

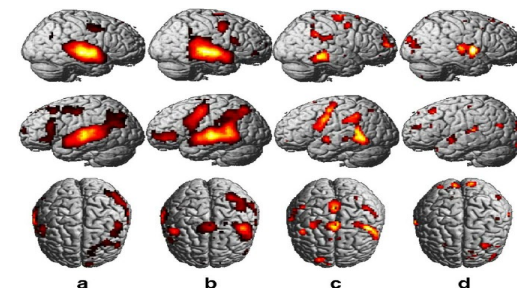
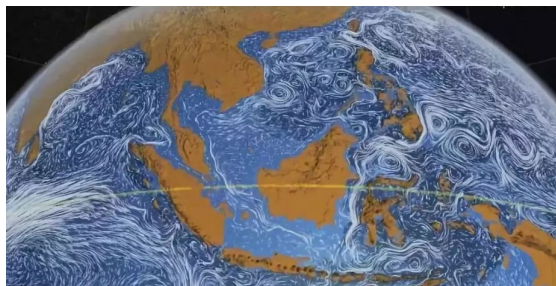
Streaming Factor Trajectory Learning for Temporal Tensor Decomposition

Shikai Fang, Xin Yu, Shibo Li, Zheng Wang, Robert Kirby, Shandian Zhe

NeurIPS 2023

Github Repo: <https://github.com/xuangu-fang/Streaming-Factor-Trajectory-Learning>

- Multi-dim Array
- Represent interactions of multiple objects/entities
- Each entry: (index1, index2..)-> value



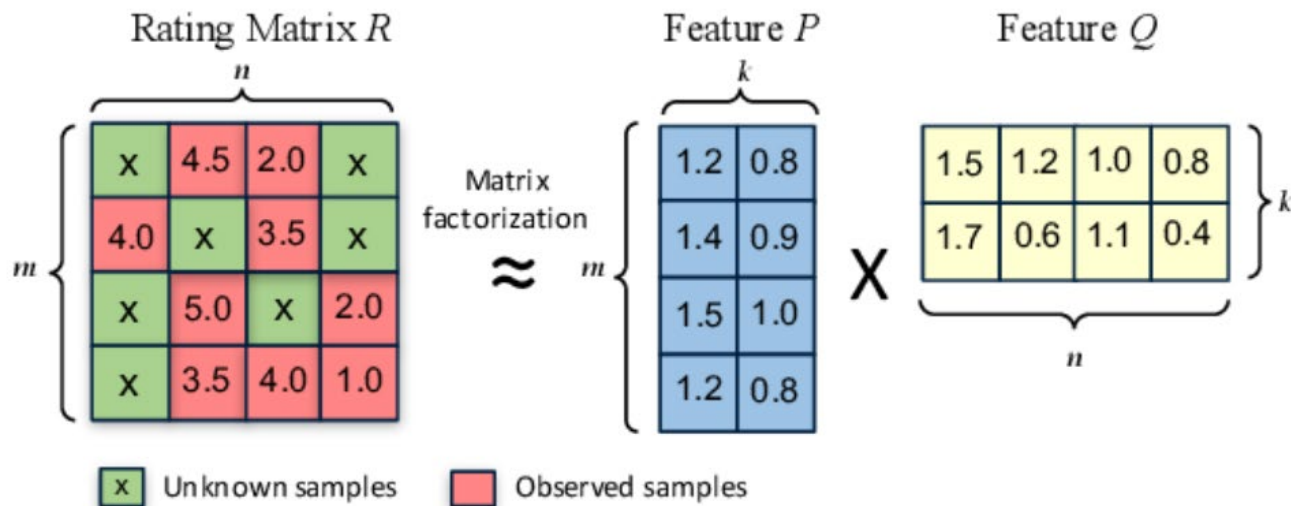
(location, region, time, climate)

(user, movie, episode)

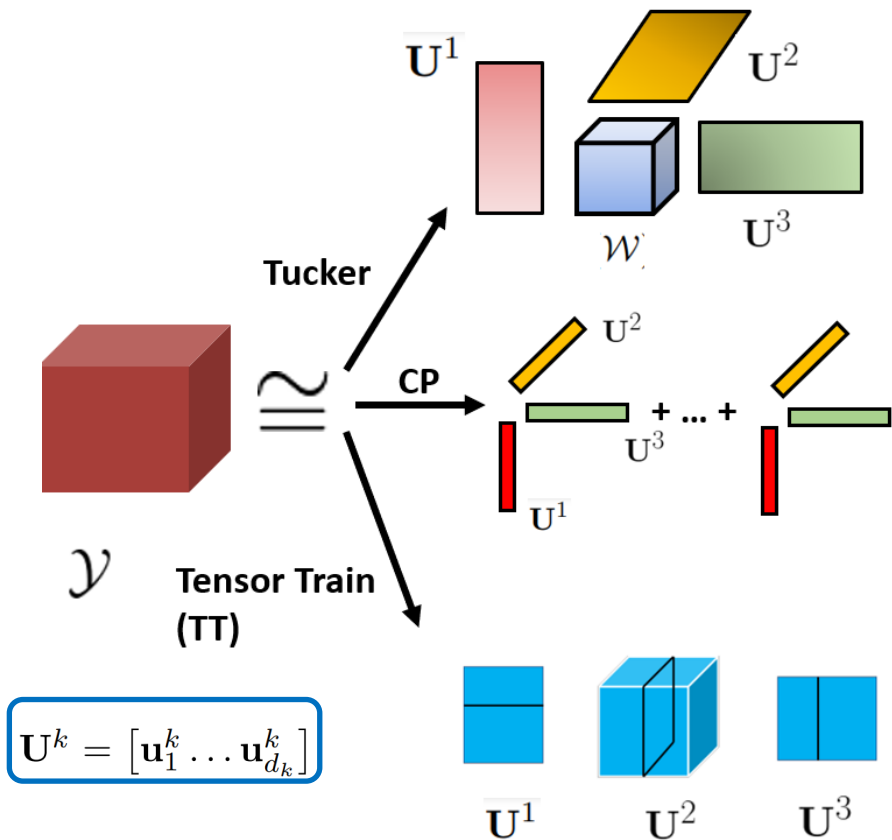
(subject, voxel, electrode)

Tensor Decomposition

- Learning representation of latent factors
- Simple case: **Collaborative Filtering** (Matrix Factorization)



CP / Tucker / TT Decomposition

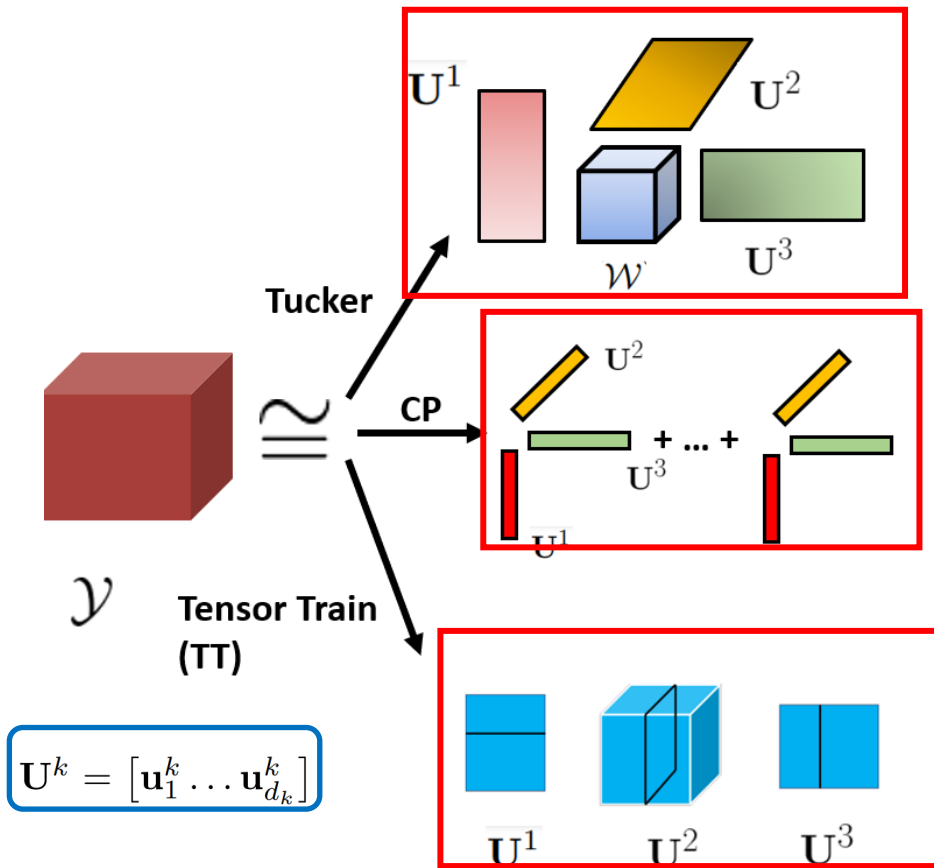


$$y_{\mathbf{i}} \approx \text{vec}(\mathcal{W})^\top \left(\mathbf{u}_{i_1}^1 \otimes \dots \otimes \mathbf{u}_{i_K}^K \right)$$

$$y_{\mathbf{i}} \approx \left(\mathbf{u}_{i_1}^1 \circ \dots \circ \mathbf{u}_{i_K}^K \right)^\top \mathbf{1}$$

$$y_{\mathbf{i}} \approx \underbrace{\mathbf{u}_{i_1}^1}_{r_0 \times r_1} \underbrace{\mathbf{u}_{i_2}^2}_{r_1 \times r_2} \dots \underbrace{\mathbf{u}_{i_K}^K}_{r_{K-1} \times r_K}$$

Bayesian Tensor Decomposition



Deterministic factors:
Vectors

Probabilistic factors :
Distributions

$$\mathcal{N}(\mathbf{u}_s^k \mid \mathbf{u}_s^k, v\mathbf{I})$$

Uncertainty counts!

Example: Bayesian CP Decomposition

Joint Distribution

Priors of factors and noise

$$p(\{y_i\}_{i \in S}, \mathcal{U}, \tau) = \text{Gam}(\tau \mid a_0, b_0) \prod_{k=1}^K \prod_{s=1}^{d_k} \mathcal{N}(\mathbf{u}_s^k \mid \mathbf{m}_s^k, v\mathbf{I})$$

$$\prod_{i \in S} \mathcal{N}(y_i \mid \mathbf{1}^\top (\mathbf{u}_{i_1}^1 \circ \dots \circ \mathbf{u}_{i_K}^K), \tau^{-1})$$

CP likelihood

Approx. Posterior

$$p(\mathcal{U}, \tau \mid \{y_i\}_{i \in S}) \approx q(\mathcal{U}, \tau) = q(\tau) \prod_{k=1}^K \prod_{s=1}^{d_k} q(\mathbf{u}_s^k)$$

Inference!

$$= \text{Gamma}(\tau \mid a^*, b^*) \prod_{k=1}^K \prod_{s=1}^{d_k} \mathcal{N}(\mathbf{u}_s^k \mid \boldsymbol{\mu}_s^{k*}, \boldsymbol{\Sigma}_s^{k*})$$

Variational inference(VI):

- minimize the **KL divergence** of exact posterior and approx. posterior

$$\mathcal{L} = \int q^*(\mathcal{U}, \tau) \log \frac{p(\{y_i\}_{i \in S_t} | \mathcal{U}, \tau) q(\mathcal{U}, \tau)}{q^*(\mathcal{U}, \tau)} d\mathcal{U} d\tau$$

Expectation propagation(EP):

- **moment match** when L(reverse order) has **closed form solution**

Assumed density filtering(ADF):

- **moment match** when having **tractable normalization term**

Reparameterization trick (SVI):

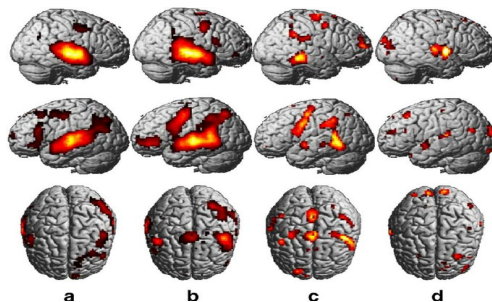
- **Probabilistic version SGD** when L is **totally intractable**

...(sampling based methods)

- Tensor-valued time series
- Represent time-varying and high-order interactions



(region, site, weather)



(subject, voxel, electrode)



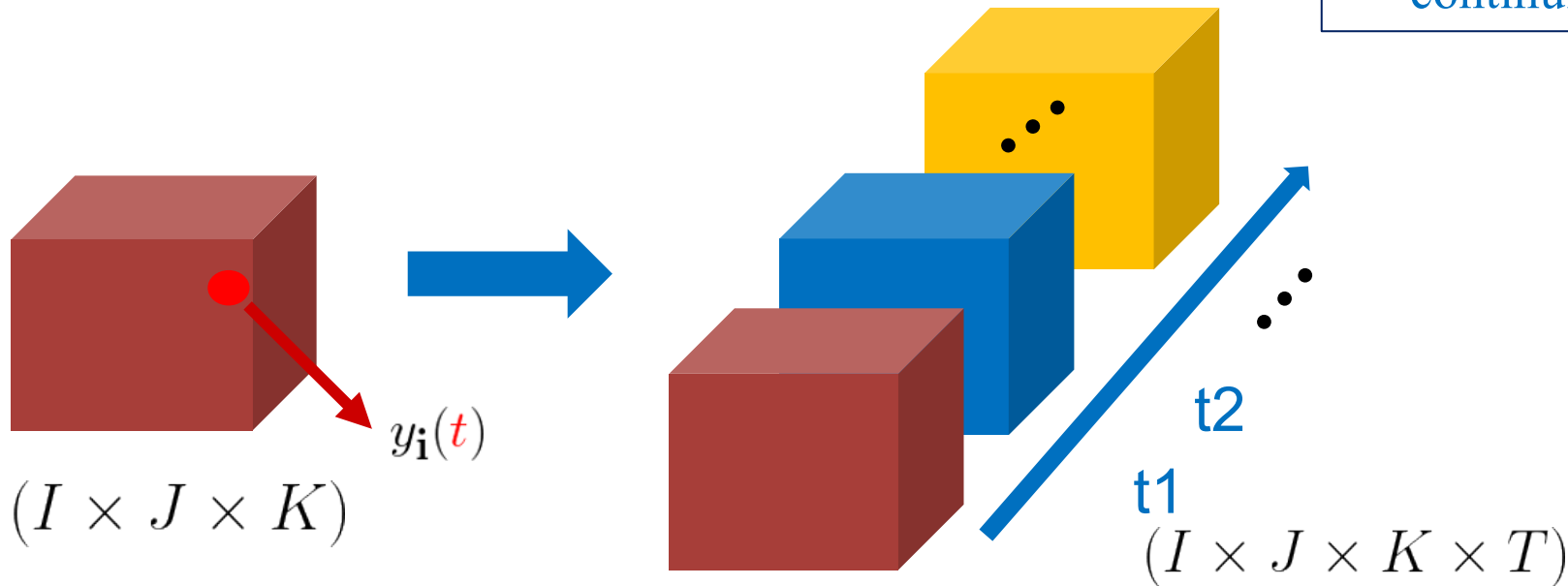
(user, user, location, message)

Tensor structure are evolving through time!

Temporal Tensor: trival solutions

- Drop timestamps
- Augment tensor with **discrete time mode**

- Too Sparse!
- Loss temporal continuity



Most existing work[1][2][3]:

Evolving weights + **Static Factors**

$$y_{\mathbf{i}}(t) \approx \mathbf{w}(t)^\top \left(\mathbf{u}_{i_1}^1 \circ \dots \circ \mathbf{u}_{i_K}^K \right)$$

- **Over-Simplistic!**
- **Evolving factors dominate in many cases**

[1] Zhang, Yanqing, et al. "Dynamic tensor recommender systems." *The Journal of Machine Learning Research* 22.1 (2021): 3032-3066.

[2] Fang, Shikai, et al. "Bayesian Continuous-Time Tucker Decomposition." *International Conference on Machine Learning*. PMLR, 2022.

[3] Li, Shibo, et al.. "Decomposing Temporal High-Order Interactions via Latent ODEs." *International Conference on Machine Learning*. PMLR, 2022.

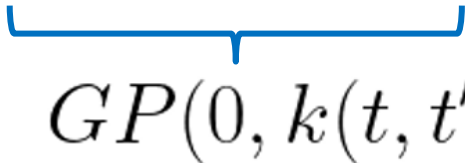
SFTL(ours): Learning functional factor trajectories!

CP:
$$y_{\mathbf{i}}(t) \approx \mathbf{w}^{\top} \left(\mathbf{u}_{i_1}^1(t) \circ \dots \circ \mathbf{u}_{i_K}^K(t) \right)^{\top}$$

Tucker:
$$y_{\mathbf{i}}(t) \approx \text{vec}(\mathcal{W})^{\top} \left(\mathbf{u}_{i_1}^1(t) \otimes \dots \otimes \mathbf{u}_{i_K}^K(t) \right)$$

Gaussian Process(GP): Bayesian Functional Prior

$$y_{\mathbf{i}}(t) \approx \mathbf{w}^{\top} \left(\mathbf{u}_{i_1}^1(t) \circ \dots \circ \mathbf{u}_{i_K}^K(t) \right)^{\top}$$



 $GP(0, k(t, t'))$

- Non-parametric: strong and flexible
- **Non-scalable : $O(N^3)$ time cost**

Linear-Cost GP with Chain-Structure

Temporal GPs

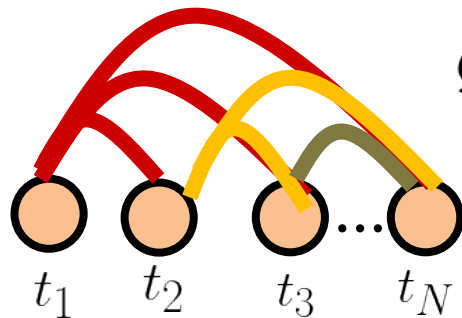
Temporal States:

LTI-SDE

$$\frac{d\gamma_{\mathbf{r}}(t)}{dt} = \mathbf{F}\gamma_{\mathbf{r}} + \mathbf{L}\xi(t)$$

discrete form

State Space Model

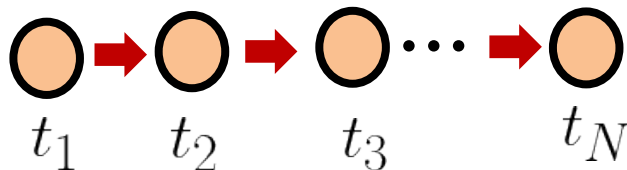


$$\mathcal{GP}(\mathbf{w}_{\mathbf{r}} \mid \mathbf{0}, \mathbf{K}_{\mathbf{r}}(t, t'))$$

Space: $\mathcal{O}(N^2)$

Time: $\mathcal{O}(N^3)$

$$p(\gamma_{\mathbf{r}}(t_{n+1}) \mid \gamma_{\mathbf{r}}(t_n)) = \mathcal{N}(\gamma_{\mathbf{r}}(t_{n+1}) \mid \mathbf{A}_n \gamma_{\mathbf{r}}(t_n), \mathbf{Q}_n)$$



Space: $\mathcal{O}(N)$

Time: $\mathcal{O}(N)$

Joint Prob of SFTL (CP)

noise

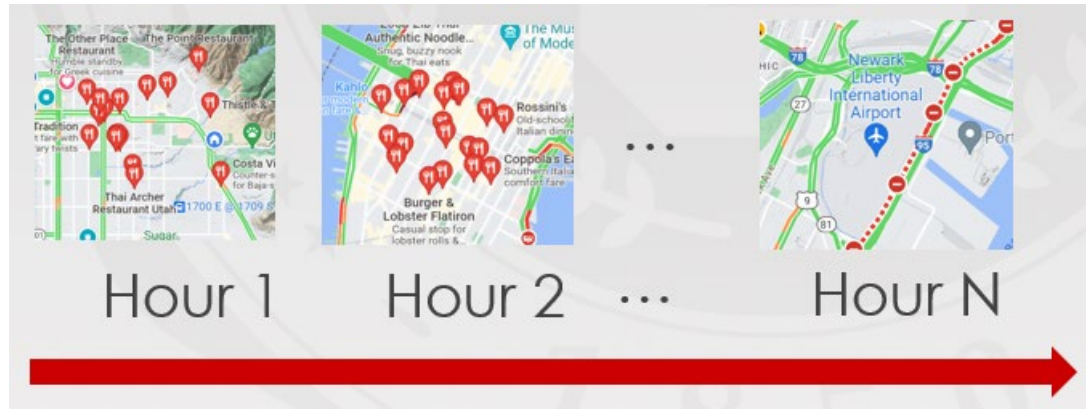
states of factor trajectory

$$\begin{aligned}
 p(\{\mathbf{h}_{j,k}^m\}, \tau, \mathbf{y}) &= \boxed{p(\tau)} \prod_{m=1}^M \prod_{j=1}^{d_m} \boxed{p(\mathbf{h}_{j,1}^m) \prod_{k=1}^{c_j^m-1} p(\mathbf{h}_{j,k+1}^m | \mathbf{h}_{j,k}^m)} \\
 &\cdot \prod_{n=1}^N \boxed{\mathcal{N}(y_n | \mathbf{1}^\top (\mathbf{u}_{\ell_{n1}}^1(t_n) \circ \dots \circ \mathbf{u}_{\ell_{nM}}^M(t_n)), \tau^{-1})}.
 \end{aligned}$$

likelihood

Inference Challenges

- One observation \longleftrightarrow interweaving of multi-trajectories
- Datasets are collected in **fast-generating & streaming** manner



Security / Privacy / Safety-constrained scenarios

- Data are **NOT** allowed to be stored/revisited



- **Offline learning** of classical tensor models:

$$\min_{U \in \mathbb{R}^{m \times d}, V \in \mathbb{R}^{n \times d}} \sum_{(i,j) \in \text{obs}} (A_{ij} - \langle U_i, V_j \rangle)^2 + w_0 \sum_{(i,j) \notin \text{obs}} (\langle U_i, V_j \rangle)^2$$

- Need to collect **all observations** at first
- Then go through data for **multi-epochs** to update

Can not fit **streaming** tensor data!

- **Online learning/ Streaming Inference:**
data come, model update, data drops
- Principle: Incremental version of Bayes' rule:

Posterior on old data

$$p(\boldsymbol{\theta} \mid \mathcal{D}_{\text{old}} \cup \mathcal{D}_{\text{new}}) \propto p(\boldsymbol{\theta} \mid \mathcal{D}_{\text{old}}) p(\mathcal{D}_{\text{new}} \mid \boldsymbol{\theta})$$

Posterior on all data

Likelihood on current model

Streaming Inference \Leftrightarrow Kalman Filter!

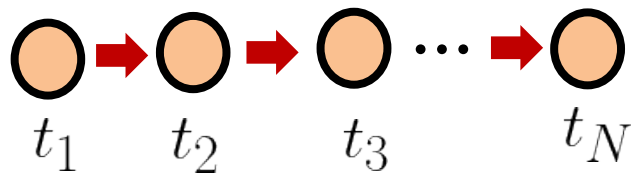
Incremental version of Bayes' rule:

$$p(\boldsymbol{\theta} \mid \mathcal{D}_{\text{old}} \cup \mathcal{D}_{\text{new}}) \propto p(\boldsymbol{\theta} \mid \mathcal{D}_{\text{old}}) p(\mathcal{D}_{\text{new}} \mid \boldsymbol{\theta})$$

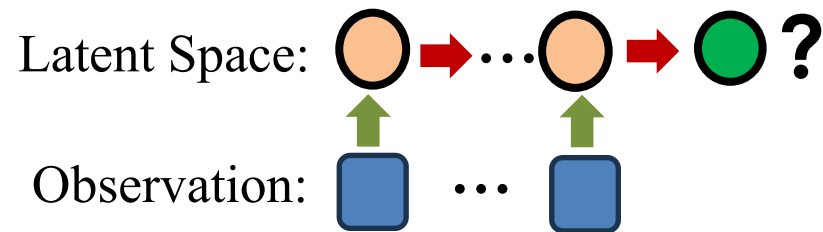
+



Chain-Structure of factor trajectory



Kalman Filter!

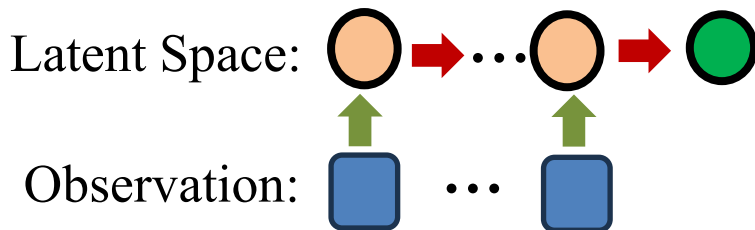


Last Challenge:

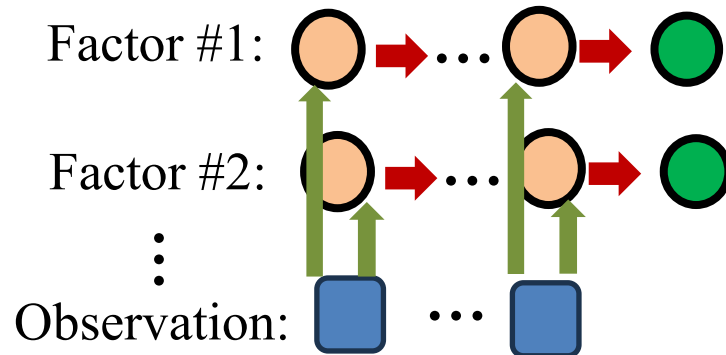
?

- **One** observation \longleftrightarrow States of **multi**-trajectories

- Classical Kalman Filter



- SFTL



- Conditional Moment-Matching

$$\mathcal{N}(y | \mathbf{1}^\top (\mathbf{u}_{\ell_1}^1(t_{n+1}) \circ \dots \circ \mathbf{u}_{\ell_M}^M(t_{n+1}))) \propto \prod_{m=1}^M \mathcal{N}(\mathbf{u}_{\ell_m}^m(t_{n+1}) | \boldsymbol{\gamma}_{\ell_m}^m, \boldsymbol{\Sigma}_{\ell_m}^m) \text{Gam}(\tau | \alpha_\ell, \omega_\ell)$$

Observation

Msgs to **multi**-trajectories

- Message Passing (**Expectation Propagation**)

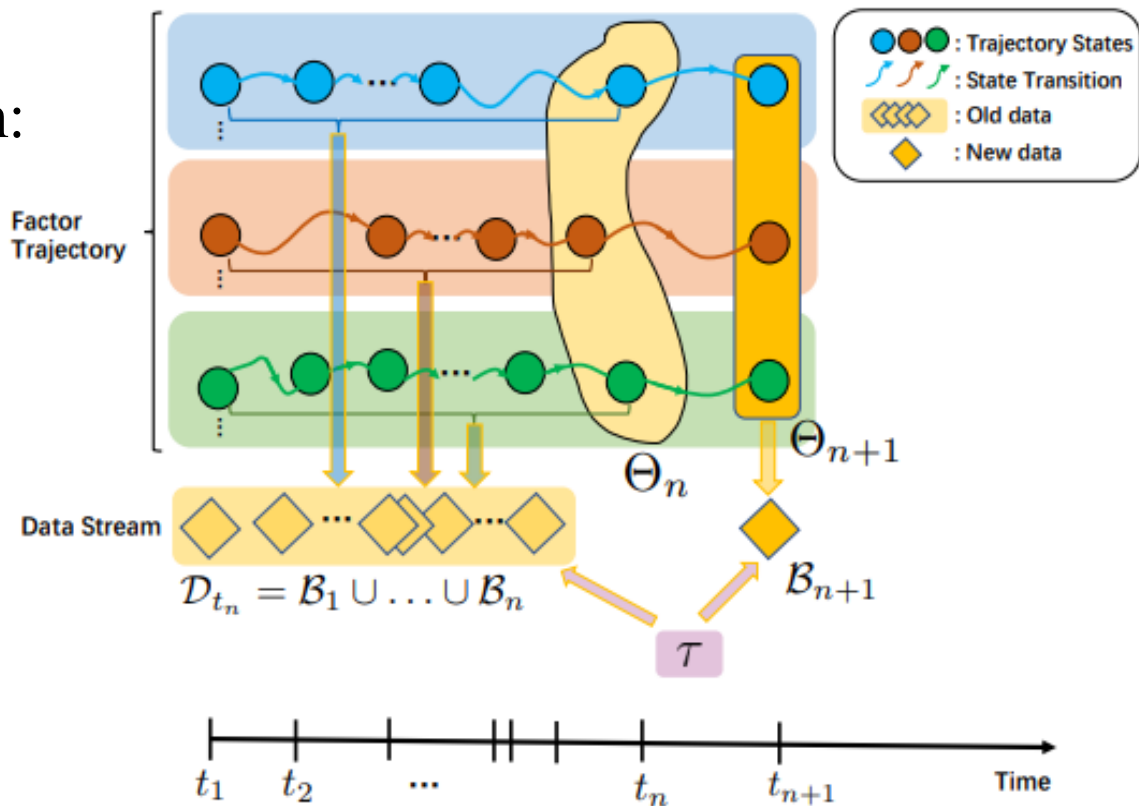
$$\mathbb{E}_q[\mathbf{v}_\ell^m] = \mathbb{E}_q[\mathbf{v}_{\ell_1}^1] \circ \dots \circ \mathbb{E}_q[\mathbf{v}_{\ell_{m-1}}^{m-1}] \circ \mathbb{E}_q[\mathbf{v}_{\ell_{m+1}}^{m+1}] \circ \dots \circ \mathbb{E}_q[\mathbf{v}_{\ell_M}^M],$$

$$\mathbb{E}_q[\mathbf{v}_\ell^m (\mathbf{v}_\ell^m)^\top] = \mathbb{E}_q[\mathbf{v}_{\ell_1}^1 (\mathbf{v}_{\ell_1}^1)^\top] \circ \dots \circ \mathbb{E}_q[\mathbf{v}_{\ell_{m-1}}^{m-1} (\mathbf{v}_{\ell_{m-1}}^{m-1})^\top]$$

$$\circ \mathbb{E}_q[\mathbf{v}_{\ell_{m+1}}^{m+1} (\mathbf{v}_{\ell_{m+1}}^{m+1})^\top] \circ \dots \circ \mathbb{E}_q[\mathbf{v}_{\ell_M}^M (\mathbf{v}_{\ell_M}^M)^\top].$$

Efficient **Online** algorithm:

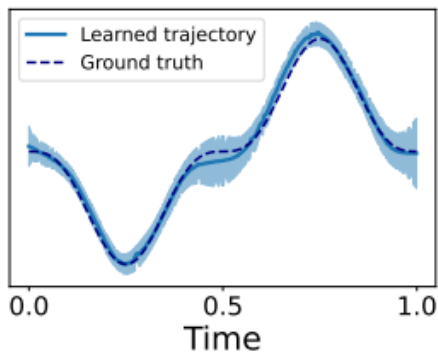
- **Closed-form Updates**
- **Conditional Moment Match**
- **Kalman Filter**
+ (RTS Smoother)



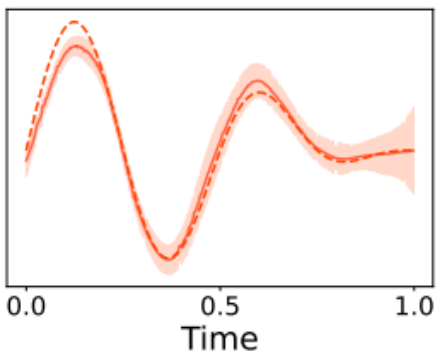
Algorithm 1 Streaming Factor Trajectory Learning (SFTL)

- 1: **Input:** kernel hyper-parameters $a, \rho, \nu = p + \frac{1}{2}$ ($p \in \{0, 1, 2, \dots\}$)
 - 2: $n \leftarrow 0$
 - 3: **while** Receiving a new batch of entreis \mathcal{B}_{n+1} **do**
 - 4: **if** $n = 0$ **then**
 - 5: Set $a_n = a_0, b_n = b_0, \hat{\boldsymbol{\mu}}_{j,c_{j,n+1}}^m = \mathbf{0}$, and $\hat{\mathbf{V}}_{j,c_{j,n+1}}^m = \bar{\mathbf{P}}_\infty$ in (8).
 - 6: Goto 9.
 - 7: **end if**
 - 8: Retrieve the involved preceding factor states $\Theta_n = \{\mathbf{h}_{j,c_{j,n}}^m | j \in \mathcal{I}_{n+1}^m\}_m$ and their running posterior, $p(\Theta_n, \tau | \mathcal{D}_{t_n}) \approx \text{Gam}(\tau | a_n, b_n) \prod_{m=1}^M \prod_{j \in \mathcal{I}_{n+1}^m} \mathcal{N}(\mathbf{h}_{j,c_{j,n}}^m | \hat{\boldsymbol{\mu}}_{j,c_{j,n}}^m, \hat{\mathbf{V}}_{j,c_{j,n}}^m)$.
 - 9: According to (8) and (9), use conditional Expectation Propagation to calculate the running posterior of the current factor states, $p(\Theta_{n+1}, \tau | \mathcal{D}_{t_{n+1}}) \approx \text{Gam}(\tau | a_{n+1}, b_{n+1}) \prod_{m=1}^M \prod_{j \in \mathcal{I}_{n+1}^m} \mathcal{N}(\mathbf{h}_{j,c_{j,n+1}}^m | \hat{\boldsymbol{\mu}}_{j,c_{j,n+1}}^m, \hat{\mathbf{V}}_{j,c_{j,n+1}}^m)$.
 - 10: **if** Needed **then**
 - 11: Run RTS smoothing on any factor state chain $\{\mathbf{h}_{j,k}^m | k = 1, 2, \dots\}$ of interest.
 - 12: **end if**
 - 13: $n \leftarrow n + 1$
 - 14: **end while**
 - 15: Run RTS smoothing for every factor state chain $\{\mathbf{h}_{j,k}^m | k = 1, 2, \dots\}$.
 - 16: **Return:** $\{q(\mathbf{h}_{j,k}^m | \mathcal{D}) | k = 1, 2, \dots\}_{1 \leq m \leq M, 1 \leq j \leq d_m}, q(\tau | \mathcal{D})$, where \mathcal{D} is all the data received.
-

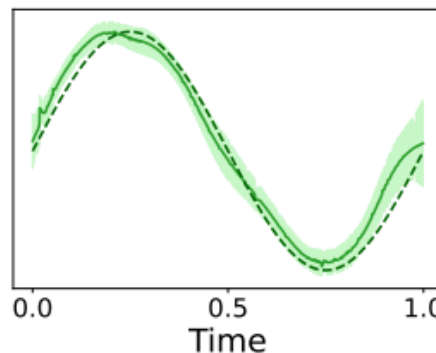
$$y_{(i,j)}(t) \sim \mathcal{N}(u_i^1(t)u_j^2(t), 0.05)$$



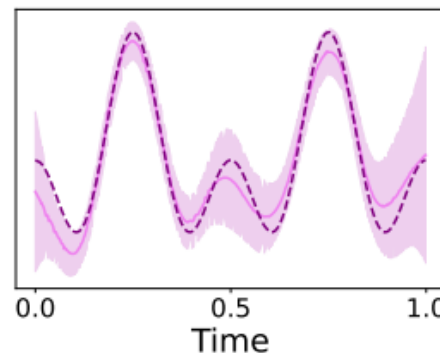
(a) $u_1^1(t)$



(b) $u_2^1(t)$

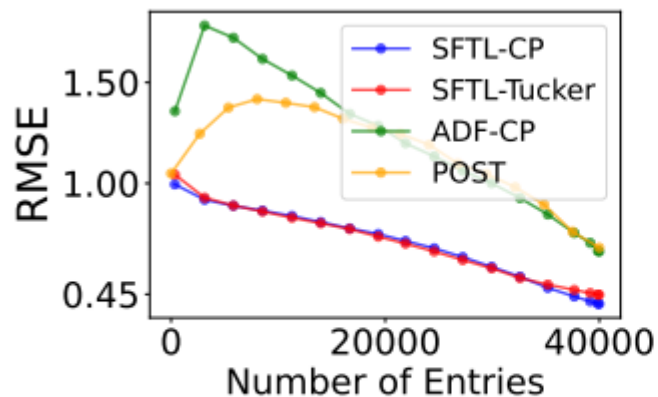


(c) $u_1^2(t)$

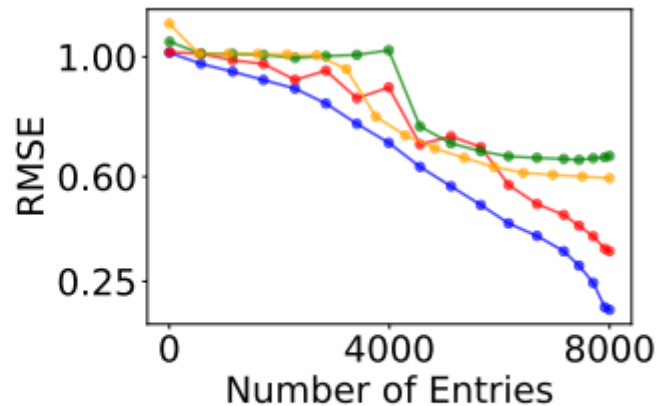


(d) $u_2^2(t)$

Experiments: Online Prediction



(a) *FitRecord*



(b) *ServerRoom*

Experiments: Prediction Accuracy

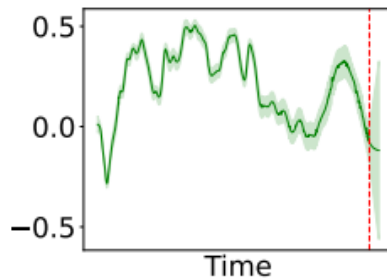
- **Online beat Offline!**

	RMSE	<i>FitRecord</i>
Static	PTucker	0.656 ± 0.147
	Tucker-ALS	0.846 ± 0.005
	CP-ALS	0.882 ± 0.017
	CT-CP	0.664 ± 0.007
	CT-GP	0.604 ± 0.004
	BCTT	0.518 ± 0.007
	NONFAT	0.503 ± 0.002
	THIS-ODE	0.526 ± 0.004
Stream	POST	0.696 ± 0.019
	ADF-CP	0.648 ± 0.008
	BASS-Tucker	0.976 ± 0.024
	SFTL-CP	0.424 ± 0.014
	SFTL-Tucker	0.430 ± 0.010

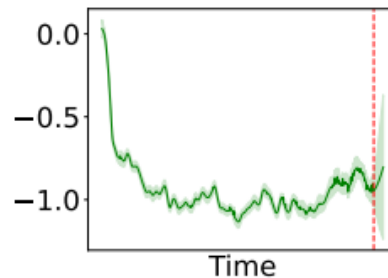
	MAE	
Static	PTucker	0.369 ± 0.009
	Tucker-ALS	0.615 ± 0.006
	CP-ALS	0.642 ± 0.012
	CT-CP	0.46 ± 0.004
	CT-GP	0.414 ± 0.001
	BCTT	0.355 ± 0.005
	NONFAT	0.341 ± 0.001
	THIS-ODE	0.363 ± 0.004
Stream	POST	0.478 ± 0.014
	ADF-CP	0.449 ± 0.006
	BASS	0.772 ± 0.031
	SFTL-CP	0.242 ± 0.006
	SFTL-Tucker	0.246 ± 0.001

Experiments: Learned Trajectories

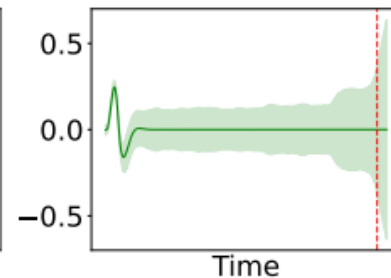
- Reasonable uncertainty
- Indicate tensor rank?



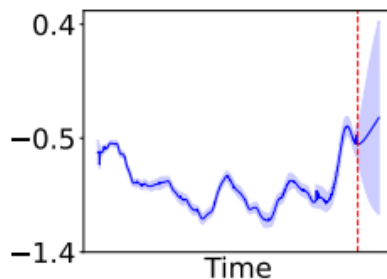
(a) $u_{1,1}^1(t)$



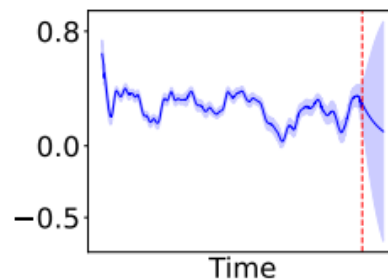
(b) $u_{1,2}^1(t)$



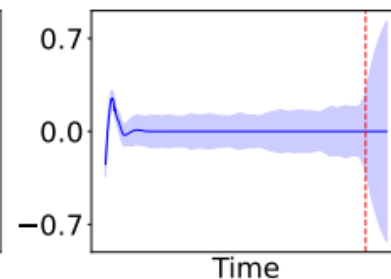
(c) $u_{1,3}^1(t)$



(d) $u_{2,1}^2(t)$



(e) $u_{2,2}^2(t)$



(f) $u_{2,3}^2(t)$

Thanks.

Github Repo

