# DETERMINANTAL POINT PROCESSES FOR NATURAL LANGUAGE PROCESSING

Jennifer Gillenwater
Joint work with Alex Kulesza and Ben Taskar

## Motivation & background on DPPs

# Motivation & background on DPPs Large-scale settings

# Motivation & background on DPPs Large-scale settings Structured summarization

Motivation & background on DPPs

Large-scale settings

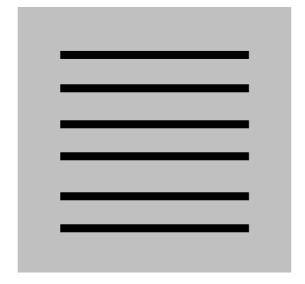
Structured summarization

Other potential NLP applications

# MOTIVATION & BACKGROUND



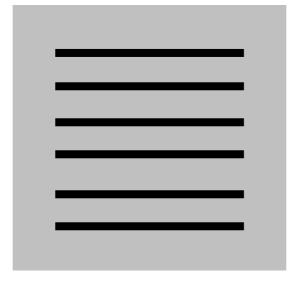






Quality:

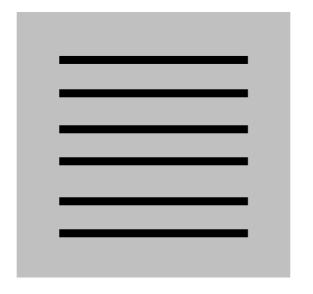
relevance to the topic





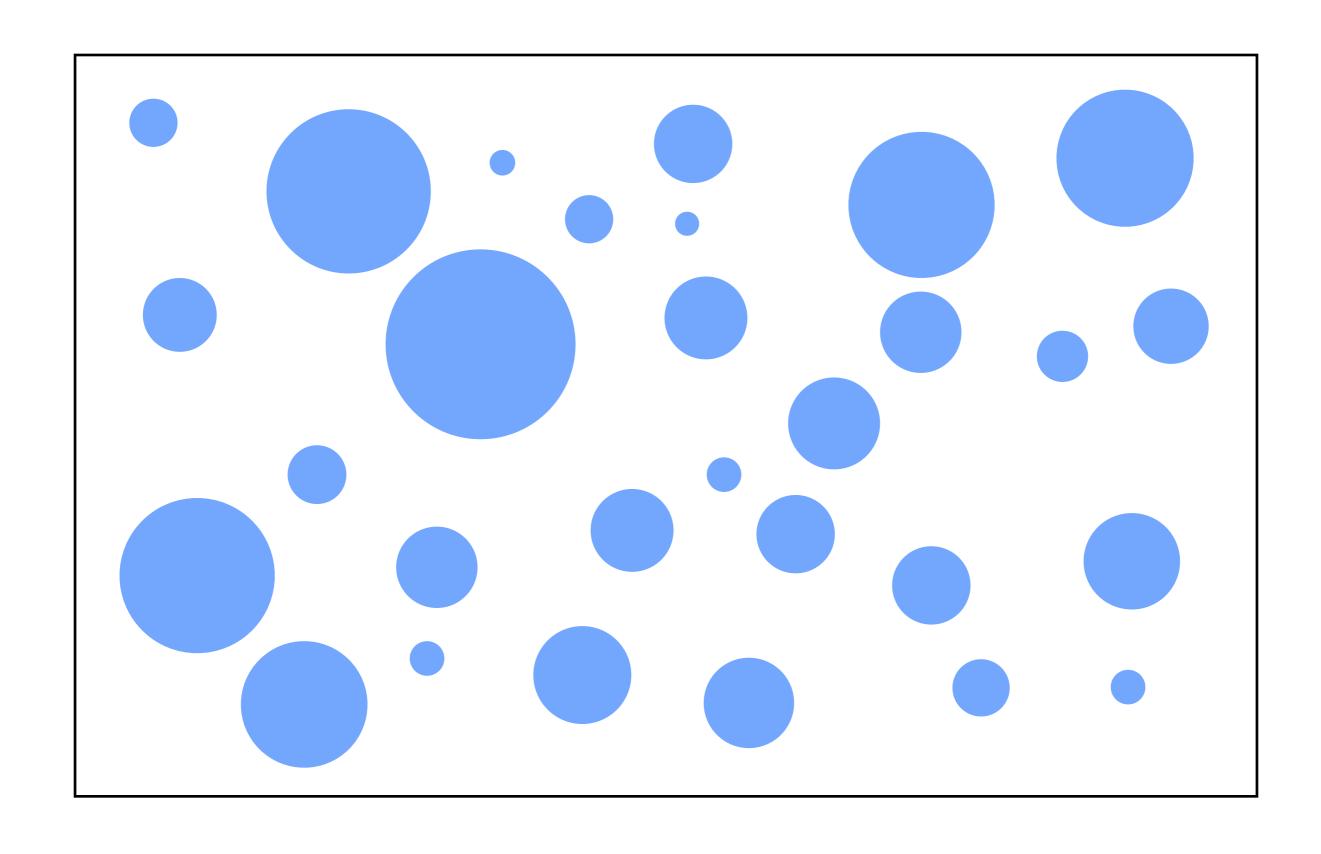
#### Quality:

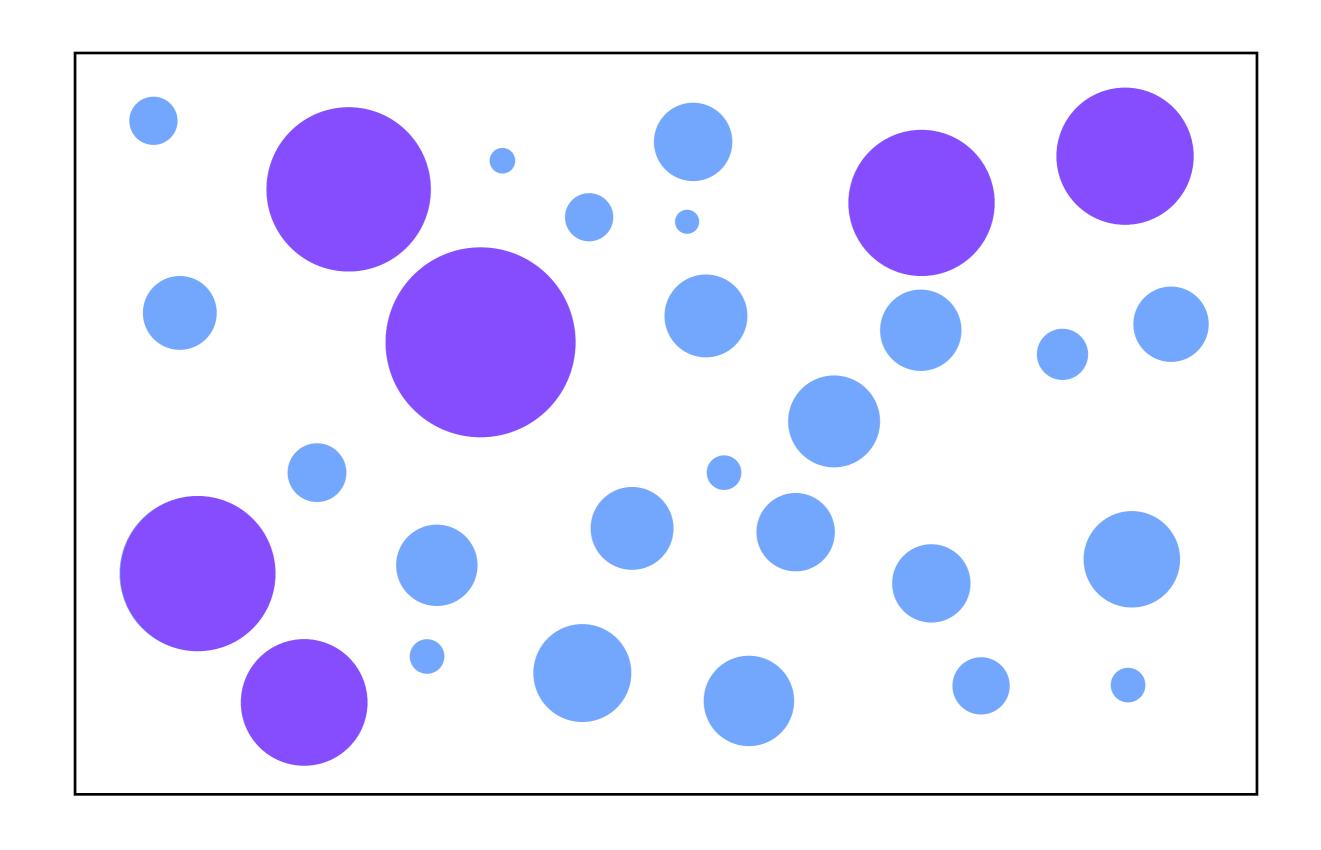
relevance to the topic

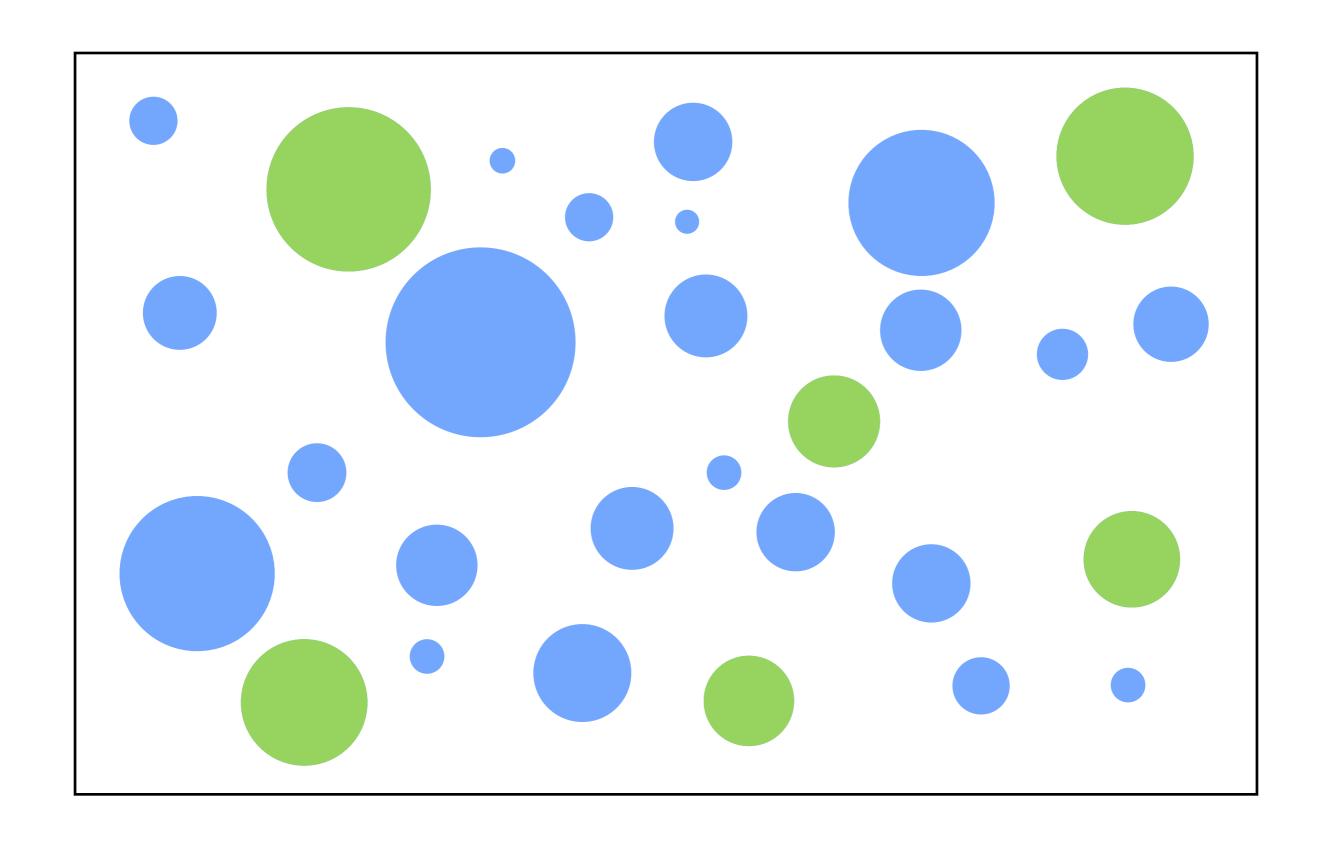


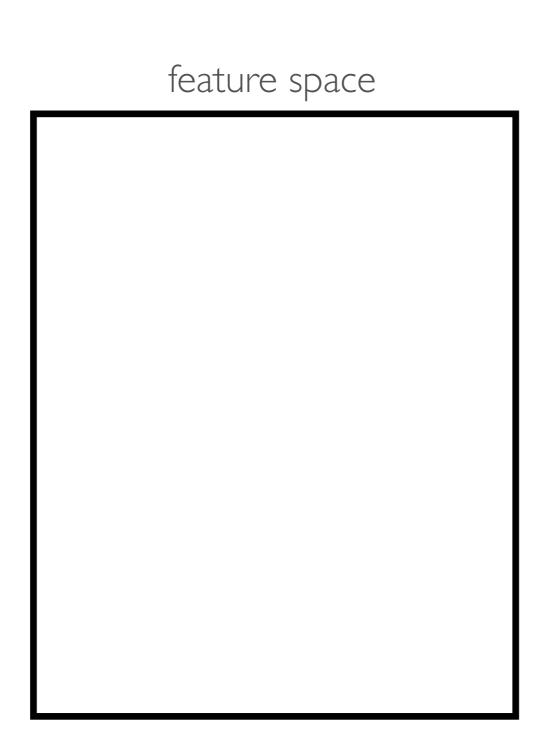
#### **Diversity**:

coverage of core ideas

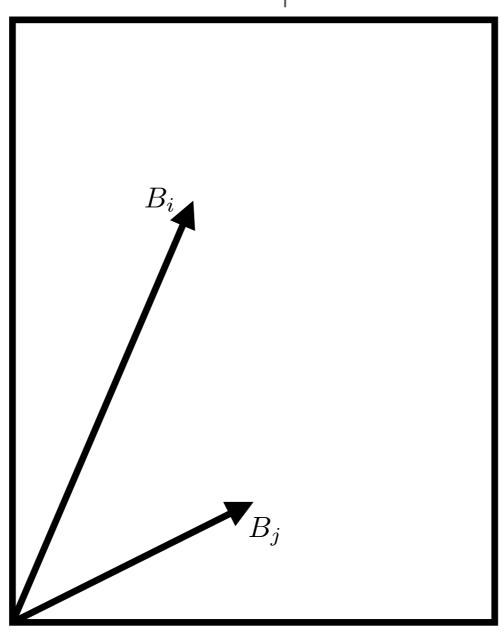




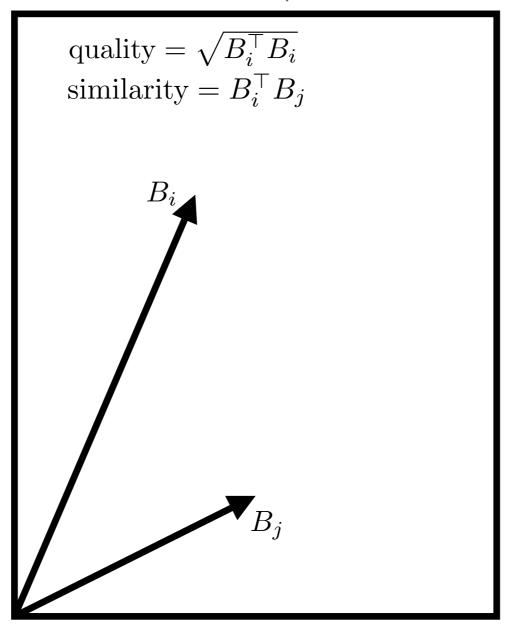




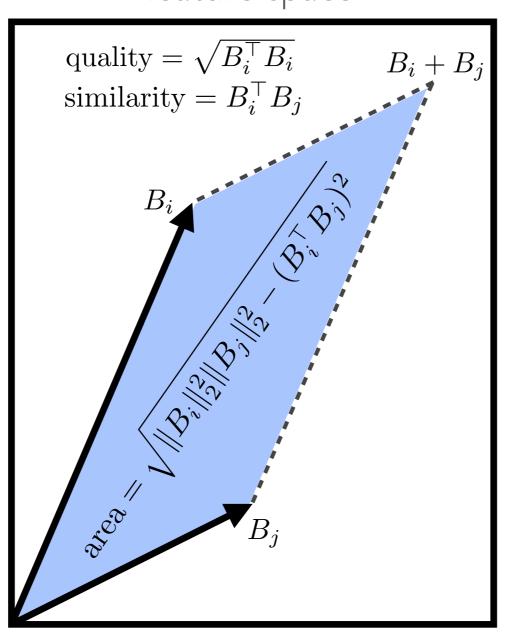
feature space



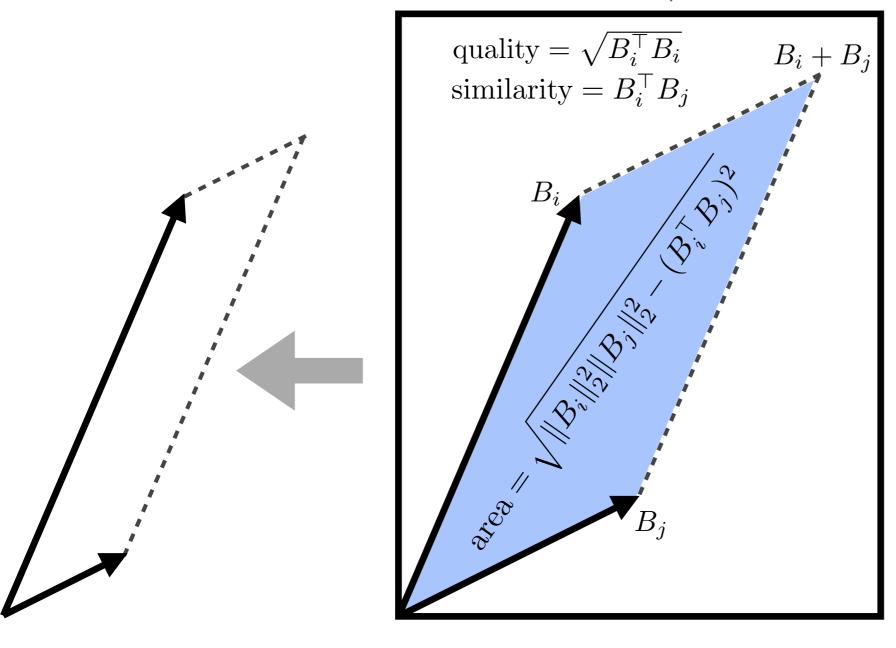
#### feature space

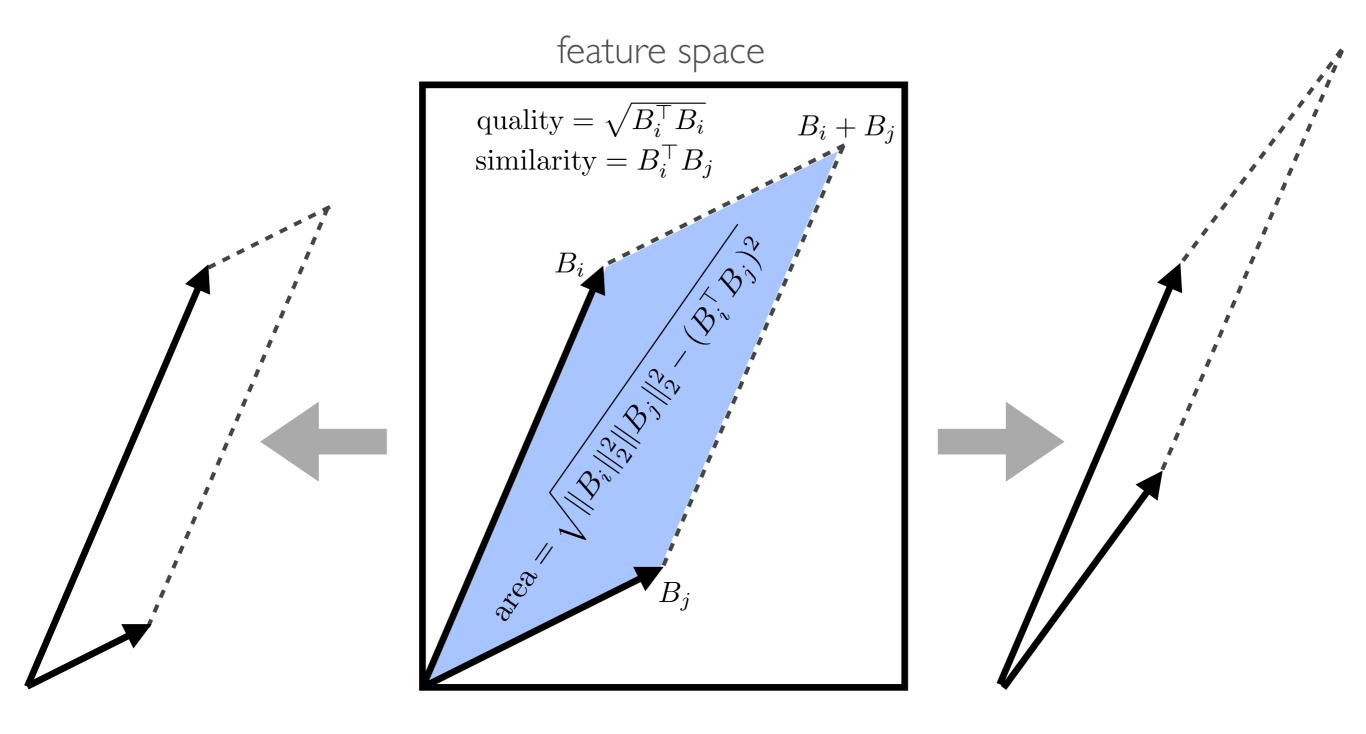


#### feature space



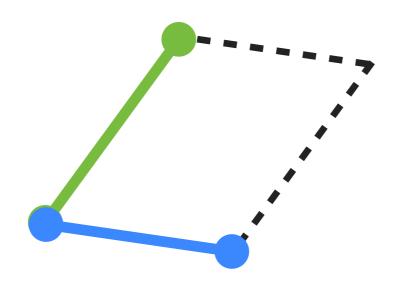




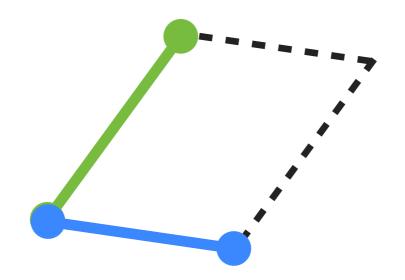


area = 
$$\sqrt{\|B_i\|_2^2 \|B_j\|_2^2 - (B_i^\top B_j)^2}$$

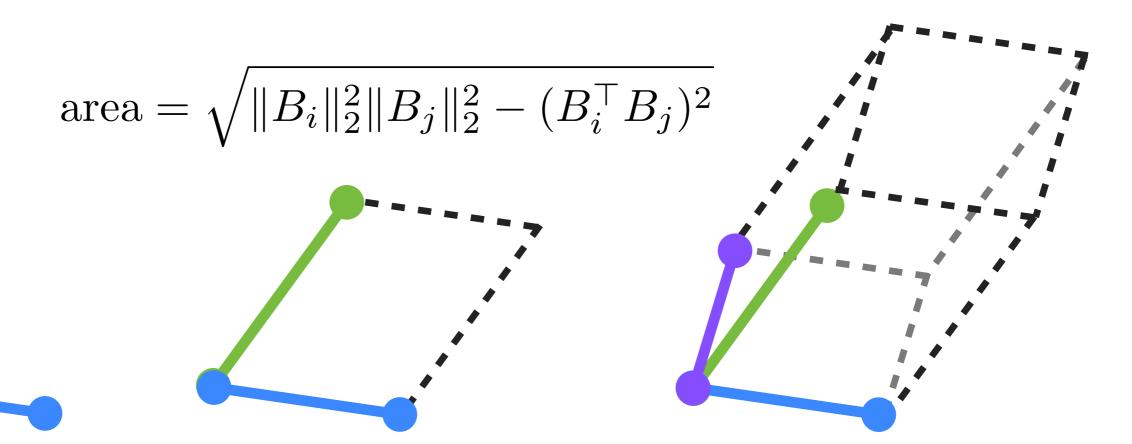
area = 
$$\sqrt{\|B_i\|_2^2 \|B_j\|_2^2 - (B_i^\top B_j)^2}$$



area = 
$$\sqrt{\|B_i\|_2^2 \|B_j\|_2^2 - (B_i^\top B_j)^2}$$

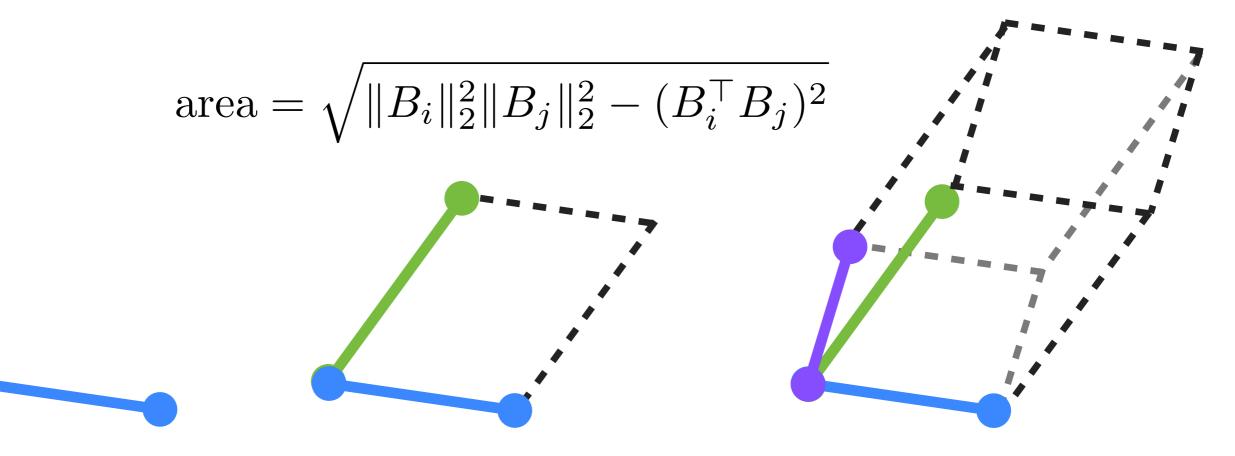


length =  $||B_i||_2$ 



length =  $||B_i||_2$ 

 $volume = base \times height$ 



length = 
$$||B_i||_2$$

 $volume = base \times height$ 

$$vol(B) = height \times base$$
  
=  $||B_1||_2 vol(proj_{\perp B_1}(B_{2:N}))$ 

#### AREA AS A DET

area = 
$$\sqrt{\|B_i\|_2^2 \|B_j\|_2^2 - (B_i^\top B_j)^2}$$

#### AREA AS A DET

area = 
$$\sqrt{\|B_i\|_2^2 \|B_j\|_2^2 - (B_i^\top B_j)^2}$$

$$= \det \left( \begin{array}{cc} ||B_i||_2^2 & B_i^{\top} B_j \\ B_i^{\top} B_j & ||B_j||_2^2 \end{array} \right)^{\frac{1}{2}}$$

#### AREA AS A DET

area = 
$$\sqrt{\|B_i\|_2^2 \|B_j\|_2^2 - (B_i^\top B_j)^2}$$

$$= \det \left( \begin{array}{cc} ||B_i||_2^2 & B_i^{\top} B_j \\ B_i^{\top} B_j & ||B_j||_2^2 \end{array} \right)^{\frac{1}{2}}$$

$$= \det \left( \begin{array}{c} -B_i - \\ -B_j - \end{array} \right)^{\frac{1}{2}}$$

#### VOLUME AS A DET

$$\operatorname{vol}(B_{\{i,j\}}) = \det\left(\begin{array}{c} B_i - \\ B_j - \\ \end{array}\right)^{\frac{1}{2}}$$

#### VOLUME AS A DET

$$\operatorname{vol}(B_{\{i,j\}}) = \det\left(\begin{array}{c} B_i - \\ B_j - \end{array}\right)^{\overline{2}}$$

$$\operatorname{vol}(B) = \det \left( \begin{array}{c} -B_1 - \\ \vdots \\ -B_N - \end{array} \right)^{\frac{1}{2}}$$

$$\operatorname{vol}(B)^2 = \det(B^{\mathsf{T}}B) = \det(L)$$

#### COMPLEX STATISTICS

#### COMPLEX STATISTICS



#### COMPLEX STATISTICS







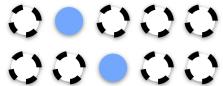
00000











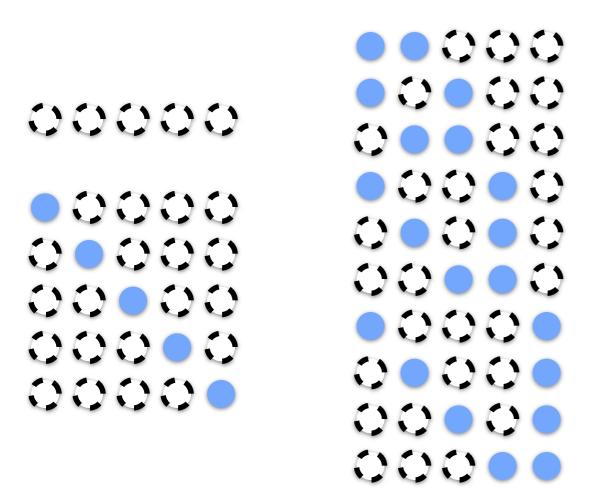


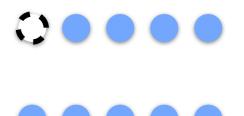




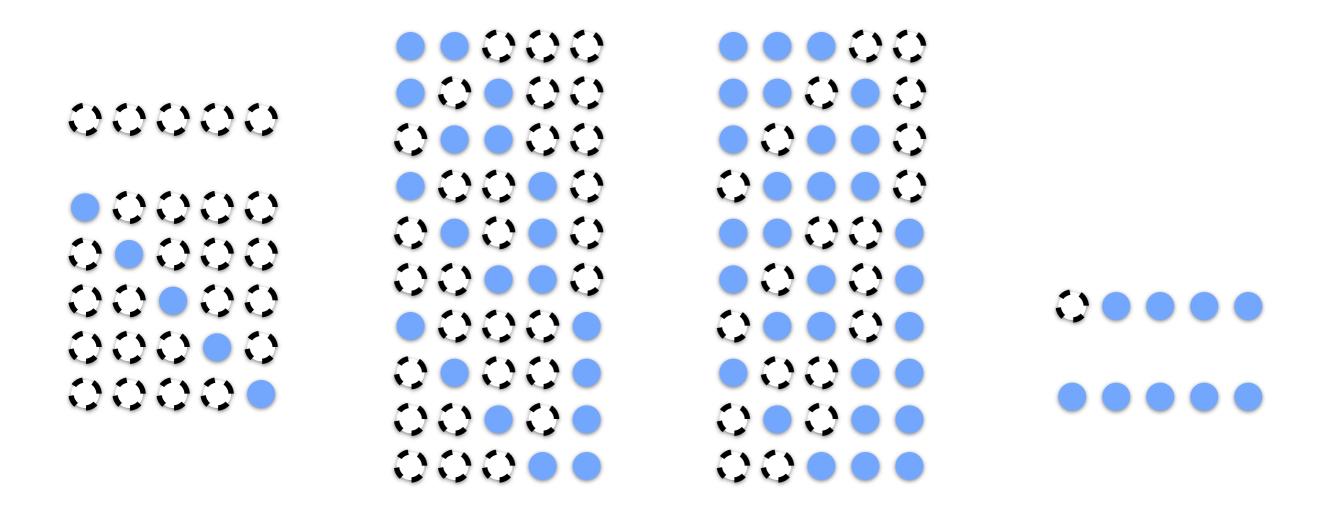




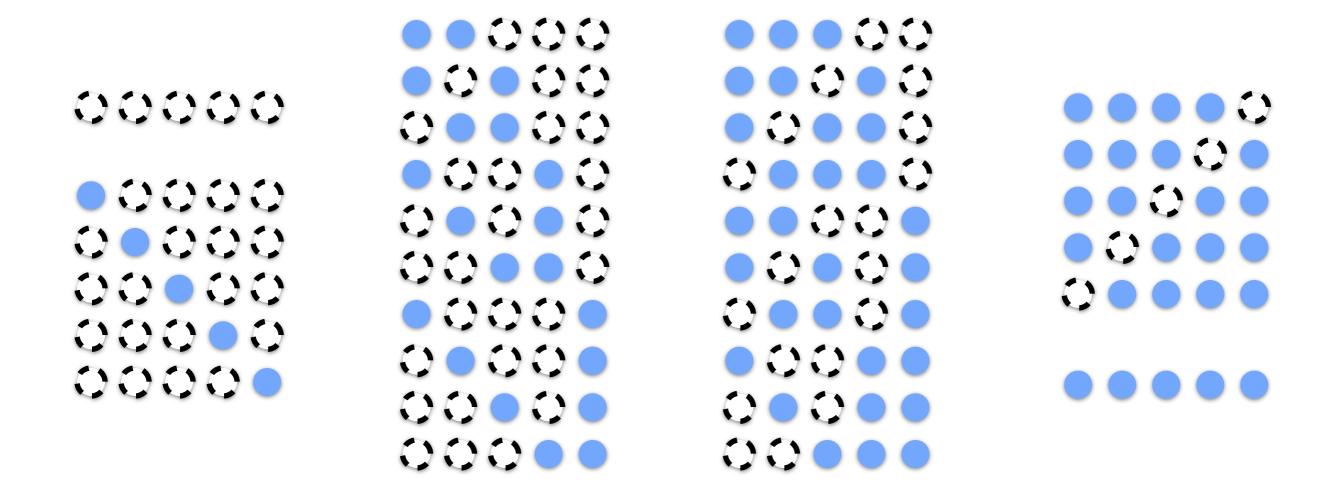


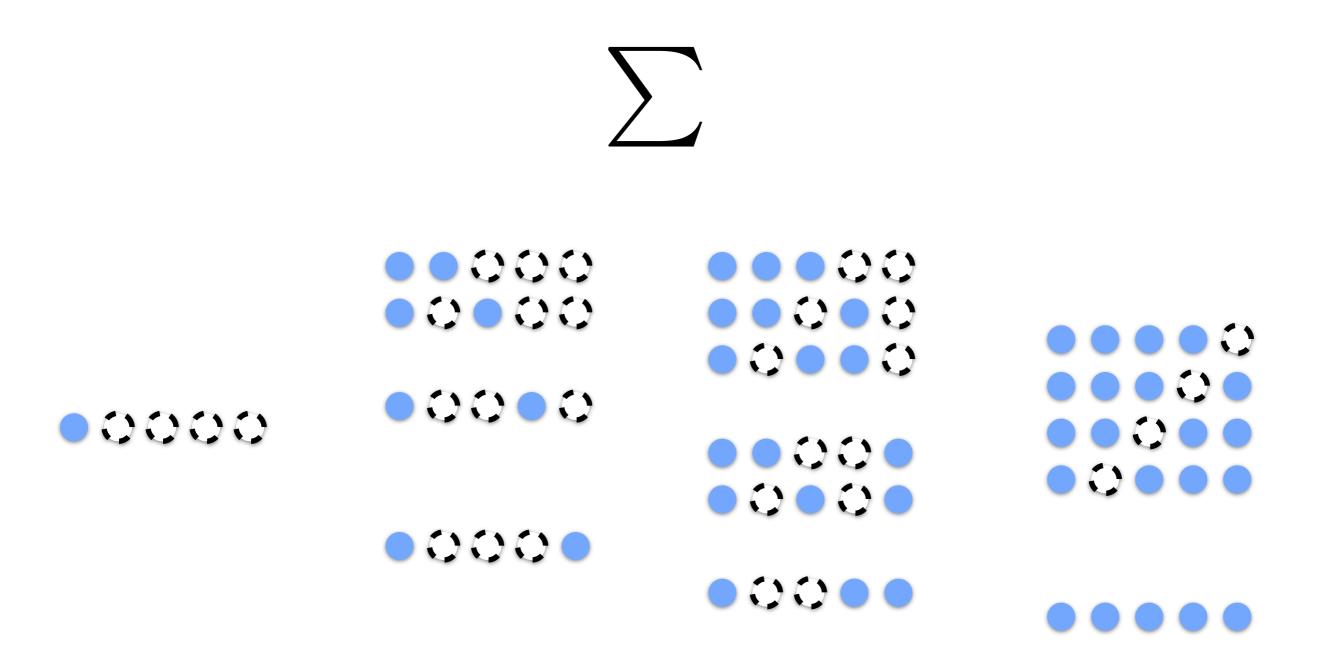




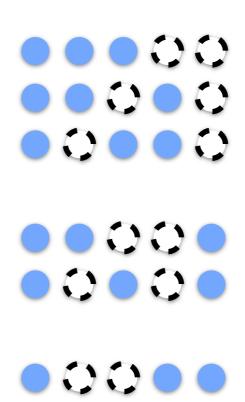




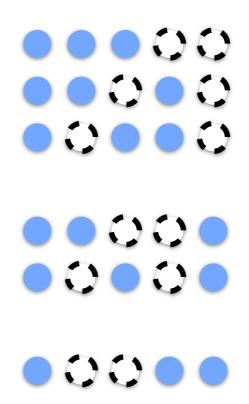


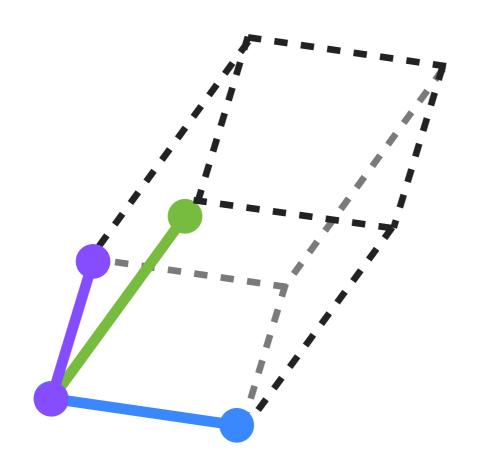




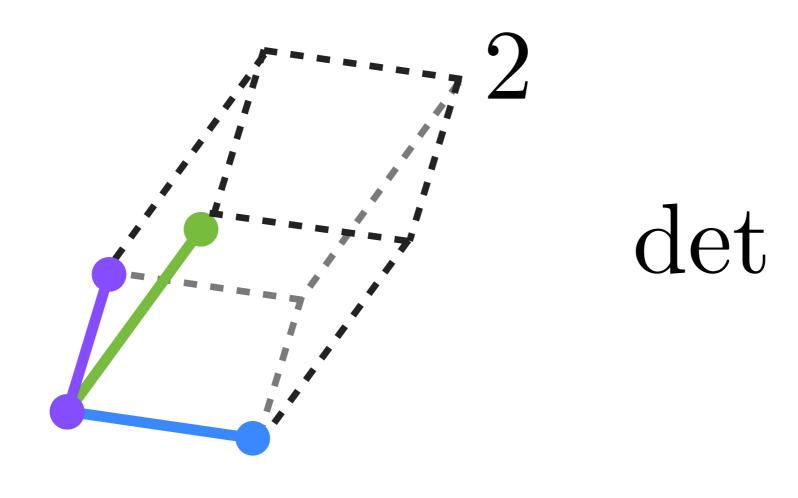


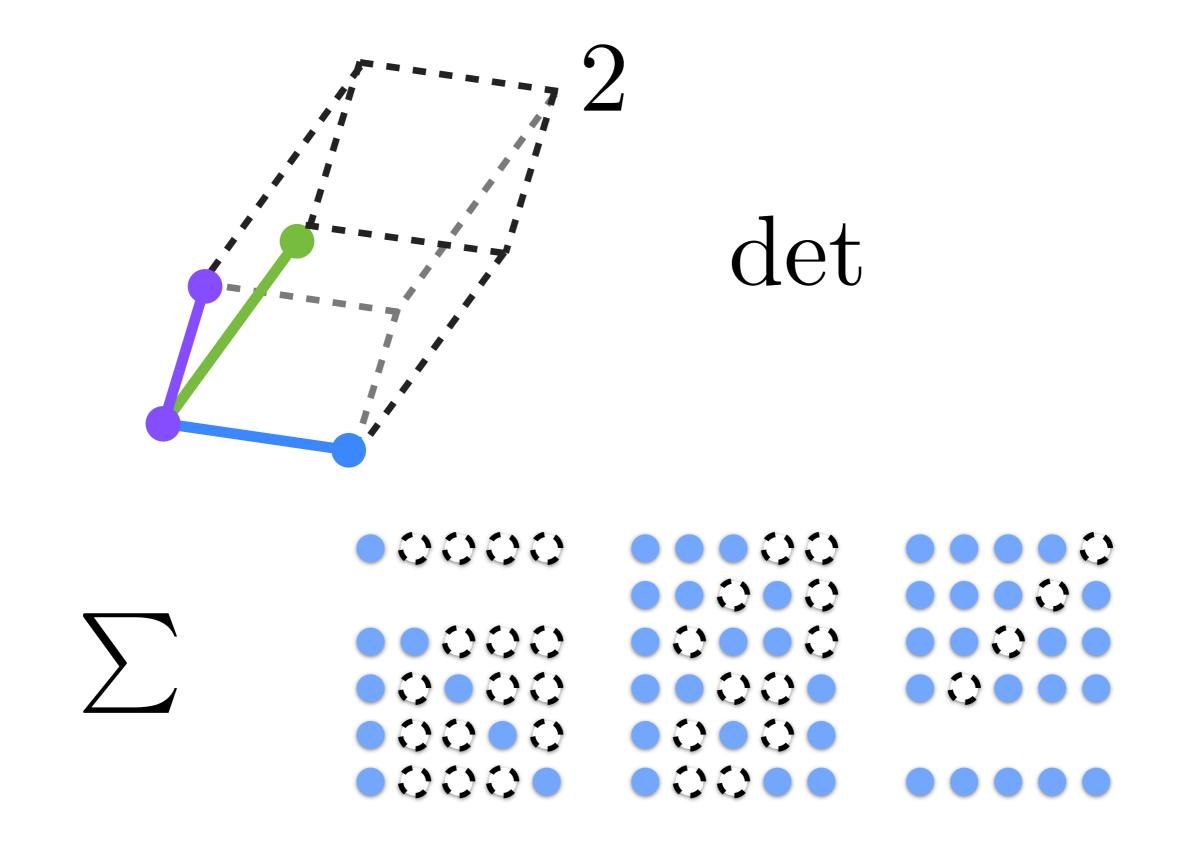
 $N \text{ items} \implies 2^N \text{ sets}$ 

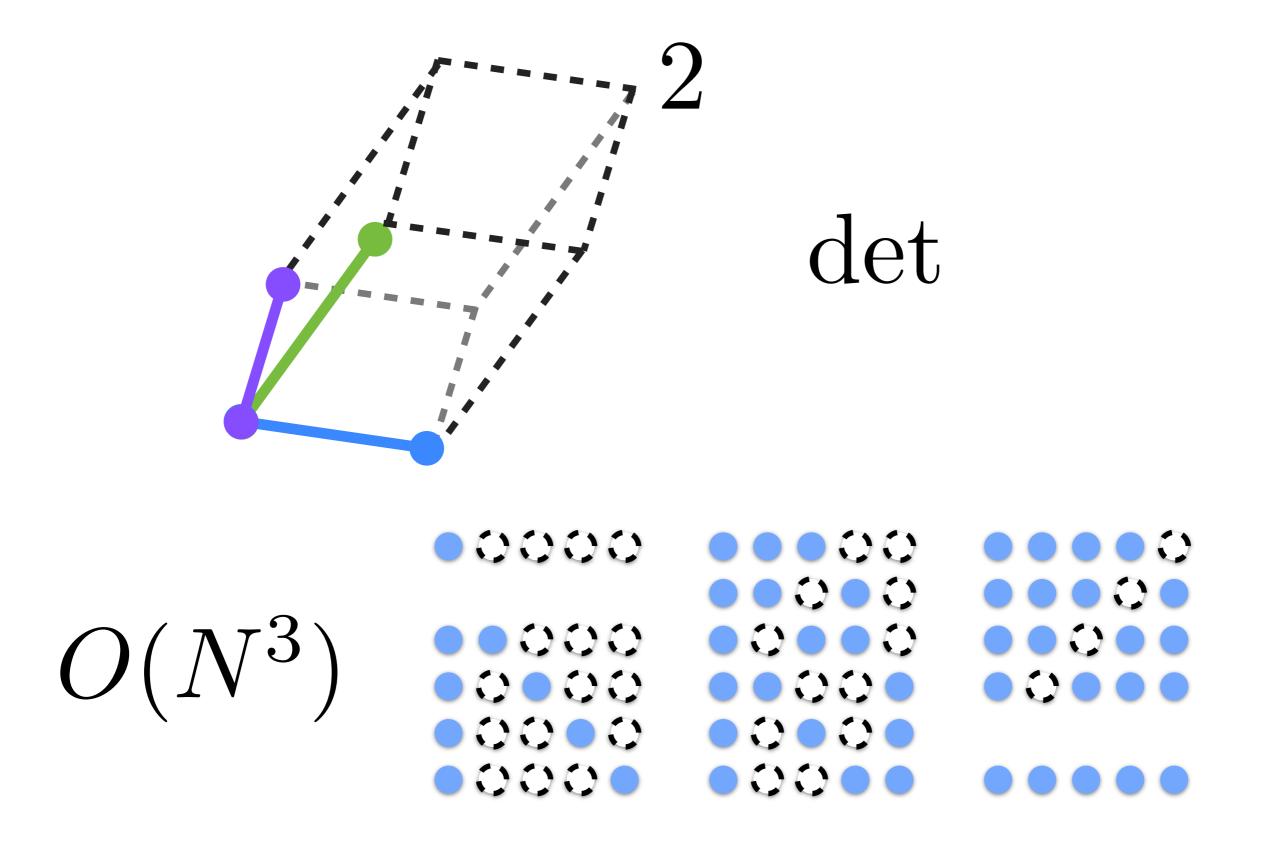




 $\det^{\frac{1}{2}}$ 

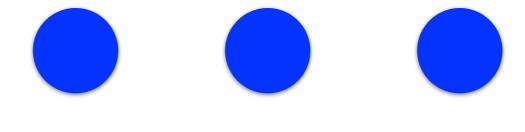






$$\mathcal{Y} = \{1, \dots, N\}$$

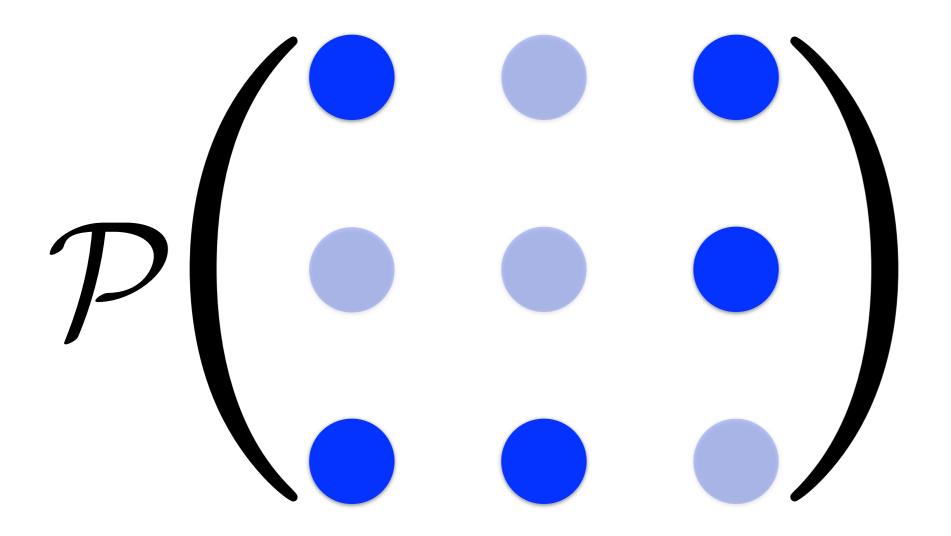
$$\mathcal{Y} = \{1, \dots, N\}$$



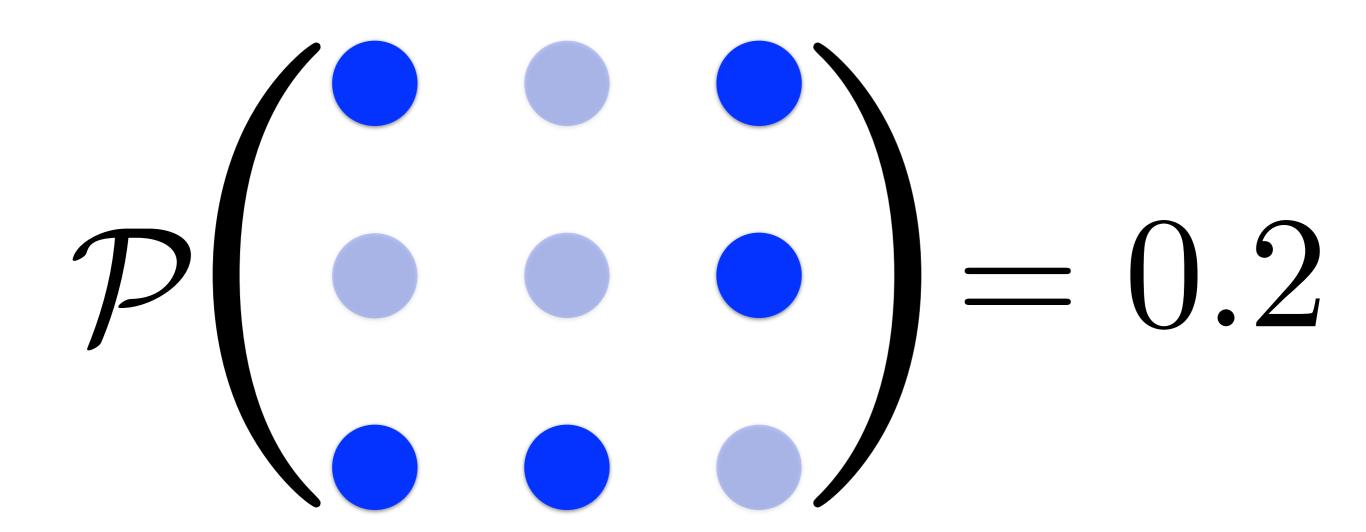




$$\mathcal{Y} = \{1, \dots, N\}$$



$$\mathcal{Y} = \{1, \dots, N\}$$



$$\mathcal{P}(\{2,3,5\}) \propto$$

$$\mathcal{P}(\{2,3,5\}) \propto$$

$$\mathcal{P}(\{2,3,5\}) \propto$$

$L_{11}$	$L_{12}$	$L_{13}$	$L_{14}$	$L_{15}$
$L_{21}$	$ L_{22} $	$L_{23}$	$L_{24}$	$L_{25}$
$L_{31}$	$L_{32}$	$L_{33}$	$L_{34}$	$L_{35}$
$L_{41}$	$ L_{42} $	$ L_{43} $	$L_{44}$	$ L_{45} $
$L_{51}$	$oxed{L_{52}}$	$L_{53}$	$L_{54}$	$oxed{L_{55}}$

$$\mathcal{P}(\{2,3,5\}) \propto$$

$$L_{22}$$
  $L_{23}$   $L_{25}$ 

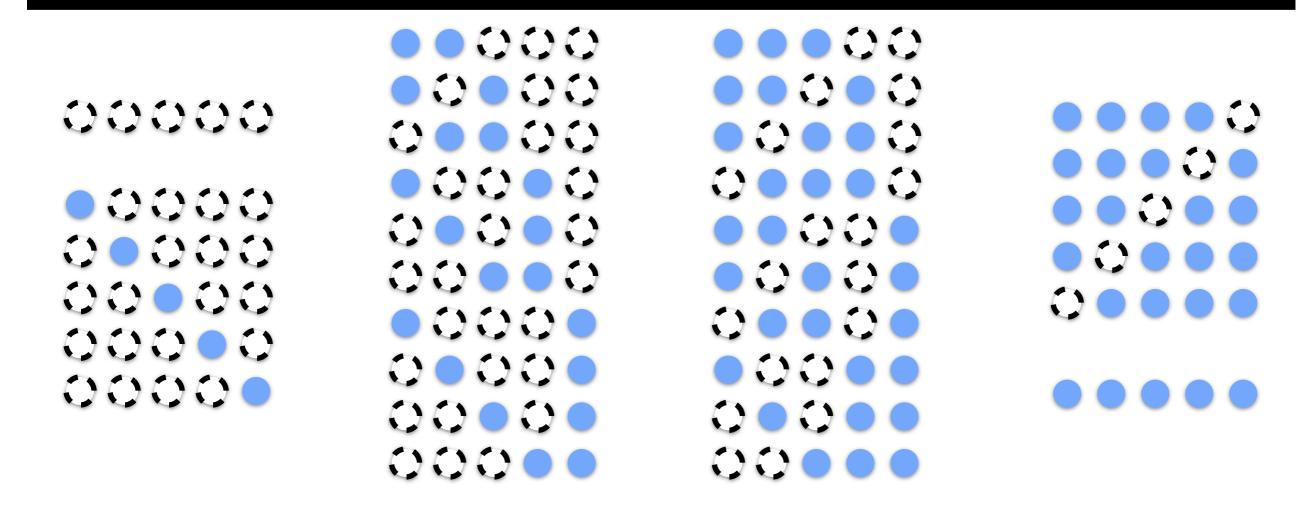
$$L_{32}$$
  $L_{33}$   $L_{35}$ 

$$L_{52}$$
  $L_{53}$   $L_{55}$ 

$$\mathcal{P}(\{2,3,5\}) \propto \det egin{pmatrix} L_{22} & L_{23} & L_{25} \ L_{32} & L_{33} & L_{35} \ L_{52} & L_{53} & L_{55} \ \end{pmatrix}$$

$$\mathcal{P}(\{2,3,5\}) = \det$$

 $egin{pmatrix} L_{22} & L_{23} & L_{25} \ L_{32} & L_{33} & L_{35} \ L_{52} & L_{53} & L_{55} \ \end{pmatrix}$ 



$$\mathcal{P}(\{2,3,5\}) = \det \begin{pmatrix} L_{22} & L_{23} & L_{25} \\ L_{32} & L_{33} & L_{35} \\ L_{52} & L_{53} & L_{55} \end{pmatrix}$$

$$\det(L+I)$$

Normalizing:  $\mathcal{P}_L(\mathbf{Y}=Y)$ 

Normalizing:  $\mathcal{P}_L(\mathbf{Y}=Y)$ 

Marginalizing:  $\mathcal{P}(Y \subseteq \mathbf{Y})$ 

Normalizing: 
$$\mathcal{P}_L(\mathbf{Y}=Y)$$

Marginalizing: 
$$\mathcal{P}(Y \subseteq \mathbf{Y})$$

Conditioning: 
$$\mathcal{P}_L(\mathbf{Y} = B \mid A \subseteq \mathbf{Y})$$
  
 $\mathcal{P}_L(\mathbf{Y} = B \mid A \cap \mathbf{Y} = \emptyset)$ 

Normalizing:  $\mathcal{P}_L(\mathbf{Y}=Y)$ 

Marginalizing:  $\mathcal{P}(Y \subseteq \mathbf{Y})$ 

Conditioning:  $\mathcal{P}_L(\mathbf{Y} = B \mid A \subseteq \mathbf{Y})$  $\mathcal{P}_L(\mathbf{Y} = B \mid A \cap \mathbf{Y} = \emptyset)$ 

Sampling:  $\mathbf{Y} \sim \mathcal{P}_L$ 

Normalizing:  $\mathcal{P}_L(\mathbf{Y}=Y)$ 

Marginalizing:  $\mathcal{P}(Y \subseteq \mathbf{Y})$ 

Conditioning:  $\mathcal{P}_L(\mathbf{Y} = B \mid A \subseteq \mathbf{Y})$  $\mathcal{P}_L(\mathbf{Y} = B \mid A \cap \mathbf{Y} = \emptyset)$ 

Sampling:  $\mathbf{Y} \sim \mathcal{P}_L$ 

 $O(N^3)$ 

# LARGE-SCALE SETTINGS

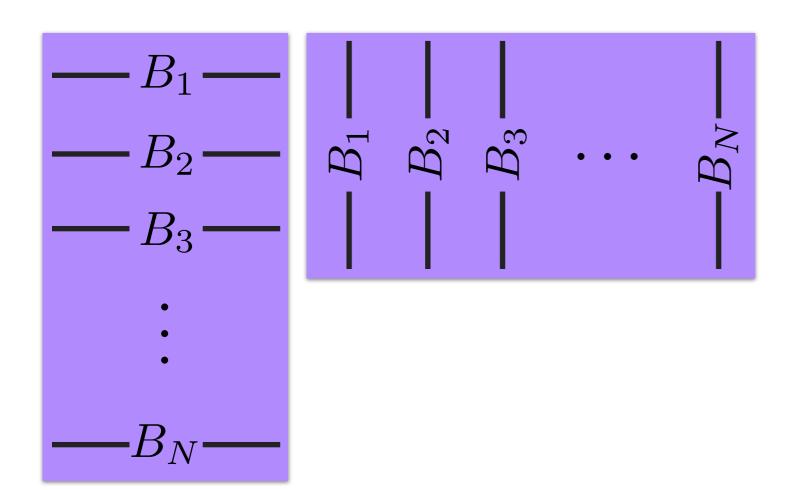
# DUAL KERNEL

KULESZA AND TASKAR (NIPS 2010)

## DUAL KERNEL

KULESZA AND TASKAR (NIPS 2010)

L



## DUAL KERNEL

KULESZA AND TASKAR (NIPS 2010)

L

### DUAL KERNEL

KULESZA AND TASKAR (NIPS 2010)

C

### DUAL KERNEL

KULESZA AND TASKAR (NIPS 2010)

C

#### DUAL KERNEL

KULESZA AND TASKAR (NIPS 2010)

C

$$L = V \Lambda V^{\top}$$

$$C = \hat{V} \Lambda \hat{V}^{\top}$$

$$L = V\Lambda V^{\top} \left\langle V = B^{\top} \hat{V} \Lambda^{-\frac{1}{2}} \right\rangle C = \hat{V} \Lambda \hat{V}^{\top}$$

$$L = V\Lambda V^{\top} \left\langle V = B^{\top} \hat{V} \Lambda^{-\frac{1}{2}} \right\rangle C = \hat{V} \Lambda \hat{V}^{\top}$$

Normalizing 
$$\sum_{Y} \det(L_Y)$$
  $O(D^3)$ 

$$L = V\Lambda V^{\top} \left\langle \begin{array}{c} V = B^{\top} \hat{V} \Lambda^{-\frac{1}{2}} \\ \end{array} \right\rangle C = \hat{V} \Lambda \hat{V}^{\top}$$

Normalizing 
$$\sum_{Y} \det(L_Y)$$
  $O(D^3)$ 

Marginalizing & Conditioning

$$O(D^3 + D^2k^2)$$

$$L = V\Lambda V^{\top} \left\langle V = B^{\top} \hat{V} \Lambda^{-\frac{1}{2}} \right\rangle C = \hat{V} \Lambda \hat{V}^{\top}$$

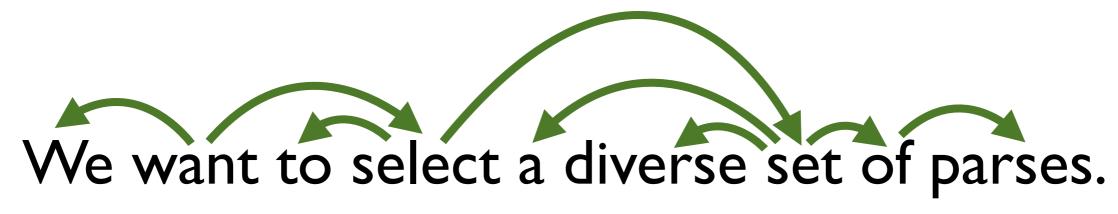
Normalizing 
$$\sum_{Y} \det(L_Y)$$
  $O(D^3)$ 

1arginalizing & Conditioning 
$$O(D^3 + D^2 k^2)$$

Sampling 
$$\mathbf{Y} \sim \mathcal{P}_L$$
  $O(ND^2k)$ 

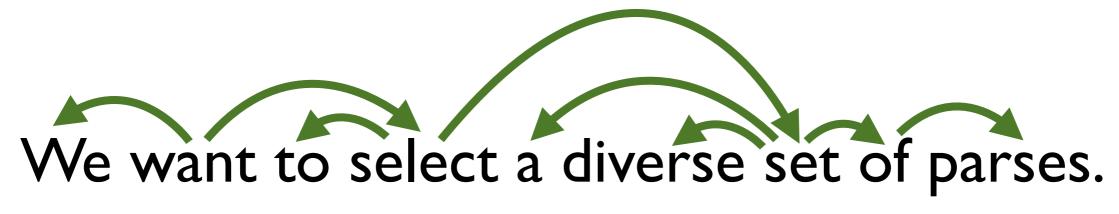
# EXPONENTIAL N

#### EXPONENTIAL N

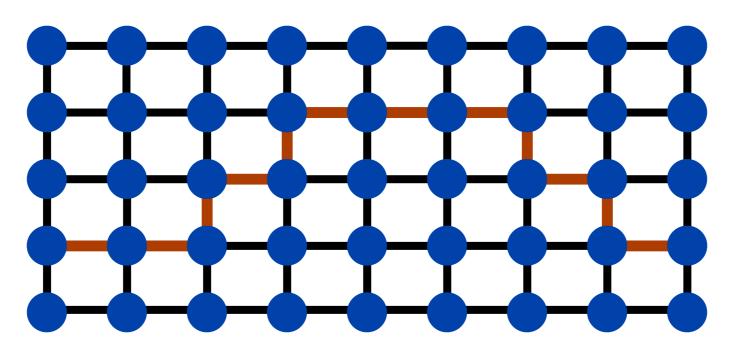


 $N = O(\{\text{sentence length}\}^{\{\text{sentence length}\}})$ 

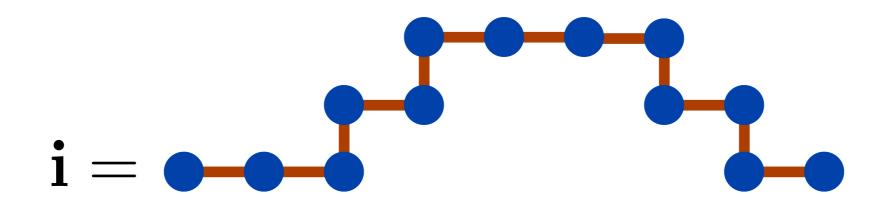
#### EXPONENTIAL N

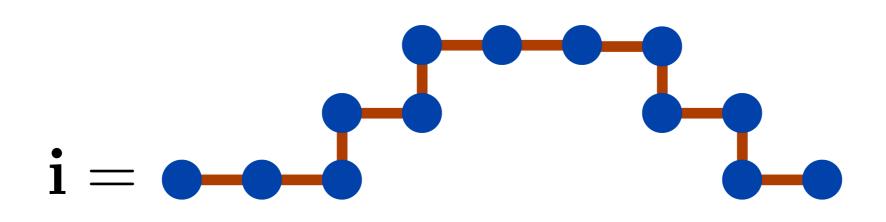


 $N = O(\{\text{sentence length}\}^{\{\text{sentence length}\}})$ 



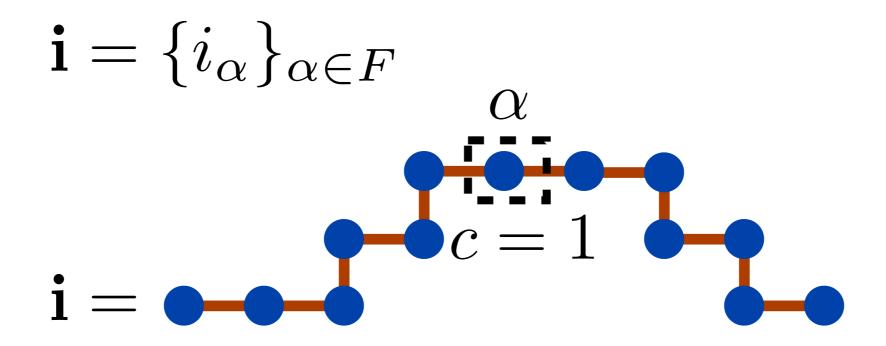
 $N = O(\{\text{node degree}\}^{\{\text{path length}}\})$ 





$$B_{\mathbf{i}} = q(\mathbf{i})\phi(\mathbf{i})$$

$$\uparrow \qquad \uparrow$$
quality similarity



$$B_{\mathbf{i}} = q(\mathbf{i})\phi(\mathbf{i})$$

$$\mathbf{\uparrow} \qquad \mathbf{\uparrow}$$
quality similarity

$$\mathbf{i} = \{i_{\alpha}\}_{\alpha \in F}$$

$$\mathbf{i} = \mathbf{i}_{\alpha}\}_{\alpha \in F}$$

$$\mathbf{i} = \mathbf{i}_{\alpha}$$

$$B_{\mathbf{i}} = q(\mathbf{i})\phi(\mathbf{i})$$

$$\mathbf{\uparrow} \qquad \mathbf{\uparrow}$$
quality similarity

$$\mathbf{i} = \{i_{\alpha}\}_{\alpha \in F}$$

$$\mathbf{i} = \{i_{\alpha}\}_{\alpha \in F}$$

$$\mathbf{i} = \{i_{\alpha}\}_{\alpha \in F}$$

$$B_{\mathbf{i}} = \left| \prod_{\alpha \in F} q(i_{\alpha}) \right| \phi(\mathbf{i})$$

$$\mathbf{i} = \{i_{\alpha}\}_{\alpha \in F}$$

$$\alpha$$

$$\mathbf{i} = \mathbf{i}$$

$$\mathbf{i} = \mathbf{i}$$

$$\mathbf{j}$$

$$\mathbf{i} = \mathbf{i}$$

$$\mathbf{j}$$

$$\mathbf{i} = \mathbf{i}$$

$$\mathbf{j}$$

$$\mathbf{j}$$

$$\mathbf{i} = \mathbf{j}$$

$$\mathbf{j}$$

$$\mathbf{i} = \{i_{\alpha}\}_{\alpha \in F}$$

$$\mathbf{i} = \mathbf{i}$$

$$\mathbf{i$$

$$M = \frac{\alpha}{c = 2} \qquad R = \boxed{\phantom{a}}$$

$$B_{\mathbf{i}} = \left[ \prod_{\alpha \in F} q(i_{\alpha}) \right] \left[ \sum_{\alpha \in F} \phi(i_{\alpha}) \right]$$

$$\mathbf{Y} \sim \mathcal{P}_L \quad O(D^2k^3 + Dk^2M^cR)$$

$$M = \frac{\alpha}{c = 2} \qquad R = \boxed{\phantom{a}}$$

$$B_{\mathbf{i}} = \left[ \prod_{\alpha \in F} q(i_{\alpha}) \right] \left[ \sum_{\alpha \in F} \phi(i_{\alpha}) \right]$$

$$\mathbf{Y} \sim \mathcal{P}_L \quad O(D^2 k^3 + Dk^2 M^c R)$$
 $M^c R = 4^2 * 12 = 192 \ll N = 4^{12} = 16,777,216$ 

# LARGE FEATURE SETS?

#### LARGE FEATURE SETS?

N = # of items

Large Exponential

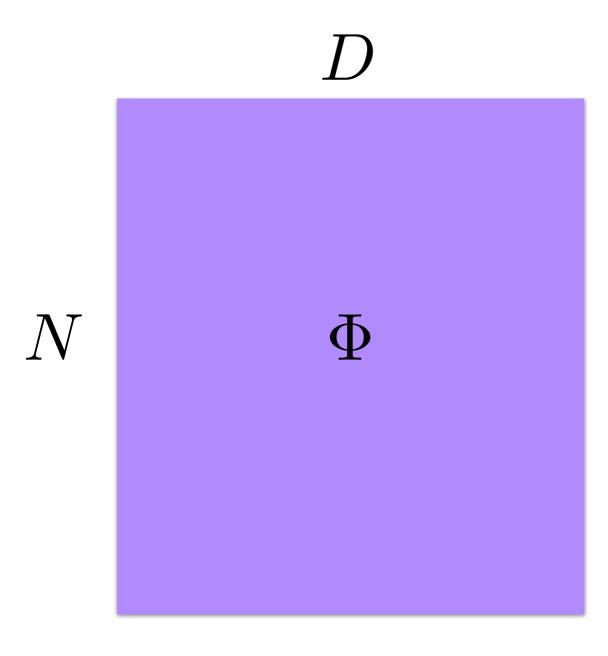
Small dual dual structure

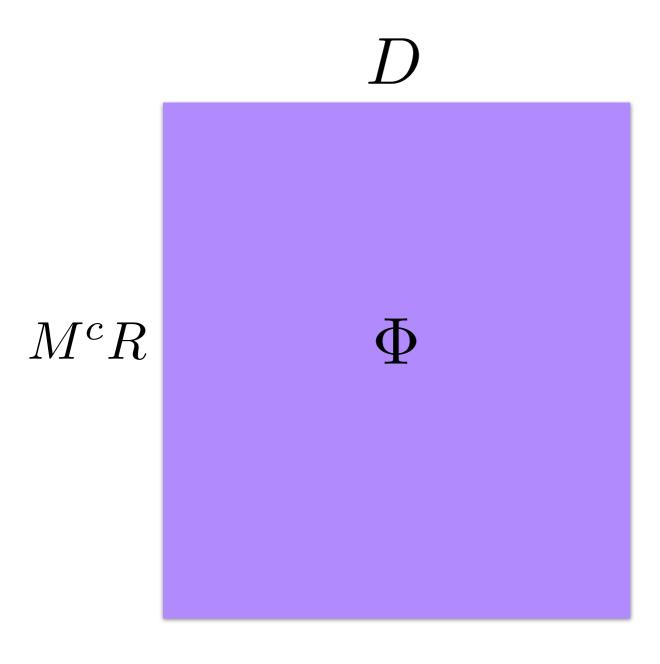
#

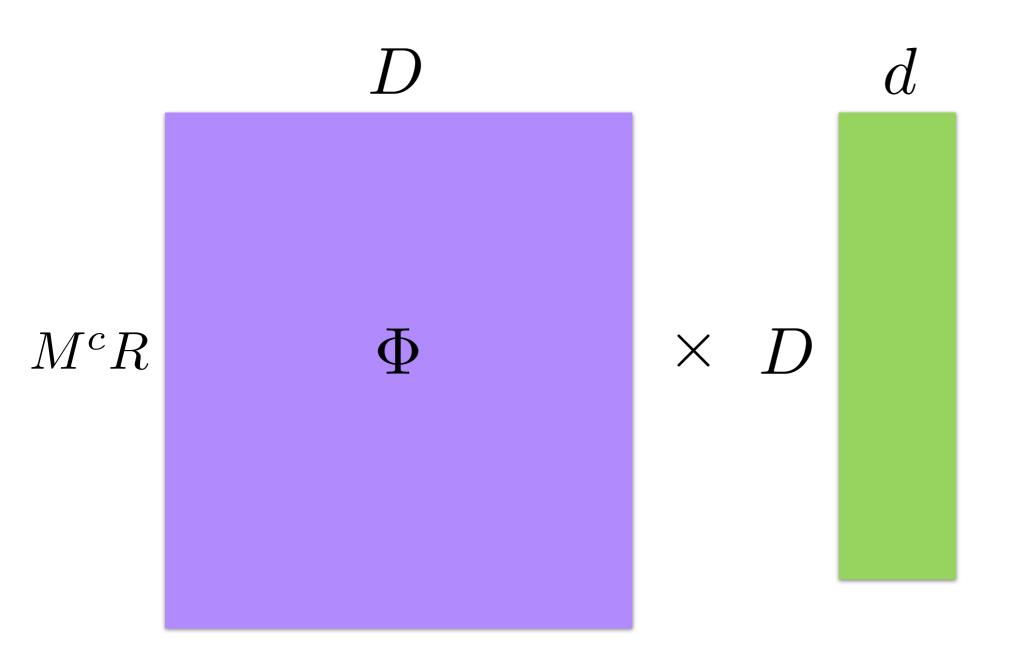
### LARGE FEATURE SETS?

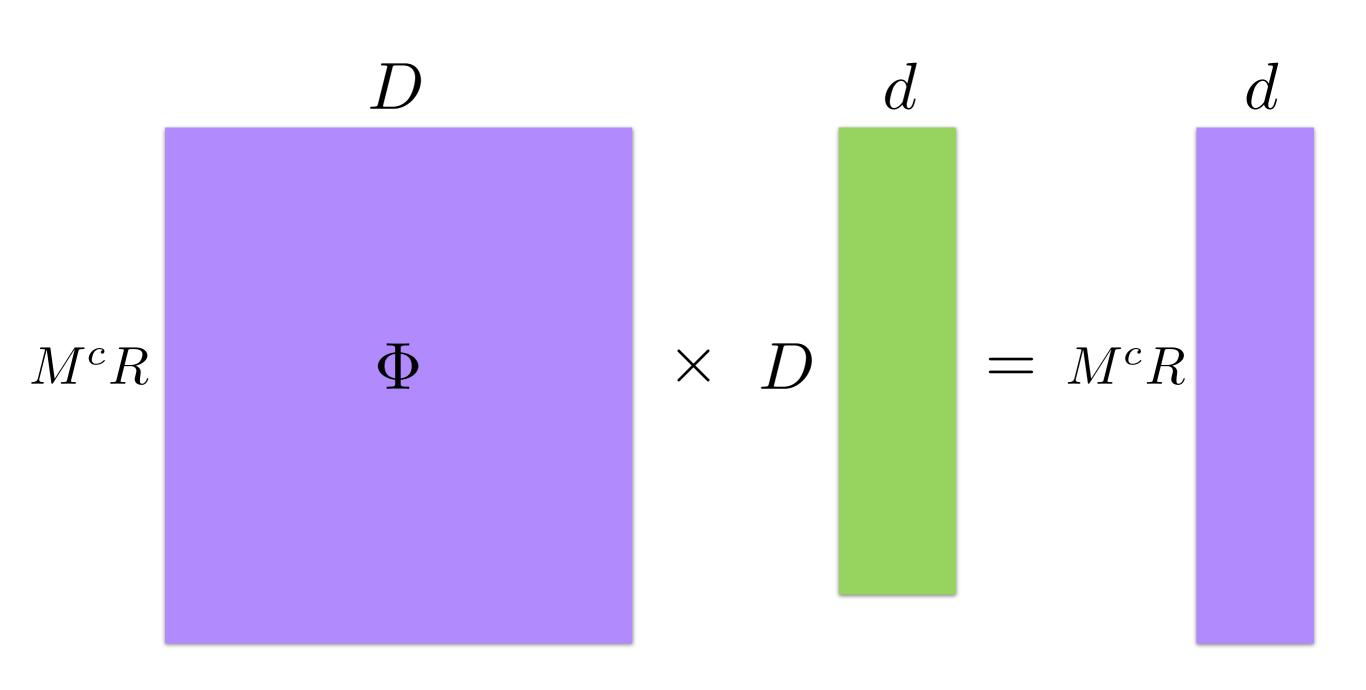
N = # of items

.es		Large	Exponential
of features	Small	dual	dual + structure
D = #  of	Large	?	









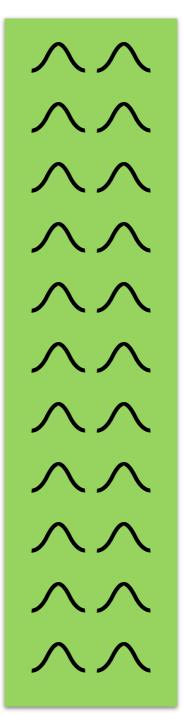
GILLENWATER, KULESZA, AND TASKAR (EMNLP 2012)

d

D

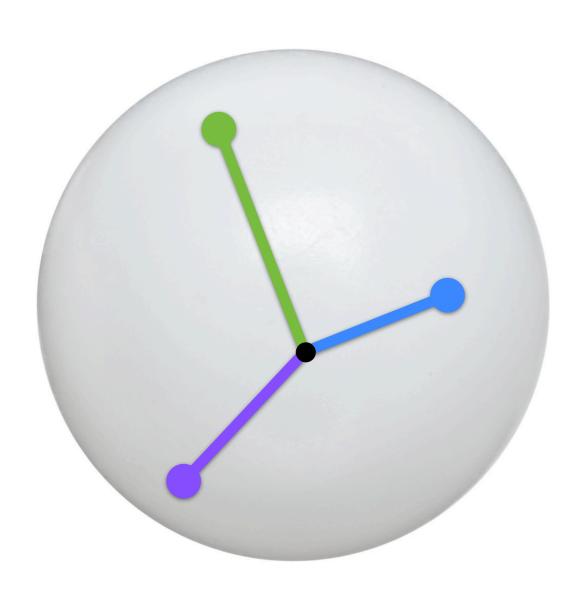
GILLENWATER, KULESZA, AND TASKAR (EMNLP 2012)

d

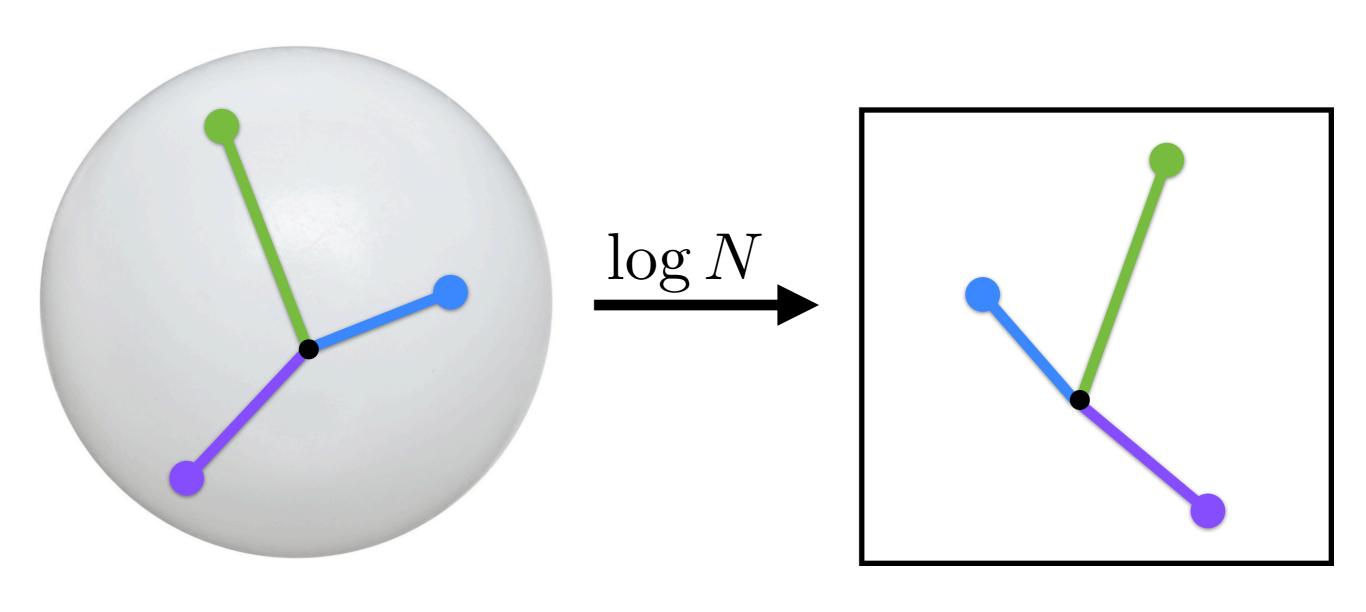


 $\int$ 

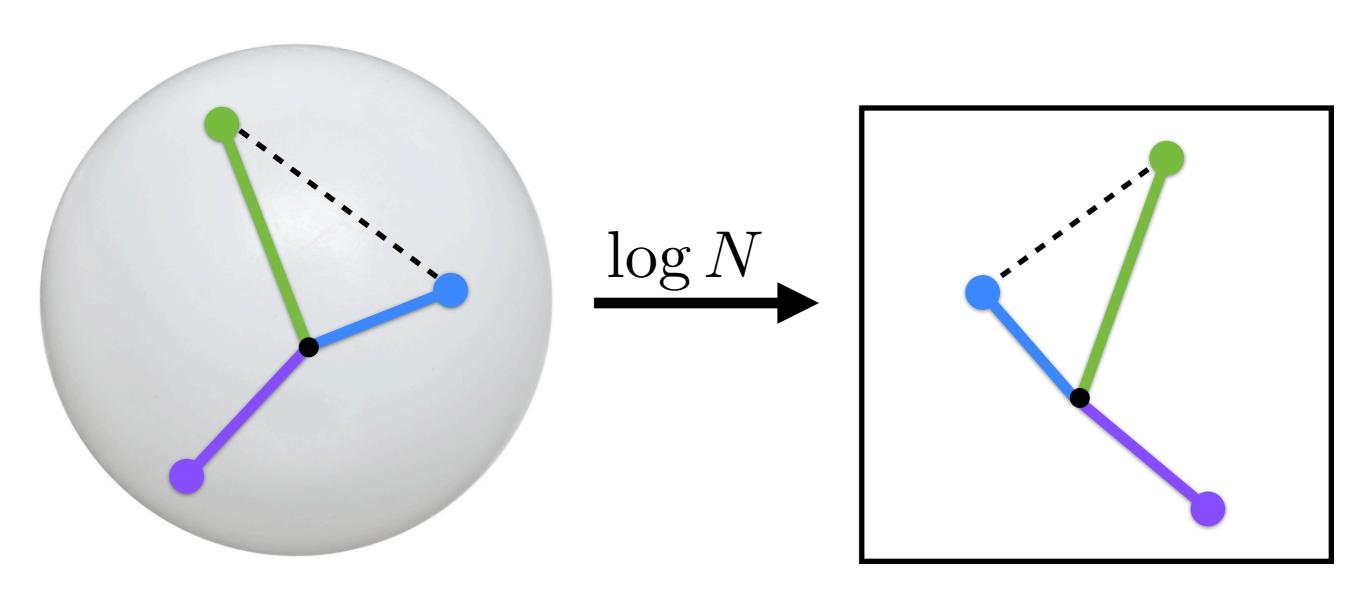
JOHNSON AND LINDENSTRAUSS (1984)



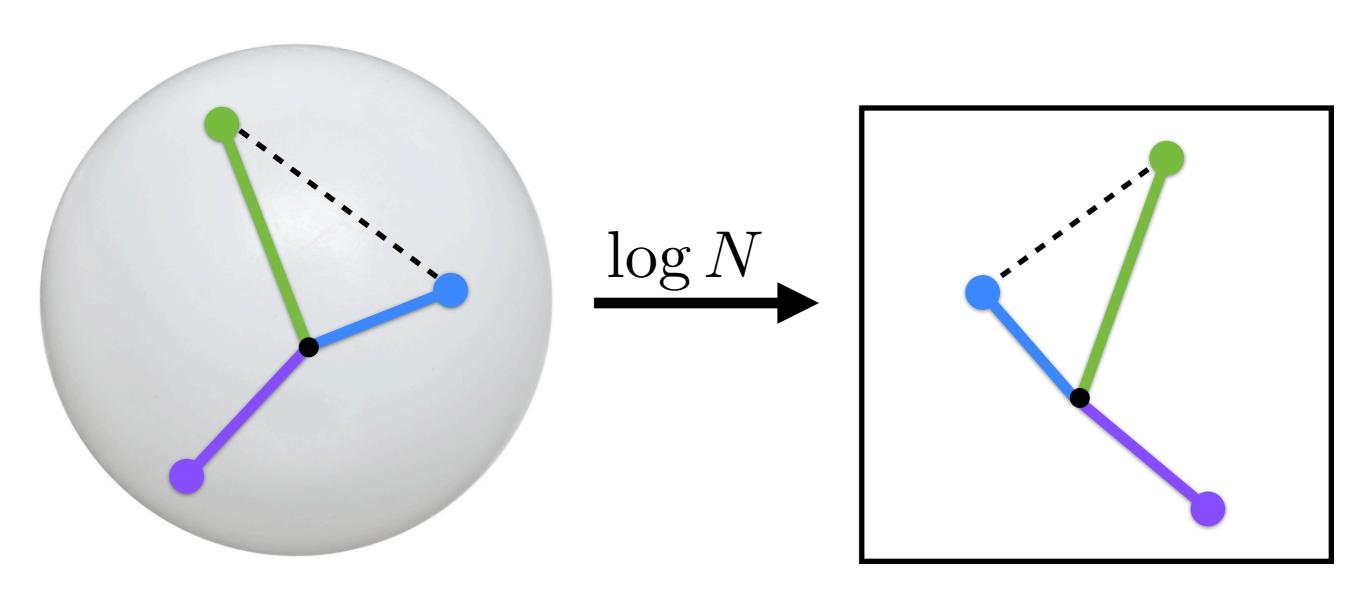
JOHNSON AND LINDENSTRAUSS (1984)



JOHNSON AND LINDENSTRAUSS (1984)

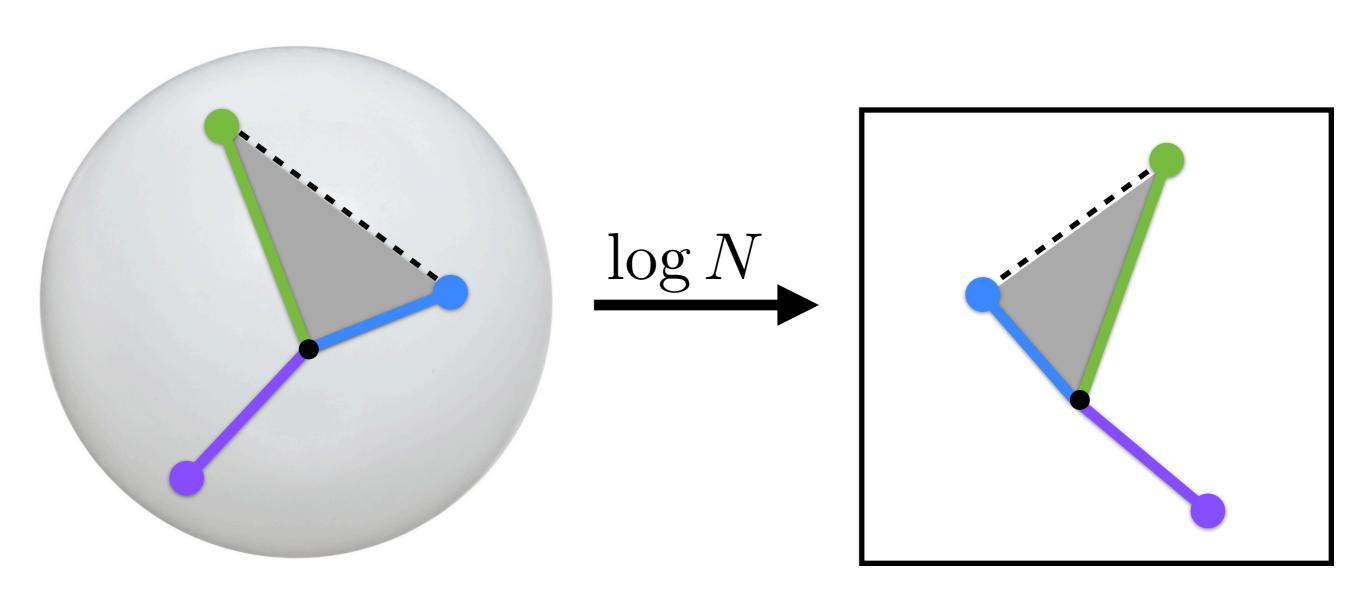


MAGEN AND ZOUZIAS (2008)



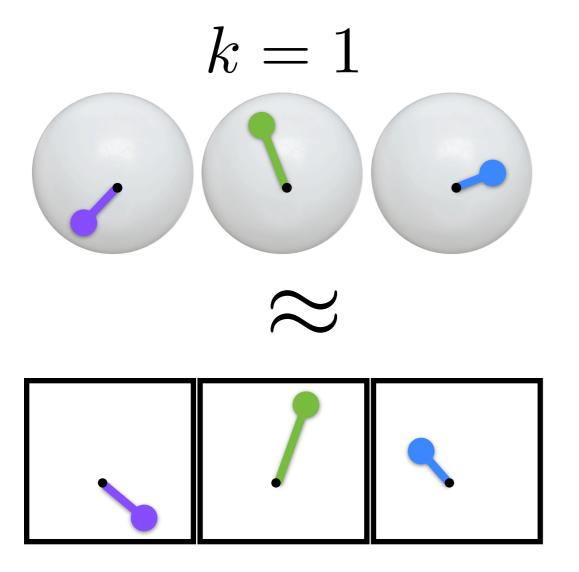
## VOLUME PRESERVATION

MAGEN AND ZOUZIAS (2008)

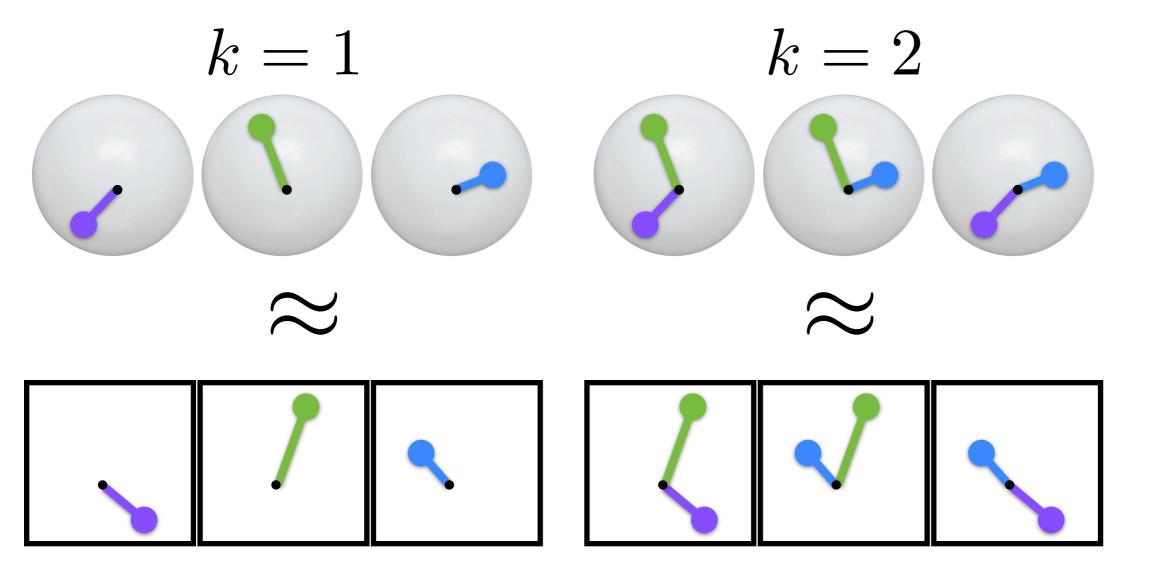


$$vol^2 = det$$

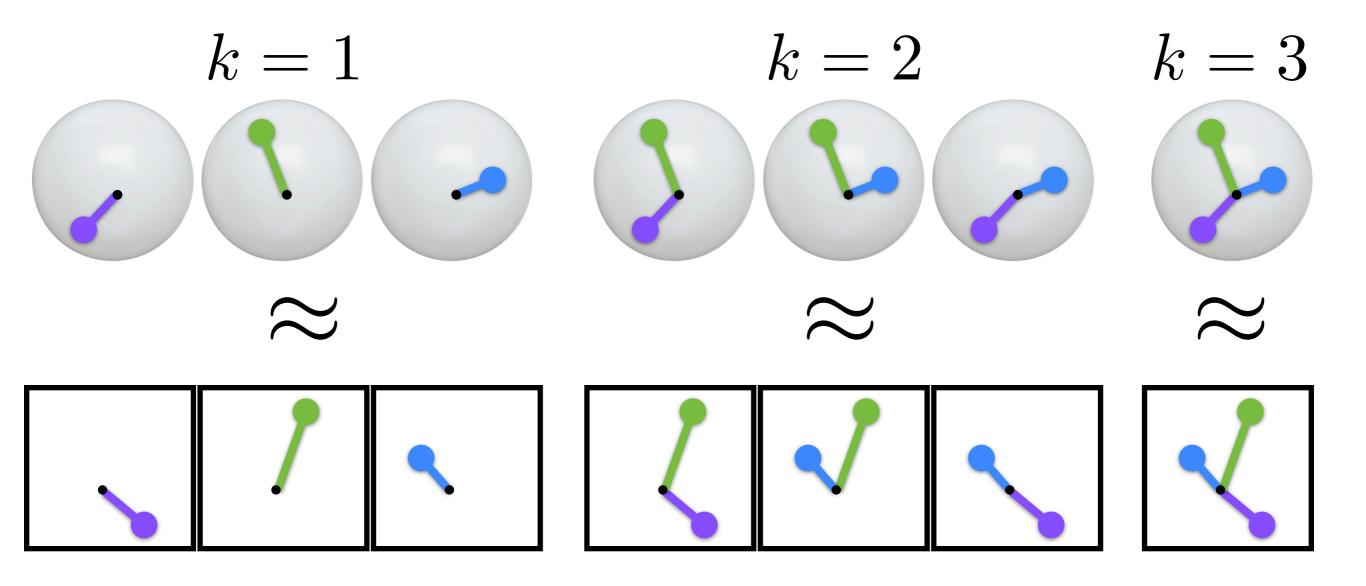
$$vol^2 = det$$



$$vol^2 = det$$



$$vol^2 = det$$



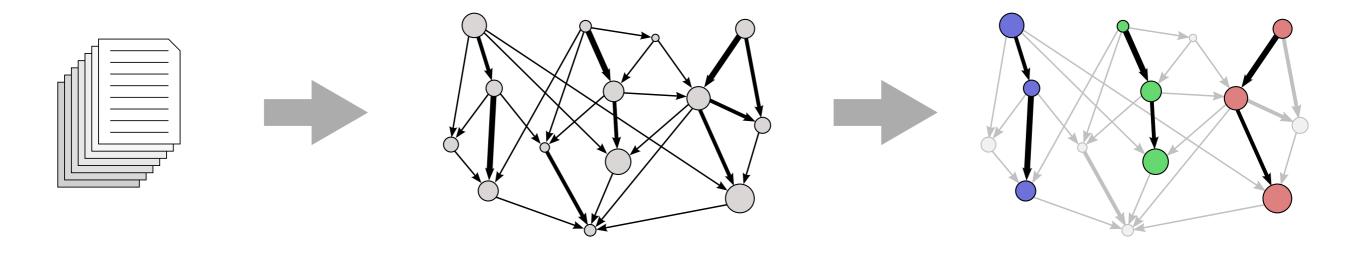
$$d = O\left(\max\left\{\frac{k}{\epsilon}, \frac{\log(1/\delta) + \log(N)}{\epsilon^2} + k\right\}\right)$$

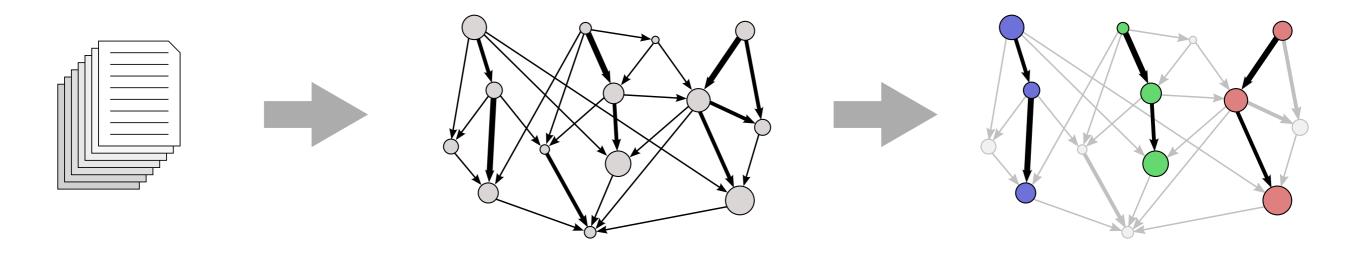
subset size total # of items
$$d = O\left(\max\left\{\frac{\frac{k}{\epsilon}}{\epsilon}, \frac{\log(1/\delta) + \log(N)}{\epsilon^2} + k\right\}\right)$$

subset size total # of items
$$d = O\left(\max\left\{\frac{\frac{k}{\epsilon}}{\epsilon}, \frac{\log(1/\delta) + \log(N)}{\epsilon^2} + k\right\}\right)$$
w.p.  $1 - \delta: \|\mathcal{P}^k - \tilde{\mathcal{P}}^k\|_1 \le e^{6k\epsilon} - 1$ 

subset size total # of items 
$$d = O\left(\max\left\{\frac{k}{\epsilon}, \frac{\log(1/\delta) + \log(N)}{\epsilon^2} + k\right\}\right)$$
 w.p.  $1 - \delta: \|\mathcal{P}^k - \tilde{\mathcal{P}}^k\|_1 \le e^{6k\epsilon} - 1$ 

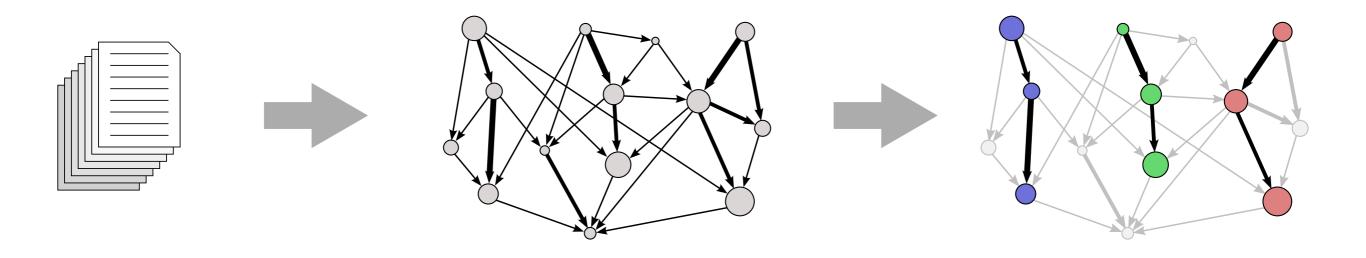
## STRUCTURED SUMMARIZATION



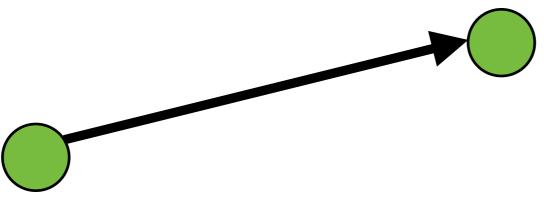




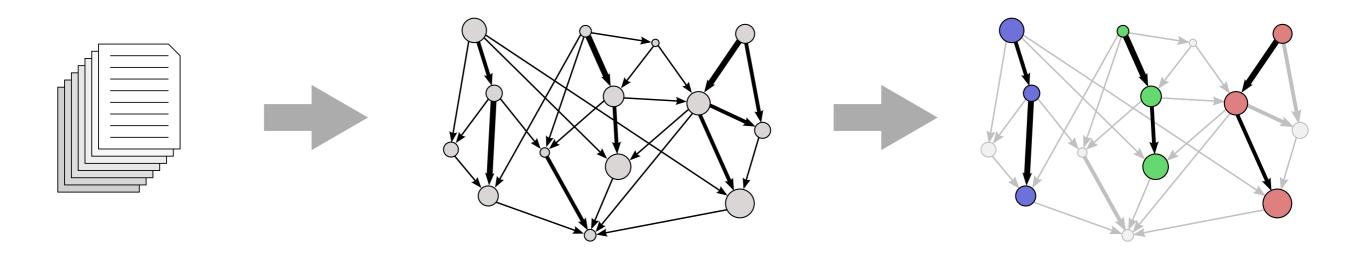
March 28: Health officials confirm Ebola outbreak in Guinea's capital



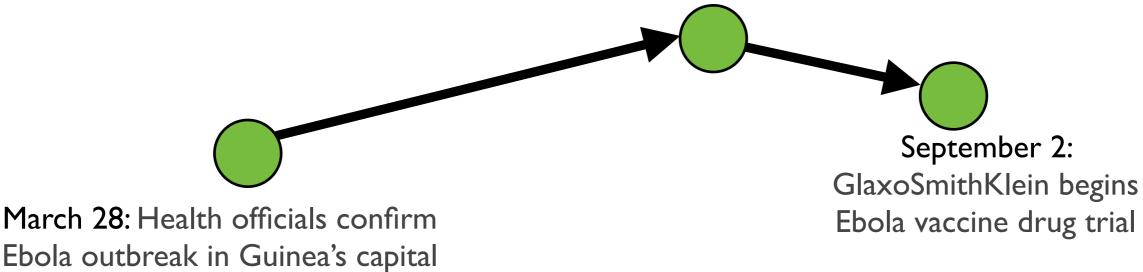
August 8: World Health Organization declares Ebola epidemic an international health emergency

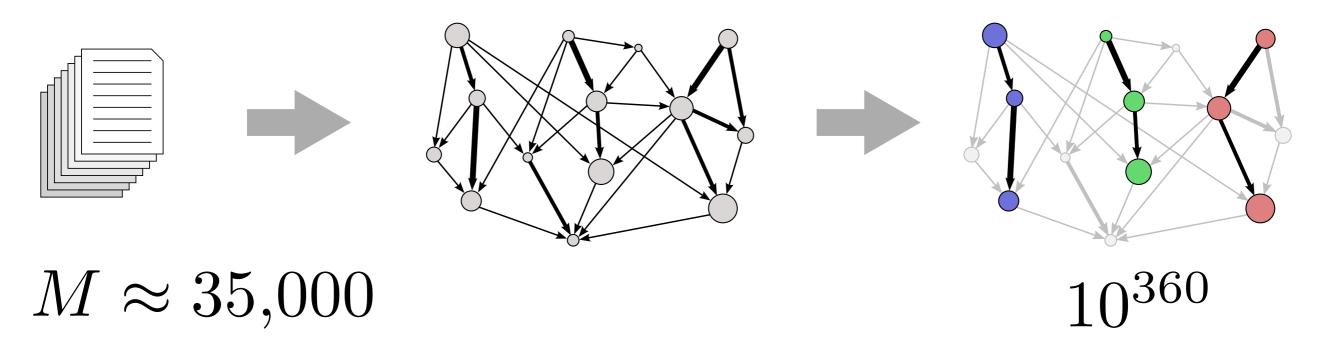


March 28: Health officials confirm Ebola outbreak in Guinea's capital

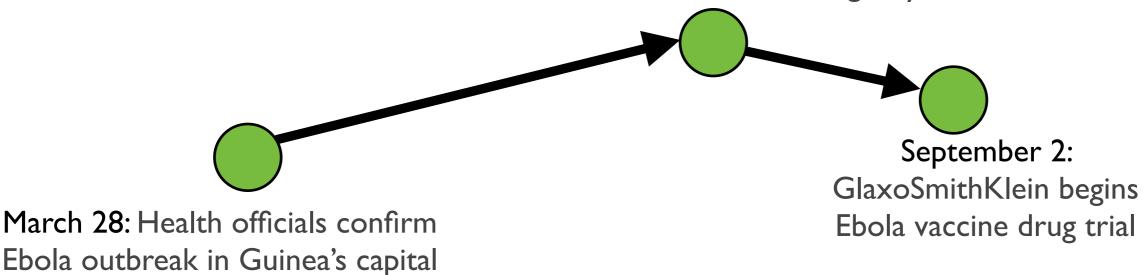


August 8: World Health Organization declares Ebola epidemic an international health emergency





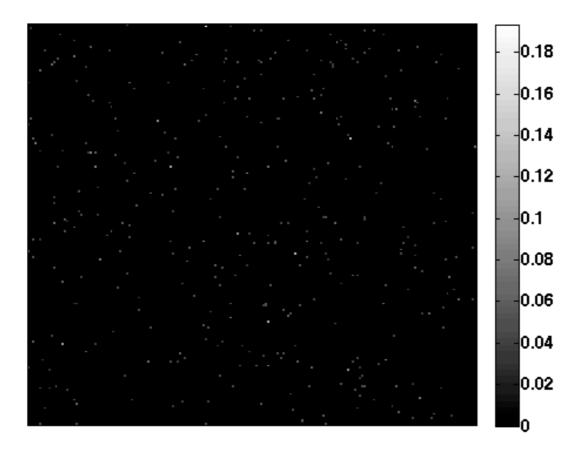
August 8: World Health Organization declares Ebola epidemic an international health emergency



## PROJECTING NEWS FEATURES

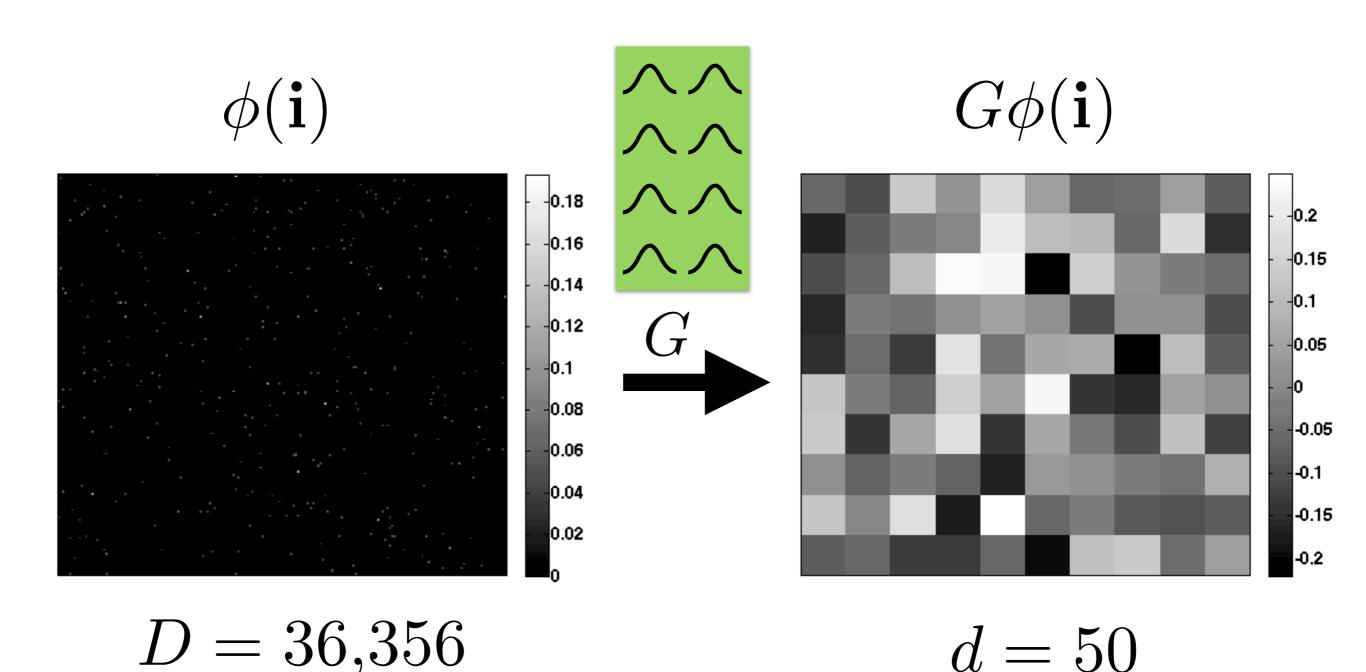
## PROJECTING NEWS FEATURES

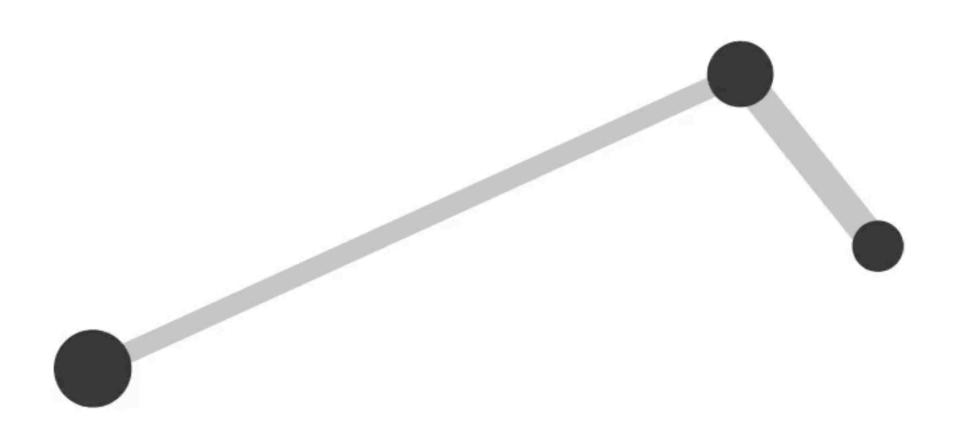
$$\phi(\mathbf{i})$$

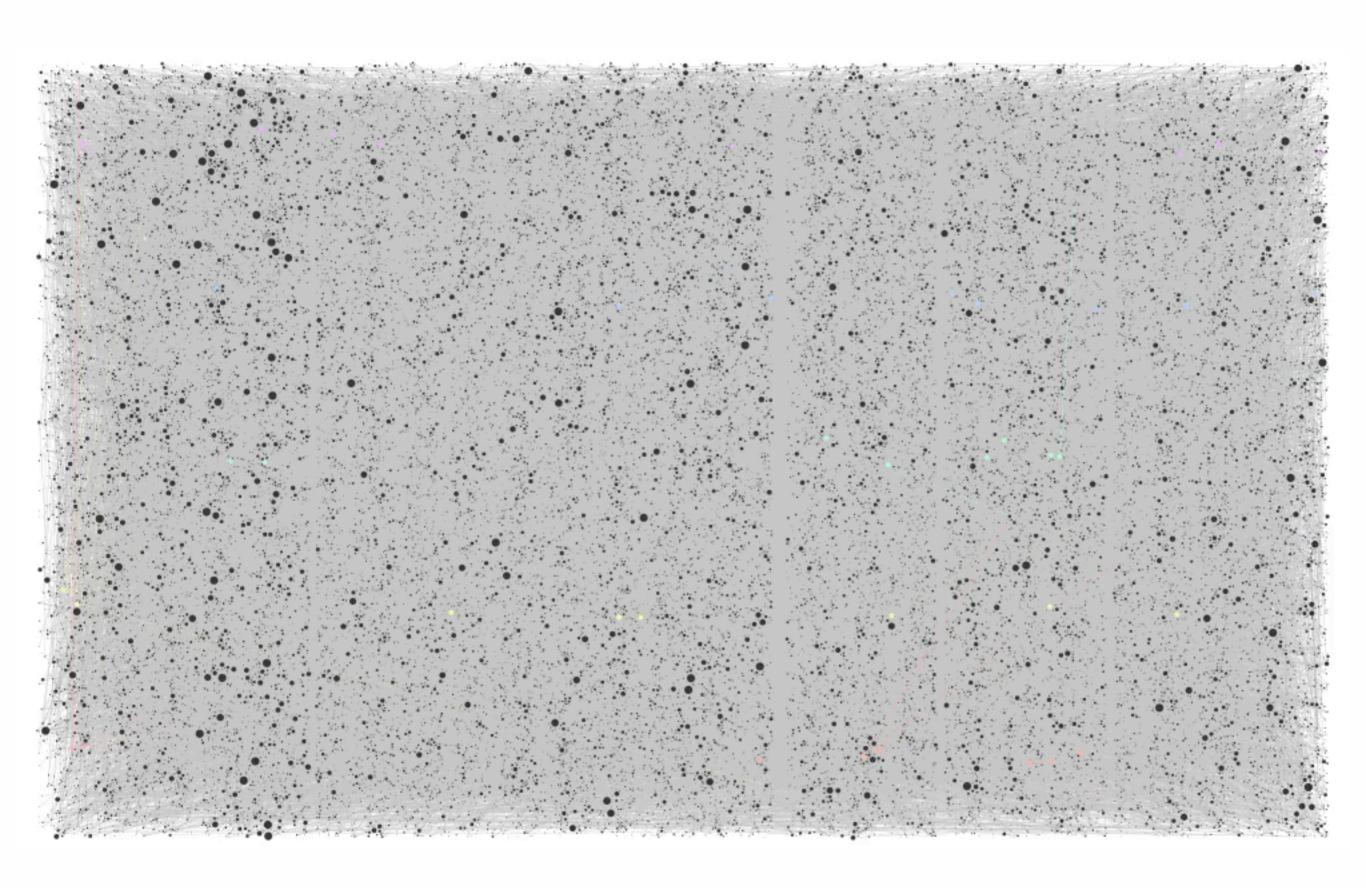


$$D = 36,356$$

## PROJECTING NEWS FEATURES







## DPPTHREADS

#### **DPPTHREADS**

iraq iraqi killed baghdad arab marines deaths forces

social tax security democrats rove accounts

owen nominees senate democrats judicial filibusters

israel palestinian iraqi israeli gaza abbas baghdad

pope vatican church parkinson

Jan 08 Jan 28 Feb 17 Mar 09 Mar 29 Apr 18 May 08 May 28 Jun 17

### DPPTHREADS



- Feb 24: Parkinson's Disease Increases Risks to Pope
- Feb 26: Pope's Health Raises Questions About His Ability to Lead
- Mar 13: Pope Returns Home After 18 Days at Hospital
- Apr 01: Pope's Condition Worsens as World Prepares for End of Papacy
- Apr 02: Pope, Though Gravely III, Utters Thanks for Prayers
- **Apr 18**: Europeans Fast Falling Away from Church
- Apr 20: In Developing World, Choice [of Pope] Met with Skepticism
- May 18: Pope Sends Message with Choice of Name

System	
ROUGE-1F	
R-SU4F	
Coherence	

System	k-means	
ROUGE-1F	16.5	
R-SU4F	3.76	
Coherence	2.73	

System	k-means	DTM	
ROUGE-1F	16.5	14.7	
R-SU4F	3.76	3.44	
Coherence	2.73	3.2	

System	k-means	DTM	DPP
ROUGE-1F	16.5	14.7	17.2
R-SU4F	3.76	3.44	3.98
Coherence	2.73	3.2	3.3

System	k-means	DTM	DPP
ROUGE-1F	16.5	14.7	17.2
R-SU4F	3.76	3.44	3.98
Coherence	2.73	3.2	3.3
Runtime (s)	626	19,434	252

# OTHER POTENTIAL NLP APPLICATIONS

 Parser: simple model with local features defines basic scores for all possible parse trees

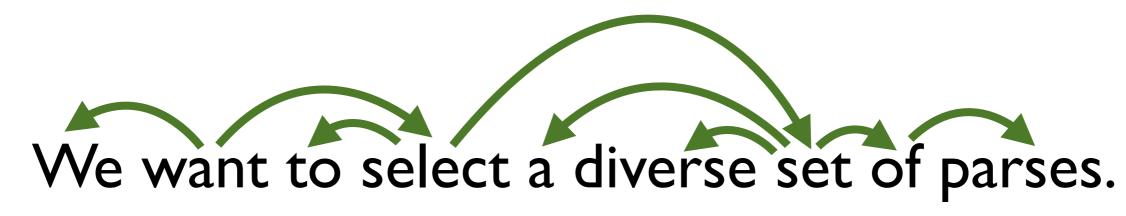
- Parser: simple model with local features defines basic scores for all possible parse trees
- Re-ranker: more complex model with non-local features provides more refined scores

- Parser: simple model with local features defines basic scores for all possible parse trees
- Re-ranker: more complex model with non-local features provides more refined scores
- Typical pipeline: find the k highest-scoring parses under the simple model, then score these k with the more complex model and output the best

- Parser: simple model with local features defines basic scores for all possible parse trees
- Re-ranker: more complex model with non-local features provides more refined scores
- Typical pipeline: find the k highest-scoring parses under the simple model, then score these k with the more complex model and output the best
- Issue: the k may be largely redundant, so reranker does not get to consider significantly different parses

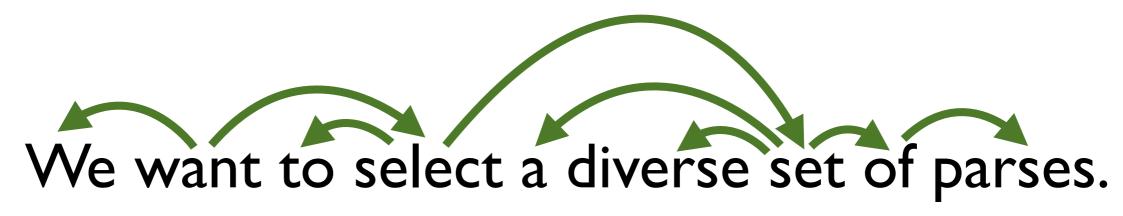
# IDEA: USE DPPS FOR SELECTING RE-RANKER INPUT

## IDEA: USE DPPS FOR SELECTING RE-RANKER INPUT



 $N = O(\{\text{sentence length}\}^{\{\text{sentence length}\}})$ 

# IDEA: USE DPPS FOR SELECTING RE-RANKER INPUT

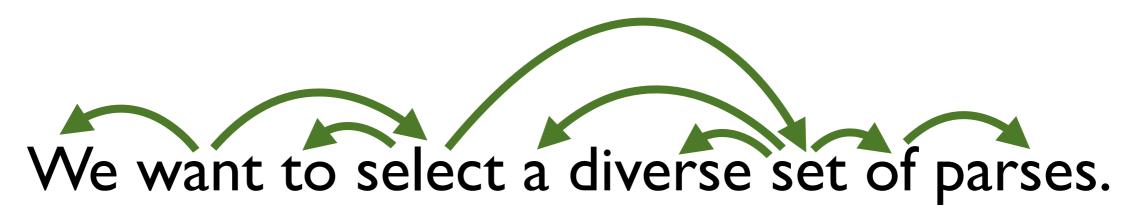


 $N = O(\{\text{sentence length}\} \{\text{sentence length}\})$ 

#### Quality:

standard parser scores

## IDEA: USE DPPS FOR SELECTING RE-RANKER INPUT



 $N = O(\{\text{sentence length}\} \{\text{sentence length}\})$ 

#### Quality:

standard parser scores

#### **Diversity**:

edge lengths, POS pairs, etc.

 Goal: identify all possible senses of ambiguous words (e.g. river bank vs bank deposit)

- Goal: identify all possible senses of ambiguous words (e.g. river bank vs bank deposit)
- Typical approach: unsupervised clustering, cluster centers represent word senses

- Goal: identify all possible senses of ambiguous words (e.g. river bank vs bank deposit)
- Typical approach: unsupervised clustering, cluster centers represent word senses
- Why DPPs fit: can re-express finding cluster centers as the problem of finding a high-quality, diverse set

- Goal: identify all possible senses of ambiguous words (e.g. river bank vs bank deposit)
- Typical approach: unsupervised clustering, cluster centers represent word senses
- Why DPPs fit: can re-express finding cluster centers as the problem of finding a high-quality, diverse set

#### Quality:

centrality (density of points nearby)

- Goal: identify all possible senses of ambiguous words (e.g. river bank vs bank deposit)
- Typical approach: unsupervised clustering, cluster centers represent word senses
- Why DPPs fit: can re-express finding cluster centers as the problem of finding a high-quality, diverse set

Quality:

centrality (density of points nearby)

**Diversity**:

same as standard WSI features

## QUESTIONS?