

# Learning about Language with Normalizing Flows

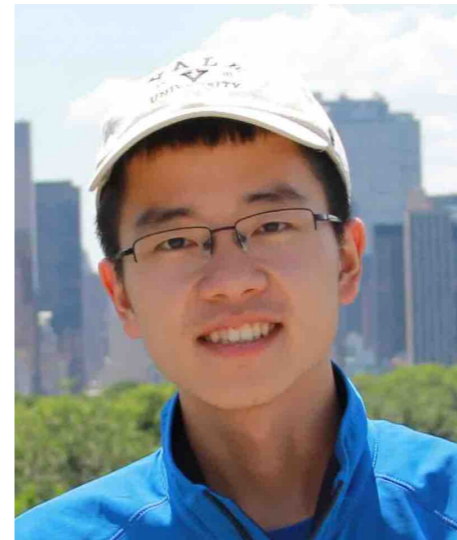
**Graham Neubig**

Language Technologies Institute, Carnegie Mellon University

Chunting Zhou



Junxian He



Di Wang, Xuezhe Ma, Daniel Spokoyny, Taylor Berg-Kirkpatrick

# Learning about Language?

# Learning about Language?

- **Syntactic structure**

# Learning about Language?

- **Syntactic structure**

The cat sat on a green wall

# Learning about Language?

- **Syntactic structure**

The cat sat on a green wall

Parts-of-speech: DT NN VBD IN DT JJ NN

# Learning about Language?

- **Syntactic structure**

The cat sat on a green wall

Parts-of-speech: DT NN VBD IN DT JJ NN

Dependency: 

# Learning about Language?

- **Syntactic structure**

The cat sat on a green wall

Parts-of-speech: DT NN VBD IN DT JJ NN

Dependency: 

- **Cross-lingual correspondences**

# Learning about Language?

- Syntactic structure

The cat sat on a green wall

Parts-of-speech: DT NN VBD IN DT JJ NN

Dependency: 

- Cross-lingual correspondences

a cat green on sat the wall

のは上壁猫緑座った



# Learning about Language?

- Syntactic structure

The cat sat on a green wall

Parts-of-speech: DT NN VBD IN DT JJ NN

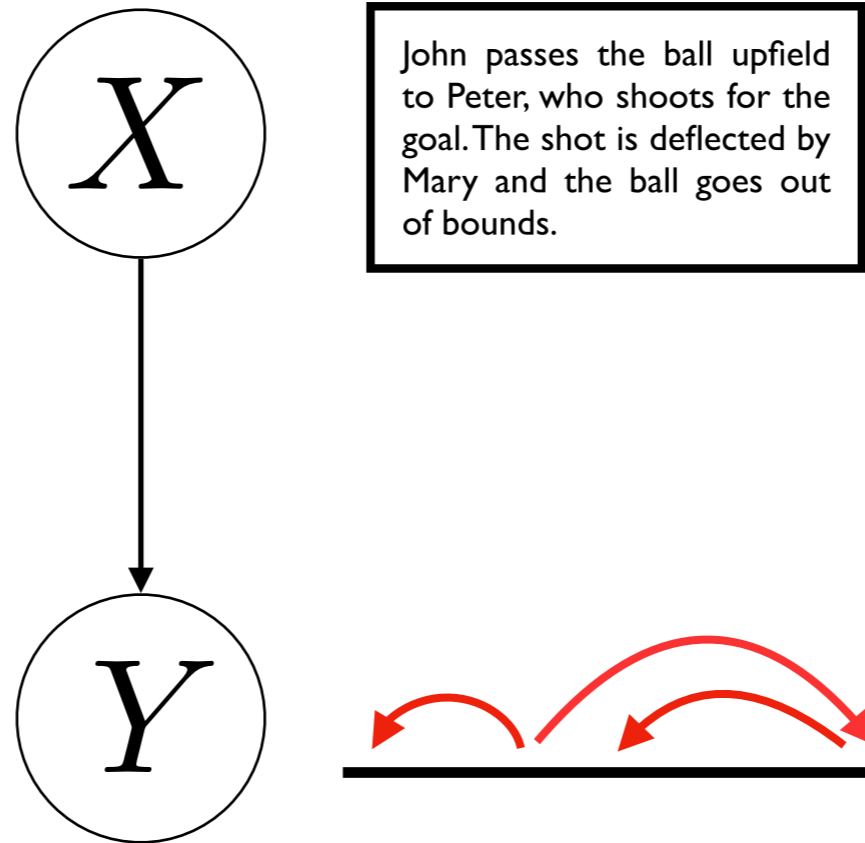
Dependency: 

- Cross-lingual correspondences

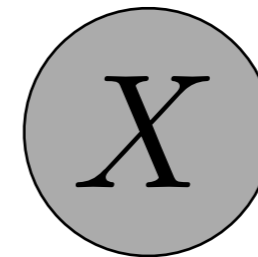
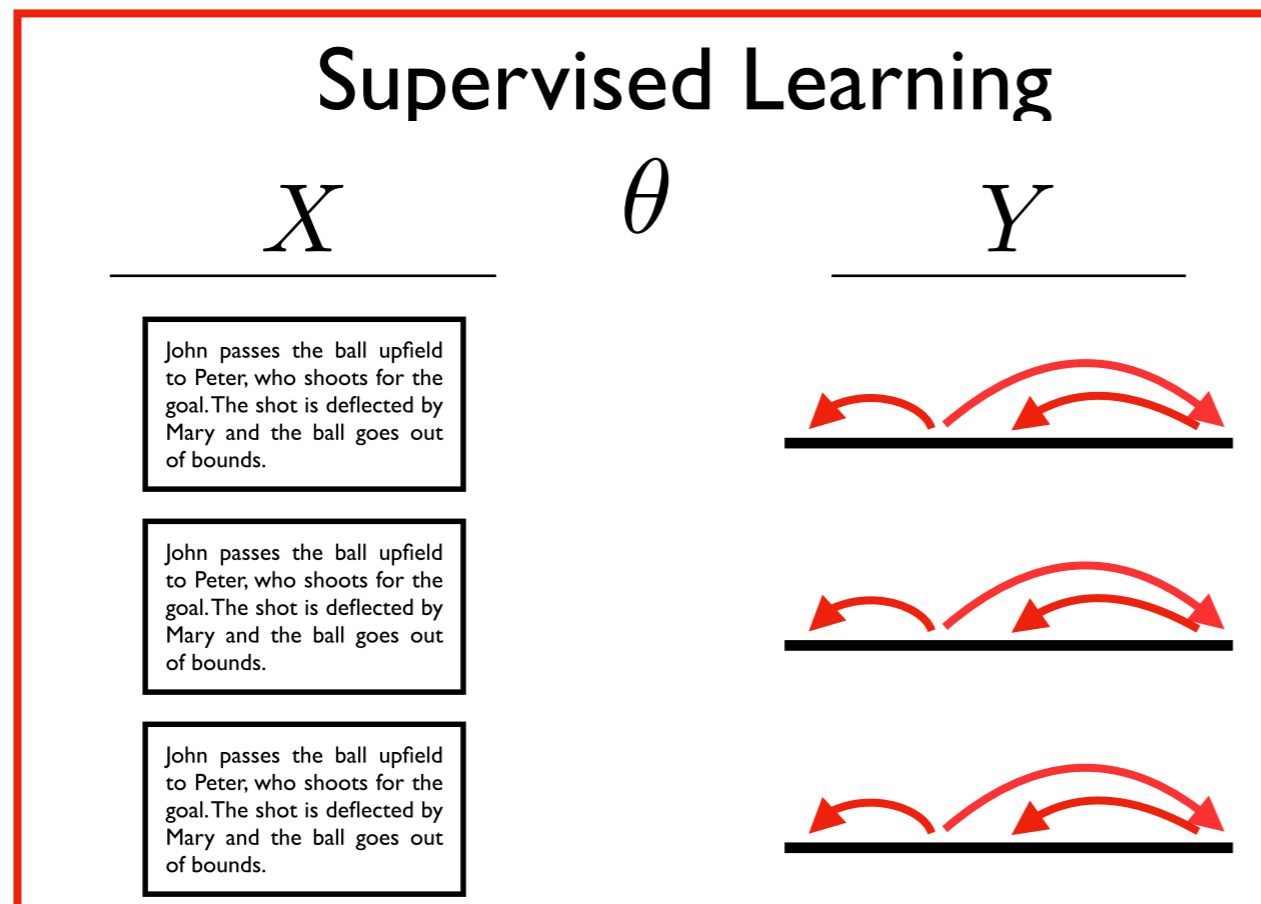
a cat green on sat the wall

の は 上 壁 猫 緑 座った

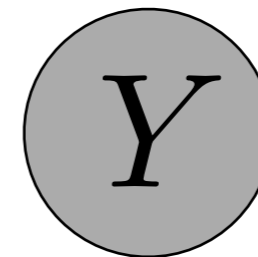
# Supervised Approaches



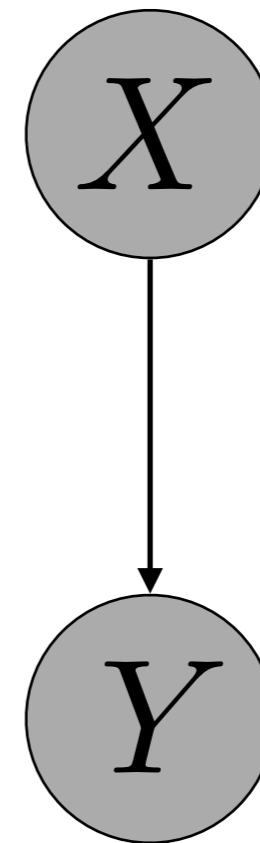
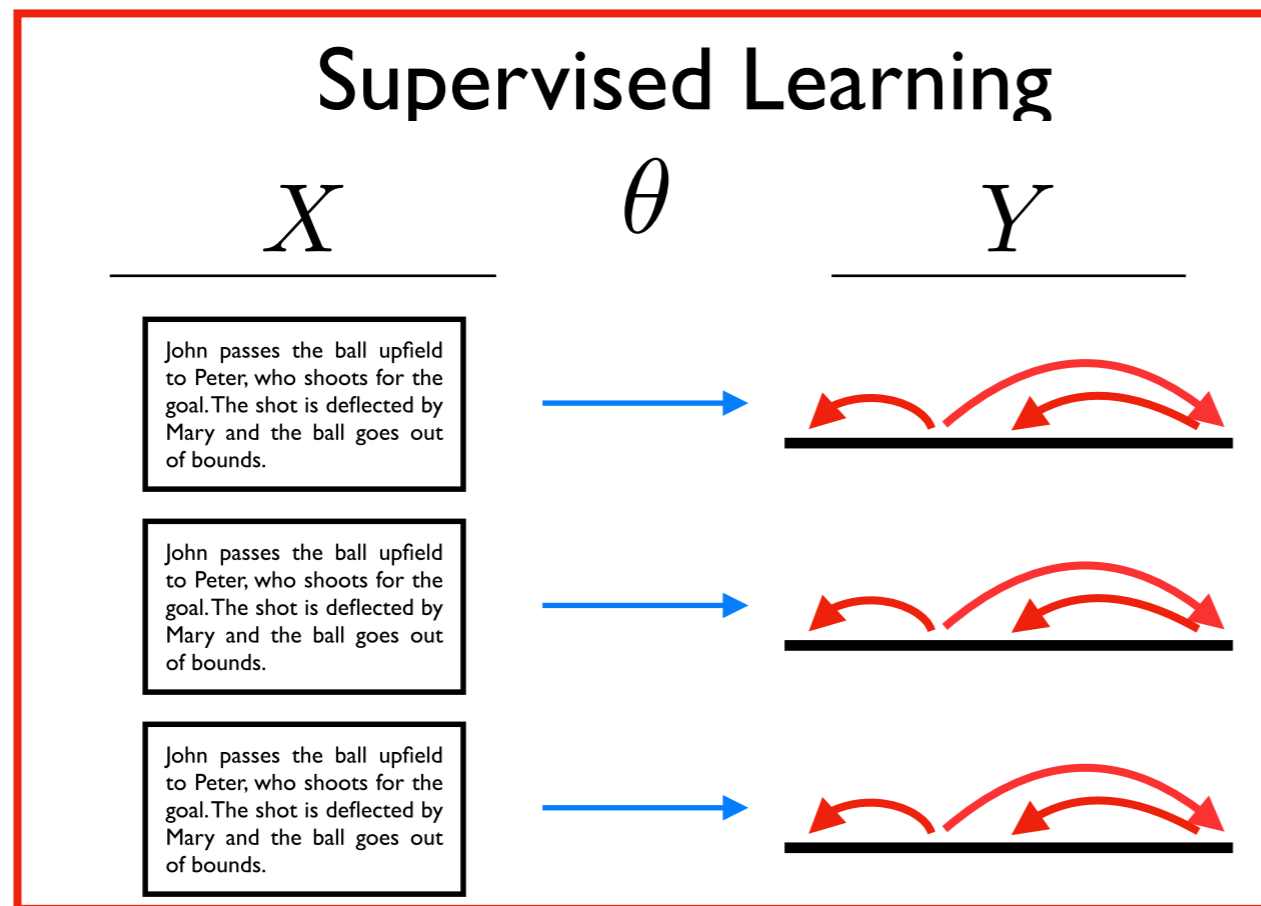
# Supervised Approaches



John passes the ball upfield to Peter, who shoots for the goal. The shot is deflected by Mary and the ball goes out of bounds.



# Supervised Approaches



John passes the ball upfield to Peter, who shoots for the goal. The shot is deflected by Mary and the ball goes out of bounds.

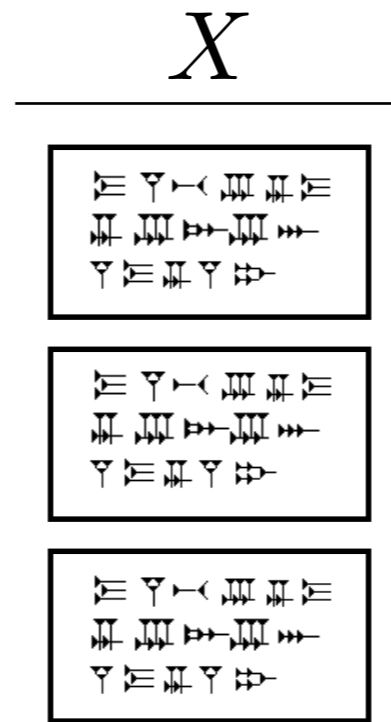


# Unsupervised Approaches

$$X$$

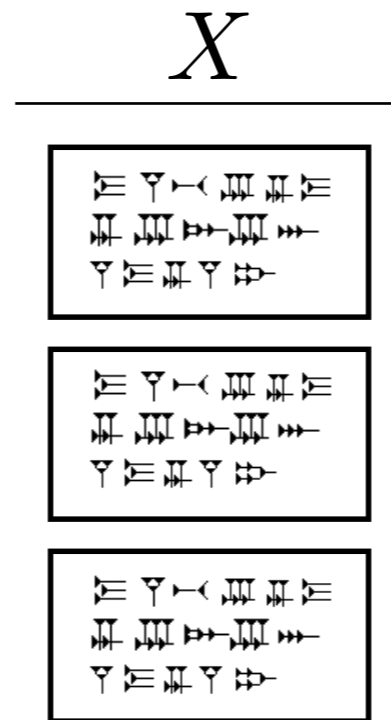

---

# Unsupervised Approaches



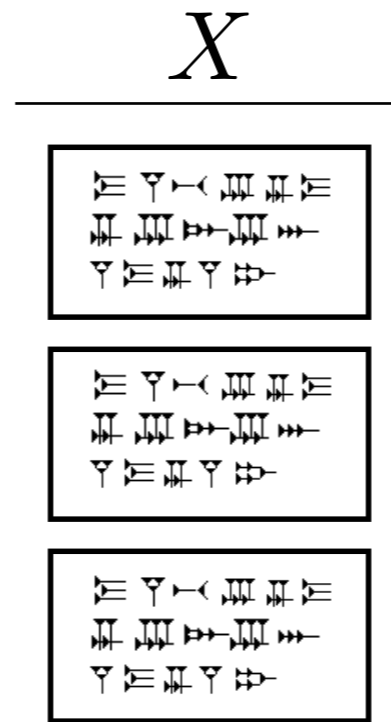
- Learning language models  $P(X)$

# Unsupervised Approaches



- Learning language models  $P(X)$
- Learning continuous features from language models (e.g. word2vec, skipthought, BERT)

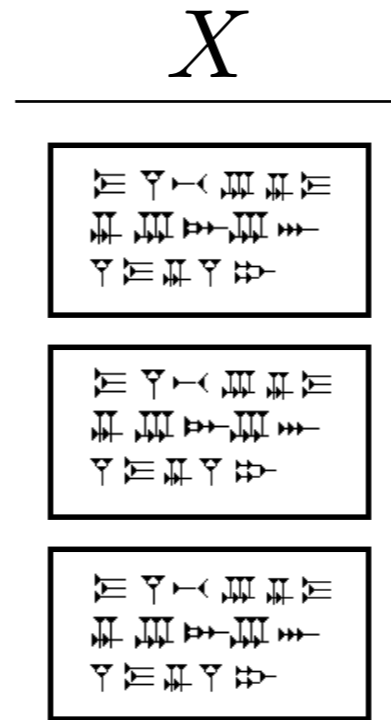
# Unsupervised Approaches



- Learning language models  $P(X)$
- Learning continuous features from language models (e.g. word2vec, skipthought, BERT)
- But how do we turn this into **interpretable structure**?

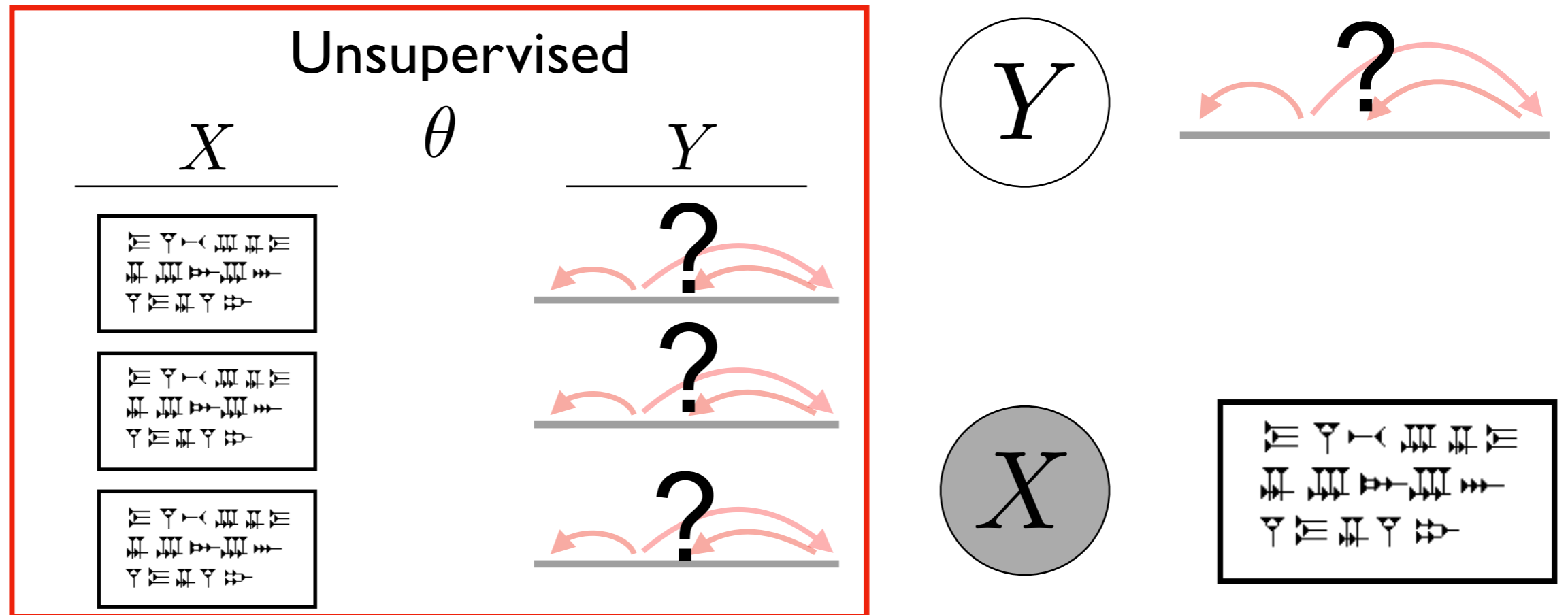


# Unsupervised Approaches

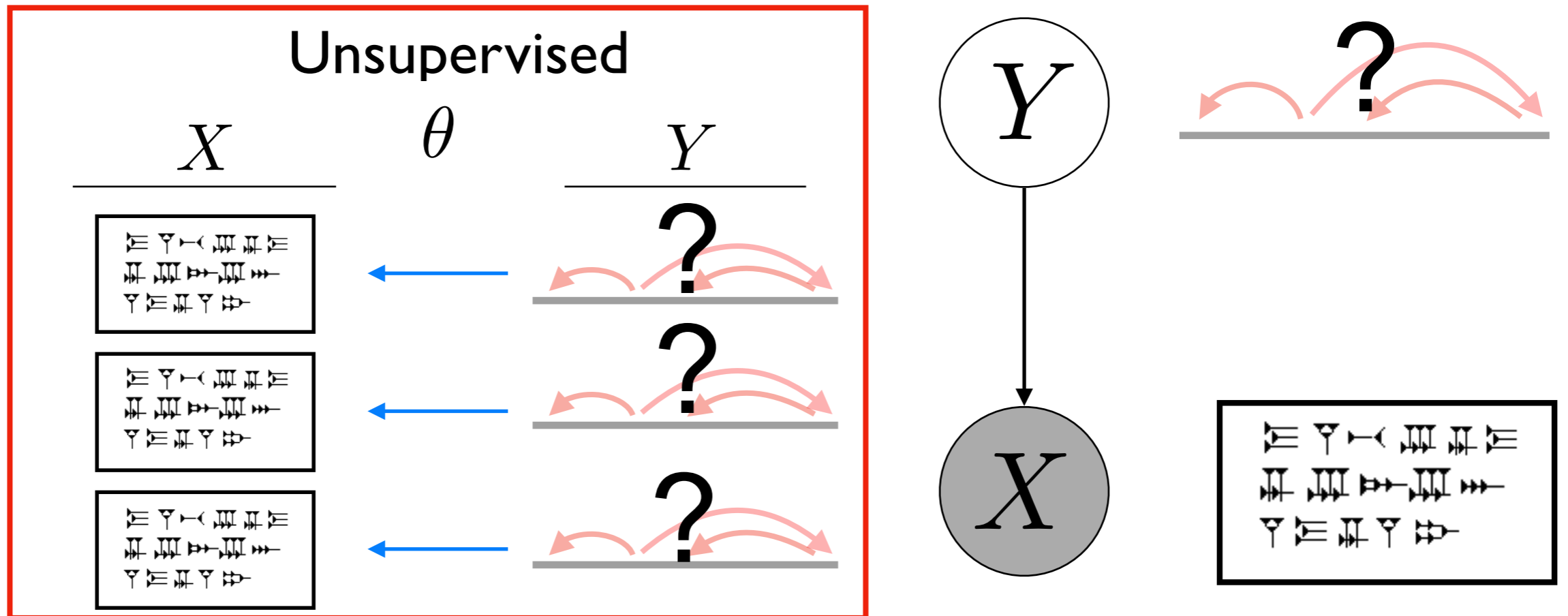


- Learning language models  $P(X)$
- Learning continuous features from language models (e.g. word2vec, skipthought, BERT)
- But how do we turn this into **interpretable structure?**
- How do we do it while **taking advantage of continuous features?**

# Latent Variable Approaches



# Latent Variable Approaches

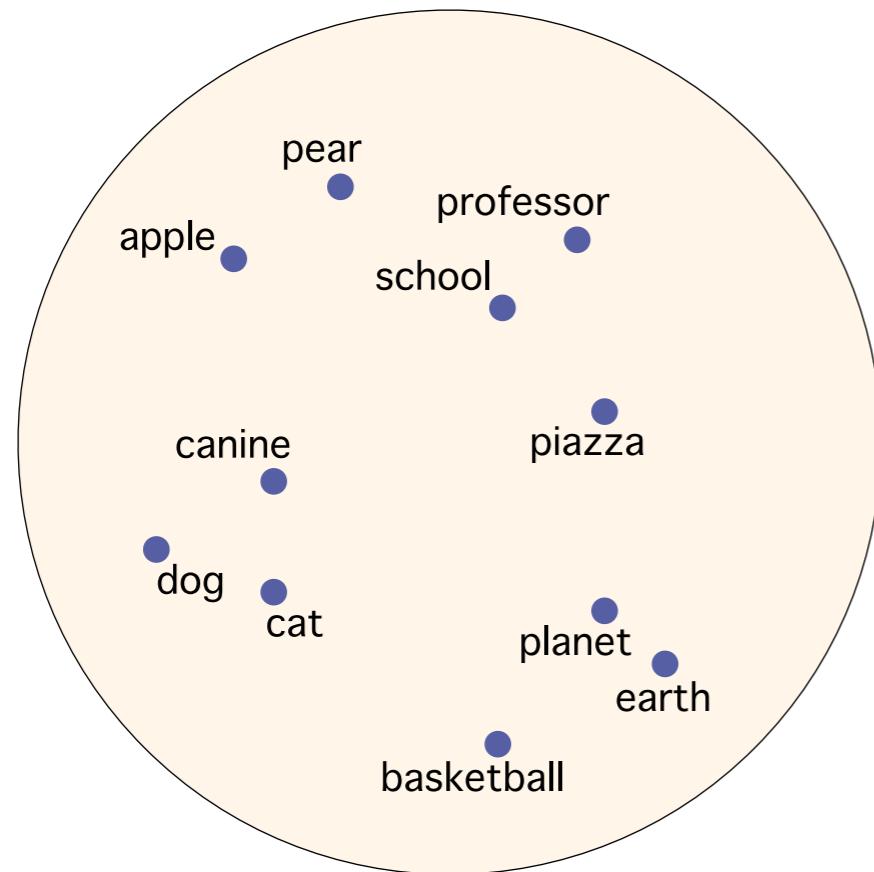


# Density Matching for Bilingual Word Embedding

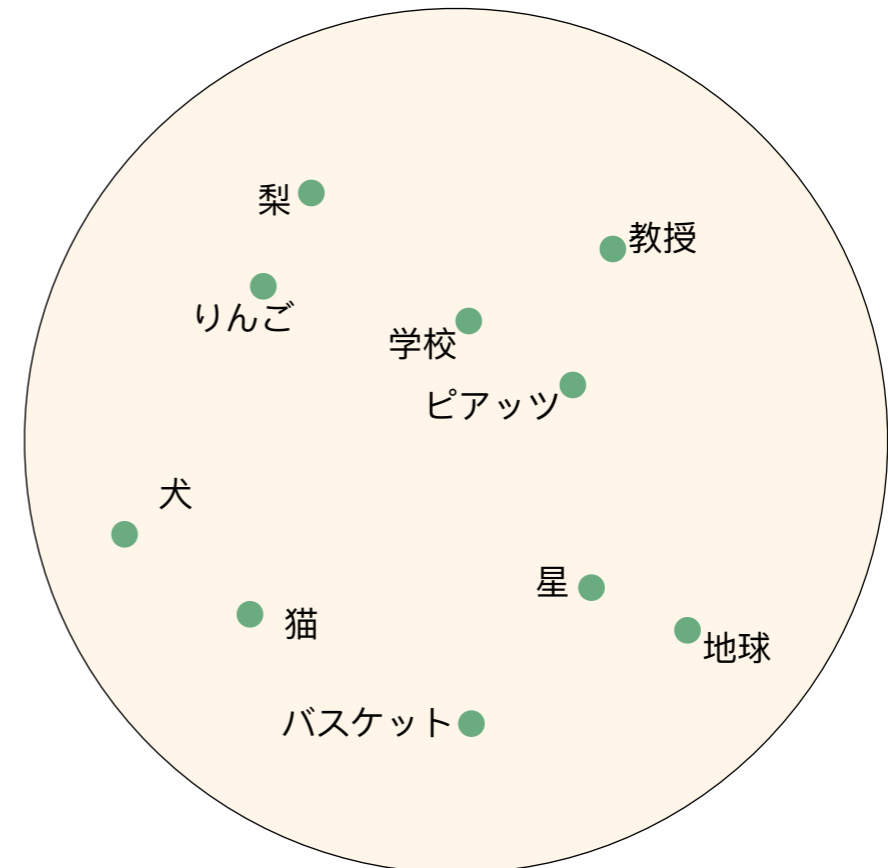
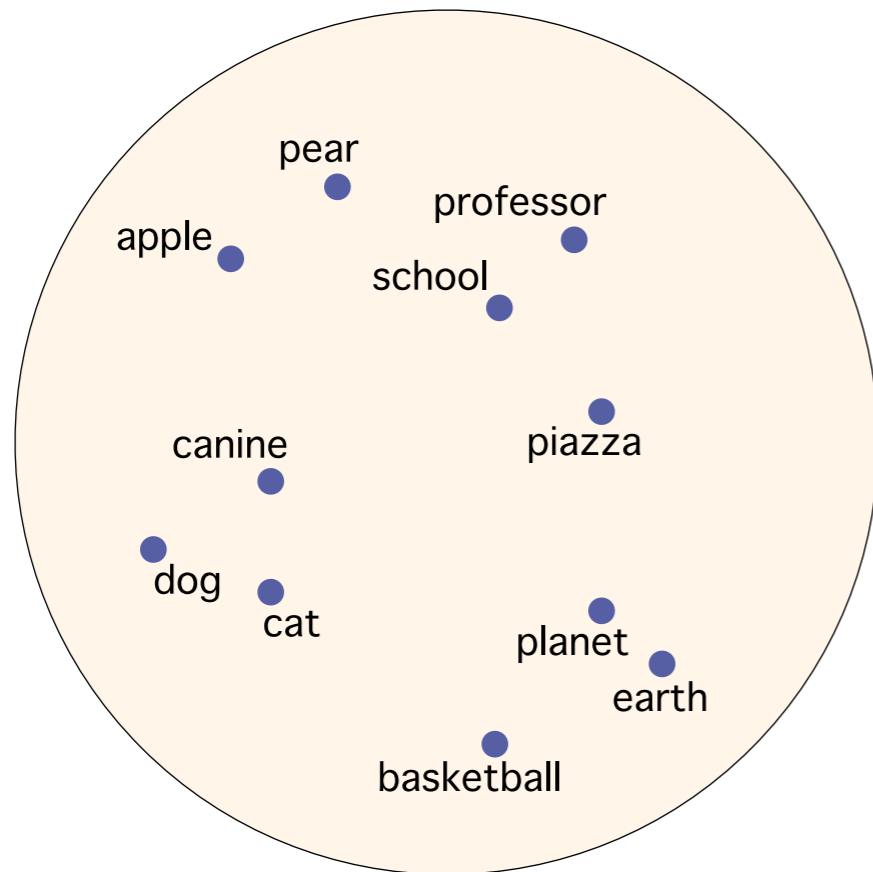
Chunting Zhou, Xuezhe Ma, Di Wang, Graham Neubig  
(NAACL 2019)

# Bilingual Word Embedding

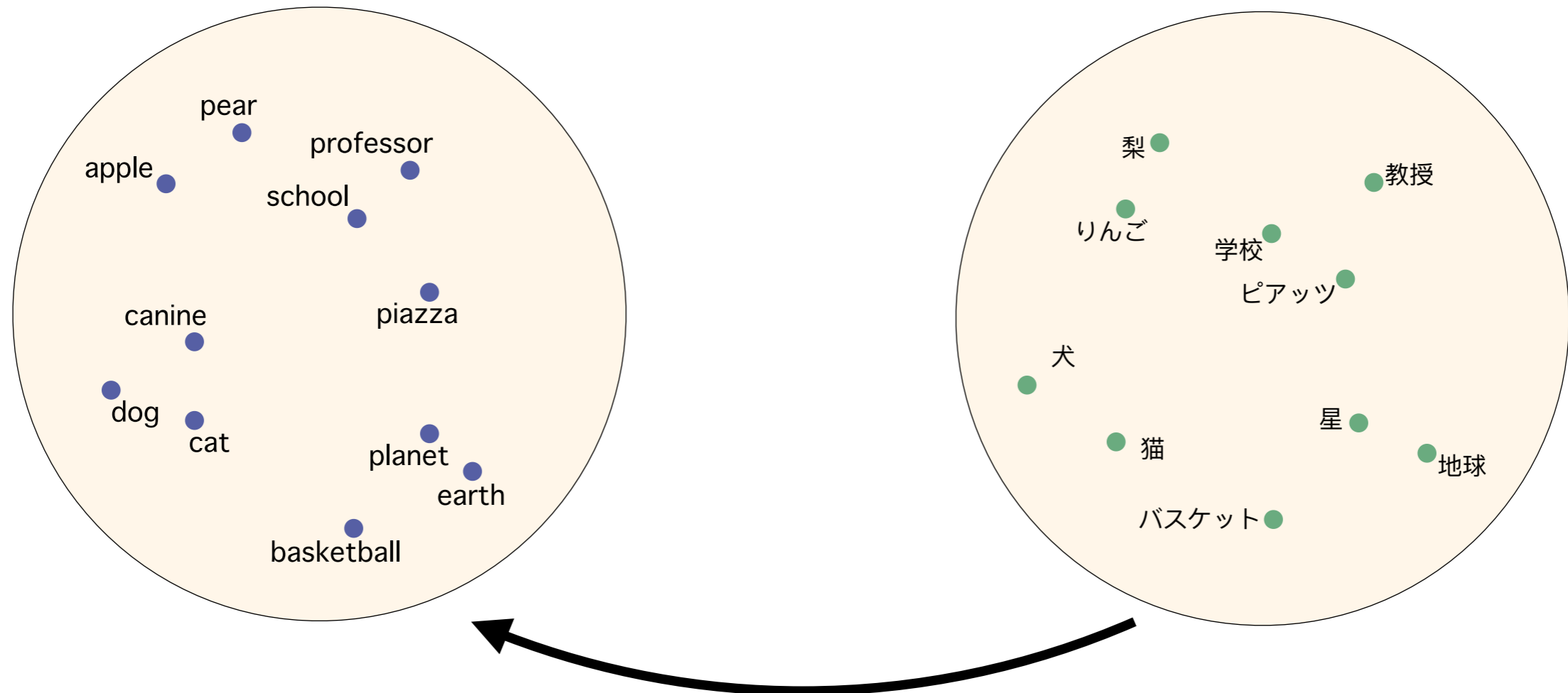
# Bilingual Word Embedding



# Bilingual Word Embedding



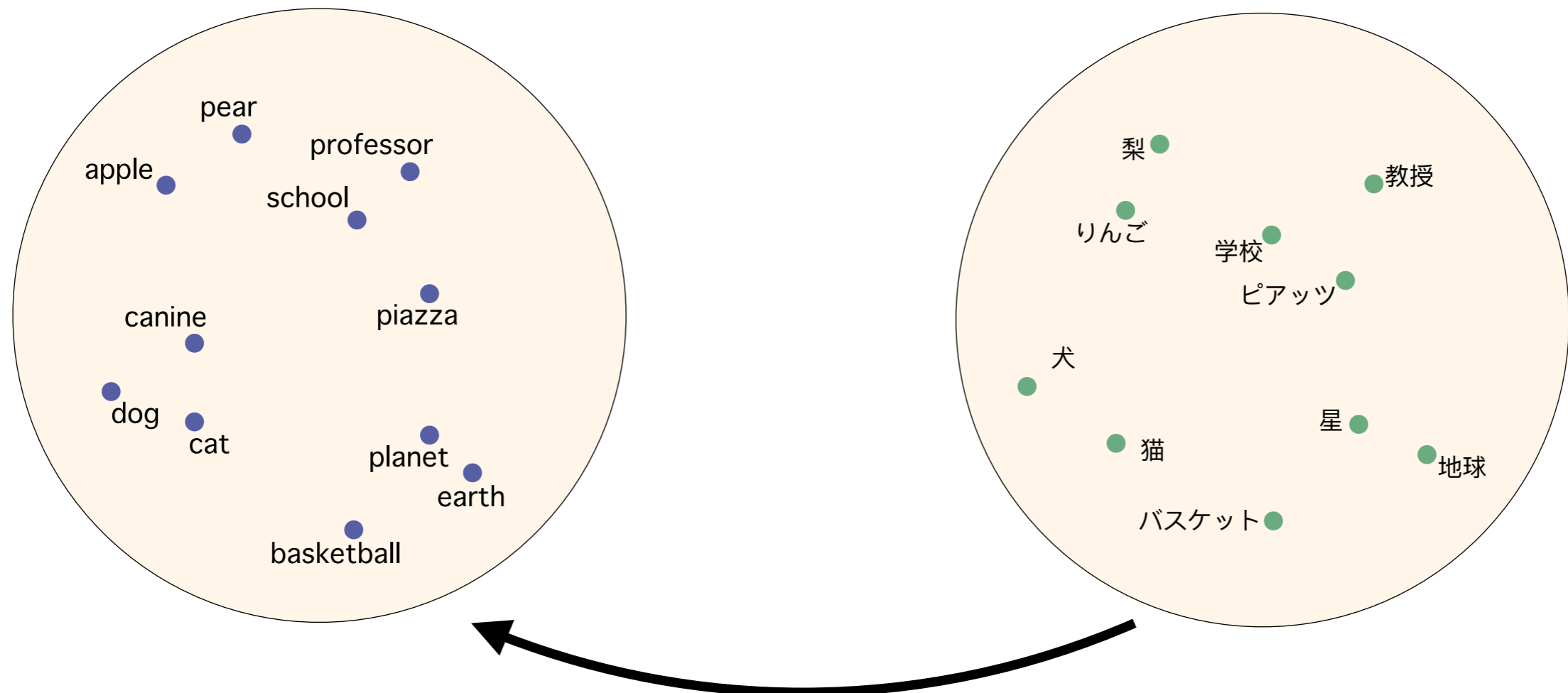
# Bilingual Word Embedding



- Map word embeddings from different languages into a single vector space

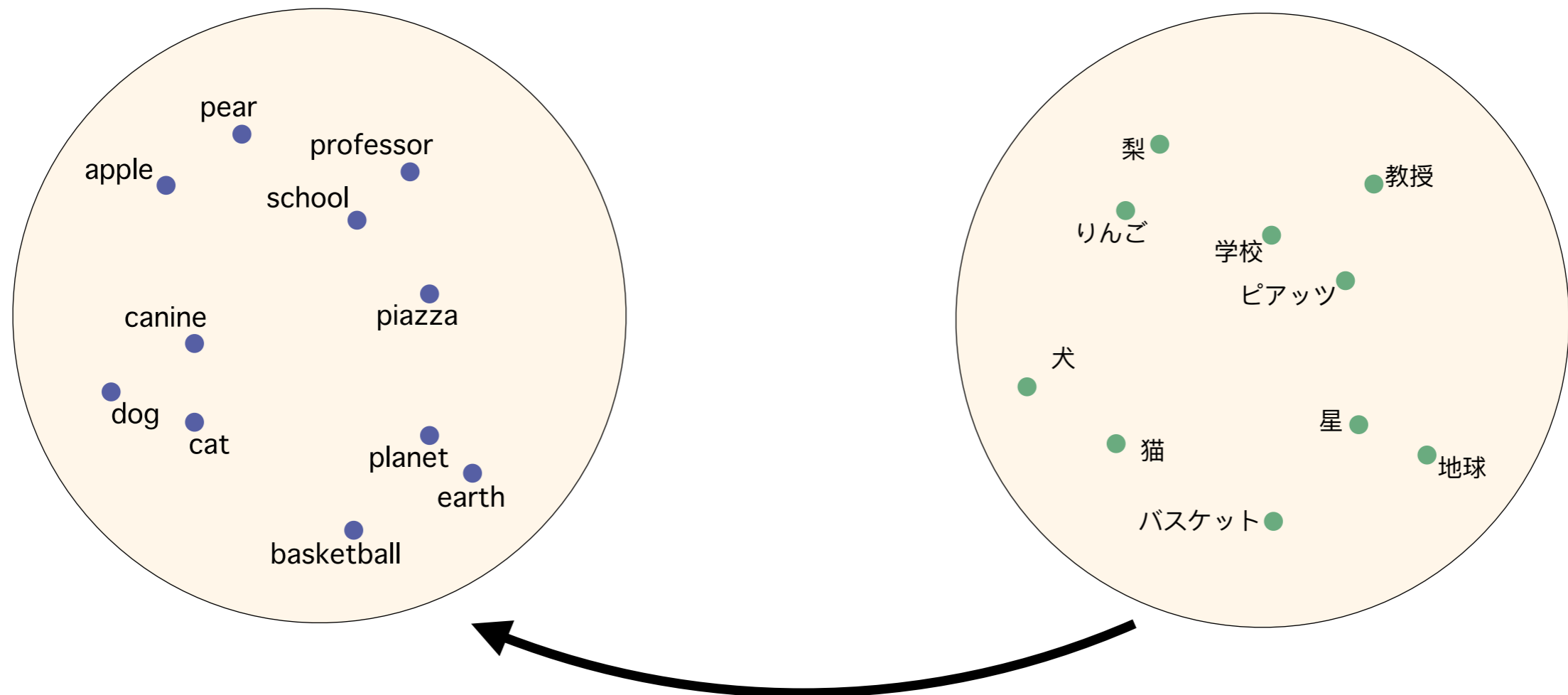


# Bilingual Word Embedding



- Map word embeddings from different languages into a single vector space
  - Cross-lingual transfer

# Bilingual Word Embedding



- Map word embeddings from different languages into a single vector space
  - Cross-lingual transfer
  - Cross-lingual NLP tasks

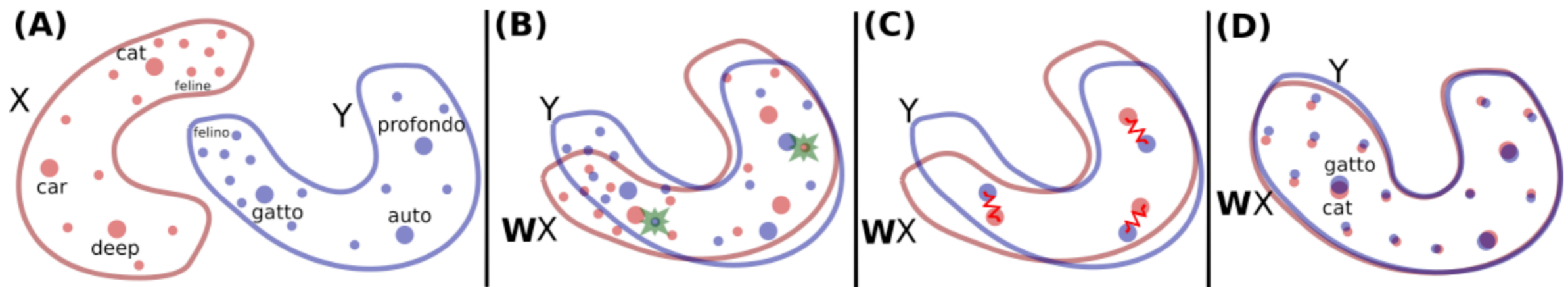
# Previous Work on Unsupervised BWE

# Previous Work on Unsupervised BWE

- Unsupervised methods of minimization some form of distance between distributions of discrete vector sets:

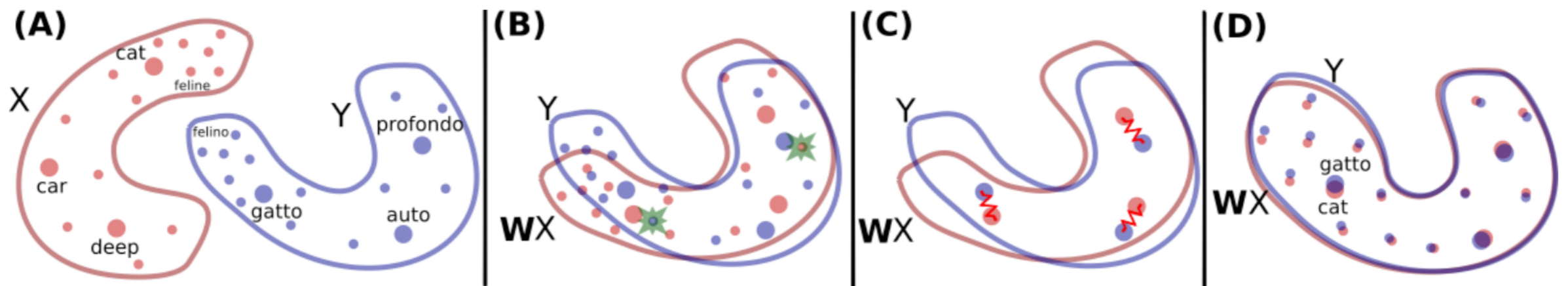
# Previous Work on Unsupervised BWE

- Unsupervised methods of minimization some form of distance between distributions of discrete vector sets:



# Previous Work on Unsupervised BWE

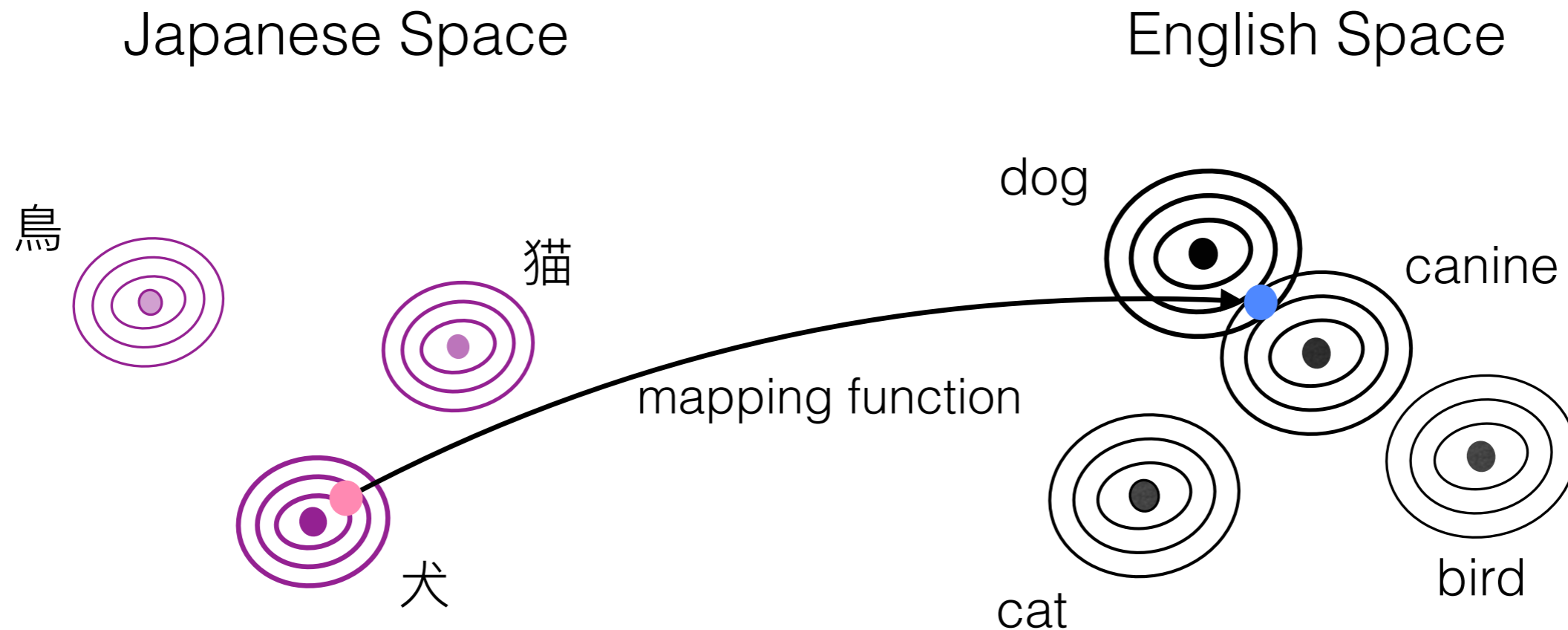
- Unsupervised methods of minimization some form of distance between distributions of discrete vector sets:



- No direct probabilistic interpretation, not a "typical" unsupervised generative model

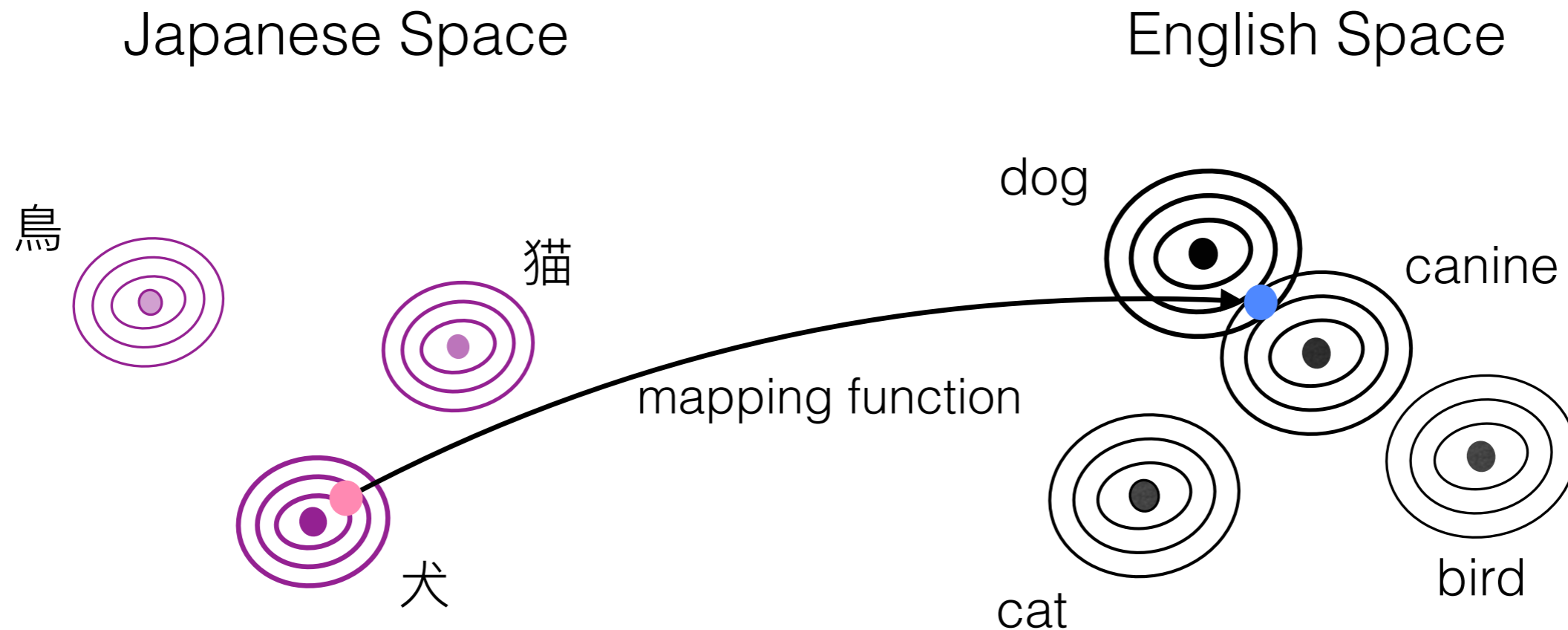
# Density Mapping for Bilingual Word Embedding (DeMa-BWE)

# Density Mapping for Bilingual Word Embedding (DeMa-BWE)





# Density Mapping for Bilingual Word Embedding (DeMa-BWE)



- Mapping function is learned with **normalizing flow**

# Normalizing Flows

# Normalizing Flows

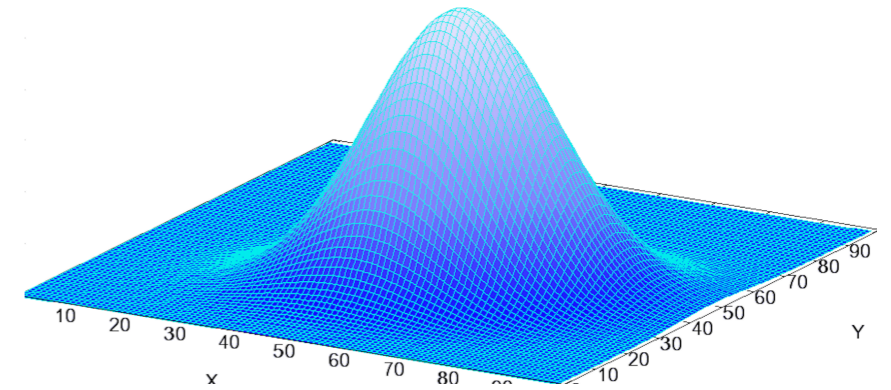


$$X \sim P(X)$$

# Normalizing Flows



$$X \sim P(X)$$



$$Z \sim N(0, I)$$

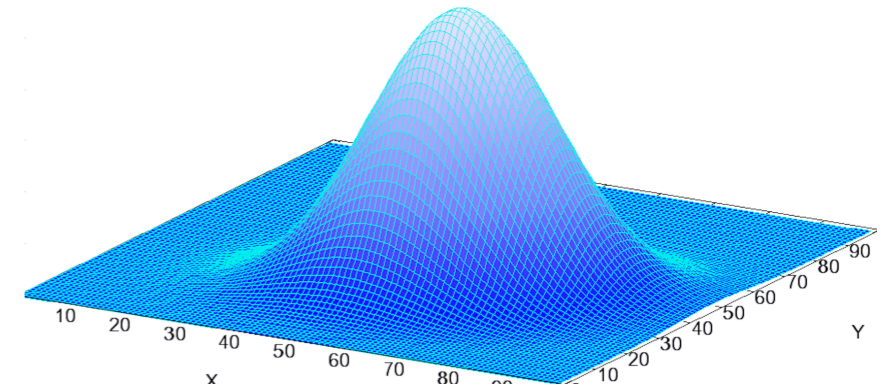
# Normalizing Flows



$$X \sim P(X)$$

$$X = f_{\theta}^{-1}(Z)$$

$$Z = f_{\theta}(X)$$



$$Z \sim N(0, I)$$

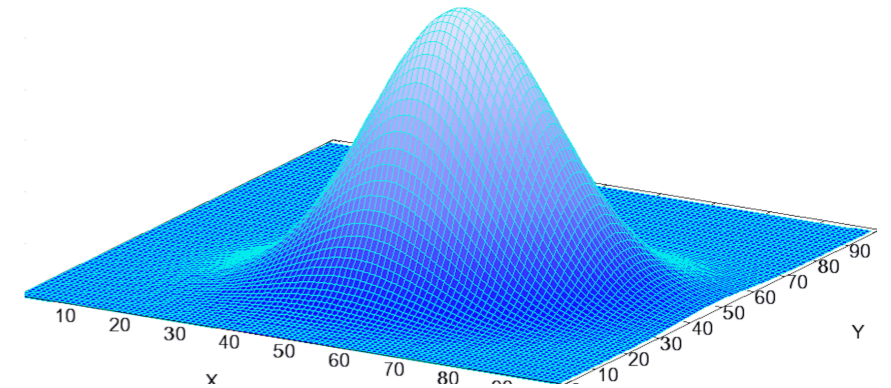
# Normalizing Flows



$$X \sim P(X)$$

$$X = f_{\theta}^{-1}(Z)$$

$$Z = f_{\theta}(X)$$



$$Z \sim N(0, I)$$

**Change of variable formula:**

$$p_{\theta}(x) = p_Z(f_{\theta}(x)) \left| \det\left(\frac{\partial f_{\theta}(x)}{\partial x}\right) \right|$$

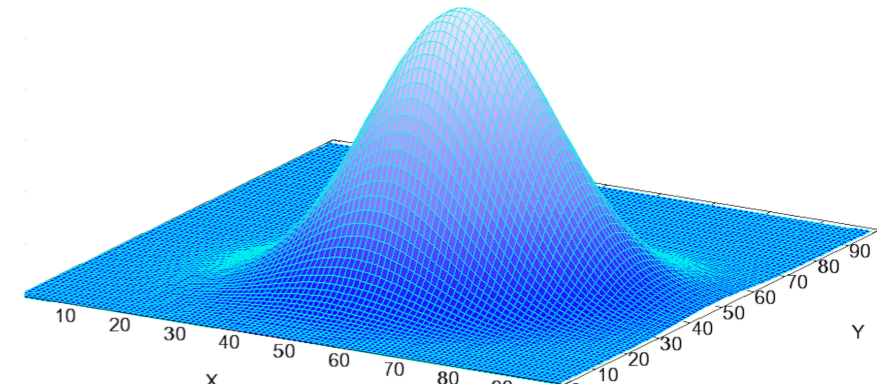
# Normalizing Flows



$$X \sim P(X)$$

$$X = f_{\theta}^{-1}(Z)$$

$$Z = f_{\theta}(X)$$



$$Z \sim N(0, I)$$

**Change of variable formula:**

$$p_{\theta}(x) = p_Z(f_{\theta}(x)) \left| \det \left( \frac{\partial f_{\theta}(x)}{\partial x} \right) \right|$$

Intuitively, prevents degenerative mapping of everything to zero vector

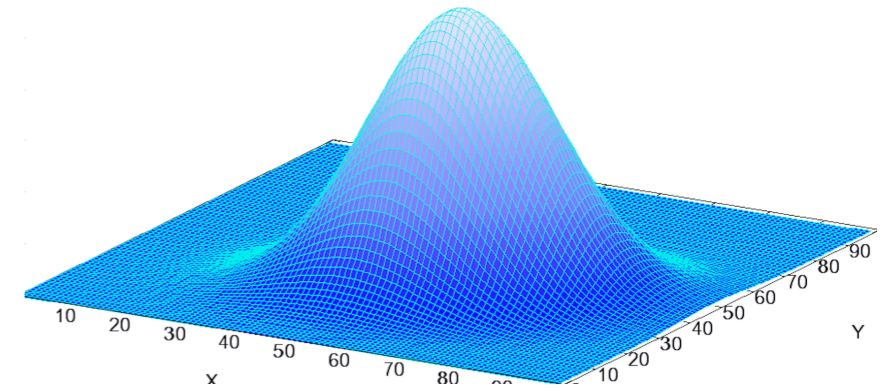
# Normalizing Flows



$$X \sim P(X)$$

$$X = f_{\theta}^{-1}(Z)$$

$$Z = f_{\theta}(X)$$



$$Z \sim N(0, I)$$

**Change of variable formula:**

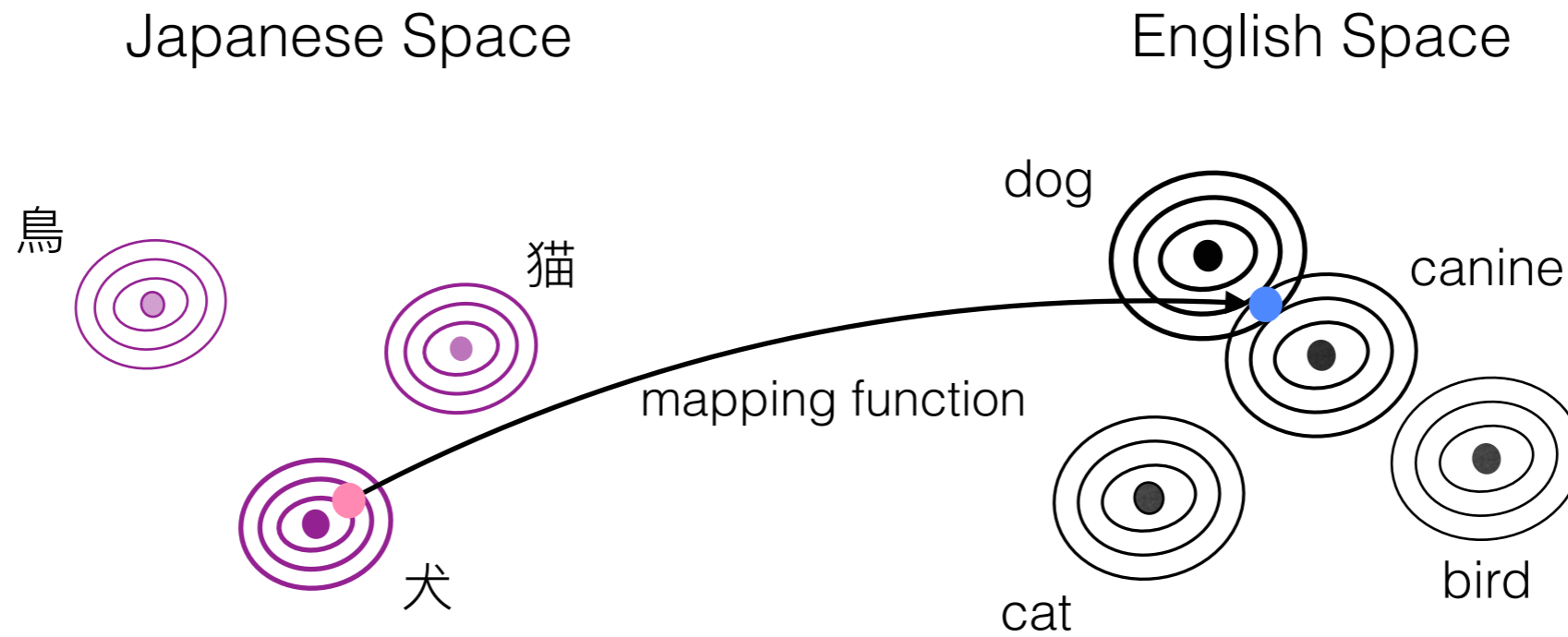
$$p_{\theta}(x) = p_Z(f_{\theta}(x)) \left| \det \left( \frac{\partial f_{\theta}(x)}{\partial x} \right) \right|$$

Intuitively, prevents degenerative mapping of everything to zero vector

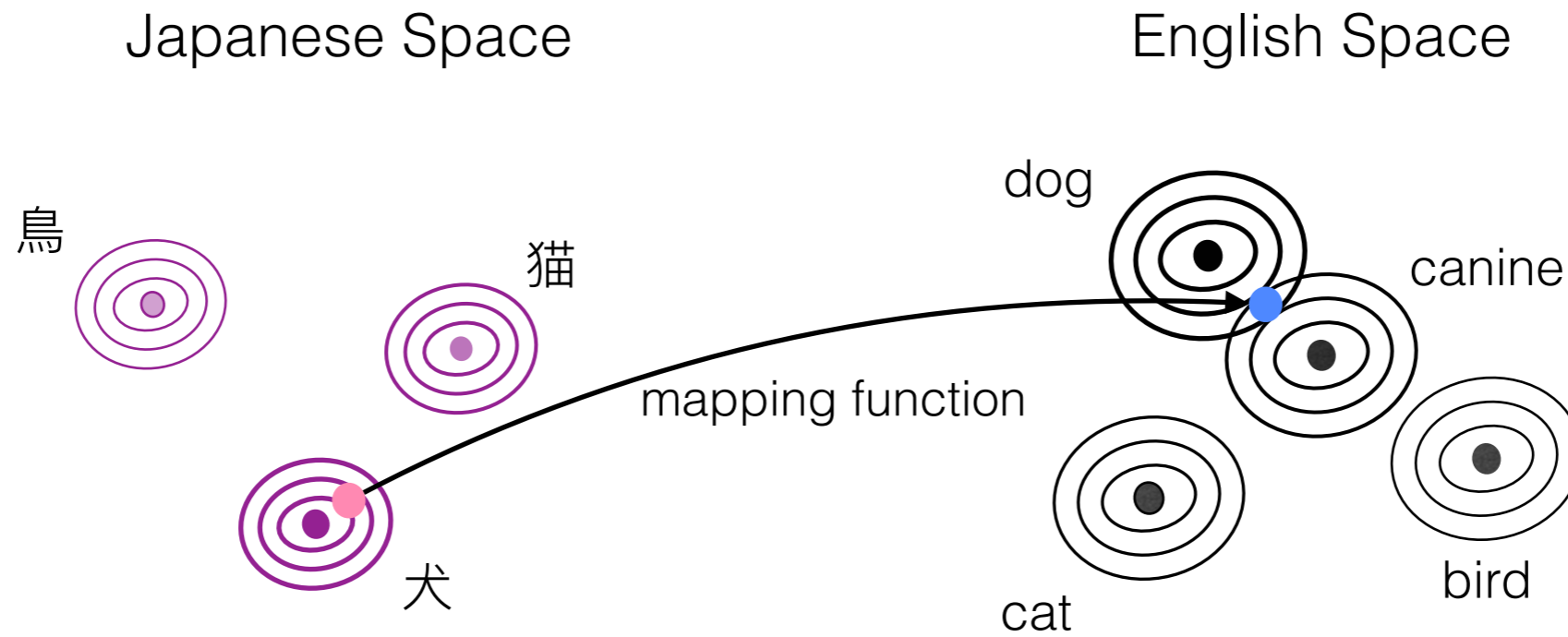
**Normalizing Flow:** A series of such invertible transformations  $f$



# DeMa-BWE: Preliminaries

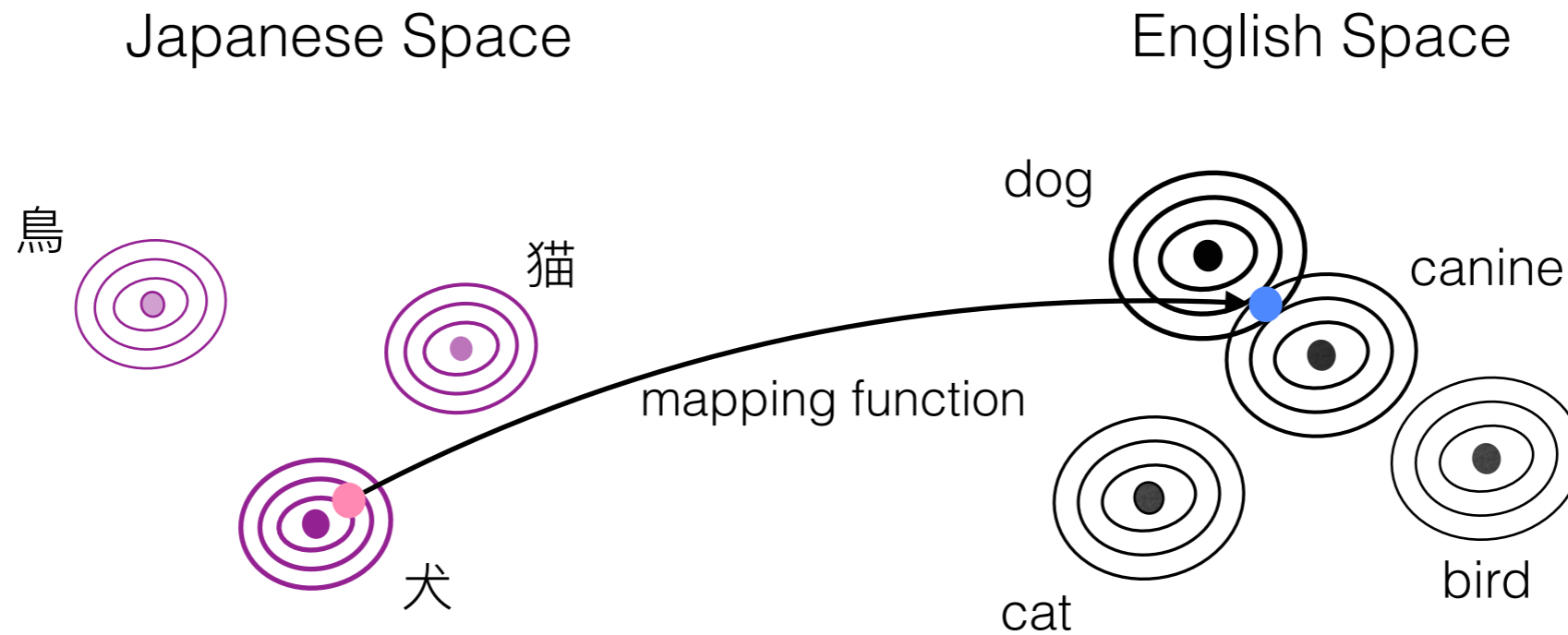


# DeMa-BWE: Preliminaries



Notations:

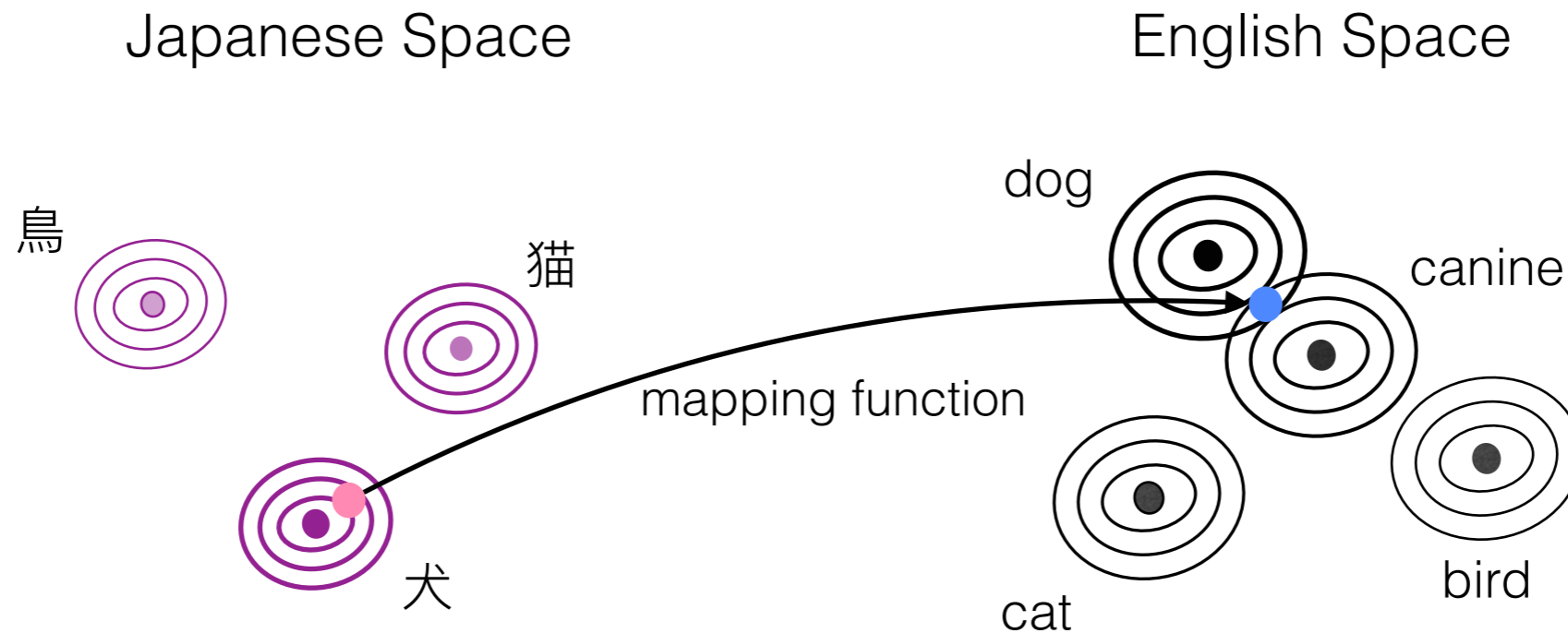
# DeMa-BWE: Preliminaries



Notations:

$\mathbf{x} \in \mathbb{R}^d$ ,  $\mathbf{y} \in \mathbb{R}^d$  : denote vectors in the src and tgt embedding space

# DeMa-BWE: Preliminaries

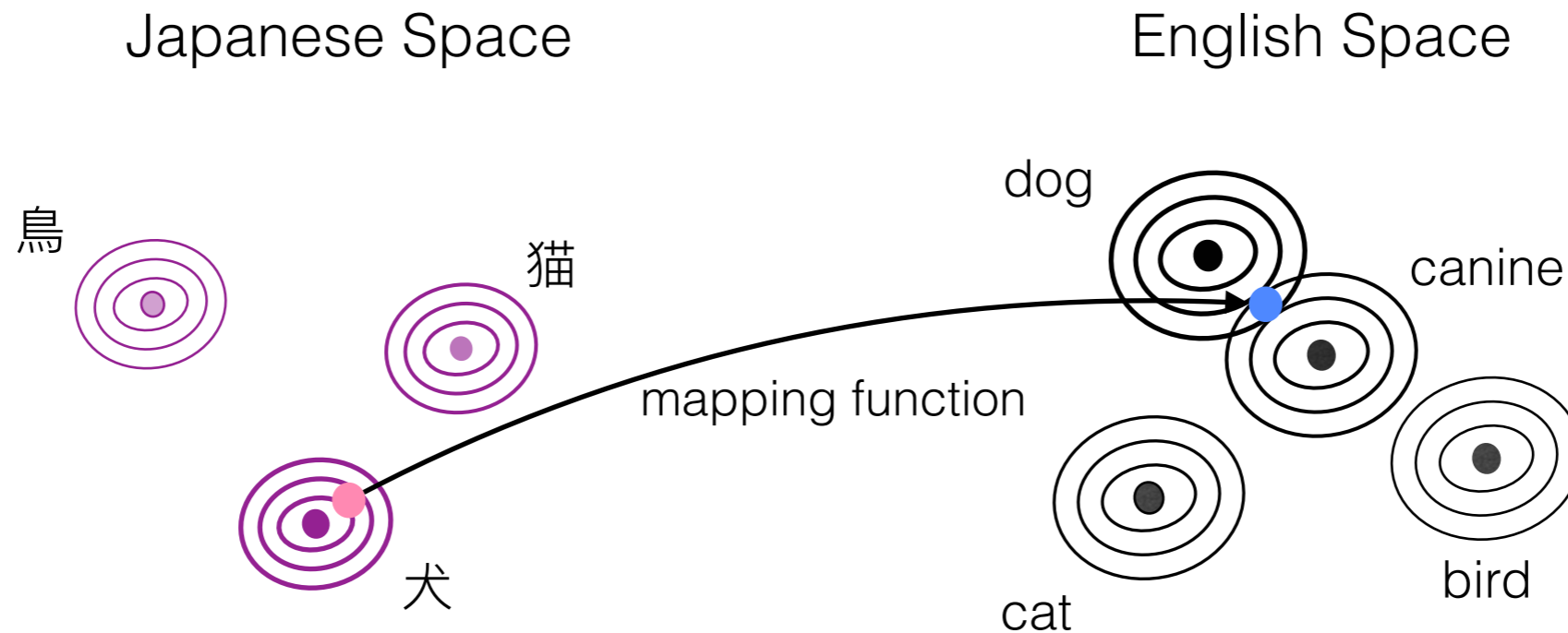


## Notations:

$\mathbf{x} \in \mathbb{R}^d$ ,  $\mathbf{y} \in \mathbb{R}^d$  : denote vectors in the src and tgt embedding space

$x_i$ ,  $y_j$  : denote an actual word in src and tgt vocabularies

# DeMa-BWE: Preliminaries



## Notations:

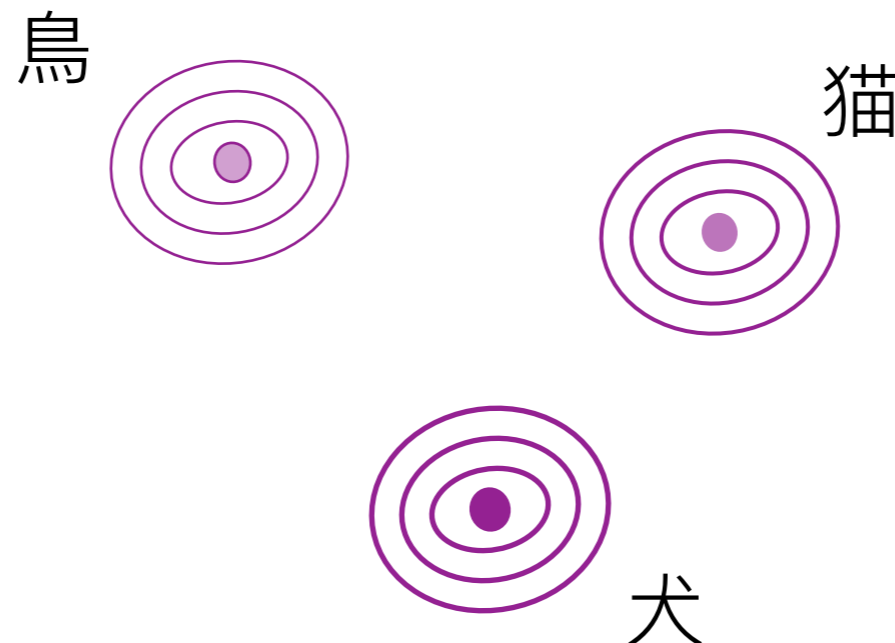
$\mathbf{x} \in \mathbb{R}^d$ ,  $\mathbf{y} \in \mathbb{R}^d$  : denote vectors in the src and tgt embedding space

$x_i$ ,  $y_j$  : denote an actual word in src and tgt vocabularies

$f_{xy}$ ,  $f_{yx}$  : denote src->tgt, and tgt-src mapping functions

# Prior Distribution

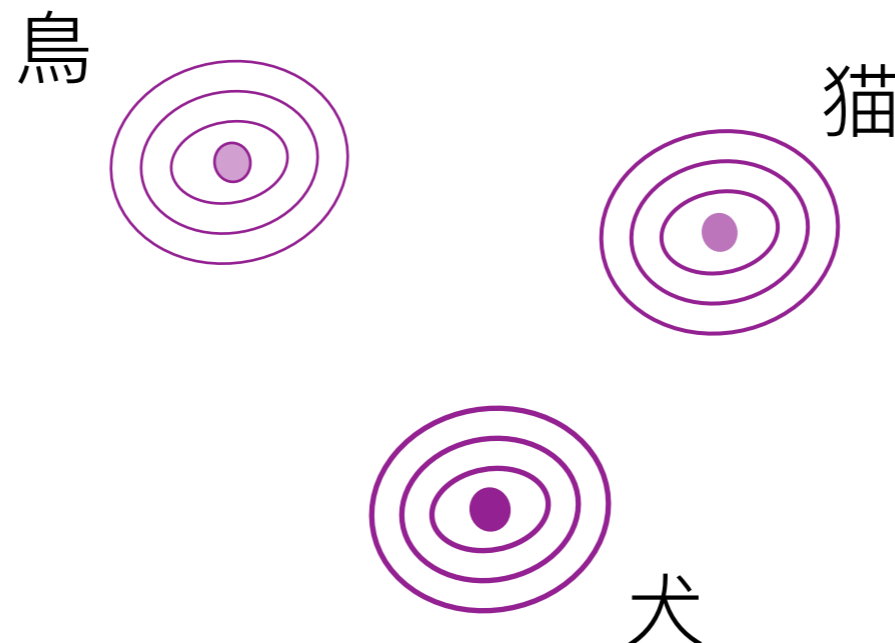
- Assumption on the monolingual word embedding space: Gaussian mixture model



# Prior Distribution

- Assumption on the monolingual word embedding space: Gaussian mixture model

$$p(\mathbf{x}) = \sum_{i \in \{1, \dots, N_x\}} \pi(x_i) \tilde{p}(\mathbf{x} | x_i)$$

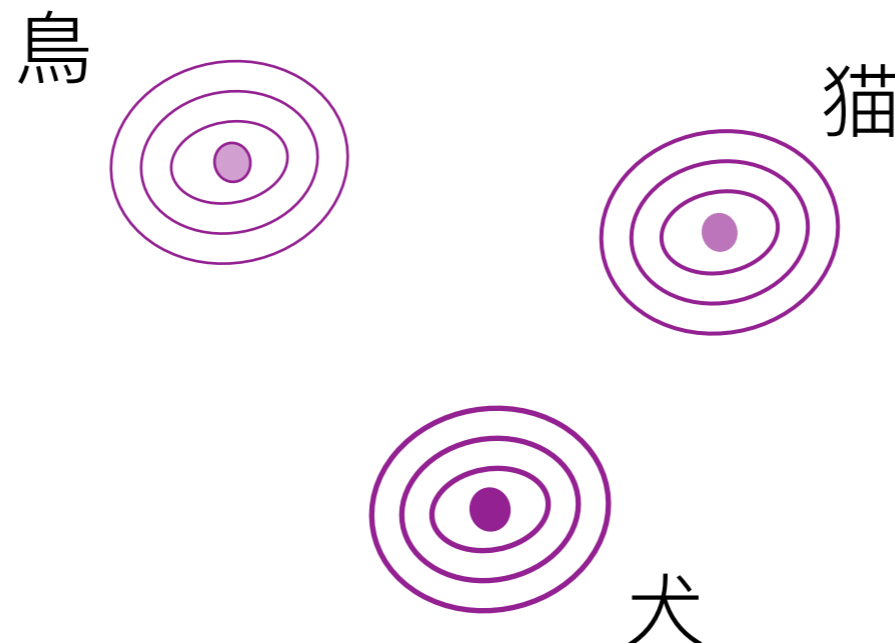


# Prior Distribution

- Assumption on the monolingual word embedding space: Gaussian mixture model

$$p(\mathbf{x}) = \sum_{i \in \{1, \dots, N_x\}} \pi(x_i) \tilde{p}(\mathbf{x} | x_i)$$

$$\tilde{p}(\mathbf{x} | x_i) = \mathcal{N}(\mathbf{x} | \mathbf{x}_i, \sigma_x^2 \mathbf{I})$$





# DeMa-BWE: Density Matching

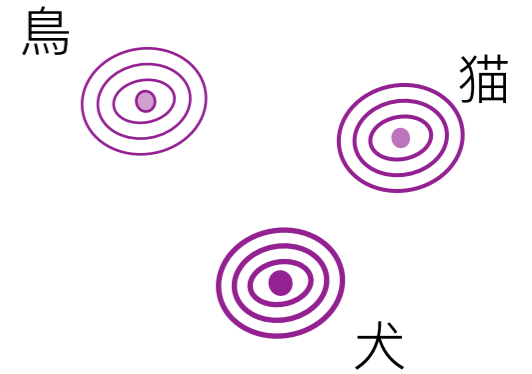
# DeMa-BWE: Density Matching

- Sampling a continuous vector from the GMM

# DeMa-BWE: Density Matching

- Sampling a continuous vector from the GMM

$$x_i \sim \pi(x_i) \quad \mathbf{x} \sim \tilde{p}(\mathbf{x}|x_i)$$

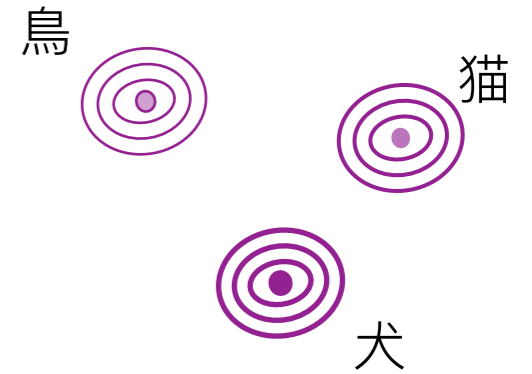


# DeMa-BWE: Density Matching

- Sampling a continuous vector from the GMM

$$x_i \sim \pi(x_i) \quad \mathbf{x} \sim \tilde{p}(\mathbf{x}|x_i)$$

- Apply the mapping function  $f_{xy}$  to obtain the transformed vector in the target space.



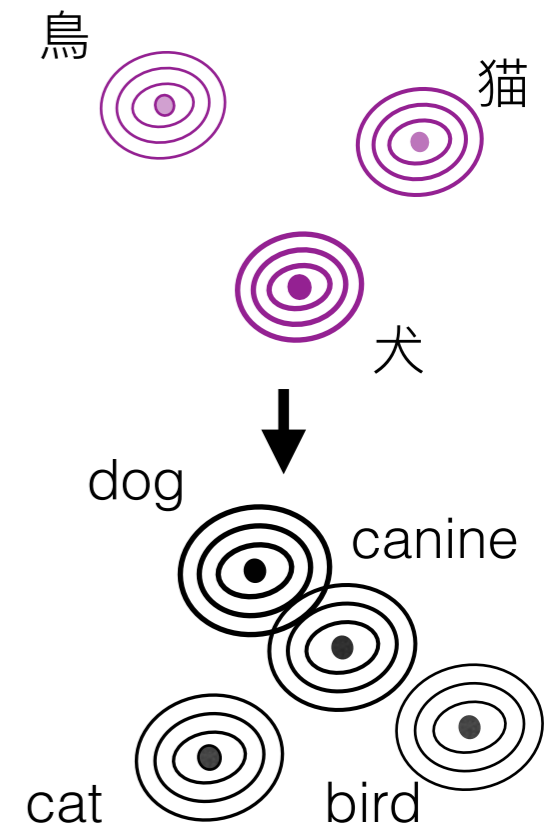
# DeMa-BWE: Density Matching

- Sampling a continuous vector from the GMM

$$x_i \sim \pi(x_i) \quad \mathbf{x} \sim \tilde{p}(\mathbf{x}|x_i)$$

- Apply the mapping function  $f_{xy}$  to obtain the transformed vector in the target space.

$$f_{xy}(\cdot) = \mathbf{W}_{xy} \cdot$$



# DeMa-BWE: Density Matching

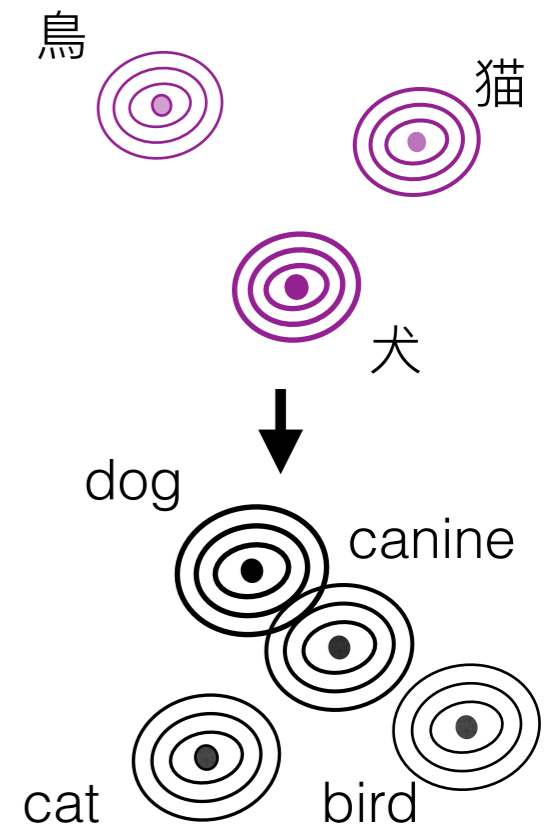
- Sampling a continuous vector from the GMM

$$x_i \sim \pi(x_i) \quad \mathbf{x} \sim \tilde{p}(\mathbf{x}|x_i)$$

- Apply the mapping function  $f_{xy}$  to obtain the transformed vector in the target space.

$$f_{xy}(\cdot) = \mathbf{W}_{xy} \cdot$$

- Computing the density of  $x$  in the mapped target space



# DeMa-BWE: Density Matching

- Sampling a continuous vector from the GMM

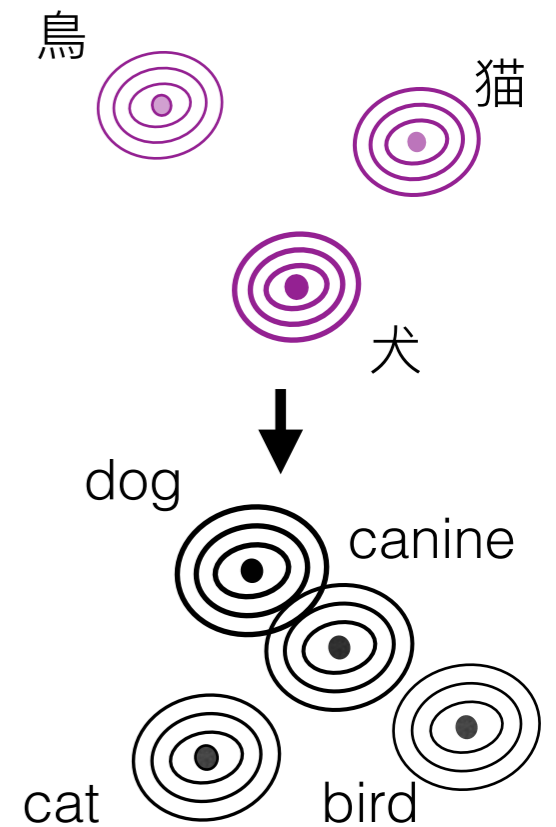
$$x_i \sim \pi(x_i) \quad \mathbf{x} \sim \tilde{p}(\mathbf{x}|x_i)$$

- Apply the mapping function  $f_{xy}$  to obtain the transformed vector in the target space.

$$f_{xy}(\cdot) = \mathbf{W}_{xy} \cdot$$

- Computing the density of  $\mathbf{x}$  in the mapped target space

$$\log p(\mathbf{x}; \mathbf{W}_{xy}) = \log p(\mathbf{y}) + \log |\det(\mathbf{W}_{xy})|$$



# DeMa-BWE: Density Matching

- Sampling a continuous vector from the GMM

$$x_i \sim \pi(x_i) \quad \mathbf{x} \sim \tilde{p}(\mathbf{x}|x_i)$$

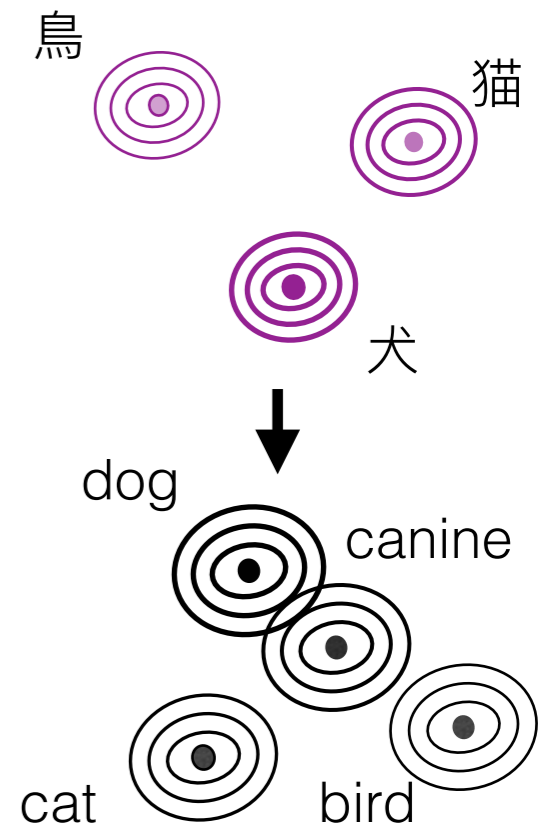
- Apply the mapping function  $f_{xy}$  to obtain the transformed vector in the target space.

$$f_{xy}(\cdot) = \mathbf{W}_{xy} \cdot$$

- Computing the density of  $\mathbf{x}$  in the mapped target space

$$\log p(\mathbf{x}; \mathbf{W}_{xy}) = \log p(\mathbf{y}) + \log |\det(\mathbf{W}_{xy})|$$

$$\text{minimize: } \text{KL}(p(\mathbf{x}) || p(\mathbf{x}; \mathbf{W}_{xy}))$$





# DeMa-BWE: Density Matching

- Sampling a continuous vector from the GMM

$$x_i \sim \pi(x_i) \quad \mathbf{x} \sim \tilde{p}(\mathbf{x}|x_i)$$

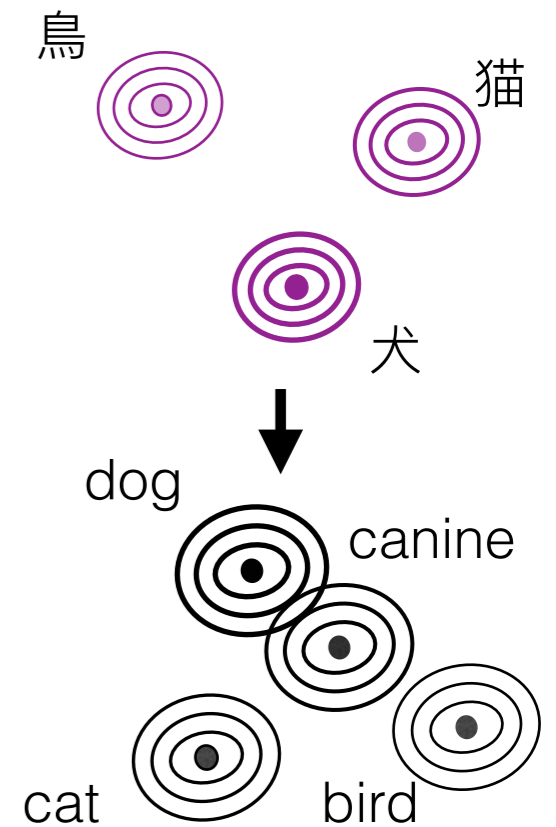
- Apply the mapping function  $f_{xy}$  to obtain the transformed vector in the target space.

$$f_{xy}(\cdot) = \mathbf{W}_{xy} \cdot$$

- Computing the density of  $\mathbf{x}$  in the mapped target space

$$\log p(\mathbf{x}; \mathbf{W}_{xy}) = \log p(\mathbf{y}) + \log |\det(\mathbf{W}_{xy})|$$

- Objective: minimize:  $\text{KL}(p(\mathbf{x}) || p(\mathbf{x}; \mathbf{W}_{xy}))$



# DeMa-BWE: Density Matching

- Sampling a continuous vector from the GMM

$$x_i \sim \pi(x_i) \quad \mathbf{x} \sim \tilde{p}(\mathbf{x}|x_i)$$

- Apply the mapping function  $f_{xy}$  to obtain the transformed vector in the target space.

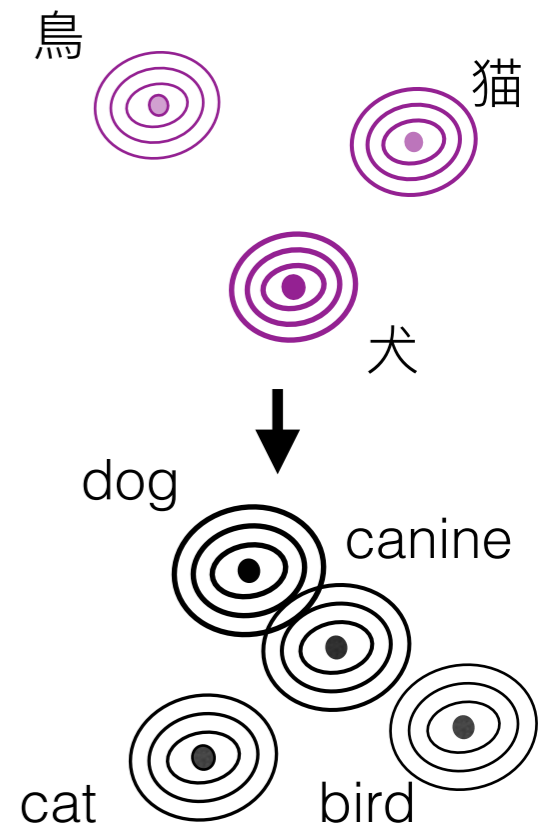
$$f_{xy}(\cdot) = \mathbf{W}_{xy} \cdot$$

- Computing the density of  $\mathbf{x}$  in the mapped target space

$$\log p(\mathbf{x}; \mathbf{W}_{xy}) = \log p(\mathbf{y}) + \log |\det(\mathbf{W}_{xy})|$$

- Objective: minimize:  $\text{KL}(p(\mathbf{x}) || p(\mathbf{x}; \mathbf{W}_{xy}))$

$$\mathcal{L}_{xy} = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\log p(\mathbf{y}) + \log |\det(\mathbf{W}_{xy})|]$$



# Method Details

# Method Details

- **Weak Orthogonality Constraint:** Try to make sure that the transformation is close to orthogonal

# Method Details

- **Weak Orthogonality Constraint:** Try to make sure that the transformation is close to orthogonal

$$\mathcal{L}_{bt} = \mathbb{E}_{x_i \sim \pi(x_i), \mathbf{x} \sim \tilde{p}(\mathbf{x}|x_i)} [g(\mathbf{W}_{yx} \mathbf{W}_{xy} \mathbf{x}, \mathbf{x})] + \mathbb{E}_{y_j \sim \pi(y_j), \mathbf{y} \sim \tilde{p}(\mathbf{y}|y_j)} [g(\mathbf{W}_{xy} \mathbf{W}_{yx} \mathbf{y}, \mathbf{y})]$$

# Method Details

- **Weak Orthogonality Constraint:** Try to make sure that the transformation is close to orthogonal

$$\mathcal{L}_{bt} = \mathbb{E}_{x_i \sim \pi(x_i), \mathbf{x} \sim \tilde{p}(\mathbf{x}|x_i)} [g(\mathbf{W}_{yx} \mathbf{W}_{xy} \mathbf{x}, \mathbf{x})] + \mathbb{E}_{y_j \sim \pi(y_j), \mathbf{y} \sim \tilde{p}(\mathbf{y}|y_j)} [g(\mathbf{W}_{xy} \mathbf{W}_{yx} \mathbf{y}, \mathbf{y})]$$

- **Weak Supervision w/ Identical Strings:** Take advantage of the fact that identical strings are usually the same word in both languages

# Method Details

- **Weak Orthogonality Constraint:** Try to make sure that the transformation is close to orthogonal

$$\mathcal{L}_{bt} = \mathbb{E}_{x_i \sim \pi(x_i), \mathbf{x} \sim \tilde{p}(\mathbf{x}|x_i)} [g(\mathbf{W}_{yx} \mathbf{W}_{xy} \mathbf{x}, \mathbf{x})] + \mathbb{E}_{y_j \sim \pi(y_j), \mathbf{y} \sim \tilde{p}(\mathbf{y}|y_j)} [g(\mathbf{W}_{xy} \mathbf{W}_{yx} \mathbf{y}, \mathbf{y})]$$

- **Weak Supervision w/ Identical Strings:** Take advantage of the fact that identical strings are usually the same word in both languages

$$\mathcal{L}_{sup} = \sum_{v \in \mathcal{W}_{id}} g(\mathbf{v}_x \mathbf{W}_{xy}^T, \mathbf{v}_y) + g(\mathbf{v}_y \mathbf{W}_{yx}^T, \mathbf{v}_x)$$

# Method Details

- **Weak Orthogonality Constraint:** Try to make sure that the transformation is close to orthogonal

$$\mathcal{L}_{bt} = \mathbb{E}_{x_i \sim \pi(x_i), \mathbf{x} \sim \tilde{p}(\mathbf{x}|x_i)} [g(\mathbf{W}_{yx} \mathbf{W}_{xy} \mathbf{x}, \mathbf{x})] + \mathbb{E}_{y_j \sim \pi(y_j), \mathbf{y} \sim \tilde{p}(\mathbf{y}|y_j)} [g(\mathbf{W}_{xy} \mathbf{W}_{yx} \mathbf{y}, \mathbf{y})]$$

- **Weak Supervision w/ Identical Strings:** Take advantage of the fact that identical strings are usually the same word in both languages

$$\mathcal{L}_{sup} = \sum_{v \in \mathcal{W}_{id}} g(\mathbf{v}_x \mathbf{W}_{xy}^T, \mathbf{v}_y) + g(\mathbf{v}_y \mathbf{W}_{yx}^T, \mathbf{v}_x)$$

- **Alignment Selection Methods:** Use cross-domain similarity local scaling (CSLS)



# Method Details

- **Weak Orthogonality Constraint:** Try to make sure that the transformation is close to orthogonal

$$\mathcal{L}_{bt} = \mathbb{E}_{x_i \sim \pi(x_i), \mathbf{x} \sim \tilde{p}(\mathbf{x}|x_i)} [g(\mathbf{W}_{yx} \mathbf{W}_{xy} \mathbf{x}, \mathbf{x})] + \mathbb{E}_{y_j \sim \pi(y_j), \mathbf{y} \sim \tilde{p}(\mathbf{y}|y_j)} [g(\mathbf{W}_{xy} \mathbf{W}_{yx} \mathbf{y}, \mathbf{y})]$$

- **Weak Supervision w/ Identical Strings:** Take advantage of the fact that identical strings are usually the same word in both languages

$$\mathcal{L}_{sup} = \sum_{v \in \mathcal{W}_{id}} g(\mathbf{v}_x \mathbf{W}_{xy}^T, \mathbf{v}_y) + g(\mathbf{v}_y \mathbf{W}_{yx}^T, \mathbf{v}_x)$$

- **Alignment Selection Methods:** Use cross-domain similarity local scaling (CSLS)

$$\text{CSLS}(\mathbf{x}', \mathbf{y}) = 2\cos(\mathbf{x}', \mathbf{y}) - r_T(\mathbf{x}') - r_S(\mathbf{y})$$

# Experiments

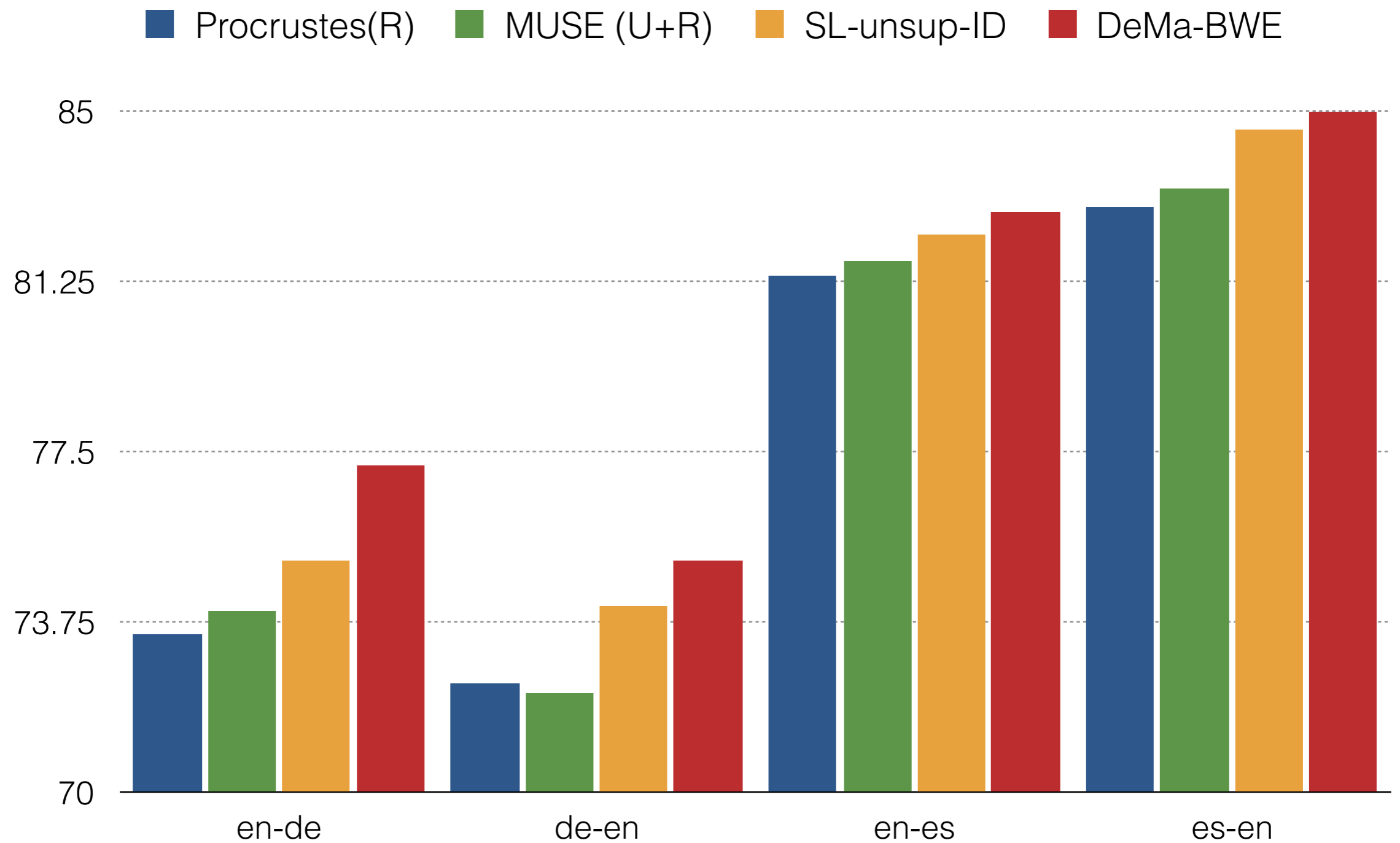
- **Dataset and Tasks**

- Bilingual Lexicon Induction Task: MUSE dataset (Conneau et al., 2017)
- Cross-lingual Word Similarity Task: SemEval 2017

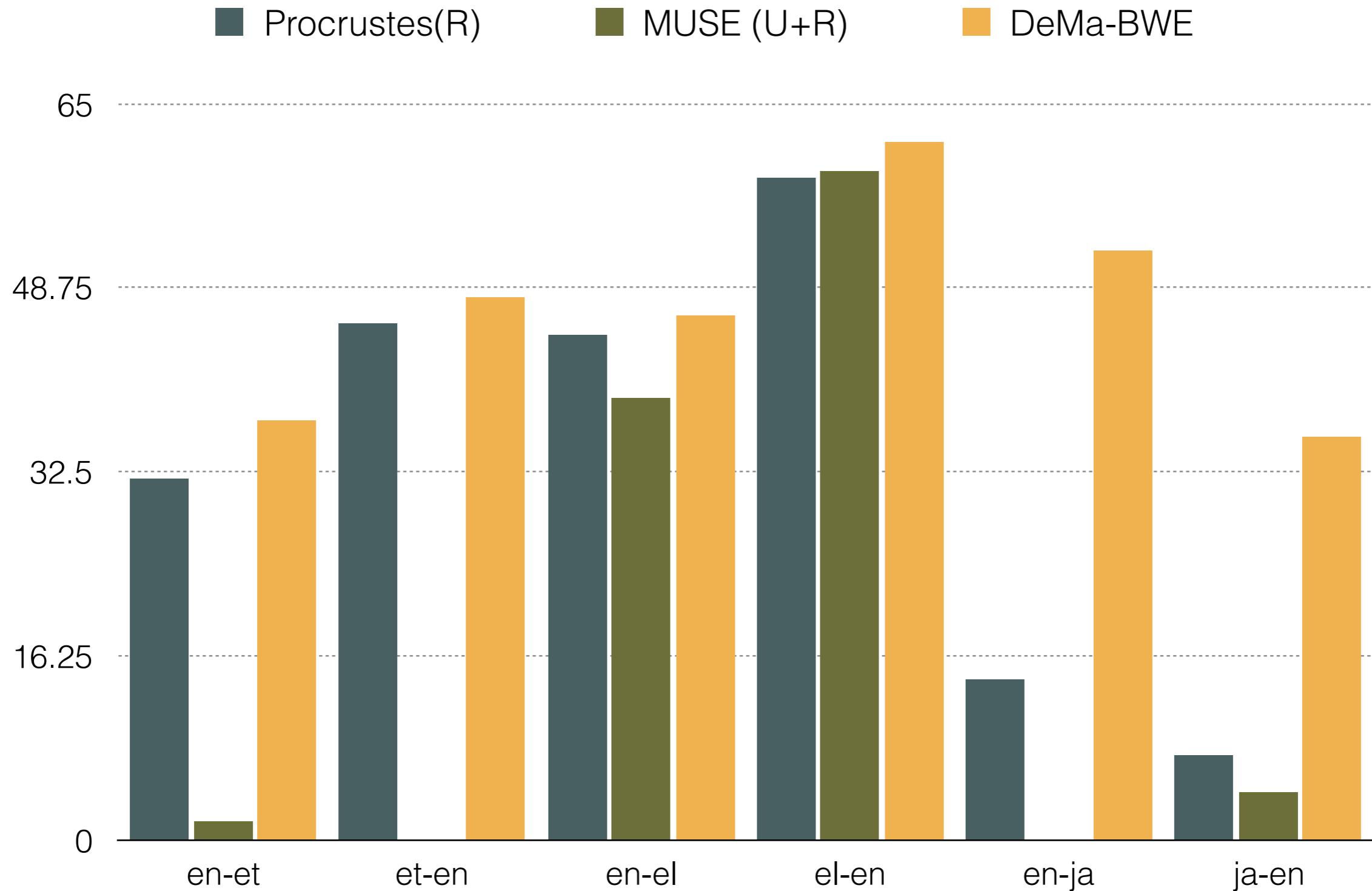
- **Languages**

- Baseline languages: en - es, de, fr, ru, zh, ja
- Morphologically rich languages: en - et, fi, el, hu, pl, tr

# Main Results on BLI (close languages)



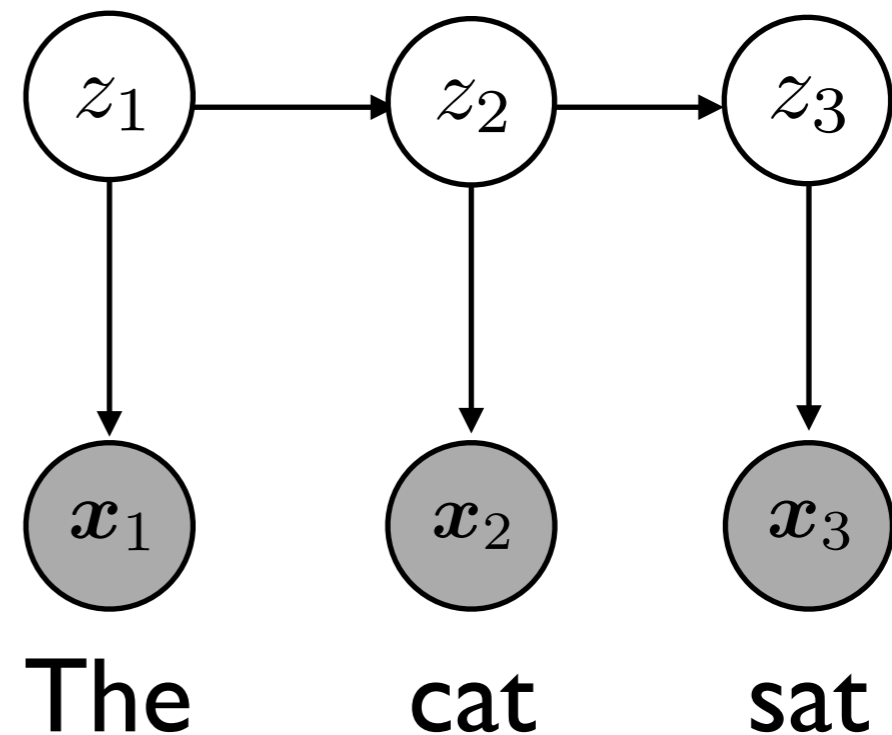
# Main Results on BLI (distant languages)



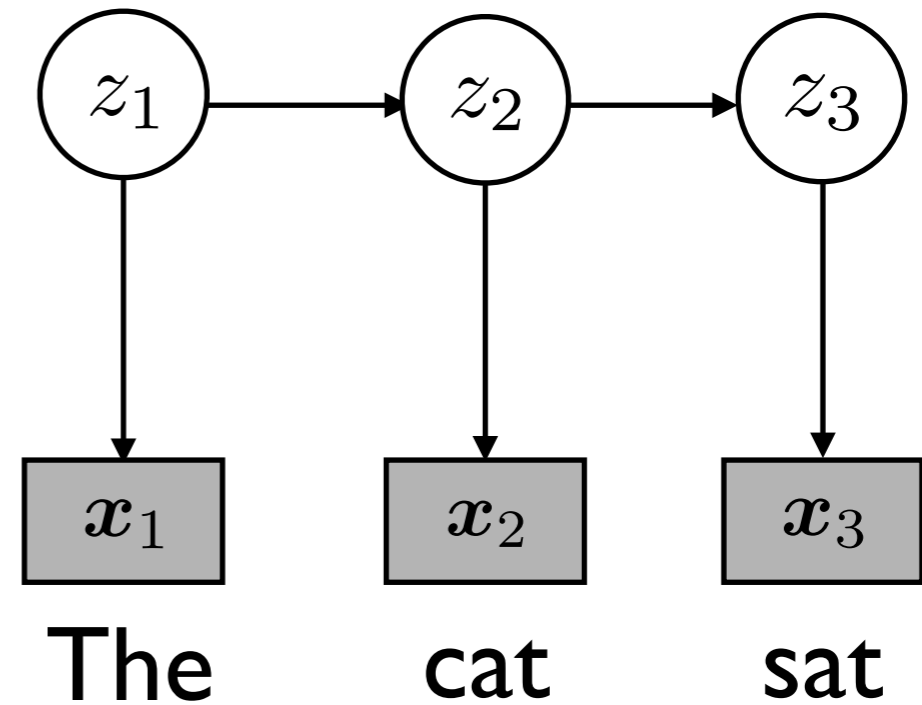
# Unsupervised Learning of Syntactic Structure w/ Invertible Neural Projections

Junxian He, Graham Neubig, Taylor Berg-Kirkpatrick  
(EMNLP 2018)

# HMM for Part-of-Speech Induction



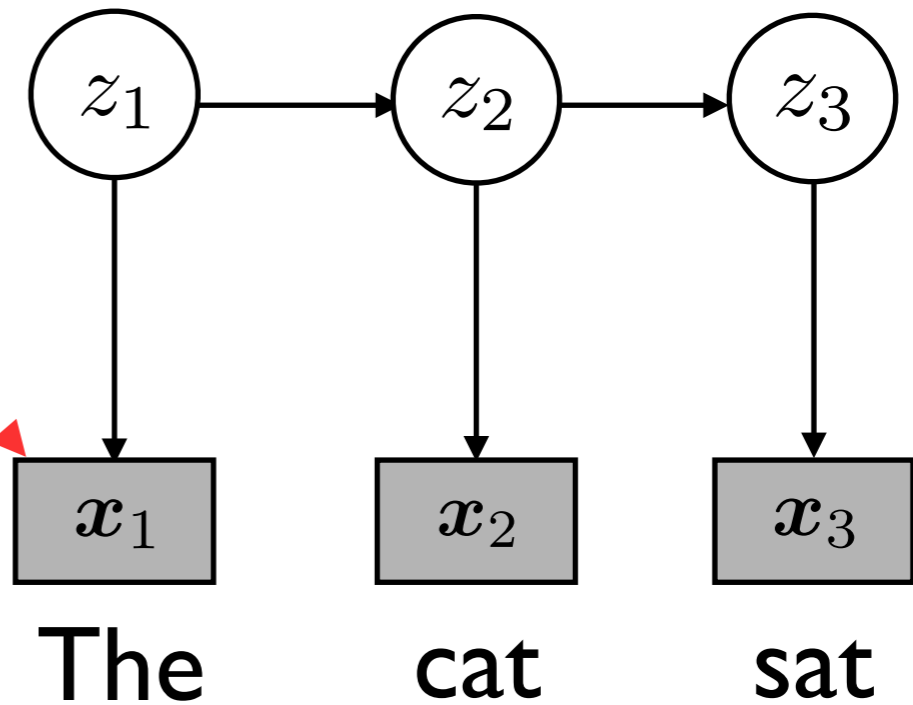
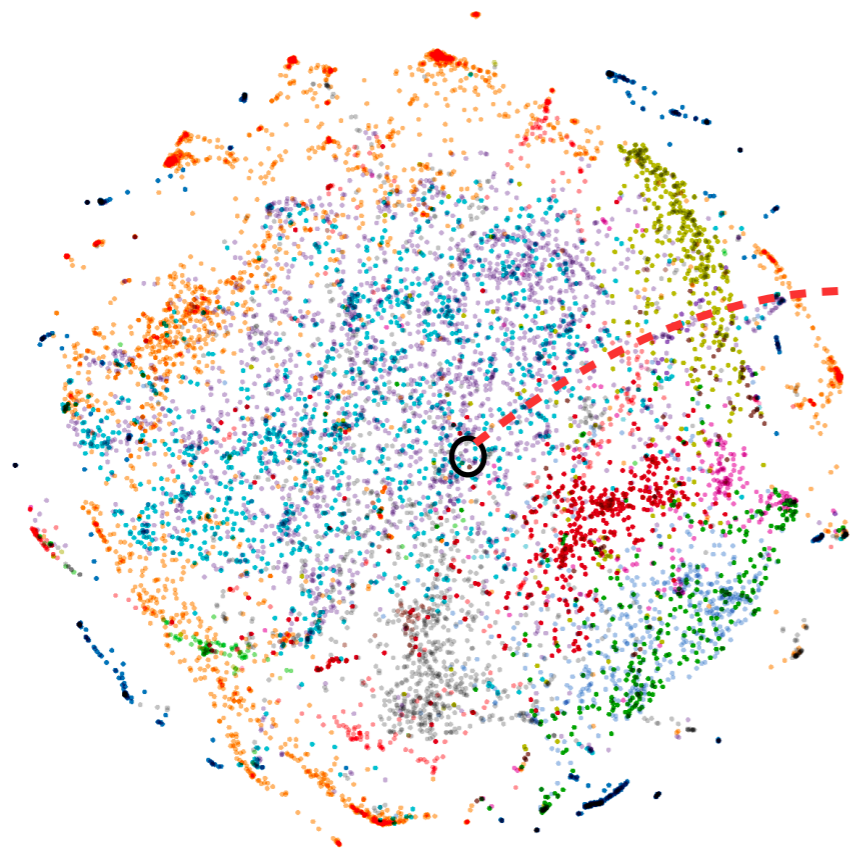
# Gaussian HMM for POS Induction



$$\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i})$$

[Lin et al. 2015]

# Gaussian HMM for POS Induction

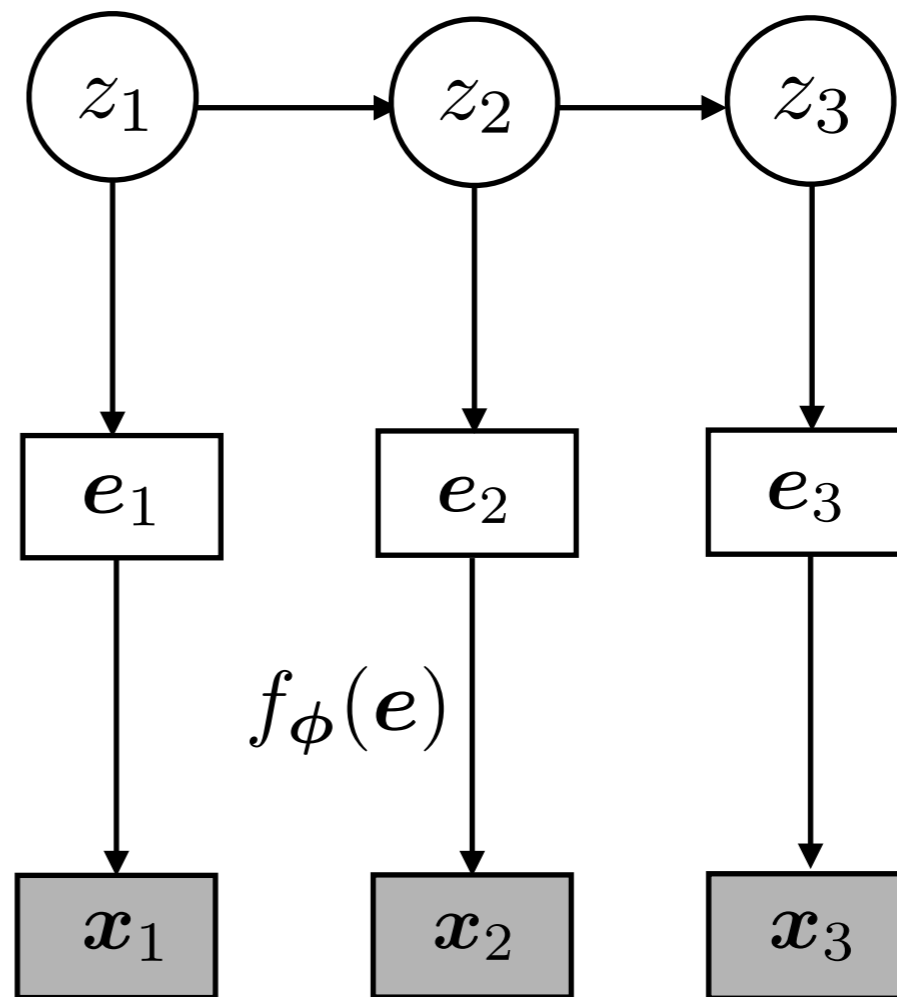


$$\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i})$$

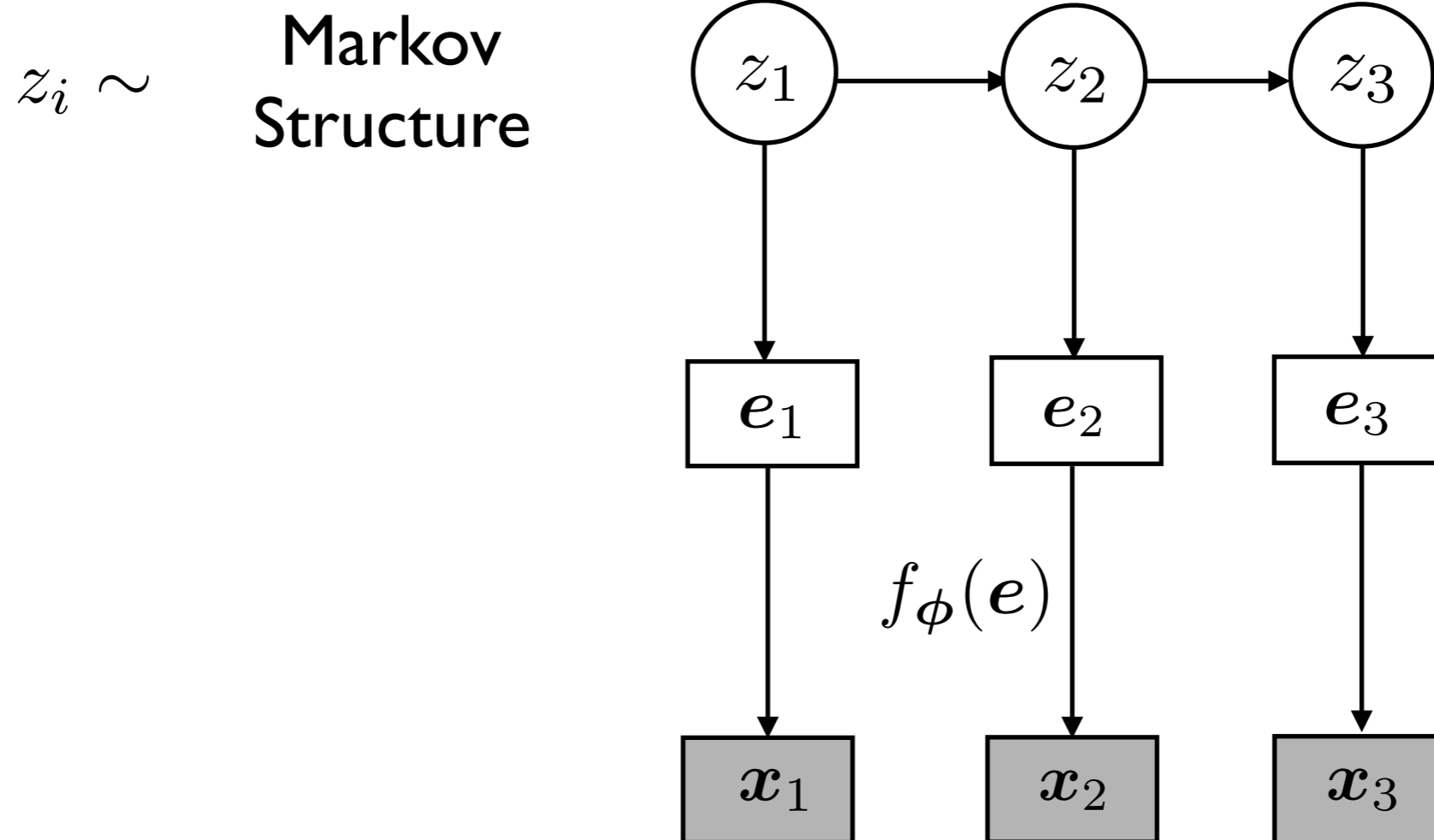
[Lin et al. 2015]



# Latent Embeddings w/ Neural Projection



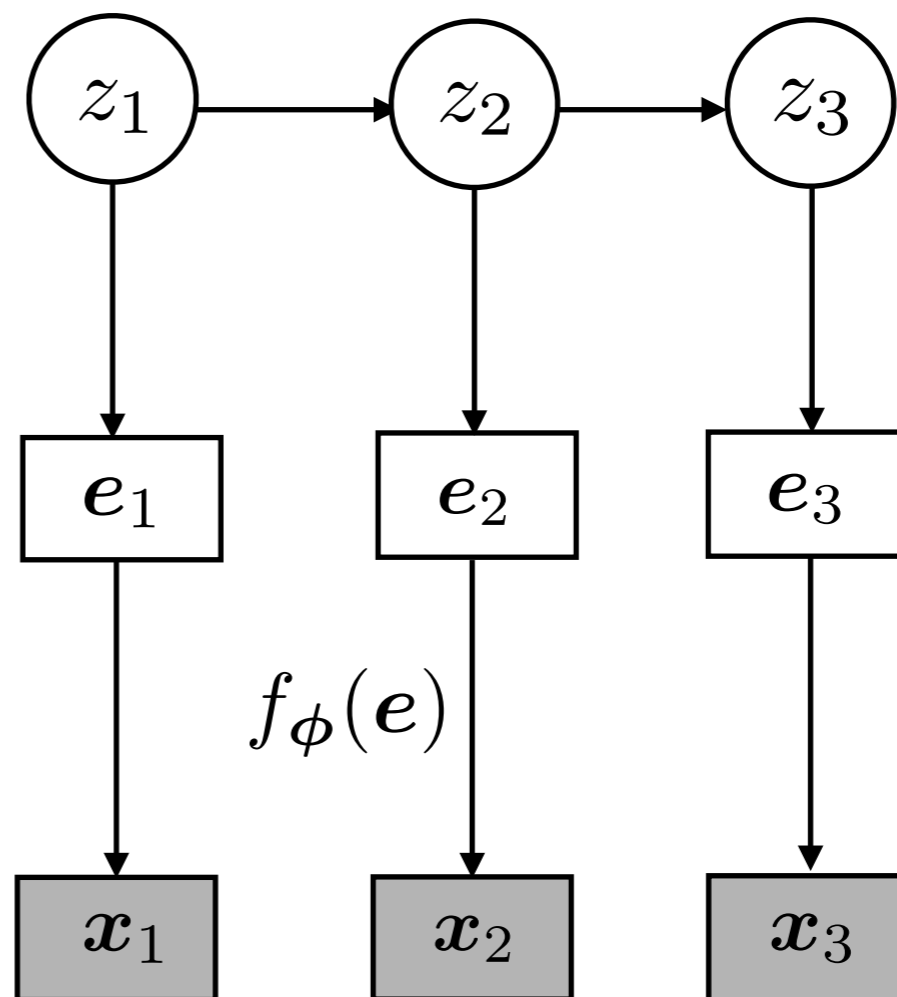
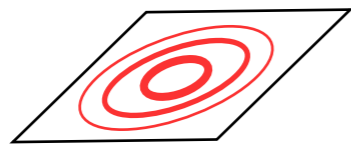
# Latent Embeddings w/ Neural Projection



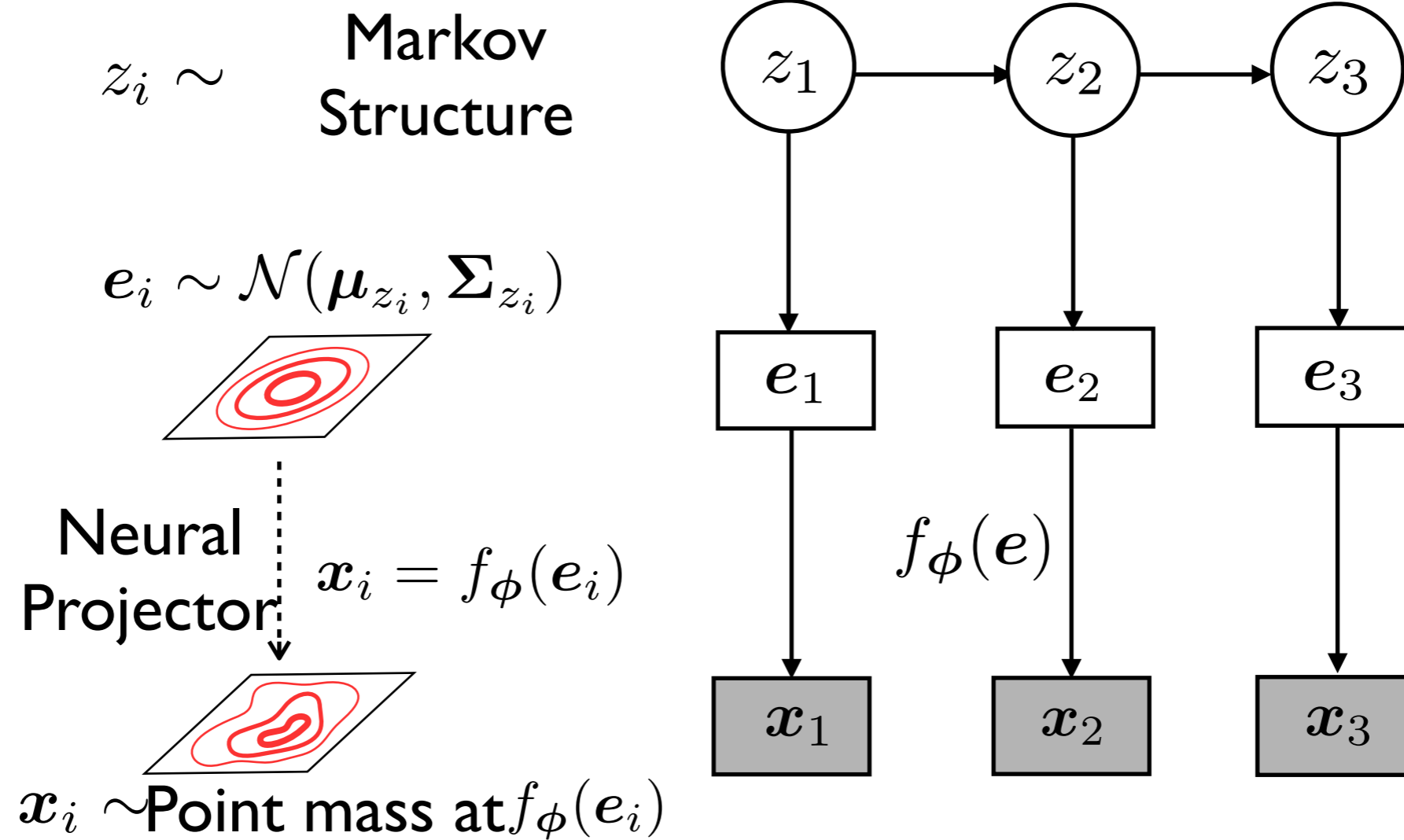
# Latent Embeddings w/ Neural Projection

$z_i \sim$  Markov Structure

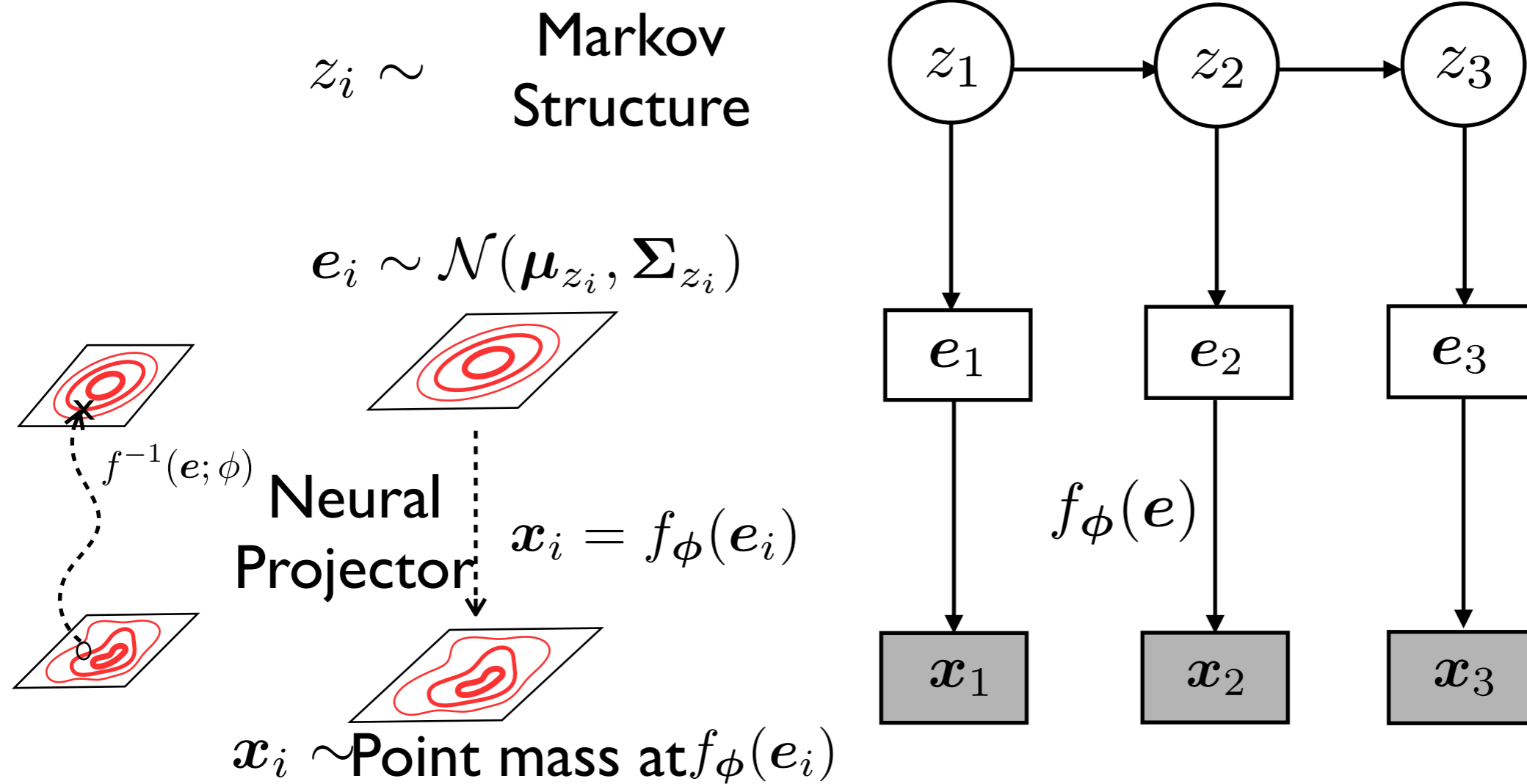
$$e_i \sim \mathcal{N}(\mu_{z_i}, \Sigma_{z_i})$$



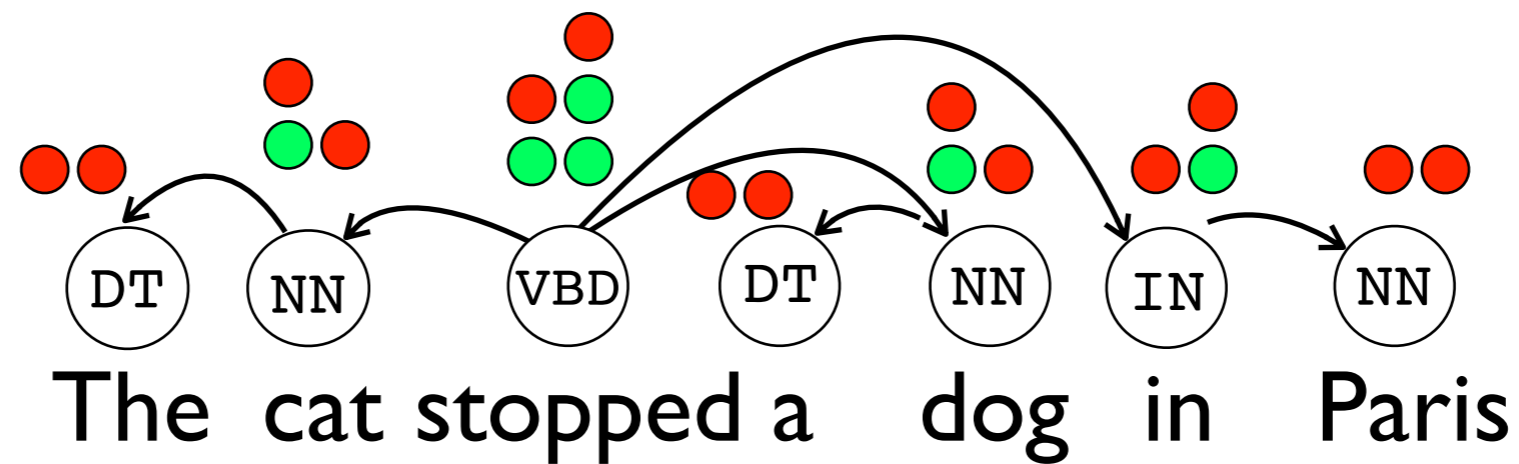
# Latent Embeddings w/ Neural Projection



# Latent Embeddings w/ Neural Projection

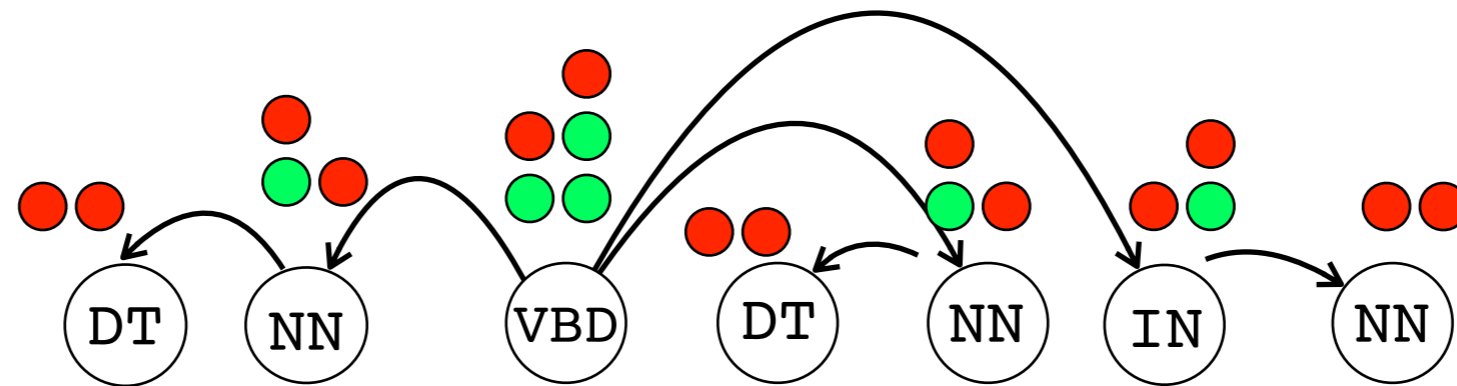


# Dependency Model with Valence

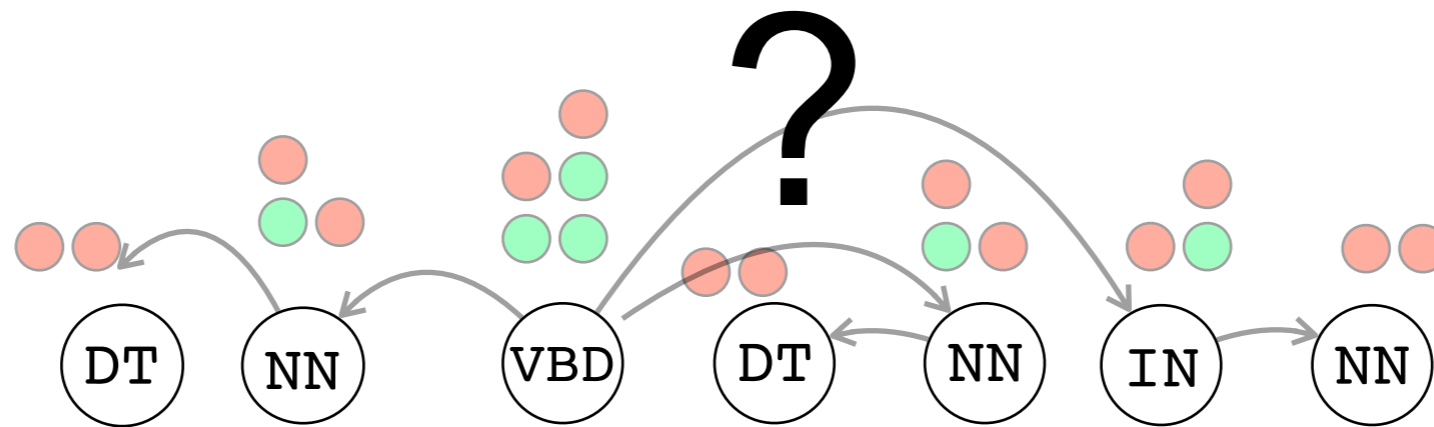


[Klein and Manning 2004]

# Dependency Model with Valence

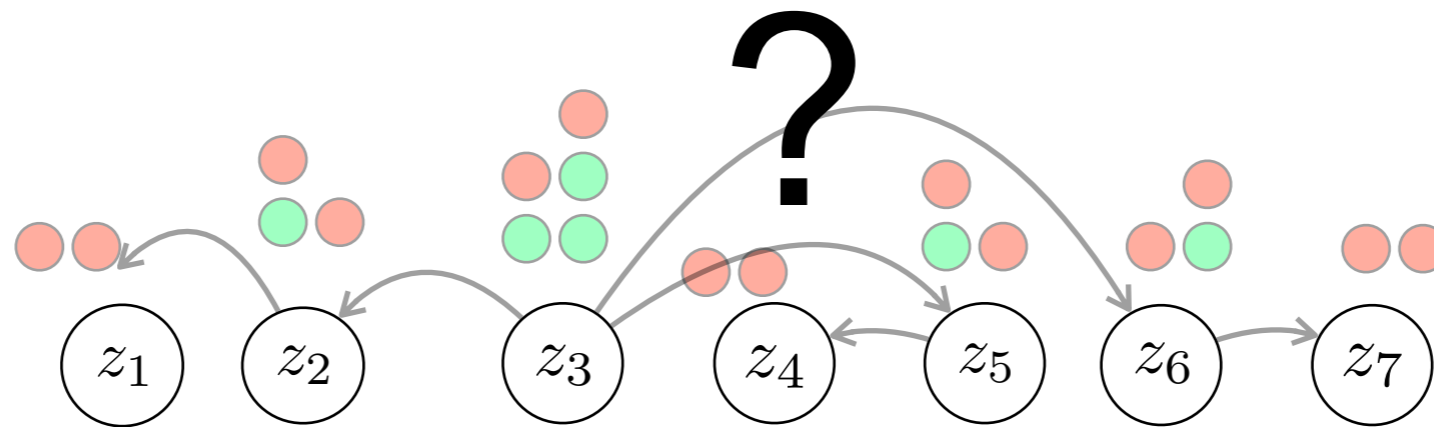


[Klein and Manning 2004]

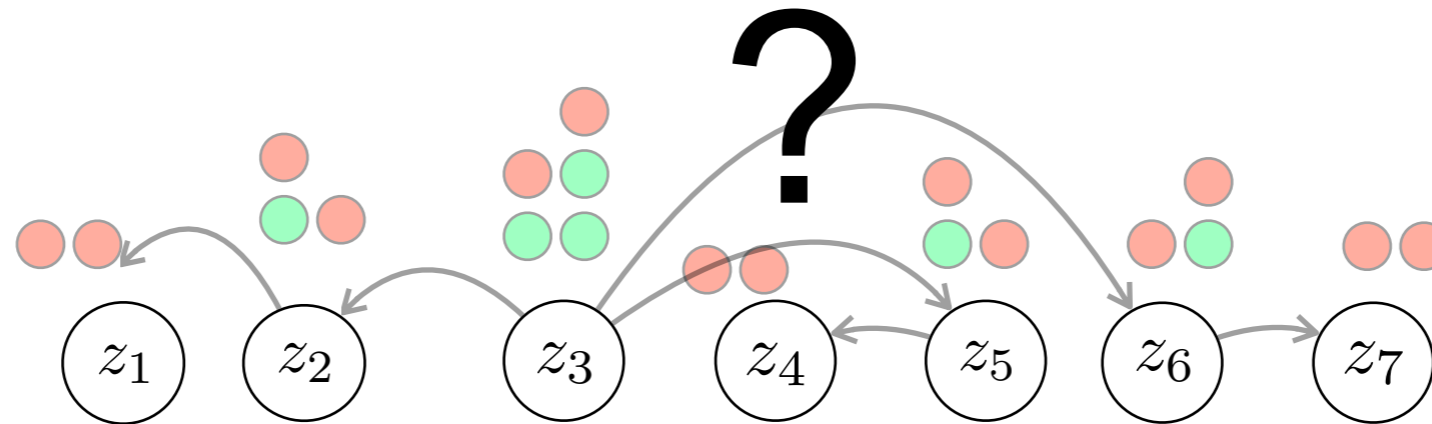
Dependency Parse Induction  
from POS



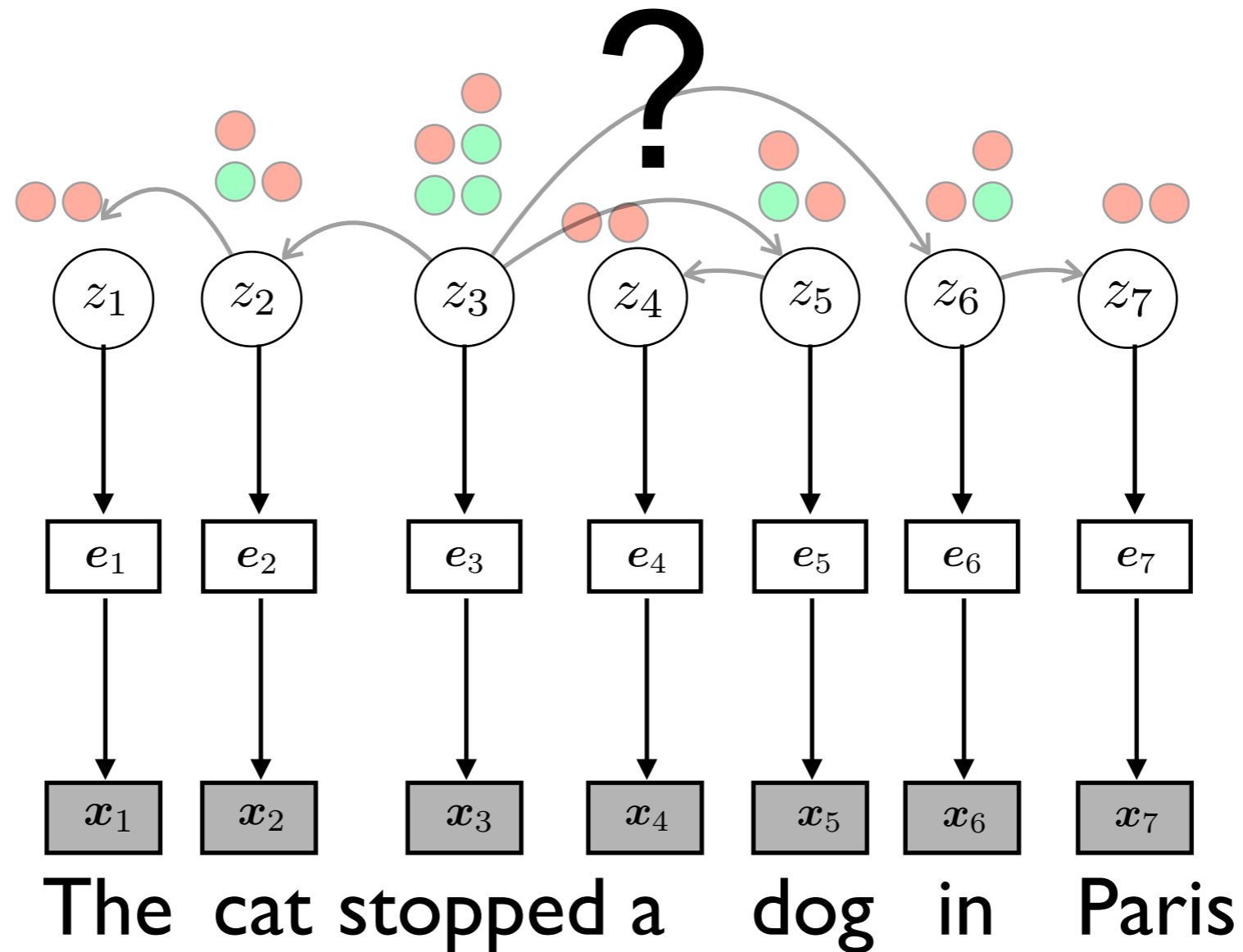
# Grammar Induction from Raw Text



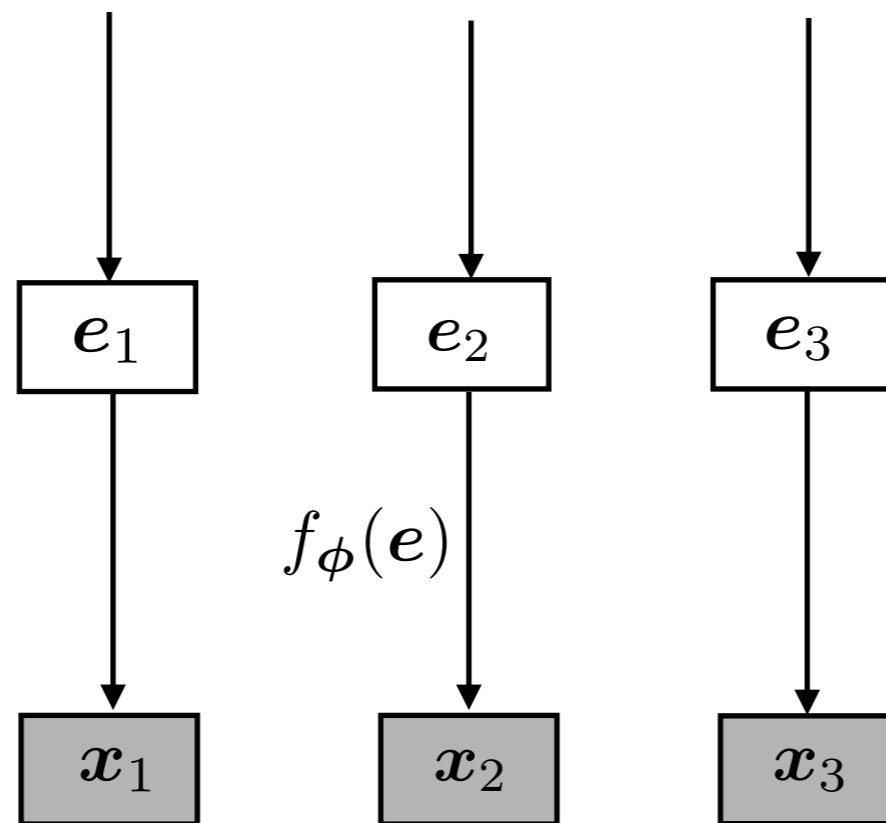
# Grammar Induction from Raw Text



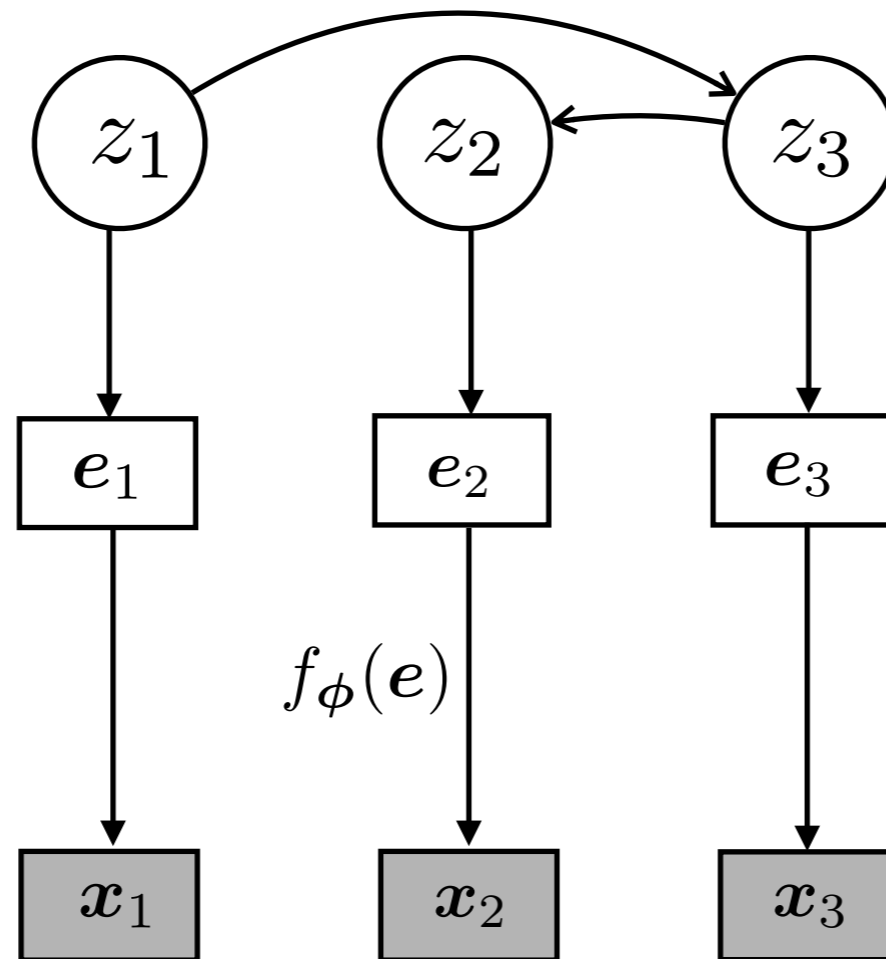
# Grammar Induction from Raw Text



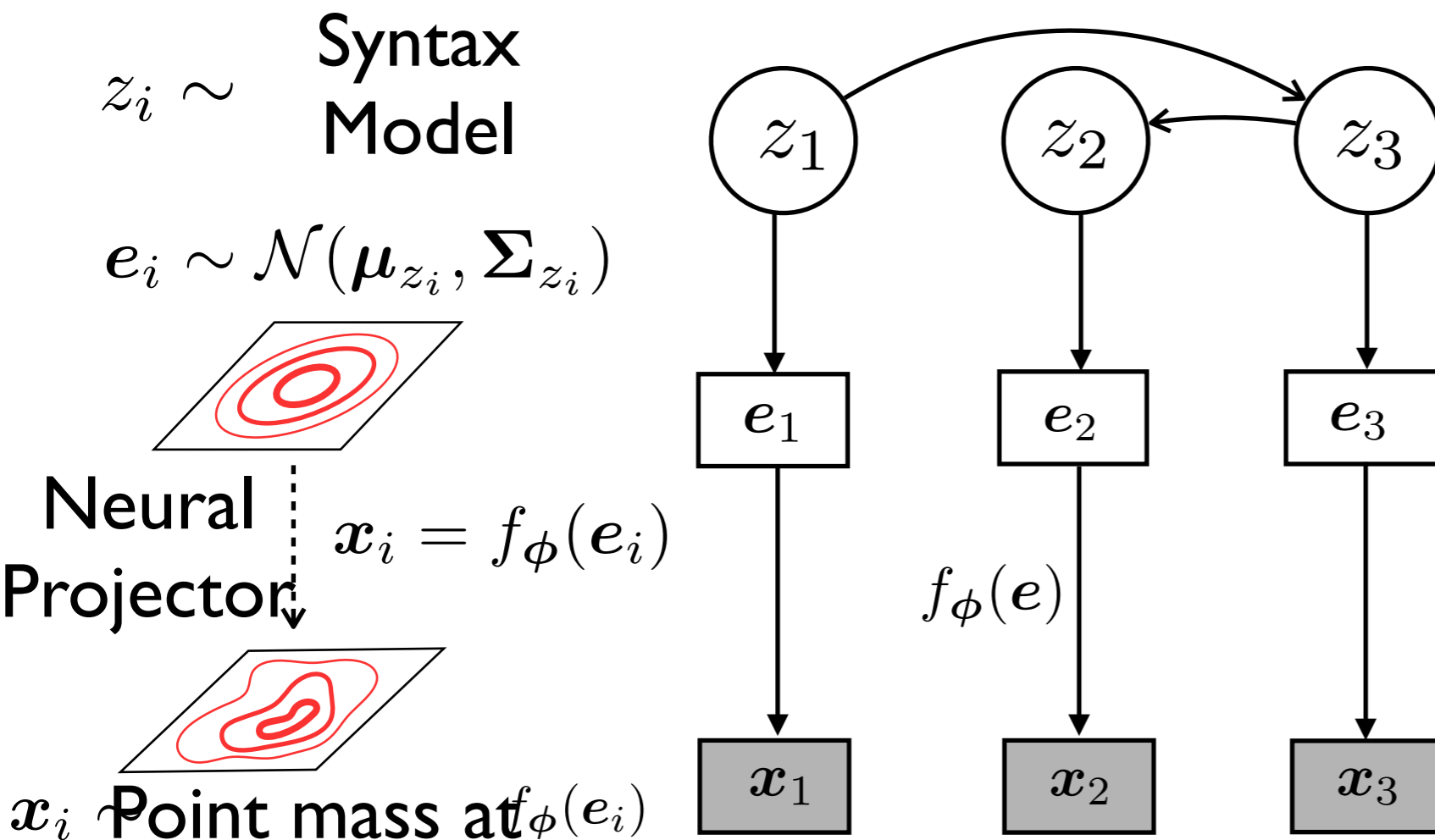
# Latent Embeddings w/ Neural Projection



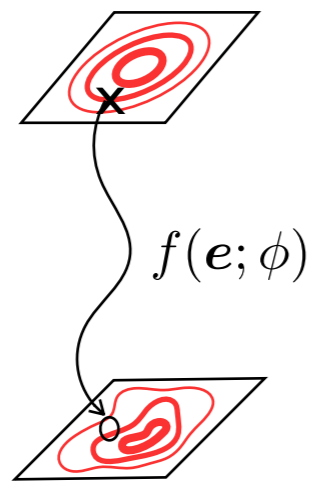
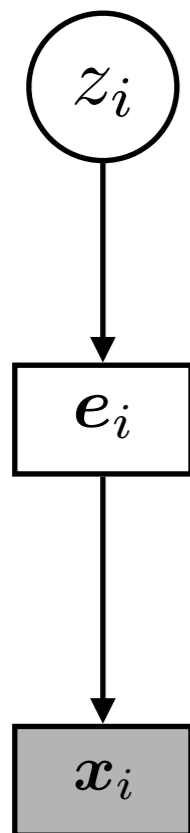
# Latent Embeddings w/ Neural Projection



# Latent Embeddings w/ Neural Projection

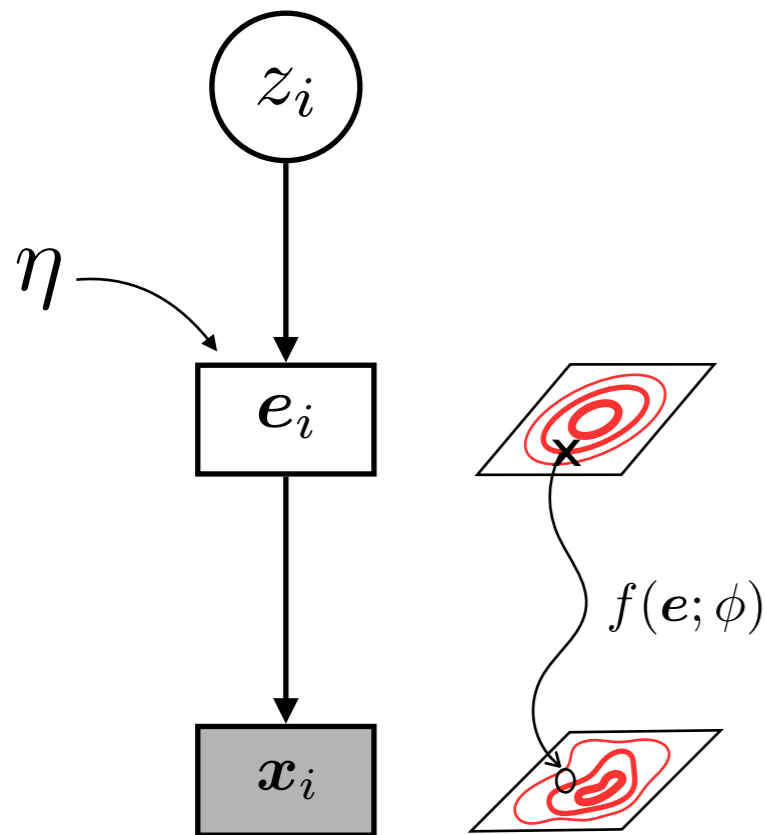


# Learning and Inference



$$p(\mathbf{x}_i | z_i; \eta, \phi)$$

# Learning and Inference

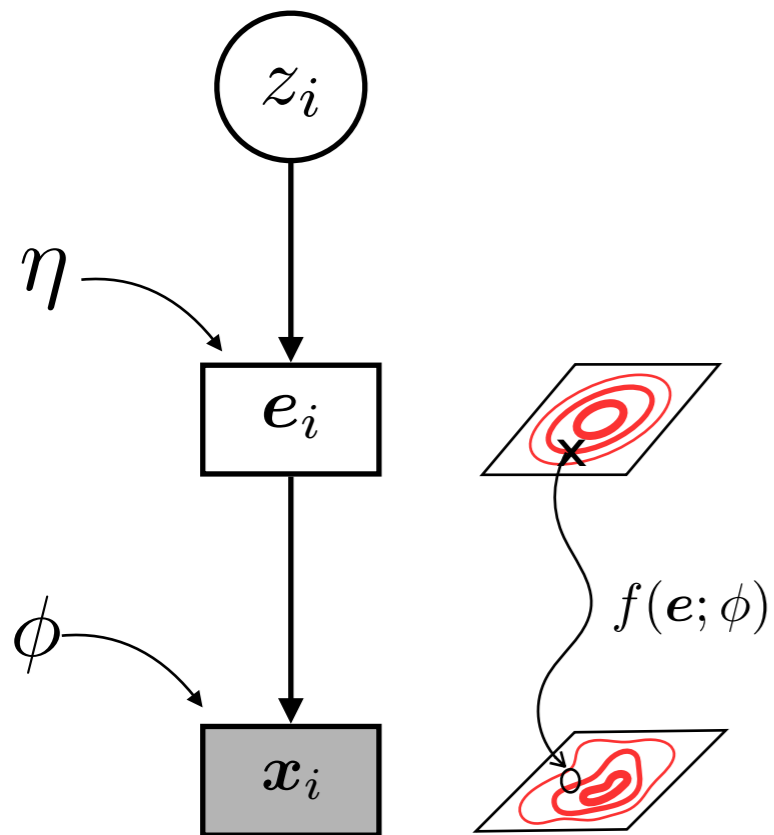


$$p(\mathbf{x}_i | z_i; \eta, \phi)$$

Gaussian embedding parameters



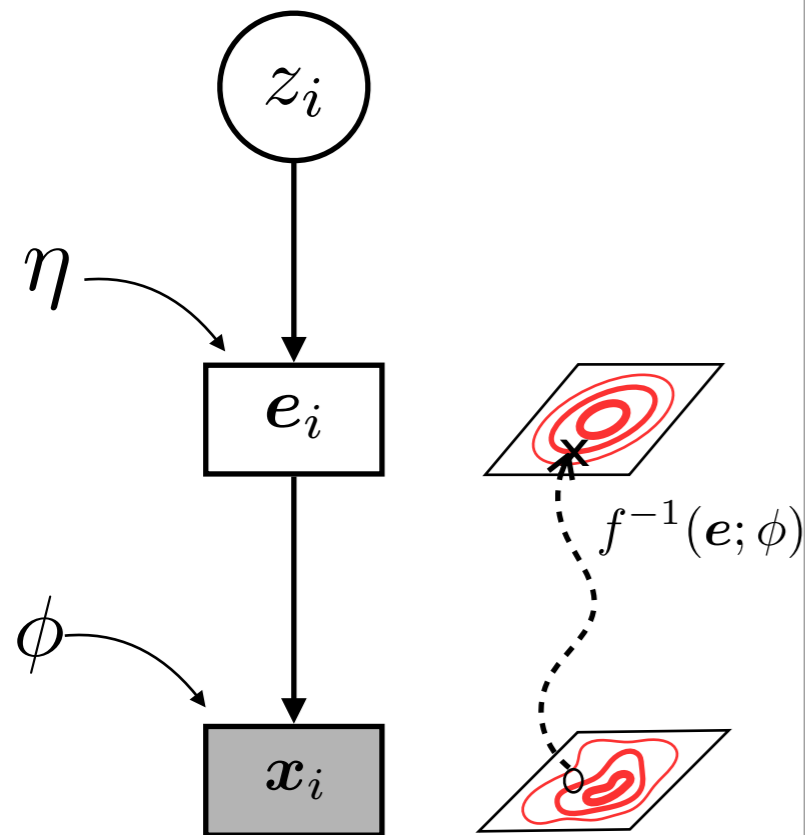
# Learning and Inference



$$p(\mathbf{x}_i | z_i; \eta, \phi)$$

Projection parameters

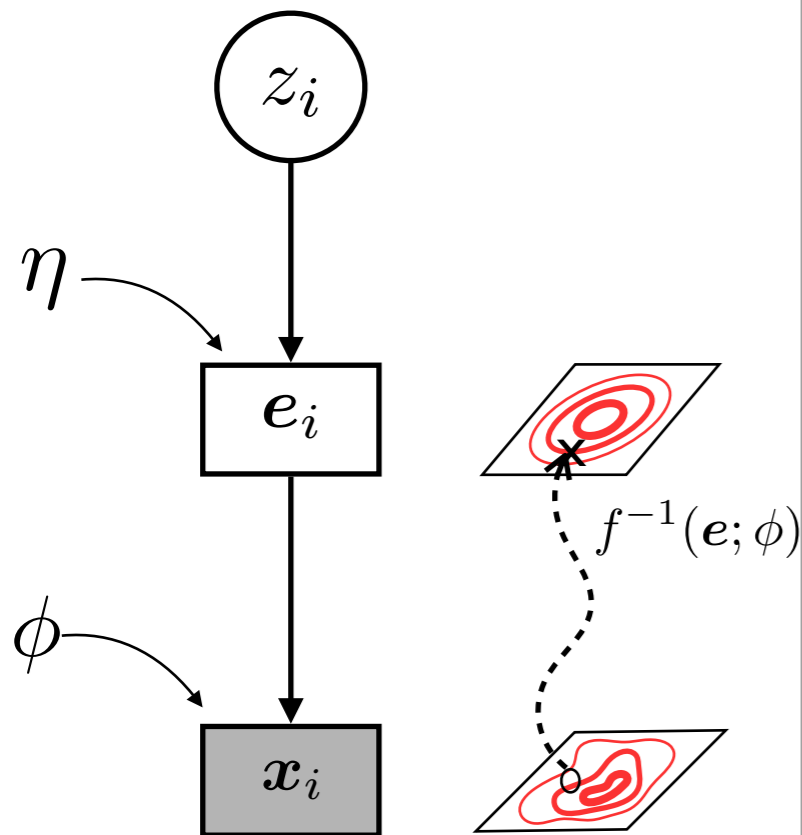
# Learning and Inference



$\dim(\mathbf{x}) = \dim(\mathbf{e})$  and  $f$  is invertible

$$p(\mathbf{x}_i | z_i; \eta, \phi)$$

# Learning and Inference

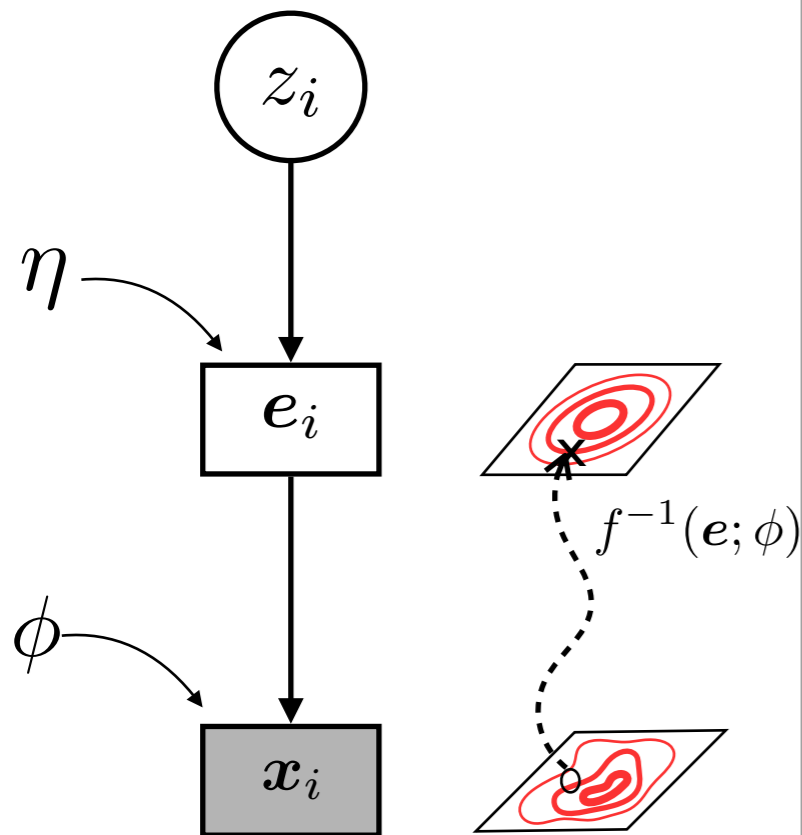


$\dim(\mathbf{x}) = \dim(\mathbf{e})$  and  $f$  is invertible

$$p(\mathbf{x}_i | z_i; \eta, \phi)$$

$$= p(f_\phi^{-1}(\mathbf{x}_i) | z_i; \eta) \left| \det \frac{\partial f^{-1}}{\partial \mathbf{x}_i} \right|$$

# Learning and Inference



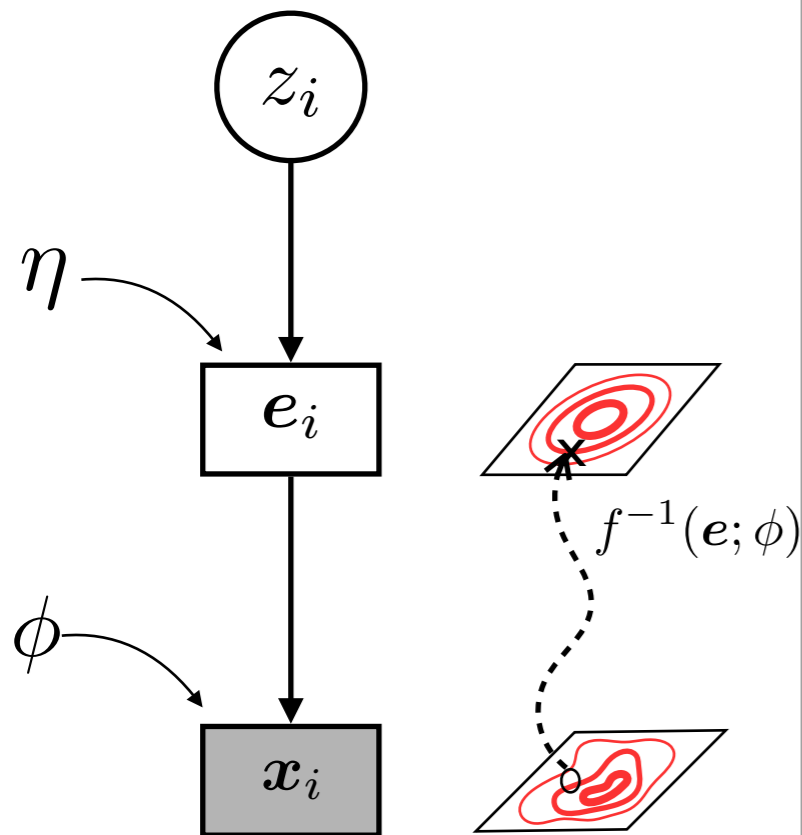
$\dim(\mathbf{x}) = \dim(\mathbf{e})$  and  $f$  is invertible

$$p(\mathbf{x}_i | z_i; \eta, \phi)$$

$$= p(f_\phi^{-1}(\mathbf{x}_i) | z_i; \eta) \left| \det \frac{\partial f^{-1}}{\partial \mathbf{x}_i} \right|$$

Determinant of Jacobian matrix

# Learning and Inference



$\dim(\mathbf{x}) = \dim(\mathbf{e})$  and  $f$  is invertible

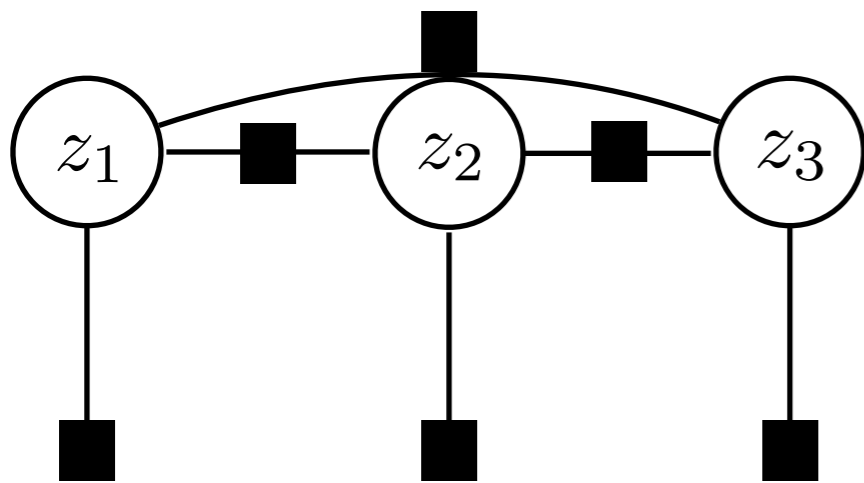
$$p(\mathbf{x}_i | z_i; \eta, \phi)$$

$$= p(f_\phi^{-1}(\mathbf{x}_i) | z_i; \eta) \left| \det \frac{\partial f^{-1}}{\partial \mathbf{x}_i} \right|$$

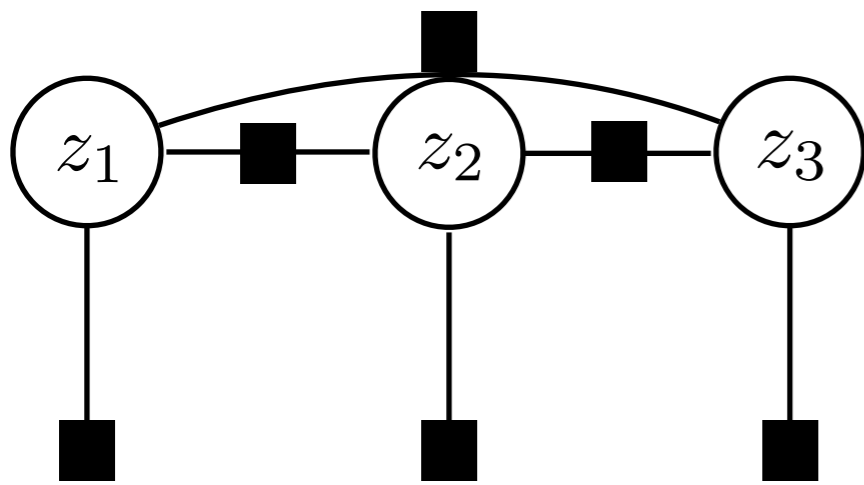
Gaussian distribution

Determinant of Jacobian matrix

# Learning and Inference

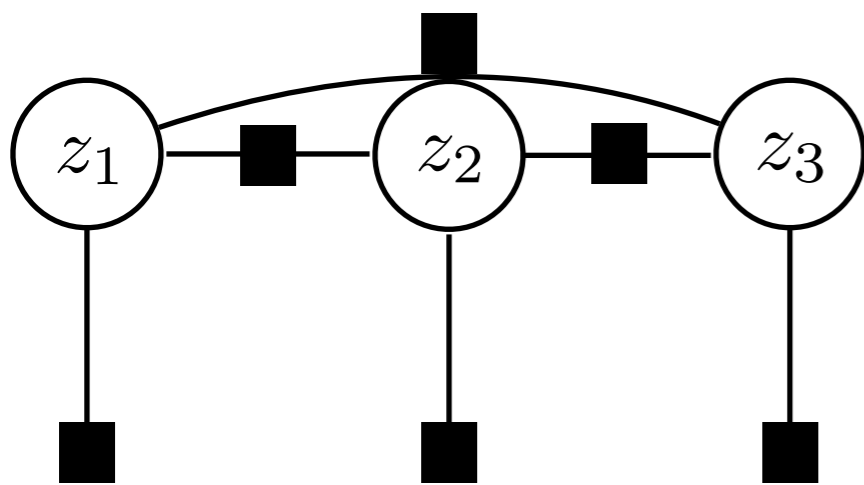


# Learning and Inference



$$p(f_{\phi}^{-1}(\mathbf{x}_i) | z_i; \eta) \left| \det \frac{\partial f^{-1}}{\partial \mathbf{x}_i} \right|$$

# Learning and Inference



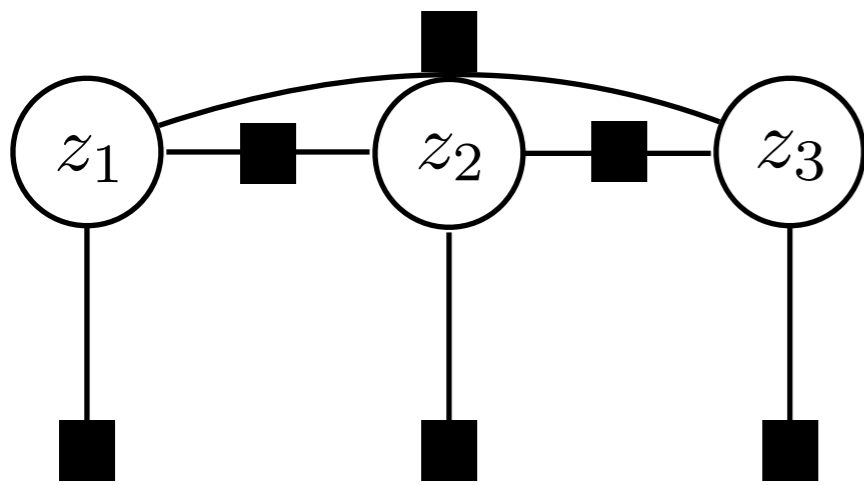
$$p(f_{\phi}^{-1}(\mathbf{x}_i) | z_i; \eta) \left| \det \frac{\partial f^{-1}}{\partial \mathbf{x}_i} \right|$$

Example of Markov prior

$$\log p(\mathbf{x}) = \log p_{\text{GHMM}}(f_{\phi}^{-1}(\mathbf{x}))$$



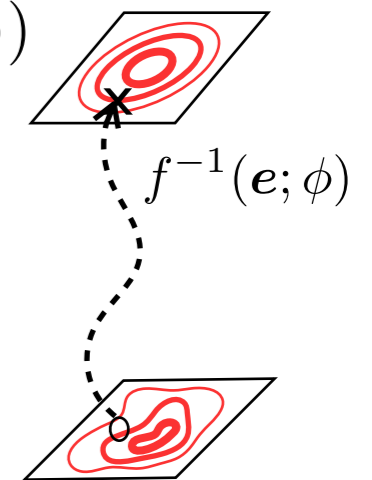
# Learning and Inference



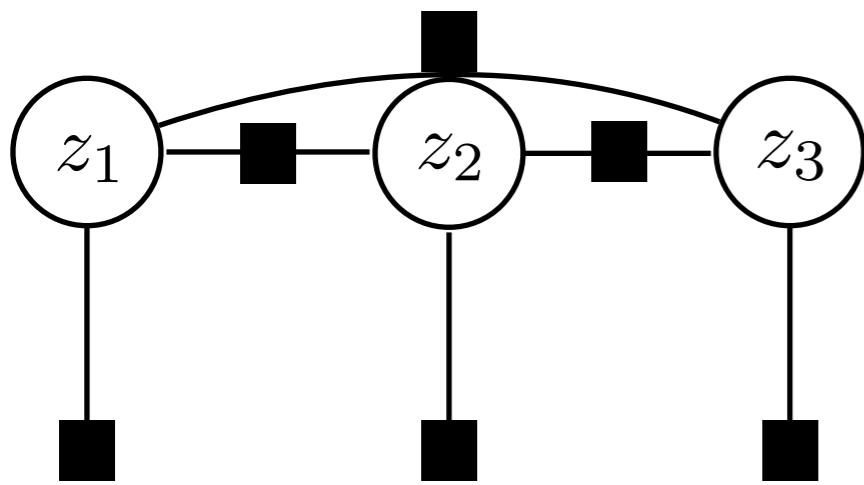
$$p(f_{\phi}^{-1}(\mathbf{x}_i) | z_i; \eta) \left| \det \frac{\partial f^{-1}}{\partial \mathbf{x}_i} \right|$$

Example of Markov prior

$$\log p(\mathbf{x}) = \log p_{\text{GHMM}}(f_{\phi}^{-1}(\mathbf{x}))$$



# Learning and Inference

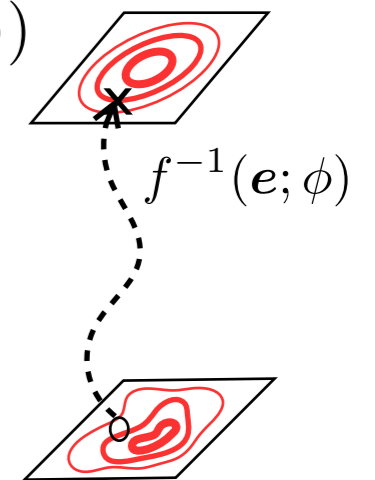


$$p(f_{\phi}^{-1}(\mathbf{x}_i) | z_i; \eta) \left| \det \frac{\partial f^{-1}}{\partial \mathbf{x}_i} \right|$$

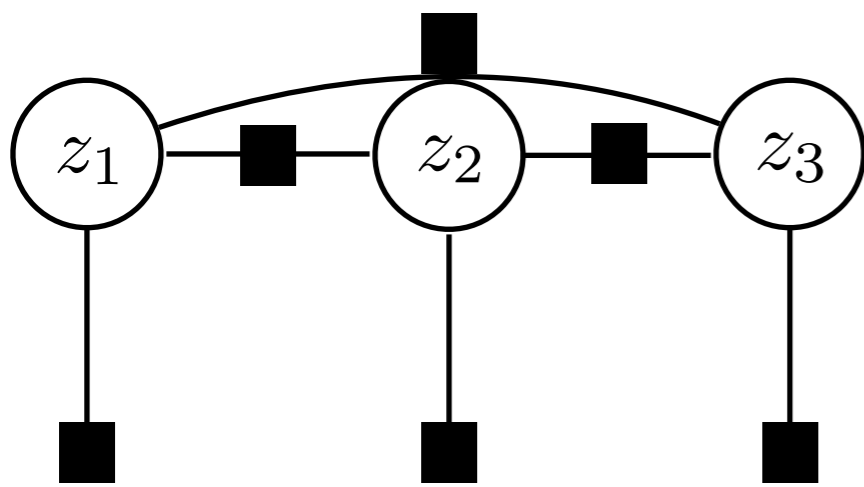
## Example of Markov prior

$$\log p(\mathbf{x}) = \log p_{\text{GHMM}}(f_{\phi}^{-1}(\mathbf{x}))$$

$$+ \sum \log \left| \det \frac{\partial f_{\phi}^{-1}}{\partial \mathbf{x}_i} \right|$$



# Learning and Inference



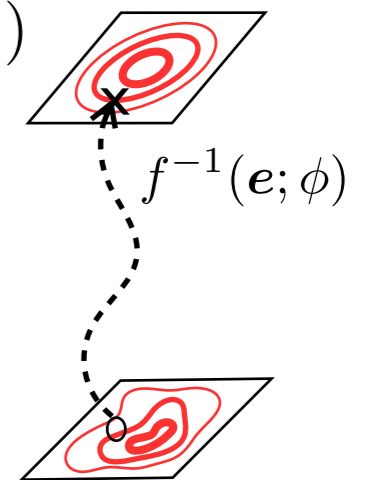
$$p(f_{\phi}^{-1}(\mathbf{x}_i) | z_i; \eta) \left| \det \frac{\partial f^{-1}}{\partial \mathbf{x}_i} \right|$$

## Example of Markov prior

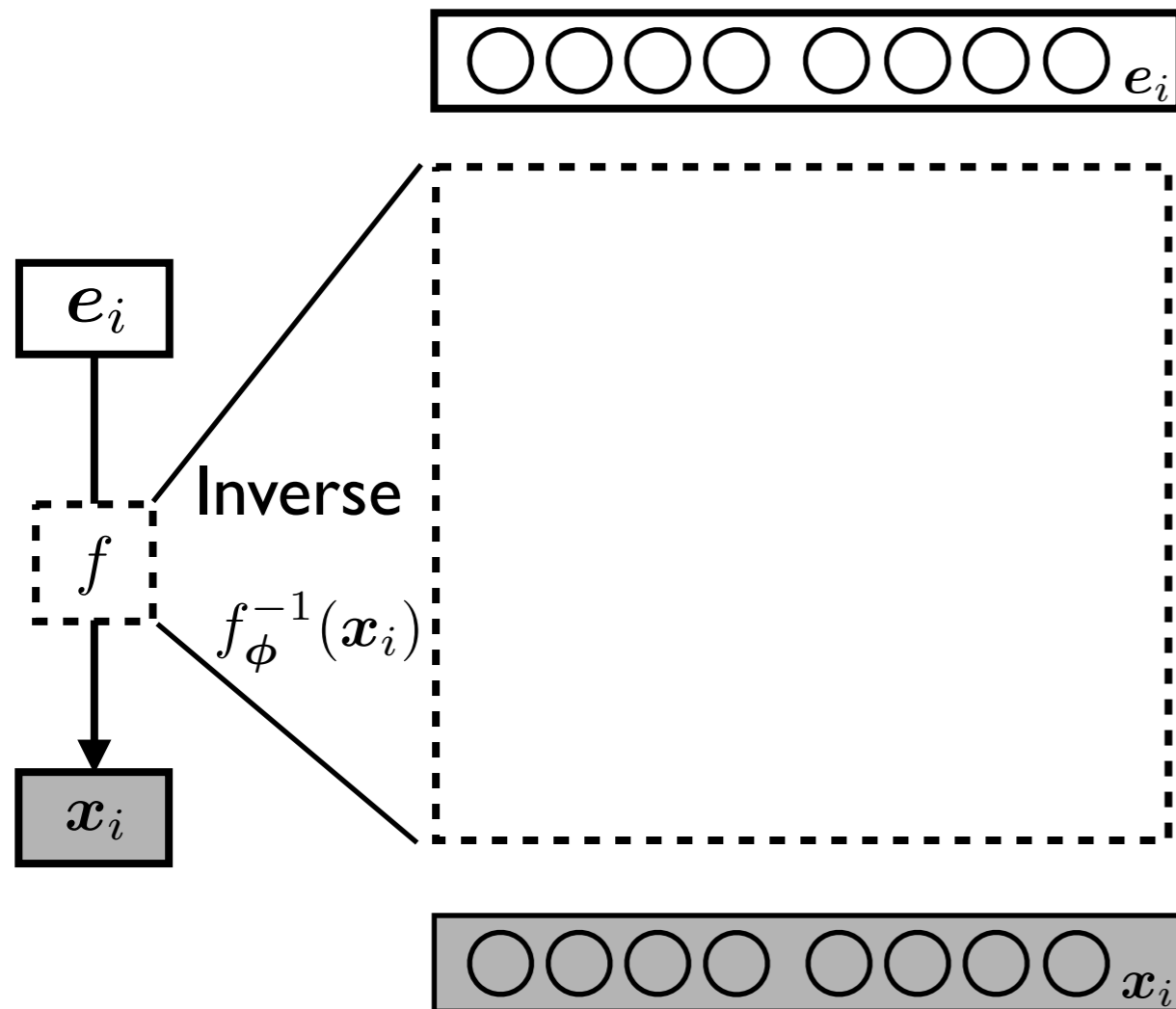
$$\log p(\mathbf{x}) = \log p_{\text{GHMM}}(f_{\phi}^{-1}(\mathbf{x}))$$

$$+ \sum \log \left| \det \frac{\partial f_{\phi}^{-1}}{\partial \mathbf{x}_i} \right|$$

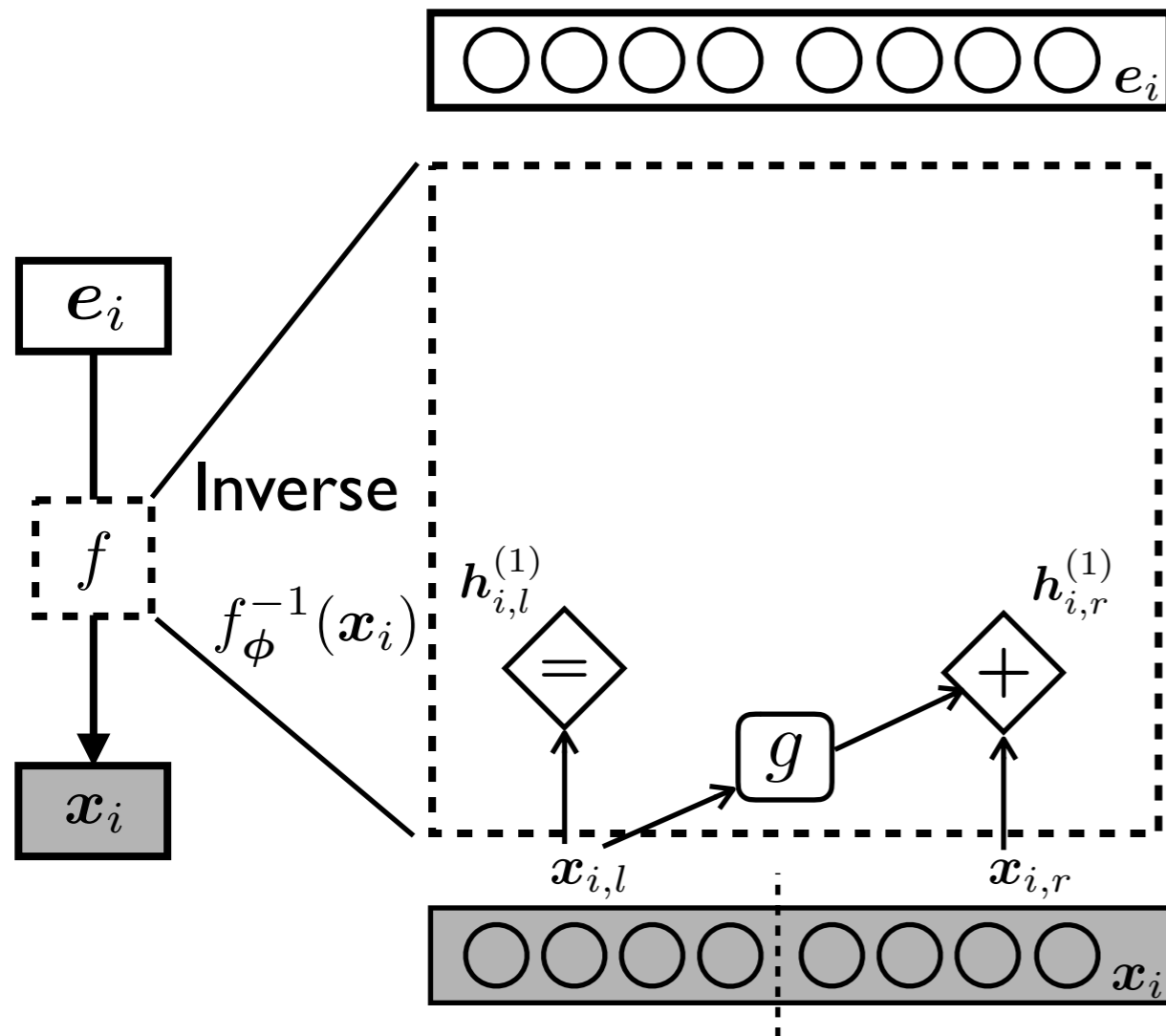
$-\infty$  when  $f$  is not invertible



# Learning with Inverse Projection



# Learning with Inverse Projection

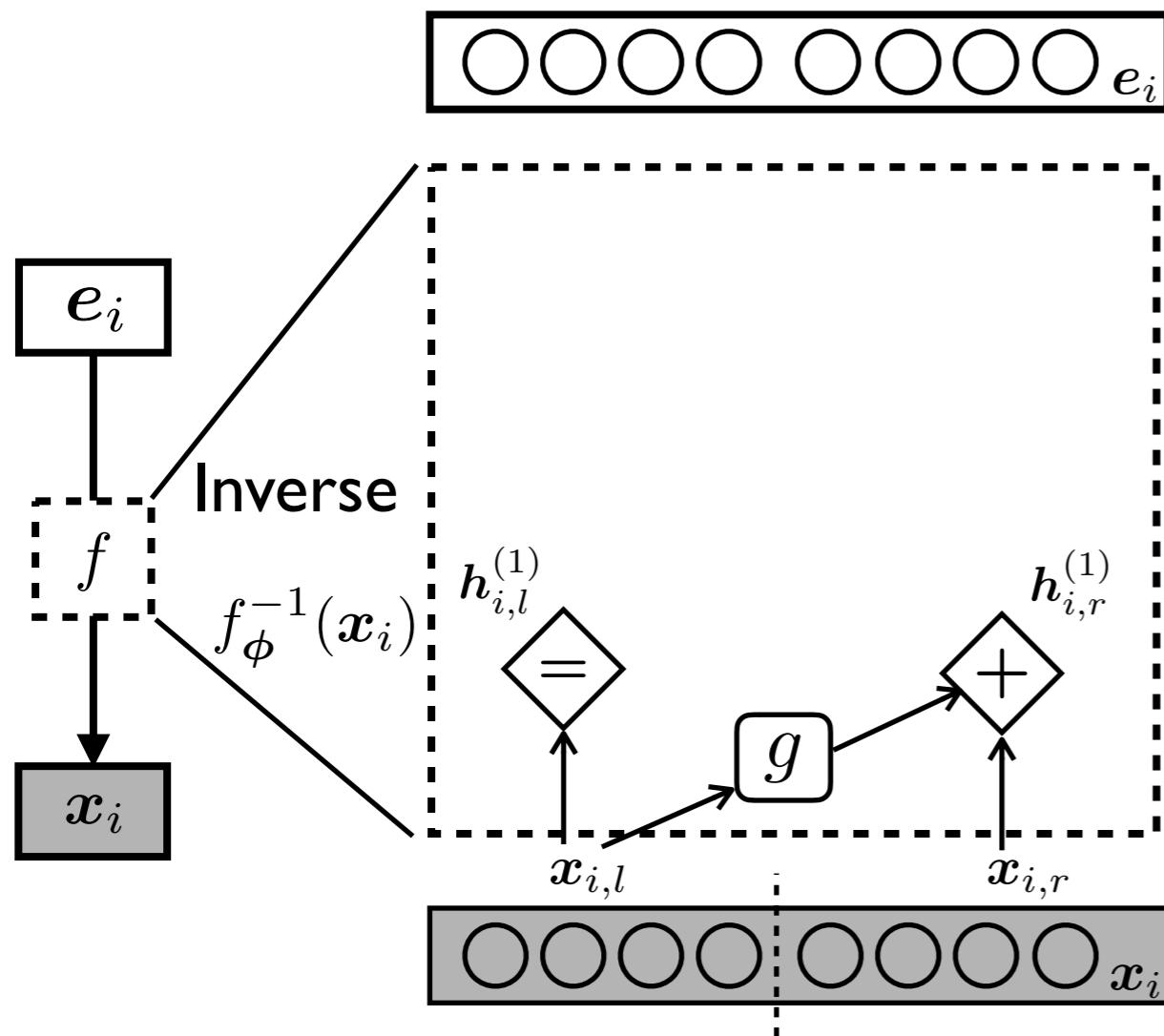


$$h_{i,l}^{(1)} = \mathbf{x}_{i,l}$$

$$h_{i,r}^{(1)} = \mathbf{x}_{i,r} + g(\mathbf{x}_{i,l})$$

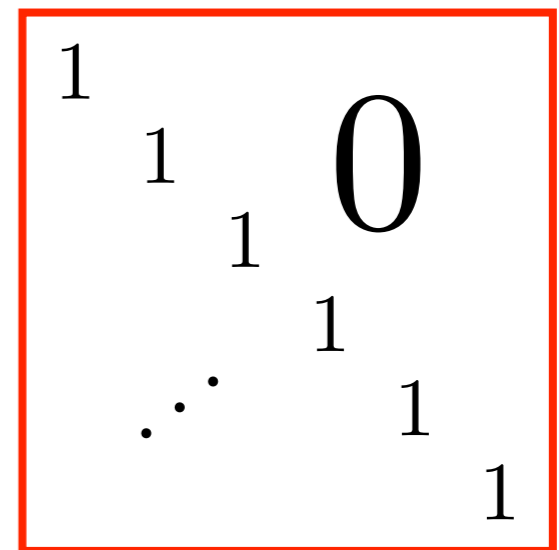
[Dinh et al. 2014]

# Learning with Inverse Projection



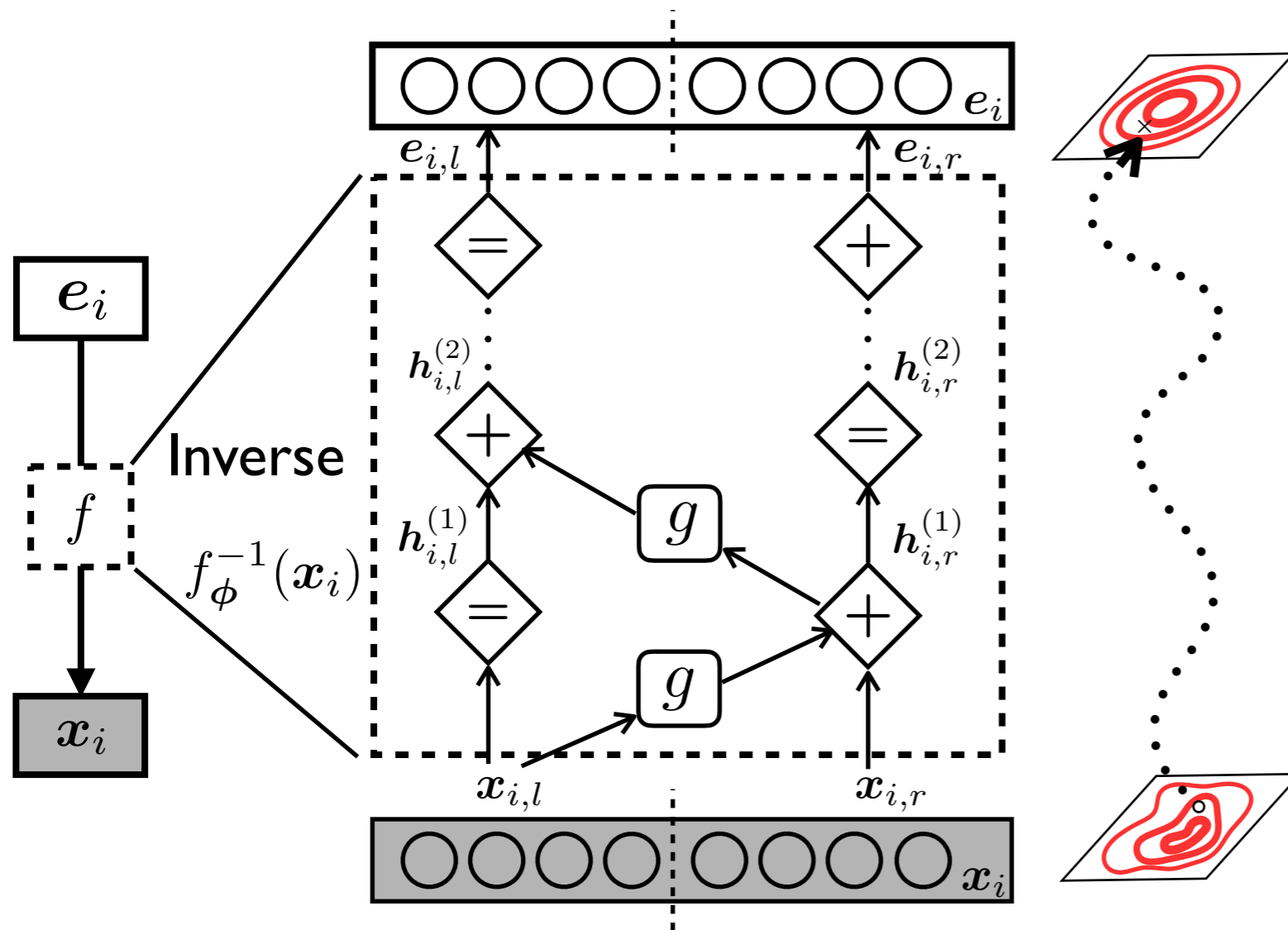
$$h_{i,l}^{(1)} = x_{i,l}$$

$$h_{i,r}^{(1)} = x_{i,r} + g(x_{i,l})$$



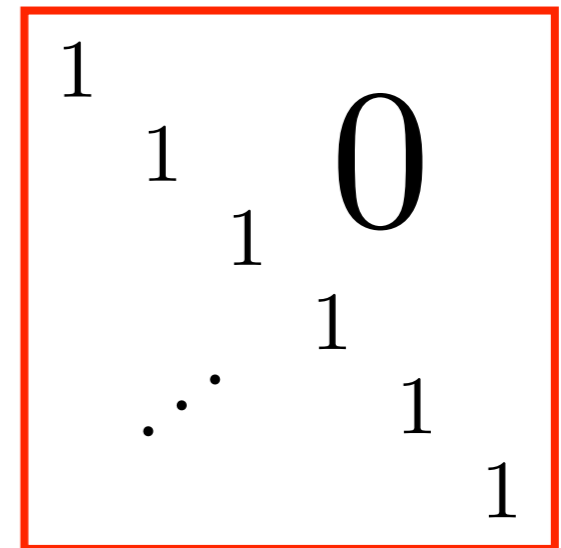
[Dinh et al. 2014]

# Learning with Inverse Projection



$$h_{i,l}^{(1)} = \mathbf{x}_{i,l}$$

$$h_{i,r}^{(1)} = \mathbf{x}_{i,r} + g(\mathbf{x}_{i,l})$$



[Dinh et al. 2014]

# Experiments

- Dataset: English Penn Treebank
- POS tagging

Trained and tested on whole PTB

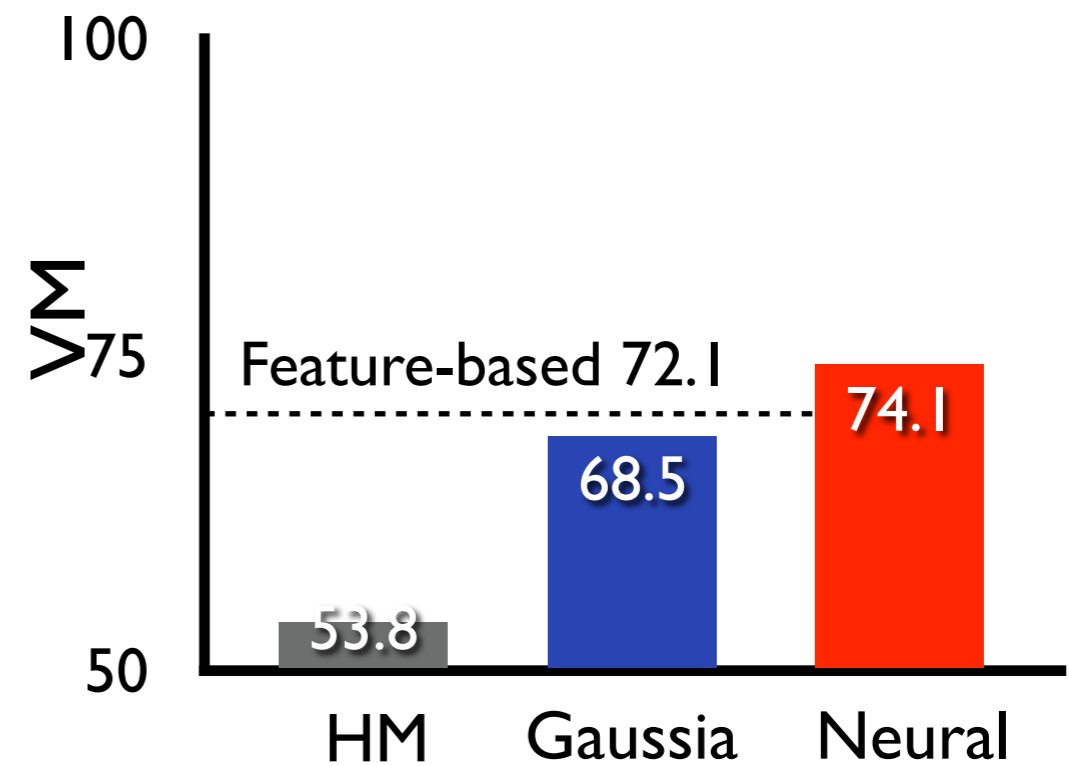
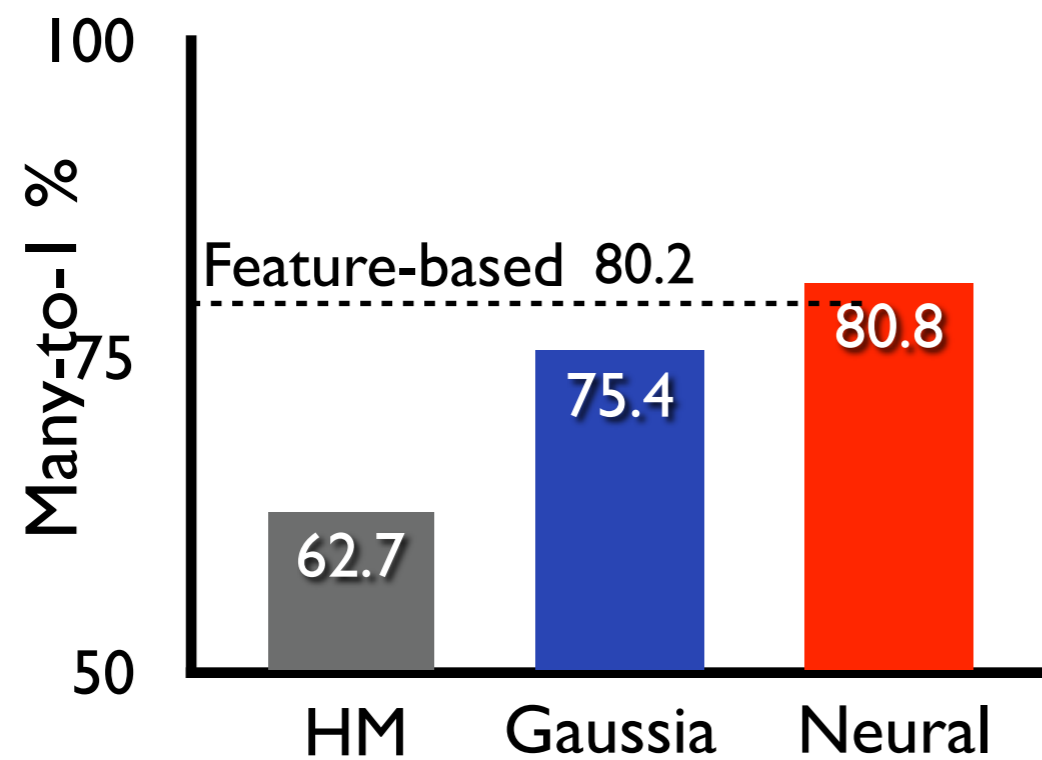
- Grammar induction

Trained on sentences of length  $\leq 10$  in section 2-21

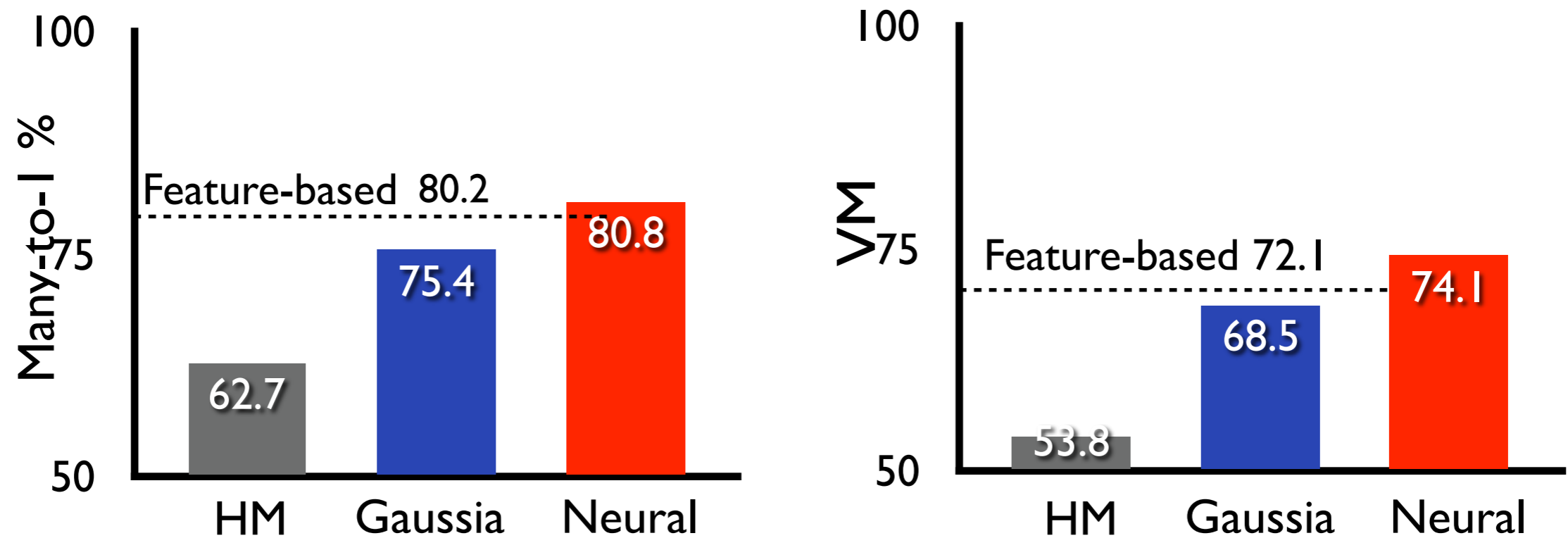
Tested on sentences in section 23



# Part-of-speech Induction

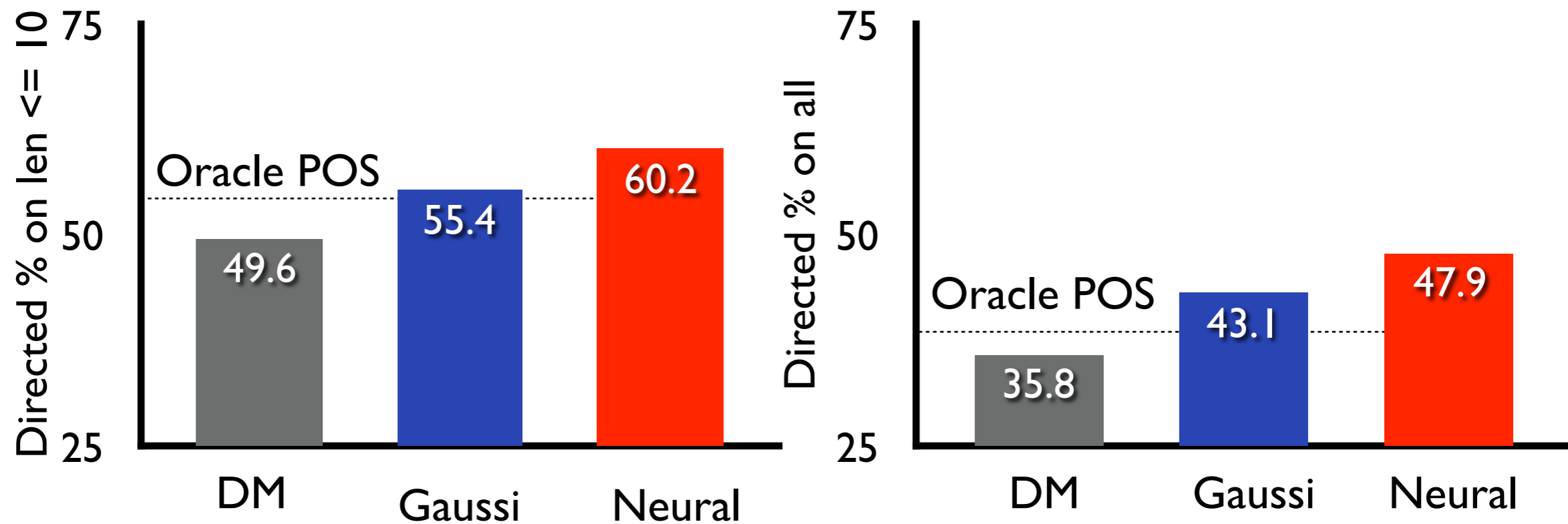


# Part-of-speech Induction

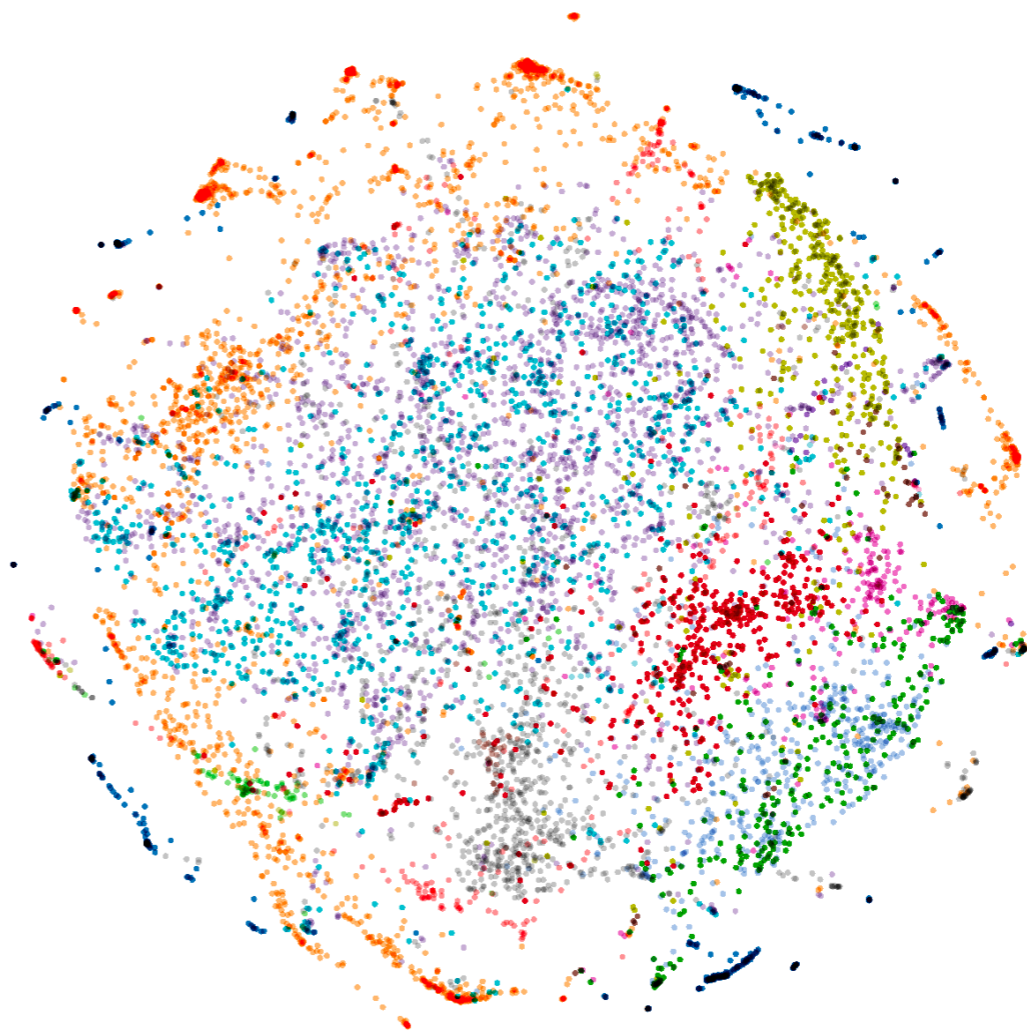


Outperform feature-based SOTA

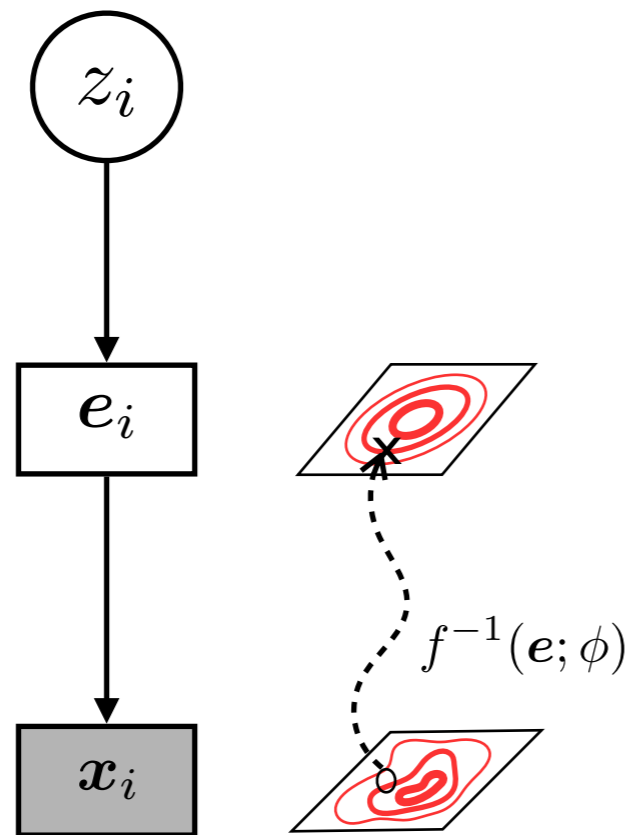
# Dependency Parse Induction



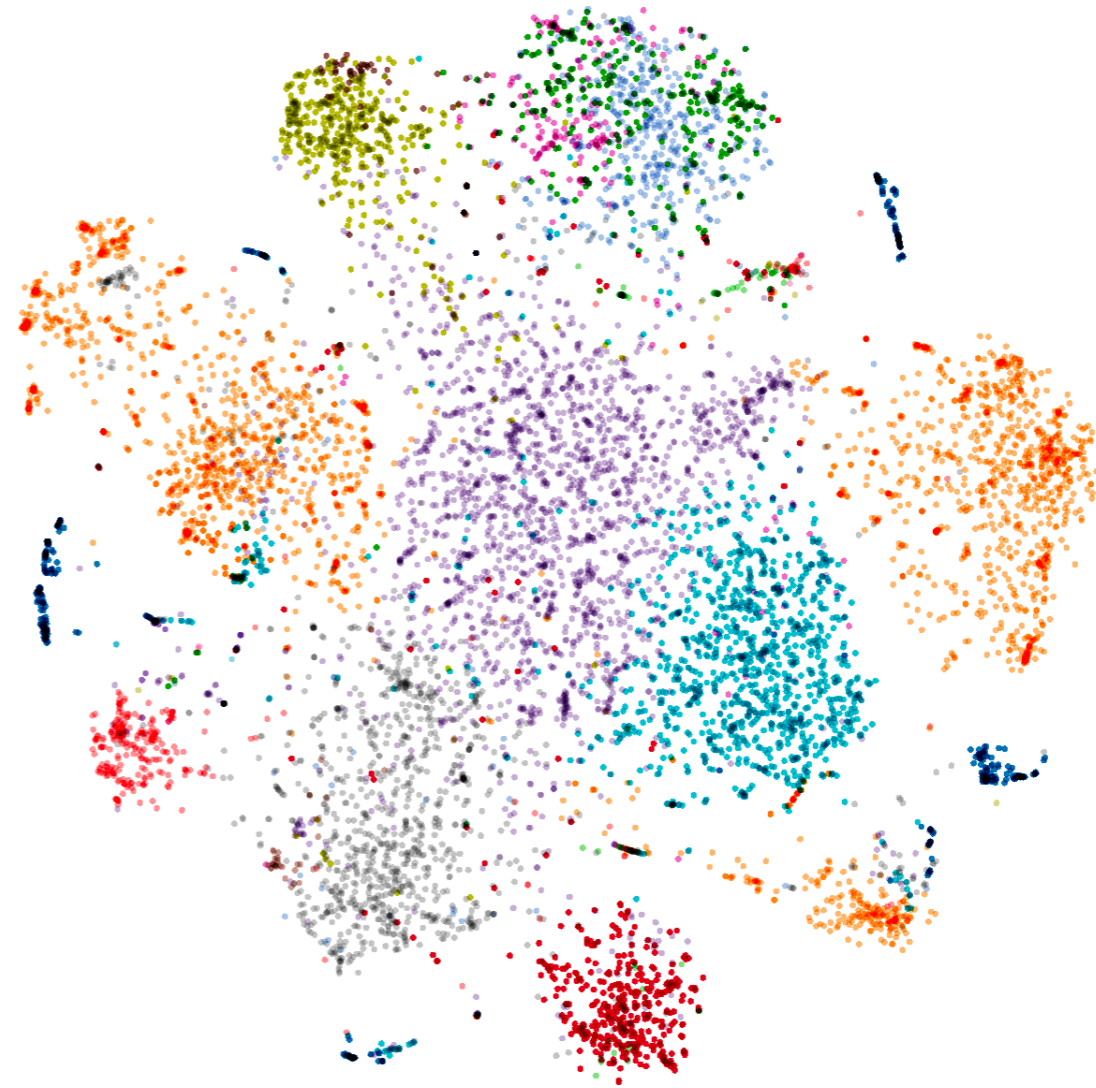
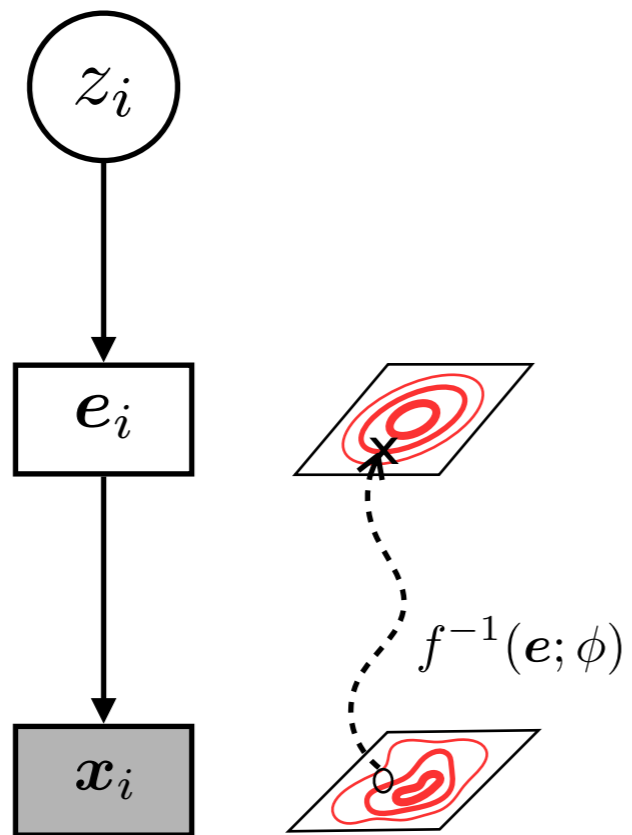
# Original Embedding Space



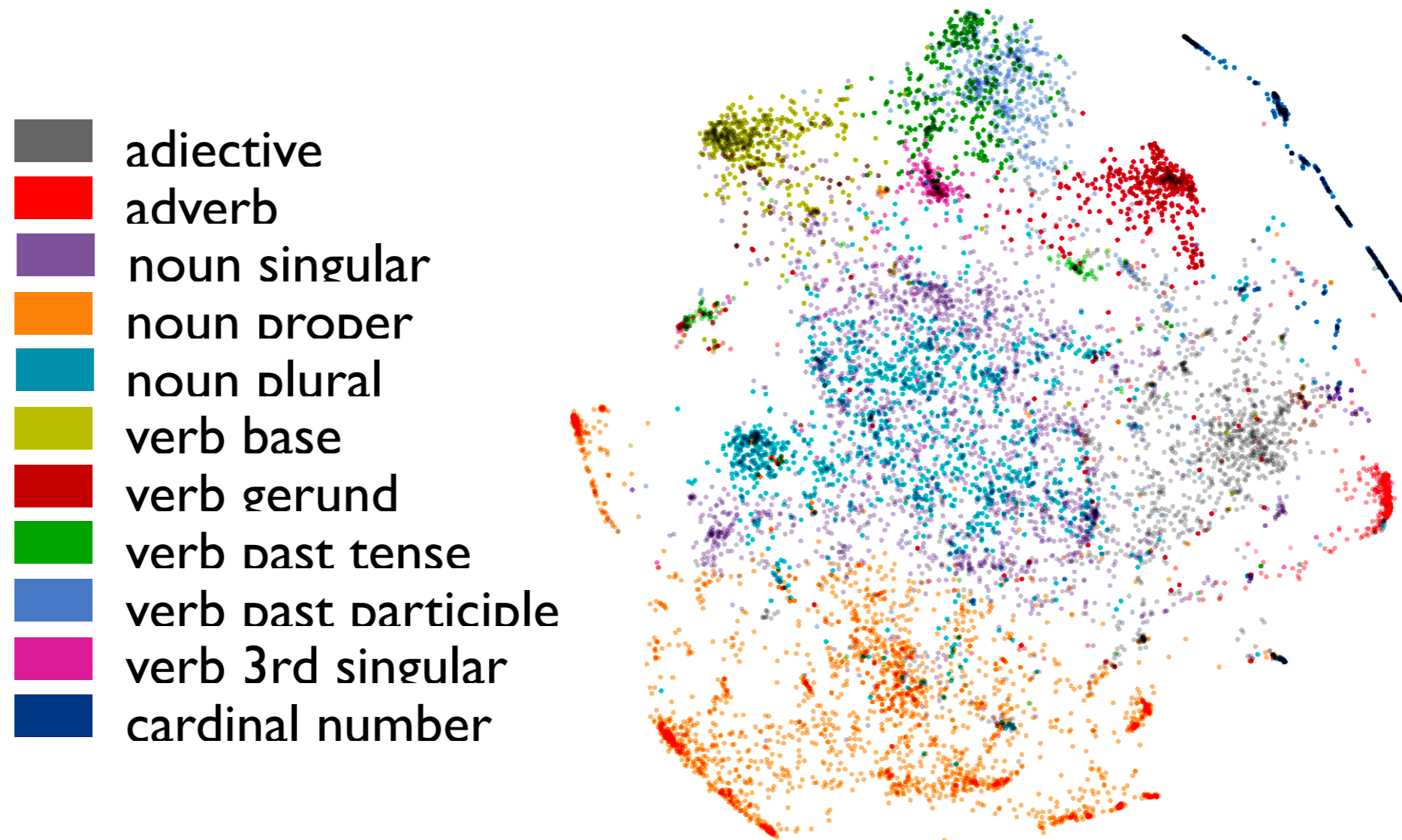
# Projected Embedding Space w/ Markov Prior



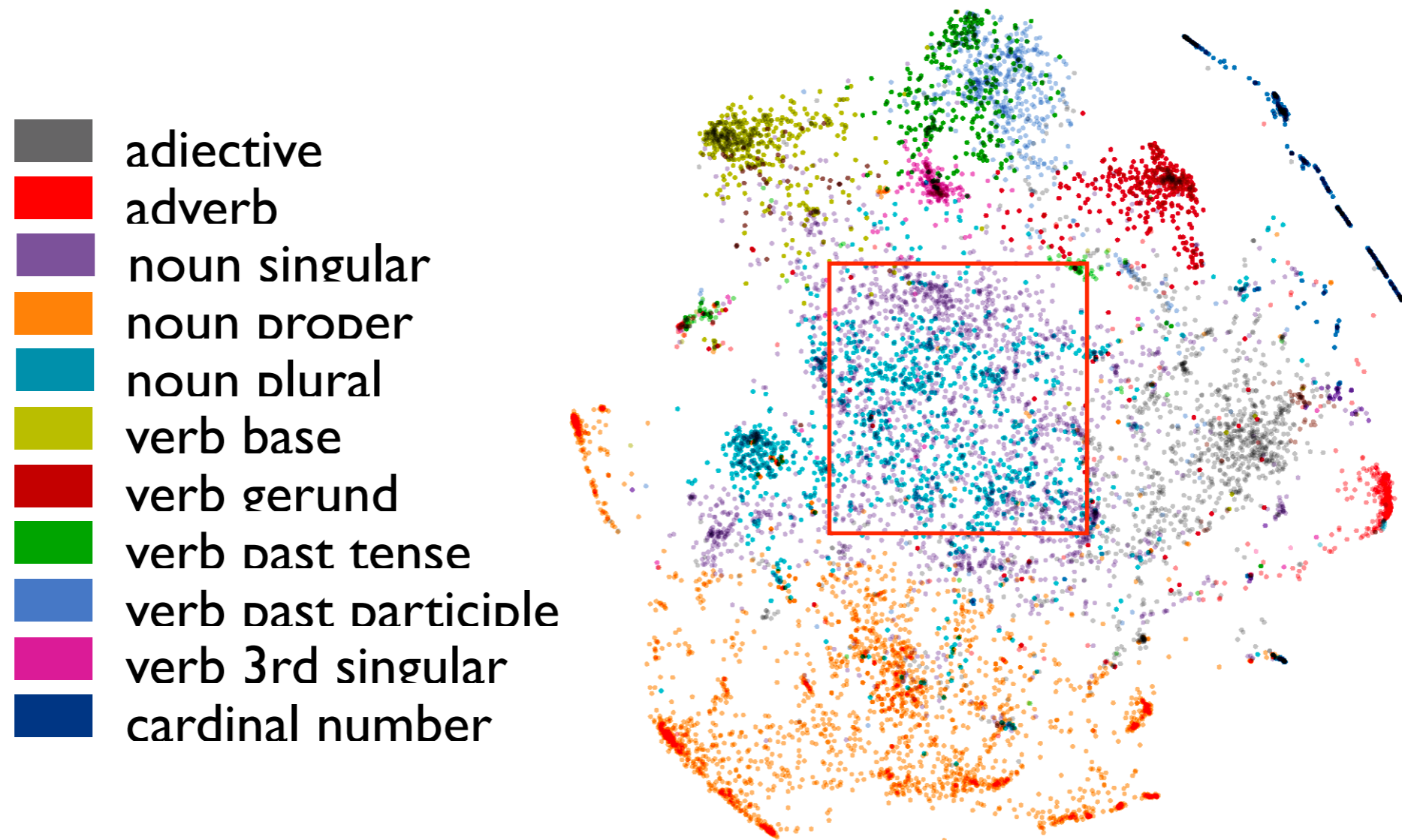
# Projected Embedding Space w/ Markov Prior



# Projected Embedding Space w/ DMV Prior

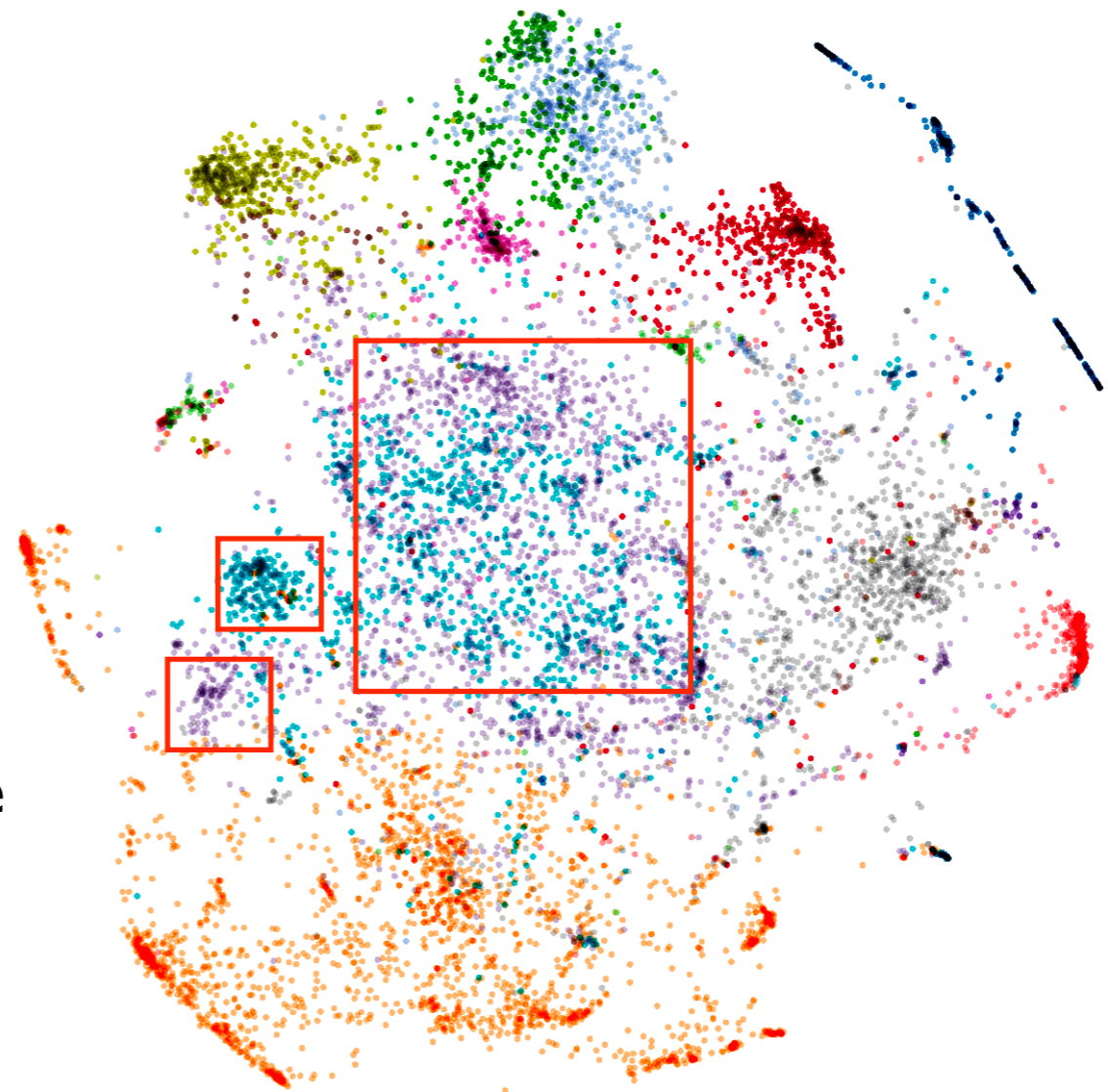


# Projected Embedding Space w/ DMV Prior

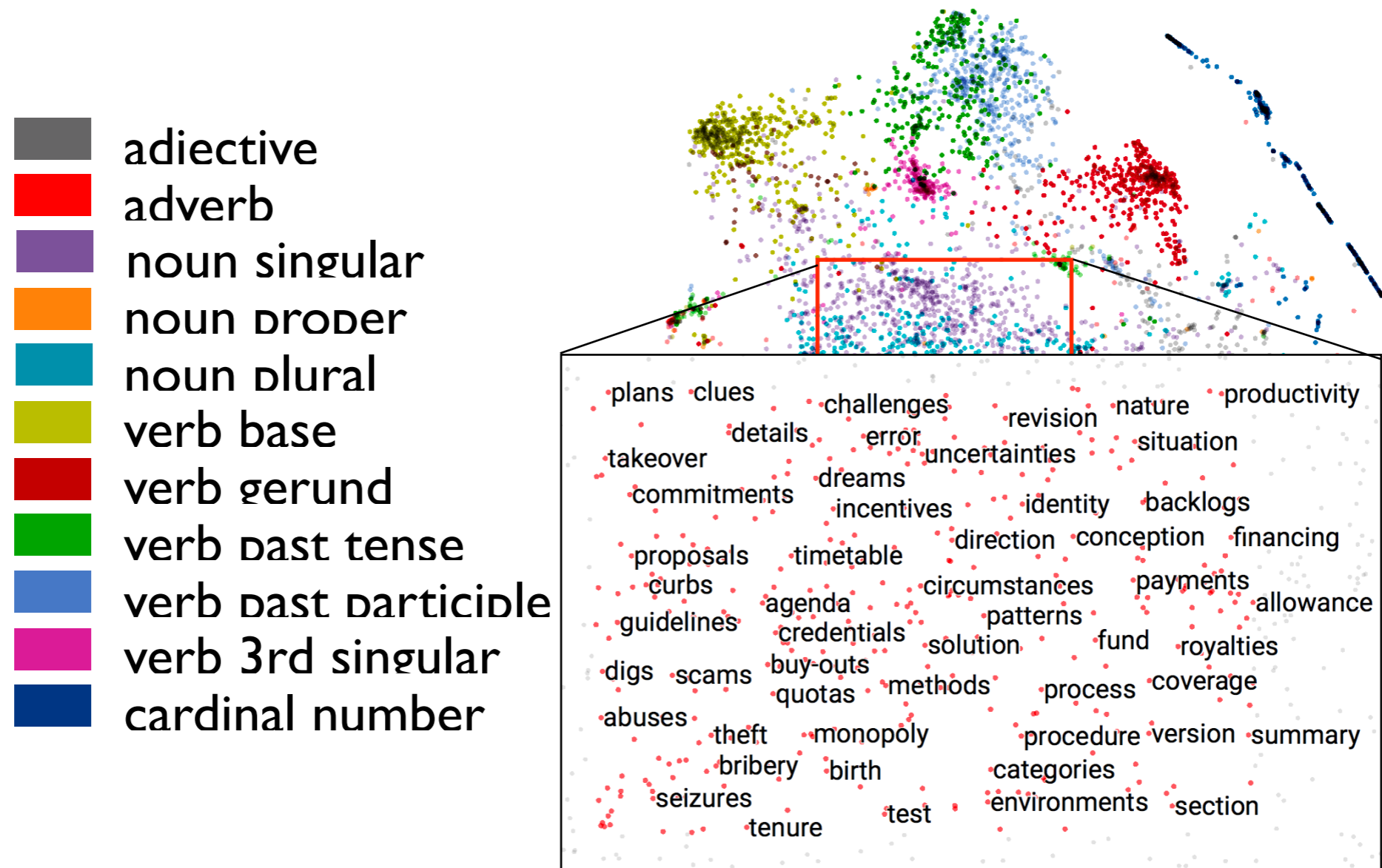




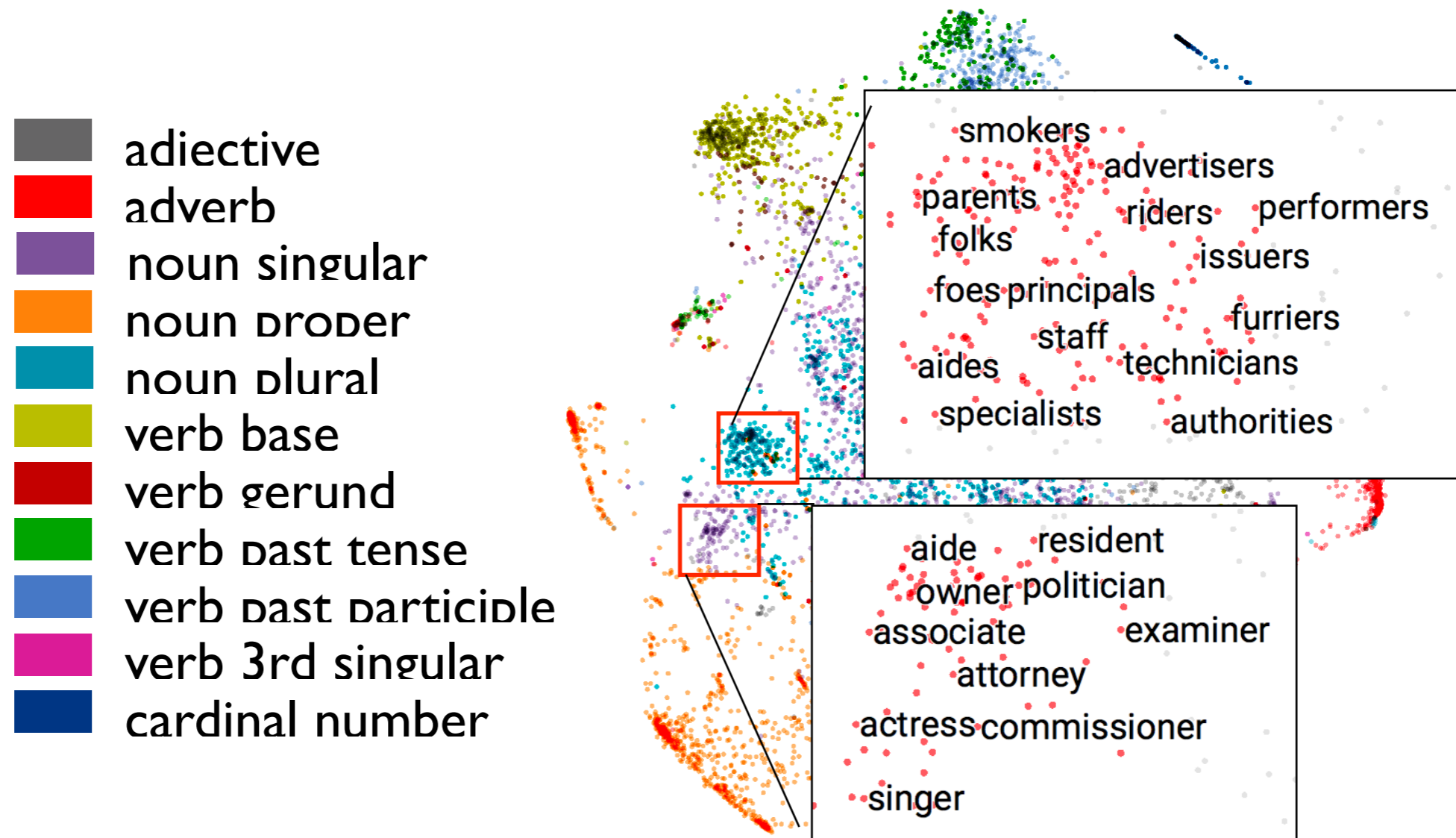
# Projected Embedding Space w/ DMV Prior



# Projected Embedding Space w/ DMV Prior



# Projected Embedding Space w/ DMV Prior



# Conclusion

# Conclusion

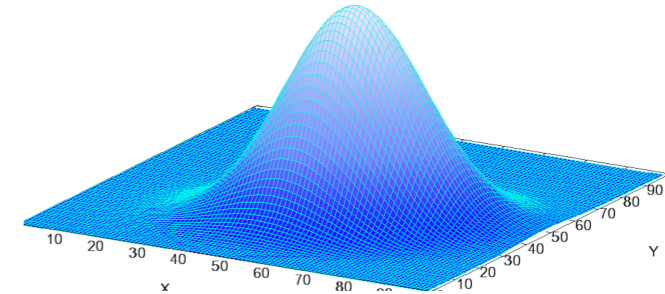
# Conclusion

- Normalizing flows for unsupervised learning



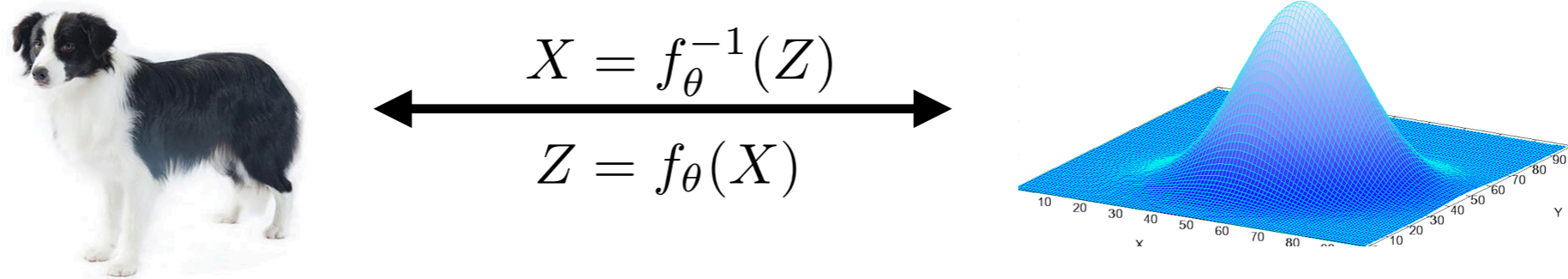
$$X = f_{\theta}^{-1}(Z)$$

$$Z = f_{\theta}(X)$$

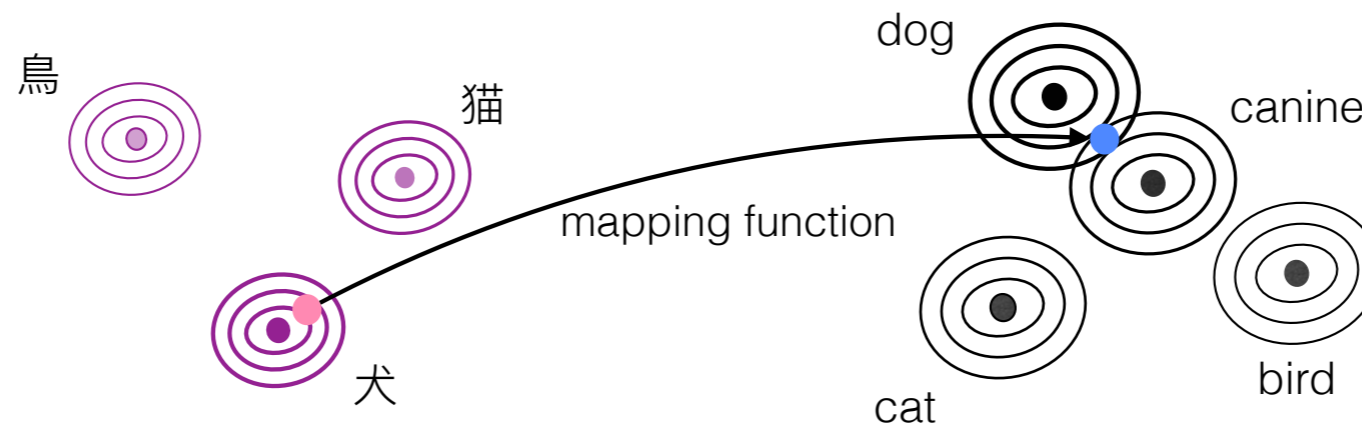


# Conclusion

- Normalizing flows for unsupervised learning

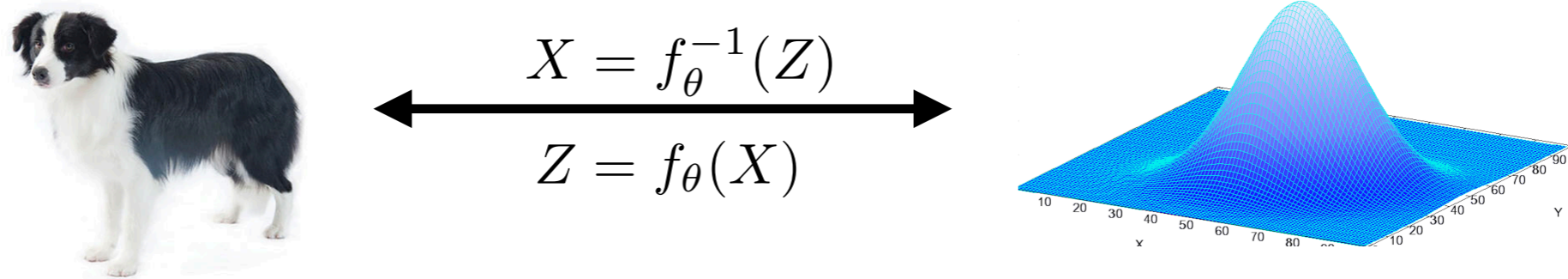


- Learning of bilingual lexicons

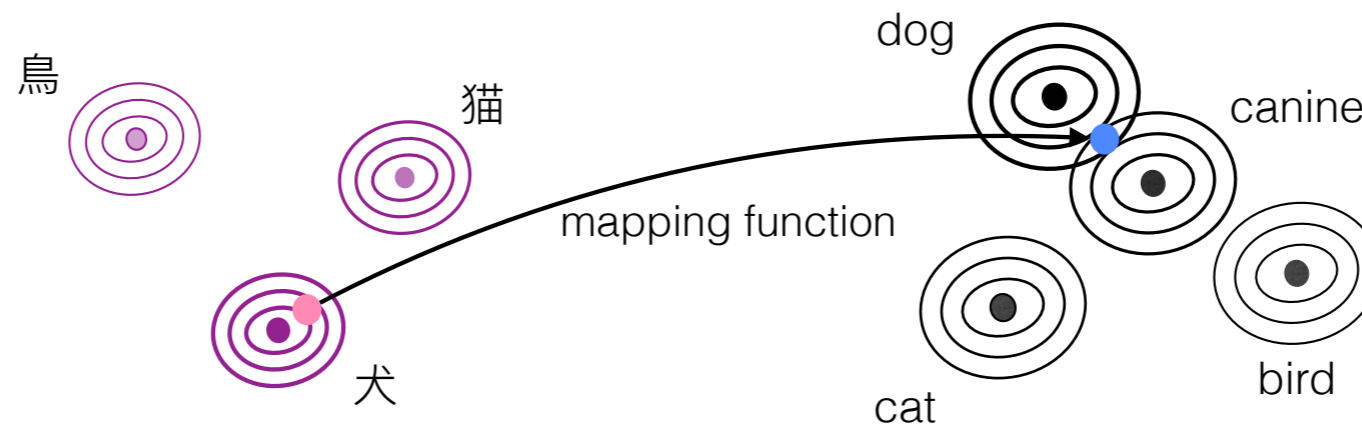


# Conclusion

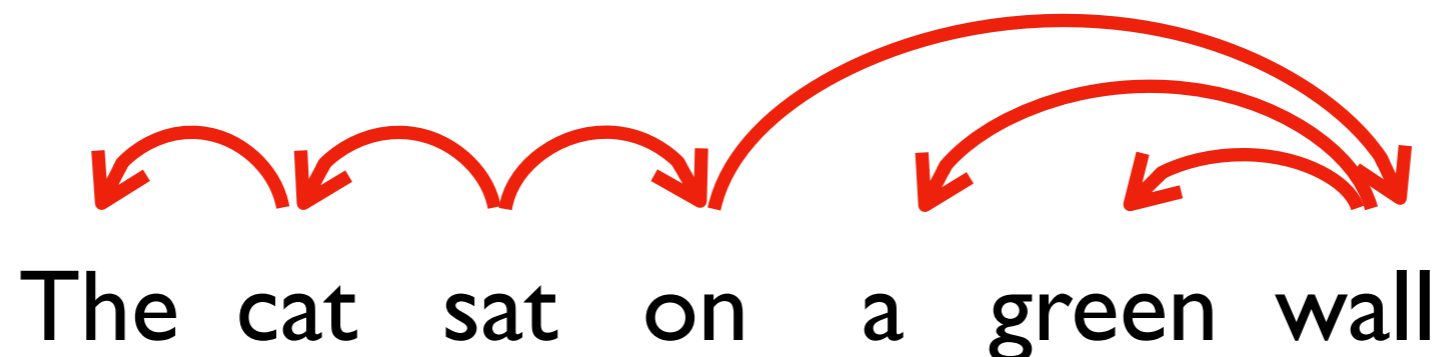
- Normalizing flows for unsupervised learning



- Learning of bilingual lexicons



- Learning of syntactic structure





# Thank You! Questions?

**DeMa-BWE**



<https://github.com/violet-zct/DeMa-BWE>

The cat sat on a green wall



<https://github.com/jxhe/struct-learning-with-flow>