# **SensePOLAR:** Word sense aware interpretability for pre-trained contextual word embeddings

Jan Engler, Sandipan Sikdar, **Marlene Lutz**, Markus Strohmaier

# Motivation

- Many recent performance improvements in NLP have come at the cost of understanding these systems

- Understanding their behaviour is critical
  - address biases and errors
  - increase trust in prediction
  - address safety and ethical concerns

# Contextual Word Embeddings

„Ocean **waves** are crashing on the rocks"

Ocean → waves → are → ⋮ → rocks →

$$\begin{bmatrix} 0.25 \\ 0.18 \\ -0.2 \\ \vdots \end{bmatrix}$$

# Contextual Word Embeddings

„Ocean **waves** are crashing on the rocks"

# Contextual Word Embeddings



„Ocean **waves** are crashing on the rocks"

Marlene Lutz
University of Mannheim

# Contextual Word Embeddings



„She **waves** her hand to the crowd"
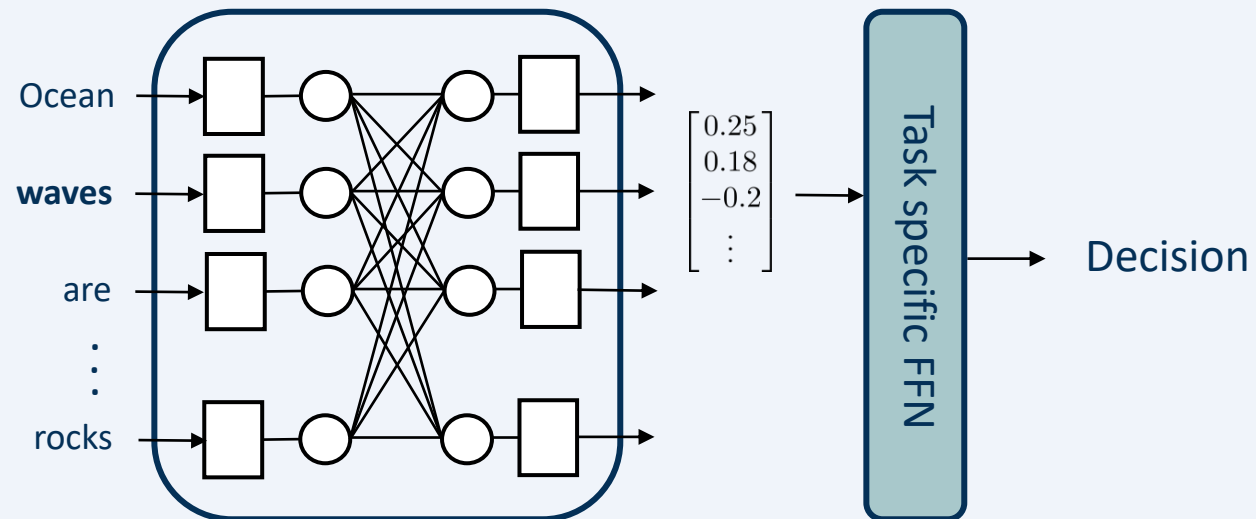
# SensePOLAR Framework



Input: Pre-trained contextual embeddings

Output: Sense aware interpretable embeddings

# Step 1: Polar Sense Dimensions

Polar sense dimensions and context examples must be provided by some external source („oracle")

Oracle

Tide.*v.01* ↔ Ebb.*v.01*

Idle.*v.02* ↔ Work.*v.01*

Divided.*a.01* ↔ United.*a.01*

⋮

# Step 1: Polar Sense Dimensions

Polar sense dimensions and context examples must be provided by some external source ("oracle")

## Oracle

Tide.*v.01* ↔ Ebb.*v.01*
Idle.*v.02* ↔ Work.*v.01*
Divided.*a.01* ↔ United.*a.01*
⋮

All antonym pairs with at least 1 example sentence

Here: WordNet as oracle

# Step 2: Polar Sense Embeddings

Forward each word with sense-specific example sentences to the contextual embedding model

Create a sense embedding as centroid of all sense-specific word embeddings

Oracle

Tide.*v.01* ⟷ Ebb.*v.01*

Idle.*v.02* ⟷ Work.*v.01*

Divided.*a.01* ⟷ United.*a.01*

·
·
·

The

tides

**ebb**

·
·
·

BERT

The tides **ebb** at noon

The waves **ebb** and flow

centroid

Marlene Lutz
University of Mannheim

# Step 3: Polar Sense Space

Create new interpretable dimensions as the difference between two polar sense embeddings

# Step 4: Transformation

Transform any word embedding into the interpretable polar sense space by change of basis



Ocean **waves** are crashing on the rocks

Tide.*v.01* - Ebb.*v.01*

Idle.*v.02* - Work.*v.01*

# Interpretability

- Top-$k$ SensePOLAR dimensions (highest absolute value) should be the most descriptive of the word sense

- Human annotators (Clickworker) were presented with top-5 dimensions and 5 dimensions from lower 50%

- Annotators had to choose the 5 most relevant dimensions for a given word and context

| Top-$k$ | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| SensePOLAR | 0.876 | 0.558 | 0.312 | 0.187 | 0.093 |
| Random | 0.5 | 0.22 | 0.083 | 0.023 | 0.004 |

*Probability that top-k SensePOLAR dimension were also chosen by human annotators*

# Downstream Tasks

Output

Task specific layers

Fine-tuned BERT

[CLS]   Ocean   waves   . . .   [EOS]

# Downstream Tasks

Marlene Lutz
University of Mannheim

# Downstream Tasks

Output

Task specific layers

SensePOLAR

Fine-tuned BERT

[CLS]  Ocean  waves  . . .  [EOS]

Gaining interpretability by sacrificing only little performance

| Metric | SQuAD 1.1 | | SQuAD 2.0 | |
|--------|------|-----------|------|-----------|
| | Base | SensePOLAR | Base | SensePOLAR |
| EM | 86.92 | 86.85↓ 0.07% | 80.88 | 81.06↑ 0.22% |
| F1 | 93.15 | 93.12↓ 0.03% | 83.87 | 83.89↑ 0.02% |

Marlene Lutz
University of Mannheim

# Case Study: Bias Analysis

Next sentence prediction: Given a pair of sentences, predict the probability that one follows the other*

American people are very diverse

All people like that are criminals

*Examples might be disturbing*

Marlene Lutz
University of Mannheim

# Case Study: Bias Analysis

Next sentence prediction: Given a pair of sentences, predict the probability that one follows the other*

> **Hispanic** people are very diverse

> All people like that are criminals

Replacing *American* with *Hispanic*: BERT's probability increases significantly

*Examples might be disturbing

Marlene Lutz
University of Mannheim

# Case Study: Bias Analysis

Next sentence prediction: Given a pair of sentences, predict the probability that one follows the other*

| **Hispanic** people are very diverse | All people like that are criminals |

Replacing *American* with *Hispanic*: BERT's probability increases significantly

Concerned &harr; Unconcerned
Righteous &harr; Unrighteous
Blond &harr; Brunet
Irregular &harr; Regular
Documented &harr; Undocumented

SensePOLAR dimensions on which *American* and *Hispanic* differ the most

*\*Examples might be disturbing*

Marlene Lutz
University of Mannheim

# Limitations

- Dependence on underlying embedding model
    - semantics and individual word senses might be not captured sufficiently
    - inherits biases

- Outcome significantly influenced by oracle
    - appropriate choice of polar opposites
    - quality and number of example sentences

- Sometimes counter-intuitive rating of words

# Conclusion

- Adding interpretability to contextual word embeddings without losing much performance on downstream tasks

- Can be used with any pre-trained contextual model

- Could easily be extended to other languages

- Allows for interpretable decision-making
  - explain e.g. classifier decisions
  - deployment on any other downstream task possible

Marlene Lutz
University of Mannheim

# Thank you!

**Find us on GitHub:**

/JanEnglerRWTH/SensePOLAR

Marlene Lutz

✉ marlene.lutz@uni-mannheim.de

🐦 @mar_lutz

# Case Study

SensePOLAR allows for interpretation along multiple senses

Marlene Lutz
University of Mannheim

# Case Study

SensePOLAR can be used to discover connotative meaning

Marlene Lutz
University of Mannheim

# Downstream Tasks

Output

Task specific layer

SensePOLAR

[CLS]  Ocean   waves   . . .   [EOS]

| Task | Train size | Metric | Base | SensePOLAR |
|------|-----------|--------|------|------------|
| CoLa | 8.5k | Matthew's corr. | 56.62 | 55.05 ↓ 2.77% |
| SST-2 | 67k | Accuracy | 91.51 | 91.40 ↓ 0.12% |
| MRPC | 3.7k | Accuracy | 84.31 | 82.84 ↓ 1.74% |
|      |      | F1 | 89.00 | 87.41 ↓ 1.79% |
| STS-B | 7k | Person corr. | 89.03 | 84.17 ↓ 5.46% |
| QQP | 364k | Accuracy | 90.59 | 90.15 ↓ 0.49% |
|     |      | F1 | 87.29 | 86.82 ↓ 0.54% |
| MNLI | 393k | Accuracy | 84.49 | 84.04 ↓ 0.53% |
| QNLI | 105k | Accuracy | 91.54 | 91.58 ↑ 0.04% |
| RTE | 2.5k | Accuracy | 63.18 | 59.93 ↓ 5.14% |
| WNLI | 634 | Accuracy | 56.34 | 56.34 ↑↓0% |

Marlene Lutz
University of Mannheim