



# Many Processors, Little Time: MCMC for Partitions via Optimal Transport Couplings Brian L. Trippe

Postdoctoral Research Fellow

Columbia University Department of Statistics





Tin D. Nguyen Tamara Broderick

Example: Cluster cells based on gene expression [1]

- How large are the clusters (cell types)?
- Which cells are of the same type?



Example: Cluster cells based on gene expression [1]

- How large are the clusters (cell types)?
- Which cells are of the same type?
   Intractability of Bayesian posterior → MCMC



Example: Cluster cells based on gene expression [1]

- How large are the clusters (cell types)?
- Which cells are of the same type? Intractability of Bayesian posterior  $\rightarrow$  MCMC





BISCUIT (F-score: 0.91)





Example: Cluster cells based on gene expression [1]

- How large are the clusters (cell types)?
- Which cells are of the same type? Intractability of Bayesian posterior  $\rightarrow$  MCMC



Challenge: MCMC takes a long time!







Example: Cluster cells based on gene expression [1]

- How large are the clusters (cell types)?
- Which cells are of the same type? Intractability of Bayesian posterior  $\rightarrow$  MCMC





BISCUIT (F-score: 0.91)





Challenge: MCMC takes a long time! Goal: Use computational parallelism to accelerate MCMC

Example: Cluster cells based on gene expression [1]

- How large are the clusters (cell types)?
- Which cells are of the same type? Intractability of Bayesian posterior  $\rightarrow$  MCMC





BISCUIT (F-score: 0.91)





Challenge: MCMC takes a long time! Goal: Use computational parallelism to accelerate MCMC This work: Optimal transport couplings make this possible

### Roadmap



#### • Parallelizing MCMC with Couplings:

- Background & Notation
- The Label Switching Problem
- We Frame Gibbs Sampling as Markov Chain on Partitions
- Our Optimal Transport Coupling
- Big-O Analysis Demonstrates Fast Computation
- Improved Estimation Error and Intervals with OTC over Naïve Parallelism in Practice

Set-up & challenge of burn-in bias:

• Want to compute  $H^* = \int h(X) p_X(X) dX$ 

Set-up & challenge of burn-in bias:

- Want to compute  $H^* = \int h(X)p_X(X)dX$
- Choose <u>Markov chain</u>:  $X_0, X_1, \dots \rightsquigarrow p_X$

- Set-up & challenge of burn-in bias:
- Want to compute  $H^* = \int h(X)p_X(X)dX$
- Choose <u>Markov chain</u>:  $X_0, X_1, \dots \rightsquigarrow p_X$
- Problem: for any  $T<\infty\,,\mathbb{E}[h(X_T)]\neq H^*$

- Set-up & challenge of burn-in bias:
- Want to compute  $H^* = \int h(X) p_X(X) dX$
- Choose <u>Markov chain</u>:  $X_0, X_1, \dots \rightsquigarrow p_X$
- Problem: for any  $T<\infty\,,\mathbb{E}[h(X_T)]\neq H^*$
- "Naïve parallelism" gives little help!

Set-up & challenge of burn-in bias:

- Want to compute  $H^* = \int h(X) p_X(X) dX$
- Choose <u>Markov chain</u>:  $X_0, X_1, \dots \rightsquigarrow p_X$
- Problem: for any  $T < \infty$ ,  $\mathbb{E}[h(X_T)] \neq H^*$
- "Naïve parallelism" gives little help!

Unbiased MCMC Set-Up (simplified)<sup>†</sup> Use "coupled chain" ( $Y_0, Y_1, ...$ ) where: 1.  $Y_t \sim X_t$ 2.  $Y_t \sim X_t$  for  $t > \tau$  ("monting time")

2.  $Y_{t-1} = X_t$  for  $t > \tau$  ("meeting time")

Set-up & challenge of burn-in bias:

- Want to compute  $H^* = \int h(X)p_X(X)dX$
- Choose <u>Markov chain</u>:  $X_0, X_1, \dots \rightsquigarrow p_X$
- Problem: for any  $T < \infty$  ,  $\mathbb{E}[h(X_T)] \neq H^*$
- "Naïve parallelism" gives little help!

Unbiased MCMC Set-Up (simplified)<sup>†</sup> Use "coupled chain" ( $Y_0, Y_1, ...$ ) where: 1.  $Y_t \sim X_t$ 2.  $Y_{t-1} = X_t$  for  $t > \tau$  ("meeting time") Then:  $\mathbb{E}[h(X_{T'}) + \sum_{t=T'+1}^{\tau} h(X_t) - h(Y_{t-1})] = H^*$ 

$$\frac{\langle X_T' \rangle + \langle L_{t=T'+1} n \langle X_t \rangle - n \langle T_{t-1} \rangle}{\text{unbiased estimate}} = n$$

Set-up & challenge of burn-in bias:

- Want to compute  $H^* = \int h(X)p_X(X)dX$
- Choose <u>Markov chain</u>:  $X_0, X_1, \dots \rightsquigarrow p_X$
- Problem: for any  $T<\infty\,,\mathbb{E}[h(X_T)]\neq H^*$
- "Naïve parallelism" gives little help!

Unbiased MCMC Set-Up (simplified)<sup>†</sup> Use "coupled chain" ( $Y_0, Y_1, ...$ ) where: 1.  $Y_t \sim X_t$ 2.  $Y_{t-1} = X_t$  for  $t > \tau$  ("meeting time") Then:  $\mathbb{E}[h(X_{T'}) + \sum_{t=T'+1}^{\tau} h(X_t) - h(Y_{t-1})] = H^*$ 

• Reduce error <u>with parallelism</u>

Set-up & challenge of burn-in bias:

- Want to compute  $H^* = \int h(X)p_X(X)dX$
- Choose <u>Markov chain</u>:  $X_0, X_1, \dots \rightsquigarrow p_X$
- Problem: for any  $T < \infty$  ,  $\mathbb{E}[h(X_T)] \neq H^*$
- "Naïve parallelism" gives little help!

Unbiased MCMC Set-Up (simplified)<sup>†</sup> Use "coupled chain" ( $Y_0, Y_1, ...$ ) where: 1.  $Y_t \sim X_t$ 2.  $Y_{t-1} = X_t$  for  $t > \tau$  ("meeting time") Then:  $\mathbb{E}[h(X_{T'}) + \sum_{t=T'+1}^{\tau} h(X_t) - h(Y_{t-1})] = H^*$ 

• Reduce error <u>with parallelism</u>

 $H^* = \lim_{T \to \infty} \mathbb{E}[h(X_T)]$ 

Set-up & challenge of burn-in bias:

- Want to compute  $H^* = \int h(X) p_X(X) dX$
- Choose <u>Markov chain</u>:  $X_0, X_1, \dots \rightsquigarrow p_X$
- Problem: for any  $T<\infty\,,\mathbb{E}[h(X_T)]\neq H^*$
- "Naïve parallelism" gives little help!

Unbiased MCMC Set-Up (simplified)<sup>†</sup> Use "coupled chain" ( $Y_0, Y_1, ...$ ) where: 1.  $Y_t \sim X_t$ 2.  $Y_{t-1} = X_t$  for  $t > \tau$  ("meeting time") Then:

$$\mathbb{E}\left[\frac{h(X_{T'}) + \sum_{t=T'+1}^{\tau} h(X_t) - h(Y_{t-1})}{\text{unbiased estimate}}\right] = H^*$$

$$\begin{split} H^* &= \lim_{T \to \infty} \mathbb{E}[h(X_T)] & \text{Telescopic sum} & \text{Reduce error with parallelism} \\ &= \lim_{T \to \infty} \mathbb{E}\left[h(X_{T'}) + \sum_{t=T'+1}^{T} h(X_t) - h(X_{t-1})\right] \end{split}$$

Set-up & challenge of burn-in bias:

- Want to compute  $H^* = \int h(X)p_X(X)dX$
- Choose <u>Markov chain</u>:  $X_0, X_1, \dots \rightsquigarrow p_X$
- Problem: for any  $T<\infty\,,\mathbb{E}[h(X_T)]\neq H^*$
- "Naïve parallelism" gives little help!

Unbiased MCMC Set-Up (simplified)<sup>†</sup> Use "coupled chain" ( $Y_0, Y_1, ...$ ) where: 1.  $Y_t \sim X_t$ 2.  $Y_{t-1} = X_t$  for  $t > \tau$  ("meeting time") Then:

$$\mathbb{E}\left[\frac{h(X_{T'}) + \sum_{t=T'+1}^{\tau} h(X_t) - h(Y_{t-1})}{\text{unbiased estimate}}\right] = H^*$$

• Reduce error <u>with parallelism</u>

$$H^* = \lim_{T \to \infty} \mathbb{E}[h(X_T)]$$
  
= 
$$\lim_{T \to \infty} \mathbb{E}\left[h(X_{T'}) + \sum_{t=T'+1}^T h(X_t) - h(X_{t-1})\right]$$

Set-up & challenge of burn-in bias:

- Want to compute  $H^* = \int h(X)p_X(X)dX$
- Choose <u>Markov chain</u>:  $X_0, X_1, \dots \rightsquigarrow p_X$
- Problem: for any  $T < \infty$  ,  $\mathbb{E}[h(X_T)] \neq H^*$
- "Naïve parallelism" gives little help!

Unbiased MCMC Set-Up (simplified)<sup>†</sup> Use "coupled chain" ( $Y_0, Y_1, ...$ ) where: 1.  $Y_t \sim X_t$ 2.  $Y_{t-1} = X_t$  for  $t > \tau$  ("meeting time") Then:

$$\mathbb{E}\left[\frac{h(X_{T'}) + \sum_{t=T'+1}^{\tau} h(X_t) - h(Y_{t-1})}{\text{unbiased estimate}}\right] = H^*$$

$$H^* = \lim_{T \to \infty} \mathbb{E}[h(X_T)]$$
$$= \lim_{T \to \infty} \mathbb{E}\left[h(X_T') + \sum_{t=T'+1}^T h(X_t) - \frac{h(Y_{t-1})}{P(t-1)}\right] \longleftarrow Y_{t-1} \sim X_{t-1}$$

Set-up & challenge of burn-in bias:

- Want to compute  $H^* = \int h(X) p_X(X) dX$
- Choose <u>Markov chain</u>:  $X_0, X_1, \dots \rightsquigarrow p_X$
- Problem: for any  $T<\infty\,,\mathbb{E}[h(X_T)]\neq H^*$
- "Naïve parallelism" gives little help!

Unbiased MCMC Set-Up (simplified)<sup>†</sup> Use "coupled chain" ( $Y_0, Y_1, ...$ ) where: 1.  $Y_t \sim X_t$ 2.  $Y_{t-1} = X_t$  for  $t > \tau$  ("meeting time") Then:

$$\mathbb{E}\left[\frac{h(X_{T'}) + \sum_{t=T'+1}^{\tau} h(X_t) - h(Y_{t-1})}{\text{unbiased estimate}}\right] = H^*$$

• Reduce error <u>with parallelism</u>



Set-up & challenge of burn-in bias:

- Want to compute  $H^* = \int h(X) p_X(X) dX$
- Choose <u>Markov chain</u>:  $X_0, X_1, \dots \rightsquigarrow p_X$
- Problem: for any  $T < \infty$  ,  $\mathbb{E}[h(X_T)] \neq H^*$
- "Naïve parallelism" gives little help!

Unbiased MCMC Set-Up (simplified)<sup>†</sup> Use "coupled chain" ( $Y_0, Y_1, ...$ ) where: 1.  $Y_t \sim X_t$ 2.  $Y_{t-1} = X_t$  for  $t > \tau$  ("meeting time") Then:

$$\mathbb{E}\left[\frac{h(X_{T'}) + \sum_{t=T'+1}^{\tau} h(X_t) - h(Y_{t-1})}{\text{unbiased estimate}}\right] = H^*$$

• Reduce error <u>with parallelism</u>

$$H^* = \lim_{T \to \infty} \mathbb{E}[h(X_T)]$$

$$= \lim_{T \to \infty} \mathbb{E}\left[h(X_{T'}) + \sum_{t=T'+1}^{T} h(X_t) - h(Y_{t-1})\right] \longleftarrow Y_{t-1} \sim X_{t-1}$$

$$= \mathbb{E}\left[h(X_{T'}) + \sum_{t=T'+1}^{\tau} h(X_t) - h(Y_{t-1})\right] \longleftarrow \text{Only finitely many non-zero terms.}$$

How do we apply this to clustering problems?

### Unbiased MCMC Set-Up

Use coupled chains such that

1. 
$$Y_t \sim X_t$$
  
2.  $Y_{t-1} = X_t$  for  $t > \tau$ 

Choices for Clustering Applications:

• Transition kernel for  $X_t$ 

### Unbiased MCMC Set-Up

Use coupled chains such that

1.  $Y_t \sim X_t$ 2.  $Y_{t-1} = X_t$  for  $t > \tau$ 

#### Choices for Clustering Applications:

• Transition kernel for  $X_t \rightarrow$  Gibbs

### Unbiased MCMC Set-Up

Use coupled chains such that

1. 
$$Y_t \sim X_t$$
  
2.  $Y_{t-1} = X_t$  for  $t > \tau$ 

#### Choices for Clustering Applications:

- Transition kernel for  $X_t \xrightarrow{\phantom{*}}$  Gibbs
- Coupling that meets quickly

Unbiased MCMC Set-Up Use coupled chains such that 1.  $Y_t \sim X_t$ 

2. 
$$Y_{t-1} = X_t$$
 for  $t > \tau$ 

Estimate (1 per processor)  $h(X_{T'}) + \sum_{t=T'+1}^{\tau} h(X_t) - h(Y_{t-1})$ Usual MCMC Bias estimate correction

Choices for Clustering Applications:

- Transition kernel for  $X_t \xrightarrow{\phantom{*}}$  Gibbs
- Coupling that meets quickly

Unbiased MCMC Set-Up Use coupled chains such that 1.  $Y_t \sim X_t$ 

2. 
$$Y_{t-1} = X_t$$
 for  $t > \tau$ 

Estimate (1 per processor)  $h(X_{T'}) + \sum_{t=T'+1}^{\tau} h(X_t) - h(Y_{t-1})$ Usual MCMC Bias estimate correction

#### Choices for Clustering Applications:

- Transition kernel for  $X_t \xrightarrow{\phantom{*}}$  Gibbs
- Coupling that meets quickly
  - large  $\tau \xrightarrow{}$  high variance

Unbiased MCMC Set-Up Use coupled chains such that 1.  $Y_t \sim X_t$ 

2. 
$$Y_{t-1} = X_t$$
 for  $t > \tau$ 



- Transition kernel for  $X_t \xrightarrow{\phantom{*}}$  Gibbs
- Coupling that meets quickly
  - large  $\tau \rightarrow$  high variance
  - not addressed by existing work!



Unbiased MCMC Set-Up Use coupled chains such that 1.  $Y_t \sim X_t$ 

2. 
$$Y_{t-1} = X_t$$
 for  $t > \tau$ 

#### Choices for Clustering Applications:

- Transition kernel for  $X_t \xrightarrow{\phantom{*}}$  Gibbs
- Coupling that meets quickly
  - large  $\tau \xrightarrow{}$  high variance
  - not addressed by existing work!

**Challenge:** the "label-switching" problem Equivalent re-labelings impede mixing









Unbiased MCMC Set-Up Use coupled chains such that 1.  $Y_t \sim X_t$ 

2. 
$$Y_{t-1} = X_t$$
 for  $t > \tau$ 

#### Choices for Clustering Applications:

- Transition kernel for  $X_t \xrightarrow{\phantom{*}}$  Gibbs
- Coupling that meets quickly
  - large  $\tau \xrightarrow{}$  high variance
  - not addressed by existing work!

**Challenge:** the "label-switching" problem Equivalent re-labelings impede mixing



Key idea: Develop a coupling that is agnostic to the labeling







### Coupling Gibbs Over Partitions via Optimal Transport

# Coupling Gibbs Over Partitions via Optimal Transport

We frame Gibbs samplers as over <u>partitions</u> instead of over labelings

- $X \sim p_{\Pi}(\cdot)$  is a random partition (e.g.  $X = \{\{1,3\}, \{2\}\}\}$ )
- Define X(-n) as leaving out n (e.g.  $X(-1) = \{\{2\}, \{3\}\}\}$ )

- We frame Gibbs samplers as over <u>partitions</u> instead of over labelings
- $X \sim p_{\Pi}(\cdot)$  is a random partition (e.g.  $X = \{\{1,3\},\{2\}\}\}$ )
- Define X(-n) as leaving out n (e.g.  $X(-1)=\{\{2\},\{3\}\})$
- Gibbs transition kernel:  $X_{t+1} \sim p_{\Pi \mid \Pi(-n)} (\cdot \mid X_t(-n))$

- $X \sim p_{\Pi}(\cdot)$  is a random partition (e.g.  $X = \{\{1,3\}, \{2\}\}\}$ )
- Define X(-n) as leaving out n (e.g.  $X(-1) = \{\{2\}, \{3\}\}\}$ )
- Gibbs transition kernel:  $X_{t+1} \sim p_{\Pi \mid \Pi(-n)} (\cdot \mid X_t(-n)) = \sum_k a_k \delta_{\pi^k}(\cdot)$

- $X \sim p_{\Pi}(\cdot)$  is a random partition (e.g.  $X = \{\{1,3\},\{2\}\}\}$ )
- Define X(-n) as leaving out n (e.g.  $X(-1) = \{\{2\}, \{3\}\}\}$ )
- Gibbs transition kernel:  $X_{t+1} \sim p_{\Pi \mid \Pi(-n)} (\cdot \mid X_t(-n)) = \sum_k a_k \delta_{\pi^k}(\cdot)$  $Y_t \sim p_{\Pi \mid \Pi(-n)} (\cdot \mid Y_{t-1}(-n)) = \sum_{k'} b_{k'} \delta_{\nu^{k'}}(\cdot)$

- $X \sim p_{\Pi}(\cdot)$  is a random partition (e.g.  $X = \{\{1,3\}, \{2\}\}\}$ )
- Define X(-n) as leaving out n (e.g.  $X(-1)=\{\{2\},\{3\}\})$
- Gibbs transition kernel:  $X_{t+1} \sim p_{\Pi|\Pi(-n)} (\cdot |X_t(-n)) = \sum_k a_k \delta_{\pi^k}(\cdot)$  $Y_t \sim p_{\Pi|\Pi(-n)} (\cdot |Y_{t-1}(-n)) = \sum_{k'} b_{k'} \delta_{\nu^{k'}}(\cdot)$

To couple  $X_{t+1}$  and  $Y_t$ , we use optimal transport:  $(X_{t+1}, Y_t) \sim \boldsymbol{\gamma}^*(p(\cdot | X_t), p(\cdot | Y_{t-1}))$ 

- $X \sim p_{\Pi}(\cdot)$  is a random partition (e.g.  $X = \{\{1,3\}, \{2\}\}\}$ )
- Define X(-n) as leaving out n (e.g.  $X(-1)=\{\{2\},\{3\}\})$
- Gibbs transition kernel:  $X_{t+1} \sim p_{\Pi|\Pi(-n)} (\cdot |X_t(-n)) = \sum_k a_k \delta_{\pi^k}(\cdot)$  $Y_t \sim p_{\Pi|\Pi(-n)} (\cdot |Y_{t-1}(-n)) = \sum_{k'} b_{k'} \delta_{\nu^{k'}}(\cdot)$

To couple  $X_{t+1}$  and  $Y_t$ , we use optimal transport:  $(X_{t+1}, Y_t) \sim \boldsymbol{\gamma}^* (p(\cdot | X_t), p(\cdot | Y_{t-1}))$ 

Strategy for  $\gamma^*$ : make  $X_{t+1}$  and  $Y_t$  as <u>close</u> as possible

- $X \sim p_{\Pi}(\cdot)$  is a random partition (e.g.  $X = \{\{1,3\}, \{2\}\}\}$ )
- Define X(-n) as leaving out n (e.g.  $X(-1)=\{\{2\},\{3\}\})$
- Gibbs transition kernel:  $X_{t+1} \sim p_{\Pi|\Pi(-n)} (\cdot |X_t(-n)) = \sum_k a_k \delta_{\pi^k}(\cdot)$  $Y_t \sim p_{\Pi|\Pi(-n)} (\cdot |Y_{t-1}(-n)) = \sum_{k'} b_{k'} \delta_{\nu^{k'}}(\cdot)$

To couple  $X_{t+1}$  and  $Y_t$ , we use optimal transport:  $(X_{t+1}, Y_t) \sim \boldsymbol{\gamma}^* (p(\cdot | X_t), p(\cdot | Y_{t-1}))$ Strategy for  $\boldsymbol{\gamma}^*$ : make  $X_{t+1}$  and  $Y_t$  as <u>close</u> as possible Need a metric:

- $X \sim p_{\Pi}(\cdot)$  is a random partition (e.g.  $X = \{\{1,3\}, \{2\}\}\}$ )
- Define X(-n) as leaving out n (e.g.  $X(-1) = \{\{2\}, \{3\}\}\}$ )
- Gibbs transition kernel:  $X_{t+1} \sim p_{\Pi|\Pi(-n)} (\cdot |X_t(-n)) = \sum_k a_k \delta_{\pi^k}(\cdot)$  $Y_t \sim p_{\Pi|\Pi(-n)} (\cdot |Y_{t-1}(-n)) = \sum_{k'} b_{k'} \delta_{\nu^{k'}}(\cdot)$

To couple  $X_{t+1}$  and  $Y_t$ , we use optimal transport:  $(X_{t+1}, Y_t) \sim \boldsymbol{\gamma}^* (p(\cdot | X_t), p(\cdot | Y_{t-1}))$ Strategy for  $\boldsymbol{\gamma}^*$ : make  $X_{t+1}$  and  $Y_t$  as close as possible

Need a metric: Use adjacency matrix  $\rightarrow$  Hamming distance

- $X \sim p_{\Pi}(\cdot)$  is a random partition (e.g.  $X = \{\{1,3\}, \{2\}\}\}$ )
- Define X(-n) as leaving out n (e.g.  $X(-1) = \{\{2\}, \{3\}\}\}$ )
- Gibbs transition kernel:  $X_{t+1} \sim p_{\Pi|\Pi(-n)} (\cdot |X_t(-n)) = \sum_k a_k \delta_{\pi^k}(\cdot)$  $Y_t \sim p_{\Pi|\Pi(-n)} (\cdot |Y_{t-1}(-n)) = \sum_{k'} b_{k'} \delta_{\nu^{k'}}(\cdot)$

To couple  $X_{t+1}$  and  $Y_t$ , we use optimal transport:  $(X_{t+1}, Y_t) \sim \boldsymbol{\gamma}^*(p(\cdot | X_t), p(\cdot | Y_{t-1}))$ 

Strategy for  $\gamma^*$ : make  $X_{t+1}$  and  $Y_t$  as <u>close</u> as possible Need a metric: Use adjacency matrix  $\rightarrow$  Hamming distance

OT  
Problem:  
s.t. 
$$\gamma \ge 0, \sum_{k} \gamma(\pi^{k}, \nu^{k'}) = b^{k'}, \sum_{k'} \gamma(\pi^{k}, \nu^{k'}) = a^{k}$$

Tosh and Dasgupta [2014]; Rand [1971]; Nguyen, Trippe, Broderick [2022]

- $X \sim p_{\Pi}(\cdot)$  is a random partition (e.g.  $X = \{\{1,3\}, \{2\}\}\}$ )
- Define X(-n) as leaving out n (e.g.  $X(-1) = \{\{2\}, \{3\}\}\}$ )
- Gibbs transition kernel:  $X_{t+1} \sim p_{\Pi|\Pi(-n)} (\cdot |X_t(-n)) = \sum_k a_k \delta_{\pi^k}(\cdot)$  $Y_t \sim p_{\Pi|\Pi(-n)} (\cdot |Y_{t-1}(-n)) = \sum_{k'} b_{k'} \delta_{\nu^{k'}}(\cdot)$

To couple  $X_{t+1}$  and  $Y_t$ , we use optimal transport:  $(X_{t+1}, Y_t) \sim \boldsymbol{\gamma}^* (p(\cdot | X_t), p(\cdot | Y_{t-1}))$ 

Strategy for  $\gamma^*$ : make  $X_{t+1}$  and  $Y_t$  as <u>close</u> as possible Need a metric: Use adjacency matrix  $\rightarrow$  Hamming distance



- $X \sim p_{\Pi}(\cdot)$  is a random partition (e.g.  $X = \{\{1,3\}, \{2\}\}\}$ )
- Define X(-n) as leaving out n (e.g.  $X(-1) = \{\{2\}, \{3\}\}\}$ )
- Gibbs transition kernel:  $X_{t+1} \sim p_{\Pi \mid \Pi(-n)} (\cdot \mid X_t(-n)) = \sum_k a_k \delta_{\pi^k}(\cdot)$  $Y_t \sim p_{\Pi \mid \Pi(-n)} (\cdot \mid Y_{t-1}(-n)) = \sum_{k'} b_{k'} \delta_{\eta k'} (\cdot)$

To couple  $X_{t+1}$  and  $Y_t$ , we use optimal transport:  $(X_{t+1}, Y_t) \sim \boldsymbol{\gamma}^* (p(\cdot | X_t), p(\cdot | Y_{t-1}))$ 

Strategy for  $\gamma^*$ : make  $X_{t+1}$  and  $Y_t$  as <u>close</u> as possible Need a metric: Use adjacency matrix  $\rightarrow$  Hamming distance

OT  
Problem:  
s.t. 
$$\gamma \ge 0$$
,  $\sum_{k} \sum_{k'} \gamma(\pi^{k}, \nu^{k'}) d_{\text{Hamming}}(\pi^{k}, \nu^{k'})$   
s.t.  $\gamma \ge 0$ ,  $\sum_{k} \gamma(\pi^{k}, \nu^{k'}) = b^{k'}$ ,  $\sum_{k'} \gamma(\pi^{k}, \nu^{k'}) = a^{k}$  is  
by construction, does not suffer from label switching!  
Tosh and Descurta [2014]: Rend [1971]: Nouven Trippe Broderick [2022]

Tosh and Dasgupta [2014]; Rand [1971]; Nguyen, Trippe, Broderick [2022]

- $X \sim p_{\Pi}(\cdot)$  is a random partition (e.g.  $X = \{\{1,3\}, \{2\}\}\}$ )
- Define X(-n) as leaving out n (e.g.  $X(-1) = \{\{2\}, \{3\}\}\}$ )
- Gibbs transition kernel:  $X_{t+1} \sim p_{\Pi|\Pi(-n)} (\cdot |X_t(-n)) = \sum_k a_k \delta_{\pi^k}(\cdot)$  $Y_t \sim p_{\Pi|\Pi(-n)} (\cdot |Y_{t-1}(-n)) = \sum_{k'} b_{k'} \delta_{\nu^{k'}}(\cdot)$

To couple  $X_{t+1}$  and  $Y_t$ , we use optimal transport:  $(X_{t+1}, Y_t) \sim \boldsymbol{\gamma}^* (p(\cdot | X_t), p(\cdot | Y_{t-1}))$ 

Strategy for  $\gamma^*$ : make  $X_{t+1}$  and  $Y_t$  as <u>close</u> as possible Need a metric: Use adjacency matrix  $\rightarrow$  Hamming distance

OT  

$$\gamma^{*} = \inf_{\gamma} \sum_{k} \sum_{k'} \gamma(\pi^{k}, \nu^{k'}) d_{\text{Hamming}}(\pi^{k}, \nu^{k'})$$
Problem:  
s.t.  $\gamma \ge 0, \sum_{k} \gamma(\pi^{k}, \nu^{k'}) = b^{k'}, \sum_{k'} \gamma(\pi^{k}, \nu^{k'}) = a^{k} \inf_{i \in \mathbb{T}} 5000$ 
• By construction, does not suffer from label switching!  
• We prove:  $\gamma^{*}$  permits unbiased estimation  
Tosh and Dasgupta [2014]; Rand [1971]; Nguyen, **Trippe**, Broderick [2022]

- Single-cell clustering Dirichlet process mixture model
- Run many pairs of coupled chains → compute meeting time distribution





OT vs. label-based  $10^{0}$ Single-cell clustering – Dirichlet process mixture model • Run many pairs of coupled Meeting-time chains  $\rightarrow$  compute meeting survival function  $10^{-2}$ time distribution (lower is better) CommonRNG • Consider label-based couplings: Maximal - Common RNG (Gibbs, 2004) OT - Maximal coupling (Jerrum, 1998)  $10^{3}$  $10^{1}$ 

Meeting Time (Sweeps)

**Coupling meeting-time** 

### Roadmap



- Parallelizing MCMC with Couplings:
  - Background & Notation
  - The Label Switching Problem
- We Frame Gibbs Sampling as Markov Chain on Partitions
- Our Optimal Transport Coupling
- Big-O Analysis Demonstrates Fast Computation
- Improved Estimation Error and Intervals with OTC over Naïve Parallelism in Practice

Can we compute our coupling fast enough?

• If coupling is too time intensive, we might prefer single chains

Can we compute our coupling fast enough?

- If coupling is too time intensive, we might prefer single chains
- Naïve computation of  $d_{\text{Hamming}}(\cdot, \cdot) \rightarrow O(N^2)$  time (let alone OT problem)

#### Can we compute our coupling fast enough?

- If coupling is too time intensive, we might prefer single chains
- Naïve computation of  $d_{\text{Hamming}}(\cdot, \cdot) \rightarrow O(N^2)$  time (let alone OT problem)

We show: can compute coupling in  $O(K^3\log K)$  amortized time!

• K: # of clusters --- typically fixed or  $O(\log N)$ 

#### Can we compute our coupling fast enough?

- If coupling is too time intensive, we might prefer single chains
- Naïve computation of  $d_{\text{Hamming}}(\cdot, \cdot) \rightarrow O(N^2)$  time (let alone OT problem)

We show: can compute coupling in  $O(K^3\log K)$  amortized time!

- K: # of clusters --- typically fixed or  $O(\log N)$
- Bottleneck : Orlin's algorithm in OT problem (but fast in practice)

#### Can we compute our coupling fast enough?

- If coupling is too time intensive, we might prefer single chains
- Naïve computation of  $d_{\text{Hamming}}(\cdot, \cdot) \rightarrow O(N^2)$  time (let alone OT problem)

We show: can compute coupling in  $O(K^3\log K)$  amortized time!

- K: # of clusters --- typically fixed or  $O(\log N)$
- Bottleneck : Orlin's algorithm in OT problem (but fast in practice)
- Compute cost dominated by marginal kernel in practice

### Can we compute our coupling fast enough?

- If coupling is too time intensive, we might prefer single chains
- Naïve computation of  $d_{\text{Hamming}}(\cdot, \cdot) \rightarrow O(N^2)$  time (let alone OT problem)

We show: can compute coupling in  $O(K^3\log K)$  amortized time!

- K: # of clusters --- typically fixed or  $O(\log N)$
- Bottleneck : Orlin's algorithm in OT problem (but fast in practice)
- Compute cost dominated by marginal kernel in practice

Additional challenge: Higher variance than single chains

### Can we compute our coupling fast enough?

- If coupling is too time intensive, we might prefer single chains
- Naïve computation of  $d_{\text{Hamming}}(\cdot, \cdot) \rightarrow O(N^2)$  time (let alone OT problem)

We show: can compute coupling in  $O(K^3\log K)$  amortized time!

- K: # of clusters --- typically fixed or  $O(\log N)$
- Bottleneck : Orlin's algorithm in OT problem (but fast in practice)
- Compute cost dominated by marginal kernel in practice

Additional challenge: Higher variance than single chains

- How many processors are needed?
- Previous works do not compare to naïve use of parallelism

- Coupled chains: aggregate estimates from multiple pairs of chains
- Naïve parallelism (baseline): average (biased) estimates from single chains

- Coupled chains: aggregate estimates from multiple pairs of chains
- Naïve parallelism (baseline): average (biased) estimates from single chains

Single-cell clustering (largest component proportion)



- Coupled chains: aggregate estimates from multiple pairs of chains
- Naïve parallelism (baseline): average (biased) estimates from single chains



• Further improvement with robust estimators (clipping outliers)

- Coupled chains: aggregate estimates from multiple pairs of chains
- Naïve parallelism (baseline): average (biased) estimates from single chains



• Further improvement with robust estimators (clipping outliers)

• Each process gives an i.i.d. sample  $\rightarrow$  use standard errors to form confidence intervals

- Each process gives an i.i.d. sample  $\rightarrow$  use standard errors to form confidence intervals
- Correct coverage with many processors

- Each process gives an i.i.d. sample  $\rightarrow$  use standard errors to form confidence intervals
- Correct coverage with many processors **Single-cell clustering**



- Each process gives an i.i.d. sample  $\rightarrow$  use standard errors to form confidence intervals
- Correct coverage with many processors **Single-cell clustering**



• Analogous "intervals" from single chains do not cover the estimand

- Each process gives an i.i.d. sample  $\rightarrow$  use standard errors to form confidence intervals
- Correct coverage with many processors



• Analogous "intervals" from single chains do not cover the estimand



Contact: tdn@mit.edu, btrippe@mit.edu, tamarab@mit.edu

#### Main References:

"Many processors, little time: MCMC for partitions via optimal transport couplings." Tin Nguyen, **Brian Trippe** & Tamara Broderick in *AISTATS (*2022)

"Optimal transport couplings of Gibbs samplers on partitions for unbiased estimation" **Brian Trippe\***, Tin Nguyen\* & Tamara Broderick in *AABI (*2021) [\*equal contribution]

**References Cited:** Jacob, O'Leary & Atchadé. "Unbiased Markov chain Monte Carlo methods with couplings." JRSS-B (2020); Tosh & Dasgupta. "Lower bounds for the Gibbs sampler over mixtures of Gaussians." ICML (2014).; Rand "Objective criteria for the evaluation of clustering methods." JASA (1971) 9/9