

# Mechanistic Interpretability

GPT-2 small  
WT 2024/25

Frederick Riemenschneider



21.11.2024

`https://cdn.openai.com/better-language-models/  
language_models_are_unsupervised_multitask_  
learners.pdf`

- 124M parameters (BERT base: 110M)
- 40GB pre-training data (BERT base: 16GB)

# Experiment Setup

- English Web Treebank (Universal Dependencies)
- first 1000 sentences
  - 25533 words
- random split, test size = 0.2
- logistic regression







