# Interpretable Structure Induction Via **Sparse Attention**

Ben Peters    Instituto de Telecomunicações

→ **Vlad Niculae**    IT

André Martins    IT & Unbabel

# Sparse linear models are more interpretable...

| LogL | Collocation | Sense |
|------|-------------|-------|
| 10.12 | *plant* growth | ⇒ A |
| 9.68 | car (within ±$k$ words) | ⇒ B |
| 9.64 | *plant* height | ⇒ A |
| 9.61 | union (within ±$k$ words) | ⇒ B |
| 9.54 | equipment (within ±$k$ words) | ⇒ B |
| 9.51 | assembly *plant* | ⇒ B |
| 9.50 | nuclear *plant* | ⇒ B |
| 9.31 | flower (within ±$k$ words) | ⇒ A |
| 9.24 | job (within ±$k$ words) | ⇒ B |
| 9.03 | fruit (within ±$k$ words) | ⇒ A |
| 9.02 | *plant* species | ⇒ A |
| ... | ... | |

**Final decision list for *plant* (abbreviated)**

Decision list from Yarowsky (1995)

# **Sparse** linear models are more interpretable... but we use bigger models today!



Decision list from Yarowsky (1995)

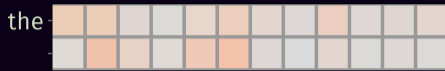# Neural Attention Mechanisms

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

# Neural Attention Mechanisms



$p = 0.0003\%$

$p = 0.15\%$

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

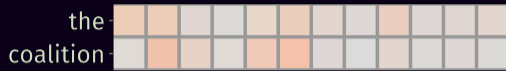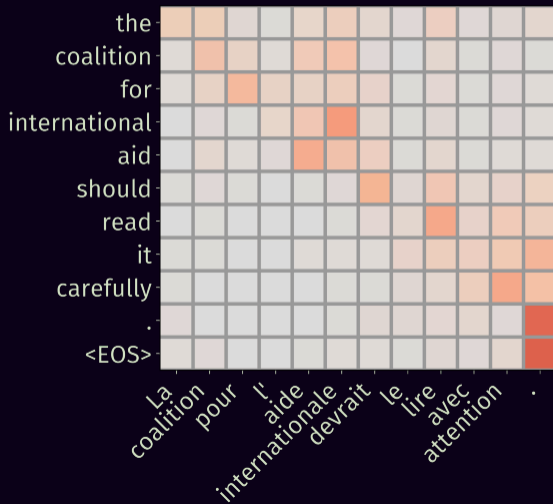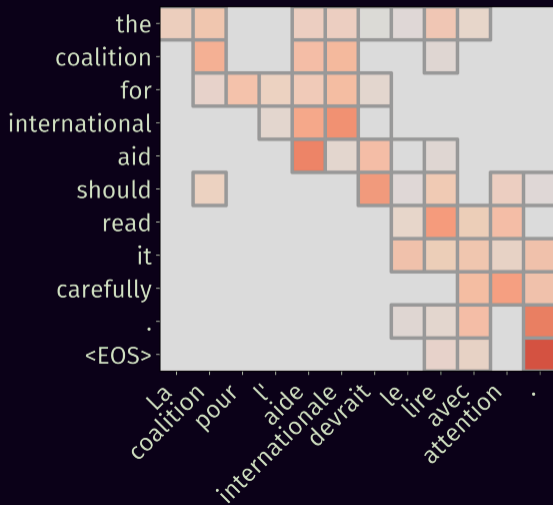# Neural Attention Mechanisms

# Neural Attention Mechanisms

# Neural Attention Mechanisms

# Neural Attention Mechanisms

# Sparse Neural Attention



sparsemax
(Martins & Astudillo, 2016)

# Sparse Neural Attention



$p$ = 0.0%!

sparsemax
(Martins & Astudillo, 2016)

# Structured & Sparse Attention



fusedmax
(Niculae & Blondel, 2017)

# Structured & Sparse Attention



exactly equal!

fusedmax
(Niculae & Blondel, 2017)

# Smoothed Max Operators

$$\text{softmax}(\boldsymbol{\theta}) = \boldsymbol{p}$$

# Smoothed Max Operators

$$\Pi_\Omega(\boldsymbol{\theta}) = \arg\max_{\boldsymbol{p} \in \triangle} \boldsymbol{p}^\top \boldsymbol{\theta} - \Omega(\boldsymbol{p})$$
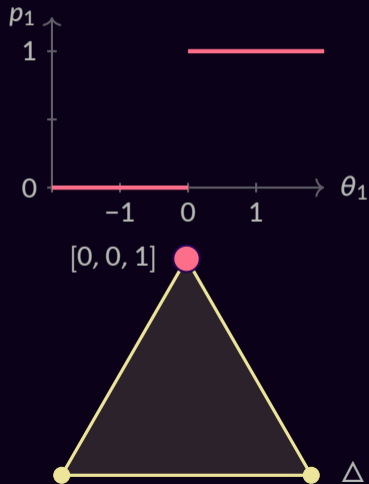
●    argmax: $\Omega(\boldsymbol{p}) = 0$

# Smoothed Max Operators

$$\Pi_\Omega(\boldsymbol{\theta}) = \underset{\boldsymbol{p} \in \triangle}{\arg\max}\, \boldsymbol{p}^\top \boldsymbol{\theta} - \Omega(\boldsymbol{p})$$

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

# Smoothed Max Operators

$$\Pi_\Omega(\boldsymbol{\theta}) = \arg\max_{\boldsymbol{p} \in \triangle} \boldsymbol{p}^\top \boldsymbol{\theta} - \Omega(\boldsymbol{p})$$

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

- ●   softmax:   $\Omega(\boldsymbol{p}) = \sum_j p_j \log p_j$

- ● sparsemax:   $\Omega(\boldsymbol{p}) = \tfrac{1}{2}\|\boldsymbol{p}\|_2^2$

# Smoothed Max Operators

$$\Pi_\Omega(\boldsymbol{\theta}) = \arg\max_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^\top\boldsymbol{\theta} - \Omega(\boldsymbol{p})$$



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

# Constrained Attention

$$\arg\max_{\substack{\boldsymbol{p}\in\triangle \\ \boldsymbol{a}\leq\boldsymbol{p}\leq\boldsymbol{b}}} \boldsymbol{p}^{\top}\boldsymbol{\theta} - \Omega_1(\boldsymbol{p})$$

$$= \arg\max_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^{\top}\boldsymbol{\theta} - \underbrace{\Omega(\boldsymbol{p})}_{:=\Omega_1 + \mathsf{Id}_{[a,b]}}$$



*e.g.*, fertility constraints for NMT

(Kreutzer & Martins, 18)
(Malaviya et al, 18)

# Structured Attention & Graphical Models

# Structured Attention & Graphical Models

- ● **argmax** $\arg\max\limits_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^\top\boldsymbol{\theta}$

- ● **softmax** $\arg\max\limits_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^\top\boldsymbol{\theta} + \mathrm{H}(\boldsymbol{p})$

- ● **sparsemax** $\arg\max\limits_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^\top\boldsymbol{\theta} - \frac{1}{2}\|\boldsymbol{p}\|^2$



$\triangle$

**argmax** $\underset{p \in \triangle}{\arg\max} \, p^\top \theta$

**MAP** $\underset{\mu \in \mathcal{M}}{\arg\max} \, \mu^\top \eta$

**softmax** $\underset{p \in \triangle}{\arg\max} \, p^\top \theta + \mathsf{H}(p)$

**marginals** $\underset{\mu \in \mathcal{M}}{\arg\max} \, \mu^\top \eta + \widetilde{\mathsf{H}}(\mu)$

**sparsemax** $\underset{p \in \triangle}{\arg\max} \, p^\top \theta - \frac{1}{2}\|p\|^2$

**SparseMAP** $\underset{\mu \in \mathcal{M}}{\arg\max} \, \mu^\top \eta - \frac{1}{2}\|\mu\|^2$

# Structured Attention for Alignments

NLI

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

input

(P, H)

output

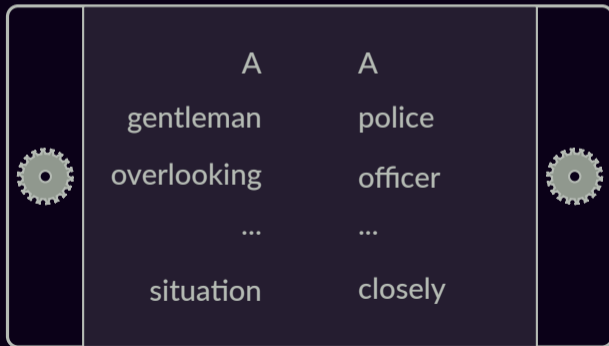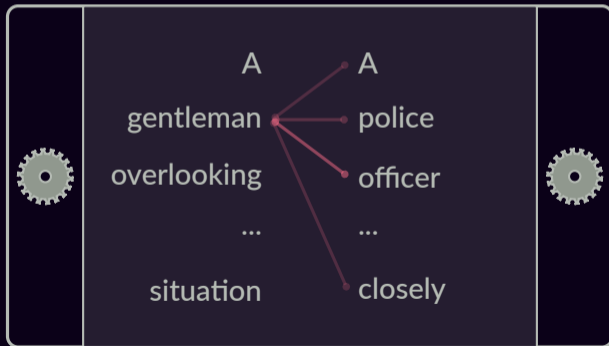| A | A | |
|---|---|---|
| gentleman | police | entails |
| overlooking | officer | contradicts |
| ... | ... | neutral |
| situation | closely | |

Model: ESIM (Chen, 16)

# 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 (Chen, 16)

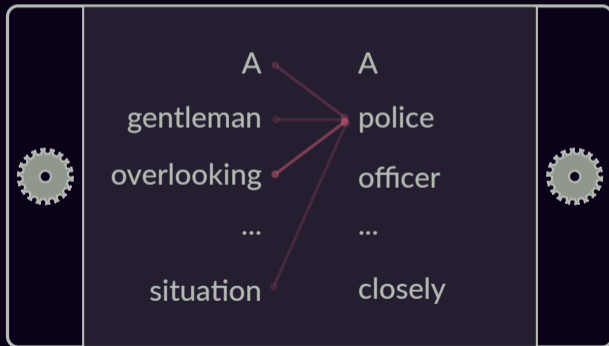output

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)



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

output

entails

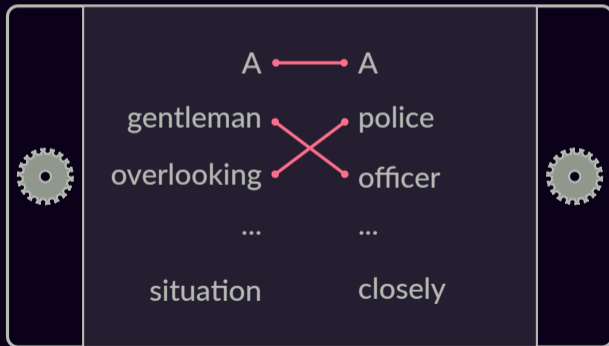contradicts

neutral

Model: ESIM (Chen, 16)

# Structured Attention for Alignments

NLI

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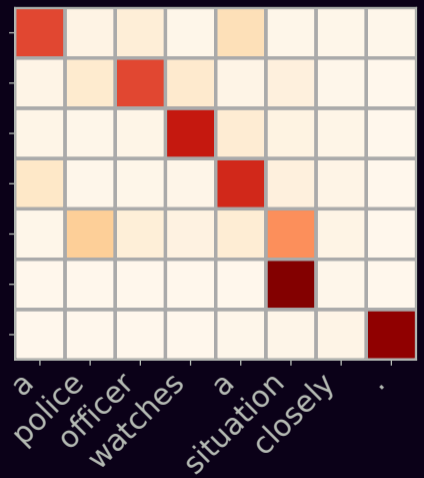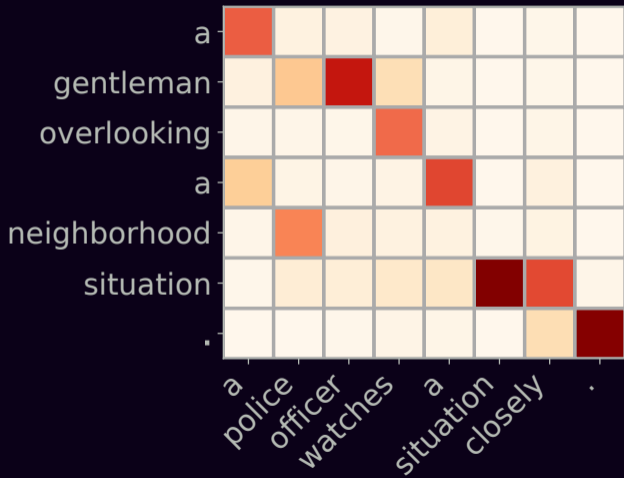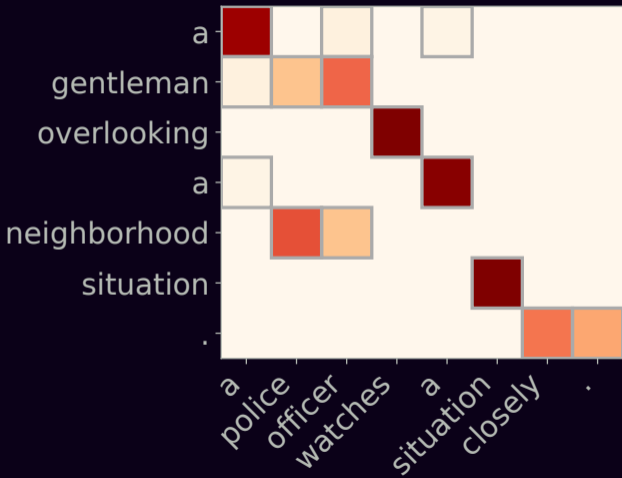
input

(P, H)



output

A ——— A

gentleman        police

overlooking        officer

...        ...

situation        closely
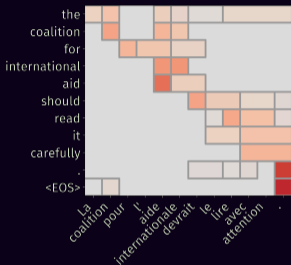
entails

contradicts

neutral

Proposed model: global matching
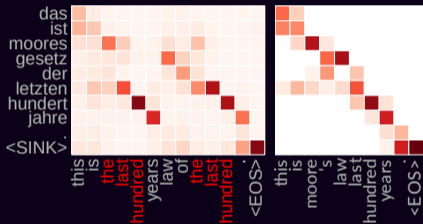
# Summary: Neural attention with...



**structured sparsity**
(*e.g.* fusedmax)

**constraints**
(*e.g.* csparsemax — fertility)

**structure**
(*e.g.* SparseMAP alignments)

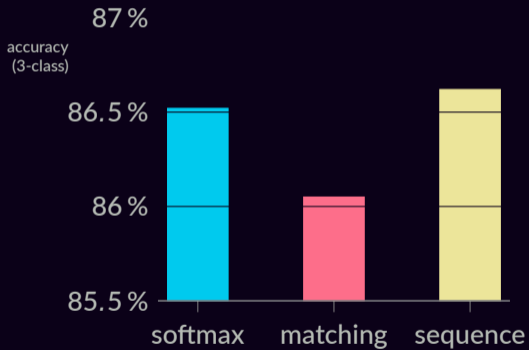and dynamic computation graphs with structured latent variables! (Friday 15:36 in 3B)

✉ vlad@vene.ro  ⌂ github.com/vene/sparsemap  🏠 https://vene.ro  🐦 @vnfrombucharest

# Acknowledgements

Some icons by Dave Gandy and Freepik via flaticon.com.

# Extra slides

**SNLI**

accuracy (3-class)

87 %

86.5 %

86 %

85.5 %

softmax  matching  sequence

**MultiNLI**

76.5 %

76 %

75.5 %

75 %

softmax  matching  sequence