# GPT-1 and GPT-2 Review

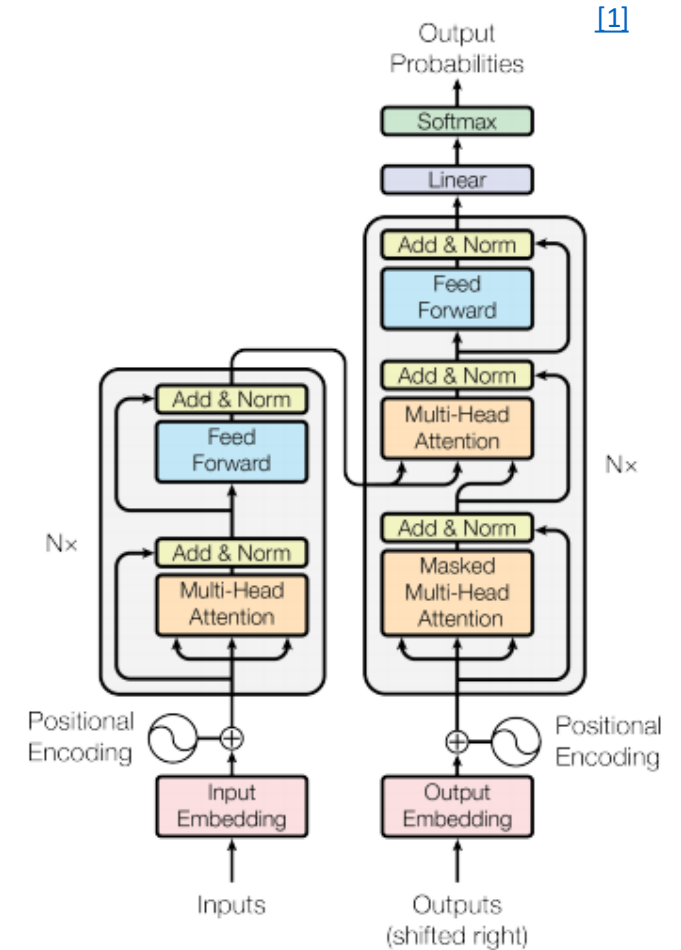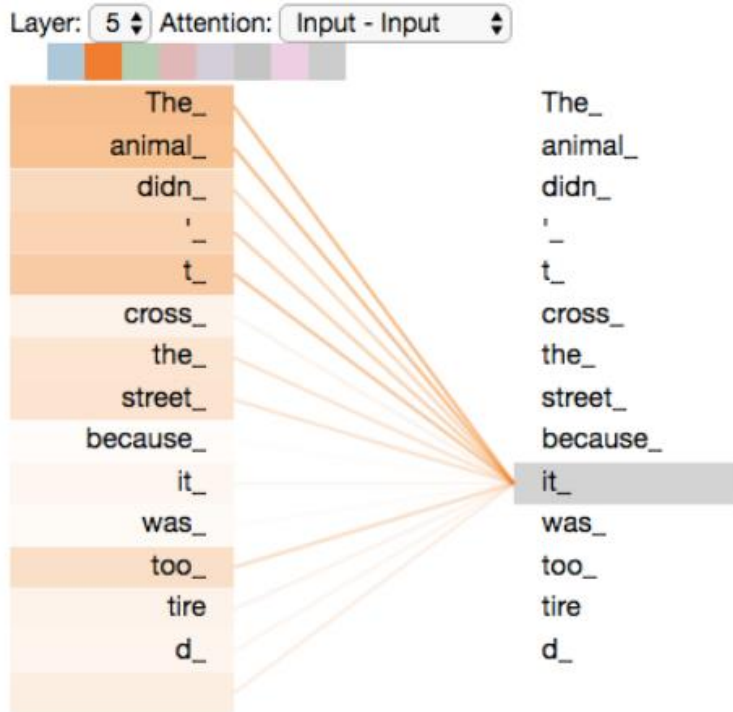Click to add text

Amin Saied

2021-01-15

Background

# Transformers: Life before GPT-1

- Sequence-to-sequence model
- Evolution of RNNs
- Review:
  - Self-attention
  - Multi-headed attention
  - Encoder/decoder

[1] Attention is all you need - Vaswani, Shazeer, Parmar, Uszkoreit, Jones, Gomez, Kaiser, Polosukhin

# Breaking down self-attention

"The animal didn't cross the street because it was too tired"



Layer: 5 ⬍ Attention: Input - Input ⬍

| | |
|---|---|
| The_ | The_ |
| animal_ | animal_ |
| didn_ | didn_ |
| '_ | '_ |
| t_ | t_ |
| cross_ | cross_ |
| the_ | the_ |
| street_ | street_ |
| because_ | because_ |
| it_ | it_ |
| was_ | was_ |
| too_ | too_ |
| tire | tire |
| d_ | d_ |

$$Z := softmax\left(\frac{Q \cdot K^T}{\sqrt{d}}\right) \cdot V$$

|  | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \bullet k_1 = 112$ | $q_1 \bullet k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

The illustrated transformer –Jay Alammar (Blog)
Attention is all you need - Vaswani, Shazeer, Parmar, Uszkoreit, Jones, Gomez, Kaiser, Polosukhin

# Breaking down **multi-headed** self-attention

1) This is our input sentence*

2) We embed each word*

Thinking Machines

**X**

* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

**R**

3) Split into 8 heads. We multiply X or R with weight matrices

$W_0^Q$
$W_0^K$
$W_0^V$

$W_1^Q$
$W_1^K$
$W_1^V$

...

$W_7^Q$
$W_7^K$
$W_7^V$

4) Calculate attention using the resulting Q/K/V matrices

$Q_0$
$K_0$
$V_0$

$Q_1$
$K_1$
$V_1$

...

$Q_7$
$K_7$
$V_7$

5) Concatenate the resulting Z matrices, then multiply with weight matrix $W^O$ to produce the output of the layer

$Z_0$

$Z_1$

...

$Z_7$

$W^O$

$Z$

$Z_0$ $Z_1$ $Z_2$ $Z_3$ $Z_4$ $Z_5$ $Z_6$ $Z_7$

Layer: 5 Attention: Input - Input

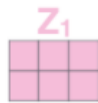The_ · animal_ · didn_ · '_ · t_ · cross_ · the_ · street_ · because_ · it_ · was_ · too_ · tire · d_

The_ · animal_ · didn_ · '_ · t_ · cross_ · the_ · street_ · because_ · it_ · was_ · too_ · tire · d_

ENCODER #1

Add & Normalize

Feed Forward    Feed Forward

Add & Normalize

Self-Attention

POSITIONAL ENCODING ⊕    ⊕

$x_1$ Thinking    $x_2$ Machines

The illustrated transformer —Jay Alammar (Blog)
Attention is all you need - Vaswani, Shazeer, Parmar, Uszkoreit, Jones, Gomez, Kaiser, Polosukhin

# Encoder / decoder

Decoding time step: 1 (2) 3 4 5 6          OUTPUT          I

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT          Je          suis          étudiant          PREVIOUS OUTPUTS          I

K_encdec   V_encdec

ENCODERS          DECODERS

Linear + Softmax

**Encoder**
- Full-context
- (Contextual) Word-embeddings

**Decoder**
- Left-context (masking)
- => predict next word

The illustrated transformer — Jay Alammar (Blog)
Attention is all you need - Vaswani, Shazeer, Parmar, Uszkoreit, Jones, Gomez, Kaiser, Polosukhin

# Glossary

- Tokens
- Attention (self-attention, multi-headed)
- Transformer
- Encoder / decoder

GPT 1

# GPT-1

Improving Language Understanding by Generative Pre-Training – Radford et al

**Key takeaways**

- Semi-supervised learning with transformers
  - Pretraining / finetuning
- Decoder-only architecture
- Simplified approach to transfer learning

=>

- SOTA in 9/12 tasks studied

# Language modelling (unsupervised approach)

**3.1  Unsupervised pre-training**

Given an unsupervised corpus of tokens $\mathcal{U} = \{u_1, \ldots, u_n\}$, we use a standard language modeling objective to maximize the following likelihood:

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \ldots, u_{i-1}; \Theta) \qquad (1)$$

where $k$ is the size of the context window, and the conditional probability $P$ is modeled using a neural network with parameters $\Theta$. These parameters are trained using stochastic gradient descent [51].

- Different from above: this is unsupervised!

- Training data:
  - Data: Edon Lulzim Zhegrova (born 31 March 1999) is a Kosovan professional footballer who plays as a right winger for Swiss club Basel
  - Input: Edon Lulzim Zhegrova (born 31 March 1999) is a Kosovan professional
  - Output: Edon Lulzim Zhegrova (born 31 March 1999) is a Kosovan professional footballer

- Jargon: Auto-regressive language modelling

- Transfer learning in NLP!

# Decoder-only architecture

- Based on previous work [2] using decoder-only transformer to generate Wikipedia articles

- Key-insight [2]: convert seq-to-seq task into language modelling task
  - Seq-to-seq: $(x_1, \ldots, x_m) \mapsto (y_1, \ldots, y_n)$
  - LM: $(x_1, \ldots, x_m, \delta, y_1, \ldots, y_n)$, where $\delta$=separator token

$$p(w^1, \ldots, w^{n+\eta}) = \prod_{j=1}^{n+\eta} p(w^i | w^1, \ldots, w^{j-1})$$

- [1]: Semi-supervised approach!

[2] P. J. Liu, M. Saleh, E. Pot, B. Goodrich, R. Sepassi, L. Kaiser, and N. Shazeer. Generating wikipedia by summarizing long sequences. ICLR, 2018.

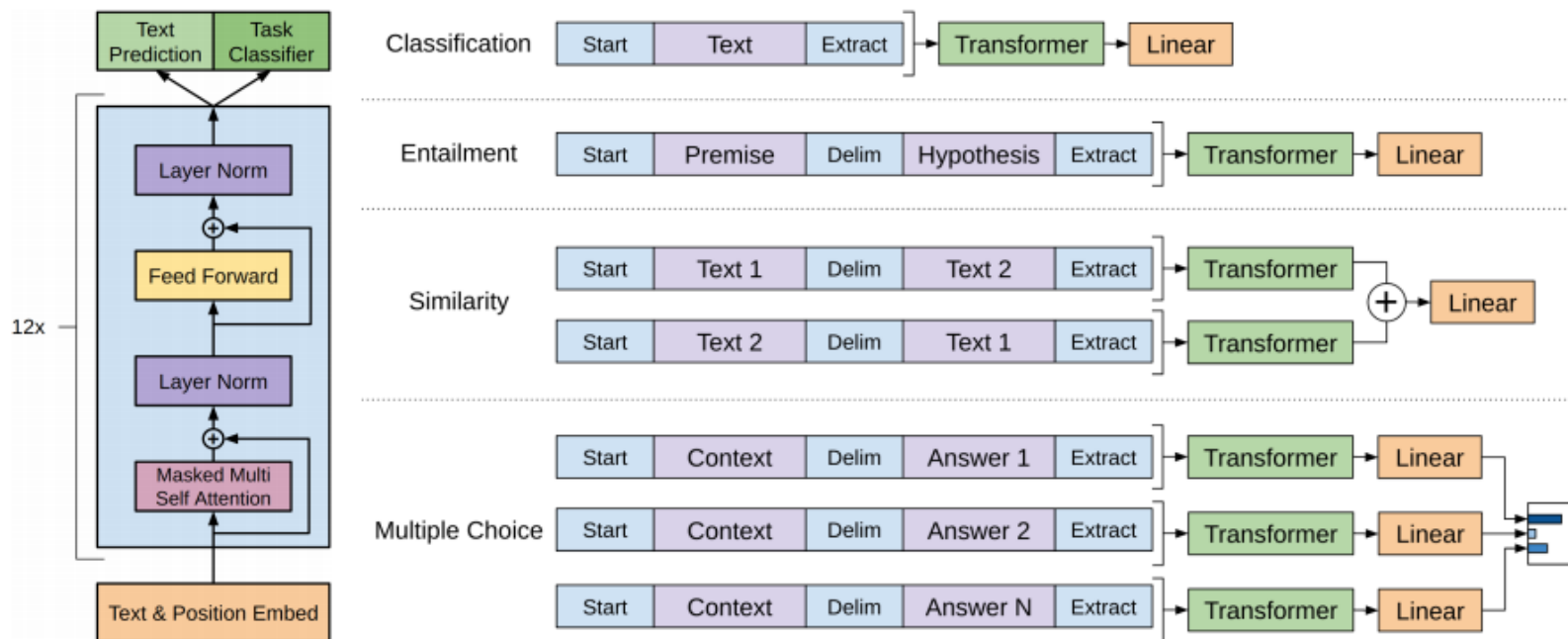[1] Improving Language Understanding by Generative Pre-Training – Radford et al

## 3.2 Supervised fine-tuning

After training the model with the objective in Eq. 1, we adapt the parameters to the supervised target task. We assume a labeled dataset $\mathcal{C}$, where each instance consists of a sequence of input tokens, $x^1, \ldots, x^m$, along with a label $y$. The inputs are passed through our pre-trained model to obtain the final transformer block's activation $h_l^m$, which is then fed into an added linear output layer with parameters $W_y$ to predict $y$:

$$P(y|x^1, \ldots, x^m) = \mathtt{softmax}(h_l^m W_y). \tag{3}$$

This gives us the following objective to maximize:

$$L_2(\mathcal{C}) = \sum_{(x,y)} \log P(y|x^1, \ldots, x^m). \tag{4}$$



Improving Language Understanding by Generative Pre-Training – Radford et al

# Glossary

- Tokens
- Attention (self-attention, multi-headed)
- Transformer
- Encoder / decoder
- **Pretrain / Finetune**
- **Language modelling**
- **Auto-regressive**

# Experimental Results

**Model specifications**   Our model largely follows the original transformer work [62]. We trained a 12-layer decoder-only transformer with masked self-attention heads (768 dimensional states and 12 attention heads). For the position-wise feed-forward networks, we used 3072 dimensional inner states. We used the Adam optimization scheme [27] with a max learning rate of 2.5e-4. The learning rate was increased linearly from zero over the first 2000 updates and annealed to 0 using a cosine schedule. We train for 100 epochs on minibatches of 64 randomly sampled, contiguous sequences of 512 tokens. Since layernorm [2] is used extensively throughout the model, a simple weight initialization of $N(0, 0.02)$ was sufficient. We used a bytepair encoding (BPE) vocabulary with 40,000 merges [53] and residual, embedding, and attention dropouts with a rate of 0.1 for regularization. We also employed a modified version of L2 regularization proposed in [37], with $w = 0.01$ on all non bias or gain weights. For the activation function, we used the Gaussian Error Linear Unit (GELU) [18]. We used learned position embeddings instead of the sinusoidal version proposed in the original work. We use the *ftfy* library[2] to clean the raw text in BooksCorpus, standardize some punctuation and whitespace, and use the *spaCy* tokenizer.[3]

- 12 layer decoder
- 768 dim hidden states
- 12 attention heads (multi-headed attention)

| Method | Classification | | Semantic Similarity | | | GLUE |
|---|---|---|---|---|---|---|
| | CoLA (mc) | SST2 (acc) | MRPC (F1) | STSB (pc) | QQP (F1) | |
| Sparse byte mLSTM [16] | - | **93.2** | - | - | - | - |
| TF-KLD [23] | - | - | **86.0** | - | - | - |
| ECNU (mixed ensemble) [60] | - | - | - | 81.0 | - | - |
| Single-task BiLSTM + ELMo + Attn [64] | 35.0 | 90.2 | 80.2 | 55.5 | 66.1 | 64.8 |
| Multi-task BiLSTM + ELMo + Attn [64] | 18.9 | 91.6 | 83.5 | 72.8 | 63.3 | 68.9 |
| Finetuned Transformer LM (ours) | **45.4** | 91.3 | 82.3 | **82.0** | 70.3 | 72.8 |

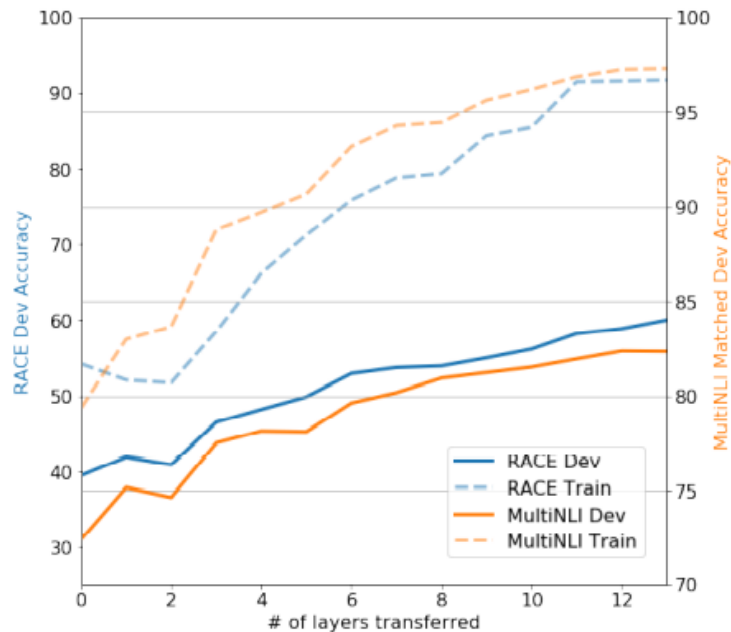| Task | Datasets |
|---|---|
| Natural language inference | SNLI [5], MultiNLI [66], Question NLI [64], RTE [4], SciTail [25] |
| Question Answering | RACE [30], Story Cloze [40] |
| Sentence similarity | MSR Paraphrase Corpus [14], Quora Question Pairs [9], STS Benchmark [6] |
| Classification | Stanford Sentiment Treebank-2 [54], CoLA [65] |

Question Answering

| Method | Story Cloze | RACE-m | RACE-h | RACE |
|---|---|---|---|---|
| val-LS-skip [55] | 76.5 | - | - | - |
| Hidden Coherence Model [7] | 77.6 | - | - | - |
| Dynamic Fusion Net [67] (9x) | - | 55.6 | 49.4 | 51.2 |
| BiAttention MRU [59] (9x) | - | 60.2 | 50.3 | 53.3 |
| Finetuned Transformer LM (ours) | **86.5** | **62.9** | **57.4** | **59.0** |

Natural language inference

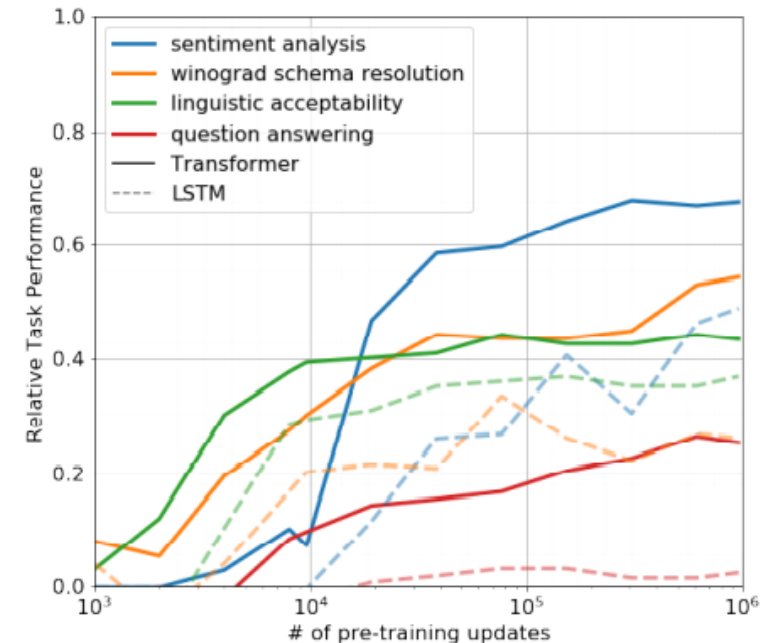| Method | MNLI-m | MNLI-mm | SNLI | SciTail | QNLI | RTE |
|---|---|---|---|---|---|---|
| ESIM + ELMo [44] (5x) | - | - | 89.3 | - | - | - |
| CAFE [58] (5x) | 80.2 | 79.0 | 89.3 | - | - | - |
| Stochastic Answer Network [35] (3x) | 80.6 | 80.1 | - | - | - | - |
| CAFE [58] | 78.7 | 77.9 | 88.5 | 83.3 | | |
| GenSen [64] | 71.4 | 71.3 | - | - | 82.3 | 59.2 |
| Multi-task BiLSTM + Attn [64] | 72.2 | 72.1 | - | - | 82.1 | **61.7** |
| Finetuned Transformer LM (ours) | **82.1** | **81.4** | **89.9** | **88.3** | **88.1** | 56.0 |

# Details

- Augmented objective function in finetuning

- More layers is better!

- Zero-shot

We additionally found that including language modeling as an auxiliary objective to the fine-tuning helped learning by (a) improving generalization of the supervised model, and (b) accelerating convergence. This is in line with prior work [50, 43], who also observed improved performance with such an auxiliary objective. Specifically, we optimize the following objective (with weight $\lambda$):

$$L_3(\mathcal{C}) = L_2(\mathcal{C}) + \lambda * L_1(\mathcal{C}) \tag{5}$$

Overall, the only extra parameters we require during fine-tuning are $W_y$, and embeddings for delimiter tokens (described below in Section 3.3).



Improving Language Understanding by Generative Pre-Training – Radford et al

GPT 2

# GPT-2

- Current paradigm => "narrow learners"
  - Don't generalize well to out-of-distribution data
  - Hypothesis: Single task training
- Idea: Use LM and zero-shot => "general learners"
- + Make your models huge ☺

- P(output|input) → P(output|input, task)
  - (translate to french, english text, french text)
  - (answer the question, document, question, answer)

Language Models are Unsupervised Multitask Learners – Radford et al

# WebText

- Common Crawl: big but low-quality
  - Don't' use
- WebText:
  - Outbound links from Reddit (with karma >= 3)
  - 45 million links
  - 40 GB of text
  - (Removed Wikipedia to avoid conflicts with other datasets)

# Zero-shot Language Modelling

**Language Models are Unsupervised Multitask Learners**

| | LAMBADA (PPL) | LAMBADA (ACC) | CBT-CN (ACC) | CBT-NE (ACC) | WikiText2 (PPL) | PTB (PPL) | enwik8 (BPB) | text8 (BPC) | WikiText103 (PPL) | 1BW (PPL) |
|---|---|---|---|---|---|---|---|---|---|---|
| SOTA | 99.8 | 59.23 | 85.7 | 82.3 | 39.14 | 46.54 | 0.99 | 1.08 | 18.3 | **21.8** |
| 117M | **35.13** | 45.99 | **87.65** | **83.4** | **29.41** | 65.85 | 1.16 | 1.17 | 37.50 | 75.20 |
| 345M | **15.60** | 55.48 | **92.35** | **87.1** | **22.76** | 47.33 | 1.01 | **1.06** | 26.37 | 55.72 |
| 762M | **10.87** | **60.12** | **93.45** | **88.0** | **19.93** | **40.31** | **0.97** | 1.02 | 22.05 | 44.575 |
| 1542M | **8.63** | **63.24** | **93.30** | **89.05** | **18.34** | **35.76** | **0.93** | **0.98** | **17.48** | 42.16 |

*Table 3.* Zero-shot results on many datasets. No training or fine-tuning was performed for any of these results. PTB and WikiText-2 results are from (Gong et al., 2018). CBT results are from (Bajgar et al., 2016). LAMBADA accuracy result is from (Hoang et al., 2018) and LAMBADA perplexity result is from (Grave et al., 2016). Other results are from (Dai et al., 2019).

# Zero-shot Downstream



- Promising and impressive (compared to expectations)
- But far from SOTA

# Example: Natural Questions

- Top 30 most-confident answers
- Question: did these show up in the training data?

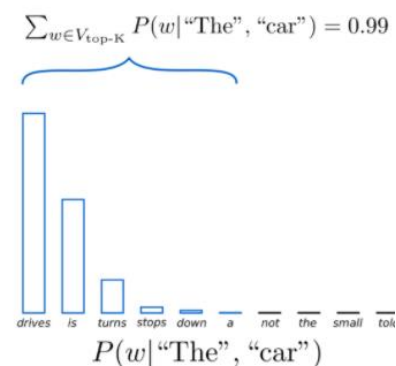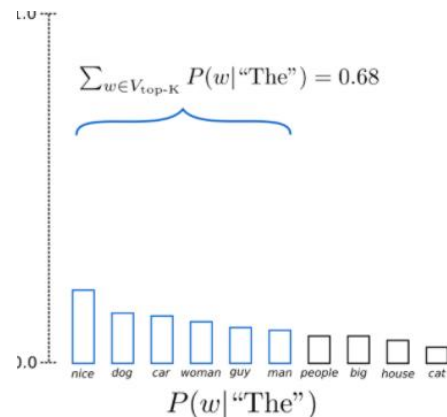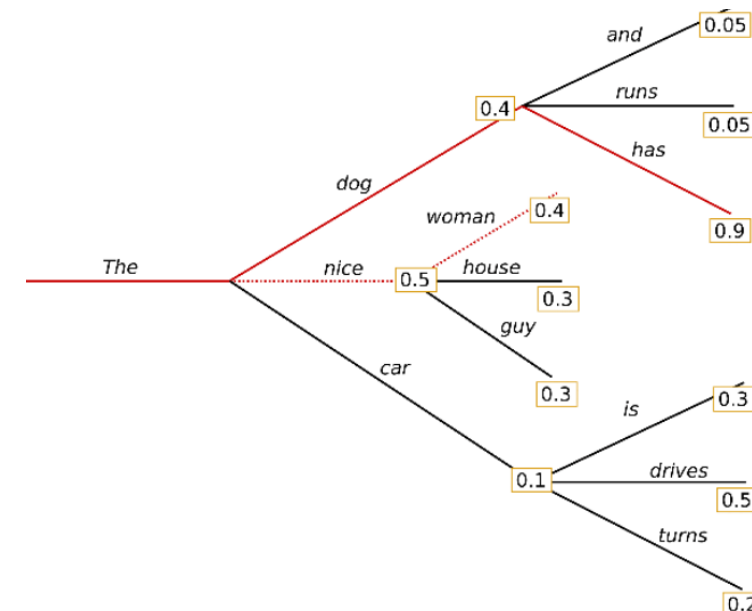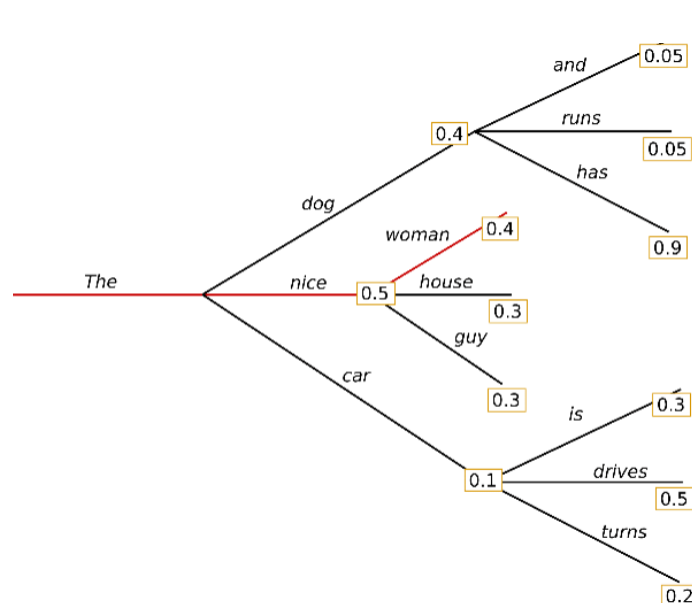| Question | Generated Answer | Correct | Probability |
|---|---|---|---|
| Who wrote the book the origin of species? | Charles Darwin | ✓ | 83.4% |
| Who is the founder of the ubuntu project? | Mark Shuttleworth | ✓ | 82.0% |
| Who is the quarterback for the green bay packers? | Aaron Rodgers | ✓ | 81.1% |
| Panda is a national animal of which country? | China | ✓ | 76.8% |
| Who came up with the theory of relativity? | Albert Einstein | ✓ | 76.4% |
| When was the first star wars film released? | 1977 | ✓ | 71.4% |
| What is the most common blood type in sweden? | A | ✗ | 70.6% |
| Who is regarded as the founder of psychoanalysis? | Sigmund Freud | ✓ | 69.3% |
| Who took the first steps on the moon in 1969? | Neil Armstrong | ✓ | 66.8% |
| Who is the largest supermarket chain in the uk? | Tesco | ✓ | 65.3% |
| What is the meaning of shalom in english? | peace | ✓ | 64.0% |
| Who was the author of the art of war? | Sun Tzu | ✓ | 59.6% |
| Largest state in the us by land mass? | California | ✗ | 59.2% |
| Green algae is an example of which type of reproduction? | parthenogenesis | ✗ | 56.5% |
| Vikram samvat calender is official in which country? | India | ✓ | 55.6% |
| Who is mostly responsible for writing the declaration of independence? | Thomas Jefferson | ✓ | 53.3% |
| What us state forms the western boundary of montana? | Montana | ✗ | 52.3% |
| Who plays ser davos in game of thrones? | Peter Dinklage | ✗ | 52.1% |
| Who appoints the chair of the federal reserve system? | Janet Yellen | ✗ | 51.5% |
| State the process that divides one nucleus into two genetically identical nuclei? | mitosis | ✓ | 50.7% |
| Who won the most mvp awards in the nba? | Michael Jordan | ✗ | 50.2% |
| What river is associated with the city of rome? | the Tiber | ✓ | 48.6% |
| Who is the first president to be impeached? | Andrew Johnson | ✓ | 48.3% |
| Who is the head of the department of homeland security 2017? | John Kelly | ✓ | 47.0% |
| What is the name given to the common currency to the european union? | Euro | ✓ | 46.8% |
| What was the emperor name in star wars? | Palpatine | ✓ | 46.5% |
| Do you have to have a gun permit to shoot at a range? | No | ✓ | 46.4% |
| Who proposed evolution in 1859 as the basis of biological development? | Charles Darwin | ✓ | 45.7% |
| Nuclear power plant that blew up in russia? | Chernobyl | ✓ | 45.7% |
| Who played john connor in the original terminator? | Arnold Schwarzenegger | ✗ | 45.2% |

# Generalization vs Memorization

- Bloom filters with 8-grams => estimate overlap
  - Given Datasets A, B.
  - Question: What is the percentage of 8-grams from A that are also in B?
- Interesting: 1BW has overlap of ~13% with its own training set…
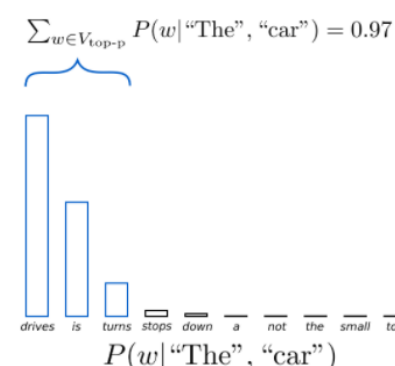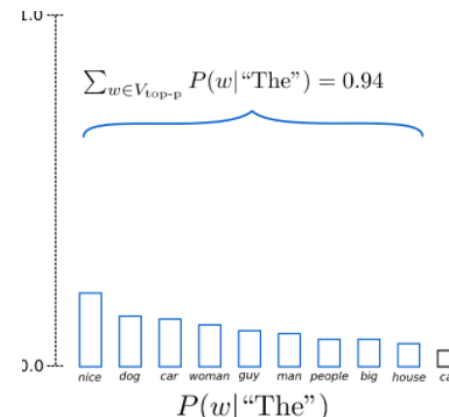- TL;DR – WebText has low or no overlap with the datasets used in the studies

# Text generation from LMs

- Greedy
- Beam search
- Sampling
  - Top-k sampling
  - Top-p sampling



Top-k

Top-p

**Context (human-written):** In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

**GPT-2:** The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

Pérez and the others then ventured further into the valley. "By the time we reached the top of one peak, the water looked blue, with some crystals on top," said Pérez.

Pérez and his friends were astonished to see the unicorn herd. These creatures could be seen from the air without having to move too much to see them – they were so close they could touch their horns.

While examining these bizarre creatures the scientists discovered that the creatures also spoke some fairly regular English. Pérez stated, "We can see, for example, that they have a common 'language,' something like a dialect or dialectic."

Dr. Pérez believes that the unicorns may have originated in Argentina, where the animals were believed to be descendants of a lost race of people who lived there before the arrival of humans in those parts of South America.

While their origins are still unclear, some believe that perhaps the creatures were created when a human and a unicorn met each other in a time before human civilization. According to Pérez, "In South America, such incidents seem to be quite common."

However, Pérez also pointed out that it is likely that the only way of knowing for sure if unicorns are indeed the descendants of a lost alien race is through DNA. "But they seem to be able to communicate in English quite well, which I believe is a sign of evolution, or at least a change in social organization," said the scientist.