



THE UNIVERSITY OF UTAH

Bayesian Continuous-Time Tucker Decomposition

Shikai Fang, Akil Narayan, Robert M. Kirby,
Shandian Zhe

Presenter: Shikai Fang

School of computing, The University of Utah

For ICML 2022



Outline

1. Background
2. Motivation
3. Dynamic Tucker-Core via SDE
4. Message-Passing Inference: SDE
Discretization+ Conditional Moment Matching
5. Experiments on Real-world Data



THE UNIVERSITY OF UTAH

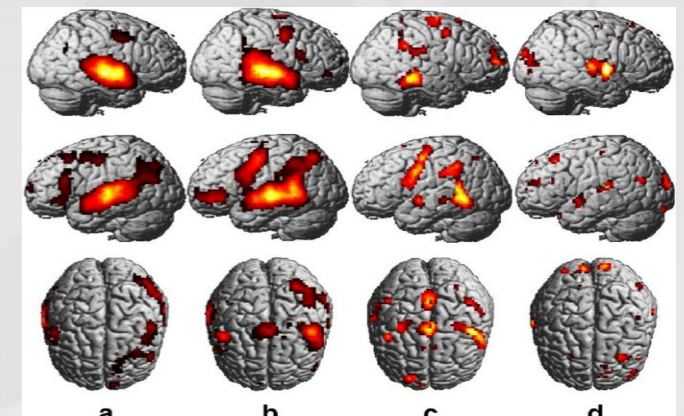
Tensor Data: Widely Used High-Order Data Structures to Represent Interactions of Multiple Objects/Entities



(user, movie, episode)



(user, advertisement, page-section)



(subject, voxel, electrode)



(user, item, online-store)



(user, user, location, message-type)

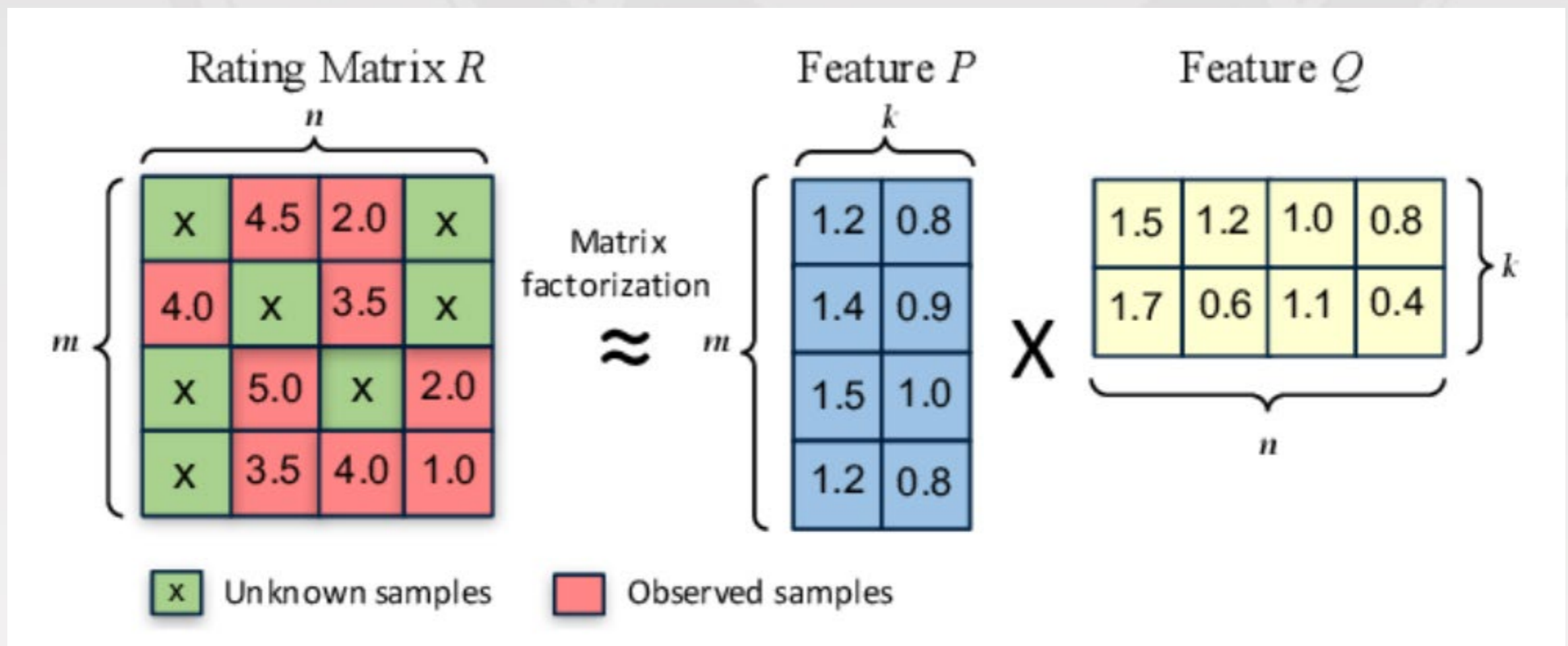


(patient, gene, condition)



Tensor Decomposition: estimate latent factors to reconstruct tensor with observed entries

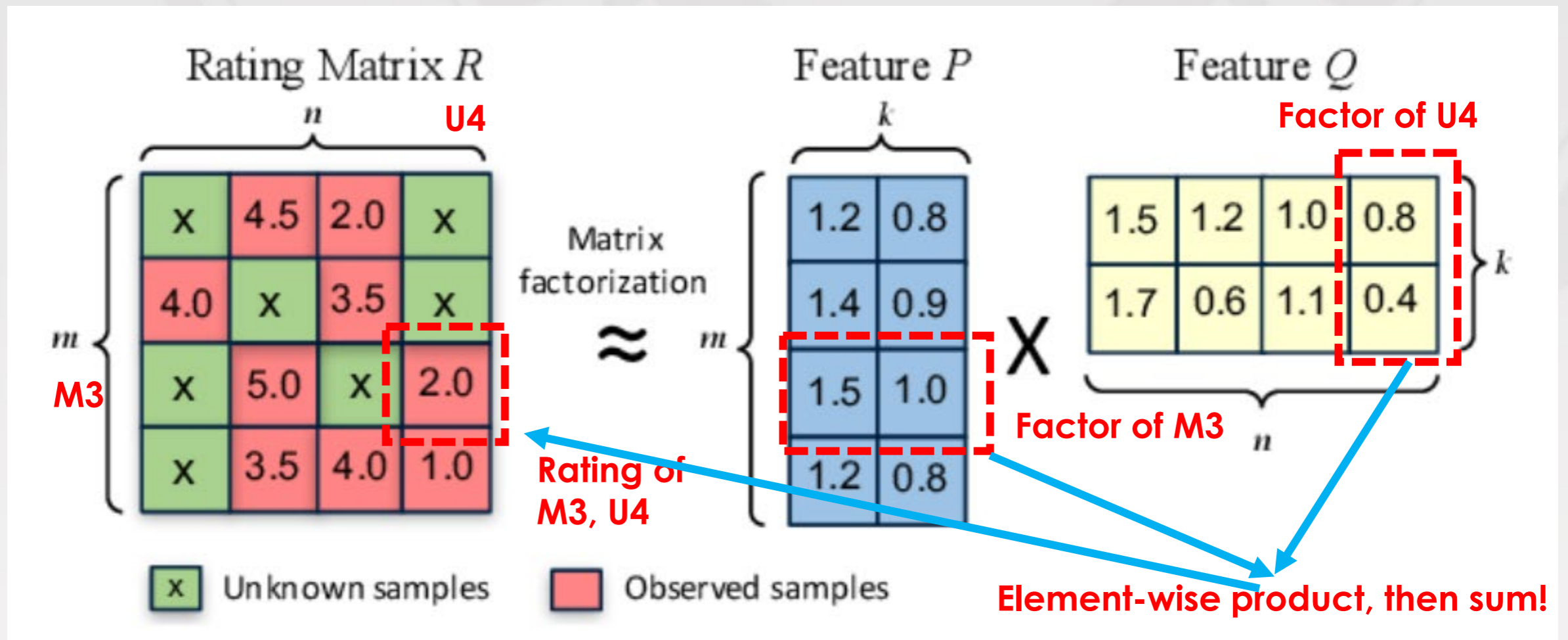
- Simple case:
Collaborative Filtering (Matrix Factorization)





Tensor Decomposition: estimate latent factors to reconstruct tensor with observed entries

- Simple case:
Collaborative Filtering (Matrix Factorization)



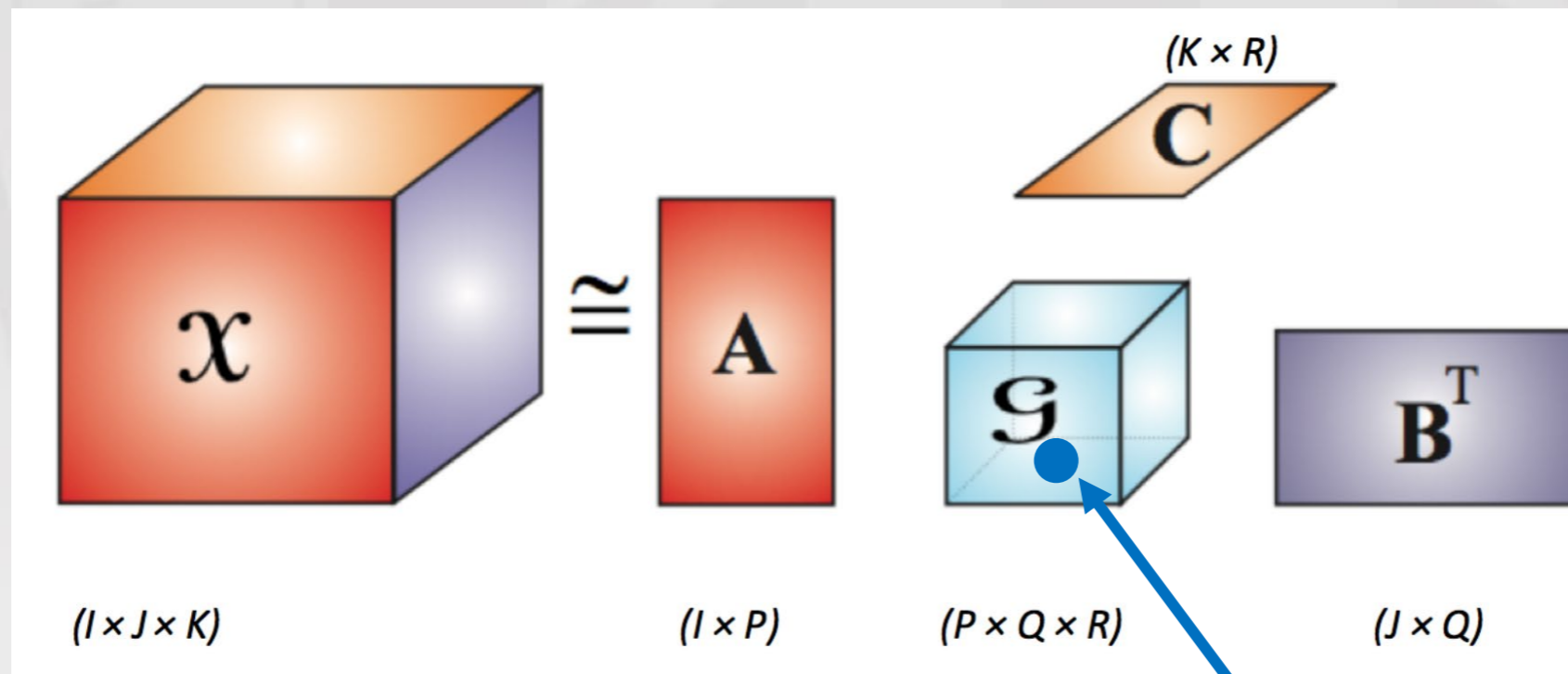


Tucker Decomposition

- 2-D matrix \Rightarrow N-D tensor
- Element-wise interaction \Rightarrow all possible interactions



Tucker Decomposition



One interaction weight

One Interaction of latent factors

Element-wise form for a K-mode tensor \mathcal{Y} :

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

$$= \sum_{r_1=1}^{R_1} \dots \sum_{r_K=1}^{R_K} \left[w_{(r_1, \dots, r_K)} \cdot \prod_{k=1}^K u_{i_k, r_k}^k \right]$$

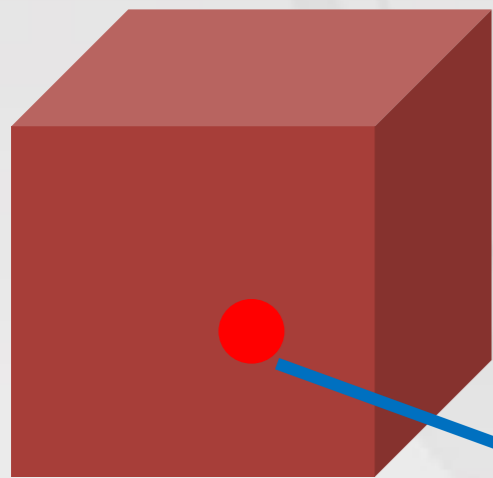


Challenge: Temporal info in Tensor

What about each entry is time-dependent?

Straightforward Solution:

- Drop time or
- Augment tensor with **time-step mode**

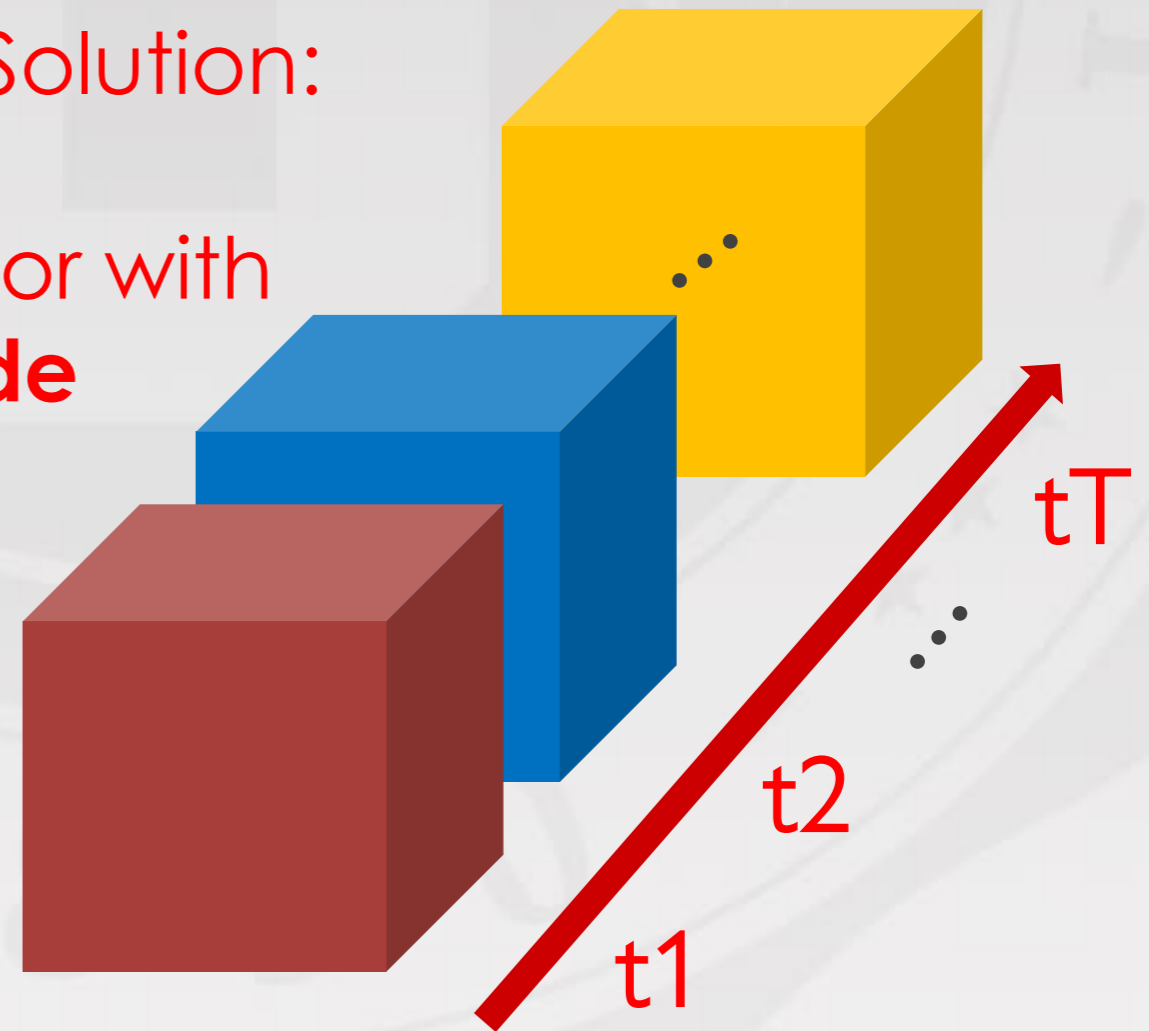


$$X_{ijk}(t)$$

$$(I \times J \times K)$$

Problem:

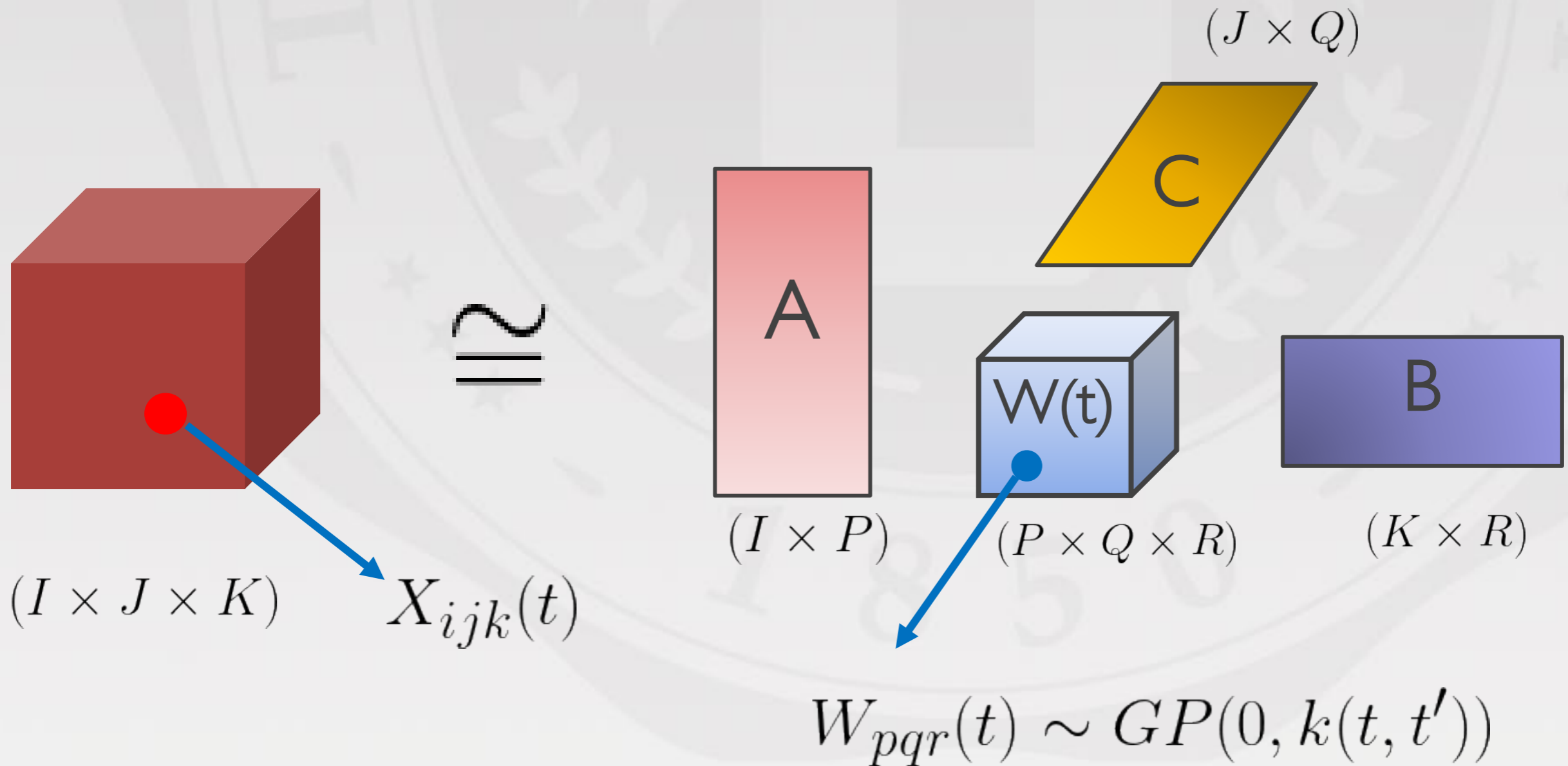
1. Too Sparse
2. Ignore the temporary continuity



$$(I \times J \times K \times T)$$



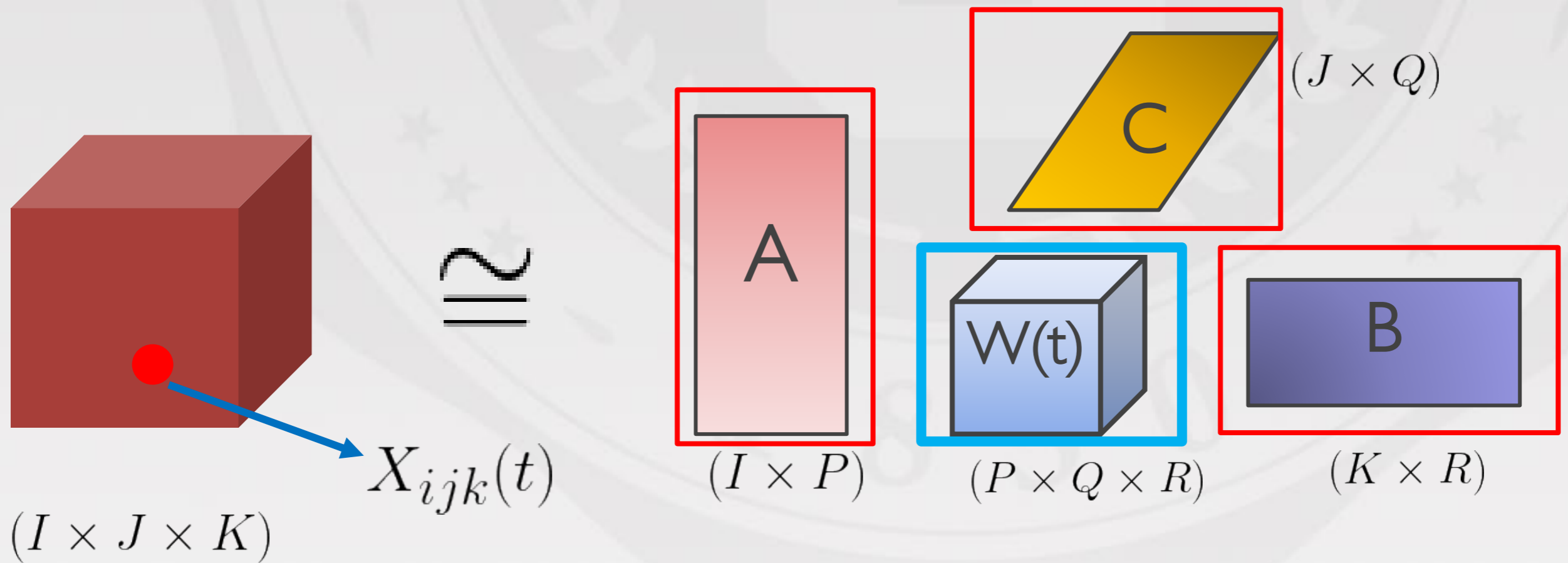
Our Solution: Modeling Dynamic Tucker Core by Temporal Gaussian Processes





High-level Motivation:

Decouple the **representation learning of factors**
and the **capture of dynamic pattern**





Joint Probability:

$$p(\mathcal{U}, \{\mathbf{w}_r\}_r, \tau, \mathbf{y}) =$$

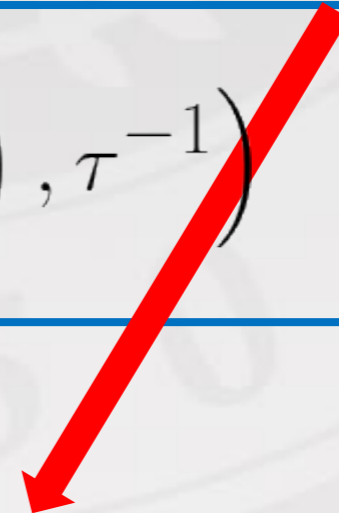
$$\text{Gam}(\tau \mid b_0, c_0) \prod_{k=1}^K \prod_{j=1}^{d_k} \mathcal{N}(\mathbf{u}_j^k \mid \mathbf{0}, \mathbf{I}) \times \prod_{\mathbf{r}=(1,\dots,1)}^{(R_1,\dots,R_K)} \mathcal{N}(\mathbf{w}_r \mid \mathbf{0}, \mathbf{K}_r) \times$$

Priors of factors and noise

Temporal GPs on Tucker Core

$$\prod_{n=1}^N \mathcal{N}(y_n \mid \text{vec}(\mathcal{W}(t_n))^\top (\mathbf{u}_{i_{n_1}}^1 \otimes \dots \otimes \mathbf{u}_{i_{n_K}}^K), \tau^{-1})$$

Gaussian Likelihood



Computational challenge: $\mathbf{O}(N^3)$ cost of full GPs



THE UNIVERSITY OF UTAH

To **avoid low-rank/sparse approx.** (low quality),
but enjoy **linear-cost inference of full GPs,**

We apply a crucial fact:

Temporal GPs



with stationary kernel

Linear Time-Invariant(LTI) SDE



discrete form on $(t_1, t_2 \dots t_N)$

State Space Model (Gauss Markov Chain)

Can be solved by
**Kalman filtering &
RTS Smoothing in $O(N)^*$**

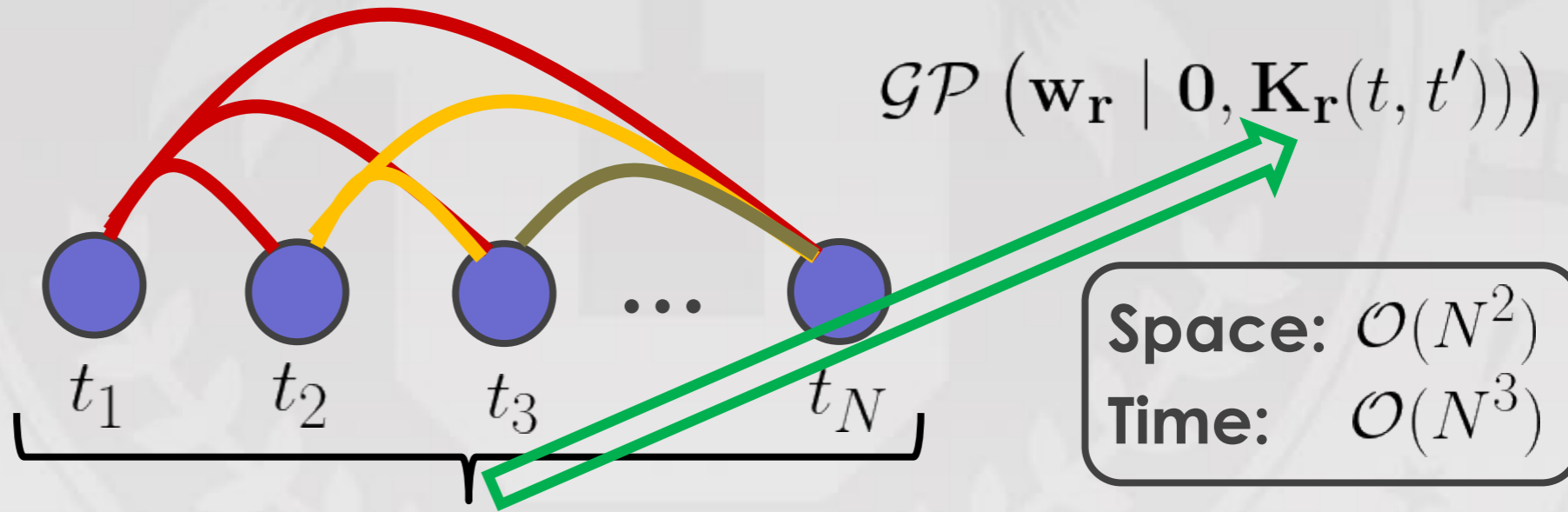


*: it holds with linear emission/observed likelihood, if with non-linear, we could apply non-linear filter and smoothing



Illustration of computation cost:

Temporal GPs



Temporal States:

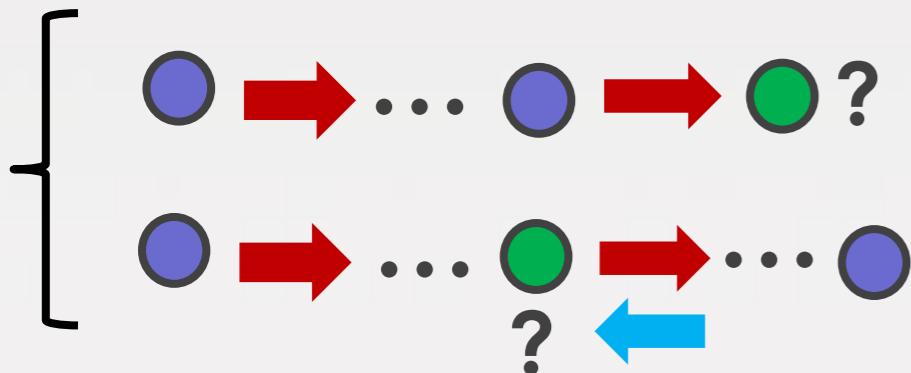
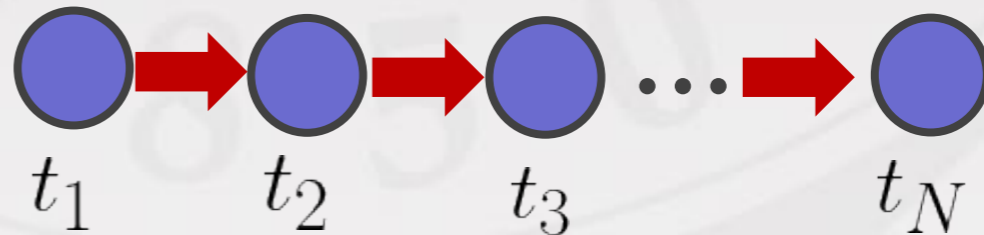
LTI-SDE

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

discrete form

State Space Model

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



Kalman Filter

RTS Smoothing

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



Specifically,

Temporal GPs

$$\prod_{\mathbf{r}=(1,\dots,1)}^{(R_1,\dots,R_K)} \mathcal{N}(\mathbf{w}_{\mathbf{r}} \mid \mathbf{0}, \mathbf{K}_{\mathbf{r}})$$



with stationary kernel
(e.g., Matern25)

$$\begin{cases} \gamma_{\mathbf{r}}(t) = \left(w_{\mathbf{r}}, \frac{dw_{\mathbf{r}}}{dt} \right)^{\top} \\ \frac{d\gamma_{\mathbf{r}}(t)}{dt} = \mathbf{F}\gamma_{\mathbf{r}} + \mathbf{L}\xi(t) \end{cases}$$

Linear Time-Invariant(LTI) SDE



discrete form on $(t_1, t_2 \dots t_N)$

$$\begin{cases} p(\gamma_{\mathbf{r}}(t_1)) = \mathcal{N}(\gamma_{\mathbf{r}}(t_1) \mid \mathbf{0}, \mathbf{P}_{\infty}) \\ p(\gamma_{\mathbf{r}}(t_{n+1}) \mid \gamma_{\mathbf{r}}(t_n)) = \mathcal{N}(\gamma_{\mathbf{r}}(t_{n+1}) \mid \mathbf{A}_n\gamma_{\mathbf{r}}(t_n), \mathbf{Q}_n) \end{cases}$$

State Space Model (Gauss Markov Chain)



Recall: linear-cost solver - KF, RTS



Reformulate **Tucker core** with **State Space Priors**

$$p(\bar{\gamma}_1) \prod_{n=1}^{N-1} p(\bar{\gamma}_{n+1} \mid \bar{\gamma}_n)$$

We post **Gaussian-Gamma Approx.** to fit each data-llk

$$\mathcal{N}\left(y_n \mid (\mathbf{H}\bar{\gamma}_n)^\top \left(\mathbf{u}_{i_{n1}}^1 \otimes \dots \otimes \mathbf{u}_{i_{nK}}^K\right), \tau^{-1}\right) \approx$$

$$Z_n \prod_{k=1}^K \mathcal{N}\left(\mathbf{u}_{i_{nk}}^k \mid \mathbf{m}_{i_{nk}}^{k,n}, \mathbf{V}_{i_{nk}}^{k,n}\right) \cdot \text{Gam}(\tau \mid b_n, c_n)$$

Approx. Msg of Factors & noise

$$\times \mathcal{N}(\mathbf{H}\bar{\gamma}_n \mid \boldsymbol{\beta}_n, \mathbf{S}_n)$$

Approx. Msg of SDE states /Tucker core

Substitute these into joint prob.



The proposed approx. posterior is:

$$q(\mathcal{U}, \{\bar{\gamma}_n\}, \tau) \propto \prod_{k=1}^K \prod_{j=1}^{d_k} \mathcal{N}(\mathbf{u}_j^k \mid \mathbf{0}, \mathbf{I}) \text{Gam}(\tau \mid b_0, c_0)$$

Standard moment match? Infeasible!

$$\prod_{n=1}^N \prod_{k=1}^K \mathcal{N}(\mathbf{u}_{i_{n_k}}^k \mid \mathbf{m}_{i_{n_k}}^{k,n}, \mathbf{V}_{i_{n_k}}^{k,n}) \text{Gam}(\tau \mid b_n, c_n)$$

$$p(\bar{\gamma}_1) \mathcal{N}(\mathbf{H}\bar{\gamma}_1 \mid \boldsymbol{\beta}_1, \mathbf{S}_1) \prod_{n=1}^{N-1} p(\bar{\gamma}_{n+1} \mid \bar{\gamma}_n) \mathcal{N}(\mathbf{H}\bar{\gamma}_n \mid \boldsymbol{\beta}_n, \mathbf{S}_n)$$

SDE states: Solve by KF and RTS

Apply conditional moment matching and delta method!



- **Conditional Moment Match**

$$\mathbb{E}_{\tilde{p}}[\phi(\boldsymbol{\eta}_n)] = \mathbb{E}_{\tilde{p}(\Theta \setminus \eta_n)} \left[\mathbb{E}_{\tilde{p}(\boldsymbol{\eta}_n | \Theta \setminus \eta_n)} [\phi(\boldsymbol{\eta}) | \Theta \setminus \eta_n] \right]$$

- **Delta method:**

$$\mathbb{E}_q(\Theta \setminus \eta_n) [\boldsymbol{\rho}_n] \approx \rho_n \left(\mathbb{E}_q [\Theta \setminus \eta_n] \right)$$

Enable **tractable moment matching**

to update **approx. probability terms** under
Expectation Propagation (EP) framework



Algorithm 1 BCTT

Input: $\mathcal{D} = \{(\mathbf{i}_1, t_1, y_1), \dots, (\mathbf{i}_N, t_N, y_N)\}$, kernel hyper-parameters l, σ^2

Initialize approximation terms in (10) for each likelihood.

repeat

Run KF and RTS smoothing to compute each $q(\bar{\gamma}_n)$

for $n = 1$ **to** N **in parallel do**

Simultaneously update $\mathcal{N}(\mathbf{H}\bar{\gamma}_n | \beta_n, \mathbf{S}_n)$,
 $\text{Gam}(\tau | b_n, c_n)$ and $\left\{ \mathcal{N} \left(\mathbf{u}_{i_{n_k}}^k | \mathbf{m}_{i_{n_k}}^{k,n}, \mathbf{V}_{i_{n_k}}^{k,n} \right) \right\}_k$
in (10) with conditional moment matching and multi-variate delta method.

end for

until Convergence

Return: $\{q(\mathcal{W}(t_n))\}_{n=1}^N, \{q(\mathbf{u}_j^k)\}_{1 \leq k \leq K, 1 \leq j \leq d_k}, q(\tau)$

Time cost: $\mathcal{O}(N\bar{R})$ **Space cost:** $\mathcal{O} \left(N \left(\bar{R}^2 + \sum_{k=1}^K R_k^2 \right) \right)$

➔ Kalman Filtering (forward)
 ➔ RTS Smoothing (backward)
 ➔ Conditional Moment Matching (parallel)

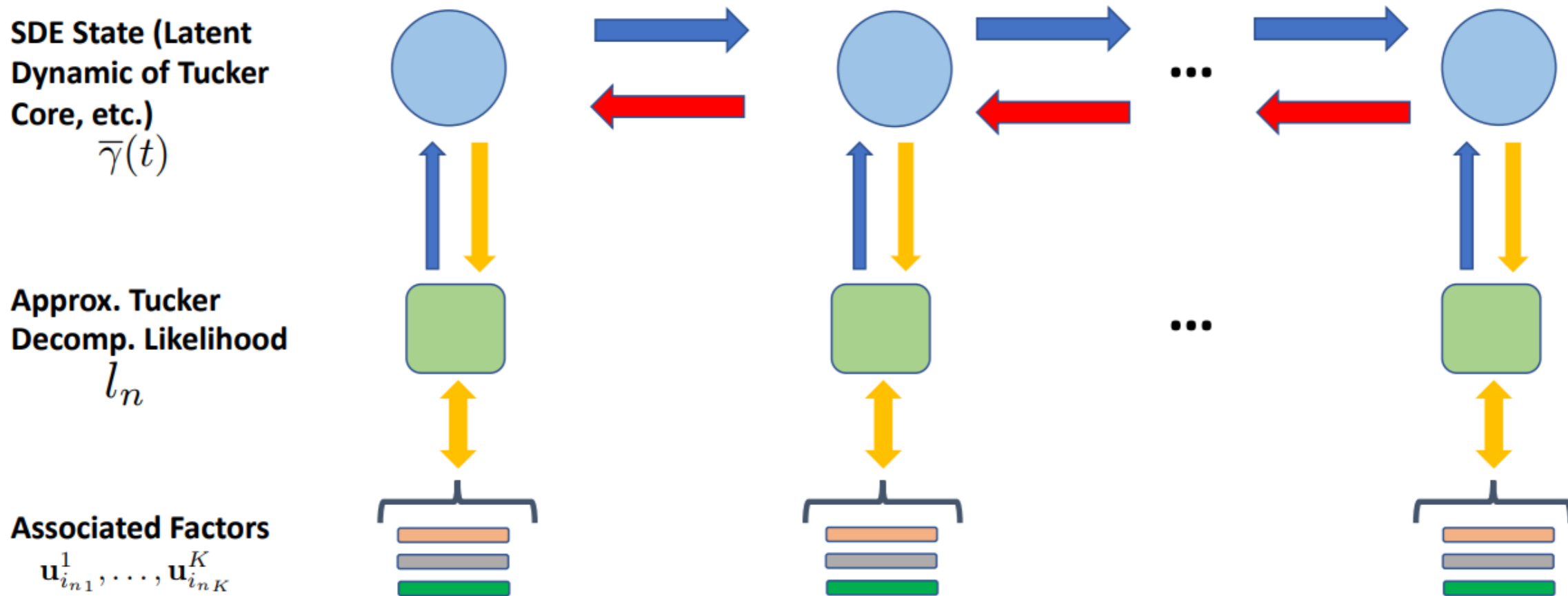
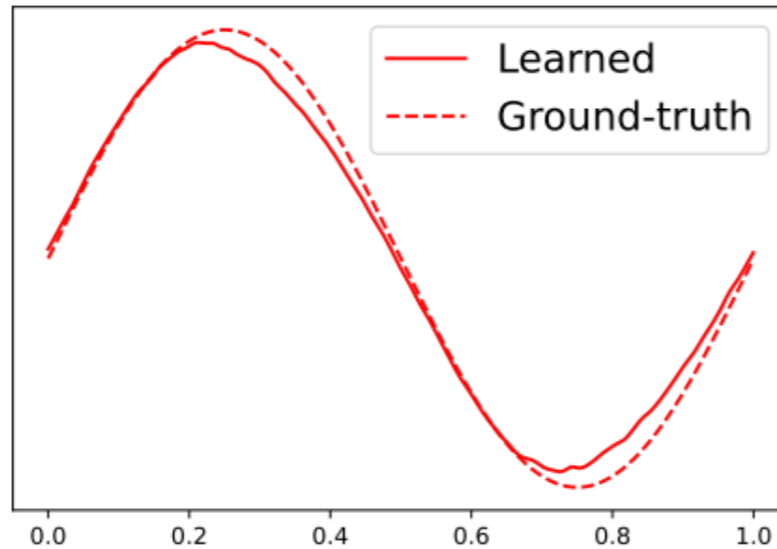


Figure 1. Graphical illustration of the message-passing inference algorithm.

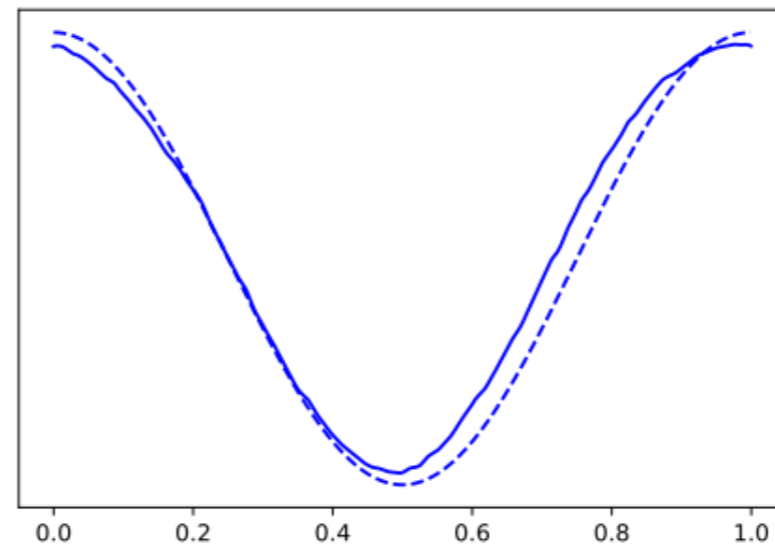


Can BCTT capture the temporal patterns in tensor?

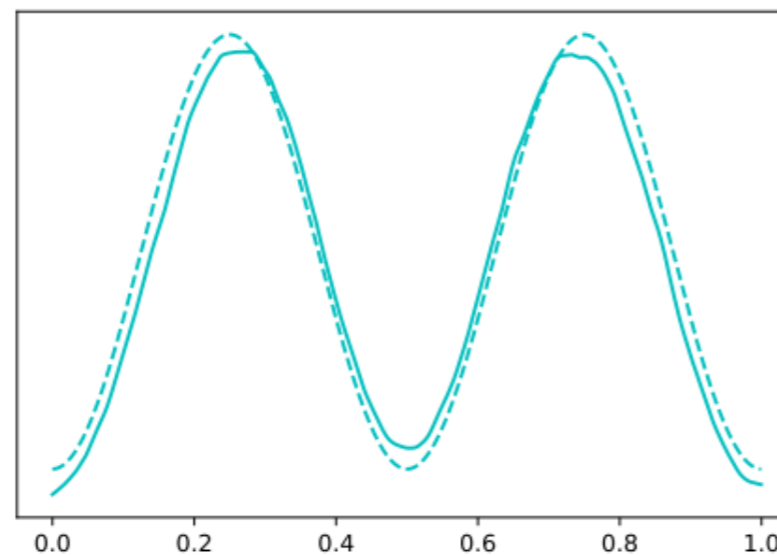
- Exp on **simulation data**
- Plot the dynamics of Tucker core



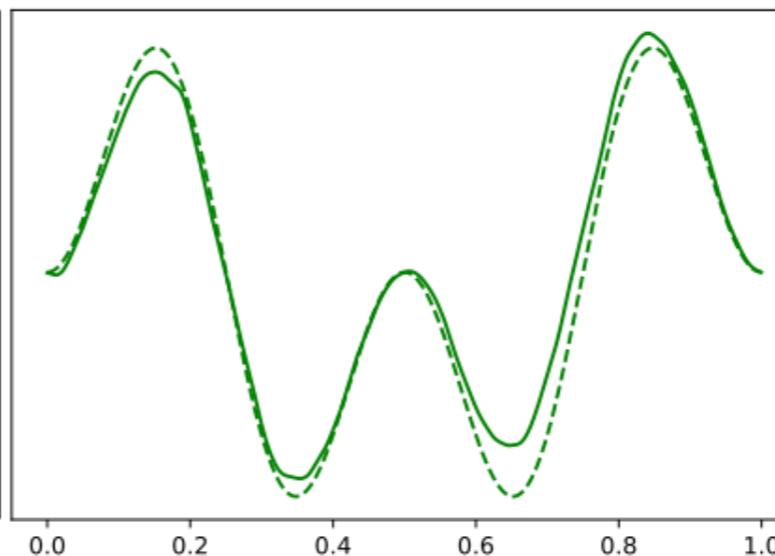
(a) $w_{(1,1)}(t)$



(b) $w_{(1,2)}(t)$



(c) $w_{(2,1)}(t)$



(d) $w_{(2,2)}(t)$



Can BCTT capture the temporal patterns in tensor?

- Exp on **real-world data (DBLP dataset)**
- Scatter low-rank structures of Tucker core

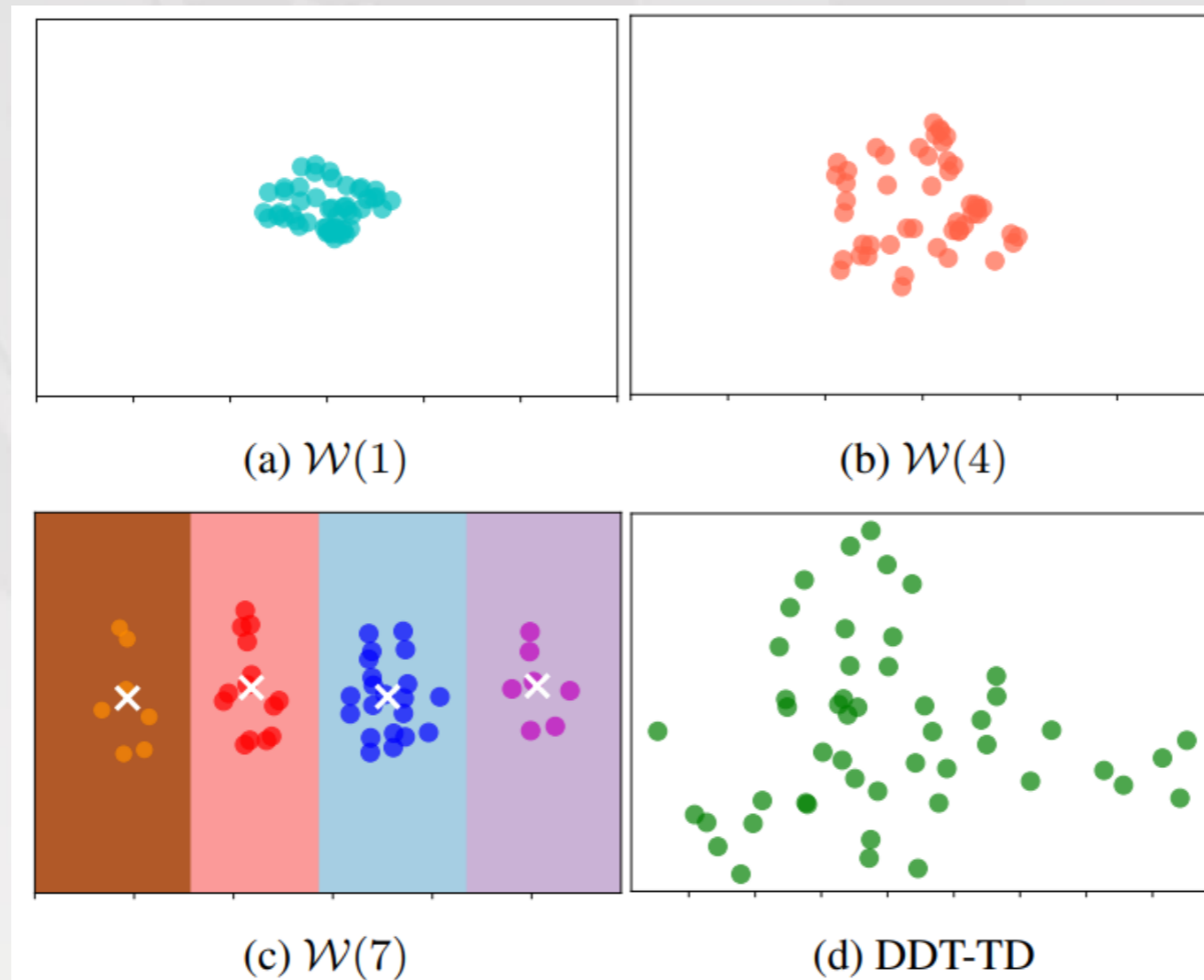


Figure 4. The structures of learned tensor-core at different time points by BCTT (a-c) and the static tensor-score learned by dynamic discrete-time Tucker decomposition (DDT-TD).



Prediction with BCTT

- Prediction performance of BCTT on 3 real-world data

| RMSE | <i>MovieLens</i> | <i>AdsClicks</i> | <i>DBLP</i> |
|-------------|----------------------|----------------------|----------------------|
| CT-CP | 1.113 ± 0.004 | 1.337 ± 0.013 | 0.240 ± 0.007 |
| CT-GP | 0.949 ± 0.008 | 1.422 ± 0.008 | 0.227 ± 0.009 |
| DT-GP | 0.963 ± 0.008 | 1.436 ± 0.015 | 0.227 ± 0.007 |
| DDT-GP | 0.957 ± 0.008 | 1.437 ± 0.010 | 0.225 ± 0.006 |
| DDT-CP | 1.022 ± 0.003 | 1.420 ± 0.020 | 0.245 ± 0.004 |
| DDT-TD | 1.059 ± 0.006 | 1.401 ± 0.022 | 0.232 ± 0.09 |
| BCTT | 0.922 ± 0.002 | 1.322 ± 0.012 | 0.214 ± 0.009 |

MAE

| | | | |
|-------------|----------------------|----------------------|----------------------|
| CT-CP | 0.788 ± 0.004 | 0.787 ± 0.006 | 0.105 ± 0.001 |
| CT-GP | 0.714 ± 0.004 | 0.891 ± 0.011 | 0.092 ± 0.004 |
| DT-GP | 0.722 ± 0.008 | 0.893 ± 0.008 | 0.084 ± 0.003 |
| DDT-GP | 0.720 ± 0.003 | 0.894 ± 0.009 | 0.083 ± 0.001 |
| DDT-CP | 0.755 ± 0.002 | 0.901 ± 0.011 | 0.114 ± 0.002 |
| DDT-TD | 0.742 ± 0.006 | 0.866 ± 0.012 | 0.101 ± 0.001 |
| BCTT | 0.698 ± 0.002 | 0.777 ± 0.016 | 0.084 ± 0.001 |

(a) $R = 3$

| RMSE | <i>MovieLens</i> | <i>AdsClicks</i> | <i>DBLP</i> |
|-------------|----------------------|----------------------|----------------------|
| CT-CP | 1.165 ± 0.008 | 1.324 ± 0.013 | 0.263 ± 0.006 |
| CT-GP | 0.965 ± 0.019 | 1.410 ± 0.015 | 0.227 ± 0.007 |
| DT-GP | 0.949 ± 0.007 | 1.425 ± 0.015 | 0.225 ± 0.008 |
| DDT-GP | 0.948 ± 0.005 | 1.421 ± 0.012 | 0.220 ± 0.006 |
| DDT-CP | 1.141 ± 0.007 | 1.623 ± 0.013 | 0.282 ± 0.011 |
| DDT-TD | 0.944 ± 0.003 | 1.453 ± 0.035 | 0.312 ± 0.072 |
| BCTT | 0.895 ± 0.007 | 1.304 ± 0.018 | 0.202 ± 0.009 |

MAE

| | | | |
|-------------|----------------------|----------------------|----------------------|
| CT-CP | 0.835 ± 0.006 | 0.792 ± 0.007 | 0.128 ± 0.001 |
| CT-GP | 0.717 ± 0.012 | 0.883 ± 0.016 | 0.092 ± 0.002 |
| DT-GP | 0.714 ± 0.005 | 0.886 ± 0.012 | 0.084 ± 0.001 |
| DDT-GP | 0.707 ± 0.004 | 0.882 ± 0.015 | 0.082 ± 0.003 |
| DDT-CP | 0.843 ± 0.003 | 1.082 ± 0.013 | 0.141 ± 0.004 |
| DDT-TD | 0.712 ± 0.002 | 0.903 ± 0.024 | 0.221 ± 0.047 |
| BCTT | 0.679 ± 0.001 | 0.785 ± 0.010 | 0.080 ± 0.001 |

(b) $R = 7$



THE UNIVERSITY OF UTAH

Thanks for attention Q&A Time

Presenter' email: shikai.fang@utah.edu

Focus: Bayesian machine learning, tensor learning