



# Adversarial Scrubbing of Demographic Information for Text Classification

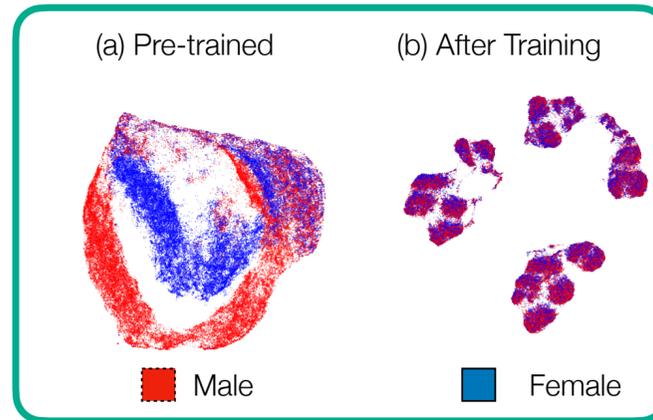
Somnath Basu Roy Chowdhury, Sayan Ghosh, Yiyuan Li, Junier B. Oliva, Shashank Srivastava, Snigdha Chaturvedi

{somnath, sayghosh, yiyuanli, joliva, ssvivastava, snigdha}@cs.unc.edu

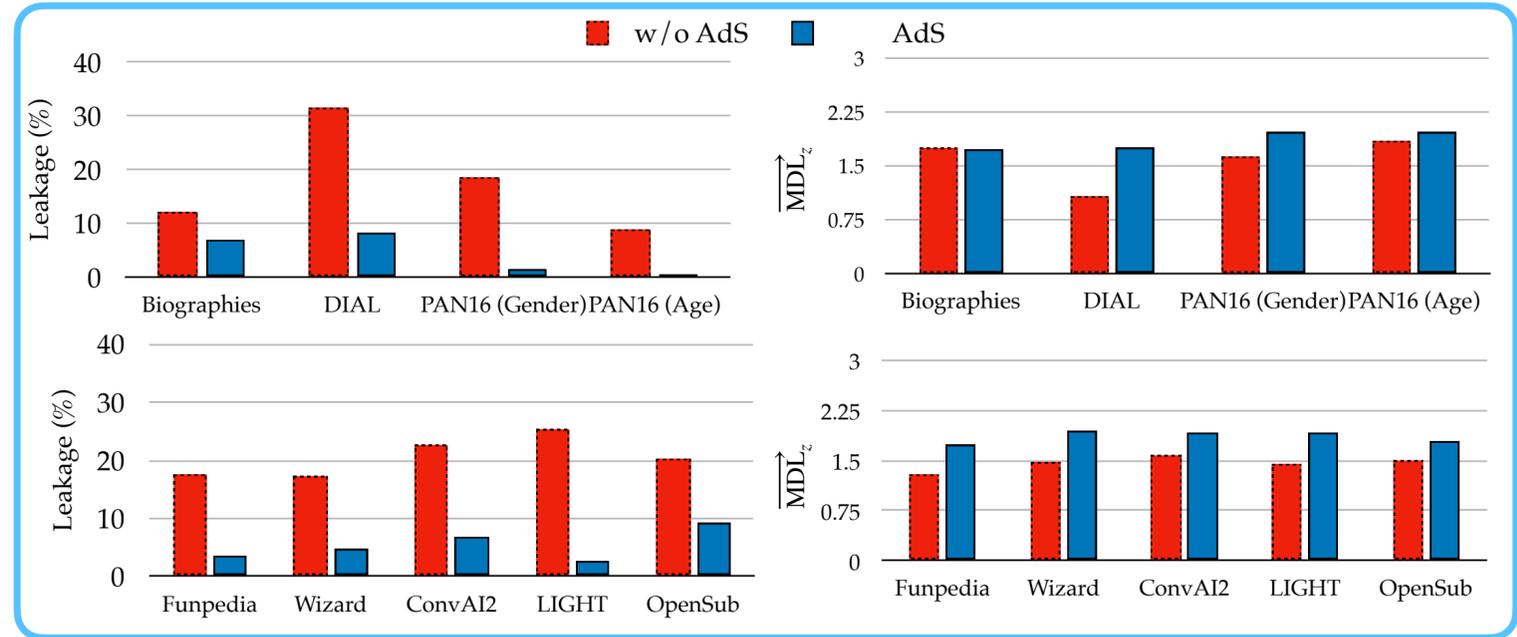
## Problem Statement

- Text classification systems can be **biased** towards certain demographic groups
- Natural language is highly indicative of demographic attributes
- ML model representation **encode sensitive information** even without having direct access to them
- We propose AdS to learn **fair representations** during training text classification systems

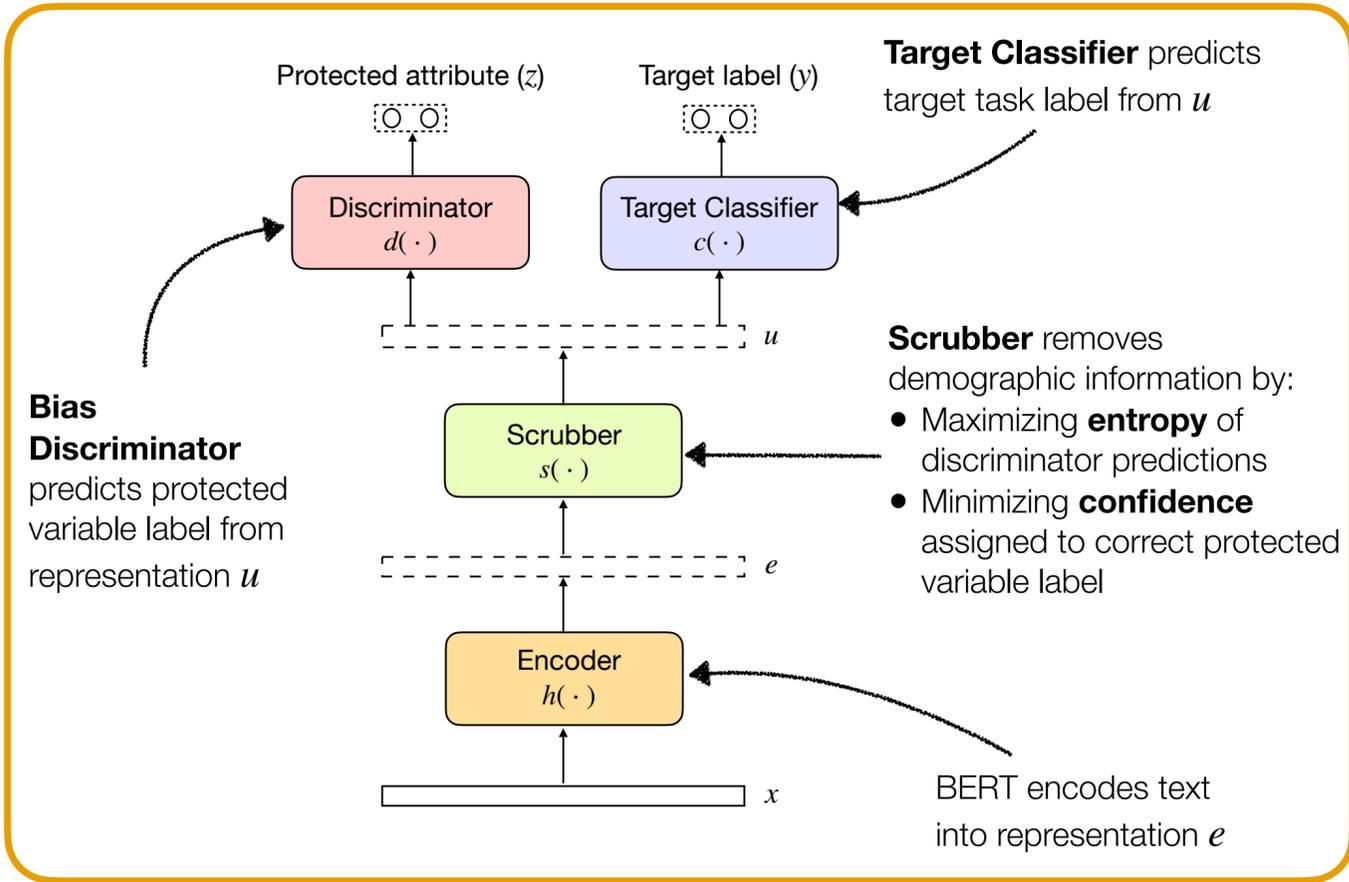
## Visualization



## Performance



## Adversarial Scrubber



## Theoretical Results

- Minimizing Scrubber loss increases the Bias discriminator loss
- Discriminator loss doesn't decrease when Encoder, Scrubber and Bias discriminator are updated
- For a strong Bias discriminator, the encoder and scrubber converge on the target task w/o leaking demographic information

## Probing

- We evaluate representation  $u$  using off-the shelf sklearn classifier and evaluate the metrics below:
- Protected variable **accuracy** (lower value is desirable)
  - Minimum description length (**MDL**) of protected variable (high value desired)
  - Target variable accuracy (high value is desirable)

## Conclusion

- We propose **Adversarial Scrubber** (AdS) to learn fair representation during text classification
- Theoretical and empirical analysis show AdS converges on target task **without leaking** sensitive information
- Effectively evaluate information content using MDL