# Attention is not Explanation

Mark Jacobsen

# Assumptions and Research Questions

**Questions:**
Does attention provide model transparency?
Are attended-to features responsible for outputs?
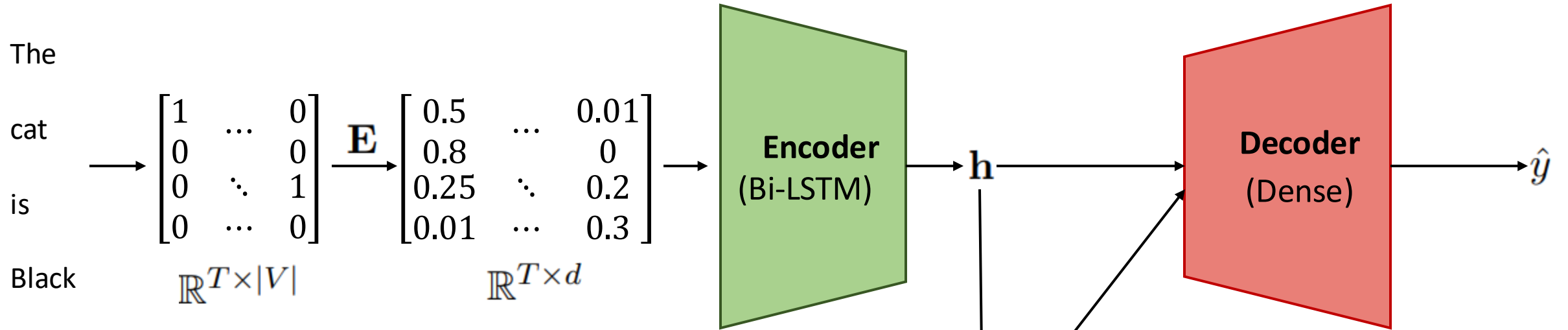
Authors claim: No

**If yes:**
1. Attention weights should correlate with feature importance methods.
   - Gradient-based methods
   - Leave-one-out
2. Alternative attention weight configurations should yield corresponding changes in prediction.

Experiments

# Attention / Model Architecture

The

cat

$$\begin{bmatrix} 1 & & 0 \\ 0 & \cdots & 0 \\ 0 & \ddots & 1 \\ 0 & \cdots & 0 \end{bmatrix} \xrightarrow{\mathbf{E}} \begin{bmatrix} 0.5 & \cdots & 0.01 \\ 0.8 & & 0 \\ 0.25 & \ddots & 0.2 \\ 0.01 & \cdots & 0.3 \end{bmatrix}$$

is

Black

$$\mathbb{R}^{T \times |V|} \qquad\qquad \mathbb{R}^{T \times d}$$

**Encoder**
(Bi-LSTM)

$\mathbf{h}$

**Decoder**
(Dense)

$\hat{y}$

$$\hat{\boldsymbol{\alpha}} = \mathrm{softmax}(\phi(\mathbf{h}, \mathbf{Q})) \in \mathbb{R}^T$$

1. $\phi(\mathbf{h}, \mathbf{Q}) = \dfrac{\mathbf{hQ}}{\sqrt{m}}$

2. $\phi(\mathbf{h}, \mathbf{Q}) = \mathbf{v}^T \tanh(\mathbf{W_1}\mathbf{h} + \mathbf{W_2}\mathbf{Q})$

# Experiments

Mark Jacobsen

# Datasets / NLP Tasks

**Sentiment Analysis**:
- Stanford Sentiment Treebank (SST)
- IMDB Large Movie Revies Corpus

**Other Binary Text Classification**:
- Twitter Adverse Drug Reaction dataset
- 20 Newsgroups (Hockey vs Basketball)
- AG News Corpus
- ...

**Natural Language Inference**:
- SNLI Dataset

**Question Answering**:
- CNN News Articles
- bAbI

| Dataset | $|V|$ | Avg. length | Train size | Test size | Test performance |
|---|---|---|---|---|---|
| SST | 16175 | 19 | 3034 / 3321 | 863 / 862 | 0.81 |
| IMDB | 13916 | 179 | 12500 / 12500 | 2184 / 2172 | 0.88 |
| ADR Tweets | 8686 | 20 | 14446 / 1939 | 3636 / 487 | 0.61 |
| 20 Newsgroups | 8853 | 115 | 716 / 710 | 151 / 183 | 0.94 |
| AG News | 14752 | 36 | 30000 / 30000 | 1900 / 1900 | 0.96 |
| Diabetes (MIMIC) | 22316 | 1858 | 6381 / 1353 | 1295 / 319 | 0.79 |
| Anemia (MIMIC) | 19743 | 2188 | 1847 / 3251 | 460 / 802 | 0.92 |
| CNN | 74790 | 761 | 380298 | 3198 | 0.64 |
| bAbI (Task 1 / 2 / 3) | 40 | 8 / 67 / 421 | 10000 | 1000 | 1.0 / 0.65 / 0.64 |
| SNLI | 20982 | 14 | 182764 / 183187 / 183416 | 3219 / 3237 / 3368 | 0.78 |

# Feature Importance Correlation

1. Attention weights should correlate with feature importance methods.

**Algorithm 1** Feature Importance Computations

$$\mathbf{h} \leftarrow \mathrm{Enc}(\mathbf{x}), \; \hat{\alpha} \leftarrow \mathrm{softmax}(\phi(\mathbf{h}, \mathbf{Q}))$$
$$\hat{y} \leftarrow \mathrm{Dec}(\mathbf{h}, \alpha)$$

Model Computation

$$g_t \leftarrow \left| \sum_{w=1}^{|V|} \mathbb{1}[\mathbf{x}_{tw} = 1] \frac{\partial y}{\partial \mathbf{x}_{tw}} \right|, \; \forall t \in [1, T]$$
$$\tau_g \leftarrow \mathrm{Kendall}\text{-}\tau(\alpha, g)$$

Gradient Based Feature Importance
Kendall Correlation

$$\Delta \hat{y}_t \leftarrow \mathrm{TVD}(\hat{y}(\mathbf{x}_{-t}), \hat{y}(\mathbf{x})), \; \forall t \in [1, T]$$
$$\tau_{loo} \leftarrow \mathrm{Kendall}\text{-}\tau(\alpha, \Delta \hat{y})$$

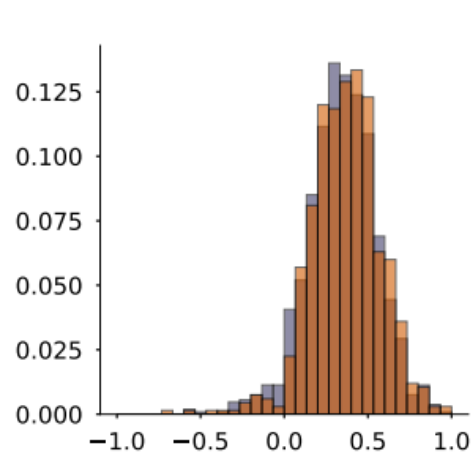Leave-One-Out Feature Importance
Kendall Correlation

UNIVERSITY OF CAMBRIDGE

# Feature Importance Correlation

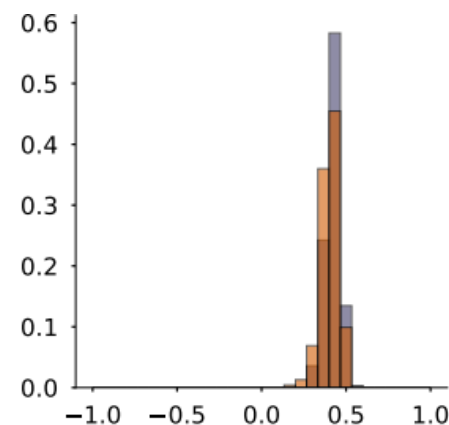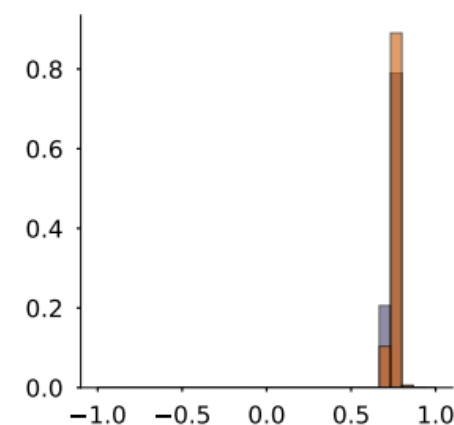| Dataset | Class | Gradient (BiLSTM) $\tau_g$ | | Gradient (Average) $\tau_g$ | | Leave-One-Out (BiLSTM) $\tau_{loo}$ | |
|---|---|---|---|---|---|---|---|
| | | Mean $\pm$ Std. | Sig. Frac. | Mean $\pm$ Std. | Sig. Frac. | Mean $\pm$ Std. | Sig. Frac. |
| SST | 0 | $0.34 \pm 0.21$ | 0.48 | $0.61 \pm 0.20$ | 0.87 | $0.27 \pm 0.19$ | 0.33 |
| | 1 | $0.36 \pm 0.21$ | 0.49 | $0.60 \pm 0.21$ | 0.83 | $0.32 \pm 0.19$ | 0.40 |
| IMDB | 0 | $0.44 \pm 0.06$ | 1.00 | $0.67 \pm 0.05$ | 1.00 | $0.34 \pm 0.07$ | 1.00 |
| | 1 | $0.43 \pm 0.06$ | 1.00 | $0.68 \pm 0.05$ | 1.00 | $0.34 \pm 0.07$ | 0.99 |
| ADR Tweets | 0 | $0.47 \pm 0.18$ | 0.76 | $0.73 \pm 0.13$ | 0.96 | $0.29 \pm 0.20$ | 0.44 |
| | 1 | $0.49 \pm 0.15$ | 0.85 | $0.72 \pm 0.12$ | 0.97 | $0.44 \pm 0.16$ | 0.74 |
| 20News | 0 | $0.07 \pm 0.17$ | 0.37 | $0.79 \pm 0.07$ | 1.00 | $0.06 \pm 0.15$ | 0.29 |
| | 1 | $0.21 \pm 0.22$ | 0.61 | $0.75 \pm 0.08$ | 1.00 | $0.20 \pm 0.20$ | 0.62 |
| AG News | 0 | $0.36 \pm 0.13$ | 0.82 | $0.78 \pm 0.07$ | 1.00 | $0.30 \pm 0.13$ | 0.69 |
| | 1 | $0.42 \pm 0.13$ | 0.90 | $0.76 \pm 0.07$ | 1.00 | $0.43 \pm 0.14$ | 0.91 |
| Diabetes | 0 | $0.42 \pm 0.05$ | 1.00 | $0.75 \pm 0.02$ | 1.00 | $0.41 \pm 0.05$ | 1.00 |
| | 1 | $0.40 \pm 0.05$ | 1.00 | $0.75 \pm 0.02$ | 1.00 | $0.45 \pm 0.05$ | 1.00 |
| Anemia | 0 | $0.47 \pm 0.05$ | 1.00 | $0.77 \pm 0.02$ | 1.00 | $0.46 \pm 0.05$ | 1.00 |
| | 1 | $0.46 \pm 0.06$ | 1.00 | $0.77 \pm 0.03$ | 1.00 | $0.47 \pm 0.06$ | 1.00 |
| CNN | Overall | $0.24 \pm 0.07$ | 0.99 | $0.50 \pm 0.10$ | 1.00 | $0.20 \pm 0.07$ | 0.98 |
| bAbI 1 | Overall | $0.25 \pm 0.16$ | 0.55 | $0.72 \pm 0.12$ | 0.99 | $0.16 \pm 0.14$ | 0.28 |
| bAbI 2 | Overall | $-0.02 \pm 0.14$ | 0.27 | $0.68 \pm 0.06$ | 1.00 | $-0.01 \pm 0.13$ | 0.27 |
| bAbI 3 | Overall | $0.24 \pm 0.11$ | 0.87 | $0.61 \pm 0.13$ | 1.00 | $0.26 \pm 0.10$ | 0.89 |
| SNLI | 0 | $0.31 \pm 0.23$ | 0.36 | $0.59 \pm 0.18$ | 0.80 | $0.16 \pm 0.26$ | 0.20 |
| | 1 | $0.33 \pm 0.21$ | 0.38 | $0.58 \pm 0.19$ | 0.80 | $0.36 \pm 0.19$ | 0.44 |
| | 2 | $0.31 \pm 0.21$ | 0.36 | $0.57 \pm 0.19$ | 0.80 | $0.34 \pm 0.20$ | 0.40 |

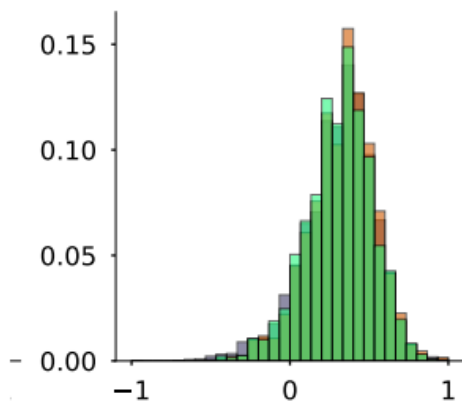# Gradient Feature Importance Correlation
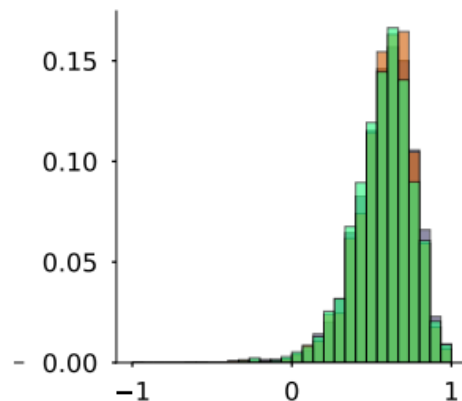


(a) SST (BiLSTM)

(b) SST (Average)
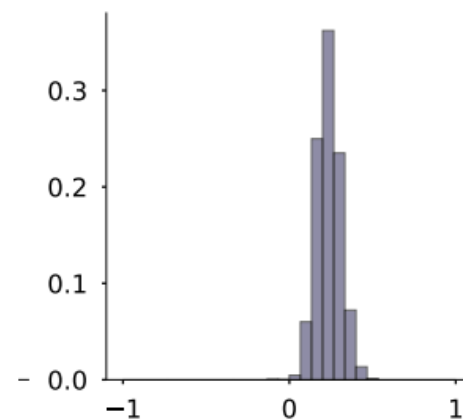
(c) Diabetes (BiLSTM)
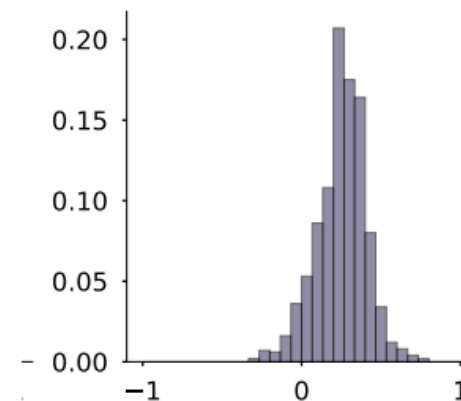
(d) Diabetes (Average)

(e) SNLI (BiLSTM)

(f) SNLI (Average)

(g) CNN-QA (BiLSTM)

(h) BAbI 1 (BiLSTM)

# Attention Changes

2. Alternative attention weight configurations should yield corresponding changes in prediction.

---

**Algorithm 2** Permuting attention weights

---

$\mathbf{h} \leftarrow \mathrm{Enc}(\mathbf{x}), \hat{\alpha} \leftarrow \mathrm{softmax}(\phi(\mathbf{h}, \mathbf{Q}))$

$\hat{y} \leftarrow \mathrm{Dec}(\mathbf{h}, \hat{\alpha})$

**for** $p \leftarrow 1$ to $100$ **do**

$\quad \alpha^p \leftarrow \mathrm{Permute}(\hat{\alpha})$

$\quad \hat{y}^p \leftarrow \mathrm{Dec}(\mathbf{h}, \alpha^p) \qquad \triangleright \text{Note} : \mathbf{h} \text{ is not changed}$

$\quad \Delta \hat{y}^p \leftarrow \mathrm{TVD}[\hat{y}^p, \hat{y}]$

**end for**

$\Delta \hat{y}^{med} \leftarrow \mathrm{Median}_p(\Delta \hat{y}^p)$

---

Model Computation
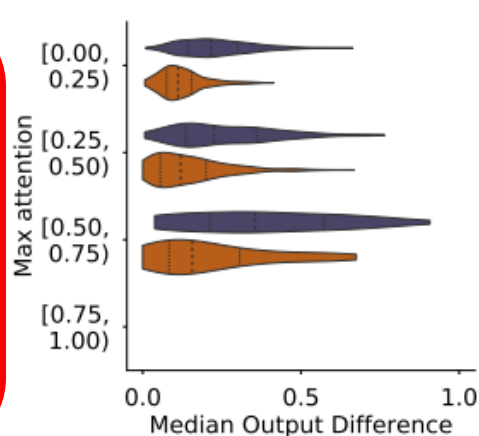
Random Attention Permutation
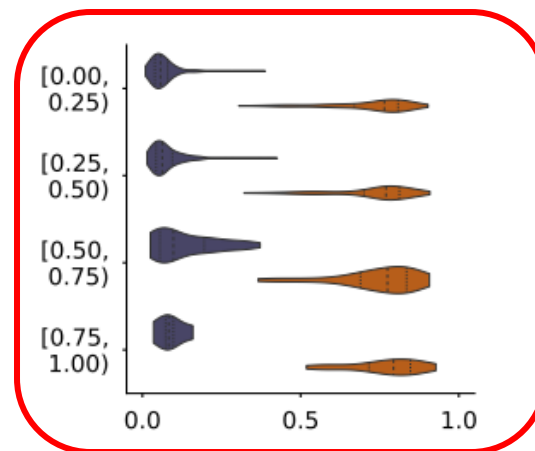Compute Output Difference

Median Output Difference
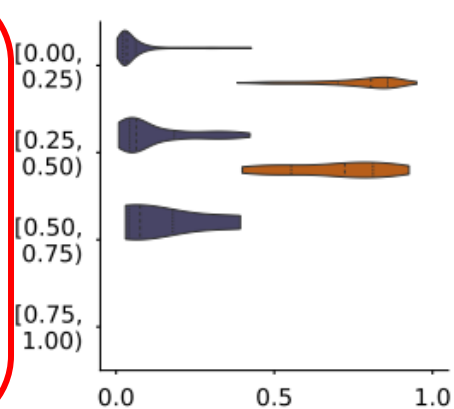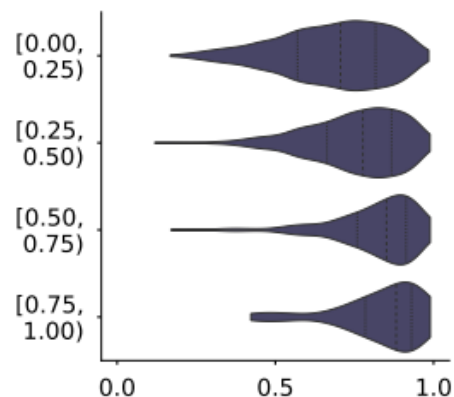
# Random Attention Permutation
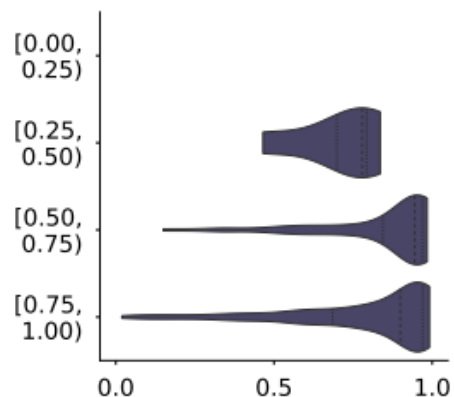


(a) SST (BiLSTM)

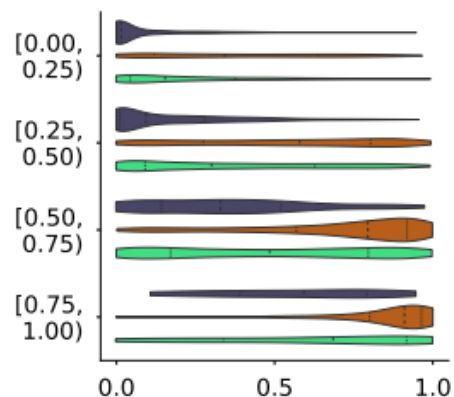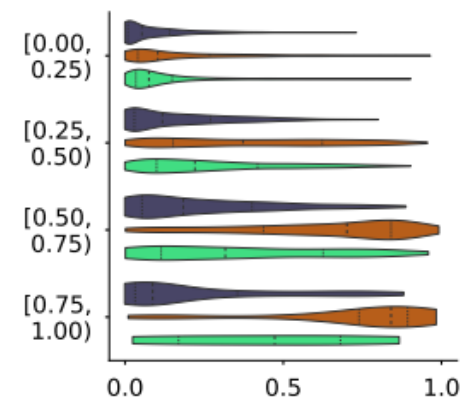(b) SST (CNN)

(c) Diabetes (BiLSTM)

(d) Diabetes (CNN)

(e) CNN-QA (BiLSTM)

(f) bAbI 1 (BiLSTM)

(g) SNLI (BiLSTM)

(h) SNLI (CNN)

# Adversarial Attention

2. Alternative attention weight configurations should yield corresponding changes in prediction.

**Algorithm 3** Finding adversarial attention weights

$\mathbf{h} \leftarrow \mathrm{Enc}(\mathbf{x}), \hat{\alpha} \leftarrow \mathrm{softmax}(\phi(\mathbf{h}, \mathbf{Q}))$

$\hat{y} \leftarrow \mathrm{Dec}(\mathbf{h}, \hat{\alpha})$

$\alpha^{(1)}, ..., \alpha^{(k)} \leftarrow \mathrm{Optimize\ Eq\ 1}$

**for** $i \leftarrow 1$ **to** $k$ **do**

$\quad \hat{y}^{(i)} \leftarrow \mathrm{Dec}(\mathbf{h}, \alpha^{(i)}) \qquad \triangleright \mathbf{h}$ is not changed

$\quad \Delta\hat{y}^{(i)} \leftarrow \mathrm{TVD}[\hat{y}, \hat{y}^{(i)}]$

$\quad \Delta\alpha^{(i)} \leftarrow \mathrm{JSD}[\hat{\alpha}, \alpha^{(i)}]$

**end for**

$\epsilon\text{-max JSD} \leftarrow \max_i \mathbb{1}[\Delta\hat{y}^{(i)} \leq \epsilon]\Delta\alpha^{(i)}$
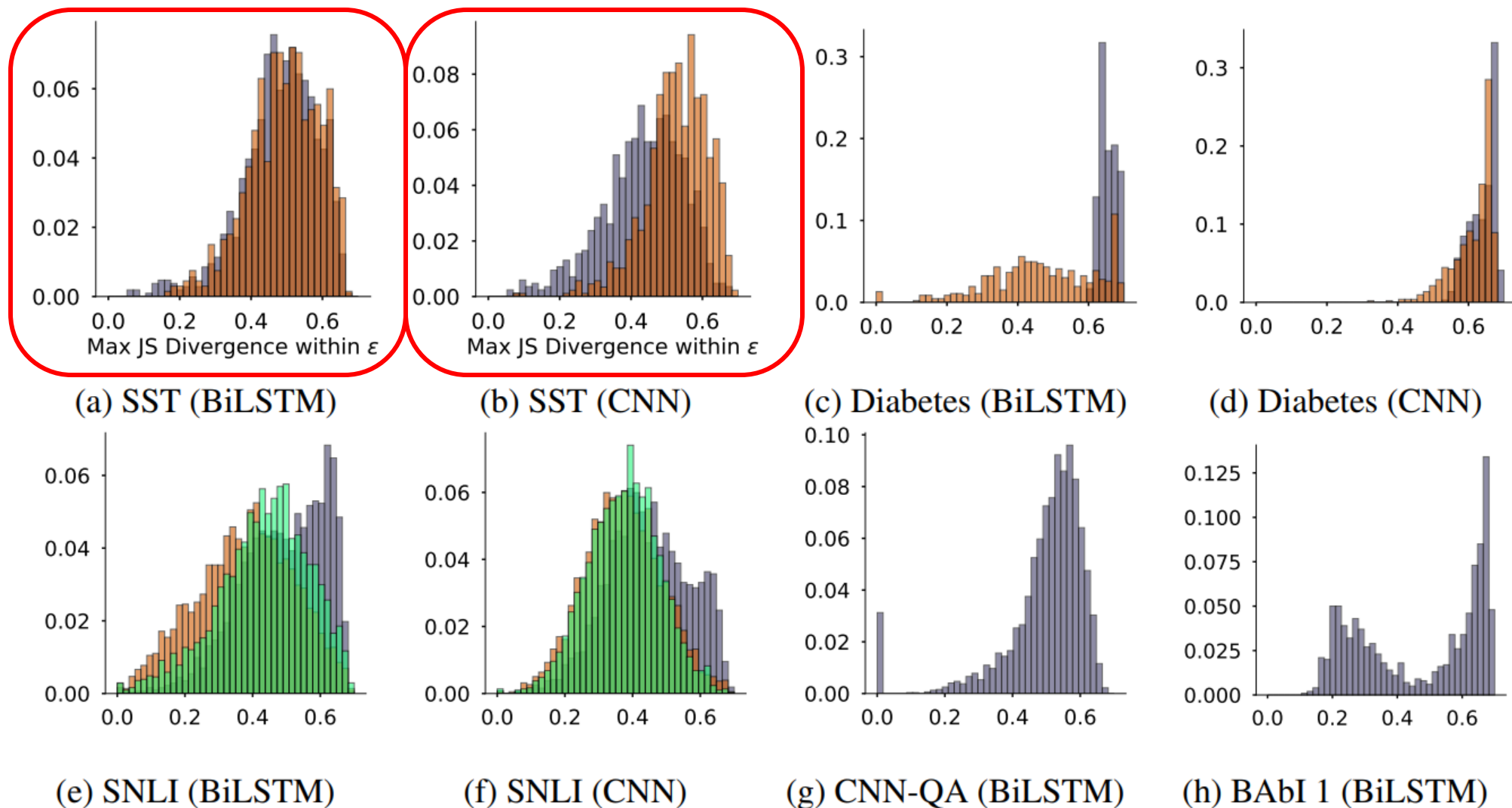
Model Computation

Adversarially Optimize Attention
Compute difference in
- Output
- Attention Weight

Find maximum attention
difference with maximum
threshold output difference

UNIVERSITY OF
CAMBRIDGE

# Adversarial Attention



(a) SST (BiLSTM)  (b) SST (CNN)  (c) Diabetes (BiLSTM)  (d) Diabetes (CNN)

(e) SNLI (BiLSTM)  (f) SNLI (CNN)  (g) CNN-QA (BiLSTM)  (h) BAbI 1 (BiLSTM)

# Results

**Prior Assumptions**

**1.** Attention weights should correlate with feature importance methods.
- Gradient-based methods
- Leave-one-out

Weak correlation

**2.** Alternative attention weight configurations should yield corresponding changes in prediction.

Negligible Changes

→ Attention is Not Explanation

# Discussion

Mark Jacobsen

# Pros / Cons

**Positive:**
- Clear research question and experiments for evaluation

- Various datasets and empirical evidence to backup claims

- They use previously established definitions of "explainability" / "transparency"
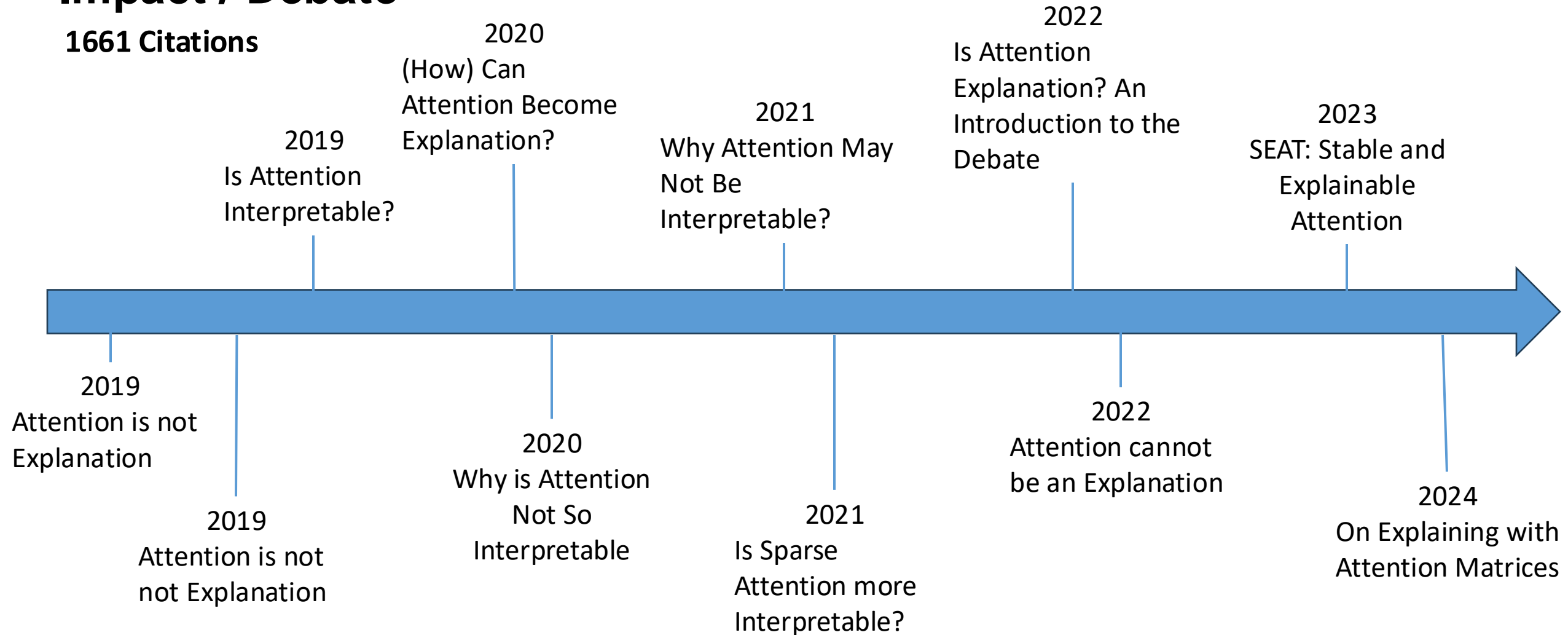  - Still up to debate

**Negative:**
- Questionable Assumptions
  - Using feature attribution as ground-truth
  - Explanations need to be exclusive

- Strong focus on binary classification and LSTMs

- Is changing the attention weights valid?

- Their own results sometimes show that attention can be explanation

# Impact / Debate

**1661 Citations**

2019
Is Attention
Interpretable?

2020
(How) Can
Attention Become
Explanation?

2021
Why Attention May
Not Be
Interpretable?

2022
Is Attention
Explanation? An
Introduction to the
Debate

2023
SEAT: Stable and
Explainable
Attention

2019
Attention is not
Explanation

2019
Attention is not
not Explanation

2020
Why is Attention
Not So
Interpretable

2021
Is Sparse
Attention more
Interpretable?

2022
Attention cannot
be an Explanation

2024
On Explaining with
Attention Matrices

# Sources

[1] Jain, S., & Wallace, B. C. (2019). Attention is not explanation. arXiv preprint arXiv:1902.10186. Retrieved from https://arxiv.org/abs/1902.10186.

[2] Wiegreffe, S., & Pinter, Y. (2019). Attention is not not explanation. arXiv preprint arXiv:1908.04626. Retrieved from https://arxiv.org/abs/1908.04626.

[3] Bibal, A., Cardon, R., Alfter, D., Wilkens, R., Wang, X., François, T., & Watrin, P. (2022). Is attention explanation? An introduction to the debate. In S. Muresan, P. Nakov, & A. Villavicencio (Eds.), Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) (pp. 3889–3900). Association for Computational Linguistics. https://aclanthology.org/2022.acl-long.269/.

UNIVERSITY OF CAMBRIDGE