# Structured Attention
# &
# Differentiable Dynamic Programming

Mathieu Blondel

Staff Research Scientist

NTT, Kyoto, Japan

January 4th, 2019

# Outline

1. Structured attention

2. Differentiable dynamic programming

# Outline

1. Str[...] differentiable

**differentiable max and argmax operators!**

2. Differe[...]le dynamic programming

# Outline

1. Structured attention

# Sequence to sequence with attention

Bahdanau et al., ICLR, 2015

# Sequence to sequence with attention



encoder

$h_1$     $h_2$     $h_3$

input

$w_1$     $w_2$     $w_3$

La     coalition     internationale

$H = \text{encode}(W)$

$W = \text{lookup}(\text{words})$

Bahdanau et al., ICLR, 2015

# Sequence to sequence with attention



aggregated vector

$a_1$

encoder

$h_1$   $h_2$   $h_3$

input

$w_1$   $w_2$   $w_3$

La   coalition   internationale

$\theta_t = H q_{t-1}$ # attn scores

$p_t = \mathbf{softmax}(\theta_t)$ # attn proba

$a_t = p_t^\top H$ # aggregated vector

$H = \mathrm{encode}(W)$

$W = \mathrm{lookup(words)}$

Bahdanau et al., ICLR, 2015

# Sequence to sequence with attention



decoder

$q_1$

aggregated
vector

$a_1$

encoder

$h_1$ $h_2$ $h_3$

input

$w_1$ $w_2$ $w_3$

La coalition internationale

$$q_t = \text{decode}(a_t, y_{t-1}, q_{t-1})$$

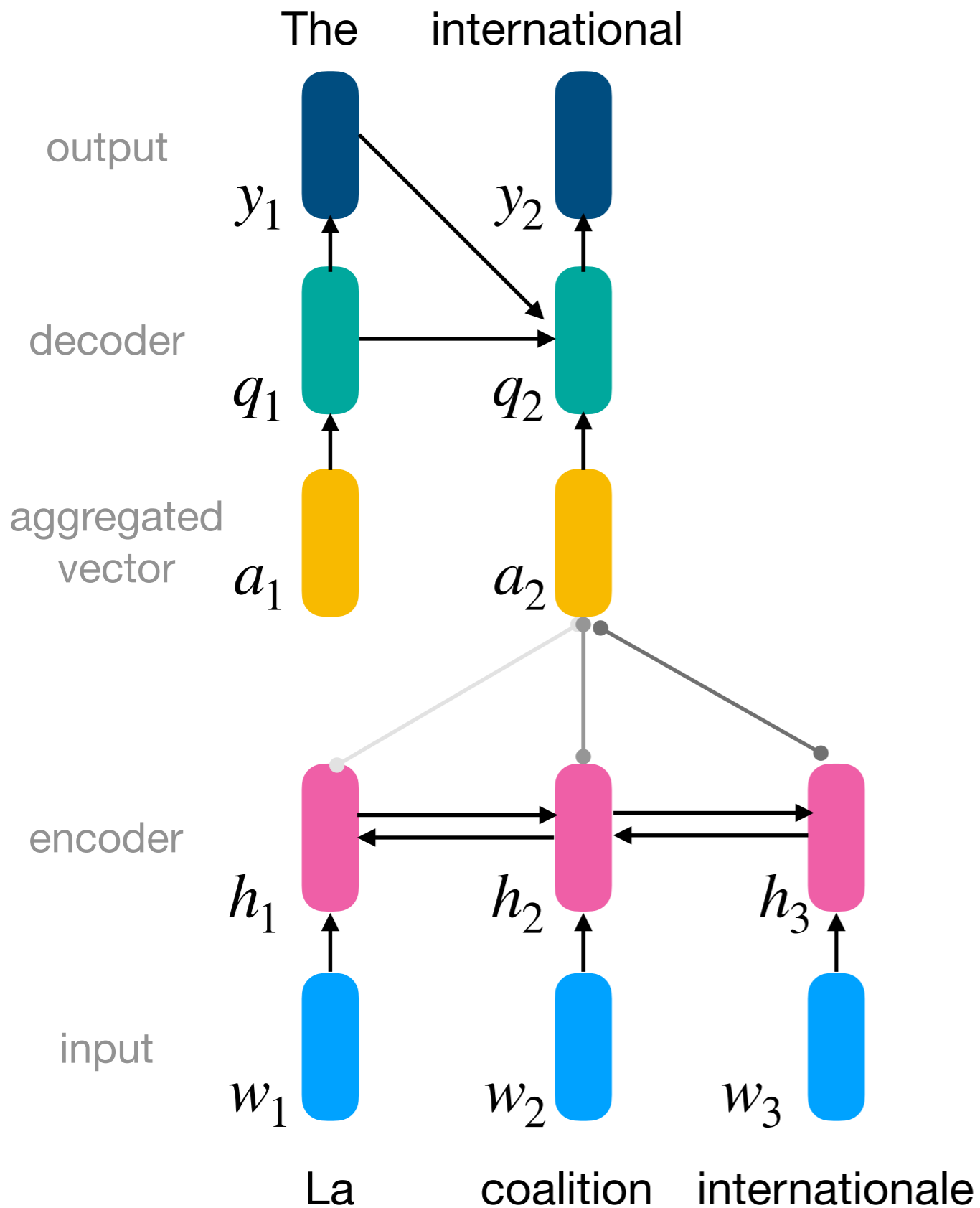$$\theta_t = H q_{t-1} \quad \text{\# attn scores}$$

$$p_t = \textbf{softmax}(\theta_t) \quad \text{\# attn proba}$$

$$a_t = p_t^\top H \quad \text{\# aggregated vector}$$

$$H = \text{encode}(W)$$

$$W = \text{lookup}(words)$$

Bahdanau et al., ICLR, 2015

# Sequence to sequence with attention



$$y_t = Mq_t + b$$

$$q_t = \text{decode}(a_t, y_{t-1}, q_{t-1})$$

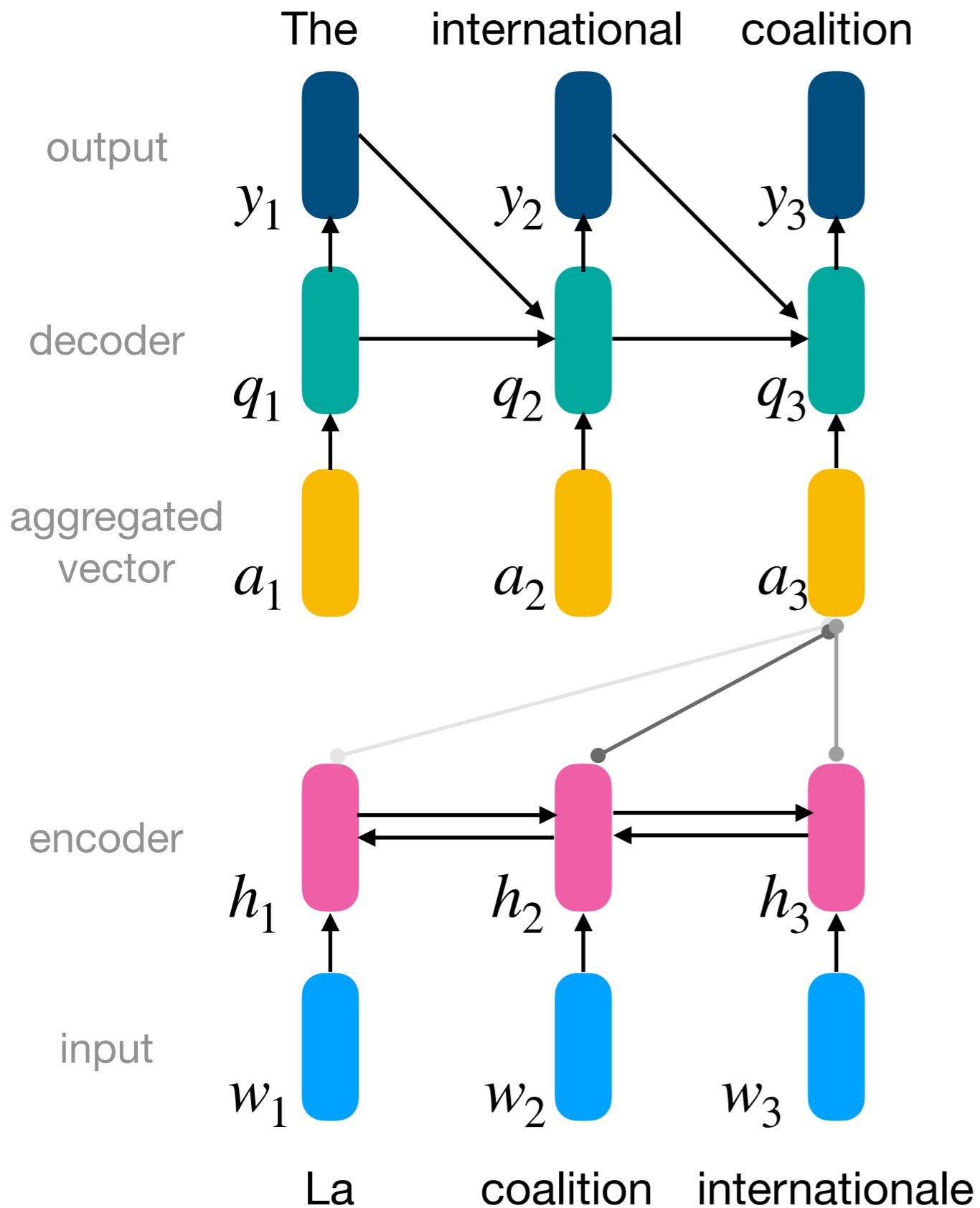$$\theta_t = Hq_{t-1} \quad \text{\# attn scores}$$

$$p_t = \textbf{softmax}(\theta_t) \quad \text{\# attn proba}$$

$$a_t = p_t^{\top}H \quad \text{\# aggregated vector}$$

$$H = \text{encode}(W)$$

$$W = \text{lookup(words)}$$

Bahdanau et al., ICLR, 2015

# Sequence to sequence with attention

The    international

output    $y_1$    $y_2$

decoder    $q_1$    $q_2$

aggregated vector    $a_1$    $a_2$

encoder    $h_1$    $h_2$    $h_3$

input    $w_1$    $w_2$    $w_3$

La    coalition    internationale

$$y_t = Mq_t + b$$

$$q_t = \text{decode}(a_t, y_{t-1}, q_{t-1})$$

$\theta_t = Hq_{t-1}$ # attn scores

$p_t = \textbf{softmax}(\theta_t)$ # attn proba

$a_t = p_t^\top H$ # aggregated vector

$$H = \text{encode}(W)$$

$$W = \text{lookup}(\text{words})$$

Bahdanau et al., ICLR, 2015

# Sequence to sequence with attention



The   international   coalition

output   $y_1$   $y_2$   $y_3$

decoder   $q_1$   $q_2$   $q_3$

aggregated vector   $a_1$   $a_2$   $a_3$

encoder   $h_1$   $h_2$   $h_3$

input   $w_1$   $w_2$   $w_3$

La   coalition   internationale

$$y_t = Mq_t + b$$

$$q_t = \text{decode}(a_t, y_{t-1}, q_{t-1})$$

$$\theta_t = Hq_{t-1} \quad \text{\# attn scores}$$

$$p_t = \textbf{softmax}(\theta_t) \quad \text{\# attn proba}$$

$$a_t = p_t^\top H \quad \text{\# aggregated vector}$$

$$H = \text{encode}(W)$$

$$W = \text{lookup}(words)$$

Bahdanau et al., ICLR, 2015

# Softmax attention

$$\mathbf{softmax}(\theta) \triangleq \frac{\exp(\theta)}{\sum_{i=1}^{m} \exp(\theta_i)}$$

Bahdanau et al., ICLR, 2015
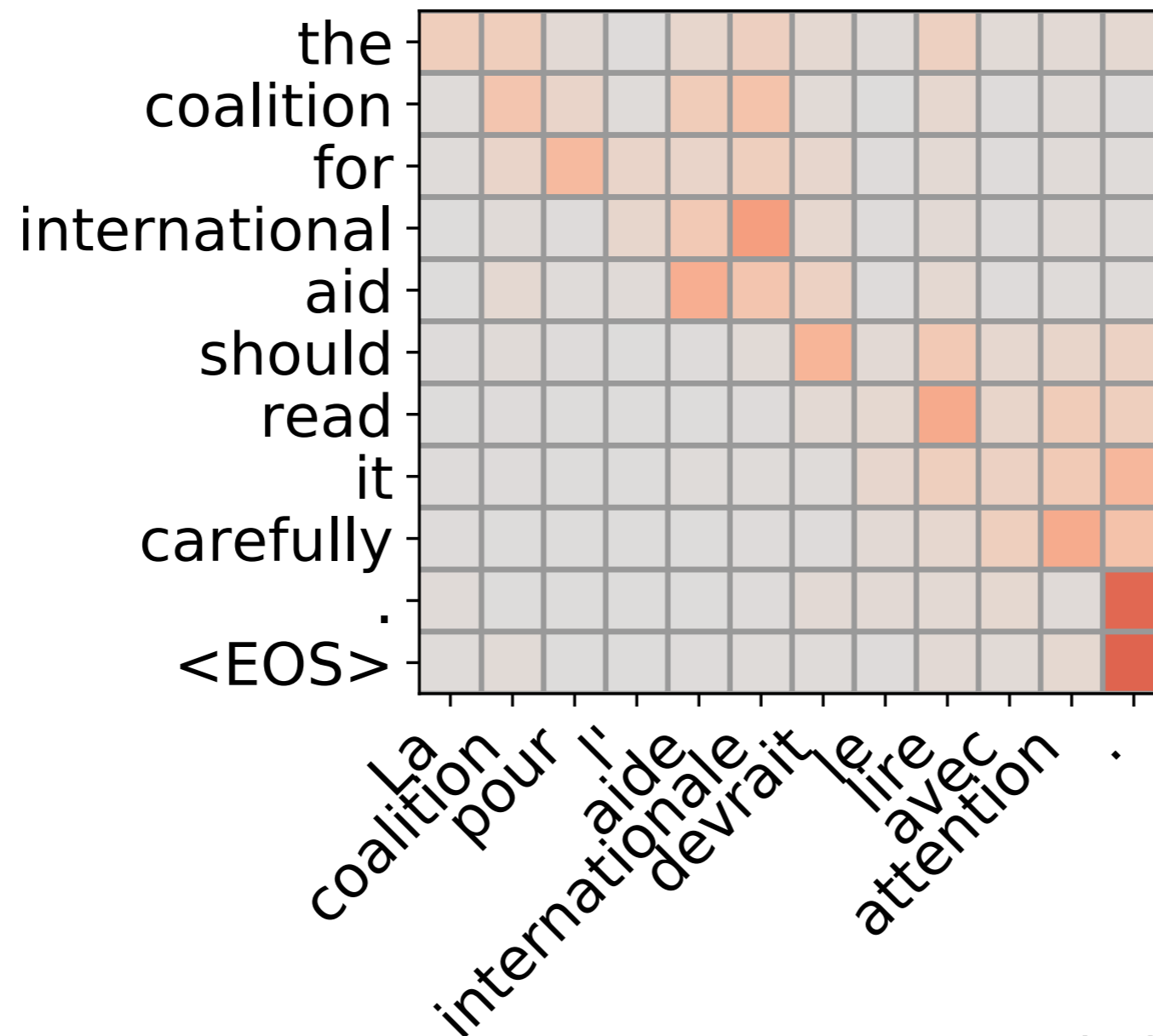
# Softmax attention

$$\mathbf{softmax}(\theta) \triangleq \frac{\exp(\theta)}{\sum_{i=1}^{m} \exp(\theta_i)}$$

the 

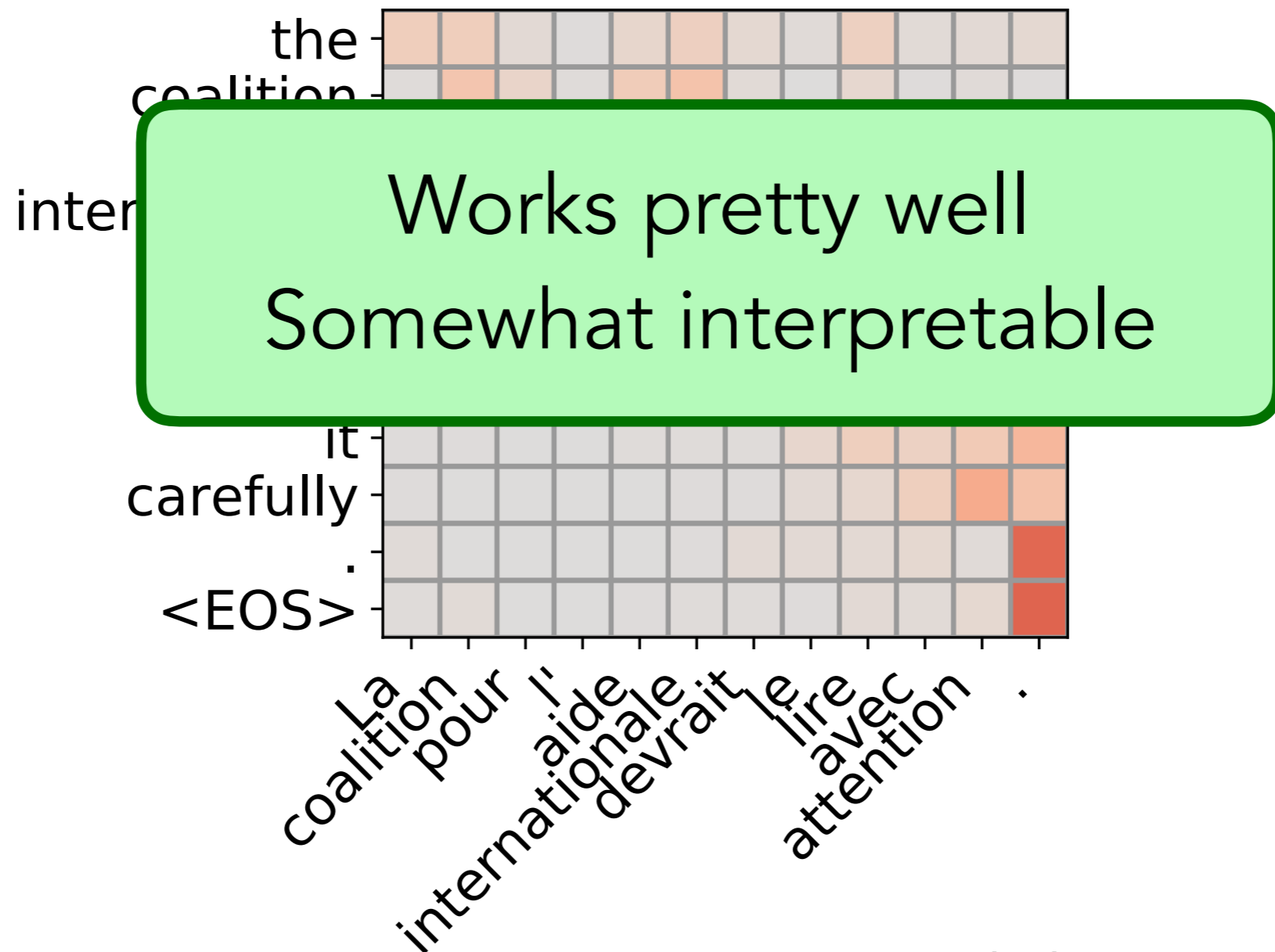La coalition pour l' aide internationale devrait le lire avec attention .

Bahdanau et al., ICLR, 2015

# Softmax attention

$$\textbf{softmax}(\theta) \triangleq \frac{\exp(\theta)}{\sum_{i=1}^{m} \exp(\theta_i)}$$



Bahdanau et al., ICLR, 2015

# Softmax attention

$$\mathbf{softmax}(\theta) \triangleq \frac{\exp(\theta)}{\sum_{i=1}^{m} \exp(\theta_i)}$$



Bahdanau et al., ICLR, 2015

# Softmax attention

$$\mathbf{softmax}(\theta) \triangleq \frac{\exp(\theta)}{\sum_{i=1}^{m} \exp(\theta_i)}$$



Bahdanau et al., ICLR, 2015

# Softmax attention

$$\mathbf{softmax}(\theta) \triangleq \frac{\exp(\theta)}{\sum_{i=1}^{m} \exp(\theta_i)}$$



the
coalition
inter...
it
carefully
.
<EOS>

La coalition pour l' aide internationale devrait le lire avec attention .

Works pretty well
Somewhat interpretable

Bahdanau et al., ICLR, 2015

# Softmax attention

$$\mathbf{softmax}(\theta) \triangleq \frac{\exp(\theta)}{\sum_{i=1}^{m} \exp(\theta_i)}$$



Works pretty well
Somewhat interpretable

Pays attention to all words
(**dense**)

Bahdanau et al., ICLR, 2015

# Sparsemax attention

$$\mathbf{sparsemax}(\theta) \triangleq \arg\min_{p\in\triangle^m} \|p - \theta\|^2$$



Martins & Astudillo, ICML, 2016

# Sparsemax attention

$$\textbf{sparsemax}(\theta) \triangleq \arg\min_{p \in \triangle^m} \|p - \theta\|^2$$



**Sparse**

More interpretable

Martins & Astudillo, ICML, 2016

# Sparsemax attention

$$\textbf{sparsemax}(\theta) \triangleq \arg\min_{p \in \triangle^m} \|p - \theta\|^2$$



**Sparse**
More interpretable

Can we **generalize** it?
Can we incorporate **structure**?

Martins & Astudillo, ICML, 2016

# Fusedmax attention (proposed)

**fusedmax**$(\theta) \triangleq$ **???**



Niculae & Blondel, NIPS 2017

# Fusedmax attention (proposed)

**fusedmax**$(\theta) \triangleq$ **???**



Sparse

**Adjacent grouping**

Good prior / Inductive bias
(encourage peeking at entire blocks of words)

Niculae & Blondel, NIPS 2017

# Our contributions

# Our contributions

- A principled framework for **differentiable argmax** operators

# Our contributions

- A principled framework for **differentiable argmax** operators

  - Recovers softmax and sparsemax as special cases

Niculae & Blondel, NIPS 2017

# Our contributions

- A principled framework for **differentiable argmax** operators

  - Recovers softmax and sparsemax as special cases

  - Enables to construct new operators easily

# Our contributions

- A principled framework for **differentiable argmax** operators

  - Recovers softmax and sparsemax as special cases

  - Enables to construct new operators easily

- Efficient **forward** and **backward** computations for **fusedmax**

# Our contributions

- A principled framework for **differentiable argmax** operators

  - Recovers softmax and sparsemax as special cases

  - Enables to construct new operators easily

- Efficient **forward** and **backward** computations for **fusedmax**

- Extensive experiments on NMT and sentence summarization

Niculae & Blondel, NIPS 2017

# From argmax to softmax

$$i^\star \in \arg\max_{i \in [m]} \theta_i$$

# From argmax to softmax

$$\mathbf{argmax}(\theta) \triangleq e_{i^\star} \qquad i^\star \in \arg\max_{i \in [m]} \theta_i$$

One-hot representation
of integer argmax

# From argmax to softmax

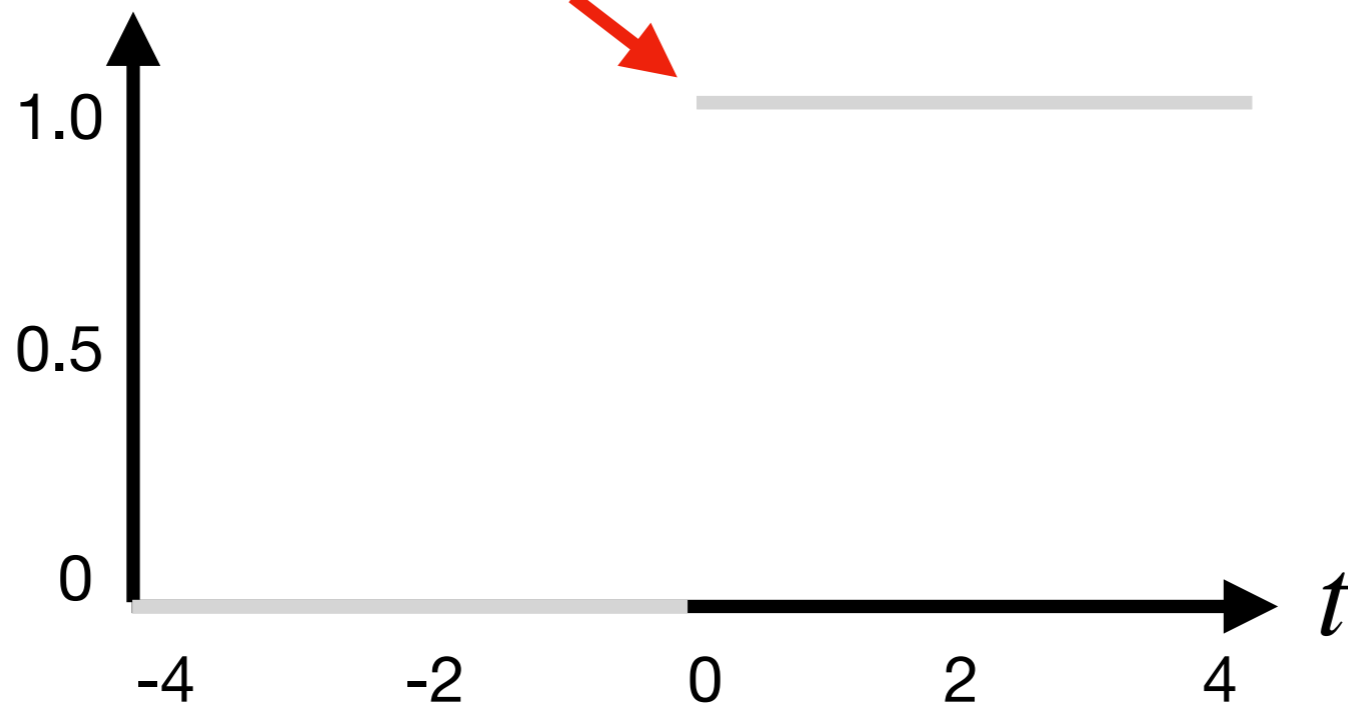Function from
$\mathbb{R}^m$ to $\{e_1, \ldots, e_m\}$

$\mathbf{argmax}(\theta) \triangleq e_{i^\star}$     $i^\star \in \arg\max_{i \in [m]} \theta_i$

One-hot representation
of integer argmax

# From argmax to softmax

Function from
$R^m$ to $\{e_1, ..., e_m\}$

$$\mathbf{argmax}(\theta) \triangleq e_{i^\star} \qquad i^\star \in \arg\max_{i \in [m]} \theta_i$$

One-hot representation
of integer argmax

Discontinuous
function



— $\mathbf{argmax}([t,0])_1$

# From argmax to softmax

Function from $R^m$ to $\{e_1, \ldots, e_m\}$

$$\mathbf{argmax}(\theta) \triangleq e_{i^\star} \qquad i^\star \in \arg\max_{i \in [m]} \theta_i$$

One-hot representation of integer argmax

Discontinuous function



usual sigmoid function when m = 2

━━━ $\mathbf{argmax}([t,0])_1$ ━━━ $\mathbf{softmax}([t,0])_1$

Should really be called soft *argmax*

# From argmax to softmax

Function from $R^m$ to $\{e_1, \ldots, e_m\}$

$$\mathbf{argmax}(\theta) \triangleq e_{i^\star} \qquad i^\star \in \arg\max_{i \in [m]} \theta_i$$

One-hot representation

Where does the softmax come from?

Can we generalize it?



—— $\mathbf{argmax}([t,0])_1$    —— $\mathbf{softmax}([t,0])_1$

Should really be called soft *argmax*

# Differentiable argmax: a variational perspective

$$\textbf{argmax}(\theta) = \underset{p \in \{e_1, \ldots, e_m\}}{\arg\max} \langle p, \theta \rangle$$

$$\mathbf{argmax}(\theta) = \underset{p \in \{e_1, \ldots, e_m\}}{\arg\max} \langle p, \theta \rangle$$

$$= \underset{p \in \triangle^m}{\arg\max} \langle p, \theta \rangle$$

$$\textbf{argmax}(\theta) = \underset{p \in \{e_1, \ldots, e_m\}}{\arg\max} \langle p, \theta \rangle$$

$$= \underset{p \in \triangle^m}{\arg\max} \langle p, \theta \rangle$$

Fundamental theorem
of linear programming
(Dantzig, 1955)



$p^\star$

$\theta$

$e_2$          $e_3$

**unregularized (Ω=0)**

# Differentiable argmax: a variational perspective

$$\mathbf{argmax}(\theta) = \underset{p \in \{e_1, \ldots, e_m\}}{\arg\max} \langle p, \theta \rangle \qquad \mathbf{argmax}_{\Omega}(\theta) \triangleq \underset{p \in \triangle^m}{\arg\max} \langle p, \theta \rangle - \Omega(p)$$

$$= \underset{p \in \triangle^m}{\arg\max} \langle p, \theta \rangle$$

**Strongly-convex** regularizer

Fundamental theorem
of linear programming
(Dantzig, 1955)

$p^\star$

$\theta$

$e_2$   $e_3$

**unregularized (Ω=0)**

# Differentiable argmax: a variational perspective

Introduce regularization

$$\mathbf{argmax}(\theta) = \underset{p\in\{e_1,\ldots,e_m\}}{\arg\max}\ \langle p,\theta\rangle \qquad \mathbf{argmax}_\Omega(\theta) \triangleq \underset{p\in\triangle^m}{\arg\max}\langle p,\theta\rangle-\Omega(p)$$

**Strongly-convex** regularizer

$$= \underset{p\in\triangle^m}{\arg\max}\langle p,\theta\rangle$$

Fundamental theorem
of linear programming
(Dantzig, 1955)

Move solution
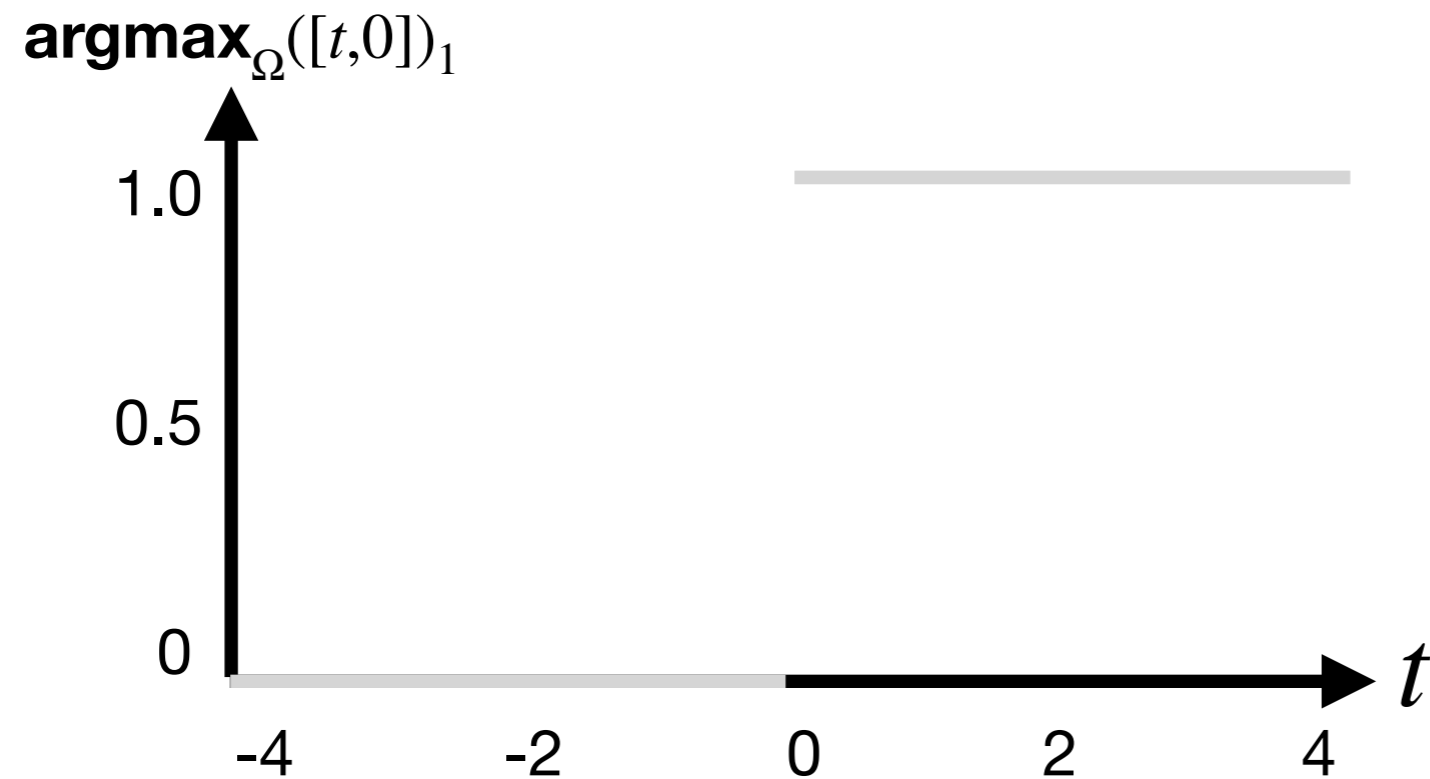away from the simplex vertices
(spread probability mass)

**unregularized (Ω=0)**

**regularized**

# Examples

# Examples

**Unregularized**

$$\Omega(p) = 0$$



**argmax**$_\Omega([t,0])_1$

1.0

0.5

0

-4    -2    0    2    4    $t$

**argmax**

$e_2$    $e_3$

— **argmax**$([t,0])_1$

# Examples

**Unregularized**

$$\Omega(p) = 0$$

**Shannon (negative) entropy**

$$\Omega(p) = \sum_i p_i \log p_i$$



**argmax**$([t,0])_1$   **softmax**$([t,0])_1$
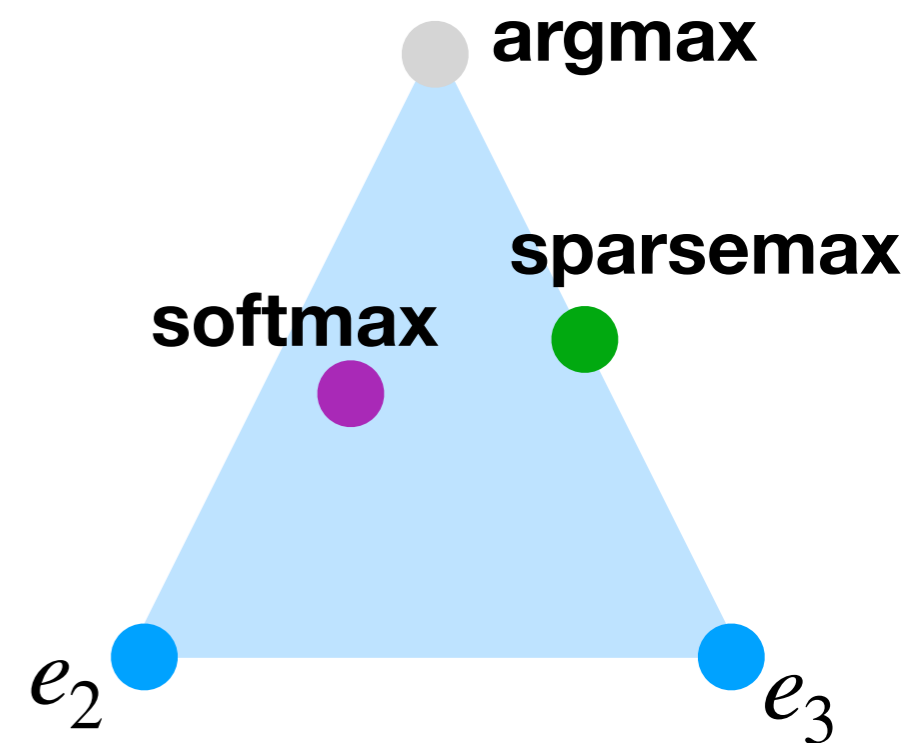
# Examples

**Unregularized**

$$\Omega(p) = 0$$

**Shannon (negative) entropy**

$$\Omega(p) = \sum_i p_i \log p_i$$
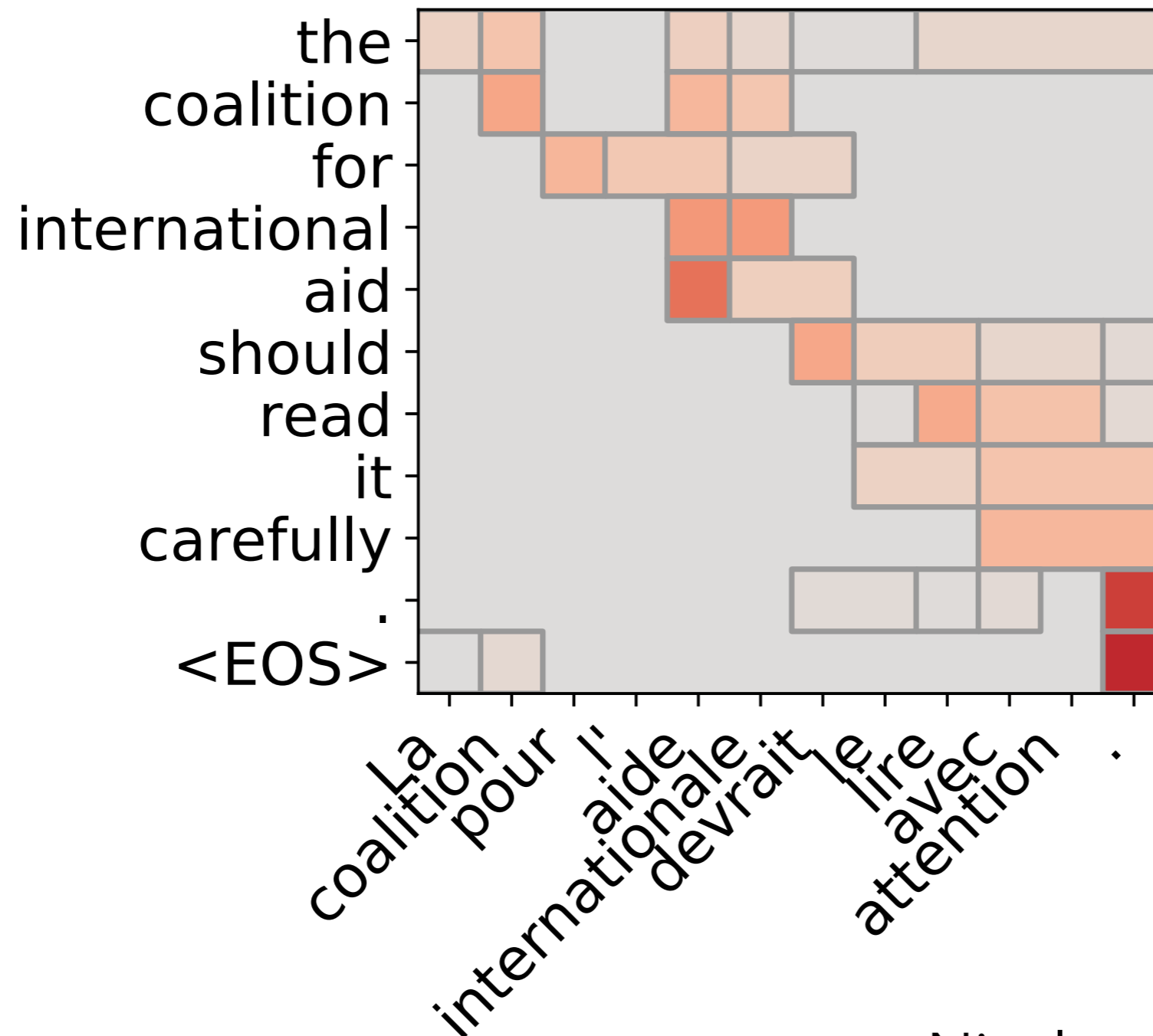
**Squared norm**

$$\Omega(p) = \frac{1}{2}\|p\|^2$$



**argmax** $([t,0])_1$  **softmax** $([t,0])_1$  **sparsemax** $([t,0])_1$

# Fusedmax attention

$$\mathbf{fusedmax}(\theta) = \mathbf{argmax}_{\Omega}(\theta)$$



Niculae & Blondel, NIPS 2017

# Fused Lasso (a.k.a. 1d total variation)

$$\mathbf{prox}_{TV}(x) \triangleq \arg\min_{y \in \mathbb{R}^m} \|x - y\|^2 + \lambda \sum_{i=1}^{m-1} |y_{i+1} - y_i|$$

$x$

$y^\star$

Total variation signal denoising

# Fusedmax attention

We choose 
$$\Omega(p) \triangleq \underbrace{\frac{1}{2}\|p\|^2}_{\text{sparsemax}} + \underbrace{\lambda \sum_{i=1}^{m-1} |p_{i+1} - p_i|}_{\text{fused lasso}}$$

# Fusedmax attention

We choose 

**sparsemax**          **fused lasso**

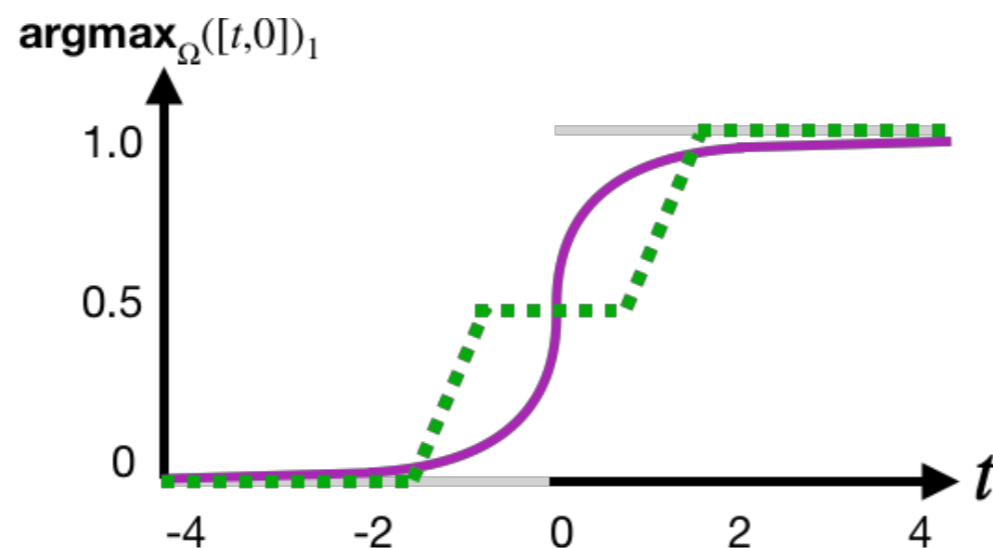$$\Omega(p) \triangleq \frac{1}{2}\|p\|^2 + \lambda \sum_{i=1}^{m-1} |p_{i+1} - p_i|$$

leading to

$$\textbf{fusedmax}(\theta) \triangleq \arg\min_{p \in \triangle^m} \frac{1}{2}\|p - \theta\|^2 + \lambda \sum_{i=1}^{m-1} |p_{i+1} - p_i|$$
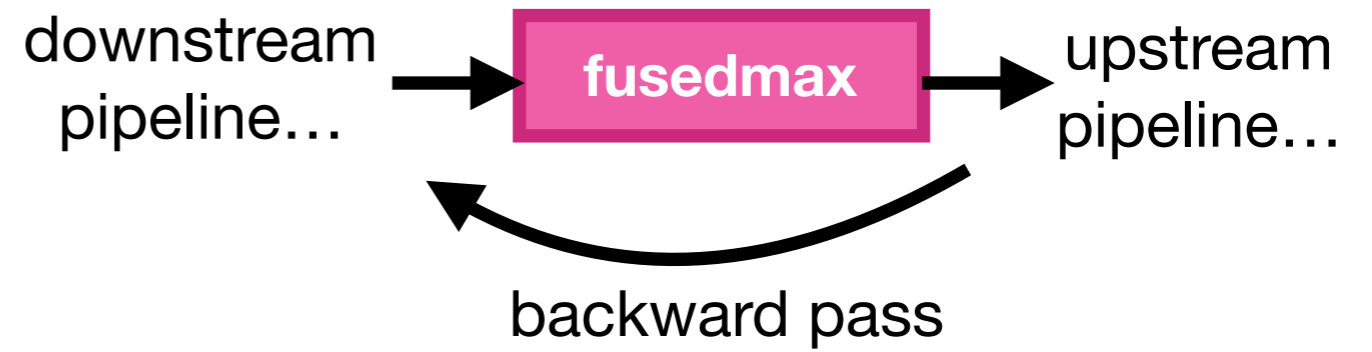


$\textbf{argmax}([t,0])_1$     $\textbf{softmax}([t,0])_1$     $\textbf{fusedmax}([t,0])_1$
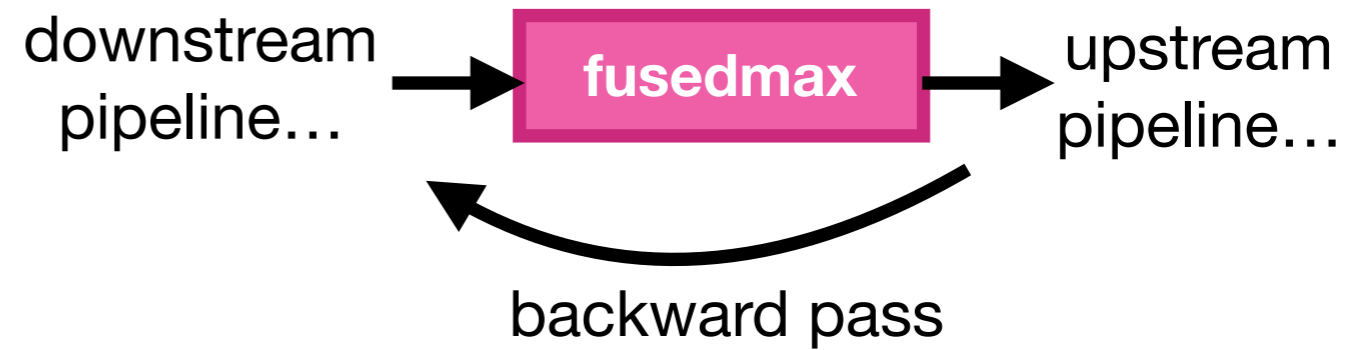
# Fusedmax: computation

How to compute forward and backward passes?

downstream pipeline... → **fusedmax** → upstream pipeline...

backward pass

# Fusedmax: computation

How to compute forward and backward passes?

downstream pipeline… → **fusedmax** → upstream pipeline…
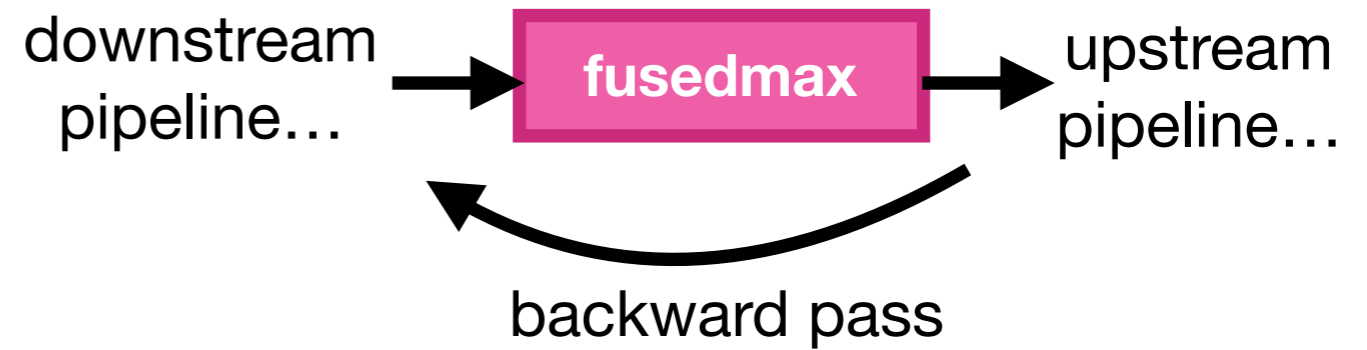
backward pass

Proposition (Niculae & Blondel, 2017)

$$\textbf{fusedmax} = \textbf{sparsemax} \circ \textbf{prox}_{TV}$$

# Fusedmax: computation

How to compute forward and backward passes?

downstream pipeline… → **fusedmax** → upstream pipeline…
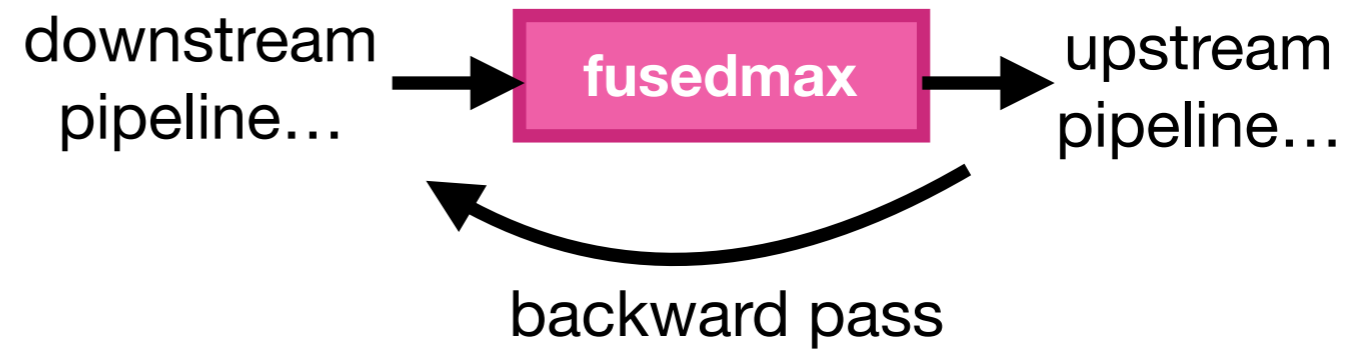
backward pass

Proposition (Niculae & Blondel, 2017)

**Not true for all regularizers!**

$$\mathbf{fusedmax} = \mathbf{sparsemax} \circ \mathbf{prox}_{TV}$$

# Fusedmax: computation

How to compute forward and backward passes?

downstream pipeline…  →  **fusedmax**  →  upstream pipeline…

backward pass

## Proposition (Niculae & Blondel, 2017)

**Not true for all regularizers!**

$$\textbf{fusedmax} = \textbf{sparsemax} \circ \textbf{prox}_{TV}$$

|  | sparsemax | prox$_{TV}$ |
|---|---|---|
| forward | Michelot, 1986 | Condat, 2013 |
| backward (Jacobian) | Martins & Atstudillo, 2016 | ? |

# Jacobian of **prox**$_{TV}$

# Jacobian of **prox**$_{TV}$

<u>Proposition (Niculae & Blondel, 2017)</u>

$|G_1|$  $\qquad$  $|G_7|$

...

the
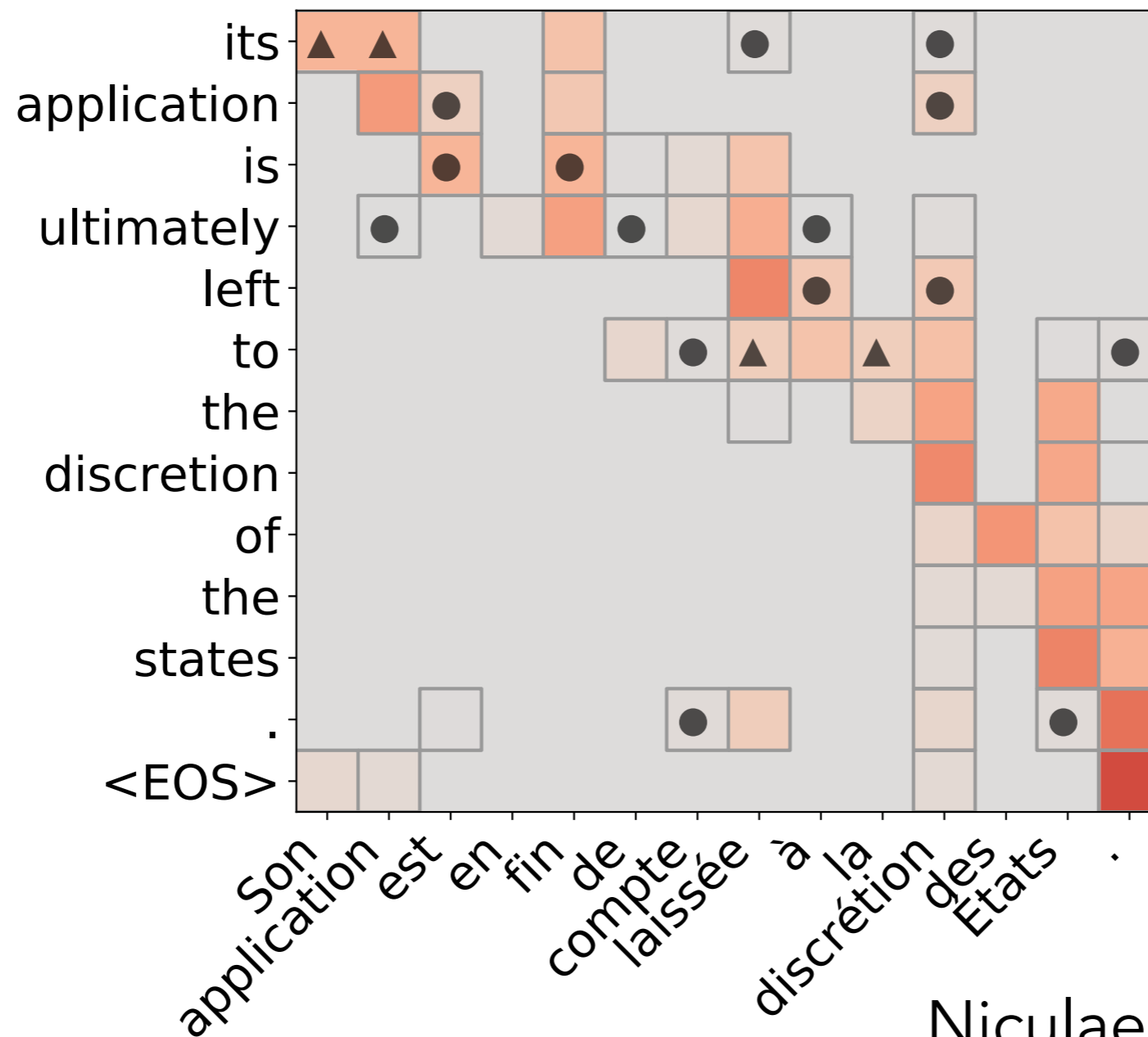
output of
**fusedmax**$(\theta)$
(forward pass)

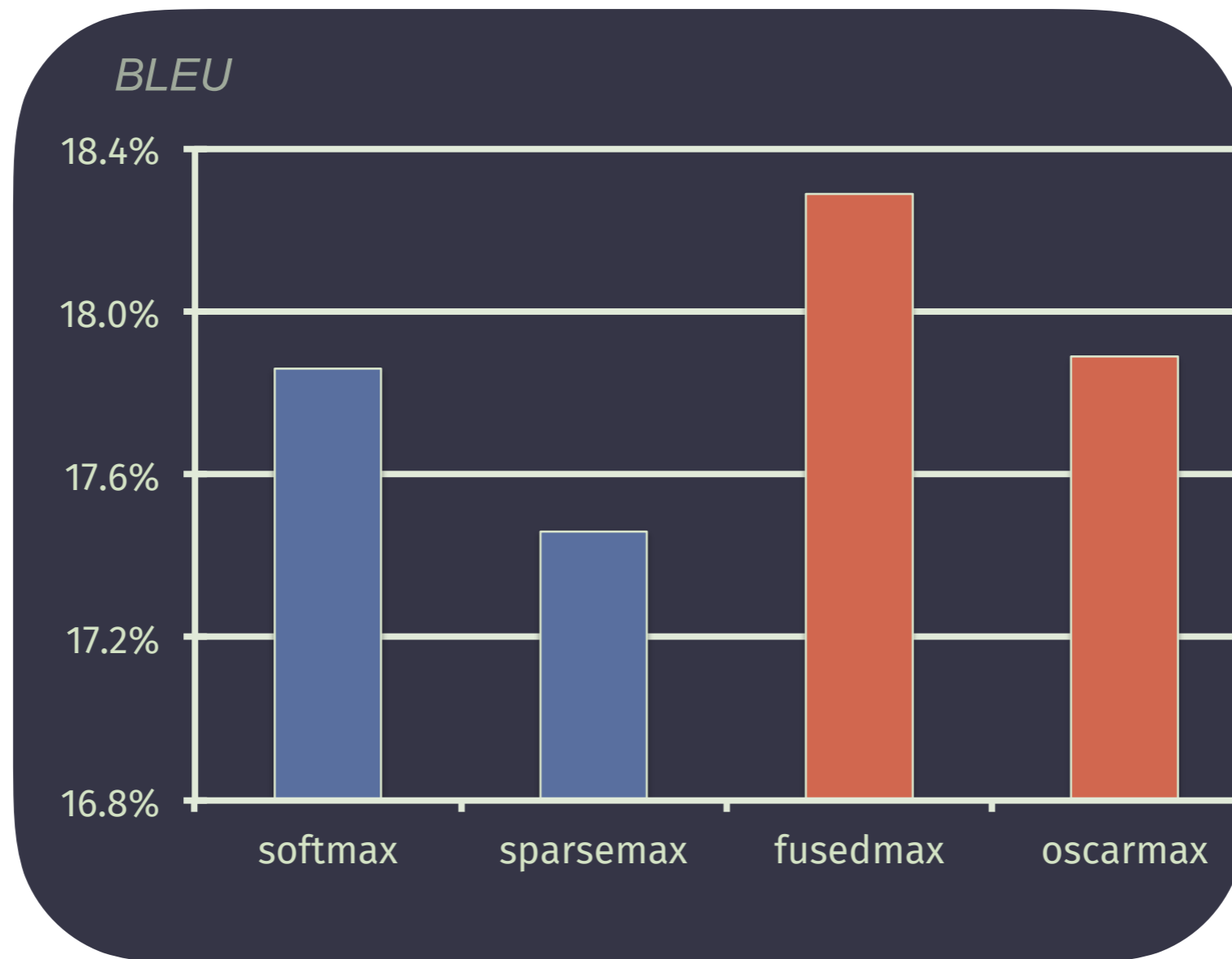# Jacobian of **prox**$_{TV}$

<u>Proposition (Niculae & Blondel, 2017)</u>

# Oscarmax attention

$$\mathbf{oscarmax}(\theta) \triangleq \arg\min_{p \in \triangle^m} \frac{1}{2}\|p - \theta\|^2 + \lambda \sum_{i<j} \max\{\,|p_i|, |p_j|\,\}$$



Niculae & Blondel, NIPS 2017

# Neural Machine Translation
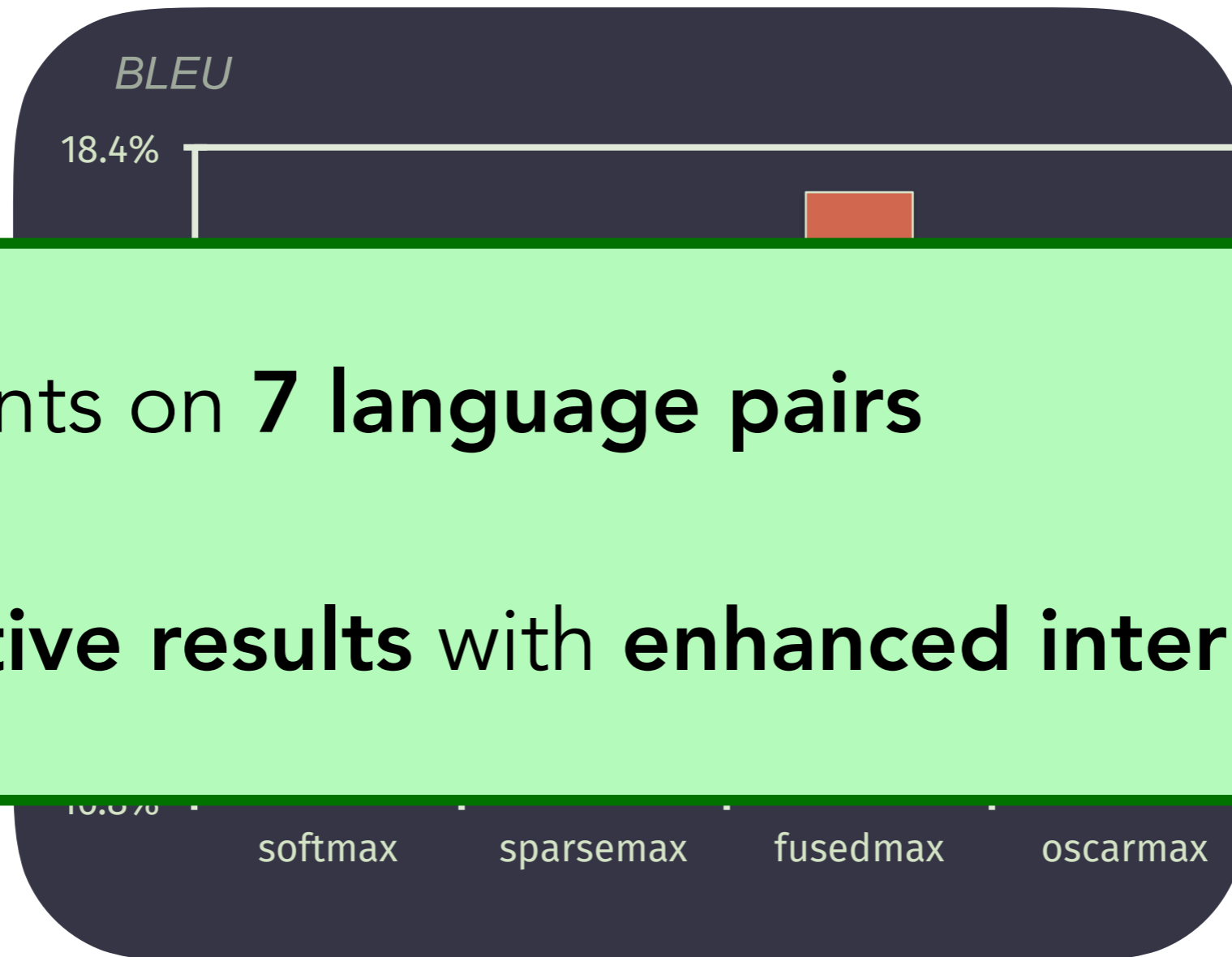
Romanian-English



Experiments based on Open-NMT
using WMT16 dataset

# Neural Machine Translation

Romanian-English



BLEU

18.4%

Experiments based on Open-NMT
using WMT16 dataset

. Experiments on **7 language pairs**

. **Competitive results** with **enhanced interpretability**!

softmax    sparsemax    fusedmax    oscarmax
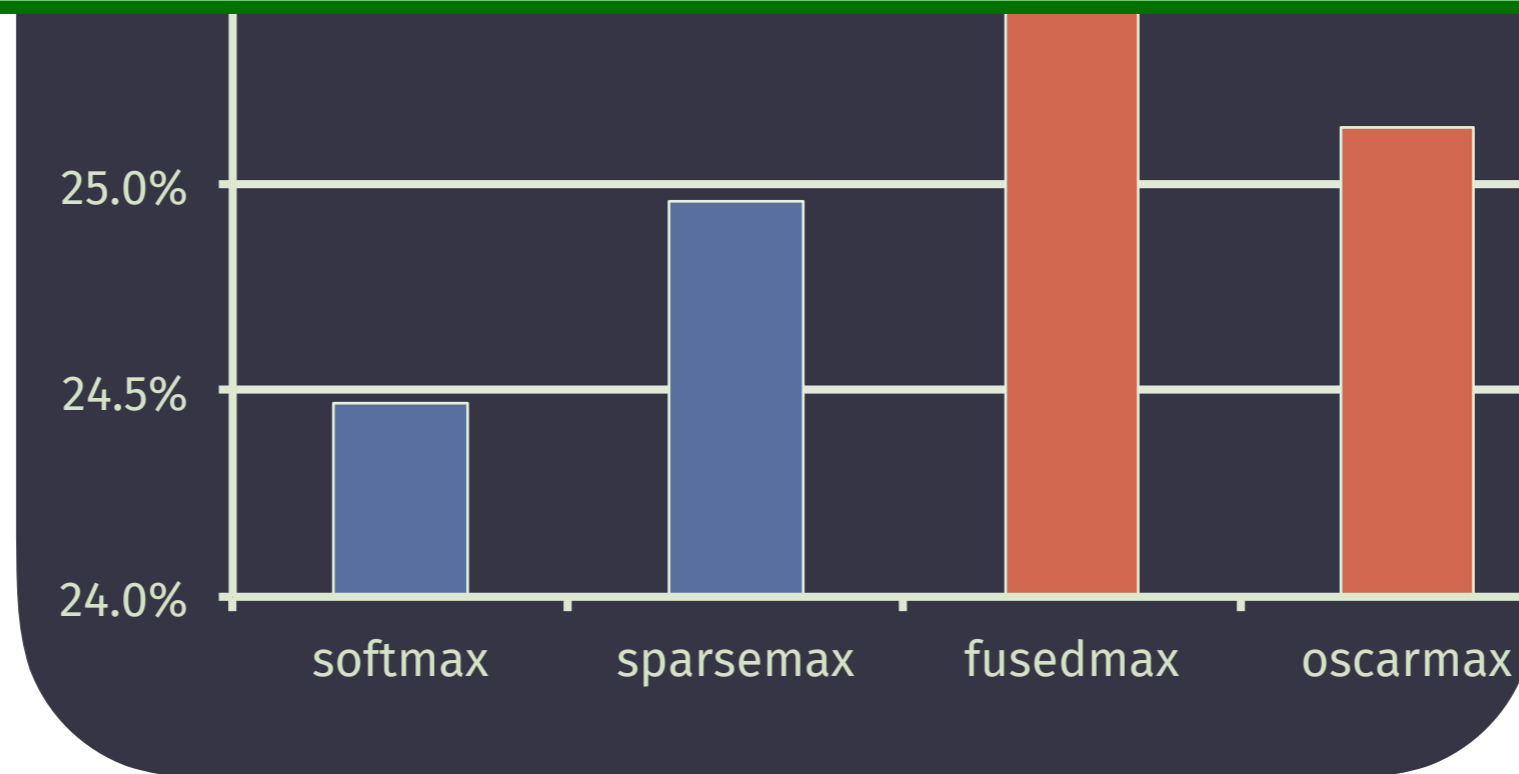
fusedmax

# Sentence summarization



ROUGE-L

Experiments based on Open-NMT
using the Gigaword sentence summarization dataset

# Sentence summarization

- **Significant accuracy improvement**

- Greatly **enhanced interpretability**



Experiments based on Open-NMT
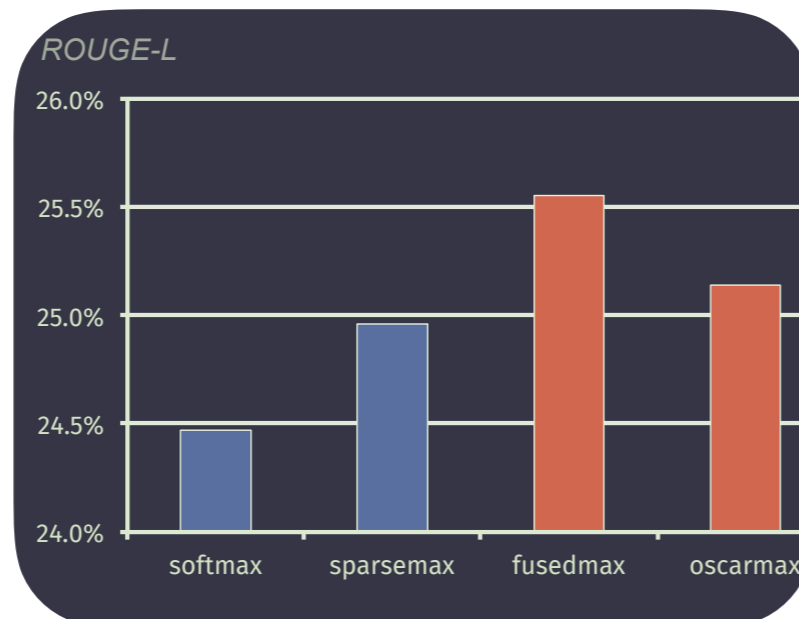using the Gigaword sentence summarization dataset

# Summary so far

## Principled framework for differentiable argmax operators

$$\mathbf{argmax}_{\Omega}(\theta) \triangleq \arg\max_{p \in \triangle^m} \langle p, \theta \rangle - \Omega(p)$$

| mechanism | regularization Ω |
|-----------|------------------|
| softmax | Shannon's neg-entropy |
| sparsemax | squared norm |
| fusedmax | squared norm + fused lasso |

## New interpretable attention mechanisms



## Faster training by leveraging sparsity

| attention | time per epoch |
|-----------|----------------|
| softmax | 1h 26m 40s $\pm$ 51s |
| sparsemax | 1h 24m 21s $\pm$ 54s |
| fusedmax | 1h 23m 58s $\pm$ 50s |
| oscarmax | 1h 23m 19s $\pm$ 50s |

## Great accuracy on various applications

# Outline

1. Structured attention

2. Differentiable dynamic programming

# Soft Viterbi algorithm: sequence tagging



one path in the DAG = one possible tag sequence

# Soft Viterbi algorithm: sequence tagging



Unregularized

Entropic regularization
⇔ Linear-chain CRF

(Lafferty et al., 2001)

one path in the DAG = one possible tag sequence

# Soft Viterbi algorithm: sequence tagging



Unregularized

Entropic regularization
⇔ Linear-chain CRF

(Lafferty et al., 2001)

Quadratic regularization
(Mensch & Blondel, 2018)
**Sparse sequence distribution**

one path in the DAG = one possible tag sequence

# Soft DTW: time series alignment

DTW = Dynamic Time Warping  [Sakoe & Chiba, 1978]
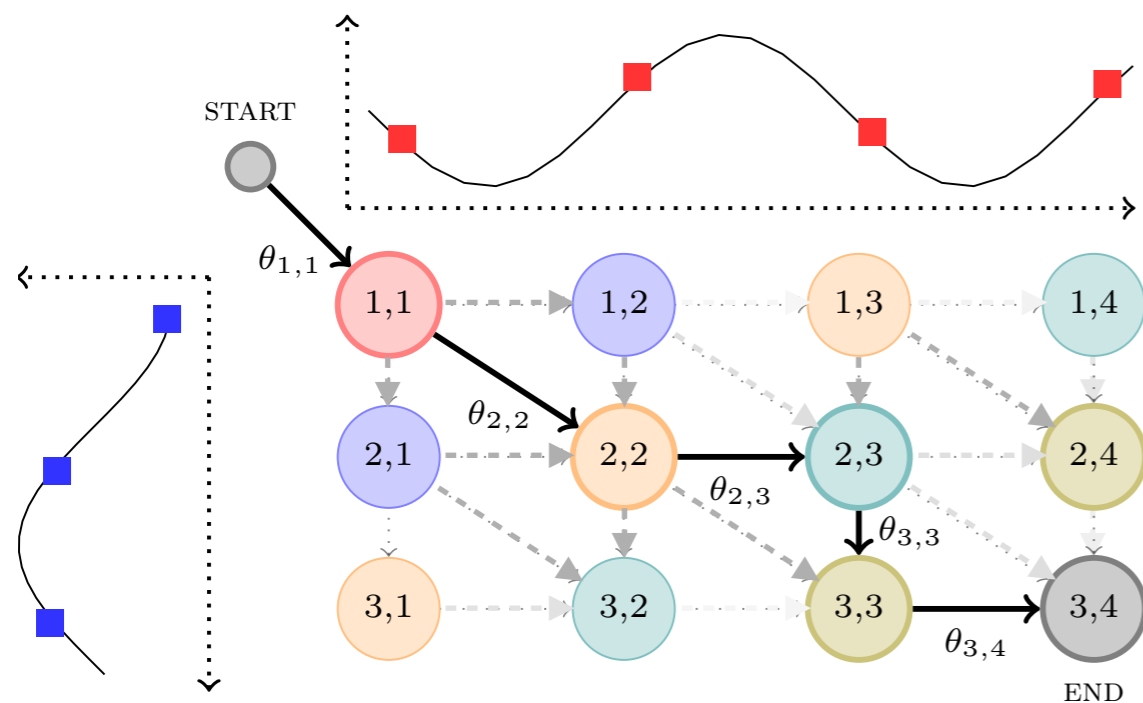
# Soft DTW: time series alignment

DTW = Dynamic Time Warping  [Sakoe & Chiba, 1978]



Entropic regularization
⇔ Soft-DTW

(Cuturi & Blondel, 2017)

one path in the DAG
=
one possible **monotonic** time-series alignment

# Soft DTW: time series alignment

DTW = Dynamic Time Warping  [Sakoe & Chiba, 1978]



Entropic regularization
⇔ Soft-DTW

(Cuturi & Blondel, 2017)

Quadratic regularization
(Mensch & Blondel, 2018)
**Sparse alignment distribution**

one path in the DAG
=
one possible **monotonic** time-series alignment

Entropic regularization
(Cuturi & Blondel, 2017)

Hard solution (DTW alignment)   Soft solution (**expected alignment** $\mathbb{E}_p[Y]$)

# Expected Alignment (Path)



Entropic regularization
(Cuturi & Blondel, 2017)

Quadratic regularization
(Mensch & Blondel, 2018)

Hard solution (DTW alignment)   Soft solution (**expected alignment** $\mathbb{E}_p[Y]$)

$$\mathbf{MAP}(\theta) \triangleq \arg\max_{y \in \mathcal{Y} \subseteq \mathbb{R}^m} \langle y, \theta \rangle$$

$\theta$

$$\mathbf{MAP}(\theta) \triangleq \arg\max_{y \in \mathscr{Y} \subseteq \mathbb{R}^m} \langle y, \theta \rangle = \arg\max_{y \in \text{conv}(\mathscr{Y})} \langle y, \theta \rangle$$



$\theta$

$\text{conv}(\mathscr{Y})$

Marginal polytope
(Wainwright & Jordan, 2008)

$$\mathbf{MAP}(\theta) \triangleq \arg\max_{y \in \mathcal{Y} \subseteq \mathbb{R}^m} \langle y, \theta \rangle = \arg\max_{y \in \mathrm{conv}(\mathcal{Y})} \langle y, \theta \rangle$$



$\theta$

$\mathrm{conv}(\mathcal{Y})$

Marginal polytope
(Wainwright & Jordan, 2008)

**Can be computed efficiently by
dynamic programming
in the case of DAGs (no cycle)**

# MAP inference: Highest-scoring Structure

$$\mathbf{MAP}(\theta) \triangleq \arg\max_{y \in \mathcal{Y} \subseteq \mathbb{R}^m} \langle y, \theta \rangle = \arg\max_{y \in \mathrm{conv}(\mathcal{Y})} \langle y, \theta \rangle$$

$\theta$

$\mathrm{conv}(\mathcal{Y})$

Marginal polytope
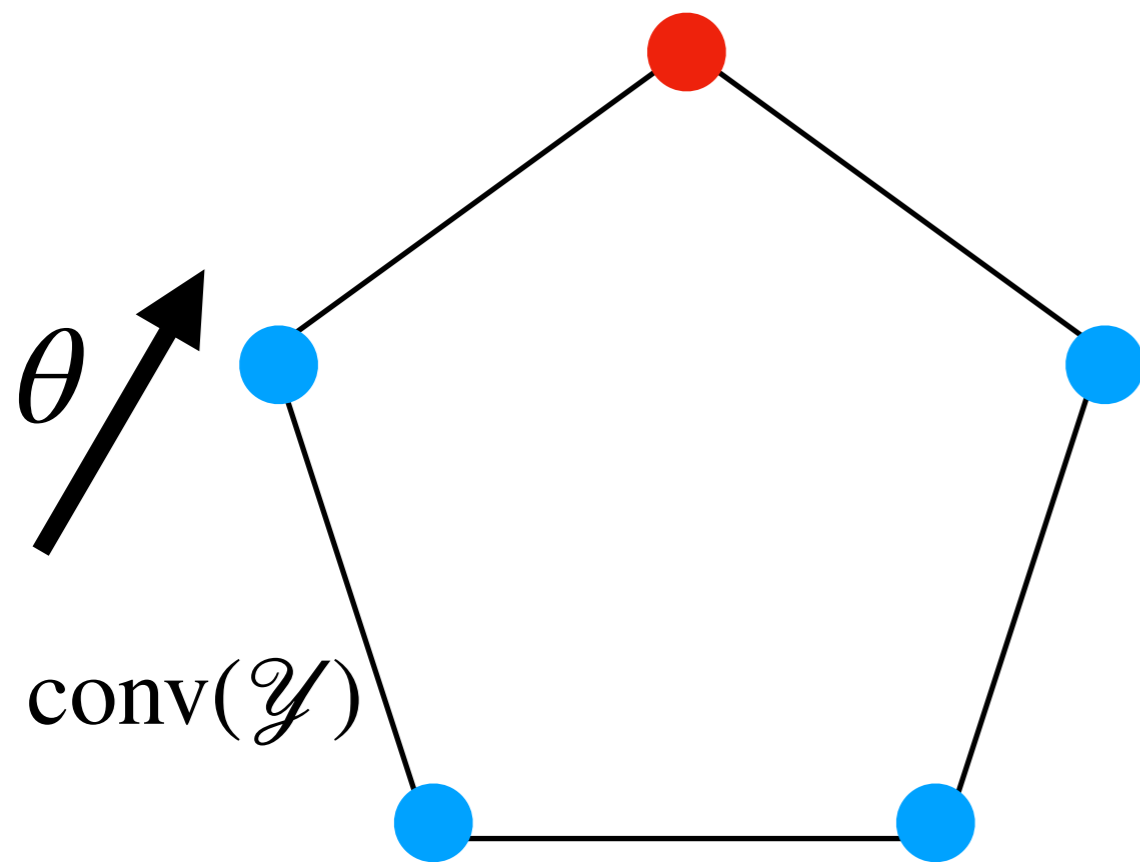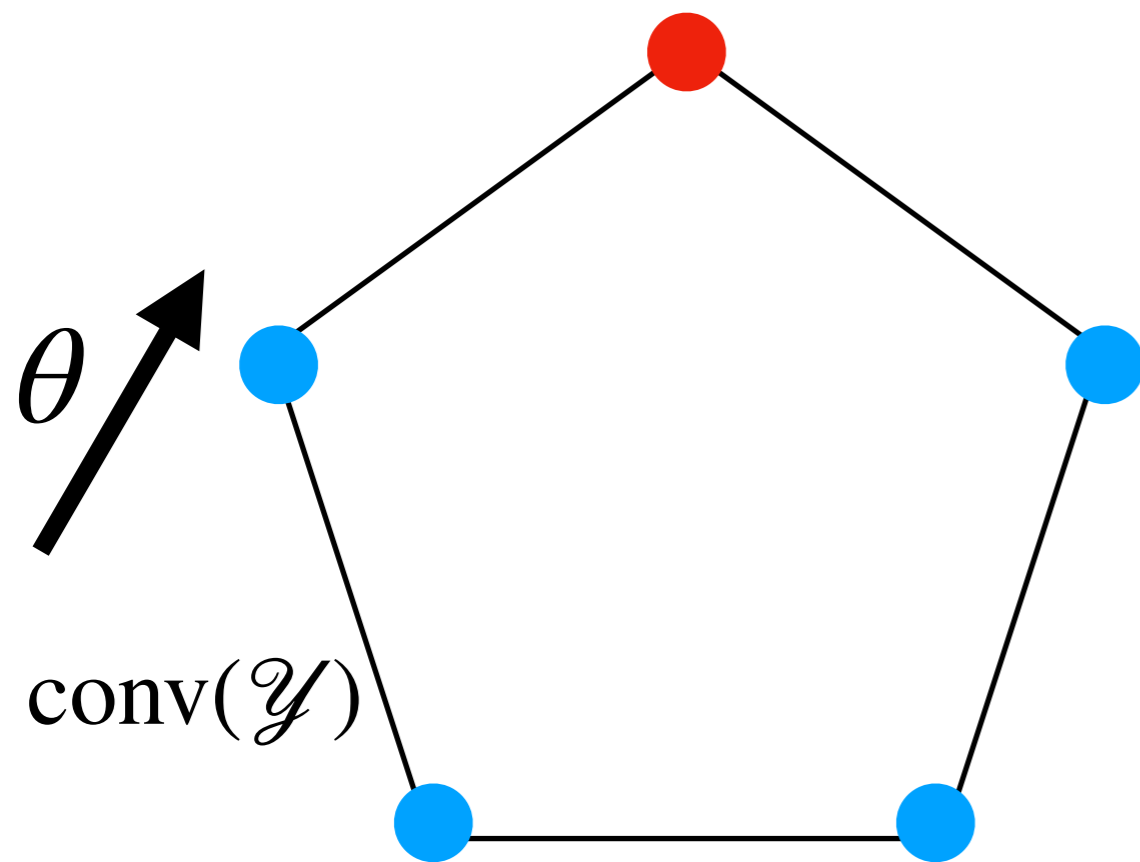(Wainwright & Jordan, 2008)

Can be computed efficiently by
dynamic programming
in the case of DAGs (no cycle)

MAP(θ) is a discontinuous function

# Bellman's recursion

Best value in state i up to time t

$$v_{t,i}(\theta) = \max_j v_{t-1,j}(\theta) + \theta_{t,i,j}$$

$t-1$        $t$        $t+1$

$v_{t-1,1}(\theta)$    $\theta_{t,1,1}$    $v_{t,1}(\theta)$

$\theta_{t,1,2}$

$\cdots$ $v_{t-1,2}(\theta)$               $\cdots$

$\theta_{t,1,3}$

$v_{t-1,2}(\theta)$

Forward pass

# DP value and optimality

# DP value and optimality

Optimality: $\mathrm{DP}(\theta) = \max_{y \in \mathscr{Y}} \langle y, \theta \rangle \in \mathbb{R}$



Forward pass

# Maintaining back pointers

$$b_{t,i}(\theta) = \arg\max_{j} \; v_{t-1,j}(\theta) + \theta_{t,i,j} \in [S]$$



$t-1$        $t$        $t+1$

$b_{t,1}(\theta)$    $v_{t,1}(\theta)$

$b_{t,2}(\theta)$    $v_{t,2}(\theta)$

$b_{t,3}(\theta)$    $v_{t,3}(\theta)$

Forward pass

# Backtracking

Optimal path equals $\mathbf{MAP}(\theta) = \arg\max_{y \in \mathcal{Y} \subseteq \mathbb{R}^m} \langle y, \theta \rangle$



Start from end state and follow back pointers

# Marginal inference: Expected Structure

Gibbs distribution

$$\textbf{marginal}(\theta) \triangleq \mathbb{E}_p[Y] \qquad p = \textbf{softmax}\left((\langle y, \theta \rangle)_{y \in \mathscr{Y}}\right) \in \triangle^{|\mathscr{Y}|}$$



$\theta$

$\textbf{marginal}(\theta)$

$\text{conv}(\mathscr{Y})$

Marginal polytope
(Wainwright & Jordan, 2008)

# Marginal inference: Expected Structure

Gibbs distribution

$$\mathbf{marginal}(\theta) \triangleq \mathbb{E}_p[Y] \qquad p = \mathbf{softmax}\left(\left(\langle y, \theta \rangle\right)_{y \in \mathcal{Y}}\right) \in \triangle^{|\mathcal{Y}|}$$

**Differentiable** but **completely dense**

(always in the interior of the polytope)



$\theta$

**marginal**$(\theta)$

$\mathrm{conv}(\mathcal{Y})$

Marginal polytope
(Wainwright & Jordan, 2008)

# Marginal inference: Expected Structure

Gibbs distribution

$$\mathbf{marginal}(\theta) \triangleq \mathbb{E}_p[Y] \qquad p = \mathbf{softmax}\left((\langle y, \theta \rangle)_{y \in \mathcal{Y}}\right) \in \triangle^{|\mathcal{Y}|}$$

**Differentiable** but **completely dense**
(always in the interior of the polytope)

**Computation:** change semiring

$$x \to e^x \quad (\max, +) \to (+, \times)$$

$\theta$

**marginal**$(\theta)$

$\mathrm{conv}(\mathcal{Y})$

Marginal polytope
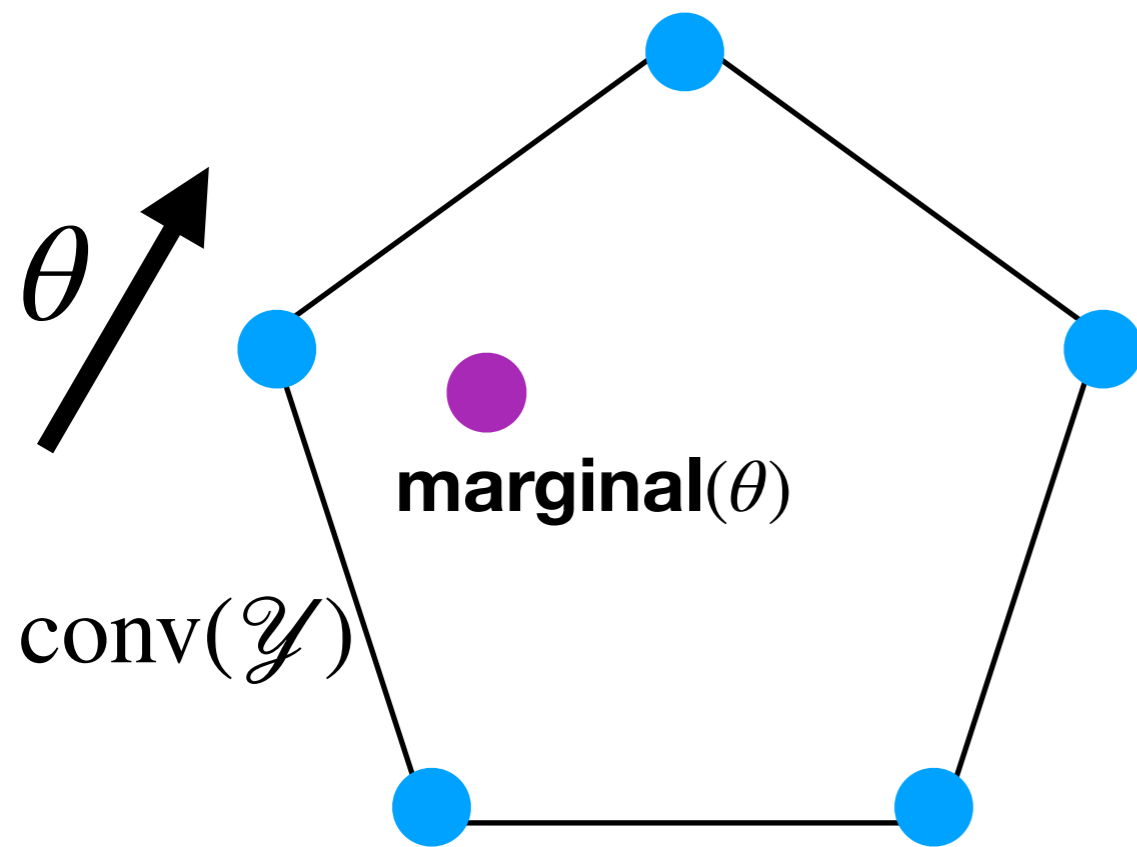(Wainwright & Jordan, 2008)

# Marginal inference: Expected Structure

Gibbs distribution

$$\mathbf{marginal}(\theta) \triangleq \mathbb{E}_p[Y] \qquad p = \mathbf{softmax}\left(\left(\langle y, \theta \rangle\right)_{y \in \mathcal{Y}}\right) \in \triangle^{|\mathcal{Y}|}$$



$\theta$

conv($\mathcal{Y}$)

**marginal**$(\theta)$

Marginal polytope
(Wainwright & Jordan, 2008)

**Differentiable** but **completely dense**
(always in the interior of the polytope)

**Computation:** change semiring

$$x \rightarrow e^x \quad (\max, +) \rightarrow (+, \times)$$

Viterbi → Forward-Backward
CKY → Inside-Outside
DTW → Soft-DTW
max-sum → sum-product (BP)

# Sparse marginal inference?

$$\textbf{marginal}_{\Omega}(\theta) \triangleq \mathbb{E}_p[Y] \quad p = \textbf{argmax}_{\Omega}\left((\langle y, \theta \rangle)_{y \in \mathscr{Y}}\right) \in \triangle^{|\mathscr{Y}|}$$



$\theta$

$\textbf{marginal}_{\Omega}(\theta)$

$\mathrm{conv}(\mathscr{Y})$

# Sparse marginal inference?

$$\mathbf{marginal}_{\Omega}(\theta) \triangleq \mathbb{E}_p[Y] \quad p = \mathbf{argmax}_{\Omega}\left((\langle y, \theta \rangle)_{y \in \mathcal{Y}}\right) \in \triangle^{|\mathcal{Y}|}$$

Can we use $\Omega(p) = \dfrac{1}{2}\|p\|^2$ ?



$\theta$

$\mathbf{marginal}_{\Omega}(\theta)$

$\mathrm{conv}(\mathcal{Y})$

# Sparse marginal inference?

$$\mathbf{marginal}_{\Omega}(\theta) \triangleq \mathbb{E}_p[Y] \quad p = \mathbf{argmax}_{\Omega}\left((\langle y, \theta \rangle)_{y \in \mathcal{Y}}\right) \in \triangle^{|\mathcal{Y}|}$$

Can we use $\Omega(p) = \frac{1}{2}\|p\|^2$ ?



$\theta$

$\mathbf{marginal}_{\Omega}(\theta)$

$\mathrm{conv}(\mathcal{Y})$

**No longer a semiring change in general**

# Sparse marginal inference?

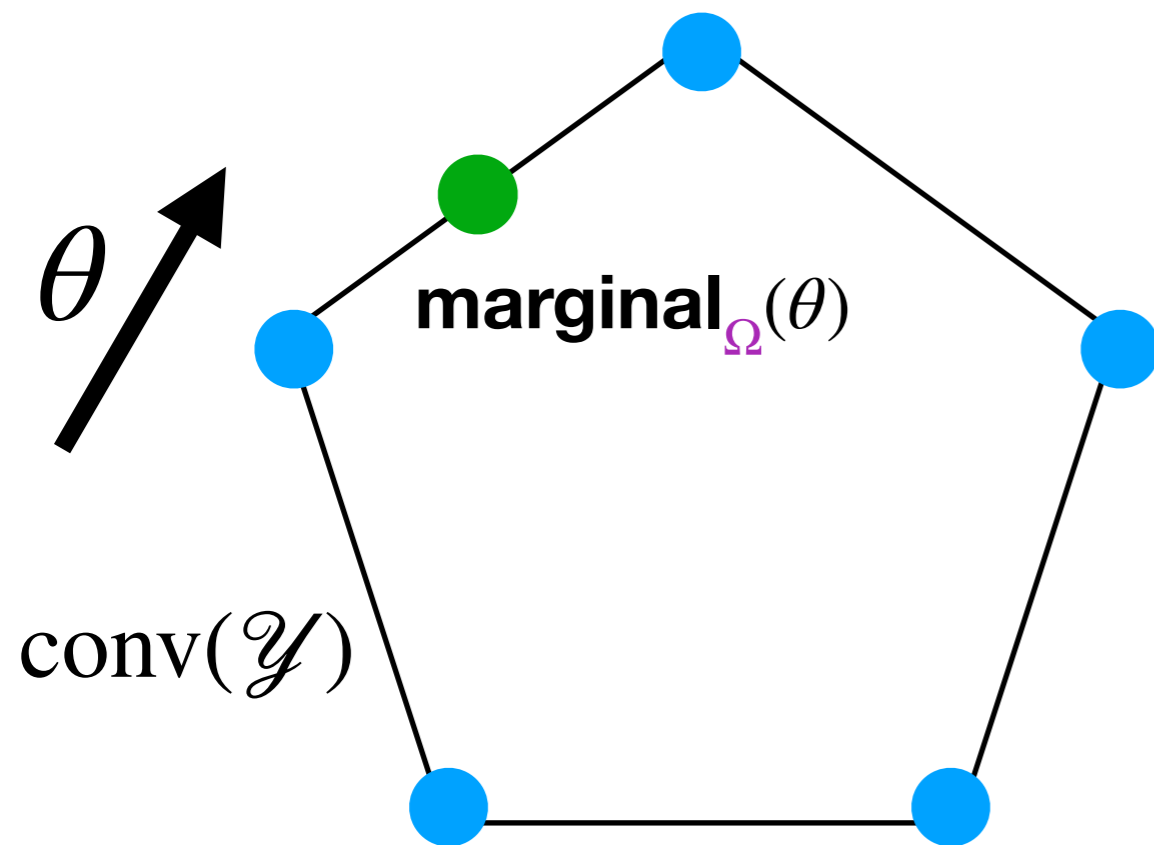$$\mathbf{marginal}_{\Omega}(\theta) \triangleq \mathbb{E}_p[Y] \quad p = \mathbf{argmax}_{\Omega}\left((\langle y, \theta \rangle)_{y \in \mathcal{Y}}\right) \in \triangle^{|\mathcal{Y}|}$$

Can we use $\Omega(p) = \dfrac{1}{2}\|p\|^2$ ?

**No longer a semiring change in general**

**Difficult to compute exactly**

$\theta$

$\mathbf{marginal}_{\Omega}(\theta)$

$\text{conv}(\mathcal{Y})$

# Our proposal for differentiable DP

Mensch & Blondel, ICML 2018

# Our proposal for differentiable DP

- Based on the novel viewpoint of **smoothed max operators**

# Our proposal for differentiable DP

- Based on the novel viewpoint of **smoothed max operators**

- Works for any **shortest path** problem over a **DAG**

# Our proposal for differentiable DP

- Based on the novel viewpoint of **smoothed max operators**

- Works for any **shortest path** problem over a **DAG**

- Enjoys **same big-O complexity** as regular DP

# Our proposal for differentiable DP

- Based on the novel viewpoint of **smoothed max operators**

- Works for any **shortest path** problem over a **DAG**

- Enjoys **same big-O complexity** as regular DP

- **Sparse solutions** when using quadratic regularization

# Our proposal for differentiable DP

- Based on the novel viewpoint of **smoothed max operators**

- Works for any **shortest path** problem over a **DAG**

- Enjoys **same big-O complexity** as regular DP

- **Sparse solutions** when using quadratic regularization

- **Probabilistic interpretation**

# Our proposal for differentiable DP

- Based on the novel viewpoint of **smoothed max operators**

- Works for any **shortest path** problem over a **DAG**

- Enjoys **same big-O complexity** as regular DP

- **Sparse solutions** when using quadratic regularization

- **Probabilistic interpretation**

- **Unified** and **numerically stable** implementation (computations directly in log-domain!)

# Smoothed max operators

# Smoothed max operators

Recall the definition of differentiable **argmax** operator

$$\mathbf{argmax}_{\Omega}(\theta) \triangleq \arg\max_{p \in \triangle^m} \langle p, \theta \rangle - \Omega(p) \in \triangle^m$$

# Smoothed max operators

Recall the definition of differentiable **argmax** operator

$$\mathbf{argmax}_{\Omega}(\theta) \triangleq \arg\max_{p \in \triangle^m} \langle p, \theta \rangle - \Omega(p) \in \triangle^m$$

Similarly we define the smoothed **max** operator (Nesterov, 2005)

$$\max_{\Omega}(\theta) \triangleq \max_{p \in \triangle^m} \langle p, \theta \rangle - \Omega(p) \in \mathbb{R}$$

# Smoothed max operators

Recall the definition of differentiable **argmax** operator

$$\mathbf{argmax}_{\Omega}(\theta) \triangleq \arg\max_{p \in \triangle^m} \langle p, \theta \rangle - \Omega(p) \in \triangle^m$$

Similarly we define the smoothed **max** operator (Nesterov, 2005)

$$\max_{\Omega}(\theta) \triangleq \max_{p \in \triangle^m} \langle p, \theta \rangle - \Omega(p) \in \mathbb{R}$$

From the duality between smoothness and strong convexity

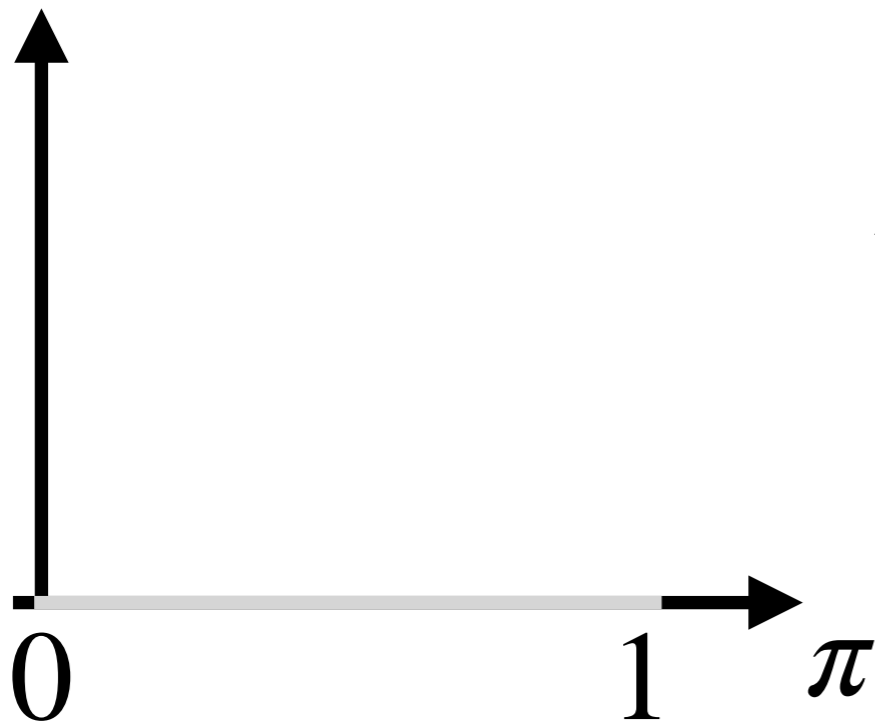Strongly convex $\Omega$ over $\triangle$ $\iff$ Smooth $\max_{\Omega}$

# Examples

$$\max_\Omega(\theta) \triangleq \max_{p \in \triangle^m} \langle p, \theta \rangle - \Omega(p)$$

**Unregularized**

$$\Omega(p) = 0$$

**Regularization**

$-\Omega([\pi, 1 - \pi])$

**Smoothed max**

$\max_\Omega([t, 0])$

0          1     $\pi$

$t$

# Examples

$$\max{}_{\Omega}(\theta) \triangleq \max_{p \in \triangle^m} \langle p, \theta \rangle - \Omega(p)$$

**Unregularized**

**Shannon (negative) entropy**

$$\Omega(p) = 0$$

$$\Omega(p) = \sum_i p_i \log p_i$$

**Regularization**

**Smoothed max**

$-\Omega([\pi, 1 - \pi])$

$\max{}_{\Omega}([t, 0])$

# Examples

$$\max_{\Omega}(\theta) \triangleq \max_{p \in \triangle^m} \langle p, \theta \rangle - \Omega(p)$$
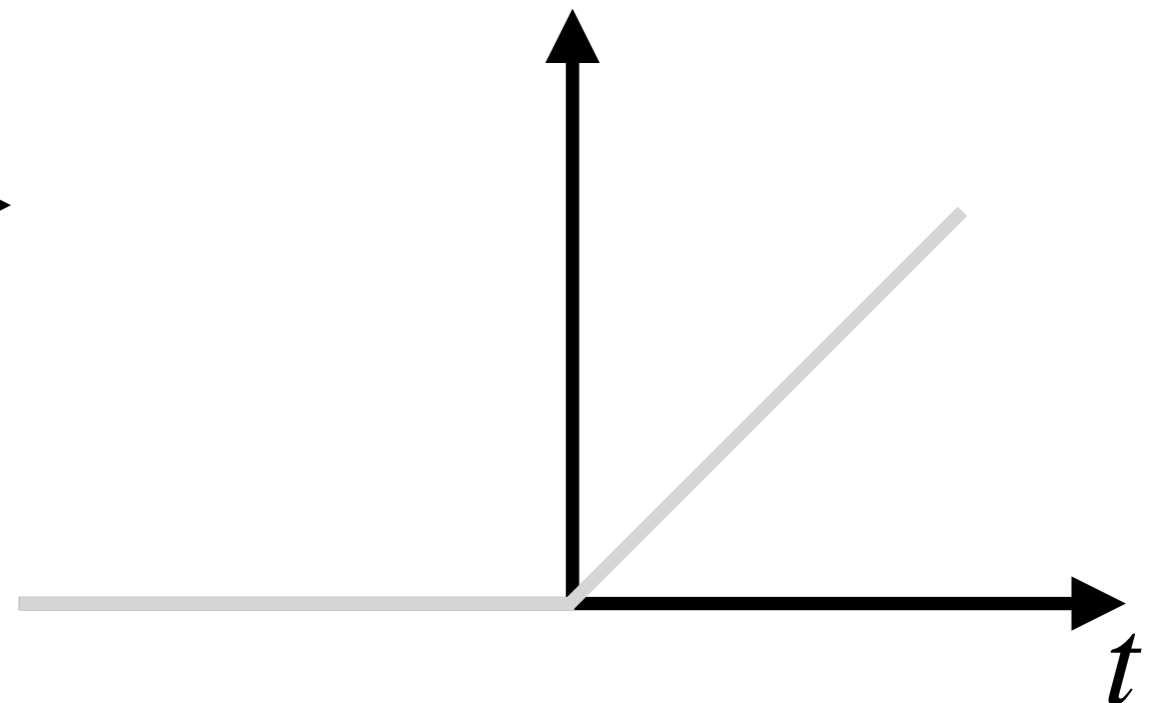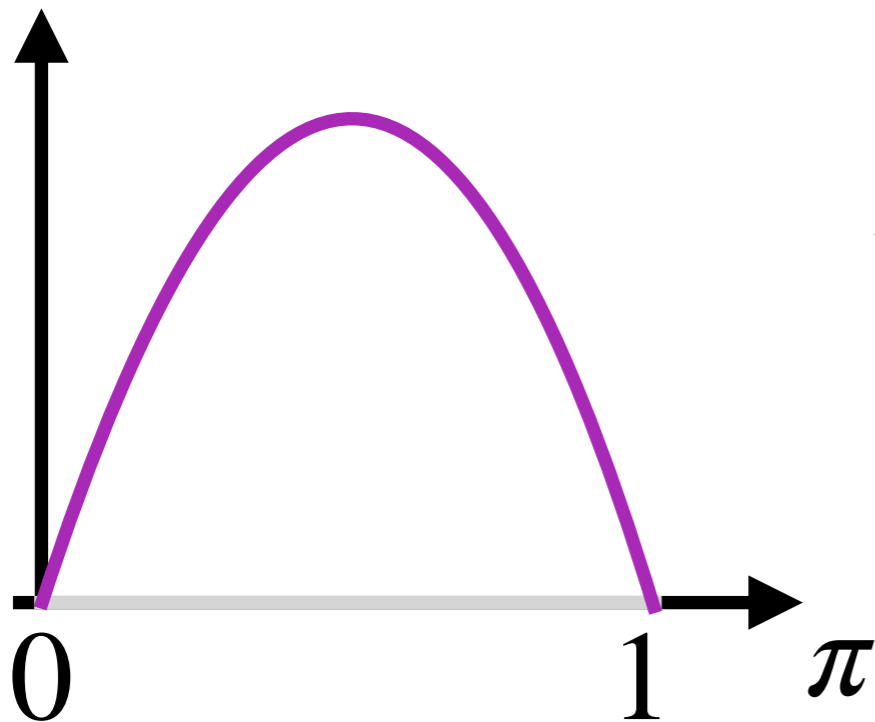
**Unregularized**

$$\Omega(p) = 0$$

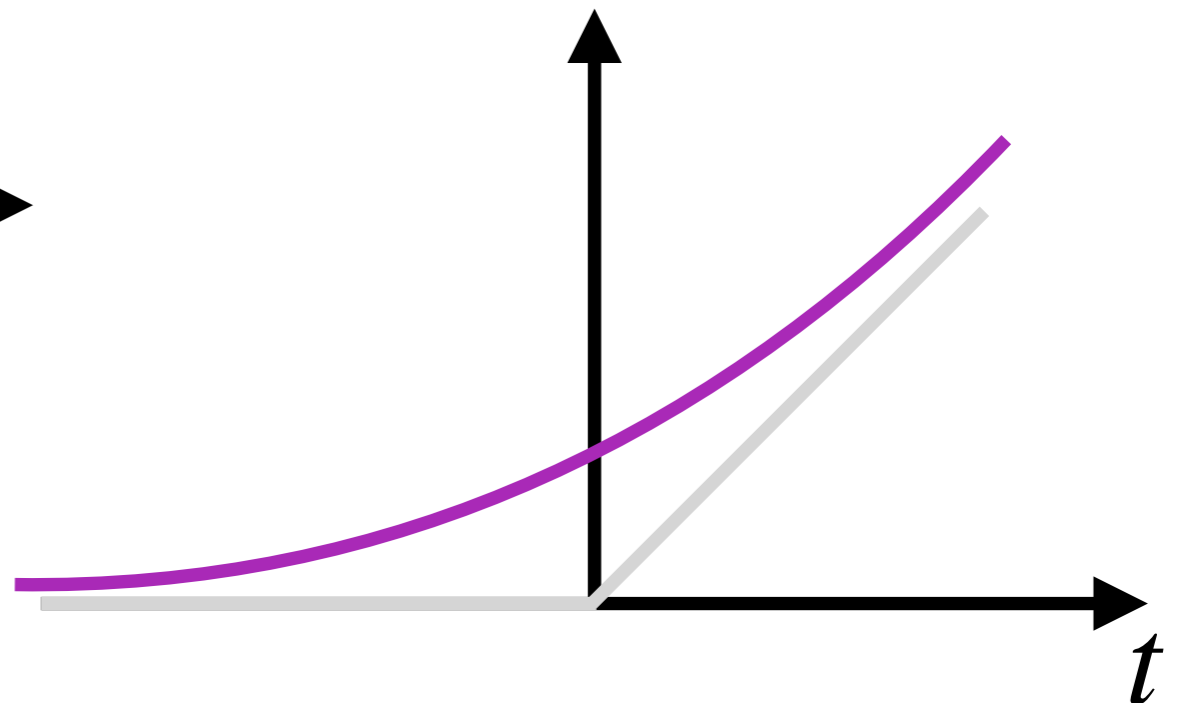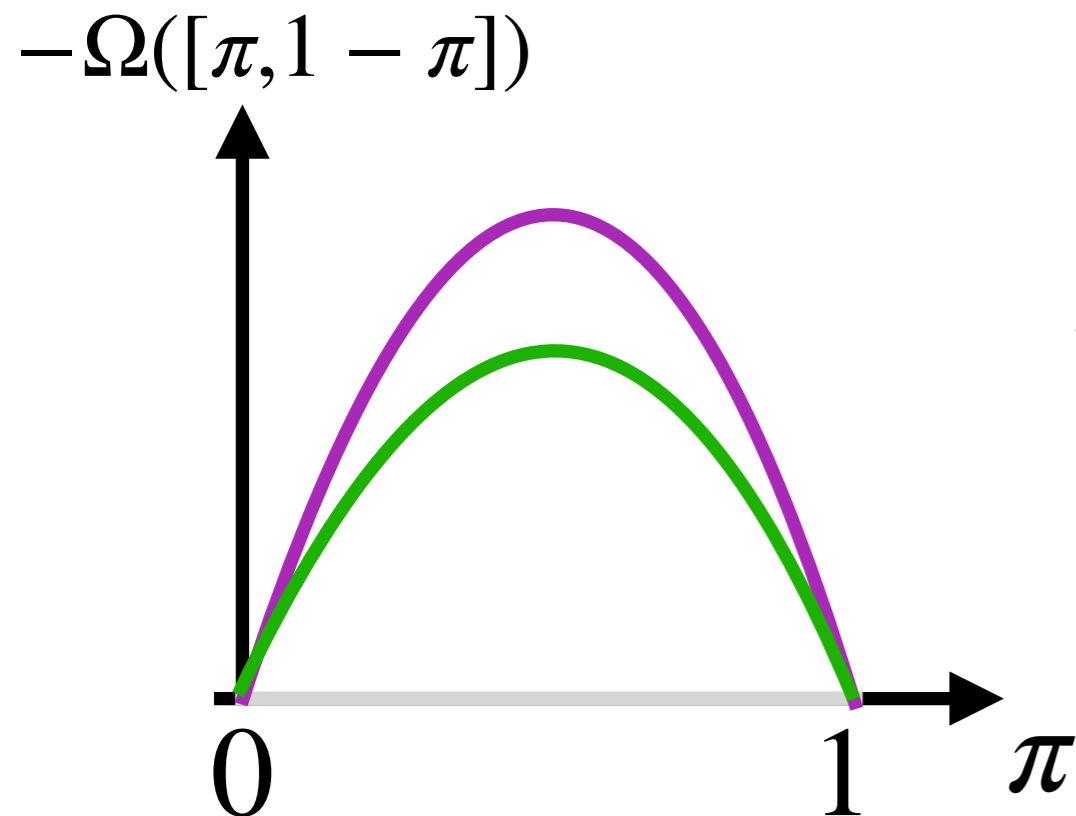**Shannon (negative) entropy**

$$\Omega(p) = \sum_i p_i \log p_i$$

**Gini (negative) index**

$$\Omega(p) = \frac{1}{2}(\|p\|^2 - 1)$$

**Regularization**

$-\Omega([\pi, 1 - \pi])$

**Smoothed max**

$\max_{\Omega}([t, 0])$



**Exactly 0**

# **Smoothed** Bellman's recursion

$$v_{t,i}(\theta) = \max_{\Omega}((v_{t-1,j}(\theta) + \theta_{t,i,j})_{j \in [S]})$$

(max, +)

↓

(max$_{\Omega}$, +)

$t - 1$    $t$    $t + 1$



$v_{t-1,1}(\theta)$    $\theta_{t,1,1}$    $v_{t,1}(\theta)$

$\theta_{t,1,2}$

... $v_{t-1,2}(\theta)$    ...

$\theta_{t,1,3}$

$v_{t-1,2}(\theta)$

Forward pass

# **Smoothed** DP value

# **Smoothed** DP value

$$\mathrm{DP}_{\Omega}(\theta) \leq \max_{\Omega}((\langle y, \theta \rangle)_{y \in \mathcal{Y}})$$

# Probabilistic backpointers

Replace back pointers with **distribution** over states

$$q_{t,i}(\theta) = \mathbf{argmax}_\Omega((v_{t-1,j}(\theta) + \theta_{t,i,j})_{j \in [S]}) \in \triangle^S$$



$t-1$     $t$     $t+1$

$q_{t,1}(\theta)$   $v_{t,1}(\theta)$

$q_{t,2}(\theta)$   $v_{t,2}(\theta)$

$q_{t,3}(\theta)$   $v_{t,3}(\theta)$

Forward pass

# Random walk

Random walk (finite Markov chain) defines
a distribution p over paths



Each time step t has its own transition matrix $Q_t \in \mathbb{R}^{S \times S}$

# Random walk

Sampling is easy.

How to compute **expectation** $\mathbb{E}_p[Y]$ ?

$T$



START

END

Each time step t has its own transition matrix $Q_t \in \mathbb{R}^{S \times S}$

# Gradient = Expected path

$$\nabla \mathrm{DP}_\Omega(\theta) = \mathbb{E}_p[Y] \in \mathrm{conv}(\mathscr{Y})$$

# Gradient = Expected path

$$\nabla \text{DP}_\Omega(\theta) = \mathbb{E}_p[Y] \in \text{conv}(\mathcal{Y})$$

Can compute $\mathbb{E}_p[Y]$ at the same cost as computing $\text{DP}_\Omega(\theta)$ by **backpropagation**

# Gradient = Expected path

<u>Proposition (Mensch & Blondel, 2018)</u>   (See also Eisner, 2016)

$$\nabla \text{DP}_\Omega(\theta) = \mathbb{E}_p[Y] \in \text{conv}(\mathcal{Y})$$

Can compute $\mathbb{E}_p[Y]$ at the same cost as computing $\text{DP}_\Omega(\theta)$ by **backpropagation**

**Intractable sum
if computed naively**

For $\Omega$ = negative entropy, we have

$$\nabla \text{DP}_\Omega(\theta) = \mathbb{E}_p[Y] = \frac{\sum_{y \in \mathcal{Y}} \exp\langle y, \theta\rangle y}{Z(\theta)}$$

# Backpropagation

$$E \triangleq \mathbb{E}_p[Y] \qquad \mathbf{e}_{t,\cdot,j} = \mathbf{q}_{t+1,\cdot,j} \circ (\mathbf{e}_{t+1,\cdot,j}^\top \mathbf{1})$$

# Backpropagation

$$E \triangleq \mathbb{E}_p[Y] \qquad \mathbf{e}_{t,\cdot,j} = \mathbf{q}_{t+1,\cdot,j} \circ (\mathbf{e}_{t+1,\cdot,j}^\top \mathbf{1})$$



Up to **12x faster** in our experiments compared to PyTorch's autodiff

Backpropagation

# Theoretical results

# Theoretical results

1. $DP_\Omega(\theta)$ is convex

   Proof uses that $x \leq y \Rightarrow \max_\Omega(x) \leq \max_\Omega(y)$

# Theoretical results

1. $DP_{\Omega}(\theta)$ is convex

   Proof uses that $x \leq y \Rightarrow \max_{\Omega}(x) \leq \max_{\Omega}(y)$

2. Approximation error

   N: #nodes in DAG
   L, U: constants that depend on $\Omega$

   $$(N-1) \, L \leq DP_{\Omega}(\theta) - DP(\theta) \leq (N-1) \, U$$

# Theoretical results

1. $DP_\Omega(\theta)$ is convex

   Proof uses that $x \leq y \Rightarrow \max_\Omega(x) \leq \max_\Omega(y)$

2. Approximation error

   N: #nodes in DAG
   L, U: constants that depend on $\Omega$

$$(N-1)\, L \leq DP_\Omega(\theta) - DP(\theta) \leq (N-1)\, U$$

3. $DP_\Omega(\theta) = \max_\Omega((\langle y, \theta \rangle)_{y \in \mathcal{Y}}) \quad \Leftrightarrow \quad \Omega = -H$ (Shannon's negentropy)

   Proof reduces to showing that $\max_{-H}$ is the only $\max_\Omega$ supporting
   **associativity**, i.e., $\max_{-H}(x, \max_{-H}(y, z)) = \max_{-H}(\max_{-H}(x, y), z)$

# Structured prediction losses

# Structured prediction losses

## Training time

Structured perceptron loss (Collins, 2002)

$$\max_{y \in \mathcal{Y}} \langle \theta, y \rangle - \langle \theta, y_{\text{true}} \rangle$$

# Structured prediction losses

## Training time

Structured perceptron loss (Collins, 2002)

$$\max_{y \in \mathcal{Y}} \langle \theta, y \rangle - \langle \theta, y_{\text{true}} \rangle = \text{DP}(\theta) - \langle \theta, y_{\text{true}} \rangle$$

# Structured prediction losses

## Training time

Structured perceptron loss (Collins, 2002)

$$\max_{y \in \mathcal{Y}} \langle \theta, y \rangle - \langle \theta, y_{\text{true}} \rangle = \mathrm{DP}(\theta) - \langle \theta, y_{\text{true}} \rangle$$

Smoothed loss (proposed)

$$\mathrm{DP}_{\Omega}(\theta) - \langle \theta, y_{\text{true}} \rangle$$

# Structured prediction losses

## Training time

Structured perceptron loss (Collins, 2002)

$$\max_{y \in \mathcal{Y}} \langle \theta, y \rangle - \langle \theta, y_{\text{true}} \rangle = \text{DP}(\theta) - \langle \theta, y_{\text{true}} \rangle$$

Smoothed loss (proposed)

$$\text{DP}_{\Omega}(\theta) - \langle \theta, y_{\text{true}} \rangle$$

Entropic regularization $\rightarrow$ CRF loss
Quadratic regularization $\rightarrow$ new loss

# Structured prediction losses

## Training time

Structured perceptron loss (Collins, 2002)

$$\max_{y \in \mathcal{Y}} \langle \theta, y \rangle - \langle \theta, y_{\text{true}} \rangle = \text{DP}(\theta) - \langle \theta, y_{\text{true}} \rangle$$

Smoothed loss (proposed)

$$\text{DP}_{\Omega}(\theta) - \langle \theta, y_{\text{true}} \rangle$$

Entropic regularization $\rightarrow$ CRF loss

Quadratic regularization $\rightarrow$ new loss

## Test time

MAP solution

$$\arg\max_{y \in \mathcal{Y} \subseteq \mathbb{R}^m} \langle y, \theta \rangle$$

# Structured prediction losses

## Training time

Structured perceptron loss (Collins, 2002)

$$\max_{y \in \mathcal{Y}} \langle \theta, y \rangle - \langle \theta, y_{\text{true}} \rangle = \text{DP}(\theta) - \langle \theta, y_{\text{true}} \rangle$$

Smoothed loss (proposed)

$$\text{DP}_{\Omega}(\theta) - \langle \theta, y_{\text{true}} \rangle$$

Entropic regularization → CRF loss
Quadratic regularization → new loss

## Test time

MAP solution

Expected solution

$$\arg\max_{y \in \mathcal{Y} \subseteq \mathbb{R}^m} \langle y, \theta \rangle \qquad \nabla \text{DP}_{\Omega}(\theta) = \mathbb{E}_p[Y]$$

# Structured prediction losses

## Training time

Structured perceptron loss (Collins, 2002)

$$\max_{y \in \mathcal{Y}} \langle \theta, y \rangle - \langle \theta, y_{\text{true}} \rangle = \text{DP}(\theta) - \langle \theta, y_{\text{true}} \rangle$$

Smoothed loss (proposed)

$$\text{DP}_\Omega(\theta) - \langle \theta, y_{\text{true}} \rangle$$

Entropic regularization $\rightarrow$ CRF loss

Quadratic regularization $\rightarrow$ new loss

## Test time

MAP solution

$$\arg\max_{y \in \mathcal{Y} \subseteq \mathbb{R}^m} \langle y, \theta \rangle$$

Expected solution

$$\nabla \text{DP}_\Omega(\theta) = \mathbb{E}_p[Y]$$

Ranking

Sort by probability
(sparse case)

# NER experiments

| S-ORG | O | B-PER | E-PER | O | O | O | O | S-LOC |
|-------|---|-------|-------|---|---|---|---|-------|
| **Apple** | CEO | **Tim** | **Cook** | introduces | new | iphone | in | **Cupertino**. |

Tags: {Location, Organization, Person, Misc}  x  {Singleton, Begin, Inside, End}

# NER experiments

S-ORG    O    B-PER   E-PER     O      O      O    O    S-LOC

**Apple** CEO **Tim Cook** introduces new iphone in **Cupertino**.

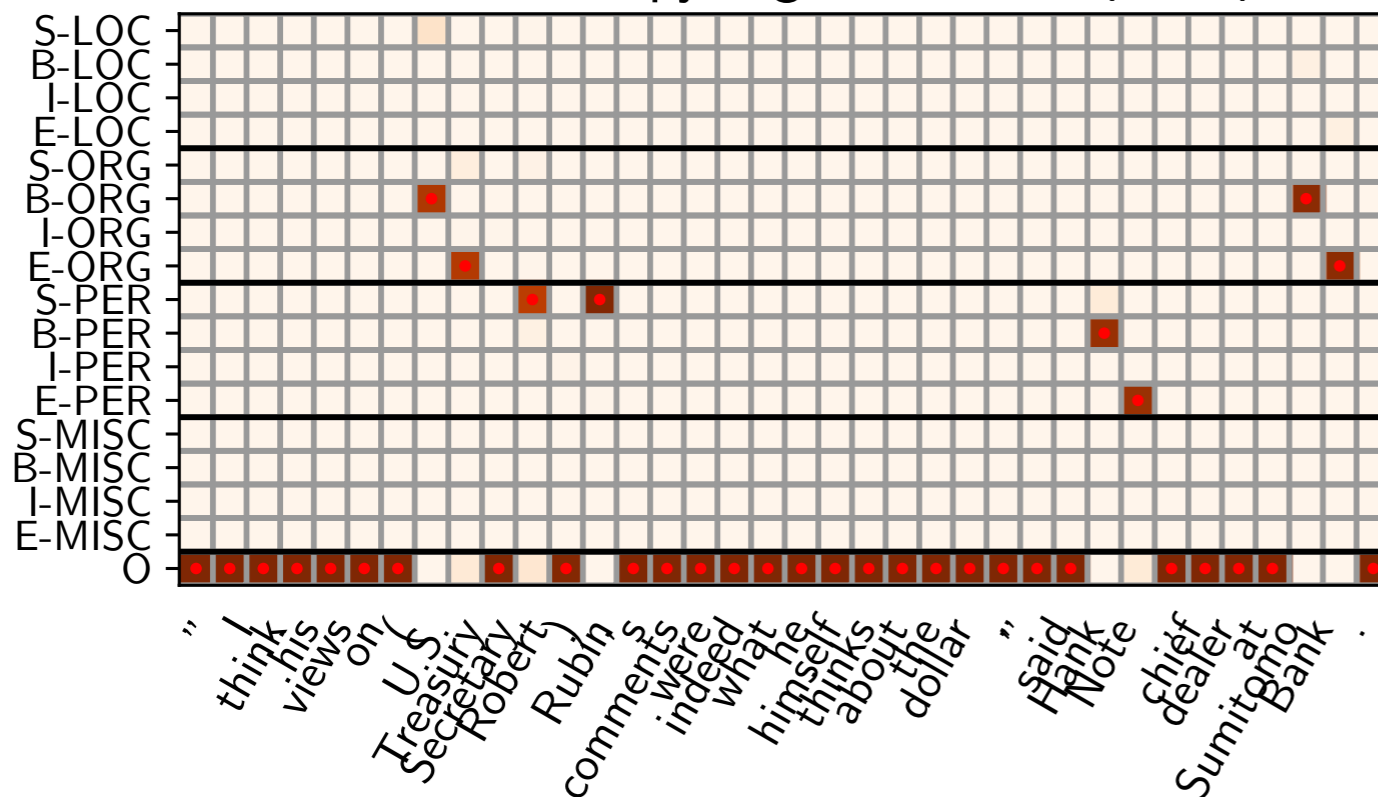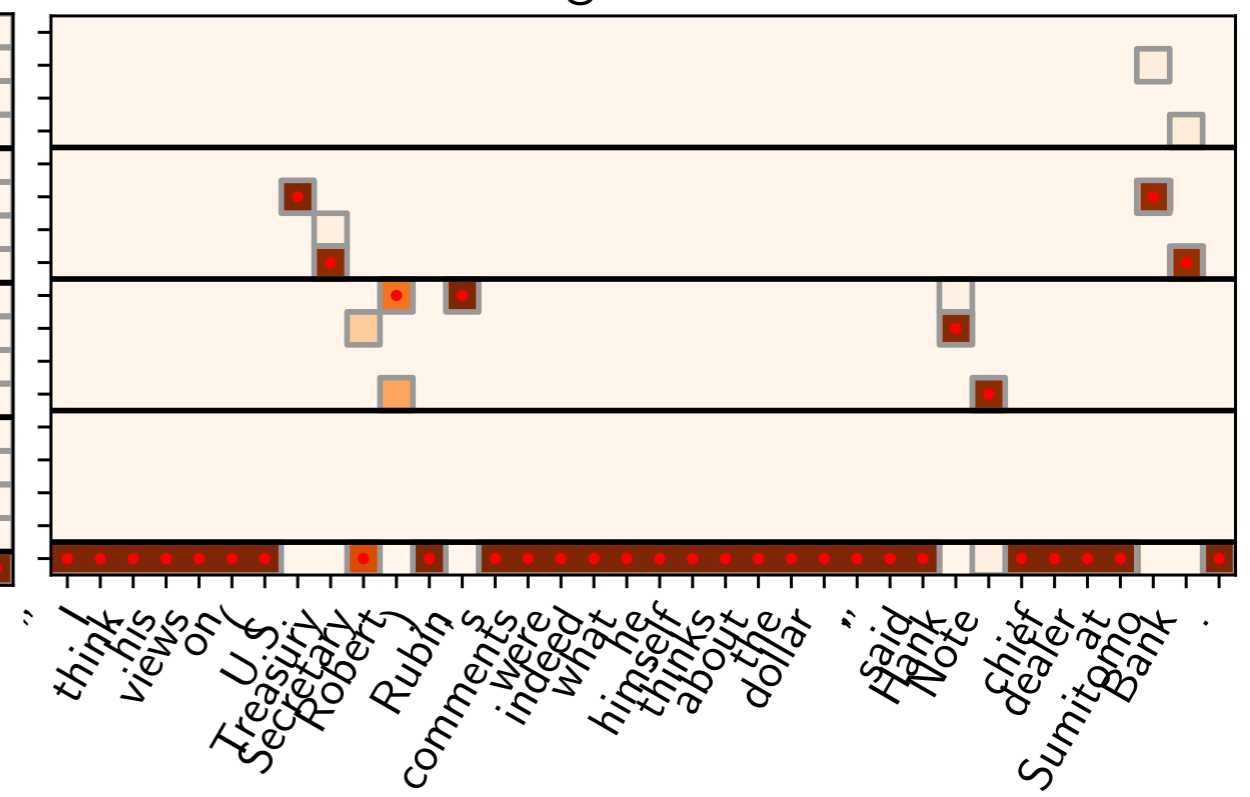Tags: {Location, Organization, Person, Misc}  x  {Singleton, Begin, Inside, End}

## Examples of predicted soft assignments at test time



Entropy regularization (CRF)                        L2 regularization

# NER experiments

F$_1$ score comparison on CoNLL03 NER datasets

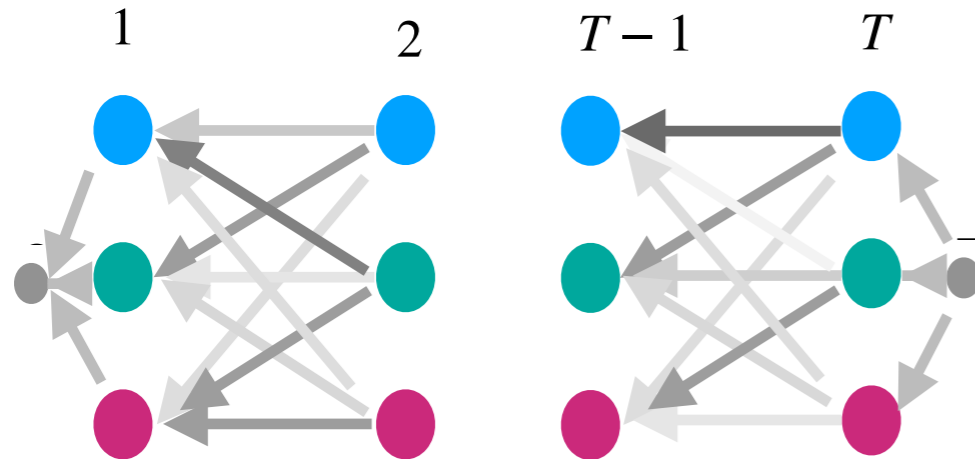|  | English | Spanish | German | Dutch |
|---|---|---|---|---|
| CRF loss (Entropy) | 90.80 | 86.68 | 77.35 | 87.56 |
| Squared norm | 90.86 | 85.51 | 76.01 | 86.58 |
| Lample et al 2016 (CRF loss) | 90.96 | 85.75 | 78.76 | 81.74 |

# NER experiments

F$_1$ score comparison on CoNLL03 NER datasets

. Competitive results with other losses

. **Fast convergence** at train time thanks to **smoothness**

. **Sparse probabilistic model** available at test time!

| Squared norm | 90.86 | 85.51 | 76.01 | 86.58 |
|---|---|---|---|---|
| Lample et al 2016 (CRF loss) | 90.96 | 85.75 | 78.76 | 81.74 |

# Summary of second part

**Smoothing induces a random walk**



a distribution over paths in the DAG

**Gradient = Expected path**

$$\nabla \mathrm{DP}_\Omega(\theta) = \mathbb{E}_p[Y]$$

computed efficiently by backprop

**Entropic regularization = CRF**



**L2 regularization = new sparse model**

# Conclusion

# Conclusion

- The log-sum-exp and softmax are ubiquitous in deep learning

# Conclusion

- The log-sum-exp and softmax are ubiquitous in deep learning

- $\max_\Omega$ and $\text{argmax}_\Omega$ operators provide **drop-in replacement** for them with **sparse** and/or **structured** outputs

# Conclusion

- The log-sum-exp and softmax are ubiquitous in deep learning

- $\max_\Omega$ and $\text{argmax}_\Omega$ operators provide **drop-in replacement** for them with **sparse** and/or **structured** outputs

- Induce a **probabilistic perspective**

# Conclusion

- The log-sum-exp and softmax are ubiquitous in deep learning

- max$_\Omega$ and argmax$_\Omega$ operators provide **drop-in replacement** for them with **sparse** and/or **structured** outputs

- Induce a **probabilistic perspective**

- Many more potential applications to explore