# A Weak Characterization of Slow Variables in Stochastic Dynamical Systems

Andreas Bittracher[1(✉)] and Christof Schütte[1,2]

[1] Department of Mathematics and Computer Science, Freie Universität Berlin, Berlin, Germany
bittracher@mi.fu-berlin.de, Christof.Schuette@fu-berlin.de
[2] Zuse Institute Berlin, Berlin, Germany

**Abstract.** We present a novel characterization of slow variables for continuous Markov processes that provably preserve the slow timescales. These slow variables are known as reaction coordinates in molecular dynamical applications, where they play a key role in system analysis and coarse graining. The defining characteristics of these slow variables is that they parametrize a so-called transition manifold, a low-dimensional manifold in a certain density function space that emerges with progressive equilibration of the system's fast variables. The existence of said manifold was previously predicted for certain classes of metastable and slow-fast systems. However, in the original work, the existence of the manifold hinges on the pointwise convergence of the system's transition density functions towards it. We show in this work that a convergence in average with respect to the system's stationary measure is sufficient to yield reaction coordinates with the same key qualities. This allows one to accurately predict the timescale preservation in systems where the old theory is not applicable or would give overly pessimistic results. Moreover, the new characterization is still constructive, in that it allows for the algorithmic identification of a good slow variable. The improved characterization, the error prediction and the variable construction are demonstrated by a small metastable system.

## 1 Introduction

The ability and practice to perform all-atom molecular simulations of more and more complex biochemical systems has led to an unprecedented increase in the available amount of dynamical data about those systems. This has exponentiated the importance to identify good chemical reaction coordinates (RCs), low-dimensional observables of the full system that are associated with the relevant, often slowly-progressing sub-processes. For one, a meaningful RC permits insight into the essential mechanisms and parameters of a reaction, by acting as a filter

for the overwhelming complexity of the data. As an example, computing the free energy (also known as the potential of mean force) along such a coordinate is typically used for identifying energy barriers and associated transition states [10,37]. RCs are also essential for the development of accurate reduced dynamical models. The Mori-Zwanzig formalism and related schemes [18,27,42,44] can be used to derive approximate closed equations of motion of the dynamics projected onto the image space of the RC. Depending on the chosen RC, the essential dynamical properties of the reduced model — such as transition rates between reactant and product — may or may not resemble those of the original system [43]. Finally, accelerated sampling schemes such as metadynamics [20], Blue Moon sampling [8] and umbrella sampling [38] also rely heavily on an accurate RC to guide them efficiently into unexplored territory.

In each of those applications, the result depends crucially on the "quality" of the RC, an elusive measure for how well the RC suits the specified task. In most cases, this quality can be brought down to how well the RC "captures the essential dynamics", in particular the rates of transitions between reactant and product state (see also [28] for an in-depth review on the effect of poorly chosen RCs on different classic rate theories). Due to this ambiguity, the search for universal and mathematically rigorous optimality criteria for RCs remains an active field of research, and numerous new approaches have been suggested during the last decade. For reactions involving one clearly defined reactant and product state, a in multiple ways ideal RC is the committor function [3,23], a one-dimensional observable that in each point describes the probability to hit the product state before returning to the reactant state. As the committor function is notoriously hard to compute, advanced numerical schemes have been developed to either approximate it efficiently [12], or find RCs that are equivalent by certain metrics [29]. Still, the computation of committor-like RCs often remains out of reach for high-dimensional systems.

For systems where the relevant behavior involves transitions between more than tw states [36], where the reaction is not adequately described by a transition between isolated states [35], or where the states are not known or cannot be computed, other optimality criteria must be employed. Here one common approach is to demand the preservation of the system's longest (equilibration) time scales under projection of the dynamics onto the RC. This leads naturally to a characterization of RCs in terms of the eigenvalues of the system's transfer operator, a widely used mathematical tool for time scale analysis in molecular dynamics and beyond [7,11,19,34,41]. It is in this setting where the authors and coworkers have previously proposed a novel mathematical framework for the characterization and numerical computation of ideal RCs [6]. The proposed theory builds on the insight that in many systems, the equilibration of the fast sub-processes over time manifests as the convergence of the system's transition density functions towards a certain low-dimensional manifold in density space, the so-called transition manifold (TM). This convergence is observed even if there is no equivalent low-dimensional structure in state space, such as a transition pathway between isolated states. Any parametrization of the TM then can in theory be used to construct an ideal RC.

The framework demands that the convergence towards the TM must occur for *all* transition density functions, i.e., for every conceivable starting state. In practice however, this rather strong condition is often violated for starting states with high potential energy, as the associated transition density functions may stay far away from any sensible candidate TM for all times. The probability to encounter these states in the canonical ensemble is however exponentially low, and thus should not contribute significantly to the shape of the RC. Indeed, the numerical methods built around parametrizing the TM are able to successfully deal with this problem by heuristically ignoring sparse outliers by tuning the manifold learning algorithm [4,5].

Still, a rigorous argument for why those outliers can be safely ignored was lacking so far, a gap that the present article aims to fill. In short, we show that the distance to the TM does not need to be uniformly low for all transition density functions, but that the distance is permitted to scale with the potential energy of the starting state. The RC received by parametrizing the TM is then of the same quality as in the uniform distance case. This extension to the TM theory will therefore allow to measure the quality of given RCs, and the numerical computation of ideal RCs in systems that been previously deemed unsuitable for the theory.

This paper is structured as follows: Section 2 reviews the time scale-based definition of good RCs. Section 3 presents the main contribution of this article, weakened but sufficient conditions for the existence of good RCs. In Sect. 4 we give an example of a metastable toy system that fulfills the relaxed but not the original reducibility condition, and demonstrate how the new characterization can improve the quality of error bounds for the dominant timescales. In Sect. 5, concluding remarks and an outlook on future work are given.

## 2    Good Reaction Coordinates

Before introducing the (generalized) transition manifold framework, we first revisit the fundamental time scale-based definition of good reaction coordinates.

### 2.1    Timescales of Molecular Dynamics

We consider a time- and space-continuous, reversible and ergodic Markov process $\mathbf{X}_t$ on a state space $\mathbb{X} \subset \mathbb{R}^n$. In a molecular dynamical system consisting of $N$ atoms, $\mathbb{X}$ often is the Euclidean space describing the three-dimensional positions of all atoms, i.e., $\mathbb{X} = \mathbb{R}^{3N}$ (or $\mathbb{X} = \mathbb{R}^{6N}$ if the atom's momenta are also included). In this case, $\mathbf{X}_t$ is typically described by a thermostated Hamiltonian dynamics or Langevin dynamics.

$\mathbf{X}_t$ is fully characterized by its stochastic transition functions $p^t(x, \cdot) : \mathbb{X} \to \mathbb{R}^+$, or, equivalently, by its family of *transfer operators* $\mathcal{T}^t : L^1_\mu \to L^1_\mu$, $t \geq 0$,

$$\mathcal{T}^t u(x) = \int_{\mathbb{X}} \frac{\rho(x')}{\rho(x)} p^t(x', x) u(x') \, dx'.$$

Here, $\rho$ is the system's (positive) stationary density, which is unique due to the ergodicity of $\mathbf{X}_t$, and $\mu$ is the associated invariant measure. Operating on $L_\mu^1$, $\mathcal{T}^t$ can be understood as the evolution operator of densities with respect to $\mu$ under the dynamics.

On $L_\mu^1$, $\mathcal{T}^t$ is a linear Markov operator, [21, Chap. 3], and in particular non-expansive. Hence, no eigenvalue of $\mathcal{T}^t$ has absolute value greater than 1. Due to the uniqueness of the stationary density, the eigenvalue $\lambda_0^t := 1$ is single; the associated unique eigenfunction is $\varphi_0 \equiv 1$.

Furthermore, $\mathcal{T}^t$ is well-defined as an operator $\mathcal{T}^t : L_\mu^p \to L_\mu^p$ for any $1 \le p \le \infty$ [2]. We understand $\mathcal{T}^t$ as an operator on $L_\mu^2$ from now on, where we will be able to exploit the additional Hilbert space structure. In particular, $\mathcal{T}^t$ is self-adjoint with respect to the inner product on $L_\mu^2$ [33], hence its point spectrum is real and therefore confined to the interval $(-1, 1]$. Note that $\mathcal{T}^t$ cannot possess the eigenvalue $-1$, as this would imply the existence of an eigenfunction $\widetilde{\varphi}_0 \neq \varphi_0$ of $\mathcal{T}^{2t}$ to eigenvalue 1. This however contradicts the uniqueness of $\varphi_0$ as the only eigenfunction to eigenvalue 1 of $\mathcal{T}^t$ for all $t$.

In the following we will always order the eigenvalues so that

$$1 = \lambda_0^t > \lambda_1^t \ge \lambda_2^t \ge \cdots .$$

The associated eigenfunctions $\varphi_i$ of $\mathcal{T}^t$ form an orthonormal basis of $L_\mu^2$. Hence, on $L_\mu^2$, $\mathcal{T}^t$ admits the decomposition

$$\mathcal{T}^t = \sum_{i=0}^\infty \lambda_i^t \langle \varphi_i, \cdot \rangle_\mu \, \varphi_i,$$

which lets us examine the behavior of $\mathbf{X}_t$ on different time scales. The $i$-th *relaxation rate*, i.e., the exponential rate with which the $i$-th eigenfunction $\varphi_i$ of $\mathcal{T}^t$ decays, is given by

$$\sigma_i = -\log(\lambda_i^t)/t, \quad i = 0, 1, 2, \ldots, \tag{1}$$

independent of $t$. These rates, as well as their inverse, the *relaxation time scales* $t_i = 1/\sigma_i, \ i = 0, 1, 2, \ldots$, measure the influence of the different $\varphi_i$ on the long time density transport under $\mathcal{T}^t$, and hence are central quantities of the system.

## 2.2  Reaction Coordinates

A reaction coordinate (RC) now is a continuous map $\xi : \mathbb{X} \to \mathbb{Y} \subset \mathbb{R}^r$, where typically $r \ll n$. Note that the term "reaction coordinate" does not imply that $\xi$ describes a reaction of some sort, it simply is a continuous map. For $y \in \mathbb{Y}$, let $\Sigma_\xi(y)$ be the $y$-level set of $\xi$, i.e.,

$$\Sigma_\xi(y) = \big\{ x \in \mathbb{X} \mid \xi(x) = y \big\}.$$

Following [22], we now define the *coordinate projection operator* $\Pi_\xi : L_\mu^1 \to L_\mu^1$ for a RC $\xi$ by

$$
\begin{aligned}
\big(\Pi_\xi u\big)(x) &= \int_{\Sigma_\xi(\xi(x))} u(x')d\mu_{\xi(x)}(x') \\
&= \frac{1}{\Gamma\big(\xi(x)\big)} \int_{\Sigma_\xi(\xi(x))} u(x')\rho(x') \det\big(\nabla\xi(x')^\mathsf{T}\nabla\xi(x')\big)^{-1/2} d\sigma_{\xi(x)}(x'),
\end{aligned}
$$

where $\Gamma(y)$ is a normalization constant given by

$$
\Gamma(y) = \int_{\Sigma_\xi(y)} \rho(x') \det\big(\nabla\xi(x')^\mathsf{T}\nabla\xi(x')\big)^{-1/2} d\sigma_y(x'),
$$

and $d\sigma_y$ denotes the surface measure on $\Sigma_\xi(y)$. $\mu_y$ can be understood as the invariant measure $\mu$ conditioned on $\Sigma_\xi(y)$, and formally is induced by the density

$$
\rho_y = \frac{\rho}{\Gamma(y)}\big[\det\big(\nabla\xi^\mathsf{T}\nabla\xi\big)\big]^{-1/2}.
$$

As $L_\mu^2 \subset L_\mu^1$ due to Hölder's inequality, $\Pi_\xi$ is defined on $L_\mu^2$ as well. Informally, $\Pi_\xi$ has the effect of averaging an input function $u$ over each level set $\Sigma_\xi(y)$ with respect to $\mu_y$.

It has been shown in [6] that $\Pi_\xi$ is indeed a projection operator. Moreover, $\Pi_\xi$ is equivalent to the Zwanzig projection operator, described in detail in [17], although the latter is typically constructed so that its image are functions over the reduced space $\mathbb{Y}$. For our presentation, however, it is advantageous to define $\Pi_\xi$ to project onto a true subspace of $L_\mu^2$ (namely the subspace of functions that are constant on each $\Sigma_\xi(y),\ y \in \mathbb{Y}$).

The *effective transfer operator* $\mathcal{T}_\xi^t : L_\mu^2 \to L_\mu^2$ associated with the RC $\xi$ is now defined by

$$
\mathcal{T}_\xi^t = \Pi_\xi \circ \mathcal{T}^t \circ \Pi_\xi.
$$

Originally considered in [42], $\mathcal{T}_\xi^t$ has been shown to again be self-adjoint and bounded in $L_\mu^2$-norm by 1 [6]. Hence, the eigenvalues $\lambda_{\xi,i}^t,\ i = 0, 1, 2, \dots$ of $\mathcal{T}_\xi^t$ are also confined to the interval $[-1, 1]$.

## 2.3   Preservation of Time Scales

Our characterization of *good* RCs — originally proposed in [6] — now revolves around the central assumption that the relevant part of the dynamics (the "reaction") occurs on the slowest time scales of $\mathbf{X}_t$. Moreover, we assume that the time scales of the reaction are well-separated from non-reactive time scales, i.e., $t_0 > t_1 \geq \cdots \geq t_d \gg t_{d+1}$ for some $d \in \mathbb{N}$. This is a sensible and commonly made assumption [26,31,32,34], as it holds true for many difference classes of chemical and molecular reactions. However, there are relevant molecular systems whose effective behavior cannot be explained by its slowest timescales alone [25,40],

and hence valid criticism of the general equivalence of the slow with the relevant time scales exist. Nevertheless, we assume that the reaction in question is associated with the $d$ dominant time scales.

With the goal of preserving the dominant time scales under projection onto the RC, and the close connection between those time scales and the dominant transfer operator eigenvalues (1), we use the following definition of good RCs:

**Definition 1 (Good reaction coordinates [6]).** Let $\lambda_i^t$, $i = 0, 1, 2, \ldots$ and $\lambda_{\xi,i}^t$, $i = 0, 1, 2, \ldots$ denote the eigenvalues of $\mathcal{T}^t$ and $\mathcal{T}_\xi^t$, respectively. Let $t_d$ be the last time scale of the system that is relevant to the reaction. Let $\varepsilon > 0$.

An RC $\xi : \mathbb{X} \to \mathbb{Y}$ is called a $\varepsilon$-good RC, if for all $t > 0$ holds

$$|\lambda_i^t - \lambda_{\xi,i}^t| \leq \varepsilon, \quad i = 0, 1, \ldots, d. \tag{2}$$

Informally, we will call $\xi$ a *good RC* if it is $\varepsilon$-good for small $\varepsilon$.

Alternatively, the following sufficient condition characterizes good RC by the projection error of the dominant eigenfunctions under $\Pi_\xi$:

**Theorem 1 ([6], Corollary 3.6).** *Let $(\lambda_i^t, \varphi_i)$, $i = 1, 2, \ldots$ denote the eigenpairs of $\mathcal{T}^t$. For any given $i$, if*

$$\|\Pi_\xi \varphi_i - \varphi_i\|_{L_\mu^2} \leq \varepsilon,$$

*then there is an eigenvalue $\lambda_{\xi,i}^t$ of $\mathcal{T}_\xi^t$ such that*

$$\left|\lambda_i^t - \lambda_{\xi,i}^t\right| \leq \frac{\varepsilon}{\sqrt{1 - \varepsilon^2}}.$$

*Remark 1.* By the above theorem, choosing the $d$ dominant eigenfunctions as the $d$ components of $\xi$ results in a "perfect" RC. However, this approach may lead to redundancy if the $\varphi_i$, $i = 1, \ldots, d$ are strongly correlated and can be parametrized by a common, lower-dimensional $\xi$. For example, a system with $d$ metastable sets along a common, one-dimensional transition pathway would possess $d$ dominant eigenfunctions, but a one-dimensional good RC that parametrizes the transition pathway (see [6, Sect. 5.2] for a detailed example).

Using eigenfunctions as RCs was also promoted by Froyland et al. [14, 15], for the special case where the timescale separation stems from a pointwise local separation of the dynamics into a slow and a fast part. Just like for the transition manifold approach presented in Sect. 3, the short-time equilibration of the dynamics again plays an important part, but unlike in our approach it is assumed to take place on certain "fast fibers" of state space. The transition manifold framework can therefore be considered a generalization of the approach of Froyland et al.

# 3   Weak Reducibility of Stochastic Systems

Definition (2) is not constructive, in that it allow one to check the quality of
a given RC, but does not indicate how to find a good RC algorithmically. To
this end, we will now derive a reducibility condition that binds the existence of
good RCs to the existence of a certain low-dimensional structure in the space of
transition density functions. This structure, called the *transition manifold*, can
be interpreted as the backbone of the essential dynamics, can be visualized, and
ultimately can be used to numerically compute good RCs.

## 3.1   Condition for Good Reaction Coordinates Based on Transfer Operator Eigenfunctions

It was shown in [6] that if for some functions $\hat\varphi_i : \mathbb{Y} \to \mathbb{R}$ the condition

$$\|\varphi_i - \hat\varphi_i \circ \xi\|_\infty \leq \varepsilon, \quad i = 0, 1, \ldots, d \tag{3}$$

holds, then $\xi$ is a $\frac{\varepsilon}{\sqrt{1-\varepsilon^2}}$-good RC by Theorem 1. In other words, if the dominant
eigenfunctions are pointwise almost constant along the level sets of $\xi$, then $\xi$ is
a good RC.

   It turns out, however, that condition (3) is unnecessarily strong. To be pre-
cise, the pointwise approximation implied by the $\|\cdot\|_\infty$-norm can be replaced
by the following weaker condition. This was already observed previously [6,
Remark 4.3], but has not been proven formally.

**Theorem 2.** *Assume that for an RC $\xi : \mathbb{X} \to \mathbb{Y}$ and some functions $\hat\varphi_i : \mathbb{Y} \to \mathbb{R}$, $i = 0, 1, \ldots, d$ holds*

$$\int_{\Sigma_\xi(y)} \left|\varphi_i(x') - \hat\varphi_i(y)\right| d\mu_y(x') \leq \varepsilon \tag{4}$$

*for all level sets $\Sigma_\xi(y)$ of $\xi$. Then*

$$\|\Pi_\xi \varphi_i - \varphi_i\|_{L^2_\mu} \leq 2\varepsilon.$$

*Remark 2.* In words, for a specific value $y \in \mathbb{Y}$, the dominant eigenfunctions $\varphi_i$
do not need to be almost constant everywhere on $\Sigma_\xi(y)$, but only the average
deviation of $\varphi_i$ from some value $\hat\varphi(y)$ along $\Sigma_\xi(y)$, weighted by $\mu_y$, must be
small. Hence, $\xi$ may be a good RC even if $\varphi_i(x')$ substantially deviates from
the value $\hat\varphi(y)$, as long as it is in regions where the measure $\mu_y$ is small. These
are precisely the regions of state space that are lowly-populated in the canonical
ensemble, and thus are statistically irrelevant.

*Proof (**Proof of Theorem** 2).* The projection error is

$$\|\Pi_\xi \varphi_i - \varphi_i\|_{L^2_\mu} \leq \|\Pi_\xi \varphi_i - (\hat\varphi_i \circ \xi)\|_{L^2_\mu} + \|(\hat\varphi_i \circ \xi) - \varphi_i\|_{L^2_\mu}.$$

For the first summand, consider

$$\big(\Pi_\xi \varphi_i\big)(x) = \int_{\Sigma_\xi(\xi(x))} \varphi_i(x') d\mu_{\xi(x)}(x')$$

$$= \int_{\Sigma_\xi(\xi(x))} \Big(\hat\varphi_i\big(\underbrace{\xi(x')}_{=\xi(x)}\big) + \varphi_i(x') - \hat\varphi_i\big(\xi(x')\big)\Big) d\mu_{\xi(x)}(x')$$

$$= \hat\varphi_i\big(\xi(x)\big) + \int_{\Sigma_\xi(\xi(x))} \Big(\varphi_i(x') - \hat\varphi_i\big(\xi(x')\big)\Big) d\mu_{\xi(x)}(x'),$$

and hence

$$\|\Pi_\xi \varphi_i - (\hat\varphi_i \circ \xi)\|_{L^2_\mu}^2 \le \int_{\mathbb{X}} \Big(\underbrace{\int_{\Sigma_\xi(\xi(x))} \big|\varphi_i(x') - \hat\varphi_i\big(\xi(x')\big)\big| d\mu_{\xi(x)}(x')}_{\le \varepsilon}\Big)^2 d\mu(x)$$

$$\le \varepsilon^2 \int_{\mathbb{X}} d\mu(x) = \varepsilon^2.$$

For the second summand, we get with the co-area formula [13]

$$\|(\hat\varphi_i \circ \xi) - \varphi_i\|_{L^2_\mu}^2 = \int_{\mathbb{Y}} \int_{\Sigma_\xi(y)} \big|\hat\varphi_i\big(\xi(x')\big) - \varphi_i(x')\big|^2 d\mu_y(x')\Gamma(y)\,dy$$

$$\le \int_{\mathbb{Y}} \Big(\underbrace{\int_{\Sigma_\xi(y)} \big|\hat\varphi_i\big(\xi(x')\big) - \varphi_i(x')\big| d\mu_y(x')}_{\le \varepsilon}\Big)^2 \Gamma(y)\,dy$$

$$\le \varepsilon^2 \int_{\mathbb{Y}} \Gamma(y)\,dy = \varepsilon^2.$$

### 3.2   Weak Reducibility and Weak Transition Manifolds

From the abstract condition (4) of good RCs, one can now derive a constructive condition for the existence of a good RC. We will also repeat the strong version of this condition, based on (3), which was originally derived in [6].

The parametrizations of certain manifolds will play a central role in our constructions. Specifically, we consider the special class of manifolds $\mathbb{M} \subset L^1$ for which a compact and connected set $\mathbb{Y} \subset \mathbb{R}^r$, as well as a homeomorphism $\mathcal{E} : \mathbb{M} \to \mathbb{Y}$ exists, such that

$$\mathbb{M} = \mathcal{E}^{-1}(\mathbb{Y}). \tag{5}$$

$\mathbb{Y}$ will later become the image space of our constructed RC.

For a fixed lag time $\tau > 0$, we now call the set of functions

$$\widetilde{\mathbb{M}} = \big\{p^\tau(x, \cdot) \mid x \in \mathbb{X}\big\} \subset L^1$$

the *fuzzy transition manifold*. Note that $\widetilde{\mathbb{M}}$ is not a manifold; the reason behind the choice of name will however soon become clear. Now, for any manifold $\mathbb{M} \subset \widehat{\mathbb{M}}$ of form (5), define the projection onto $\mathbb{M}$ by

$$\mathcal{Q} : \mathbb{X} \to \mathbb{M}, \quad x \mapsto \arg\min_{f \in \mathbb{M}} \|f - p^\tau(x, \cdot)\|_{L^2_{1/\mu}}. \tag{6}$$

**Definition 2.** We call the system *strongly $(\varepsilon, r, \tau)$-reducible*, if there exists a manifold $\mathbb{M} \subset \widetilde{\mathbb{M}}$ of form (5) so that for all $x \in \mathbb{X}$

$$\left\| \mathcal{Q}(x) - p^\tau(x, \cdot) \right\|_{L^2_{1/\mu}} \leq \varepsilon. \tag{7}$$

We call any such $\mathbb{M}$ a *strong transition manifold.*

 We call the system *weakly $(\varepsilon, r, \tau)$-reducible*, if there exists a manifold $\mathbb{M} \subset \widetilde{\mathbb{M}}$ of form (5) so that for all $x \in \mathbb{X}$

$$\int_{\Sigma_{\mathcal{Q}}(\mathcal{Q}(x))} \left\| \mathcal{Q}(x') - p^\tau(x', \cdot) \right\|_{L^2_{1/\mu}} d\mu_{\mathcal{Q}(x)}(x') \leq \varepsilon, \tag{8}$$

where $\Sigma_{\mathcal{Q}}(f)$ is the $f$-level set of $\mathcal{Q}$. We call any such $\mathbb{M}$ a *weak transition manifold.*

*Example 1.* As an illustration of the core idea behind the TM construction, we give a simple example of a metastable system with a strong TM, originally published in [5].

 Consider a two-dimensional system described by the overdamped Langevin equation

$$d\mathbf{X}_t = -\nabla V(\mathbf{X}_t)\, dt + \sqrt{2\beta^{-1}} d\mathbf{W}_t, \tag{9}$$

where $V$ is the potential energy function and $\mathbf{W}_t$ is a Wiener diffusion process scaled by the inverse temperature $\beta \in \mathbb{R}^+$. Now suppose that $V$ possesses two local energy wells, connected by a linear, one-dimensional transition path, such as in Fig. 1(left). The "reaction" in this system is the rare transition from one well to the other. Hence, an intuitively good RC is the horizontal coordinate of a point, $\xi(x) = x_1$, as it describes the progress of $x$ along the transition pathway.
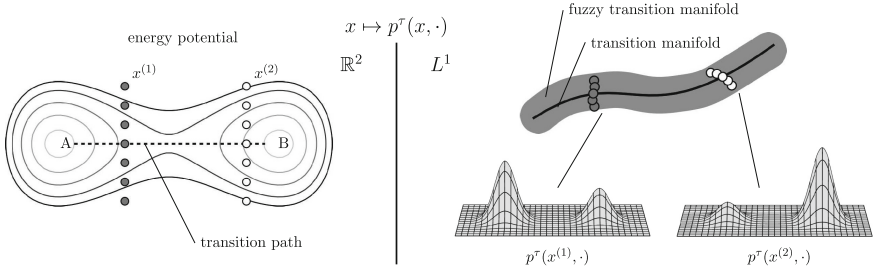
 The key insight now is that, if the lag time $\tau$ was chosen long enough for a typical trajectory to move to one of the metastable sets, then the transition densities $p^\tau(x, \cdot) \in L^1$ also essentially depend only on the progress of $x$ along the transition path. The reason is that the $p^\tau(x, \cdot)$ are essentially convex combinations of two Gaussians[1] centered in the energy minima $A$ and $B$,

$$p^\tau(x, \cdot) \approx c(x)\rho_A(\cdot) + (1 - c(x))\rho_B(\cdot)$$

with the convex factor $c(x)$ determined by the progress of the starting point $x$ along the transition path. This is represented in Fig. 1(right) by the fact that the transition densities for each gray and white starting point, respectively, concentrate around one point each in $L^1$. Hence, overall, the fuzzy TM $\widetilde{\mathbb{M}}$ concentrates around a one-dimensional manifold in $L^1$. This manifold is therefore a strong TM.

 An example of a system with only a weak TM will be discussed in detail in Sect. 4.

---

[1] To be precise, the $p^\tau(x, \cdot)$ are approximately convex combinations of the quasi-stationary densities [16] of the metastable sets, that here however resemble Gaussians.

**Fig. 1.** Illustration of the transition manifold concept for metastable systems. Left: energy potential of a two-dimensional metastable system. Right: Sketch of the (fuzzy) TM for this system. Starting points $x$ with the same progress along the transition path get mapped to approximately the same density under the map $x \mapsto p^\tau(x, \cdot)$. Geometrically, this means that the fuzzy TM concentrates around a one-dimensional manifold in $L^1$.

*Remark 3.* Note that we slightly deviate from the original definition of the transition manifold in [6] by requiring that $\mathbb{M} \subset \widetilde{\mathbb{M}}$ instead of only $\mathbb{M} \subset L^1$. Also note that $\mathcal{Q}$ is now defined on $\mathbb{X}$ and not on $\widetilde{\mathbb{M}}$ as originally in [6]. The interpretation of $\mathcal{Q}$ as "closest point projection onto $\mathbb{M}$" is still valid, however.

Condition (7) indicates whether the fuzzy TM $\widetilde{\mathbb{M}}$ clusters $\varepsilon$-closely around an actual manifold $\mathbb{M}$ with respect to the $L^2_{1/\rho}$-norm. Again, condition (8) represents a relaxation of this condition, as the integral introduces a weighting with respect to $d\mu_{\mathcal{Q}(x)}$. Informally speaking, for points $x'$ with $\rho(x') = \mathcal{O}(\varepsilon)$, a distance $\left\| \mathcal{Q}(x') - p^\tau(x', \cdot) \right\|_{L^2_{1/\mu}} = \mathcal{O}(1)$ is now permitted without violating the reducibility condition.

It was shown in [6] that strongly reducible systems possess good RCs. The following theorem now shows that weakly reducible systems still possess good RCs. It characterizes $\mathcal{Q}$ as a good "$\mathbb{M}$-valued RC" (cf. (4)):

**Theorem 3.** *Let the system be weakly $(\varepsilon, r, \tau)$-reducible. Then for each eigenpair $(\lambda_i^\tau, \varphi_i)$ of the transfer operator $\mathcal{T}^\tau$ there exists a map $\tilde{\varphi}_i : \mathbb{M} \to \mathbb{R}$ so that for all $x \in \mathbb{X}$*

$$\int_{\Sigma_{\mathcal{Q}(\mathcal{Q}(x))}} \big|\varphi_i(x') - \tilde{\varphi}_i\big(\mathcal{Q}(x')\big)\big| \, d\mu_{\mathcal{Q}(x)}(x') \leq \frac{\varepsilon}{|\lambda_i^\tau|}.$$

*Proof.* As $\mathbb{M} \subset \widetilde{\mathbb{M}}$, for $x \in \mathbb{X}$ we can choose $q(x) \in \mathbb{X}$ so that $\mathcal{Q}(x) = p^t\big(q(x), \cdot\big)$. Let $\tilde{\varphi}_i : \mathbb{M} \to \mathbb{R}$ be defined by

$$\tilde{\varphi}_i\big(\mathcal{Q}(x)\big) = \varphi_i\big(q(x)\big).$$

Then

$$\int_{\Sigma_{\mathcal{Q}(\mathcal{Q}(x))}} \big|\varphi_i(x') - \tilde{\varphi}_i\big(\mathcal{Q}(x')\big)\big| \, d\mu_{\mathcal{Q}(x)}(x') = \int_{\Sigma_{\mathcal{Q}(\mathcal{Q}(x))}} \big|\varphi_i(x') - \tilde{\varphi}_i\big(\mathcal{Q}(x)\big)\big| \, d\mu_{\mathcal{Q}(x)}(x')$$

$$= \int_{\Sigma_{\mathcal{Q}(\mathcal{Q}(x))}} \big|\varphi_i(x') - \varphi_i\big(q(x)\big)\big| \, d\mu_{\mathcal{Q}(x)}(x') =: (\star)$$

As the system is reversible, the detailed balance condition $\rho(x)p^\tau(x,x'') = \rho(x'')p^\tau(x'',x)$ holds. Hence, the eigenfunctions $\varphi_i$ of $\mathcal{T}^\tau$ have the property

$$\lambda_i^\tau \varphi_i = \mathcal{T}^\tau \varphi_i = \int_{\mathbb{X}} \frac{\rho(x'')}{\rho(x)} p^\tau(x'',\cdot)\varphi_i(x'') \, dx'' = \int_{\mathbb{X}} \varphi_i(x'')p^\tau(\cdot,x'') \, dx'',$$

and thus

$$(\star) = \int_{\Sigma_{\mathcal{Q}}(\mathcal{Q}(x))} \frac{1}{|\lambda_i^\tau|} \left| \int_{\mathbb{X}} \varphi_i(x'') \Big( p^\tau(x',x'') - p^\tau\big(q(x),x''\big) \Big) dx'' \right| d\mu_{\mathcal{Q}(x)}(x').$$

Swapping integrals gives

$$(\star) \leq \frac{1}{|\lambda_i^\tau|} \int_{\mathbb{X}} |\varphi_i(x'')| \int_{\Sigma_{\mathcal{Q}}(\mathcal{Q}(x))} \Big| p^\tau(x',x'') - p^\tau\big(q(x),x''\big) \Big| d\mu_{\mathcal{Q}(x)}(x') \, dx'',$$

and with Hölder's inequality, $\|fg\|_{L^1} \leq \|f\|_{L^2_\mu} \|g\|_{L^2_{1/\mu}}$, we get

$$\leq \frac{1}{|\lambda_i^\tau|} \underbrace{\|\varphi_i\|_{L^2_\mu}}_{=1} \left\| \int_{\Sigma_{\mathcal{Q}}(\mathcal{Q}(x))} \Big| p^\tau(x',\cdot) - p^\tau\big(q(x),\cdot\big) \Big| d\mu_{\mathcal{Q}(x)}(x') \right\|_{L^2_{1/\mu}}.$$

Applying triangle inequality and using $p^\tau\big(q(x),\cdot\big) = \mathcal{Q}(x)$ gives

$$(\star) \leq \frac{1}{|\lambda_i^\tau|} \int_{\Sigma_{\mathcal{Q}}(\mathcal{Q}(x))} \left\| p^t(x',\cdot) - p^t\big(q(x),\cdot\big) \right\|_{L^2_{1/\mu}} d\mu_{\mathcal{Q}(x)}(x')$$

$$= \frac{1}{|\lambda_i^\tau|} \int_{\Sigma_{\mathcal{Q}}(\mathcal{Q}(x))} \left\| p^t(x',\cdot) - \underbrace{\mathcal{Q}(x)}_{=\mathcal{Q}(x')} \right\|_{L^2_{1/\mu}} d\mu_{\mathcal{Q}(x)}(x').$$

By our assumption, this integral is at most $\varepsilon$. Hence,

$$(\star) \leq \frac{\varepsilon}{|\lambda_i^\tau|}.$$

As the last step, we can now construct from $\mathcal{Q}$ an $r$-dimensional RC that meets the condition (2):

**Corollary 1.** *Let the system be weakly $(\varepsilon, r, \tau)$-reducible. Let $\mathcal{E} : \mathbb{M} \to \mathbb{R}^r$ be any parametrization of the transition manifold $\mathbb{M}$. Then for the RC*

$$\xi : \mathbb{X} \to \mathbb{R}^r, \quad x \mapsto \mathcal{E}\big(\mathcal{Q}(x)\big) \tag{10}$$

*and the eigenpairs $(\lambda_i^\tau, \varphi_i)$ of $\mathcal{T}^\tau$ holds*

$$\|\Pi_\xi \varphi_i - \varphi_i\|_{L^2_\mu} \leq \frac{2\varepsilon}{|\lambda_i^\tau|}. \tag{11}$$

*Proof.* Let $\tilde{\varphi}_i : \mathbb{M} \to \mathbb{R}$ as in the proof of Theorem 3, and define $\hat{\varphi}_i : \mathbb{Y} \to \mathbb{R}$ via

$$\hat{\varphi}_i(y) := \tilde{\varphi}_i\big(\mathcal{E}^{-1}(y)\big).$$

Note that for any $x \in \mathbb{X}$ holds $\Sigma_\xi\big(\xi(x)\big) = \Sigma_\mathcal{Q}(\mathcal{Q}(x))$. Thus,

$$\int_{\Sigma_\xi(\xi(x))} \big|\varphi_i(x') - (\hat{\varphi}_i \circ \xi)(x')\big| d\mu_y(x') = \int_{\Sigma_\mathcal{Q}(\mathcal{Q}(x))} \big|\varphi_i(x') - (\tilde{\varphi}_i \circ \mathcal{Q})(x')\big| d\mu_{\mathcal{Q}(x)}(x')$$
$$\leq \frac{\varepsilon}{|\lambda_i^\tau|},$$

where the last inequality is Theorem 3. The assertion now follows from Theorem 2. □

If $(\lambda_i^\tau, \varphi_i)$ is dominant, i.e., $\lambda_i^\tau \approx 1$, then the projection error (11) is small. In that case, $\xi : x \mapsto \mathcal{E}\big(\mathcal{Q}(x)\big)$ is indeed a good RC, by Theorem 1.

*Remark 4.* Any RC of form (10) is called an *ideal RC* [6]. As in practice, however, neither the projection $\mathcal{Q}$ nor the parametrization $\mathcal{E}$ of $\mathbb{M}$ are known, this RC cannot be computed analytically. Instead, for strongly reducible systems, an approximate parametrization of $\mathbb{M}$ is computed by applying manifold learning methods to a finite sample of the fuzzy TM $\widetilde{\mathbb{M}}$ [4–6]. Our ongoing efforts to extend these techniques to the newly-identified weak reducibility condition will be discussed in the outlook in Sect. 5.

## 4    Numerical Example: A Weakly Reducible System

In order to compare the strong and weak reducibility condition, we consider a simple two-dimensional metastable system that possesses a one-dimensional RC. This system, originally considered in [22], is governed by an overdamped Langevin equation of form (9), where the potential energy function $V$ is given by
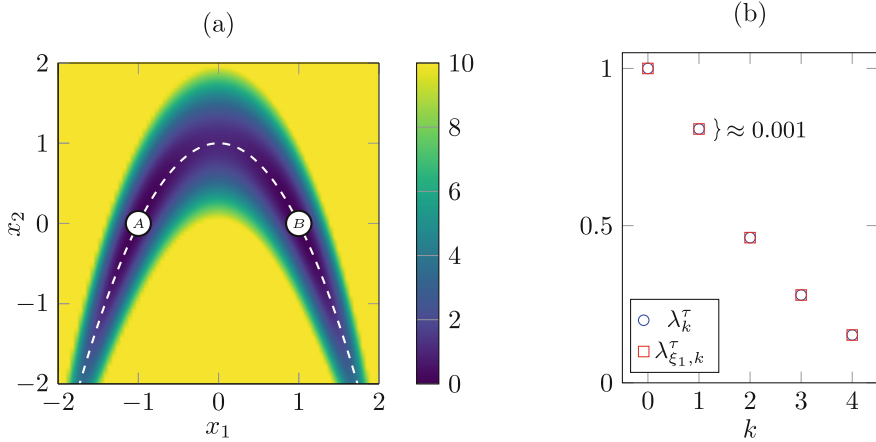
$$V(x) = (x_1^2 - 1)^2 + 10\,(x_1^2 + x_2 - 1)^2.$$

We choose the inverse temperature $\beta = 1$, and consider the system on the domain $\mathbb{X} = [-2, 2] \times [-2, 2]$ (though no boundary conditions have been enforced in the following computations). The potential $V$, depicted in Fig. 2(a), possesses two local minima in the states $A = (-1, 0)$ and $B = (1, 0)$. The reaction in question hence is the transition from the area around one minimum (without loss of generality state $A$) to the other (state $B$). The minimum energy pathway (MEP) [24], which in the zero temperature limit supports almost all reactive trajectories [30], is indicated by the white dashed line.

The spectrum of $\mathcal{T}^\tau$ for $\tau = 0.5$, computed by a Ulam method [39] from a long, equilibrated trajectory of the system, exhibits a spectral gap after $\lambda_1$ (Fig. 2(b)). The relevant reaction, i.e., the transition between the two metastable sets, is associated primarily with the process on the dominant timescale $t_1$.

The (MEP) of the potential is given by the set

$$A_{\mathrm{MEP}} = \{(x_1, x_2) \in \mathbb{X} \mid x_2 = 1 - x_1^2\}.$$

Intuitively, the manifold

**Fig. 2.** (a) Energy potential of a two-dimensional drift-diffusion system. The reaction of interest here is the transition between the two local minima. (b) Eigenvalues of the full transfer operator $\mathcal{T}^\tau$ and of the effective transfer operator $\mathcal{T}^\tau_{\xi_1}$ projected onto the computed RC $\xi_1$.

$$\mathbb{M}_{\text{MEP}} = \{p^\tau(x, \cdot) \mid x \in A_{\text{MEP}}\}$$

should constitute a good TM. This statement should come with a warning: The intuition that the MEP allows to construct a good TM is wrong in general. There are many cases where the relevant transition pathways are completely different from the MEPs of the underlying system, mainly because for finite temperatures all statistically relevant transition paths concentrate in regions not close to the MEP and only converge to the MEP in the limit of zero temperature. In the case considered herein, however, relevant transition paths concentrate around the MEP even for finite temperatures.

Before quantitatively assessing whether or not $\mathbb{M}_{\text{MEP}}$ is indeed is a good TM, we visualize the fuzzy TM of the system, i.e., the set $\widetilde{\mathbb{M}} = \{p^\tau(x, \cdot) \mid x \in \mathbb{X}\}$. As $\widetilde{\mathbb{M}}$ lies in the function space $L^1$, it first needs to be embedded into a (finite-dimensional) Euclidean space. This is done by computing the mean of every $p^\tau(x, \cdot) \in \widetilde{\mathbb{M}}$ via the function $\mathbf{m} : L^1 \to \mathbb{R}^2$,

$$\mathbf{m}(p^\tau(x, \cdot)) := \int_{\mathbb{X}} x' \, p^\tau(x, x') \, dx'. \tag{12}$$

The set $\mathbf{m}(\widetilde{\mathbb{M}})$ then serves as the Euclidean embedding[2] of $\widetilde{\mathbb{M}}$.
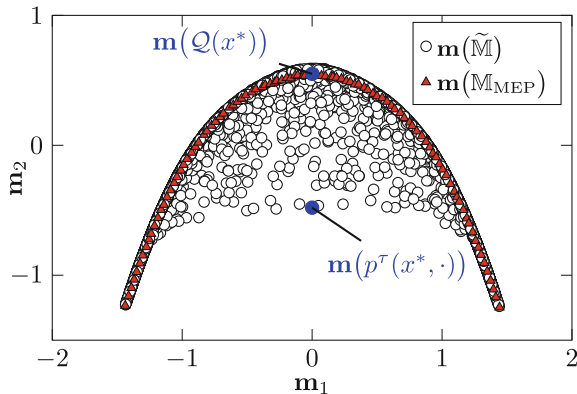
---

[2] While for general dynamics $\mathbf{m}$ is not an embedding of the fuzzy TM in the strict topological sense, we conjecture that in this system, no two transition densities $p^\tau(x_1, \cdot), p^\tau(x_2, \cdot)$ possess the same mean, and hence that $\mathbf{m}$ is homeomorphic on $\widetilde{\mathbb{M}}$ and its image. Still, we neither formally confirm this, nor assess the distortion of $\widetilde{M}$ under $\mathbf{m}$, and hence $\mathbf{m}(\widetilde{\mathbb{M}})$ as a replacement for $\widetilde{M}$ should be handled with care.

Furthermore, as $\mathbf{m}(\widetilde{\mathbb{M}})$ is an infinite set, only a finite subsample can be visualized. For this we draw a large number, specifically $N = 8000$, of starting points $\{x_1, \ldots, x_N\}$ uniformly from $\mathbb{X}$ and for each $x_k$ compute $\mathbf{m}_k := \mathbf{m}(p^\tau(x_k, \cdot))$. Here the integral in (12) is approximated via Monte Carlo quadratur, i.e., for $M \gg 1$,

$$\mathbf{m}(p^\tau(x_k, \cdot)) \approx \frac{1}{M} \sum_{l=1}^{M} z_k^{(l)}, \tag{13}$$

where the $z_k^{(l)}$ are samples of the density $p^\tau(x_k, \cdot)$. These were computed numerically by an Euler-Maruyama integrator of (9), starting in $x_k$, with a different random seed for each $l = 1, \ldots, M$.

The points $\mathbf{m}_k$ are shown in Fig. 3. We observe that most of the $\mathbf{m}_k$ lie close to a parabola-like structure, though there appear to exist systematic outliers, associated with starting points from the high energy regions in the lower part of $\mathbb{X}$. The maximum distance is assumed by the starting point $x^* = (0, -2)$. The parabola is exactly the Euclidean embedding of $\mathbb{M}_{\text{MEP}}$, which is also shown in Fig. 3.



**Fig. 3.** Euclidean embeddings via the mean embedding function $\mathbf{m}$ of the fuzzy TM $\widetilde{\mathbb{M}}$, and the TM based on the minimum energy pathway, $\mathbb{M}_{\text{MEP}}$. Shown are $N = 8000$ sample points of $\mathbf{m}(\widetilde{\mathbb{M}})$, and $N = 100$ sample points of $\mathbf{m}(\mathbb{M}_{\text{MEP}})$. $\mathbf{m}(\widetilde{\mathbb{M}})$ appears to cluster around $\mathbf{m}(\mathbb{M}_{\text{MEP}})$, except for outliers from the high energy regions below the MEP.

However, the outliers prevent $\mathbb{M}_{\text{MEP}}$ from being a good strong TM by Definition 2. To be precise, for the point $x^* = (0, -3)$, we get for the distance in (7)

$$\left\| \mathcal{Q}(x^*) - p^t(x^*, \cdot) \right\|_{L^2_{1/\mu}} \approx 2.5, \tag{14}$$

where again finite samples of $\widetilde{\mathbb{M}}$ and $\mathbb{M}_{\text{MEP}}$, and kernel density estimations of the $p^t(x, \cdot)$ were used in the computation. Using (14) as a lower bound for the

eigenvalue approximation (2) via Theorems 2 and 1 is of course worthless, hence $\mathbb{M}_{\mathrm{MEP}}$ is not a strong TM.

On the other hand, for the defining condition (8) of weak reducibility holds

$$\int_{\Sigma_{\mathcal{Q}}(\mathcal{Q}(x^*))} \big\| \mathcal{Q}(x') - p^\tau(x',\cdot) \big\|_{L^2_{1/\mu}} \, d\mu_{\mathcal{Q}(x^*)}(x') \approx 0.02 \tag{15}$$

for the problematic point $x^*$. Assuming this value is indeed an upper bound for (8), the system is weakly reducible with parameter $\varepsilon = 0.06$, and $\mathbb{M}_{\mathrm{MEP}}$ is the corresponding weak TM. The eigenvalue error for $\lambda_1^\tau$ predicted by Theorems 2 and 1 then is

$$|\lambda_1^\tau - \lambda_{\xi,1}^\tau| \le 0.06, \tag{16}$$

for any RC $\xi$ of the form (10).

To confirm this error bound, we now construct such an RC. For this, a parametrization $\mathcal{E}$ of $\mathbb{M}_{\mathrm{MEP}}$ must be chosen. Any such parametrization is sufficient, for simplicity we choose
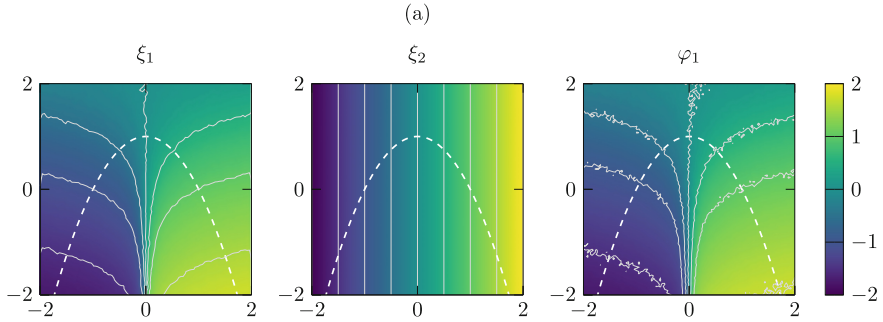
$$\mathcal{E}\big(p^\tau(x,\cdot)\big) := x_1,$$

i.e., the map of $p^\tau(x,\cdot)$ onto the first component $x_1$ of its starting point $x$. Next, the projection $\mathcal{Q}$ of $\widetilde{\mathbb{M}}$ onto the TM $\mathbb{M}_{\mathrm{MEP}}$ is required. In order to avoid the costly calculation of kernel density estimates for the large number of starting points, and to avoid the badly-conditioned scaling by the factor $1/\rho$, we replace the $L^2_{1/\rho}$ distance in (6) by the Euclidean distance between the mean-embedded densities, i.e., utilize

$$\widetilde{\mathcal{Q}}(x) = \underset{f \in \mathbb{M}_{\mathrm{MEP}}}{\arg\min} \big\| \mathbf{m}(f) - \mathbf{m}\big(p^\tau(x,\cdot)\big) \big\|_2.$$

Numerically, this projection is approximated by choosing from the 100 sample points of $\mathbf{m}(\mathbb{M}_{\mathrm{MEP}})$ that are shown in Fig. 3 the point of minimum distance from $\mathbf{m}(p^\tau(x,\cdot))$. The point $\mathbf{m}(p^\tau(x,\cdot))$ is here again computed via (13). While using the projection $\widetilde{\mathcal{Q}}$ instead of $\mathcal{Q}$ might slightly distort the computed RC, it will have a negative impact on the quality of the RC, so if the bound (16) holds for $\widetilde{\mathcal{Q}}$, it will hold for $\mathcal{Q}$ as well. Moreover, it has been shown in [5] that the $L^2_{1/\rho}$ distance is equivalent to the distance in certain embedding spaces.

The final RC is then given by $\xi_1 : x \mapsto \mathcal{E}\big(\widetilde{\mathcal{Q}}(p^\tau(x,\cdot))\big)$. By numerically evaluating $\xi_1$ at the 8000 sample points (where the $p^\tau(x,\cdot)$ are again approximated by finite samples) and interpolating the resulting values bilinearly, we receive a continuous RC on $\mathbb{X}$. Figure 4 shows the level plot of $\xi_1$. We see that the level sets of $\xi_1$ are essentially identical to those of the dominant eigenfunction $\varphi_1$, also shown in Fig. 4. This is not surprising, as $\xi_1$ is constructed to fulfill the requirements of Theorem 1, i.e., the dominant eigenfunctions are required to be almost invariant under averaging over the level sets of $\xi_1$. As there is only one dominant eigenfunction $\varphi_1$, and $\xi_1$ is also one-dimensional, this implies that the level sets of $\xi_1$ and $\varphi_1$ are almost identical. Note however that the precise ranges

**Fig. 4.** Level plots of the RCs $\xi_1$ computed by the TM method, a naively-constructed RC $\xi_2$, as well as the dominant eigenfunction $\varphi_1$ of $\mathcal{T}^\tau$. We see that the level sets of $\xi_1$ and $\varphi_1$ are essentially identical.

of $\xi$ and $\varphi_1$ are not necessarily identical, but strongly depend on the chosen parametrization $\mathcal{E}$.

The effective transfer operator $\mathcal{T}_{\xi_1}^\tau$ associated with $\xi_1$ can again be approximated by an Ulam method. Its leading eigenvalues, shown in Fig. 3(b), approximate the eigenvalues of the full transfer operator $\mathcal{T}^\tau$ very well. In particular, for the second dominant eigenvalue holds

$$|\lambda_1^\tau - \lambda_{\xi,1}^\tau| \approx 0.001.$$

As a consequence, the relaxation rate of the projected system $\xi_1(X_t)$, denoted $\sigma_{\xi,1}$ and computed from $\lambda_{\xi,1}$ via (1), also approximate the rate of the full system $\sigma_{\text{full}}$ very well; we have $\sigma_{\xi_1} \approx 0.43$, $\sigma_{\text{full}} \approx 0.43$ . In contrast, projections onto other, naively chosen RCs, such as

$$\xi_2(x) := x_1,$$

seem to systematically over-estimate the equilibration rate, hence under-estimates the metastability of the system. Specifically, we have $\sigma_{\xi_2} \approx 0.46$. Reduced models built based on $\xi_2$ would therefore run the risk of equilibrating quicker than the full model by artificially increasing the number of transitions.

That said, the difference between $|\sigma_{\xi_1} - \sigma_{\xi_2}| \approx 0.03$ is rather small, so the naive RC $\xi_2$ can already be considered quite good. The reason is that at low temperatures the dynamics concentrates near the MEP, and here for each level set of $\xi_2$ there exists a level set of $\xi_1$ that is close (in the sense that the minimum pairwise point distance is small), and the RCs are both smooth. Still, the difference is measurable, and this causes the discrepancy.

Overall, this example confirms that

(1) the RC $\xi_1$ derived from a parametrization of $\mathbb{M}_{\mathrm{MEP}}$ is good, and
(2) the error bound (16) derived from the characterization of $\mathbb{M}_{\mathrm{MEP}}$ as a weak
    TM is reasonably accurate.

## 5   Conclusion and Outlook

In this work, we derived an improved and generalized characterization of good
reaction coordinates for timescale-separated stochastic processes. We built upon
a recently developed framework that constructs good RCs from parametrizations
of the so-called transition manifold, a potentially low-dimensional manifold in the
space of probability densities. We have shown that the criteria on the underlying
system to possess such a manifold were overly strict, in the sense that certain
systems with demonstrated good reaction coordinates do not possess a transition
manifold by the old definition. We thus provided an alternative, relaxed definition
of the transition manifold that is applicable to a larger class of systems, while
still allowing the construction of good reaction coordinates.

One natural next step would be to implement the novel definition of weak
TMs into a data-driven algorithm for the identification of good RCs. Unlike
in the toy example from Sect. 4, the parametrization of the transition mani-
fold (or of a suitable candidate) is not known analytically in practice. Instead,
an approximate parametrization is identified by applying a nonlinear manifold
learning algorithm to a large sample of $\widetilde{\mathbb{M}}$ (or a suitable embedding thereof) [4].
Many manifold learning algorithms, such as the diffusion maps algorithm [9] can
be tuned to ignore outliers, which can be seen as a heuristic way weighing with
respect to the invariant measure $\mu$. A more rigorous approach however would
be to directly implement the weighted distance (8) into the diffusion maps algo-
rithm. This could be achieved by using the target measure-extension of diffusion
maps [1], which at the same time allows one to estimate the in general unknown
measure $\mu$ from data.

## References

1. Banisch, R., Trstanova, Z., Bittracher, A., Klus, S., Koltai, P.: Diffusion maps
   tailored to arbitrary non-degenerate itô processes. Appl. Comput. Harmonic Anal.
   **48**(1), 242–265 (2020)
2. Baxter, J.R., Rosenthal, J.S.: Rates of convergence for everywhere-positive Markov
   chains. Stat. Probab. Lett. **22**(4), 333–338 (1995)
3. Best, R.B., Hummer, G.: Reaction coordinates and rates from transition paths.
   Proc. Natl. Acad. Sci. **102**(19), 6732–6737 (2005)

4. Bittracher, A., Banisch, R., Schütte, C.: Data-driven computation of molecular reaction coordinates. J. Chem. Phys. **149**(15), 154103 (2018)
5. Bittracher, A., Klus, S., Hamzi, B., Koltai, P., Schütte, C.: Dimensionality reduction of complex metastable systems via kernel embeddings of transition manifolds. arXiv eprint 1904.08622 (2019)
6. Bittracher, A., Koltai, P., Klus, S., Banisch, R., Dellnitz, M., Schütte, C.: Transition manifolds of complex metastable systems: theory and data-driven computation of effective dynamics. J. Nonlinear Sci. **28**(2), 471–512 (2017)
7. Chodera, J.D., Noé, F.: Markov state models of biomolecular conformational dynamics. Curr. Opin. Struct. Biol. **25**, 135–144 (2014)
8. Ciccotti, G., Kapral, R., Vanden-Eijnden, E.: Blue moon sampling, vectorial reaction coordinates, and unbiased constrained dynamics. Chem. Phys. Chem. **6**(9), 1809–1814 (2005)
9. Coifman, R.R., Lafon, S.: Diffusion maps. Appl. Comput. Harmonic Anal. **21**(1), 5–30 (2006)
10. Daldrop, J.O., Kappler, J., Brünig, F.N., Netz, R.R.: Butane dihedral angle dynamics in water is dominated by internal friction. Proc. Natl. Acad. Sci. **115**(20), 5169–5174 (2018)
11. Dellnitz, M., Junge, O.: On the approximation of complicated dynamical behavior. SIAM J. Numer. Anal. **36**(2), 491–515 (1999)
12. Elber, R., Bello-Rivas, J.M., Ma, P., Cardenas, A.E., Fathizadeh, A.: Calculating Iso-committor surfaces as optimal reaction coordinates with milestoning. Entropy **19**(5) (2017)
13. Evans, L.C., Gariepy, R.F.: Measure Theory and Fine Properties of Functions. Chapman and Hall/CRC, New York (2015)
14. Froyland, G., Gottwald, G., Hammerlindl, A.: A computational method to extract macroscopic variables and their dynamics in multiscale systems. SIAM J. Appl. Dyn. Syst. **13**(4), 1816–1846 (2014)
15. Froyland, G., Gottwald, G.A., Hammerlindl, A.: A trajectory-free framework for analysing multiscale systems. Physica D **328**, 34–43 (2016)
16. Gesùa, G.D., Lelièvre, T., Peutreca, D.L., Nectouxa, B.: Jump markov models and transition state theory: the quasi-stationary distribution approach. Faraday Discuss. **195**, 469–495 (2016)
17. Givon, D., Kupferman, R., Stuart, A.: Extracting macroscopic dynamics: model problems and algorithms. Nonlinearity **17**(6), R55–R127 (2004)
18. Kappler, J., Daldrop, J.O., Bruenig, F.N., Boehle, M.D., Netz, R.R.: Memory-induced acceleration and slowdown of barrier crossing. J. Chem. Phys. **148**, 014903 (2018)
19. Klus, S., Nüske, F., Koltai, P., Wu, H., Kevrekidis, I., Schütte, C., Noé, F.: Data-driven model reduction and transfer operator approximation. J. Nonlinear Sci. **28**, 985–1010 (2018)
20. Laio, A., Gervasio, F.L.: Metadynamics: a method to simulate rare events and reconstruct the free energy in biophysics, chemistry and material science. Rep. Prog. Phys. **71**(12), 126601 (2008)
21. Lasota, A., Mackey, M.C.: Chaos, Fractals, and Noise: Stochastic Aspects of Dynamics, vol. 97. Springer, Berlin (2013)
22. Legoll, F., Lelièvre, T.: Effective dynamics using conditional expectations. Nonlinearity **23**(9), 2131 (2010)
23. Ma, A., Dinner, A.R.: Automatic method for identifying reaction coordinates in complex systems. J. Phys. Chem. B **109**(14), 6769–6779 (2005)

24. Maragliano, L., Fischer, A., Vanden-Eijnden, E., Ciccotti, G.: String method in collective variables: minimum free energy paths and isocommittor surfaces. J. Phys. Chem. **125**(2), 024106 (2006)
25. McGibbon, R.T., Husic, B.E., Pande, V.S.: Identification of simple reaction coordinates from complex dynamics. J. Phys. Chem. **146**(4), 44109 (2017)
26. Noé, F., Nüske, F.: A variational approach to modeling slow processes in stochastic dynamical systems. Multiscale Model. Simul. **11**(2), 635–655 (2013)
27. Pavliotis, G.A., Stuart, A.M.: Multiscale Methods: Averaging and Homogenization. Springer, New York (2008)
28. Peters, B.: Reaction coordinates and mechanistic hypothesis tests. Annu. Rev. Phys. Chem. **67**(1), 669–690 (2016). PMID: 27090846
29. Peters, B., Trout, B.L.: Obtaining reaction coordinates by likelihood maximization. J. Chem. Phys. **125**(5), 054108 (2006)
30. Ren, W.: Higher order string method for finding minimum energy paths. Commun. Math. Sci. **1**(2), 377–384 (2003)
31. Sarich, M., Noé, F., Schütte, C.: On the approximation quality of Markov state models. Multiscale Model. Simul. **8**(4), 1154–1177 (2010)
32. Sarich, M., Schütte, C.: Approximating selected non-dominant timescales by Markov state models. Commun. Math. Sci. **10**(3), 1001–1013 (2012)
33. Schervish, M.J., Carlin, B.P.: On the convergence of successive substitution sampling. J. Comput. Graph. Stat. **1**(2), 111–127 (1992)
34. Schütte, C., Fischer, A., Huisinga, W., Deuflhard, P.: A direct approach to conformational dynamics based on hybrid Monte Carlo. J. Comput. Phys. **151**(1), 146–168 (1999)
35. Sengupta, U., Carballo-Pacheco, M., Strodel, B.: Automated markov state models for molecular dynamics simulations of aggregation and self-assembly. J. Chem. Phys. **150**(11), 115101 (2019)
36. Sirur, A., De Sancho, D., Best, R.B.: Markov state models of protein misfolding. J. Chem. Phys. **144**(7), 075101 (2016)
37. Smith, P.E.: The alanine dipeptide free energy surface in solution. J. Chem. Phys. **111**(12), 5568–5579 (1999)
38. Torrie, G.M., Valleau, J.P.: Nonphysical sampling distributions in Monte Carlo free-energy estimation: umbrella sampling. J. Comput. Phys. **23**(2), 187–199 (1977)
39. Ulam, S.: A Collection of Mathematical Problems. Interscience Tracts in Pure and Applied Mathemtics, vol. 8. Interscience Publishers, New York (1960)
40. Wedemeyer, W.J., Welker, E., Scheraga, H.A.: Proline cis-trans isomerization and protein folding. Biochemistry **41**(50), 14637–14644 (2002)
41. Williams, M.O., Kevrekidis, I.G., Rowley, C.W.: A data-driven approximation of the koopman operator: extending dynamic mode decomposition. J. Nonlinear Sci. **25**(6), 1307–1346 (2015)
42. Zhang, W., Hartmann, C., Schütte, C.: Effective dynamics along given reaction coordinates, and reaction rate theory. Faraday Discuss. **195**, 365–394 (2016)
43. Zhang, W., Schuette, C.: Reliable approximation of long relaxation timescales in molecular dynamics. Entropy **19**(7), 367 (2017)
44. Zwanzig, R.: Memory effects in irreversible thermodynamics. Phys. Rev. **124**, 983–992 (1961)