Attention & Transformer LMs



Timeline: Language Modeling and Vector Semantics



Timeline: Language Modeling and Vector Semantics

1913 Markov: Probability that next letter would be vowel or consonant.

1948



Timeline: Language Modeling and Vector Semantics

1913 Markov: Probability that next letter would be vowel or consonant.

Shannon: A Mathematical Theory of Communication (first digital language model)

Osgood: *The* Measurement of Meaning

1948

Switzer: Vector Space Models

Deerwater: Indexing by Latent Semantic Analysis (LSA) Bengio:

1980

Language Models **Vector Semantics**

LMs + Vectors

~logarithmic scale

natural language 2003 Blei et al.: [LDA Top 2010

Neural-net

embeddings

based

These (or similar) are Brown et al.: Class-based ngrai behind almost all state-of-the-art modern NLP systems

GPT

RoBERTA

GPT4.5

Mikolov: word2vec

ELMO 2018

Collobert and Weston: A unified architecture for natural language BERT processing: Deep neural networks...

Jelinek et al. (IBM): Language Models for Speech Recognition

Timeline: Language Modeling and Vector Semantics 1913 Markov: Probability that next letter would be vowel or consonant. 1948 Shannon: A Mathematical Theory of Communication (first digital language model) Jelinek et al. (IBM): Language Models f<u>or Speech Recognitio</u>n 1980 These (or similar) are Brown et al.: Class-based ngrai Osgood: *The* behind almost all Measurement **Robustly Optimized** state-of-the-art of Meaning **BERTransformers** modern NLP systems Deerwater: Pretraining Approch Switzer: Vector Mikolov: word2vec Indexing b Space Models **Generative Pretrained** Semantic. (LSA) Transformers anc GPT Bengio: Weston: A unified Language Models RoBERTA tecture for Vector Semantics **Bidirectional Transformers** BERT LMs + Vectors anunyuuge embeddings processing: Deep ~logarithmic scale neural networks... GPT4.5

Transformers

(Advances in neural information processing systems, 2017)



Transformers

(Advances in neural information processing systems, 2017)

Attention is all you need

A Vaswani, N Shazeer, N Parmar... - Advances in neural ..., 2017 - proceedings.neurips.cc

... to attend to all positions in the decoder up to and including that position. We need to prevent

... We implement this inside of scaled dot-product attention by masking out (setting to -∞) ...

☆ Save 57 Cite Cited by 173107 Related articles All 73 versions ≫

Attention Is All You Need Jakob Uszkoreit* Google Research usz@google.com Niki Parmar* Google Research nikip@google.com Noam Shazeer* Google Brain Lukasz Kaiser* noam@google.com Ashish Vaswani* Google Brain lukaszkaiser@google.com Google Brain avaswani@google.com Aidan N. Gomez* University of Toronto aidan@cs.toronto.edu Llion Jones* Google Research Illia Polosukhin* ‡ illia.polosukhin@gmail.com llion@google.com on models are based on complex recurrent or include an encoder and a decoder. The best Abstract he encoder and decoder through an attention simple network architecture, the Transformer, sms, dispensing with recurrence and convolutions nachine translation tasks show these models to ing more parallelizable and requiring significantly achieves 28.4 BLEU on the WMT 2014 Englishnproving over the existing best results, including transtation approving over the existing test results, including task, we single model state-of-the-art BLEU score of 41.0 after urous survey and the training costs of the

Corman translatio less time





- self-attention multi-headed attention positional embeddings residual links
- intuition from translation
 key, queries, values
 similarity score

Attention: Motivated from Translation

As an optimization problem (Eisenstein, 2018):

$$\hat{\boldsymbol{w}}^{(t)} = \operatorname*{argmax}_{\boldsymbol{w}^{(t)}} \Psi(\boldsymbol{w}^{(s)}, \boldsymbol{w}^{(t)})$$

Why?

- \$40billion/year industry
- A center piece of many genres of science fiction
- A fairly "universal" problem:
 - Language understanding
 - Language generation
- Societal benefits of intercultural communication

Why?

- \$40billion/year industry
- A center piece of many genres of science fiction
- A fairly "universal" problem:
 - Language understanding
 - Language generation
- Societal benefits of intercultural communication



Why?

- \$40billion/year industry
- A center piece of many genres of science fiction
- A fairly "universal" problem:
 - Language understanding
 - Language generation
- Societal benefits of intercultural communication



(Douglas Adams)

Why Neural Network Approach works? (Manning, 2018)

- Joint end-to-end training: learning all parameters at once.
- Exploiting distributed representations (embeddings)
- Exploiting variable-length context
- High quality generation from deep decoders stronger language models (even when wrong, make sense)

Recurrent Neural Network



Recurrent Neural Network



RNN: Encoder







essentially a language model conditioned on the final state from the encoder.





Encoder-Decoder

Challenge:

The ball was kicked by kayla.

• Long distance dependency when translating:



Kayla kicked the ball.

Encoder-Decoder

Challenge:

The ball was kicked by kayla.

• Long distance dependency when translating:



Long Distance / Out of order dependencies

 $oldsymbol{h}_{m-1}^{(s,2)}$

 $lackslash h_{m-1}^{(s,1)}$.

 $oldsymbol{x}_{m-1}^{(s)}$

















 $\overline{\alpha_{h_i \to s}} = \operatorname{softmax}(\psi(h_i, s))$

$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \to s_n} s_n$$



$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \to s_n} s_n$$



$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \to s_n} s_n$$



$$\alpha_{h_i \to s} = \operatorname{softmax}(\psi(h_i, s))$$

$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \to s_n} s_n$$



Iternative Scoring Functions

$$\psi_{add}(h_i, s) = v^T \tanh(W_h h_i + W_s s])$$

 $\psi_{dp}(h_i, s) = s^T h_i$
 $\psi_{mult}(h_i, s) = s^T W h_i$


If variables are standardized, matrix multiply produces a similarity score.

Alternative Scoring Function $\psi_{add}(h_i, s) = v^T \tan(W_h h_i + W_s s])$ $\psi_{dp}(h_i, s) = s^T h_i$ $\psi_{mult}(h_i, s) = s^T W h_i$



$$\alpha_{h_i \to s} = \operatorname{softmax}(\psi(h_i, s))$$

$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \to s_n} s_n$$



A useful abstraction is to make the vector attended to (the "value vector", Z) separate than the "key vector" (s).

$$\psi_{mult}(h_i, s) = s^T W h_i$$

$$\alpha_{h_i \to s} = \operatorname{softmax}(\psi(h_i, s))$$

$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \to s_n} z_n$$



A useful abstraction is to make the vector attended to (the "value vector", Z) separate than the "key vector" (s).

$$\psi_{mult}(h_i, s) = s^T W h_i$$

$$\alpha_{h_i \to s} = \operatorname{softmax}(\psi(h_i, s))$$

$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \to s_n} z_n$$



Attention as weighting a value based on a query and key:



(Eisenstein, 2018)

















Evolution of Sequence Modeling

RNNs -> GRU/LSTMs-RNN <u>-> LSTMS with Attention</u> -> Attention without RNN



- self-attention multi-headed attention positional embeddings residual links
- intuition from translation key, queries, values similarity score











A weighted combination of other words' vectors.



The Transformer's Heart: Self-Attention







scaling parameter $\psi_{dp}(q,k) = (qk^t)\sigma$

RNN Limitation: Losing Track of Long Distance Dependencies

The horse which was raced past the barn tripped .



RNN Limitation: Losing Track of Long Distance Dependencies





RNN: Limitation: Not parallelizable



Contextual Word Vectors

Person A

How are you?

I feel *fine* –even *great*!

Person B

My life is a *great* mess! I'm having a very hard time being happy.

What is going on? Earlier, I played the game Yahtzee with my partner. I could not get that die to roll a 1! Now I'm lying on my bed for a rest.

My business *partner* was *lying* to me. He was trying to *game* the system and *played* me. I think I am going to *die* –he left and now I have to pay the *rest* of his *fine*.

RNN: Limitation: Not parallelizable



The Transformer: Motivation

- Capture long-distance dependencies
- Preserving sequential distances / periodicity
- Capture multiple relationships
- Easy to parallelize -- don't need sequential processing.

Introducing the Transformer

Attention Is All You Need			
Ashish Vaswani* Google Brain avaswani@google.com	Noam Shazeer* Google Brain noam@google.com	Niki Parmar* Google Research nikip@google.com	Jakob Uszkoreiť Google Research usz@google.com
Llion Jones*	Aidan N. Gomez*	† Łuka	asz Kaiser*
llion@google.com	aidan@cs.toronto.	edu lukaszkai:	ser@google.com
	illia.polosukhin Abstrac	@gmail.com	
The dominant seque convolutional neura performing models mechanism. We pr based solely on atten entirely. Experime be superior in qualif less time to train. O to-German translati ensembles, by over 2 our model establishe training for 3.5 day	ence transduction model al networks that include a also connect the encode opose a new simple net tion mechanisms, dispens nts on two machine tran y while being more para Dur model achieves 28.4 on task, improving over 2 BLEU. On the WMT 20 is a new single-model stat s on eight GPUs, a small o literature	s are based on comple an encoder and a deco er and decoder through work architecture, the sing with recurrence and islation tasks show the llelizable and requiring BLEU on the WMT 3 the existing best resu 14 English-to-French tr e-of-the-art BLEU scor fraction of the trainin	x recurrent or der. The best h an attention Transformer, d convolutions ese models to g significantly 2014 English- ilts, including anslation task, e of 41.0 after g costs of the



Introducing the Transformer



Encoder



Encoder: Input Embedding



Input Embedding

Original Sentence

Tokenization

Input IDs (embedding lookup: position in the vocab -FIXED)

Embeddings (vector of size d_{model}= 512 or 1024 or ... LEARNED)
Encoder: Positional Encoding



Positional Encoding

Original Sentence (tokens)

e.g.

Embeddings (vector of size d_{model}= 512 or 1024 or ... Learned)

$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{model}})$$

Positional Embedding

(vector of size d_{model}= 512 or 1024 or … Can be Learned or Flxed)

Encoder: Multi-Head Attention



Specs for Transformer's Self-Attention



Specs for Transformer's Self-Attention



Specs for Transformer's Self-Attention

$egin{attn_head} \operatorname{attn}(W_q^th,W_k^th,W_v^th) \ & ext{attn}(q,k,v) = \operatorname{softmax}\left(rac{qk^t}{\sqrt{d_k}} ight)v \end{array}$

Two Additions:

1. Linearly parameterized Q, K, and V:

•
$$q = W_{q}^{T}b_{i}$$

• $k = W_{k}^{T}b_{i}$
• $v = W_{v}^{T}b_{i}$

Scaling parameter for score
$$\psi_{dp}(q,k) = (qk^t) \stackrel{\searrow}{\sigma},$$
where $\sigma = rac{1}{\sqrt{d_k}}$ $\psi_{dp}(q,k) = rac{(qk^t)}{\sqrt{d}}$

Self-Attention in PyTorch

```
import nn.functional as f
class SelfAttention(nn.Module):
    def __init__(self, h_dim:int):
        self.Q = nn.Linear(h_dim, h_dim) #1 head
        self.K = nn.Linear(h_dim, h_dim)
        self.V = nn.Linear(h_dim, h_dim)
```

```
def forward(hidden_states:torch.Tensor):
    v = self.V(hidden_states)
    k = self.K(hidden_states)
    q = self.Q(hidden_states)
    attn_scores = torch.matmul(q, k.T)
    attn_probs = f.Softmax(attn_scores)
```

```
context = torch.matmul(attn_probs, v)
return context
```

```
\psi_{dp}(q,k) = (qk^t) \sigma
```

Linear layer: $W^T X$

One set of weights for each of K, Q, and V

Self-Attention in PyTorch

```
import nn.functional as f
class SelfAttention(nn.Module):
    def init (self, h dim:int):
        self.Q = nn.Linear(h dim, h dim) #1 head
        self.K = nn.Linear(h dim, h dim)
        self.V = nn.Linear(h dim, h dim)
        self.dropout = nn.dropout(p=0.1)
    def forward(hidden states:torch.Tensor):
        v = self.V(hidden states)
        k = self.K(hidden_states)
        q = self.Q(hidden_states)
        attn scores = torch.matmul(q, k.T)
        attn probs = f.Softmax(attn scores)
        attn probs = self.dropout(attn probs)
        context = torch.matmul(attn probs, v)
        return context
```

```
\psi_{dp}(\mathbf{q},k) = (qk^t)\sigma
```

Linear layer: $W^T X$

One set of weights for each of K, Q, and V

Self-Attention in PyTorch



Multi-headed Attention

Single-headed (standard self-attention)

Limitation (thus far): Can't capture multiple types of dependencies between words.



Multi-headed Attention

Limitation (thus far): Can't capture multiple types of dependencies between words. Solution: Multi-head attention



The Transformer: Multi-headed Attention



Transformer's Multi-headed Attention



Transformer's Residual Links



Transformer's Residual Links

Heat matrix of attention weights from one layer (rows) to the next (cols).



without residuals

Output=LayerNorm(h+MHAtten(h))





with residuals















Timeline: Language Modeling and Vector Semantics

1913 Markov: Probability that next letter would be vowel or consonant.

> 1948 Shannon: A Mathematical Theory of Communication (first digital language model)

Osgood: *The* Measurement of Meaning

> Switzer: Vector Space Models

Language Models **Vector Semantics**

LMs + Vectors

~logarithmic scale



Timeline: Language Modeling and Vector Semantics

1913 Markov: Probability that next letter would be vowel or consonant.

> 1948 Shannon: A Mathematical Theory of Communication (first digital language model)

Osgood: *The* Measurement of Meaning

> Switzer: Vector Space Models

Language Models **Vector Semantics**

LMs + Vectors

~logarithmic scale



BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova Google AI Language {jacobdevlin,mingweichang,kentonl,kristout}@google.com

Abstract

We introduce a new language representation model called **BERT**, which stands for Bidirectional Encoder Representations from Transformers. Unlike recent language representation models (Peters et al., 2018a; Radford et al., 2018), BERT is designed to pretrain deep bidirectional representations from unlabeled text by jointly conditioning on both left and right context in all layers. As a result, the pre-trained BERT model can be finetuned with just one additional output layer to create state-of-the-art models for a wide range of tasks, such as question answering and language inference, without substantial taskspecific architecture modifications.

BERT is conceptually simple and empirically powerful. It obtains new state-of-the-art results on eleven natural language processing tasks, including pushing the GLUE score to 80.5% (7.7% point absolute improvement), MultiNLI accuracy to 86.7% (4.6% absolute improvement), SQuAD v1.1 question answering Test F1 to 93.2 (1.5 point absolute improvement) and SQuAD v2.0 Test F1 to 83.1 (5.1 point absolute improvement).

1 Introduction

Language model pre-training has been shown to be effective for improving many natural language processing tasks (Dai and Le, 2015; Peters et al., 2018a; Radford et al., 2018; Howard and Ruder, 2018). These include sentence-level tasks such as natural language inference (Bowman et al., 2015; Williams et al., 2018) and paraphrasing (Dolan and Brockett, 2005), which aim to predict the relationships between sentences by analyzing them holistically, as well as token-level tasks such as named entity recognition and question answering, where models are required to produce fine-grained output at the token level (Tjong Kim Sang and De Meulder, 2003; Rajpurkar et al., 2016). There are two existing strategies for applying pre-trained language representations to downstream tasks: feature-based and fine-tuning. The feature-based approach, such as ELMO (Peters et al., 2018a), uses task-specific architectures that include the pre-trained representations as additional features. The fine-tuning approach, such as the Generative Pre-trained Transformer (OpenAI GPT) (Radford et al., 2018), introduces minimal task-specific parameters, and is trained on the downstream tasks by simply fine-tuning all pretrained parameters. The two approaches share the same objective function during pre-training, where they use unidirectional language models to learn general language representations.

We argue that current techniques restrict the power of the pre-trained representations, especially for the fine-tuning approaches. The major limitation is that standard language models are unidirectional, and this limits the choice of architectures that can be used during pre-training. For example, in OpenAI GPT, the authors use a left-toright architecture, where every token can only at tend to previous tokens in the self-attention layers of the Transformer (Vaswani et al., 2017). Such restrictions are sub-optimal for sentence-level tasks, and could be very harmful when applying finetuning based approaches to token-level tasks such as question answering, where it is crucial to incorporate context from both directions.

In this paper, we improve the fine-tuning based approaches by proposing BERT: Bidirectional Encoder Representations from Transformers. BERT alleviates the previously mentioned unidirectionality constraint by using a "masked language model" (MLM) pre-training objective, inspired by the Cloze task (Taylor, 1953). The masked language model randomly masks some of the tokens from the input, and the objective is to predict the original vocabulary id of the masked

4171

Proceedings of NAACL-HLT 2019, pages 4171–4186 Minneapolis, Minnesota, June 2 - June 7, 2019. ©2019 Association for Computational Linguistics

LMs + Vectors

~logarithmic scale

Ietter would be vowel or consonant. *atical Theory of Communication* (first digital language model) *atical Theory of Communication* (for the state of the stat

tent

vsis

Bengio:

based

Neural-net

embeddings

Brown et al.: *Class-based ngram models of* 2003 *natural language*

Modeling and **Vector Semantics**

Blei et al.: [LDA Topic Modeling] 2010 Mikolov: word2vec 2018 Collobert and

Weston: A unified architecture for natural language BERT processing: Deep neural networks...

XLNet RoBERTA

BERT Rediscovers the Classical NLP Pipeline

Ian Tenney¹ Dipanjan Das¹ Ellie Pavlick^{1,2} ¹Google Research ²Brown University {iftenney, dipanjand, epavlick}@google.com

Abstract

Pre-trained text encoders have rapidly advanced the state of the art on many NLP tasks. We focus on one such model, BERT, and aim to quantify where linguistic information is captured within the network. We find that the model represents the steps of the traditional NLP pipeline in an interpretable and localizable way, and that the regions responsible for each step appear in the expected sequence: POS tagging, parsing, NER, semantic roles, then coreference. Qualitative analysis reveals that the model can and often does adjust this pipeline dynamically, revising lowerlevel decisions on the basis of disambiguating information from higher-level representations. of the network directly, to assess whether there exist localizable regions associated with distinct types of linguistic decisions. Such work has produced evidence that deep language models can encode a range of syntactic and semantic information (e.g. Shi et al., 2016; Belinkov, 2018; Tenney et al., 2019), and that more complex structures are represented hierarchically in the higher layers of the model (Peters et al., 2018b; Blevins et al., 2018).

We build on this latter line of work, focusing on the BERT model (Devlin et al., 2019), and use a suite of probing tasks (Tenney et al., 2019) derived from the traditional NLP pipeline to quantify where specific types of linguistic information are

Bengio:

based

Neural-net

embeddings

lationships between sentences by analyzing then holistically, as well as token-level tasks such a named entity recognition and question answering where models are required to produce fine-grainer output at the token level (Tjong Kim Sang an De Meulder, 2003; Raiourkar et al., 2016).

rectionality constraint by using a "masked lan guage model" (MLM) pre-training objective, in spired by the Cloze task (Taylor, 1953). Th masked language model randomly masks some or the tokens from the input, and the objective is to predict the original vocabulary id of the maske.

171

Proceedings of NAACL-HLT 2019, pages 4171–4186 finneapolis, Minnesota, June 2 - June 7, 2019. ©2019 Association for Computational Linguistics

LMs + Vectors

logarithmic scale

and Vector Semantics

wel or consonant.

ommunication (first digital language model) anguage Models for Speech Recognition

lass-based ngram models of atural language

et al.: [LDA Topic Modeling]

2010 Mikolov: *word2vec*

2018

Collobert and Weston: A unified architecture for natural language BERT processing: Deep neural networks...

XLNet RoBERTA

GPT3

Jacob Devlin N

Abstract

Ve introduce a new languag ion model called BERT, whi Bidirectional Encoder Represe Transformers. Unlike recent la entation models (Peters et al. ord et al., 2018), BERT is de rain deep bidirectional represe nlabeled text by jointly conditi eft and right context in all lay ult, the pre-trained BERT mod uned with just one additional o create state-of-the-art mode ange of tasks, such as question. anguage inference, without su

SERT is conceptually simple an ownerful. It obtains new state alts on eleven natural langua asks, including pushing the G MultiNLI accuracy to 86.7% (mprovement), SQuAD v1.1 (mprovement) and SQuAD v2.0 T orverment) and SQuAD v2.0 T

Introduction

anguage model pre-training ha e effective for improving many rocessing tasks (Dai and Le, 2C 018a; Radford et al., 2018; Ho 018). These include sentence-le atural language inference (Bow Villiams et al., 2018) and para di Brockett, 2005), which aim ationships between sentences by olistically, as well as token-leu ande antity recognition and que here models are required to proutput at the token level (Tjon uptut at the token level (Tjon uptut at the token level (Tjon

> Pro Minneapolis, Minnesota, J

Journalism Quarterly

DEVOTED TO RESEARCH STUDIES IN THE FIELD OF MASS COMMUNICATIONS

FALL 1953

"Cloze Procedure": A New Tool For Measuring Readability

BY WILSON L. TAYLOR*

Here is the first comprehensive statement of a research method and its theory which were introduced briefly during a workshop at the 1953 AEJ convention. Included are findings from three pilot studies and two experiments in which "cloze procedure" results are compared with those of two readability formulas.

"CLOZE PROCEDURE" IS A NEW PSYchological tool for measuring the effectiveness of communication. The method is straightforward; the data are easily quantifiable; the findings seem to stand up.

At the outset, this tool was looked on as a new approach to "readability." It was so used in three pilot studies and two experiments, the main findings of which are reported here.

*The writer is particularly obligated to Prof. Charles E. Osgood, University of Illinois, and Melvin R. Marks, Personnel Research Section, A.G.O., Department of the Army, for instigating and assisting in the series of efforts that yielded the notion of "cloze procedure." Both are experimental psychologists. Among others who have advised, encouraged or otherwise aided are these of the University of Illinois: Prof. Lee J. Cronbach, educational psychologist and statistician; Dean Wilbur Schramm, Division of Communications; Prof. Charles E. Swanson, Institute of Communications Research, and George R. Klare, psychologist, both of whom have authored articles on readability; and several journalism teachers who lent their classes. Kalmer E. Stordahl and Clifford M. Christensen, until recently research associates of the Institute, also contributed. First, the results of the new method were repeatedly shown to conform with the results of the Flesch and Dale-Chall devices for estimating readability. Then the scope broadened, and cloze procedure was pitted against those standard formulas.

If future research substantiates the results so far, this tool seems likely to have a variety of applications, both theoretical and practical, in other fields involving communication functions.

THE "CLOZE UNIT"

At the heart of the procedure is a functional unit of measurement tentatively dubbed a "cloze." It is pronounced like the verb "close" and is derived from "closure." The last term is one gestalt psychology applies to the human tendency to complete a familiar but not-quite-finished pattern—to "see" a broken circle as a whole one, for example, by mentally closing up the gaps.

embeddings

415

ng and Vector Semantics

e vowel or consonant.

of Communication (first digital language model)

1): Language Models for Speech Recognition

I.: Class-based ngram models of natural language

Blei et al.: [LDA Topic Modeling]

2010 Mikolov: word2vec 2018 Collobert and Weston: A unified architecture for natural language processing: Deep neural networks...

[~]logarithmic scale

Task: Estimate $P(w_i | w_1, ..., w_{i-1}, w_{i+1}, ..., w_n)$:P(masked word given context)P(with | He ate the cake < M > the fork) = ?

Task: Estimate $P(w_i | w_1, ..., w_{i-1}, w_{i+1}, ..., w_n)$:P(masked word given context)P(with | He ate the cake < M > the fork) = ?



Task: Estimate $P(w_i | w_1, ..., w_{i-1}, w_{i+1}, ..., w_n)$:P(masked word given context)P(with | He ate the cake < M > the fork) = ?



Task: Estimate $P(w_i | w_1, ..., w_{i-1}, w_{i+1}, ..., w_n)$:P(masked word given context)P(with | He ate the cake < M > the fork) = ?

Sequence (He, at, the, cake,<MASK>, the, fork)



What is the masked word in the sequence?



Transformer Language Models: Uses multiple layers of a transformer



layer k: (used for language modeling)

layer k-1: (taken as contextual embedding)

layers 1 to k-2: (compose embeddings with context)

layer 0: (input: word-type embeddings)

sentence (sequence) input:

(Kjell, Kjell, and Schwartz, 2023)

<u>Auto-encoder (MLM):</u>

- Connections go both directions.
- Task is predict word in middle: p(wi| ..., pwi-2, wi-1, wi+1, wi+2...)
- Better for:
 - \circ embeddings
 - fine-tuning (transfer learning)



<u>Auto-encoder (MLM):</u>

- Connections go both directions.
- Task is predict word in middle: p(wi| ..., pwi-2, wi-1, wi+1, wi+2...)
- Better for:
 - \circ embeddings
 - fine-tuning (transfer learning)

<u>Auto-regressor</u> (generator):

- Connections go forward only
- Task is predict word next word: p(wi| wi-1, wi-2, ...)
- Better for:
 - generating text
 - zero-shot learning





<u> Auto-encoder (MLM):</u>

- Connections go both directions.
- Task is predict word in middle: p(wi| ..., pwi-2, wi-1, wi+1, wi+2...)
- Better for:
 - embeddings
 - fine-tuning (transfer learning)

Auto-regressor (generator):

- Connections go forward only
- Task is predict word next word: p(wi| wi-1, wi-2, ...)
- Better for:
 - generating text
 - zero-shot learning







Bert: Attention by Layers

https://colab.research.google.com/drive/1vIOJ1IhdujVjfH857hvYKIdKPTD9Kid8

(Vig, 2019)
Hugging Face or AllenNLP

https://github.com/huggingface/transformers

```
#example for getting embeddings
from transformers import BertModel, PreTrainedTokenizerFast, pipeline
```

```
bert_tokenizer = PreTrainedTokenizerFast.from_pretrained('google-bert/bert-base-uncased')
bert_model = BertModel.from_pretrained('google-bert/bert-base-uncased')
pipe = pipeline('feature-extraction', model=bert_model, tokenizer=bert_tokenizer)
emb = pipe(text)
print(emb[0][0])
```

https://docs.allennlp.org/v2.10.1/api/modules/transformer/transformer_module/

Transformer (as of 2017)

"WMT-2014" Data Set. BLEU scores:



Transformers as of 2024

General Language Understanding Evaluations:

https://gluebenchmark.com/leaderboard

https://super.gluebenchmark.com/leaderboard/

ChatGPT

B

ChatGPT is an artificial intelligence chatbot developed by OpenAI and launched in November 2022. It is built on top of OpenAI's GPT-3.5 and GPT-4 families of large language models and has been fine-tu...



Transformers as of 2023



BERT Performance: e.g. Question Answering

GLUE scores evolution over 2018-2019



https://rajpurkar.github.io/SQuAD-explorer/

The Transformer: Take Away

Challenges to sequential representation learning

- Capture long-distance dependencies Self-attention treats far away words similar to those close.
- Preserving sequential distances / periodicity
 Positional embeddings encode distances/periods.
- Capture multiple relationships *Multi-headed attention enables multiple compositions*.
- Easy to parallelize -- don't need sequential processing. Entire layer can be computed at once. Is only matrix multiplications + standardizing.

Part 3: Applying Transformer LMs

Foundational

Applied









Large Training Corpus

softmax for LM:

layer k: (used for language modeling)

layer k-1: (taken as contextual embedding)

layers 1 to k-2: (compose embeddings with context)

layer 0: (input: word-type embeddings)

sentence (sequence) input:



softmax for LM:

layer k: (used for language modeling)

layer k-1: (taken as contextual embedding)

layers 1 to k-2: (compose embeddings with context)

layer 0: (input: word-type embeddings)

sentence (sequence) input:

New Continued Training Corpus



softmax for LM:

layer k: (used for language modeling)

layer k-1: (taken as contextual embedding)

layers 1 to k-2: (compose embeddings with context)

layer 0: (input: word-type embeddings)

sentence (sequence) input:

Task Promptse.g. What topic is this about? "Last night, the
Seawolves won the game." answer: sports





Task Fine-Tuning





Large Training Corpus

(used for language modeling)

SV

TIAX TOT

layer k-1: (taken as contextual embedding)

layers 1 to k-2: (compose embeddings with context)

layer 0: (input: word-type embeddings)

sentence (sequence) input:

Task Fine-Tuning

classifier or regressor: (e.g. sentiment, topic classification, etc.)

optional layer(s) for task:

layer k-1: (taken as contextual embedding)

layers 1 to k-2: (compose embeddings with context)

layer 0: (input: word-type embeddings)

sentence (sequence) input:



Large Training Corpus





Contextual Embeddings: for Supervised ML; for Similarity (unsup)



New Corpus

softmax for LM:

layer k: (used for language modeling)

layer k-1: (taken as contextual embedding)

layers 1 to k-2: (compose embeddings with context)

layer 0: (input: word-type embeddings)

sentence (sequence) input:

Contextual Embeddings: for Supervised ML; for Similarity (unsup)





New Corpus



layer k-1: (taken as contextual embedding)

layers 1 to k-2: (compose embeddings with context)

layer 0: (input: word-type embeddings)

sentence (sequence) input:

Contextual Embeddings: for Supervised ML

classifier or regressor: (e.g. sentiment, topic classification, etc.)

layer(s) for task:

layer k-1: (taken as contextual embedding)

layers 1 to k-2: (compose embeddings with context)

layer 0: (input: word-type embeddings)

sentence (sequence) input:



Contextual Embeddings: for Similarity (unsup)

classifier or regressor: (e.g. sentiment, topic classification, etc.)

layer(s) for task:

layer k-1: (taken as contextual embedding)

layers 1 to k-2: (compose embeddings with context)

layer 0: (input: word-type embeddings)

sentence (sequence) input:







RAG, Few-Shot, Zero-Shot

softmax for LM:

layer k: (used for language modeling)

layer k-1: (taken as contextual embedding)

layers 1 to k-2: (compose embeddings with context)

layer 0: (input: word-type embeddings)

sentence (sequence) input:

Answer(s)

No training! The model is frozen

Zero shot = Prompt has no examples, just prompting directly for the task, without answer.

Few shot = Prompt has a few examples of the task with answer, then prompting for the task without answer.

RAG = Using other NLP techniques to retrieve relevant information to include in the prompt (retrieval approach can use other models).

Task Prompts

e.g. What topic is this about? "Last night, the Seawolves won the game." answer: sports



Supplemental Review Material



simpler version





softmax for LM:

layer k: (used for language modeling)

layer k-1: (taken as contextual embedding)

layers 1 to k-2: (compose embeddings with context)

layer 0: (input: word-type embeddings)

sentence (sequence) input:

Decoder



Decoder: Cross Attention



Decoder: Masked Multi-Head Attention

