

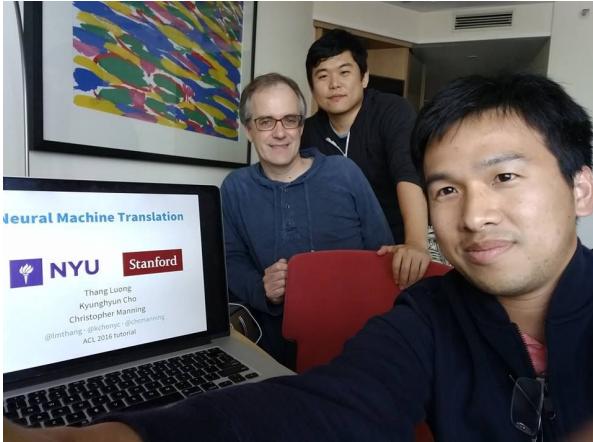
Transformer & Question Answering

Thang Luong

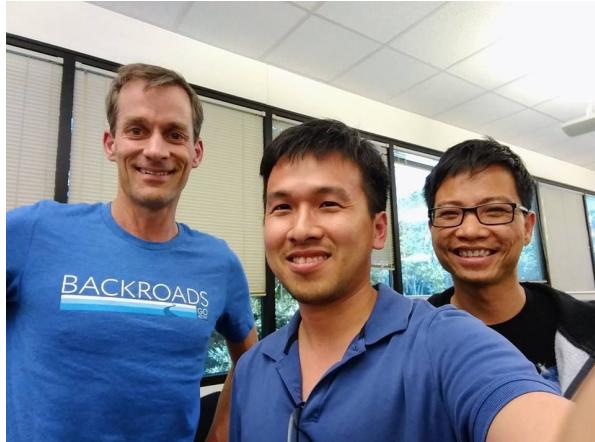
VietAI Guest Lecture

About me

- Research Scientist at Google Brain; PhD from Stanford



Chris Manning & Kyunghyun Cho
(4-hr tutorial rehearsal at ACL'16)



Jeff Dean & Quoc Le
(Google office on July 4th, 2017)



Vietnamese @ Google Brain
Farewell Hung Bui
DeepMind → VinAI, March 2019

Why Transformer?



Why Transformer?

Best NLP Model Ever? Google BERT Sets New Standards in 11 Language Tasks

 Synced in SyncedReview [Follow](#)
Oct 16, 2018 · 3 min read *



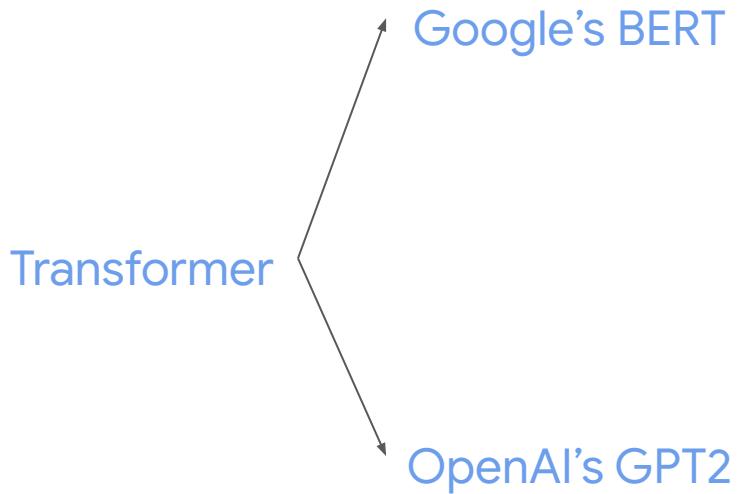
Google's BERT



Transformer



Why Transformer?



Best NLP Model Ever? Google BERT Sets New Standards in 11 Language Tasks

Synced in SyncedReview Follow
Oct 16, 2018 · 3 min read *



VietAI



ABOUT PROGRESS RESOURCES BLOG

Better Language Models and Their Implications

We've trained a large-scale unsupervised language model which generates coherent paragraphs of text, achieves state-of-the-art performance on many language modeling benchmarks, and performs rudimentary reading comprehension, machine translation, question answering, and summarization—all without task-specific training.

Outline

- Transformer
 - Self-attention
 - Architecture
- Question Answering (BiDAF, QANet)
- Beyond

Let's translate “Apple” into your language



Apple is loved by many people; it is good for your **health**



Apple is loved by many people; it is good in its **design**

Needs to capture long-range dependency



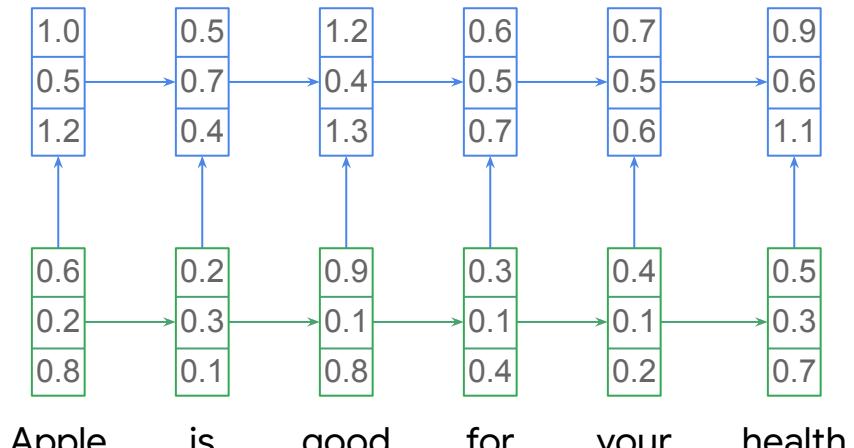
Apple is loved by many people; it is good for your **health**



Apple is loved by many people; it is good in its **design**

RNNs/LSTMs

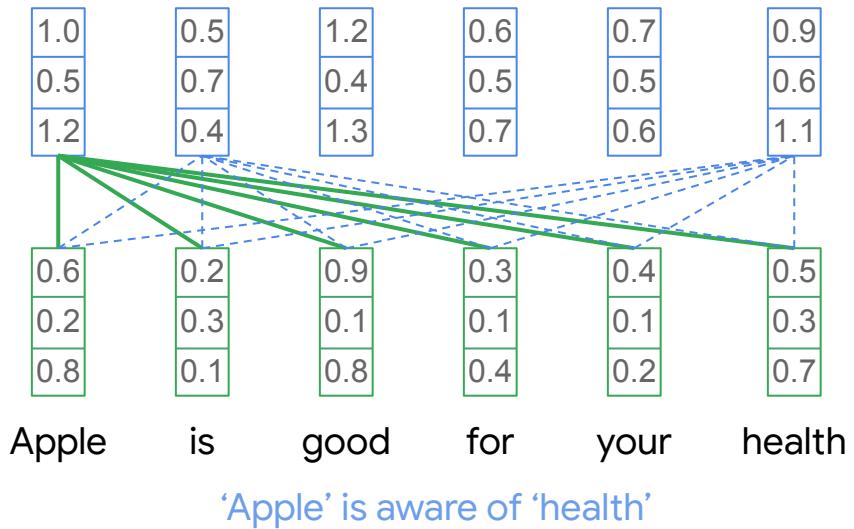
Sequential: slow & difficult to capture long-range dependency.



'Apple' is not aware of 'health'

Self-attention

All words connect to one another



(Simplified) Self-attention



The diagram illustrates the calculation of attention weights for each word in the sequence "Apple is good for your health".

Query Q: A vertical vector with values [1.0, 0.5, 1.2].

Keys K: A horizontal sequence of vectors, each representing a word from the input sequence:

- "Apple": [0.6, 0.2, 0.8]
- "is": [0.2, 0.3, 0.1]
- "good": [0.9, 0.1, 0.8]
- "for": [0.3, 0.1, 0.4]
- "your": [0.4, 0.1, 0.2]
- "health": [0.5, 0.3, 0.7]

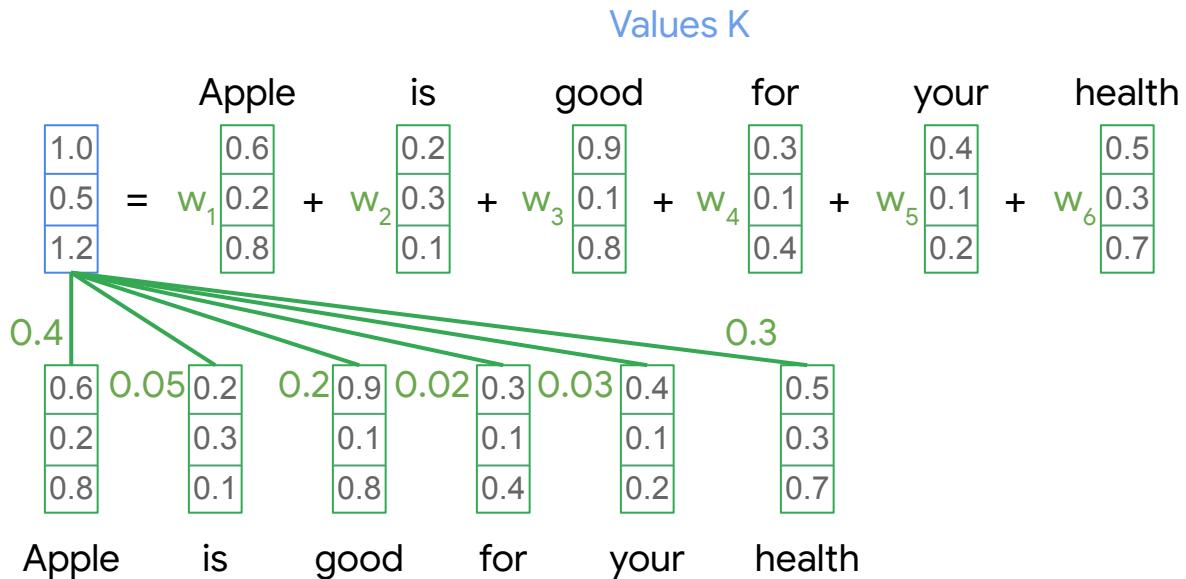
Weights W: A set of weight vectors $w_1, w_2, w_3, w_4, w_5, w_6$ that map the Query Q vector to the Keys K vectors.

The attention weights are calculated as the dot product of Query Q and Keys K, scaled by the softmax of the weights w_i :

$$\text{Attention Weight} = \frac{\text{dot}(Q, K)}{\text{softmax}(w_i)}$$

Word	Key Vector	Weight w_i	Attention Weight
Apple	[0.6, 0.2, 0.8]	w_1	$\frac{\text{dot}(Q, K)}{\text{softmax}(w_1)}$
is	[0.2, 0.3, 0.1]	w_2	$\frac{\text{dot}(Q, K)}{\text{softmax}(w_2)}$
good	[0.9, 0.1, 0.8]	w_3	$\frac{\text{dot}(Q, K)}{\text{softmax}(w_3)}$
for	[0.3, 0.1, 0.4]	w_4	$\frac{\text{dot}(Q, K)}{\text{softmax}(w_4)}$
your	[0.4, 0.1, 0.2]	w_5	$\frac{\text{dot}(Q, K)}{\text{softmax}(w_5)}$
health	[0.5, 0.3, 0.7]	w_6	$\frac{\text{dot}(Q, K)}{\text{softmax}(w_6)}$

(Simplified) Self-attention



Self-attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled factor = sqrt(hidden size)

Query Q		Keys K						
	Apple	Apple	is	good	for	your	health	
$w_1, w_2, w_3, w_4, w_5, w_6$ = softmax ($\begin{matrix} 0.6 \\ 0.2 \\ 0.8 \end{matrix}$, $\begin{matrix} 0.6 \\ 0.2 \\ 0.8 \end{matrix}$	$\begin{matrix} 0.2 \\ 0.3 \\ 0.1 \end{matrix}$	$\begin{matrix} 0.9 \\ 0.1 \\ 0.8 \end{matrix}$	$\begin{matrix} 0.3 \\ 0.1 \\ 0.4 \end{matrix}$	$\begin{matrix} 0.4 \\ 0.1 \\ 0.2 \end{matrix}$	$\begin{matrix} 0.5 \\ 0.3 \\ 0.7 \end{matrix}$)

Values V						
	Apple	is	good	for	your	health
$\begin{matrix} 1.0 \\ 0.5 \\ 1.2 \end{matrix}$ = w_1	$\begin{matrix} 0.6 \\ 0.2 \\ 0.8 \end{matrix}$	$\begin{matrix} 0.2 \\ 0.3 \\ 0.1 \end{matrix}$	$\begin{matrix} 0.9 \\ 0.1 \\ 0.8 \end{matrix}$	$\begin{matrix} 0.3 \\ 0.1 \\ 0.4 \end{matrix}$	$\begin{matrix} 0.4 \\ 0.1 \\ 0.2 \end{matrix}$	$\begin{matrix} 0.5 \\ 0.3 \\ 0.7 \end{matrix}$

Multi-headed Self-attention

One parameter set per head i

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\begin{aligned} \text{MultiHead}(Q, K, V) &= \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \\ \text{where } \text{head}_i &= \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \end{aligned}$$

Query Q	Keys K						
Apple	Apple		is	good	for	your	health
$w_1, w_2, w_3, w_4, w_5, w_6 = \text{softmax} \left(\begin{matrix} 0.6 & 0.2 & 0.8 \\ 0.2 & 0.3 & 0.1 \\ 0.8 & 0.1 & 0.8 \\ 0.6 & 0.2 & 0.8 \\ 0.2 & 0.3 & 0.1 \\ 0.8 & 0.1 & 0.8 \end{matrix} \right)$	0.6	0.2	0.9	0.3	0.4	0.5	
	0.2	0.3	0.1	0.1	0.1	0.2	
	0.8	0.1	0.8	0.4	0.2	0.3	

Values V							
Apple	is	good	for	your	health		
1.0	0.6	0.2	0.9	0.3	0.4	0.5	
0.5	0.2	0.3	0.1	0.1	0.1	0.3	
1.2	0.8	0.1	0.8	0.4	0.2	0.7	
$= w_1 + w_2 + w_3 + w_4 + w_5 + w_6$							



Tensor2Tensor

[tensorflow / tensor2tensor](#)

Tensor2Tensor

[pypi package](#) 1.13.4 [issues](#) 406 open [contributions](#) welcome [chat on gitter](#) License Apache 2.0 build error Run on FloydHub

Tensor2Tensor, or T2T for short, is a library of deep learning models and datasets designed to make deep learning more accessible and accelerate ML research. T2T is actively used and maintained by researchers and engineers within the Google Brain team and a community of users. We're eager to collaborate with you too, so feel free to [open an issue on GitHub](#) or send along a pull request (see our [contribution doc](#)). You can chat with us on [Gitter](#) and join the [T2T Google Group](#).

Quick Start

This [iPython notebook](#) explains T2T and runs in your browser using a free VM from Google, no installation needed. Alternatively, here is a one-command version that installs T2T, downloads MNIST, trains a model and evaluates it:

```
pip install tensor2tensor && t2t-trainer \
--generate_data \
--data_dir=~/t2t_data \
--output_dir=~/t2t_train/mnist \
--problem=image_mnist \
--model=shake_shake \
--hparams_set=shake_shake_quick \
--train_steps=1000 \
--eval_steps=100
```

<https://github.com/tensorflow/tensor2tensor>

Watch 427

Star 8,142

Fork 2,081

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 Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

<https://arxiv.org/abs/1706.03762>

Tensor2Tensor: Multi-headed Self-attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad \text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

```
def scaled_dot_product_attention_simple(q, k, v, bias, name=None):
    """Scaled dot-product attention. One head. One spatial dimension.
```

Args:

- q: a Tensor with shape [batch, length_q, depth_k]
- k: a Tensor with shape [batch, length_kv, depth_k]
- v: a Tensor with shape [batch, length_kv, depth_v]
- bias: optional Tensor broadcastable to [batch, length_q, length_kv]
- name: an optional string

https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/common_attention.py#L4390
https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/common_attention.py#L5985

Tensor2Tensor: Multi-headed Self-attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad \text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

```

scalar = tf.rsqrt(tf.to_float(common_layers.shape_list(q)[2]))
logits = tf.matmul(q * scalar, k, transpose_b=True)
if bias is not None:
    logits += bias
weights = tf.nn.softmax(logits, name="attention_weights")
if common_layers.should_generate_summaries():
    tf.summary.image(
        "attention", tf.expand_dims(tf.pow(weights, 0.2), 3), max_outputs=1)
return tf.matmul(weights, v)
  
```

https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/common_attention.py#L4390
https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/common_attention.py#L5985

Tensor2Tensor: Multi-headed Self-attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

```

scalar = tf.rsqrt(tf.to_float(common_layers.shape_list(q)[2]))
logits = tf.matmul(q * scalar, k, transpose_b=True)
if bias is not None:
    logits += bias
weights = tf.nn.softmax(logits, name="attention_weights")
if common_layers.should_generate_summaries():
    tf.summary.image(
        "attention", tf.expand_dims(tf.pow(weights, 0.2), 3), max_outputs=1)
return tf.matmul(weights, v)

```

https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/common_attention.py#L4390
https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/common_attention.py#L5985

Tensor2Tensor: Multi-headed Self-attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

```

scalar = tf.rsqrt(tf.to_float(common_layers.shape_list(q)[2]))
logits = tf.matmul(q * scalar, k, transpose_b=True)

if bias is not None:
    logits += bias
weights = tf.nn.softmax(logits, name="attention_weights")
if common_layers.should_generate_summaries():
    tf.summary.image(
        "attention", tf.expand_dims(tf.pow(weights, 0.2), 3), max_outputs=1)
return tf.matmul(weights, v)

```

https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/common_attention.py#L4390
https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/common_attention.py#L5985

Tensor2Tensor: Multi-headed Self-attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

```

scalar = tf.rsqrt(tf.to_float(common_layers.shape_list(q)[2]))
logits = tf.matmul(q * scalar, k, transpose_b=True)
if bias is not None:
    logits += bias
weights = tf.nn.softmax(logits, name="attention_weights")
if common_layers.should_generate_summaries():
    tf.summary.image(
        "attention", tf.expand_dims(tf.pow(weights, 0.2), 3), max_outputs=1)
return tf.matmul(weights, v)
  
```

https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/common_attention.py#L4390
https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/common_attention.py#L5985

Tensor2Tensor: Multi-headed Self-attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

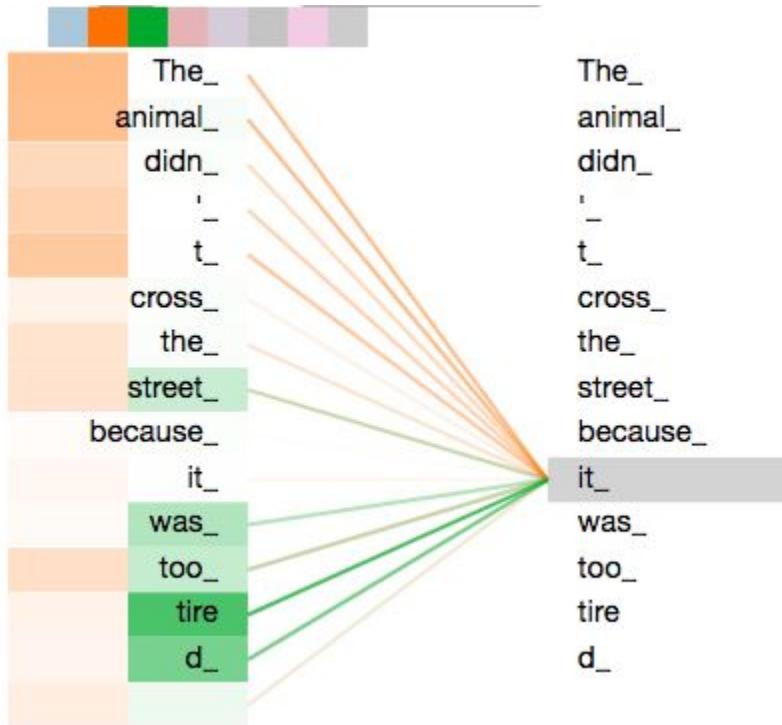
$$\begin{aligned} \text{MultiHead}(Q, K, V) &= \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \\ \text{where } \text{head}_i &= \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \end{aligned}$$

```

scalar = tf.rsqrt(tf.to_float(common_layers.shape_list(q)[2]))
logits = tf.matmul(q * scalar, k, transpose_b=True)
if bias is not None:
    logits += bias
weights = tf.nn.softmax(logits, name="attention_weights")
if common_layers.should_generate_summaries():
    tf.summary.image(
        "attention", tf.expand_dims(tf.pow(weights, 0.2), 3), max_outputs=1)
return tf.matmul(weights, v)
  
```

https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/common_attention.py#L4390
https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/common_attention.py#L5985

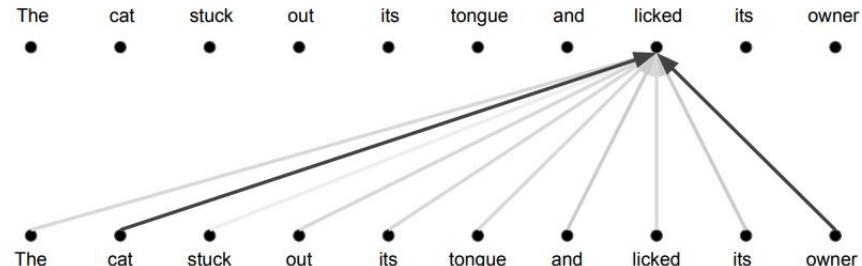
Visualization of multi-head attention



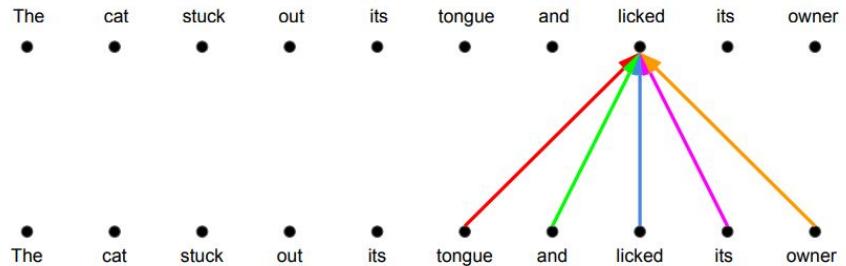
The Illustrated Transformer. <https://jalammar.github.io/illustrated-transformer>

Attention vs. Convolution

Attention: a weighted average



Convolution:
Different linear transformations by relative position.



Content from <http://web.stanford.edu/class/cs224n/slides/cs224n-2019-lecture13-contextual-representations.pdf>

Attention is cheap!

FLOPs

Self-Attention	$O(\text{length}^2 \cdot \text{dim})$	$= 4 \cdot 10^9$
RNN (LSTM)	$O(\text{length} \cdot \text{dim}^2)$	$= 16 \cdot 10^9$
Convolution	$O(\text{length} \cdot \text{dim}^2 \cdot \text{kernel_width})$	$= 6 \cdot 10^9$

length=1000 dim=1000 kernel_width=3

Content from <http://web.stanford.edu/class/cs224n/slides/cs224n-2019-lecture13-contextual-representations.pdf>

Outline

- Transformer
 - Self-attention
 - Architecture
- Question Answering (BiDAF, QANet)
- Beyond

Self-attention as a layer

N vectors in → N vectors out



Transformer

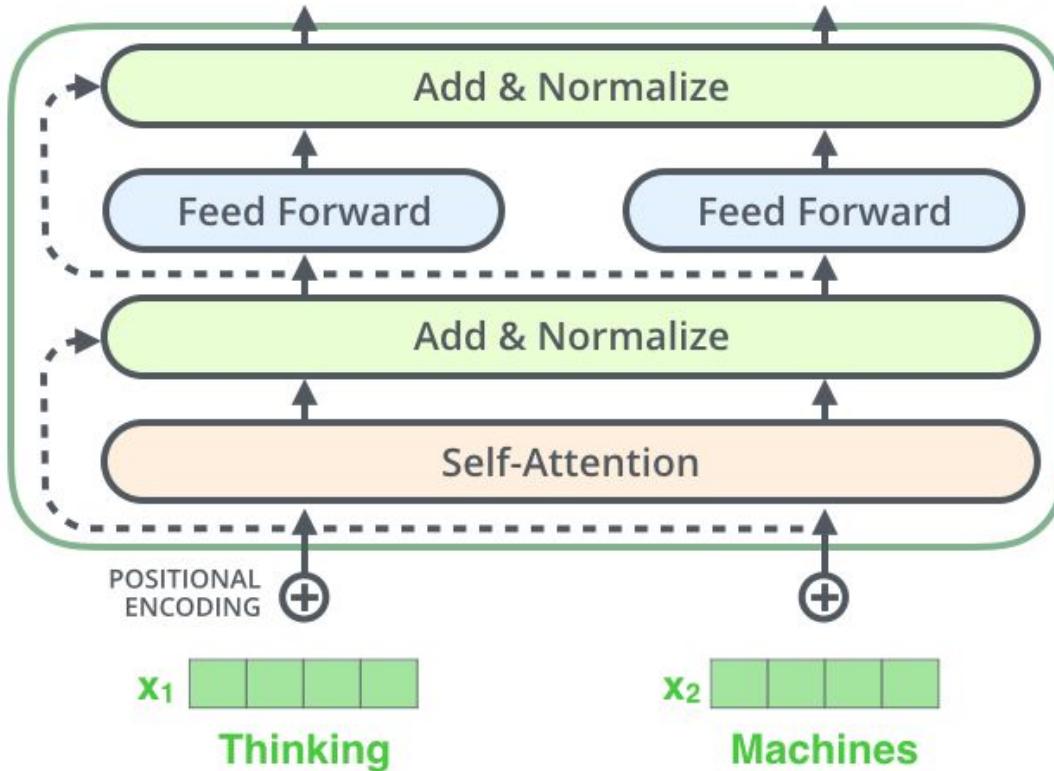


Image from <https://jalammar.github.io/illustrated-transformer>

Transformer

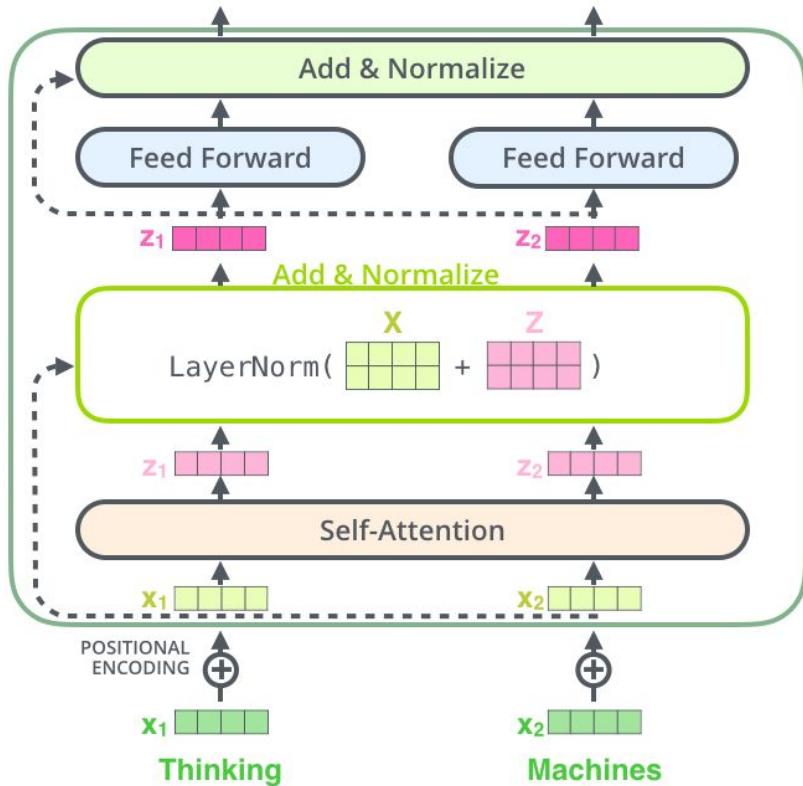


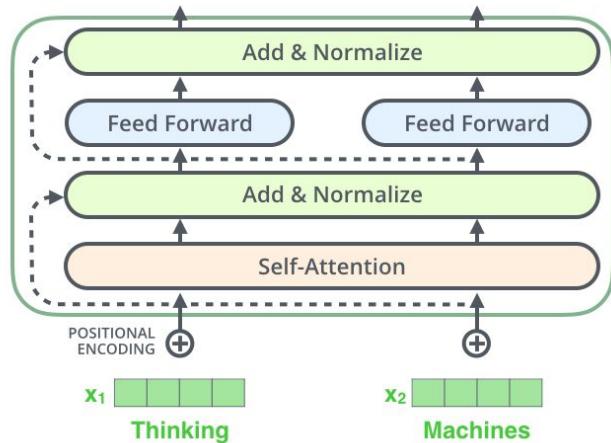
Image from <https://jalammar.github.io/illustrated-transformer>

Transformer



VietAI

```
for layer in range(hparams.num_encoder_layers or hparams.num_hidden_layers):
    with tf.variable_scope("layer_%d" % layer):
        with tf.variable_scope("self_attention"):
            y = common_attention.multihead_attention(
                common_layers.layer_preprocess(x, hparams),
                training=(hparams.get("mode", tf.estimator.ModeKeys.TRAIN)
                          == tf.estimator.ModeKeys.TRAIN))
            x = common_layers.layer_postprocess(x, y, hparams)
        with tf.variable_scope("ffn"):
            y = transformer_ffn_layer(
                common_layers.layer_preprocess(x, hparams),
                hparams,
                pad_remover,
                conv_padding="SAME",
                nonpadding_mask=nonpadding,
                losses=losses)
            x = common_layers.layer_postprocess(x, y, hparams)
```



<https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/models/transformer.py#L185>

https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/transformer_layers.py#L128

Transformer

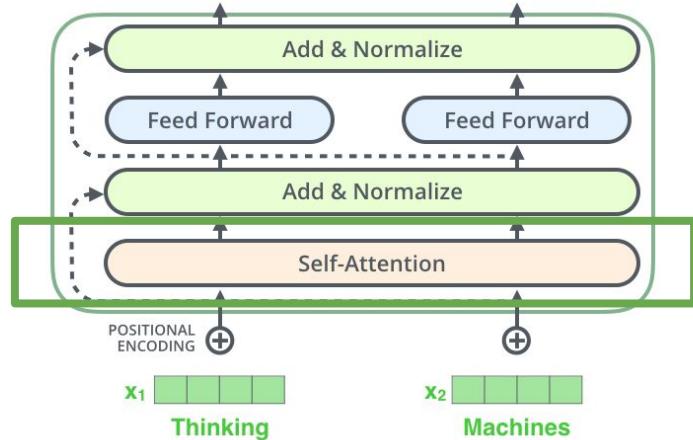


VietAI

```
for layer in range(hparams.num_encoder_layers or hparams.num_hidden_layers):
    with tf.variable_scope("layer_%d" % layer):
        with tf.variable_scope("self_attention"):
            y = common_attention.multihead_attention(
                common_layers.layer_preprocess(x, hparams),
                training=(hparams.get("mode", tf.estimator.ModeKeys.TRAIN)
                          == tf.estimator.ModeKeys.TRAIN))
            x = common_layers.layer_postprocess(x, y, hparams)
        with tf.variable_scope("ffn"):
            y = transformer_ffn_layer(
                common_layers.layer_preprocess(x, hparams),
                hparams,
                pad_remover,
                conv_padding="SAME",
                nonpadding_mask=nonpadding,
                losses=losses)
            x = common_layers.layer_postprocess(x, y, hparams)
```

<https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/models/transformer.py#L185>

https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/transformer_layers.py#L128

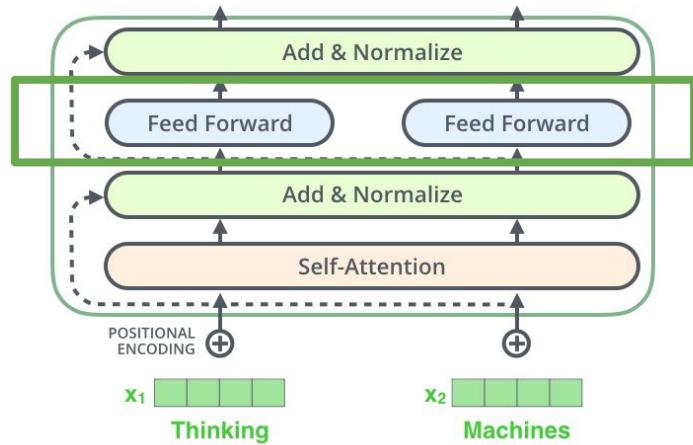


Transformer



VietAI

```
for layer in range(hparams.num_encoder_layers or hparams.num_hidden_layers):
    with tf.variable_scope("layer_%d" % layer):
        with tf.variable_scope("self_attention"):
            y = common_attention.multihead_attention(
                common_layers.layer_preprocess(x, hparams),
                training=(hparams.get("mode", tf.estimator.ModeKeys.TRAIN)
                          == tf.estimator.ModeKeys.TRAIN))
            x = common_layers.layer_postprocess(x, y, hparams)
        with tf.variable_scope("ffn"):
            y = transformer_ffn_layer(
                common_layers.layer_preprocess(x, hparams),
                hparams,
                pad_remover,
                conv_padding="SAME",
                nonpadding_mask=nonpadding,
                losses=losses)
            x = common_layers.layer_postprocess(x, y, hparams)
```



<https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/models/transformer.py#L185>

https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/transformer_layers.py#L128

Transformer



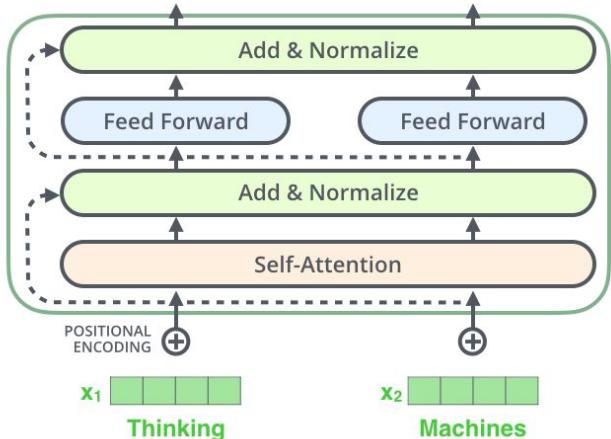
VietAI

```
for layer in range(hparams.num_encoder_layers or hparams.num_hidden_layers):
    with tf.variable_scope("layer_%d" % layer):
        with tf.variable_scope("self_attention"):
            Pre-process: x = LayerNorm(x)

            y = common_attention.multihead_attention(
                common_layers.layer_preprocess(x, hparams),
                training=(hparams.get("mode", tf.estimator.ModeKeys.TRAIN)
                          == tf.estimator.ModeKeys.TRAIN))
            x = common_layers.layer_postprocess(x, y, hparams)

        with tf.variable_scope("ffn"):
            y = transformer_ffn_layer(
                common_layers.layer_preprocess(x, hparams),
                hparams,
                pad_remover,
                conv_padding="SAME",
                nonpadding_mask=nonpadding,
                losses=losses)
            Post-process: x = dropout(y) + x

            x = common_layers.layer_postprocess(x, y, hparams)
```



<https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/models/transformer.py#L185>

https://github.com/tensorflow/tensor2tensor/blob/master/tensor2tensor/layers/transformer_layers.py#L128

Transformer for Translation

3 types of attentions:

- Self-attention Encoder
- (Masked) Self-attention Decoder
- Attention Decoder → Encoder

Parallelism:

- Training: Encoder & Decoder
- Inference: Encoder

Outline

- Transformer
- Question Answering (BiDAF, QANet)
- Beyond

Question answering is Google's core task!

what is the most popular question answering dataset?

All Images News Videos Shopping More Settings Tools

About 66,700,000 results (0.69 seconds)

The Stanford Question Answering Dataset

<https://rajpurkar.github.io/SQuAD-explorer/> ▾

Stanford Question Answering Dataset (SQuAD) is a new reading comprehension dataset, consisting of questions posed by crowdworkers on a set of Wikipedia ...
You've visited this page many times. Last visit: 8/31/18

The Stanford Question Answering Dataset - mlx - Pranav Rajpurkar

<https://rajpurkar.github.io/mlx/qa-and-squad/> ▾

Apr 3, 2017 - The **questions** had to be entered in a text field, and the **answers** highlighted in the passage.
To guide the workers, we had examples of **good** ...

WikiTableQuestions: a Complex Real-World Question Understanding ...

<https://nlp.stanford.edu/.../wikitablequestions-a-complex-real-world-question-understa...> ▾

Feb 11, 2016 - Natural language **question** understanding has been one of the **most** important ...
However, even the best **question** answering systems today still face two ... In the WikiTableQuestions **dataset**, each **question** comes with a table from Wikipedia. ... phrases must be inferred using the context and **common** sense.

Have to read to find answer

General QA systems

Q: How many of Warsaw's inhabitants spoke Polish in 1933?

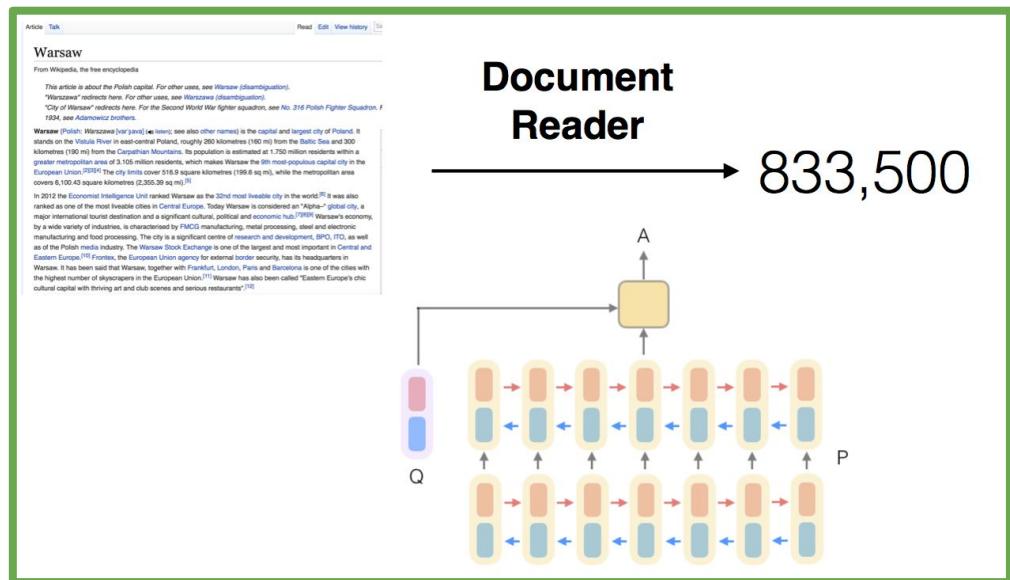
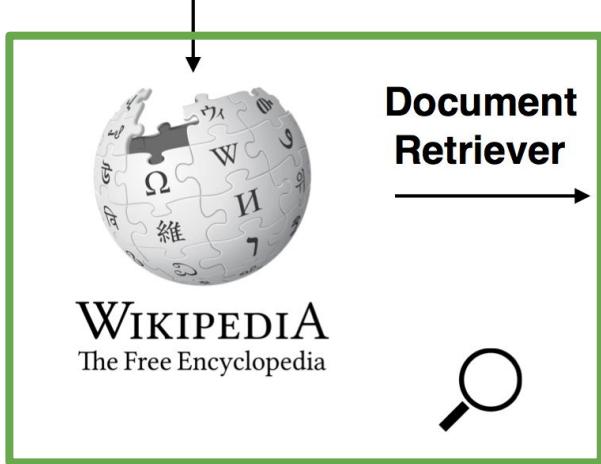
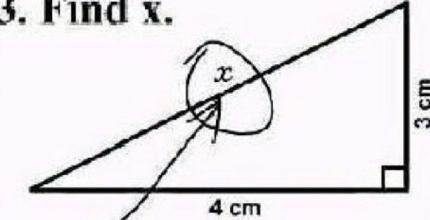


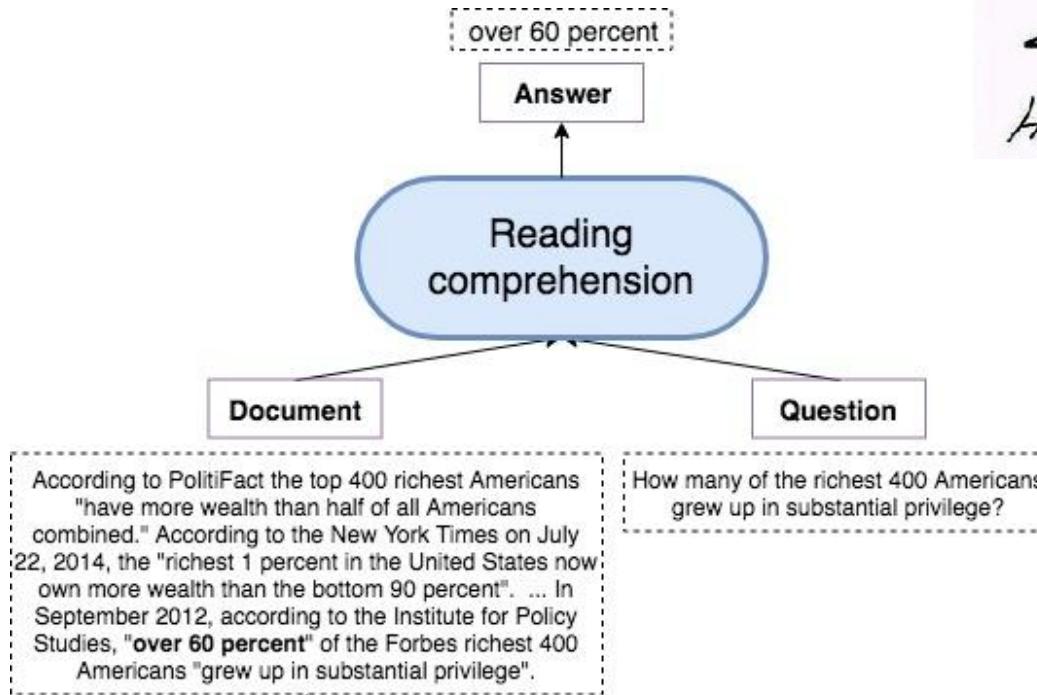
Image from <https://github.com/facebookresearch/DrQA>

3. Find x.



Here it is

Reading Comprehension



Stanford Question Answering Dataset (**SQuAD 1.1**) <https://rajpurkar.github.io/SQuAD-explorer/>

- **100K+ (question, document, answer)** over 500+ Wikipedia articles

Vietnamese Reading Comprehension System

Sample BERT-based Reading Comprehension System Demonstration

Paragraph

Google là một công ty Internet có trụ sở tại Hoa Kỳ, được thành lập vào năm 1998. Sản phẩm chính của công ty này là công cụ tìm kiếm Google, được nhiều người đánh giá là công cụ tìm kiếm hữu ích và mạnh mẽ nhất trên Internet. Trụ sở của Google tên là "Googleplex" tại Mountain View, California. Giám đốc công ty là Larry Page, một trong hai người sáng lập ra công ty. Tên "Google" là một lỗi chính tả của từ googol, bằng 10¹⁰⁰. Google chọn tên này để thể hiện sứ mệnh của công ty để sắp xếp số lượng thông tin khổng lồ trên mạng. Googleplex, tên của trụ sở Google, có nghĩa là 10googol.

Question

giám đốc của Google là ai?

Paragraph must be in English or Vietnamese.

Answer: Larry Page, - Confidence: 99.0 %

Answer: Larry - Confidence: 0.0 %

Answer: Larry Page, một trong hai - Confidence: 0.0 %

Question should end with "?"

Choose model type

Vietnamese

Run

SQuAD: 100+ teams since Aug'16

Industry

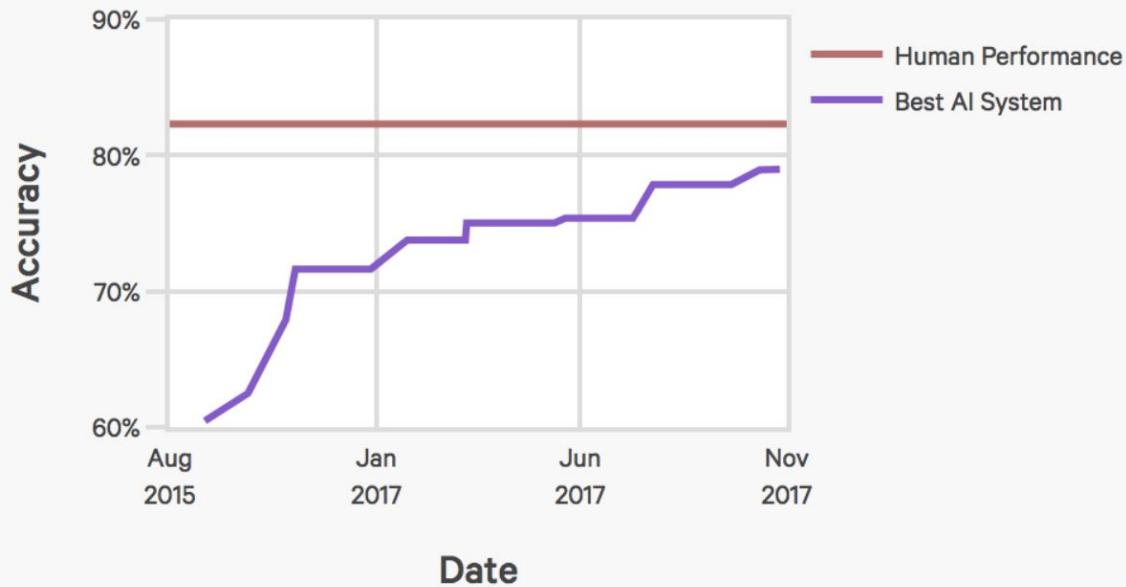


Academia



Until 2017

Question Answering, SQuAD v1.1

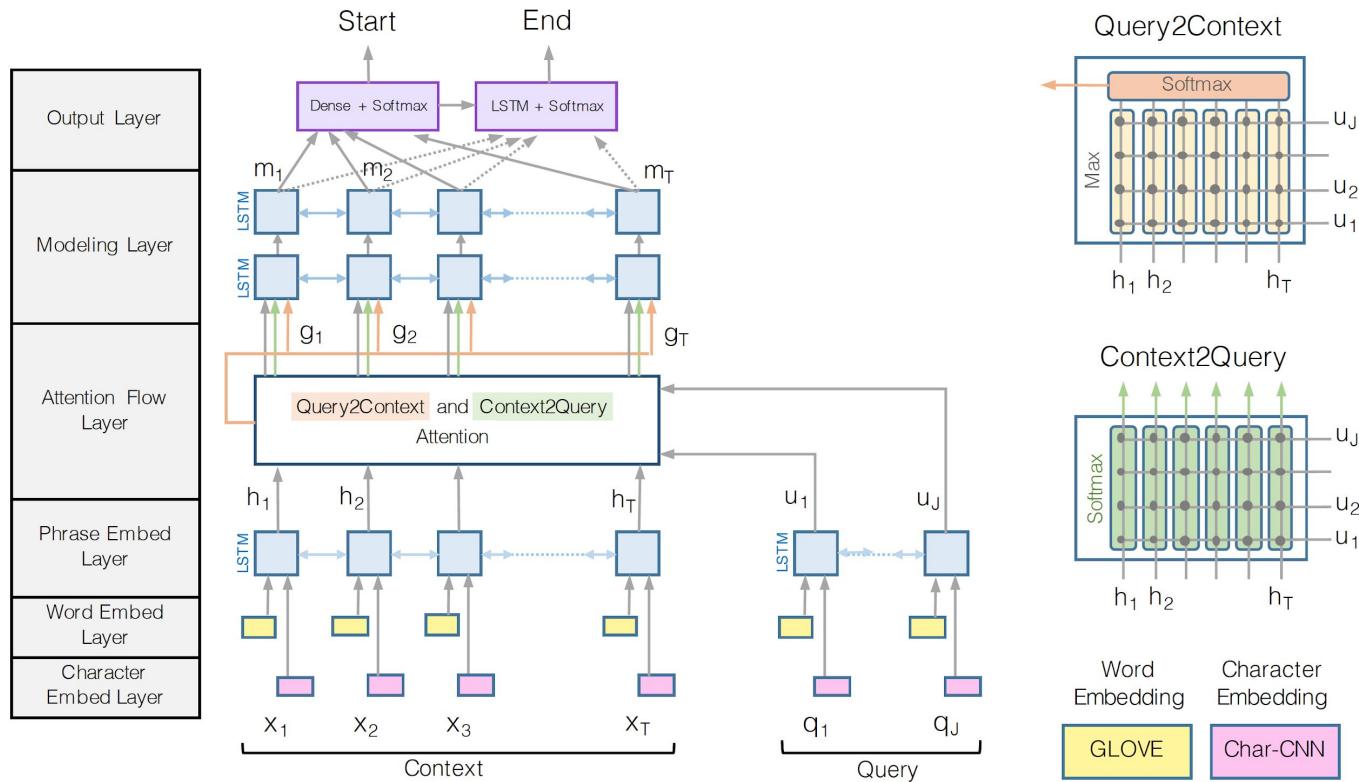


Source: stanford-qa.com

AIINDEX.ORG 

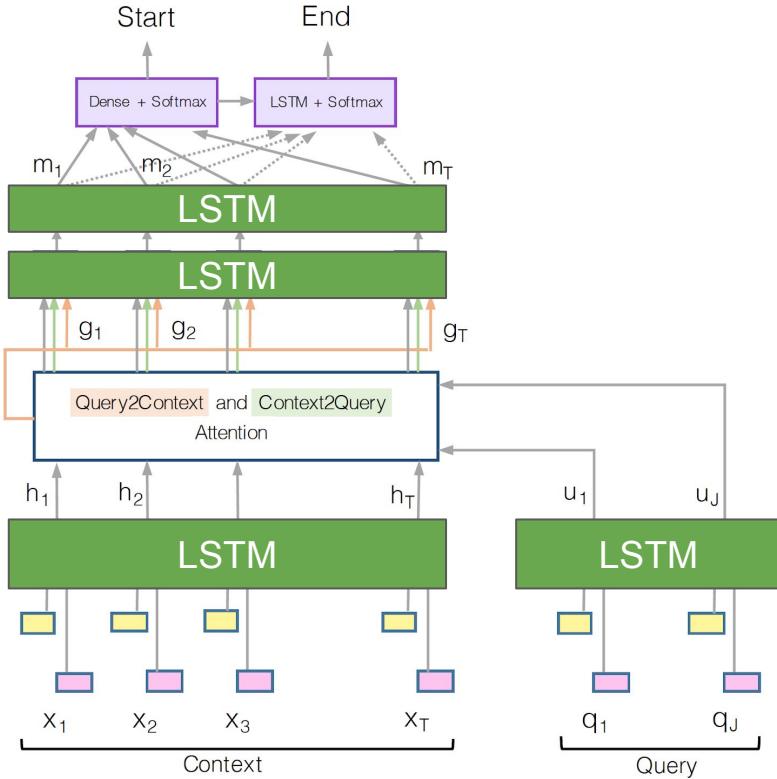
BiDAF -- LSTM

(Seo et al., ICLR'17)
<https://allenai.github.io/bi-att-flow/>



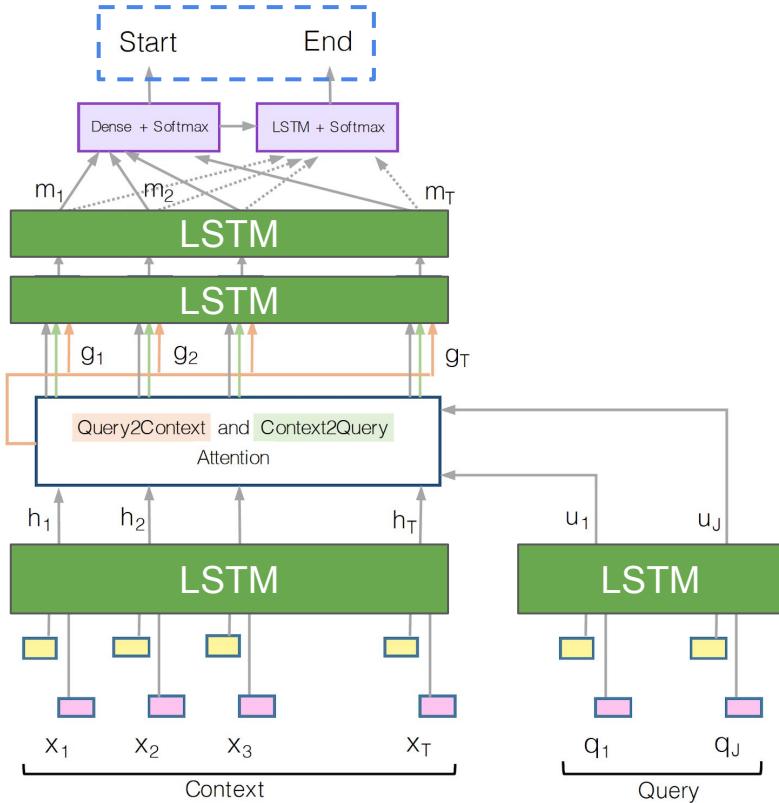
BiDAF -- LSTM

(Seo et al., ICLR'17)
<https://allenai.github.io/bi-att-flow/>



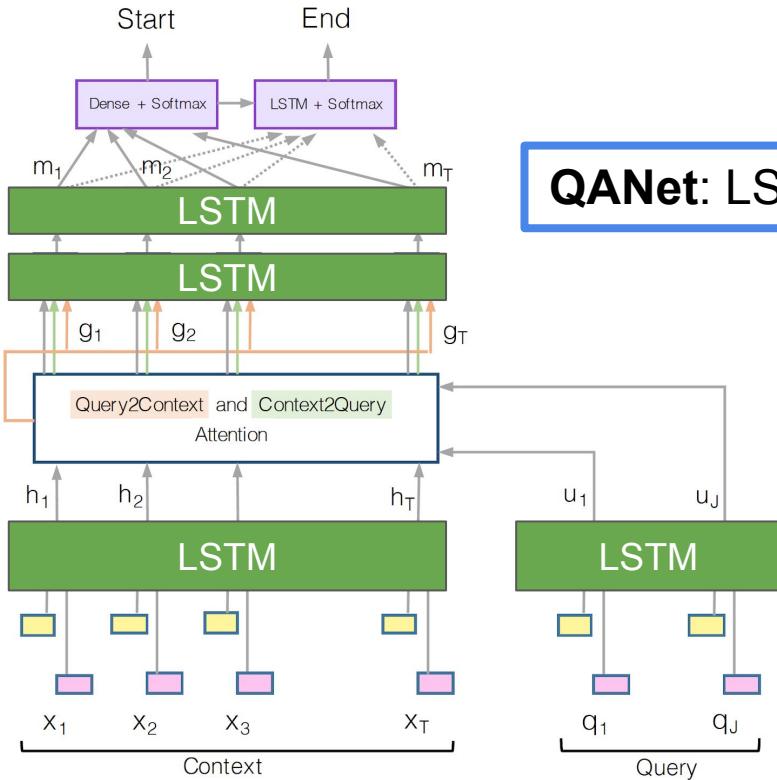
BiDAF -- LSTM

Predict start / end positions instead of words



BiDAF -- LSTM

(Seo et al., ICLR'17)
<https://allenai.github.io/bi-att-flow/>



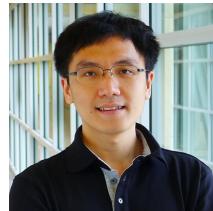
QANet: LSTM is slow, let's replace it!

SQuAD1.1 Leaderboard

Since the release of SQuAD1.0, the community has made rapid progress, with the best models now rivaling human performance on the task. Here are the ExactMatch (EM) and F1 scores evaluated on the test set of v1.1.

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar et al. '16)	82.304	91.221
1 Jul 11, 2018	QANet (ensemble) Google Brain & CMU	84.454	90.490
2 Jul 08, 2018	r-net (ensemble) Microsoft Research Asia	84.003	90.147
3 Jun 20, 2018	MARS (ensemble) YUANFUDAO research NLP	83.982	89.796
4 Mar 19, 2018	QANet (ensemble) Google Brain & CMU	83.877	89.737
5 Jun 21, 2018	MARS (single model) YUANFUDAO research NLP	83.122	89.224
6 Mar 06, 2018	QANet (ensemble) Google Brain & CMU	82.744	89.045
7 May 09, 2018	MARS (single model) YUANFUDAO research NLP	82.587	88.880
7 Feb 19, 2018	Reinforced Mnemonic Reader + A2D (ensemble model) Microsoft Research Asia & NUDT	82.849	88.764
7 Jun 20, 2018	QANet (single) Google Brain & CMU	82.471	89.306
7 Jan 22, 2018	Hybrid AoA Reader (ensemble) Joint Laboratory of HIT and iFLYTEK Research	82.482	89.281
7 Jan 03, 2018	r-net+ (ensemble) Microsoft Research Asia	82.650	88.493
7 Jan 05, 2018	SLQA+ (ensemble) Alibaba iDST NLP	82.440	88.607

#1 Mar-Aug 2018



Adams
Wei Yu



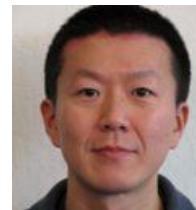
David Dohan



Thang Luong



Rui Zhao



Kai Chen



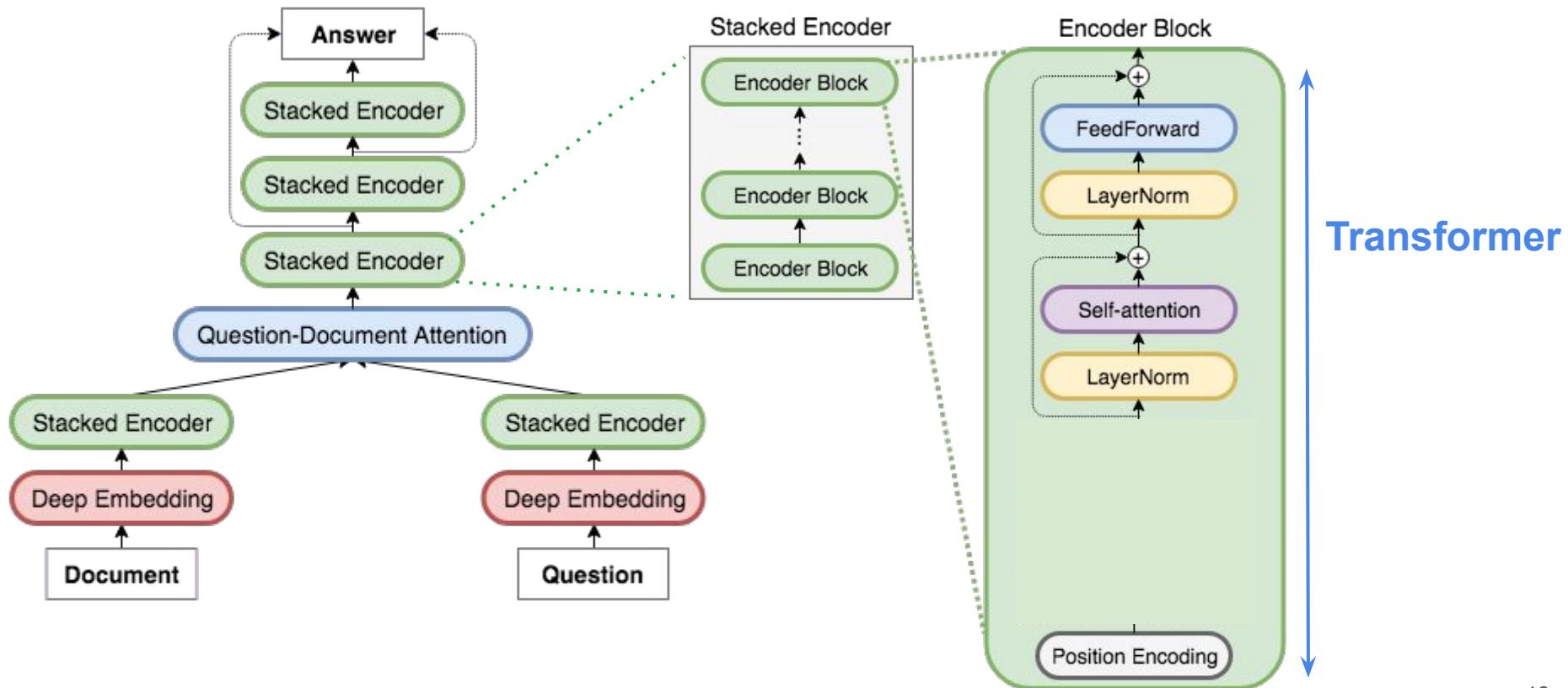
Mohammad
Norouzi



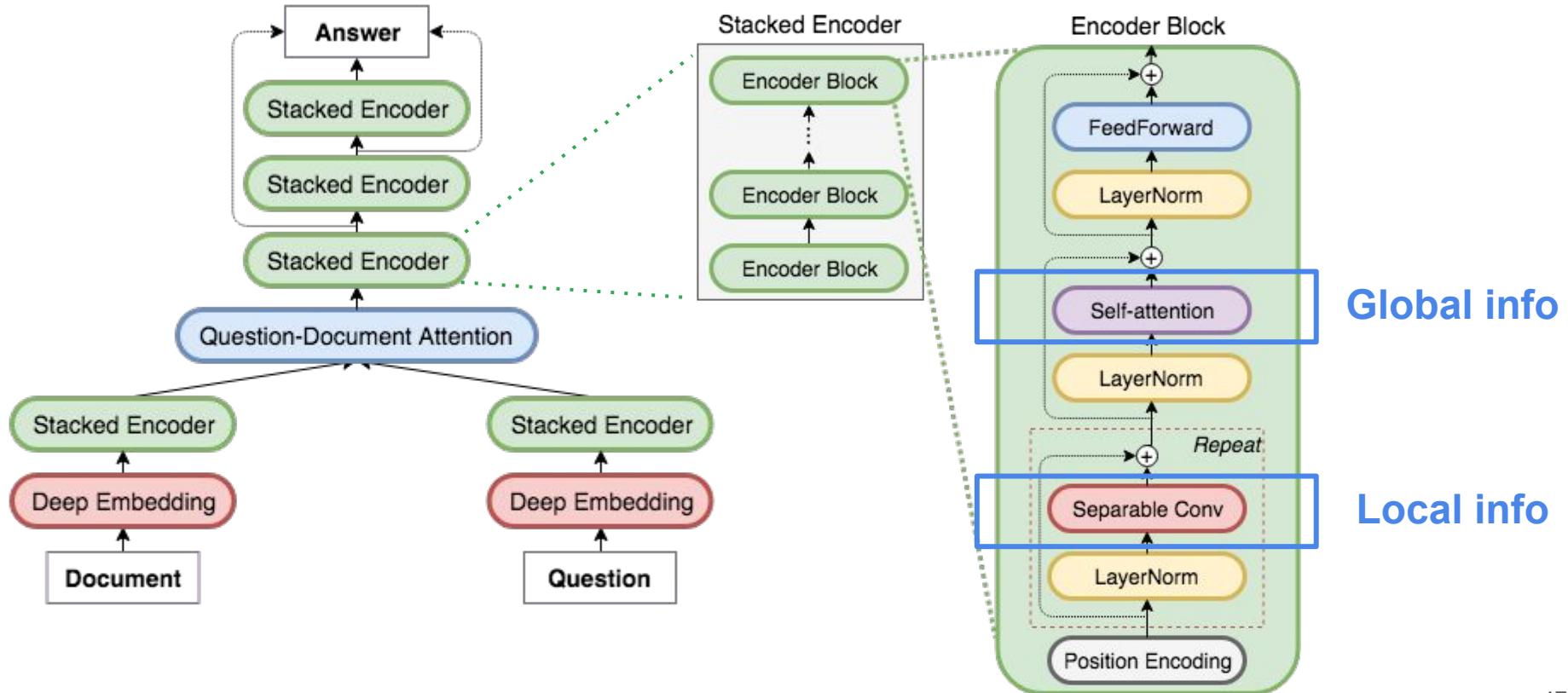
Quoc Le

Combining Local Convolution & Global Self-Attention for Reading Comprehension
ICLR'18, <https://arxiv.org/abs/1804.09541>

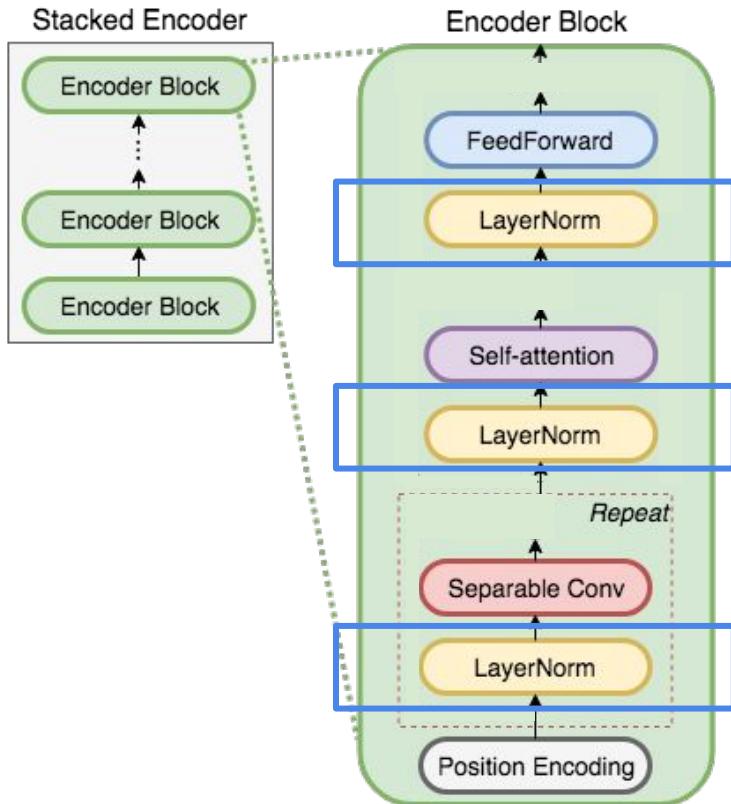
QANet - Inspired by Transformer



QANet - Inspired by Transformer



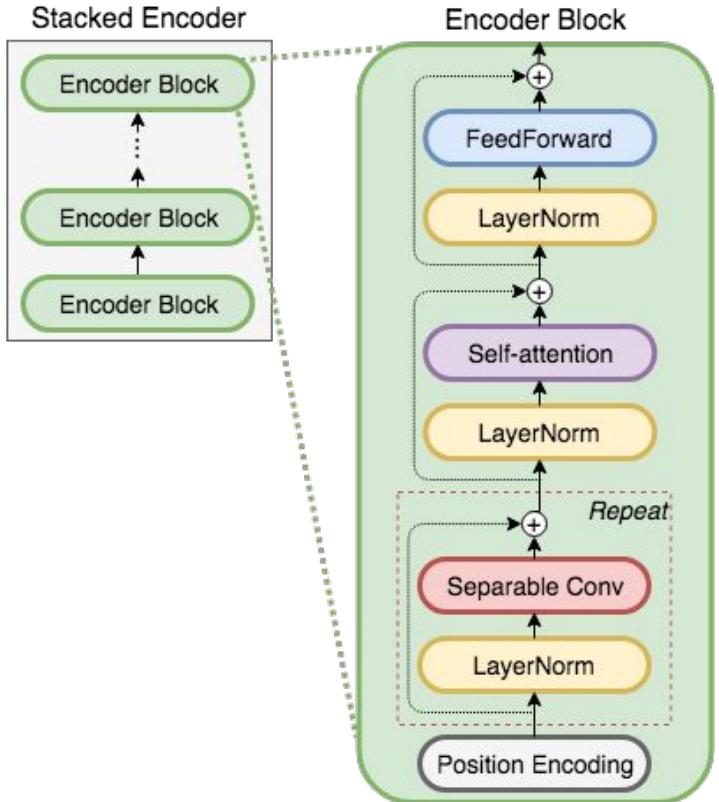
Lots of Regularizations



- **LayerNorm everywhere**

Layer Normalization (Ba et al., arXiv'16)
<https://arxiv.org/pdf/1607.06450.pdf>

Lots of Regularizations



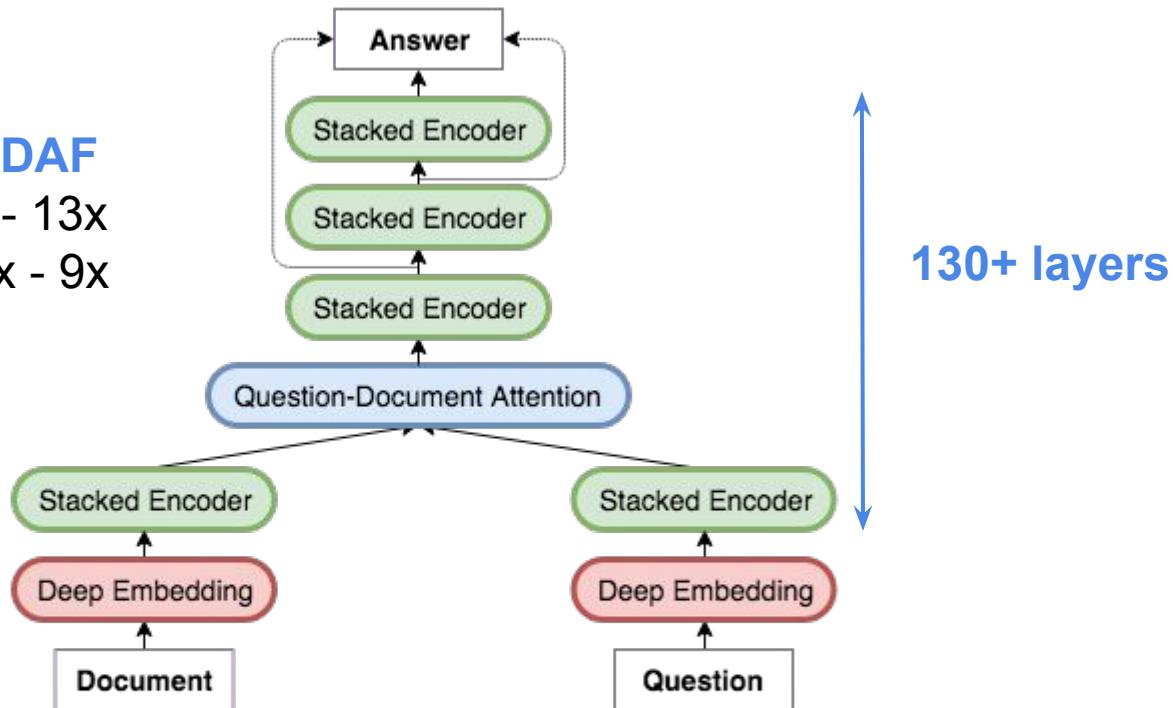
- **LayerNorm everywhere**
- **Residual Connections everywhere**

ResNet (He et al., CVPR'16)

<https://arxiv.org/pdf/1512.03385.pdf>

QANet: fast & accurate

- Compared to BiDAF
 - Training: 3x - 13x
 - Inference: 4x - 9x



Previously: 29 conv layers (Conneau et al., 2017) <https://arxiv.org/abs/1606.01781>

QANet was #1 (Mar-Aug 2018)



Jeff Dean

@JeffDean

Following



You have questions? QANet, a new stacked architecture designed to do well on question answering tasks, has answers.



Thang Luong @lmthang

Happy to announce our QANet models, #1 on @stanfordnlp question answering dataset (SQuAD). 3 ideas: deep & fast arch (130+ layers), data augmentation, transfer learning. Joint work /w @AdamsYu @dmdohan @oahziur, Quoc Le, et al. See our ...

Show this thread

6:09 PM - 25 Apr 2018

Outline

- Transformer
- Question Answering (BiDAF, QANet)
- Beyond

BERT: fully surpassing human (Oct 2018)



Thang Luong

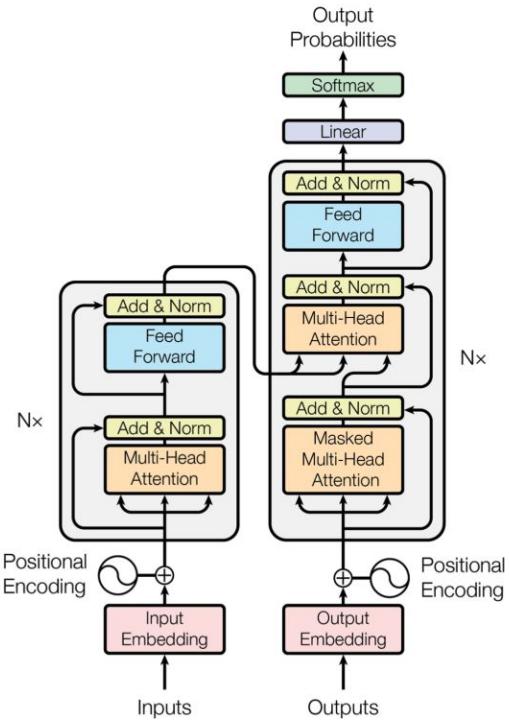
@lmthang

A new era of NLP has just begun a few days ago: large **pretraining** models (**Transformer** [24 layers, 1024 dim, 16 heads]) + massive compute is all you need. BERT from [@GoogleAI](#): SOTA results on everything arxiv.org/abs/1810.04805. Results on SQuAD are just mind-blowing. Fun time ahead!

SQuAD1.1 Leaderboard

Since the release of SQuAD1.0, the community has made rapid progress, with the best models now rivaling human performance on the task. Here are the ExactMatch (EM) and F1 scores evaluated on the test set of v1.1.

Rank	Model	EM	F1
	Human Performance <i>Stanford University (Rajpurkar et al. '16)</i>	82.304	91.221
1	BERT (ensemble) Google A.I.	87.433	93.160
	Oct 05, 2018		



A successful transfer learning story



WIKIPEDIA
The Free Encyclopedia



Pretraining: large, unlabeled data

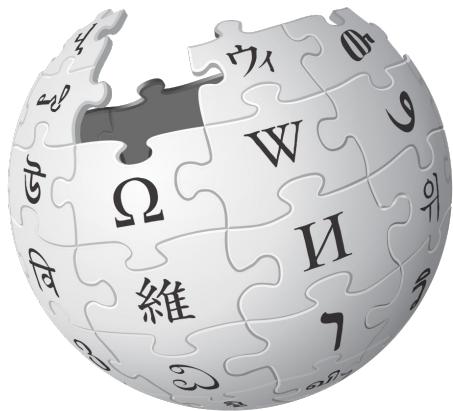
Transfer Learning: small, labeled data

The first recorded travels by Europeans to China and back date from this time. The most famous traveler of the period was the Venetian Marco Polo, whose account of his trip to "Cambaluc," the capital of the Great Khan, and of life there astounded the people of Europe. The account of his travels, *Il milione* (or, *The Million*, known in English as the *Travels of Marco Polo*), appeared about the year 1299. Some argue over the accuracy of Marco Polo's accounts due to the lack of mentioning the Great Wall of China, tea houses, which would have been a prominent sight since Europeans had yet to adopt a tea culture, as well the practice of foot binding by the women in capital of the Great Khan. Some suggest that Marco Polo acquired much of his knowledge **through contact with Persian traders** since many of the places he named were in Persian.

How did some suspect that Polo learned about China instead of by actually visiting it?

Answer: **through contact with Persian traders**

BERT → SQuAD



WIKIPEDIA

The Free Encyclopedia

BERT pretrained on
English Wikipedia (~4M articles)

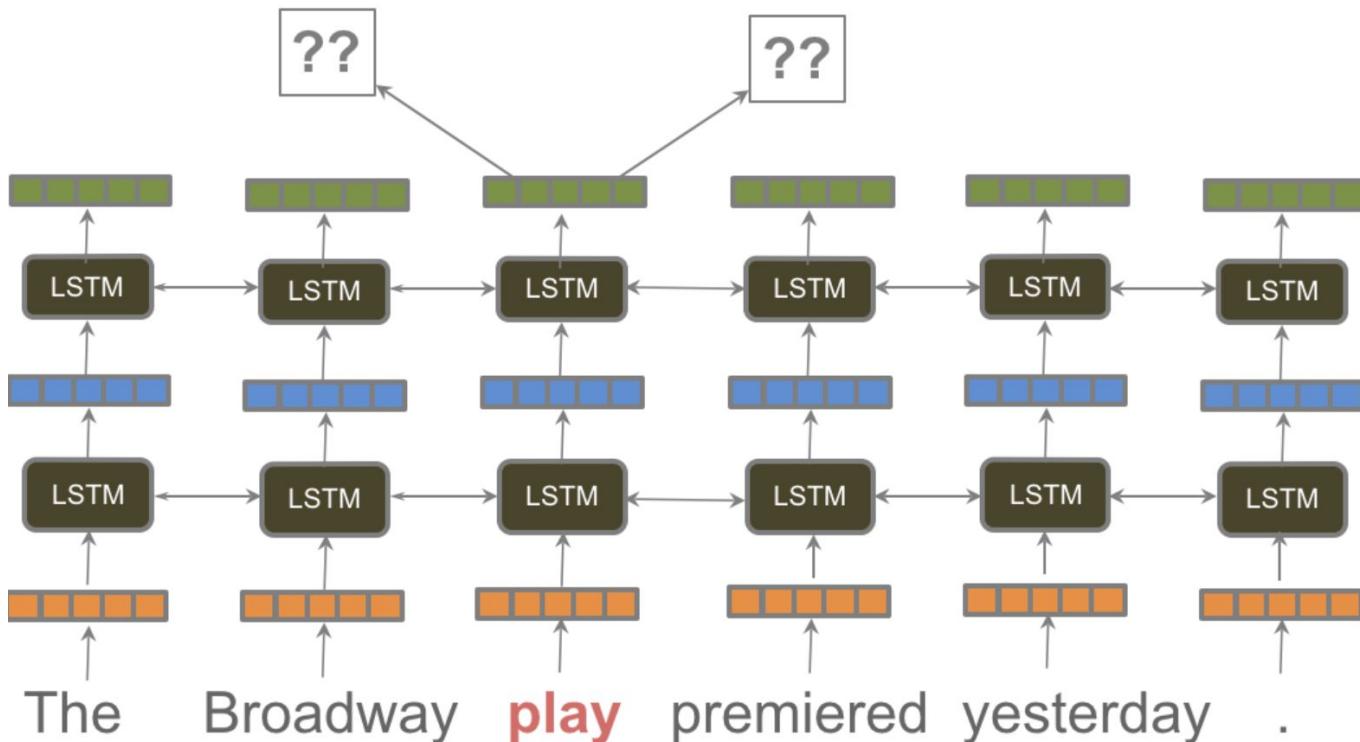
The first recorded travels by Europeans to China and back date from this time. The most famous traveler of the period was the Venetian Marco Polo, whose account of his trip to "Cambaluc," the capital of the Great Khan, and of life there astounded the people of Europe. The account of his travels, *Il milione* (or, *The Million*, known in English as the *Travels of Marco Polo*), appeared about the year 1299. Some argue over the accuracy of Marco Polo's accounts due to the lack of mentioning the Great Wall of China, tea houses, which would have been a prominent sight since Europeans had yet to adopt a tea culture, as well the practice of foot binding by the women in capital of the Great Khan. Some suggest that Marco Polo acquired much of his knowledge **through contact with Persian traders** since many of the places he named were in Persian.

How did some suspect that Polo learned about China instead of by actually visiting it?

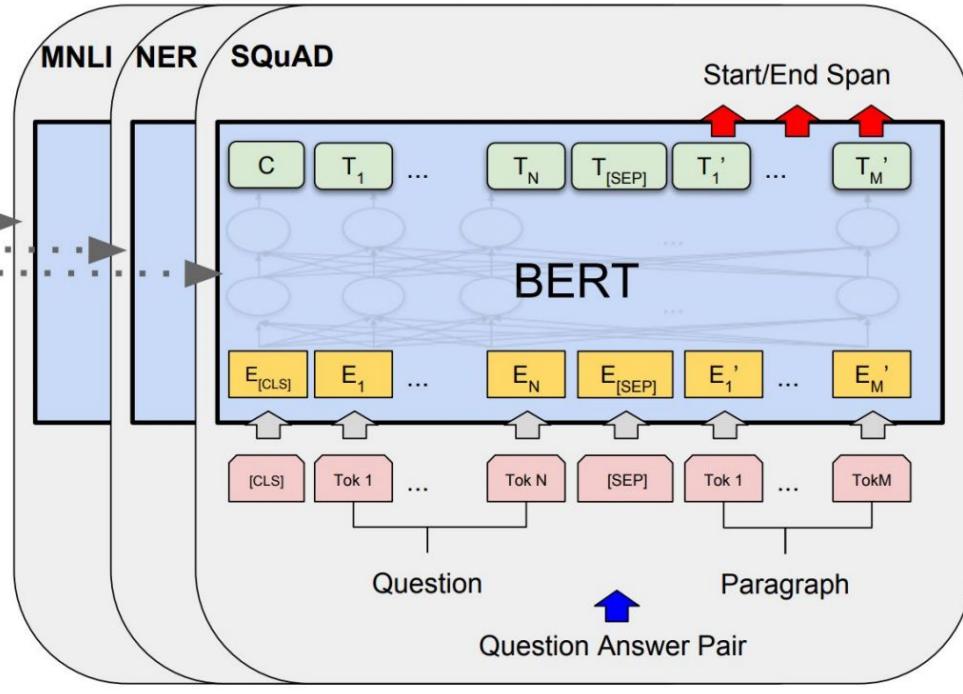
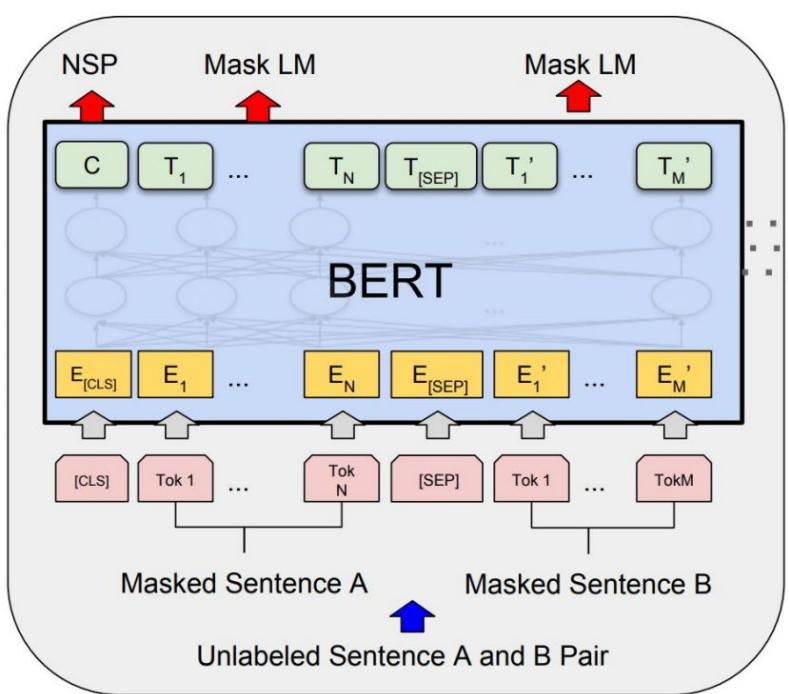
Answer: **through contact with Persian traders**

SQuAD: 500 articles
Stanford Question Answering Dataset

Unsupervised Learning for Text



BERT for SQuAD



Content from <https://nlp.stanford.edu/seminar/details/jdevlin.pdf>

BERT for SQuAD

1.1

What was another term used for the oil crisis?

Ground Truth Answers: first oil shock shock shock first oil
 shock shock
 Prediction: shock

The 1973 **oil crisis** began in October 1973 when the members of the Organization of Arab Petroleum Exporting Countries (OAPEC, consisting of the Arab members of OPEC plus Egypt and Syria) proclaimed an **oil embargo**. By the end of the embargo in March 1974, the price of **oil** had risen from US\$3 per barrel to nearly \$12 globally; US prices were significantly higher. The embargo caused an **oil crisis**, or "shock", with many short- and long-term effects on global politics and the global economy. It was later called the "**first oil shock**", followed by the 1979 **oil crisis**, termed the "second **oil shock**".

What action did the US begin that started the second oil shock?

Ground Truth Answers: <No Answer>

Prediction: <No Answer>

2.0

The 1973 **oil crisis** began in October 1973 when the members of the Organization of Arab Petroleum Exporting Countries (OAPEC, consisting of the Arab members of OPEC plus Egypt and Syria) proclaimed an **oil embargo**. By the end of the embargo in March 1974, the price of **oil** had risen from US\$3 per barrel to nearly \$12 globally; US prices were significantly higher. The embargo caused an **oil crisis**, or "shock", with many short- and long-term effects on global politics and the global economy. It was later called the "**first oil shock**", followed by the 1979 **oil crisis**, termed the "second **oil shock**".

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar et al. '16)	82.304	91.221
1	BERT (ensemble) Google AI Language https://arxiv.org/abs/1810.04805	87.433	93.160
Oct 05, 2018			
2	BERT (single model) Google AI Language https://arxiv.org/abs/1810.04805	85.083	91.835
Oct 05, 2018			
2	nlNet (ensemble) Microsoft Research Asia	85.954	91.677
Sep 26, 2018			
5	nlNet (single model) Microsoft Research Asia	83.468	90.133
Sep 09, 2018			
3	QANet (ensemble) Google Brain & CMU	84.454	90.490
Jul 11, 2018			

Content from <https://nlp.stanford.edu/seminar/details/jdevlin.pdf>

Final: Let's work on something for Vietnamese!



Binhvq News Corpus

Thông tin cơ bản:

- Trích xuất từ khoảng 14.896.998 bài báo trên internet bao gồm các báo:

2Sao, ANTG, ANTT, ANTD, ATGT, BVPL, BizLIVE, Biên Phòng, Bnews, Báo Văn hóa, Bộ GTVT, Bộ KHTT, Bộ Ngoại Giao, Bộ Nội Vụ, Bộ Tài Chính, Bộ VHTTDL, CAND, CAT Cartimes, Chính Phủ, Công Luận, Công Lý, Công Thương, DNVN, Doanh Nghiệp, Dân Em Đẹp, GD&ĐT, GTVT, Gia Đình Mới, Gia Đình VN, Giao Thông, Giáo Dục VN, GB&Xt Hà Tĩnh, Hải Quan, ICTNews, Infonet, KTNT, KTĐT, Khỏe 365, Khỏe Plus, Khỏe Plus Kiểm Sát, Kiểm sát, Ký Nguyên Số, Lao Động, LĐTBQ, MT&CS, Mặt Trận, Một Thế Giới Nghe Nhìn VN, Nghệ An, Ngày Nay, Người Lành Bão, Người Tiêu Dùng, Người Đô Thị, Nhân Dân, Nông Nghiệp, NB&DS, PC World, PL&XH, PLO, PNNews, PNSK, PetroTimes, Pháp Luật Plus, Pháp Luật VN, Phụ Nữ VN, Quốc Hội, Quốc Hội TV, QĐND, SGGP, SC Seatiners, Sài Gòn Tiếp Thị, TBDN, TBKTSG, TG&VN, TGT, TH&PL, TNMT, TTOL, TTXVN Thanh Niên, Thành Tra, TheLEADER, Thương Gia, Thủ Giới Trẻ, Thủ Giới Xe, Tin Tức TTXVN, Tiền Phong, TuanVietNam, Tuyên Giáo, Tuổi Trẻ TD, Tài Chính, Tạp chí Xây dựng Đảng, Tạp chí công sản, Tổ Quốc, VEF, VNCA, VNEWS, VOV, VTC, Vietnam Finance, VietnamNet, VietnamPlus, VnEconomy, VnMedia, Văn Hiến, Văn Hóa Xây Dựng Đảng, Zing, Ôtô - xe máy, Ôtô Xe Máy, ĐCSVN, ĐS&PL, ĐTCK, Đại Đoàn Kết Việt, Đầu Thủ, Đầu Tư, Đời Sống Plus

<https://github.com/binhvq/news-corpus>

Sample BERT-based Reading Comprehension System Demonstration

Paragraph

Google là một công ty Internet có trụ sở tại Hoa Kỳ, được thành lập vào năm 1995. Sản phẩm chính của công ty này là công cụ tìm kiếm Google, được nhiều người đánh giá là công cụ tìm kiếm hữu ích và mạnh mẽ nhất trên Internet. Trụ sở của Google tên là "Googleplex" tại Mountain View, California. Giám đốc công ty là Larry Page, một trong hai người sáng lập ra công ty. Tên "Google" là một lỗi chính tả của từ googol, bằng 10¹⁰⁰. Google chọn tên này để thể hiện sứ mệnh của công ty để sắp xếp số lượng thông tin khổng lồ trên mạng. Googleplex, tên của trụ sở Google, có nghĩa là 10googol.

Question

giám đốc của Google là ai?

Paragraph must be in English or Vietnamese.

Answer: Larry Page, - Confidence: 99.0 %

Answer: Larry - Confidence: 0.0 %

Answer: Larry Page, một trong hai - Confidence: 0.0 %

Question should end with ?

Choose model type

Vietnamese

Run

Thank you!

External Slides for BERT

Jacob Devlin's slides:

<https://nlp.stanford.edu/seminar/details/jdevlin.pdf>

ICML slides 40-51:

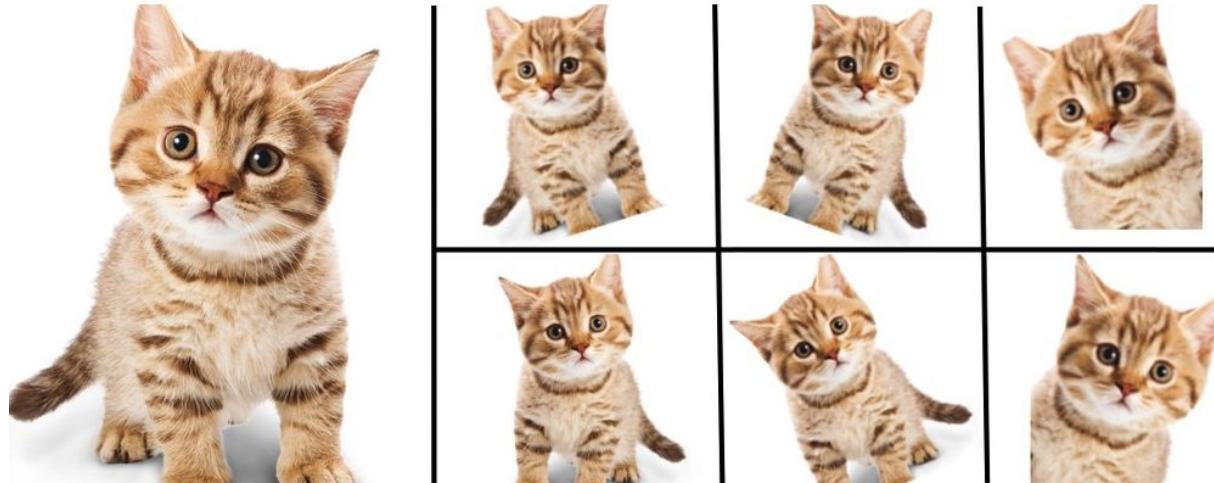
<https://media.neurips.cc/Conferences/NIPS2018/Slides/graves-deeplearning2.pdf>

The Illustrated BERT.

<http://jalammar.github.io/illustrated-bert/>

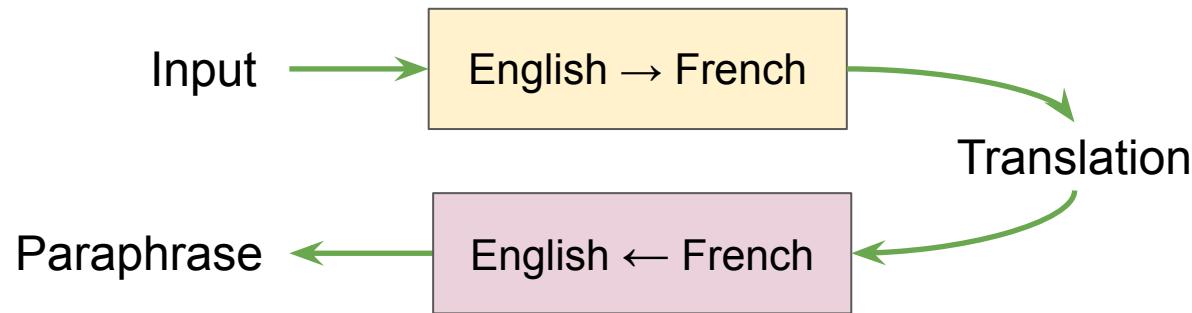
QANet Idea #2: Data Augmentation for Text

Data augmentation: popular in vision & speech, but not NLP



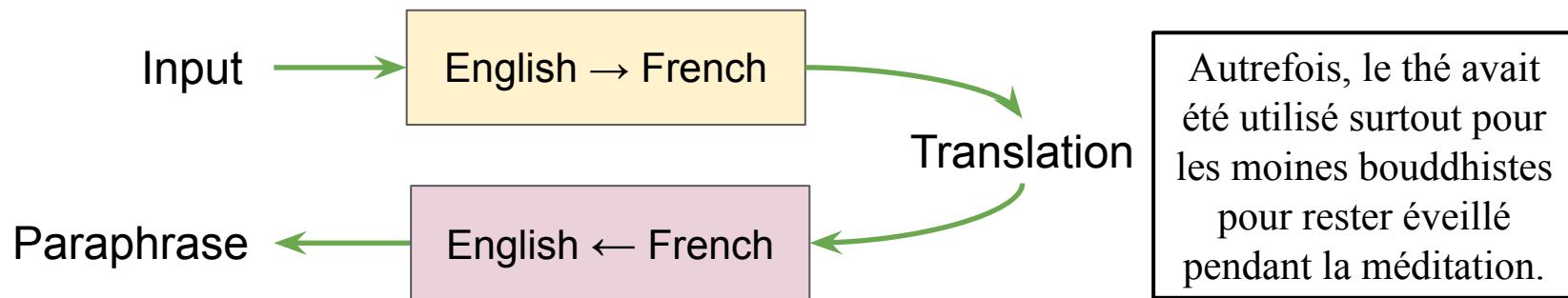
Enlarge your Dataset

More data with NMT back-translation



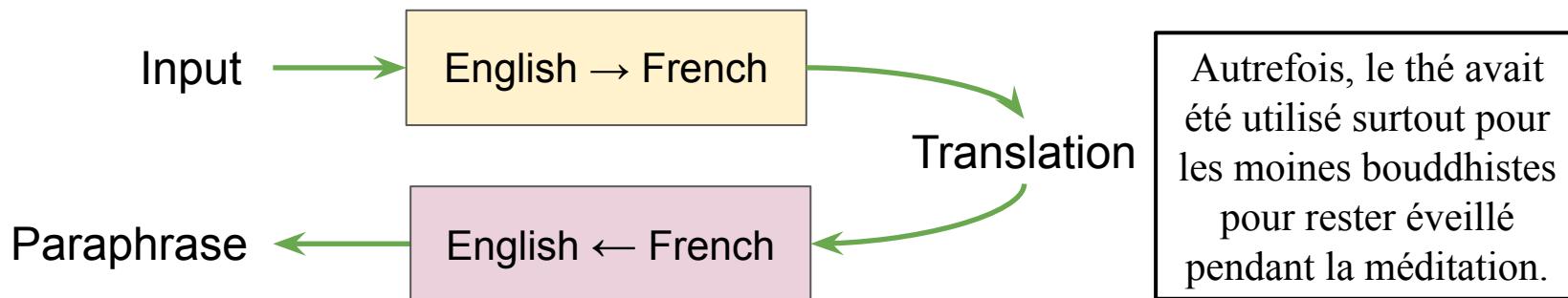
More data with NMT back-translation

Previously, tea had been used primarily for Buddhist monks to stay awake during meditation.



More data with NMT back-translation

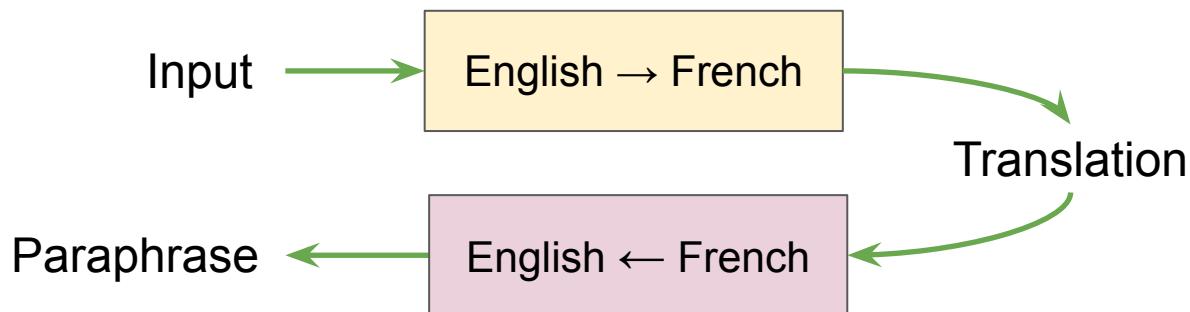
Previously, tea had been used primarily for Buddhist monks to stay awake during meditation.



In the past, tea was used mostly for Buddhist monks to stay awake during the meditation.

More data with NMT back-translation

Previously, tea had been used primarily for Buddhist monks to stay awake during meditation.



In the past, tea was used mostly for Buddhist monks to stay awake during the meditation.

- More data
 - (Input, *label*)
 - (Paraphrase, *label*)

Applicable to virtually any NLP tasks!

QANet augmentation

Original

en-de-en

Question

What individuals live at Fatima House at Notre Dame?

Context

... **Retired priests and brothers** reside in Fatima House (a former retreat center), Holy Cross House, as well as Columba Hall near the Grotto. ...

Answer

Retired priests and brothers

Context

... In the Fatima House (a former retreat), the **tired priests and the brothers** are found in the magnificent church of Saint Cross, the Columbian Hall in the vicinity of the cave. ...

Answer

tired priests and the brothers

QANet augmentation

Original

en-fr-en

Question

What individuals live at Fatima House at Notre Dame?

Context

... **Retired priests and brothers** reside in Fatima House (a former retreat center), Holy Cross House, as well as Columba Hall near the Grotto. ...

Answer

Retired priests and brothers

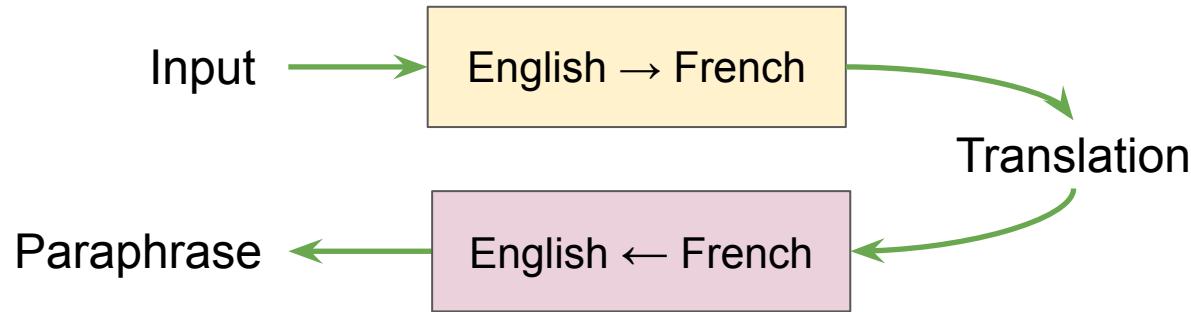
Context

... **Priests and retired brothers** reside at Fatima House (former public centre), the Holy Cross House and the Columbayn Hall near the Cave. ...

Answer

Priests and retired brothers

QANet augmentation



Use 2 language pairs: *English-French*, *English-German*. **3x data**.

Improvement: +1.1 F1

Failure cases

Question

Where by mass is oxygen a major part?

Context

... Oxygen constitutes 49.2 % of the Earth ' s crust by mass and is the major component of the **world ' s oceans** (**88.8 %** by mass) . Oxygen gas is the second most common component of the Earth ' s atmosphere , taking up 20.8 % of its volume and 23.1 % of its mass (some 1015 tonnes) ...

Ground truths

"world's oceans", "the world's oceans", "oceans"

Ensemble Answers

model 0: 88.8%
model 1: 88.8%
model 2: world's oceans
model 3: 88.8%
model 4: the world's oceans
model 5: 88.8%
model 6: 88.8%
model 7: 88.8%
model 8: world's oceans
model 9: world's oceans
model 10: 88.8%
model 11: 88.8%
model 12: world's oceans
model 13: world's oceans

Failure cases

Ensemble Answers

Question

Who else did Tesla make the acquaintance of in 1886?

Context

... In late 1886 Tesla met **Alfred S. Brown** , a Western Union superintendent , and New York attorney **Charles F. Peck** ... Together in April 1887 they formed the Tesla Electric Company with an agreement that profits from generated patents would go $\frac{1}{3}$ to Tesla , $\frac{1}{3}$ to Peck and Brown , and $\frac{1}{3}$ to fund development .

Ground truths

'Charles F. Peck', 'New York attorney Charles F. Peck'

model 0: Alfred S. Brown
model 1: Alfred S. Brown
model 2: Alfred S. Brown
model 3: Alfred S. Brown
model 4: Alfred S. Brown
model 5: Alfred S. Brown
model 6: Alfred S. Brown
model 7: Alfred S. Brown
model 8: Alfred S. Brown
model 9: Alfred S. Brown
model 10: Alfred S. Brown
model 11: Alfred S. Brown, a Western Union superintendent, and New York attorney Charles F. Peck
model 12: Alfred S. Brown
model 13: Peck and Brown

Học AI thế nào?

A typical day for researchers on Google's Brain Team

Posted Sep 15, 2017 by John Mannes (@JohnMannes)

- Everyone universally spends a lot of time reading papers on arXiv that are relevant to their research and collaborating with colleagues. Sara Hooker, a resident on the team, likes to chat with colleagues over breakfast — and lunch and dinner — to keep up to date on what other researchers are doing to solve similar problems.
- Chat with colleagues: *very important!*

<https://techcrunch.com/2017/09/15/a-typical-day-for-researchers-on-googles-brain-team/>

Cập nhật thông tin Deep Learning



- Twitter: most deep learning researchers.



Home



/MachineLearning @slashML · 2h

🔥 Latest Deep Learning OCR with Keras and Supervisely in 15 minutes



[P] 🔥 Latest Deep Learning OCR with Keras and... • r/... reddit.com

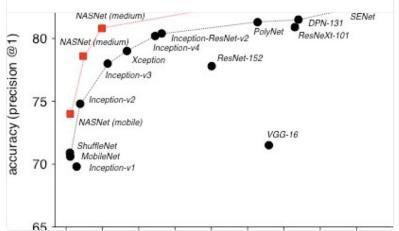


Thang Luong @lmthang · 2h

Neural architecture search work from @GoogleBrain enables state-of-the-art vision models of varying sizes depending on your need!

Google Research @googleresearch

AutoML for large scale image classification and object detection
goo.gl/F6QXbd





Cập nhật thông tin Deep Learning

- Twitter: most deep learning researchers.
- Newsletters
 - Import AI: <https://jack-clark.net/>
 - The Wild Week in AI <https://www.getrevue.co/profile/wildml>

Import AI: #66: Better synthetic images heralds a Fake News future, GraphCore shows why AI chips are about to get very powerful, simulated rooms for better reinforcement learning

Welcome to Import AI, [subscribe here](#).

NVIDIA achieves remarkable quality on synthetic images:

...AKA: The Fake Photo News Tidalwave Cometh

...NVIDIA researchers have released 'Progressive Growing of GANS for Improved QUality, Stability, and Variation' – a training methodology for creating better synthetic images with generative adversarial networks

News

Theano Development Ends

GROUPS.GOOGLE.COM — Share

The MILA group will stop the development of **Theano**, one of the first and most popular Deep Learning frameworks. Many people have started their Deep Learning journey with Theano, which has been responsible for early breakthroughs and served as inspiration for frameworks like Tensorflow. Theano will continue to be available, but will not be maintained.

The Wild Week in AI

October 2 - Issue #63 - [View online](#)

The Wild Week in AI is a weekly AI & Deep Learning newsletter curated by [@dennybritz](#).



Home



/MachineLearning @slashML · 2h

Latest Deep Learning OCR with Keras and Supervisely in 15 minutes



[P] 🔥 Latest Deep Learning OCR with Keras and... • r/... reddit.com

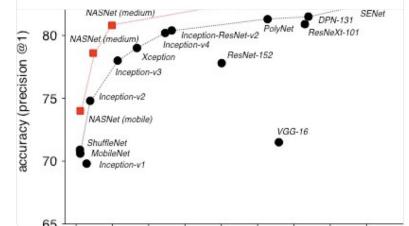


Thang Luong @lmthang · 2h

Neural architecture search work from [@GoogleBrain](#) enables state-of-the-art vision models of varying sizes depending on your need!

Google Research @googleresearch

AutoML for large scale image classification and object detection
goo.gl/F6QXBd



Tài liệu

- Deep learning book
- Online courses:
 - Coursera & Udacity
 - Stanford CS224N (language), CS231N (vision)
- Github
 - Learn from others & create your own

