# Understanding Attention Training via Output Relevance

Charlie Snell, Ruiqi Zhong, Jacob Steinhardt, Dan Klein

# Outline

- Prior work on text classification:

  - Standardly trained attention $\approx$ Explanations .

  - Can train models to attend to irrelevant words, same accuracy.

?

- Our work studies how the attention evolves at training time, for text classification and machine translation.
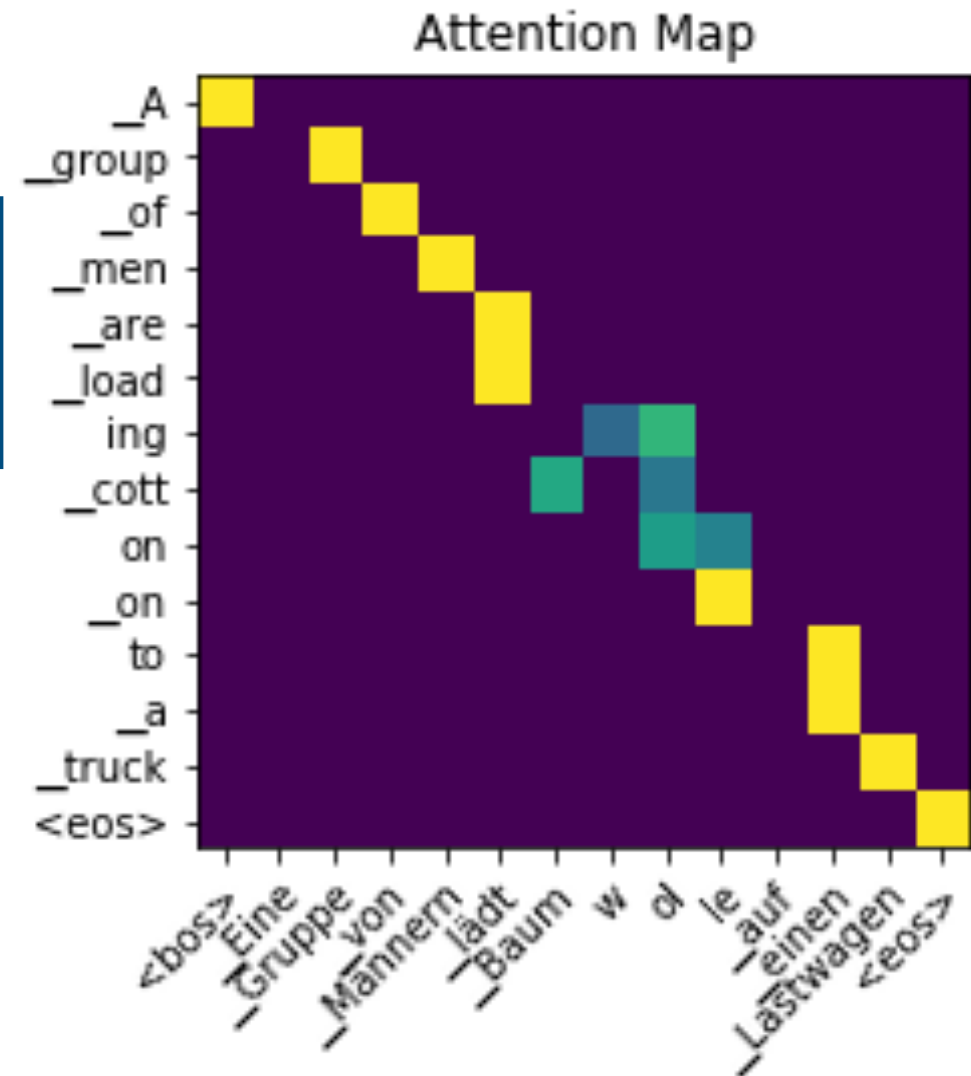
# Attention Distribution over Inputs
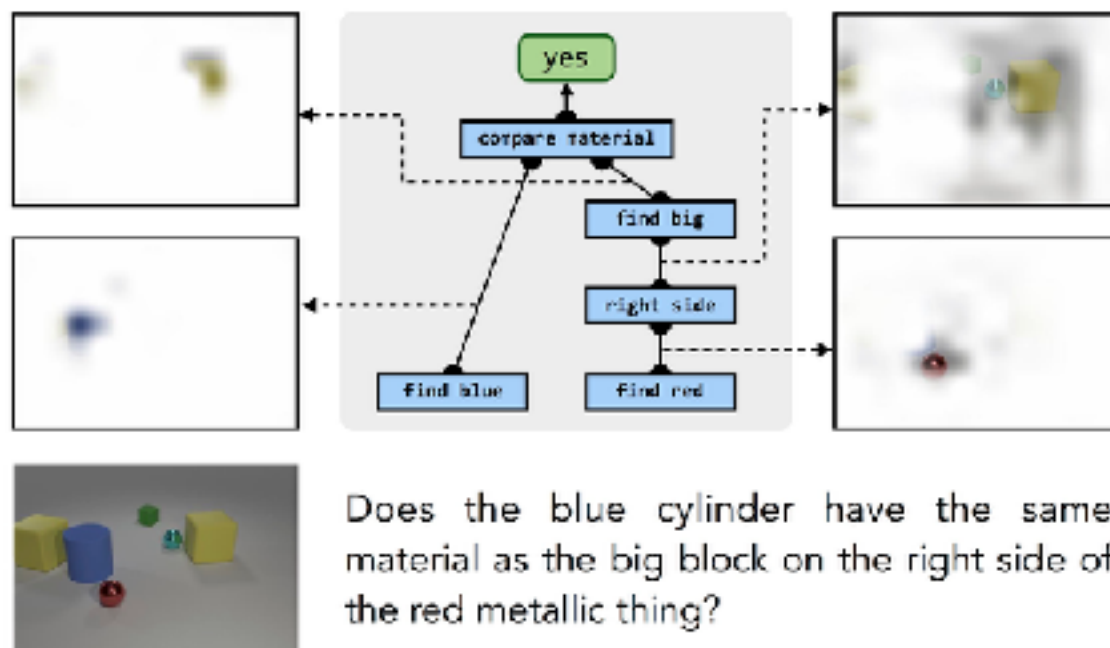
Sentiment Classification

*Negative*

*after 15 minutes watching the movie i was asking myself what to do leave the theater sleep or try to keep watching the movie to see if there was anything worth i finally watched the movie what a waste of time maybe i am not a 5 years old kid anymore*

**Seems to provide explanations !!**

German -> English Translation

Does the blue cylinder have the same material as the big block on the right side of the red metallic thing?
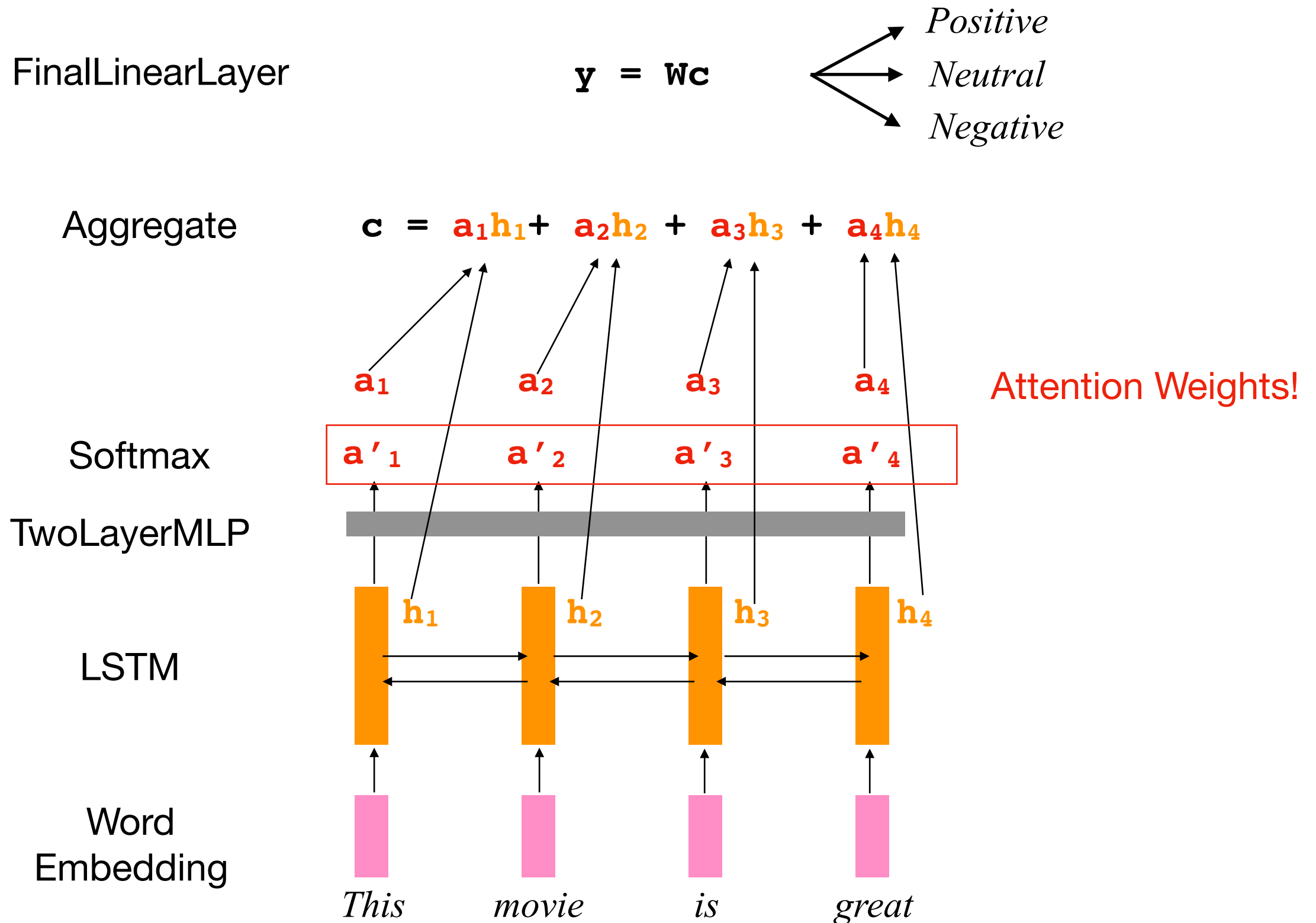
Visual Question Answering

# Attention ≈ Explanation?

- Attention is not Explanation (Jain & Wallace 2019)

- Is Attention Interpretable? (Serrano & Smith 2019)

- Attention is not not Explanation (Wiegreffe & Pinter 2019)

- Attention Interpretability Across Tasks (Shikhar et al. 2019)

- Understanding Attention Training via Output Relevance (Snell et al. 2020)

**A series of work trying to understand attention mechanism.**
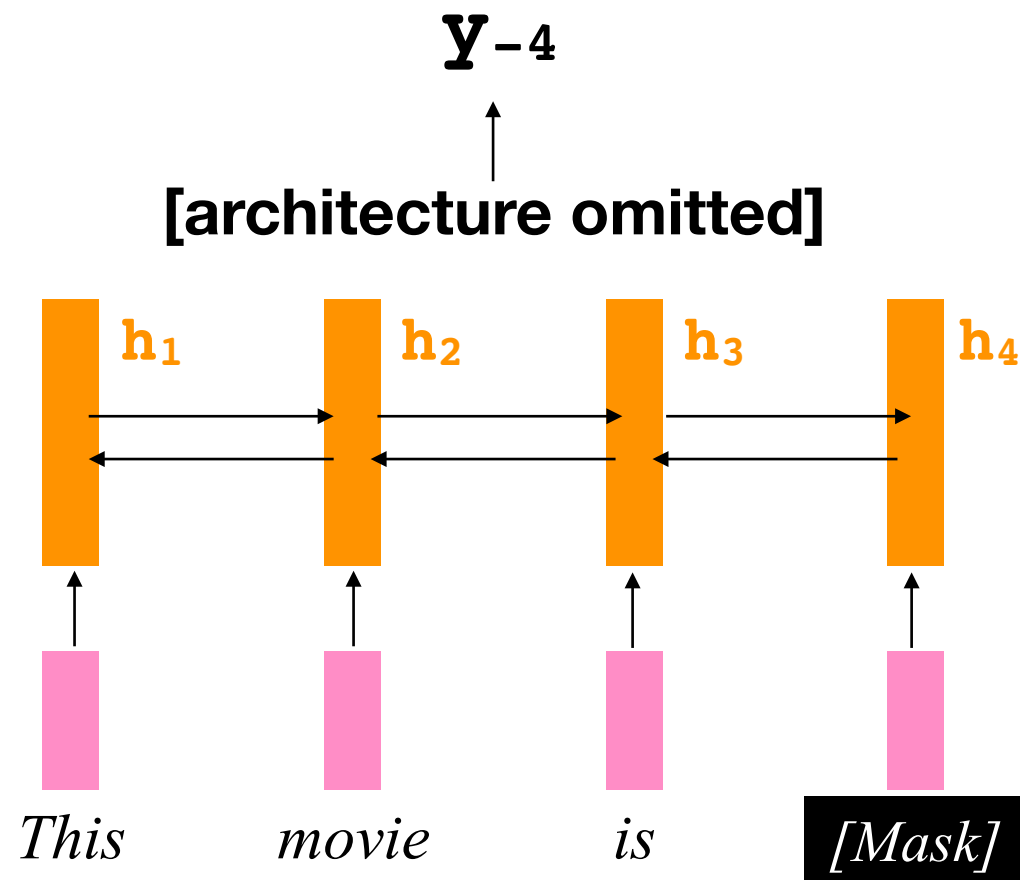
# Prior Works on Classification

- If the model is trained in a standard way, attention weights correlate strongly with individual token influence.

- Uniform attention gives the same accuracy (attention does not matter) .

- Models can be trained to attend to irrelevant words, *without harming accuracy*.

# Text Classification

FinalLinearLayer

$$y = Wc$$

*Positive*
*Neutral*
*Negative*

Aggregate

$$c = a_1h_1 + a_2h_2 + a_3h_3 + a_4h_4$$

$a_1$ $a_2$ $a_3$ $a_4$   Attention Weights!

Softmax

$a'_1$ $a'_2$ $a'_3$ $a'_4$

TwoLayerMLP

LSTM

$h_1$ $h_2$ $h_3$ $h_4$

Word Embedding

*This* *movie* *is* *great*

# Attention ≈ Explanations

$\mathbf{y_{-4}}$

[architecture omitted]

$\mathbf{h_1}$   $\mathbf{h_2}$   $\mathbf{h_3}$   $\mathbf{h_4}$

*This*   *movie*   *is*   *[Mask]*

Leave-one-out influence:
$$\Delta_4 \ := \ y \ - \ y_{-4}$$

|  | *This* | *movie* | *is* | *great* |  |
|---|---|---|---|---|---|
| Attention | $\mathbf{a_1}$ | $\mathbf{a_2}$ | $\mathbf{a_3}$ | $\uparrow\mathbf{a_4}$ | Strongly Correlates |
| Influence | $\Delta_1$ | $\Delta_2$ | $\Delta_3$ | $\uparrow\Delta_4$ | |

# Train "Deceptive" Attention

Call a standard model P.
Now train the deceptive model Q.

$$\mathcal{L}(\mathcal{M}_P, \mathcal{M}_Q) = TVD(\hat{y_P}, \hat{y_Q}) - \lambda KL(a_P \| a_Q)$$

Q makes similar
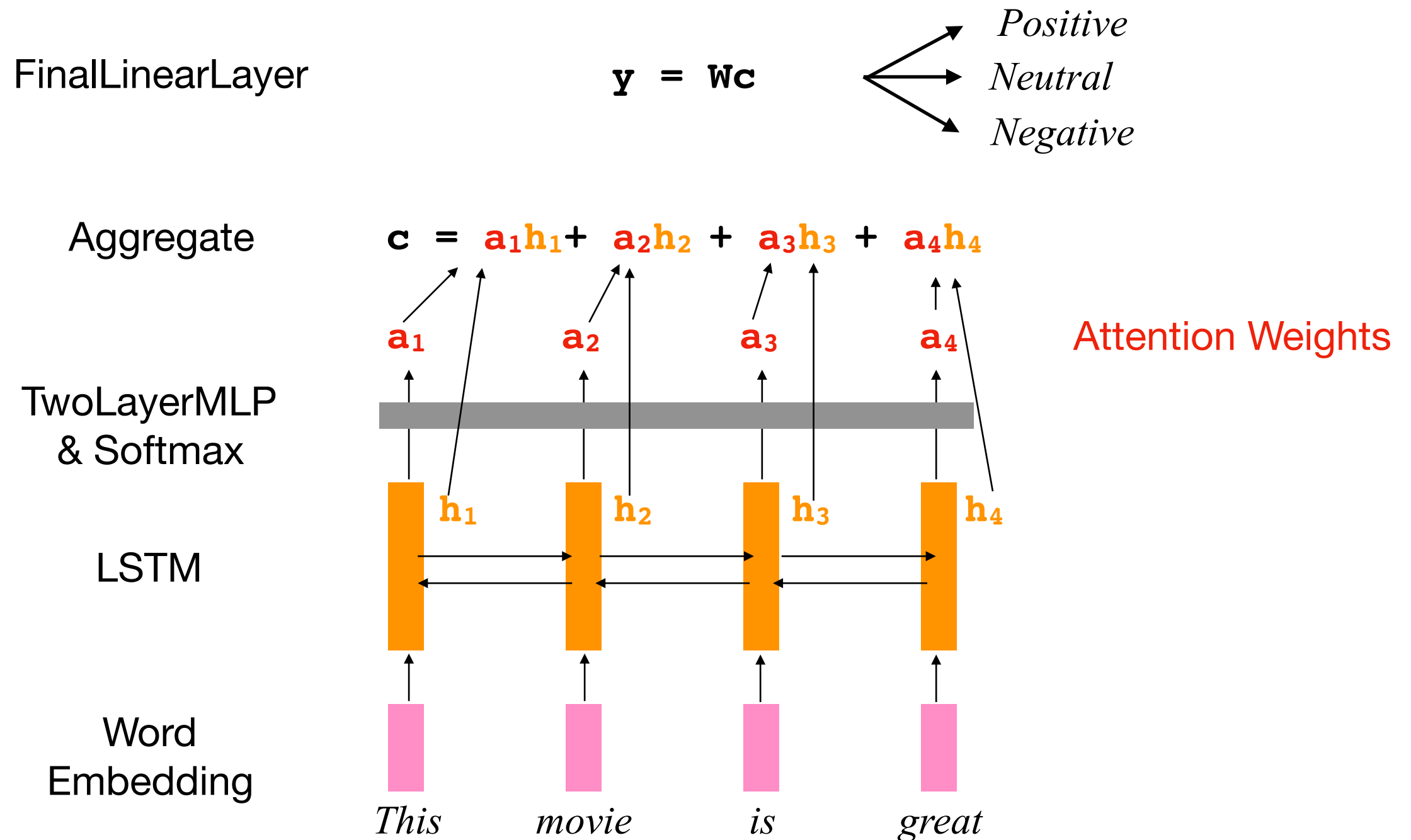predictions as P.

Q attends
differently from P.

**Attention becomes anti-correlated
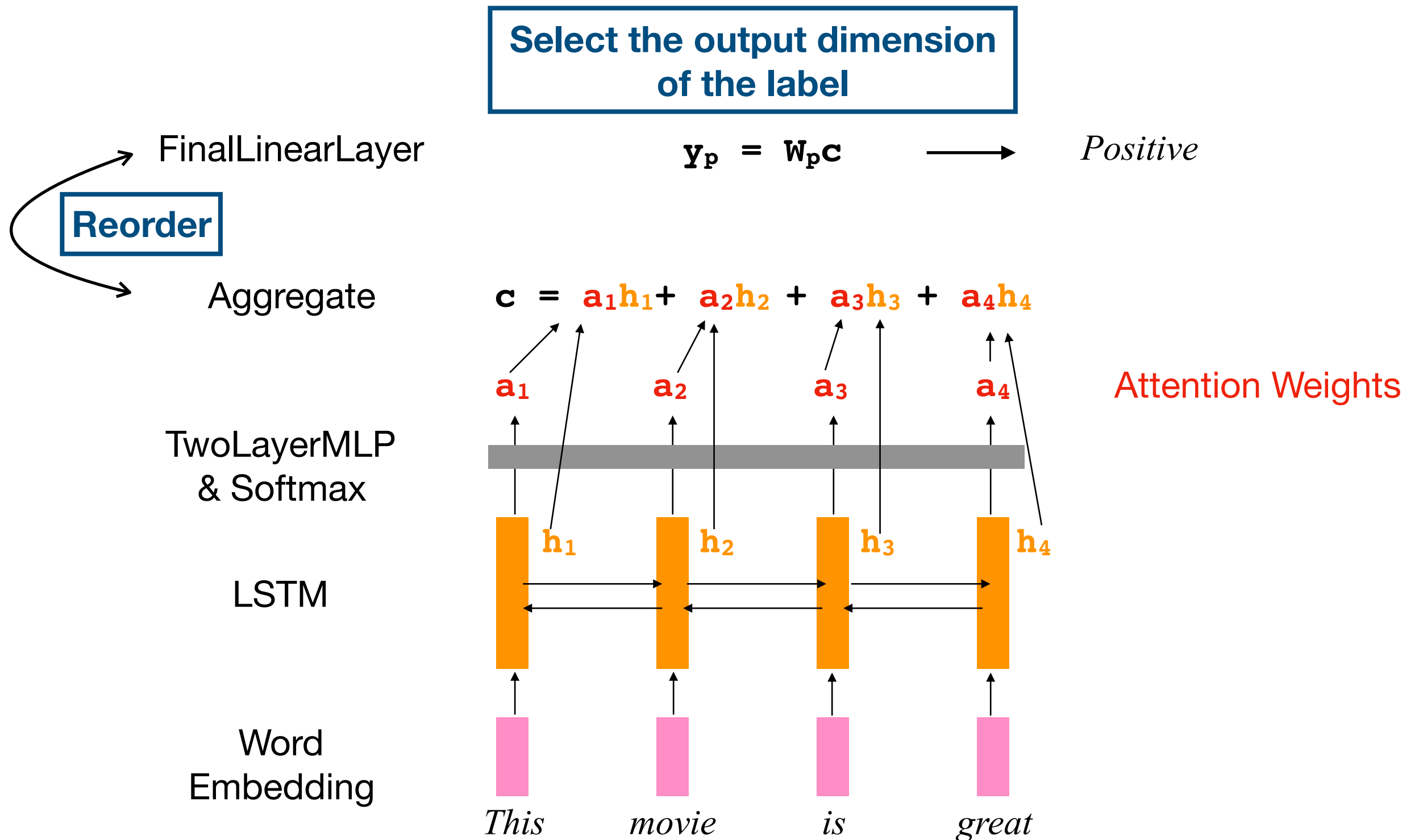with explanations !**

# Our Work:
# Understand Attention Training

- We can construct two models with the same loss but different attention weights.

- Need to open the black-box of standard training.

# Defining Output Relevance

FinalLinearLayer $\quad$ $\mathtt{y = Wc}$

*Positive*
*Neutral*
*Negative*

Aggregate $\quad$ $\mathtt{c = a_1 h_1 + a_2 h_2 + a_3 h_3 + a_4 h_4}$

$\mathtt{a_1}$ $\qquad$ $\mathtt{a_2}$ $\qquad$ $\mathtt{a_3}$ $\qquad$ $\mathtt{a_4}$ $\qquad$ Attention Weights

TwoLayerMLP & Softmax

$\mathtt{h_1}$ $\qquad$ $\mathtt{h_2}$ $\qquad$ $\mathtt{h_3}$ $\qquad$ $\mathtt{h_4}$

LSTM

Word Embedding

*This* $\qquad$ *movie* $\qquad$ *is* $\qquad$ *great*

# Defining Output Relevance

Select the output dimension of the label

FinalLinearLayer $\qquad$ $\mathbf{y_p} = \mathbf{W_p c}$ $\qquad\longrightarrow\qquad$ *Positive*

Reorder

Aggregate $\qquad$ $\mathbf{c} = \mathbf{a_1 h_1} + \mathbf{a_2 h_2} + \mathbf{a_3 h_3} + \mathbf{a_4 h_4}$

$\mathbf{a_1}$ $\qquad$ $\mathbf{a_2}$ $\qquad$ $\mathbf{a_3}$ $\qquad$ $\mathbf{a_4}$ $\qquad$ Attention Weights

TwoLayerMLP & Softmax

$\mathbf{h_1}$ $\qquad$ $\mathbf{h_2}$ $\qquad$ $\mathbf{h_3}$ $\qquad$ $\mathbf{h_4}$

LSTM

Word Embedding

*This* $\qquad$ *movie* $\qquad$ *is* $\qquad$ *great*

# Defining Output Relevance

**Training objective optimizes $y_p$**

**Output Relevance $r_4$: how much does the model associate $h_4$ with the positive label**

Aggregate    $y_p = a_1 r_1 + a_2 r_2 + a_3 r_3 + a_4 r_4 \longrightarrow$ *Positive* Logit

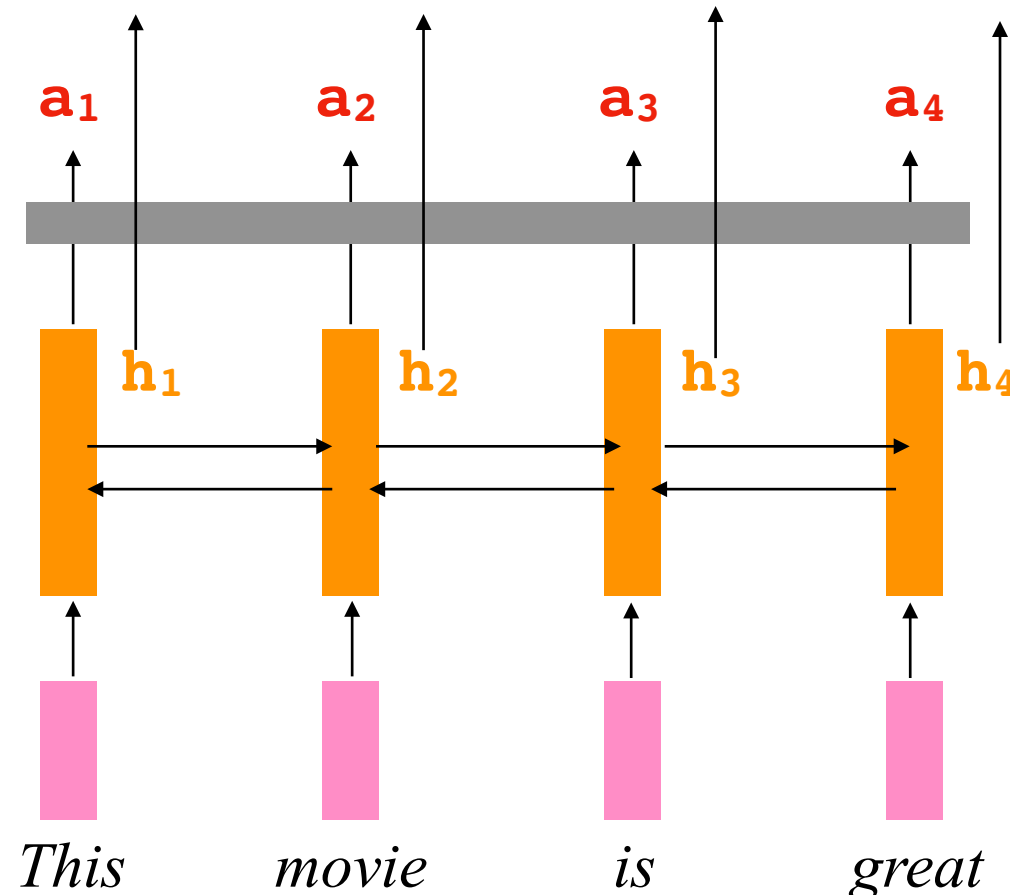FinalLinearLayer    $r_1 = W_p h_1$    $r_2 = W_p h_2$    $r_3 = W_p h_3$    $r_4 = W_p h_4$

$a_1$    $a_2$    $a_3$    $a_4$    Attention Weights

TwoLayerMLP & Softmax

$h_1$    $h_2$    $h_3$    $h_4$

LSTM

Word Embedding

*This*    *movie*    *is*    *great*

# Attention $a$ and Output Relevance $r$

Training objective maximizes $y_p$

Output Relevance $r_4$:
how much does the model associate $h_4$ with the positive label

$$y_p = a_1 r_1 + a_2 r_2 + a_3 r_3 + a_4 r_4 \longrightarrow \textit{Positive } \text{Logit}$$

- $a$ attracted to larger $r$

- What is $a$ and $r$ when initialized? => both uniform

- How will $a$ and $r$ interact? Near the beginning, $a$ remains uniform.  Under uniform attention, $r$ increases faster at "keyword" positions, then attracts $a$

# Attention **a** and Output Relevance **r**

Output Relevance $r_4$:
how much does the model associate
$h_4$ with the positive label

$$y_p = a_1 r_1 + a_2 r_2 + a_3 r_3 + a_4 r_4 \longrightarrow \textit{Positive} \text{ Logit}$$

- **a** attracted to larger **r**

- What is **a** and **r** when initialized? => both uniform

- How will **a** and **r** interact? Near the beginning, **a** remains uniform.  Under uniform attention, <u>r increases faster at "keyword" positions</u>, then attracts **a**

?

# Intuition of Increasing r

Aggregate $\qquad$ $y_p = a_1 r_1 + a_2 r_2 + a_3 r_3 + a_4 r_4 \longrightarrow$ *Positive* Logit

FinalLinearLayer $\qquad$ $r_1 = W_p h_1 \quad r_2 = W_p h_2 \quad r_3 = W_p h_3 \quad r_4 = W_p h_4$

$a_1 \qquad\qquad a_2 \qquad\qquad a_3 \qquad\qquad a_4$ $\qquad$ Attention Weights

TwoLayerMLP
& Softmax

$h_1 \qquad\qquad h_2 \qquad\qquad h_3 \qquad\qquad h_4$

LSTM

Word
Embedding

*This* $\qquad$ *movie* $\qquad$ *is* $\qquad$ *great*

# Intuition of Increasing r

Aggregate  $\mathbf{y_p} = \mathbf{0.25}(\mathbf{r_{This}} + \mathbf{r_{movie}} + \mathbf{r_{is}} + \mathbf{r_{great}}) \longrightarrow$ *Positive* Logit

FinalLinearLayer  $\mathbf{r_{This}}$   $\mathbf{r_{movie}}$   $\mathbf{r_{is}}$   $\mathbf{r_{great}}$

**.25**   **.25**   **.25**   **.25**   Attention Weights

TwoLayerMLP
& Softmax

$\mathbf{h_{This}}$   $\mathbf{h_{movie}}$   $\mathbf{h_{is}}$   $\mathbf{h_{great}}$

LSTM

Word
Embedding

*This*   *movie*   *is*   *great*

# Intuition of Increasing r

$$y_p = 0.25(r_{This} + r_{movie} + r_{is} + r_{great}) \longrightarrow \textit{Positive } \text{Logit}$$

$$y_n = -0.25(r_{This} + r_{movie} + r_{is} + r_{bad}) \longrightarrow \textit{Negative } \text{Logit}$$

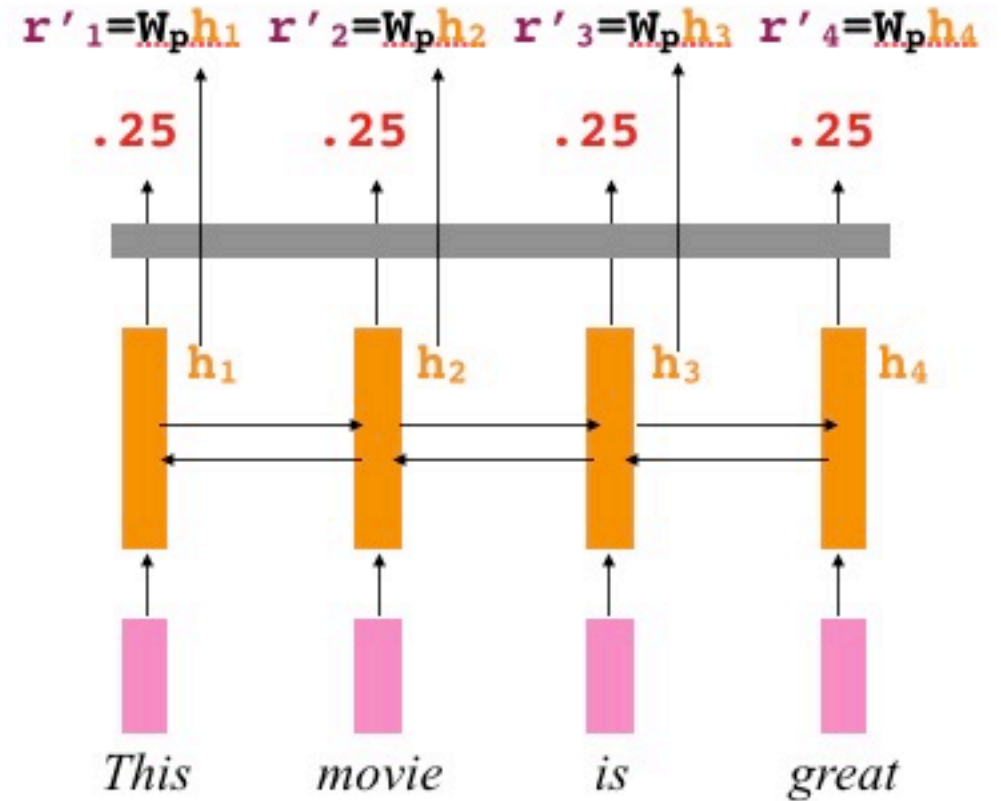- $r_{This}$, $r_{movie}$ and $r_{is}$ remains roughly unchanged (cancels out)

- $r_{great}$ increases

# Increasing r Under Uniform Attention

$$y_p = a_1r_1 + a_2r_2 + a_3r_3 + a_4r_4$$

$$y_p = .25 (r'_1 + r'_2 + r'_3 + r'_4)$$

$r_1 = W_p h_1 \quad r_2 = W_p h_2 \quad r_3 = W_p h_3 \quad r_4 = W_p h_4$

$r'_1 = W_p h_1 \quad r'_2 = W_p h_2 \quad r'_3 = W_p h_3 \quad r'_4 = W_p h_4$

$a_1 \qquad a_2 \qquad a_3 \qquad a_4$

$.25 \qquad .25 \qquad .25 \qquad .25$

$h_1 \qquad h_2 \qquad h_3 \qquad h_4$

$h_1 \qquad h_2 \qquad h_3 \qquad h_4$

*This    movie    is    great*

*This    movie    is    great*

standard training
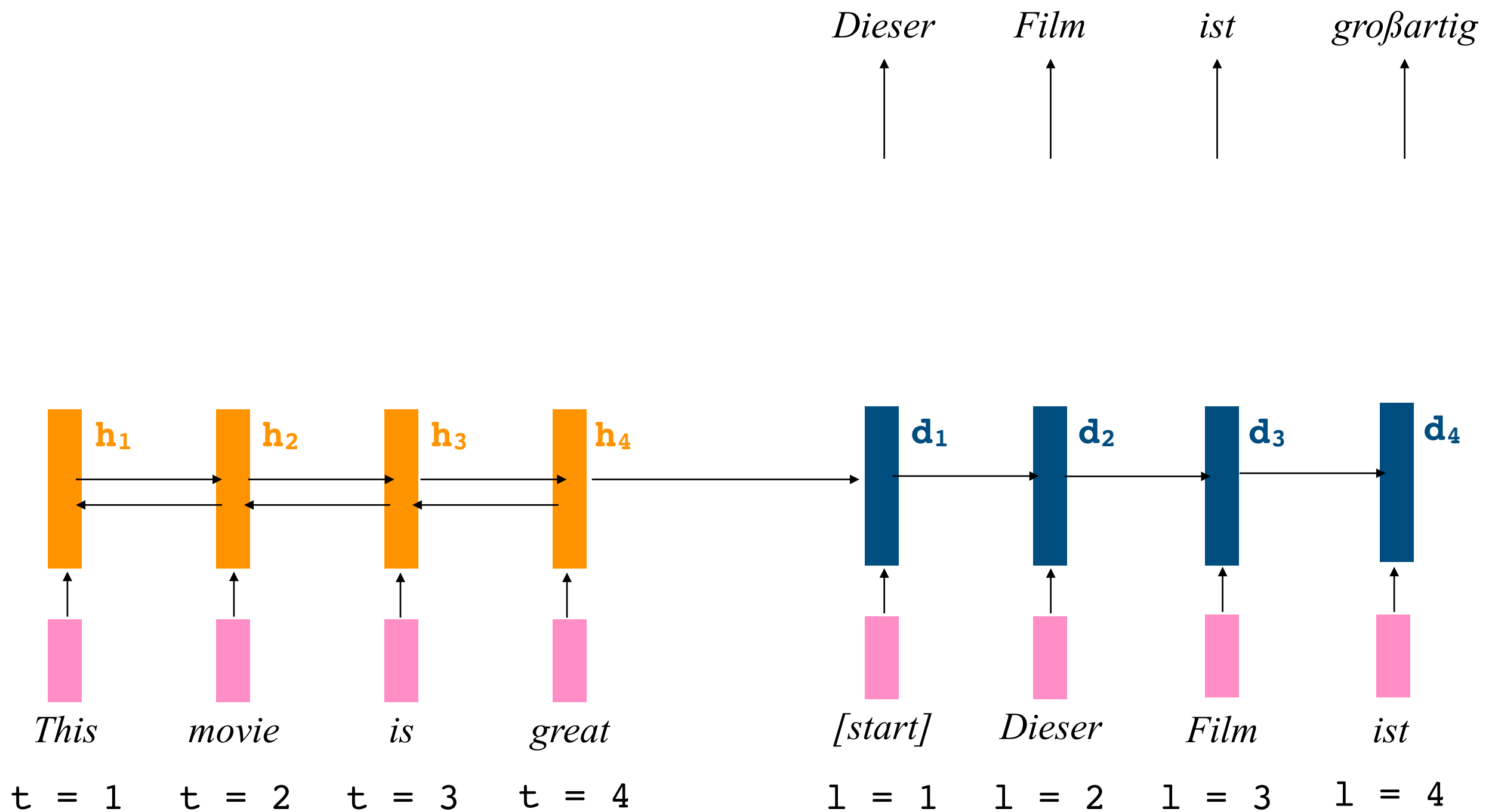
uniform attention training

- **r'** still correlates with standard attention **a**

- **r'** correlates with individual token influence.

# Seq2Seq Setup

(simplified)

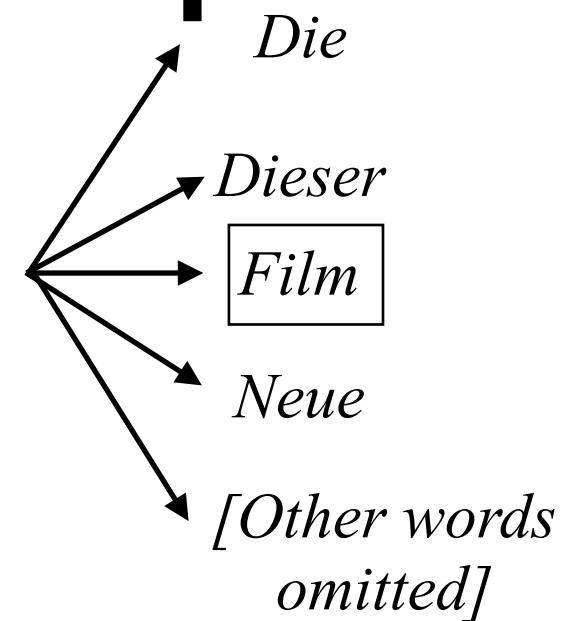*Dieser*        *Film*        *ist*        *großartig*

| $h_1$ | $h_2$ | $h_3$ | $h_4$ | | $d_1$ | $d_2$ | $d_3$ | $d_4$ |

*This*        *movie*        *is*        *great*        *[start]*        *Dieser*        *Film*        *ist*

t = 1        t = 2        t = 3        t = 4                l = 1        l = 2        l = 3        l = 4

# Seq2Seq Setup

(simplified)

$$y_2 = Wc_2$$

*Die*

*Dieser*

$\boxed{\textit{Film}}$

*Neue*

*[Other words omitted]*

**Reorder**

$$c_2 = a_{21}h_1 + a_{22}h_2 + a_{23}h_3 + a_{24}h_4$$

$a'_{21}=h_1d_2$  $a'_{22}=h_2d_2$  $a'_{23}=h_3d_2$  $a'_{24}=h_4d_2$

$h_1$  $h_2$  $h_3$  $h_4$  $d_1$  $d_2$

**To produce the second output word**

*This*  *movie*  *is*  *great*  *[start]*  *Dieser*

t = 1  t = 2  t = 3  t = 4  l = 1  l = 2
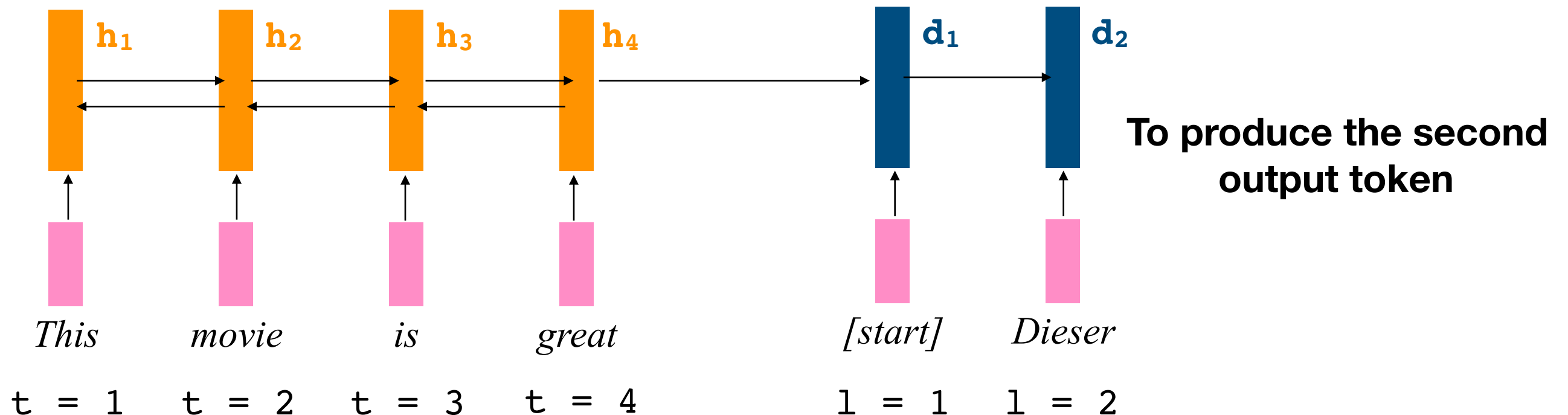
# Seq2Seq Output Relevance

(simplified)

**Training objective optimizes $y_{2,Film}$**

**Output Relevance $r_{Film, 2}$: how much does the model associate $h_2$ with the word "Film"**

$$y_{2,Film} = a_{21}r_{Film,1} + a_{22}r_{Film,2} + a_{23}r_{Film,3} + a_{24}r_{Film,4}$$

$$r_{Film,t} := W_{Film} \, h_t$$

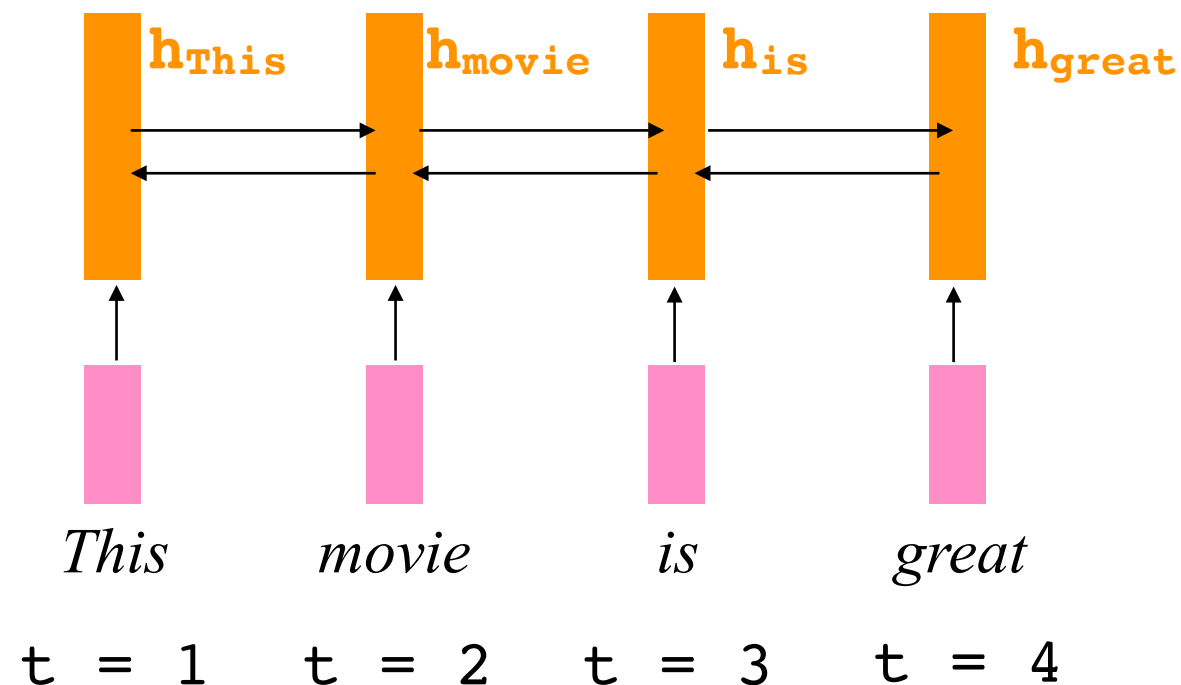$$a'_{21} = h_1 d_2 \quad a'_{22} = h_2 d_2 \quad a'_{23} = h_3 d_2 \quad a'_{24} = h_4 d_2$$



To produce the second output token

| $h_1$ | $h_2$ | $h_3$ | $h_4$ | | $d_1$ | $d_2$ |

| *This* | *movie* | *is* | *great* | | *[start]* | *Dieser* |

| t = 1 | t = 2 | t = 3 | t = 4 | | l = 1 | l = 2 |

# Intuition of Learning r, "Dictionary"

**Gradient update to maximize:**

Likelihood($l = 2$) = $.25(r_{Film, This} + r_{Film,movie} + r_{Film,is} + r_{Film,great})$

$r_{Film, This}$   $r_{Film,movie}$   $r_{Film,is}$   $r_{Film,great}$

.25      .25      .25      .25

$h_{This}$    $h_{movie}$    $h_{is}$    $h_{great}$

*This*     *movie*     *is*     *great*

t = 1    t = 2    t = 3    t = 4

# Intuition of Learning r, "Dictionary"

**Gradient update to maximize: Likelihood($l$ = 2) = $\Sigma_t$ r$_{Film,t}$**

**$t \in$ {This, movie, is, great }**

$\mathbf{r}_{Film, This}$ $\quad$ $\mathbf{r}_{Film,movie}$ $\quad$ $\mathbf{r}_{Film,is}$ $\quad$ $\mathbf{r}_{Film,great}$

.25 $\qquad$ .25 $\qquad$ .25 $\qquad$ .25

h$_{This}$ $\qquad$ h$_{movie}$ $\qquad$ h$_{is}$ $\qquad$ h$_{great}$

*This* $\qquad$ *movie* $\qquad$ *is* $\qquad$ *great*

t = 1 $\quad$ t = 2 $\quad$ t = 3 $\quad$ t = 4

# Intuition of Learning r, "Dictionary"

This movie is great
=> Dieser Film ist großartig

**Gradient update to maximize:**
**Loss(sentence) $\Sigma_l \Sigma_t r_{l,t}$**
**$t \in$ {This, movie, is, great }**
**$l \in$ {Dieser, Film, ist, großartig}**

How many times $r_{l,t}$ appears in the sum.

| r | This | Movie | Is | Great |
|---|---|---|---|---|
| Dieser | 1 | 1 | 1 | 1 |
| Film | 1 | 1 | 1 | 1 |
| Ist | 1 | 1 | 1 | 1 |
| großartig | 1 | 1 | 1 | 1 |

**Looks like there is no reason**
**for the model to learn word-to-word correspondence**

# Intuition of Learning r, "Dictionary"

This movie is great
=> Dieser Film ist großartig

**Gradient update to maximize: $\Sigma_l \, \Sigma_t \, r_{l,t}$**
**t ∈ {This, movie, is, great }**
**l ∈ {Dieser, Film, ist, großartig}**

**Gradient update to maximize: $\Sigma_l \, \Sigma_t \, r_{l,t}$**
**t ∈ {This, movie, is, bad }**
**l ∈ {Dieser, Film, ist, schlecht}**

This movie is bad
=> Dieser Film ist schlecht

add

| r | This | Movie | Is | Great | Bad |
|---|---|---|---|---|---|
| Dieser | 2 | 2 | 2 | 1 | 1 |
| Film | 2 | 2 | 2 | 1 | 1 |
| Ist | 2 | 2 | 2 | 1 | 1 |
| großartig | 1 | 1 | 1 | 1 | 0 |
| schlecht | 1 | 1 | 1 | 0 | 1 |

**Can recover r (word-to-word correspondence) from "co-occurrence".**

# Testing the "co-occurrence" Intuition

- Hypothesis: If we remove the co-occurrence statistics, $r$ cannot be learned, and hence attention $a$ fails to learn.

- Experiments on a "sequence copying task".

  - Setting 1: the model learns the copying task from a distribution of permutations of [1, 40].

    **~50% fails**

  - Setting 2: the model to learns from a distribution of length 40 array, each token is a uniform i.i.d. sample from [1, 40]. **Always successful**

```
input:  3 2 1 0 ; output:  3' 2' 1' 0'
input:  0 2 1 3 ; output:  0' 2' 1' 3'
input:  3 2 0 1 ; output:  3' 2' 0' 1'
input:  0 2 1 3 ; output:  0' 2' 1' 3'
input:  3 2 0 1 ; output:  3' 2' 0' 1'
```

```
input:  0 3 3 1 2 ; output:  2' 1' 1' 3' 0'
input:  1 2 2 3 3 ; output:  3' 0' 0' 1' 1'
input:  2 2 1 2 3 ; output:  0' 0' 3' 0' 1'
input:  3 2 1 0 0 ; output:  1' 0' 3' 2' 2'
input:  3 0 2 1 2 ; output:  1' 2' 0' 3' 0'
```

# Takeaways ...

- Interpretable attention might not be necessary to achieve high accuracy (e.g. in text classification).

- Attention is shaped by training dynamics.

- Open the blackbox of training to understand neural networks.

# References

- Jain, Sarthak, and Byron C. Wallace. "Attention is not explanation." *arXiv preprint arXiv:1902.10186* (2019).

- Wiegreffe, Sarah, and Yuval Pinter. "Attention is not not explanation." arXiv preprint arXiv:1908.04626 (2019).

- Snell, Charlie, et al. "Understanding Attention Training via Output Relevance" OpenReview Preprint (2020)