

# Sampling from Probabilistic Submodular Models

Alkis Gotovos



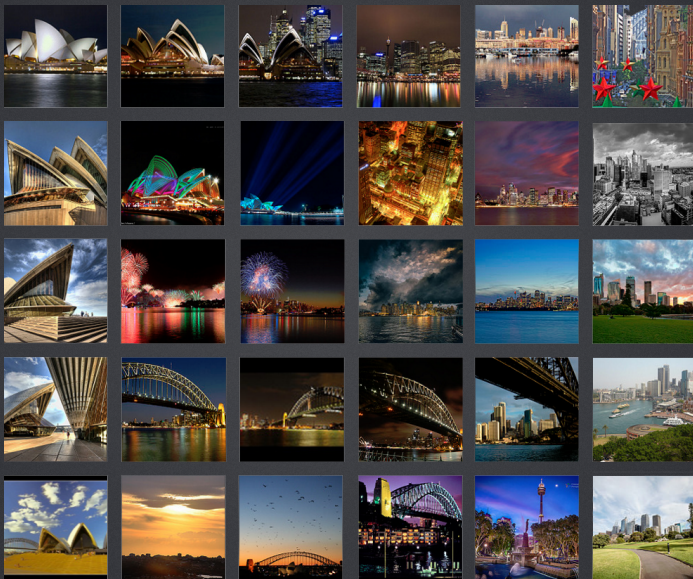
Hamed Hassani



Andreas Krause



# Image Collection Summarization



# Submodularity

- Facility location objective [Lin and Bilmes, '12] [Tschatschek et al., '14]

# Submodularity

- Facility location objective [Lin and Bilmes, '12] [Tschatschek et al., '14]
- Encourage coverage and diversity of the summary

# Submodularity

- Facility location objective [Lin and Bilmes, '12] [Tschatschek et al., '14]
- Encourage coverage and diversity of the summary
- Set of all images  $V$

# Submodularity

- Facility location objective [Lin and Bilmes, '12] [Tschachtschek et al., '14]
- Encourage coverage and diversity of the summary
- Set of all images  $V$
- For any summary  $S \subseteq V \longrightarrow F(S) \in \mathbb{R}$

# Submodularity

- Facility location objective [Lin and Bilmes, '12] [Tschachtschek et al., '14]
- Encourage coverage and diversity of the summary
- Set of all images  $V$
- For any summary  $S \subseteq V \longrightarrow F(S) \in \mathbb{R}$



# Submodularity

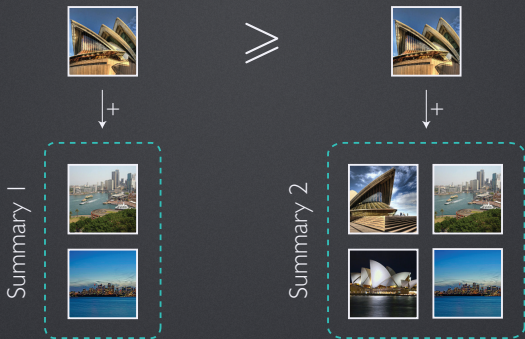
- Facility location objective [Lin and Bilmes, '12] [Tschatschek et al., '14]
- Encourage coverage and diversity of the summary
- Set of all images  $V$
- For any summary  $S \subseteq V \rightarrow F(S) \in \mathbb{R}$





# Submodularity

- Facility location objective [Lin and Bilmes, '12] [Tschatschek et al., '14]
- Encourage coverage and diversity of the summary
- Set of all images  $V$
- For any summary  $S \subseteq V \rightarrow F(S) \in \mathbb{R}$



# Submodularity

$F$  is submodular

$\iff$

$-F$  is supermodular

↓

coverage / diversity

↓

smoothness / cooperation

# Submodularity

$F$  is submodular



$-F$  is supermodular



coverage / diversity



smoothness / cooperation

- Submodular optimization is well-studied

# Submodularity

$F$  is submodular



$-F$  is supermodular



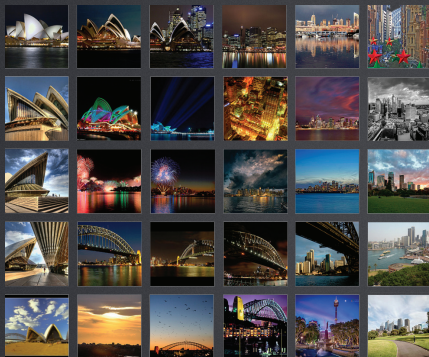
coverage / diversity



smoothness / cooperation

- Submodular optimization is well-studied
- Little existing work on probabilistic models

# Sampling Summaries

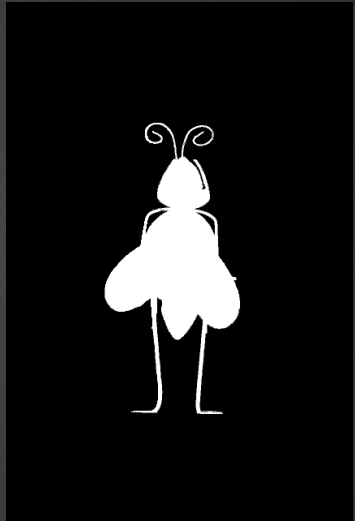


⋮

# Foreground / Background Segmentation



# Foreground / Background Segmentation

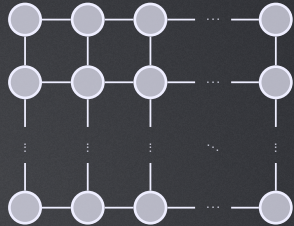


# Foreground / Background Segmentation

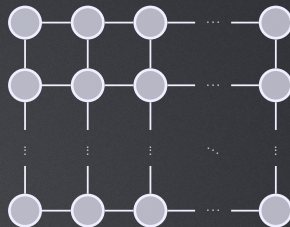




# Foreground / Background Segmentation

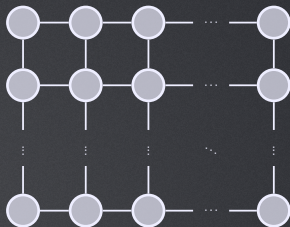


# Foreground / Background Segmentation



- Set of all pixels  $V$

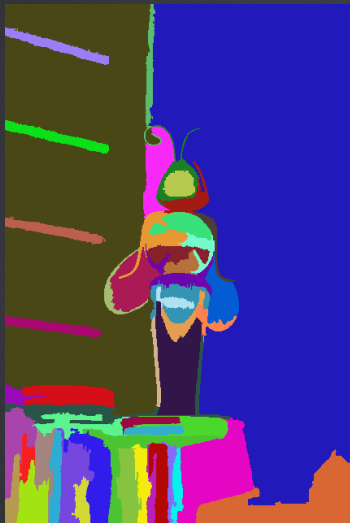
# Foreground / Background Segmentation



- Set of all pixels  $V$
- For  $S \subseteq V$  of foreground pixels,

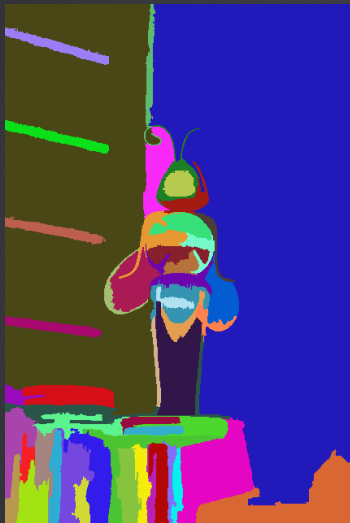
$$p(S) \propto \exp \left( \sum_{v \sim w} F_{v,w}(S) \right)$$

# Higher-order Models



Superpixel potentials [Kohli et al., '08]

# Higher-order Models



$$V = V_1 \cup V_2 \cup \dots \cup V_L$$

Superpixel potentials [Kohli et al., '08]

# Higher-order Models



$$V = V_1 \cup V_2 \cup \dots \cup V_L$$

Superpixel potentials [Kohli et al., '08]

# Higher-order Models



$$V = V_1 \cup V_2 \cup \dots \cup V_L$$

$$F_i(S) = \phi(|S \cap V_i|)$$

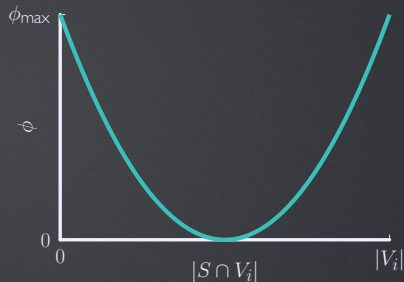
Superpixel potentials [Kohli et al., '08]

# Higher-order Models



$$V = V_1 \cup V_2 \cup \dots \cup V_L$$

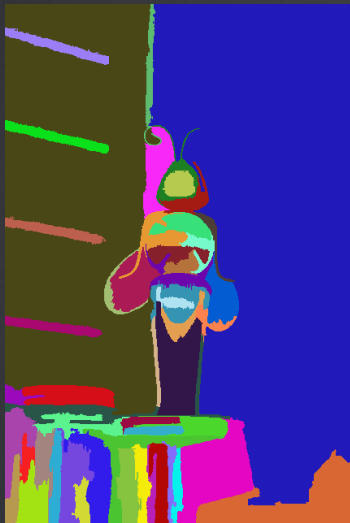
$$F_i(S) = \phi(|S \cap V_i|)$$



Superpixel potentials [Kohli et al., '08]



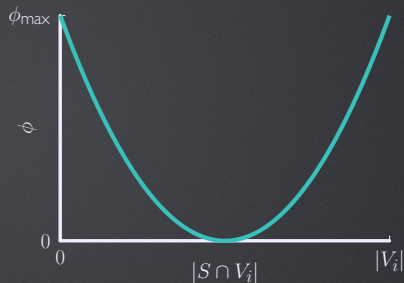
# Higher-order Models



Superpixel potentials [Kohli et al., '08]

$$V = V_1 \cup V_2 \cup \dots \cup V_L$$

$$F_i(S) = \phi(|S \cap V_i|)$$



$$p(S) \propto \exp\left(\sum_{i=1}^L F_i(S)\right)$$

# Higher-order Models [Djolonga and Krause, '15]

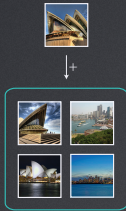
Pairwise



Higher-order



# Probabilistic Submodular Models

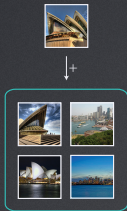


Use submodular functions  
in probabilistic models



Equip existing models with  
higher-order interactions

# Probabilistic Submodular Models

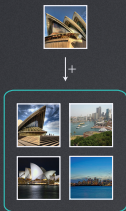


Use submodular functions  
in probabilistic models

Equip existing models with  
higher-order interactions

Probabilistic Submodular Models

# Probabilistic Submodular Models



Use submodular functions  
in probabilistic models

Equip existing models with  
higher-order interactions

Probabilistic Submodular Models

$$p(S) = \frac{1}{Z} \exp(F(S))$$

$F : 2^V \rightarrow \mathbb{R}$  is a submodular or supermodular function

# Probabilistic Submodular Models

PSMs

Markov Random Fields

---

# Probabilistic Submodular Models

PSMs

Markov Random Fields

---

- Ground set  $V$  with  $|V| = n$

# Probabilistic Submodular Models

PSMs

Markov Random Fields

---

- Ground set  $V$  with  $|V| = n$
- Sub- or supermodular function

$$F : 2^V \rightarrow \mathbb{R}$$



# Probabilistic Submodular Models

PSMs

Markov Random Fields

---

- Ground set  $V$  with  $|V| = n$
- Sub- or supermodular function

$$F : 2^V \rightarrow \mathbb{R}$$

- Distribution over subsets

$$p(S) \propto \exp(F(S))$$

---

# Probabilistic Submodular Models

PSMs

Markov Random Fields

---

- Ground set  $V$  with  $|V| = n$

- Binary random vector

$$X = (X_1, \dots, X_n)$$

- Sub- or supermodular function

$$F : 2^V \rightarrow \mathbb{R}$$

- Distribution over subsets

$$p(S) \propto \exp(F(S))$$

---

# Probabilistic Submodular Models

## PSMs

- Ground set  $V$  with  $|V| = n$

- Sub- or supermodular function

$$F : 2^V \rightarrow \mathbb{R}$$

- Distribution over subsets

$$p(S) \propto \exp(F(S))$$

## Markov Random Fields

- Binary random vector

$$X = (X_1, \dots, X_n)$$

- Set of factors

$$\phi_i : \{0, 1\}^{C_i} \rightarrow \mathbb{R}$$

# Probabilistic Submodular Models

## PSMs

- Ground set  $V$  with  $|V| = n$

- Sub- or supermodular function

$$F : 2^V \rightarrow \mathbb{R}$$

- Distribution over subsets

$$p(S) \propto \exp(F(S))$$

## Markov Random Fields

- Binary random vector

$$X = (X_1, \dots, X_n)$$

- Set of factors

$$\phi_i : \{0, 1\}^{C_i} \rightarrow \mathbb{R}$$

- Distribution over binary vectors

$$p(X) \propto \exp\left(\sum_i \phi_i(X_{C_i})\right)$$

# Probabilistic Submodular Models

## PSMs

- Ground set  $V$  with  $|V| = n$

- Sub- or supermodular function

$$F : 2^V \rightarrow \mathbb{R}$$

- Distribution over subsets

$$p(S) \propto \exp(F(S))$$

## Markov Random Fields

- Binary random vector

$$X = (X_1, \dots, X_n)$$

- Set of factors

$$\phi_i : \{0, 1\}^{\mathcal{C}_i} \rightarrow \mathbb{R}$$

- Distribution over binary vectors

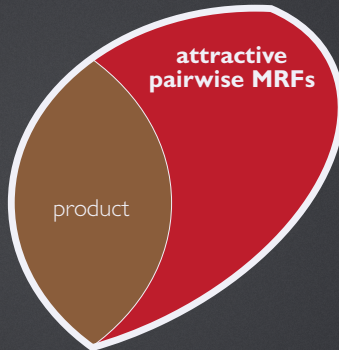
$$p(X) \propto \exp\left(\sum_i \phi_i(X_{\mathcal{C}_i})\right)$$

Model order:  $\max_i |\mathcal{C}_i|$

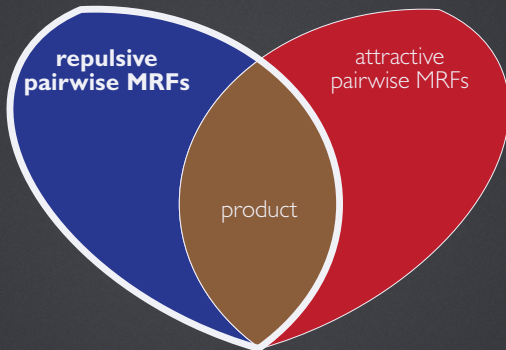
# Landscape of Models



# Landscape of Models

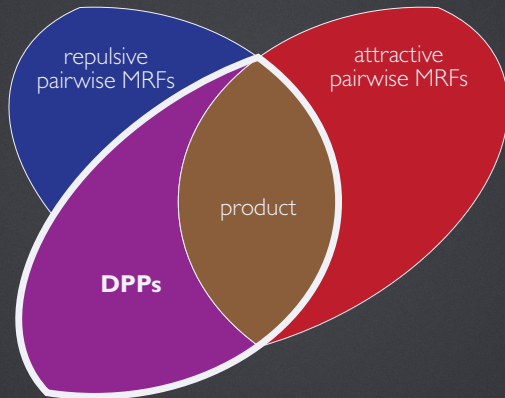


# Landscape of Models

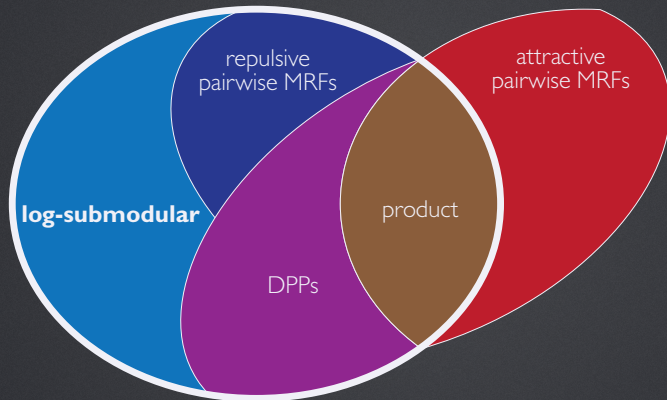




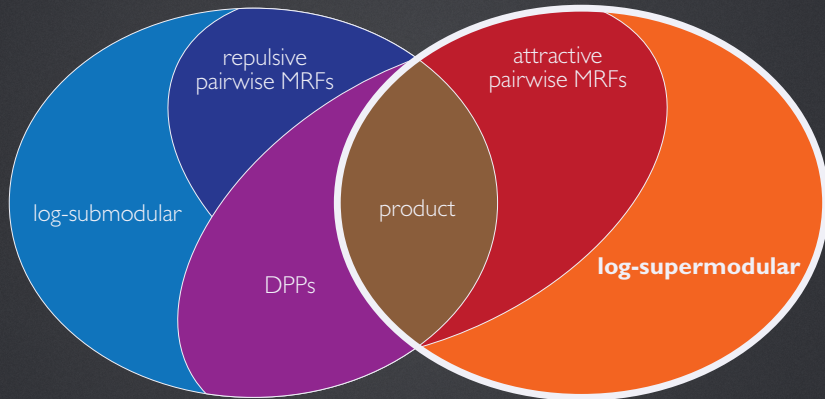
# Landscape of Models



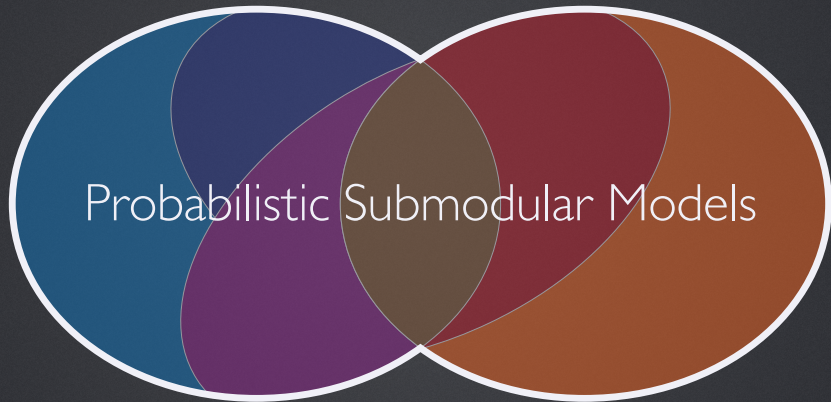
# Landscape of Models



# Landscape of Models



# Landscape of Models



# Inference

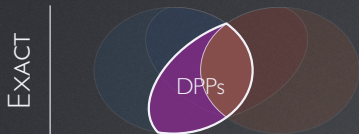
$\mathbb{P}(\text{pixel label})$



$\mathbb{P}(\text{image} \in \text{summary} \mid \text{selected})$



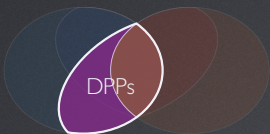
# Inference



- Tractable only for limited subclasses
- #P-hard even for Ising models

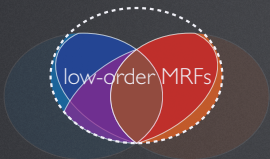
# Inference

EXACT



- Tractable only for limited subclasses
- #P-hard even for Ising models

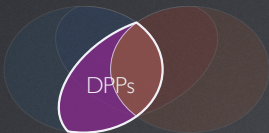
BP, MF, ...



- Extensively studied model class
- Complexity exponential in model order

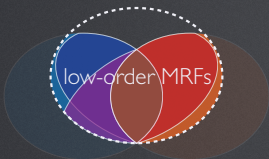
# Inference

EXACT



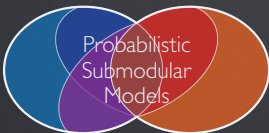
- Tractable only for limited subclasses
- #P-hard even for Ising models

BP, MF, ...



- Extensively studied model class
- Complexity exponential in model order

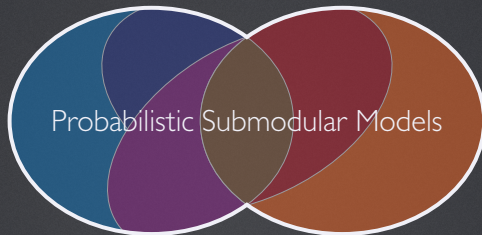
L-FIELD



- Variational approach for general PSMs  
[Djolonga and Krause, '14]



# Inference



What about sampling?

# Markov Chain Monte Carlo

- State space  $\Omega$
- Transition matrix  $P$
- Stationary distribution  $\pi$

# Markov Chain Monte Carlo

- State space  $\Omega$
- Transition matrix  $P$
- Stationary distribution  $\pi$

Markov chain  $(S_t)_{t \geq 0}$  that moves according to  $P$

# Markov Chain Monte Carlo

- State space  $\Omega$  powerset of  $V$
- Transition matrix  $P$
- Stationary distribution  $\pi$

Markov chain  $(S_t)_{t \geq 0}$  that moves according to  $P$

# Markov Chain Monte Carlo

- State space  $\Omega$                       powerset of  $V$
- Transition matrix  $P$                       Gibbs sampler
- Stationary distribution  $\pi$

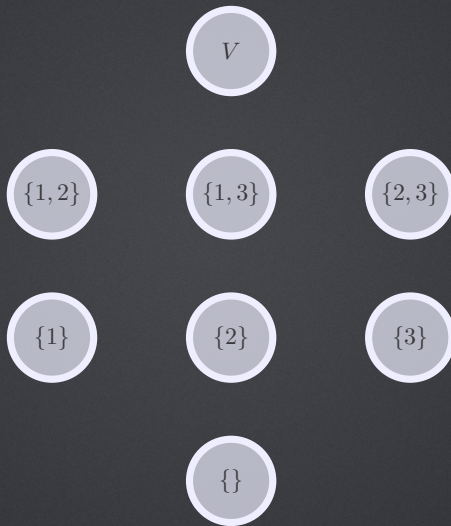
Markov chain  $(S_t)_{t \geq 0}$  that moves according to  $P$

# Markov Chain Monte Carlo

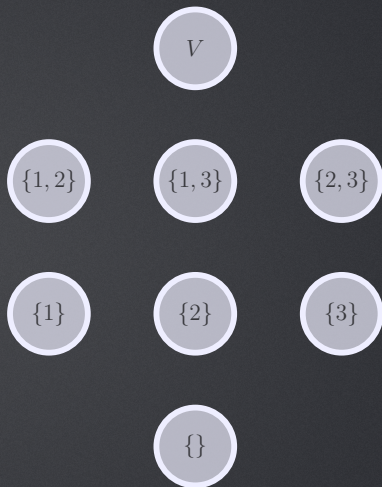
- State space  $\Omega$                       powerset of  $V$
- Transition matrix  $P$                       Gibbs sampler
- Stationary distribution  $\pi$                       PSM distribution

Markov chain  $(S_t)_{t \geq 0}$  that moves according to  $P$

# State Space of $V = \{1, 2, 3\}$



# Gibbs Sampler

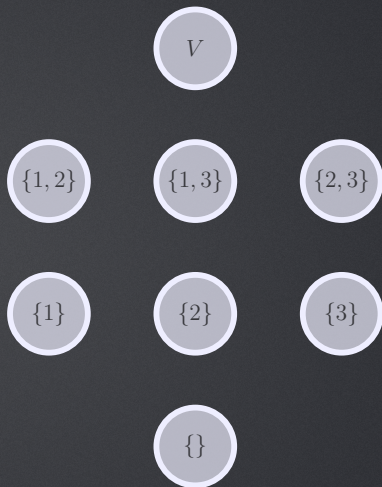




# Gibbs Sampler

Start at  $S_0$

For  $t = 1, 2, \dots$



# Gibbs Sampler

Start at  $S_0$

For  $t = 1, 2, \dots$



# Gibbs Sampler

Start at  $S_0$

For  $t = 1, 2, \dots$

- Select random  $v \in V$

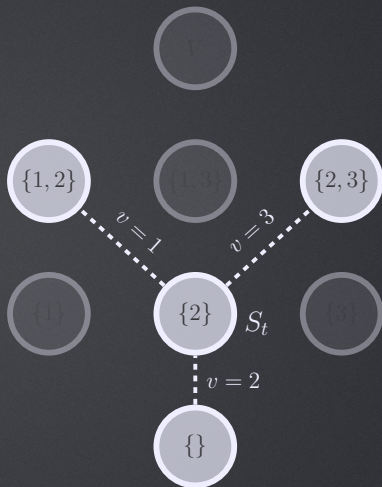


# Gibbs Sampler

Start at  $S_0$

For  $t = 1, 2, \dots$

- Select random  $v \in V$

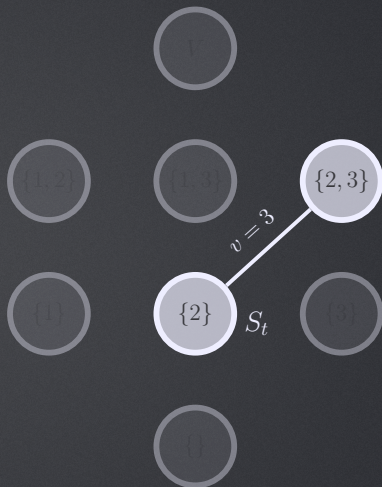


# Gibbs Sampler

Start at  $S_0$

For  $t = 1, 2, \dots$

- Select random  $v \in V$

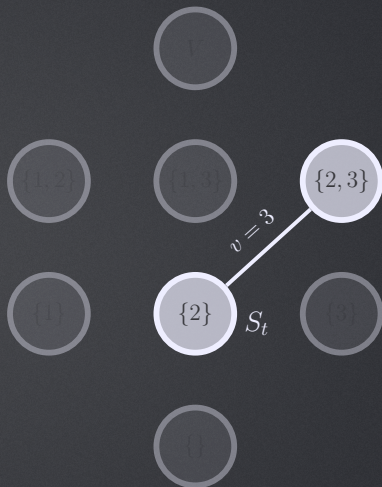


# Gibbs Sampler

Start at  $S_0$

For  $t = 1, 2, \dots$

- Select random  $v \in V$
- Compute conditional  $p_{\text{add}}$



# Gibbs Sampler

Start at  $S_0$

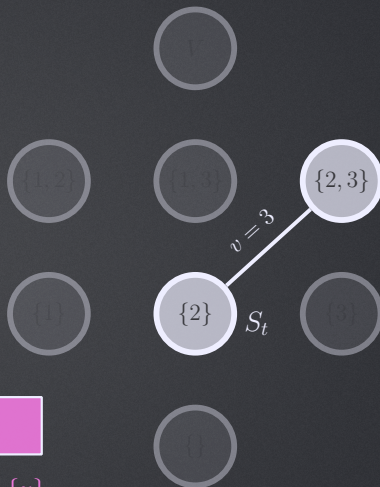
For  $t = 1, 2, \dots$

- Select random  $v \in V$
- Compute conditional  $p_{\text{add}}$
- Flip biased coin



$$S_{t+1} \leftarrow S_t \cup \{v\}$$

$$S_{t+1} \leftarrow S_t \setminus \{v\}$$



# Gibbs Sampler

Start at  $S_0$

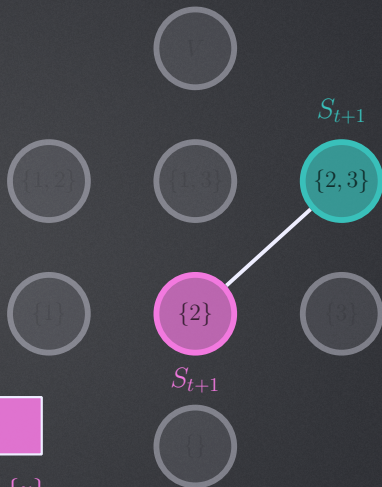
For  $t = 1, 2, \dots$

- Select random  $v \in V$
- Compute conditional  $p_{\text{add}}$
- Flip biased coin



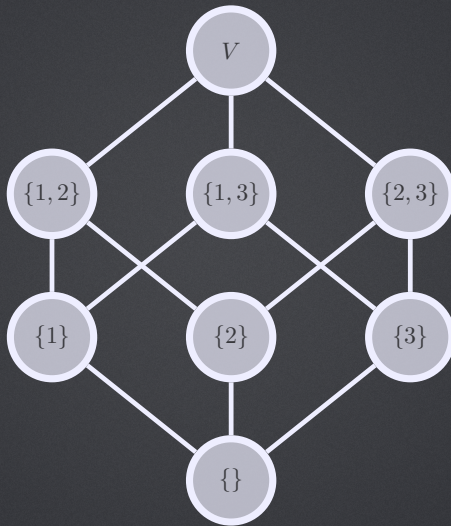
$$S_{t+1} \leftarrow S_t \cup \{v\}$$

$$S_{t+1} \leftarrow S_t \setminus \{v\}$$

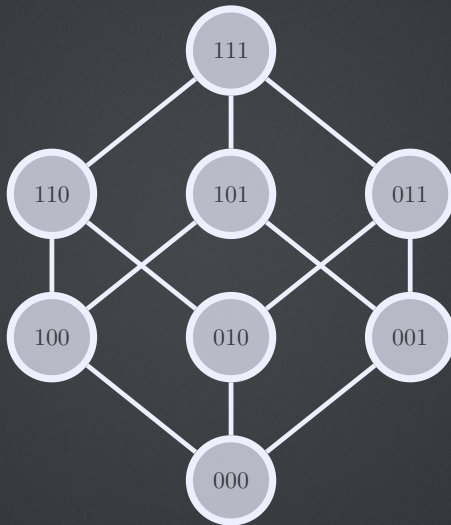




# Gibbs Sampler



# Gibbs Sampler



## Mixing time

Does the Markov chain converge?

# Mixing time

Does the Markov chain converge?

Total variation distance

$$d(t) = d_{\text{tv}}(\mathbb{P}_{S_t}, \pi)$$

# Mixing time

Does the Markov chain converge?

Total variation distance

$$d(t) = \max\{d_{\text{TV}}(\mathbb{P}_{S_t}, \pi) \mid S_0 \in \Omega\}$$

# Mixing time

Does the Markov chain converge?

Total variation distance

$$d(t) = \max\{d_{\text{TV}}(\mathbb{P}_{S_t}, \pi) \mid S_0 \in \Omega\}$$

Under mild assumptions (ergodicity),  $d(t) \xrightarrow{t \rightarrow \infty} 0$

# Mixing time

Does the Markov chain converge?

Total variation distance

$$d(t) = \max\{d_{\text{TV}}(\mathbb{P}_{S_t}, \pi) \mid S_0 \in \Omega\}$$

Under mild assumptions (ergodicity),  $d(t) \xrightarrow{t \rightarrow \infty} 0$

How long does it take to get “close enough” to  $\pi$ ?

# Mixing time

Does the Markov chain converge?

Total variation distance

$$d(t) = \max\{d_{\text{TV}}(\mathbb{P}_{S_t}, \pi) \mid S_0 \in \Omega\}$$

Under mild assumptions (ergodicity),  $d(t) \xrightarrow{t \rightarrow \infty} 0$

How long does it take to get “close enough” to  $\pi$ ?

Mixing time  $t_{\text{mix}}(\epsilon) = \min\{t \mid d(t) \leq \epsilon\}$



# Goal

- Mixing times for general PSMs are exponential in  $|V| = n$

# Goal

- Mixing times for general PSMs are exponential in  $|V| = n$
- Exponential even for pairwise models [Jerrum and Sinclair, '93]

# Goal

- Mixing times for general PSMs are exponential in  $|V| = n$
- Exponential even for pairwise models [Jerrum and Sinclair, '93]

We establish sufficient conditions for sub-exponential mixing of the Gibbs sampler on PSMs.

## Polynomial-time Mixing

$$F \text{ is modular if } F(A) + F(B) = F(A \cup B) + F(A \cap B)$$

# Polynomial-time Mixing

$$\begin{array}{l} \text{sub-} \\ F \text{ is modular if} \end{array} \quad F(A) + F(B) \stackrel{\geq}{=} F(A \cup B) + F(A \cap B)$$

# Polynomial-time Mixing

$$\begin{array}{l} \text{sub-} \\ F \text{ is modular if } \\ \text{super-} \end{array} \quad F(A) + F(B) \begin{array}{l} \geq \\ = \\ \leq \end{array} F(A \cup B) + F(A \cap B)$$

# Polynomial-time Mixing

$$\begin{array}{l} \text{sub-} \\ F \text{ is modular if } \\ \text{super-} \end{array} \quad F(A) + F(B) \begin{array}{l} \geq \\ = \\ \leq \end{array} F(A \cup B) + F(A \cap B)$$

“Distance” from modularity

$$|F(A) + F(B) - F(A \cup B) - F(A \cap B)|$$

# Polynomial-time Mixing

$$\begin{array}{l} \text{sub-} \\ F \text{ is modular if } \\ \text{super-} \end{array} \quad F(A) + F(B) \begin{array}{l} \geq \\ = \\ \leq \end{array} F(A \cup B) + F(A \cap B)$$

“Distance” from modularity

$$\zeta_F := \max_{A, B \subseteq V} |F(A) + F(B) - F(A \cup B) - F(A \cap B)|$$



# Polynomial-time Mixing

submodular

constant

modular

normalized  
monotone  
submodular

$$F(S) = c + \sum_{v \in S} m_v + f(S)$$

# Polynomial-time Mixing

submodular

constant

modular

normalized  
monotone  
submodular

$$F(S) = c + \sum_{v \in S} m_v + f(S)$$



$$\exp(F(S)) \propto$$

PSM



$$\prod_{v \in S} \exp(m_v) \times$$

product



$$\exp(f(S))$$

interactions

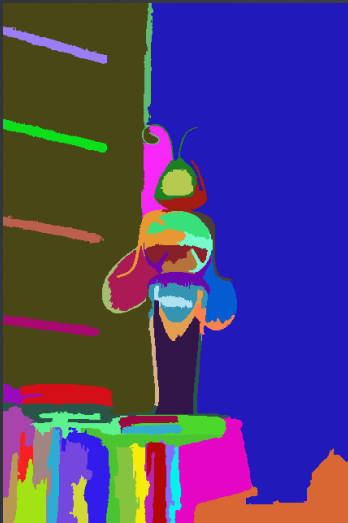
# Polynomial-time Mixing

## Theorem 1

For any submodular or supermodular set function  $F$ , the mixing time of the Gibbs sampler is bounded by

$$t_{\text{mix}}(\epsilon) = \mathcal{O} \left( n^2 \exp(\zeta_f) \log \epsilon^{-1} \right).$$

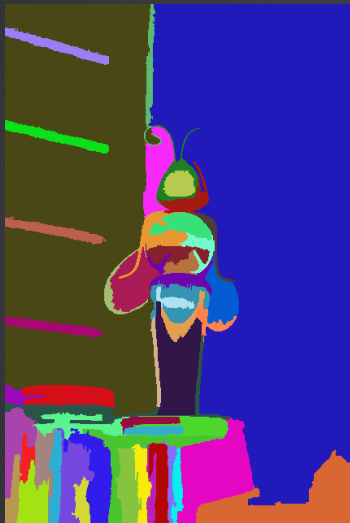
# Polynomial-time Mixing



$$F_i(S) = \phi(|S \cap V_i|)$$

$$F(S) = \sum_{i=1}^L F_i(S)$$

# Polynomial-time Mixing



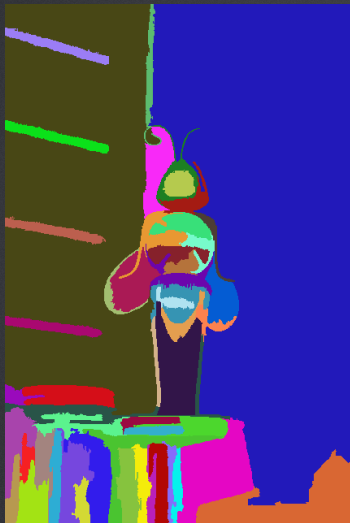
$$F_i(S) = \phi(|S \cap V_i|)$$

$$F(S) = \sum_{i=1}^L F_i(S)$$

Easy to show that

$$\zeta_f \leq L\phi_{\max}$$

# Polynomial-time Mixing



$$F_i(S) = \phi(|S \cap V_i|)$$

$$F(S) = \sum_{i=1}^L F_i(S)$$

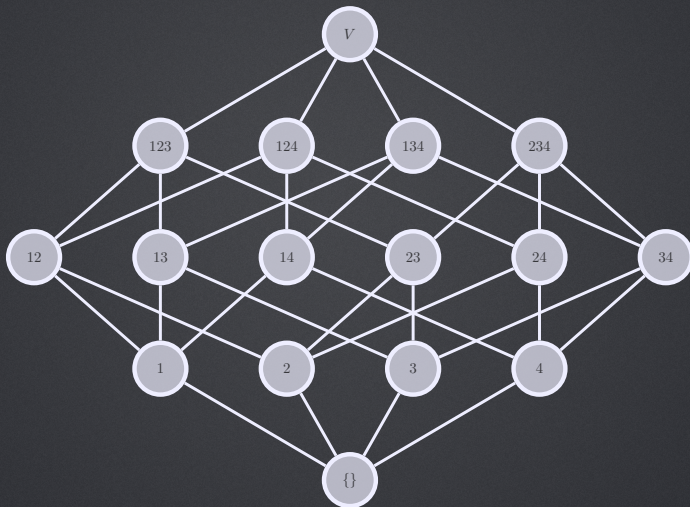
Easy to show that

$$\zeta_f \leq L\phi_{\max}$$

$$|V_i| \approx 10^5 \text{ vs. } L \approx 50$$

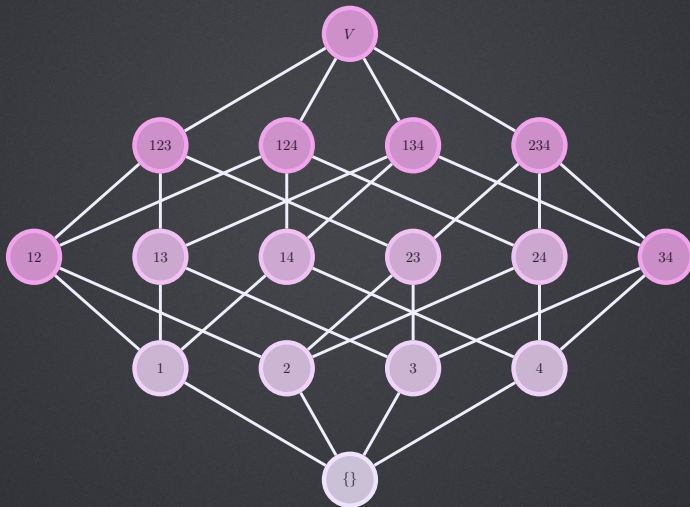
# Proof Outline

Method of canonical paths [Sinclair, '92]



# Proof Outline

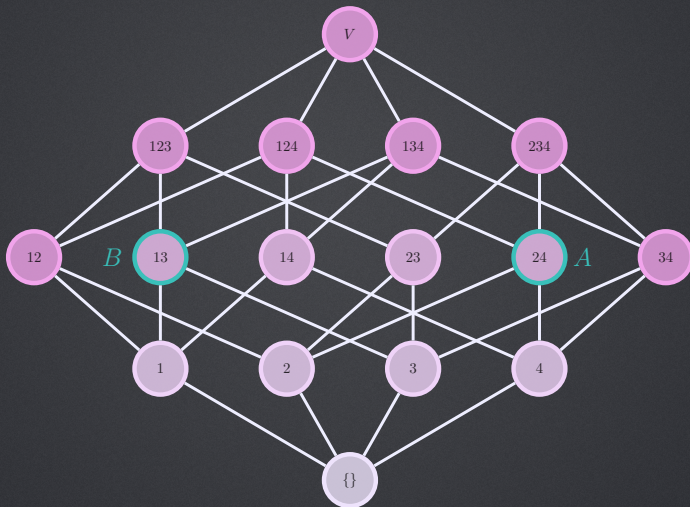
For each  $A, B \subseteq V$ , need to route  $p(A)p(B)$  amount of flow





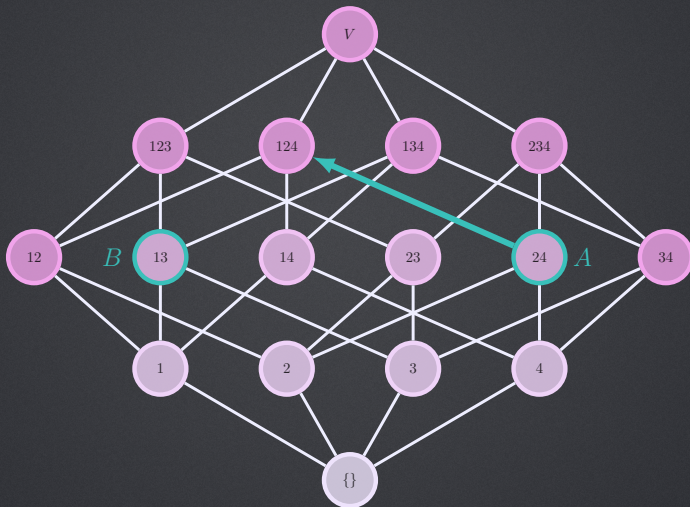
# Proof Outline

Construct a “canonical path” between each  $A, B \subseteq V$



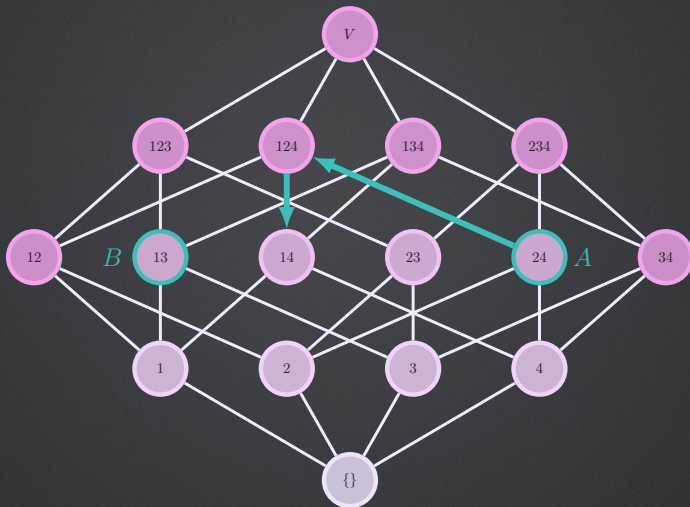
# Proof Outline

Construct a “canonical path” between each  $A, B \subseteq V$



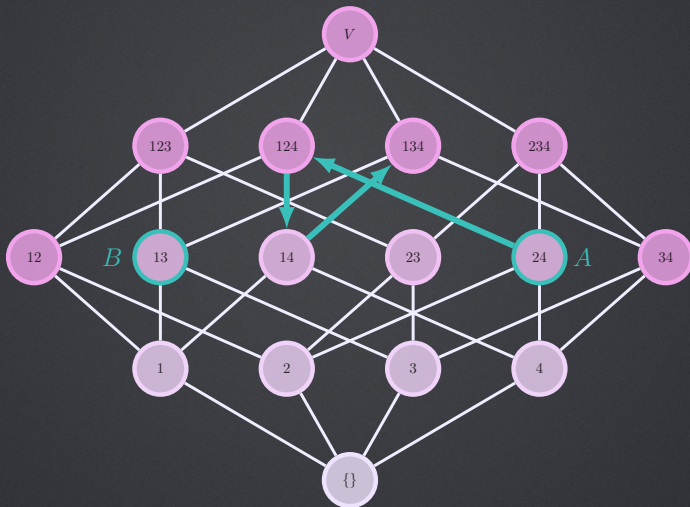
# Proof Outline

Construct a “canonical path” between each  $A, B \subseteq V$



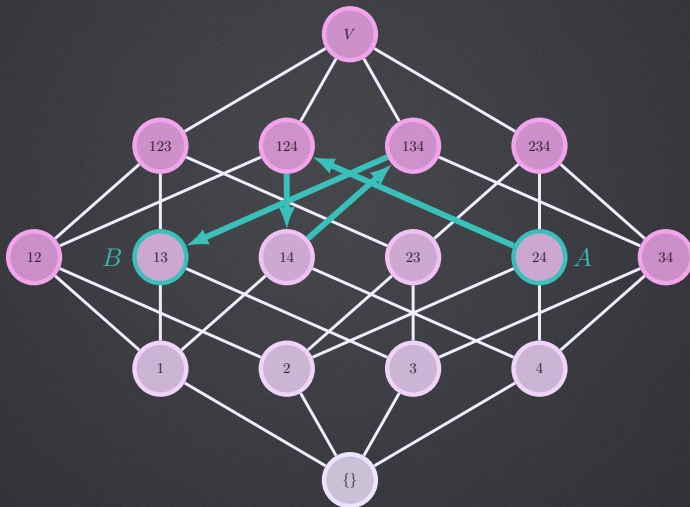
# Proof Outline

Construct a “canonical path” between each  $A, B \subseteq V$



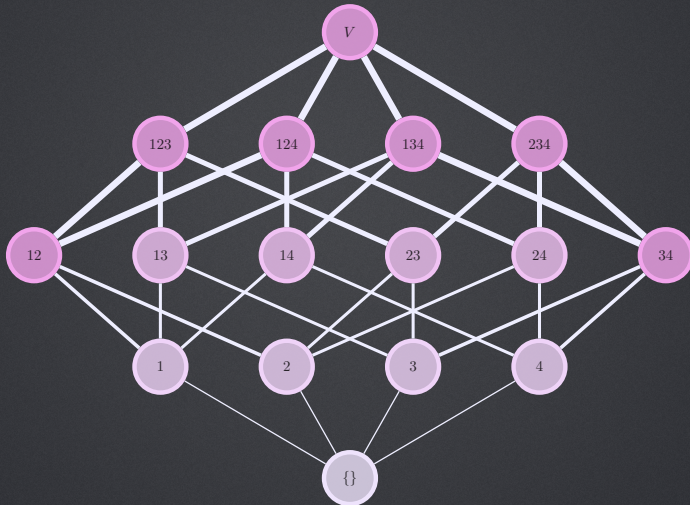
# Proof Outline

Construct a “canonical path” between each  $A, B \subseteq V$



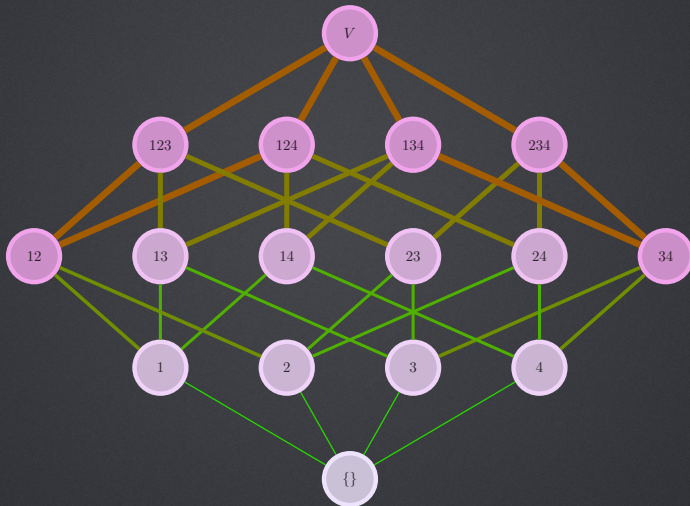
# Proof Outline

Capacity of an edge  $\sim$  transition probability of Gibbs sampler



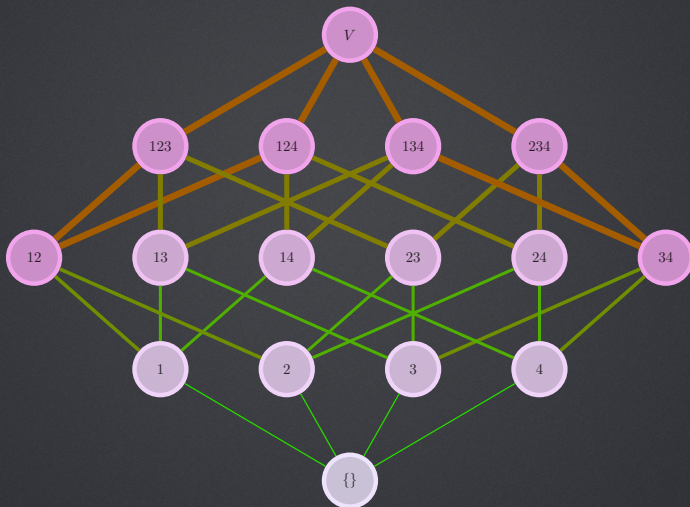
# Proof Outline

Congestion of an edge  $\sim$  (total flow through edge) / capacity



# Proof Outline

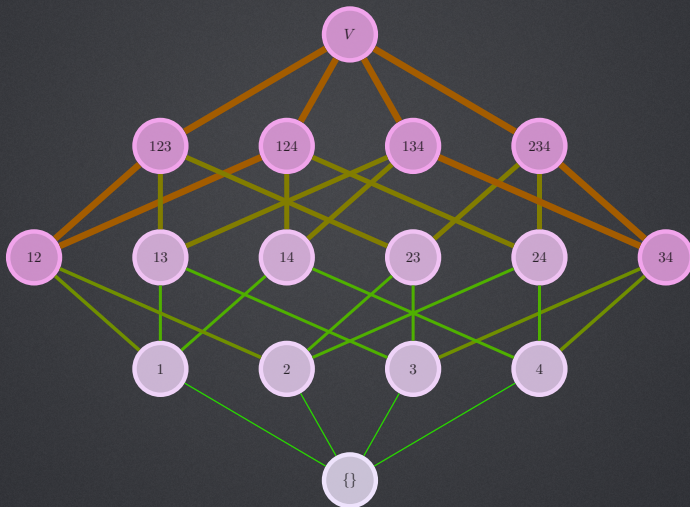
$$t_{\text{mix}}(\epsilon) = \mathcal{O}(\max\{\text{congestion}\} \log \epsilon^{-1}) \text{ [Sinclair, '92]}$$





# Proof Outline

We bound the maximum congestion of a PSM using  $\zeta_f$



# Fast Mixing

## Theorem 2

For any submodular or supermodular set function  $F$ , if  $\gamma_f < 1$ , the mixing time of the Gibbs sampler is bounded by

$$t_{\text{mix}}(\epsilon) \leq \frac{1}{1 - \gamma_f} n (\log n + \log \epsilon^{-1}).$$

- $\gamma_f$  = “maximum total influence”

# Fast Mixing

## Theorem 2

For any submodular or supermodular set function  $F$ , if  $\gamma_f < 1$ , the mixing time of the Gibbs sampler is bounded by

$$t_{\text{mix}}(\epsilon) \leq \frac{1}{1 - \gamma_f} n (\log n + \log \epsilon^{-1}).$$

- $\gamma_f$  = “maximum total influence”
- Simple way to bound  $\gamma_f$ , if  $f(S) = \sum_i f_i(S)$

# Fast Mixing

## Theorem 2

For any submodular or supermodular set function  $F$ , if  $\gamma_f < 1$ , the mixing time of the Gibbs sampler is bounded by

$$t_{\text{mix}}(\epsilon) \leq \frac{1}{1 - \gamma_f} n (\log n + \log \epsilon^{-1}).$$

- $\gamma_f$  = “maximum total influence”
- Simple way to bound  $\gamma_f$ , if  $f(S) = \sum_i f_i(S)$
- Closely related to Dobrushin uniqueness conditions, and influence matrix norms [Dyer et al., '09]

# Fast Mixing

## Theorem 2

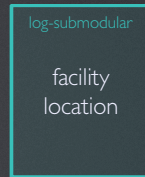
For any submodular or supermodular set function  $F$ , if  $\gamma_f < 1$ , the mixing time of the Gibbs sampler is bounded by

$$t_{\text{mix}}(\epsilon) \leq \frac{1}{1 - \gamma_f} n (\log n + \log \epsilon^{-1}) .$$

- $\gamma_f$  = “maximum total influence”
- Simple way to bound  $\gamma_f$ , if  $f(S) = \sum_i f_i(S)$
- Closely related to Dobrushin uniqueness conditions, and influence matrix norms [Dyer et al., '09]
- Similar theorem by [Rebeschini and Karbasi, '15]

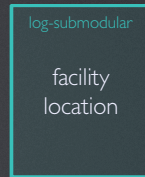
# Evaluation

Compare against variational approach [Djolonga and Krause, '14]



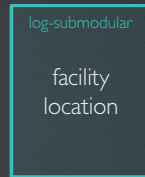
# Evaluation

Compare against variational approach [Djolonga and Krause, '14]



# Evaluation

Compare against variational approach [Djolonga and Krause, '14]

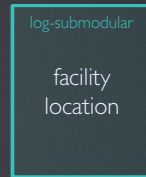


- Compute  $p(v | S)$



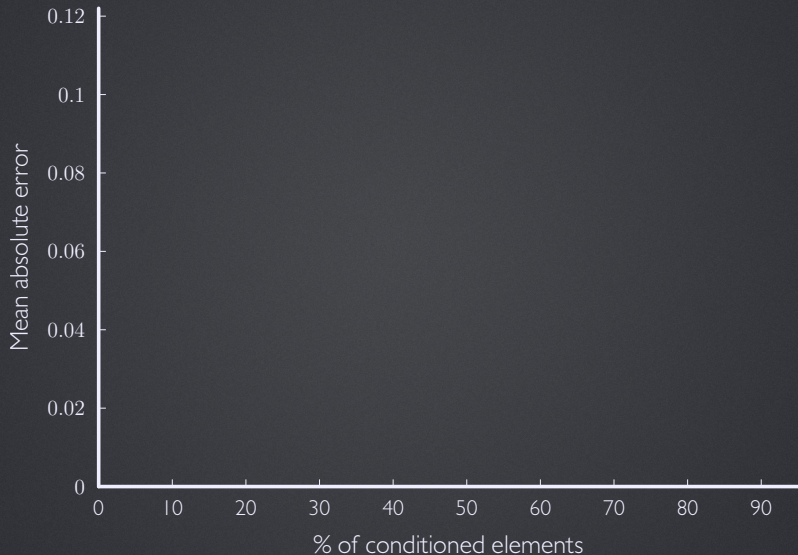
# Evaluation

Compare against variational approach [Djolonga and Krause, '14]

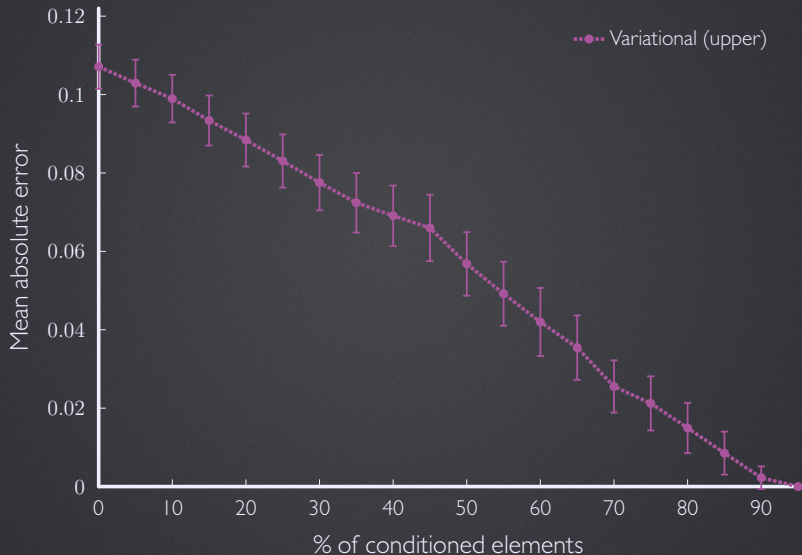


- Compute  $p(v \mid S)$
- $|V| = 20 \longrightarrow$  compare to exact marginals

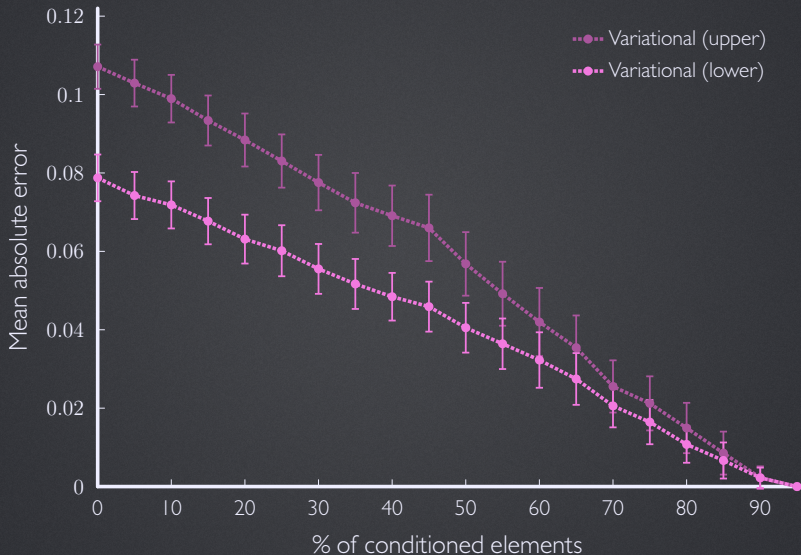
# Evaluation



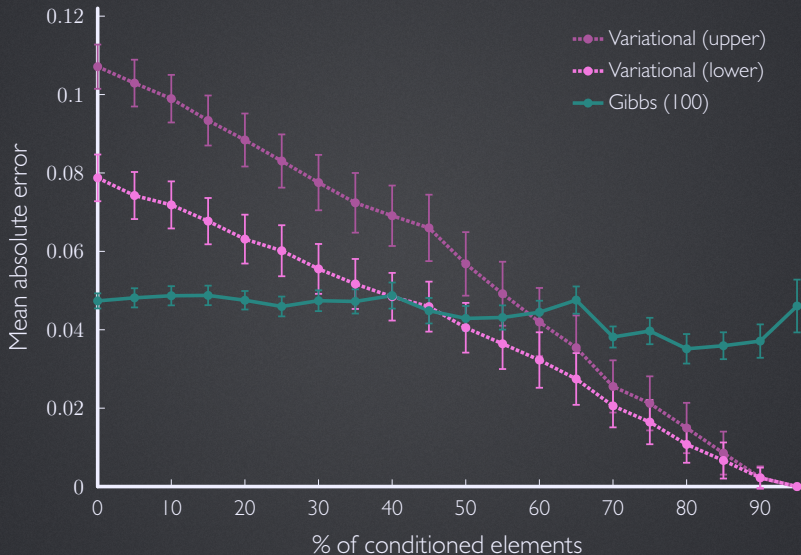
# Evaluation



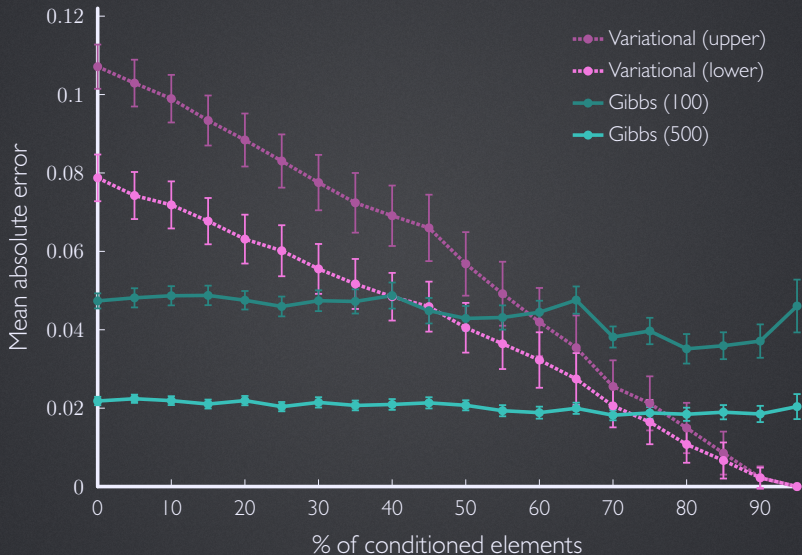
# Evaluation



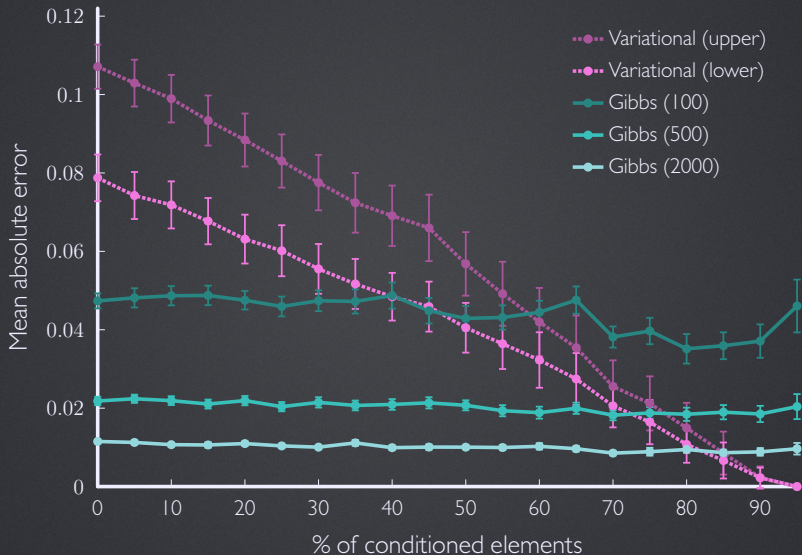
# Evaluation



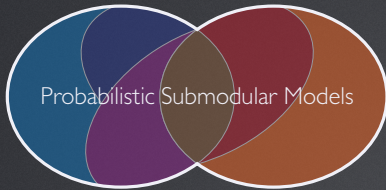
# Evaluation



# Evaluation

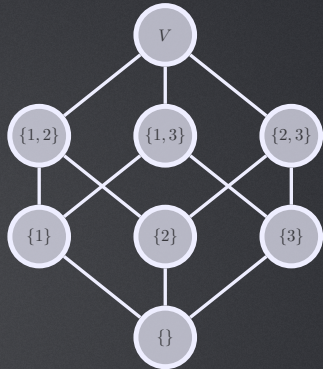
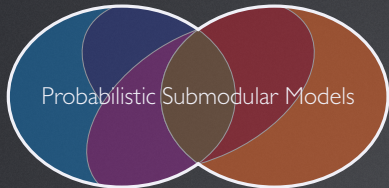


# Conclusion

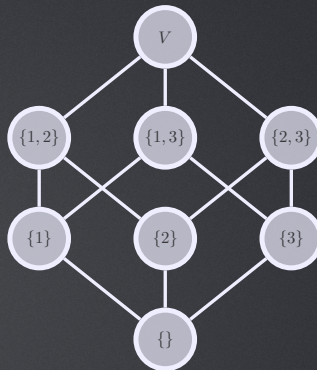
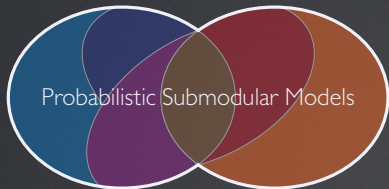




# Conclusion

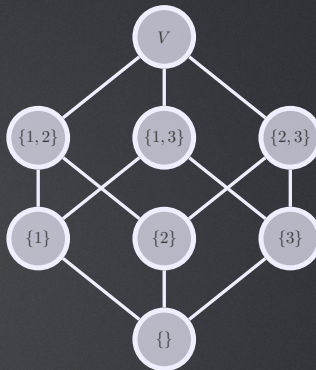
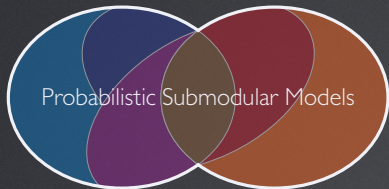


# Conclusion



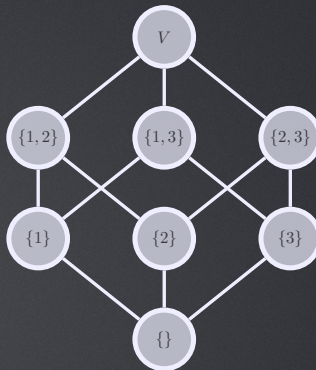
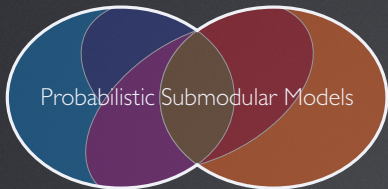
- Identify higher-order models amenable to efficient inference

# Conclusion



- Identify higher-order models amenable to efficient inference
- First indications that sub-/supermodularity can lead to faster mixing

# Conclusion



- Identify higher-order models amenable to efficient inference
- First indications that sub-/supermodularity can lead to faster mixing

Poster #70

# Backup I

Start at  $S_0$

For  $t = 1, 2, \dots$

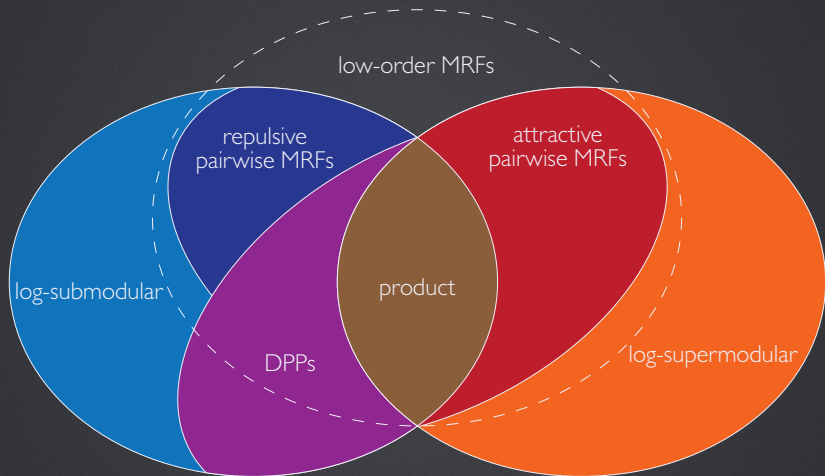
- Select random  $v \in V$
- $\Delta \leftarrow F(S_t \cup \{v\}) - F(S_t \setminus \{v\})$
- $p_{\text{add}} \leftarrow e^{\Delta} / (1 + e^{\Delta})$
- Flip biased coin



$$S_{t+1} \leftarrow S_t \cup \{v\}$$

$$S_{t+1} \leftarrow S_t \setminus \{v\}$$

## Backup II



## Backup III

### Theorem I

For any set function  $F$ , the mixing time of the Gibbs sampler is bounded by

$$t_{\text{mix}}(\epsilon) = \mathcal{O} \left( n^2 \exp(2\zeta_F) \log \epsilon^{-1} \right).$$

For any **submodular or supermodular** set function  $F$ , the mixing time of the Gibbs sampler is bounded by

$$t_{\text{mix}}(\epsilon) = \mathcal{O} \left( n^2 \exp(\zeta_f) \log \epsilon^{-1} \right).$$