

# Contextual word representations: ELECTRA

(Efficiently Learning an Encoder that Classifies Token Replacements Accurately)

Christopher Potts

Stanford Linguistics

CS224u: Natural language understanding



# Addressing the known limitations with BERT

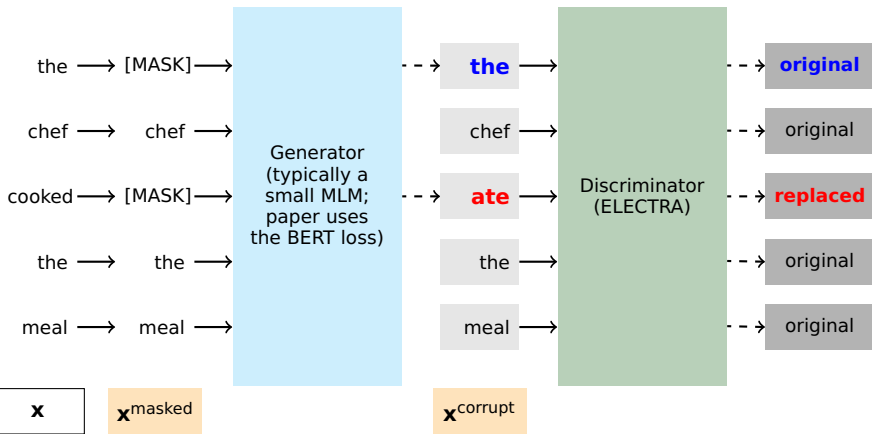
1. Devlin et al. (2019:§5): admirably detailed but still partial ablation studies and optimization studies.
2. Devlin et al. (2019): “The first [downside] is that we are creating a mismatch between pre-training and fine-tuning, since the [MASK] token is never seen during fine-tuning.”
3. Devlin et al. (2019): “The second downside of using an MLM is that only 15% of tokens are predicted in each batch”
4. Yang et al. (2019): “BERT assumes the predicted tokens are independent of each other given the unmasked tokens, which is oversimplified as high-order, long-range dependency is prevalent in natural language”

# Core model structure (Clark et al. 2019)

Random sample of  
 $\approx 15\%$  of tokens masked

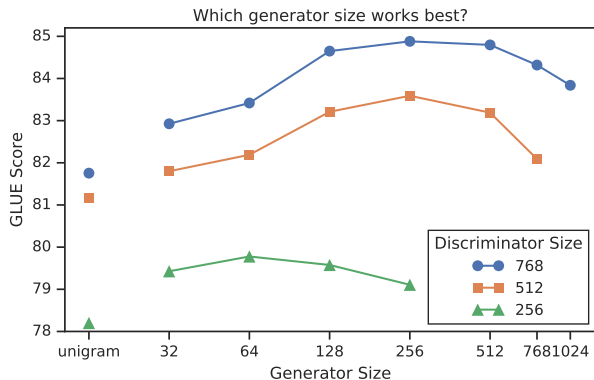
Masked tokens replaced  
 proportional to Gen-  
 erator probabilities

Loss:  
 Generator +  
 $\lambda$  ELECTRA



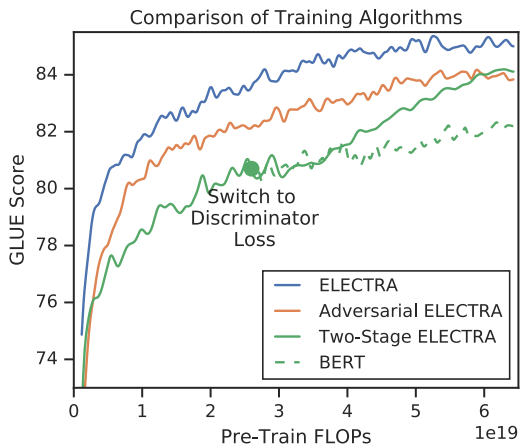
# Generator/Discriminator relationships

Where Generator and Discriminator are the same size, they can share Transformer parameters, and more sharing is better. However, the best results come from having a Generator that is small compared to the Discriminator:



Clark et al. 2019, Figure 3

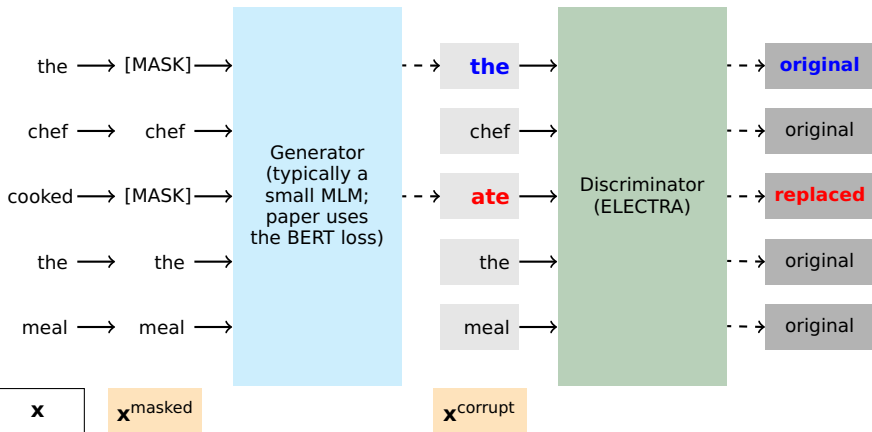
# Efficiency



Clark et al. 2019, Figure 3

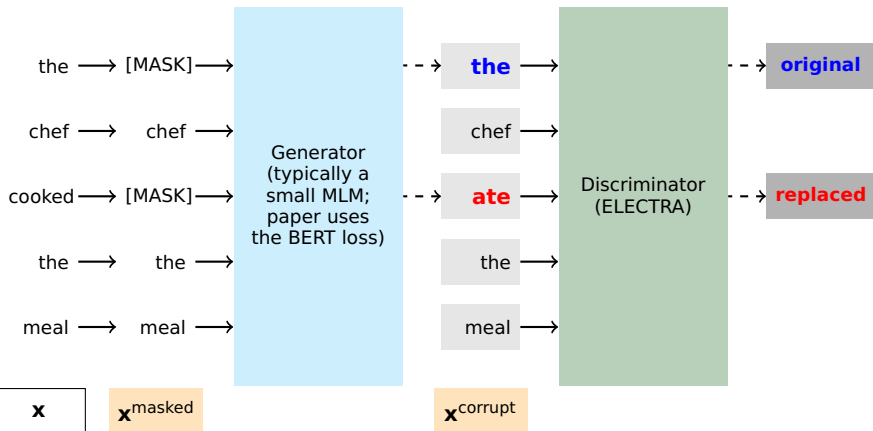
# ELECTRA efficiency analyses

## Full ELECTRA



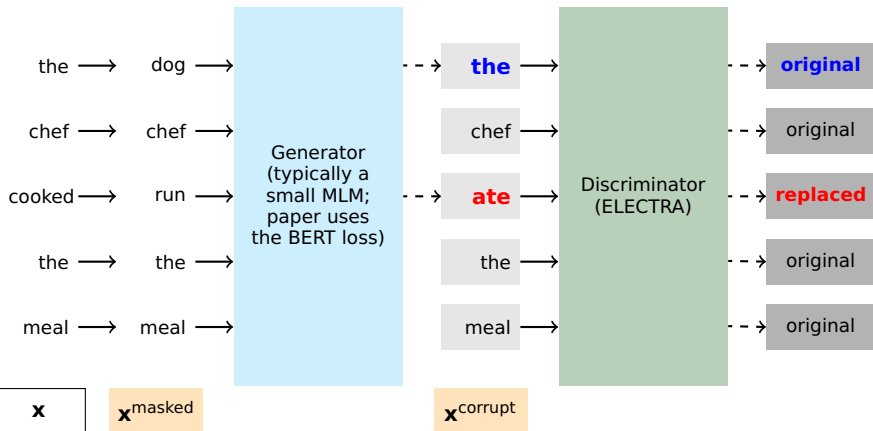
# ELECTRA efficiency analyses

## ELECTRA 15%



# ELECTRA efficiency analyses

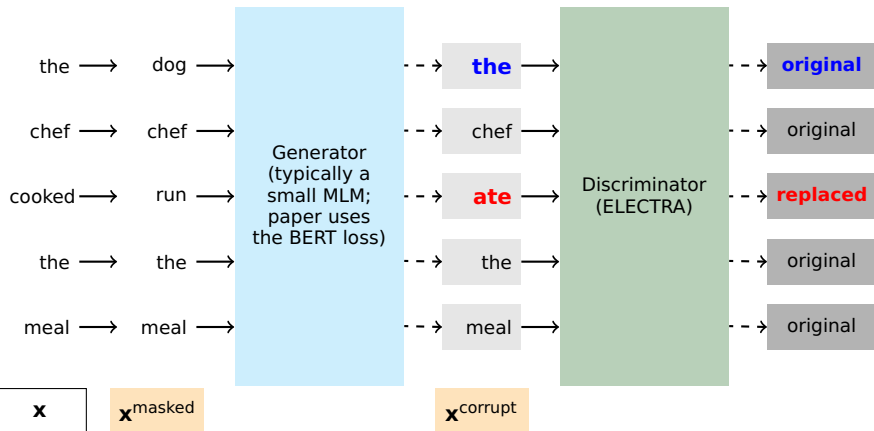
## Replace MLM





# ELECTRA efficiency analyses

## All-tokens MLM



# ELECTRA efficiency analyses

Model	GLUE score
<b>ELECTRA</b>	<b>85.0</b>
All-tokens MLM	84.3
Replace MLM	82.4
ELECTRA 15%	82.4
BERT	82.2

# ELECTRA model releases

Available from the [project site](#):

Model	Layers	Hidden Size	Params	GLUE test
Small	12	256	14M	77.4
Base	12	768	110M	82.7
Large	24	1024	335M	85.2

‘Small’ is the model designed to be “quickly trained on a single GPU”.

# References I

- Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. 2019. Electra: Pre-training text encoders as discriminators rather than generators. In *International Conference on Learning Representations*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. [XLNet: Generalized autoregressive pretraining for language understanding](#). In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 5753–5763. Curran Associates, Inc.