

Transition manifolds of complex metastable systems

Theory and data-driven computation of effective dynamics

Andreas Bittracher¹, Péter Koltai¹, Stefan Klus¹, Ralf Banisch¹,
Michael Dellnitz², and Christof Schütte^{1,3}

¹Department of Mathematics and Computer Science, Freie Universität Berlin, Germany

²Department of Mathematics, Paderborn University, Germany

³Zuse Institute Berlin, Germany

Abstract

We consider complex dynamical systems showing metastable behavior but no local separation of fast and slow time scales. The article raises the question of whether such systems exhibit a low-dimensional manifold supporting its effective dynamics. For answering this question, we aim at finding nonlinear coordinates, called reaction coordinates, such that the projection of the dynamics onto these coordinates preserves the dominant time scales of the dynamics. We show that, based on a specific reducibility property, the existence of good low-dimensional reaction coordinates preserving the dominant time scales is guaranteed. Based on this theoretical framework, we develop and test a novel numerical approach for computing good reaction coordinates. The proposed algorithmic approach is fully local and thus not prone to the curse of dimension with respect to the state space of the dynamics. Hence, it is a promising method for data-based model reduction of complex dynamical systems such as molecular dynamics.

1. Introduction

With the advancement of computing power, we are able to simulate and analyze more and more complicated and high-dimensional models of dynamical systems, ranging from astronomical scales for the simulation of galaxies, over planetary and continental scales for climate and weather prediction, down to molecular and sub-atomistic scales via, e.g., Molecular Dynamics (MD) simulations aimed at gaining insight into complex biological processes. Particular aspects of such processes, however, can often be described by much simpler means than the full process, thus *reducing* the full dynamics to some *essential* behavior or *effective dynamics* in terms of some essential observables of the system.

Extracting these observables and the related effective dynamics from a dynamical system, though, is one of the most challenging problems in computational modeling [31].

One prominent example of dynamical reduction is arguably given by a variety of multiscale systems with explicit fast-slow time scale separation, mostly singularly perturbed systems, where either the fast component is considered in a quasi-stationary regime (i.e. the slow components are fixed and assumed not to change for the observation period), or the effective behavior of the fast components is injected into the slow processes, e.g. by averaging or homogenization [60]. Much of the recent attention has been directed to the case where the deduction of the slow (or fast) effective dynamics is not possible by purely analytic means, due to the lack of an analytic description of the system, or because the complexity of the system renders this task unfeasible [31, 32, 12, 21, 56, 71, 14, 79, 36]. However, all of these approaches still depend on some local form of time scale separation between the “fast” and the “slow” components of the dynamics.

The focus of this work is on specific multiscale systems *without* local dynamical slow-fast time scale separation, but for which a reduction to an effective dynamical behavior supported on some low-dimensional manifold is still possible. The dynamical property lying at the heart of our approach is that there is a time scale separation in the *global kinetic* behavior of the process, as opposed to the aforementioned slow-fast behavior encoded in the *local dynamics*. Here, global kinetic behavior means that the multiple scales show up if we consider the *Fokker–Planck equation* associated with the dynamics, say $\dot{u} = \mathcal{L}u$, where the Fokker–Planck operator \mathcal{L} will have several small eigenvalues, while the rest of its spectrum is significantly larger. Such dynamical systems exhibit *metastable* behavior and the slow time scales are the time scales of statistical relaxation between the main metastable sets, while there is no time scale gap for the local dynamics within each of the metastable regions [6, 70].

Global time scale separation induced by metastability has been analyzed for deterministic [18] and stochastic dynamical systems [68, 33] for more than a decade. A typical trajectory of a metastable dynamical system will spend most time within the metastable sets, while rare transitions between these sets happen as sudden “jumps” roughly along low-dimensional *transition pathways* that connect the metastable sets [16, 58, 26]. For an example, see Figure 1.

The tool to describe the global kinetic behavior of a metastable system is the so-called *transfer operator* (the evolution operator of the Fokker–Planck equation), which acts on functions on the state space. The time scale separation we rely on here implies a spectral gap for this operator. This fact has been exploited to find low-dimensional representations of the global kinetics in form of Markov chains whose (discrete) states represent the metastable sets while the transition probabilities between the states approximate the jump statistics between the sets on long time scales. Under the name “Markov State Models” (MSM), this approach has led to a variety of methods [7, 70] with broad application, e.g., in molecular dynamics, cf. [68, 59, 69, 10]. This reduction comes with a price: Since the relaxation kinetics is described just by jumps between the metastable sets in a (finite) discrete state space, any information about the transition process and its dynamical features is lost. A variety of approaches have been developed for complementing the MSM approach appropriately [53], but a continuous (in time and

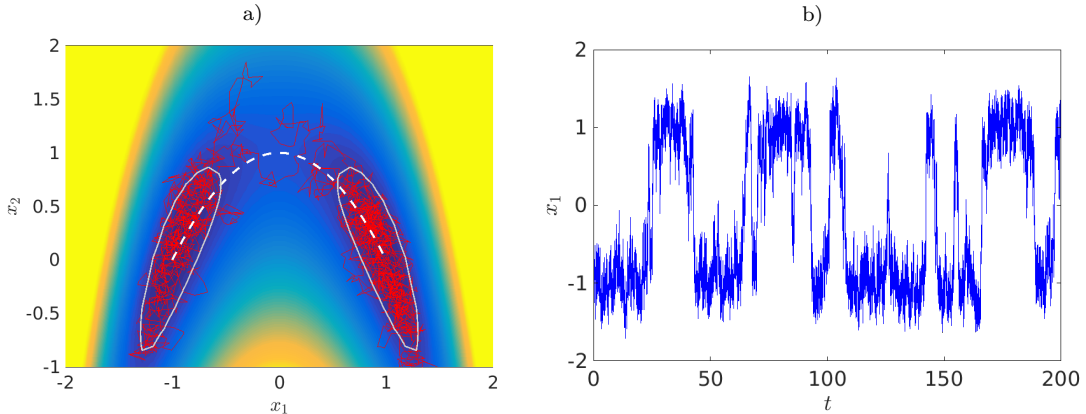


Figure 1: a) Curved double-well potential with two metastable sets (areas encircled by light grey lines) around the global minima $(-1, 0)$ and $(1, 0)$. In a typical trajectory (red line), transitions between the metastable sets are rare events and generally happen along the transition path (white dashed line). b) The x_1 -component of a longer trajectory that shows multiple rare transitions (or events).

space) low-dimensional effective description based on MSMs allowing to understand the transition mechanism is infeasible.

In another branch of the literature, again heavily influenced by molecular dynamics applications, model reduction techniques have been developed that assume the existence of a low-dimensional *reaction coordinate* or *order parameter* in order to construct an effective dynamics or kinetics: Examples are free energy based techniques [75, 42], trajectory-based sampling techniques [27, 2, 54, 61], methods based on diffusive processes [4, 84, 60], and many more that rely on the assumption that the reaction coordinates are known. The problem of actually *constructing* good reaction coordinates remains an area of ongoing research [46], to which this paper contributes. Typically, reaction coordinates are either postulated using system specific expert knowledge [8, 72], an approximation to the dominant eigenfunctions of the transfer operator is sought [70, 10, 61], or machine learning techniques are proposed [48]. Froyland et al. [31] show that these eigenfunctions are indeed optimal — in the sense of optimally representing the slow dynamics — but for high dimensional systems computational reaction coordinate identification still is often infeasible. In the context of *transition path theory* [78], the committor function is known to be an ideal [47] reaction coordinate. In [62], the authors construct a level set of the committor using support vector machines, but the computation of reaction coordinates is infeasible for high-dimensional systems. The main problem in computing reaction coordinates for high-dimensional metastable systems results from the fact that all of these algorithms try to solve a *global* problem in the entire state space that cannot be decomposed easily into purely local computations.

In this article, we elaborate on the definition, existence and algorithmic identification of reaction coordinates for metastable systems: We define reaction coordinates as a small set of *nonlinear* coordinates on which a reduced system [45, 84] can be defined

having the same dominant time scales (in terms of transfer operator eigenvalues) as the original system. We then consider a low-dimensional state space on which the reduced dynamics is a Markov process. Thus, our approach utilizes concepts and transfer operator theory developed previously, but in our case the projected transfer operator is still *infinite-dimensional*, in stark contrast to its reduction to a stochastic matrix in the MSM approach.

The contribution of this paper is twofold. First, we develop a conceptual framework that identifies *good* reaction coordinates as the ones that parameterize a low-dimensional *transition manifold* \mathbb{M} in the function space L^1 , which is the natural state space of the Fokker–Planck equation $\dot{u} = \mathcal{L}u$ associated with the dynamics. The property which defines \mathbb{M} is that, on moderate time scales $t_{\text{fast}} < t \ll t_{\text{slow}}$, the *transition density functions* of the dynamics concentrate around \mathbb{M} . We provide evidence that such an \mathbb{M} indeed exists due to metastability and the existence of transition pathways. Crucially, the dimension of \mathbb{M} is often lower than the number of dominant eigenfunctions.

Second, we present an algorithm to construct approximate reaction coordinates. Our algorithm is data-driven and *fully local*, thus circumventing the main problem of previously proposed algorithms: In order to compute the value of the desired reaction coordinate ξ at a location x in the state space \mathbb{X} , only the ability to simulate short trajectories initialized at x is needed. In particular, we assume no a priori knowledge of metastable sets, no global equilibration, and we do not need to resolve the slow time scales numerically. The algorithm is built on two pillars:

1. The simulation time scale t can be chosen a lot smaller than the dominant time scales t_{slow} of the system, such that it is feasible to simulate many short trajectories of length t .
2. We utilize *embedding techniques* inspired by the seminal work of Whitney [83] and the recent work [19] that allows one to take almost any mapping into a Euclidean space of more than twice the dimension of the manifold \mathbb{M} and to obtain a one-to-one image of it.

These two pillars together with the low-dimensionality of \mathbb{M} imply that we can represent the image of the reaction coordinate in a space with moderate (finite) dimension. Then, we can use established *manifold learning* techniques [56, 12, 71] to obtain a parametrization of the manifold in the embedding space and pull this parametrization back to the original state space, hence obtaining a reaction coordinate.

The locality of the algorithm also implies that reaction coordinates are only computed in the region of state space where sampled points are available. This is a common issue with manifold learning algorithms; here it manifests as the transition manifold being reliably learned only in regions we have good sampling coverage of. However, recently several methods have appeared in the literature that allow a fast exploration of the state space. These methods do not provide equilibrium sampling, but instead try to rapidly cover the essential part of the state space with sampling points. This can be achieved with enhanced sampling methods such as Umbrella Sampling [41, 76], Metadynamics [43, 44], Blue-Moon sampling [11], Adaptive Biasing Force method [15],

or Temperature-Accelerated Molecular Dynamics [49], as well as trajectory-based techniques like Milestoning [28], Transition Interface Sampling [55], or Forward Flux Sampling [3]. Alternatively, several techniques like the equation-free approach [36], the heterogeneous multiscale method (HMM) [23] and methods based on diffusion maps [9] have been developed to utilize short unbiased MD trajectories for extracting information that allows much larger timesteps. This can be combined with reaction coordinate based effective dynamics [84, 85].

In principle, the method we present in this article may be combined with any enhanced sampling technique in order to generate sampling points that cover a large part of the state space. For simplicity, we will use long MD trajectories to generate our sampling points, but we do not require that the points are distributed according to an equilibrium distribution.

The paper is organized as follows: Section 2 introduces transfer operators, which describe the global kinetics of the stochastic process. Based on these transfer operators, we define metastability, i.e. the existence of dominant time scales. In Section 3, we describe the model reduction techniques *Markov state modeling* and *coordinate projection* that are designed to capture the dominant time scales of metastable systems. Furthermore, we characterize *good* reaction coordinates. In the first part of Section 4, we show that our dynamical assumption ensures the existence of good reaction coordinates, then in the second part we describe our approach to compute them. Several numerical examples are given in Section 5. Concluding remarks and an outlook are provided in Section 6.

2. Transfer operators and their properties

As mentioned in the introduction, global properties of dynamical systems such as metastable sets or a partitioning into fast and slow subprocesses can be obtained using transfer operators associated with the system and their eigenfunctions. In this section, we will introduce different transfer operators needed for our considerations.

2.1. Transfer operators

In what follows, $P[\cdot \mid \mathfrak{E}]$ denotes probabilities conditioned on the event \mathfrak{E} and $E[\cdot \mid \mathfrak{E}]$ the expectation value. Furthermore, $\{\mathbf{X}_t\}_{t \geq 0}$ is a stochastic process defined on a state space $\mathbb{X} \subset \mathbb{R}^n$.

Definition 2.1 (Transition density function). *Let \mathbb{A} be any measurable set, then the transition density function $p^t: \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}_{\geq 0}$ of a time-homogeneous stochastic process $\{\mathbf{X}_t\}_{t \geq 0}$ is defined by*

$$P[\mathbf{X}_t \in \mathbb{A} \mid \mathbf{X}_0 = x] = \int_{\mathbb{A}} p^t(x, y) dy.$$

That is, $p^t(x, y)$ is the conditional probability density of $\mathbf{X}_t = y$ given that $\mathbf{X}_0 = x$.

With the aid of the transition density function, we can now define transfer operators. Note, however, that the transition density is in general not known explicitly and needs

to be estimated from simulation data. In what follows, we assume that there is a *unique* equilibrium density ϱ that is invariant under $\{\mathbf{X}_t\}_{t \geq 0}$, that is, it satisfies

$$\varrho(x) = \int_{\mathbb{X}} p^t(y, x) \varrho(y) \, dy,$$

a.e. on \mathbb{X} . Let μ denote the associated invariant measure $d\mu = \varrho \, dx$.

Definition 2.2 (Transfer operators). *Let $p \in L^1(\mathbb{X})$ be a probability density¹, $u = p/\varrho \in L^1_\mu(\mathbb{X})$ be a probability density with respect to the equilibrium density ϱ , and $f \in L^\infty(\mathbb{X})$ an observable of the system. For a given lag time t :*

(a) *The Perron–Frobenius operator $\mathcal{P}^t: L^1(\mathbb{X}) \rightarrow L^1(\mathbb{X})$ is defined by the unique linear extension of*

$$\mathcal{P}^t p(x) = \int_{\mathbb{X}} p^t(y, x) p(y) \, dy$$

to $L^1(\mathbb{X})$.

(b) *The Perron–Frobenius operator $\mathcal{T}^t: L^1_\mu(\mathbb{X}) \rightarrow L^1_\mu(\mathbb{X})$ with respect to the equilibrium density is defined by the unique linear extension of*

$$\mathcal{T}^t u(x) = \int_{\mathbb{X}} \frac{\varrho(y)}{\varrho(x)} p^t(y, x) u(y) \, dy$$

to $L^1_\mu(\mathbb{X})$.

(c) *The Koopman operator $\mathcal{K}^t: L^\infty(\mathbb{X}) \rightarrow L^\infty(\mathbb{X})$ is defined by*

$$\mathcal{K}^t f(x) = \int_{\mathbb{X}} p^t(x, y) f(y) \, dy = \mathbb{E}[f(\mathbf{X}_t) \mid \mathbf{X}_0 = x]. \quad (1)$$

All these are well-defined non-expanding operators on the respective spaces.

The equilibrium density ϱ satisfies $\mathcal{P}^t \varrho = \varrho$, that is, ϱ is an eigenfunction of \mathcal{P}^t with associated eigenvalue $\lambda_0 = 1$. The definition of \mathcal{T}^t relies on ϱ , we have $\varrho \mathcal{T}^t u = \mathcal{P}^t(u\varrho)$.

Instead of their natural domains from Definition 2.2, all our transfer operators are considered on the following Hilbert spaces: $\mathcal{P}^t: L^2_{1/\mu}(\mathbb{X}) \rightarrow L^2_{1/\mu}(\mathbb{X})$, $\mathcal{T}^t: L^2_\mu(\mathbb{X}) \rightarrow L^2_\mu(\mathbb{X})$, and $\mathcal{K}^t: L^2_\mu(\mathbb{X}) \rightarrow L^2_\mu(\mathbb{X})$. They are still well-defined non-expansive operators on these spaces [1, 67, 38].

Furthermore, we will need the notion of reversibility for our considerations. Reversibility means that the process is statistically indistinguishable from its time-reversed counterpart.

¹We denote by L^q the space (equivalence class) of q -integrable functions with respect to the Lebesgue measure. L^q_ν denotes the same space of function, now integrable with respect to the measure ν .

Definition 2.3 (Reversibility). *A system is said to be reversible if the detailed balance condition*

$$\varrho(x) p^t(x, y) = \varrho(y) p^t(y, x)$$

is satisfied for all $x, y \in \mathbb{X}$.

In what follows, we will assume that the system is reversible.

One prominent example for a class of SDEs satisfying uniqueness of the equilibrium density and reversibility is given by

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

Here, V is called the potential, β is the non-dimensionalized inverse temperature, and \mathbf{W}_t is a standard Wiener process. The process generated by (2) is ergodic and thus admits a unique positive equilibrium density, given by $\varrho(x) = \exp(-\beta V(x))/Z$, under mild growth conditions on the potential V [50, 51]. Note that the subsequent considerations hold for all stochastic processes that satisfy reversibility and ergodicity with respect to a unique positive invariant measure and are *not* limited to the class of dynamical systems given by (2). See [70] for a discussion of a variety of stochastic dynamical systems that have been considered in this context.

As a result of the detailed balance condition, the Koopman operator \mathcal{K}^t and the Perron–Frobenius operator with respect to the equilibrium density \mathcal{T}^t become identical and we obtain

$$\langle \mathcal{P}^t f, g \rangle_{1/\mu} = \langle f, \mathcal{P}^t g \rangle_{1/\mu} \quad \text{and} \quad \langle \mathcal{T}^t f, g \rangle_{\mu} = \langle f, \mathcal{T}^t g \rangle_{\mu},$$

i.e. all the transfer operators become self-adjoint on the respective Hilbert spaces from above. Here $\langle \cdot, \cdot \rangle_{\mu}$ and $\langle \cdot, \cdot \rangle_{1/\mu}$ denote the natural scalar products on the weighted spaces L_{μ}^2 and $L_{1/\mu}^2$, respectively.

2.2. Spectral decomposition

Due to the self-adjointness, the eigenvalues λ_i^t of \mathcal{P}^t and \mathcal{T}^t are real-valued and the eigenfunctions form an orthogonal basis with respect to $\langle \cdot, \cdot \rangle_{1/\mu}$ and $\langle \cdot, \cdot \rangle_{\mu}$, respectively. In what follows, we assume that the spectrum of \mathcal{T}^t is purely discrete given by (infinitely many) isolated eigenvalues. This assumption is made for the sake of simplicity. It is actually not required for the rest of our considerations; it would be sufficient to assume that the spectral radius R of the essential spectrum of \mathcal{T}^t is strictly smaller than 1, and some isolated eigenvalues of modulus larger than R exist. It has been shown that this condition is satisfied for a large class of metastable dynamical systems, see [70, Sec. 5.3] for details. For example, the process generated by (2) has purely discrete spectrum under mild growth and regularity assumptions on the potential V .

Under this condition, ergodicity implies that the dominant eigenvalue λ_0 is the only eigenvalue with absolute value 1 and we can thus order the eigenvalues so that

$$1 = \lambda_0^t > \lambda_1^t \geq \lambda_2^t \geq \dots$$

The eigenfunction of \mathcal{T}^t corresponding to $\lambda_0 = 1$ is the constant function $\varphi_0 = \mathbb{1}_{\mathbb{X}}$. Let φ_i be the normalized eigenfunctions of \mathcal{T}^t , i.e. $\langle \varphi_i, \varphi_j \rangle_{\mu} = \delta_{ij}$, then any function $f \in L^2_{\mu}(\mathbb{X})$ can be written in terms of the eigenfunctions as $f = \sum_{i=0}^{\infty} \langle f, \varphi_i \rangle_{\mu} \varphi_i$. Applying \mathcal{T}^t thus results in

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

For more details, we refer to [38] and references therein.

2.3. Implied time scales

For some $d \in \mathbb{N}$, we call the $d + 1$ dominant eigenvalues $\lambda_0^t, \dots, \lambda_d^t$ of \mathcal{T}^t the *dominant spectrum* of \mathcal{T}^t , i.e.

$$\sigma_{\text{dom}}(\mathcal{T}^t) := \{\lambda_0^t, \dots, \lambda_d^t\}.$$

Usually, d is chosen in such a way that there is a *spectral gap* after λ_d^t , i.e. $1 - \lambda_d^t \ll \lambda_d^t - \lambda_{d+1}^t$. The (*implied*) *time scales* on which the associated dominant eigenfunctions decay are given by

$$t_i = -t / \log(\lambda_i^t). \quad (3)$$

If \mathcal{T}^t is a semigroup of operators, then there are $\kappa_i \leq 0$ with $\lambda_i^t = \exp(\kappa_i t)$ such that $t_i = -\kappa_i^{-1}$ holds. Assuming there is a spectral gap, the dominant time scales satisfy $t_1 \geq \dots \geq t_d \gg t_{d+1}$. These are the time scales of the *slow* dynamical processes, also called *rare events*, which are of primary interest in applications. The other, *fast* processes are regarded as fluctuations around the relative equilibria (or *metastable states*) between which the relevant slow processes travel.

3. Projected transfer operators and reaction coordinates

The purpose of dimension reduction in molecular dynamics is to find a reduced dynamical model that captures the dominant time scales of the system correctly while keeping the model as simple as possible. In this section, we will introduce two different projections and the corresponding projected transfer operators. The goal is to find suitable projections onto the slow processes.

3.1. Galerkin projections and Markov state models

One frequently used approach to obtain a reduced model is *Markov state modeling*. The goal is to find a model that is as simple as possible and yet correctly reproduces the dominant time scales. Given a fixed $t > 0$, most authors [58, 59] refer to a Markov state model (MSM) as a matrix $T^t \in \mathbb{R}^{(d+1) \times (d+1)}$ such that

$$\sigma_{\text{dom}}(\mathcal{T}^t) \approx \sigma_{\text{dom}}(T^t), \quad (4)$$

and it has been studied in detail under which condition this can be achieved [20, 65].

There are different ways of constructing an MSM, maybe the most intuitive one is also the simplest: Let the entries of T^t be the transition rates between metastable sets. A typical molecular system with d dominant time scales will have $d + 1$ metastable sets $\mathbb{C}_1, \dots, \mathbb{C}_{d+1}$ (also called *cores*) and its dynamics is characterized by transitions between these sets and fluctuations inside the sets (see Figure 1 for an illustration). Since the fluctuations are on faster time scales, we neglect them by setting [68]

$$T_{\text{core},ij}^t = \mathbb{P}_\mu[\mathbf{X}_t \in \mathbb{C}_j \mid \mathbf{X}_0 \in \mathbb{C}_i], \quad (5)$$

where \mathbb{P}_μ denotes the probability measure conditioned to the initial condition \mathbf{X}_0 being distributed according to μ . Thus, $T_{\text{core},ij}^t$ is the probability that the process in equilibrium jumps to the metastable set \mathbb{C}_j in time t , given that it started in the metastable set \mathbb{C}_i . Note that (5) can be equivalently rewritten as

$$T_{\text{core},ij}^t = \frac{\langle \mathcal{T}^t \mathbb{1}_{\mathbb{C}_i}, \mathbb{1}_{\mathbb{C}_j} \rangle_\mu}{\langle \mathbb{1}_{\mathbb{C}_i}, \mathbb{1}_{\mathbb{C}_i} \rangle_\mu}, \quad (6)$$

where $\mathbb{1}_{\mathbb{C}_i}$ is the characteristic function of the set \mathbb{C}_i .

Equation (6) readily suggests that T_{core}^t is a projection of the transfer operator \mathcal{T}^t , namely its *Galerkin projection* onto the space spanned by the characteristic functions $\mathbb{1}_{\mathbb{C}_1}, \dots, \mathbb{1}_{\mathbb{C}_{d+1}}$ [68].

Definition 3.1 (Galerkin projection). *Given a set of basis functions $\psi_1, \dots, \psi_m \in L_\mu^2(\mathbb{X})$, let $\mathbb{V} := \text{span}\{\psi_1, \dots, \psi_m\}$ and $\psi := (\psi_1, \dots, \psi_m)^\top$. The projection to \mathbb{V} or, equivalently, to ψ , $\Pi_\psi = \Pi_\psi: L_\mu^2(\mathbb{X}) \rightarrow \mathbb{V}$ is defined as*

$$\langle \Pi_\psi f - f, g \rangle_\mu = 0 \quad \forall f \in L_\mu^2(\mathbb{X}), \forall g \in \mathbb{V}.$$

The residual projection is given by $\Pi_\psi^\perp = \text{Id} - \Pi_\psi$, where Id is the identity. The Galerkin projection of \mathcal{T}^t to \mathbb{V} is given by the linear operator $T^t: \mathbb{V} \rightarrow \mathbb{V}$ satisfying

$$\langle \mathcal{T}^t f - T^t f, g \rangle_\mu = 0 \quad \forall f, g \in \mathbb{V}.$$

Equivalently, $T^t = \Pi_\psi \mathcal{T}^t$. We also denote the extension of T^t to the whole $L_\mu^2(\mathbb{X})$, given by $\Pi_\psi \mathcal{T}^t \Pi_\psi$, by T^t . Furthermore, we denote the matrix representation of T^t with respect to the basis (ψ_0, \dots, ψ_d) by T^t as well. Either it will be clear from the context which of the objects T^t is meant or it will not matter; e.g., the dominant spectrum is the same for all of them.

We see that T_{core} is the matrix representation of the Galerkin projection with respect to the basis functions $\langle \mathbb{1}_{\mathbb{C}_i}, \mathbb{1}_{\mathbb{C}_i} \rangle_\mu^{-1} \mathbb{1}_{\mathbb{C}_i}$, $i = 1, \dots, d + 1$. More general MSMs can be built by Galerkin projections of the transfer operator to spaces spanned by other — not necessarily piecewise constant — basis functions [80, 69, 82, 37, 38, 61, 57]. However, in some of these methods, one also often loses the interpretation of the entries of the matrix T^t as probabilities.

Ultimately, the best MSM in terms of approximation quality in (4) is given by the Galerkin projection of \mathcal{T}^t onto the space spanned by its dominant eigenfunctions φ_0, \dots ,

φ_d . This space is invariant under \mathcal{T}^t since $\mathcal{T}^t \varphi_i = \lambda_i^t \varphi_i$ and the dominant eigenvalues (and hence the time scales) are the same for the MSM and for \mathcal{T}^t . Due to the curse of dimensionality, however, the computation of the eigenfunctions φ_i is in general infeasible for high-dimensional problems.

Remark 3.2. There are quantitative results assessing the error in (4) of the MSM in terms of the projection errors $\|\Pi_\psi^\perp \varphi_i\|_{L_\mu^2}$, $i = 0, \dots, d$, cf. [70, Section 5.3]. One can obtain a weaker, but similar result from our Lemma 3.5 in the next section.

3.2. Coordinate projections and effective transfer operators

While the MSMs from above successfully reproduce the dominant time scales of the original system, they often discard all other information about the system, such as the transition paths between metastable sets. Minimal coordinates that describe these transitions are called *reaction coordinates* and reducing the dynamics onto these coordinates yields *effective dynamics* [45, 84]. The goal of the previous section — namely to retain the dominant time scales of the original dynamics in a reduced model — can now be reformulated for this lower-dimensional effective dynamics or, equivalently, for its (effective) transfer operator.

Let $\xi: \mathbb{X} \rightarrow \mathbb{R}^k$ be a C^1 function, where $k \leq n$. Let $\mathbb{L}_z = \{x \in \mathbb{X} \mid \xi(x) = z\}$ be the z -level set of ξ . The so-called *coarea formula* [29, Section 3.2], which can be considered as a nonlinear variant of Fubini's theorem, splits integrals over \mathbb{X} into consecutive integrals over level sets of ξ and then over the range of ξ . For $f \in L_\mu^2(\mathbb{X})$, we have²

$$\int_{\mathbb{X}} f(x) d\mu(x) = \int_{\xi(\mathbb{X})} \int_{\mathbb{L}_z} f(x') \varrho(x') \det(\nabla \xi(x')^\top \nabla \xi(x'))^{-1/2} d\sigma_z(x') dz, \quad (7)$$

where $z = \xi(x)$ and σ_z is the surface measure on \mathbb{L}_z . The *coordinate projection*, defined next, averages a given function along the level sets of a coordinate function ξ .

Definition 3.3 (Coordinate projection). *For $f \in L_\mu^2(\mathbb{X})$, we define*

$$P_\xi f(x) = \int_{\mathbb{L}_z} f(x') d\mu_z(x') \quad (8)$$

$$= \frac{1}{\Gamma(z)} \int_{\mathbb{L}_z} f(x') \varrho(x') \det(\nabla \xi(x')^\top \nabla \xi(x'))^{-1/2} d\sigma_z(x'), \quad (9)$$

where μ_z is a probability measure on \mathbb{L}_z with density $\frac{\varrho}{\Gamma(z)} \det(\nabla \xi^\top \nabla \xi)^{-1/2}$ with respect to σ_z . Here, $\Gamma(z)$ is just the normalization constant so that μ_z becomes a probability measure. The residual projection is given by $P_\xi^\perp = \text{Id} - P_\xi$.

To get a better feeling for the action of P_ξ , note that $P_\xi f(x)$ is the expectation of $f(\mathbf{x}')$ with respect to μ conditional to $\xi(\mathbf{x}') = \xi(x)$, i.e.

$$P_\xi f(x) = \mathbf{E}_\mu [f(\mathbf{x}') \mid \xi(\mathbf{x}') = \xi(x)].$$

²The coarea formula holds for L^1 functions, but $L_\mu^2 \subset L_\mu^1$, since μ is a probability measure (i.e., it is finite).

Or, in other words, μ_z is the marginal of μ conditional to $\xi(x) = z$. Note, in particular, that $P_\xi f$ is itself a function on \mathbb{X} , but it is constant on the level sets of ξ , and thus let us set $\widehat{P_\xi f}(\xi(x)) = P_\xi f(x)$ for $x \in \mathbb{L}_{\xi(x)}$. It follows from the coarea formula (7) and (9) that

$$\int_{\mathbb{X}} f(x) d\mu(x) = \int_{\xi(\mathbb{X})} \Gamma(z) \widehat{P_\xi f}(z) dz. \quad (10)$$

Next, we state some properties of the coordinate projection.

Proposition 3.4. *The coordinate projection has the following properties.*

- (a) P_ξ is a linear projection, i.e. $P_\xi^2 = P_\xi$.
- (b) P_ξ is self-adjoint with respect to $\langle \cdot, \cdot \rangle_\mu$.
- (c) $P_\xi: L_\mu^2(\mathbb{X}) \rightarrow L_\mu^2(\mathbb{X})$ is orthogonal, hence non-expansive, i.e. $\|P_\xi f\|_{L_\mu^2} \leq \|f\|_{L_\mu^2}$.

Proof. See Appendix A. □

We use the coordinate projection to describe the dynamics-induced propagation of reduced distributions with respect to the variable ξ . To this end, we define the *effective transfer operator* $\mathcal{T}_\xi^t: L_\mu^2(\mathbb{X}) \rightarrow L_\mu^2(\mathbb{X})$ by

$$\mathcal{T}_\xi^t = P_\xi \mathcal{T}^t P_\xi. \quad (11)$$

We immediately obtain from the self-adjointness of \mathcal{T}^t (see Section 2) and Proposition 3.4 (b) that \mathcal{T}_ξ^t is a self-adjoint operator on $L_\mu^2(\mathbb{X})$. Moreover, $\|\mathcal{T}^t\|_{L_\mu^2} \leq 1$ and Proposition 3.4 (c) imply that $\|\mathcal{T}_\xi^t\|_{L_\mu^2} \leq 1$. Thus, the spectrum of the effective transfer operator lies in the interval $[-1, 1]$, too.

Returning to the purpose of these constructions, we call ξ a *good reaction coordinate* if

$$\sigma_{\text{dom}}(\mathcal{T}^t) \approx \sigma_{\text{dom}}(\mathcal{T}_\xi^t). \quad (12)$$

While the previously introduced Markov state model T^t obtained by the Galerkin projection was approximating the dominant spectrum of the original transfer operator by a finite-dimensional operator (i.e. a matrix), the effective transfer operator still acts on an infinite-dimensional space. The reduction lies in the fact that \mathcal{T}^t operates on functions over $\mathbb{X} \subseteq \mathbb{R}^n$, but the effective transfer operator \mathcal{T}_ξ^t operates *essentially* on functions over $\xi(\mathbb{X}) \subset \mathbb{R}^k$, although we embed those into \mathbb{X} through the level sets of ξ .

As mentioned above, a Galerkin projection of the transfer operator onto its dominant eigenfunctions is a perfect MSM. In the same vein, we ask here how we can characterize a good reaction coordinate. We can make use of the following general result.

Lemma 3.5. *Let \mathbb{H} be a Hilbert space with scalar product $\langle \cdot, \cdot \rangle$, and associated norm $\|\cdot\|$, let $Q: \mathbb{H} \rightarrow \mathbb{H}$ be some orthogonal projection on a linear subspace of \mathbb{H} , with $Q^\perp = \text{Id} - Q$. Let $T: \mathbb{H} \rightarrow \mathbb{H}$ be a self-adjoint non-expansive linear operator, and u with $\|u\| = 1$ its eigenvector, i.e., $Tu = \lambda u$ for some $\lambda \in \mathbb{R}$. If $\|Q^\perp u\| < \varepsilon$, then $T_Q := QTQ$ has an eigenvalue $\lambda_Q \in \mathbb{R}$ with $|\lambda - \lambda_Q| < \varepsilon/\sqrt{1 - \varepsilon^2}$.*

Proof. Using $Q = \text{Id} - Q^\perp$, we have

$$T_Q Qu = QT \underbrace{QQ}_{=Q} u = QTu - \underbrace{QTQ^\perp}_{=-\zeta} u = \lambda Qu + \zeta,$$

where $\|\zeta\| \leq \|Q^\perp u\| < \varepsilon$ since Q and T are non-expanding. Thus, $u' := Qu/\|Qu\|$ satisfies $T_Q u' = \lambda u' + \zeta/\|Qu\|$, and the orthogonality of Q gives $\|Qu\| > \sqrt{1 - \varepsilon^2}$. Now, any orthogonal projection is self-adjoint, as is shown in the proof of Proposition 3.4, hence the operator QTQ is self-adjoint, too, and thus normal. From the theory of pseudospectra for normal operators [77, Theorems 2.1, 2.2, and §4], we know that if $\|T_Q u' - \lambda u'\| < \varepsilon/\sqrt{1 - \varepsilon^2}$, then T_Q has an eigenvalue $\lambda_Q \in \mathbb{R}$ in the $\varepsilon/\sqrt{1 - \varepsilon^2}$ -neighborhood of λ . \square

With $\mathbb{H} = L_\mu^2$, $Q = P_\xi$, and $T = \mathcal{T}^t$ we immediately obtain the following result.

Corollary 3.6. *As before, let λ_i^t and φ_i , $i = 0, \dots, d$, denote the dominant eigenvalues and eigenfunctions of \mathcal{T}^t , respectively. For any given i , if $\|P_\xi^\perp \varphi_i\|_{L_\mu^2} < \varepsilon$, then there is an eigenvalue $\tilde{\lambda}_i^t$ of \mathcal{T}_ξ^t with $|\lambda_i^t - \tilde{\lambda}_i^t| < \varepsilon/\sqrt{1 - \varepsilon^2}$.*

Corollary 3.6 implies that if the projection error of *all dominant eigenfunctions* is small, then ξ is a good reaction coordinate in the sense of (12). Very similar results are available for approximation of the eigenvalues of the infinitesimal generator of the Fokker–Planck equation associated with the transfer operator if the dynamical system under consideration is continuous in time [85].

Under which conditions is the projection error small? Let us consider the case where there are $\tilde{\varphi}_i : \mathbb{R}^k \rightarrow \mathbb{R}$, $i = 1, \dots, d$, such that $\varphi_i(x) = \tilde{\varphi}_i(\xi(x))$. We then say that φ_i is a function of ξ or that ξ parametrizes φ_i . If ξ parametrizes φ_i perfectly, the projection error obviously vanishes. Thus, trivially, by choosing $\xi = \varphi = (\varphi_1, \dots, \varphi_d)^\top$, we obtain a perfect reaction coordinate since with $\tilde{\varphi}_i(z) := z_i$ with $\varphi_i = \tilde{\varphi}_i \circ \xi$. However, the eigenfunctions are *global* objects, i.e., their computation is prohibitive in high dimensions. Since we are aiming at computing a reaction coordinate, we have to answer the question of whether there is a reaction coordinate ξ that can be evaluated based on local computations only while it parametrizes the dominant eigenfunctions of \mathcal{T}^t well enough such that it leads to a small projection error. We will see next that this question can be answered by utilizing a common property of most metastable systems: The transitions between the metastable sets happen along so-called *reaction pathways*, which imply the existence of *transition manifolds* in the space of transition densities. A *suitable* parametrization of this manifold results in a parametrization of the dominant eigenfunctions with a small error.

4. Identifying good reaction coordinates

The goal is now to find a reaction coordinate ξ that is as low-dimensional as possible and results in a good projected transfer operator in the sense of (12). As we saw in the previous section, the condition $\|P_\xi^\perp \varphi_i\|_{L_\mu^2} \approx 0$ is sufficient. Thus, the idea to numerically

seek ξ that parametrizes the dominant eigenfunctions of \mathcal{T}^t in the $\|\cdot\|_{L^2_\mu}$ -norm seems natural since this would lead to small projection error $\|P_\xi^\perp \varphi_i\|_{L^2_\mu}$.

In fact, eigenfunctions of transfer operators have been used before to compute reduced dynamics and reaction coordinates: In [31], methods to decompose multiscale systems into fast and slow processes and to project the dynamics onto these subprocesses based on eigenfunctions of the Koopman operator \mathcal{K}^t are proposed. In [52], the dominant eigenfunctions of the transfer operator \mathcal{T}^t , which due to the assumed reversibility of the system is identical to \mathcal{K}^t , are shown to be good reaction coordinates. Also, comitor functions (introduced in Appendix B), which are closely related to the dominant eigenfunctions, have been used as reaction coordinates in [22, 47].

However, we propose a fundamentally different path in defining and finding reaction coordinates, as working with dominant eigenfunctions has two major disadvantages:

1. The eigenproblem is *global*. Thus if we wish to learn the value of an eigenfunction φ_i at only one location $x \in \mathbb{X}$, we need an approximation of the transfer operator \mathcal{T}_t that has to be accurate on all of \mathbb{X} . The computational effort to construct such an approximation grows exponential with $\dim(\mathbb{X})$, this is the *curse of dimensionality*. There have been attempts to mitigate this [80, 35, 81], but we aim to circumvent this problem entirely. Given two points $x, y \in \mathbb{X}$, we will decide whether $\xi(x)$ is close to $\xi(y)$ or not by using only local computations around x and y (i.e. samples from the transition densities $p^t(x, \cdot)$ and $p^t(y, \cdot)$ for moderate t).
2. The number of dominant eigenfunctions ($d + 1$) equals the number of metastable states, and this number can be much larger than the dimension of the transition manifold. This fact is illustrated in Example 4.1 below.

Example 4.1. Let us consider a diffusion process of the form (2) with the circular multi-well potential shown in Figure 2. Choosing a temperature that is not high enough for the central potential barrier to be overcome easily, transitions between the wells typically happen in the vicinity of a one-dimensional reaction pathway, the unit circle. The number of dominant eigenfunctions, however, corresponds to the number of wells. Nevertheless, projecting the system onto the unit circle would retain the dominant time scales of the system, cf. Section 5. △

4.1. Parametrization of dominant eigenfunctions

If the $(d + 1)$ dominant eigenfunctions do not depend fully on the phase space \mathbb{X} , a lower-dimensional and ultimately easier to find reaction coordinate suffices for keeping the eigenvalue approximation error (12) small. It is easy to see that if there exists a function $\xi: \mathbb{X} \rightarrow \mathbb{R}^k$ for some k so that the eigenfunctions φ are constant on the level sets of ξ , i.e., there exist functions $\tilde{\varphi}_i: \mathbb{R}^k \rightarrow \mathbb{R}$, $i = 1, \dots, d$ such that $\varphi_i = \tilde{\varphi}_i \circ \xi$, then the projection error $\|P_\xi^\perp \varphi_i\|_{L^2_\mu}$ is zero. A quantitative generalization of this is the statement that if the eigenfunctions φ_i are *almost constant* on level sets of a ξ , then the projection error is small.

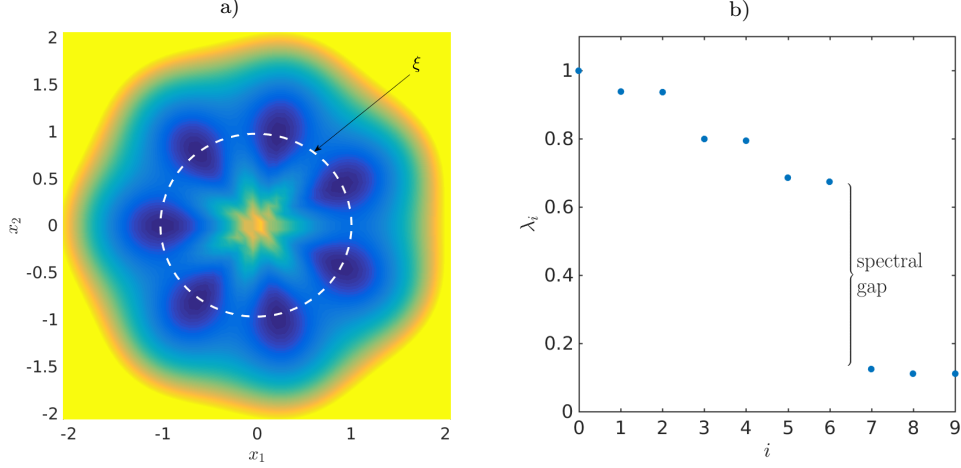


Figure 2: a) Potential with seven wells and thus seven dominant eigenvalues, but only a one-dimensional reaction coordinate. The reaction pathway is marked by a dashed white line. b) Dominant eigenvalues of \mathcal{T}^t for $t = 0.1$ and $\beta = 0.5$. The spectral gap is clearly visible.

Lemma 4.2. *Assume that there exists a function $\xi: \mathbb{X} \rightarrow \mathbb{R}^k$ for some k and functions $\tilde{\varphi}_i: \mathbb{R}^k \rightarrow \mathbb{R}$, $i = 1, \dots, d$, with*

$$|\varphi_i(x) - \tilde{\varphi}_i(\xi(x))| \leq \varepsilon \quad \forall x \in \mathbb{X}. \quad (13)$$

Then $\|P_\xi^\perp \varphi_i\|_{L_\mu^2} \leq 2\varepsilon$.

Proof. Assuming (13) holds, there exists a function $c_i: \mathbb{R} \rightarrow \mathbb{R}$ with $c_i(x) \leq 1 \quad \forall x \in \mathbb{X}$ so that

$$\varphi_i(x) = \tilde{\varphi}_i(\xi(x)) + c_i(x)\varepsilon.$$

Thus, we have

$$\begin{aligned} P_\xi \varphi_i(x) &= \int_{\mathbb{L}_{\xi(x)}} \left(\tilde{\varphi}_i(\xi(x')) + c_i(x')\varepsilon \right) d\mu_{\xi(x)}(x') \\ &= \tilde{\varphi}_i(\xi(x)) + \varepsilon \int_{\mathbb{L}_{\xi(x)}} c_i(x') d\mu_{\xi(x)}(x'). \end{aligned}$$

For the projection error, we then obtain

$$\begin{aligned} \|P_\xi \varphi_i - \varphi_i\|_{L_\mu^2} &\leq \|P_\xi \varphi_i - \tilde{\varphi}_i \circ \xi\|_{L_\mu^2} + \|\tilde{\varphi}_i \circ \xi - \varphi_i\|_{L_\mu^2} \\ &\leq 2\varepsilon. \end{aligned} \quad \square$$

Remark 4.3. From the proof we see that the pointwise condition (13) can be replaced by the much weaker condition

$$\int_{\mathbb{L}_z} |\varphi_i(x') - \tilde{\varphi}_i(\xi(x'))| d\mu_z(x') \leq \varepsilon,$$

for all level sets \mathbb{L}_z of ξ .

From here on, we address the following two central questions:

- (Q1) *In which dynamical situations can we expect to find low-dimensional reaction coordinates?*
- (Q2) *How can we computationally exploit the properties of the dynamics to obtain reaction coordinates?*

Let us start with the first question. We will address the second question in Section 4.2 and Section 4.3. Experience shows [24, 63, 25, 70] that transitions between metastable states tend to happen along so-called *reaction pathways*, which is the low-dimensional dynamical backbone in the high-dimensional state space, connecting the metastable states via saddle points of the potential V [30].

From now on, we observe the system at an intermediate time scale $t_{\text{slow}} \gg t \gg t_{\text{fast}}$ (where t_{slow} and t_{fast} are the implied time scales t_d , t_{d+1} from Section 2.3) and thus assume that the process \mathbf{X}_t has already left the transition region (if it started there), equilibrated to a quasi-stationary distribution inside some metastable wells, but has not had enough time to equilibrate *globally*. At this time scale, starting in some $x \in \mathbb{X}$, the transition density $p^t(x, \cdot)$ is observed to approximately depend only on progress along these reaction paths; see Figure 3 for an illustration. This means that the density $p^t(x, \cdot)$ on the fiber perpendicular to the transition pathway is approximately the same as $p^t(x^*, \cdot)$ for some x^* on the transition pathway. As this pathway is low-dimensional, this means that the image $\overline{\mathcal{Q}}(\mathbb{X})$ of the map

$$\overline{\mathcal{Q}}(x) := p^t(x, \cdot)$$

is almost a low-dimensional manifold in $L^1(\mathbb{X})$.

The existence of this low-dimensional structure in the space of probability densities is exactly the assumption we need to ensure that the dominant eigenfunctions are low-dimensionally parametrizable, and thus that a low-dimensional reaction coordinate ξ exists. This assumption is made precise in Definition 4.4. To summarize, we will see that ξ is a good reaction coordinate if $p^t(x, \cdot) \approx p^t(y, \cdot)$ for $\xi(x) = \xi(y)$.

Definition 4.4. *[(ε, r)-reducibility and transition manifold] We call the process \mathbf{X}_t (ε, r)-reducible, if there exists a smooth closed r -dimensional manifold $\mathbb{M} \subset L^2_{1/\mu} \subset L^1(\mathbb{X})$ such that for $t_{\text{fast}} \ll t \ll t_{\text{slow}}$ and all $x \in \mathbb{X}$*

$$\min_{f \in \mathbb{M}} \|f - p^t(x, \cdot)\|_{L^2_{1/\mu}} \leq \varepsilon \quad (14)$$

holds. We call \mathbb{M} the transition manifold and the map $\mathcal{Q}: \mathbb{X} \rightarrow \mathbb{M}$,

$$\mathcal{Q}(x) := \arg \min_{f \in \mathbb{M}} \|p^t(x, \cdot) - f\|_{L^2_{1/\mu}} \quad (15)$$

the mapping onto the transition manifold. We can set $\mathbb{M} = \text{cl}(\mathcal{Q}(\mathbb{X}))$, where $\text{cl}(\mathbb{Y})$ denotes the closure of the set \mathbb{Y} .³

³If it is necessary to break ties in (15), we can do so by taking any of the minimizers. The mapping $x \mapsto p^t(x, \cdot)$ can be shown to be smooth [5, Theorem C.1], hence $\mathcal{Q}(\mathbb{X})$ is a smooth manifold satisfying (14).

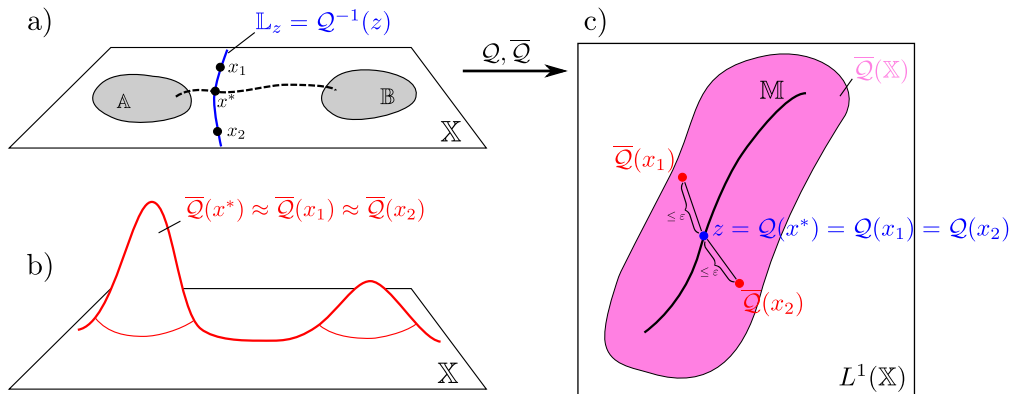


Figure 3: a) and b) The transition densities $\bar{Q}(x_1)$ and $\bar{Q}(x_2)$ are “similar” to $\bar{Q}(x^*)$ for some x^* on the transition path (dashed line) that connects the metastable sets \mathbb{A} and \mathbb{B} . c) The mapping \mathcal{Q} can be thought of as mapping all points that are “similar” under \bar{Q} to the same point in $L^1(\mathbb{X})$. The image of \mathcal{Q} thus forms a r -dimensional manifold in $L^1(\mathbb{X})$.

Remark 4.5. While it is natural to motivate (ε, r) -reducibility by the existence of reaction pathways in phase space, it is not strictly necessary. There exist stochastic systems without low-dimensional reaction pathways whose densities still quickly converge to a transition manifold in L^1 . Future work includes the identification of necessary and sufficient conditions for the existence of transition manifolds (see the first point in the conclusions). We also further elaborate on the connection between reaction pathways and transition manifolds in Appendix B.

Remark 4.6. We recall from Section 2 that the Perron–Frobenius operator \mathcal{P}^t is also naturally defined on the space $L^2_{1/\mu}$ [67]. Further, with the Dirac distribution centered in $x \in \mathbb{X}$, denoted by δ_x , we formally have $p^t(x, \cdot) = \mathcal{P}^t \delta_x$. Hence, the choice of norm in Definition 4.4 is natural. It should also be noted that since μ is a probability measure, the Hölder inequality yields $\|f\|_{L^1_\mu} \leq \|f\|_{L^2_\mu}$. Using this we have

$$\|f\|_{L^1} = \|f/\varrho\|_{L^1_\mu} \leq \|f/\varrho\|_{L^2_\mu} = \|f\|_{L^2_{1/\mu}},$$

which shows that if $p^t(x, \cdot)$ and $p^t(y, \cdot)$ are close in the $L^2_{1/\mu}$ norm, they are also close in the L^1 norm. We require the closeness of the respective $p^t(x, \cdot)$ in the $L^2_{1/\mu}$ norm for our theoretical considerations below, but otherwise we will think of them as functions in L^1 .

Note that we only need to evolve the system at hand for a moderate time $t \ll t_{\text{slow}}$, which has to be merely sufficiently large to damp out the fast fluctuations in the metastable states. This will be an important point later, allowing for numerical tractability.

Next, we show that (ε, r) -reducibility implies that dominant eigenfunctions are almost constant on the level sets of \mathcal{Q} .

Lemma 4.7. *If \mathbf{X}_t is (ε, r) -reducible, then for an eigenfunction φ_i of \mathcal{T}^t with $\|\varphi_i\|_{L_\mu^2} = 1$ and points $x, y \in \mathbb{X}$ with $\mathcal{Q}(x) = \mathcal{Q}(y)$ we have*

$$|\varphi_i(x) - \varphi_i(y)| \leq \frac{2\varepsilon}{|\lambda_i|}.$$

Proof. First note that for the transition densities $p^t(x, \cdot)$, $p^t(y, \cdot)$ it holds that

$$\begin{aligned} \|p^t(x, \cdot) - p^t(y, \cdot)\|_{L_{1/\mu}^2} &\leq \|p^t(x, \cdot) - \mathcal{Q}(x)\|_{L_{1/\mu}^2} + \|\mathcal{Q}(x) - p^t(y, \cdot)\|_{L_{1/\mu}^2} \\ &= \|p^t(x, \cdot) - \mathcal{Q}(x)\|_{L_{1/\mu}^2} + \|\mathcal{Q}(y) - p^t(y, \cdot)\|_{L_{1/\mu}^2} \leq 2\varepsilon. \end{aligned} \quad (16)$$

With this we can show the assertion:

$$\lambda_i \varphi_i(x) = \mathcal{T}^t \varphi_i(x) = \mathcal{K}^t \varphi_i(x) = \int_{\mathbb{X}} \varphi_i(x') p^t(x, x') dx'.$$

Applying (16), for some function $e \in L_{1/\mu}^2(\mathbb{X})$ with $\|e\|_{L_{1/\mu}^2} \leq 2\varepsilon$, we get

$$\begin{aligned} \lambda_i \varphi_i(x) &= \int_{\mathbb{X}} \varphi_i(x') (p^t(y, x') + e(x')) dx' \\ &= \int_{\mathbb{X}} \varphi_i(x') p^t(y, x') dx' + \int_{\mathbb{X}} \varphi_i(x') \frac{e(x')}{\varrho(x')} d\mu(x') \\ &= \lambda_i \varphi_i(y) + \int_{\mathbb{X}} \varphi_i(x') \frac{e(x')}{\varrho(x')} d\mu(x'), \end{aligned}$$

where in the last equation, we again used that due to reversibility $\mathcal{K}^t = \mathcal{T}^t$ and that φ_i is an eigenfunction. Thus for the difference, we have

$$\begin{aligned} |\varphi(x) - \varphi(y)| &= \frac{1}{|\lambda_i|} \left| \int_{\mathbb{X}} \varphi_i(x') \frac{e(x')}{\varrho(x')} d\mu(x') \right| \\ &\leq \frac{1}{|\lambda_i|} \underbrace{\|\varphi_i\|_{L_\mu^2}}_{=1} \underbrace{\|e/\varrho\|_{L_\mu^2}}_{=\|e\|_{L_{1/\mu}^2}} \leq \frac{2\varepsilon}{|\lambda_i|}. \end{aligned} \quad \square$$

Assuming that the eigenfunctions are normalized (which we do from now on), i.e., $\|\varphi_i\|_{L_\mu^2} = 1$, and that ε is sufficiently small, Lemma 4.7 implies that the dominant eigenfunctions (i.e., $|\lambda_i| \approx 1$) are almost constant on the level sets of \mathcal{Q} . This can now be used to show that the φ_i are not fully dependent on \mathbb{X} , but only on the level sets of \mathcal{Q} (up to a small error), in a sense similar to Lemma 4.2.

Corollary 4.8. *Let \mathbf{X}_t be (ε, r) -reducible. Then there exists a function $\tilde{\varphi}_i: \mathbb{M} \rightarrow \mathbb{R}$ such that*

$$|\varphi_i(x) - \tilde{\varphi}_i(\mathcal{Q}(x))| \leq \frac{\varepsilon}{|\lambda_i|}.$$

Proof. Fix $x \in \mathbb{X}$, and let $z = \mathcal{Q}(x)$. Define the function $\tilde{\varphi}_i$ by

$$\tilde{\varphi}_i(\mathcal{Q}(x)) := \frac{1}{2} \left(\inf_{\mathcal{Q}(y)=z} \varphi_i(y) + \sup_{\mathcal{Q}(y)=z} \varphi_i(y) \right).$$

Since by Lemma 4.7 it holds that $|\varphi_i(x) - \varphi_i(y)| \leq \frac{2\varepsilon}{|\lambda_i|}$ if $\mathcal{Q}(x) = \mathcal{Q}(y)$, we have that

$$\left| \sup_{\mathcal{Q}(y)=z} \varphi_i(y) - \inf_{\mathcal{Q}(y)=z} \varphi_i(y) \right| \leq \frac{2\varepsilon}{|\lambda_i|},$$

thus our choice of $\tilde{\varphi}_i$ gives

$$|\varphi_i(x) - \tilde{\varphi}_i(\mathcal{Q}(x))| \leq \frac{\varepsilon}{|\lambda_i|}.$$

□

4.2. Embedding the transition manifold

In light of Corollary 4.8, one could say that \mathcal{Q} is an “ \mathbb{M} -valued reaction coordinate”. However, as we have no access to \mathbb{M} so far, and a \mathbb{R}^k -valued reaction coordinate is more intuitive, we aim to obtain a more useful representation of the transition manifold through *embedding* it into a finite, possibly low-dimensional Euclidean space.

We will see that we are very free in the choice of the embedding mapping, even though the manifold \mathbb{M} is not known explicitly (we only assumed that it exists). To achieve this, we will use an infinite-dimensional variant of the *weak Whitney embedding theorem* [66, 83], which, roughly speaking, states that “almost every bounded linear map from $L^1(\mathbb{X})$ to \mathbb{R}^{2r+1} will be one-to-one on \mathbb{M} and its image”. We first specify what we mean by “almost every” in the context of bounded linear maps, following the notions of Sauer et al. [66].

Definition 4.9 (Prevalence). *A Borel subset \mathbb{S} of a normed linear space \mathbb{V} is called prevalent if there is a finite-dimensional subspace \mathbb{E} of \mathbb{V} such that for each $v \in \mathbb{V}$, $v + e$ belongs to \mathbb{S} for (Lebesgue) almost every e in \mathbb{E} .*

As the infinite-dimensional embedding theorem from Hunt et al. [34] is applicable not only to smooth manifolds, but to arbitrary subsets $\mathbb{A} \subset \mathbb{V}$ of fractal dimension, it uses the concepts of *box covering dimension* $\dim_B(\mathbb{A})$ and *thickness exponent* $\tau(\mathbb{A})$ from fractal geometry. Intuitively, $\dim_B(\mathbb{A})$ describes the exponent of the growth rate in the number of boxes of decreasing side length that are needed to cover \mathbb{A} , and $\tau(\mathbb{A})$ describes how well \mathbb{A} can be approximated using only finite-dimensional linear subspaces of \mathbb{V} . As these concepts coincide with the traditional measure of dimensionality in our setting, we will not go into detail here and point to [34] for a precise definition.

The general infinite-dimensional embedding theorem reads:

Theorem 4.10 ([34, Theorem 3.9]). *Let \mathbb{V} be a Banach space and $\mathbb{A} \subset \mathbb{V}$ be a compact set with box-counting dimension d and thickness exponent τ . Let $k > 2d$ be an integer, and let α be a real number with*

$$0 < \alpha < \frac{k - 2d}{k(1 + \tau)}.$$

Then for almost every (in the sense of prevalence) bounded linear function $\mathcal{E} : \mathbb{V} \rightarrow \mathbb{R}^k$ there exists $C > 0$ such that for all $x, y \in \mathbb{A}$,

$$C \|\mathcal{E}(x) - \mathcal{E}(y)\|_2^\alpha \geq \|x - y\|_2, \quad (17)$$

where $\|\cdot\|_2$ denotes the Euclidean 2-norm.

Note that (17) implies Hölder continuity of \mathcal{E}^{-1} on $\mathcal{E}(\mathbb{A})$ and in particular that \mathcal{E} is one-to-one on \mathbb{A} and its image. Using that the box counting dimension of a smooth r -dimensional manifold \mathbb{K} is simply r and that the thickness exponent is bounded from above by the box-counting dimension, thus $0 \leq \tau(\mathbb{K}) \leq r$, see [34], we get the following infinite-dimensional embedding theorem for smooth manifolds.

Corollary 4.11. *Let \mathbb{V} be a Banach space, let $\mathbb{K} \subset \mathbb{V}$ be a smooth manifold of dimension r and let $k > 2r$. Then almost every (in the sense of prevalence) bounded linear function $\mathcal{E} : \mathbb{V} \rightarrow \mathbb{R}^k$ is one-to-one on \mathbb{K} and its image in \mathbb{R}^k .*

Thus, since the transition manifold \mathbb{M} is assumed to be a smooth r -dimensional manifold in $L^1(\mathbb{X})$, an arbitrarily chosen bounded linear map $\mathcal{E} : L^1(\mathbb{X}) \rightarrow \mathbb{R}^{2r+1}$ can be assumed to be one-to-one on \mathbb{M} and its image. In particular, $\mathcal{E}(\mathbb{M})$ is again an r -dimensional manifold (although not necessarily smooth). With this insight, we can now construct a reaction coordinate in Euclidean space:

Corollary 4.12. *Let \mathbf{X}_t be (ε, r) -reducible and let $\mathcal{E} : L^1(\mathbb{X}) \rightarrow \mathbb{R}^{2r+1}$ be one-to-one on \mathbb{M} and its image. Define $\xi : \mathbb{R}^n \rightarrow \mathbb{R}^{2r+1}$ by*

$$\xi(x) := \mathcal{E}(\mathcal{Q}(x)). \quad (18)$$

Then there exists a function $\hat{\varphi}_i : \mathbb{R}^{2r+1} \rightarrow \mathbb{R}$ so that

$$|\varphi_i(x) - \hat{\varphi}_i(\xi(x))| \leq \frac{\varepsilon}{|\lambda_i|}. \quad (19)$$

Proof. As \mathcal{E} is one-to-one on \mathbb{M} and its image, it is invertible on $\mathcal{E}(\mathbb{M})$. With $\tilde{\varphi}_i$ chosen as in the proof of Corollary 4.8, define $\hat{\varphi}_i : \mathcal{E}(\mathbb{M}) \rightarrow \mathbb{R}$ by

$$\hat{\varphi}_i(\hat{z}) := \tilde{\varphi}_i(\mathcal{E}^{-1}(\hat{z})). \quad (20)$$

Then

$$|\varphi_i(x) - \hat{\varphi}_i(\xi(x))| = |\varphi_i(x) - \tilde{\varphi}_i(\mathcal{Q}(x))| \stackrel{\text{Cor. 4.8}}{\leq} \frac{\varepsilon}{|\lambda_i|}. \quad \square$$

Since $\widehat{\mathbb{M}} := \mathcal{E}(\mathbb{M})$ is an r -dimensional manifold, ξ is effectively an r -dimensional reaction coordinate. Thus, if the right-hand side of (19) is small, the φ_i are “almost parametrizable” by the r -dimensional reaction coordinate ξ . Using Lemma 4.2, we immediately see that this results in a small projection error $\|P_{\xi}^{\perp}\varphi_i\|$, and due to Corollary 3.6 in a good transfer operator approximation; hence ξ is a good reaction coordinate.

The reaction coordinate ξ remains an “ideal” case, because we have no access to the map Q and hence to \mathbb{M} , only to $\overline{Q}(x) = p^t(x, \cdot) \approx Q(x)$. We summarize the construction of the ideal reaction coordinate ξ in Figure 4.

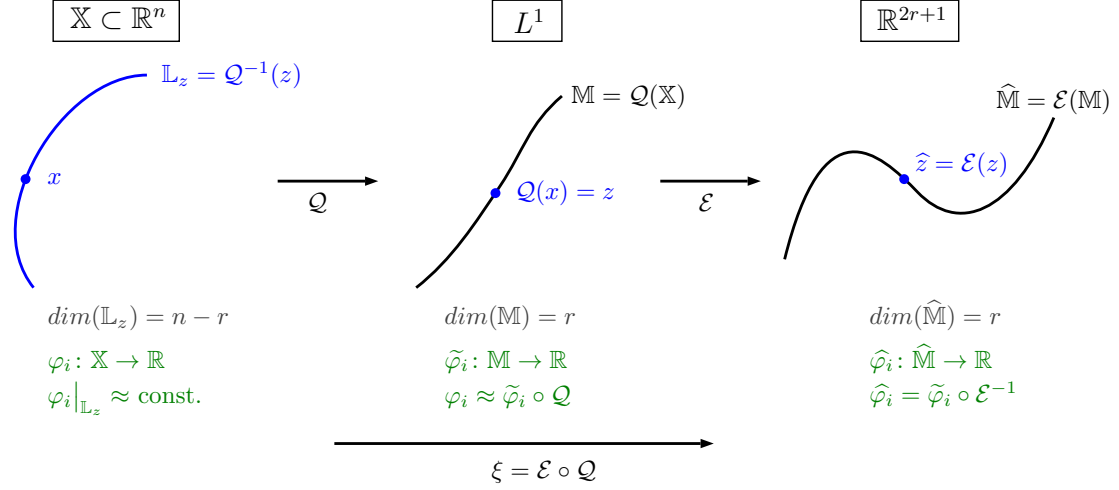


Figure 4: Summary of the construction of the ideal reaction coordinate ξ .

Remark 4.13. The recent work of Dellnitz et al. [19] uses similar embedding techniques to identify finite-dimensional objects in the state space of infinite-dimensional dynamical systems. They utilize the infinite-dimensional delay-embedding theorem of Robinson [64], a generalization of the well-known *Takens embedding theorem* [73], to compute finite-dimensional attractors of delay differential equations by established subdivision techniques [17].

4.3. Numerical approximation of the reaction coordinate

Approximate embedding of the transition manifold. We now elaborate how to construct a good reaction coordinate $\bar{\xi}$ numerically. To use the central definition (18) in practice, two points have to be addressed:

1. How to choose the embedding \mathcal{E} ?
2. How to deal with the fact that we do not know Q ?

For the choice of \mathcal{E} , we restrict ourselves to linear maps of the form

$$\mathcal{E}(f) := \begin{pmatrix} \langle f, \eta_1 \rangle \\ \vdots \\ \langle f, \eta_{2r+1} \rangle \end{pmatrix}, \quad (21)$$

with arbitrarily chosen linearly independent functions $\eta_i \in L^\infty(\mathbb{X})$, where $\langle f, \eta_i \rangle = \int f \eta_i$. In practice, we will choose the $\eta_i : \mathbb{X} \rightarrow \mathbb{R}$ as linear functions themselves, i.e. $\eta_i(x) = a_i^\top x$ for some, usually randomly drawn, $a_i \in \mathbb{R}^n$. Note that then $\eta_i \notin L^\infty$, but this is not a problem because we will embed the functions $f = p^t(x, \cdot)$, and $p^t(x, y)$ can be shown to decay exponentially as $\|y\|_2 \rightarrow \infty$, cf. [5, Theorem C.1]. Thus, $\langle f, \eta_i \rangle$ will exist. For linearly independent η_i , these maps are still generic in the sense of the Whitney embedding theorem, and thus still embed the transition manifold \mathbb{M} .

A natural choice for the approximation of the unknown map \mathcal{Q} is the mapping to the transition probability density,

$$\bar{\mathcal{Q}} : x \mapsto p^t(x, \cdot), \quad (22)$$

as $\|\bar{\mathcal{Q}}(x) - p^t(x, \cdot)\|_{L^2_{1/\mu}} \leq \varepsilon$. With this, we consider

$$\mathcal{E}(\bar{\mathcal{Q}}(x)) = \mathcal{E}(p^t(x, \cdot)) = \begin{pmatrix} \langle p^t(x, \cdot), \eta_1 \rangle \\ \vdots \\ \langle p^t(x, \cdot), \eta_{2r+1} \rangle \end{pmatrix} \stackrel{(1)}{=} \begin{pmatrix} \mathcal{K}^t \eta_1(x) \\ \vdots \\ \mathcal{K}^t \eta_{2r+1}(x) \end{pmatrix}. \quad (23)$$

The values on the right-hand side can in turn be approximated by a Monte Carlo quadrature, using only short-time trajectories of the original dynamics:

$$\mathcal{K}^t \eta_i(x) = \mathbb{E}[\eta_i(\mathbf{X}_t) \mid \mathbf{X}_0 = x] \approx \frac{1}{M} \sum_{m=1}^M \eta_i(\Phi_t^{(m)}(x)), \quad (24)$$

where the $\Phi_t^{(m)}(x)$ are independent realizations of \mathbf{X}_t with starting point $\mathbf{X}_0 = x$, in practice realized by a stochastic integrator (e.g. Euler–Maruyama).

The computationally infeasible reaction coordinate ξ . Note that $\mathcal{E} \circ \bar{\mathcal{Q}}$ is not yet an r -dimensional reaction coordinate, since $\bar{\mathcal{Q}}(\mathbb{X})$ is only approximately an r -dimensional manifold; more precisely, it lies in the ε -neighborhood of an r -dimensional submanifold \mathbb{M} of L^1 . Hence, also $\mathcal{E}(\bar{\mathcal{Q}}(\mathbb{X}))$ is only approximately an r -dimensional manifold; see the magenta regions in Figure 5.

The question now is how we can reduce $\mathcal{E} \circ \bar{\mathcal{Q}}$ to an r -dimensional *good* reaction coordinate. Since we know from above that $\xi = \mathcal{E} \circ \mathcal{Q}$ is a good reaction coordinate, let us see what would be needed to construct it.

The property of ξ that we want is that it is constant along level sets \mathbb{L}_z of \mathcal{Q} , i.e., $\xi|_{\mathbb{L}_z} = \text{const}$ (because this implies that it is a good reaction coordinate, cf. Corollary 4.12). Hence, if we could identify $\hat{\mathbb{M}}$ as an r -dimensional manifold in \mathbb{R}^{2r+1} , we

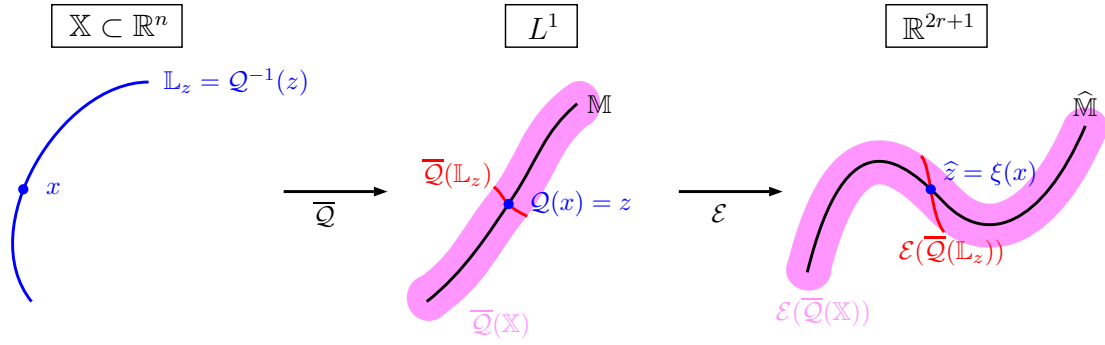


Figure 5: How to make a good r -dimensional reaction coordinate out of $\mathcal{E} \circ \overline{\mathcal{Q}}$? Given the analysis from the previous section, we would like to make the level sets \mathbb{L}_z of \mathcal{Q} also the level sets of $\bar{\xi}$ (red line segment). Unfortunately, we have no access to these.

would project $\mathcal{E}(\overline{\mathcal{Q}}(x))$ along $\mathcal{E}(\overline{\mathcal{Q}}(\mathbb{L}_z))$ onto $\hat{\mathbb{M}}$ — assuming that $\hat{\mathbb{M}}$ and $\mathcal{E}(\overline{\mathcal{Q}}(\mathbb{L}_z))$ intersect in \mathbb{R}^{2r+1} — to obtain $\xi(x)$ as the resulting point (see Figure 5, where we would project along the red line on the right). Unfortunately, we have no access to \mathcal{Q} (not to mention that $\hat{\mathbb{M}}$ and $\mathcal{E}(\overline{\mathcal{Q}}(\mathbb{L}_z))$ need not intersect in \mathbb{R}^{2r+1}) and hence to its level sets \mathbb{L}_z . Thus, this strategy seems infeasible.

A computationally feasible reaction coordinate. What helps us at this point is that there is a certain amount of arbitrariness in the definition of \mathcal{Q} . Recalling Definition 4.4, what we are given is $\overline{\mathcal{Q}}$, and we construct $\mathcal{Q}(x)$ as a projection of $\overline{\mathcal{Q}}(x)$ onto the r -dimensional manifold \mathbb{M} by the closest-point projection \mathcal{Q}' ; i.e., $\mathcal{Q} = \mathcal{Q}' \circ \overline{\mathcal{Q}}$. This choice of \mathcal{Q}' is convenient, because we can show

$$\|\overline{\mathcal{Q}}(x) - \overline{\mathcal{Q}}(y)\|_{L^2_{1/\mu}} \leq 2\varepsilon \quad \text{for every } \mathcal{Q}(x) = \mathcal{Q}(y) \text{ (i.e., on level sets of } \mathcal{Q}\text{)}, \quad (25)$$

which is used in Lemma 4.7. Other choices of \mathcal{Q}' could, however, yield a similarly practicable $\mathcal{O}(\varepsilon)$ -bound in (25). Our strategy will be to choose a specific r -dimensional reaction coordinate $\bar{\xi}$ and to show that in general it can be expected to be a good reaction coordinate.

Let us recall that, by assumption, the set $\overline{\mathcal{Q}}(\mathbb{X})$ is contained in the ε -neighborhood of an unknown smooth r -dimensional manifold $\mathbb{M} \subset L^1(\mathbb{X})$. Thus, a generic smooth map $\mathcal{E}: L^1(\mathbb{X}) \rightarrow \mathbb{R}^{2r+1}$ will embed \mathbb{M} into \mathbb{R}^{2r+1} , forming a diffeomorphism from \mathbb{M} to $\hat{\mathbb{M}}$. Thus, \mathcal{E} is going to map $\overline{\mathcal{Q}}(\mathbb{X})$ to an $\mathcal{O}(\varepsilon)$ -neighborhood of $\hat{\mathbb{M}}$. This means, the r -dimensional manifold structure of $\hat{\mathbb{M}}$ should still be detectable and can be identified with standard manifold learning tools. We use the diffusion maps algorithm (see Section 4.4 below), which gives us a map $\Psi: \mathbb{R}^{2r+1} \rightarrow \mathbb{R}^r$ (the diffusion map). Then we define $\bar{\xi}$ as

$$\bar{\xi} := \Psi \circ \mathcal{E} \circ \overline{\mathcal{Q}}. \quad (26)$$

This is depicted on the right-hand side of Figure 6, where the red dashed line shows the level set $\hat{\mathbb{L}}_{\hat{z}} = \{z \in \mathbb{R}^{2r+1} : \Psi(z) = \Psi(\hat{z})\}$.

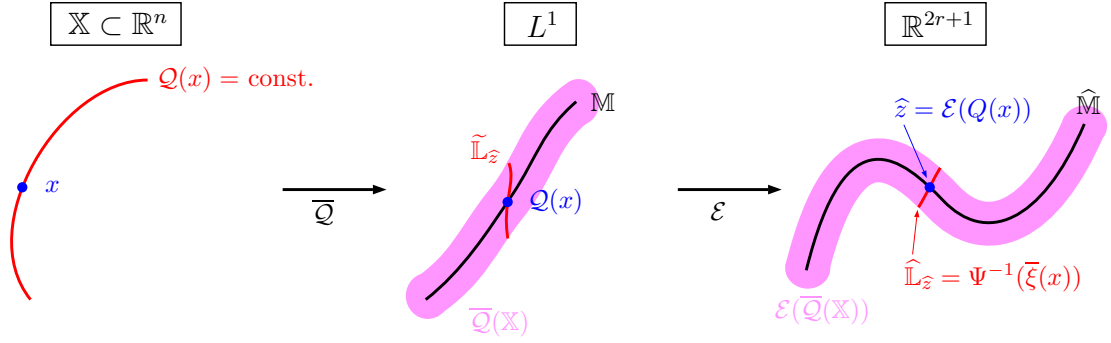


Figure 6: The realized reaction coordinate $\bar{\xi}$.

Next, we consider the set $\tilde{\mathbb{L}}_{\hat{z}} := \mathcal{E}^{-1}(\hat{\mathbb{L}}_{\hat{z}}) \cap \overline{\mathcal{Q}}(\mathbb{X})$. It holds that $\tilde{\mathbb{L}}_{\hat{z}} = \{\overline{\mathcal{Q}}(x) \mid \bar{\xi}(x) = \Psi(\hat{z})\}$. Recall that $\mathcal{E}: \mathbb{M} \rightarrow \hat{\mathbb{M}}$ is one-to-one, thus $\tilde{\mathbb{L}}_{\hat{z}}$ intersects \mathbb{M} in exactly one point. We define this one point as $\mathcal{Q}(x)$, and thus \mathcal{Q}' is the projection onto \mathbb{M} along $\tilde{\mathbb{L}}_{\hat{z}}$. We see that \mathcal{Q} is well-defined and that $\mathcal{Q}(x) = \mathcal{Q}(y) \Leftrightarrow \bar{\xi}(x) = \bar{\xi}(y)$.

At this point we assume that \mathcal{E}^{-1} is sufficiently well-behaved in a neighborhood of $\hat{\mathbb{M}}$, it does not “distort transversality” of intersections, such that the diameter of $\tilde{\mathbb{L}}_{\hat{z}}$ is $\mathcal{O}(\varepsilon)$ with a moderate constant in $\mathcal{O}(\cdot)$. We will investigate a formal justification of this fact in a future work, here we assume it holds true, and we will see in the numerical experiments that the assumption is justified. This assumption implies that $\|\overline{\mathcal{Q}}(x) - \overline{\mathcal{Q}}(y)\|_{L^2_{1/\mu}} = \mathcal{O}(\varepsilon)$ for $\mathcal{Q}(x) = \mathcal{Q}(y)$, i.e. for $\bar{\xi}(x) = \bar{\xi}(y)$. Now, however, Lemma 4.7 implies that φ_i is almost constant (up to an error $\mathcal{O}(\varepsilon)$) on level sets of $\bar{\xi}$, which, in turn, by Lemma 4.2 and Corollary 3.6 shows that $\bar{\xi}$ is a good reaction coordinate.

4.4. Identification of $\hat{\mathbb{M}}$ through Manifold Learning

In this section, we describe how to identify $\hat{\mathbb{M}}$ numerically. The task is as follows: Given that we have computed $\mathcal{E}(\overline{\mathcal{Q}}(x_i)) = \hat{z}_i \in \mathbb{R}^{2r+1}$ for a number of sample points $\{x_i\}_{i=1}^\ell \subset \mathbb{X}$, we would like to identify the r -dimensional manifold $\hat{\mathbb{M}}$, noting the points $\mathcal{E}(\overline{\mathcal{Q}}(x_i))$ are in a $\mathcal{O}(\varepsilon)$ -neighborhood of $\hat{\mathbb{M}}$ (see Section 4.3). Additionally, we would like an r -dimensional coordinate function $\Psi: \mathbb{R}^{2r+1} \rightarrow \mathbb{R}^r$ that parameterizes $\hat{\mathbb{M}}$ (so that the level sets of Ψ are transversal to $\hat{\mathbb{M}}$).

This is a default setting for which manifold learning algorithms can be applied. Many standard methods exist; we name multidimensional scaling [40, 39], Isomap [74], and diffusion maps [13] as a few of the most prominent examples. Because of its favorable properties, we choose the diffusion maps algorithm here and summarize it briefly for our setting in what follows. For details, the reader is referred to [13, 56, 12, 71].

Given sample points $\{\hat{z}_i\}_{i=1}^\ell \subset \mathbb{R}^{2r+1}$, diffusion maps proceeds by constructing a similarity matrix $W \in \mathbb{R}^{\ell \times \ell}$ with

$$W_{ij} = h \left(\frac{\|\hat{z}_i - \hat{z}_j\|_2^2}{\sigma} \right),$$

where $\|\cdot\|_2$ is the Euclidean norm in \mathbb{R}^{2r+1} , $\sigma > 0$ is a scale factor, and $h : \mathbb{R} \rightarrow \mathbb{R}_+$ is a kernel function which is most commonly chosen as $h(x) = \exp(-x)1_{x \leq R}$ with a suitably chosen cutoff R that sparsifies W and ensures that only local distances enter the construction. With D being the diagonal matrix containing the row sums of W , the similarity matrix is then normalized to give $\tilde{W} = D^{-1}WD^{-1}$. Finally, the stochastic matrix $P = \tilde{D}^{-1}\tilde{W}$ is constructed, where \tilde{D} is the diagonal matrix containing the row sums of \tilde{W} . P is similar to the symmetric matrix $\tilde{D}^{-1/2}\tilde{W}\tilde{D}^{-1/2}$, thus it has an orthonormal basis of eigenvectors $\{\psi_i\}_{i=0}^{\ell-1}$ with real eigenvalues γ_i . Since P is also stochastic, $|\gamma_i| \leq 1$. The diffusion map is then given by

$$\Psi : \mathbb{R}^{2r+1} \rightarrow \mathbb{R}^r, \quad \Psi(\hat{z}) = (\gamma_1\psi_1(\hat{z}), \dots, \gamma_r\psi_r(\hat{z}))^\top. \quad (27)$$

Using properties of the Laplacian eigenproblem on $\hat{\mathbb{M}}$, one can show that Ψ indeed parameterizes the r -dimensional manifold $\hat{\mathbb{M}}$ for suitably chosen σ [13].

Remark 4.14. The diffusion maps algorithm will only reliably identify $\hat{\mathbb{M}}$ based on the neighborhood relations between the embedded sample points z_i , if the points cover all parts of $\hat{\mathbb{M}}$ sufficiently well. In particular, as $p^t(x, \cdot)$ and thus $(\mathcal{E} \circ \mathcal{Q})(x)$ vary strongly with x traversing the transition regions, a good coverage of those regions is required.

For the various low-dimensional academical examples Section 5, this is ensured by choosing the x_i to be a dense grid of points in \mathbb{X} . For the high-dimensional example in Section 5.2, the evaluation points are generated as a subsample from a long equilibrated trajectory, essentially sampling μ . Both of these ad-hoc methods are likely to be unapplicable in realistic high-dimensional systems with very long equilibration times. However, as we mentioned in the introduction, there exist multiple statistical and dynamical approaches to this common problem of quickly sampling the relevant parts of phase space, including the transition regions. Each of these sampling methods can be easily integrated into our proposed algorithm as a pre-processing step.

Fundamentally though, the central idea of our method does not depend crucially on the applicability of diffusion maps. Rather, the latter can be considered an optional post-processing step. Using the $2r + 1$ -dimensional reaction coordinate

$$\bar{\bar{\xi}} := \mathcal{E} \circ \bar{\mathcal{Q}},$$

i.e. (26) without the manifold learning step, may in practice already represent a sufficient dimensionality reduction.

In addition, situations may occur where the a priori generation of evaluation points is not possible or desired. One of the final goals and currently work in progress is the construction of an accelerated integration scheme that generates significant evaluation points and their reaction coordinate value “on the fly”. This is related to the effective dynamics mentioned in fifth point of the conclusion. However, this also requires us to be able to evaluate the reaction coordinate at isolated points, independent of each other, and thus also necessitates the use of the above $\bar{\bar{\xi}}$ instead of $\bar{\xi}$.

5. Numerical Examples

Based on the results from the previous sections, we propose the following algorithm to compute reaction coordinates numerically:

1. Let $x_i, i = 1, \dots, \ell$, be the points for which we would like to evaluate $\bar{\xi}$. Here, we assume the points satisfy the requirements addressed in Remark 4.14.
2. Choose linearly independent functions $\eta_j \in L^\infty(\mathbb{X}), j = 1, \dots, 2r+1$. The essential boundedness of the η_j is not necessary, but $|\eta_j(x)|$ should not grow faster than a polynomial as $\|x\|_2 \rightarrow \infty$.
3. In each point x_i , start M simulations of length t and estimate $\mathcal{E}_j(\bar{\mathcal{Q}}(x_i))$ using (23) and (24), to obtain the point $\hat{z}_i \in \mathbb{R}^{2r+1}$. We discuss the appropriate choice of M and t in Section 5.1.
4. Apply the diffusion maps technique from Section 4.4 for the point cloud $\{\hat{z}_i\}_{i=1}^\ell$, and obtain $\Psi : \mathbb{R}^{2r+1} \rightarrow \mathbb{R}^r$, a parametrization of the point in its r essential directions of variation.
5. By (27), we define the reaction coordinate as

$$\bar{\xi} : x_i \mapsto \Psi(\hat{z}_i). \quad (28)$$

The numerical effort of this algorithm depends strongly on the third step. Given ℓ evaluation points, and a choice of M trajectories per point, the cost is mainly given by $M \cdot \ell \cdot c(t)$, where $c(t)$ is the effort of a single numerical realization of the dynamics up to time t . The high-dimensional phase space only enters the algorithm as the domain of the observables η_j . The cost of evaluating those typically very simple functions⁴ at the $M \cdot \ell$ end points of the trajectory is negligible. The cost of the method is thus essentially independent of n .

In order to demonstrate the efficacy of our method, we compute the reaction coordinates for three representative problems, namely a simple curved double-well potential, a multi-well potential defined on a circle, both in low and high dimensions, and two slightly different quadruple-well potentials stressing the difference between a one- and a two-dimensional reaction coordinate.

5.1. Curved double-well potential

As a first verification, we consider a system with an analytically known reaction coordinate that is then used for comparison. Consider the two-dimensional drift-diffusion process (2) with potential

$$V(x_1, x_2) = (x_1^2 - 1)^2 + 2(x_1^2 + x_2 - 1)^2$$

⁴In our examples, we used linear functions with great success.

and inverse temperature $\beta = 0.5$. This potential already served as a motivational example for the nature of reaction coordinates in the introduction and is shown in Figure 1. The system possesses two metastable sets around the minima $(-1, 0)^\top$ and $(1, 0)^\top$, which are connected by the transition path $\{x \in \mathbb{R}^2 \mid x_2 = 1 - x_1^2\}$. The implied time scales, defined in (3), can be computed from the eigenvalues using a standard Ulam-type Galerkin discretization [37, 38] of the transfer operator \mathcal{T}^t and are shown in Figure 7a⁵. We observe a significant gap between t_1 and t_2 and thus identify t_1 as the last slow and t_2 as the first fast time scale. Choosing the lag time $t = 2$ then satisfies $t_{\text{slow}} > t > t_{\text{fast}}$. A visual inspection of a typical trajectory of length t starting in one of the two metastable sets as shown in Figure 7b confirms that the respective set is sampled, yet a transition to the other set is a rare event.

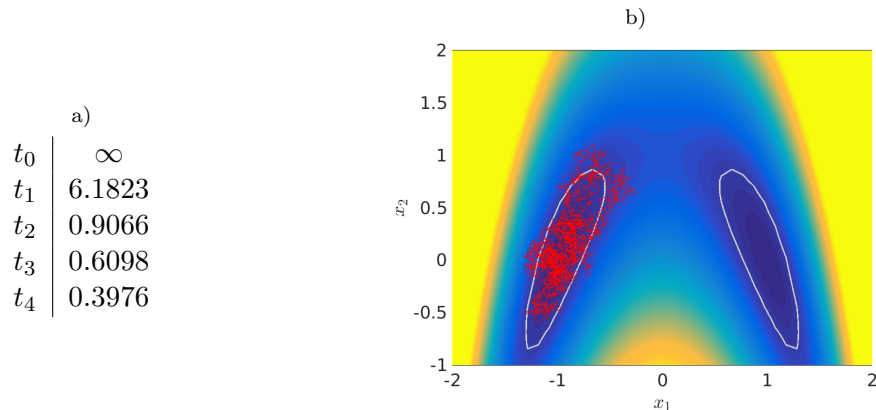


Figure 7: a) Implied time scales of the double-well system. b) Trajectory of length $t = 1$.

The low dimension of the system allows us to compute the reaction coordinate on a full regular grid over the phase space. We choose a 40×30 grid in the rectangular region $[-2, 2] \times [-1, 2]$ and denote the set of grid points by $\bar{\mathbb{X}}$. For this system, we expect a one-dimensional transition path and thus a one-dimensional reaction coordinate ξ . That is, $r = 1$ and $2r + 1 = 3$. Thus, we choose three linear observables in our embedding function (21), e.g.,

$$\begin{aligned}
 \eta_1(x_1, x_2) &= -0.2630 x_1 - 0.3186 x_2, \\
 \eta_2(x_1, x_2) &= -0.2246 x_1 + 0.0969 x_2, \\
 \eta_3(x_1, x_2) &= 0.1564 x_1 + 0.0783 x_2,
 \end{aligned} \tag{29}$$

whose coefficients were drawn uniformly from $[-1, 1]$. The expectation value in (23) is approximated by a Monte Carlo quadrature using $M = 10^5$ sample trajectories for each grid point, cf. (24). The parameter M was chosen such that the error in (24), commonly defined as the variance of the Monte Carlo sum, is sufficiently low. The

⁵In realistic, high-dimensional systems, the computation of the dominant eigenvalues using grid-based methods is likely infeasible. In these situations, the implied time scales have to be estimated, for example using standard Markov State Model techniques [7].

resulting embedding of the grid points x into \mathbb{R}^3 is shown in Figure 8. The transition path seems to be already parametrized well by the individual components of $\mathcal{E} \circ \bar{\mathcal{Q}}$.

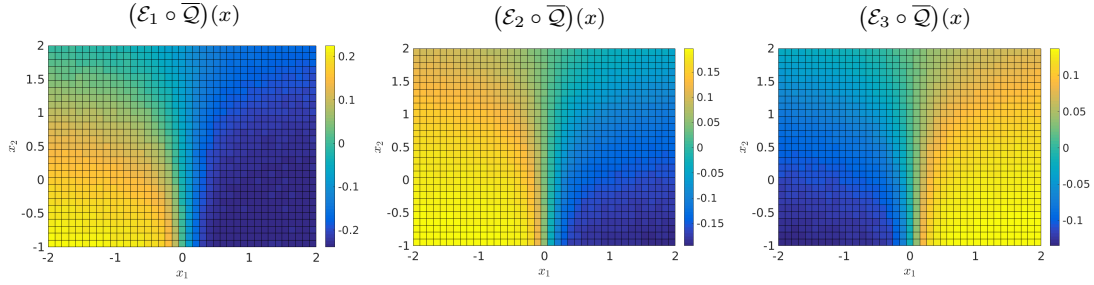


Figure 8: The individual components of the embedding $\mathcal{E} \circ \bar{\mathcal{Q}}$ on the grid points $x \in \bar{\mathbb{X}}$.

For this example, the image of \mathbb{X} under $\mathcal{E} \circ \bar{\mathcal{Q}}$ should form a compact neighborhood of the one-dimensional manifold $\mathcal{E}(\mathbb{M})$, as described in Section 4.3. The one-dimensional structure in $\mathcal{E}(\bar{\mathcal{Q}}(\bar{\mathbb{X}}))$ is clearly visible, see Figure 9a. To identify the one-dimensional coordinate along this set the diffusion map algorithm is used. Let $\Psi_1 : (\mathcal{E} \circ \bar{\mathcal{Q}})(\bar{\mathbb{X}}) \rightarrow \mathbb{R}$ denote the first diffusion map coordinate on the embedded grid points, also visualized in Figure 9a. The final reaction coordinate, shown in Figure 9b, is then given by

$$\bar{\xi}(x) := \Psi_1((\mathcal{E} \circ \bar{\mathcal{Q}})(x)), \quad x \in \bar{\mathbb{X}}.$$

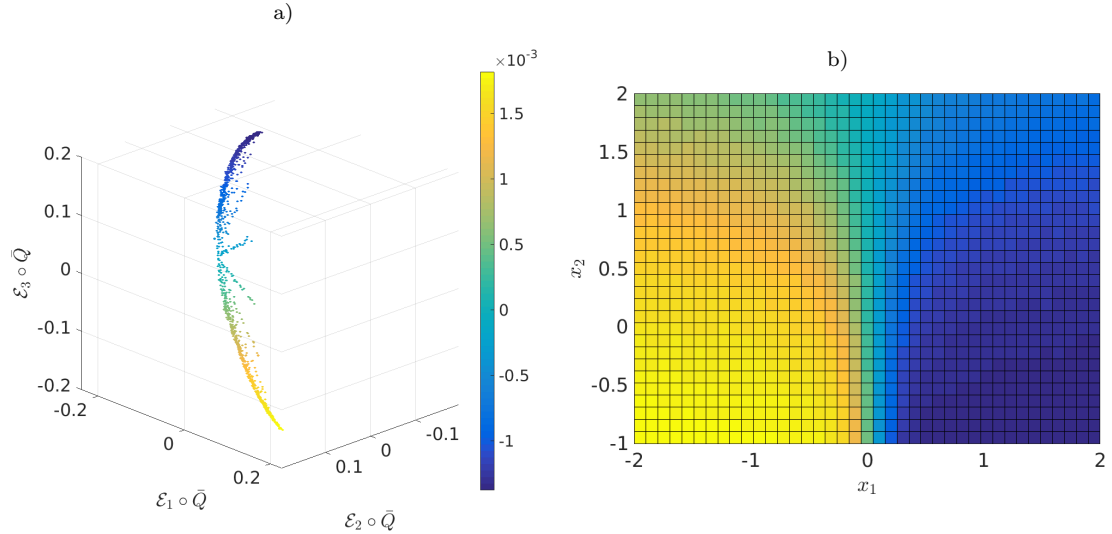


Figure 9: a) The embedded grid points colored according to the first diffusion map coordinate. b) The final reaction coordinate $\bar{\xi}$.

Legoll and Lelièvre [45] show that the effective dynamics based on the reaction coordinate

$$\xi^*(x) = x_1 \exp(-2x_2)$$

accurately reproduces the long-time dynamics of the full process — although they do not use dominant eigenvalues of the transfer operator in their argumentation. It is easy to verify that the level sets of ξ^* traverse the transition path orthogonally. Figure 10 compares the level sets of $\bar{\xi}$ and ξ^* . While the two reaction coordinates have different absolute values, their contour lines coincide well. As the projection operator P_ξ only depends on the level sets of ξ , the projected transfer operators \mathcal{T}_ξ^t and $\mathcal{T}_{\xi^*}^t$ should be similar as well.

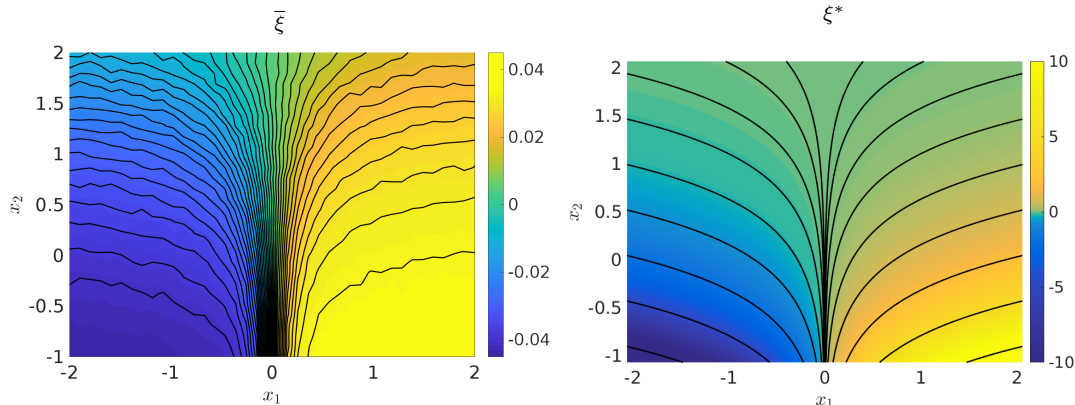


Figure 10: Selected contour lines (black) of the newly identified reaction coordinate $\bar{\xi}$ and the reference reaction coordinate ξ^* .

Projected eigenvalue error. To conclude this example, we compute the dominant spectrum of the projected transfer operator and compare it to the spectrum of the full transfer operator. To discretize \mathcal{T}_ξ^t , we use a simple Ulam-type discretization scheme based on a long equilibrated trajectory of the full dynamics. Recall from Section 3.2 that, although \mathcal{T}_ξ^t formally acts as an operator on functions over \mathbb{X} , it is constant along level sets of $\bar{\xi}$, and thus can be treated as an operator on functions over \mathbb{R}^r . For completeness, we state the rough outline of an algorithm that we used to approximate \mathcal{T}_ξ^t . An introduction to Ulam- and other Galerkin-type discretization schemes for transfer operators can be found, e.g., in [37].

1. Compute points $\bar{\mathbb{X}} := \{\Phi_{(k\tau)}x_0 \mid k = 1, \dots, N\}$, a discrete trajectory with step size τ of the full phase space dynamics that adequately samples the invariant density ϱ .
2. Compute the reaction coordinate $\bar{\xi}$ on the points $\bar{\mathbb{X}}$.
3. Divide the neighborhood of $\bar{\xi}(\bar{\mathbb{X}})$ into boxes or other suitable discretization elements $\{\mathbb{A}_1, \dots, \mathbb{A}_N\}$ and sample the boxes from the trajectory, i.e. compute

$$\bar{\mathbb{X}}_i := \{x \in \bar{\mathbb{X}} \mid \bar{\xi}(x) \in \mathbb{A}_i\} .$$

- Count the time- t -transitions within $\bar{\mathbb{X}}$ between the boxes (where t is a multiple of τ), i.e. compute the matrix

$$(T_{\bar{\xi}}^t)_{ij} := \#\{x \in \bar{\mathbb{X}}_i \mid \Phi_t x \in \bar{\mathbb{X}}_j\}.$$

- After row-normalization, the eigenvalues of $T_{\bar{\xi}}^t$ approximate the point spectrum of $\mathcal{T}_{\bar{\xi}}^t$.

Remark 5.1. Note that the equilibrated trajectory $\bar{\mathbb{X}}$ is typically unavailable for more complex systems. In practice, one would replace steps 1 and 2 by directly computing a reduced trajectory $\bar{\mathbb{Z}} = \{z_1, \dots, z_N\} \subset \mathbb{R}^r$ whose statistics approximate that of $\xi(\bar{\mathbb{X}})$. The formulation of a reduced numerical integration scheme to realize this is currently work in progress (see the fifth point in the conclusions).

For our example system, we compute $\bar{\mathbb{X}}$ as a $N = 10^6$ step trajectory with step size $\tau = 10^{-2}$ using the Euler-Maruyama scheme. However, to reduce the numerical effort, $\bar{\xi}$ is computed only on a subsample of $\bar{\mathbb{X}}$ (10^4 points) and extended to $\bar{\mathbb{X}}$ by nearest-neighbor interpolation. On $\bar{\mathbb{X}}$, the image of the $\bar{\xi}$ is contained in the interval $[-0.04, 0.04]$, which we discretize into $M = 40$ subintervals of equal length. The spectrum of the full transfer operator \mathcal{T}^t was computed using the standard Ulam method over a 40×30 uniform box discretization of the domain $[-2, 2] \times [-1, 2]$. With the choice $t = 1$ for the lag time, the spectral gap is clearly visible.

We observe in Figure 11 that the eigenvalues of $\mathcal{T}_{\bar{\xi}}^t$ and \mathcal{T}^t are in excellent agreement. Not only the dominant eigenvalues λ_0, λ_1 are approximated well (as predicted by Lemma 3.5), but also the further subdominant eigenvalues that are not covered by our theory. In particular, the reaction coordinate $\bar{\xi}$ provides a better approximation to the spectrum of \mathcal{T}^t than other, manually chosen reaction coordinates: Figure 11 also shows the eigenvalues of the projected transfer operator associated with the reaction coordinates

$$\zeta_1(x) := x_1 \quad \text{and} \quad \zeta_2(x) := x_1 + x_2.$$

We see that these are consistently outperformed by the computed reaction coordinate $\bar{\xi}$ (although it appears that ζ_1 already is quite a good reaction coordinate).

5.2. Circular potential

Let us now compute the reaction coordinates for the multi-well diffusion process described in Example 4.1. The corresponding k -well potential is defined as

$$V(x) = \cos(k \arctan(x_2, x_1)) + 10 \left(\sqrt{x_1^2 + x_2^2} - 1 \right)^2.$$

We use $k = 7$, for which the potential is shown in Figure 2a. The potential as well as the dominant eigenvalues of the corresponding transfer operator clearly indicate the existence of seven metastable sets, yet a typical longtime trajectory, shown in Figure 12a,

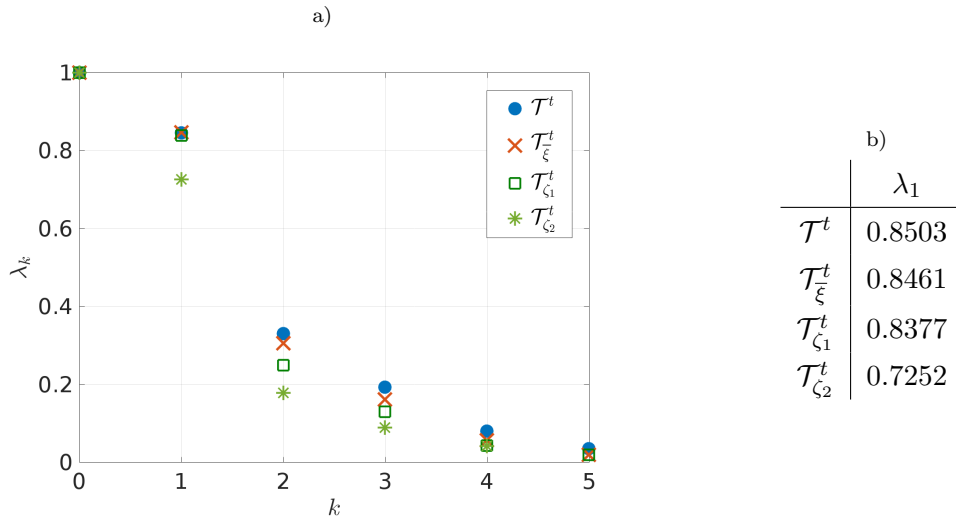


Figure 11: a) Comparison of the two dominant and first four non-dominant eigenvalues of the full transfer operator \mathcal{T}^t and the projected transfer operators $\mathcal{T}_{\bar{\xi}}^t, \mathcal{T}_{\zeta_1}^t, \mathcal{T}_{\zeta_2}^t$. b) Detailed comparison of the second eigenvalue of the various transfer operators.

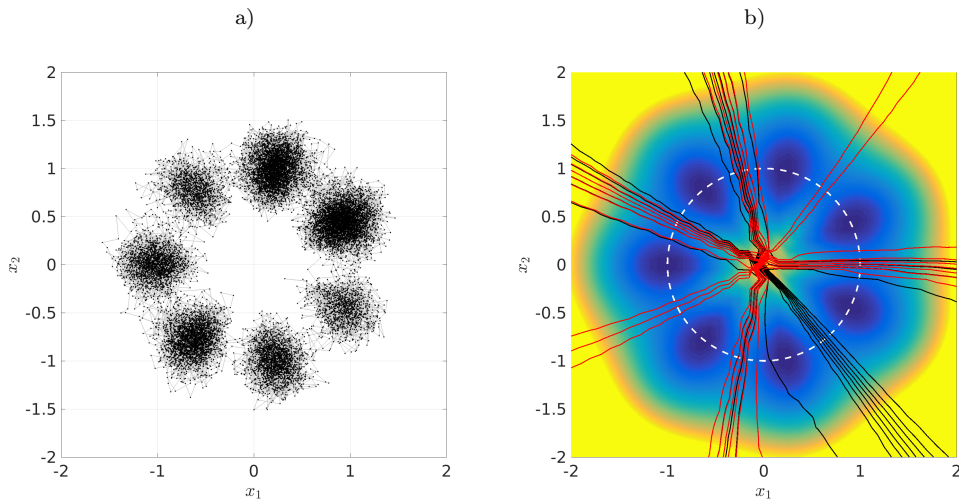


Figure 12: a) Longtime trajectory of the diffusion process with the circular seven-well potential. b) The contour lines of $\bar{\xi}_1$ (black) and $\bar{\xi}_2$ (red) show that $\bar{\xi}$ is almost constant on the metastable sets, but resolves the transition regions well.

suggests a one-dimensional transition path, the unit circle \mathbb{B}_1 . We demonstrate that with our method, a reaction coordinate of minimal dimension can be computed.

We again choose the inverse temperature $\beta = 0.5$ and perform the same analysis as in the previous subsection. For this system, a time scale gap between $t_6 \approx 1.53$ and $t_7 \approx 0.05$ can be found. We thus choose the intermediate time scale $t = 0.1$. Since we again expect a one-dimensional transition path, the three observables (29) are used for the embedding of \mathbb{M} . We use the grid points of a 40×40 grid, denoted again by $\overline{\mathbb{X}}$, over the region $[-2, 2] \times [-2, 2]$ as our test points.

The individual components of the embedding $\mathcal{E} \circ \overline{\mathcal{Q}}$ are shown in Figure 13. The embedded grid points, seen as the individual points in Figure 14a, seem to concentrate around a one-dimensional circular manifold and thus reveal the one-dimensional nature of the reaction coordinate. Although slightly unintuitive, the diffusion maps algorithm now identifies *two* significant diffusion map components, as shown in Figure 14a. The reason is that the circular manifold cannot be embedded into \mathbb{R}^1 , so that a two-component coordinate is necessary to parametrize it. Figure 12b shows some contour lines (of equidistant values) of the two components of $\bar{\xi}$. We see that $\bar{\xi}$ is almost constant on the seven metastable sets, but resolves the transition regions well.

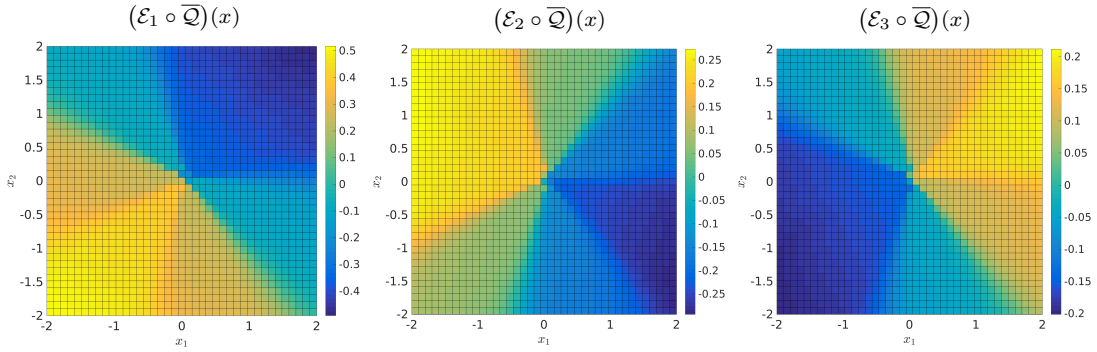


Figure 13: The individual components of the embedding $\mathcal{E} \circ \overline{\mathcal{Q}}$ on the grid points $x \in \overline{\mathbb{X}}$.

Parametrization of the dominant eigenfunctions. Next, we experimentally investigate how well the dominant eigenfunctions φ_i of \mathcal{T}^t can be parametrized by the numerically computed reaction coordinate $\bar{\xi}$. If the eigenfunctions are almost functions of $\bar{\xi}$, then by Lemma 4.2 and Corollary 3.6 the reaction coordinate is suitable to reproduce *all the dominant time scales*. To this end, we compute the dominant eigenfunctions φ_j , $j = 0, \dots, d$ by the Ulam-type Galerkin method (as in the previous example), and plot $\varphi_j(x_i)$ against $\bar{\xi}(x_i)$. Note that due to the reasons discussed above, the range of $\bar{\xi}$ is a one-dimensional manifold in \mathbb{R}^2 . If φ_j can be parametrized by $\bar{\xi}$, we expect that $\varphi_j(x_{i_1}) \approx \varphi_j(x_{i_2})$, whenever $\bar{\xi}(x_{i_1}) \approx \bar{\xi}(x_{i_2})$. The result is shown in Figure 15. We clearly see the functional dependency of the first seven (i.e., the dominant) eigenfunctions on the reaction coordinate.

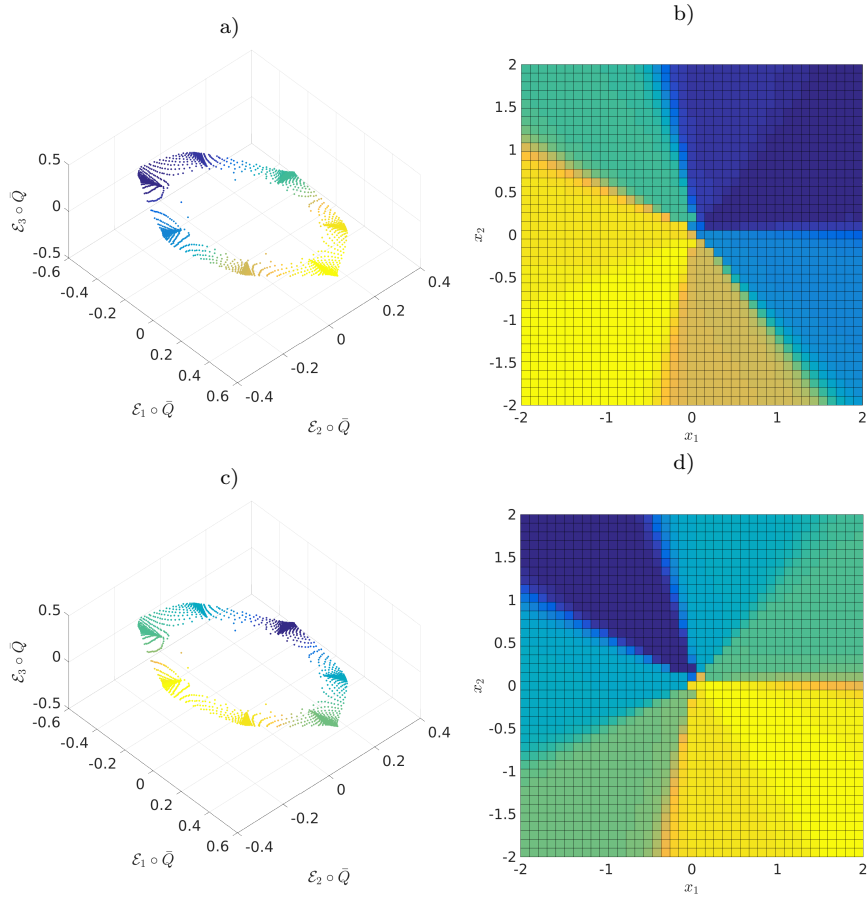


Figure 14: Left column: The embedded grid points $\mathcal{E}(\bar{Q}(\bar{\mathbb{X}}))$. The coloring shows the a) first and c) second significant diffusion map on the points. Right column: The b) first and d) second components of the final reaction coordinate $\bar{\xi}$.

Circular potential in higher dimensions. The identification of reaction coordinates is not limited to two dimensions. To show that our method can effectively find the reaction coordinates in high-dimensional systems, we extend the 7-well potential to ten dimensions by adding a quadratic term in x_3, \dots, x_{10} :

$$V(x) = \cos(7 \arctan(x_2, x_1)) + 10 \left(\sqrt{x_1^2 + x_2^2} - 1 \right)^2 + 10 \sum_{j=3}^{10} x_j^2.$$

We expect the one-dimensional circle $\{x \in \mathbb{R}^{10} \mid x_1^2 + x_2^2 = 1, x_j = 0, j = 3, \dots, 10\}$ to be the transition path and accordingly choose a three-dimensional linear observable $\eta(x) = A \cdot x$, $A \in \mathbb{R}^{3 \times 10}$, where the coefficients A_{ij} were again drawn uniformly from $[-1, 1]$.

In ten dimensions, the computation of the reaction coordinate on all points of a regular grid is no longer possible due to the curse of dimensionality, and neither is the visual-

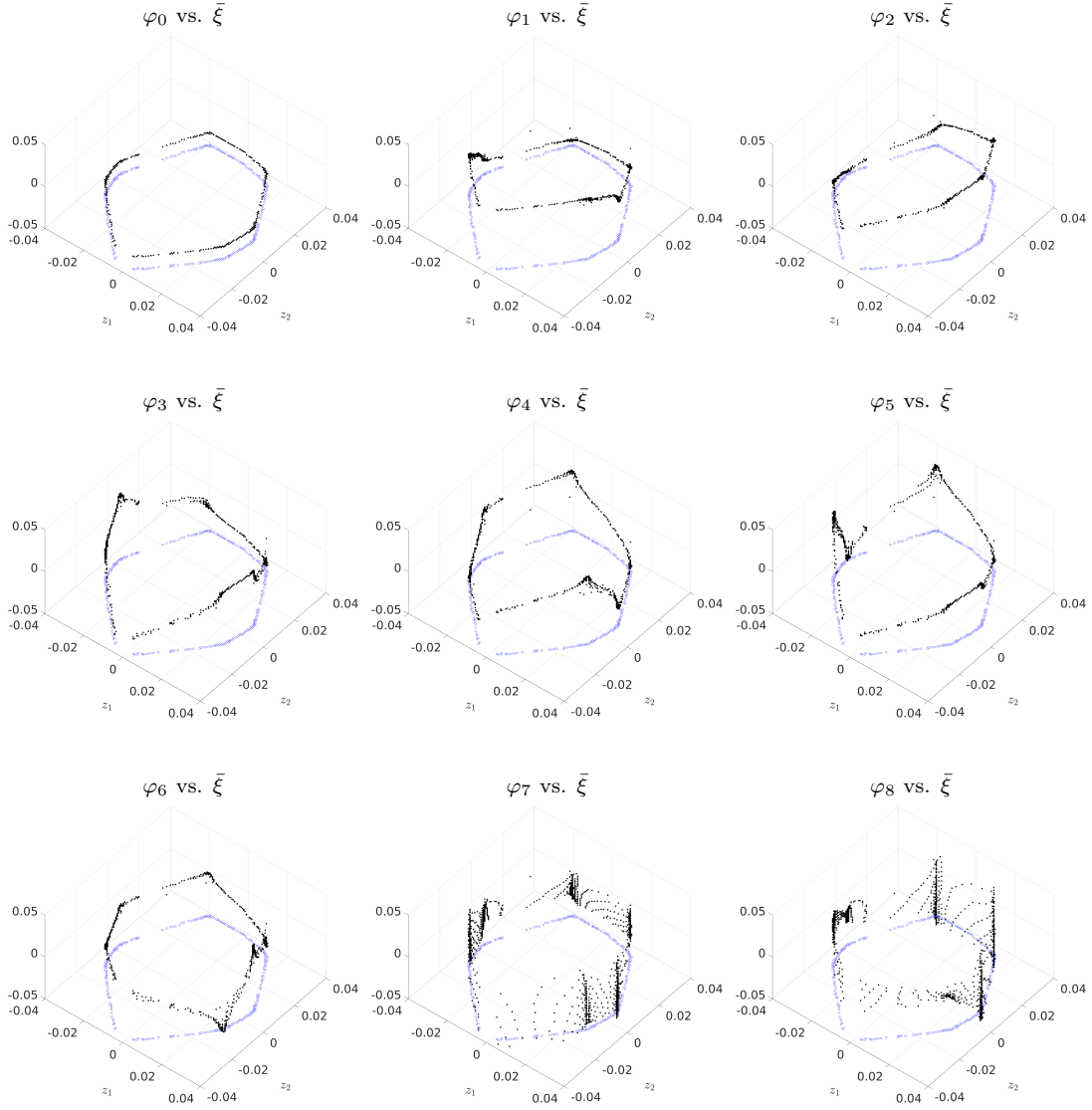


Figure 15: Black dots: The values of the first nine eigenfunctions of \mathcal{T}^t plotted against $\bar{\xi}(x_i)$, $x_i \in \bar{\mathbb{X}}$. The blue markers indicate the $\bar{\xi}(x_i)$ in the bottom plane. The seven dominant eigenfunctions (φ_0 to φ_6) seem to have a smooth dependency on $\bar{\xi}$. In contrast, the values of the non-dominant φ_7 and φ_8 vary substantially over individual level sets of $\bar{\xi}$.

ization of this grid. Instead, we compute $\bar{\xi}$ on 10^5 points sampled from the invariant measure and plot only the first three coordinates. Let this point cloud be called $\bar{\mathbb{X}}$.

Performing the standard procedure, i.e. embedding $\bar{\mathbb{X}}$ into \mathbb{R}^3 and identifying the one-dimensional core using diffusion maps, a two-component reaction coordinate is identified. Coloring the first three dimensions of $\bar{\mathbb{X}}$ by $\bar{\xi}$ (Figure 16a,b), we see that the expected reaction pathway is indeed parametrized. This pathway as well as the seven metastable states can also be recognized in a plot of the components of $\bar{\xi}(\bar{\mathbb{X}})$ plotted against each other, indicating that the information about the dominant eigenfunctions, thus the long-time jump process, is indeed retained by $\bar{\xi}$.

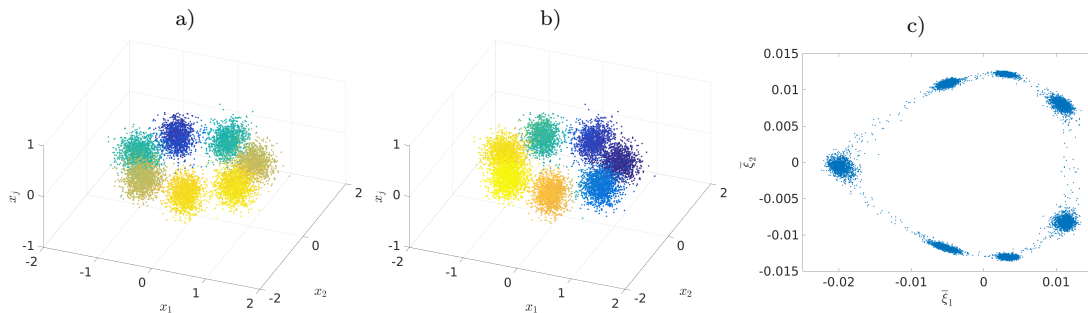


Figure 16: a) & b) The two components $\bar{\xi}_1$ and $\bar{\xi}_2$ on the sampling points $\bar{\mathbb{X}}$. The picture shows the first three dimensions of x , but is qualitatively the same when replacing x_3 by x_j , $j = 4, \dots, 10$. c) The values of $\bar{\xi}_1$ and $\bar{\xi}_2$ on $\bar{\mathbb{X}}$ plotted against each other.

5.3. Two quadruple well potentials

Our theory is based on the existence of an r -dimensional transition manifold \mathbb{M} in $L^1(\mathbb{X})$ around which the transition probability functions concentrate. In Appendix B, we argued that the existence of an r -dimensional transition path suffices to ensure the existence of \mathbb{M} . Here we illustrate how the existence of the transition path is reflected in the embedding procedure.

For this we consider the “hilly” and “flat” quadruple well potentials

$$V_1(x) = (x_1^2 - 1)^2 + (x_2^2 - 1)^2 + 5 \exp(-5(x_1^2 + x_2^2))$$

and

$$V_2(x) = 1 - \exp(-10((x_1 - 1)^2 + (x_2 - 1)^2)^2) - \exp(-10((x_1 - 1)^2 + (x_2 + 1)^2)^2) \\ - \exp(-10((x_1 + 1)^2 + (x_2 + 1)^2)^2) - \exp(-10((x_1 + 1)^2 + (x_2 - 1)^2)^2).$$

Both systems possess metastable sets around the four minima $(\pm 1, \pm 1)$, but V_1 confines its dynamics outside of the metastable sets onto a one-dimensional transition path, whereas V_2 does not impose such restrictions on the dynamics (see Figure 17). For both potentials the time $t = 1$ lies inside the slow-fast time scale gap. Assuming a

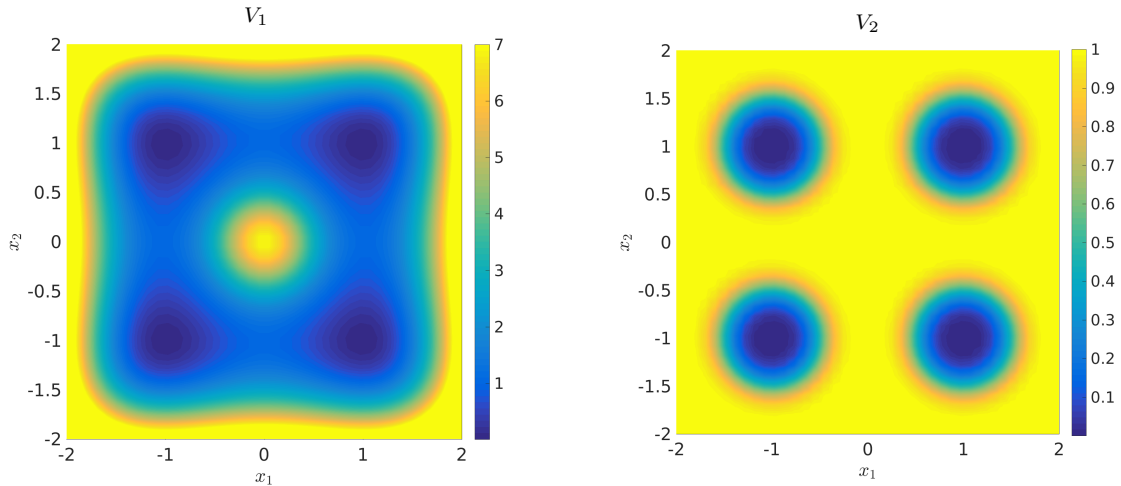


Figure 17: The two quad-well potentials V_1 and V_2 possess qualitatively different transition regions.

one-dimensional transition manifold (wrongfully for V_2), we use the three linear observables (29). A 40×40 grid on $[-2, 2] \times [-2, 2]$ is used as evaluation points for $\bar{\xi}$. The embedding of these points by $\mathcal{E} \circ \bar{Q}$ can be seen in Figure 18. We observe a one-dimensional structure in the case of the “hilly” potential V_1 , whereas the embedding points of the “flat” potential V_2 lie on a seemingly two-dimensional manifold. As these embeddings are approximately one-to-one with the respective transition manifolds \mathbb{M} , we conclude that in the case of V_1 the manifold \mathbb{M} must be one-dimensional, whereas for V_2 it is two-dimensional.

6. Conclusion

Our main contributions in this paper are:

- (a) We developed a mathematical framework to characterize good reaction coordinates for stochastic dynamical systems showing metastable behavior but no local separation of fast and slow time scales.
- (b) We showed the existence of good low-dimensional reaction coordinates under certain *dynamical* assumptions on the system.
- (c) We proposed an algorithmic approach to numerically identify good reaction coordinates and the associated low-dimension transition manifold based on local evaluation of short trajectories of the system only.

Our numerical examples show how the procedure works, that it can be used in higher dimensions, and the examples give further evidence that the dynamical assumptions from (b) are valid in many realistic cases. The application of our approach to relevant biomolecular problems, e.g. in protein folding, is ongoing work.

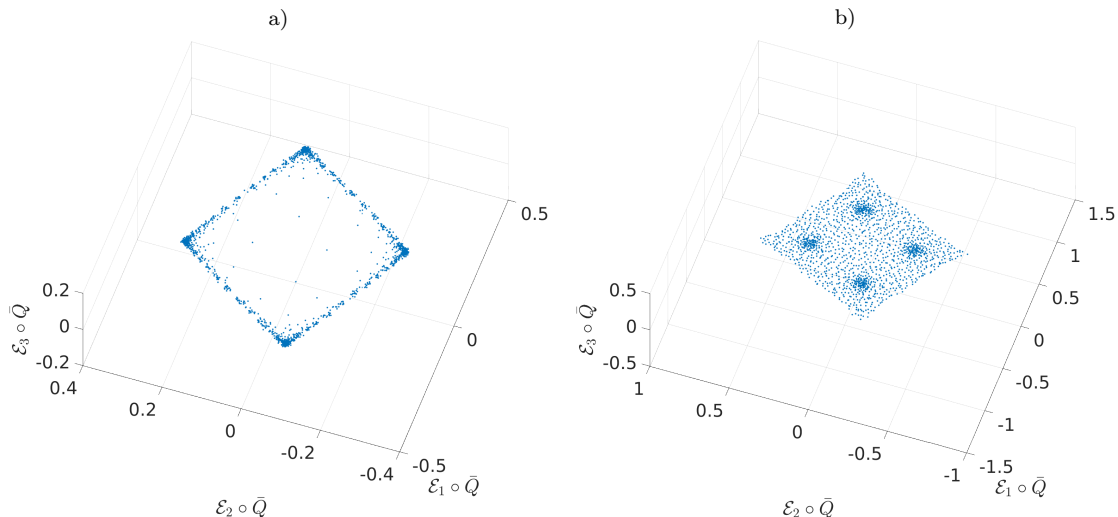


Figure 18: Embedding of the grid points for the a) “hilly” and b) “flat” four well potential. A one-dimensional structure is only visible in a), i.e. in the presence of a one-dimensional transition path.

Apart from the application to actual molecular systems, there are several open questions and challenges, which we will address in the future:

- A rigorous mathematical justification for the dynamical assumption in Definition 4.4 in terms of the potential V and the noise intensity β^{-1} in (2) would be desirable. This seems to be a demanding task, as the interplay between potential landscape and the thermal forcing is nontrivial. For $\beta^{-1} \rightarrow 0$ the problem can be handled by large deviation approaches; however, understanding increasing β^{-1} is challenging: the strength of noise increases, and additional transitions between metastable sets become more probable, as the barriers in the potential landscape become less significant, and thus the reaction coordinate may increase in dimension.
- Also related to the previous point, the choice of the correct lag time t is crucial. Choosing the time too small, the concentration of the transition densities near a low-dimensional manifold in L^1 may not have happened yet, but a too large lag time has severe consequences for the numerical expenses. If no expert knowledge of a proper lag time t is available, it has to be identified in a pre-processing step, for example using Markov State Model techniques [7].
- As discussed in the last part of Section 4.3 and in Figure 6, we need the embedding \mathcal{E} not to distort transversality close to the transition manifold \mathbb{M} too much, such that the realized reaction coordinate $\bar{\xi}$ is indeed a good one. Theoretical bounds shall be developed. This problem seems to be coupled with the problem of how to control the condition number of the embedding and its numerical realization.
- The dimension r of the reaction coordinate may not be known in advance, hence we

need an algorithmic strategy to identify this on the fly. Fortunately, once the sampling has been made, the evaluation of the embedding mapping \mathcal{E} , and finding intrinsic coordinates on the set of data points embedded in \mathbb{R}^k has a negligible numerical effort, hence different embedding dimensions k can be probed via (21). Theorem 4.10 suggests that if the identified dimension of the reaction coordinate is smaller than $k/2$, then a reaction coordinate of sufficient dimension has been found.

- To benefit from the dimensionality reduction of the reaction coordinate ξ , the dynamics that generates the reduced transfer operator \mathcal{T}_ξ^t has to be described in closed form. We are planning to employ techniques based on the Kramers–Moyal extension [84] to again receive an SDE for a stochastic process on \mathbb{R}^r .
- The embedding mapping \mathcal{E} is evaluated by Monte Carlo quadrature (24). Although Monte Carlo quadrature is known to have a convergence rate independent of the underlying dimension n of \mathbb{X} , there is still an impact of the dimension on the practical accuracy. This we shall investigate as well.

Acknowledgements

This research has been partially funded by Deutsche Forschungsgemeinschaft (DFG) through grant CRC 1114 “Scaling Cascades in Complex Systems”, Project B03 “Multilevel coarse graining of multi-scale problems”, and by the Einstein Foundation Berlin (Einstein Center ECMath).

References

- [1] J. R. Baxter and J. S. Rosenthal. Rates of convergence for everywhere-positive Markov chains. *Statistics & probability letters*, 22(4):333–338, 1995.
- [2] N. B. Becker, R. J. Allen, and P. R. ten Wolde. Non-stationary forward flux sampling. *J. Chem. Phys.*, 136(17):174118, 2012.
- [3] N. B. Becker, R. J. Allen, and P. R. ten Wolde. Non-stationary forward flux sampling. *The Journal of chemical physics*, 136(17):05B607, 2012.
- [4] R. B. Best and G. Hummer. Coordinate-dependent diffusion in protein folding. *Proceedings of the National Academy of Sciences*, 107(3):1088–1093, 2010.
- [5] A. Bittracher, P. Koltai, and O. Junge. Pseudogenerators of spatial transfer operators. *SIAM Journal on Applied Dynamical Systems*, 14(3):1478–1517, 2015.
- [6] A. Bovier, V. Gayrard, and M. Klein. Metastability in reversible diffusion processes II. Precise asymptotics for small eigenvalues. *J. Eur. Math. Soc.*, 7:69–99, 2002.
- [7] G. R. Bowman, V. S. Pande, and F. Noé, editors. *An Introduction to Markov State Models and Their Application to Long Timescale Molecular Simulation*, volume 797 of *Advances in Experimental Medicine and Biology*. Springer, 2014.

- [8] C. J. Camacho and D. Thirumalai. Kinetics and thermodynamics of folding in model proteins. *Proceedings of the National Academy of Sciences*, 90(13):6369–6372, 1993.
- [9] E. Chiavazzo, R. R. Coifman, R. Covino, C. W. Gear, A. S. Georgiou, G. Hummer, and I. G. Kevrekidis. iMapD: intrinsic map dynamics exploration for uncharted effective free energy landscapes. *arXiv preprint arXiv:1701.01513*, 2016.
- [10] J. D. Chodera and F. Noé. Markov state models of biomolecular conformational dynamics. *Current Opinion in Structural Biology*, 25:135 – 144, 2014.
- [11] G. Ciccotti, R. Kapral, and E. Vanden-Eijnden. Blue moon sampling, vectorial reaction coordinates, and unbiased constrained dynamics. *ChemPhysChem*, 6(9):1809–1814, 2005.
- [12] R. R. Coifman, I. G. Kevrekidis, S. Lafon, M. Maggioni, and B. Nadler. Diffusion maps, reduction coordinates, and low dimensional representation of stochastic systems. *Multiscale Modeling & Simulation*, 7(2):842–864, 2008.
- [13] R. R. Coifman and S. Lafon. Diffusion maps. *Applied and computational harmonic analysis*, 21(1):5–30, 2006.
- [14] M. Crosskey and M. Maggioni. ATLAS: A geometric approach to learning high-dimensional stochastic systems near manifolds. *Multiscale Modeling & Simulation*, 15(1):110–156, 2017.
- [15] E. Darve, D. Rodríguez-Gómez, and A. Pohorille. Adaptive biasing force method for scalar and vector free energy calculations. *The Journal of chemical physics*, 128(14):144120, 2008.
- [16] C. Dellago and P. G. Bolhuis. Transition path sampling and other advanced simulation techniques for rare events. In C. Holm and K. Kremer, editors, *Advanced Computer Simulation Approaches for Soft Matter Sciences III*, pages 167–233. Springer Berlin Heidelberg, Berlin, Heidelberg, 2009.
- [17] M. Dellnitz and A. Hohmann. A subdivision algorithm for the computation of unstable manifolds and global attractors. *Numerische Mathematik*, 75(3):293–317, 1997.
- [18] M. Dellnitz and O. Junge. On the approximation of complicated dynamical behavior. *SIAM J. Num. Anal.*, 36(2):491–515, 1999.
- [19] M. Dellnitz, M. H. von Molo, and A. Zießler. On the computation of attractors for delay differential equations. *Journal of Computational Dynamics*, 3(1):93–112, 2016.
- [20] N. Djurdjevac, M. Sarich, and C. Schütte. Estimating the eigenvalue error of Markov state models. *Multiscale Model. Simul.*, 10(1):61 – 81, 2012.

- [21] C. J. Dsilva, R. Talmon, C. W. Gear, R. R. Coifman, and I. G. Kevrekidis. Data-driven reduction for a class of multiscale fast-slow stochastic dynamical systems. *SIAM Journal on Applied Dynamical Systems*, 15(3):1327–1351, 2016.
- [22] R. Du, V. S. Pande, A. Y. Grosberg, T. Tanaka, and E. S. Shakhnovich. On the transition coordinate for protein folding. *The Journal of chemical physics*, 108(1):334–350, 1998.
- [23] W. E, B. Engquist, et al. The heterogenous multiscale method. *Communications in Mathematical Sciences*, 1(1):87–132, 2003.
- [24] W. E, W. Ren, and E. Vanden-Eijnden. String method for the study of rare events. *Physical Review B*, 66(5):052301, 2002.
- [25] W. E and E. Vanden-Eijnden. Towards a theory of transition paths. *J. Stat. Phys.*, 123(3):503–523, 2006.
- [26] W. E and E. Vanden-Eijnden. Transition-path theory and path-finding algorithms for the study of rare events. *Annu. Rev. Phys. Chem.*, 61(1):391–420, 2010.
- [27] A. K. Faradjian and R. Elber. Computing time scales from reaction coordinates by milestoning. *J. Chem. Phys.*, 120:10880–10889, 2004.
- [28] A. K. Faradjian and R. Elber. Computing time scales from reaction coordinates by milestoning. *The Journal of chemical physics*, 120(23):10880–10889, 2004.
- [29] H. Federer. *Geometric measure theory*, volume 1996. Springer New York, 1969.
- [30] M. Freidlin and A. D. Wentzell. *Random perturbations of dynamical systems*. Springer, New York, 1998.
- [31] G. Froyland, G. Gottwald, and A. Hammerlindl. A computational method to extract macroscopic variables and their dynamics in multiscale systems. *SIAM Journal on Applied Dynamical Systems*, 13(4):1816–1846, 2014.
- [32] G. Froyland, G. A. Gottwald, and A. Hammerlindl. A trajectory-free framework for analysing multiscale systems. *Physica D: Nonlinear Phenomena*, 328:34–43, 2016.
- [33] W. Huisinga, S. Meyn, and C. Schütte. Phase transitions & metastability in Markovian and molecular systems. *The Annals of Applied Probability*, 14 (1):419–458, 2004.
- [34] B. Hunt and V. Kaloshin. Regularity of embeddings of infinite-dimensional fractal sets into finite-dimensional spaces. *Nonlinearity*, 12(5):1263–1275, 1999.
- [35] O. Junge and P. Koltai. Discretization of the Frobenius–Perron operator using a sparse Haar tensor basis: the sparse Ulam method. *SIAM Journal on Numerical Analysis*, 47(5):3464–3485, 2009.

- [36] I. G. Kevrekidis and G. Samaey. Equation-free multiscale computation: Algorithms and applications. *Annu. Rev. Phys. Chem.*, 60(1):321–344, 2009.
- [37] S. Klus, P. Koltai, and C. Schütte. On the numerical approximation of the Perron–Frobenius and Koopman operator. *Journal of Computational Dynamics*, 3(1):51–79, 2016.
- [38] S. Klus, F. Nüske, P. Koltai, H. Wu, I. Kevrekidis, C. Schütte, and F. Noé. Data-driven model reduction and transfer operator approximation. *ArXiv e-prints*, 2017.
- [39] J. B. Kruskal. Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. *Psychometrika*, 29(1):1–27, 1964.
- [40] J. B. Kruskal. Nonmetric multidimensional scaling: a numerical method. *Psychometrika*, 29(2):115–129, 1964.
- [41] S. Kumar, J. M. Rosenberg, D. Bouzida, R. H. Swendsen, and P. A. Kollman. The weighted histogram analysis method for free-energy calculations on biomolecules. I. The method. *Journal of Computational Chemistry*, 13(8):1011–1021, 1992.
- [42] A. Laio and F. L. Gervasio. 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.
- [43] A. Laio and F. L. Gervasio. Metadynamics: a method to simulate rare events and reconstruct the free energy in biophysics, chemistry and material science. *Reports on Progress in Physics*, 71(12):126601, 2008.
- [44] A. Laio and M. Parrinello. Escaping free-energy minima. *Proceedings of the National Academy of Sciences*, 99(20):12562–12566, 2002.
- [45] F. Legoll and T. Lelièvre. Effective dynamics using conditional expectations. *Nonlinearity*, 23(9):2131, 2010.
- [46] W. Li and M. A. Recent developments in methods for identifying reaction coordinates. *Molecular simulation*, 40(10-11), 2014.
- [47] J. Lu and E. Vanden-Eijnden. Exact dynamical coarse-graining without time-scale separation. *The Journal of chemical physics*, 141(4):07B619_1, 2014.
- [48] A. Ma and A. R. Dinner. Automatic method for identifying reaction coordinates in complex systems. *J. Phys. Chem. B*, 109:6769–6779, 2005.
- [49] L. Maragliano and E. Vanden-Eijnden. A temperature accelerated method for sampling free energy and determining reaction pathways in rare events simulations. *Chemical physics letters*, 426(1):168–175, 2006.
- [50] J. C. Mattingly and A. M. Stuart. Geometric ergodicity of some hypo-elliptic diffusions for particle motions. *Markov Process. Related Fields*, 8(2):199–214, 2002.

- [51] J. C. Mattingly, A. M. Stuart, and D. J. Higham. Ergodicity for SDEs and approximations: locally Lipschitz vector fields and degenerate noise. *Stochastic processes and their applications*, 101(2):185–232, 2002.
- [52] R. T. McGibbon, B. E. Husic, and V. S. Pande. Identification of simple reaction coordinates from complex dynamics. *The Journal of Chemical Physics*, 146(4):044109, 2017.
- [53] P. Metzner, C. Schütte, and E. Vanden-Eijnden. Transition path theory for Markov jump processes. *Multiscale Modeling and Simulation*, 7(3):1192–1219, 2009.
- [54] D. Moroni, T. van Erp, and P. Bolhuis. Investigating rare events by transition interface sampling. *Physica A*, 340:395–401, 2004.
- [55] D. Moroni, T. S. van Erp, and P. G. Bolhuis. Investigating rare events by transition interface sampling. *Physica A: Statistical Mechanics and its Applications*, 340(1):395–401, 2004.
- [56] B. Nadler, S. Lafon, R. R. Coifman, and I. G. Kevrekidis. Diffusion maps, spectral clustering and reaction coordinates of dynamical systems. *Applied and Computational Harmonic Analysis*, 21(1):113–127, 2006.
- [57] F. Noé and F. Nüske. A variational approach to modeling slow processes in stochastic dynamical systems. *Multiscale Model. Simul.*, 11(2):635–655, 2013.
- [58] F. Noé, C. Schütte, E. Vanden-Eijnden, L. Reich, and T. R. Weikl. Constructing the full ensemble of folding pathways from short off-equilibrium simulations. *Proceedings of the National Academy of Sciences*, 106:19011–19016, 2009.
- [59] V. S. Pande, K. Beauchamp, and G. R. Bowman. Everything you wanted to know about Markov state models but were afraid to ask. *Methods*, 52(1):99–105, 2010.
- [60] G. Pavliotis and A. Stuart. *Multiscale methods: averaging and homogenization*. Springer Science & Business Media, 2008.
- [61] G. Pérez-Hernández, F. Paul, T. Giorgino, G. De Fabritiis, and F. Noé. Identification of slow molecular order parameters for Markov model construction. *J. Chem. Phys.*, 139(1):015102, 2013.
- [62] Z. D. Pozun, K. Hansen, D. Sheppard, M. Rupp, K.-R. Müller, and G. Henkelman. Optimizing transition states via kernel-based machine learning. *The Journal of Chemical Physics*, 136(17):174101, 2012.
- [63] W. Ren, E. Vanden-Eijnden, P. Maragakis, and W. E. Transition pathways in complex systems: Application of the finite-temperature string method to the alanine dipeptide. *The Journal of chemical physics*, 123(13):134109, 2005.
- [64] J. C. Robinson. A topological delay embedding theorem for infinite-dimensional dynamical systems. *Nonlinearity*, 18(5):2135–2143, 2005.

- [65] M. Sarich, F. Noé, and C. Schütte. On the approximation quality of Markov state models. *Multiscale Model. Simul.*, 8(4):1154 – 1177, 2010.
- [66] T. Sauer, J. A. Yorke, and M. Casdagli. Embedology. *Journal of Statistical Physics*, 65(3):579–616, 1991.
- [67] M. J. Schervish and B. P. Carlin. On the convergence of successive substitution sampling. *Journal of Computational and Graphical statistics*, 1(2):111–127, 1992.
- [68] C. Schütte, A. Fischer, W. Huisinga, and P. Deuffhard. A direct approach to conformational dynamics based on hybrid Monte Carlo. *J. Comput. Phys.*, 151(1):146–168, 1999.
- [69] C. Schütte, F. Noé, J. Lu, M. Sarich, and E. Vanden-Eijnden. Markov state models based on milestoning. *J. Chem. Phys.*, 134(20), 2011.
- [70] C. Schütte and M. Sarich. *Metastability and Markov State Models in Molecular Dynamics: Modeling, Analysis, Algorithmic Approaches*. Courant Lecture Notes in Mathematics. American Mathematical Society, 2014.
- [71] A. Singer, R. Erban, I. G. Kevrekidis, and R. R. Coifman. Detecting intrinsic slow variables in stochastic dynamical systems by anisotropic diffusion maps. *Proceedings of the National Academy of Sciences*, 106(38):16090–160955, 2009.
- [72] N. Socci, J. N. Onuchic, and P. G. Wolynes. Diffusive dynamics of the reaction coordinate for protein folding funnels. *The Journal of chemical physics*, 104(15):5860–5868, 1996.
- [73] F. Takens. Detecting strange attractors in turbulence. In *Springer Lecture Notes in Mathematics*, volume 898, pages 366–381. Springer, 1981.
- [74] J. B. Tenenbaum, V. D. Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.
- [75] G. M. Torrie and J. P. Valleau. Nonphysical sampling distributions in Monte Carlo free-energy estimation: Umbrella sampling. *J. Comput. Phys.*, 23(2):187–199, 1977.
- [76] G. M. Torrie and J. P. Valleau. Nonphysical sampling distributions in monte carlo free-energy estimation: Umbrella sampling. *Journal of Computational Physics*, 23(2):187–199, 1977.
- [77] L. N. Trefethen and M. Embree. *Spectra and pseudospectra: the behavior of non-normal matrices and operators*. Princeton University Press, 2005.
- [78] E. Vanden-Eijnden. Transition path theory. *Computer Simulations in Condensed Matter Systems: From Materials to Chemical Biology Volume 1*, pages 453–493, 2006.

- [79] E. Vanden-Eijnden. On HMM-like integrators and projective integration methods for systems with multiple time scales. *Commun. Math. Sci.*, 5(2):495–505, 06 2007.
- [80] M. Weber. *Meshless Methods in Conformation Dynamics*. PhD thesis, FU Berlin, 2006.
- [81] M. Weber. A subspace approach to molecular Markov state models via a new infinitesimal generator., 2012. Habilitation thesis.
- [82] M. Weber, K. Fackeldey, and C. Schütte. Set-free Markov state model building. *J. Chem. Phys.*, 146:124133, 2017.
- [83] H. Whitney. Differentiable manifolds. *Annals of Mathematics*, 37(3):645–680, 1936.
- [84] W. Zhang, C. Hartmann, and C. Schütte. Effective dynamics along given reaction coordinates, and reaction rate theory. *Faraday Discussions*, 195:365–394, 2016.
- [85] W. Zhang and C. Schuette. Reliable approximation of long relaxation timescales in molecular dynamics. *Submitted to Entropy*, 2017.

A. Properties of P_ξ

Proof of Proposition 3.4. (a) This property has been shown by Zhang [84] as well, we include the short reasoning for completeness. The linearity of P_ξ is obvious. The property $P_\xi^2 = P_\xi$ follows from (8) by noting that $P_\xi f$ is constant on \mathbb{L}_z and that μ_z is a probability measure for every z .

(b) From (10) we have for $f, g \in L_\mu^2(\mathbb{X})$ that

$$\begin{aligned}
 \langle P_\xi f, g \rangle_\mu &= \int_{\mathbb{X}} P_\xi f(x) g(x) d\mu(x) \\
 &\stackrel{(10)}{=} \int_{\xi(\mathbb{X})} \Gamma(z) P_\xi(\widehat{g P_\xi f})(z) dz \\
 &\stackrel{(*)}{=} \int_{\xi(\mathbb{X})} \Gamma(z) \widehat{P_\xi f}(z) \widehat{P_\xi g}(z) dz, \tag{30}
 \end{aligned}$$

where (*) follows from the linearity of P_ξ , and the fact that $P_\xi f|_{\mathbb{L}_\xi(x)} = \text{const}$, thus $P_\xi(\widehat{g P_\xi f})(z) = \widehat{P_\xi f}(z) \widehat{P_\xi g}(z)$. Expression (30) is symmetric in f and g , hence it follows that $\langle P_\xi f, g \rangle_\mu = \langle f, P_\xi g \rangle_\mu$.

(c) We first prove that P_ξ is an orthogonal projection:

$$\langle P_\xi f, f - P_\xi f \rangle_\mu \stackrel{(b)}{=} \langle f, P_\xi f - P_\xi^2 f \rangle_\mu \stackrel{(a)}{=} \langle f, P_\xi f - P_\xi f \rangle_\mu = 0.$$

Thus,

$$\|f\|_{L_\mu^2}^2 = \|f - P_\xi f\|_{L_\mu^2}^2 + \|P_\xi f\|_{L_\mu^2}^2 \geq \|P_\xi f\|_{L_\mu^2}^2,$$

and the claim follows. □

B. On the existence of reaction coordinates

To motivate the existence of low-dimensional reaction coordinates, let us assume that the dynamics of consideration has $d + 1$ metastable regions $\mathbb{C}_0, \dots, \mathbb{C}_d \subset \mathbb{X}$. Let $\mathbb{C} = \bigcup_i \mathbb{C}_i$. For a selected lag time $t > 0$ we make the following two assumptions:

- 1) Fast local equilibration: If x is in (or close to) \mathbb{C}_i then we have

$$\mathcal{P}^t \delta_x \approx \varrho_i^{qs}$$

where ϱ_i^{qs} is the quasi-stationary density of the metastable core \mathbb{C}_i :

$$\lim_{s \rightarrow \infty} \mathbb{P} [\mathbf{X}_s = y \mid \tau_{\mathbb{C}_i} > s] = \varrho_i^{qs}(y) dy$$

with $\tau_{\mathbb{C}_i}$ being the (random) exit time from the set \mathbb{C}_i .

- 2) Slow transitions: The typical transition time to reach $\mathbb{C} \setminus \mathbb{C}_i$ when starting in \mathbb{C}_i is larger than t . In other words, t is such, that if the process \mathbf{X}_s transitions from x to some \mathbb{C}_i , it did not transition through some other \mathbb{C}_j with high probability.

These two assumptions essentially say that t is much larger than the fast time scales of the system, but smaller than the dominant time scales. It follows that, for any $x \in \mathbb{X}$,

$$\mathcal{P}^t \delta_x \approx \sum_{i=0}^d q_i(x) \varrho_i^{qs}, \quad \sum_{i=0}^d q_i(x) = 1,$$

where by assumption 2) the coefficients $q_i(x)$ are given by the committor functions

$$q_i(x) = \mathbb{P} [\mathbf{X}_t \text{ reaches } \mathbb{C}_i \text{ before } \mathbb{C} \setminus \mathbb{C}_i \mid \mathbf{X}_0 = x].$$

We say that $\mathcal{P}^t \delta_x$ is an r -dimensional structure in $L^1(\mathbb{X})$ if there is a function $\xi: \mathbb{X} \rightarrow \mathbb{R}^r$ that jointly parametrizes all the committor functions, i.e., $q_i = \tilde{q}_i \circ \xi$ with $\tilde{q}_i: \mathbb{R}^r \rightarrow \mathbb{R}$. If this is the case, then

$$\overline{\mathcal{Q}}(x) = \mathcal{P}^t \delta_x \approx \sum_{i=0}^d \tilde{q}_i(\xi(x)) \varrho_i^{qs} =: \mathcal{Q}(x)$$

and clearly $\dim(\mathcal{Q}(\mathbb{X})) \leq r$ since $\dim(\xi(\mathbb{X})) = r$. Moreover, $r \leq d$ since we can explicitly construct $\xi: \mathbb{X} \rightarrow \mathbb{R}^d$ as $\xi = (q_1, \dots, q_d)$. This obviously parameterizes q_1, \dots, q_d , and it also parameterizes q_0 since $q_0 = 1 - \sum_{i=1}^d q_i$.

However, r may also be smaller than d . As an example, consider the potential with 4 minima shown in Figure 19 on the left. At low temperatures, the ‘‘hilly’’ potential energy landscape confines all transitions between the minima $\mathbb{C}_0, \dots, \mathbb{C}_3$ to a narrow region close to the red square connecting the four minima. Figure 19 shows the level sets of q_0 , the level sets of the other committors are given by the rotational symmetry of the problem. All four committors can be jointly parameterized by a single coordinate

ξ which describes clockwise movement along the red square and is constant orthogonal to it. Therefore, $r = 1$. Figure 19 on the right shows the situation with a “flat” energy landscape. Transition paths between the minima are no longer confined to a one-dimensional structure, and the committor level sets are more complicated. We can no longer parameterize all four committors with a single coordinate ξ , so $r > 1$. On the other hand, $\dim(\mathbb{X}) = 2$ so $r = 2$.

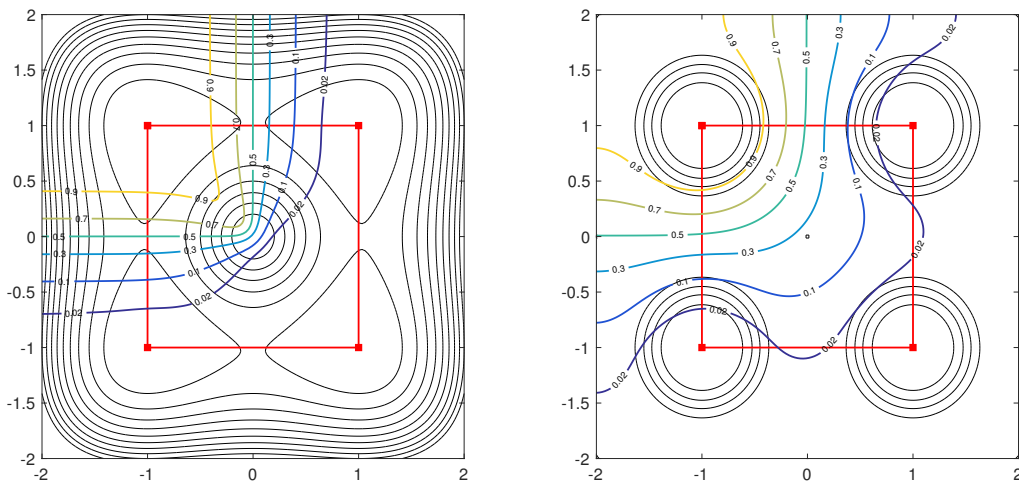


Figure 19: A potential energy landscape with four minima (black contours) and level sets of q_1 (colored contours). Left: The “hilly” landscape structure confines transition pathways to a narrow region close to the red square connecting the four minima. As a result, all committor level sets are orthogonal to this main transition path. Right: “Flat” landscape structure with more complicated level sets of the committors.

This structural difference of the potentials can also be seen when applying our algorithm to construct the reaction coordinate $\bar{\xi}$, see Figure 18 and Section 5.3.