

Interpretable Structure Induction Via **Sparse Attention**

Ben Peters Instituto de Telecomunicações


→ **Vlad Niculae** IT

André Martins IT & Unbabel

Sparse linear models are more interpretable...

Final decision list for <i>plant</i> (abbreviated)		
LogL	Collocation	Sense
10.12	<i>plant</i> growth	\Rightarrow A
9.68	car (within $\pm k$ words)	\Rightarrow B
9.64	<i>plant</i> height	\Rightarrow A
9.61	union (within $\pm k$ words)	\Rightarrow B
9.54	equipment (within $\pm k$ words)	\Rightarrow B
9.51	assembly <i>plant</i>	\Rightarrow B
9.50	nuclear <i>plant</i>	\Rightarrow B
9.31	flower (within $\pm k$ words)	\Rightarrow A
9.24	job (within $\pm k$ words)	\Rightarrow B
9.03	fruit (within $\pm k$ words)	\Rightarrow A
9.02	<i>plant</i> species	\Rightarrow A
...	...	

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

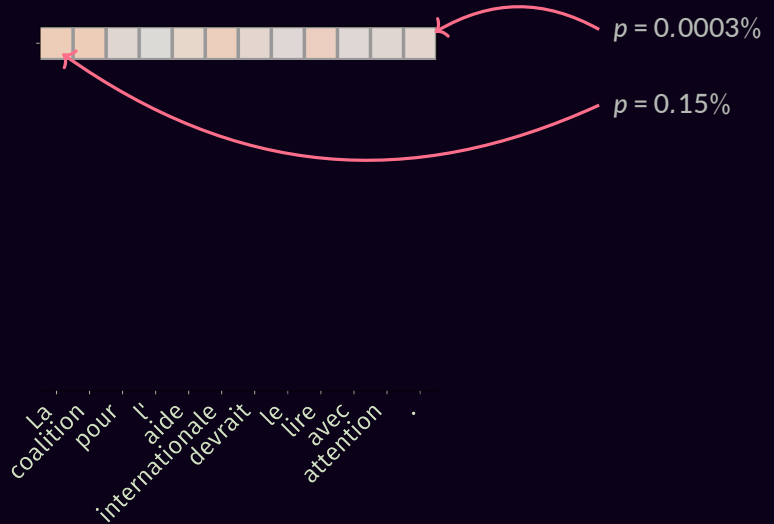
Final decision		transformer_step_200000.pt Properties		abbreviated)	
LogL	Co	Basic	Permissions	Sense	
10.12	pla	 Name: transformer_step_200000.pt Type: Program (application/octet-stream) Size: 1,3 GB (1283080693 bytes) Parent Folder: /media/hdd/transf Accessed: qui 25 out 2018 15:18:23 WEST Modified: sex 12 out 2018 00:13:32 WEST		⇒ A	
9.68	car			⇒ B	
9.64	pla			⇒ A	
9.61	uni			⇒ B	
9.54	equ			⇒ B	
9.51	ass			⇒ B	
9.50	nuc			⇒ B	
9.31	flow			⇒ A	
9.24	job			⇒ B	
9.03	fruit (within 10 words)			⇒ A	
9.02	plant species			⇒ A	
...	...				

(Decision list from Yarowsky (1995))

Neural Attention Mechanisms

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

Neural Attention Mechanisms



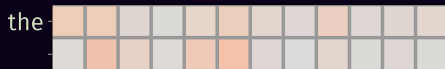
Neural Attention Mechanisms

the



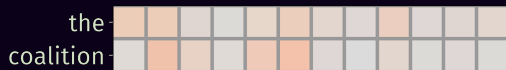
La
coalition
pour
l'
aide
internationale
devrait
le
lire
avec
attention
.

Neural Attention Mechanisms



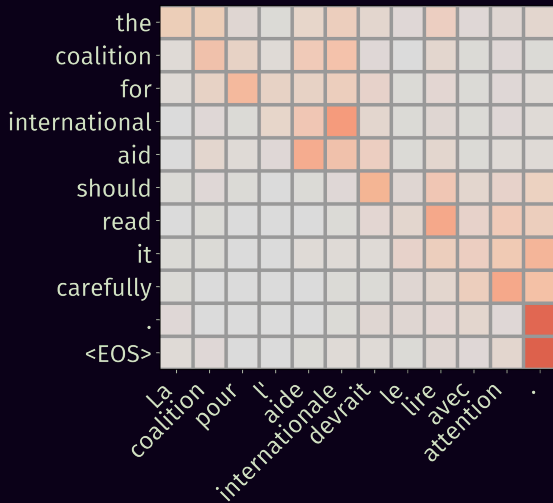
La coalition pour l'aide internationale devrait le lire avec attention .

Neural Attention Mechanisms

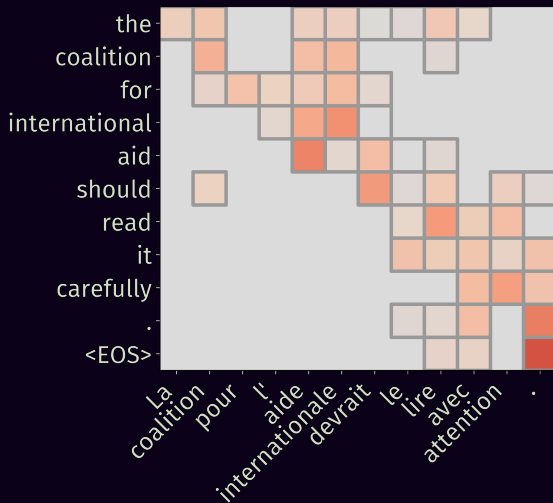


La
coalition
pour
l'
aide
internationale
devrait
le
lire
avec
attention
.

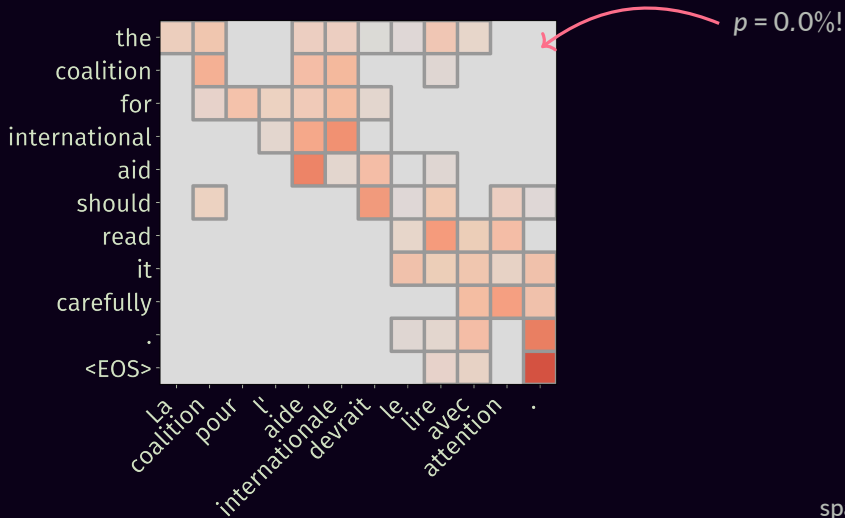
Neural Attention Mechanisms



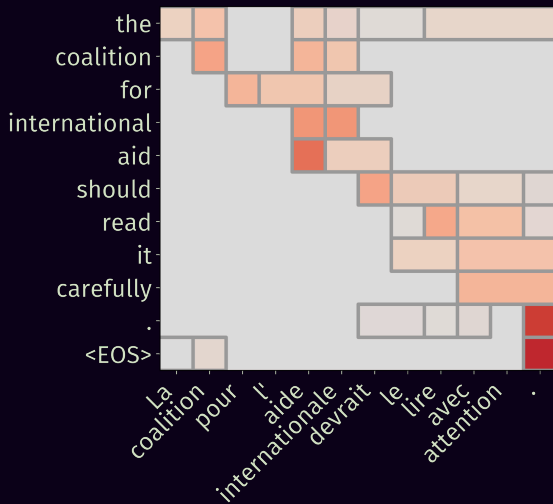
Sparse Neural Attention



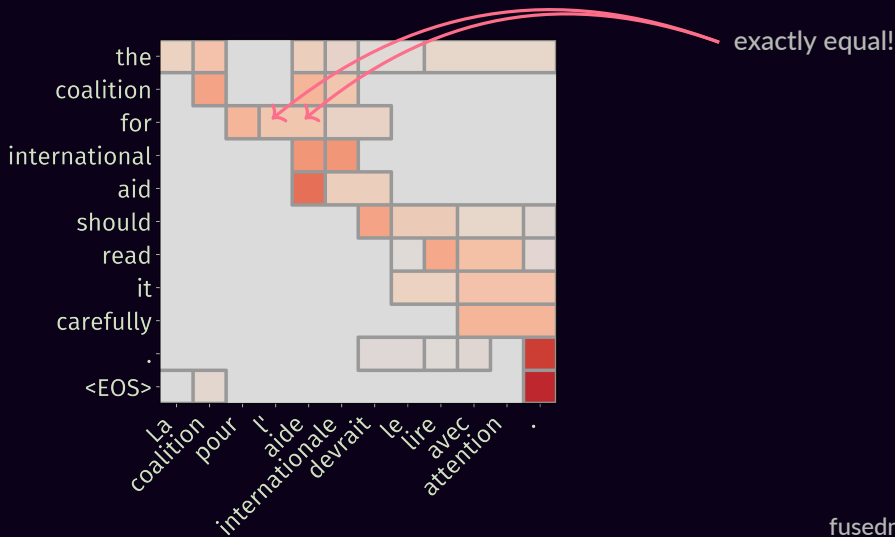
Sparse Neural Attention



Structured & Sparse Attention

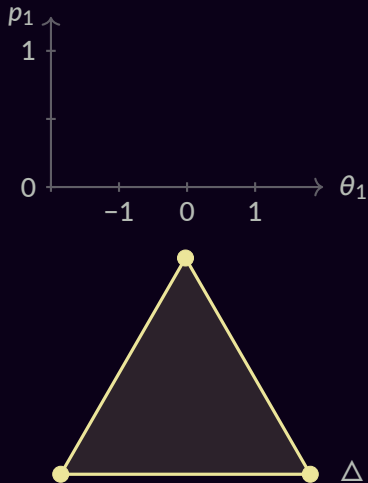


Structured & Sparse Attention



Smoothed Max Operators

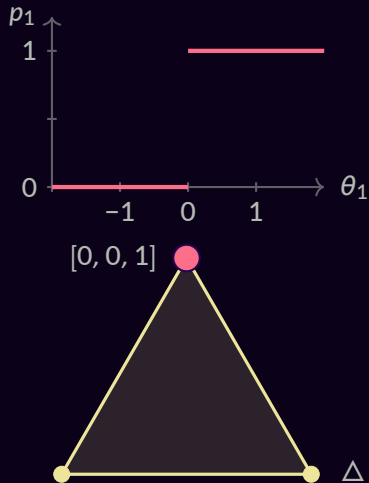
$$\Pi_{\Omega}(\boldsymbol{\theta}) = \arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^{\top} \boldsymbol{\theta} - \Omega(\boldsymbol{p})$$



Smoothed Max Operators

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

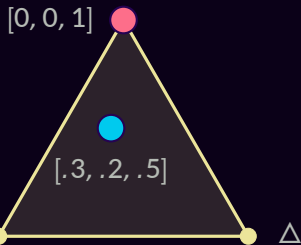
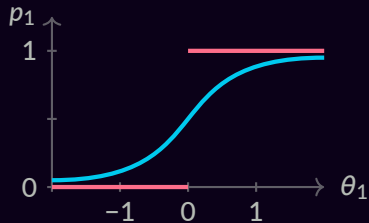
- argmax: $\Omega(\mathbf{p}) = 0$



Smoothed Max Operators

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

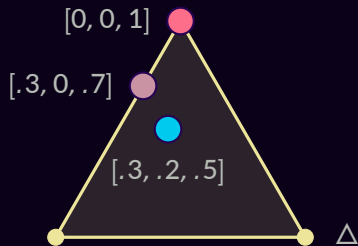
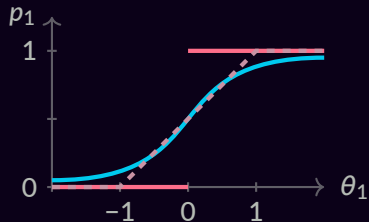
- argmax: $\Omega(\mathbf{p}) = 0$
- softmax: $\Omega(\mathbf{p}) = \sum_j p_j \log p_j$



Smoothed Max Operators

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

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



Smoothed Max Operators

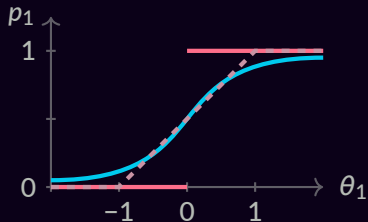
$$\Pi_{\Omega}(\boldsymbol{\theta}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^{\top} \boldsymbol{\theta} - \Omega(\mathbf{p})$$

● argmax: $\Omega(\mathbf{p}) = 0$

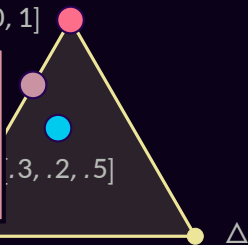
● softmax: $\Omega(\mathbf{p}) = \sum_j p_j \log p_j$

● sparsemax: $\Omega(\mathbf{p}) =$

Unlike lasso,
sparse attention needs $\frac{\partial \Pi_{\Omega}}{\partial \boldsymbol{\theta}}$



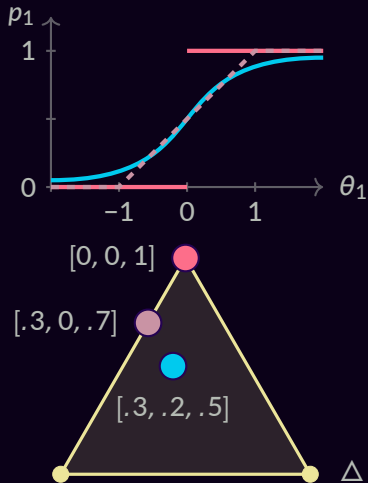
$[0, 0, 1]$



Smoothed Max Operators

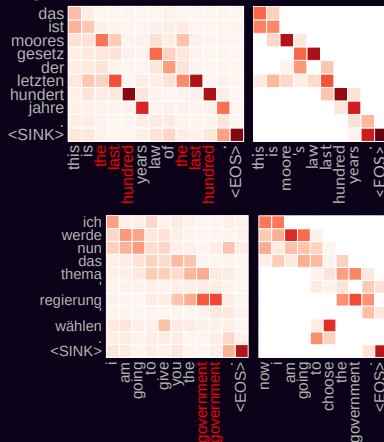
$$\Pi_{\Omega}(\boldsymbol{\theta}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^{\top} \boldsymbol{\theta} - \Omega(\mathbf{p})$$

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



Constrained Attention

e.g., fertility constraints for NMT



$$\begin{aligned} & \arg \max_{\substack{p \in \Delta \\ a \leq p \leq b}} \mathbf{p}^\top \boldsymbol{\theta} - \Omega_1(\mathbf{p}) \\ &= \arg \max_{p \in \Delta} \mathbf{p}^\top \boldsymbol{\theta} - \underbrace{\Omega(\mathbf{p})}_{:= \Omega_1 + \text{Id}_{[a,b]}} \end{aligned}$$

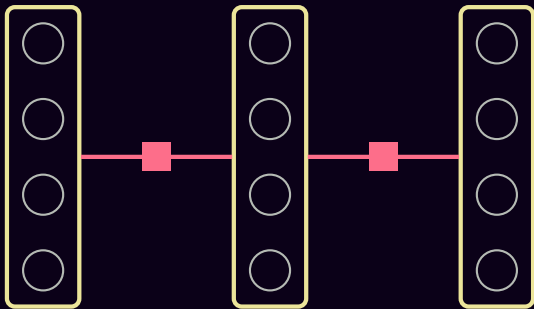
(Kreutzer & Martins, 18)

(Malaviya et al, 18)

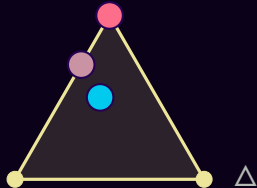
Structured Attention & Graphical Models



Structured Attention & Graphical Models



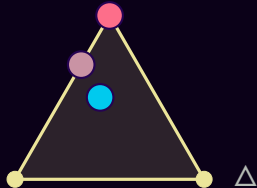
- **argmax** $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \boldsymbol{\theta}$
- **softmax** $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \boldsymbol{\theta} + H(\mathbf{p})$
- **sparsemax** $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \boldsymbol{\theta} - 1/2 \|\mathbf{p}\|^2$



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

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

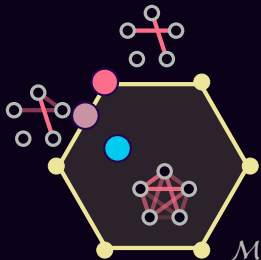
● **sparsemax** $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \boldsymbol{\theta} - 1/2 \|\mathbf{p}\|^2$



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

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

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



Structured Attention for Alignments

NLI

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

input

(P, H)

	A	A	
	gentleman	police	
	overlooking	officer	
	
	situation	closely	

output



entails

contradicts

neutral

(Model: ESIM)

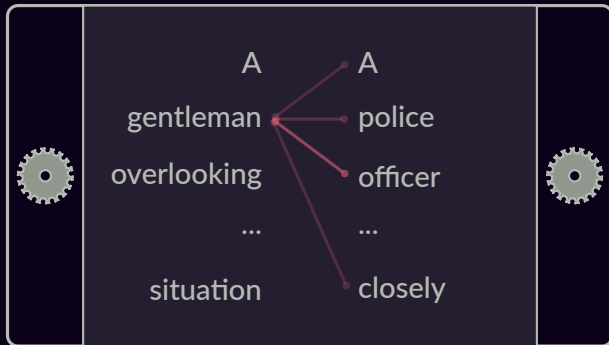
Structured Attention for Alignments

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

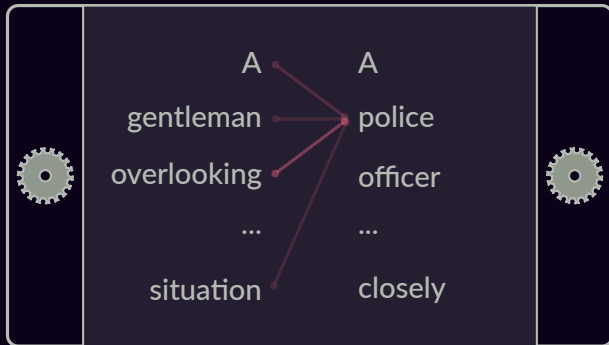
Structured Attention for Alignments

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

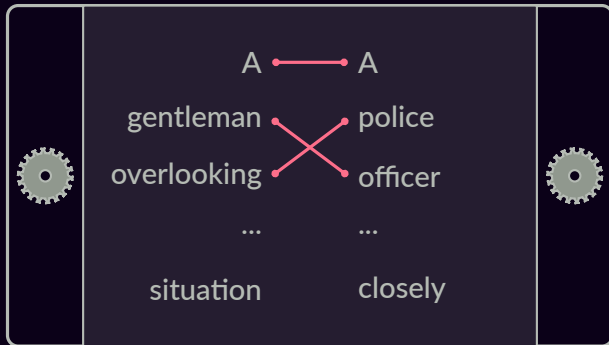
Structured Attention for Alignments

NLI

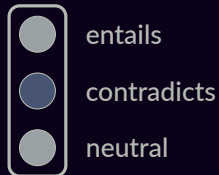
premise: A gentleman overlooking a neighborhood situation.
hypothesis: A police officer watches a situation closely.

input

(P, H)

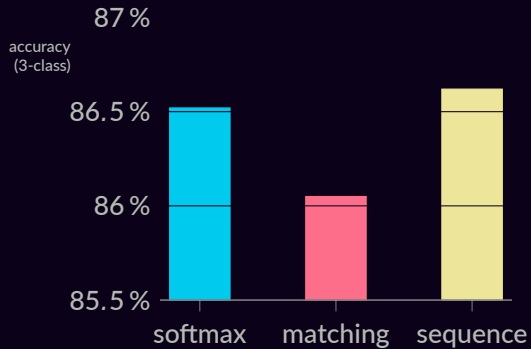


output

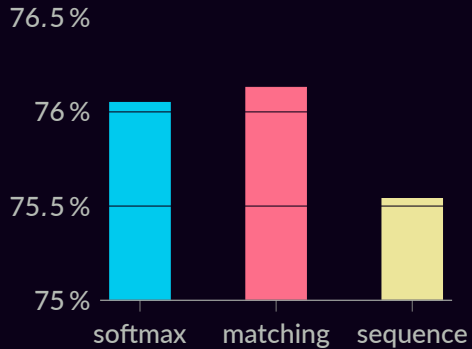


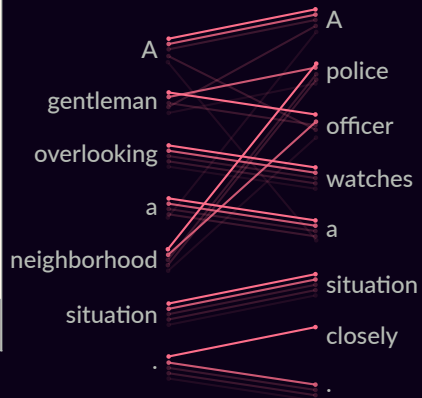
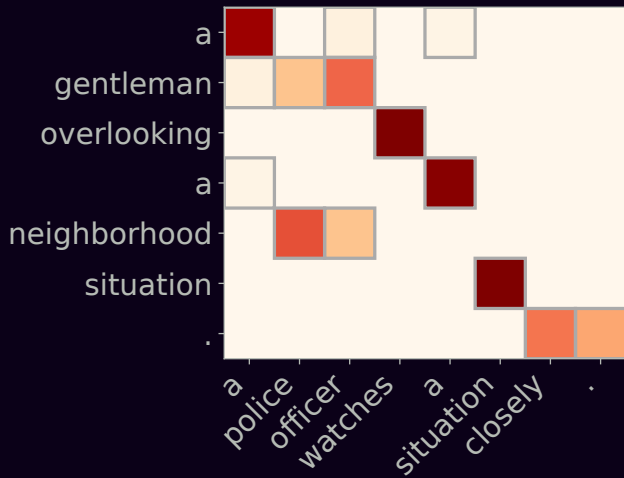
(Proposed model: global matching)

SNLI

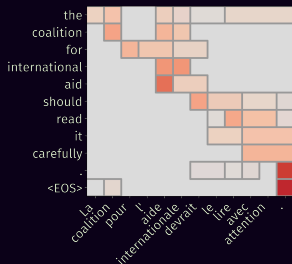


MultiNLI

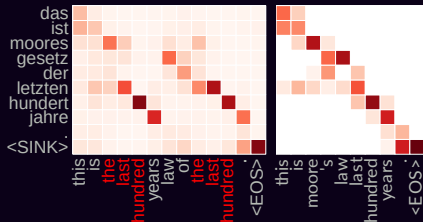




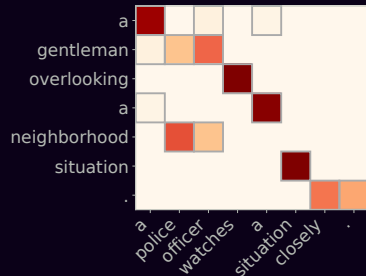
Summary: Neural attention with...



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)

Acknowledgements



This work was supported by the European Research Council (ERC StG DeepSPIN 758969) and by the Fundação para a Ciência e Tecnologia through contract UID/EEA/50008/2013.

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

Extra slides

