



# Learning with Sparse Latent Structure

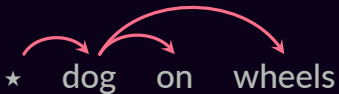
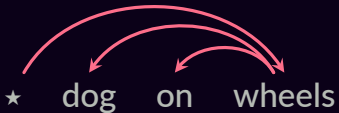
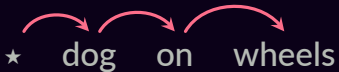
**Vlad Niculae**

Instituto de Telecomunicações

Work with: André Martins, Claire Cardie, Mathieu Blondel

# Structured Inference

...

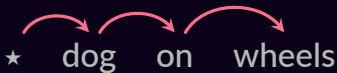


...

# Structured Inference

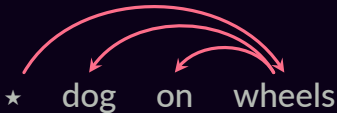
...

VERB    PREP    NOUN  
dog    on    wheels



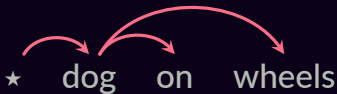
dog    hond  
on    op  
wheels    wielen

NOUN    PREP    NOUN  
dog    on    wheels



dog    hond  
on    op  
wheels    wielen

NOUN    DET    NOUN  
dog    on    wheels



dog    hond  
on    op  
wheels    wielen

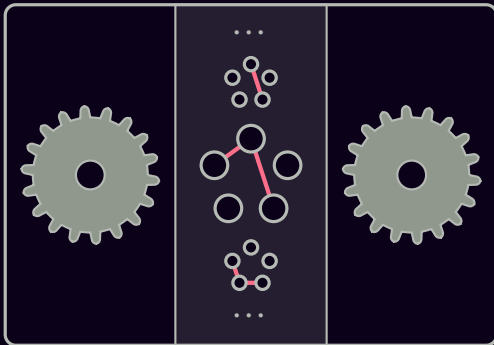
...

# Structured Inference

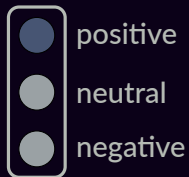


# Latent Structured Inference

input



output



\*record scratch\*

\*freeze frame\*

**How to select an item  
from a set?**

# How to select an item from a set?



...



# How to select an item from a set?

$c_1$

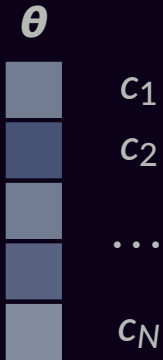
$c_2$

...

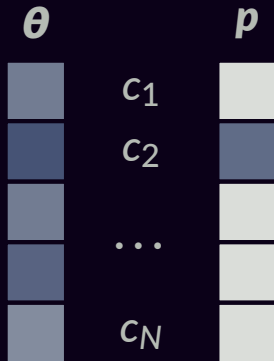
$c_N$



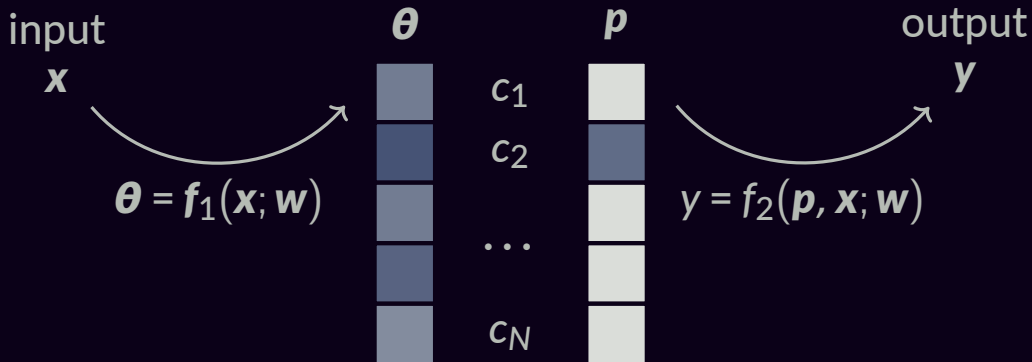
# How to select an item from a set?



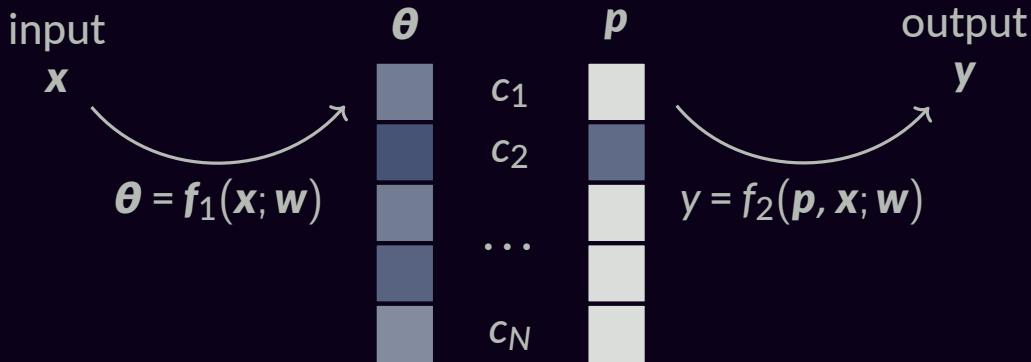
# How to select an item from a set?



# How to select an item from a set?

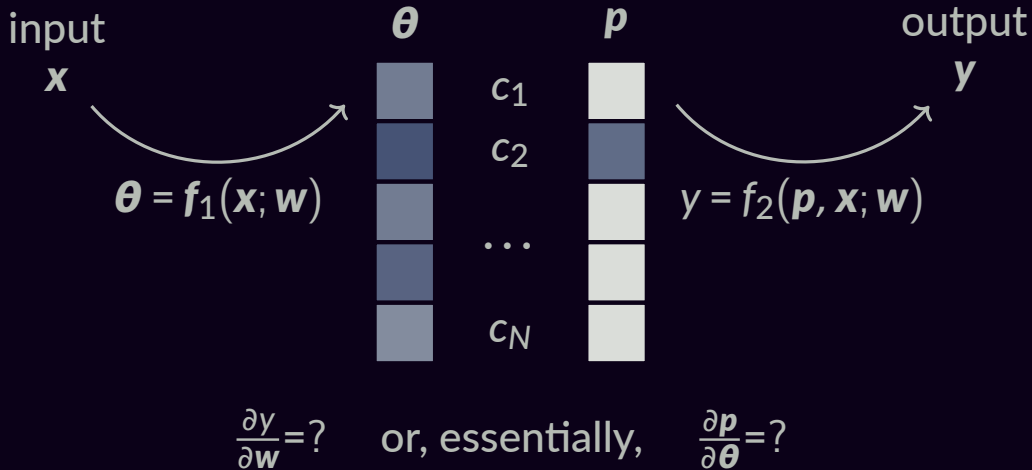


# How to select an item from a set?

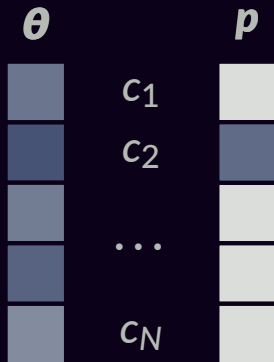


$$\frac{\partial y}{\partial \mathbf{w}} = ?$$

# How to select an item from a set?

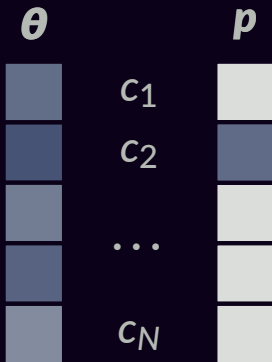


# Argmax



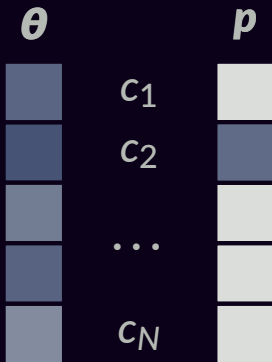
$$\frac{\partial p}{\partial \theta} = ?$$

# Argmax



$$\frac{\partial p}{\partial \theta} = ?$$

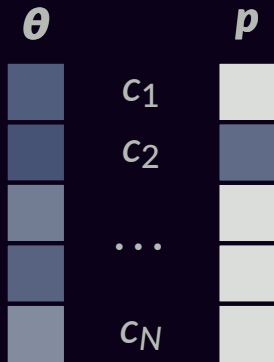
# Argmax



$$\frac{\partial p}{\partial \theta} = ?$$

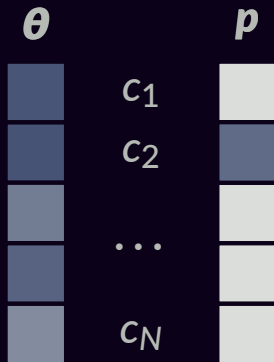


# Argmax



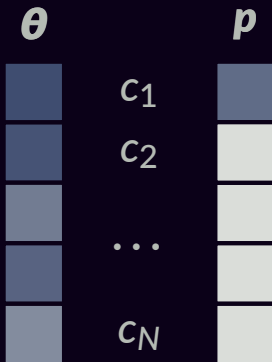
$$\frac{\partial p}{\partial \theta} = ?$$

# Argmax



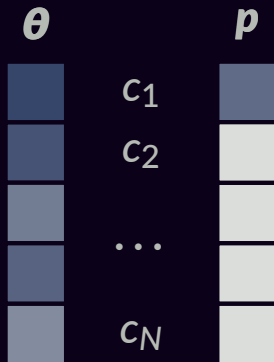
$$\frac{\partial p}{\partial \theta} = ?$$

# Argmax



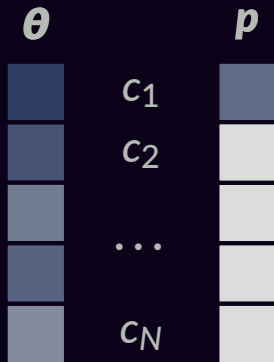
$$\frac{\partial p}{\partial \theta} = ?$$

# Argmax



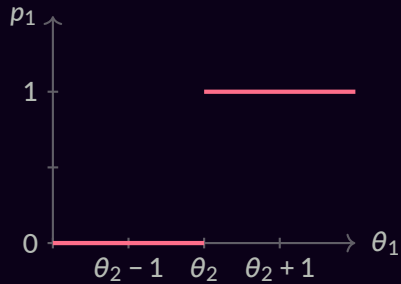
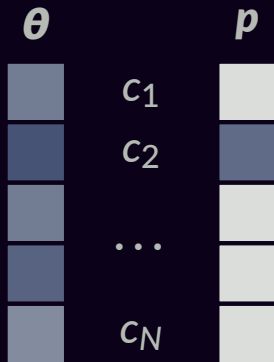
$$\frac{\partial p}{\partial \theta} = ?$$

# Argmax



$$\frac{\partial p}{\partial \theta} = ?$$

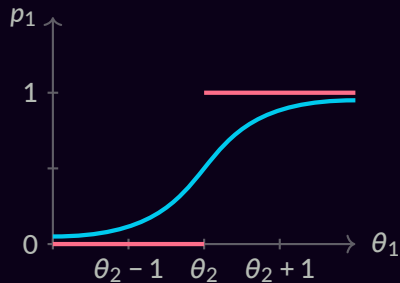
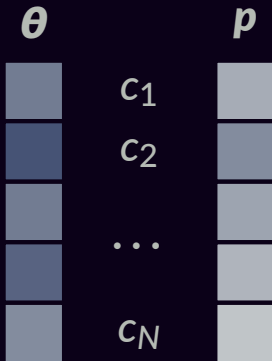
# Argmax



$$\frac{\partial p}{\partial \theta} = \mathbf{0}$$

# Argmax vs. Softmax

$$p_j = \exp(\theta_j)/Z$$



$$\frac{\partial \mathbf{p}}{\partial \boldsymbol{\theta}} = \text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top$$

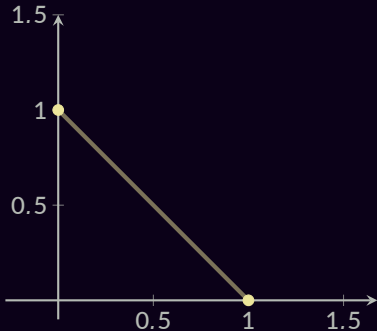
# Variational Form of Argmax

$$\Delta = \{\mathbf{p} \in \mathbb{R}^N : \mathbf{p} \geq \mathbf{0}, \mathbf{1}^\top \mathbf{p} = 1\}$$



# Variational Form of Argmax

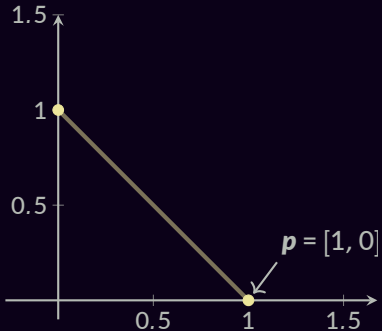
$$\Delta = \{p \in \mathbb{R}^N : p \geq 0, \mathbf{1}^\top p = 1\}$$



$N = 2$

# Variational Form of Argmax

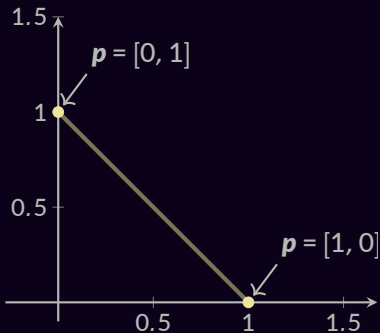
$$\Delta = \{p \in \mathbb{R}^N : p \geq 0, \mathbf{1}^\top p = 1\}$$



$N = 2$

# Variational Form of Argmax

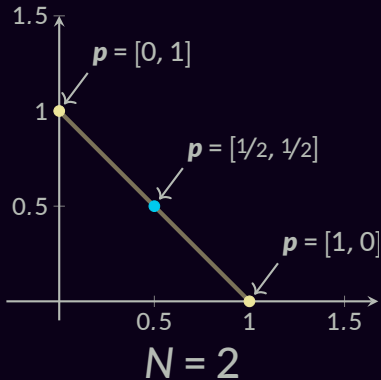
$$\Delta = \{ \mathbf{p} \in \mathbb{R}^N : \mathbf{p} \geq \mathbf{0}, \mathbf{1}^\top \mathbf{p} = 1 \}$$



$N = 2$

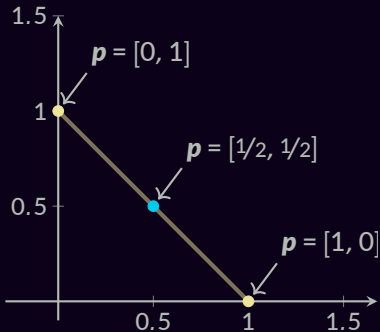
# Variational Form of Argmax

$$\Delta = \{p \in \mathbb{R}^N : p \geq 0, \mathbf{1}^\top p = 1\}$$

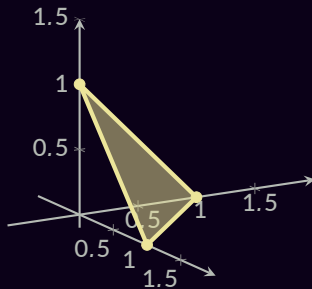


# Variational Form of Argmax

$$\Delta = \{p \in \mathbb{R}^N : p \geq 0, \mathbf{1}^\top p = 1\}$$



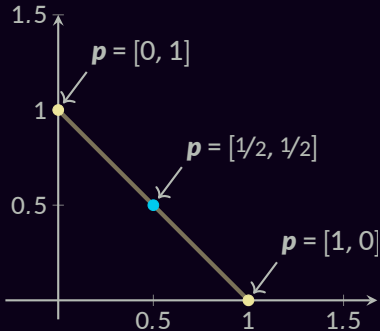
$N = 2$



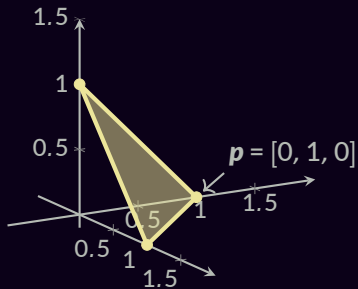
$N = 3$

# Variational Form of Argmax

$$\Delta = \{p \in \mathbb{R}^N : p \geq 0, \mathbf{1}^\top p = 1\}$$



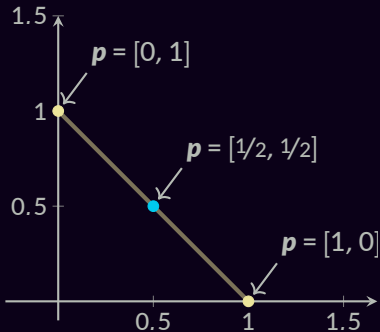
$N = 2$



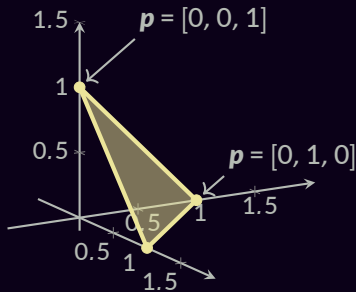
$N = 3$

# Variational Form of Argmax

$$\Delta = \{p \in \mathbb{R}^N : p \geq 0, \mathbf{1}^\top p = 1\}$$



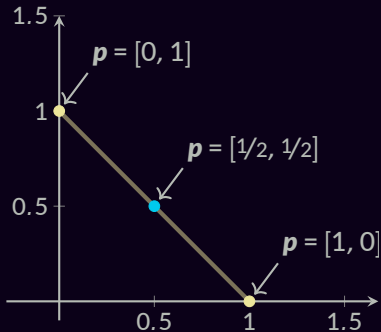
$N = 2$



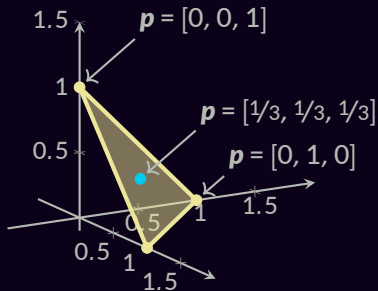
$N = 3$

# Variational Form of Argmax

$$\Delta = \{p \in \mathbb{R}^N : p \geq 0, \mathbf{1}^\top p = 1\}$$



$N = 2$



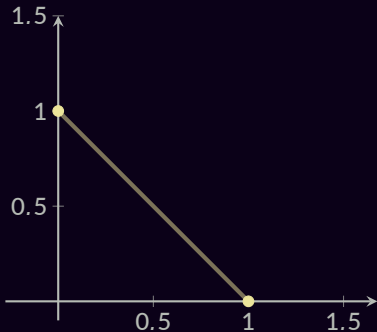
$N = 3$



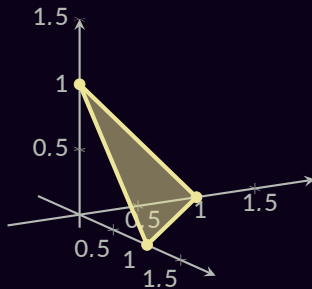
# Variational Form of Argmax

$$\max_j \theta_j = \max_{p \in \Delta} p^T \theta$$

Fundamental Thm. Lin. Prog.  
(Dantzig et al, 55)



$N = 2$

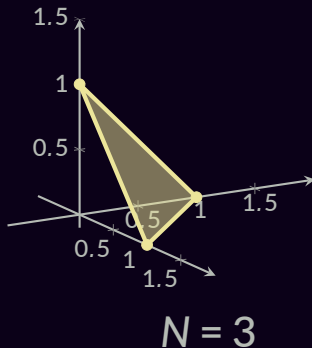
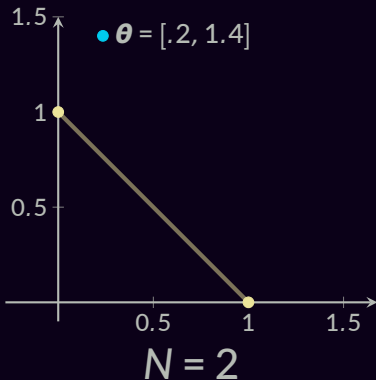


$N = 3$

# Variational Form of Argmax

$$\max_j \theta_j = \max_{p \in \Delta} p^T \theta$$

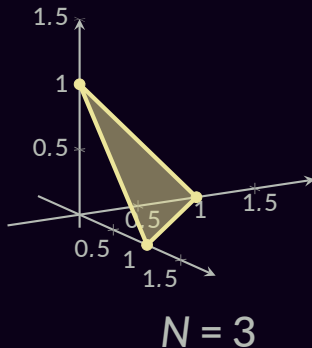
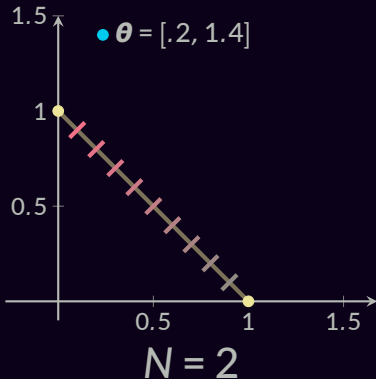
Fundamental Thm. Lin. Prog.  
(Dantzig et al, 55)



# Variational Form of Argmax

$$\max_j \theta_j = \max_{p \in \Delta} p^T \theta$$

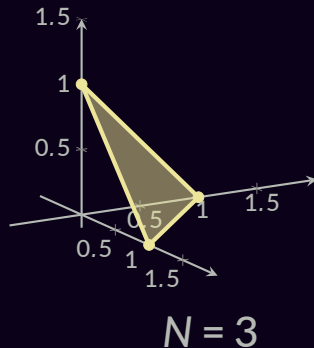
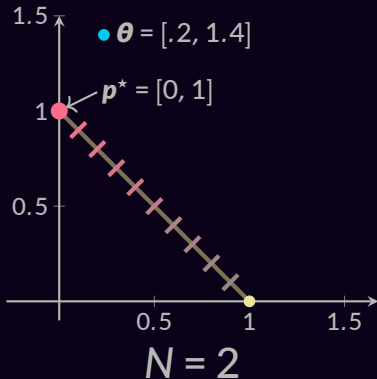
Fundamental Thm. Lin. Prog.  
(Dantzig et al, 55)



# Variational Form of Argmax

$$\max_j \theta_j = \max_{p \in \Delta} p^T \theta$$

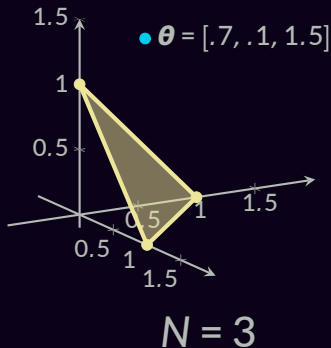
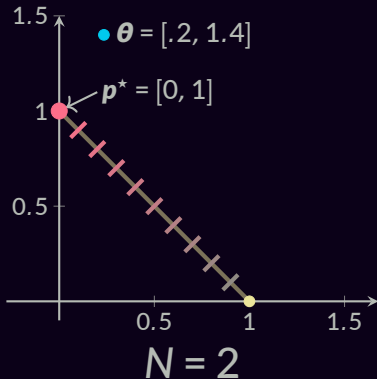
Fundamental Thm. Lin. Prog.  
(Dantzig et al, 55)



# Variational Form of Argmax

$$\max_j \theta_j = \max_{p \in \Delta} p^T \theta$$

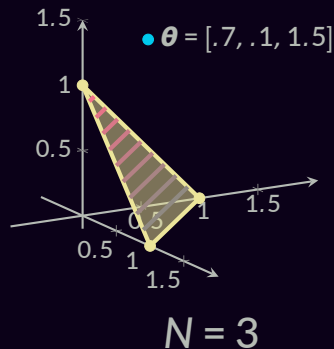
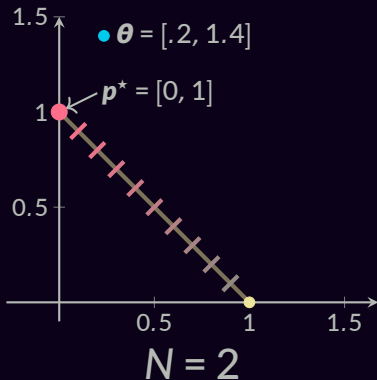
Fundamental Thm. Lin. Prog.  
(Dantzig et al, 55)



# Variational Form of Argmax

$$\max_j \theta_j = \max_{p \in \Delta} p^T \theta$$

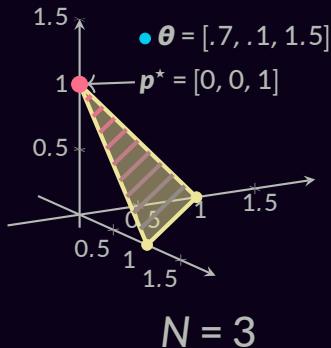
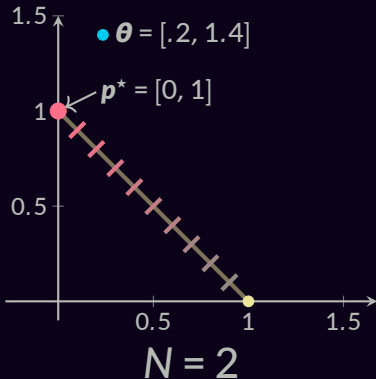
Fundamental Thm. Lin. Prog.  
(Dantzig et al, 55)



# Variational Form of Argmax

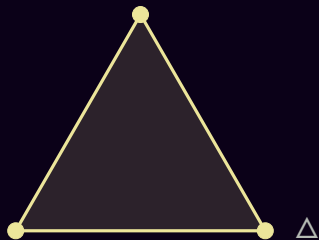
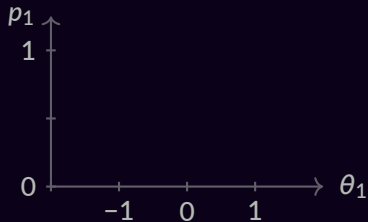
$$\max_j \theta_j = \max_{p \in \Delta} p^T \theta$$

Fundamental Thm. Lin. Prog.  
(Dantzig et al, 55)



# Smoothed Max Operators

$$\max_{\Omega}(\boldsymbol{\theta}) = \max_{\mathbf{p} \in \Delta} \mathbf{p}^{\top} \boldsymbol{\theta} - \Omega(\mathbf{p})$$

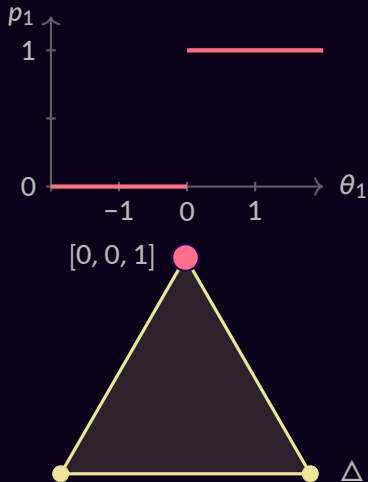




# Smoothed Max Operators

$$\max_{\Omega}(\boldsymbol{\theta}) = \max_{\mathbf{p} \in \Delta} \mathbf{p}^{\top} \boldsymbol{\theta} - \Omega(\mathbf{p})$$

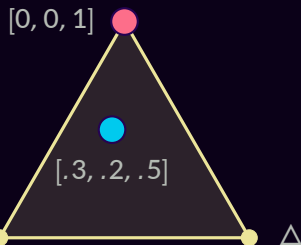
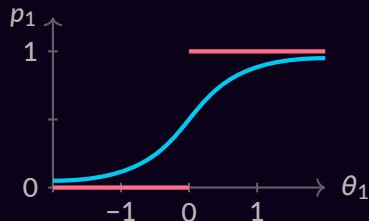
- argmax:  $\Omega(\mathbf{p}) = 0$



# Smoothed Max Operators

$$\max_{\Omega}(\boldsymbol{\theta}) = \max_{\mathbf{p} \in \Delta} \mathbf{p}^{\top} \boldsymbol{\theta} - \Omega(\mathbf{p})$$

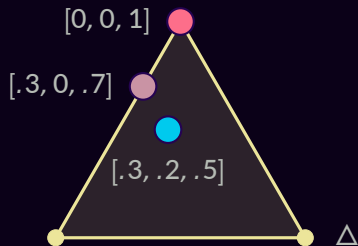
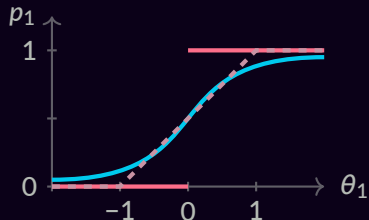
- argmax:  $\Omega(\mathbf{p}) = 0$
- softmax:  $\Omega(\mathbf{p}) = \sum_j p_j \log p_j$

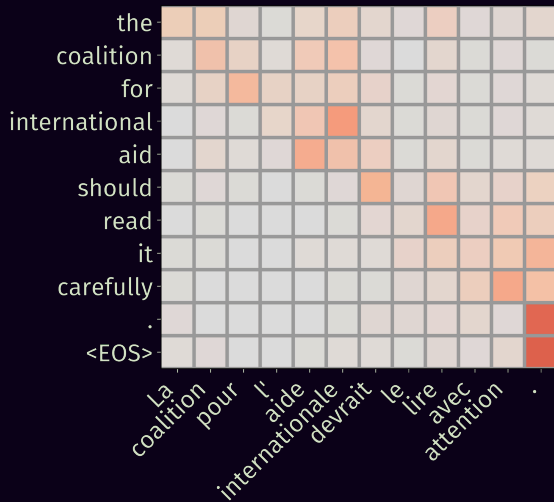


# Smoothed Max Operators

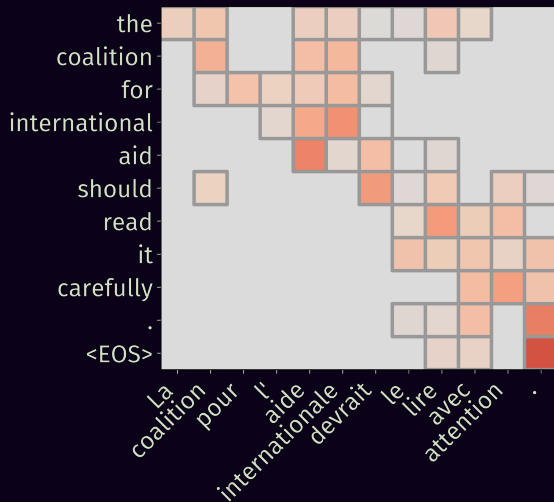
$$\max_{\Omega}(\boldsymbol{\theta}) = \max_{\mathbf{p} \in \Delta} \mathbf{p}^{\top} \boldsymbol{\theta} - \Omega(\mathbf{p})$$

- argmax:  $\Omega(\mathbf{p}) = 0$
- softmax:  $\Omega(\mathbf{p}) = \sum_j p_j \log p_j$
- sparsemax:  $\Omega(\mathbf{p}) = 1/2 \|\mathbf{p}\|_2^2$

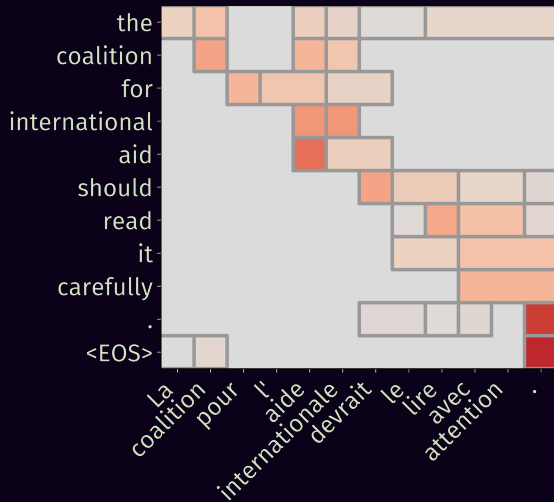




softmax



sparsemax

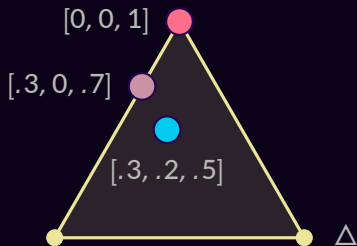
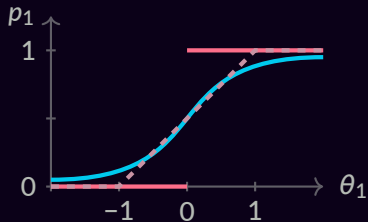


fusedmax ?!

# Smoothed Max Operators

$$\max_{\Omega}(\boldsymbol{\theta}) = \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \boldsymbol{\theta} - \Omega(\mathbf{p})$$

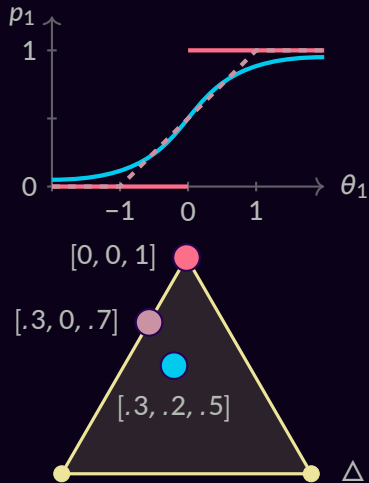
- argmax:  $\Omega(\mathbf{p}) = 0$
- softmax:  $\Omega(\mathbf{p}) = \sum_j p_j \log p_j$
- sparsemax:  $\Omega(\mathbf{p}) = 1/2 \|\mathbf{p}\|_2^2$



# Smoothed Max Operators

$$\max_{\Omega}(\boldsymbol{\theta}) = \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \boldsymbol{\theta} - \Omega(\mathbf{p})$$

- argmax:  $\Omega(\mathbf{p}) = 0$
- softmax:  $\Omega(\mathbf{p}) = \sum_j p_j \log p_j$
- sparsemax:  $\Omega(\mathbf{p}) = 1/2 \|\mathbf{p}\|_2^2$
- fusedmax:  $\Omega(\mathbf{p}) = 1/2 \|\mathbf{p}\|_2^2 + \sum_j |p_j - p_{j-1}|$
- oscarmax:  $\Omega(\mathbf{p}) = 1/2 \|\mathbf{p}\|_2^2 + \sum_{i,j} \max(p_i, p_j)$



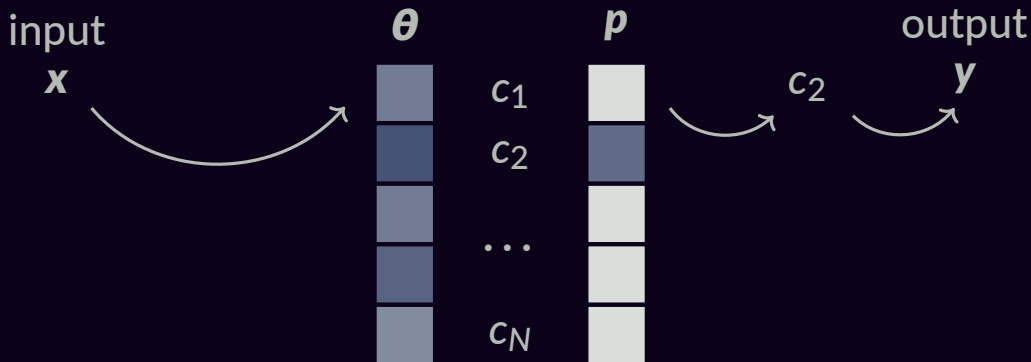


# Structured Inference

finally

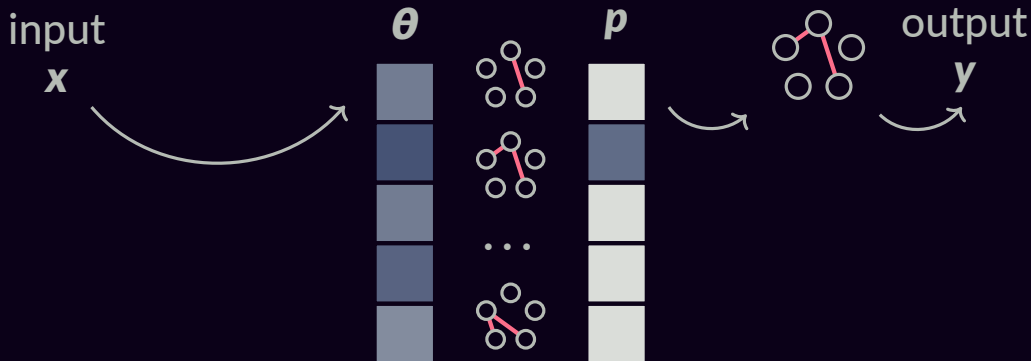
# Structured Inference

is essentially a (very high-dimensional) argmax



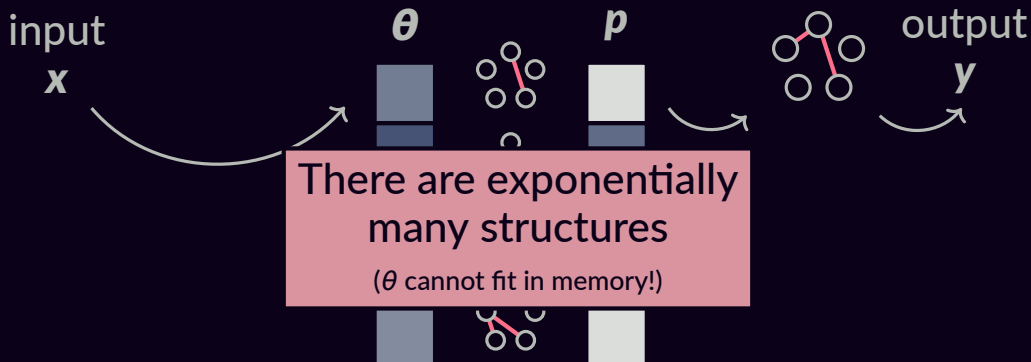
# Structured Inference

is essentially a (very high-dimensional) argmax



# Structured Inference

is essentially a (very high-dimensional) argmax



# Factorization Into Parts

$$\theta = A^{\top} \eta$$

# Factorization Into Parts

$$\theta = A^T \eta$$

★ dog on wheels

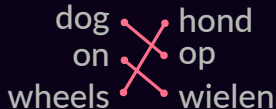


★→dog	1	0	0	$\eta = \begin{bmatrix} .1 \\ .2 \\ -.1 \\ .3 \\ .8 \\ .1 \\ -.3 \\ .2 \\ -.1 \end{bmatrix}$		
on→dog	0	1	1			
wheels→dog	0	0	0			
★→on	0	1	1			
dog→on	1	...	0		0	...
wheels→on	0	0	0		0	
★→wheels	0	0	0		0	
dog→wheels	0	1	0		0	
on→wheels	1	0	1		1	

$A =$

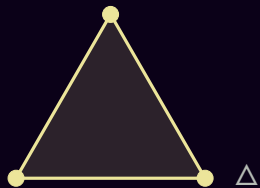
# Factorization Into Parts

$$\theta = A^T \eta$$

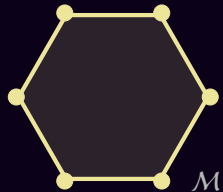
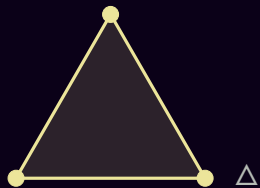


$$A = \begin{array}{l} \star \rightarrow \text{dog} \\ \text{on} \rightarrow \text{dog} \\ \text{wheels} \rightarrow \text{dog} \\ \hline \star \rightarrow \text{on} \\ \text{dog} \rightarrow \text{on} \\ \text{wheels} \rightarrow \text{on} \\ \hline \star \rightarrow \text{wheels} \\ \text{dog} \rightarrow \text{wheels} \\ \text{on} \rightarrow \text{wheels} \end{array} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \\ \hline 0 & 1 & 1 \\ 1 & \dots & 0 & 0 & \dots \\ 0 & 0 & 0 \\ \hline 0 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \eta = \begin{bmatrix} .1 \\ .2 \\ -.1 \\ \hline .3 \\ .8 \\ .1 \\ \hline -.3 \\ .2 \\ -.1 \end{bmatrix}$$

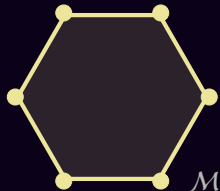
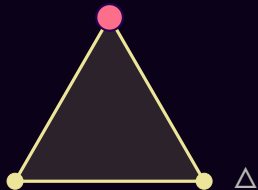
$$A = \begin{array}{l} \text{dog} - \text{hond} \\ \text{dog} - \text{op} \\ \text{dog} - \text{wielen} \\ \hline \text{on} - \text{hond} \\ \text{on} - \text{op} \\ \text{on} - \text{wielen} \\ \hline \text{wheels} - \text{hond} \\ \text{wheels} - \text{op} \\ \text{wheels} - \text{wielen} \end{array} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \\ \hline 0 & 0 & 0 \\ 1 & \dots & 0 & 0 & \dots \\ 0 & 1 & 1 \\ \hline 0 & 1 & 0 \\ 0 & 0 & 0 \\ 1 & 0 & 1 \end{bmatrix} \eta = \begin{bmatrix} .1 \\ .2 \\ -.1 \\ \hline .3 \\ .8 \\ .1 \\ \hline -.3 \\ .2 \\ -.1 \end{bmatrix}$$



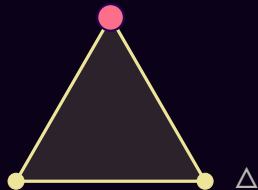




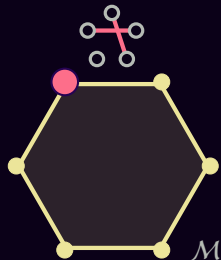
●  $\operatorname{argmax}_{p \in \Delta} p^\top \theta$



• **argmax**  $\arg \max_{p \in \Delta} p^\top \theta$



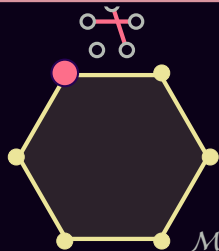
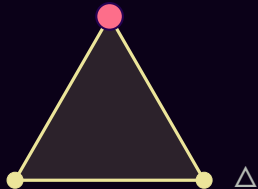
• **MAP**  $\arg \max_{\mu \in \mathcal{M}} \mu^\top \eta$



•  $\operatorname{argmax}_{p \in \Delta} p^T \theta$

•  $\operatorname{MAP}_{\mu \in \mathcal{M}} \mu^T \eta$

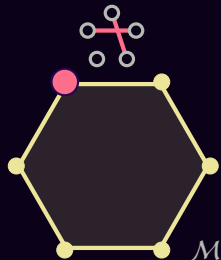
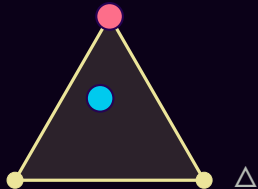
e.g. dependency parsing  $\rightarrow$  max. spanning tree  
matching  $\rightarrow$  **the Hungarian algorithm**



● **argmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \boldsymbol{\theta}$

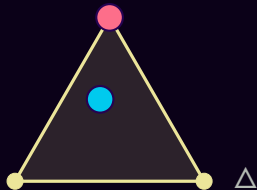
● **softmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \boldsymbol{\theta} + H(\mathbf{p})$

● **MAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta}$



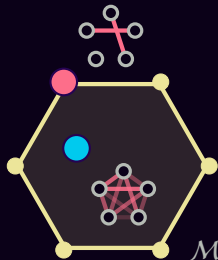
● **argmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \boldsymbol{\theta}$

● **softmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \boldsymbol{\theta} + H(\mathbf{p})$



● **MAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta}$

● **marginals**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} + \tilde{H}(\boldsymbol{\mu})$



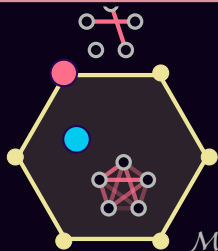
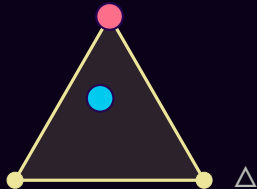
- **argmax**  $\arg \max_{p \in \Delta} p^\top \theta$

- **softmax**  $\arg \max_{p \in \Delta} p^\top \theta + H(p)$

- **MAP**  $\arg \max_{\mu \in \mathcal{M}} \mu^\top \eta$

- **marginals**  $\arg \max_{\mu \in \mathcal{M}} \mu^\top \eta + \tilde{H}(\mu)$

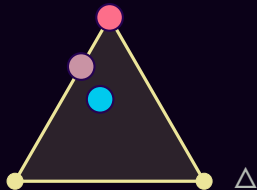
e.g. dependency parsing → **the Matrix-Tree theorem**  
 matching → **#P-complete!** (Valiant, 79)



● **argmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \boldsymbol{\theta}$

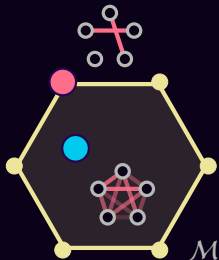
● **softmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \boldsymbol{\theta} + H(\mathbf{p})$

● **sparsemax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \boldsymbol{\theta} - 1/2 \|\mathbf{p}\|^2$



● **MAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta}$

● **marginals**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} + \tilde{H}(\boldsymbol{\mu})$

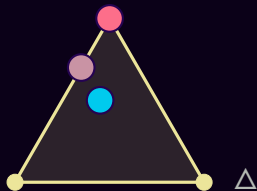




● **argmax**  $\arg \max_{p \in \Delta} p^\top \theta$

● **softmax**  $\arg \max_{p \in \Delta} p^\top \theta + H(p)$

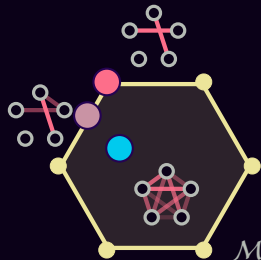
● **sparsemax**  $\arg \max_{p \in \Delta} p^\top \theta - 1/2 \|p\|^2$



● **MAP**  $\arg \max_{\mu \in \mathcal{M}} \mu^\top \eta$

● **marginals**  $\arg \max_{\mu \in \mathcal{M}} \mu^\top \eta + \tilde{H}(\mu)$

● **SparseMAP**  $\arg \max_{\mu \in \mathcal{M}} \mu^\top \eta - 1/2 \|\mu\|^2$



# SparseMAP Inference Solution

$$\boldsymbol{\mu}^{\star} = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

$$= \begin{array}{c} \circ \\ \circ \text{---} \circ \\ \circ \end{array} = .6 \begin{array}{c} \circ \\ \circ \text{---} \circ \\ \circ \end{array} + .4 \begin{array}{c} \circ \\ \circ \text{---} \circ \\ \circ \end{array}$$

$$= \mathbf{A} \mathbf{p}^{\star} \text{ with very sparse } \mathbf{p}^{\star} \in \Delta^N$$

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

**Greedy Conditional Gradient  
(Frank-Wolfe) algorithms**

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

**Greedy Conditional Gradient  
(Frank-Wolfe) algorithms**

- ▶ select a new corner of  $\mathcal{M}$

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

## Greedy Conditional Gradient (Frank-Wolfe) algorithms

- ▶ select a new corner of  $\mathcal{M}$
- ▶ update the (sparse) coefficients of  $\boldsymbol{p}$

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

## Greedy Conditional Gradient (Frank-Wolfe) algorithms

- ▶ select a new corner of  $\mathcal{M}$
- ▶ update the (sparse) coefficients of  $\boldsymbol{p}$ 
  - ▶ Update rules: vanilla, away-step, pairwise

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

## Greedy Conditional Gradient (Frank-Wolfe) algorithms

- ▶ select a new corner of  $\mathcal{M}$
- ▶ update the (sparse) coefficients of  $\boldsymbol{p}$ 
  - ▶ Update rules: vanilla, away-step, pairwise
  - ▶ Quadratic objective:  
**Active Set** (Min-Norm Point)



# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

## Greedy Conditional Gradient (Frank-Wolfe)

- ▶ select a new corner
- ▶ update the (sparse,  $\boldsymbol{\mu}^{\star}$ )
  - ▶ Update rules: vanilla, away-step, pairwise
  - ▶ Quadratic objective:  
**Active Set** (Min-Norm Point)

Active Set achieves  
**finite & linear** convergence!

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

## Greedy Conditional Gradient (Frank-Wolfe) algorithms

- ▶ select a new corner of  $\mathcal{M}$
- ▶ update the (sparse) coefficients of  $\boldsymbol{p}$ 
  - ▶ Update rules: vanilla, away-step, pairwise
  - ▶ Quadratic objective:  
**Active Set** (Min-Norm Point)

## Backward pass

$\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}$  is sparse

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

## Greedy Conditional Gradient (Frank-Wolfe) algorithms

- ▶ select a new corner of  $\mathcal{M}$
- ▶ update the (sparse) coefficients of  $\boldsymbol{p}$ 
  - ▶ Update rules: vanilla, away-step, pairwise
  - ▶ Quadratic objective:  
**Active Set** (Min-Norm Point)

## Backward pass

$\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}$  is sparse

computing  $\left(\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}\right)^{\top} d\boldsymbol{y}$   
takes  $O(\dim(\boldsymbol{\mu}) \text{nnz}(\boldsymbol{p}^{\star}))$

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^{\star} = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

Greedy Con  
(Frank-Wolfe)

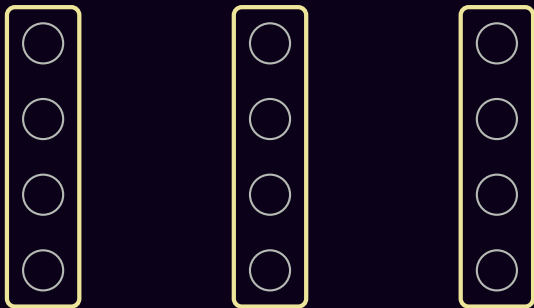
Completely modular: just add MAP pass

- ▶ select a new corner of  $\mathcal{M}$
- ▶ update the (sparse) coefficients of  $\boldsymbol{p}$ 
  - ▶ Update rules: vanilla, away-step, pairwise
  - ▶ Quadratic objective:  
**Active Set** (Min-Norm Point)

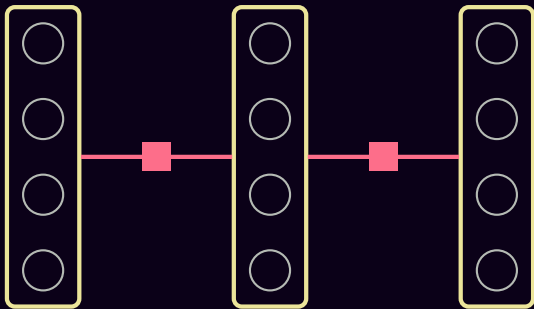
$\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}$  is sparse

computing  $\left(\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}\right)^{\top} d\boldsymbol{y}$   
takes  $O(\dim(\boldsymbol{\mu}) \text{nnz}(\boldsymbol{p}^{\star}))$

# Structured Attention & Graphical Models



# Structured Attention & Graphical Models



# Structured Attention for Alignments

NLI

premise: A gentleman overlooking a neighborhood situation.  
hypothesis: A police officer watches a situation closely.

input

(P, H)

	A	A	
	gentleman	police	
	overlooking	officer	
	...	...	
	situation	closely	

output



entails

contradicts

neutral

(Model: ESIM)

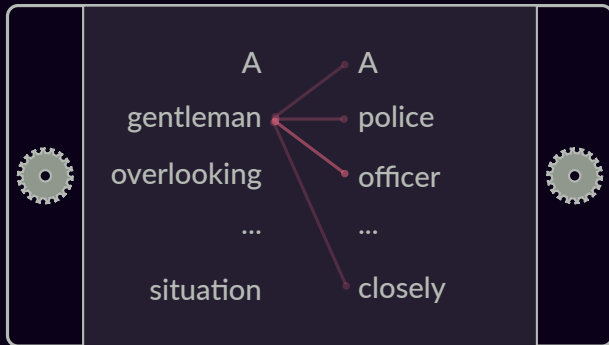
# Structured Attention for Alignments

NLI

premise: A gentleman overlooking a neighborhood situation.  
hypothesis: A police officer watches a situation closely.

input

(P, H)



(Model: ESIM)

output



entails

contradicts

neutral



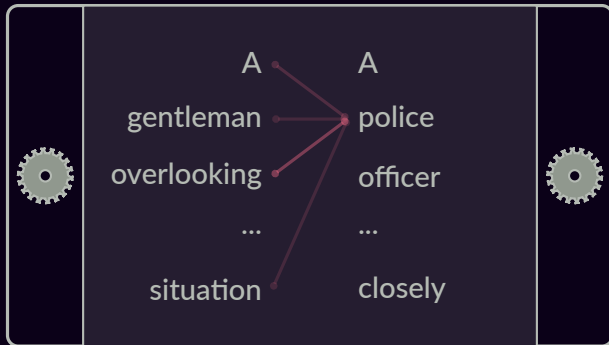
# Structured Attention for Alignments

NLI

premise: A gentleman overlooking a neighborhood situation.  
hypothesis: A police officer watches a situation closely.

input

(P, H)



(Model: ESIM)

output



entails

contradicts

neutral

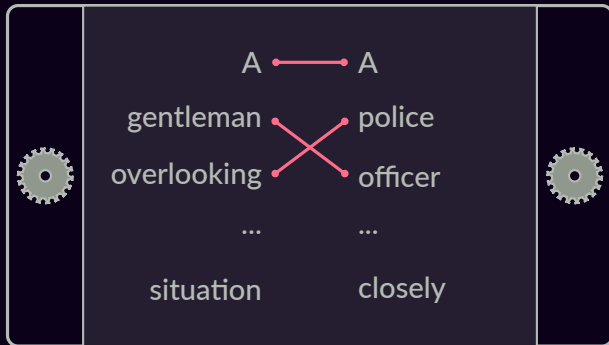
# Structured Attention for Alignments

NLI

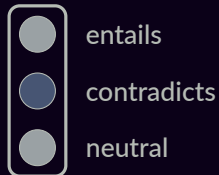
premise: A gentleman overlooking a neighborhood situation.  
hypothesis: A police officer watches a situation closely.

input

(P, H)

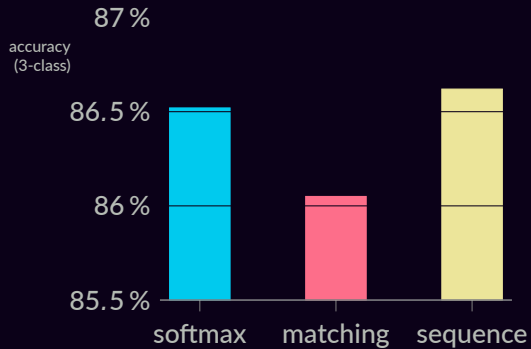


output

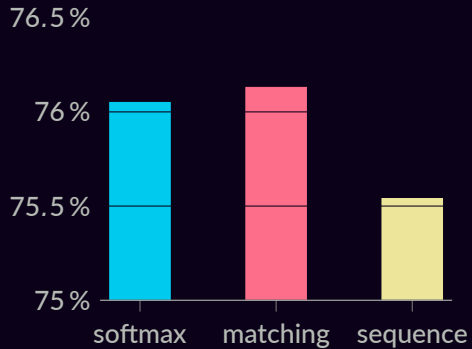


(Proposed model: global matching)

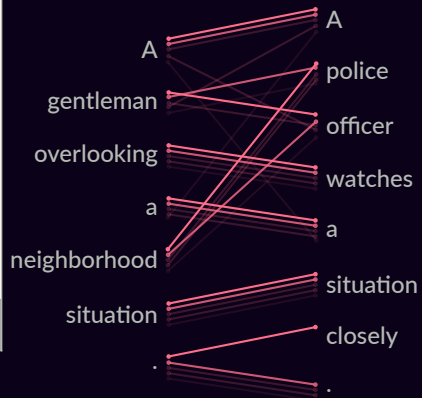
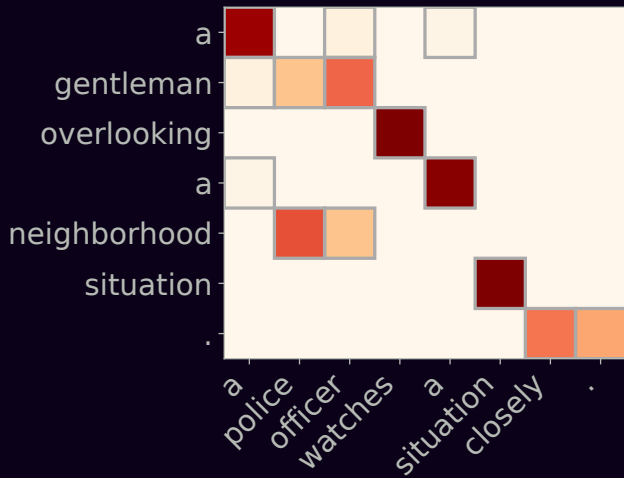
## SNLI



## MultiNLI








**Dynamically inferring  
the computation graph**

# Dependency TreeLSTM

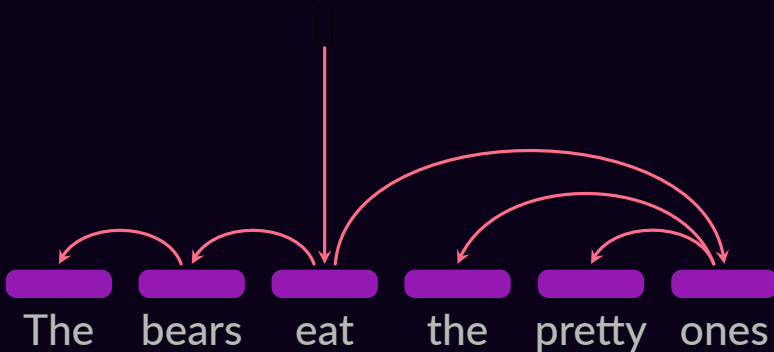
(Tai & al, 15)



The bears eat the pretty ones

# Dependency TreeLSTM

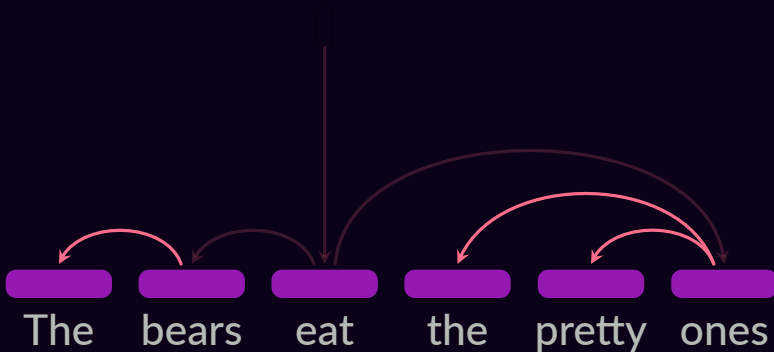
(Tai & al, 15)





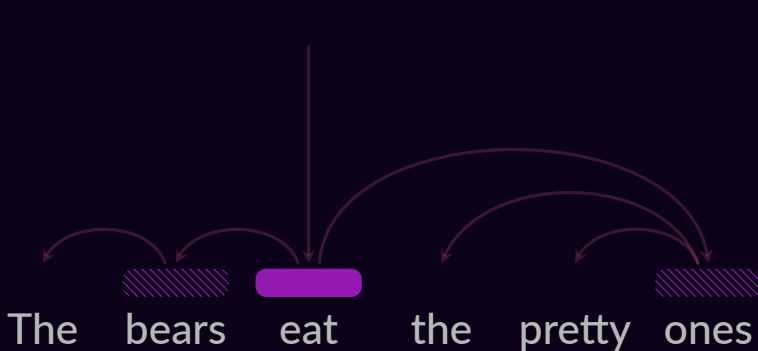
# Dependency TreeLSTM

(Tai & al, 15)



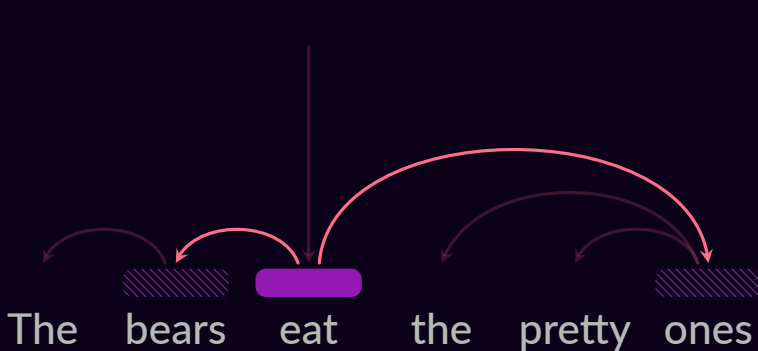
# Dependency TreeLSTM

(Tai & al, 15)



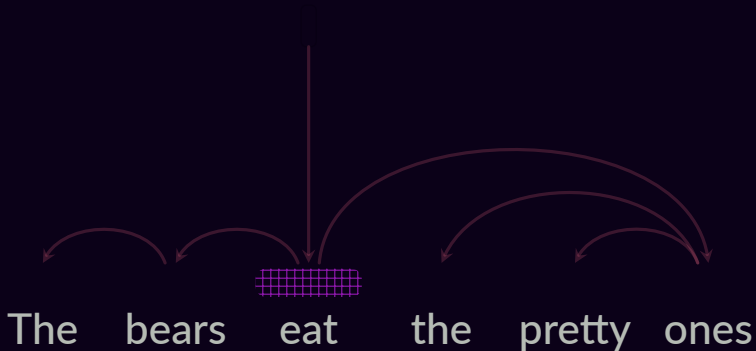
# Dependency TreeLSTM

(Tai & al, 15)



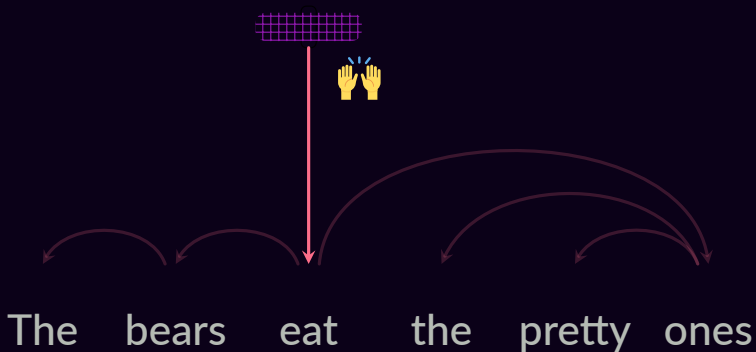
# Dependency TreeLSTM

(Tai & al, 15)



# Dependency TreeLSTM

(Tai & al, 15)

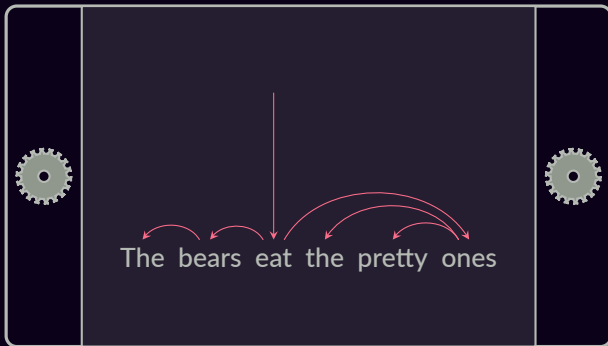


# Latent Dependency TreeLSTM

(Niculae, Martins, Cardie, 18)

input

$x$



output

$y$

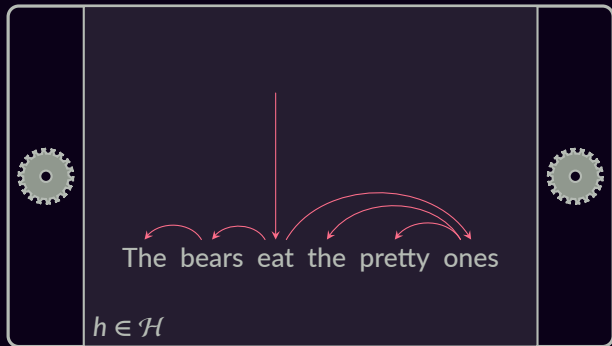
# Latent Dependency TreeLSTM

(Niculae, Martins, Cardie, 18)

$$p(y|x) = \sum_{h \in \mathcal{H}} p(y | h, x) p(h | x)$$

input

$x$



# Structured Latent Variable Models

$$p(y | x) = \sum_{h \in \mathcal{H}} p(y | h, x) p(h | x)$$




# Structured Latent Variable Models

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\phi}(y \mid h, x) p_{\pi}(h \mid x)$$

# Structured Latent Variable Models

$$p(y | x) = \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x)$$

e.g., a TreeLSTM defined by  $h$



# Structured Latent Variable Models

$$p(y | x) = \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x)$$

e.g., a TreeLSTM defined by  $h$

parsing model,  
using some score  $\pi(h; x)$

# Structured Latent Variable Models

sum over  
all possible trees

e.g., a TreeLSTM defined by  $h$

$$p(y | x) = \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x)$$

parsing model,  
using some score  $\pi(h; x)$

Exponentially large sum!

# Structured Latent Variable Models

sum over  
all possible trees

e.g., a TreeLSTM defined by  $h$

$$p(y | x) = \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x)$$

parsing model,  
using some score  $\pi(h; x)$

How to define  $p_{\pi}$ ?

idea 1

idea 2

idea 3

# Structured Latent Variable Models

sum over  
all possible trees

e.g., a TreeLSTM defined by  $h$

$$p(y | x) = \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x)$$

parsing model,  
using some score  $\pi(h; x)$

How to define  $p_{\pi}$ ?

$$\sum_{h \in \mathcal{H}}$$

idea 1

idea 2

idea 3

# Structured Latent Variable Models

sum over  
all possible trees

e.g., a TreeLSTM defined by  $h$

$$p(y | x) = \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x)$$

parsing model,  
using some score  $\pi(h; x)$

How to define  $p_{\pi}$ ?

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \pi}$$

idea 1

idea 2

idea 3

# Structured Latent Variable Models

sum over  
all possible trees

e.g., a TreeLSTM defined by  $h$

$$p(y | x) = \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x)$$

parsing model,  
using some score  $\pi(h; x)$

How to define  $p_{\pi}$ ?

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \pi}$$

idea 1  $p_{\pi}(h | x) = 1$  if  $h = h^*$  else 0

argmax

idea 2

idea 3



# Structured Latent Variable Models

sum over  
all possible trees

e.g., a TreeLSTM defined by  $h$

$$p(y | x) = \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x)$$

parsing model,  
using some score  $\pi(h; x)$


How to define  $p_{\pi}$ ?

idea 1  $p_{\pi}(h | x) = 1$  if  $h = h^*$  else 0

argmax

idea 2

idea 3

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \pi}$$


# Structured Latent Variable Models

sum over  
all possible trees

e.g., a TreeLSTM defined by  $h$

$$p(y | x) = \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x)$$

parsing model,  
using some score  $\pi(h; x)$

How to define  $p_{\pi}$ ?

idea 1  $p_{\pi}(h | x) = 1$  if  $h = h^*$  else 0

argmax

idea 2

idea 3

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \pi}$$



# Structured Latent Variable Models

sum over  
all possible trees

e.g., a TreeLSTM defined by  $h$

$$p(y | x) = \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x)$$

parsing model,  
using some score  $\pi(h; x)$

How to define  $p_{\pi}$ ?

idea 1  $p_{\pi}(h | x) = 1$  if  $h = h^*$  else 0

argmax

idea 2  $p_{\pi}(h | x) \propto \exp(\text{score}_{\pi}(h; x))$

softmax

idea 3

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \pi}$$



# Structured Latent Variable Models

sum over  
all possible trees

e.g., a TreeLSTM defined by  $h$

$$p(y | x) = \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x)$$

parsing model,  
using some score  $\pi(h; x)$

How to define  $p_{\pi}$ ?

idea 1  $p_{\pi}(h | x) = 1$  if  $h = h^*$  else 0

argmax

idea 2  $p_{\pi}(h | x) \propto \exp(\text{score}_{\pi}(h; x))$

softmax

idea 3

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \pi}$$

idea 1: 😊 (smiley face)  
idea 2: 😞 (frowny face)  
idea 3: 😊 (smiley face)

# Structured Latent Variable Models

sum over  
all possible trees

e.g., a TreeLSTM defined by  $h$

$$p(y | x) = \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x)$$

parsing model,  
using some score  $\pi(h; x)$

How to define  $p_{\pi}$ ?

idea 1  $p_{\pi}(h | x) = 1$  if  $h = h^*$  else 0

argmax

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \pi}$$



idea 2  $p_{\pi}(h | x) \propto \exp(\text{score}_{\pi}(h; x))$

softmax



idea 3

# Structured Latent Variable Models

sum over  
all possible trees

e.g., a TreeLSTM defined by  $h$

$$p(y | x) = \sum_{h \in \mathcal{H}} p_{\phi}(y | h, x) p_{\pi}(h | x)$$

parsing model,  
using some score  $\pi(h; x)$

How to define  $p_{\pi}$ ?

idea 1  $p_{\pi}(h | x) = 1$  if  $h = h^*$  else 0

argmax

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \pi}$$



idea 2  $p_{\pi}(h | x) \propto \exp(\text{score}_{\pi}(h; x))$

softmax



idea 3

SparseMAP



# SparseMAP Inference



# SparseMAP Inference

$$\text{Diagram} = .7$$

$$\text{Diagram} + .3$$

$$\text{Diagram} + 0 \text{Diagram} + \dots$$



# SparseMAP Inference

$$\begin{aligned}
 & \text{Diagram 1} = .7 \quad \text{Diagram 2} + .3 \quad \text{Diagram 3} + 0 \text{Diagram 4} + \dots \\
 p(y | x) = & .7 p_{\phi}(y | \text{Diagram 1}) + .3 p_{\phi}(y | \text{Diagram 2})
 \end{aligned}$$

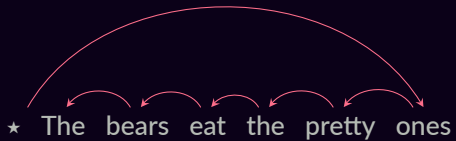
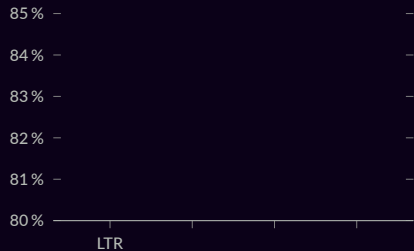
The diagrams consist of three red dots arranged horizontally. In the first diagram, a red curved arrow points from the left dot to the middle dot, and another red curved arrow points from the middle dot to the right dot. In the second diagram, a red curved arrow points from the left dot to the right dot, and a red curved arrow points from the middle dot to the right dot. In the third diagram, a red curved arrow points from the left dot to the right dot, and a red curved arrow points from the middle dot to the left dot. In the fourth diagram, a red curved arrow points from the left dot to the right dot, and a red curved arrow points from the middle dot to the left dot.

# SparseMAP Inference

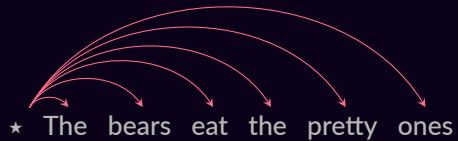
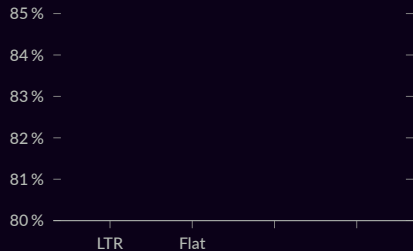
$$\begin{aligned}
 \text{Diagram 1} &= .7 & \text{Diagram 2} &+ .3 & \text{Diagram 3} &+ 0 & \text{Diagram 4} &+ \dots \\
 p(y | x) &= .7 p_{\phi}(y | \text{Diagram 1}) + .3 p_{\phi}(y | \text{Diagram 2})
 \end{aligned}$$

Diagram 1 is not a tree itself:  $p(y | x) \neq p_{\phi}(y | \text{Diagram 1})$ !

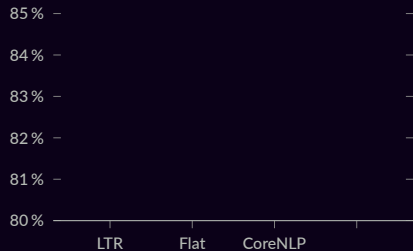




Left-to-right: regular LSTM



Flat: bag-of-words-like



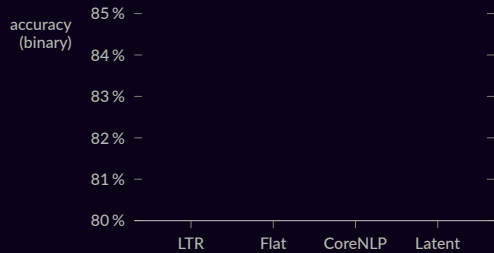
★ The bears eat the pretty ones

Diagram illustrating the dependency arcs for the sentence "The bears eat the pretty ones". The arcs show the following dependencies: "The" to "eat", "bears" to "eat", "eat" to "the", "eat" to "pretty", "eat" to "ones", "the" to "pretty", and "pretty" to "ones".

CoreNLP: off-line parser

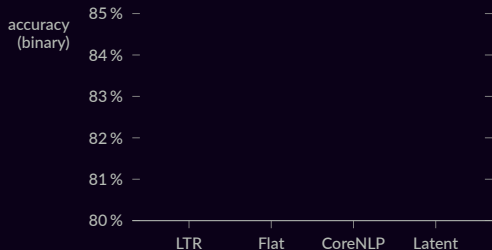


## Sentiment classification (SST)

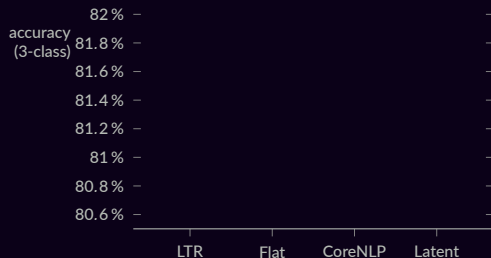




## Sentiment classification (SST)



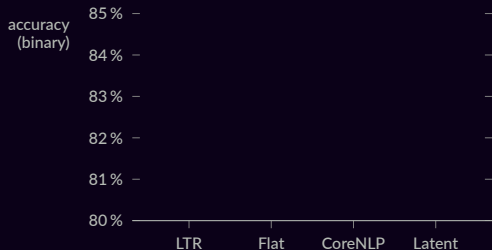
## Natural Language Inference (SNLI)



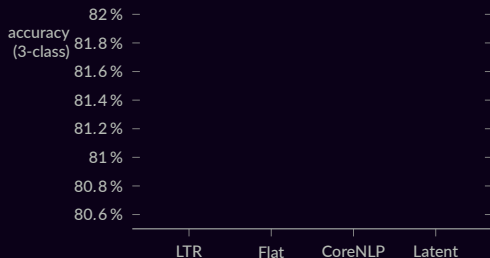
## Sentence pair classification ( $P, H$ )

$$p(y \mid P, H) = \sum_{h_P \in \mathcal{H}(P)} \sum_{h_H \in \mathcal{H}(H)} p_{\phi}(y \mid h_P, h_H) p_{\pi}(h_P \mid P) p_{\pi}(h_H \mid H)$$

## Sentiment classification (SST)



## Natural Language Inference (SNLI)

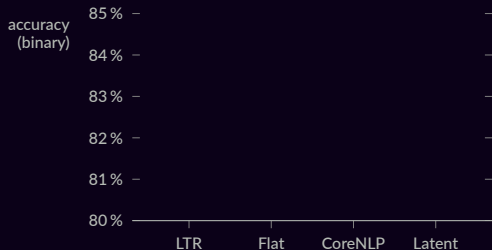


## Reverse dictionary lookup

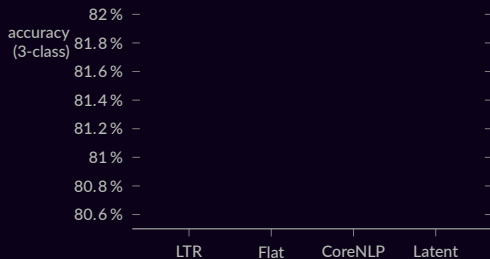
given word description, predict word embedding (Hill et al, 17)

instead of  $p(y | x)$ , we model  $\mathbb{E}_{p_{\pi}} \mathbf{g}(x) = \sum_{h \in \mathcal{H}} \mathbf{g}(x; h) p_{\pi}(h | x)$

## Sentiment classification (SST)

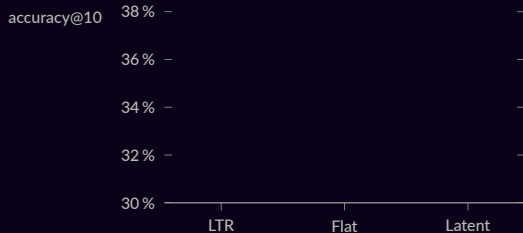


## Natural Language Inference (SNLI)

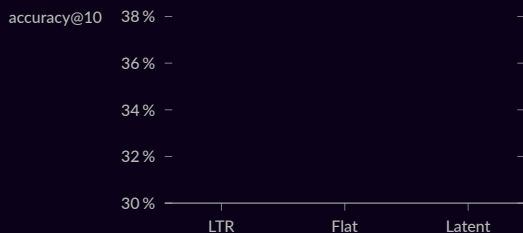


## Reverse dictionary lookup

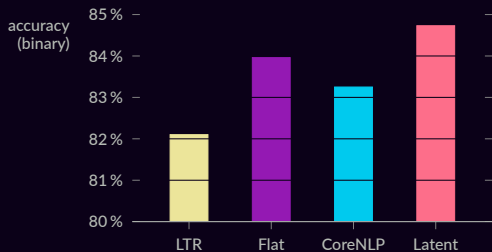
(definitions)



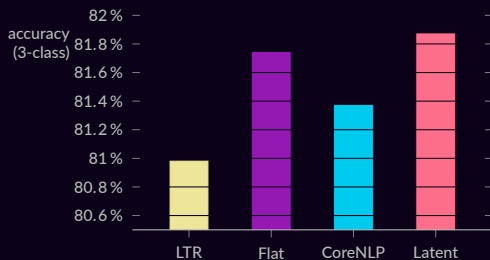
(concepts)



### Sentiment classification (SST)

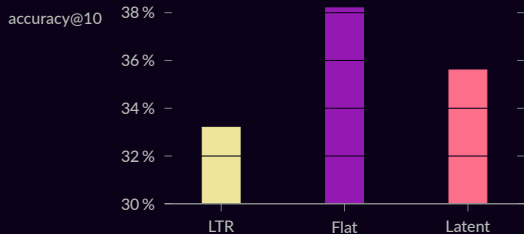


### Natural Language Inference (SNLI)

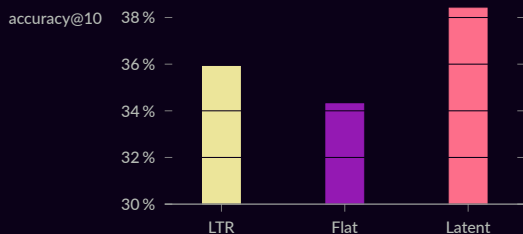


### Reverse dictionary lookup

(definitions)

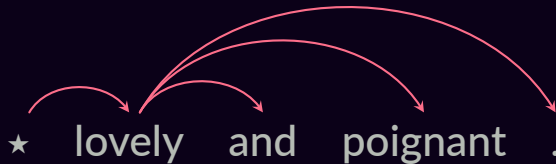


(concepts)



# Syntax vs. Composition Order

CoreNLP parse,  $p = 21.4\%$

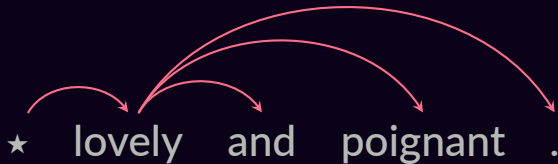


# Syntax vs. Composition Order

$p = 22.6\%$

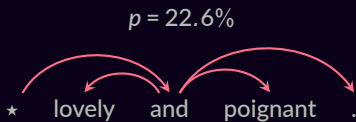


CoreNLP parse,  $p = 21.4\%$

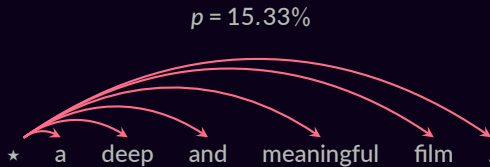


...

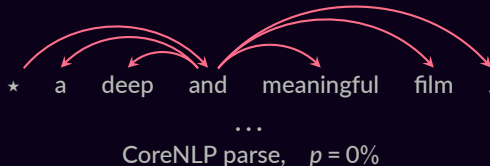
# Syntax vs. Composition Order



CoreNLP parse,  $p = 21.4\%$



$p = 15.27\%$



# Conclusions

Differentiable & sparse  
structured inference

Generic, extensible algorithms

Interpretable structured attention

Dynamically-inferred  
computation graphs

Catch us at EMNLP:

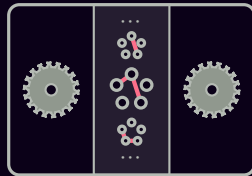
BlackboxNLP, Thursday 11:00 & EMNLP, Friday 15:36 (3B)

✉ vlad@vene.ro

🐙 [github.com/vene/sparsemap](https://github.com/vene/sparsemap)

🏠 <https://vene.ro>

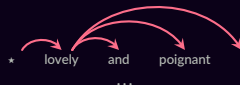
🐦 @vnfrombucharest



$p = 22.6\%$



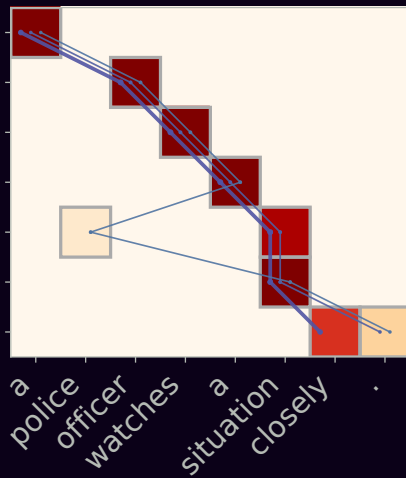
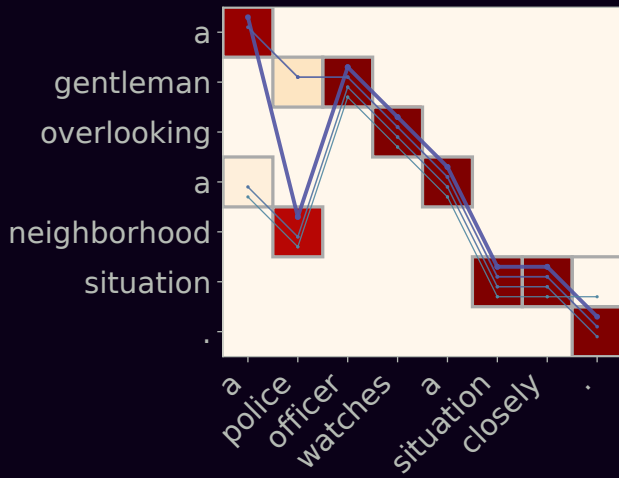
$p = 21.4\%$





**Extra slides**

Some icons by Dave Gandy and Freepik via [flaticon.com](http://flaticon.com).



# Structured Output Prediction

SparseMAP

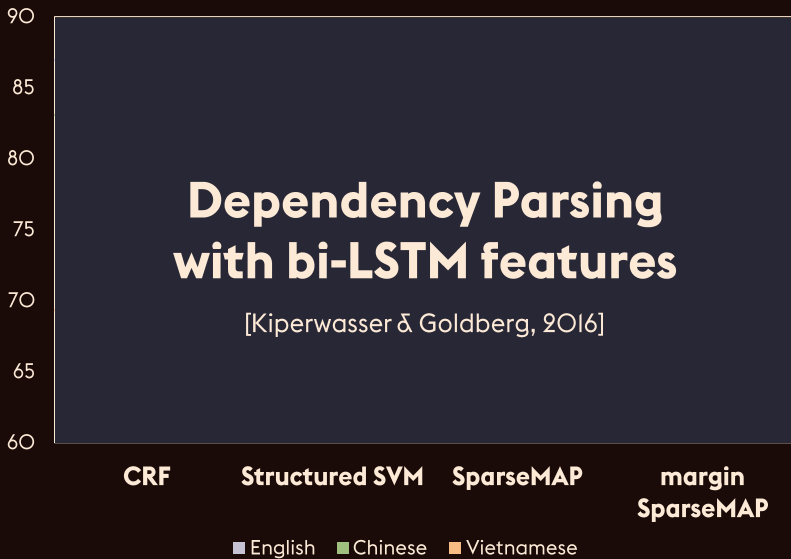
$$L_A(\boldsymbol{\eta}, \bar{\boldsymbol{\mu}}) = \max_{\boldsymbol{\mu} \in \mathcal{M}} \left\{ \boldsymbol{\eta}^\top \boldsymbol{\mu} - 1/2 \|\boldsymbol{\mu}\|^2 \right\} \\ - \boldsymbol{\eta}^\top \bar{\boldsymbol{\mu}} + 1/2 \|\bar{\boldsymbol{\mu}}\|^2$$

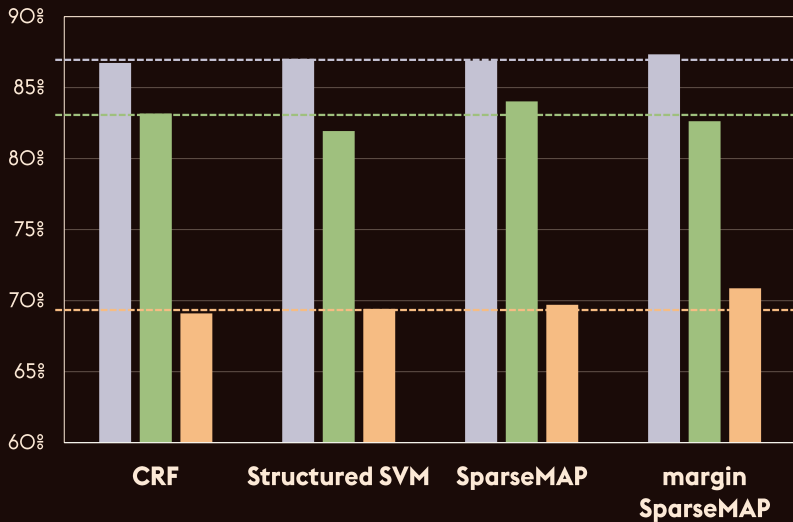
Instance of a structured Fenchel-Young loss, like CRF, SVM, etc. [Blondel, Martins, Niculae '18]

# Structured Output Prediction

$$\begin{aligned} \text{SparseMAP} \quad L_{\mathbf{A}}(\boldsymbol{\eta}, \bar{\boldsymbol{\mu}}) &= \max_{\boldsymbol{\mu} \in \mathcal{M}} \left\{ \boldsymbol{\eta}^{\top} \boldsymbol{\mu} - 1/2 \|\boldsymbol{\mu}\|^2 \right\} \\ &\quad - \boldsymbol{\eta}^{\top} \bar{\boldsymbol{\mu}} + 1/2 \|\bar{\boldsymbol{\mu}}\|^2 \\ \text{cost-SparseMAP} \quad L_{\mathbf{A}}^{\rho}(\boldsymbol{\eta}, \bar{\boldsymbol{\mu}}) &= \max_{\boldsymbol{\mu} \in \mathcal{M}} \left\{ \boldsymbol{\eta}^{\top} \boldsymbol{\mu} - 1/2 \|\boldsymbol{\mu}\|^2 + \rho(\boldsymbol{\mu}, \bar{\boldsymbol{\mu}}) \right\} \\ &\quad - \boldsymbol{\eta}^{\top} \bar{\boldsymbol{\mu}} + 1/2 \|\bar{\boldsymbol{\mu}}\|^2 \end{aligned}$$

Instance of a structured Fenchel-Young loss, like CRF, SVM, etc. [Blondel, Martins, Niculae '18]



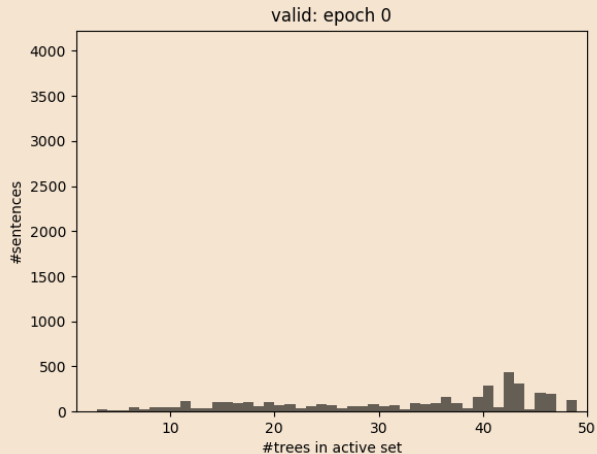
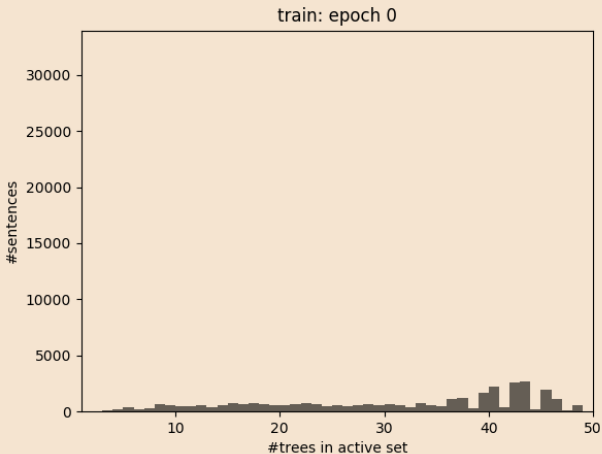


Unlabeled Accuracy (UAS)  
Universal Dependencies dataset

■ English ■ Chinese ■ Vietnamese

# Sparse Structured Output Prediction

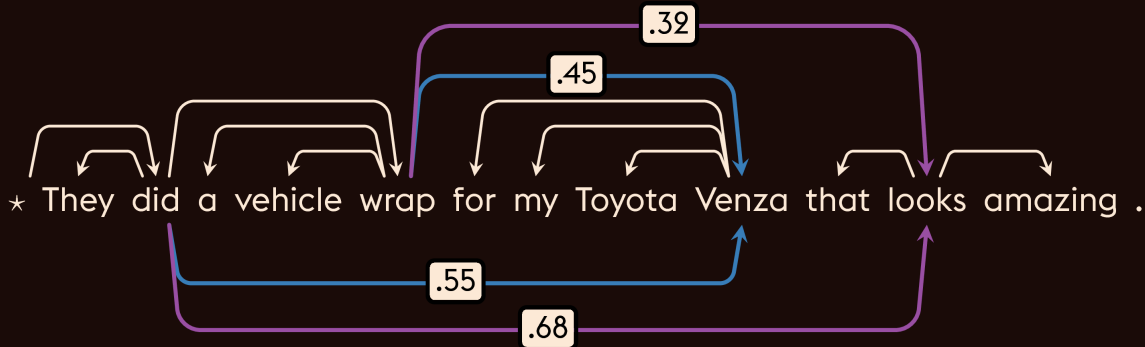
As models train, inference gets sparser!





# Sparse Structured Output Prediction

**Inference captures linguistic ambiguity!**



# Sparse Structured Output Prediction

**Inference captures linguistic ambiguity!**

