

# Neural Attention Mechanisms

Ben Peters

# Acknowledgments

Vlad Niculae made most of these slides.

# Sequence-to-Sequence With Attention

```

words = [21, 79, 14] # indices

embed = Embedding(vocab_sz, dim)
E = embed(words) # (3 × dim)
enc = LSTM(dim, dim)
H = enc(E) # (3 × dim)

dec = LSTM(2 * dim, dim)
initialize k=1, q[0], y[0]

while not done:
    s = H @ q[k - 1] # attn scores
    # s = [-.3, -1.0, 1.8]

    p = softmax(s) # attn proba
    # p = [.10, .05, .85]

    a = p @ H # (1 × dim)

    q[k] = dec(a, y[k - 1], q[k - 1])
    y[k] = W_out @ q[k] + b
    k += 1

```

*United Nations elections*

# Sequence-to-Sequence With Attention

```
words = [21, 79, 14] # indices
```

```
embed = Embedding(vocab_sz, dim)
```

```
E = embed(words) # (3 × dim)
```

```
enc = LSTM(dim, dim)
```

```
H = enc(E) # (3 × dim)
```

```
dec = LSTM(2 * dim, dim)
```

```
initialize k=1, q[0], y[0]
```

```
while not done:
```

```
    s = H @ q[k - 1] # attn scores
```

```
    # s = [-.3, -1.0, 1.8]
```

```
    p = softmax(s) # attn proba
```

```
    # p = [.10, .05, .85]
```

```
    a = p @ H # (1 × dim)
```

```
    q[k] = dec(a, y[k - 1], q[k - 1])
```

```
    y[k] = W_out @ q[k] + b
```

```
    k += 1
```

Encoder


  
 United Nations elections

# Sequence-to-Sequence With Attention

```
words = [21, 79, 14] # indices
```

```
embed = Embedding(vocab_sz, dim)
```

```
E = embed(words) # (3 × dim)
```

```
enc = LSTM(dim, dim)
```

```
H = enc(E) # (3 × dim)
```

```
dec = LSTM(2 * dim, dim)
```

```
initialize k=1, q[0], y[0]
```

```
while not done:
```

```
    s = H @ q[k - 1] # attn scores
```

```
    # s = [-.3, -1.0, 1.8]
```

```
    p = softmax(s) # attn proba
```

```
    # p = [.10, .05, .85]
```

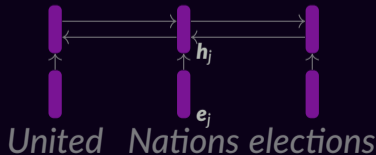
```
    a = p @ H # (1 × dim)
```

```
    q[k] = dec(a, y[k - 1], q[k - 1])
```

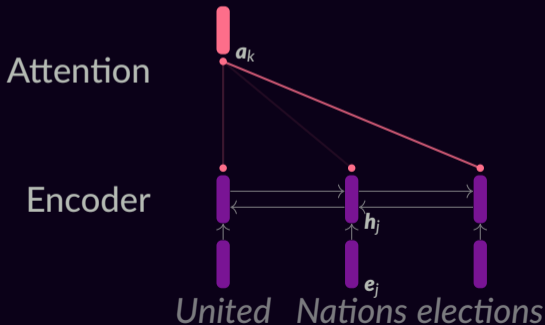
```
    y[k] = W_out @ q[k] + b
```

```
    k += 1
```

Encoder



# Sequence-to-Sequence With Attention



```
words = [21, 79, 14] # indices
```

```
embed = Embedding(vocab_sz, dim)
```

```
E = embed(words) # (3 × dim)
```

```
enc = LSTM(dim, dim)
```

```
H = enc(E) # (3 × dim)
```

```
dec = LSTM(2 * dim, dim)
```

```
initialize k=1, q[0], y[0]
```

```
while not done:
```

```
    s = H @ q[k - 1] # attn scores
```

```
    # s = [-.3, -1.0, 1.8]
```

```
    p = softmax(s) # attn proba
```

```
    # p = [.10, .05, .85]
```

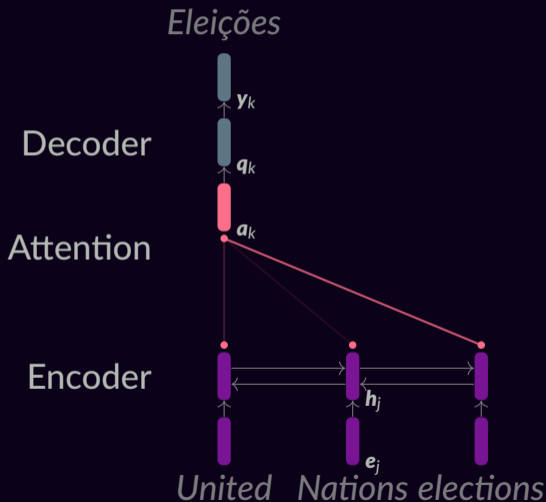
```
    a = p @ H # (1 × dim)
```

```
    q[k] = dec(a, y[k - 1], q[k - 1])
```

```
    y[k] = W_out @ q[k] + b
```

```
    k += 1
```

# Sequence-to-Sequence With Attention



```
words = [21, 79, 14] # indices
```

```
embed = Embedding(vocab_sz, dim)
```

```
E = embed(words) # (3 x dim)
```

```
enc = LSTM(dim, dim)
```

```
H = enc(E) # (3 x dim)
```

```
dec = LSTM(2 * dim, dim)
```

```
initialize k=1, q[0], y[0]
```

```
while not done:
```

```
    s = H @ q[k - 1] # attn scores
```

```
    # s = [-.3, -1.0, 1.8]
```

```
    p = softmax(s) # attn proba
```

```
    # p = [.10, .05, .85]
```

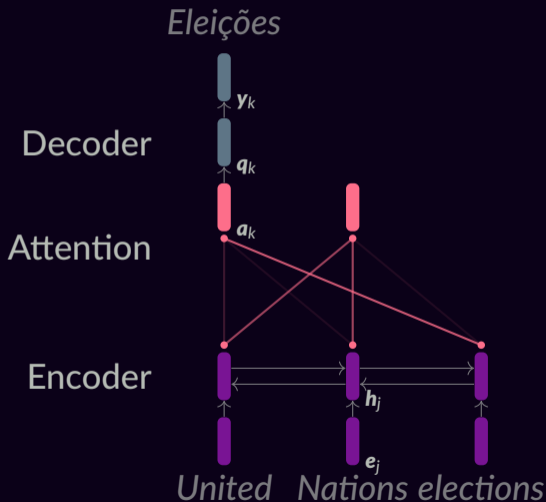
```
    a = p @ H # (1 x dim)
```

```
    q[k] = dec(a, y[k - 1], q[k - 1])
```

```
    y[k] = W_out @ q[k] + b
```

```
    k += 1
```

# Sequence-to-Sequence With Attention



```
words = [21, 79, 14] # indices
```

```
embed = Embedding(vocab_sz, dim)
```

```
E = embed(words) # (3 x dim)
```

```
enc = LSTM(dim, dim)
```

```
H = enc(E) # (3 x dim)
```

```
dec = LSTM(2 * dim, dim)
```

```
initialize k=1, q[0], y[0]
```

```
while not done:
```

```
    s = H @ q[k - 1] # attn scores
```

```
    # s = [-.3, -1.0, 1.8]
```

```
    p = softmax(s) # attn proba
```

```
    # p = [.10, .05, .85]
```

```
    a = p @ H # (1 x dim)
```

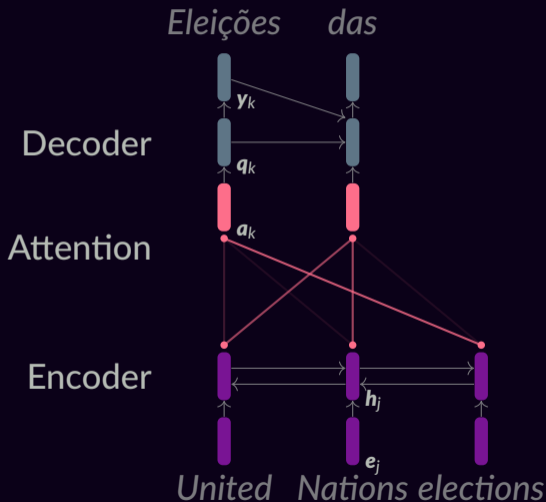
```
    q[k] = dec(a, y[k - 1], q[k - 1])
```

```
    y[k] = W_out @ q[k] + b
```

```
    k += 1
```



# Sequence-to-Sequence With Attention



```
words = [21, 79, 14] # indices
```

```
embed = Embedding(vocab_sz, dim)
```

```
E = embed(words) # (3 x dim)
```

```
enc = LSTM(dim, dim)
```

```
H = enc(E) # (3 x dim)
```

```
dec = LSTM(2 * dim, dim)
```

```
initialize k=1, q[0], y[0]
```

```
while not done:
```

```
    s = H @ q[k - 1] # attn scores
```

```
    # s = [-.3, -1.0, 1.8]
```

```
    p = softmax(s) # attn proba
```

```
    # p = [.10, .05, .85]
```

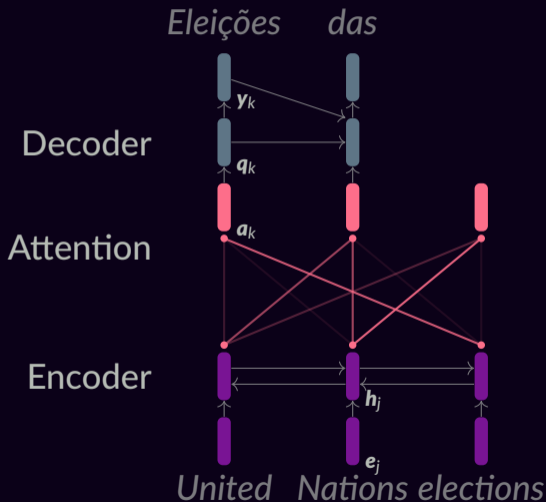
```
    a = p @ H # (1 x dim)
```

```
    q[k] = dec(a, y[k - 1], q[k - 1])
```

```
    y[k] = W_out @ q[k] + b
```

```
    k += 1
```

# Sequence-to-Sequence With Attention



```
words = [21, 79, 14] # indices
```

```
embed = Embedding(vocab_sz, dim)
```

```
E = embed(words) # (3 x dim)
```

```
enc = LSTM(dim, dim)
```

```
H = enc(E) # (3 x dim)
```

```
dec = LSTM(2 * dim, dim)
```

```
initialize k=1, q[0], y[0]
```

```
while not done:
```

```
    s = H @ q[k - 1] # attn scores
```

```
    # s = [-.3, -1.0, 1.8]
```

```
    p = softmax(s) # attn proba
```

```
    # p = [.10, .05, .85]
```

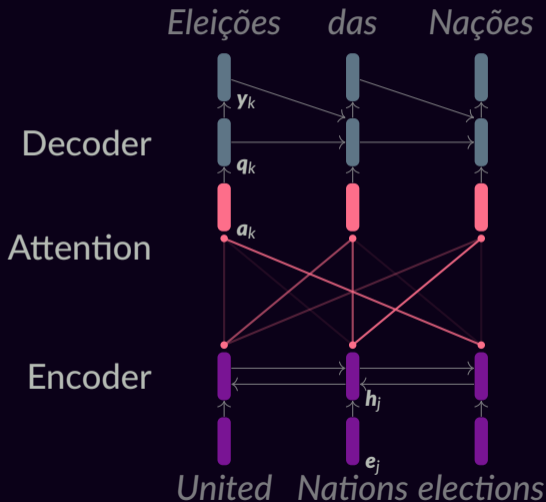
```
    a = p @ H # (1 x dim)
```

```
    q[k] = dec(a, y[k - 1], q[k - 1])
```

```
    y[k] = W_out @ q[k] + b
```

```
    k += 1
```

# Sequence-to-Sequence With Attention



```
words = [21, 79, 14] # indices
```

```
embed = Embedding(vocab_sz, dim)
```

```
E = embed(words) # (3 × dim)
```

```
enc = LSTM(dim, dim)
```

```
H = enc(E) # (3 × dim)
```

```
dec = LSTM(2 * dim, dim)
```

```
initialize k=1, q[0], y[0]
```

```
while not done:
```

```
    s = H @ q[k - 1] # attn scores
```

```
    # s = [-.3, -1.0, 1.8]
```

```
    p = softmax(s) # attn proba
```

```
    # p = [.10, .05, .85]
```

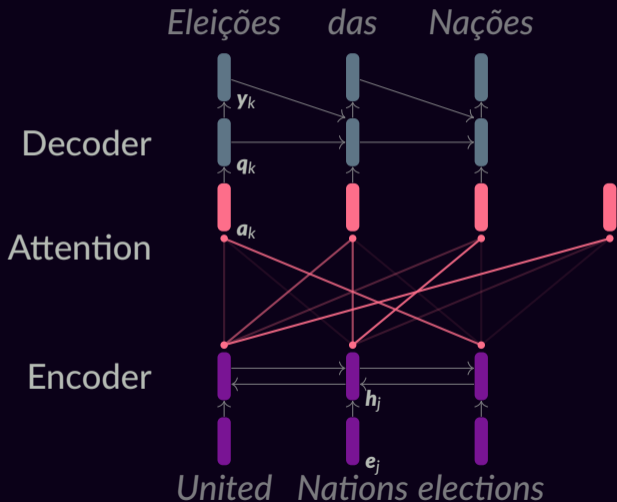
```
    a = p @ H # (1 × dim)
```

```
    q[k] = dec(a, y[k - 1], q[k - 1])
```

```
    y[k] = W_out @ q[k] + b
```

```
    k += 1
```

# Sequence-to-Sequence With Attention



```
words = [21, 79, 14] # indices
```

```
embed = Embedding(vocab_sz, dim)
```

```
E = embed(words) # (3 x dim)
```

```
enc = LSTM(dim, dim)
```

```
H = enc(E) # (3 x dim)
```

```
dec = LSTM(2 * dim, dim)
```

```
initialize k=1, q[0], y[0]
```

```
while not done:
```

```
    s = H @ q[k - 1] # attn scores
```

```
    # s = [-.3, -1.0, 1.8]
```

```
    p = softmax(s) # attn proba
```

```
    # p = [.10, .05, .85]
```

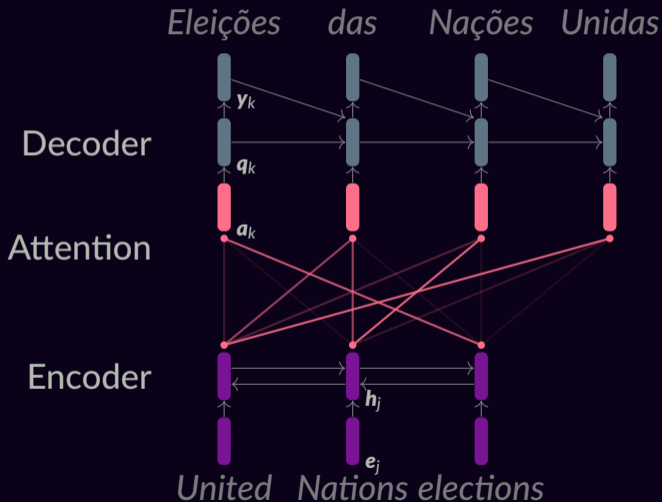
```
    a = p @ H # (1 x dim)
```

```
    q[k] = dec(a, y[k - 1], q[k - 1])
```

```
    y[k] = W_out @ q[k] + b
```

```
    k += 1
```

# Sequence-to-Sequence With Attention



```
words = [21, 79, 14] # indices
```

```
embed = Embedding(vocab_sz, dim)
```

```
E = embed(words) # (3 x dim)
```

```
enc = LSTM(dim, dim)
```

```
H = enc(E) # (3 x dim)
```

```
dec = LSTM(2 * dim, dim)
```

```
initialize k=1, q[0], y[0]
```

```
while not done:
```

```
    s = H @ q[k - 1] # attn scores
```

```
    # s = [-.3, -1.0, 1.8]
```

```
    p = softmax(s) # attn proba
```

```
    # p = [.10, .05, .85]
```

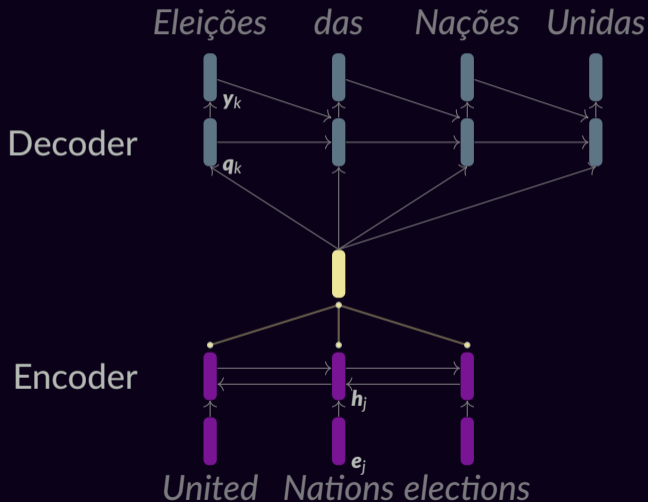
```
    a = p @ H # (1 x dim)
```

```
    q[k] = dec(a, y[k - 1], q[k - 1])
```

```
    y[k] = W_out @ q[k] + b
```

```
    k += 1
```

# Sequence-to-Sequence With Attention



```
words = [21, 79, 14] # indices
```

```
embed = Embedding(vocab_sz, dim)
```

```
E = embed(words) # (3 x dim)
```

```
enc = LSTM(dim, dim)
```

```
H = enc(E) # (3 x dim)
```

```
dec = LSTM(2 * dim, dim)
```

```
initialize k=1, q[0], y[0]
```

```
while not done:
```

```
    s = H @ q[k - 1] # attn scores
```

```
    # s = [-.3, -1.0, 1.8]
```

```
    p = softmax(s) # attn proba
```

```
    # p = [.10, .05, .85]
```

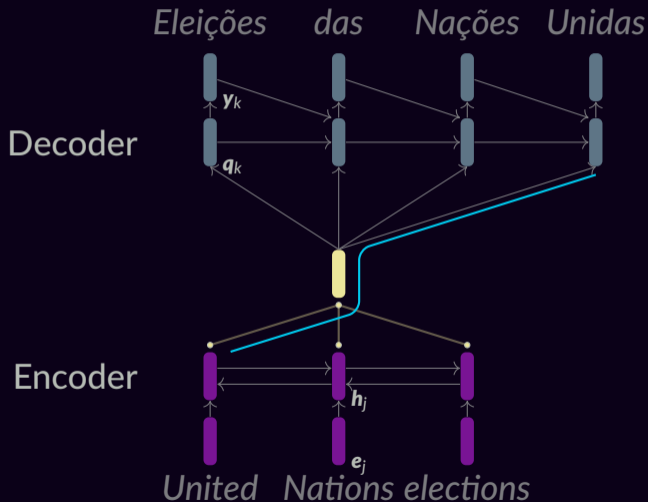
```
    a = p @ H # (1 x dim)
```

```
    q[k] = dec(a, y[k - 1], q[k - 1])
```

```
    y[k] = W_out @ q[k] + b
```

```
    k += 1
```

# Sequence-to-Sequence With Attention



```
words = [21, 79, 14] # indices
```

```
embed = Embedding(vocab_sz, dim)
```

```
E = embed(words) # (3 x dim)
```

```
enc = LSTM(dim, dim)
```

```
H = enc(E) # (3 x dim)
```

```
dec = LSTM(2 * dim, dim)
```

```
initialize k=1, q[0], y[0]
```

```
while not done:
```

```
    s = H @ q[k - 1] # attn scores
```

```
    # s = [-.3, -1.0, 1.8]
```

```
    p = softmax(s) # attn proba
```

```
    # p = [.10, .05, .85]
```

```
    a = p @ H # (1 x dim)
```

```
    q[k] = dec(a, y[k - 1], q[k - 1])
```

```
    y[k] = W_out @ q[k] + b
```

```
    k += 1
```

# In a Nutshell



# In a Nutshell

- You have one query vector (think tgt word)

# In a Nutshell

- You have one query vector (think tgt word)
- You have a lot of key vectors (one for each src word)

# In a Nutshell

- You have one query vector (think tgt word)
- You have a lot of key vectors (one for each src word)
- First: come up with a *score* between the query and each key

# In a Nutshell

- You have one query vector (think tgt word)
- You have a lot of key vectors (one for each src word)
- First: come up with a *score* between the query and each key
- Use the scores to decide where to attend

# Computing the scores

$$s_j = \sigma(\mathbf{h}_j, \mathbf{q})$$

name	$\sigma(\mathbf{h}, \mathbf{q})$	
additive	$\mathbf{v}^\top \tanh(\mathbf{W}_1 \mathbf{h} + \mathbf{W}_2 \mathbf{q})$	(Bahdanau et al., 2015)
dot-product	$\mathbf{h}^\top \mathbf{q}$	(Luong et al., 2015)
bilinear	$\mathbf{h}^\top \mathbf{W} \mathbf{q}$	(Luong et al., 2015)
scaled	$(1/\sqrt{d}) \mathbf{h}^\top \mathbf{W} \mathbf{q}$	(Vaswani et al., 2017)

```
# attention scores:  
s = H @ W_attn @ state  
# s = [-.3, -1.0, 1.8]  
  
p = softmax(s)  
# p = [.10, .05, .85]
```

\*record scratch\*

\*freeze frame\*

1. How to select an item  
from a set?

# How to select an item from a set?

*United*

*Nations*

*Elections*



# How to select an item from a set?

$c_1$

$c_2$

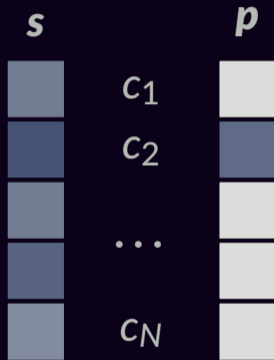
...

$c_N$

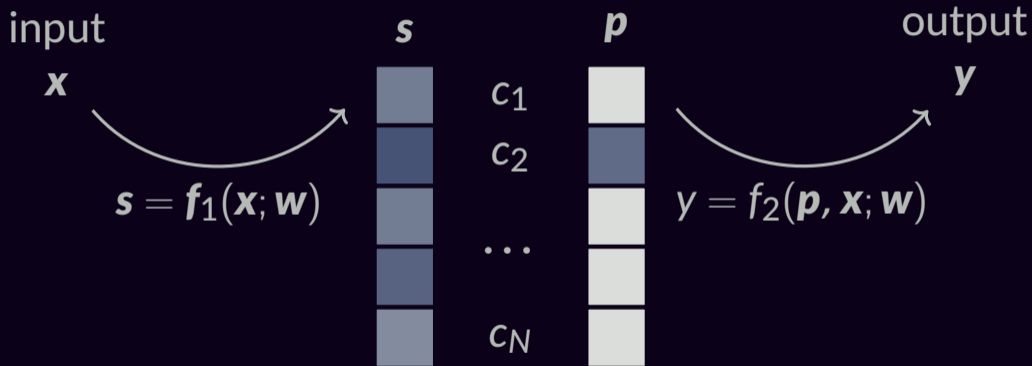
# How to select an item from a set?



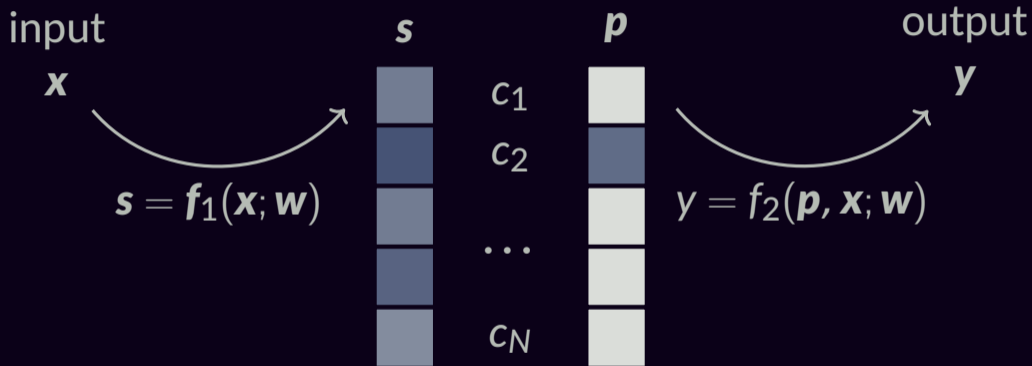
# How to select an item from a set?



# How to select an item from a set?

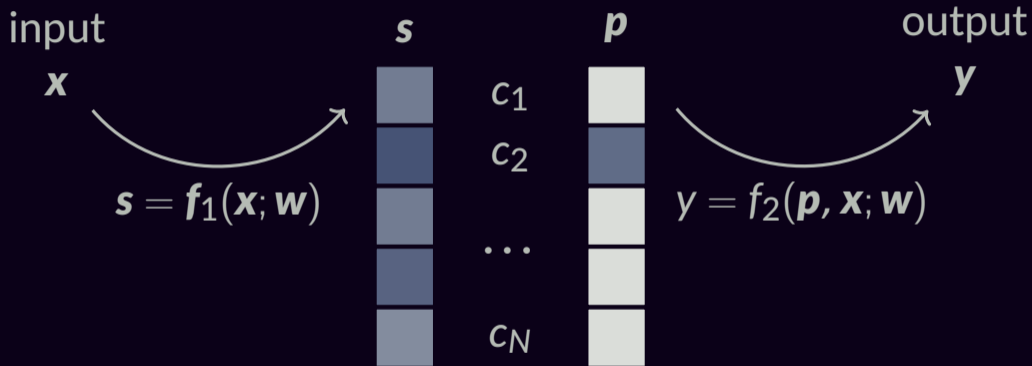


# How to select an item from a set?



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

# How to select an item from a set?



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

or, essentially,

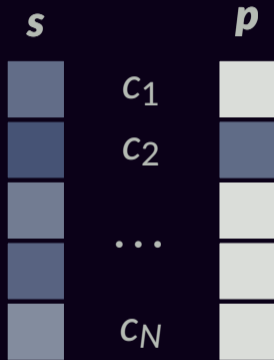
$$\frac{\partial \mathbf{p}}{\partial \mathbf{s}} = ?$$

# Winner Takes It All



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

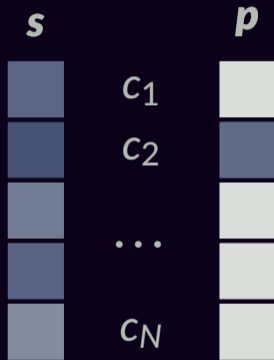
# Winner Takes It All



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

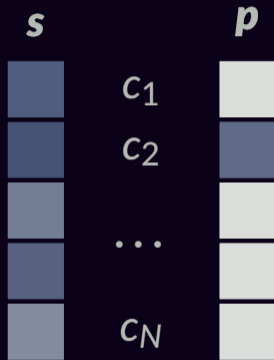


# Winner Takes It All



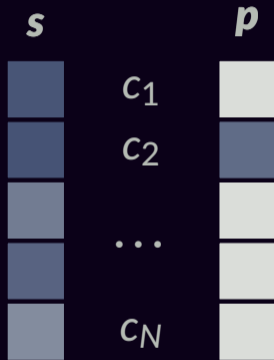
$$\frac{\partial p}{\partial s} = ?$$

# Winner Takes It All



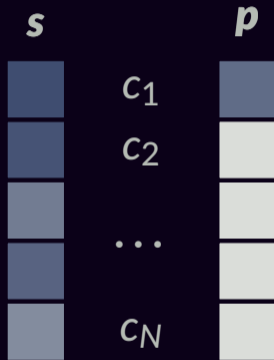
$$\frac{\partial p}{\partial s} = ?$$

# Winner Takes It All



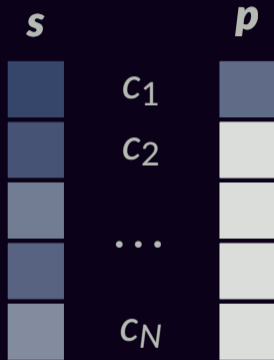
$$\frac{\partial p}{\partial s} = ?$$

# Winner Takes It All



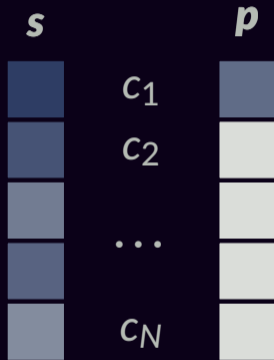
$$\frac{\partial p}{\partial s} = ?$$

# Winner Takes It All



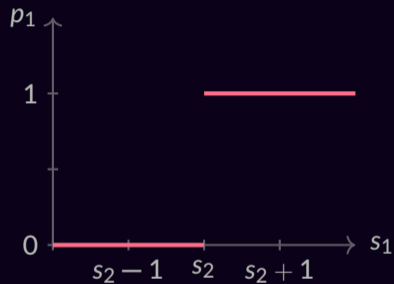
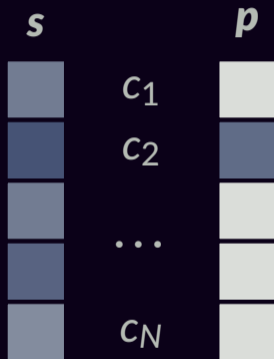
$$\frac{\partial p}{\partial s} = ?$$

# Winner Takes It All



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

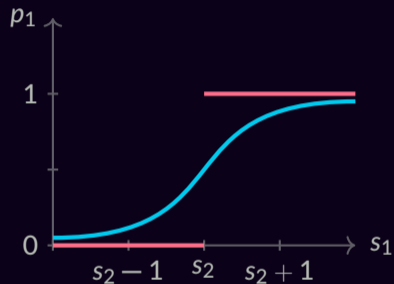
# Argmax



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

# Argmax vs. Softmax

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



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



# An Inconvenient Truth

We need gradients.  
So you can't just take the highest-scoring item.

But...

That **does not** mean that softmax is the only choice.

# But...

That **does not** mean that softmax is the only choice.

- Softmax attention is **dense**

# But...

That **does not** mean that softmax is the only choice.

- Softmax attention is **dense**
- Argmax is completely **sparse**

# But...

That **does not** mean that softmax is the only choice.

- Softmax attention is **dense**
- Argmax is completely **sparse**
- Maybe there's some kind of middle ground?

# Sparsemax

$$\begin{aligned}\text{sparsemax}(\mathbf{s}) &= \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{s} - 1/2 \|\mathbf{p}\|_2^2 \\ &= \arg \min_{\mathbf{p} \in \Delta} \|\mathbf{p} - \mathbf{s}\|_2^2\end{aligned}$$

**Computation:**

$$\mathbf{p}^* = [\mathbf{s} - \tau \mathbf{1}]_+$$

$$s_i > s_j \Rightarrow p_i \geq p_j$$

$O(d)$  via partial sort

(Held et al., 1974; Brucker, 1984; Condat, 2016)

**Backward pass:**

$$\mathbf{J}_{\text{sparsemax}} = \text{diag}(\mathbf{s}) - \frac{1}{|\mathcal{S}|} \mathbf{s} \mathbf{s}^\top$$

$$\text{where } \mathcal{S} = \{j : p_j^* > 0\},$$

$$s_j = \mathbb{1}[j \in \mathcal{S}]$$

(Martins and Astudillo, 2016)

# $\alpha$ -entmax

$$\pi_{H_{\alpha}^t}(\mathbf{s}) := \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s} + H_{\alpha}^t(\mathbf{p})$$

Solution has the form:

$$\pi_{H_{\alpha}^t}(\mathbf{s}) = [(\alpha - 1)\mathbf{s} - \tau\mathbf{1}]_+^{1/\alpha - 1}$$

## Algorithms:

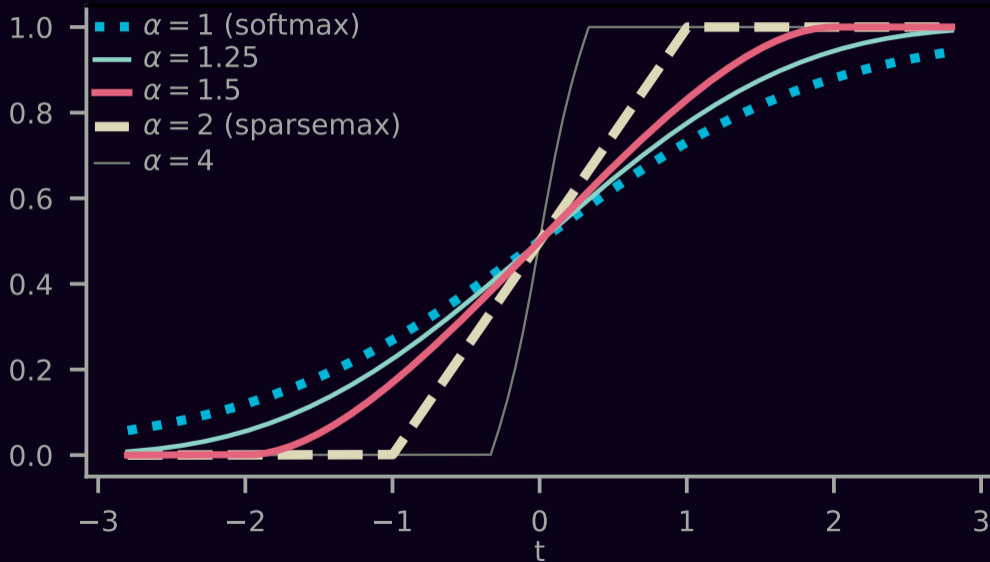
### *bisection*

- approximate; bracket  $\tau \in [\tau_{\text{lo}}, \tau_{\text{hi}}]$
- gain 1 bit per  $O(d)$  iteration
- `float32` has 23 mantissa bits

### *sort-based*

- exact algorithm,  $O(d \log d)$
- available only for  $\alpha \in \{1.5, 2\}$
- huge speed-up from partial sorting when expecting sparse solutions

$$\pi_\alpha([t, 0])_1$$





# Structural Biases?

Imagine you are an  $en \rightarrow pt$  model.

# Structural Biases?

Imagine you are an  $en \rightarrow pt$  model.

I am going to the store  $\rightarrow$

# Structural Biases?

Imagine you are an en $\rightarrow$ pt model.

I am going to the store  $\rightarrow$  Vou à loja

# Structural Biases?

Imagine you are an en→pt model.

I am going to the store → Vous à la boutique

When you generate “Vous”, where do you attend?

# Structural Biases?

Imagine you are an en→pt model.

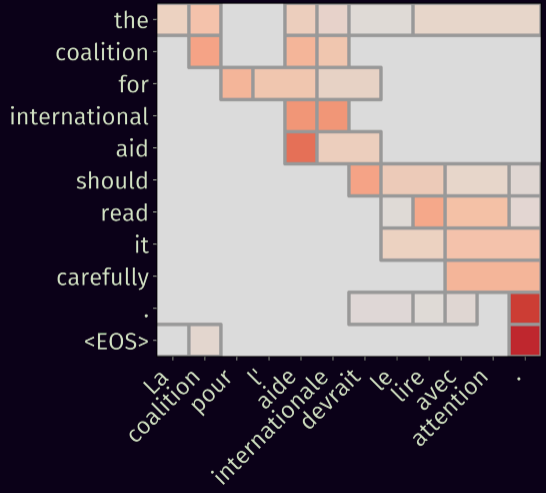
I am going to the store → Vou à loja

When you generate “Vou”, where do you attend?

I am going to the store → Vou à loja

# Enter fusedmax

Basic idea: penalize weight differences between adjacent positions.



fusedmax

# Fusedmax

$$\text{fusedmax}(\mathbf{s}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s} - \frac{1}{2} \|\mathbf{p}\|_2^2 - \sum_{2 \leq j \leq d} |p_j - p_{j-1}|$$

$$= \arg \min_{\mathbf{p} \in \Delta} \|\mathbf{p} - \mathbf{s}\|_2^2 + \sum_{2 \leq j \leq d} |p_j - p_{j-1}|$$

$$\text{prox}_{\text{fused}}(\mathbf{s}) = \arg \min_{\mathbf{p} \in \mathbb{R}^d} \|\mathbf{p} - \mathbf{s}\|_2^2 + \sum_{2 \leq j \leq d} |p_j - p_{j-1}|$$

**Proposition:**  $\text{fusedmax}(\mathbf{s}) = \text{sparsemax}(\text{prox}_{\text{fused}}(\mathbf{s}))$

(Niculae and Blondel, 2017)



# Constrained Attention

Another structural idea: limit how much attention a word gets.

# Constrained Attention

$$\text{softmax}(\mathbf{z}) = \arg \min_{\mathbf{p} \in \Delta^D} \text{KL}(\mathbf{p} \parallel \text{softmax}(\mathbf{z}))$$

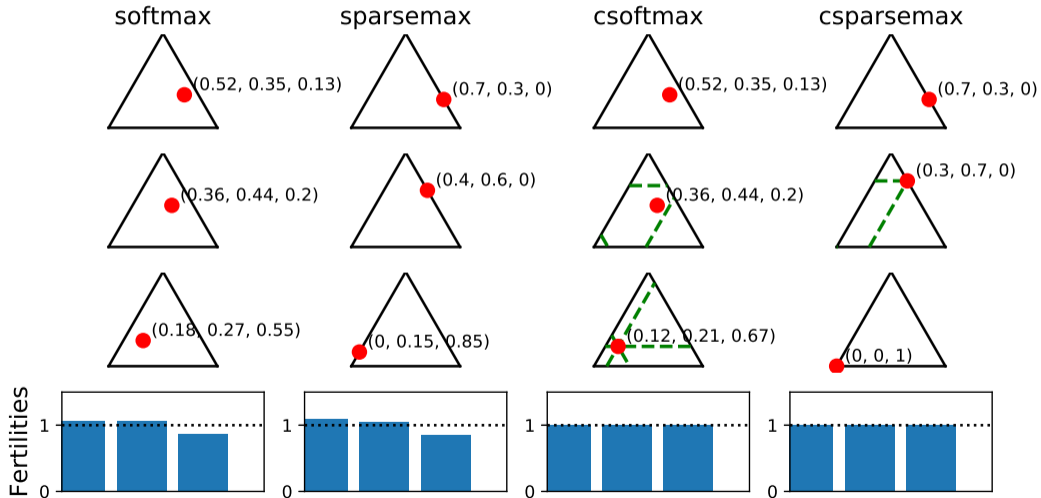
# Constrained Attention

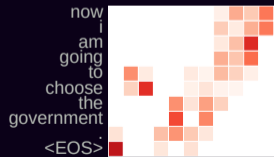
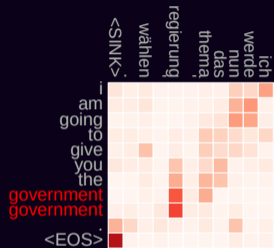
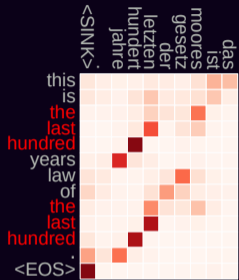
$$\text{softmax}(\mathbf{z}) = \arg \min_{\mathbf{p} \in \Delta^D} \text{KL}(\mathbf{p} \parallel \text{softmax}(\mathbf{z}))$$

$$\text{csoftmax}(\mathbf{z}, \mathbf{u}) = \arg \min_{\mathbf{p} \in \Delta^D \text{ s.t. } \mathbf{p} \leq \mathbf{u}} \text{KL}(\mathbf{p} \parallel \text{softmax}(\mathbf{z}))$$

$\mathbf{u}$ : max probability for each index

# Example: Source Sentence with Three Words





## 2. Transformers

# Beyond seq2seq

The spirit of *attention mechanisms* reaches far:

- Key-Value Attention
- Multi-head Attention
- Self-Attention

What do you get when you put these together?



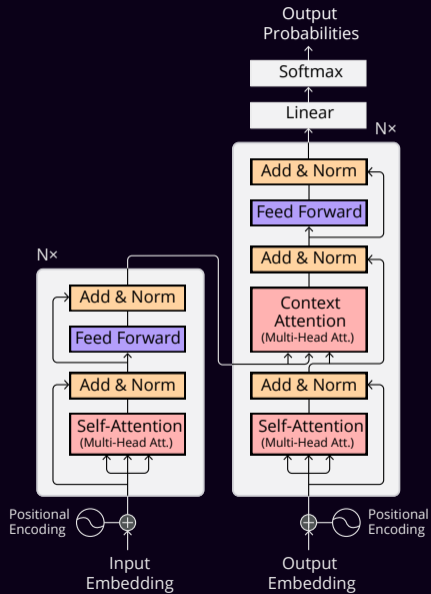


A line of toys?

A line of toys?

A film franchise starring Shia LaBeouf?

A line of toys?  
A film franchise starring Shia LaBeouf?  
*Wrong.*



# Key-Value Attention

idea: the objects we average (*values*)  
and the objects used to compute scores (*keys*)  
don't need to be identical!

$$s_j = \mathbf{h}_j^\top \mathbf{q}$$

$$\mathbf{u} = \text{softmax}(\mathbf{s})^\top \mathbf{H}$$

$$s_j = \mathbf{k}_j^\top \mathbf{q}$$

$$\mathbf{u} = \text{softmax}(\mathbf{s})^\top \mathbf{V}$$

# Multi-head Attention

idea: compute  $k$  different attention averages,  
& concatenate the outputs.

$$s_j = \mathbf{k}_j^\top \mathbf{q}$$

$$\mathbf{u} = \text{softmax}(\mathbf{s})^\top \mathbf{V}$$

$$\mathbf{u} = \text{softmax}(\mathbf{K} @ \mathbf{q}) @ \mathbf{V}$$

$$s_j^{(i)} = (\mathbf{W}_k^{(i)} \mathbf{k}_j)^\top (\mathbf{W}_q^{(i)} \mathbf{q})$$

$$\mathbf{u}^{(i)} = \text{softmax}(\mathbf{s}^{(i)})^\top (\mathbf{V} \mathbf{W}_v^{(i)})$$

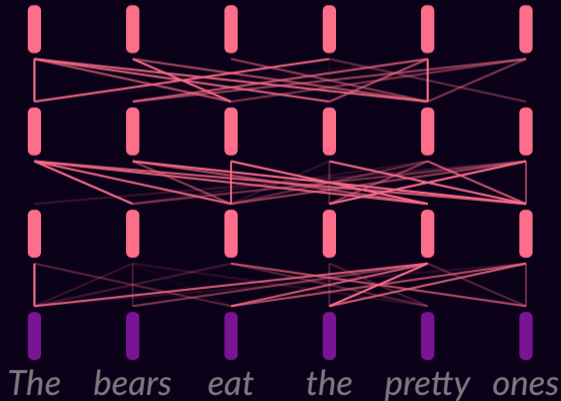
$$\mathbf{u} = [\mathbf{u}^{(1)}; \dots; \mathbf{u}^{(k)}]$$

```
for i in range(num_heads):
    Ki = K @ Wk[i].t()
    Vi = V @ Wv[i].t()
    qi = q @ Wq[i].t()
    ui = softmax(Ki @ qi) @ Vi
u = concat(ui)
```

# Self-attention

Attention as an *encoder layer*

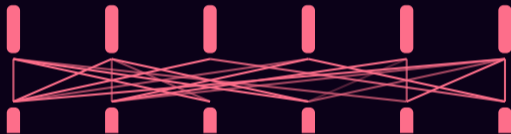
...



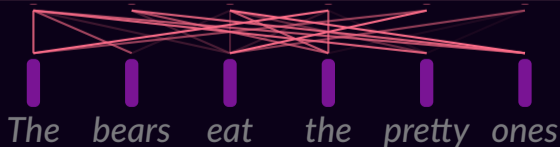
# Self-attention

Attention as an *encoder layer*

...



**Transformer** (Vaswani et al., 2017): very deep self-attention replacing LSTMs in encoder & decoder





# Could you say that again?

Each self-attention head uses key-value attention:

- Take a sequence of embeddings  $X = \mathbf{x}_1 \dots \mathbf{x}_n$ ,  $\mathbf{x}_i \in \mathbb{R}^D$

# Could you say that again?

Each self-attention head uses key-value attention:

- Take a sequence of embeddings  $X = \mathbf{x}_1 \dots \mathbf{x}_n$ ,  $\mathbf{x}_i \in \mathbb{R}^D$
- queries:  $Q = XW_q$ ,  $\mathbf{q}_i \in \mathbb{R}^{D_q}$

# Could you say that again?

Each self-attention head uses key-value attention:

- Take a sequence of embeddings  $X = \mathbf{x}_1 \dots \mathbf{x}_n$ ,  $\mathbf{x}_i \in \mathbb{R}^D$
- queries:  $Q = XW_q$ ,  $\mathbf{q}_i \in \mathbb{R}^{D_q}$
- keys:  $K = XW_k$ ,  $\mathbf{k}_i \in \mathbb{R}^{D_q}$

# Could you say that again?

Each self-attention head uses key-value attention:

- Take a sequence of embeddings  $X = \mathbf{x}_1 \dots \mathbf{x}_n$ ,  $\mathbf{x}_i \in \mathbb{R}^D$
- queries:  $Q = XW_q$ ,  $\mathbf{q}_i \in \mathbb{R}^{D_q}$
- keys:  $K = XW_k$ ,  $\mathbf{k}_i \in \mathbb{R}^{D_k}$
- values:  $V = XW_v$ ,  $\mathbf{v}_i \in \mathbb{R}^{D_v}$

# Could you say that again?

Each self-attention head uses key-value attention:

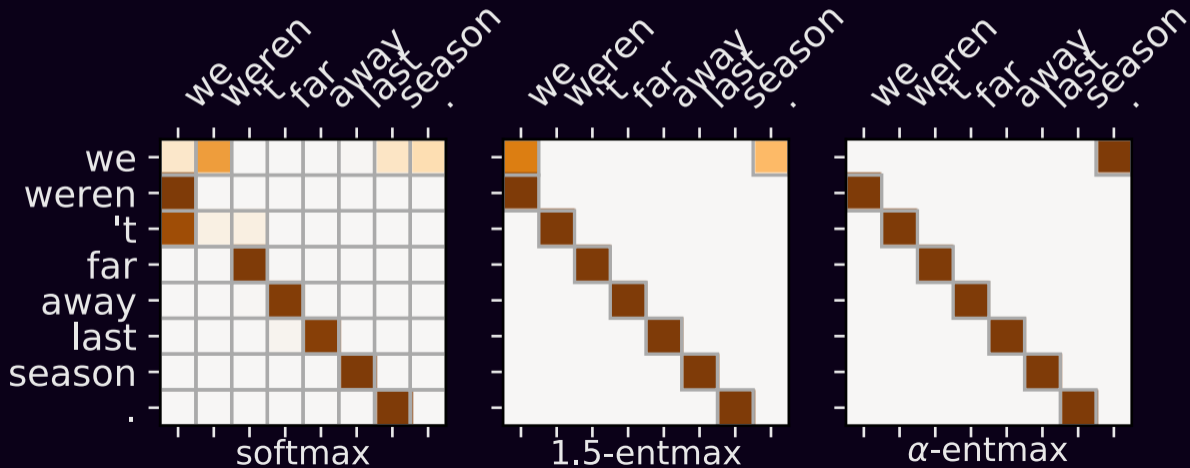
- Take a sequence of embeddings  $X = \mathbf{x}_1 \dots \mathbf{x}_n$ ,  $\mathbf{x}_i \in \mathbb{R}^D$
- queries:  $Q = XW_q$ ,  $\mathbf{q}_i \in \mathbb{R}^{D_q}$
- keys:  $K = XW_k$ ,  $\mathbf{k}_i \in \mathbb{R}^{D_k}$
- values:  $V = XW_v$ ,  $\mathbf{v}_i \in \mathbb{R}^{D_v}$

Important: each head has different  $W_q, W_k, W_v$

$$\text{softmax}\left(\frac{QK^T}{\sqrt{D}}\right)V$$

So what do heads learn?

# Previous Position









# Replacing Recurrence

Self-attention is the *only* place where positions interact.

What do we gain over RNN-based models?

What do we lose?

# References I

- Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio (2015). “Neural machine translation by jointly learning to align and translate”. In: *Proc. of ICLR*.
- Brucker, Peter (1984). “An  $O(n)$  algorithm for quadratic knapsack problems”. In: *Operations Research Letters* 3.3, pp. 163–166.
- Cheng, Jianpeng, Li Dong, and Mirella Lapata (2016). “Long Short-Term Memory-Networks for Machine Reading”. In: *Proc. of EMNLP*.
- Condat, Laurent (2016). “Fast projection onto the simplex and the  $\ell_1$  ball”. In: *Mathematical Programming* 158.1-2, pp. 575–585.
- Correia, Gonçalo M, Vlad Niculae, and André FT Martins (2019). “Adaptively Sparse Transformers”. In: *Proc. EMNLP-IJCNLP (to appear)*.
- Graves, Alex, Greg Wayne, and Ivo Danihelka (2014). “Neural Turing Machines”. In: *arXiv preprint arXiv:1410.5401*.
- Held, Michael, Philip Wolfe, and Harlan P Crowder (1974). “Validation of subgradient optimization”. In: *Mathematical Programming* 6.1, pp. 62–88.
- Luong, Minh-Thang, Hieu Pham, and Christopher D Manning (2015). “Effective approaches to attention-based neural machine translation”. In: *Proc. of EMNLP*.
- Malaviya, Chaitanya, Pedro Ferreira, and André F. T. Martins (2018). “Sparse and constrained attention for neural machine translation”. In: *Proc. of ACL*.
- Martins, André FT and Ramón Fernandez Astudillo (2016). “From softmax to sparsemax: A sparse model of attention and multi-label classification”. In: *Proc. of ICML*.

# References II

- Martins, André FT and Julia Kreutzer (2017). “Learning What’s Easy: Fully Differentiable Neural Easy-First Taggers”. In: *Proc. of EMNLP*, pp. 349–362.
- Niculescu, Vlad and Mathieu Blondel (2017). “A regularized framework for sparse and structured neural attention”. In: *Proc. of NeurIPS*.
- Peters, Ben, Vlad Niculescu, and André FT Martins (2019). “Sparse Sequence-to-Sequence Models”. In: *Proc. ACL*.
- Sukhbaatar, Sainbayar, Jason Weston, and Rob Fergus (2015). “End-to-end memory networks”. In: *Proc. of NeurIPS*.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin (2017). “Attention Is All You Need”. In: *Proc. of NeurIPS*.
- Voita, Elena, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov (July 2019). “Analyzing Multi-Head Self-Attention: Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned”. In: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, pp. 5797–5808. doi: 10.18653/v1/P19-1580. url: <https://www.aclweb.org/anthology/P19-1580>.