

On Identification & Learning of Structured Latent Representations

Yingzhen Li

yingzhen.li@imperial.ac.uk

Starting with an “Ancient” Example...

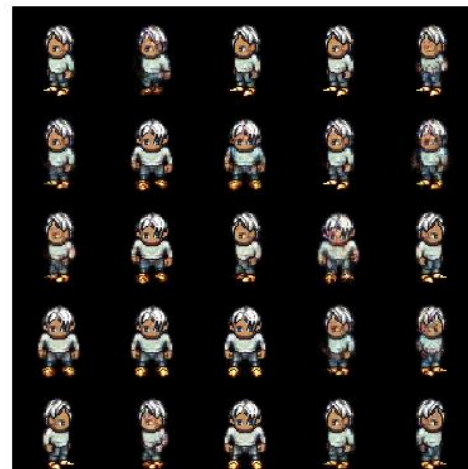
Disentangle the representation in unsupervised fashion:

- Static information (e.g., content, style)
- Temporal information (e.g., movement)

Note: no attribute labels, learned purely in an unsupervised manner.



data

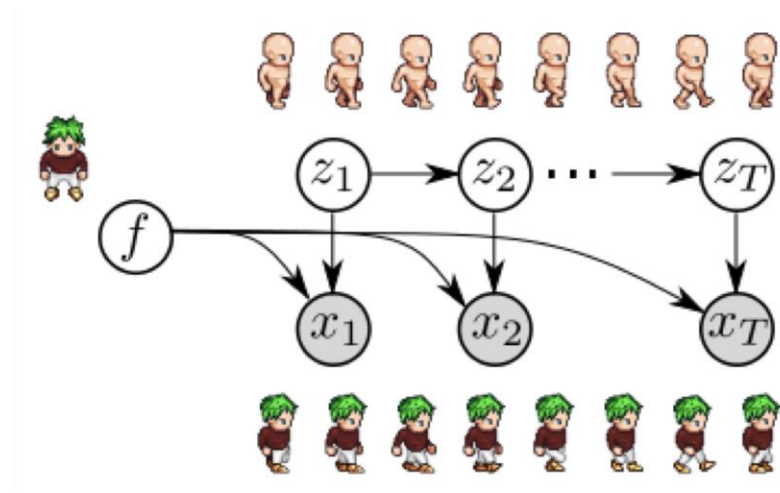


Generated (fix content)



Generated (fix dynamics)

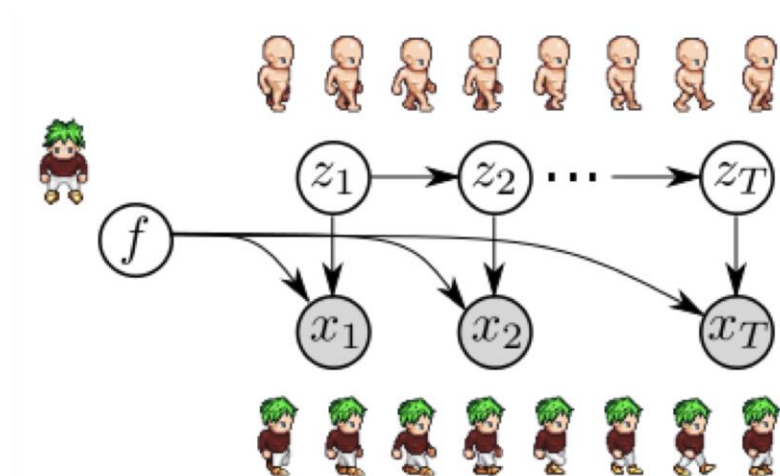
Disentangled Sequential Autoencoder



Idea:

- Build a probabilistic graphical model with f = “content” and $z_{1:T}$ = “dynamics”
- Use LSTMs to parameterise $p(z_t|z_{<t})$ and CNNs (+LSTM) to parameterise $p(x_t|f, z_t)$
- Train the model on observational data

Powerful Neural Networks Can “Cheat”

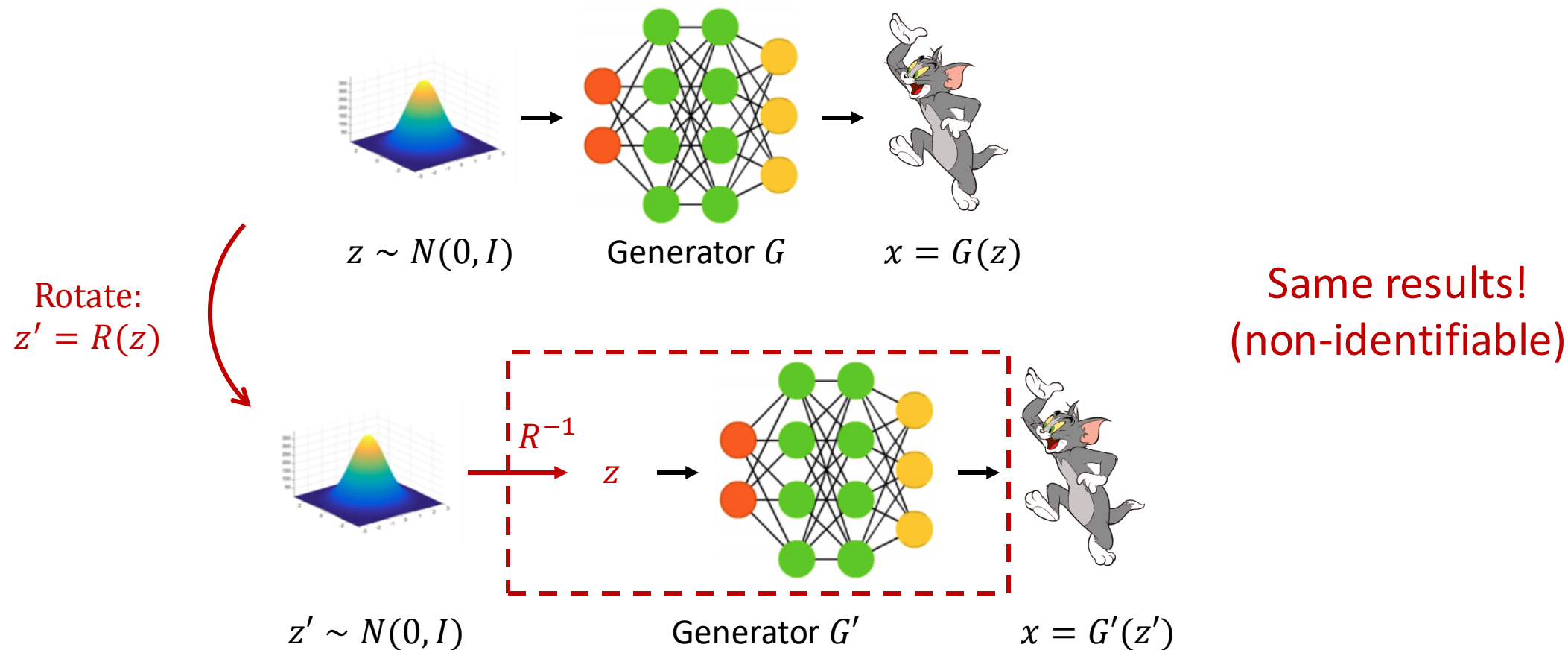


My solution back then:
Graduate student descent

Cheat in the following ways:

- The LSTM hidden cells can learn to “copy” the states
⇒ z_t captures content info
- The f variable can learn the initial condition for a deterministic dynamical system
⇒ f captures movement info

Powerful Neural Networks Can “Cheat”



Identifiability in Deep Generative Models

Workflow of causal discovery based on **identifiable DGMs**:

- Write down the model assumptions
 - E.g. $Z = g_\theta(\epsilon_1), X = f_\theta(Z) + \epsilon_2, f_\theta, g_\theta$ can be neural networks
 - This defines a model $p_\theta(X) = \int p_\theta(X|z)p_\theta(z)dz$ with parameters θ
 - Z is unobserved
- Show identifiability
 - i.e. $p_\theta(X) = p_{\theta'}(X) \Leftrightarrow f_\theta \cong f_{\theta'}, g_\theta \cong g_{\theta'}$
- Fit the model to data, and do model checking
 - If pass: use the fitted model to answer representation learning questions

Some Important Notes

- **Identifiability Proofs \neq Learning/Estimation Guarantees**
 - Assuming no model error
 - Assuming usage of consistent estimators e.g., MLE
 - Assuming abundant (e.g., infinite) amount of data
 - Assuming global optimum
- **Then why should I care about this?**
 - Identifiability as a fundamental concept of statistical inference
 - Pre-requisite for analyzing consistency of estimation
 - “Should I trust my discovery results when my deep generative model fits the data?”
 - Being able to store knowledge \neq Being able to use knowledge
 - Structured representation makes downstream use of features much easier

On the Identifiability of Switching Dynamic Models

ICML 2024

Carles Balsells Rodas¹



Yixin Wang²



Yingzhen Li¹



¹Imperial College London ²University of Michigan

Motivating Switching Dynamic Models

Regime-switching behaviour in time-series data:

“regimes”

- Complex behaviour due to switching between different **dynamical patterns**
- Within the same regime:
 - The dynamics may be stationary
 - The causal dependencies may be the same across time steps
 - The causal structure may be the same across time steps



Switching Dynamical Systems

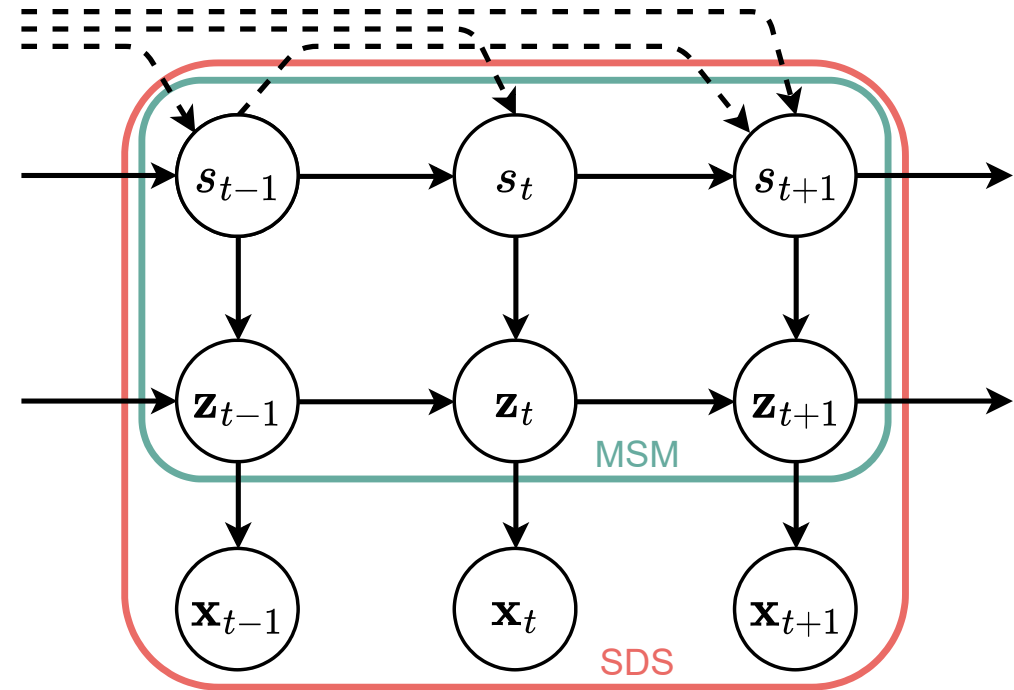
- Observations: $\mathbf{x}_t \in \mathbb{R}^n$
- Continuous latent variables: $\mathbf{z}_t \in \mathbb{R}^m$
- Discrete latent variables: $s_t \in \{1, \dots, K\}$

Switching autoregressive prior dynamics

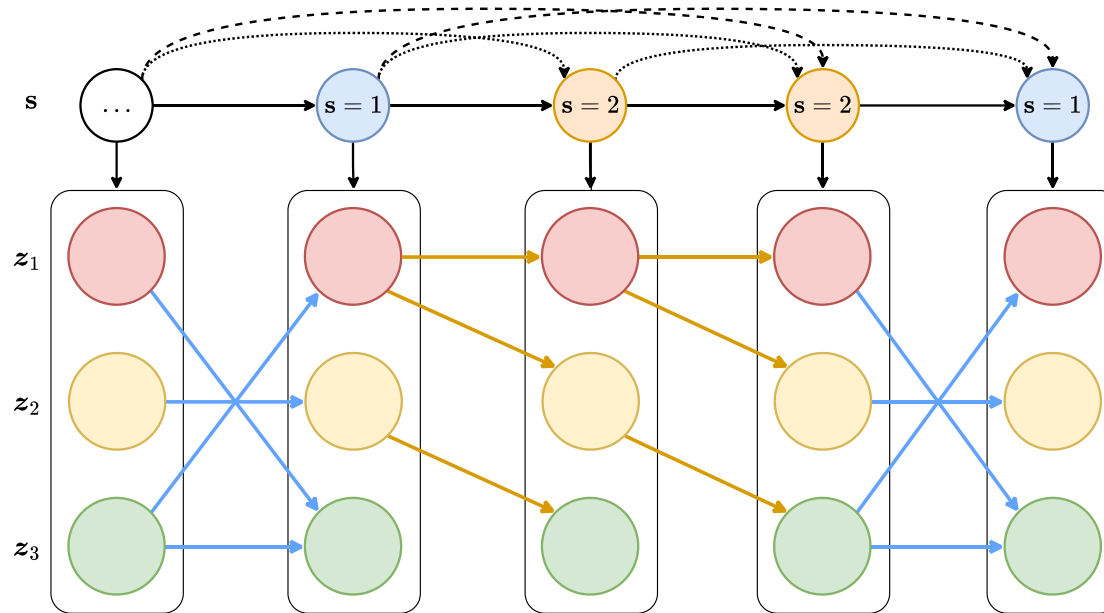
$$p_{\theta}(\mathbf{z}_{1:T}) = \sum_{\mathbf{s}_{1:T}} p_{\theta}(\mathbf{s}_{1:T}) p_{\theta}(\mathbf{z}_{1:T} | \mathbf{s}_{1:T})$$

$$p_{\theta}(\mathbf{z}_{1:T} | \mathbf{s}_{1:T}) = p_{\theta}(\mathbf{z}_1 | s_1) \prod_{t=2}^T p_{\theta}(\mathbf{z}_t | \mathbf{z}_{t-1}, s_t)$$

$$p_{\theta}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}, \mathbf{s}_{1:T}) = p_{\theta}(\mathbf{s}_{1:T}) p_{\theta}(\mathbf{z}_{1:T} | \mathbf{s}_{1:T}) \prod_{t=1}^T p_{\theta}(\mathbf{x}_t | \mathbf{z}_t)$$



Markov Switching Models



- Discrete latents: $s_t \in \{1, \dots, K\}$
- Transitions are conditionally stationary and conditional first-order Markov.

$$p(\mathbf{z}_t | \mathbf{z}_{t-1}, \dots, \mathbf{z}_1, s_t) = p(\mathbf{z}_t | \mathbf{z}_{t-1}, s_t)$$

Q: Under which **conditions** are the transitions **identifiable** from observations?

Identifiable Markov Switching Models

Theorem 1:

The following conditions render the **Markov Switching Model identifiable**¹ up to permutations:

1. *Unique indexing* for the states

$$i \neq j \iff p(\mathbf{z}_t | \mathbf{z}_{t-1}, s_t = i) \neq p(\mathbf{z}_t | \mathbf{z}_{t-1}, s_t = j)$$

2. The transition is Gaussian, and the mean and covariance are analytic in \mathbf{z}_{t-1}

$$p(\mathbf{z}_t | \mathbf{z}_{t-1}, s_t) = \mathcal{N}(\mathbf{z}_t; \mathbf{m}(\mathbf{z}_{t-1}, s_t), \Sigma(\mathbf{z}_{t-1}, s_t))$$

$\mathbf{m}(\mathbf{z}_{t-1}, s_t)$ and $\Sigma(\mathbf{z}_{t-1}, s_t)$ can be neural networks with analytic activation functions!

¹ We refer to identifiability by means of the function form rather than parameters.

Identifiable Markov Switching Models

Proof strategy: Frame the problem as a **finite mixture problem** over paths:

$$p(\mathbf{z}_{1:T}) = \sum_{i=1}^{K^T} c_i p(\mathbf{z}_{1:T} | s_1 = k_1^i, \dots, s_T = k_T^i), \quad c_i = p(\mathbf{s}_{1:T} = \{k_1^i, \dots, k_T^i\}), \quad k_t^i \in \{1, \dots, K\}$$

1. Identifiability requires linear independence over the family $\left\{ p(\mathbf{z}_{1:T} | \mathbf{s}_{1:T}) \right\}$ [1].
2. Start with $\left\{ p(\mathbf{z}_{1:2} | \mathbf{s}_{1:2}) \right\}$, then prove $T > 2$ by induction.
3. Show conditions for $p(\mathbf{z}_t | \mathbf{z}_{t-1}, s_t)$ s.t. $\left\{ p(\mathbf{z}_t | \mathbf{z}_{t-1}, s_t) p(\mathbf{z}_{1:t-1} | \mathbf{s}_{1:t-1}) \right\}$ are linearly independent. (non-parametric case)
4. Work out conditions for the Gaussian case:

$$p(\mathbf{z}_t | \mathbf{z}_{t-1}, s_t) = \mathcal{N}(\mathbf{z}_t; \mathbf{m}(\mathbf{z}_{t-1}, s_t), \Sigma(\mathbf{z}_{t-1}, s_t)) \implies \text{analytic in } \mathbf{z}_{t-1}.$$

Identifiable Markov Switching Models

Why do we need Gaussian and analytic transition functions?

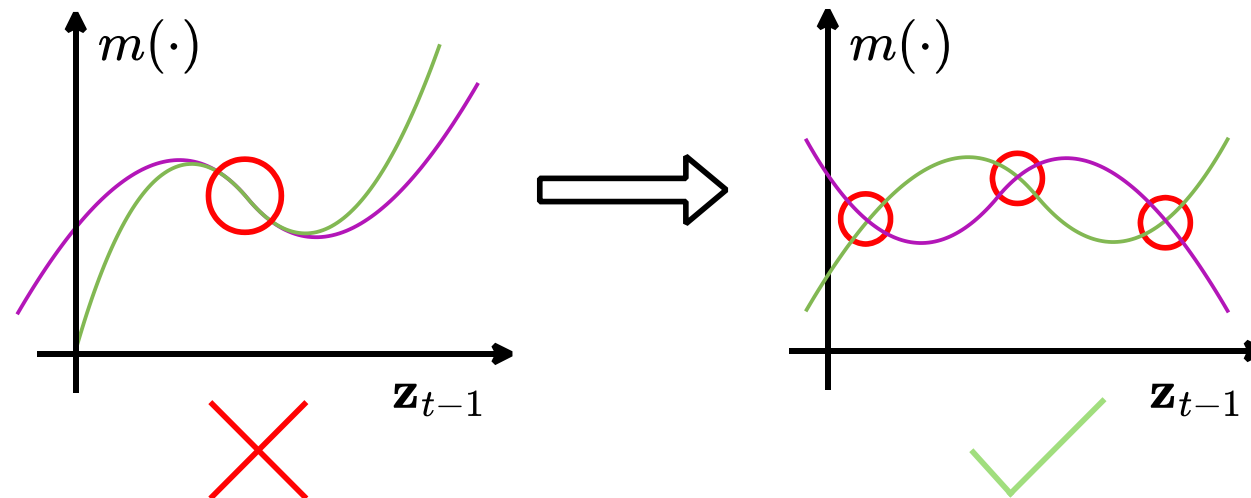
Think about linear independence of $\left\{ p(\mathbf{z}_{1:t-2}, \mathbf{z}_{t-1} | \mathbf{s}_{1:t-1}) p(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{s}_t) \right\}$

- Gaussians

$$\boldsymbol{\mu}_1 \neq \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_1 \neq \boldsymbol{\Sigma}_2 \iff p(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), p(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \text{ are linearly independent.}$$

- Analytic functions:

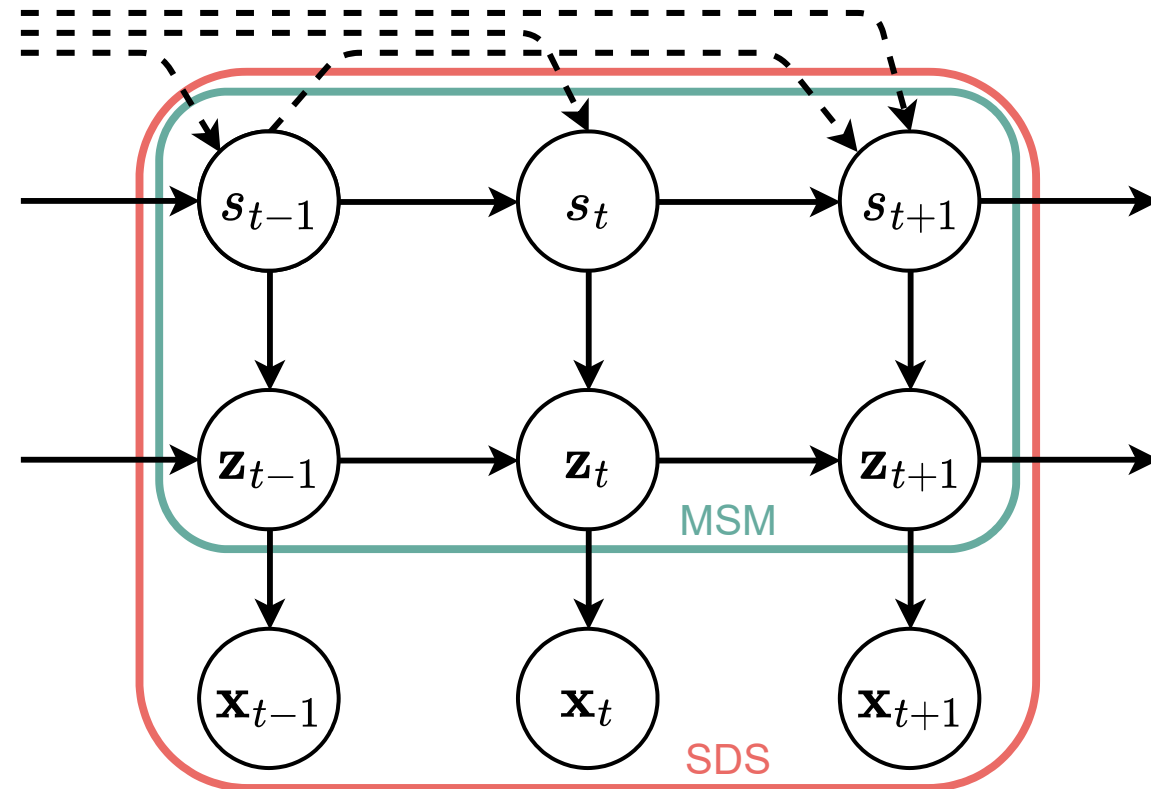
$$\mathbf{f}(\mathbf{x}) = \mathbf{g}(\mathbf{x}), \mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^d, \mu(\mathcal{X}) \neq 0 \iff \mathbf{f}(\mathbf{x}) = \mathbf{g}(\mathbf{x})$$



Identifiable Switching Dynamical Systems

- Frame the problem as an extension of iMSM and Kivva et al. (2022).

$$\mathbf{x}_t = f(\mathbf{z}_t) + \boldsymbol{\epsilon}_t, \quad \mathbf{z}_{1:T} \sim p(\mathbf{z}_{1:T}), \quad \boldsymbol{\epsilon}_t \sim p_{\boldsymbol{\epsilon}}, \quad \mathbf{x}_t \in \mathbb{R}^n, \mathbf{z}_t \in \mathbb{R}^m \quad n \geq m$$

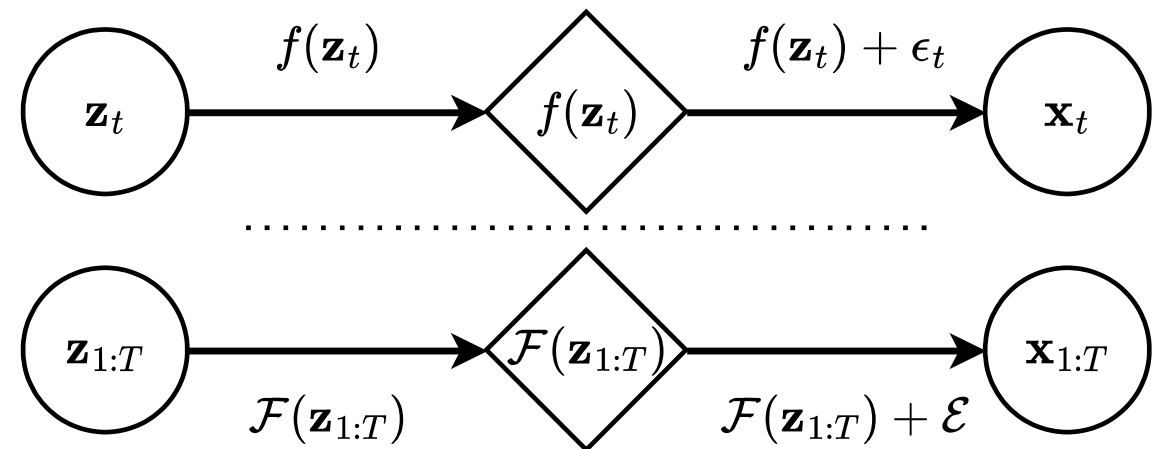


Identifiable Switching Dynamical Systems

Theorem 2: Assume an identifiable MSM.

- If f , is piece-wise linear and weakly-injective,
 - MSM is identifiable up to affine transformations.
- If f , is continuous, piece-wise linear, and injective,
 - MSM and f are identifiable up to affine transformations.

Proof sketch:



Kivva et al (2022)

Khemakhem et al. (2020)

$$\mathcal{F}(\mathbf{z}_{1:T}) = (f(\mathbf{z}_1), \dots, f(\mathbf{z}_T))^T$$

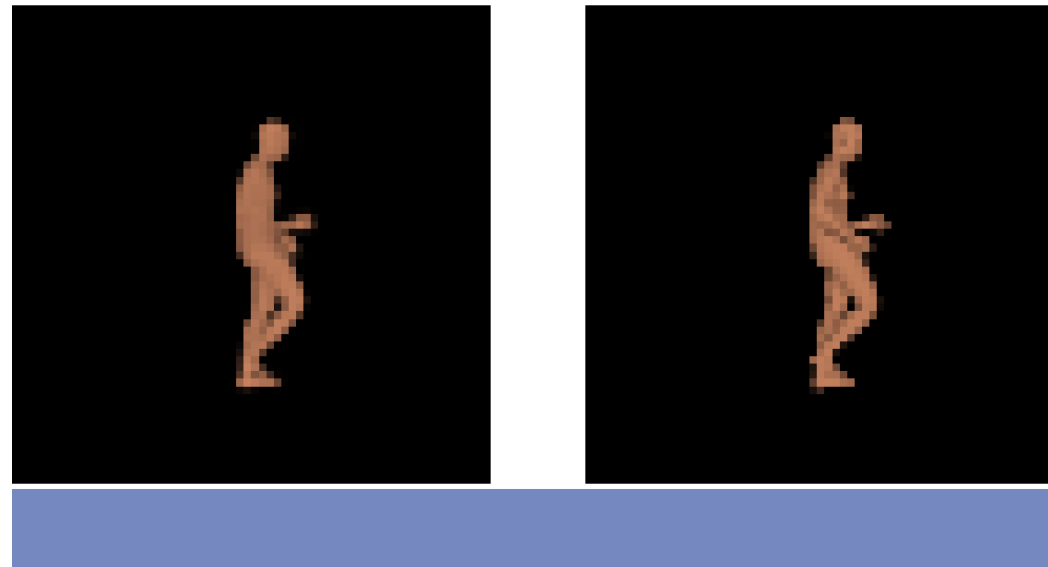
Experiments

- High-dimensional video sequences of synthetic meshes from CMU mocap data.



Reconstruction

Ground Truth

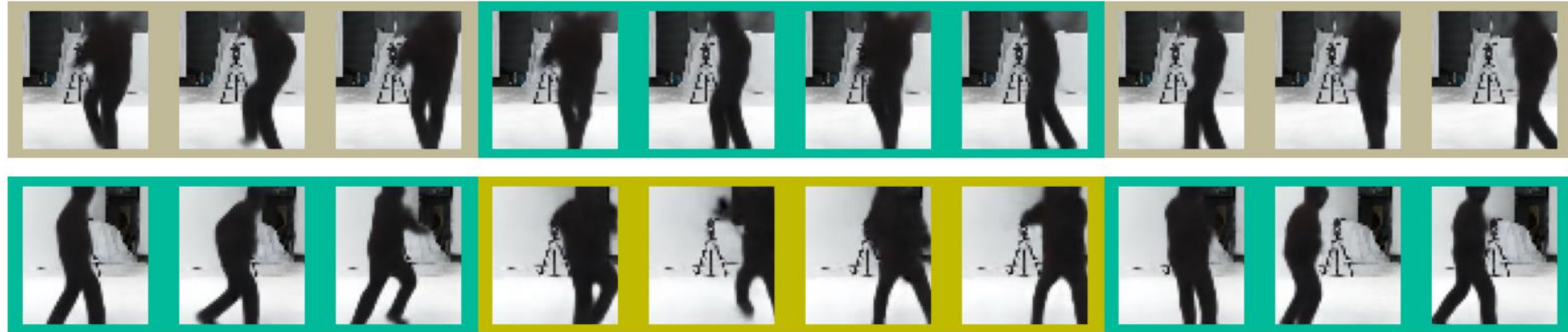


Check paper!



Experiments

- High-dimensional video sequences from the AIST Dance DB – Hip-hop.



Reconstruction

Ground Truth

Check paper!



Summary

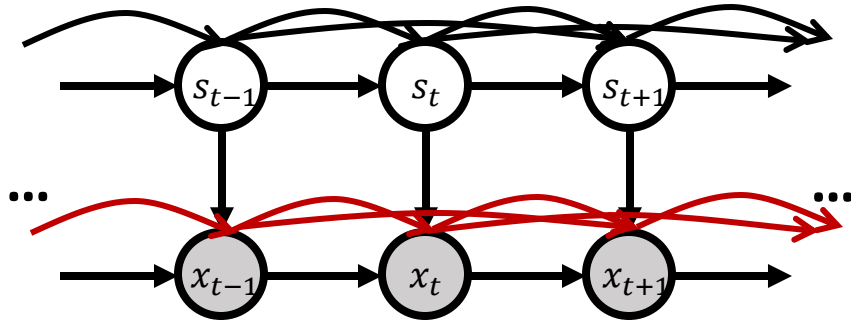
- We establish identifiability for Switching Dynamical Systems.
- Assumptions directly linked to deep generative models.
- Our proof is a major result beyond classical HMM theory.

Check paper!

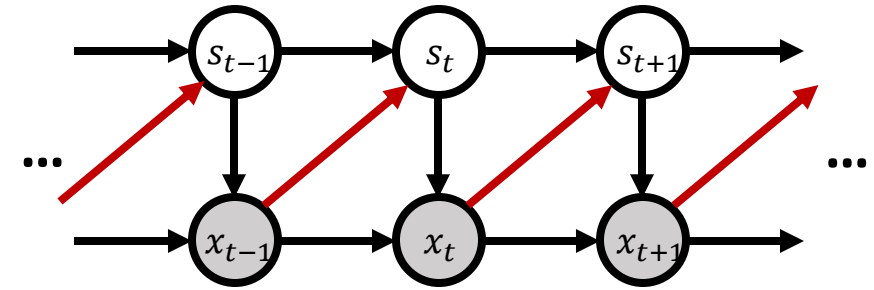


Work In-Progress & Future Extensions

- Go for more structured distributions:



higher-order Markov transitions
(e.g., neuroscience & climate data)



observation-dependent state transitions
(critical in e.g., model-based RL)

Some Important Notes (Again)

- **Identifiability Proofs \neq Learning/Estimation Guarantees**
 - Assuming no model error
 - Assuming usage of consistent estimators e.g., MLE
 - Assuming abundant (e.g., infinite) amount of data
 - Assuming global optimum
- **Can we say something closer to the practices?**

Some Questions That We Can Ask

$$Z = g_{\theta}(\epsilon_1), X = f_{\theta}(Z) + \epsilon_2, f_{\theta}, g_{\theta} \text{ can be neural networks}$$

- In VAE (i.e., DGMs learned with amortized VI) context:

Estimation consistency

If θ' maximizes ELBO
with data $p_d(x) = p_{\theta}(x)$,

Then would $f_{\theta} \cong f_{\theta'}, g_{\theta} \cong g_{\theta}'$?

- No model error
- ~~Maximum likelihood estimation~~
- Infinite amount of data
- Global optimum

Symmetry in learned representations

If θ, θ' both maximize ELBO,

Then would $f_{\theta} \cong f_{\theta'}, g_{\theta} \cong g_{\theta}'$?

- ~~No model error~~
- ~~Maximum likelihood estimation~~
- Infinite amount of data
- Global optimum

Symmetry in Learned Representations?

$$Z = g_\theta(\epsilon_1), X = f_\theta(Z) + \epsilon_2, \epsilon_2 \sim N(0, \sigma^2 I)$$

- A typical strategy for proving identifiability for DGMs:

Given: $p_\theta(x) = p_{\theta'}(x)$ and assume invertibility/injectiveness of f_θ & $f_{\theta'}$

- From noisy observations to noiseless observations
 - $p_\theta(x) = p_{\theta'}(x) \Rightarrow f_\theta(g_\theta(\epsilon_1))$ and $f_{\theta'}(g_{\theta'}(\epsilon_1))$ are equally distributed
- Removing the "volume" term (with further assumptions)
 - $\log p_\theta(z = f_\theta^{-1}(x)) + \log \left| \frac{df_\theta^{-1}(x)}{dx} \right| = \log p_{\theta'}(z = f_{\theta'}^{-1}(x)) + \log \left| \frac{df_{\theta'}^{-1}(x)}{dx} \right|$
- Analyze the equivalence class $f_\theta \cong f_{\theta'}, g_\theta \cong g_{\theta'}$

e.g., $\Phi(f_\theta^{-1}(x)) = A\Phi(f_{\theta'}^{-1}(x)) + b$ if g_θ induces an ExpFam distribution

Symmetry in Learned Representations?

- Let's look at the ELBO:

$$\begin{aligned} ELBO(x, \theta, \phi) &:= \underbrace{E_{q_\phi(z|x)}[\log p_\theta(x|z)]}_{= E_{q_\phi(z|x)} \left[\frac{1}{\sigma^2} \|x - f_\theta(z)\|_2^2 \right] + \mathcal{C}(\sigma)} + \underbrace{H[q_\phi(z|x)]}_{\text{Entropy (volume)}} + \underbrace{E_{q_\phi(z|x)}[\log p_\theta(z)]}_{\text{Stats. for inferred latent representations e.g., using ExpFam: } = E_{q_\phi(z|x)} [\langle \lambda_0, \Phi(z) \rangle] - Z(\lambda)} \end{aligned}$$

- Differences?
 - Encoder: $f_\theta^{-1}(x)$ vs $q_\phi(z|x)$
 - Additional reconstruction error term

Symmetry in Learned Representations?

$$Z = g_{\theta}(\epsilon_1), X = f_{\theta}(Z) + \epsilon_2, \epsilon_2 \sim N(0, \sigma^2 I)$$

- Possible strategy for analyzing symmetry in learned representations:

Given: $ELBO(x, \theta, \phi) = ELBO(x, \theta', \phi')$ are **optimal solutions**,

- Assume invertibility/injectiveness of f_{θ} & $f_{\theta'}$,

- Assume one of the following scenarios:

- Near deterministic regime: $\sigma \rightarrow 0$
- Optimally learned output noise variance: optimizing and share σ
- Auxiliary info available for prior only: use $p(z|u)$ instead of $p(z)$
- (each scenario needs different assumptions on $q_{\phi}(z|x)$)

$$E_{q_{\phi}(z|x)}[\log p_{\theta}(z)] = E_{q_{\phi'}(z|x)}[\log p_{\theta'}(z)]$$

Symmetry in Learned Representations?

$$E_{q_{\phi}(z|x)}[\log p_{\theta}(z)] = E_{q_{\phi'}(z|x)}[\log p_{\theta'}(z)]$$


- Why looking into this term?
 - We use q to extract latent structure in practice!
 - Example: exponential family prior and “conjugate” approximate posterior:

$$p_{\theta}(z) = \exp[\langle \lambda(\theta), \Phi(z) \rangle] - Z(\lambda(\theta))$$
$$q_{\phi}(z|x) = \exp[\langle \lambda(\phi, x), \Phi(z) \rangle] - Z(\lambda(\phi, x))$$

$$\Rightarrow E_{q_{\phi}(z|x)}[\Phi(z)] = AE_{q_{\phi'}(z|x)}[\Phi(z)] + b$$

$$\Rightarrow E_{q_{\phi^*}(z|x)}[\Phi(z)] = E_{p_{\theta^*}(z)}[\Phi(z)]$$

Take Away

- Identifiability: a fundamental question of structural representation learning
 - i.e. should you expect the learned representations to match the "structures" in data
 - The field has quite substantial advances since ~2020
 - **Our work provides strong theoretical results for (deep) switching dynamical systems**
- We need to bring identifiability theory closer to practice
 - Option 1: Getting closer to maximum likelihood
 - Option 2: Accept bias/errors in current tech, and analyze them
- Challenge : Theory for DDPMs & auto-regressive LLMs?
 - No learning for probabilistic representations
 - Studying symmetries within neural networks?

THANK YOU!

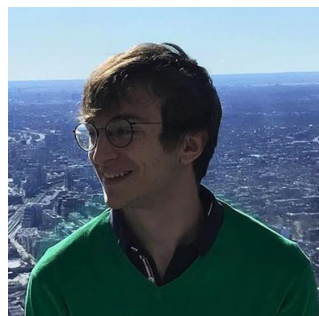
Questions? Ask now, or email:
yingzhen.li@imperial.ac.uk



PS: hiring post-docs in "AI for Chemistry"

<https://arxiv.org/abs/2305.15925>

Thanks to my awesome collaborators:



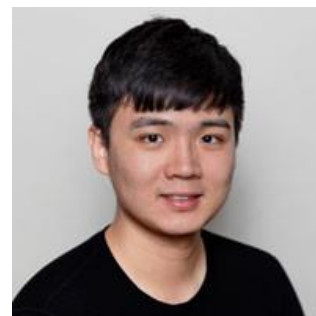
Carles Balsells-Rodas



Yixin Wang



Pedro Mediano



Ruibo Tu



Hedvig Kjellström

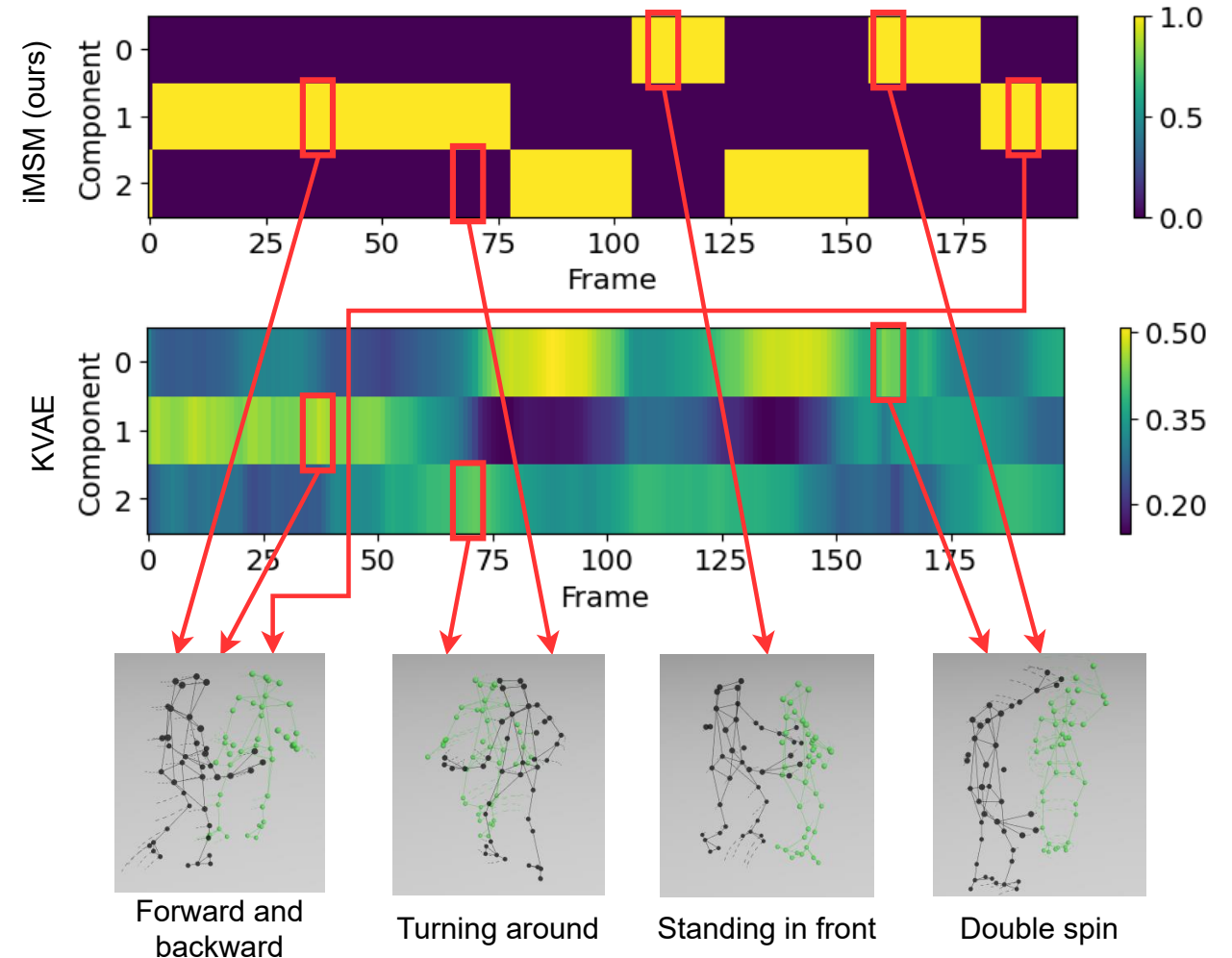


Stephan Mandt

Additional Slides

Identifiability in Markov Switching Models

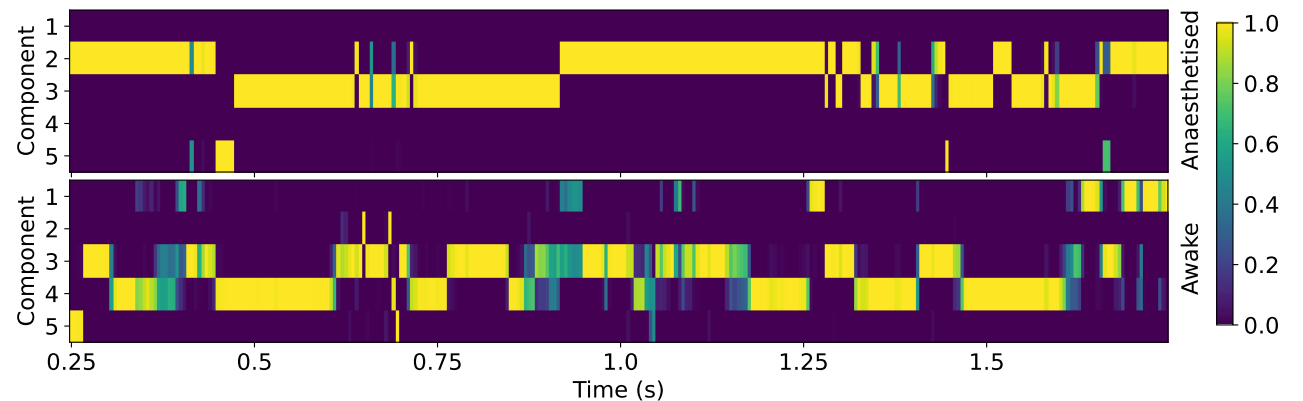
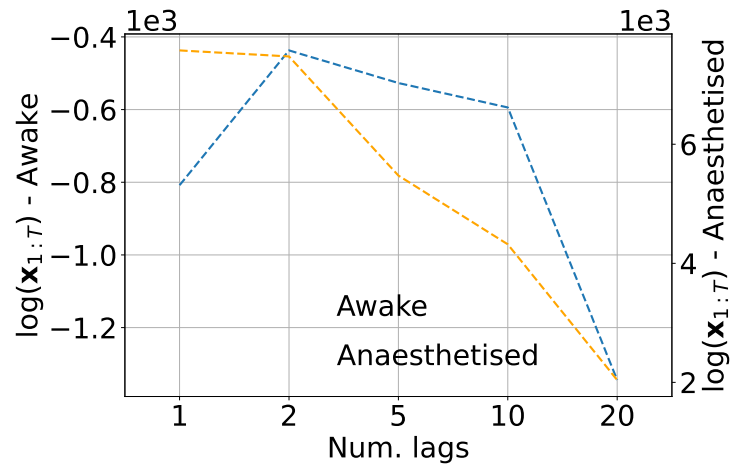
- Experiment: discovering dancing patterns
 - Data: CMU mocap
 - Estimation: Generalised EM
 - DL Baseline: KalmanVAE



Higher-Order Switching Dynamics

Neuro Activity Data Analysis:

- Recorded from Monkeys in (a) normal awake, and (b) induced anaesthetized status
- Idea: understand neuro activity by segmenting recorded signal into “regimes”

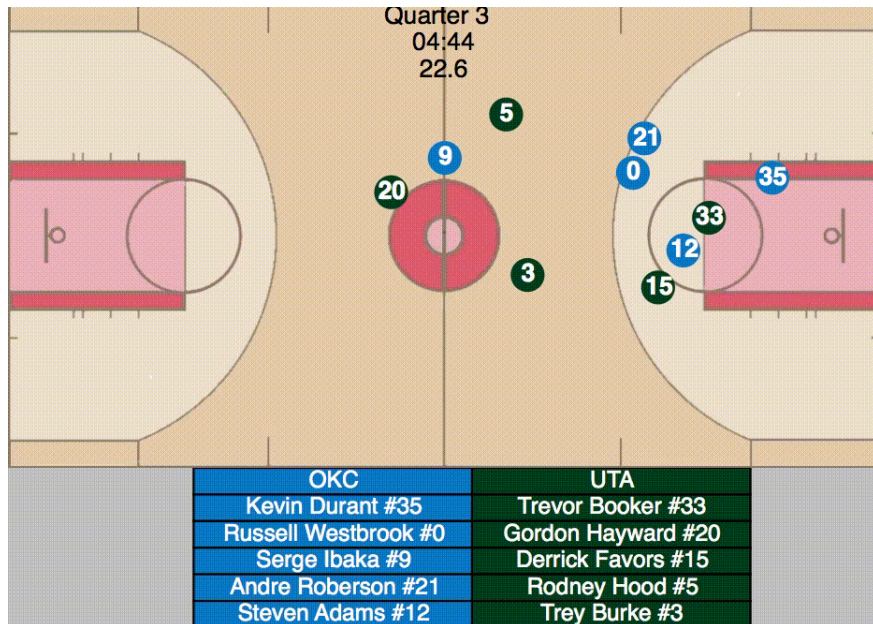


State-Dependent Causal Inference (SDCI)

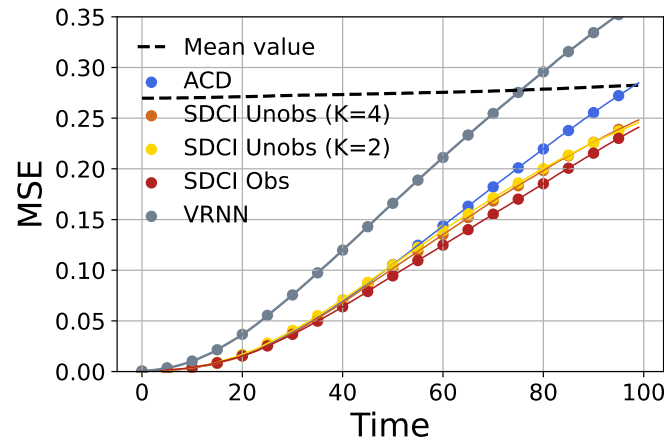
Causal discovery & sequence modelling for non-stationary time series:

Dataset: NBA player trajectories

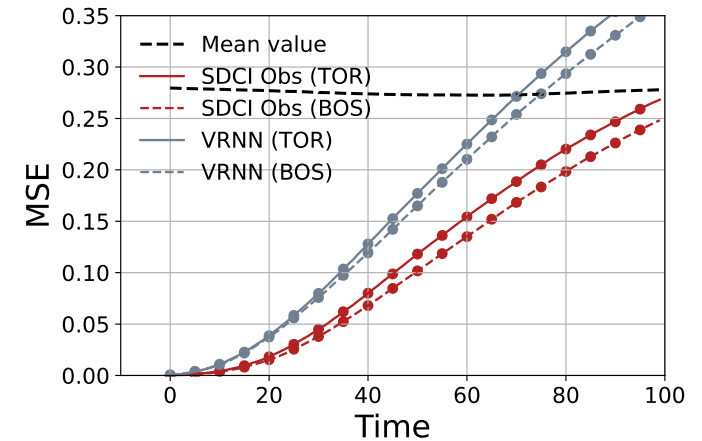
- multi-agent
- non-stationary



Forecasting error:



Train on **full data**



Train on **Boston Celtics** only

Learned hidden state visualisation:

