# SparseMAP:
# DIFFERENTIABLE SPARSE STRUCTURED INFERENCE

Presented by **Vlad Niculae**

Joint work with André FT Martins

Mathieu Blondel

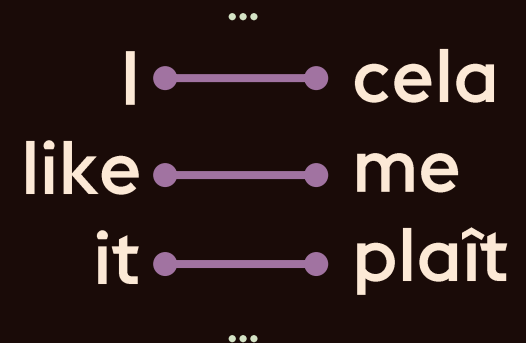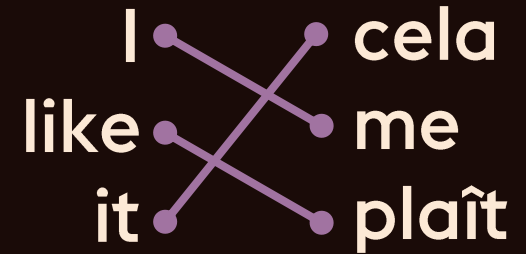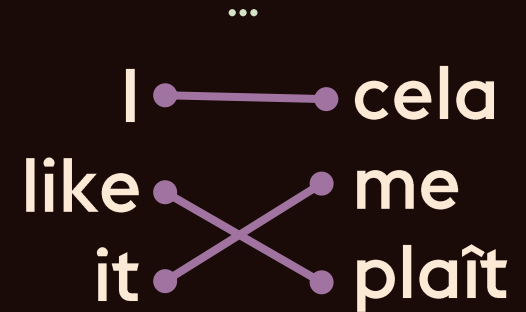Claire Cardie

poster #66 tonight
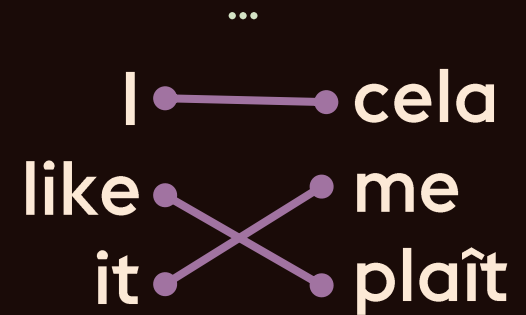
github.com/vene/sparsemap

# Structured Inference

...

* I like it

* I like it

...

* I like it

...

# Structured Inference

# Structured Inference

# Structured Inference

# (Latent) Structured Inference

input

output

# Deriving SparseMAP

# Structured Inference as argmax

input

# argmax



input

$x$   $p$

$\partial p\ /\ \partial x$ ?

# argmax → softmax
$$p_i = \exp x_i \,/\, Z$$

x    p

input

$$\partial p \,/\, \partial x \,?$$

$$x \overset{\in R^k}{=} A^\top \eta \overset{\in R^d}{} \qquad k \gg d$$

$$x = A^\top \eta$$



\* I like it

$$x = A^\top \eta$$



A=
| | | |
|---|---|---|
| 1 | O | O |
| O | 1 | 1 |
| O | O | O |
| O | 1 | 1 |
| 1 ... | O | O ... |
| O | O | O |
| O | O | O |
| O | 1 | O |
| 1 | O | 1 |

* → I
like → I
it → I
_____
* → like
I → like
it → like
_____
* → it
I → it
like → it

η=
| |
|---|
| .3 |
| .8 |
| -.5 |
| .2 |
| -.1 |
| -.2 |
| .7 |
| .6 |
| .1 |

$$x = A^\top \eta$$

$(0, 1, 0)$

$\mu = Ap$

$(1, 0, 0)$    $\Delta$

$M$

argmax $\langle x, p \rangle$
  s.t.  $p \in \Delta$

**MAP**

argmax $\langle \eta, \mu \rangle$
  s.t.  $\mu \in M$

$p^* = e_i$ where i = argmax($x$)

~

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

(0, 1, 0)

(.5, .3, .2)

$\Delta$

$\mu = Ap$

M

# MAP inference:

## Maximum spanning tree (Chu-Liu/Edmonds)



$$A = \begin{bmatrix} 1 & \bigcirc & \bigcirc \\ \bigcirc & 1 & 1 \\ \bigcirc & \bigcirc & \bigcirc \\ \hdashline \bigcirc & 1 & 1 \\ 1 & \dots & \bigcirc & \bigcirc & \dots \\ \bigcirc & \bigcirc & \bigcirc \\ \hdashline \bigcirc & \bigcirc & \bigcirc \\ \bigcirc & 1 & \bigcirc \\ 1 & \bigcirc & 1 \end{bmatrix} \begin{matrix} * & \to I \\ like & \to I \\ it & \to I \\ \hdashline * & \to like \\ I & \to like \\ it & \to like \\ \hdashline * & \to it \\ I & \to it \\ like & \to it \end{matrix}$$

$$\eta = \begin{bmatrix} .3 \\ .8 \\ -.5 \\ \hdashline .2 \\ -.1 \\ -.2 \\ \hdashline .7 \\ .6 \\ .1 \end{bmatrix}$$

## Hungarian algorithm



$$A = \begin{bmatrix} 1 & \bigcirc & \bigcirc \\ \bigcirc & 1 & \bigcirc \\ \bigcirc & \bigcirc & 1 \\ \hdashline \bigcirc & \bigcirc & 1 \\ 1 & \dots & \bigcirc & \bigcirc & \dots \\ \bigcirc & 1 & \bigcirc \\ \hdashline \bigcirc & 1 & \bigcirc \\ \bigcirc & \bigcirc & 1 \\ 1 & \bigcirc & \bigcirc \end{bmatrix} \begin{matrix} I & - cela \\ I & - me \\ I & - plait \\ \hdashline like & - cela \\ like & - me \\ like & - plait \\ \hdashline it & - cela \\ it & - me \\ it & - plait \end{matrix}$$

$$\eta = \begin{bmatrix} .3 \\ .8 \\ -.5 \\ \hdashline .2 \\ -.1 \\ -.2 \\ \hdashline .7 \\ .6 \\ .1 \end{bmatrix}$$

argmax $\langle x, p \rangle$
    s.t.   $p \in \Delta$

**MAP**

argmax $\langle \eta, \mu \rangle$
    s.t.   $\mu \in M$

argmax $\langle x, p \rangle + H(p)$
    s.t.   $p \in \Delta$

**Marginal**

argmax $\langle \eta, \mu \rangle + \tilde{H}(\mu)$
    s.t.   $\mu \in M$

argmax $\langle x, p \rangle - \frac{1}{2} \|Ap\|^2$
    s.t.   $p \in \Delta$

**SparseMAP**

argmax $\langle \eta, \mu \rangle - \frac{1}{2} \|\mu\|^2$
    s.t.   $\mu \in M$

$x = A^\top \eta$

$\mu = Ap$

$(0, 1, 0)$

$(.6, .4, 0)$

$(.5, .3, .2)$

$\Delta$

$M$

# Efficiently Computing SparseMAP

$$\text{argmax } \langle \, \eta, \mu \, \rangle - \tfrac{1}{2} \|\mu\|^2$$
$$\text{s.t. } \quad \mu \in M$$

QP with exponentially many vertices!

## Forward Pass:

Active Set algorithm

only accesses $M$
through MAP calls

linear **& finite**
convergence

## Backward Pass:

$$\frac{\partial \mu^*}{\partial \eta}$$

Linear in dim(M)
and in # selected structures

# Sparse Latent Structure

# Natural Language Inference

Prem:   A gentleman overlooking a neighborhood situation.

Hypo:   A police officer watches a situation closely.

(P, H)

entailment
contradiction
neither

# Natural Language Inference

Prem:   A gentleman overlooking a neighborhood situation.

Hypo:   A police officer watches a situation closely.



(P, H)

A gentleman overlooking ...

A police officer watches ...

entailment
contradiction
neither

Model: ESIM [Chen & al, 2017]

# Natural Language Inference

Prem:   A gentleman overlooking a neighborhood situation.

Hypo:   A police officer watches a situation closely.



(P, H)

*A gentleman overlooking ...*

*A police officer watches ...*

entailment
contradiction
neither

Model: ESIM [Chen & al, 2017]

# Natural Language Inference

# SNLI

87%
86%
85%
84%
83%

softmax  matching  sequence

# MultiNLI

76.5%
76.0%
75.5%
75.0%
74.5%
74.0%

softmax  matching  sequence

(3-class accuracy)

# Natural Language Inference with Linear Assignment

# Sparse Structured Output Prediction

# Sparse Structured Output Prediction

scores          gold structure

**SparseMAP loss**

$$L_A(\eta, \vec{\mu}) = \max_{\mu \in M} \left\{ \langle \eta, \mu \rangle - \tfrac{1}{2} \|\mu\|^2 \right\} - \langle \eta, \vec{\mu} \rangle + \tfrac{1}{2} \|\mu\|^2$$

cost (as in structured SVM)

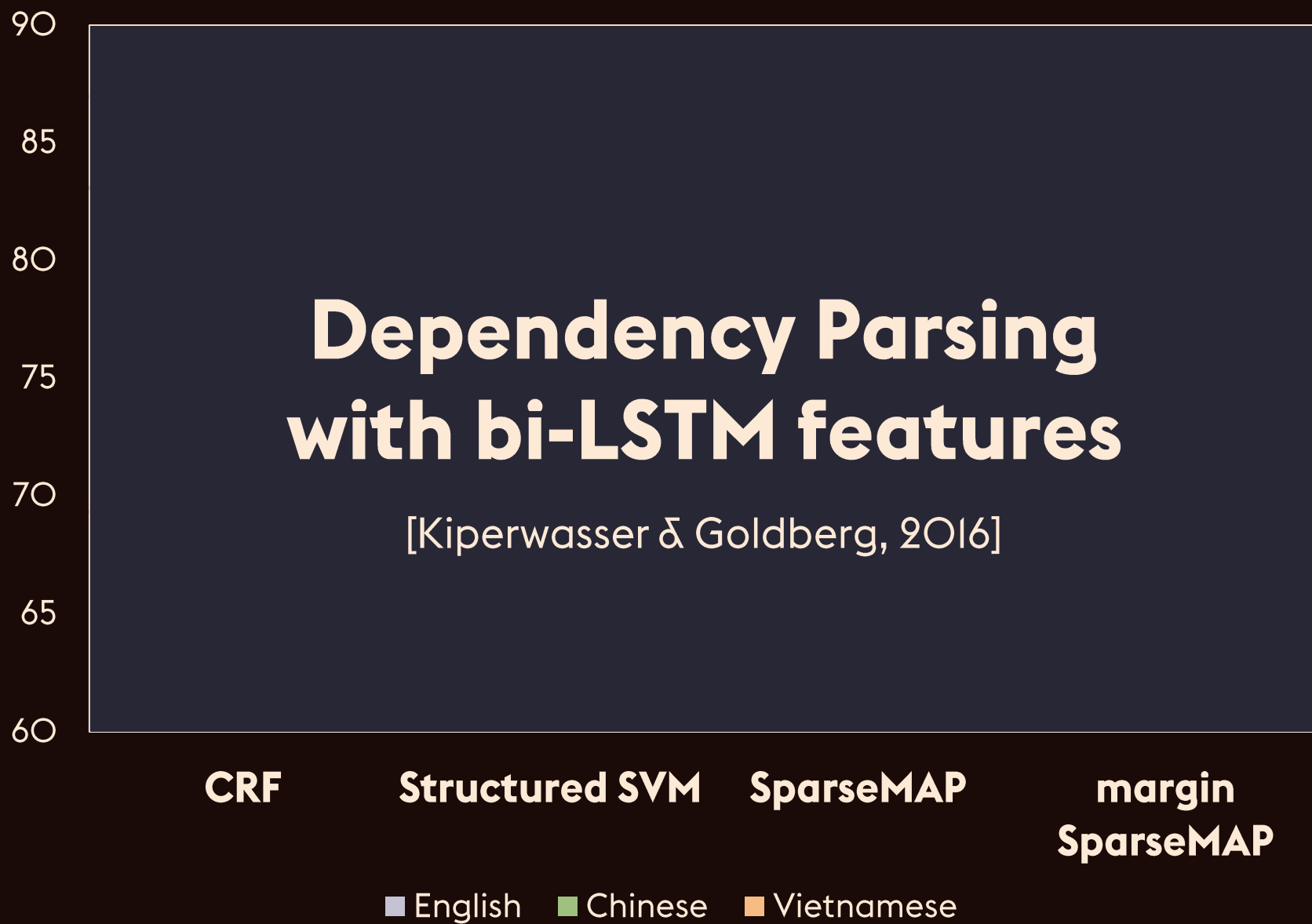**margin-SparseMAP loss**
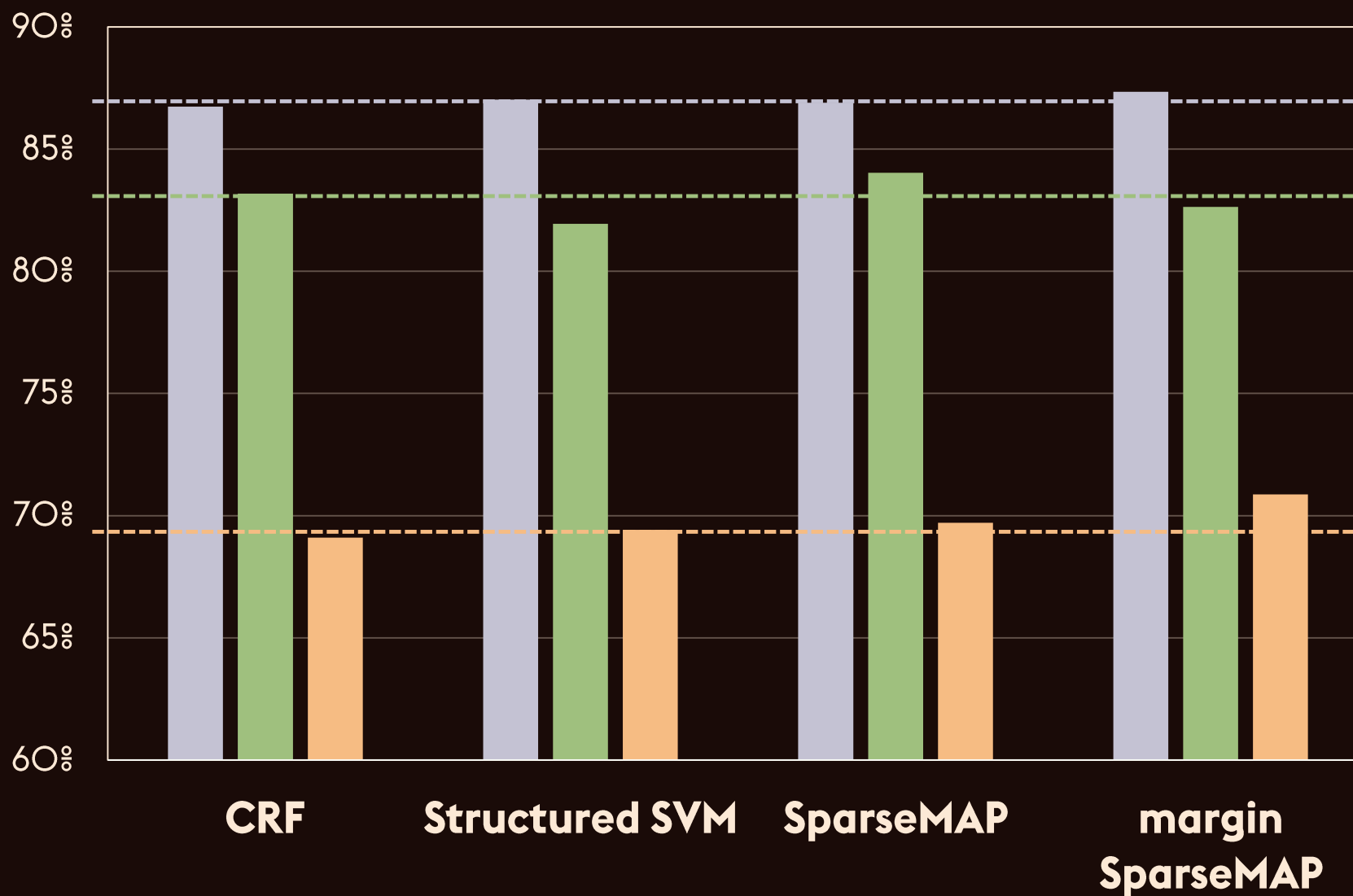
$$L_A^\rho(\eta, \vec{\mu}) = \max_{\mu \in M} \left\{ \langle \eta, \mu \rangle - \tfrac{1}{2} \|\mu\|^2 + \rho(\mu, \bar{\mu}) \right\} - \langle \eta, \vec{\mu} \rangle + \tfrac{1}{2} \|\mu\|^2$$

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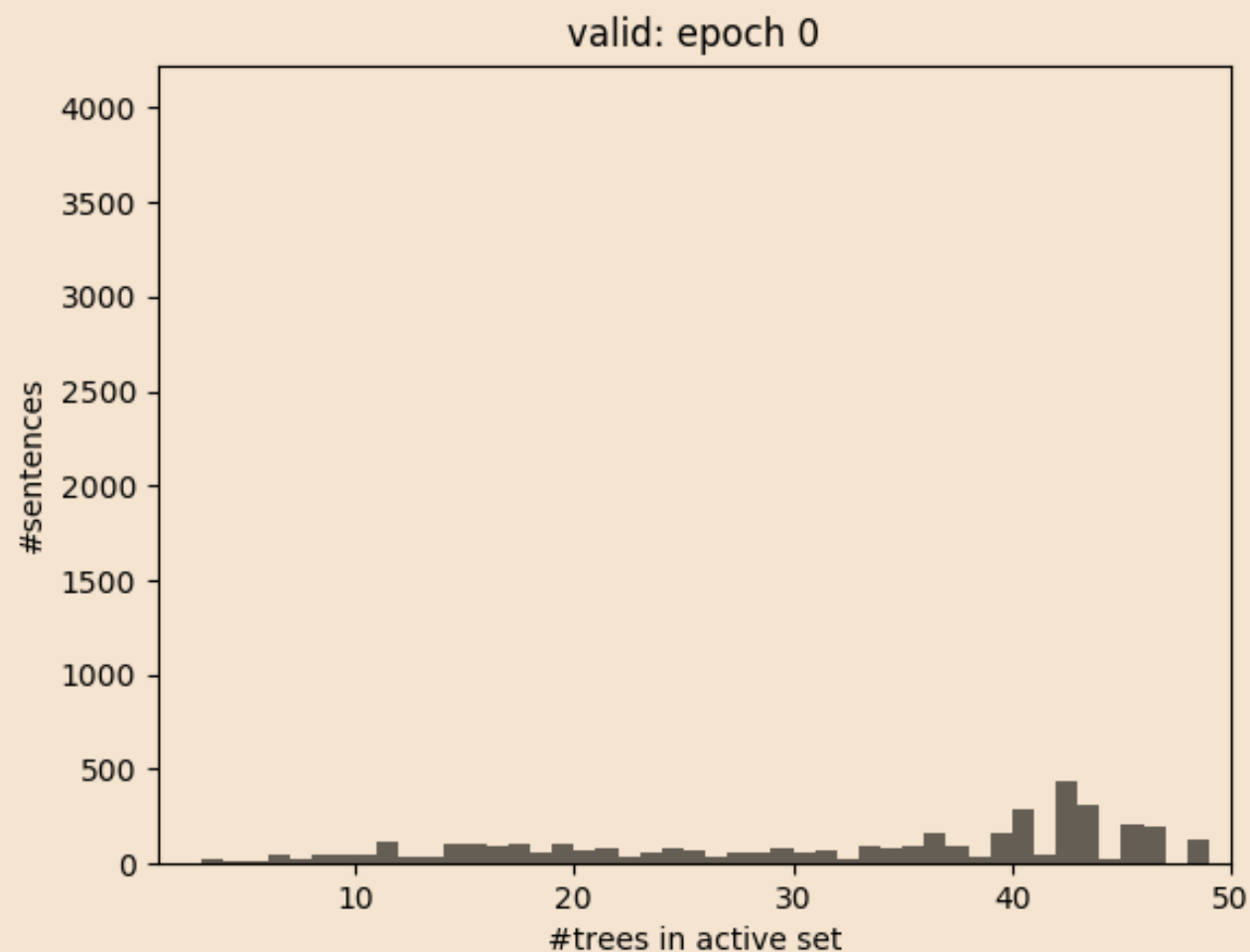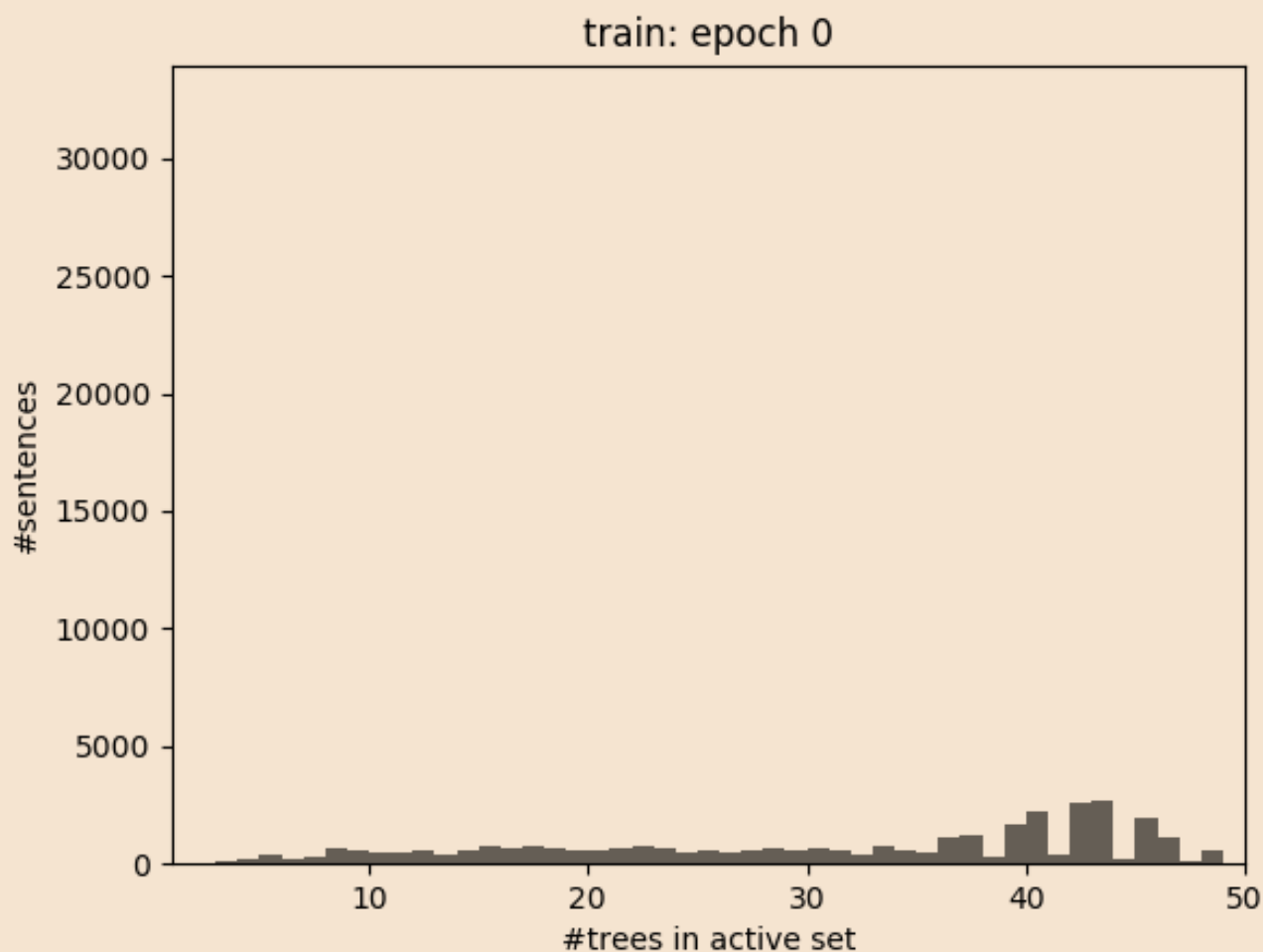
[Blondel, Martins, Niculae '18]

# Dependency Parsing
# with bi-LSTM features

[Kiperwasser & Goldberg, 2016]

90

85

80

75

70

65

60

CRF    Structured SVM    SparseMAP    margin
SparseMAP

■ English    ■ Chinese    ■ Vietnamese

Unlabeled Accuracy (UAS)
Universal Dependencies dataset

CRF    Structured SVM    SparseMAP    margin SparseMAP

English    Chinese    Vietnamese

(0, 1, 0)

(.6, .4, 0)

(.5, .3, .2)

Δ

M

## Sparse Latent Structure

A gentleman overlooking ...

A police officer watches ...

a
gentleman
overlooking
a
neighborhood
situation
.

a
police
officer
watches
a
situation
closely
.

## Sparse Structured Output Prediction

.24

.76

⋆ the broccoli looks browned around the edges .