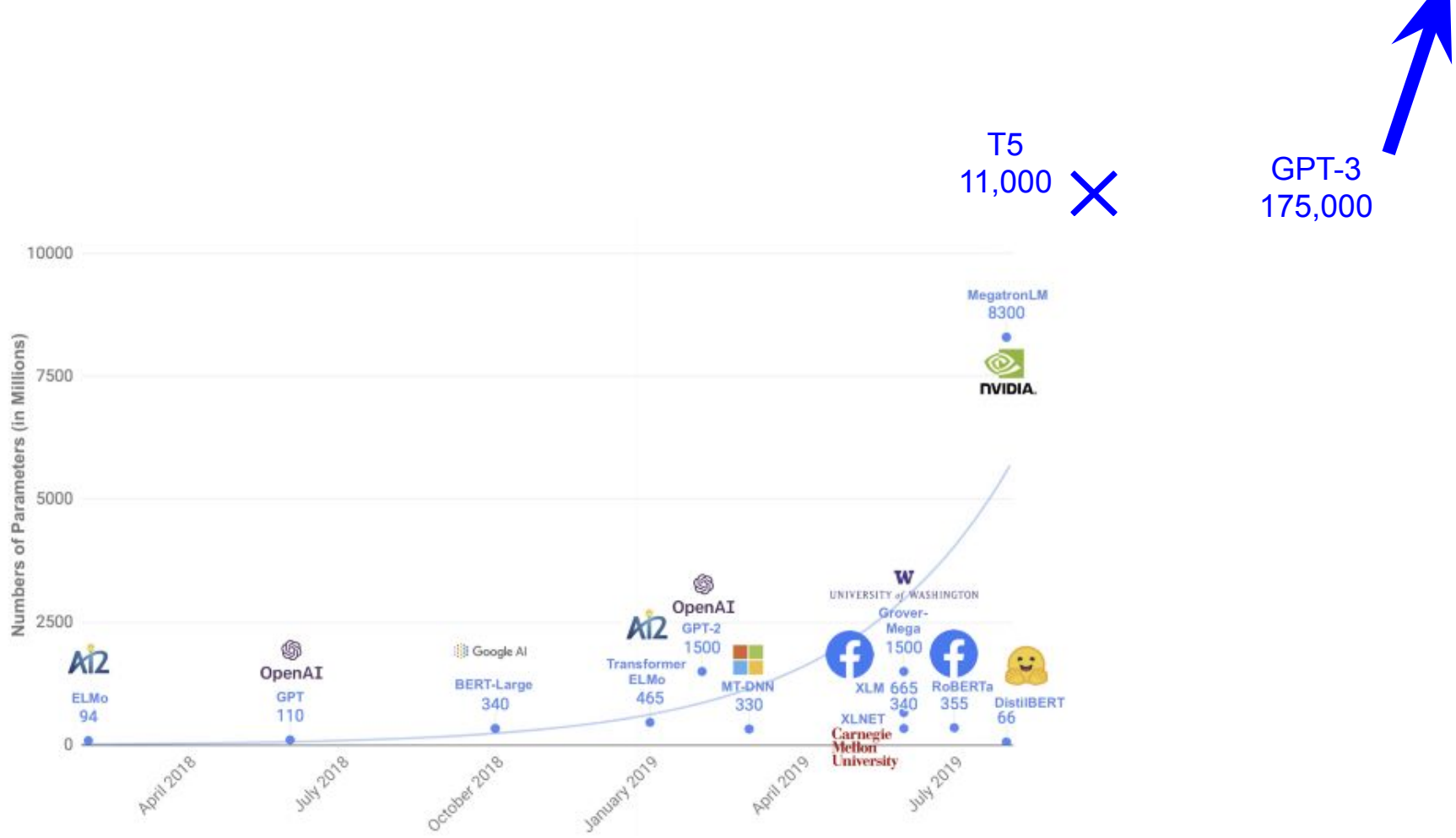


# Compressing BERT - An Evaluation and Combination of Methods

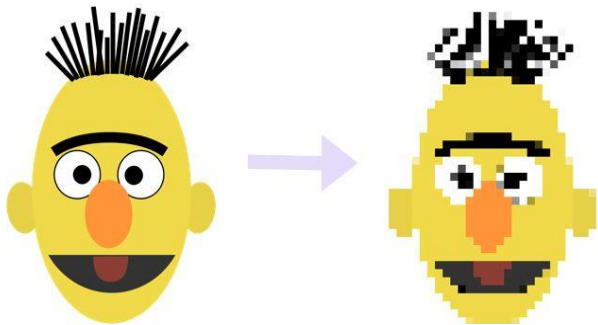
Sebastian Jäger

Model compression aims to reduce the number of parameters without or with as little performance loss as possible.

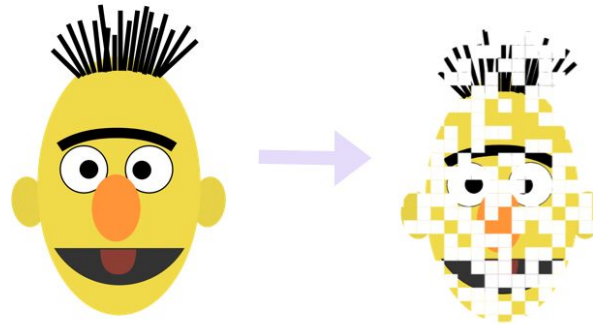
Cristian Bucila, Rich Caruana, and Alexandru Niculescu-Mizil. “Model Compression”. In: Proceedings of the Twelfth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD). 2006, pp. 535–541. doi: 10.1145/1150402.1150464.



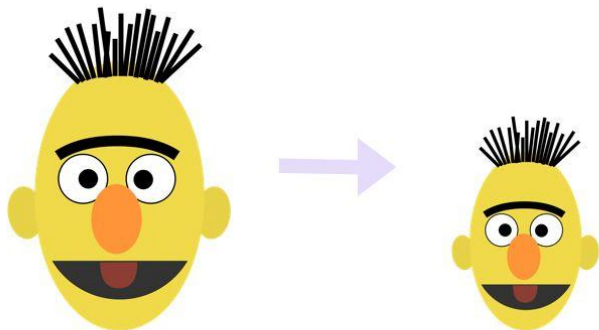
## Quantization



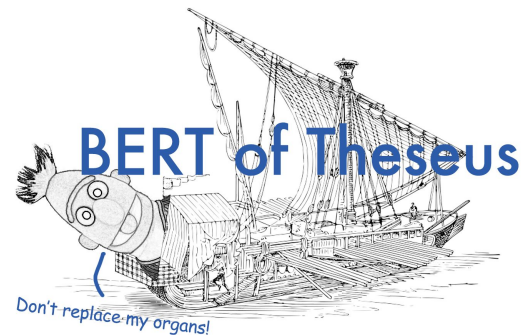
## Pruning



## Knowledge Distillation



## Theseus Compression



Images:

<https://blog.rasa.com/compressing-bert-for-faster-prediction-2/>

<https://raw.githubusercontent.com/JetRunner/BERT-of-Theseus/master/bert-of-theseus.png>

“Using a domain-specific BERT compressed with Theseus Compression helps to speed up experiment iterations for downstream task models.”

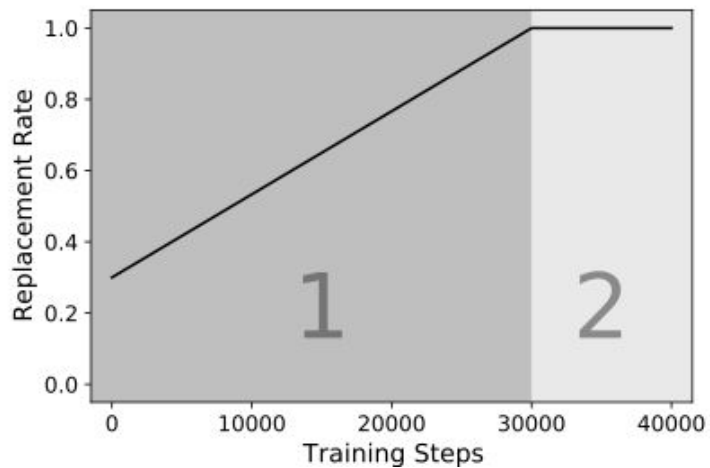
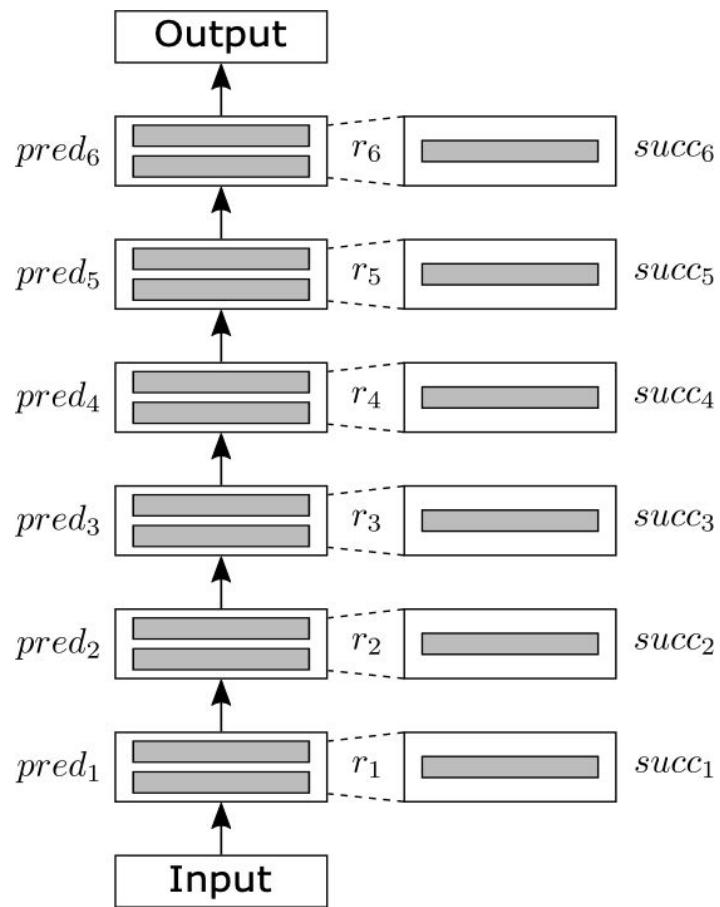
# Tasks and Datasets

## **In-hospital Mortality Prediction**

- MIMIC-III

## **German Hate Speech Detection**

- GermEval 19
- NOHATE + NOHATE LM



# Results of Hyperparameter Analysis

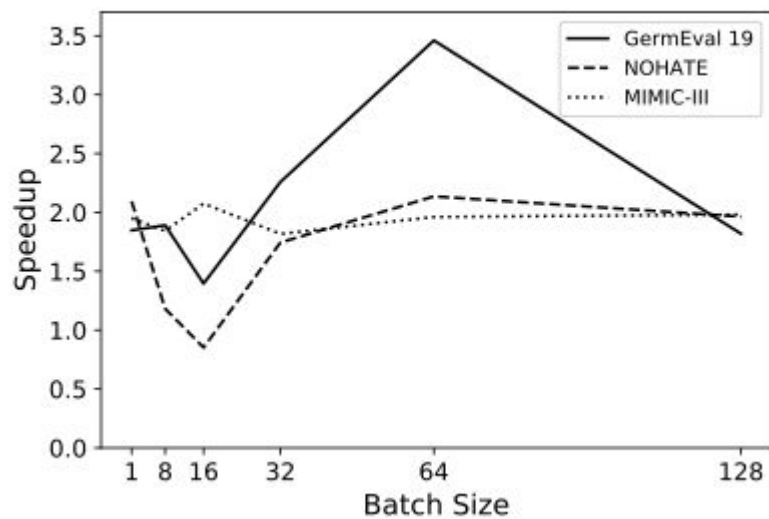
- Starting replacement rate around 0.1 to 0.3
- Module replacement for 85% of compression steps
- Both module replacement and successor fine-tuning are important
- Learning rate plays minor role in the compression performance



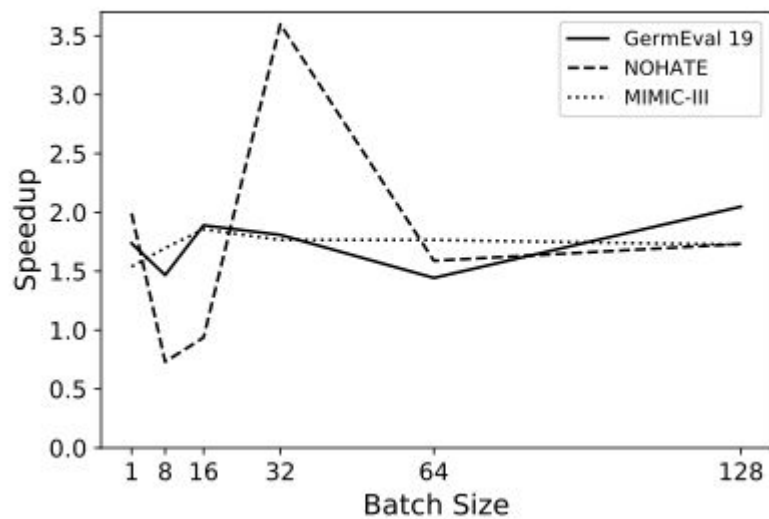
Model Settings	Retained Performance	Speedup Fine-tuning
NOHATE LM + GermEval 19	0.9750	2.50x
NOHATE LM + NOHATE	0.9889	2.16x
MIMIC-III LM + MIMIC-III	0.9947	2.14x

# Prediction Speedup

## CPU



## GPU



# Qualitative Error Analysis

100 FPs and FNs for baseline and compressed models.

- **GermEval 19:** Change of -5 and introduced 2 new error classes
- **NOHATE:** More stable, mostly changes within  $\pm 1$ , max +3
- **MIMIC-III:** All changes within  $\pm 1$

# Summary of Results

- Compression faster than plain LM fine-tuning
- Compression ratio: 1.67x
- Retained performance: up to 99%
- Speedup downstream training: >2x
- Speedup inference: 1.94x on CPU, 1.73x on GPU

# Conclusion and Limitation

- + Theseus Compression speeds up experiments
- + Compressed models are faster and smaller
  
- Theseus Compression always reduces prediction performance

# Thank You!

Questions?