

Contextual word representations: Overview

Christopher Potts

Stanford Linguistics

CS224u: Natural language understanding



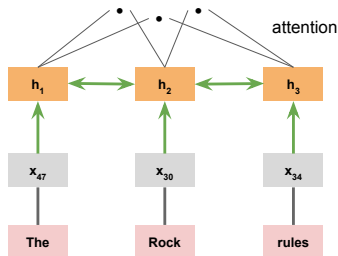
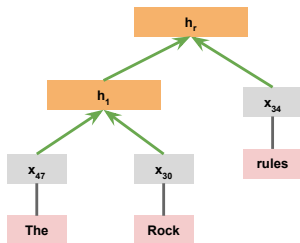
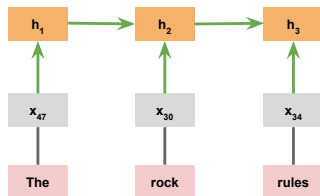
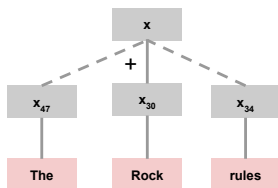
Associated materials

- Notebook: `finetuning.ipynb`
- Smith 2019
- Transformers
 1. Vaswani et al. 2017
 2. Alexander Rush: The Annotated Transformer [[link](#)]
- Hugging Face transformers: [project site](#)
- BERT: Devlin et al. 2019; [project site](#)
- RoBERTa: Liu et al. 2019; [project site](#)
- ELECTRA: Clark et al. 2019; [project site](#)

Word representations and context

- The vase broke.
 - Dawn broke.
 - The news broke.
 - Sandy broke the world record.
 - Sandy broke the law.
 - The burgler broke into the house.
 - The newscaster broke into the movie broadcast.
 - We broke even.
- flat tire/beer/note/surface
 - throw a party/fight/ball/fit
- A crane caught a fish.
 - A crane picked up the steel beam.
 - I saw a crane.
- Are there typos? I didn't see any.
 - Are there bookstores downtown? I didn't see any.

Model structure and linguistic structure



Guiding idea: Attention

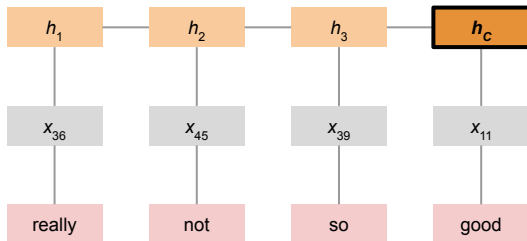
classifier $y = \mathbf{softmax}(\tilde{h}W + b)$

attention combo $\tilde{h} = \tanh([\kappa; h_C]W_\kappa)$

context $\kappa = \mathbf{mean}([\alpha_1 h_1, \alpha_2 h_2, \alpha_3 h_3])$

attention weights $\alpha = \mathbf{softmax}(\tilde{\alpha})$

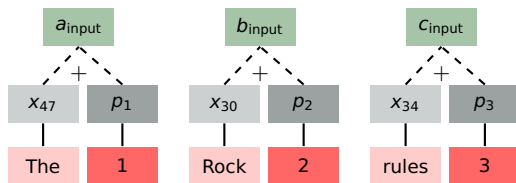
scores $\tilde{\alpha} = \begin{bmatrix} h_C^\top h_1 & h_C^\top h_2 & h_C^\top h_3 \end{bmatrix}$



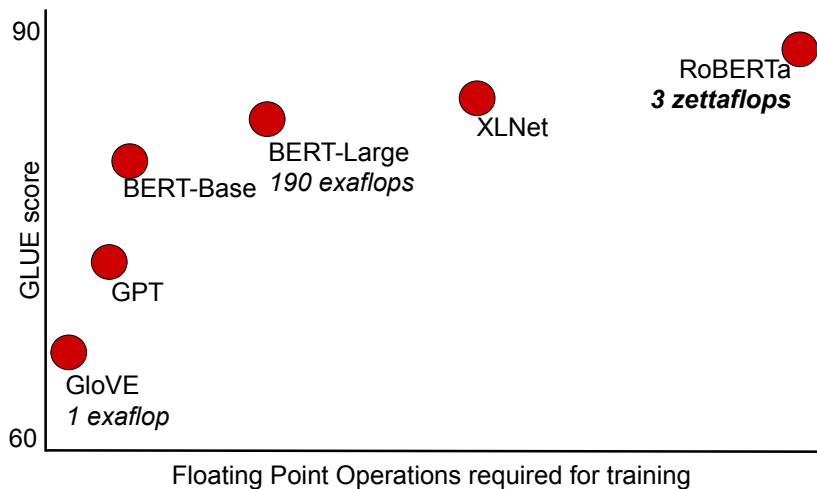
Guiding idea: Word pieces

```
[1]: from transformers import BertTokenizer
[2]: tokenizer = BertTokenizer.from_pretrained('bert-base-cased')
[3]: tokenizer.tokenize("This isn't too surprising.")
[3]: ['This', 'isn', "'", 't', 'too', 'surprising', '.']
[4]: tokenizer.tokenize("Encode me!")
[4]: ['En', '##code', 'me', '!']
[5]: tokenizer.tokenize("Snuffleupagus?")
[5]: ['S', '##nu', '##ffle', '##up', '##agu', '##s', '?']
[6]: tokenizer.vocab_size
[6]: 28996
```

Guiding idea: Positional encoding

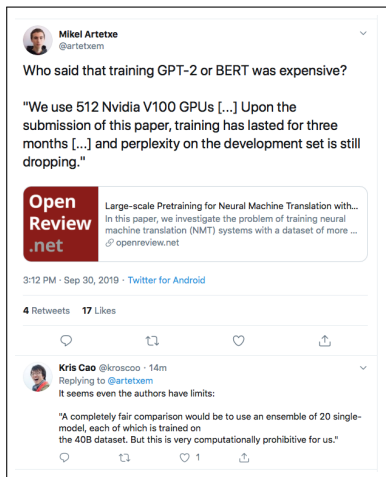


Current issues and efforts



Clark et al. 2019

Current issues and efforts



Mikel Artetxe @artetxem

Who said that training GPT-2 or BERT was expensive?

"We use 512 Nvidia V100 GPUs [...] Upon the submission of this paper, training has lasted for three months [...] and perplexity on the development set is still dropping."

Open Review .net

Large-scale Pretraining for Neural Machine Translation with...
In this paper, we investigate the problem of training neural machine translation (NMT) systems with a dataset of more ...
openreview.net

3:12 PM · Sep 30, 2019 · Twitter for Android

4 Retweets 17 Likes

Kris Cao @kroscoo · 14m
Replying to @artetxem
It seems even the authors have limits:

"A completely fair comparison would be to use an ensemble of 20 single-model, each of which is trained on the 40B dataset. But this is very computationally prohibitive for us."

<https://twitter.com/artetxem/status/1178794889229864962>

Current issues and efforts

Consumption	CO₂e (lbs)
Air travel, 1 person, NY↔SF	1984
Human life, avg, 1 year	11,023
American life, avg, 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000
Training one model (GPU)	
NLP pipeline (parsing, SRL)	39
w/ tuning & experiments	78,468
Transformer (big)	192
w/ neural arch. search	626,155

Table 1: Estimated CO₂ emissions from training common NLP models, compared to familiar consumption.¹

Current issues and efforts



Transformers

[Back to home](#)

All Models and checkpoints

Also check out our list of [Community contributors](#) and [Organizations](#)

Search models...

Tags: All

Sort: Default

Filter by model tags

All

PyTorch

TensorFlow

French 🇫🇷

German 🇩🇪

Dutch 🇳🇱

Italian 🇮🇹

Spanish 🇪🇸

Swedish 🇸🇪

Finnish 🇫🇮

Greek 🇬🇷

Turkish 🇹🇷

Arabic 🇸🇦

Chinese 🇨🇳

Malay 🇲🇾

Polish 🇵🇱

Esperanto

Multilingual 🌐

<https://huggingface.co>

Current issues and efforts

Compressing Large-Scale Transformer-Based Models: A Case Study on BERT

Prakhar Ganesh¹, Yao Chen¹, Xin Lou¹, Mohammad Ali Khan¹, Yin Yang²,
Deming Chen³, Marianne Winslett³, Hassan Sajjad^{4,2} and Preslav Nakov^{4,2}

¹Advanced Digital Sciences Center

²Hamad Bin Khalifa University

³University of Illinois at Urbana-Champaign

⁴Qatar Computing Research Institute

{prakhar.g, yao.chen, lou.xin, mohammad.k}@adsc-create.edu.sg,
{yyang, hsajjad, pnakov}@hbku.edu.qa, {dchen, winslett}@illinois.edu

Mitchell A. Gordon

About Blog Bookshelf

All The Ways You Can Compress BERT

Nov 18, 2019

Model compression reduces redundancy in a trained neural network. This is useful, since BERT barely fits on a GPU (BERT-Large does not) and definitely won't fit on your smart phone. Improved memory and inference speed efficiency can also save costs at scale.

<http://mitchgordon.me/>

Current issues and efforts

A Primer in BERTology: What we know about how BERT works

Anna Rogers, Olga Kovaleva, Anna Rumshisky

Department of Computer Science, University of Massachusetts Lowell
Lowell, MA 01854

{arogers, okovalev, arum}@cs.uml.edu

Some other Transformer-based models

- SBERT (**S**entence-**B**ERT; Reimers and Gurevych 2019)
- **G**enerative **P**re-trained **T**ransformer
 - ▶ GPT (Radford et al. 2018)
 - ▶ GPT-2 (Radford et al. 2019)
 - ▶ GPT-3 (Brown et al. 2020)
- XLNet (**X**tra **L**ong **T**ransformer: Yang et al. 2019)
- T5 (**T**ext-**T**o-**T**ext **T**ransfer **T**ransformer; Raffel et al. 2019)
- BART: Devlin et al. 2019

References I

- T. Brown, B. Mann, Nick Ryder, Melanie Subbiah, J. Kaplan, P. Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, G. Krüger, Tom Henighan, R. Child, Aditya Ramesh, D. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, E. Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, J. Clark, Christopher Berner, Sam McCandlish, A. Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. *ArXiv*, abs/2005.14165.
- 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.
- Prakhar Ganesh, Yao Chen, Xin Lou, Mohammad Ali Khan, Yin Yang, Deming Chen, Marianne Winslett, Hassan Sajjad, and Preslav Nakov. 2020. Compressing large-scale Transformer-based models: A case study on BERT. *ArXiv:2002.11985*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. ROBERTa: A robustly optimized BERT pretraining approach. *ArXiv:1907.11692*.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training. Ms, OpenAI.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.
- Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in bertology: What we know about how bert works. *ArXiv:2002.12327*.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Neural machine translation of rare words with subword units](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Noah A. Smith. 2019. Contextual word representations: A contextual introduction. *ArXiv:1902.06006v2*.
- Emma Strubell, Ananya Ganesh, and Andrew McCallum. 2019. [Energy and policy considerations for deep learning in NLP](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3645–3650, Florence, Italy. Association for Computational Linguistics.

References II

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc.
- 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.