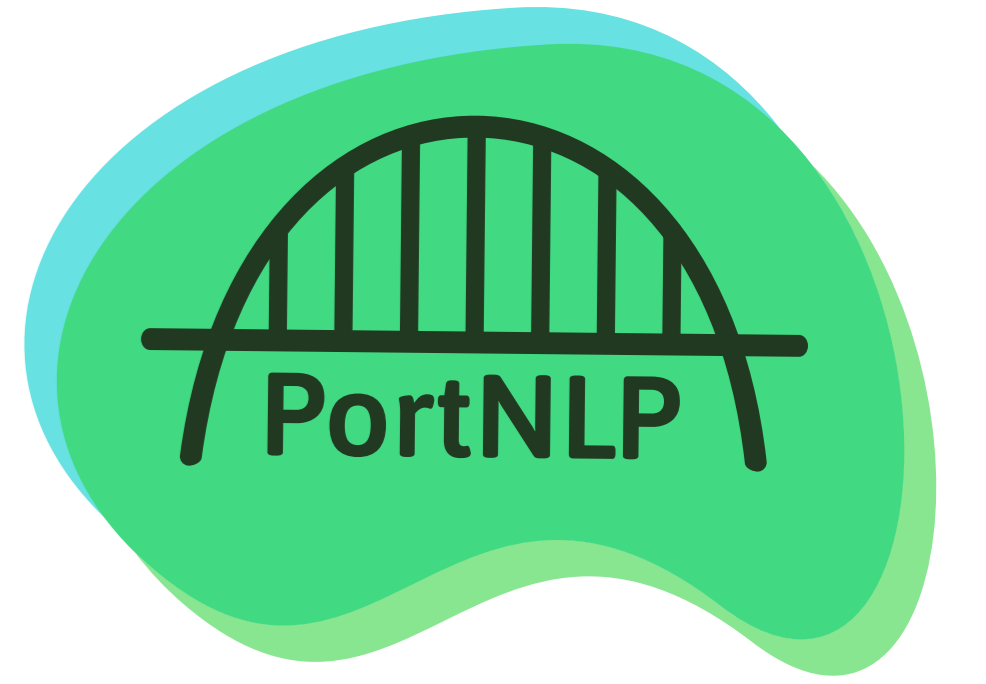# Estimating Semantic Similarity between In-Domain and Out-of-Domain Samples

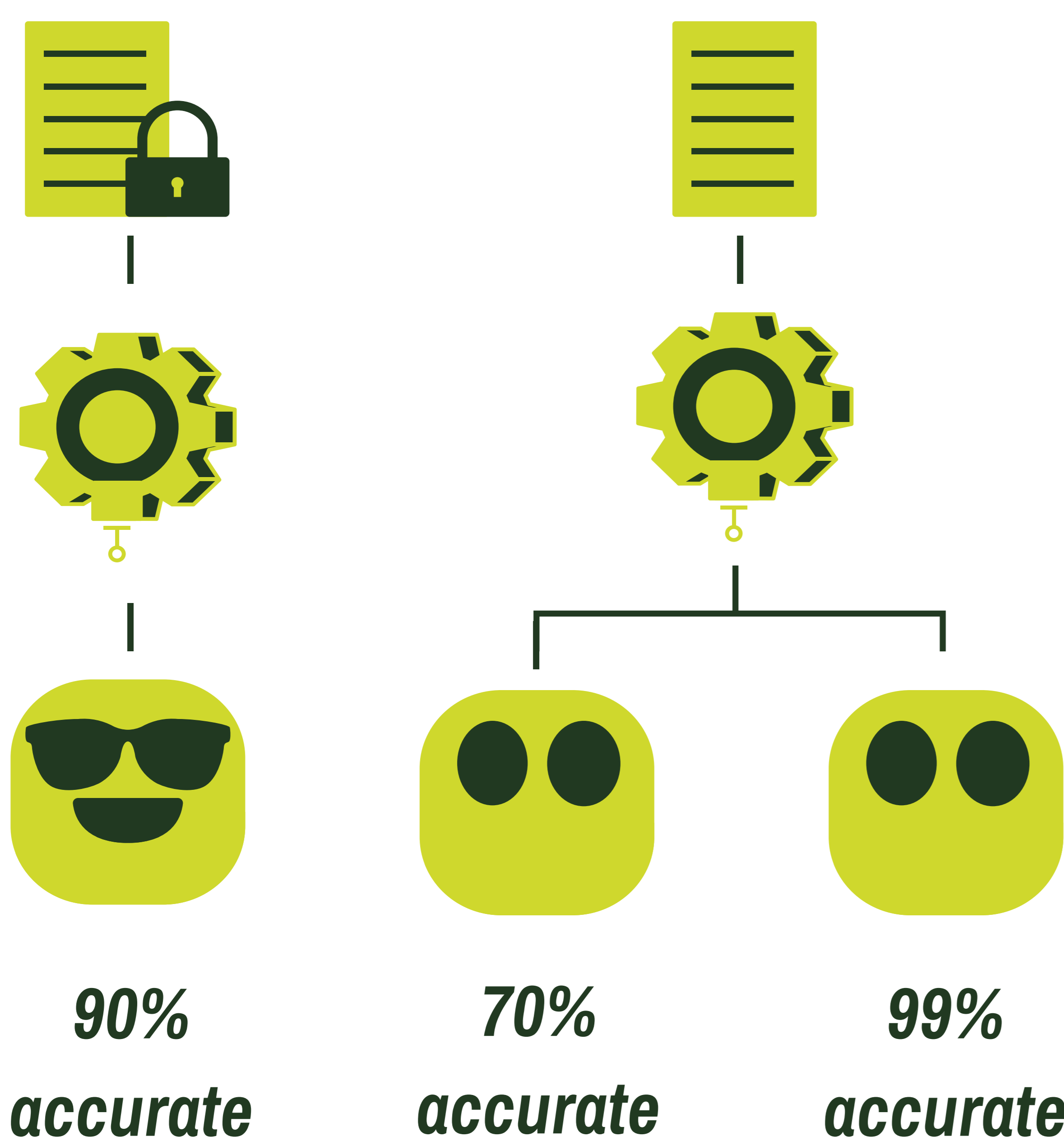**Rhitabrat Pokharel and Ameeta Agrawal**
*PortNLP Lab, Department of Computer Science, Portland State University*
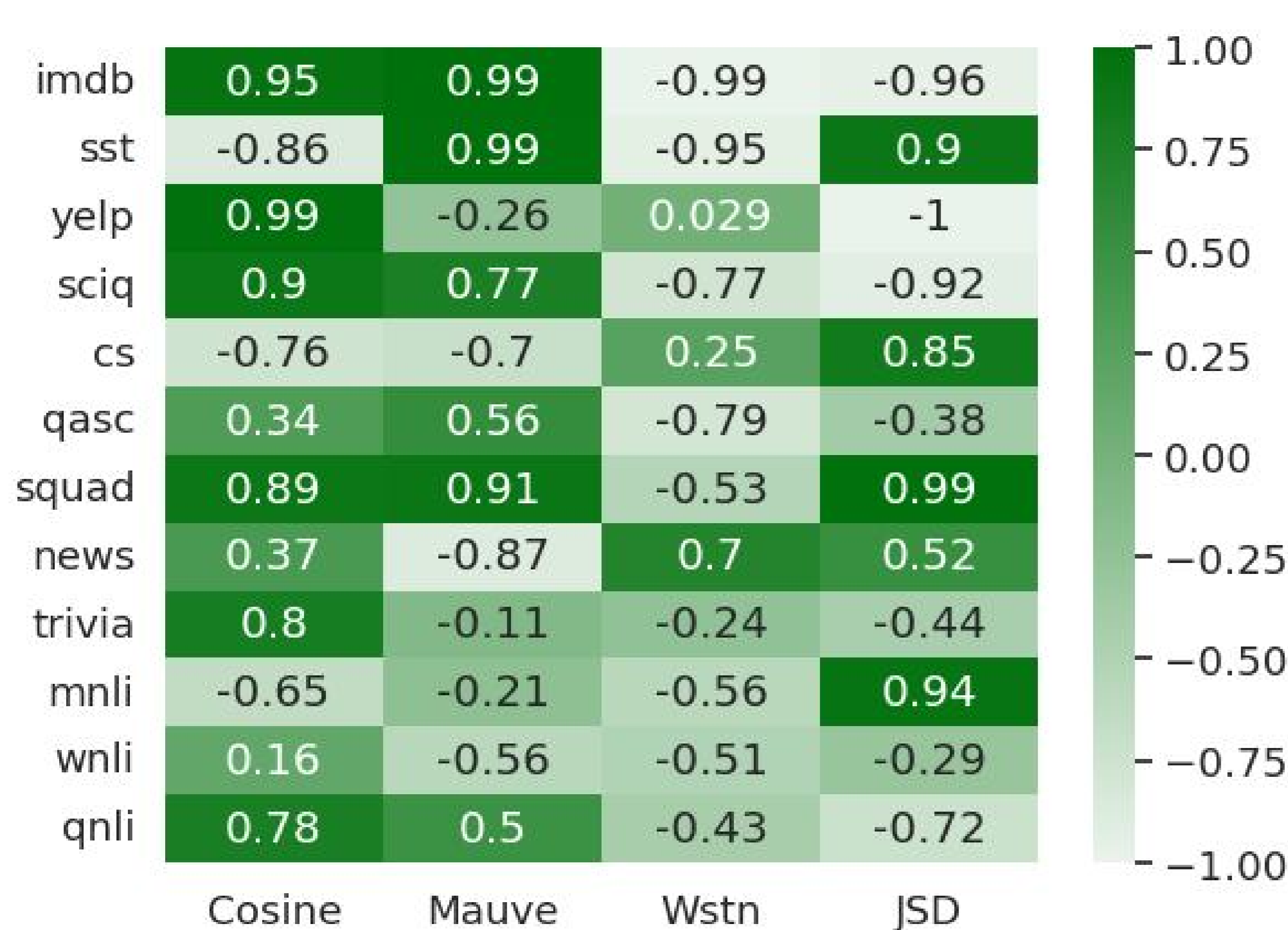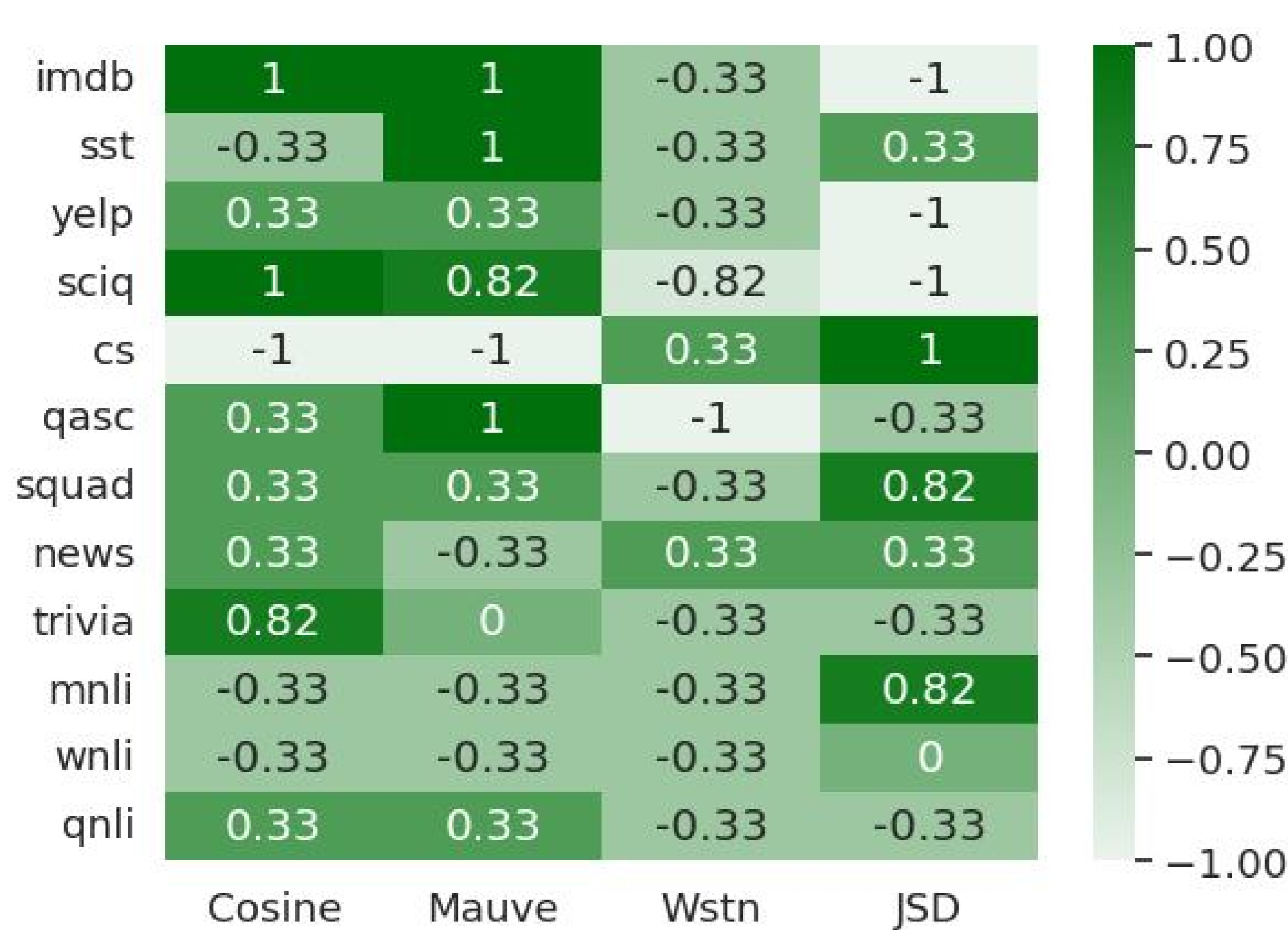*SEM, ACL 2023, Toronto*

## MOTIVATION

Models that demonstrate strong performance on carefully curated test/train sets may not necessarily showcase equivalent levels of effectiveness on real-world datasets.



90% accurate    70% accurate    99% accurate

*Out-of-domain (OOD) vs Out-of-distribution (OODist)*

## RESULTS



| | Cosine | Mauve | Wstn | JSD |
|---|---|---|---|---|
| imdb | 1 | 1 | -0.33 | -1 |
| sst | -0.33 | 1 | -0.33 | 0.33 |
| yelp | 0.33 | 0.33 | -0.33 | -1 |
| sciq | 1 | 0.82 | -0.82 | -1 |
| cs | -1 | -1 | 0.33 | 1 |
| qasc | 0.33 | 1 | -1 | -0.33 |
| squad | 0.33 | 0.33 | -0.33 | 0.82 |
| news | 0.33 | -0.33 | 0.33 | 0.33 |
| trivia | 0.82 | 0 | -0.33 | -0.33 |
| mnli | -0.33 | -0.33 | -0.33 | 0.82 |
| wnli | -0.33 | -0.33 | -0.33 | 0 |
| qnli | 0.33 | 0.33 | -0.33 | -0.33 |

| | Cosine | Mauve | Wstn | JSD |
|---|---|---|---|---|
| imdb | 0.95 | 0.99 | -0.99 | -0.96 |
| sst | -0.86 | 0.99 | -0.95 | 0.9 |
| yelp | 0.99 | -0.26 | 0.029 | -1 |
| sciq | 0.9 | 0.77 | -0.77 | -0.92 |
| cs | -0.76 | -0.7 | 0.25 | 0.85 |
| qasc | 0.34 | 0.56 | -0.79 | -0.38 |
| squad | 0.89 | 0.91 | -0.53 | 0.99 |
| news | 0.37 | -0.87 | 0.7 | 0.52 |
| trivia | 0.8 | -0.11 | -0.24 | -0.44 |
| mnli | -0.65 | -0.21 | -0.56 | 0.94 |
| wnli | 0.16 | -0.56 | -0.51 | -0.29 |
| qnli | 0.78 | 0.5 | -0.43 | -0.72 |

*Wstn and cosine show the most consistent correlation*

## OOD vs OODist

- Data from a related but different domain[1] (Amazon vs Twitter sentiment)
- Different datasets for the same task[2] (SST, IMDb, and Yelp for sentiment classification)
- Data collected at a different time[3] maybe under different settings
- Datasets that are not in the training set[4]

## DATASETS

| Datasets | Task |
|---|---|
| IMDb, SST2, Yelp | Sentiment Analysis |
| SCIQ, Commonsense, QASC | MCQ |
| SQUAD, News, Trivia | Extractive Question Answering |
| MNLI, WNLI, QNLI | Natural Language Inference |

For each of train, validation (when available), and test sets, we **downsampled** to the size of the smallest dataset.

## METHODOLOGY



Dataset X    Dataset Y    $X_{train}$    $Y_{train}$    Train (BERT)

Semantic Similarity

**Cosine, Mauve, Wasserstein, JSD**    Test    **Accuracy**

Correlated?

## CONCLUSION

- Wasserstein could be a potential metric for determining OOD samples
- Model does not always perform worse on OOD samples

## CODE & PAPER



pokharel@pdx.edu

[1] Dai, Wenyuan, Gui-Rong Xue, Qiang Yang, and Yong Yu. "Co-clustering based classification for out-of-domain documents." [2] Chrysostomou, George, and Nikolaos Aletras. "An empirical study on explanations in out-of-domain settings." [3] Ovadia, Yaniv, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. "Can you trust your model's uncertainty? evaluating predictive uncertainty under dataset shift." [4] Lin, Bill Yuchen, Sida Wang, Xi Victoria Lin, Robin Jia, Lin Xiao, Xiang Ren, and Wen-tau Yih. "On continual model refinement in out-of-distribution data streams."