

– IN5550 –

Neural Methods in Natural Language Processing

Sustainability 2: Towards Green(er) NLP

Lilja Øvrelid

(With thanks to Jeremy Barnes)



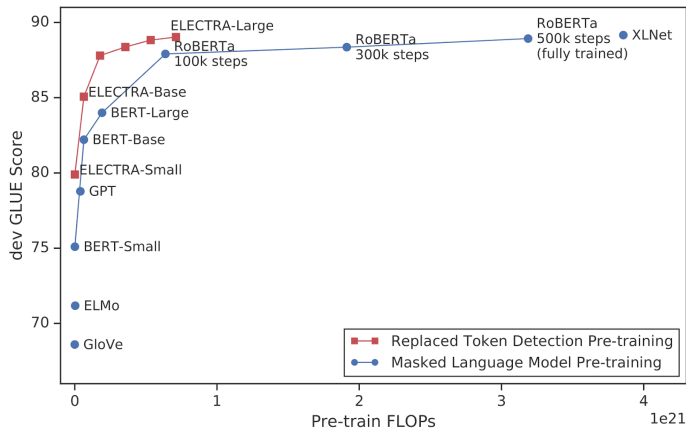


How can we mitigate the negative effects of large LMs?

- ▶ Enhance reporting of computational budgets
- ▶ Efficiency as a core evaluation metric

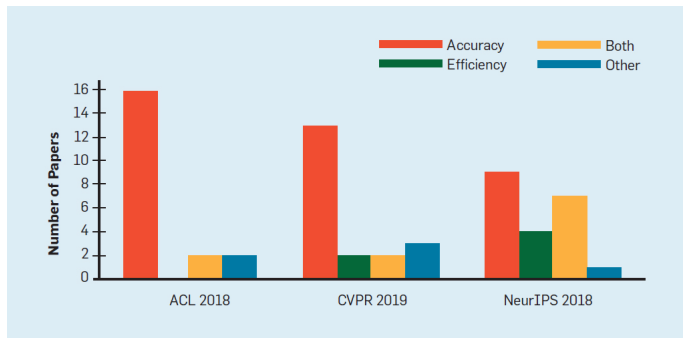
* from Schwartz et al (2020): Green AI

Improved reporting



* from Clark et al (2020): ELECTRA

Beyond accuracy



* from Schwartz et al (2020): Green AI



What is made more efficient?

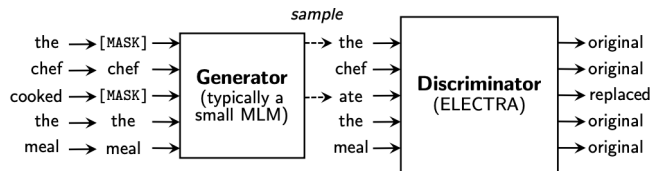
- ▶ Training
- ▶ Inference
- ▶ Model selection

How do we measure it?

- ▶ Space
- ▶ Time
- ▶ Energy

ELECTRA (Clark et al, 2020)

- ▶ Modifies pre-training objective
- ▶ Trained to distinguish "real" input tokens vs "fake" input tokens generated by another neural network
- ▶ Strong results even when trained on a single GPU.



* from Clark et al (2020): ELECTRA



To avoid retraining lots of models

- ▶ We can share the trained models
- ▶ Nordic Language Processing Laboratory (NLPL) is a good example
- ▶ Huggingface
- ▶ But it's important to get things right
 - ▶ METADATA!!!
 - ▶ same format for all models



What if we can reduce the size of these giant models?



What if we can reduce the size of these giant models?

- ▶ Often, overparameterized transformer models lead to better performance, even with less data



What if we can reduce the size of these giant models?

- ▶ Often, overparameterized transformer models lead to better performance, even with less data
- ▶ Lottery-ticket hypothesis: for large enough models, there is a small chance that random initialization will lead to a submodel that already has good weights for the task

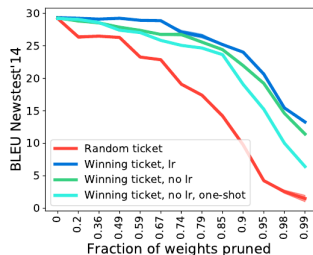


What if we can reduce the size of these giant models?

- ▶ Often, overparameterized transformer models lead to better performance, even with less data
- ▶ Lottery-ticket hypothesis: for large enough models, there is a small chance that random initialization will lead to a submodel that already has good weights for the task
- ▶ But interestingly, you can often remove a large number of the parameters for only a small decrease in performance

What if we can reduce the size of these giant models?

- ▶ Often, overparameterized transformer models lead to better performance, even with less data
- ▶ Lottery-ticket hypothesis: for large enough models, there is a small chance that random initialization will lead to a submodel that already has good weights for the task
- ▶ But interestingly, you can often remove a large number of the parameters for only a small decrease in performance



Reduce model size?



Reduce model size?

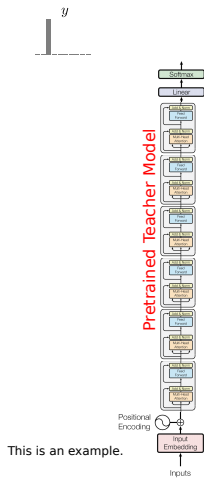


	Compression	Performance	Speedup	Model	Evaluation	
Distillation	BERT-base (Devlin et al., 2019)	×1	100%	×1	BERT ₁₂	All GLUE tasks, SQuAD
	BERT-small	×3.8	91%	–	BERT ₄ [†]	All GLUE tasks
	DistilBERT (Sanh et al., 2019)	×1.5	90% [‡]	×1.6	BERT ₆	All GLUE tasks, SQuAD
	BERT ₆ -PKD (Sun et al., 2019a)	×1.6	98%	×1.9	BERT ₆	No WNLI, CoLA, STS-B; RACE
	BERT ₃ -PKD (Sun et al., 2019a)	×2.4	92%	×3.7	BERT ₃	No WNLI, CoLA, STS-B; RACE
	Aguilar et al. (2019), Exp. 3	×1.6	93%	–	BERT ₆	CoLA, MRPC, QQP, RTE
	BERT-48 (Zhao et al., 2019)	×62	87%	×77	BERT ₁₂ ^{+†}	MNLI, MRPC, SST-2
	BERT-192 (Zhao et al., 2019)	×5.7	93%	×22	BERT ₁₂ ^{+†}	MNLI, MRPC, SST-2
	TinyBERT (Jiao et al., 2019)	×7.5	96%	×9.4	BERT ₄ [†]	No WNLI; SQuAD
	MobileBERT (Sun et al., 2020)	×4.3	100%	×4	BERT ₂₄ [†]	No WNLI; SQuAD
	PD (Turc et al., 2019)	×1.6	98%	×2.5 [‡]	BERT ₆ [†]	No WNLI, CoLA and STS-B
	WaLDORf (Tian et al., 2019)	×4.4	93%	×9	BERT ₈	SQuAD
	MiniLM (Wang et al., 2020b)	×1.65	99%	×2	BERT ₆	No WNLI, STS-B, MNLI _{mm} ; SQuAD
	MiniBERT (Tsai et al., 2019)	×6 ^{**}	98%	×27 ^{**}	mBERT ₃ [†]	CoNLL-18 POS and morphology
	BiLSTM-soft (Tang et al., 2019)	×110	91%	×434 [‡]	BiLSTM ₁	MNLI, QQP, SST-2
Quantization	Q-BERT-MP (Shen et al., 2019)	×13	98% [¶]	–	BERT ₁₂	MNLI, SST-2, CoNLL-03, SQuAD
	BERT-QAT (Zafrir et al., 2019)	×4	99%	–	BERT ₁₂	No WNLI, MNLI; SQuAD
	GOBO (Zadeh and Moshovos, 2020)	×9.8	99%	–	BERT ₁₂	MNLI
Pruning	McCarley et al. (2020), ff2	×2.2 [‡]	98% [‡]	×1.9 [‡]	BERT ₂₄	SQuAD, Natural Questions
	RPP (Guo et al., 2019)	×1.7 [‡]	99% [‡]	–	BERT ₂₄	No WNLI, STS-B; SQuAD
	Soft MvP (Sanh et al., 2020)	×33	94% [¶]	–	BERT ₁₂	MNLI, QQP, SQuAD
	IMP (Chen et al., 2020), rewind 50%	×1.4–2.5	94–100%	–	BERT ₁₂	No MNLI-mm; SQuAD
Other	ALBERT-base (Lan et al., 2020)	×9	97%	–	BERT ₁₂ [†]	MNLI, SST-2
	ALBERT-xxlarge (Lan et al., 2020)	×0.47	107%	–	BERT ₁₂ [†]	MNLI, SST-2
	BERT-of-Theseus (Xu et al., 2020)	×1.6	98%	×1.9	BERT ₆	No WNLI
	PoWER-BERT (Goyal et al., 2020)	N/A	99%	×2–4.5	BERT ₁₂	No WNLI; RACE

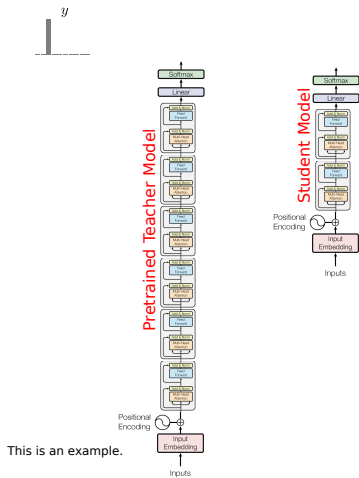
From Rogers et al. (2020) A Primer in BERTology: What We Know About How BERT Works.

Reduce model size?

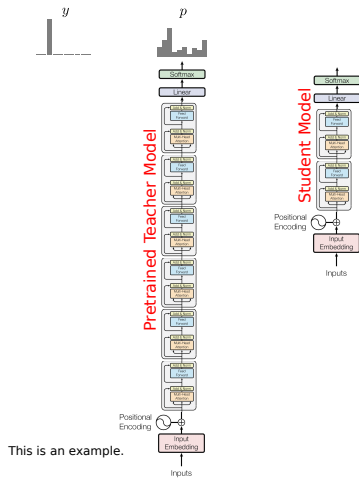




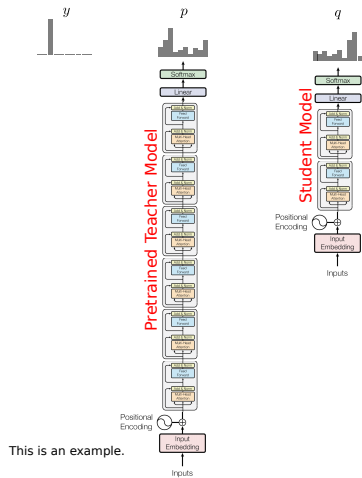
Model distillation



Model distillation



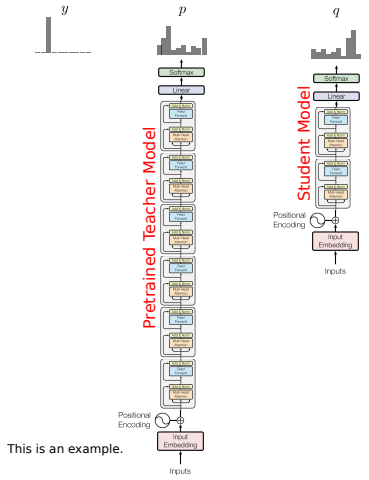
Model distillation



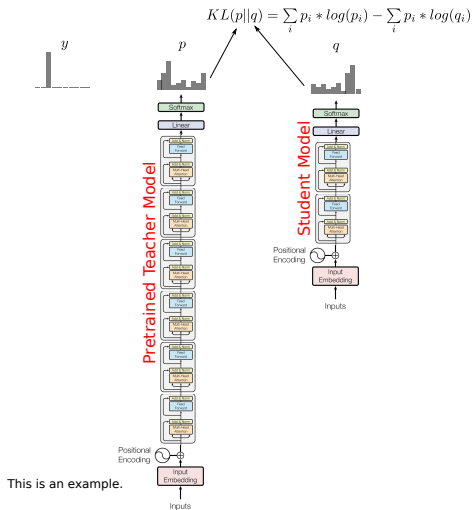
Model distillation



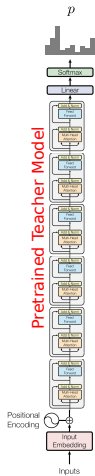
$$KL(p||q) = \sum_i p_i * \log(p_i) - \sum_i p_i * \log(q_i)$$



Model distillation



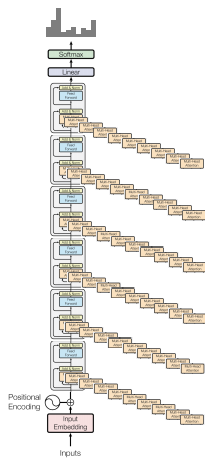
Head pruning



Head pruning



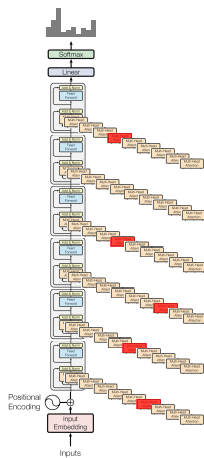
Performance = 93.7



Head pruning



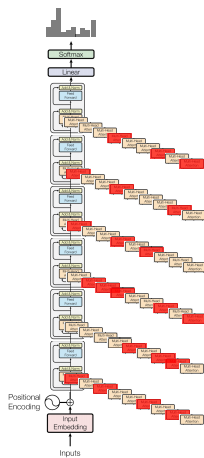
Performance = 93.7



Head pruning



Performance = 93.0



But how can we NLPers contribute to sustainability?

But how can we NLPers contribute to sustainability?

- ▶ When possible use pre-trained models

But how can we NLPers contribute to sustainability?

- ▶ When possible use pre-trained models
- ▶ If you train a strong model, similarly make it available to the community.

But how can we NLPers contribute to sustainability?

- ▶ When possible use pre-trained models
- ▶ If you train a strong model, similarly make it available to the community.
- ▶ Try to reduce the amount of hyperparameter tuning we do (for example, by working with models that are more robust to hyperparameters)

But how can we NLPers contribute to sustainability?

- ▶ When possible use pre-trained models
- ▶ If you train a strong model, similarly make it available to the community.
- ▶ Try to reduce the amount of hyperparameter tuning we do (for example, by working with models that are more robust to hyperparameters)
- ▶ Improved reporting, take into account efficiency

But how can we NLPers contribute to sustainability?

- ▶ When possible use pre-trained models
- ▶ If you train a strong model, similarly make it available to the community.
- ▶ Try to reduce the amount of hyperparameter tuning we do (for example, by working with models that are more robust to hyperparameters)
- ▶ Improved reporting, take into account efficiency
 - ▶ budget/performance curves

But how can we NLPers contribute to sustainability?

- ▶ When possible use pre-trained models
- ▶ If you train a strong model, similarly make it available to the community.
- ▶ Try to reduce the amount of hyperparameter tuning we do (for example, by working with models that are more robust to hyperparameters)
- ▶ Improved reporting, take into account efficiency
 - ▶ budget/performance curves

Model	Train / Infer FLOPs	Speedup	Params	Train Time + Hardware	GLUE
ELMo	3.3e18 / 2.6e10	19x / 1.2x	96M	14d on 3 GTX 1080 GPUs	71.2
GPT	4.0e19 / 3.0e10	1.6x / 0.97x	117M	25d on 8 P6000 GPUs	78.8
BERT-Small	1.4e18 / 3.7e9	45x / 8x	14M	4d on 1 V100 GPU	75.1
BERT-Base	6.4e19 / 2.9e10	1x / 1x	110M	4d on 16 TPUv3s	82.2
ELECTRA-Small	1.4e18 / 3.7e9	45x / 8x	14M	4d on 1 V100 GPU	79.9
50% trained	7.1e17 / 3.7e9	90x / 8x	14M	2d on 1 V100 GPU	79.0
25% trained	3.6e17 / 3.7e9	181x / 8x	14M	1d on 1 V100 GPU	77.7
12.5% trained	1.8e17 / 3.7e9	361x / 8x	14M	12h on 1 V100 GPU	76.0
6.25% trained	8.9e16 / 3.7e9	722x / 8x	14M	6h on 1 V100 GPU	74.1
ELECTRA-Base	6.4e19 / 2.9e10	1x / 1x	110M	4d on 16 TPUv3s	85.1