

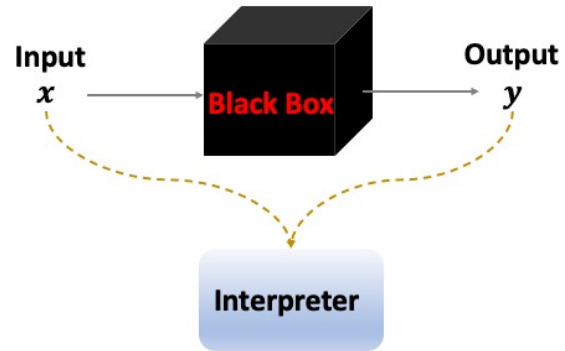
CS 4501/6501 Interpretable Machine Learning

Post-hoc explanations: beyond feature-level

Hanjie Chen, Yangfeng Ji
Department of Computer Science
University of Virginia
{hc9mx, yangfeng}@virginia.edu

Explaining Black-box Model

Model-agnostic



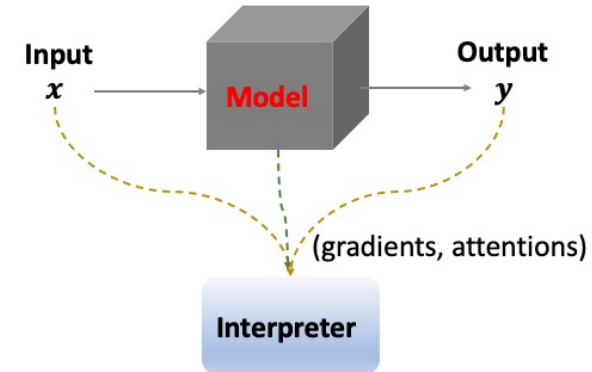
LIME

- Applicable to any black-box models
- Computational complexity
- Work well on traditional models (e.g., CNN), but not on complex DNN

SHAP

- Applicable to any black-box models
- Computational complexity
- Best performance (empirically)

Model-dependent



IG

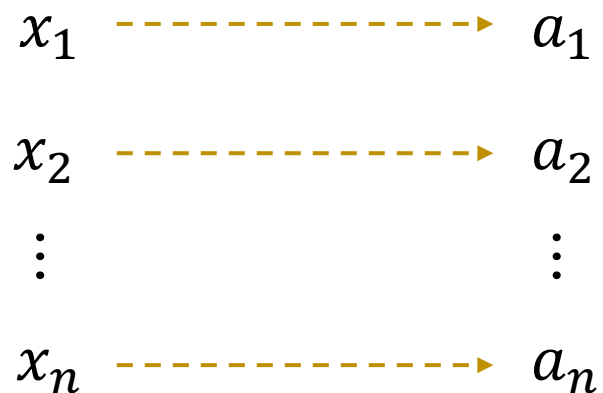
- Require access to model gradients
- Simple, fast
- Work well on both traditional and DNN models

Attention

- Rely on attention mechanism
- Simple, fast (no additional computation)
- Not clear (much debate)

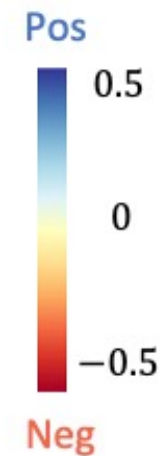
Single Feature-level Explanation

Input features Importance



Explanation

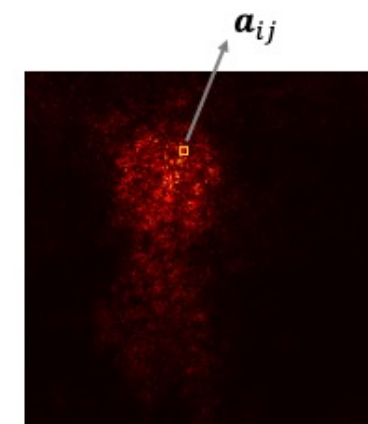
$a_1 = 0.11$	a
$a_2 = 0.46$	clever
$a_3 = 0.01$	piece
$a_4 = -0.02$	of
$a_5 = 0.06$	cinema



Input



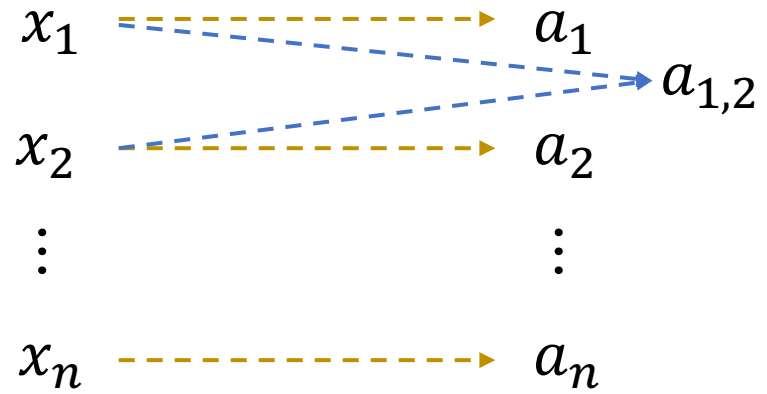
Explanation



Composite Feature-level Explanation

Input features

Importance

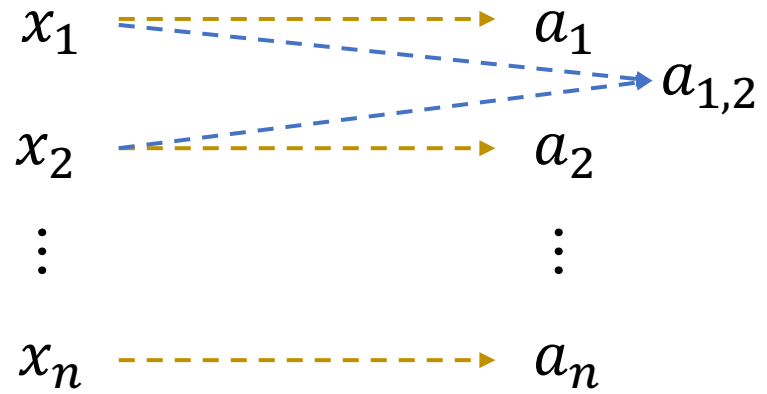


When single features have interactions, it is critical to know the importance of the composite feature composed with these single features

Composite Feature-level Explanation

Input features

Importance



Example 1

Input
(prediction: **positive**)

not

a

bad

journey

Explanation

not

a

bad

journey

Pos

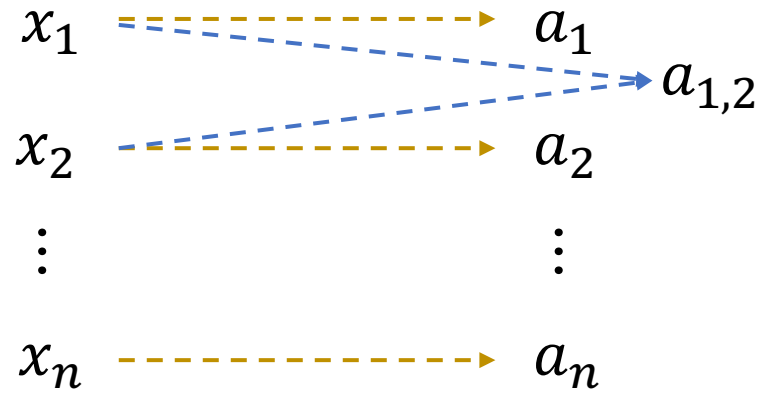
Neg

Why does the model think
"journey" as positive?

Composite Feature-level Explanation

Input features

Importance



Example 1

Input
(prediction: **positive**)

not

a

bad

journey

Composite feature
importance

Explanation

not

a

bad

journey

not bad

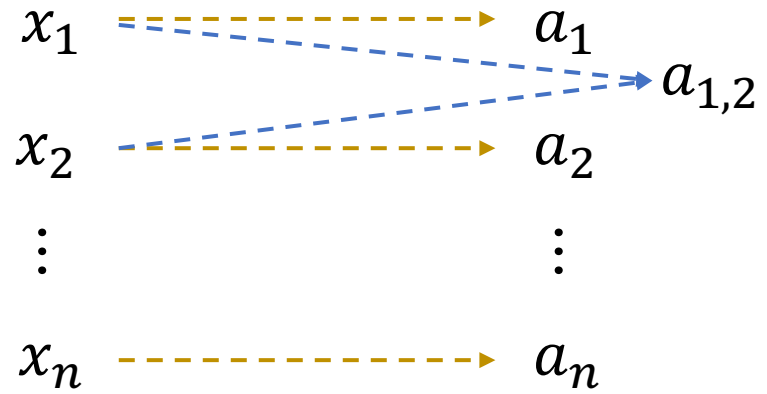
Pos

Neg

Composite Feature-level Explanation

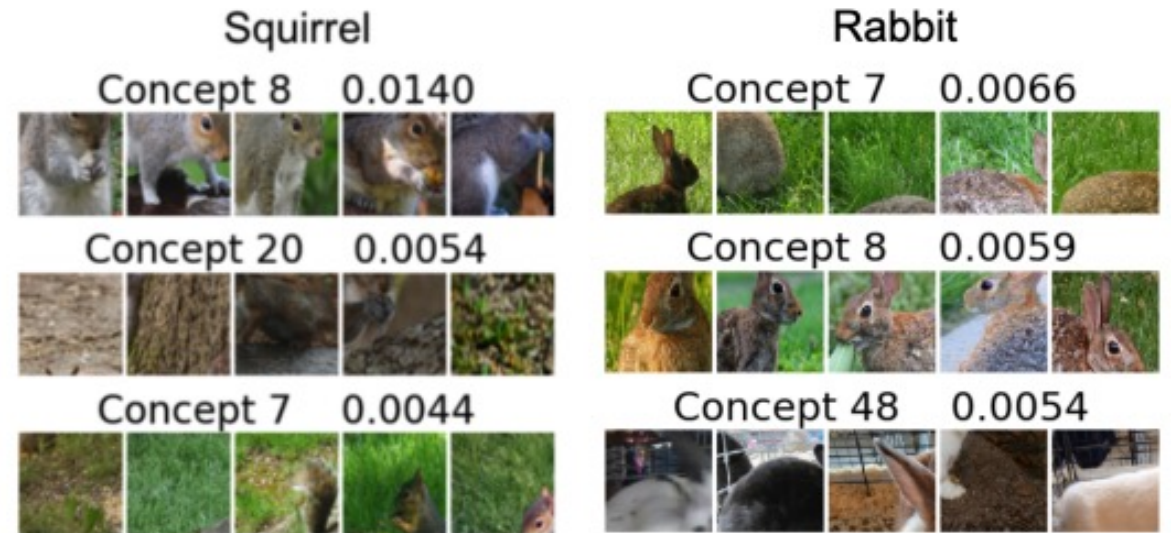
Input features

Importance



Example 2

Composite feature importance is more intuitive for human understanding



(Yeh et al., 2020)

Beyond Feature Attribution

- Contextual Decomposition (CD)
- Hierarchical Explanation via Divisive Generation (HEDGE)

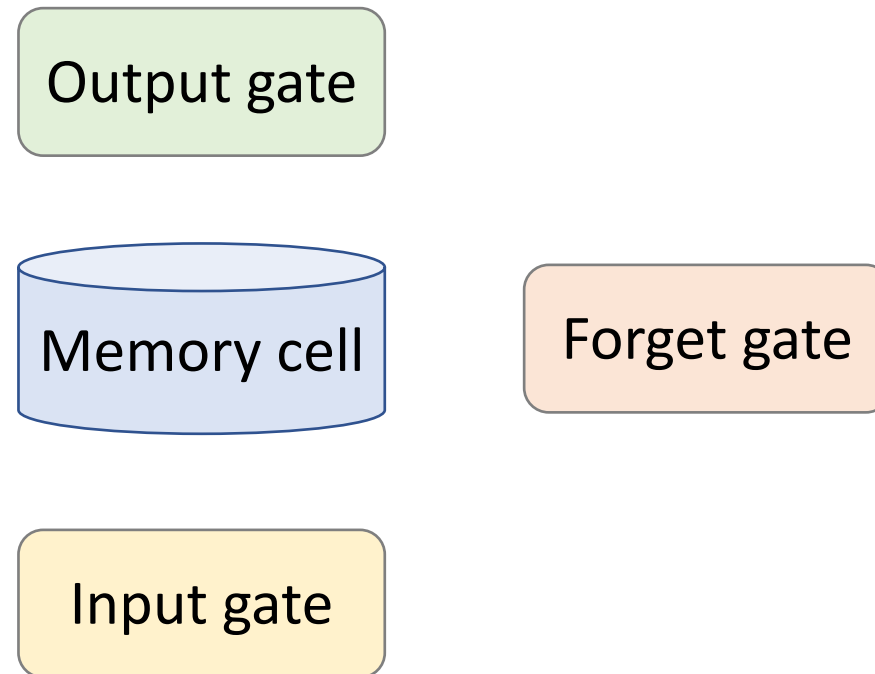
Beyond Word Importance: Contextual Decomposition to Extract Interactions From LSTMs

W. James Murdoch, Peter J. Liu, Bin Yu

(ICLR, 2018)

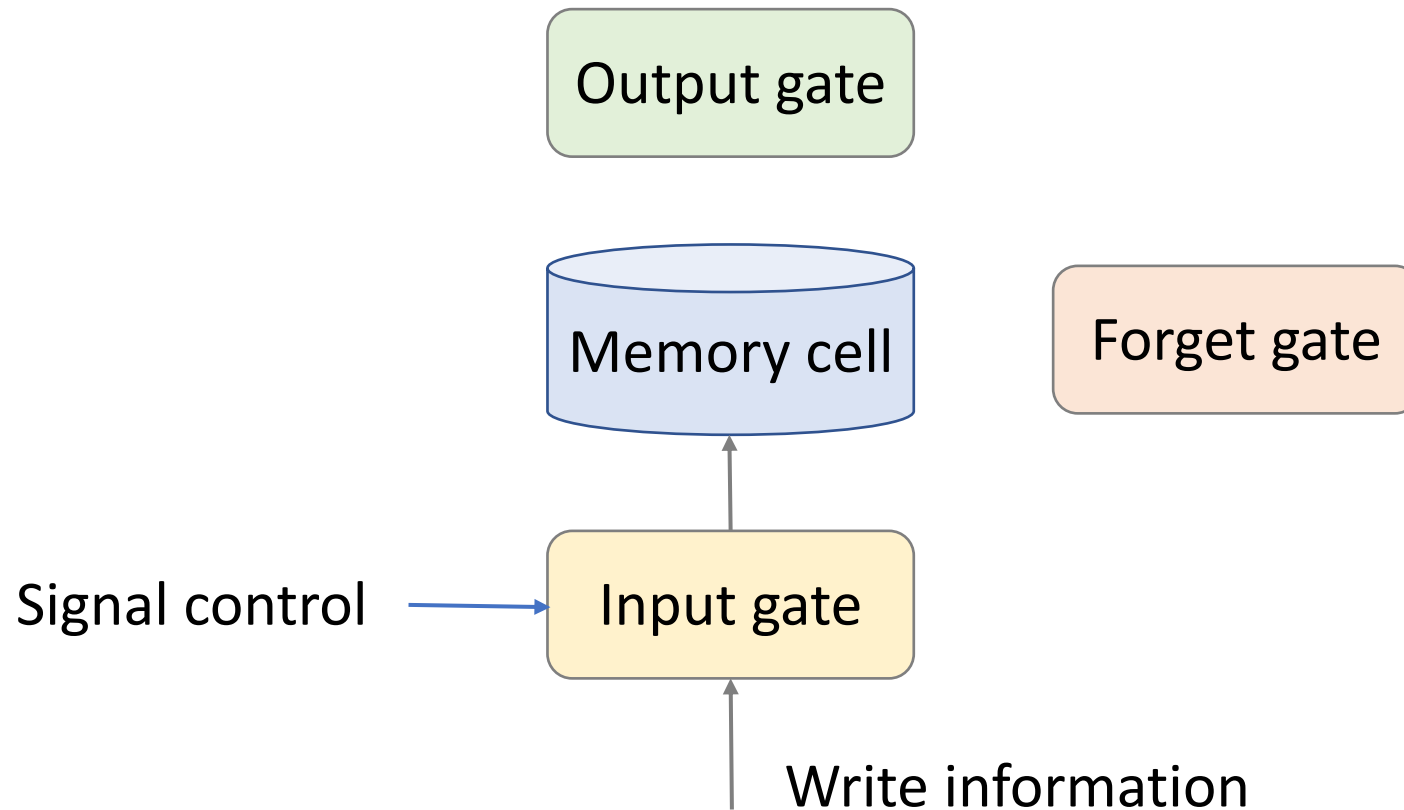
LSTM

- Long Short-term Memory Network (LSTM) [Hochreiter and Schmidhuber, 1997]



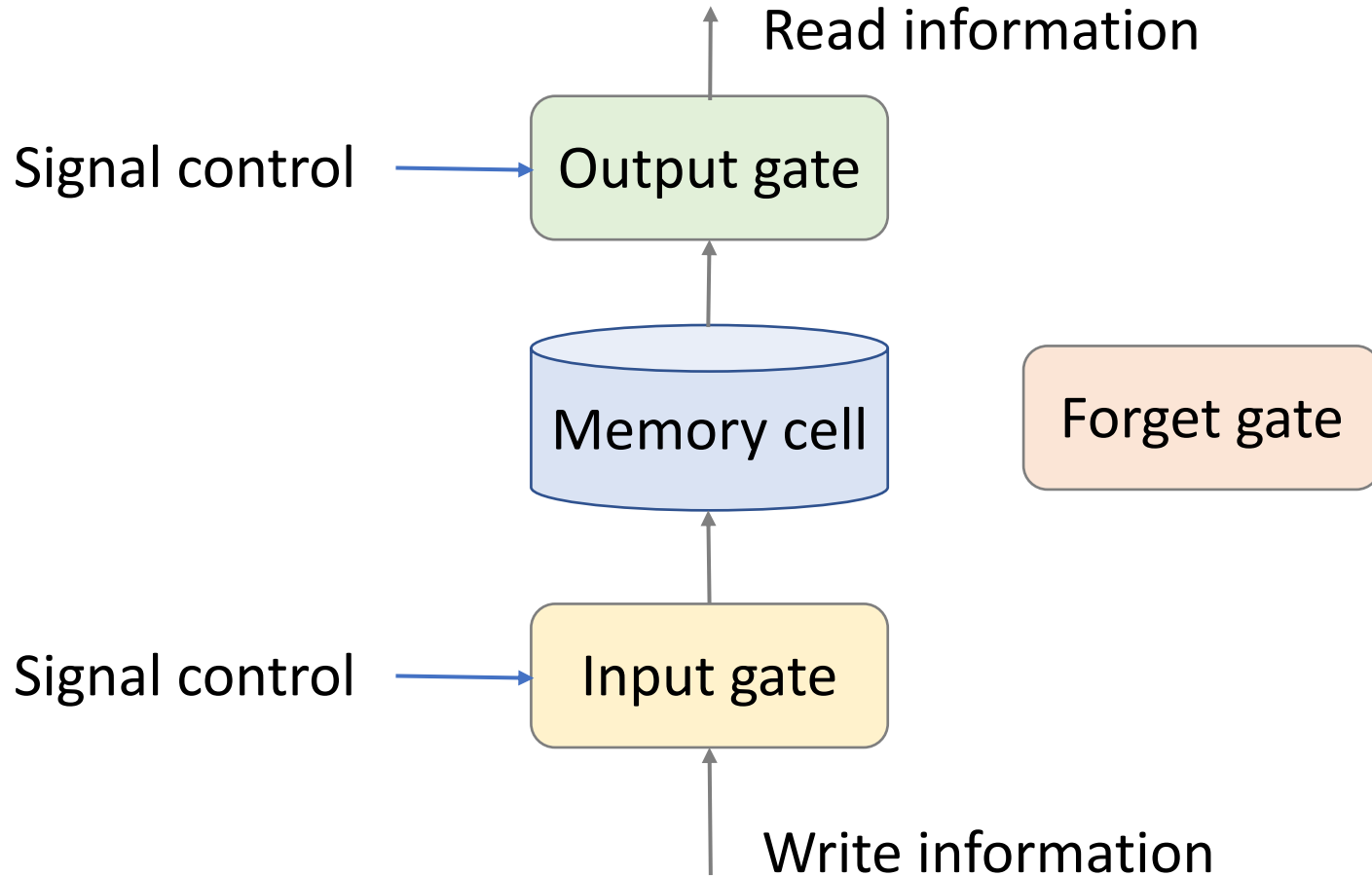
LSTM

- Long Short-term Memory Network (LSTM) [Hochreiter and Schmidhuber, 1997]



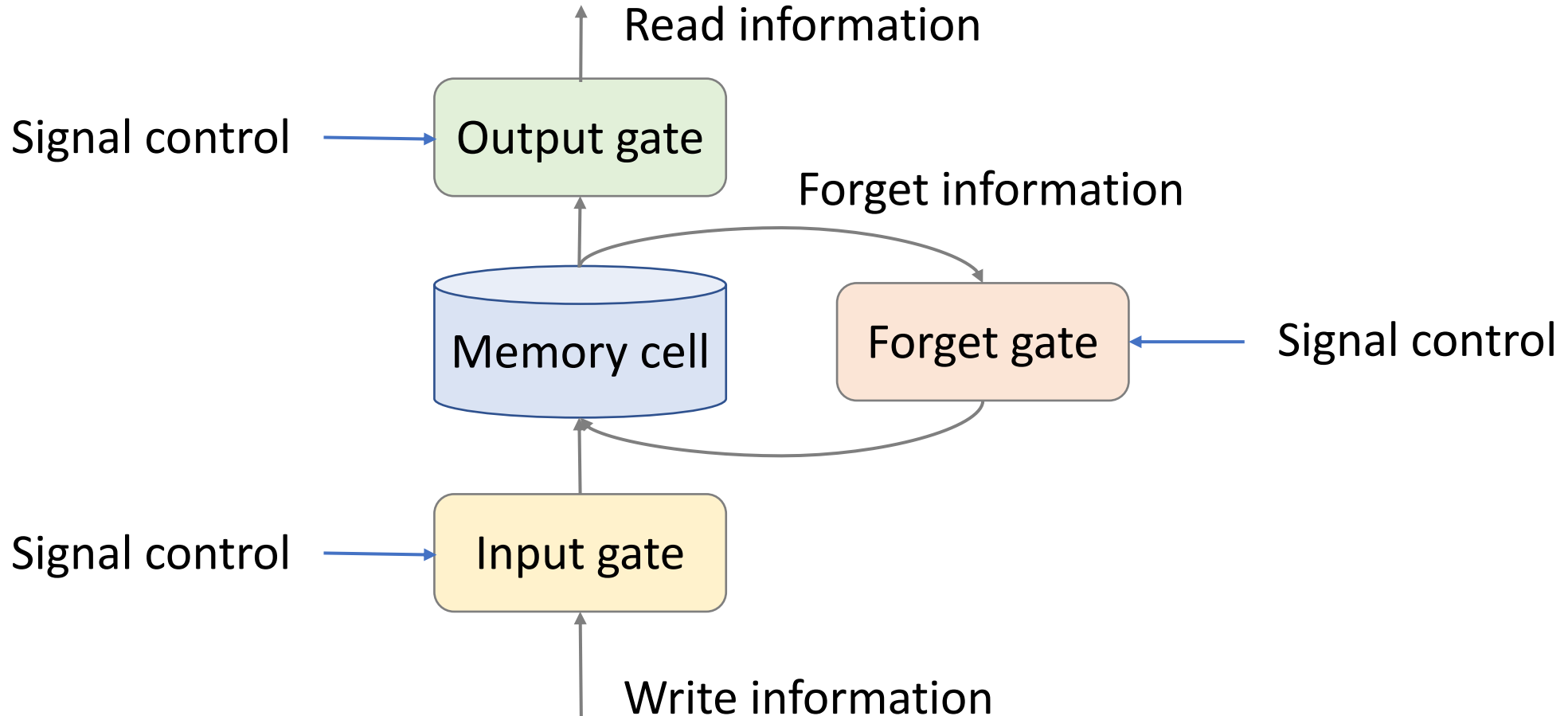
LSTM

- Long Short-term Memory Network (LSTM) [Hochreiter and Schmidhuber, 1997]



LSTM

- Long Short-term Memory Network (LSTM) [Hochreiter and Schmidhuber, 1997]



LSTM

- Long Short-term Memory Network (LSTM) [Hochreiter and Schmidhuber, 1997]

Input word embeddings: $x_1, \dots, x_T \in \mathbb{R}^{d_1}$

Cell: $c_t \in \mathbb{R}^{d_2}$ $(h_0 = c_0 = 0)$

State vector: $h_t \in \mathbb{R}^{d_2}$

LSTM

- Long Short-term Memory Network (LSTM) [Hochreiter and Schmidhuber, 1997]

Input word embeddings: $x_1, \dots, x_T \in \mathbb{R}^{d_1}$

Cell: $c_t \in \mathbb{R}^{d_2}$ ($h_0 = c_0 = 0$)

State vector: $h_t \in \mathbb{R}^{d_2}$

(Output gate) $o_t = \sigma(W_o x_t + V_o h_{t-1} + b_o)$

(Forget gate) $f_t = \sigma(W_f x_t + V_f h_{t-1} + b_f)$

(Input gate) $i_t = \sigma(W_i x_t + V_i h_{t-1} + b_i)$

x_t : current input

h_{t-1} : previous output

LSTM

- Long Short-term Memory Network (LSTM) [Hochreiter and Schmidhuber, 1997]

Input word embeddings: $x_1, \dots, x_T \in \mathbb{R}^{d_1}$

Cell: $c_t \in \mathbb{R}^{d_2}$ $(h_0 = c_0 = 0)$

State vector: $h_t \in \mathbb{R}^{d_2}$

(Output gate) $o_t = \sigma(W_o x_t + V_o h_{t-1} + b_o)$

(Forget gate) $f_t = \sigma(W_f x_t + V_f h_{t-1} + b_f)$

(Input gate) $i_t = \sigma(W_i x_t + V_i h_{t-1} + b_i)$

W, V, b are model parameters

LSTM

- Long Short-term Memory Network (LSTM) [Hochreiter and Schmidhuber, 1997]

Input word embeddings: $x_1, \dots, x_T \in \mathbb{R}^{d_1}$

Cell: $c_t \in \mathbb{R}^{d_2}$ ($h_0 = c_0 = 0$)

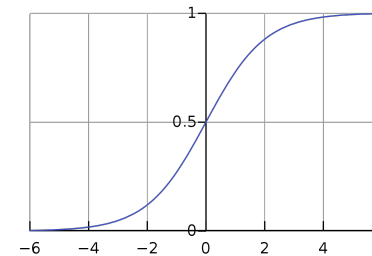
State vector: $h_t \in \mathbb{R}^{d_2}$

(Output gate) $o_t = \sigma(W_o x_t + V_o h_{t-1} + b_o)$

(Forget gate) $f_t = \sigma(W_f x_t + V_f h_{t-1} + b_f)$

(Input gate) $i_t = \sigma(W_i x_t + V_i h_{t-1} + b_i)$

$\sigma(\cdot)$: sigmoid function



Range: 0 to 1

LSTM

- Long Short-term Memory Network (LSTM) [Hochreiter and Schmidhuber, 1997]

Input word embeddings: $x_1, \dots, x_T \in \mathbb{R}^{d_1}$

Cell: $c_t \in \mathbb{R}^{d_2}$ ($h_0 = c_0 = 0$)

State vector: $h_t \in \mathbb{R}^{d_2}$

(Output gate) $o_t = \sigma(W_o x_t + V_o h_{t-1} + b_o)$

(Forget gate) $f_t = \sigma(W_f x_t + V_f h_{t-1} + b_f)$

(Input gate) $i_t = \sigma(W_i x_t + V_i h_{t-1} + b_i)$

$$g_t = \tanh(W_g x_t + V_g h_{t-1} + b_g)$$

$$c_t = f_t \odot c_{t-1} + \boxed{i_t \odot g_t} \text{ Information written into the cell}$$

$$h_t = o_t \odot \tanh(c_t)$$

LSTM

- Long Short-term Memory Network (LSTM) [Hochreiter and Schmidhuber, 1997]

Input word embeddings: $x_1, \dots, x_T \in \mathbb{R}^{d_1}$

Cell: $c_t \in \mathbb{R}^{d_2}$ ($h_0 = c_0 = 0$)

State vector: $h_t \in \mathbb{R}^{d_2}$

(Output gate) $o_t = \sigma(W_o x_t + V_o h_{t-1} + b_o)$

(Forget gate) $f_t = \sigma(W_f x_t + V_f h_{t-1} + b_f)$

(Input gate) $i_t = \sigma(W_i x_t + V_i h_{t-1} + b_i)$

$$g_t = \tanh(W_g x_t + V_g h_{t-1} + b_g)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t \quad \text{Information left in the cell after forgetting}$$

$$h_t = o_t \odot \tanh(c_t)$$

LSTM

- Long Short-term Memory Network (LSTM) [Hochreiter and Schmidhuber, 1997]

Input word embeddings: $x_1, \dots, x_T \in \mathbb{R}^{d_1}$

Cell: $c_t \in \mathbb{R}^{d_2}$ $(h_0 = c_0 = 0)$

State vector: $h_t \in \mathbb{R}^{d_2}$

(Output gate) $o_t = \sigma(W_o x_t + V_o h_{t-1} + b_o)$

(Forget gate) $f_t = \sigma(W_f x_t + V_f h_{t-1} + b_f)$

(Input gate) $i_t = \sigma(W_i x_t + V_i h_{t-1} + b_i)$

$$g_t = \tanh(W_g x_t + V_g h_{t-1} + b_g)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

$$h_t = o_t \odot \tanh(c_t) \text{ Current output}$$

LSTM

- Long Short-term Memory Network (LSTM) [Hochreiter and Schmidhuber, 1997]

Input word embeddings: $x_1, \dots, x_T \in \mathbb{R}^{d_1}$

Cell: $c_t \in \mathbb{R}^{d_2}$ ($h_0 = c_0 = 0$)

State vector: $h_t \in \mathbb{R}^{d_2}$

(Output gate) $o_t = \sigma(W_o x_t + V_o h_{t-1} + b_o)$

(Forget gate) $f_t = \sigma(W_f x_t + V_f h_{t-1} + b_f)$

(Input gate) $i_t = \sigma(W_i x_t + V_i h_{t-1} + b_i)$

$$g_t = \tanh(W_g x_t + V_g h_{t-1} + b_g)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

$$h_t = o_t \odot \tanh(c_t)$$

$t = 1, \dots, T$
→

Probability distribution

$$p = \text{softmax}(W h_T)$$

Question?

CD

- Contextual Decomposition

An arbitrary phrase: x_q, \dots, x_r ($1 \leq q \leq r \leq T$)

Decompose each c_t and h_t into a sum of two contributions

$$\begin{aligned} h_t &= \beta_t + \gamma_t & \beta_t, \beta_t^c &: \text{contributions made solely by the given phrase} \\ c_t &= \beta_t^c + \gamma_t^c & \gamma_t, \gamma_t^c &: \text{contributions involving elements outside of the phrase} \end{aligned}$$

Goal: compute the contribution of the phrase to model prediction

CD

- Contextual Decomposition

An arbitrary phrase: x_q, \dots, x_r ($1 \leq q \leq r \leq T$)

Decompose each c_t and h_t into a sum of two contributions

$$\begin{aligned} h_t &= \beta_t + \gamma_t & \beta_t, \beta_t^c &: \text{contributions made solely by the given phrase} \\ c_t &= \beta_t^c + \gamma_t^c & \gamma_t, \gamma_t^c &: \text{contributions involving elements outside of the phrase} \end{aligned}$$

$$p = \text{softmax}(Wh_T) \longrightarrow p = \text{softmax}(W\underline{\beta}_T + W\gamma_T)$$

the phrase's contribution
to the LSTM's prediction

Goal: compute the contribution of the phrase to model prediction

CD

- Contextual Decomposition

An arbitrary phrase: x_q, \dots, x_r ($1 \leq q \leq r \leq T$)

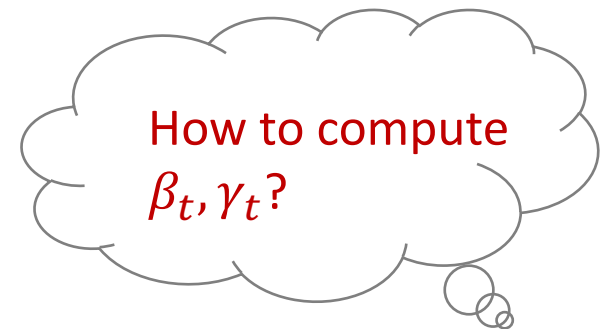
Decompose each c_t and h_t into a sum of two contributions

$$\begin{aligned} h_t &= \beta_t + \gamma_t & \beta_t, \beta_t^c &: \text{contributions made solely by the given phrase} \\ c_t &= \beta_t^c + \gamma_t^c & \gamma_t, \gamma_t^c &: \text{contributions involving elements outside of the phrase} \end{aligned}$$

$$p = \text{softmax}(Wh_T) \longrightarrow p = \text{softmax}(W\underline{\beta}_T + W\gamma_T)$$

the phrase's contribution
to the LSTM's prediction

Goal: compute the contribution of the phrase to model prediction



CD

- Contextual Decomposition

Disambiguating interactions between gates

$$\begin{aligned}i_t &= \sigma(W_i x_t + V_i h_{t-1} + b_i) \\ &= L_\sigma(W_i x_t) + L_\sigma(V_i h_{t-1}) + L_\sigma(b_i)\end{aligned}$$

$$\begin{aligned}g_t &= \tanh(W_g x_t + V_g h_{t-1} + b_g) \\ &= L_{\tanh}(W_g x_t) + L_{\tanh}(V_g h_{t-1}) + L_{\tanh}(b_g)\end{aligned}$$

Assume we have a way
of linearizing the gates

CD

- Contextual Decomposition

Disambiguating interactions between gates

$$\begin{aligned} & i_t \odot g_t \\ &= (L_\sigma(W_i x_t) + L_\sigma(V_i h_{t-1}) + L_\sigma(b_i)) \odot (L_{\tanh}(W_g x_t) + L_{\tanh}(V_g h_{t-1}) + L_{\tanh}(b_g)) \\ &= (L_\sigma(W_i x_t) + L_\sigma(V_i \beta_{t-1}) + L_\sigma(V_i \gamma_{t-1}) + L_\sigma(b_i)) \odot (L_{\tanh}(W_g x_t) + L_{\tanh}(V_g \beta_{t-1}) \\ &+ L_{\tanh}(V_g \gamma_{t-1}) + L_{\tanh}(b_g)) \end{aligned}$$

CD

- Contextual Decomposition

Disambiguating interactions between gates

$$\begin{aligned}i_t \odot g_t &= (L_\sigma(W_i x_t) + L_\sigma(V_i h_{t-1}) + L_\sigma(b_i)) \odot (L_{\tanh}(W_g x_t) + L_{\tanh}(V_g h_{t-1}) + L_{\tanh}(b_g)) \\ &= (L_\sigma(W_i x_t) + L_\sigma(V_i \beta_{t-1}) + L_\sigma(V_i \gamma_{t-1}) + L_\sigma(b_i)) \odot (L_{\tanh}(W_g x_t) + L_{\tanh}(V_g \beta_{t-1}) \\ &\quad + L_{\tanh}(V_g \gamma_{t-1}) + L_{\tanh}(b_g))\end{aligned}$$

Cross-terms:

- solely from the phrase, e.g., $L_\sigma(V_i \beta_{t-1}) \odot L_{\tanh}(V_g \beta_{t-1})$
- from some interaction between the phrase and other factors, e.g., $L_\sigma(V_i \beta_{t-1}) \odot L_{\tanh}(V_g \gamma_{t-1})$
- purely from other factors, e.g., $L_\sigma(b_i) \odot L_{\tanh}(V_g \gamma_{t-1})$

CD

- Contextual Decomposition

Disambiguating interactions between gates

$$\begin{aligned}i_t \odot g_t &= (L_\sigma(W_i x_t) + L_\sigma(V_i h_{t-1}) + L_\sigma(b_i)) \odot (L_{\tanh}(W_g x_t) + L_{\tanh}(V_g h_{t-1}) + L_{\tanh}(b_g)) \\ &= (L_\sigma(W_i x_t) + L_\sigma(V_i \beta_{t-1}) + L_\sigma(V_i \gamma_{t-1}) + L_\sigma(b_i)) \odot (L_{\tanh}(W_g x_t) + L_{\tanh}(V_g \beta_{t-1}) \\ &\quad + L_{\tanh}(V_g \gamma_{t-1}) + L_{\tanh}(b_g))\end{aligned}$$

Cross-terms:

- solely from the phrase, e.g., $L_\sigma(V_i \beta_{t-1}) \odot L_{\tanh}(V_g \beta_{t-1})$ β_t^u
- from some interaction between the phrase and other factors, e.g., $L_\sigma(V_i \beta_{t-1}) \odot L_{\tanh}(V_g \gamma_{t-1})$
- purely from other factors, e.g., $L_\sigma(b_i) \odot L_{\tanh}(V_g \gamma_{t-1})$

CD

- Contextual Decomposition

Disambiguating interactions between gates

$$\begin{aligned} i_t \odot g_t &= (L_\sigma(W_i x_t) + L_\sigma(V_i h_{t-1}) + L_\sigma(b_i)) \odot (L_{\tanh}(W_g x_t) + L_{\tanh}(V_g h_{t-1}) + L_{\tanh}(b_g)) \\ &= (L_\sigma(W_i x_t) + L_\sigma(V_i \beta_{t-1}) + L_\sigma(V_i \gamma_{t-1}) + L_\sigma(b_i)) \odot (L_{\tanh}(W_g x_t) + L_{\tanh}(V_g \beta_{t-1}) \\ &\quad + L_{\tanh}(V_g \gamma_{t-1}) + L_{\tanh}(b_g)) \end{aligned}$$

Cross-terms:

- solely from the phrase, e.g., $L_\sigma(V_i \beta_{t-1}) \odot L_{\tanh}(V_g \beta_{t-1})$
- from some interaction between the phrase and other factors, e.g., $L_\sigma(V_i \beta_{t-1}) \odot L_{\tanh}(V_g \gamma_{t-1})$
- purely from other factors, e.g., $L_\sigma(b_i) \odot L_{\tanh}(V_g \gamma_{t-1})$ γ_t^u

CD

- Contextual Decomposition

Disambiguating interactions between gates

$$i_t \odot g_t = \beta_t^u + \gamma_t^u$$

$$f_t \odot c_{t-1} = \beta_t^f + \gamma_t^f$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

$$= \beta_t^u + \gamma_t^u + \beta_t^f + \gamma_t^f$$

$$= \beta_t^c + \gamma_t^c$$

CD

- Contextual Decomposition

Disambiguating interactions between gates

$$i_t \odot g_t = \beta_t^u + \gamma_t^u$$

$$f_t \odot c_{t-1} = \beta_t^f + \gamma_t^f$$

$$\begin{aligned} c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\ &= \beta_t^u + \gamma_t^u + \beta_t^f + \gamma_t^f \\ &= \beta_t^c + \gamma_t^c \end{aligned}$$

$$h_t = o_t \odot \tanh(c_t)$$

$$= o_t \odot \tanh(\beta_t^c + \gamma_t^c)$$

$$= o_t \odot (L_{\tanh}(\beta_t^c) + L_{\tanh}(\gamma_t^c))$$

$$= o_t \odot L_{\tanh}(\beta_t^c) + o_t \odot L_{\tanh}(\gamma_t^c)$$

$$= \beta_t + \gamma_t$$

CD

- Contextual Decomposition

Disambiguating interactions between gates

$$i_t \odot g_t = \beta_t^u + \gamma_t^u$$

$$f_t \odot c_{t-1} = \beta_t^f + \gamma_t^f$$

$$\begin{aligned} c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\ &= \beta_t^u + \gamma_t^u + \beta_t^f + \gamma_t^f \\ &= \beta_t^c + \gamma_t^c \end{aligned}$$

$$h_t = o_t \odot \tanh(c_t)$$

$$= o_t \odot \tanh(\beta_t^c + \gamma_t^c)$$

$$= o_t \odot (L_{\tanh}(\beta_t^c) + L_{\tanh}(\gamma_t^c))$$

$$= o_t \odot L_{\tanh}(\beta_t^c) + o_t \odot L_{\tanh}(\gamma_t^c)$$

$$= \beta_t + \gamma_t$$

Iteratively decomposing until we get

$$h_T = \beta_T + \gamma_T$$

$$\beta_0 = \gamma_0 = 0$$

$$\beta_t: x_t (q \leq t \leq r) \quad \gamma_t: x_t (t > r, t < q)$$

CD

- Contextual Decomposition

Disambiguating interactions between gates

$$i_t \odot g_t = \beta_t^u + \gamma_t^u$$

$$f_t \odot c_{t-1} = \beta_t^f + \gamma_t^f$$

$$\begin{aligned} c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\ &= \beta_t^u + \gamma_t^u + \beta_t^f + \gamma_t^f \\ &= \beta_t^c + \gamma_t^c \end{aligned}$$

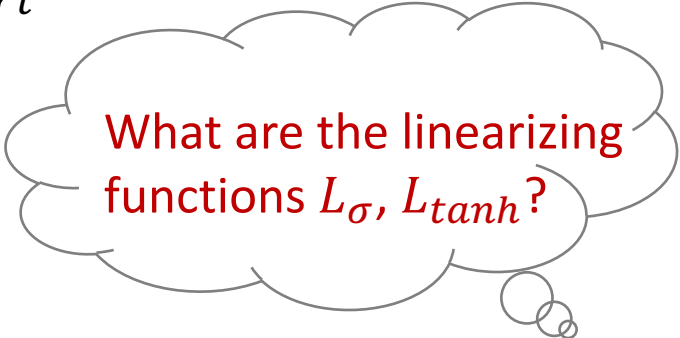
$$h_t = o_t \odot \tanh(c_t)$$

$$= o_t \odot \tanh(\beta_t^c + \gamma_t^c)$$

$$= o_t \odot (L_{\tanh}(\beta_t^c) + L_{\tanh}(\gamma_t^c))$$

$$= o_t \odot L_{\tanh}(\beta_t^c) + o_t \odot L_{\tanh}(\gamma_t^c)$$

$$= \beta_t + \gamma_t$$



What are the linearizing functions L_σ, L_{\tanh} ?

CD

- Contextual Decomposition

Linearizing activation functions (L_σ, L_{\tanh})

$$\tanh\left(\sum_{i=1}^N y_i\right) = \sum_{i=1}^N L_{\tanh}(y_i) \quad (N \leq 4)$$

CD

- Contextual Decomposition

Linearizing activation functions (L_σ, L_{\tanh})

$$\tanh\left(\sum_{i=1}^N y_i\right) = \sum_{i=1}^N L_{\tanh}(y_i) \quad (N \leq 4)$$

Telescoping sum (given a natural ordering to $\{y_i\}$)

$$L_{\tanh}(y_k) = \tanh\left(\sum_{j=1}^k y_j\right) - \tanh\left(\sum_{j=1}^{k-1} y_j\right)$$

CD

- Contextual Decomposition

Linearizing activation functions (L_σ, L_{\tanh})

$$\tanh\left(\sum_{i=1}^N y_i\right) = \sum_{i=1}^N L_{\tanh}(y_i) \quad (N \leq 4)$$

Telescoping sum (given a natural ordering to $\{y_i\}$)

$$L_{\tanh}(y_k) = \tanh\left(\sum_{j=1}^k y_j\right) - \tanh\left(\sum_{j=1}^{k-1} y_j\right)$$

$$\sum_{i=1}^N L_{\tanh}(y_i) = \tanh\left(\sum_{j=1}^N y_j\right) - \tanh\left(\sum_{j=1}^{N-1} y_j\right) + \tanh\left(\sum_{j=1}^{N-1} y_j\right) - \tanh\left(\sum_{j=1}^{N-2} y_j\right) + \dots + \tanh\left(\sum_{j=1}^2 y_j\right) - \tanh\left(\sum_{j=1}^1 y_j\right) + \tanh\left(\sum_{j=1}^1 y_j\right)$$

CD

- Contextual Decomposition

Linearizing activation functions (L_σ, L_{\tanh})

$$\tanh\left(\sum_{i=1}^N y_i\right) = \sum_{i=1}^N L_{\tanh}(y_i) \quad (N \leq 4)$$

Telescoping sum (given a natural ordering to $\{y_i\}$)

$$L_{\tanh}(y_k) = \tanh\left(\sum_{j=1}^k y_j\right) - \tanh\left(\sum_{j=1}^{k-1} y_j\right)$$

$$\sum_{i=1}^N L_{\tanh}(y_i) = \tanh\left(\sum_{j=1}^N y_j\right) - \cancel{\tanh\left(\sum_{j=1}^{N-1} y_j\right)} + \cancel{\tanh\left(\sum_{j=1}^{N-1} y_j\right)} - \cancel{\tanh\left(\sum_{j=1}^{N-2} y_j\right)} + \dots + \cancel{\tanh\left(\sum_{j=1}^2 y_j\right)} - \cancel{\tanh\left(\sum_{j=1}^1 y_j\right)} + \cancel{\tanh\left(\sum_{j=1}^1 y_j\right)}$$

CD

- Contextual Decomposition

Linearizing activation functions (L_σ, L_{\tanh})

$$\tanh\left(\sum_{i=1}^N y_i\right) = \sum_{i=1}^N L_{\tanh}(y_i) \quad (N \leq 4)$$

Telescoping sum (given a natural ordering to $\{y_i\}$)

$$L_{\tanh}(y_k) = \tanh\left(\sum_{j=1}^k y_j\right) - \tanh\left(\sum_{j=1}^{k-1} y_j\right)$$

$\{\beta_{t-1}, \gamma_{t-1}, x_t\}$ have
no clear ordering

CD

- Contextual Decomposition

Linearizing activation functions (L_σ, L_{\tanh})


$$\tanh\left(\sum_{i=1}^N y_i\right) = \sum_{i=1}^N L_{\tanh}(y_i) \quad (N \leq 4)$$

Telescoping sum (given a natural ordering to $\{y_i\}$)

$$L_{\tanh}(y_k) = \tanh\left(\sum_{j=1}^k y_j\right) - \tanh\left(\sum_{j=1}^{k-1} y_j\right)$$

All permutations: π_1, \dots, π_{M_N}
 $\pi_i^{-1}(k)$: the position of y_k in π_i

Average over
all orderings


$$L_{\tanh}(y_k) = \frac{1}{M_N} \sum_{i=1}^{M_N} \left[\tanh\left(\sum_{j=1}^{\pi_i^{-1}(k)} y_{\pi_i(j)}\right) - \tanh\left(\sum_{j=1}^{\pi_i^{-1}(k)-1} y_{\pi_i(j)}\right) \right]$$

Summary

- **LSTM**

Input word embeddings: $x_1, \dots, x_T \in \mathbb{R}^{d_1}$

$$\begin{aligned}o_t &= \sigma(W_o x_t + V_o h_{t-1} + b_o) \\f_t &= \sigma(W_f x_t + V_f h_{t-1} + b_f) \\i_t &= \sigma(W_i x_t + V_i h_{t-1} + b_i) \\g_t &= \tanh(W_g x_t + V_g h_{t-1} + b_g) \\c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\h_t &= o_t \odot \tanh(c_t)\end{aligned}$$

- **Contribution of an arbitrary phrase:**

x_q, \dots, x_r ($1 \leq q \leq r \leq T$)

$$\begin{aligned}\sigma(\cdot) &= \sum L_\sigma, \tanh(\cdot) = \sum L_{\tanh} \\h_t &= \beta_t + \gamma_t\end{aligned}$$

$t = 1, \dots, T$



$$p = \text{softmax}(W \beta_T + W \gamma_T)$$

the phrase's contribution
to the LSTM's prediction

Question?

Visualizations

Text

“used to be my favorite” (negative)

“not worth the time” (negative)

Attribution Method	Heat Map									
Gradient	used	to	be	my	favorite	not	worth	the	time	
Leave One Out (Li et al., 2016)	used	to	be	my	favorite	not	worth	the	time	
Cell decomposition (Murdoch & Szlam, 2017)	used	to	be	my	favorite	not	worth	the	time	
Integrated gradients (Sundararajan et al., 2017)	used	to	be	my	favorite	not	worth	the	time	
Contextual decomposition	used	to	be	my	favorite	not	worth	the	time	

Legend: Very Negative (red), Negative (orange), Neutral (yellow), Positive (green), Very Positive (blue)

Visualizations

The first phrase is **positive**, but the second one is **negative**

CD is the only method that accurately captures this dynamic

Attribution Method	Heat Map
Gradient	It's easy to love Robin Tunney – she's pretty and she can act – but it gets harder and harder to understand her choices.
Leave one out (Li et al., 2016)	It's easy to love Robin Tunney – she's pretty and she can act – but it gets harder and harder to understand her choices.
Cell decomposition (Murdoch & Szlam, 2017)	It's easy to love Robin Tunney – she's pretty and she can act – but it gets harder and harder to understand her choices.
Integrated gradients (Sundararajan et al., 2017)	It's easy to love Robin Tunney – she's pretty and she can act – but it gets harder and harder to understand her choices.
Contextual decomposition	It's easy to love Robin Tunney – she's pretty and she can act – but it gets harder and harder to understand her choices.

Legend Very Negative Negative Neutral Positive Very Positive

Discussion

- CD is model-dependent
- Decomposing complex DNN (e.g., transformer) is not trivial

Beyond Feature Attribution

- Contextual Decomposition (CD)
- Hierarchical Explanation via Divisive Generation (HEDGE)

Generating Hierarchical Explanations on Text Classification via Feature Interaction Detection

Hanjie Chen, Guangtao Zheng, Yangfeng Ji

(ACL, 2020)

HEDGE

Why we need hierarchical explanations?

Prediction: Negative

LIME Explanation [Ribeiro et al., 2016]

a waste of good performance

CD Explanation [Murdoch et al., 2018]

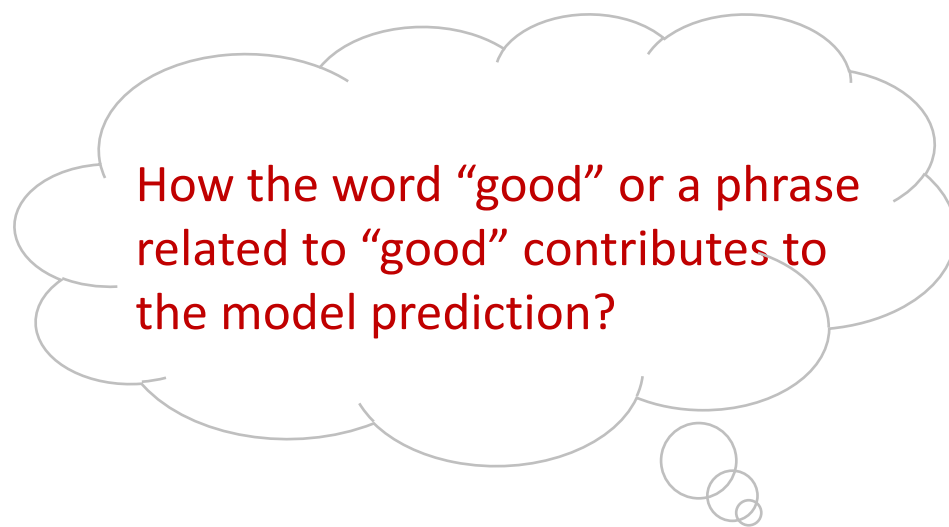
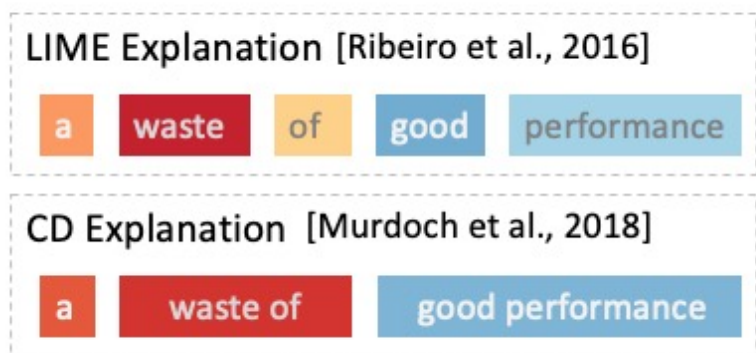
a waste of good performance



HEDGE

Why we need hierarchical explanations?

Prediction: Negative



HEDGE

Why we need hierarchical explanations?

Prediction: Negative

LIME Explanation [Ribeiro et al., 2016]

a waste of good performance

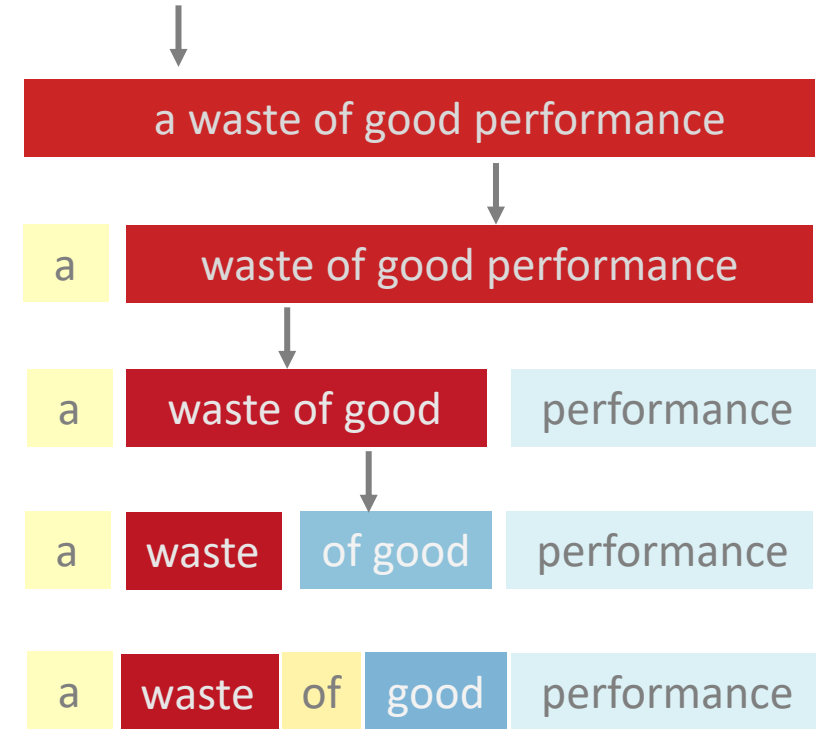
CD Explanation [Murdoch et al., 2018]

a waste of good performance



Hierarchical explanation via divisive generation (HEDGE)

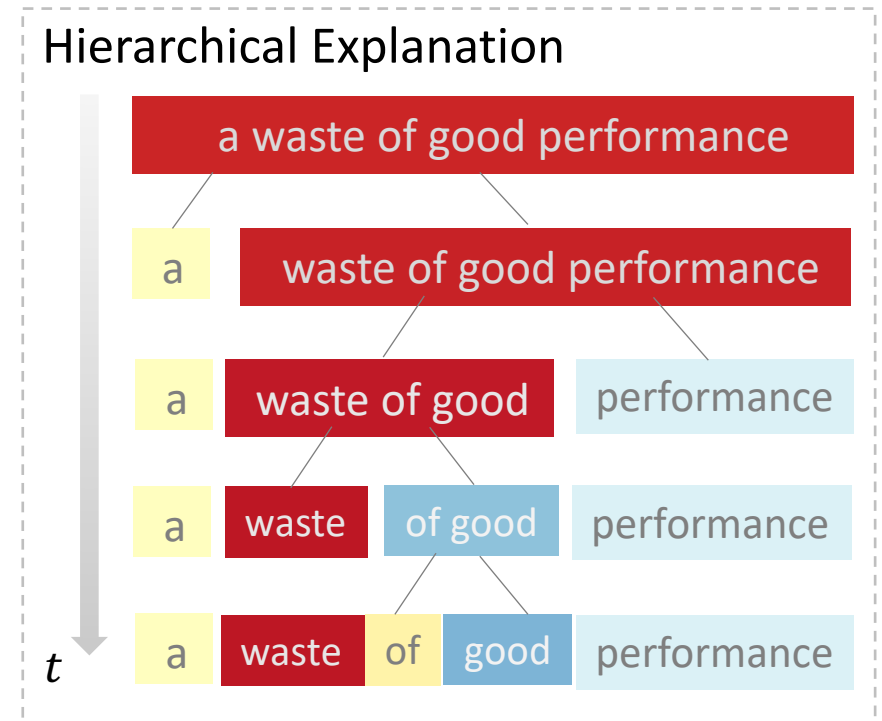
- Where is the dividing point?
- Which text segment should be split?
- How to quantify feature importance?



HEDGE

- Definition

- A text with n words: $\mathbf{x} = (x_1, \dots, x_n)$
- A text span: $\mathbf{x}_{(s_i, s_{i+1}]} = (x_{s_i+1}, \dots, x_{s_{i+1}})$
- A partition: $\mathcal{P} = \{\mathbf{x}_{(0, s_1]}, \mathbf{x}_{(s_1, s_2]}, \dots, \mathbf{x}_{(s_{P-1}, n]}\}$
- Interaction score: $\phi(\cdot, \cdot)$
- Importance score: $\psi(\cdot)$

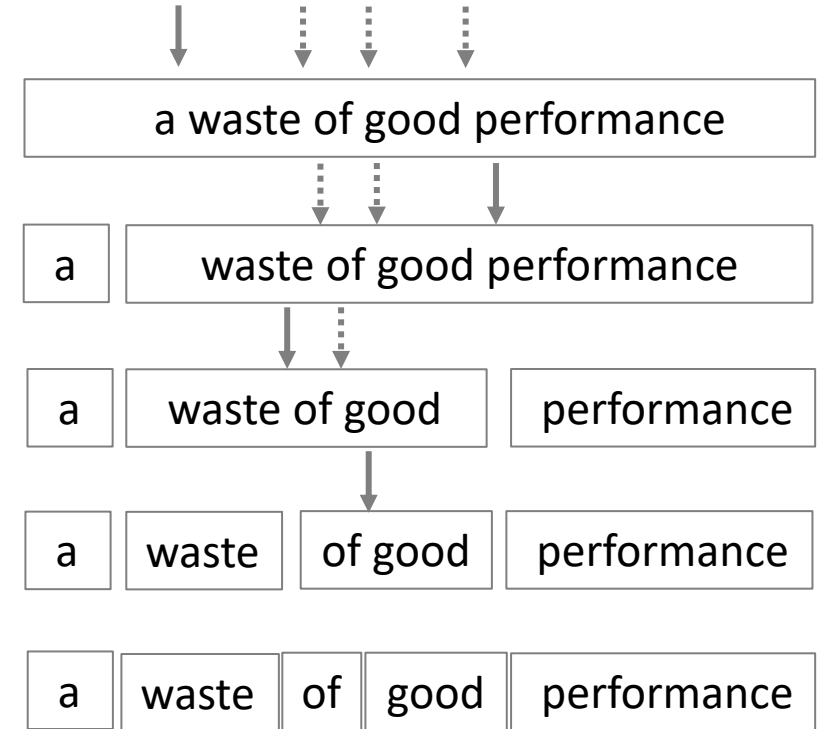


HEDGE

- Where is the dividing point?

The *local weakest* interaction point:

$$\min_{j \in (s_i, s_{i+1})} \phi(\mathbf{x}_{(s_i, j]}, \mathbf{x}_{(j, s_{i+1}]} | \mathcal{P})$$



HEDGE

- Where is the dividing point?

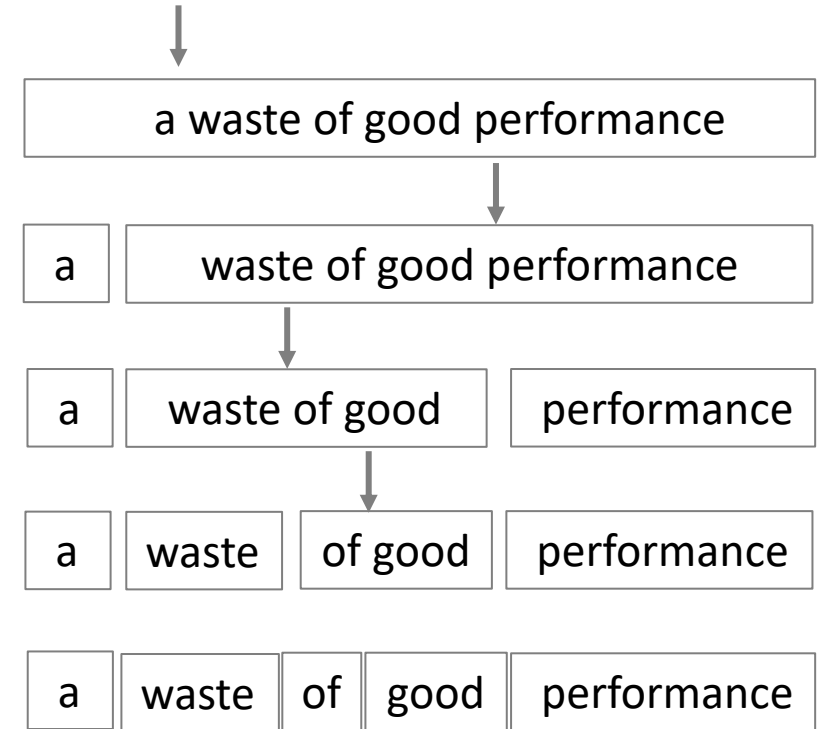
The *local weakest* interaction point:

$$\min_{j \in (s_i, s_{i+1})} \phi(\mathbf{x}_{(s_i, j]}, \mathbf{x}_{(j, s_{i+1}]} | \mathcal{P})$$

- Which text segment should be split?

The *global weakest* interaction point:

$$\min_{\mathbf{x}_{(s_i, s_{i+1}]} \in \mathcal{P}} \min_{j \in (s_i, s_{i+1})} \phi(\mathbf{x}_{(s_i, j]}, \mathbf{x}_{(j, s_{i+1}]} | \mathcal{P})$$



HEDGE

- Where is the dividing point?

The *local weakest* interaction point:

$$\min_{j \in (s_i, s_{i+1})} \phi(\mathbf{x}_{(s_i, j]}, \mathbf{x}_{(j, s_{i+1}]} | \mathcal{P})$$

- Which text segment should be split?

The *global weakest* interaction point:

$$\min_{\mathbf{x}_{(s_i, s_{i+1}]} \in \mathcal{P}} \min_{j \in (s_i, s_{i+1})} \phi(\mathbf{x}_{(s_i, j]}, \mathbf{x}_{(j, s_{i+1}]} | \mathcal{P})$$

- How to quantify feature importance?

Feature importance score: $\psi(\cdot)$



HEDGE

- Where is the dividing point?

The *local weakest* interaction point:

$$\min_{j \in (s_i, s_{i+1})} \phi(\mathbf{x}_{(s_i, j]}, \mathbf{x}_{(j, s_{i+1}]} | \mathcal{P})$$

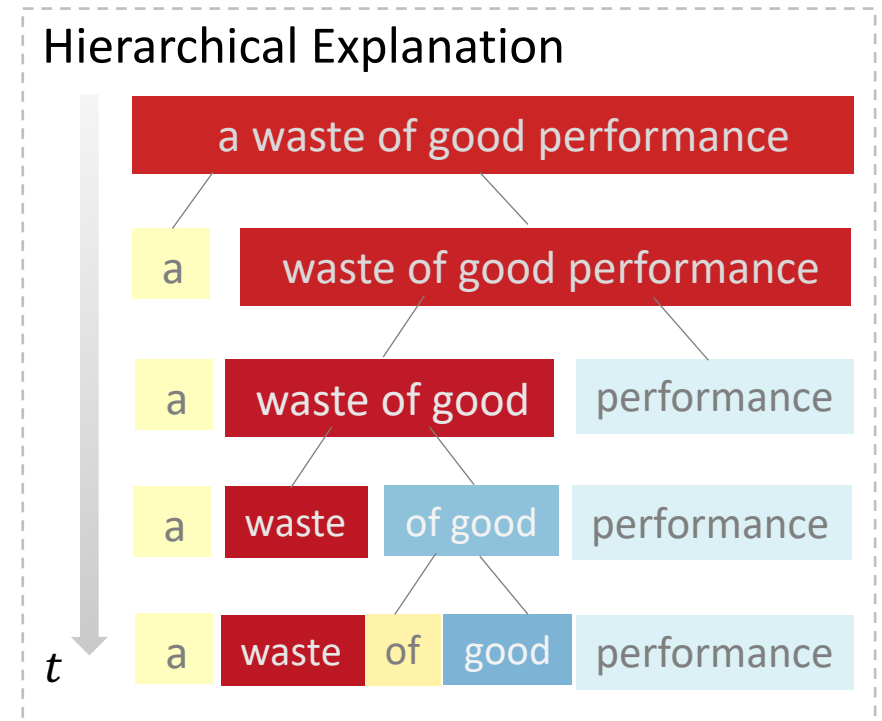
- Which text segment should be split?

The *global weakest* interaction point:

$$\min_{\mathbf{x}_{(s_i, s_{i+1}]} \in \mathcal{P}} \min_{j \in (s_i, s_{i+1})} \phi(\mathbf{x}_{(s_i, j]}, \mathbf{x}_{(j, s_{i+1}]} | \mathcal{P})$$

- How to quantify feature importance?

Feature importance score: $\psi(\cdot)$



HEDGE

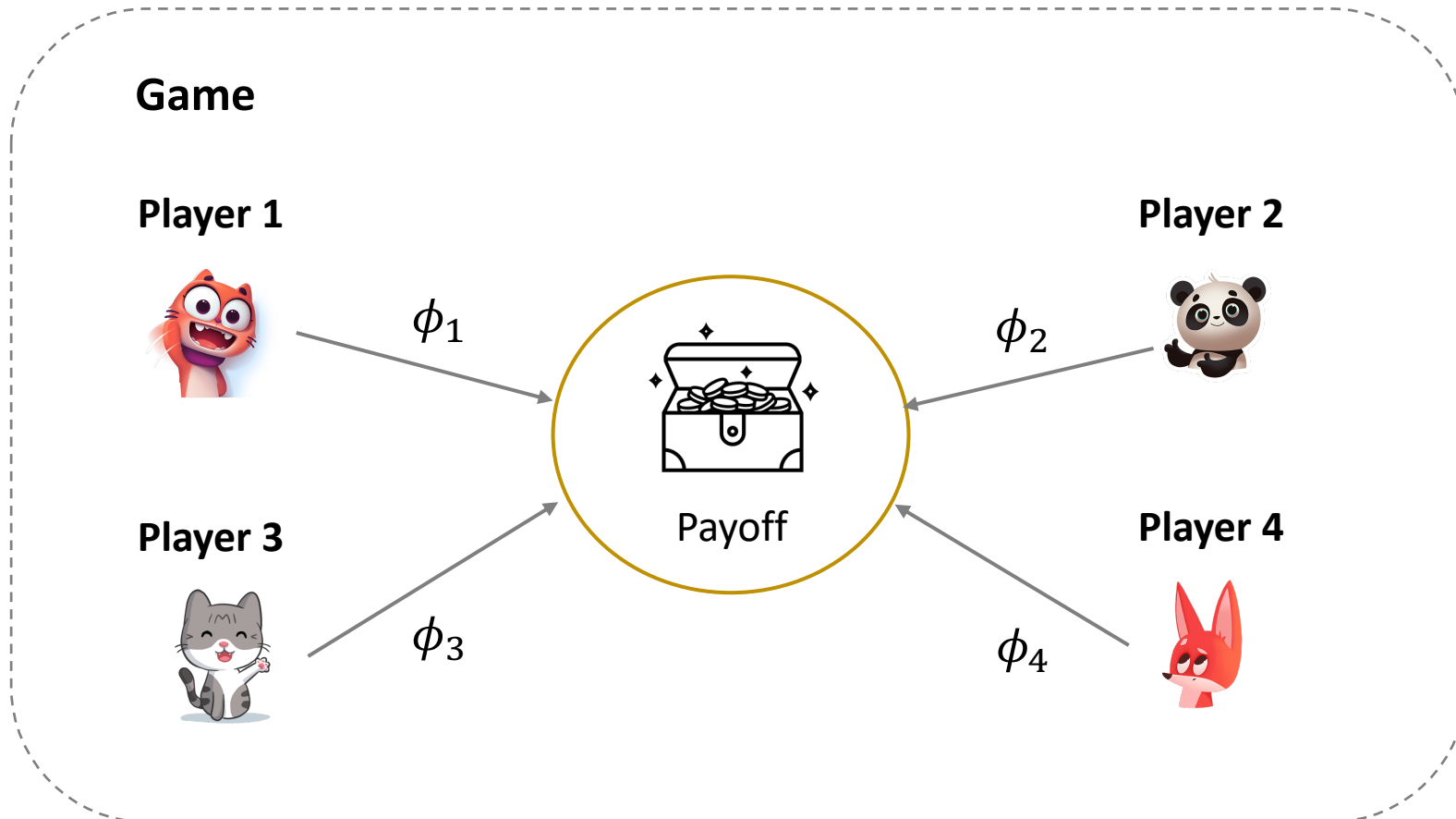
- Feature interaction score $\phi(\cdot, \cdot)$

Calculate the interaction between j_1 and j_2 via Shapley interaction index

[Fujimoto et al., 2006, Lundberg et al., 2018]

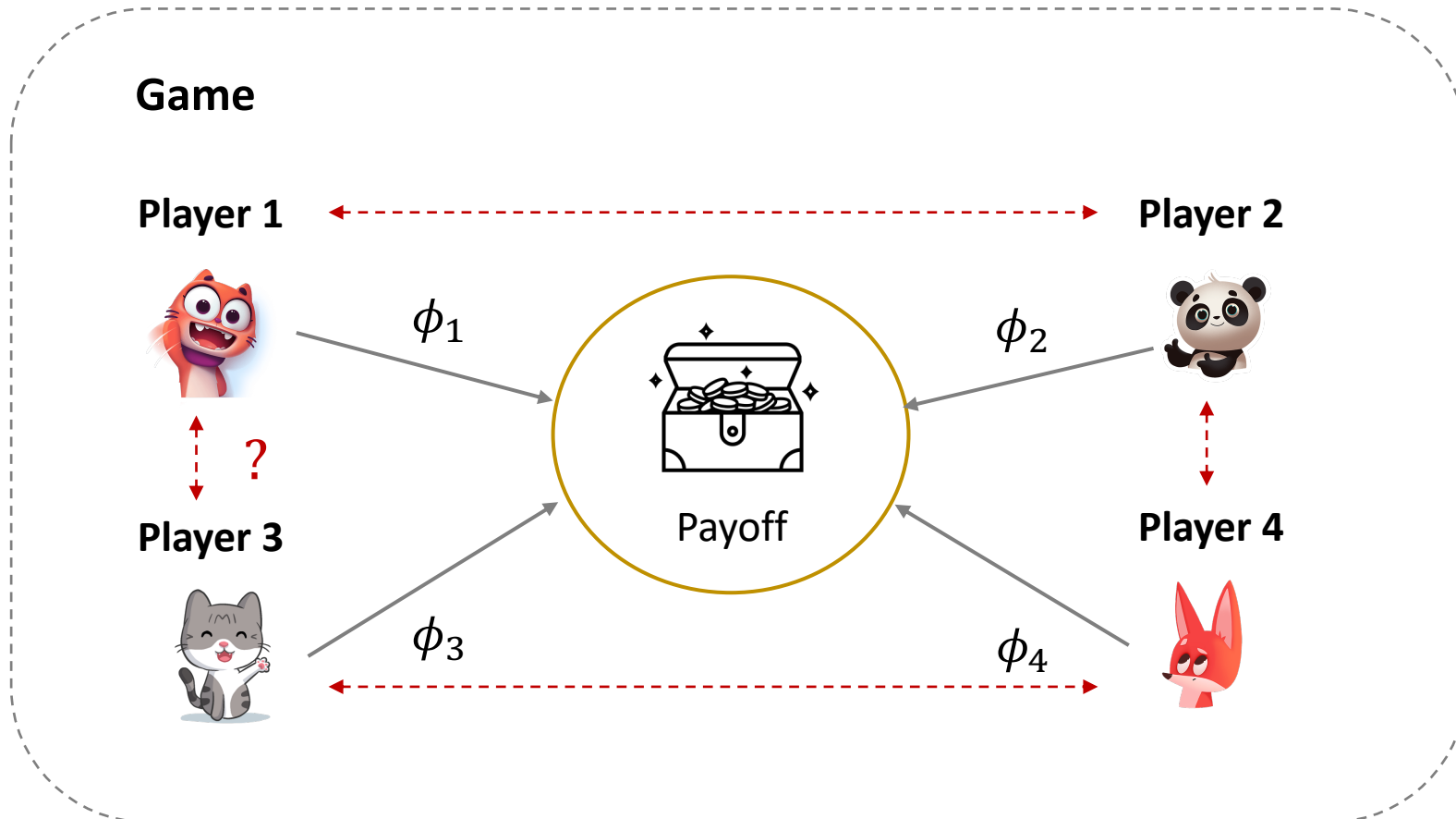
Shapley

Quantifying the contribution of each player



Shapley Interaction Index

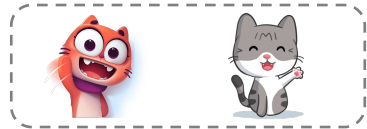
Quantifying the interaction between players



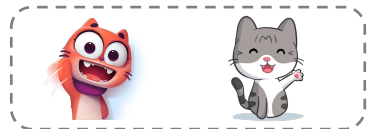
Shapley Interaction Index

Coalitions

Payoff



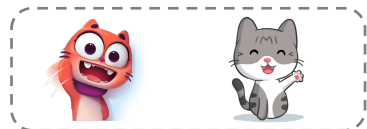
P_1



P_2



P_3



\emptyset

P_4

Shapley Interaction Index

Coalitions

Payoff



P_1

P_1'



P_2

P_2'



P_3

P_3'





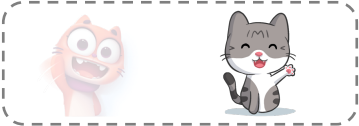




\emptyset

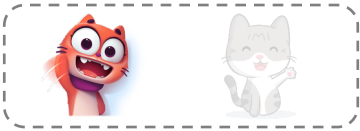

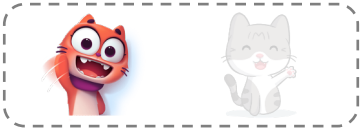

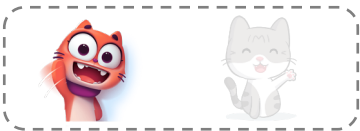

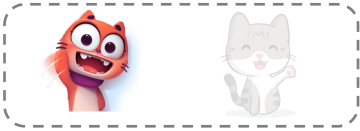
P_4

P_4'








Shapley Interaction Index

Coalitions		Payoff	Marginal contribution with player 3
		$P_1 - P_1'$	ΔP_1
		$P_2 - P_2'$	ΔP_2
		$P_3 - P_3'$	ΔP_3
	\emptyset	$P_4 - P_4'$	ΔP_4








Shapley Interaction Index

Coalitions		Payoff	Marginal contribution with player 3	Payoff
		$P_1 - P_1'$	ΔP_1	Q_1
		$P_2 - P_2'$	ΔP_2	Q_2
		$P_3 - P_3'$	ΔP_3	Q_3
	\emptyset	$P_4 - P_4'$	ΔP_4	Q_4

Shapley Interaction Index

Coalitions		Payoff	Marginal contribution with player 3	Payoff	Marginal contribution without player 3
		$P_1 - P_1'$	ΔP_1	$Q_1 - Q_1'$	ΔQ_1
		$P_2 - P_2'$	ΔP_2	$Q_2 - Q_2'$	ΔQ_2
		$P_3 - P_3'$	ΔP_3	$Q_3 - Q_3'$	ΔQ_3
	\emptyset	$P_4 - P_4'$	ΔP_4	$Q_4 - Q_4'$	ΔQ_4

Shapley Interaction Index

Coalitions		Payoff	Marginal contribution with player 3	Payoff	Marginal contribution without player 3
		$P_1 - P_1'$	ΔP_1	$Q_1 - Q_1'$	ΔQ_1
		$P_2 - P_2'$	ΔP_2	$Q_2 - Q_2'$	ΔQ_2
		$P_3 - P_3'$	ΔP_3	$Q_3 - Q_3'$	ΔQ_3
	\emptyset	$P_4 - P_4'$	ΔP_4	$Q_4 - Q_4'$	ΔQ_4

$$\phi_{1,3} = \sum \Delta P_i - \Delta Q_i$$

Shapley Interaction Index

Coalitions		Payoff	Marginal contribution with player 3	Payoff	Marginal contribution without player 3
		$P_1 - P_1'$	ΔP_1	$Q_1 - Q_1'$	ΔQ_1
		$P_2 - P_2'$	ΔP_2	$Q_2 - Q_2'$	ΔQ_2
		$P_3 - P_3'$	ΔP_3	$Q_3 - Q_3'$	ΔQ_3
	\emptyset	$P_4 - P_4'$	ΔP_4	$Q_4 - Q_4'$	ΔQ_4

$$\phi_{1,3} = \sum \Delta P_i - \Delta Q_i$$

$$\phi_{3,1} = \phi_{1,3}$$



interaction

$$\phi_{3,1} + \phi_{1,3}$$

HEDGE

- Feature interaction score $\phi(\cdot, \cdot)$

Calculate the interaction between j_1 and j_2 via Shapley interaction index

[Fujimoto et al., 2006, Lundberg et al., 2018]

$$\phi(j_1, j_2 | \mathcal{P}) = \sum_{S \subseteq \mathcal{N} \setminus \{j_1, j_2\}} \frac{|S|! (P - |S| - 1)!}{P!} \gamma(j_1, j_2, S)$$

$$\gamma(j_1, j_2, S) = \underbrace{\mathbb{E}[f(\mathbf{x}') | S \cup \{j_1, j_2\}]} - \mathbb{E}[f(\mathbf{x}') | S \cup \{j_2\}] - \underbrace{(\mathbb{E}[f(\mathbf{x}') | S \cup \{j_1\}] - \mathbb{E}[f(\mathbf{x}') | S])}$$

The influence of j_1 on the model
output with j_2 considered

without j_2 considered

HEDGE

- Feature interaction score $\phi(\cdot, \cdot)$

Calculate the interaction between j_1 and j_2 via Shapley interaction index

[Fujimoto et al., 2006, Lundberg et al., 2018]

$$\phi(j_1, j_2 | \mathcal{P}) = \sum_{S \subseteq \mathcal{N} \setminus \{j_1, j_2\}} \frac{|S|! (P - |S| - 1)!}{P!} \gamma(j_1, j_2, S)$$

$$\gamma(j_1, j_2, S) = \underbrace{\mathbb{E}[f(\mathbf{x}') | S \cup \{j_1, j_2\}]}_{\text{The influence of } j_1 \text{ on the model output with } j_2 \text{ considered}} - \underbrace{\mathbb{E}[f(\mathbf{x}') | S \cup \{j_2\}]}_{\text{without } j_2 \text{ considered}} - (\mathbb{E}[f(\mathbf{x}') | S \cup \{j_1\}] - \mathbb{E}[f(\mathbf{x}') | S])$$

The influence of j_1 on the model output with j_2 considered

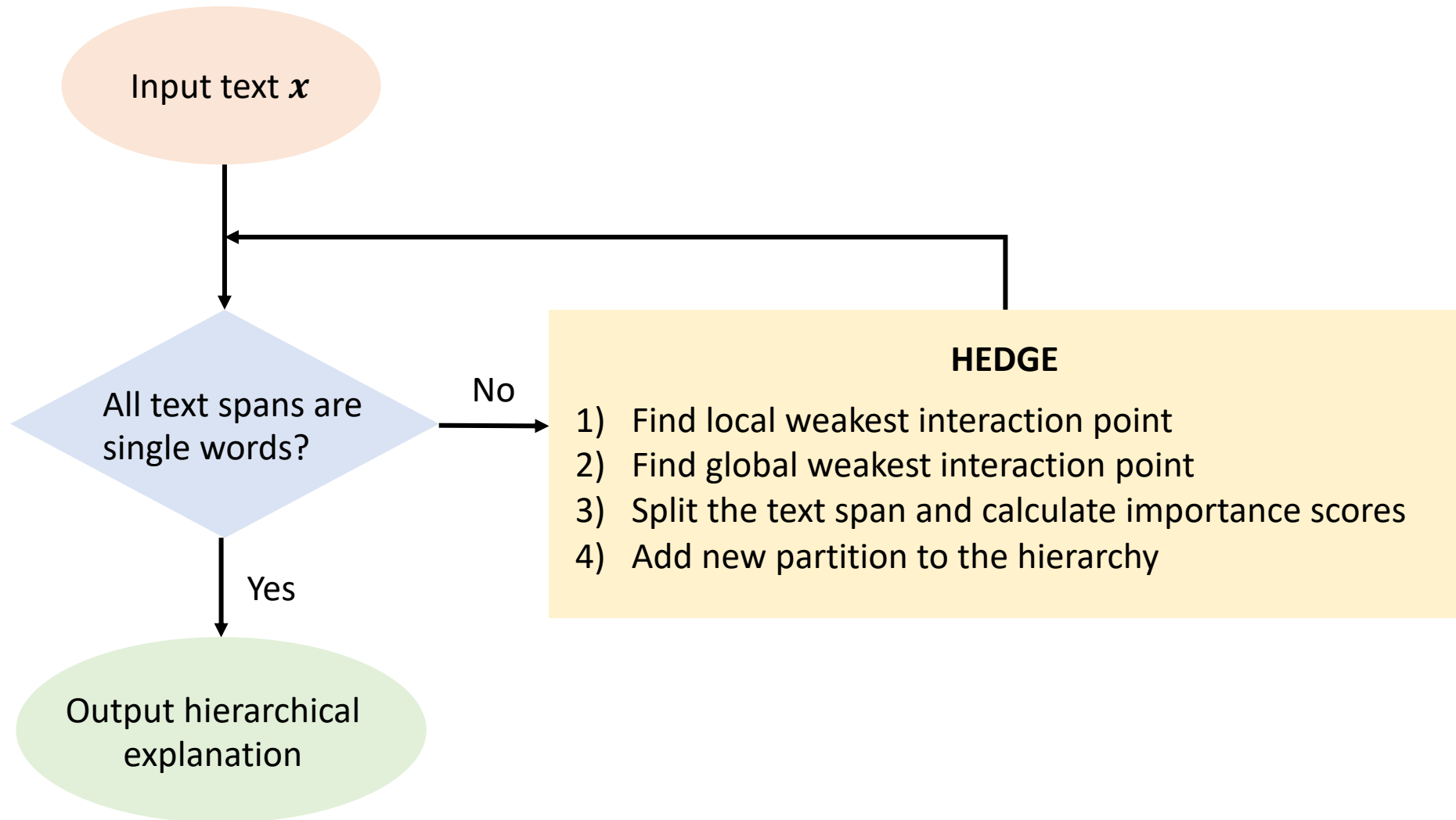
without j_2 considered

- Feature importance score $\psi(\cdot, \cdot)$

$$\psi(\mathbf{x}_{(s_i, s_{i+1})}) = \underline{f_{\hat{y}}}(\mathbf{x}_{(s_i, s_{i+1})}) - \max_{y' \neq \hat{y}, y' \in \mathcal{Y}} f_{y'}(\mathbf{x}_{(s_i, s_{i+1})})$$

Predicted label on \mathbf{x}

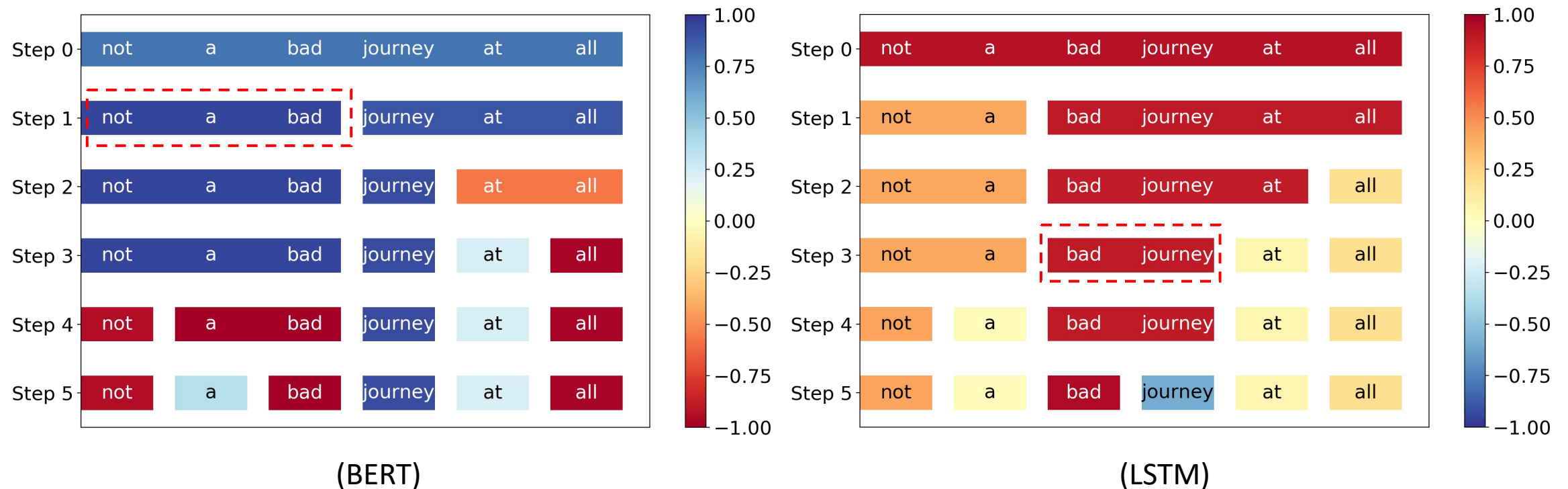
Pipeline



Question?

Qualitative Analysis

- Compare HEDGE in interpreting the LSTM and BERT model
 - BERT gives the correct prediction “positive”, while LSTM makes a wrong prediction “negative”
 - HEDGE can explain different model prediction behaviors

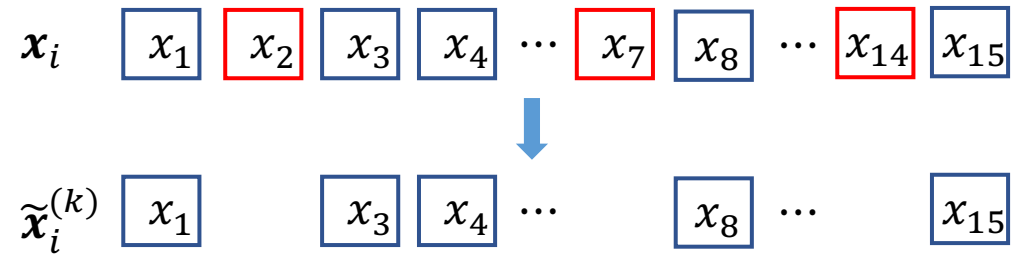


Quantitative Evaluation

- The area over the perturbation curve (AOPC) [Nguyen, 2018, Samek et al., 2016]

$$AOPC(k) = \frac{1}{N} \sum_{i=1}^N \left\{ p(\hat{y} | \mathbf{x}_i) - p(\hat{y} | \tilde{\mathbf{x}}_i^{(k)}) \right\}$$

✓ Higher AOPCs are better

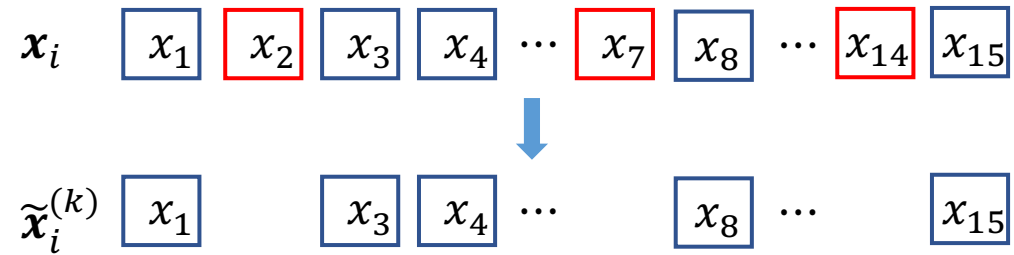


Quantitative Evaluation

- The area over the perturbation curve (AOPC) [Nguyen, 2018, Samek et al., 2016]

$$AOPC(k) = \frac{1}{N} \sum_{i=1}^N \left\{ p(\hat{y} | \mathbf{x}_i) - p(\hat{y} | \tilde{\mathbf{x}}_i^{(k)}) \right\}$$

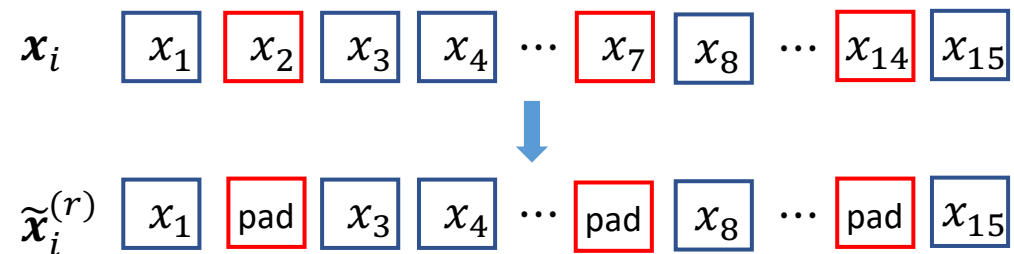
- ✓ Higher AOPCs are better



- Log-odds [Shrikumar et al., 2017, Chen et al., 2018]

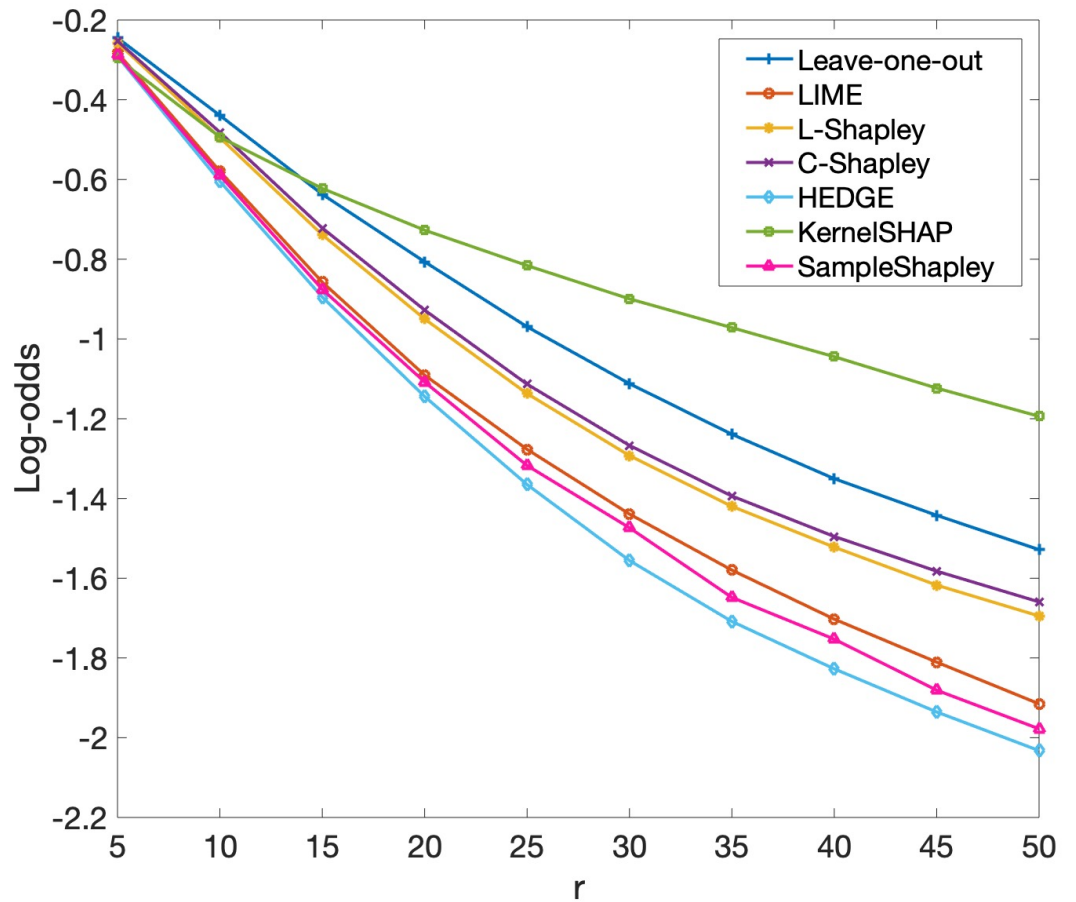
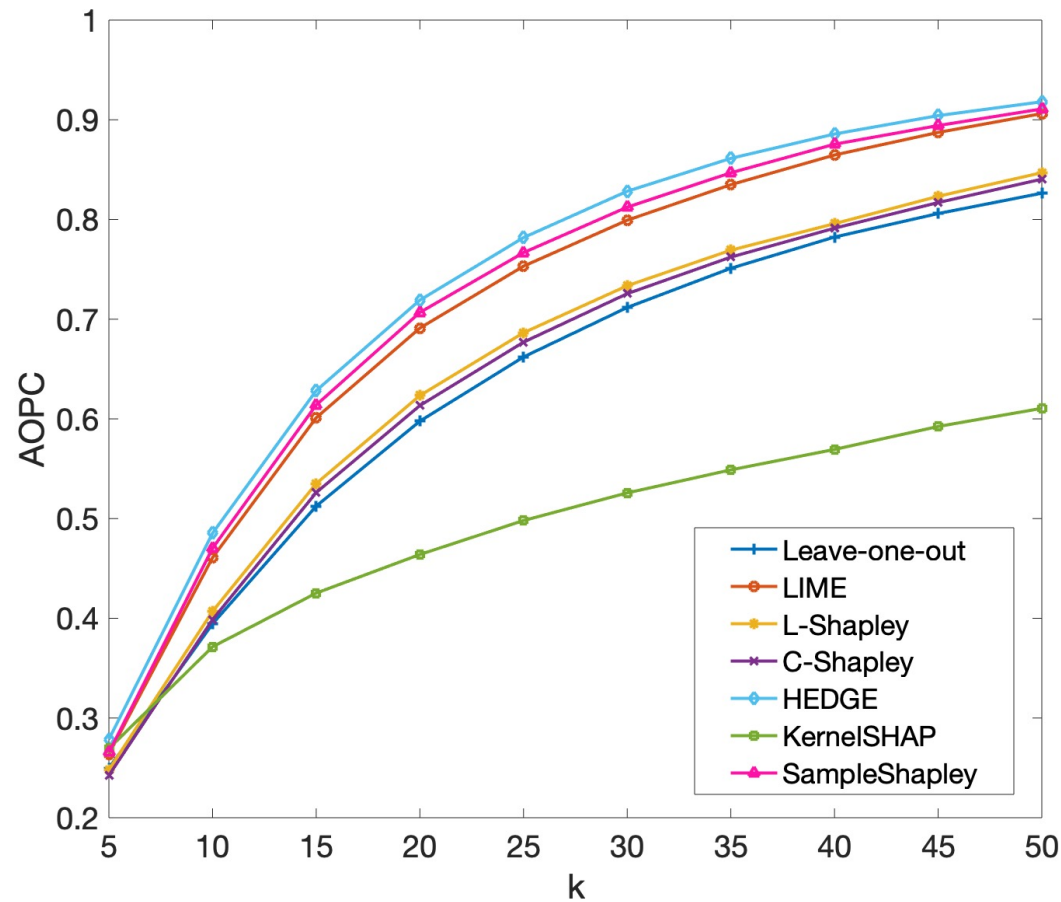
$$\text{Log-odds}(r) = \frac{1}{N} \sum_{i=1}^N \log \frac{p(\hat{y} | \tilde{\mathbf{x}}_i^{(r)})}{p(\hat{y} | \mathbf{x}_i)}$$

- ✓ Lower log-odds scores are better



Quantitative Evaluation

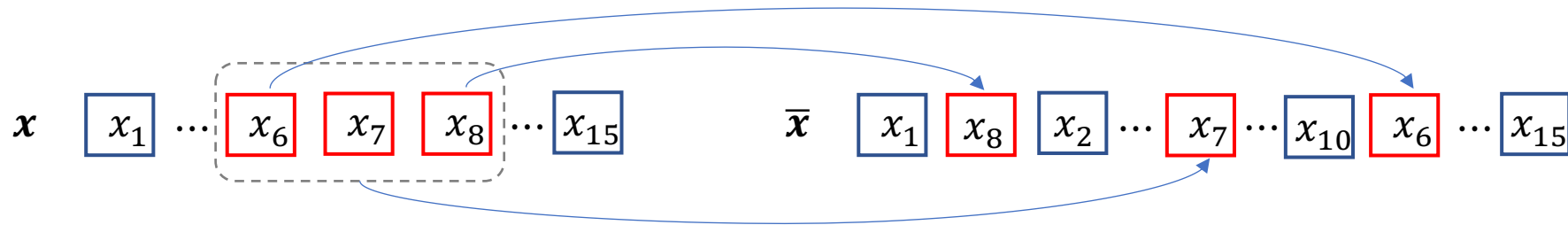
- AOPC and log-odds scores of the CNN model on the IMDB dataset
- HEDGE achieves the best performance under both evaluation metrics



Quantitative Evaluation

- Cohesion-score

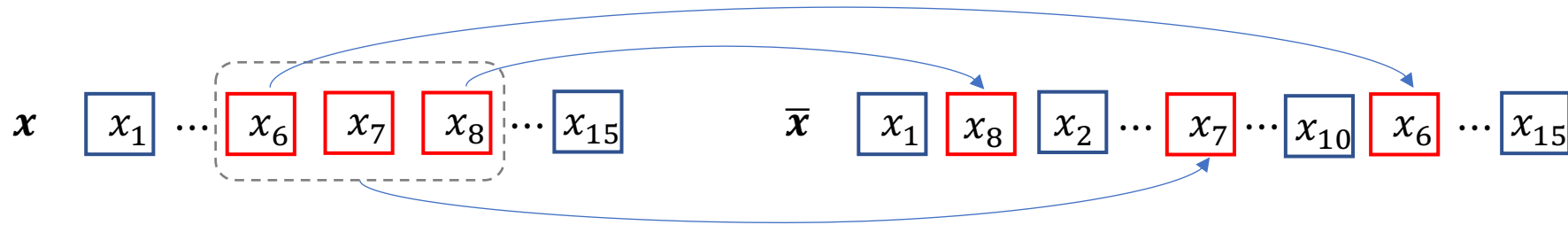
- $Cohesion - score = \frac{1}{N} \sum_{i=1}^N \frac{1}{Q} \sum_{q=1}^Q \left\{ p(\hat{y} | \mathbf{x}_i) - p(\hat{y} | \bar{\mathbf{x}}_i^{(q)}) \right\}$ ✓ Higher cohesion-scores are better



Quantitative Evaluation

- Cohesion-score

- $$\text{Cohesion - score} = \frac{1}{N} \sum_{i=1}^N \frac{1}{Q} \sum_{q=1}^Q \left\{ p(\hat{y} | \mathbf{x}_i) - p(\hat{y} | \bar{\mathbf{x}}_i^{(q)}) \right\}$$
✓ Higher cohesion-scores are better



- Results

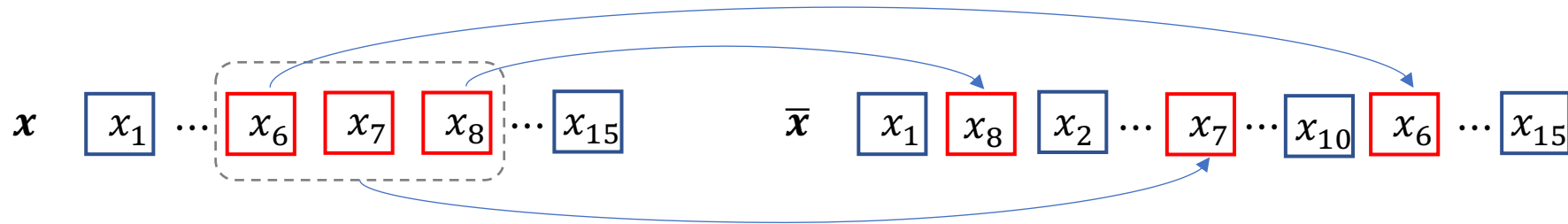
Methods	Models	Cohesion-score	
		SST	IMDB
HEDGE	CNN	0.016	0.012
	BERT	0.124	0.103
	LSTM	0.020	0.050
ACD	LSTM	0.015	0.038

- ✓ HEDGE is better at capturing feature interactions

Quantitative Evaluation

- Cohesion-score

- $$\text{Cohesion - score} = \frac{1}{N} \sum_{i=1}^N \frac{1}{Q} \sum_{q=1}^Q \left\{ p(\hat{y} | \mathbf{x}_i) - p(\hat{y} | \bar{\mathbf{x}}_i^{(q)}) \right\}$$
✓ Higher cohesion-scores are better



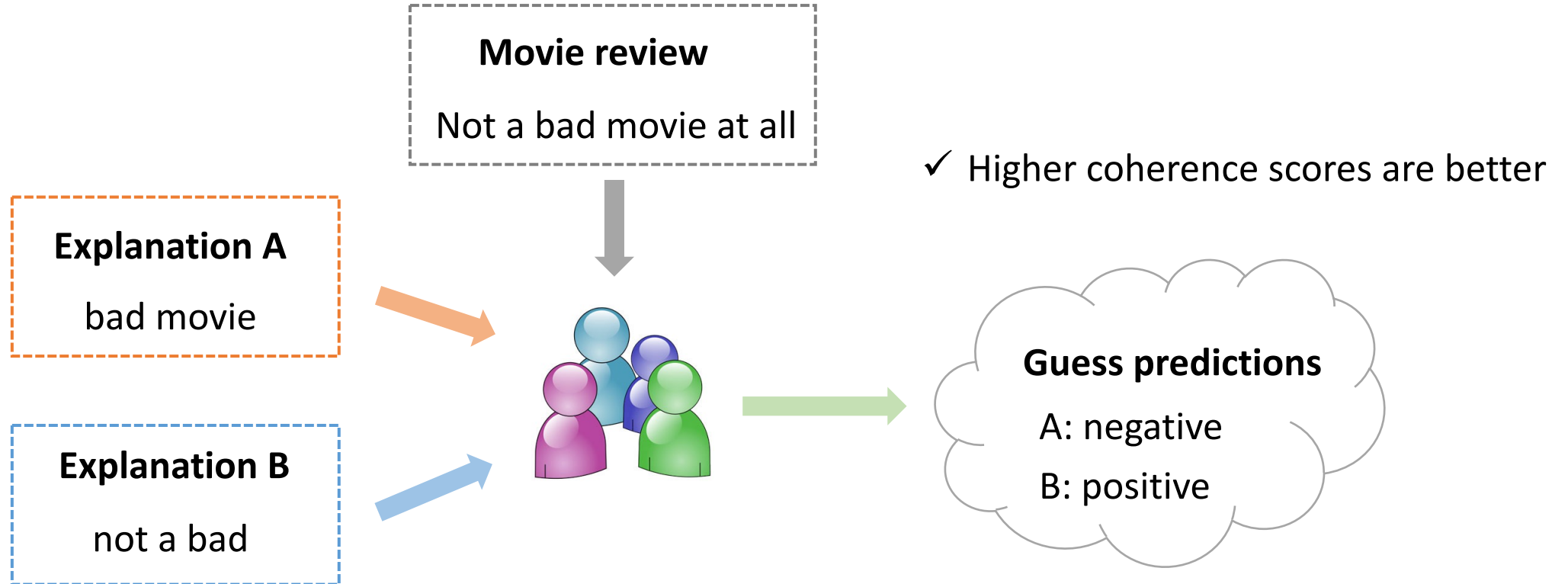
- Results

Methods	Models	Cohesion-score	
		SST	IMDB
HEDGE	CNN	0.016	0.012
	BERT	0.124	0.103
	LSTM	0.020	0.050
ACD	LSTM	0.015	0.038

✓ BERT is more sensitive to perturbations on important phrases

Human Evaluation

- Compare human annotations and model predictions



Human Evaluation

- Coherence scores of different explanation methods with LSTM model on the IMDB dataset

Methods	Coherence Score
Leave-one-out	0.82
ACD	0.68
LIME	0.85
L-Shapley	0.75
C-Shapley	0.73
KernelSHAP	0.56
SampleShapley	0.78
HEDGE	0.89

Question?

Reference

- Murdoch, W. James, Peter J. Liu, and Bin Yu. "Beyond word importance: Contextual decomposition to extract interactions from lstms." *arXiv preprint arXiv:1801.05453* (2018).
- Chen, Hanjie, Guangtao Zheng, and Yangfeng Ji. "Generating hierarchical explanations on text classification via feature interaction detection." *arXiv preprint arXiv:2004.02015* (2020).
- Yeh, Chih-Kuan, et al. "On completeness-aware concept-based explanations in deep neural networks." *Advances in Neural Information Processing Systems* 33 (2020): 20554-20565.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.
- Scott M Lundberg, Gabriel G Erion, and Su-In Lee. 2018. Consistent individualized feature attribution for tree ensembles. *arXiv preprint arXiv:1802.03888*.