

EDGAR

Jonathan Lung

EDGAR

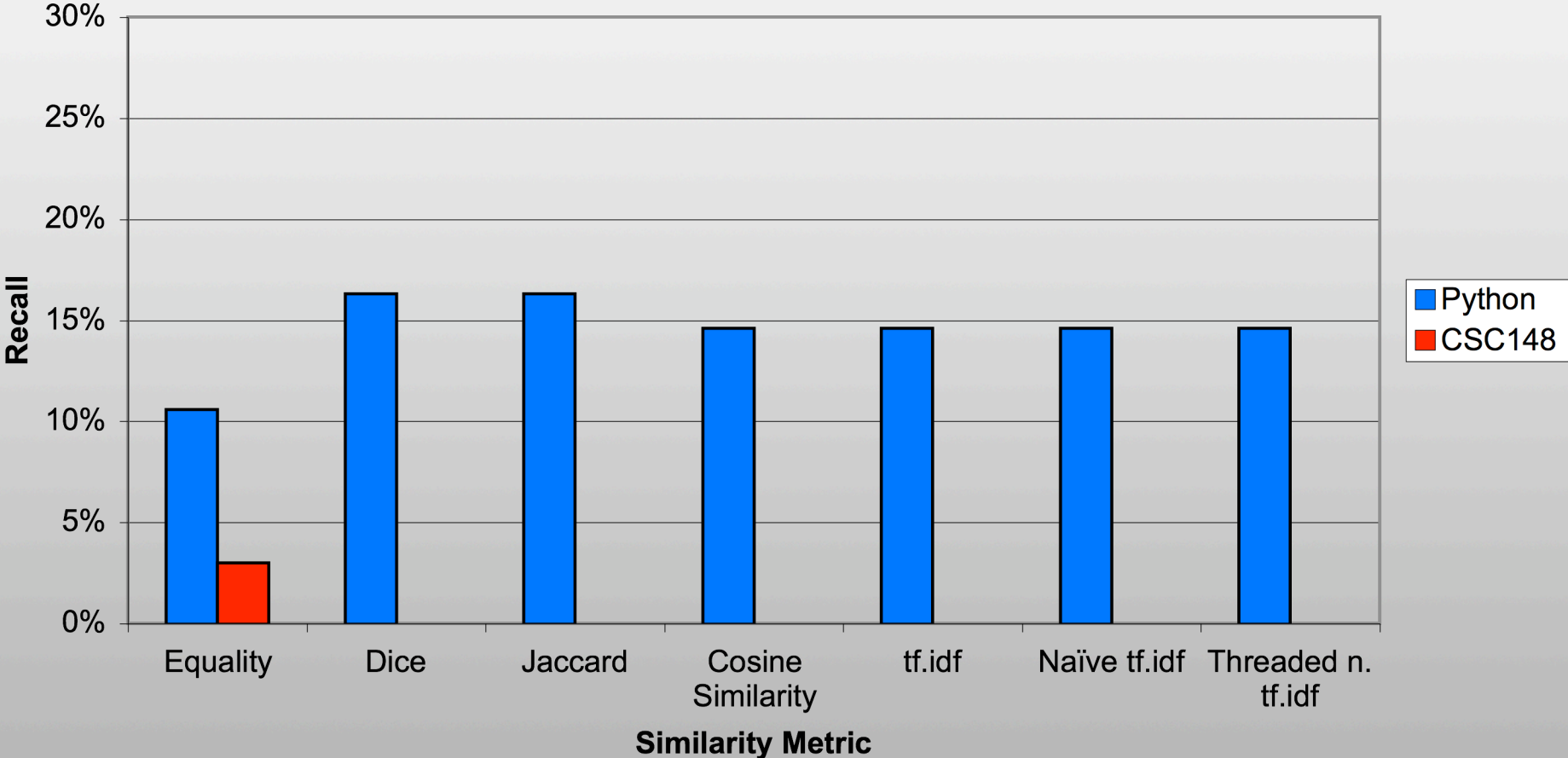
# Goals

- Reduce noise in electronic fora
  - Filtering
  - Throttling
- No human intervention
- Not domain-specific
- Can be incorporated into existing software

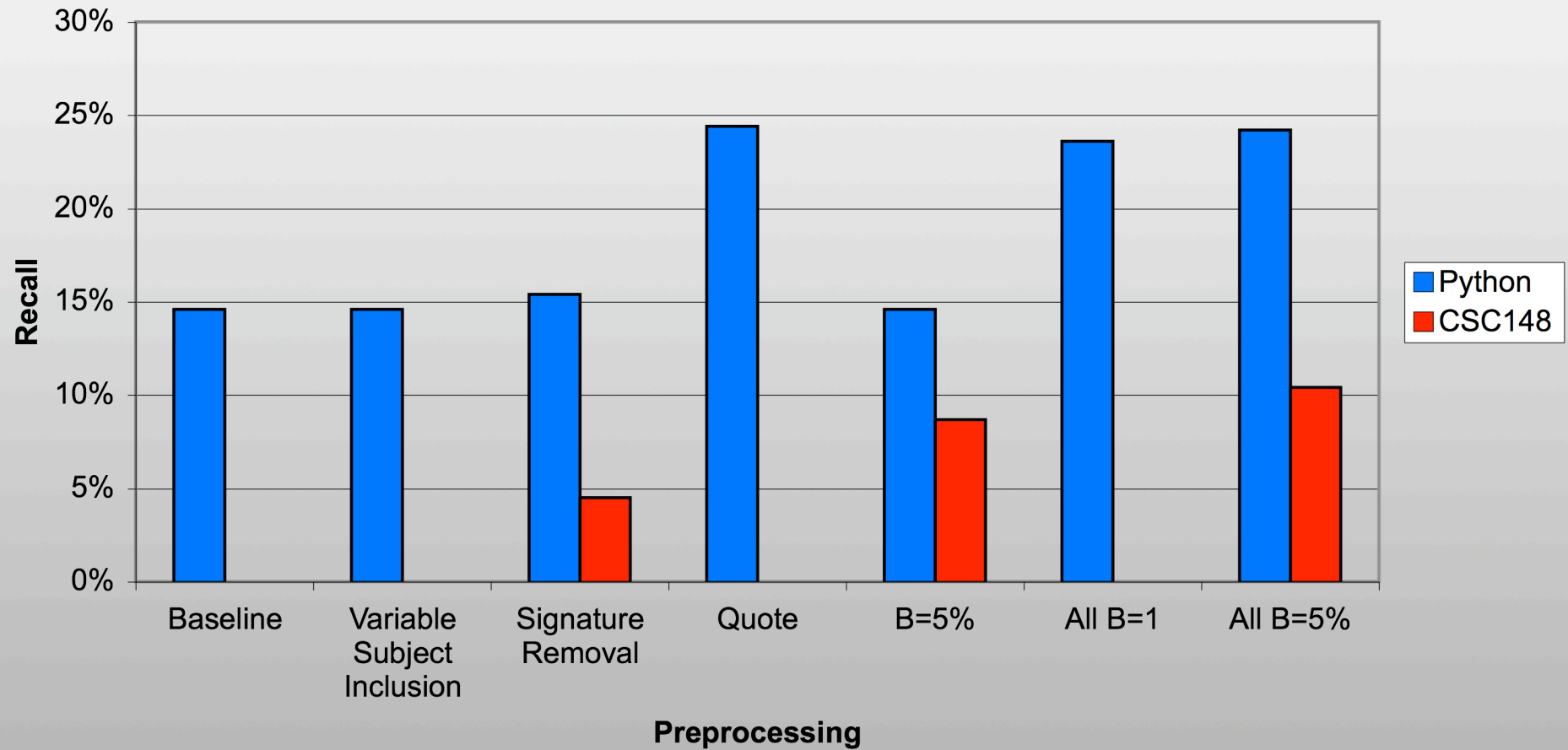
# Comparing Algorithms

- Precision and recall
- Parameter adjustment
- Corpora
  - Python mailing list
  - CSC148 newsgroups

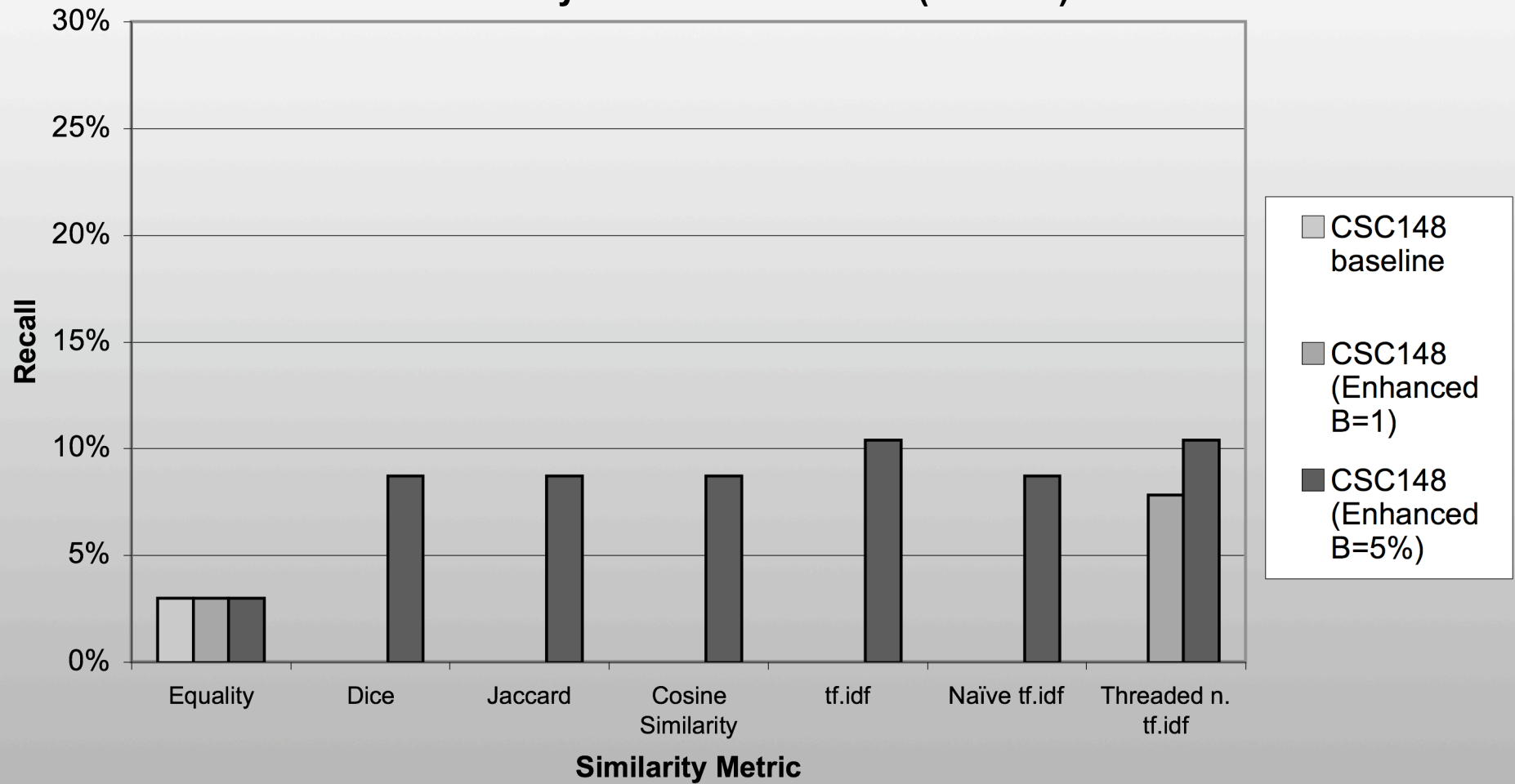
# Baseline Similarity Metric Performance



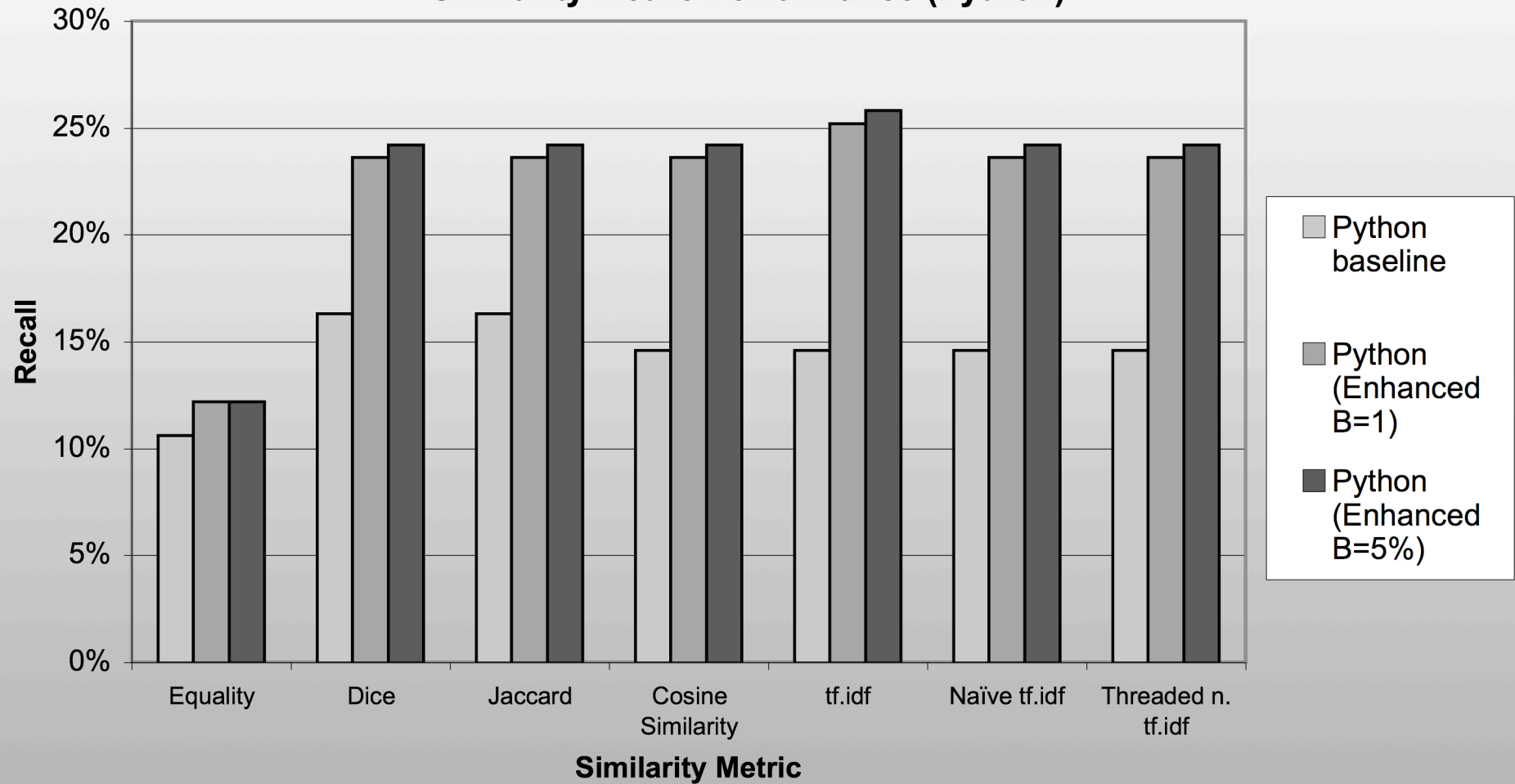
## Effects of Preprocessors on Recall



## Similarity Metric Performance (CSC148)

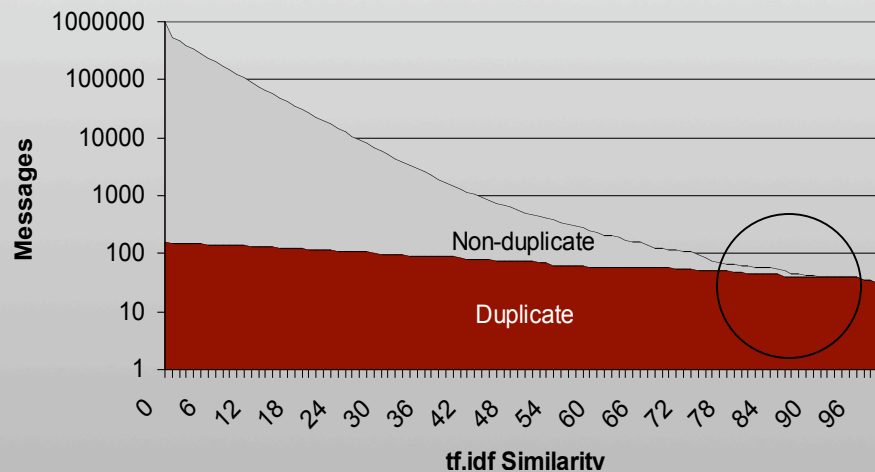


## Similarity Metric Performance (Python)

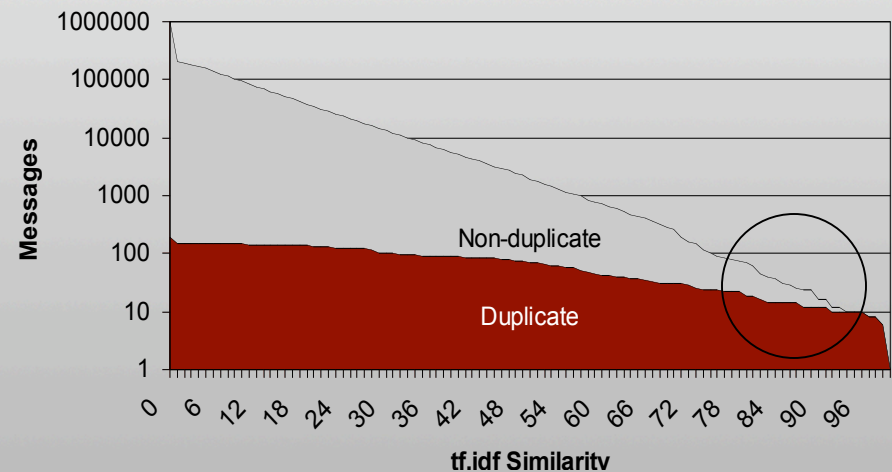


# Comparing the Corpora

Similarity vs. Proportion of Misidentified Messages (Python)



Similarity vs. Proportion of Misidentified Messages (CSC148)





# Conclusion

- Best ROI for improving pre-processors
- Room for improvement
- Many possible enhancements
- More data still required

# Questions

Goals

Comparing  
Algorithms

Enhancing  
Algorithms

Comparing  
Corpora

Conclusions

**Questions**

**EDGAR**

---