Applied Deep Learning



More on Transformers



April 14th, 2020 http://adl.miulab.tw



2 Why Transformer?





- Apply neural architecture search (NAS) on Transformer architecture
- It also proved to be efficient at smaller sizes, achieving the same quality as the original "big" Transformer with 37.6% less parameters and outperforming the Transformer by 0.7 BLEU at a mobile-friendly model size of ~7M parameters.



4











• Google recently release a model "Meena" base on the ET.

Towards a Conversational Agent that Can Chat About... Anything

Tuesday, January 28, 2020

Posted by Daniel Adiwardana, Senior Research Engineer, and Thang Luong, Senior Research Scientist, Google Research, Brain Team



7 Though achieving great improvements...

- No Recurrent Inductive Bias: The Transformer trades the recurrent inductive bias of RNN's for parallelizability. However, the recurrent inductive bias appears to be crucial for generalizing on different sequence modeling tasks of varying complexity.
- For instance, when it is necessary to model the hierarchical structure of the input, or when the distribution of input length is different during training and inference, i.e. when good length generalization is needed.

No Recurrent Inductive Bias

- "The Importance of Being Recurrent for Modeling Hierarchical Structure"
- Transformer vs LSTM
- Substitution of Control of Con

```
(d(or f)) \Box (f(and a))
(d(and(c(or d)))) #(not f)
(not(d(or(f(or c))))) \Box (not(c(and(not d))))
```



No Recurrent Inductive Bias

- The experiments showed that LSTMs are more robust and generalize better.
- This does not imply that LSTMs should always be preferred over non-recurrent architectures.
- In fact, both FAN- and CNN-based networks have proved to perform comparably or better than LSTM-based ones on a very complex task like machine translation



10— The Transformer Is Not Turing Complete

- While the Transformer executes a total number of operations that scales with the input size, the number of sequential operations is constant and independent of the input size, determined solely by the number of layers.
- Solution Assuming finite precision, this means that the Transformer cannot be computationally universal.



1 The Transformer Is Not Turing Complete

An intuitive example are functions whose execution requires the sequential processing of each input element. In this case, for any given choice of depth T, one can construct an input sequence of length N > T that cannot be processed correctly by a Transformer.





"N" input symbols

12—Lack of Conditional Computation

- The Transformer applies the same amount of computation to all inputs (as well as all parts of a single input). However, not all inputs need the same amount of computation and this can be conditioned on the complexity of the input.
- "I arrived at the **bank** after crossing the **river**."



13— Universal Transformers

- A Concurrent-Recurrent Sequence Model
- The Universal Transformer is an extension to the Transformer models which combines the parallelizability and global receptive field of the Transformer model with the recurrent inductive bias of RNNs.
- Recurrence in depth







Positions





- modulates the number of computational steps needed to process each input symbol dynamically based on a scalar pondering value that is predicted by the model at each step.
- Ises an Adaptive Computation Time (ACT) mechanism, which was originally proposed for RNNs, to enable conditional computation.





RNN vs RNN with ACT



UT with Dynamic Halting



1

- weight sharing in depth leads to better performance of UTs (compared to the standard Transformer) on very small datasets and allows the UT to be a very data efficient model
- Transition functions can be replaced

Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context

• The original Transformer: no recurrence

19

- 2nd generation (Universal Transformer): recurrence in depth
 4
- 3rd generation (**Transformer-XL**): recurrence in length

20— Transformer for LM

In language modeling, Transformers are currently implemented with a fixed-length context, i.e. a long text sequence is truncated into fixed-length segments of a few hundred characters, and each segment is processed separately.

21— Transformer for LM

- This introduces two critical limitations:
- The algorithm is not able to model dependencies that are longer than a fixed length.
- The segments usually do not respect the sentence boundaries, resulting in context fragmentation which leads to inefficient optimization.

它不僅是一個能夠處理可變長度序列的模型,在多個任 務中刷新了當前的最好性能。

²²— Transformer-XL: Segment-level Recurrence

Ouring training, the representations computed for the previous segment are fixed and cached to be reused as an extended context when the model processes the next new segment.

23 Relative Positional Encodings

- Naively applying segment-level recurrence does not work, however, because the positional encodings are not coherent when we reuse the previous segments.
- For example, consider an old segment with contextual positions [0, 1, 2, 3]. When a new segment is processed, we have positions [0, 1, 2, 3, 0, 1, 2, 3] for the two segments combined, where the semantics of each position id is incoherent through out the sequence.
- Information based on content

- Transformer-XL is up to 1,800+ times faster than a vanilla Transformer during evaluation on language modeling tasks, because no re-computation is needed.
- Transformer-XL has better performance in perplexity (more accurate at predicting a sample) on long sequences because of long-term dependency modeling, and also on short sequences by resolving the context fragmentation problem.

20 Reformer: The Efficient Transformer

- extending Transformer to even larger context windows runs into limitations. The power of Transformer comes from attention, the process by which it considers all possible pairs of words within the context window to understand the connections between them.
- So, in the case of a text of 100K words, this would require assessment of 100K x 100K word pairs, or 10 billion pairs for each step, which is impractical.

27— Reformer: The Efficient Transformer

- Since softmax is dominated by the largest elements, for each query q_i we only need to focus on the keys in K that are closest to q_i. For example, if K is of length 64K, for each q_i we could only consider a small subset of, say, the 32 or 64 closest keys.
- That is much more efficient, but how can we find the nearest neighbors among the keys?

28— Locality Sensitive Hashing (LSH)

- The problem of finding nearest neighbors quickly in highdimensional spaces can be solved by locality-sensitive hashing (LSH).
- A hashing scheme that assigns each vector x to a hash h(x) is called locality-sensitive if nearby vectors get the same hash with high probability and distant ones do not.

 Attention is then applied within these much shorter chunks (and their adjoining neighbors to cover the overflow), greatly reducing the computational load.

30— The Memory Problem

- While LSH solves the problem with attention, there is still a memory issue. A single layer of a network often requires up to a few GB of memory and usually fits on a single GPU, so even a model with long sequences could be executed if it only had one layer.
- Sut when training a multi-layer model with gradient descent, activations from each layer need to be saved for use in the backward pass. A typical Transformer model has a dozen or more layers, so memory quickly runs out if used to cache values from each of those layers.

- "The Reversible Residual Network: Backpropagation Without Storing Activations"
- recompute the input of each layer on-demand during backpropagation, rather than storing it in memory.

recompute the input of each layer on-demand during backpropagation, rather than storing it in memory.

 $y_1 = x_1 + \mathcal{F}(x_2)$ $x_2 = y_2 - \mathcal{G}(y_1)$ $y_2 = x_2 + \mathcal{G}(y_1)$ $x_1 = y_1 - \mathcal{F}(x_2)$

Reversible Transformer: F becomes an attention layer while G becomes the feed-forward layer.

34 Reformer: The Efficient Transformer

- Combines the modeling capacity of a Transformer with an architecture that can be executed efficiently on long sequences and with small memory use even for models with a large number of layers.
- The Attention Problem: Locality-Sensitive Hashing
- The Memory Problem: **Reversible Layer**

- https://ai.googleblog.com/2020/01/towards-conversationalagent-that-can.html
- https://arxiv.org/abs/1901.11117
- https://arxiv.org/pdf/1803.03585.pdf
- https://mostafadehghani.com/2019/05/05/universaltransformers/
- https://arxiv.org/pdf/1807.03819.pdf
- https://arxiv.org/pdf/1603.08983.pdf

https://medium.com/@moocaholic/adaptive-computation time-act-in-neural-networks-part-1-2a28484b53df

- https://ai.googleblog.com/2020/01/towards-conversationalagent-that-can.html
- https://ai.googleblog.com/2019/01/transformer-xlunleashing-potential-of.html
- https://arxiv.org/pdf/2001.04451.pdf
- https://arxiv.org/pdf/1707.04585.pdf

