribution using a window of samples. Here, we seek to understand how these factors impact the accuracy of the learned DWB parameters. Using the aforementioned setup in Sec. V-A, we simulate two systems with different sampling rates  $r \in \{150, 200\}$  Hz, generating 500 time series per value of  $r$ . We then run our DWB model varying the window size between  $n = 50, 100, \dots, 400, 410, \dots, 500$ . We find the optimal regularization weights by performing a grid search over the range of values as Sec. V-A-1. For each value of  $n$  and  $r$ , we set  $\lambda_x, \lambda_q$  to be the average of the results of this grid search performed for 5 randomly selected time series.

In this experiment, we remove potential confounding variables across the various window sizes. First, we use the same initialization for the pure states distributions, specifically the one generated from the spectral clustering method discussed in Sec. IV-D for the configuration of  $n = 200$ . Secondly, we ensure that the sequence of sample windows estimates the time series at the same points in time, such that regardless of  $r, n$  the windows are centered at  $\tau = [1.5, 2.0, \dots, 7.5]$ .

The results shown in Fig. 9 imply that the factors discussed in Sec. III-B that impact the ability of a sample window to accurately estimate a dynamically evolving data distribution similarly affect the ability of the DWB model to accurately learn a system’s pure states and latent states from these estimated windows. The U-shape curves in Fig. 9 of the average ground truth error as a function of window size bear strong similarities to the U-shape curves illustrating the window size tradeoff in Fig. 5 where small windows lack the samples to precisely estimate the data distribution and the accuracy of large windows suffer due to the dynamics of the systems. Furthermore, we see the same trend as in Sec. III-B where decreasing the per-sample change in distribution, shown here by increasing the sampling rate  $r$ , improves the accuracy for models with larger windows shifting this window size tradeoff and the “optimal” window size towards larger windows.

3) *Convergence:* Convergence of the optimization process to a limit point is shown in Fig. 10 for various configurations of the simulated data defined in Sec. V-A.

## B. Real World Data

In this section, we compare the performance of the non-parametric DWB model to the Gaussian DWB model [1] on univariate data. To highlight the difference between the Gaussian and non-parametric discrete quantile parameterizations, we use the regularization framework proposed in this paper for both models. We evaluate using our simulated data and two human activity datasets.

- 1) *Beep Test (BT, proprietary)*: Subjects run between two points to a metronome with increasing frequency, alternating between two states: running and standing. We use the vertical component (z-axis) of the 3-axis accelerometer, which is the dimension in which the distributions of the two pure states are best differentiated. The sensor is sampled at 100 Hz.
- 2) *Microsoft Research Human Activity (MSR, [30])*: 126 subjects perform exercises in a gym setting. Exercises vary among subjects covering strength, cardio, cross-fit, and static exercises. Each time series is truncated to five minutes. Discrete labels corresponding to activities are provided, thus we set  $K$  to the number of labeled discrete states in the truncated time series (range:  $K = 2$  to 7). We use the x-axis of the 3-axis accelerometer sampled at 50 Hz.
- 3) *Simulated Data (Sim)*: Following the data generating process outlined in Sec. V-A, we simulate 500 time series setting  $r = 200$ .

**Evaluation:** Since ground truth is not known in the real world setting, we assess performance by considering the model fit to the data, computed using the data-fit term in Eq. (11),

$$e_y = \frac{1}{N} \sum_{t=1}^N \mathcal{W}_2^2(\rho_{y_t}, \rho_{B_t}) \quad (22)$$

In the non-parametric model, this distance is computed according to Eq. (15). In the Gaussian DWB model where  $\rho_{B_t}$  is Gaussian, this distance is computed using a Monte-Carlo method using  $1e^5$  IID samples from  $\rho_{B_t}$ .

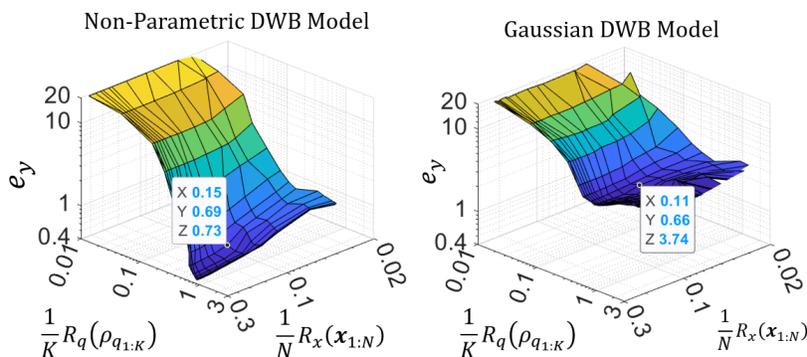


Fig. 11: **L-surface for selection of regularization weights.** L-surface for MSR time series for (left) non-parametric and (right) Gaussian DWB models varying  $\lambda_x, \lambda_q$  from  $[1e^{-5}, 1]$ . Each figure plots the magnitude of the data loss against the magnitude of the two regularization parameters on a log scale. Selecting the regularization weight according to L-surfaces amounts finding the  $\lambda_x, \lambda_q$  corresponding to the "corner" of the surface plot, the point where further increasing  $\lambda_x, \lambda_q$  (which corresponds to moving toward the back corner of the plot, decreasing  $\frac{1}{K}R_q(\rho_{q_{1:K}})$  and  $\frac{1}{N}R_x(x_{1:T})$ ) results in a sharp increase in the magnitude of the data loss.

To select regularization weights when ground truth is not available, we draw from the field of inverse problems and use L-surfaces [29]. Fig. 11 shows the L-surface for one MSR dataset plotting  $e_y$  against the magnitude of the  $R_x$  and  $R_q$  on a log scale while varying  $\lambda_x, \lambda_q$  ranging from  $1e^{-5}$  to 1. For each model configuration and dataset, we use this L-surface method to pick a set of regularization parameters based on one representative time series and apply those parameters to the rest of the dataset. These parameters are detailed in Tab. I.

Fig. 12 compares the learned pure states of the Gaussian DWB and non-parametric DWB for one example MSR time series. As seen from Fig. 12b, the distributions of the data corresponding to many of the activities are clearly multi-modal and therefore not Gaussian. Compared to the Gaussian DWB approach shown in Fig. 12c, the pure states learned using the proposed non-parametric representation shown in Fig. 12d more closely match the estimated pure states from the data. This is also reflected quantitatively in Tab. II where the values of  $e_y$  show that our non-parametric DWB model better approximates the sample windows of the time series compared to the Gaussian DWB model for each of the evaluated datasets.

We note that the magnitude of the  $e_y$ , and thus the results in Tab. II are dependent on the choice of the regularization parameters, which are chosen in a partially subjective manner. However, we can see in the L-surfaces in Fig. 11 that compared to the Gaussian model the error is significantly lower in the non-parametric model. Therefore any small subjective changes in the choice of regularization parameter would not significantly alter these conclusions.

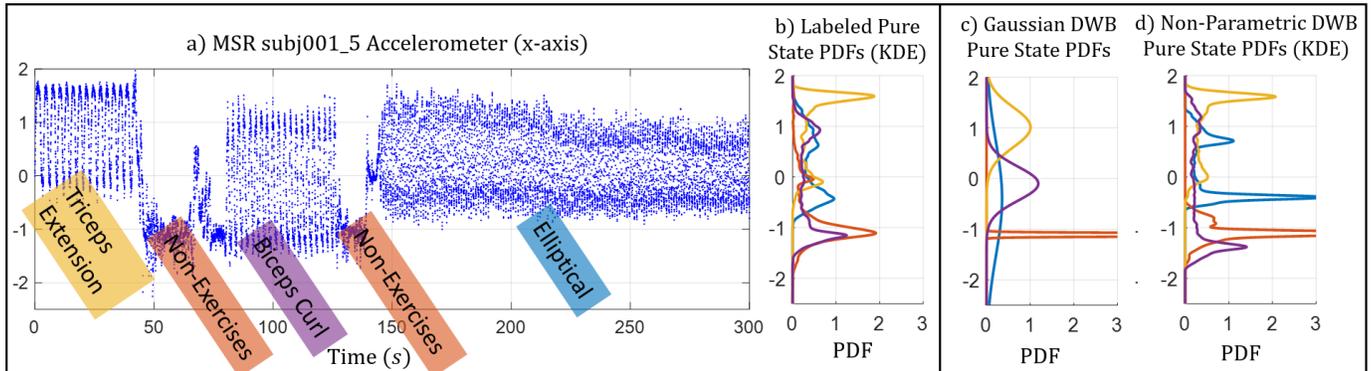


Fig. 12: **Non-parametric vs Gaussian pure states distributions for MSR data.** (a) MSR time series consisting of 4 activities over 5 minutes. (b) Pure state distributions estimated from discrete labeled data converted to a PDF using KDE with a Gaussian kernel ( $\sigma = 0.04$ ). (c) Learned pure states for Gaussian model from [1]. (d) Learned pure states for our proposed non-parametric model. The  $n$ -DQV  $\mathbf{q}_k$  for each pure state is converted to a PDF using KDE with the same Gaussian kernel centered on the values of  $\mathbf{q}_k$ . The distribution of the data in each pure state activity (blue: Elliptical, yellow: Tricep Extension, purple: Biceps Curl, Red: Non-exercise) is better captured using the non-parametric approach  $e_y = 0.73$  compared to the Gaussian  $e_y = 3.74$ .

	MSR	BT	Sim
$n$	250	100	100
NP $\lambda_x, \lambda_q$	$5e^{-2}, 5e^{-3}$	$2e^{-1}, 5e^{-2}$	$2e^{-1}1e^{-2}$
Gauss $\lambda_x, \lambda_q$	$5e^{-2}, 2e^{-3}$	$2e^{-1}, 2e^{-2}$	$1e^{-1}2e^{-3}$

TABLE I: **DWB Model configuration.** Window size and regularization weights for non-parametric (NP) and Gaussian (Gauss) DWB models for simulated (Sim) and real world (BT, MSR) datasets.

## VI. CONCLUSIONS AND FUTURE WORK

In this work, we build upon the DWB model in [1]. We present and discuss the inverse scaling relationship which captures the lack of uniqueness in the DWB model. We also consider the challenges of estimating the data distribution of a dynamically evolving time series. We then propose a temporal smoothness regularization framework to simultaneously address both of these challenges. Finally, we move beyond the Gaussian assumption of [1] by using a discrete approximation to the pure state quantile function which results in a least-square “data-fit” term in the DWB objective function. Using simulated data, we demonstrate how the two proposed regularization terms work together to achieve the desired smoothness in the time evolution of the latent state as well as the impact of window size on the accuracy of the learned model. In the real world setting of human activity analysis, we demonstrate how compared to the original Gaussian DWB model our non-parametric DWB model better characterizes the time-varying data distribution of the time series.

An important future work is to extend the univariate framework to the general non-parametric multivariate setting. Lack of closed form expressions for barycenters in the non-parametric multivariate setting as well as the “curse of dimensionality” in approximating a high-dimensional distribution from samples [31] makes this extension numerically and statistically challenging. However, fast algorithms for computing Wasserstein barycenters [32] can be effectively utilized for low dimensional problems.

Furthermore, in this work we only consider the Wasserstein barycenter to model the change in distribution as a system moves among pure states. An interesting topic for future research is to consider barycenters corresponding to alternative probability distribution distances such as the Sinkhorn divergence [33].

Our choice to approximate the quantile function of  $\rho_{q_k}$  by sampling on a uniform interval combined with the use of a Dirac-delta functions to estimate  $\rho_{y_i}$  from windows of samples lead to the finite dimensional least-squares data fit term in Eq. (16). Future work can consider alternative methods of estimating distributions from samples (e.g kernel-density estimation [34]) and quantile approximation methods (e.g linear or spline decomposition [35]) to derive alternative forms of this univariate

	MSR	BT	Sim
NP	<b>1.32</b>	<b>2.07</b>	<b>2.19</b>
Gauss	3.76	6.00	12.33

TABLE II: **Quantitative comparison of DWB model.** Average  $e_y$  across all time series for simulated (Sim) and real world (BT, MSR) datasets show clear benefits of the non-parametric DWB model compared to the Gaussian data-generating distribution model used by [1].

DWB problem.

We empirically show the monotonic convergence of the two-block cyclic descent method used in the paper. The theory in [26] indicates that we are converging to a critical point of the cost function. In future work, we can seek stronger guarantees of converging to a second-order stationary point using a proximal point block coordinate descent algorithm [36].

Finally, in this work, we also explore the issue of approximating the data distribution of a dynamically evolving system from a window of independent, non-identically distributed random variables. We demonstrate the tradeoff between the errors associated with large and small windows using simulated data to illustrate the effects of system parameters that drive the per-sample change in the distribution on the accuracy of this window estimate. Future work may consider an analytical approach to this problem using the relationship between quantiles and order statistics [37] [38] and by bounding the maximum change in distribution over the length of the window.

## REFERENCES

- [1] K. Cheng, S. Aeron, M. C. Hughes, and E. L. Miller, “Dynamical Wasserstein Barycenters for Time-series Modeling,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 27991–28003, 2021.
- [2] D. Mukhin, A. Gavrilov, A. Feigin, E. Loskutov, and J. Kurths, “Principal nonlinear dynamical modes of climate variability,” *Scientific reports*, vol. 5, no. 1, pp. 1–11, 2015, publisher: Nature Publishing Group.
- [3] S. A. Imtiaz, “A systematic review of sensing technologies for wearable sleep staging,” *Sensors*, vol. 21, no. 5, p. 1562, 2021, publisher: MDPI.
- [4] Y. Chi, T. Yang, and P. Zhang, “Dynamical Mode Recognition of Triple Flickering Buoyant Diffusion Flames: from Physical Space to Phase Space and to Wasserstein Space,” *arXiv preprint arXiv:2201.01085*, 2022.
- [5] M. C. Hughes and E. B. Sudderth, “Nonparametric discovery of activity patterns from video collections,” in *2012 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*. IEEE, 2012, pp. 25–32.
- [6] O. D. Lara and M. A. Labrador, “A survey on human activity recognition using wearable sensors,” *IEEE communications surveys & tutorials*, vol. 15, no. 3, pp. 1192–1209, 2012, publisher: IEEE.
- [7] R. Kalman, “A New Approach to Linear Filtering and Prediction Problems,” *Trans. ASME, J. Basic Eng.*, vol. 82, pp. 34–45, 1960.
- [8] M. I. Ribeiro, “Kalman and extended kalman filters: Concept, derivation and properties,” *Institute for Systems and Robotics*, vol. 43, p. 46, 2004.
- [9] E. Fox, E. Sudderth, M. Jordan, and A. Willsky, “Nonparametric Bayesian learning of switching linear dynamical systems,” *Advances in neural information processing systems*, vol. 21, 2008.
- [10] L. R. Rabiner, “A tutorial on hidden Markov models and selected applications in speech recognition,” *Proceedings of the IEEE*, vol. 77, no. 2, pp. 257–286, 1989, publisher: Ieee.
- [11] Z. Ghahramani and M. Jordan, “Factorial hidden Markov models,” *Advances in Neural Information Processing Systems*, vol. 8, 1995.
- [12] E. B. Fox, M. C. Hughes, E. B. Sudderth, and M. I. Jordan, “Joint modeling of multiple time series via the beta process with application to motion capture segmentation,” *The Annals of Applied Statistics*, vol. 8, no. 3, pp. 1281–1313, 2014, publisher: Institute of Mathematical Statistics.
- [13] R. J. McCann, “A convexity principle for interacting gases,” *Advances in mathematics*, vol. 128, no. 1, pp. 153–179, 1997, publisher: Elsevier.
- [14] M. Agueh and G. Carlier, “Barycenters in the Wasserstein space,” *SIAM Journal on Mathematical Analysis*, vol. 43, no. 2, pp. 904–924, 2011, publisher: SIAM.
- [15] T.-M. Nguyen and S. Volkov, “On a class of random walks in simplexes,” *Journal of Applied Probability*, vol. 57, no. 2, pp. 409–428, 2020, publisher: Cambridge University Press.
- [16] J. Martín-Fernández, C. Barceló-Vidal, V. Pawłowsky-Glahn, A. Buccianti, G. Nardi, and R. Potenza, “Measures of difference for compositional data and hierarchical clustering methods,” in *Proceedings of IAMG*, vol. 98, 1998, pp. 526–531, issue: 1.
- [17] F. Åström, S. Petra, B. Schmitzer, and C. Schnörr, “Image labeling by assignment,” *Journal of Mathematical Imaging and Vision*, vol. 58, no. 2, pp. 211–238, 2017, publisher: Springer.
- [18] G. Peyré, M. Cuturi, and others, “Computational optimal transport: With applications to data science,” *Foundations and Trends® in Machine Learning*, vol. 11, no. 5-6, pp. 355–607, 2019, publisher: Now Publishers, Inc.
- [19] F. Santambrogio, “Optimal transport for applied mathematicians,” *Birkhäuser, NY*, vol. 55, no. 58-63, p. 94, 2015, publisher: Springer.
- [20] N. Bonneel and H. Pfister, “Sliced Wasserstein barycenter of multiple densities,” 2013.
- [21] P. Embrechts and M. Hofert, “A note on generalized inverses,” *Mathematical Methods of Operations Research*, vol. 77, no. 3, pp. 423–432, 2013, publisher: Springer.
- [22] R. Rakestraw, “The convex cone of n-monotone functions,” *Pacific Journal of Mathematics*, vol. 43, no. 3, pp. 735–752, 1972, publisher: Mathematical Sciences Publishers.
- [23] N. Fournier and A. Guillin, “On the rate of convergence in Wasserstein distance of the empirical measure,” *Probability Theory and Related Fields*, vol. 162, no. 3, pp. 707–738, 2015, publisher: Springer.
- [24] D. Kraft, “A software package for sequential quadratic programming,” *Forschungsbericht- Deutsche Forschungs- und Versuchsanstalt für Luft- und Raumfahrt*, 1988.
- [25] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, and others, “SciPy 1.0: fundamental algorithms for scientific computing in Python,” *Nature methods*, vol. 17, no. 3, pp. 261–272, 2020, publisher: Nature Publishing Group.
- [26] L. Grippo and M. Sciandrone, “On the convergence of the block nonlinear Gauss–Seidel method under convex constraints,” *Operations research letters*, vol. 26, no. 3, pp. 127–136, 2000, publisher: Elsevier.
- [27] L. Buitinck, G. Louppe, M. Blondel, F. Pedregosa, A. Mueller, O. Grisel, V. Niculae, P. Prettenhofer, A. Gramfort, J. Grobler, and others, “API design for machine learning software: experiences from the scikit-learn project,” in *European Conference on Machine Learning and Principles and Practices of Knowledge Discovery in Databases*, 2013.
- [28] S. Aminikhanghahi and D. J. Cook, “A survey of methods for time series change point detection,” *Knowledge and information systems*, vol. 51, no. 2, pp. 339–367, 2017, publisher: Springer.
- [29] D. H. Brooks, G. F. Ahmad, R. S. MacLeod, and G. M. Maratos, “Inverse electrocardiography by simultaneous imposition of multiple constraints,” *IEEE Transactions on Biomedical Engineering*, vol. 46, no. 1, pp. 3–18, 1999, publisher: IEEE.
- [30] D. Morris, T. S. Saponas, A. Guillory, and I. Kerner, “RecoFit: using a wearable sensor to find, recognize, and count repetitive exercises,” in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 2014, pp. 3225–3234.
- [31] J. WEED and F. BACH, “Sharp asymptotic and finite-sample rates of convergence of empirical measures in Wasserstein distance,” *Bernoulli*, vol. 25, no. 4A, pp. 2620–2648, 2019.
- [32] M. Cuturi and A. Doucet, “Fast computation of Wasserstein barycenters,” in *International conference on machine learning*. PMLR, 2014, pp. 685–693.
- [33] Z. Shen, Z. Wang, A. Ribeiro, and H. Hassani, “Sinkhorn barycenter via functional gradient descent,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 986–996, 2020.

- [34] A. Harvey and V. Oryshchenko, "Kernel density estimation for time series data," *International journal of forecasting*, vol. 28, no. 1, pp. 3–14, 2012, publisher: Elsevier.
- [35] M. Daehlen and T. Lyche, "Decomposition of splines," in *Mathematical methods in computer aided geometric design II*. Elsevier, 1992, pp. 135–160.
- [36] Q. Li, Z. Zhu, and G. Tang, "Alternating minimizations converge to second-order optimal solutions," in *International Conference on Machine Learning*. PMLR, 2019, pp. 3935–3943.
- [37] H. A. David and H. N. Nagaraja, *Order statistics*. John Wiley & Sons, 2004.
- [38] P. K. Sen, "A note on order statistics for heterogeneous distributions," *The Annals of Mathematical Statistics*, vol. 41, no. 6, pp. 2137–2139, 1970, publisher: Institute of Mathematical Statistics.

APPENDIX A  
EXTENSIONS OF DWB NON-UNIQUENESS

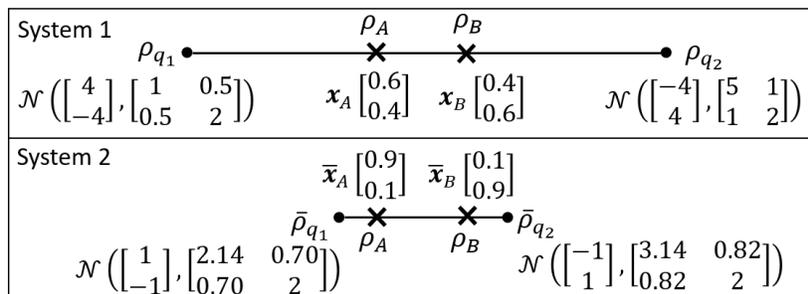


Fig. 13: **Non-uniqueness of Wasserstein barycenter parameters for multivariate Gaussians.** Both systems of pure state distributions and barycentric weights result in the same barycenters,  $\rho_A = B(\mathbf{x}_A, \rho_{q_{1:2}}) = B(\bar{\mathbf{x}}_A, \bar{\rho}_{q_{1:2}})$  and  $\rho_B = B(\mathbf{x}_B, \rho_{q_{1:2}}) = B(\bar{\mathbf{x}}_B, \bar{\rho}_{q_{1:2}})$ . Consistent with the inverse-scaling relationship discussed in Sec. III-A, compared to system 1, in system 2 the Wasserstein distance between  $\bar{\rho}_{q_1}$  and  $\bar{\rho}_{q_2}$  is diminished while the distance between  $\bar{\mathbf{x}}_A$  and  $\bar{\mathbf{x}}_B$  is increased.

While the discussion of this paper pertains mainly to the univariate case, the issue of uniqueness discussed in Sec. III also exists in the multivariate case. Consider the two dimensional systems specified in Fig. 13 of systems with two Gaussian pure states. For both systems,  $\rho_A = B(\mathbf{x}_A, \rho_{q_{1:2}}) = B(\bar{\mathbf{x}}_A, \bar{\rho}_{q_{1:2}})$  and  $\rho_B = B(\mathbf{x}_B, \rho_{q_{1:2}}) = B(\bar{\mathbf{x}}_B, \bar{\rho}_{q_{1:2}})$ .

This multivariate example also exhibits the inverse scaling relationship discussed in Sec. III-A. The Wasserstein distance between the pure states in system 1 ( $\rho_{q_1}, \rho_{q_2}$ ) is *larger* than that of system 2 ( $\bar{\rho}_{q_1}, \bar{\rho}_{q_2}$ ), however, the resulting distance between the latent states corresponding to  $\rho_A, \rho_B$  in system 1 ( $\mathbf{x}_A, \mathbf{x}_B$ ) is *smaller* than that of system 2 ( $\bar{\mathbf{x}}_A, \bar{\mathbf{x}}_B$ ).

Furthermore, we can create an example where the construction specified in Sec. III-A does not result in multiple possible values of  $\mathbf{x}_B$  and  $P_{q_{1:K}}^{-1}$  and thus cannot be used to demonstrate the non-uniqueness of the Wasserstein barycenter parameters.

Let  $\rho_{q_1} = \delta_0$  and  $\rho_{q_2} = U[0, 1]$ , where  $P_{q_1}^{-1}(\xi) = 0$  and  $P_{q_2}^{-1}(\xi) = \xi$  for  $\xi \in [0, 1]$ . Then the Wasserstein barycenter  $\rho_t = B(\mathbf{x}_t, \rho_{1:2})$  for any barycentric weight  $\mathbf{x}_t = [(1-t), t]^T$  will have distribution  $U[0, t]$  and quantile function  $P_t^{-1}(\xi) = t\xi$ . Per the construction in Sec. III-A, let us pick  $\mathbf{x}_B = [1, 0]$  to be at a vertex, thus  $P_B^{-1}(\xi) = 0$  and let  $\mathbf{x}_0 = [0.5, 0.5]^T$ , thus in  $P_0^{-1}(\xi) = 0.5\xi$ . According to Eq. (7),  $\bar{\mathbf{x}}_B$  falls of the simplex for any  $\alpha < 1$  thus making  $\alpha_0 = 1$ . Similarly, according to Eq. (8),  $\bar{P}_B^{-1}(\xi) = (1 - \alpha)t\xi$ , which breaks the monotonically increasing constraint of quantile functions for any  $\alpha > 1$ , thus restricting  $\alpha_m = 1$ . Thus, according to the construction in Sec. III-A, the set  $[\alpha_0, \alpha_m] = 1$ .

Although this example breaks the argument of the non-uniqueness of the Wasserstein barycenter parameters for this specific construction used to highlight the inverse-scaling relationship, it does not mean that the parameters of the Wasserstein barycenter are indeed unique. For example, since  $\mathbf{x}_B$  is taken to be at a vertex  $[1, 0]^T$  resulting in  $\rho_B = \rho_{q_1}$ , any valid distribution can be chosen for  $\rho_{q_2}$  and still have  $\rho_B = B([1, 0]^T, \rho_{q_{1:2}})$ .

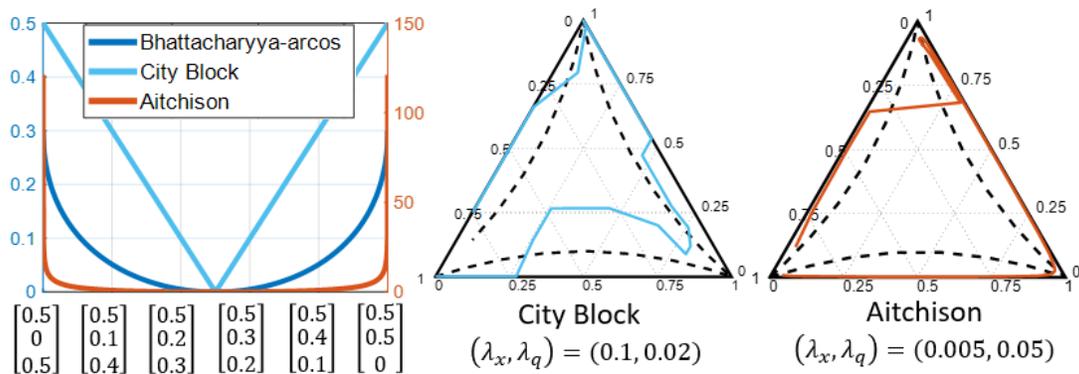


Fig. 14: **Comparison of alternative simplex distances:** (left) Example of the distance profiles for various simplex distances given by  $d^2(\mathbf{x}_0, \mathbf{x})$  where  $\mathbf{x}_0 = [0.5, 0.25, 0.25]^\top$  and  $\mathbf{x} \in [0.5, 0, 0.5]^\top, [0.5, 0.5, 0]^\top$ . The Bhattacharyya-arccos and city-block distances are plotted against the left axis while the Aitchison, which diverges at the simplex boundary, is plotted against the right axis. (Center and right) The learned simplex trajectory using the city-block and Aitchison distance in the regularizer (Eq. (12)) using the same data shown in Fig. 8 which uses the Bhattacharyya-arccos distance.  $\lambda_x, \lambda_q$  are selected via grid search using the parameters specified in Sec. V-A and  $r = 200, n = 100$ .

## APPENDIX B ALTERNATIVE SIMPLEX DISTANCES

In this work, the Bhattacharyya-arccos distance was used as a regularizer on the latent state. However, a variety of simplex distances can be used [16], the choice of which may be application dependent. We compare the Bhattacharyya-arccos, city-block<sup>3</sup>, and Aitchison distance in Fig. 14. The Aitchison distance diverges as one of the points moves towards the simplex boundary ( $\mathbf{x}$  has one or more zeros). As a result, when using this distance as the regularizer in our DWB model, the learned simplex trajectory will also avoid edges of the simplex. This is not the case for the city-block or Bhattacharyya-arccos distance which is finite for any two points on the simplex. In our simulated and real world experiments, we did not find significant differences in the ground truth error among these three simplex distances. We leave further investigation of these distance properties and their effect as regularizers to future work.

<sup>3</sup>We approximate  $|x| \approx \sqrt{x^2 + \epsilon}$  with  $\epsilon = 1e-8$