

# The BERTs family

Longer, bigger, smaller, smarter

# Pre-trained transformer types (by training method)

## Encoder only (autoencoder)

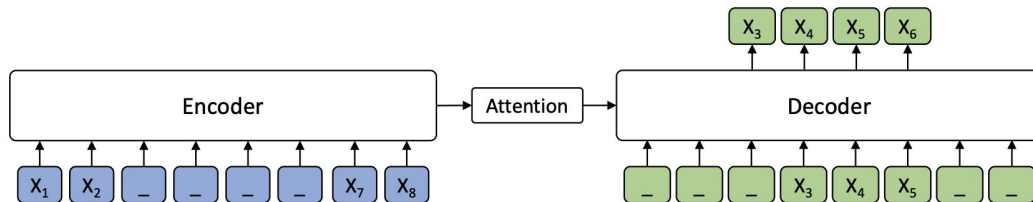
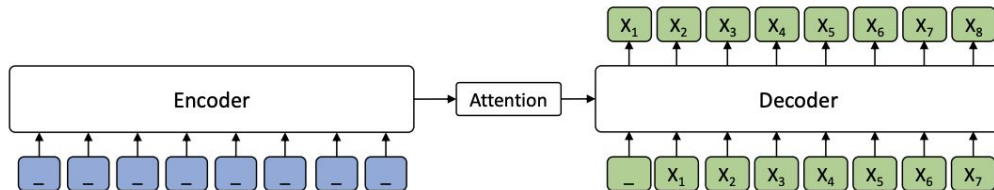
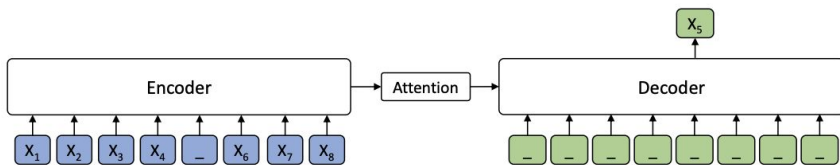
- BERT, ALBERT, RoBERTa
- Seq classification, token classification
- Masked words and predict

## Encoder-decoder (seq2seq)

- MASS, BART, T5
- Machine translation, text summary
- Mask phrases

## Decoder only (autoregressive)

- GPT, CTRL
- Text generation
- Predict next word



# CTRL

For language generation, with style control signal  
Style (embedding)

**Wikipedia** *Anarchism is* a political philosophy that advocates the abolition of all forms of hierarchy and domination, including capitalism, patriarchy, racism, sexism, heterosexism and other oppressive social structures. \n \n The term "anarchism" was coined by Pierre-Joseph Proudhon in his book "The **Books** *Anarchism is* the \n only true and practical form of Socialism. It has been said that Socialism \n will never be established until it can be proved to work in practice. The \n same may be said, with equal truth, about Anarchism. But if we are not able \n to prove its practicability at present, then let us

**Horror** *A knife* handle pulled through the open hole in the front. I jumped when the knife hit.  
  
Eyes widened in horror. Her scream was the only sound I heard besides her sobs.  
  
The spider touched her feet as it started to dig into the top of her arch. The creature's tears began to

**Reviews** *A knife* is a tool and this one does the job well.  
  
Rating: 4.0  
  
I bought these for my husband who has been using them to cut up his own meat since he got them. He says they are very sharp so be careful when you use them, but that doesn't seem like much of an issue because he's used

# Can use url as control signal

Condition on topic, dates, etc.

Links <https://www.cnn.com/2014/09/20/us-president-meets-british-pm>  
JUST WATCHED\n\nObama meets with British PM\n\nMUST WATCH\n\nStory highlights\n\nPresident Barack Obama met with Britain's Prime Minister David Cameron

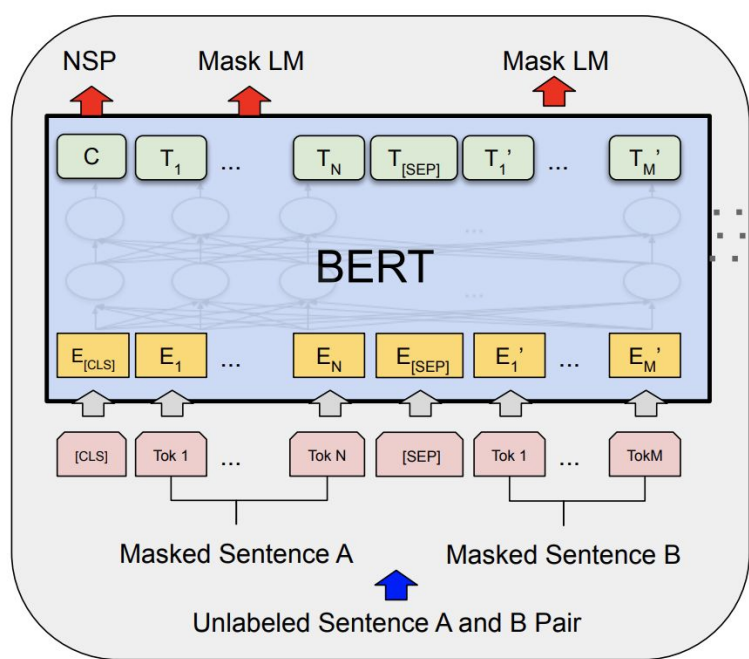
Links <https://www.cnn.com/2018/09/20/us-president-meets-british-pm>  
JUST WATCHED\n\nTrump and May meet for first time\n\nMUST WATCH\n\nWashington (CNN) President Donald Trump, who has been criticized by some in the UK over his decision to leave the European Union, met with British Prime Minister Theresa May, a White House official said on Thursday.

Links <https://www.cnn.com/style/09/20/2018/george-clooney-interview>  
George Clooney on the future of his acting career\n\nBy\n\nUpdated 10:51 AM ET, Thu September 20, 2018\n\nChat with us in Facebook Messenger. Find out what's happening in the world as it unfolds.\n\nPhotos:George Clooney, 'Ocean's 8'\n\nActor George Clooney attends a photocall for "Ocean's 8" at Grauman's Chinese Theatre on August 31, 2018, in Los Angeles.\n\n...

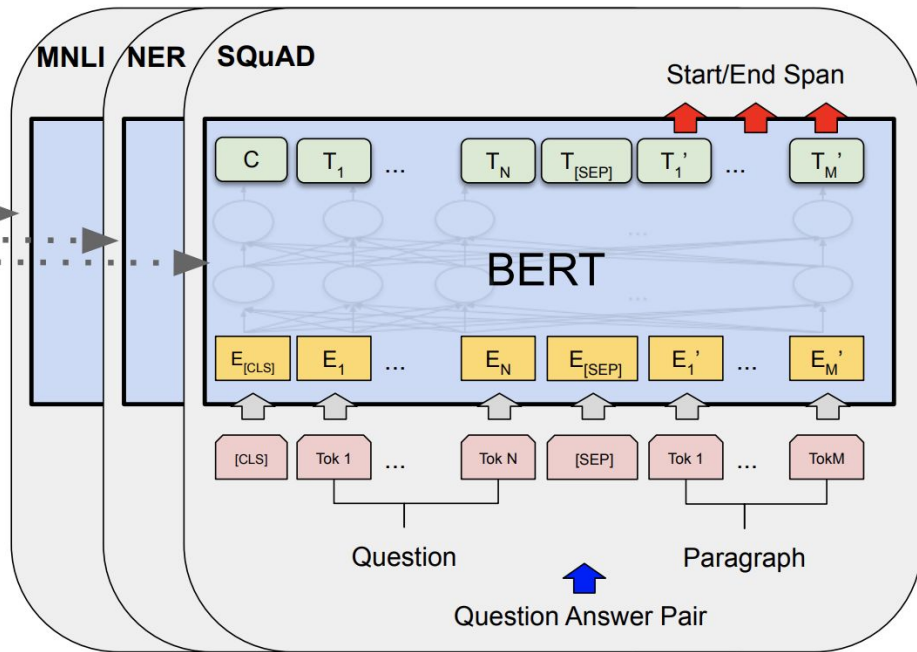
Links <https://www.cnn.com/politics/09/20/2018/george-clooney-interview>  
JUST WATCHED\n\nGeorge Clooney on the Trump administration\n\nMUST WATCH\n\n(CNN) Actor and activist George Clooney, who has been a vocal critic of President Donald Trump, said he is "ready to go back into the political arena" after his role in an anti-Trump documentary was cut from theaters this week.\n\n...

# BERT

Full transformer  
Masked LM to learn bi-directional LM  
Next sentence prediction to learn discourse



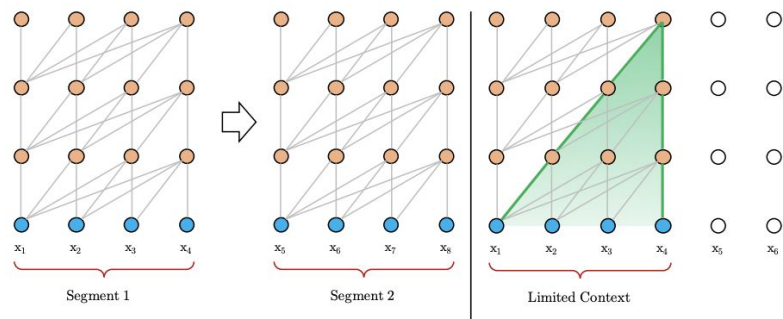
Pre-training



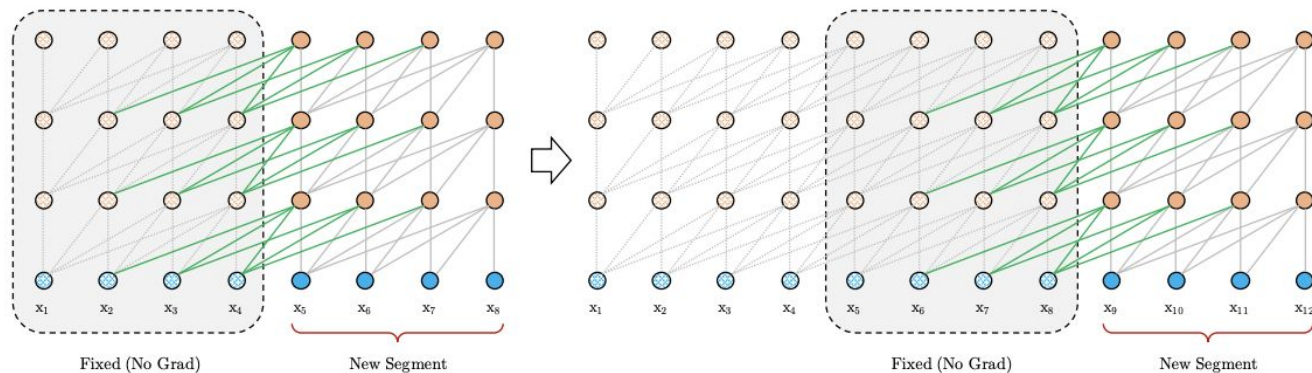
Fine-Tuning



# XL-transformer



Chunking limits the length of context



Let next chunk attention has access to previous chunk

(a) Training phase.



# XLNet

- Transformer-XL + permutation LM
- Mask LM creates mismatch between training and test.
  - Use different permutations to capture different context direction
- Add attention masks to simulate random permutation orders

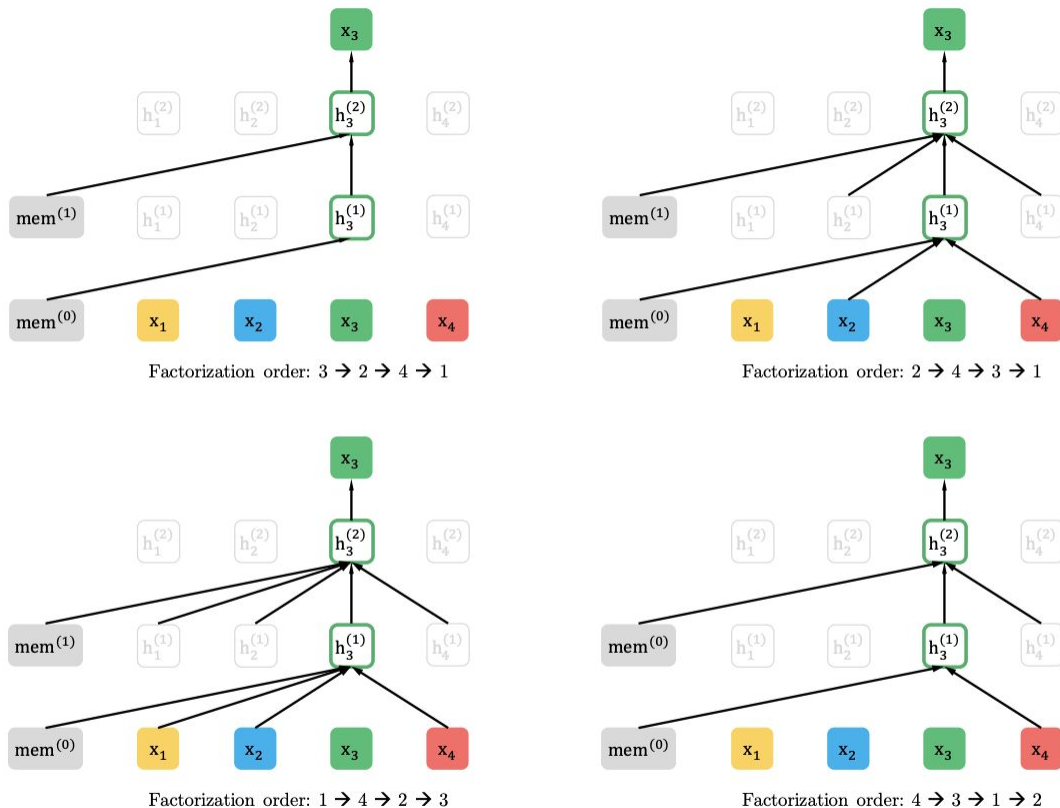


Figure 1: Illustration of the permutation language modeling objective for predicting  $x_3$  given the same input sequence  $x$  but with different factorization orders.

# Roberta (Robustly optimized BERT approach)

A trick and tuning study

Dynamic masking > static

Next sentence prediction is not optimal

Larger batch + higher learning rate

Masking	SQuAD 2.0	MNLI-m	SST-2
reference	76.3	84.3	92.8
<i>Our reimplementation:</i>			
static	78.3	84.3	92.5
dynamic	78.7	84.0	92.9

bsz	steps	lr	ppl	MNLI-m	SST-2
256	1M	1e-4	3.99	84.7	92.7
2K	125K	7e-4	<b>3.68</b>	<b>85.2</b>	<b>92.9</b>
8K	31K	1e-3	3.77	84.6	92.8

# Side note on large batch size training

## Don't Decay the Learning Rate, Increase the Batch Size

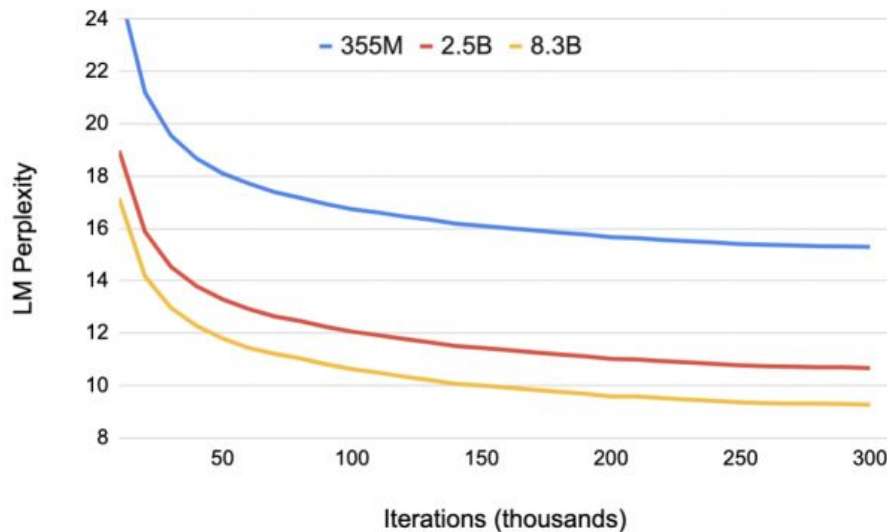
Samuel L. Smith, Pieter-Jan Kindermans, Chris Ying, Quoc V. Le

It is common practice to decay the learning rate. Here we show one can usually obtain the same learning curve on both training and test sets by instead increasing the batch size during training. This procedure is successful for stochastic gradient descent (SGD), SGD with momentum, Nesterov momentum, and Adam. It reaches equivalent test accuracies after the same number of training epochs, but with fewer parameter updates, leading to greater parallelism and shorter training times. We can further reduce the number of parameter updates by increasing the learning rate  $\epsilon$  and scaling the batch size  $B \propto \epsilon$ . Finally, one can increase the momentum coefficient  $m$  and scale  $B \propto 1/(1 - m)$ , although this tends to slightly reduce the test accuracy. Crucially, our techniques allow us to repurpose existing training schedules for large batch training with no hyper-parameter tuning. We train ResNet-50 on ImageNet to 76.1% validation accuracy in under 30 minutes.

Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour  
Don't Decay the Learning Rate, Increase the Batch Size

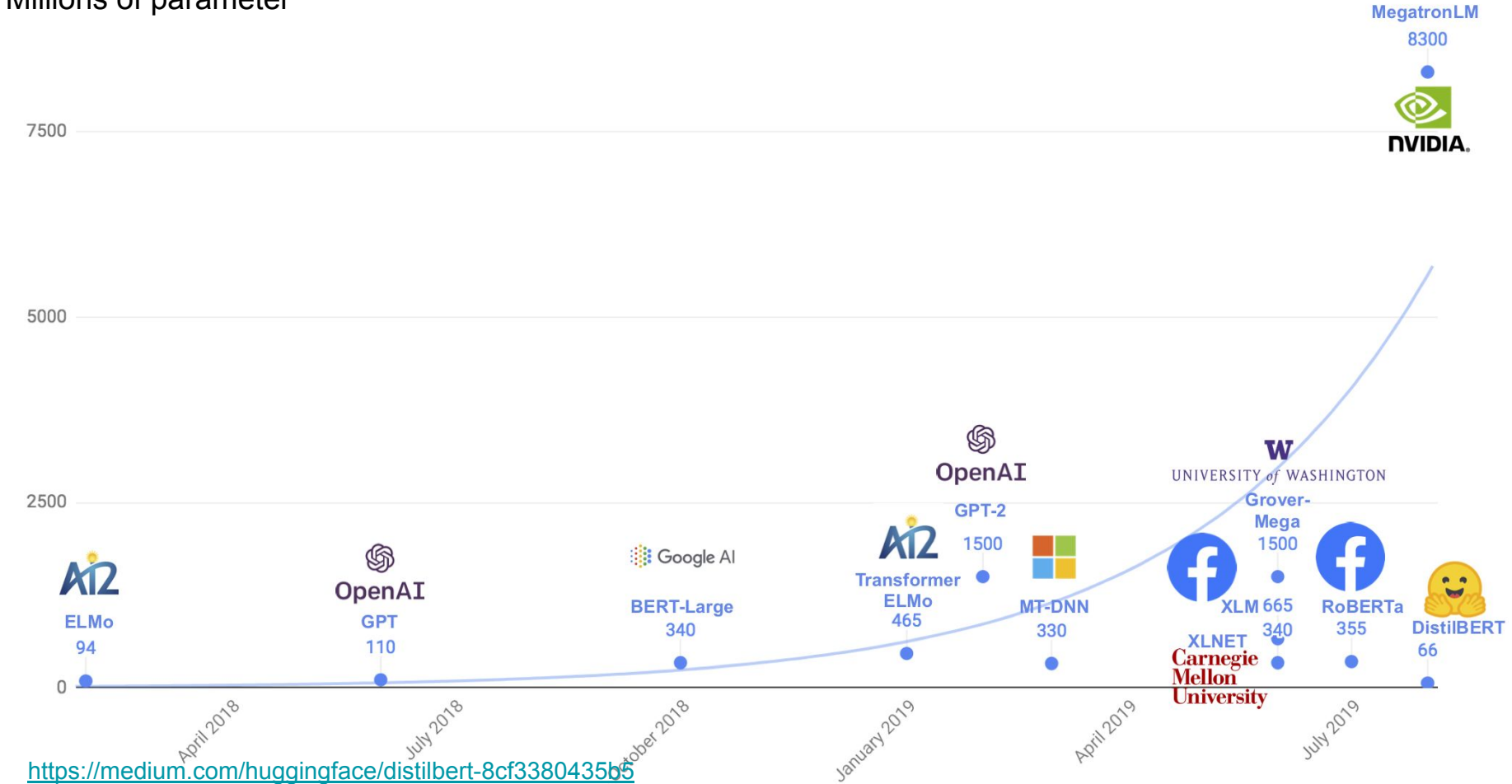
# Megatron-LM

Distributed training



*Figure 7.* Validation set perplexity. All language models are trained for 300k iterations. Larger language models converge noticeably faster and converge to lower validation perplexities than their smaller counterparts.

Millions of parameter



<https://medium.com/huggingface/distilbert-8cf3380435b5>

# Distill bert

Knowledge distillation to get smaller models

Reduce the # of transformer layers by half. Use tricks in Roberta.

Use KL-divergence between teacher and student model

“Cheaper training”

eight 16GB V100 GPUs for approximately three and a half days

	Nb of parameters (millions)	Inference Time (s)
GLUE BASELINE (ELMo + BiLSTMs)	180	895
BERT base	110	668
DistilBERT	66	410

	Macro Score	CoLA	MNLI	MNLI-MM	MRPC		QNLI	QQP		RTE	SST-2	STS-B		WNLI
		mcc	acc	acc	acc	f1	acc	acc	f1	acc	acc	pearson	spearmanr	acc
GLUE BASELINE (ELMo + BiLSTMs)	68.7	44.1	68.6 (avg)		70.8	82.3	71.1	88.0	84.3	53.4	91.5	70.3	70.5	56.3
BERT base	78.0	55.8	83.7	84.1	86.3	90.5	91.1	90.9	87.7	68.6	92.1	89.0	88.6	43.7
DistilBERT	75.2	42.5	81.6	81.1	82.4	88.3	85.5	90.6	87.7	60.0	92.7	84.5	85.0	55.6



# ALBERT

Want higher hidden units without growing the model. Factorized embedding matrix

$$V \times E \rightarrow V \times E + E \times H$$

Share attention layer parameters across layers. More stable training as a side effect.

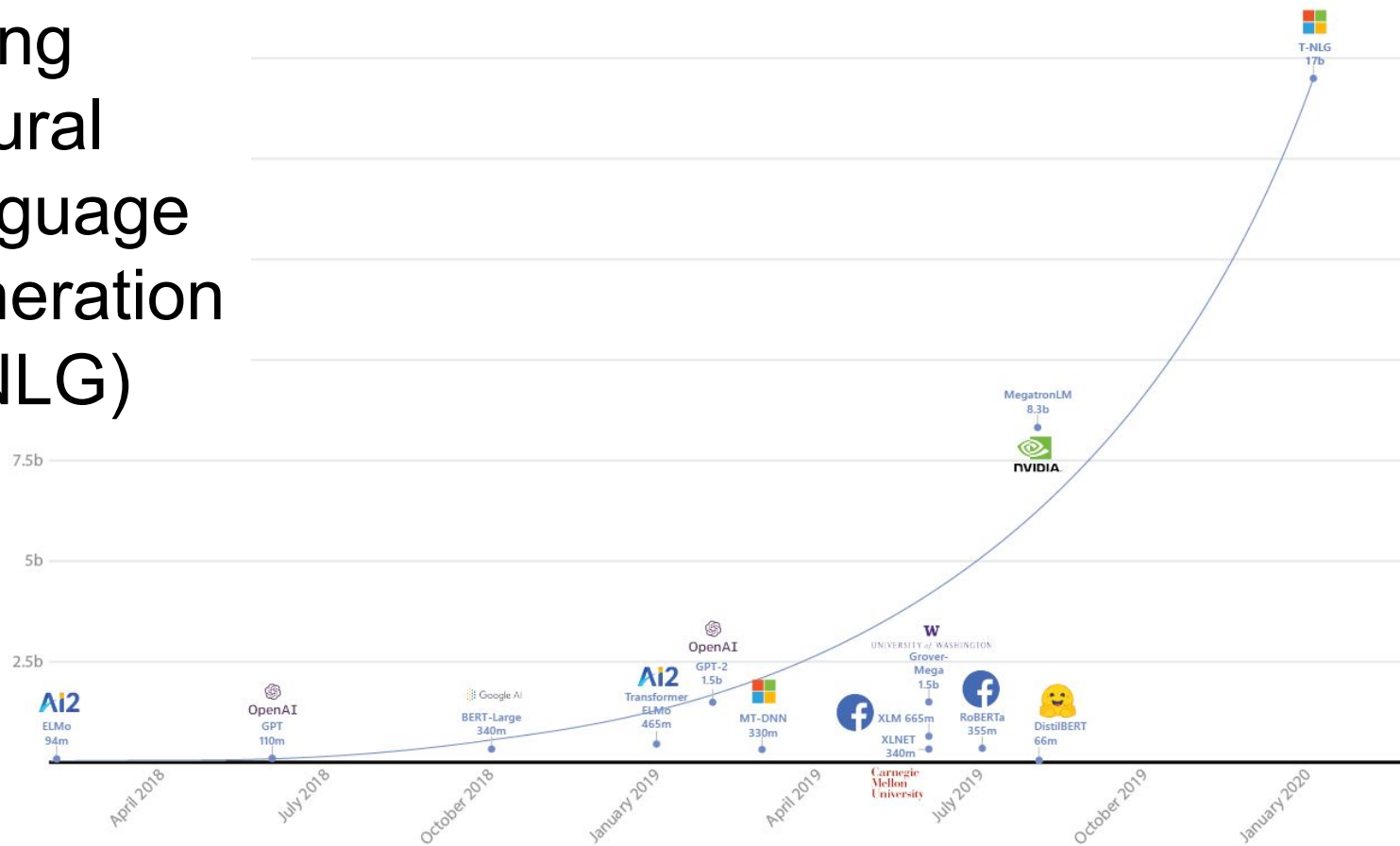
NSP is easy compared to LM tasks (multi-task imbalance)

Next sentence prediction (random sentence) -> Sentence order prediction (swapped sentence or not)

Model		Parameters	Layers	Hidden	Embedding	Parameter-sharing	Avg	Speedup
BERT	base	108M	12	768	768	False	82.1	17.7x
	large	334M	24	1024	1024	False	85.1	3.8x
	xlarge	1270M	24	2048	2048	False	76.7	1.0
ALBERT	base	12M	12	768	128	True	80.1	21.1x
	large	18M	24	1024	128	True	82.4	6.5x
	xlarge	59M	24	2048	128	True	85.5	2.4x
	xxlarge	233M	12	4096	128	True	<b>88.7</b>	1.2x

Some experiments show dropout hurt performance

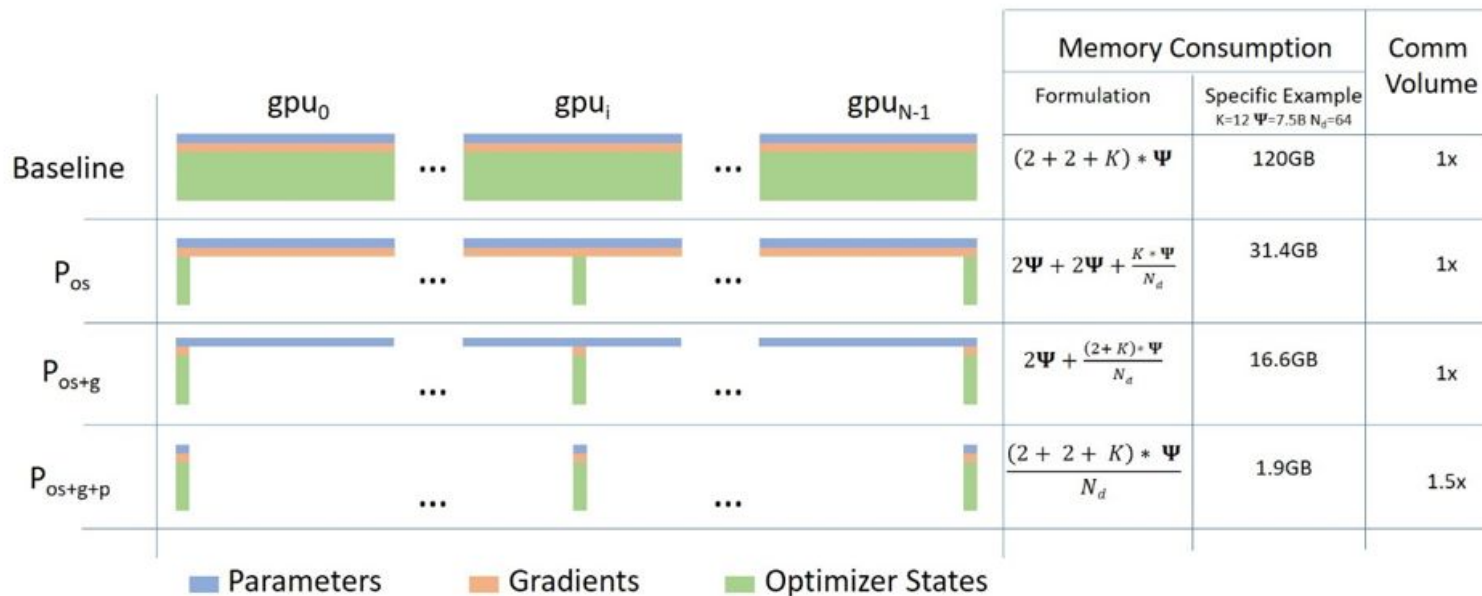
# Turing Natural Language Generation (T-NLG)



<https://www.microsoft.com/en-us/research/blog/turing-nlg-a-17-billion-parameter-language-model-by-microsoft/>

# Zero & DeepSpeed

Clever management of training parameters across GPUs and machines



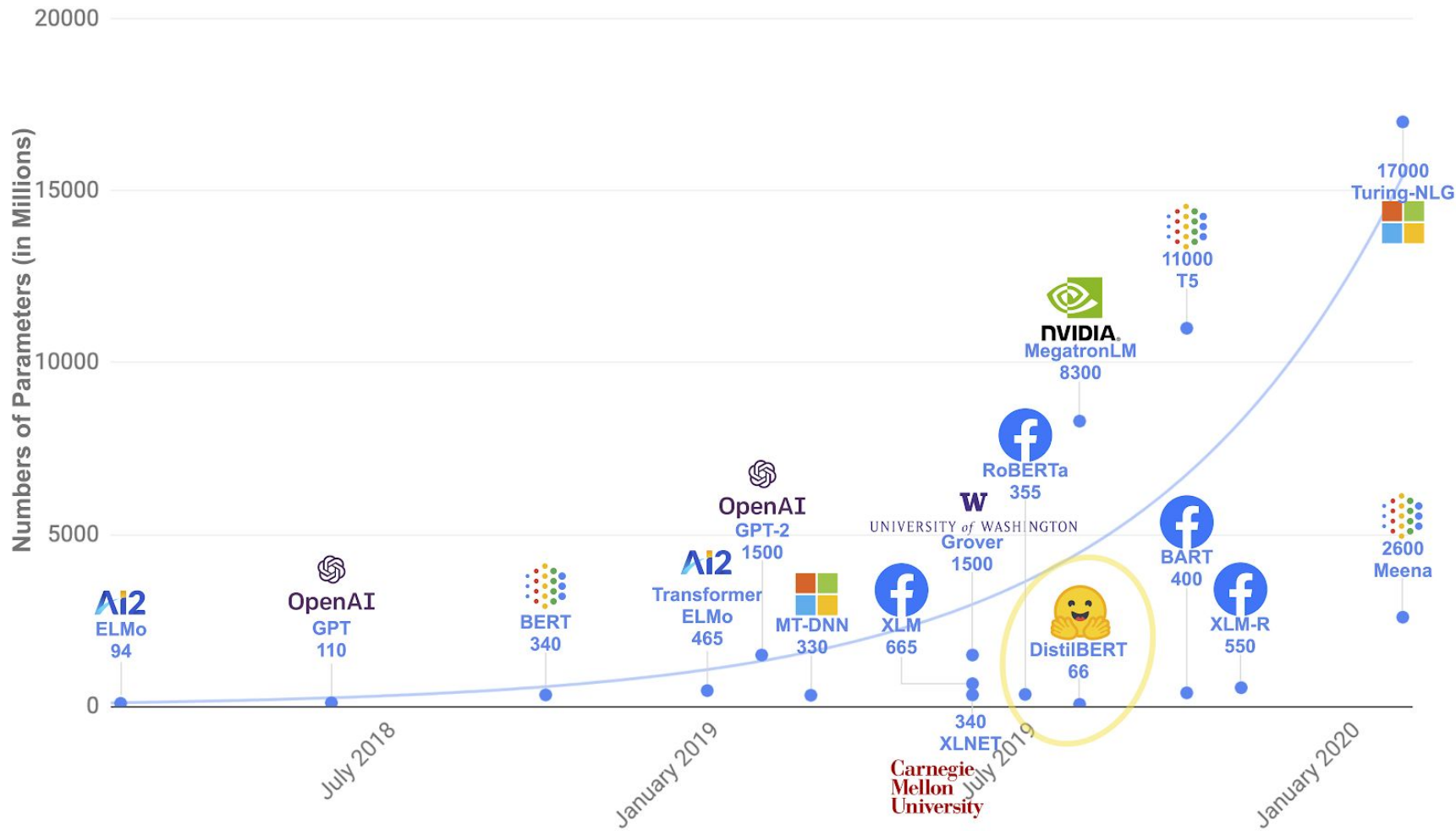
# T-NLG

	LAMBADA (acc) strict	WikiText-103 (test adj. ppl)
Open AI GPT-2 1.5B	52.66 (63.24)*	17.48
Megatron-LM 8.3B	66.51	10.81
<b>T-NLG 17B</b>	<b>67.98</b>	<b>10.21</b>

\*Open AI used additional processing (stopword filtering) to achieve higher numbers than the model achieved alone. Neither Megatron nor T-NLG use this stopwords filtering technique.

Can do Q/A by just LOTs of internet text

<b>When did WW2 end?</b>	WW2 ended in 1945.
<b>How many people live in the US?</b>	There are over 300 million people living in the US.



# Longer, bigger, smaller, smarter

XLNet

Megatron

Roberta

Distill Bert

Albert

T-NLG



# Tricks to make better transformers

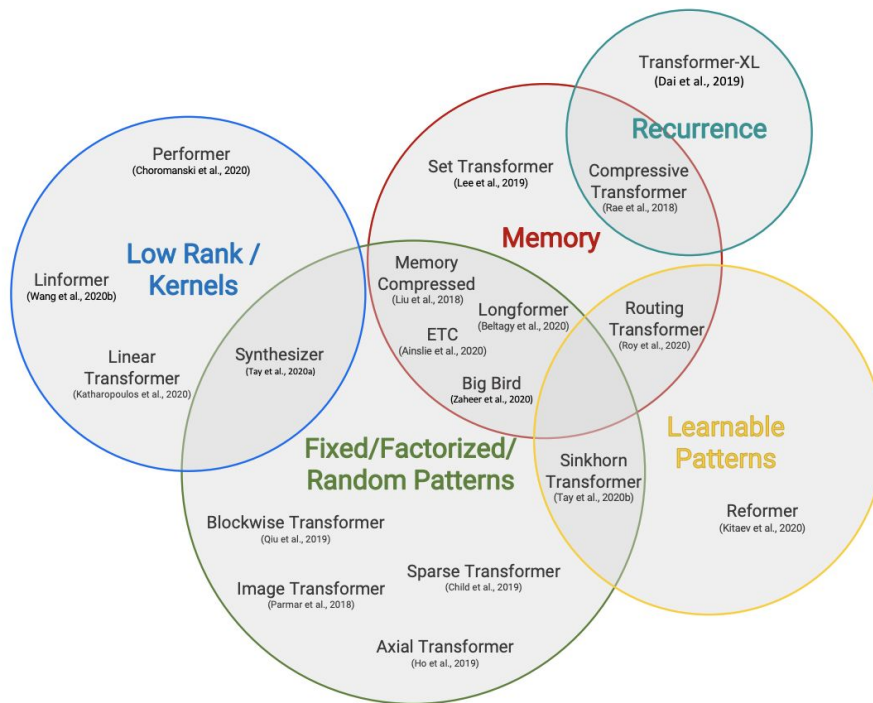
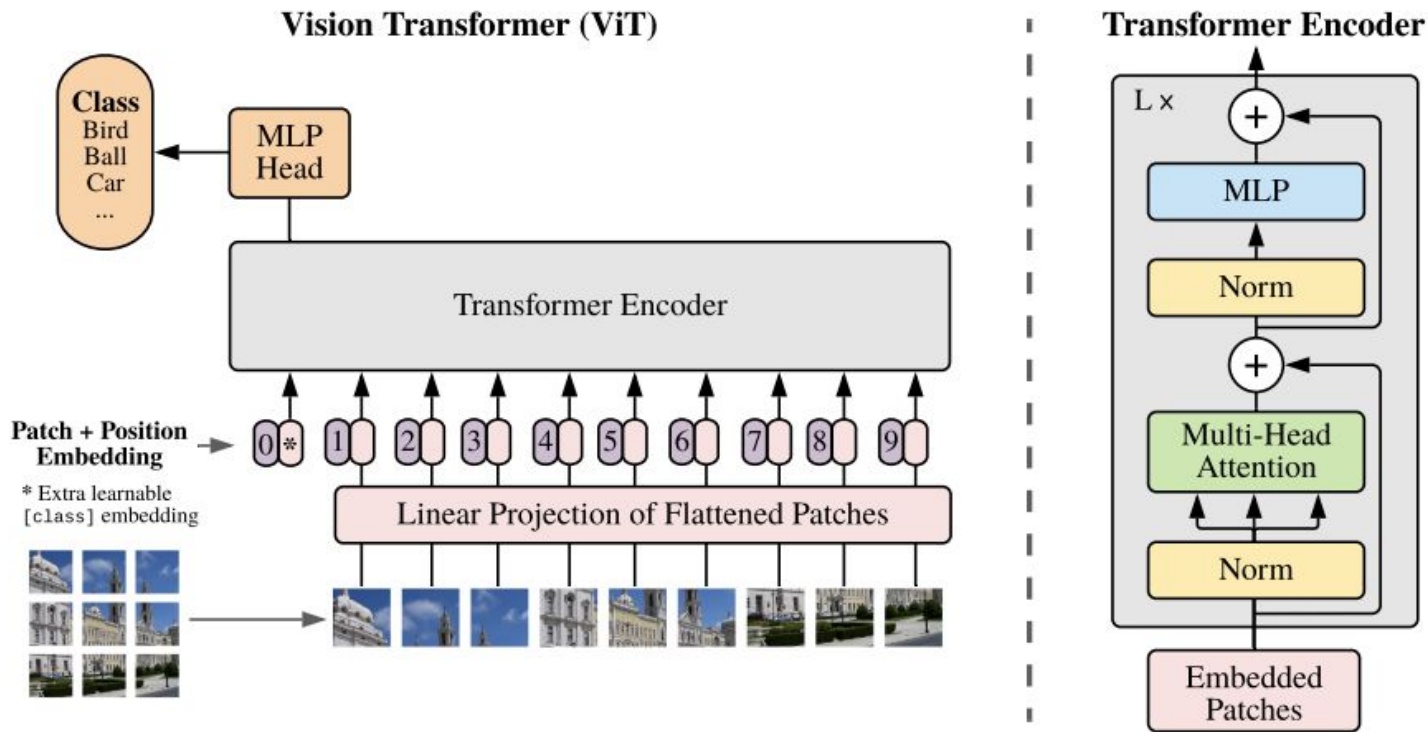


Figure 2: Taxonomy of Efficient Transformer Architectures.

# Vision transformer



AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

<https://arxiv.org/pdf/2010.11929v1.pdf>

# Dall E

GPT3 for images

Cluster image  
patches into tokens  
and combined with  
text to train an  
autoregressive  
model



(a) a tapir made of accordion.  
a tapir with the texture of an  
accordion.

(b) an illustration of a baby  
hedgehog in a christmas  
sweater walking a dog

(c) a neon sign that reads  
"backprop". a neon sign that  
reads "backprop". backprop  
neon sign

(d) the exact same cat on the  
top as a sketch on the bottom

*Figure 2.* With varying degrees of reliability, our model appears to be able to combine distinct concepts in plausible ways, create anthropomorphized versions of animals, render text, and perform some types of image-to-image translation.

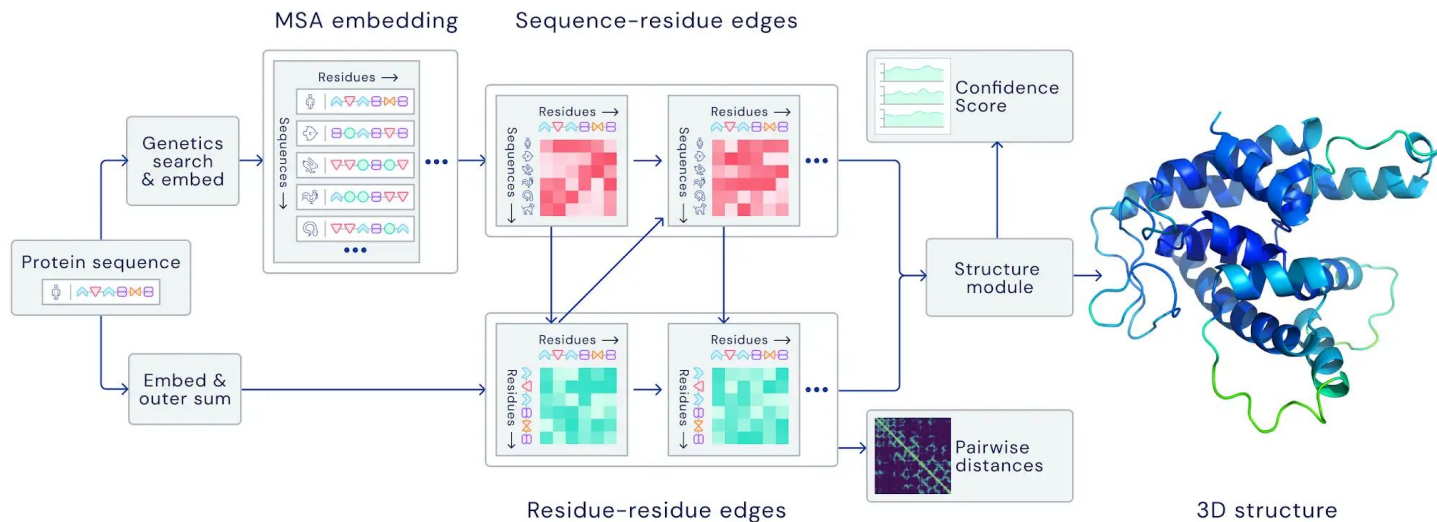
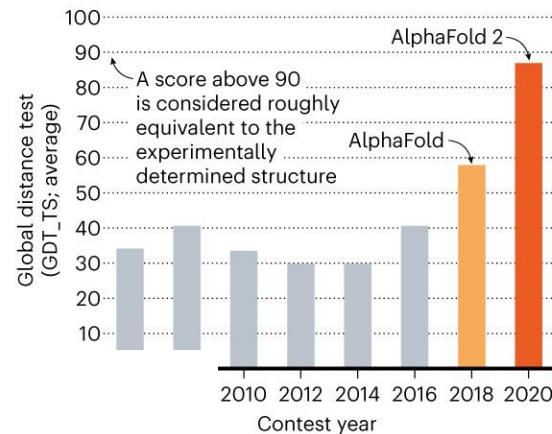
# AlphaFold2

<https://deepmind.com/blog/article/alphafold-a-solution-to-a-50-year-old-grand-challenge-in-biology>

<https://www.nature.com/articles/d41586-020-03348-4>

## STRUCTURE SOLVER

DeepMind's AlphaFold 2 algorithm significantly outperformed other teams at the CASP14 protein-folding contest — and its previous version's performance at the last CASP.



# Homework

Working with WangchanBERTa

A Roberta model

Different pretrained models with different domains/tokenizations

# Schedule

29-Mar-2021	3-Apr-2021	10	Recent Research in NLP Project Announcement + Paper Announcement
5-Apr-2021	10-Apr-2021	11	NLP Application 1 (Guest); 5-Apr-2021
12-Apr-2021	17-Apr-2021		Songkran Holiday
19-Apr-2021	24-Apr-2021	12	Paper Presentation & Progress Report
26-Apr-2021	1-May-2021	13	NLP Application 2 (Guest); 26-Apr-2021
3-May-2021	8-May-2021	14	Project Presentation



# HOW TO READ A SCIENTIFIC ARTICLE

## 2 Paper types

- Review article/tutorial

- Give insights about the field
- Useful for learning about a new field
- Read multiple to avoid the author's bias
- Title usually has “review” or “tutorial”

- Primary research article

- More details on the experiments and results

# Parts of an article

- Abstract
- Introduction
- Methods
- Results and discussion
- Conclusion
- Reference

# Things to look for before reading an article

- Publication date
- Author names
  - Previous and newer publications
- Keywords
- Acknowledgements and funding sources

## REFORMER: THE EFFICIENT TRANSFORMER

**Nikita Kitaev\***

U.C. Berkeley & Google Research  
kitaev@cs.berkeley.edu

**Łukasz Kaiser\***

Google Research  
{lukaszkaizer, levskaya}@google.com

**Anselm Levskaya**

Google Research

# Getting the big picture

- Read the abstract
- Read the introduction
  - What is the research question?
  - What is the method?
  - What had been done? How is it different from other work?
- Look at figures and results

Tip: keep track of terms you don't understand

# First reading

- Reread the introduction
- Skim methods
- Read results and discussion
  - Does the figures make sense now?
- Write on the article!



# Understanding the article

- Reread the article (until you get what you want)
- Check references for parts you don't understand
- Reread the abstract
  - Does your understanding match the abstract?
- Note down important points. This might come in handy when you write your paper/thesis!

# Evaluating the article

- Does the method make sense?
  - What are the limitations that the authors mention?
  - Are there other limitations?
  - Can it be used in other situations?
- Are the experiments legitimate?
  - The sample size is big enough?
  - What kind of dataset is used? How big?
  - The evaluation criterion is sound?
- Have these results been reproduced?
  - Look for articles that cite this paper

# ML paper checklist

- What is being done?
- How is it being done?
  - How is it different from previous work
- What is the dataset?
  - Nature of dataset
  - How many training/testing samples? How many classes/vocab size?
- Evaluation metric
  - What are the baselines?
- Practicality
  - Prone to parameter tuning?
  - Computing resource / Runtime (training and testing)

# Useful tools

- <https://scholar.google.com>

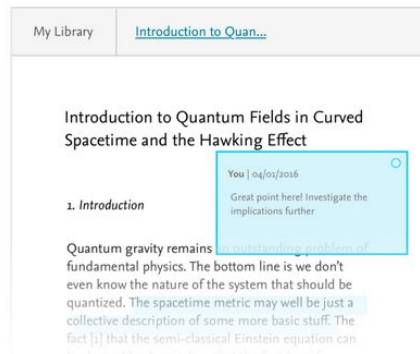
- For finding other articles by the same authors or paper that cites the article

- <https://www.mendeley.com/>

- Reference manager

## Annotate as you read

Easily add your thoughts on documents in your own library, even from mobile devices. For ease of collaboration, you can also share documents with groups of colleagues and annotate them together.



# Useful resources

Most famous NLP papers have blogs explaining them nowadays

Huggingface usually implements a paper within a couple weeks

<https://lilianweng.github.io/lil-log/>

Yannic's youtube channel

<https://www.youtube.com/channel/UCZHmQk67mSJgfCCTn7xBfew>

# Project

Group of 3-5 people

Anything text/NLP related

Must has some application component (cannot be purely basic NLP task)

Recent Research in NLP			
29-Mar-2021	3-Apr-2021	10	<a href="#">Project Announcement + Paper Announcement</a>
5-Apr-2021	10-Apr-2021	11	NLP Application 1 (Guest); 5-Apr-2021
12-Apr-2021	17-Apr-2021		Songkran Holiday
19-Apr-2021	24-Apr-2021	12	<a href="#">Paper Presentation &amp; Progress Report</a>
26-Apr-2021	1-May-2021	13	NLP Application 2 (Guest); 26-Apr-2021
3-May-2021	8-May-2021	14	<a href="#">Project Presentation</a>

# Summary

Different ways to pretrained models

Different ideas of optimizing transformers

Transformer for everything