

CARTE ML Workshop

GPT and Attention

Challenges in NLP

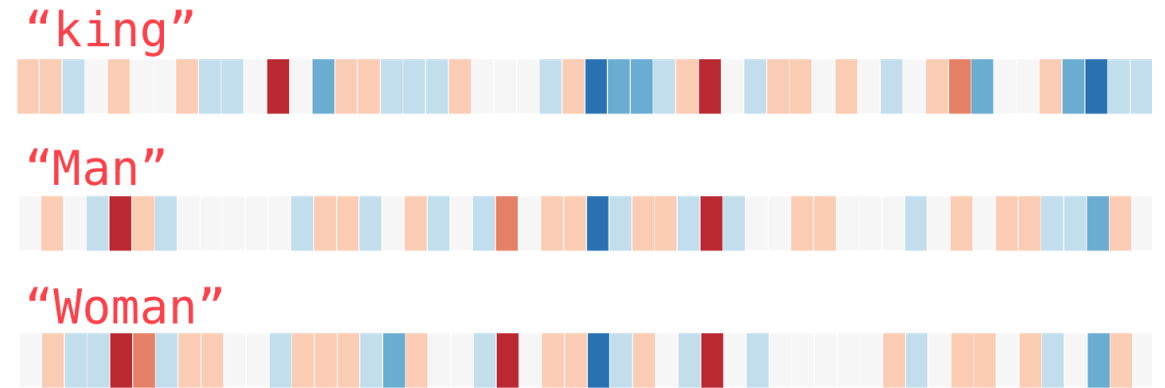
- **Ambiguity:** Words can mean different things depending on context
- **Nuances:** Languages are full of idioms, slang, cultural references, sarcasm...
- **Syntax vs Semantics:** A grammatically correct sentence might not make sense, or a grammatically incorrect one might be easy to understand
- I saw a man on the hill with the telescope
- That's a cool cat
- Colourless green ideas sleep furiously
- Me went store

Foundations of LLMs

- Fundamental goal of language modelling: next word prediction
- $P(cat \mid the \ dog \ and \ the)$
- To generate, pick the word with highest likelihood
- Early models could handle one, two words of context
- Locally coherent, but longer texts quickly lose meaning
- More context requires more complexity!

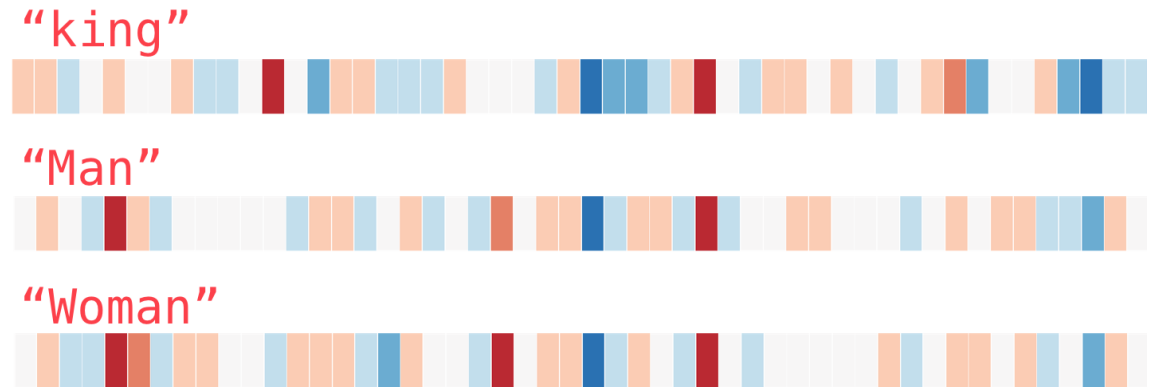
How does an LM “understand” word meaning?

- In order to predict the likelihood of a word, we must have some sense of its meaning
- Some words have similar meanings, and can easily fit in the same place
- In the same way CNNs convert an image into a set of feature maps, we can convert a word into a set of abstract linguistic features
- Word2Vec: 300 features
- GPT-3: 12,888

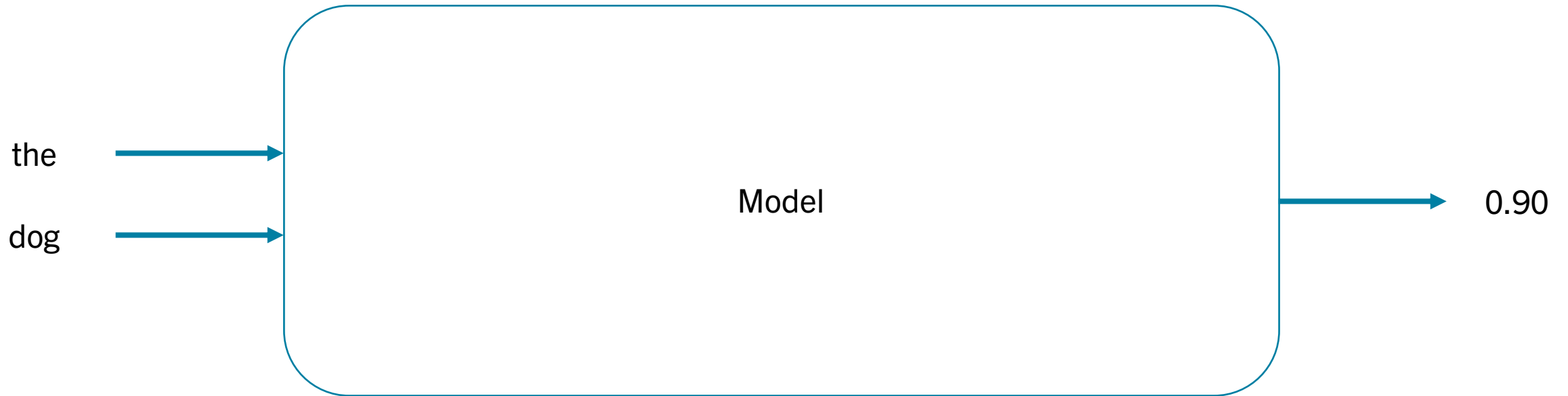


How does an LM “understand” word meaning?

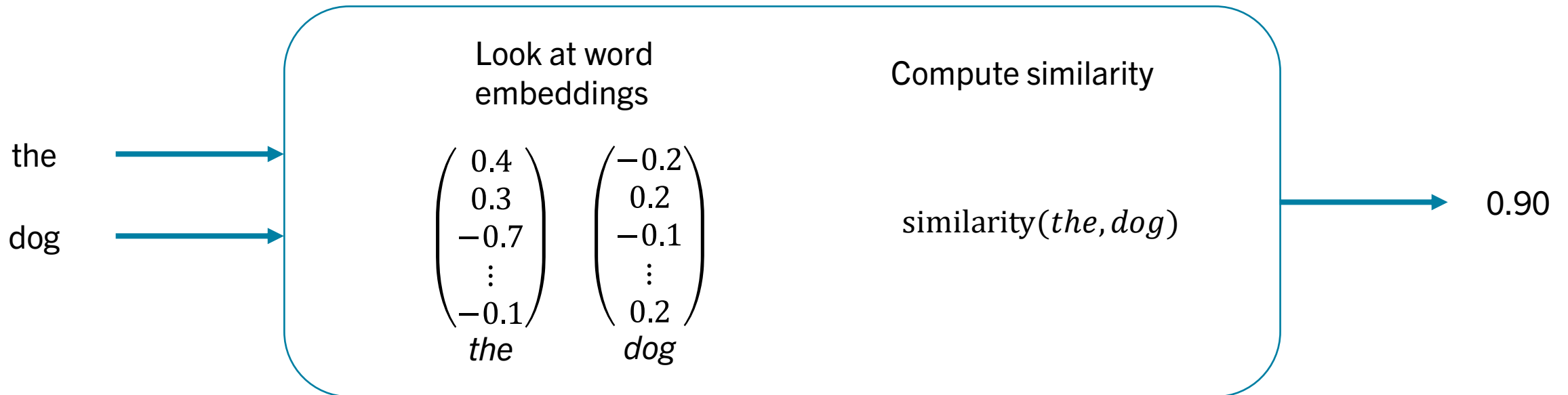
- Key concept of word embeddings: similar words should have similar vectors
- How do we accomplish this?



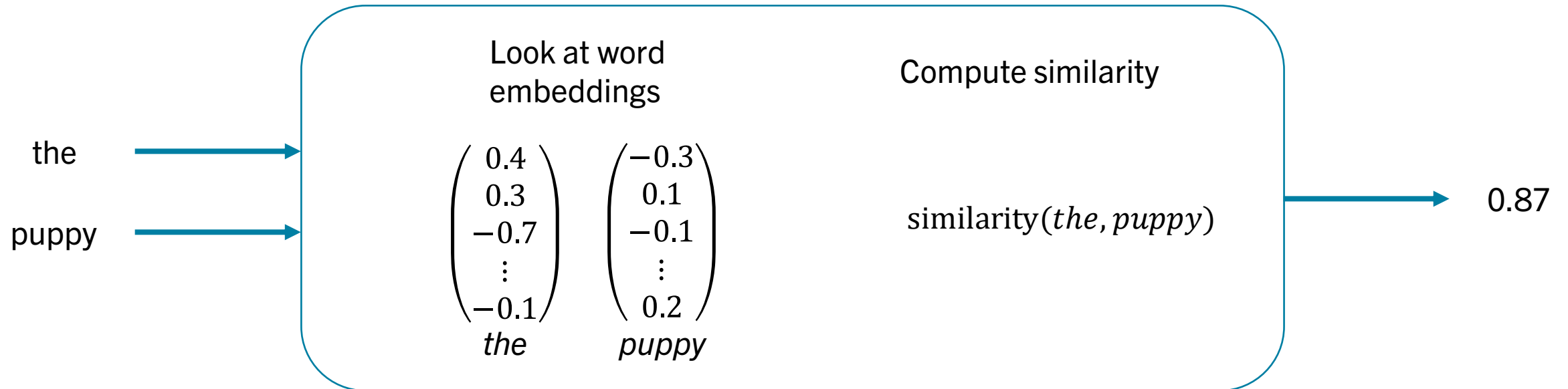
Building Word Embeddings



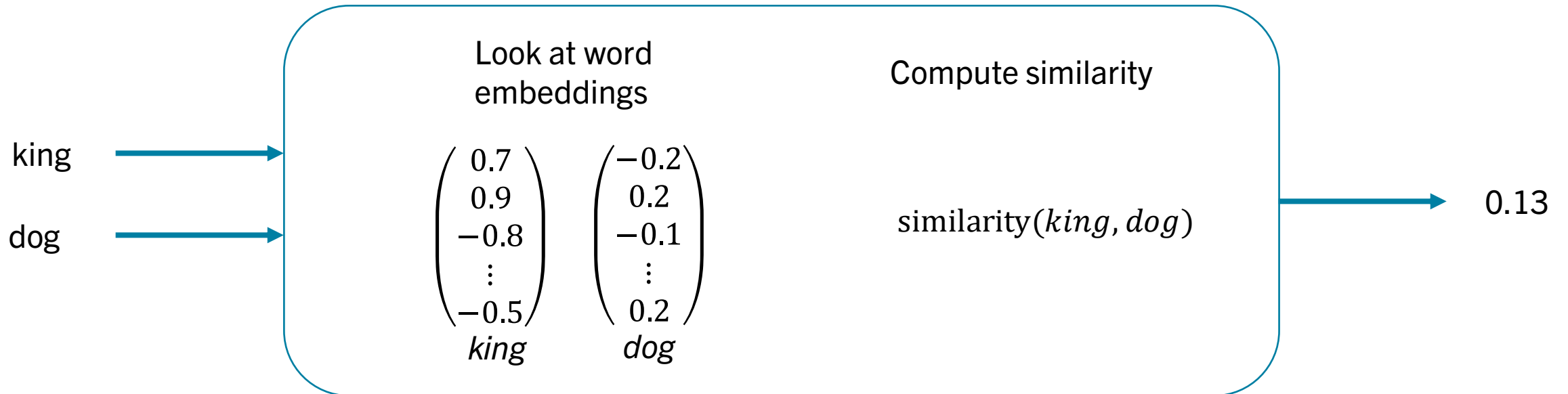
Building Word Embeddings



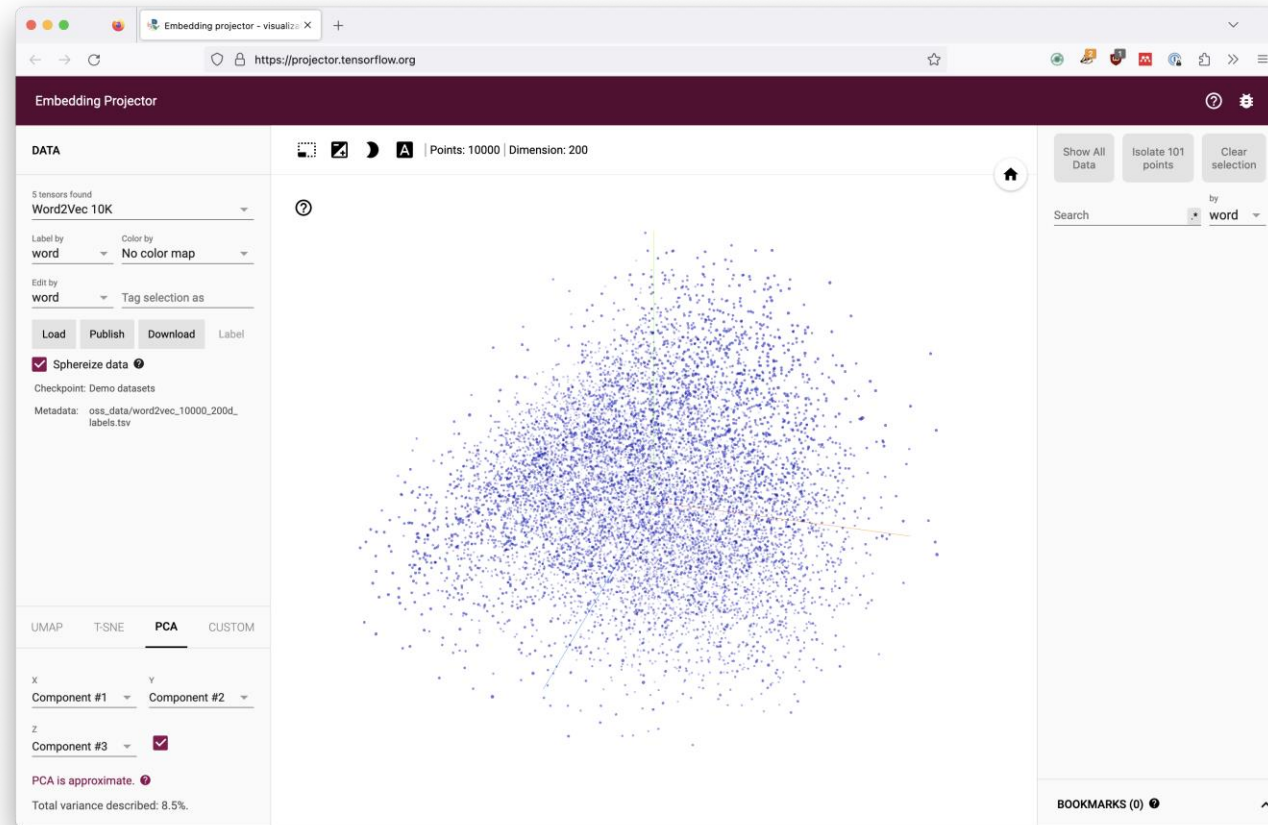
Building Word Embeddings



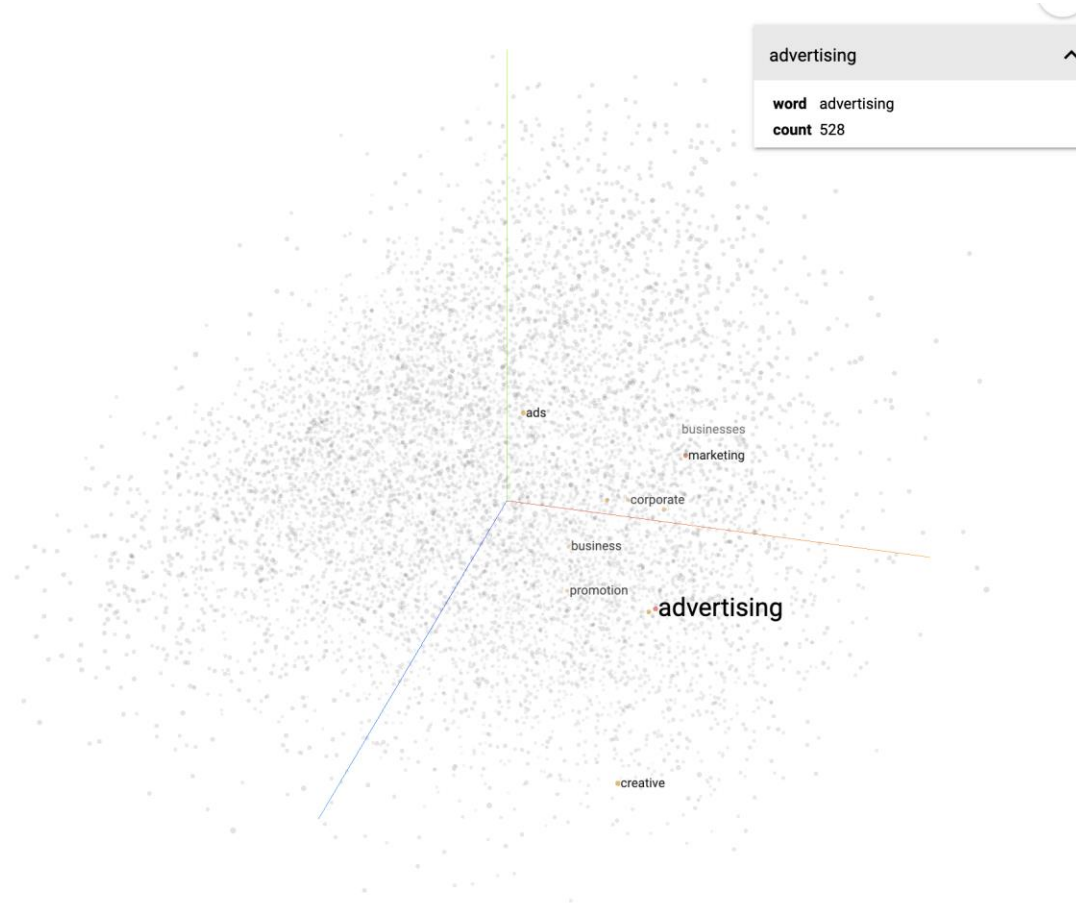
Building Word Embeddings



Results



Results



advertising ^

word	advertising
count	528

by
Search * word

neighbors ? 10

distance COSINE EUCLIDEAN

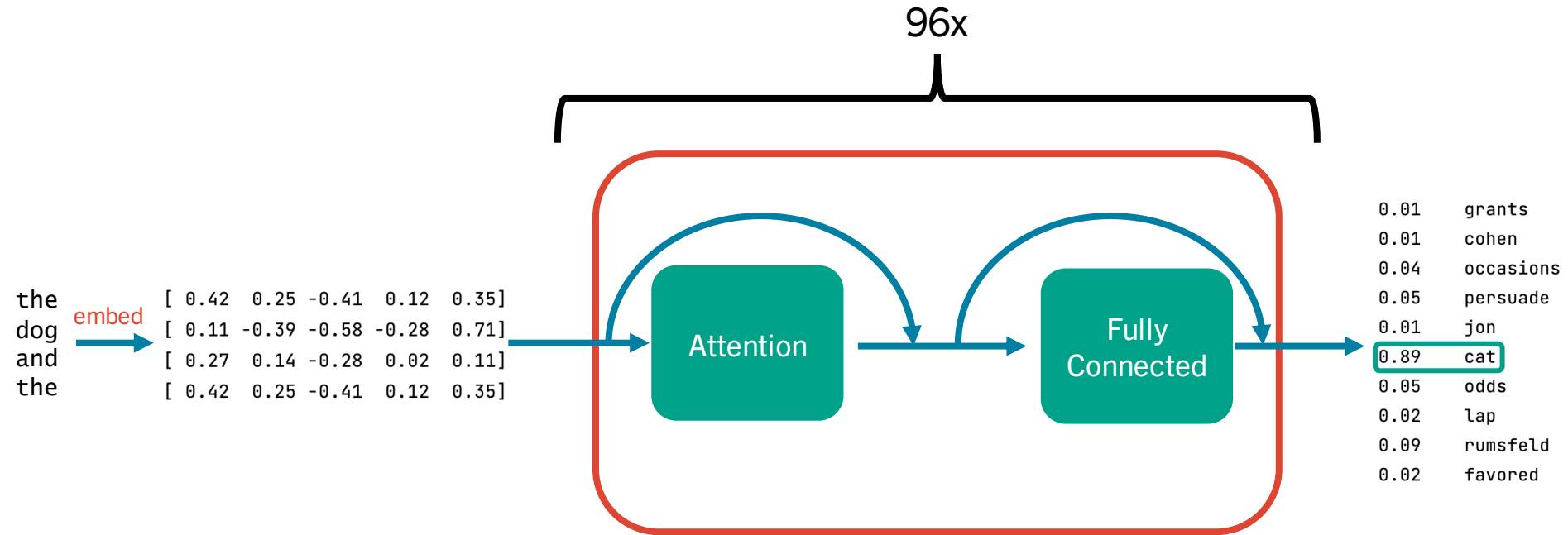
Nearest points in the original space:

marketing	0.384
corporate	0.549
media	0.572
sales	0.574
promotion	0.585
ads	0.585
business	0.604
creative	0.626
commercial	0.641
businesses	0.660

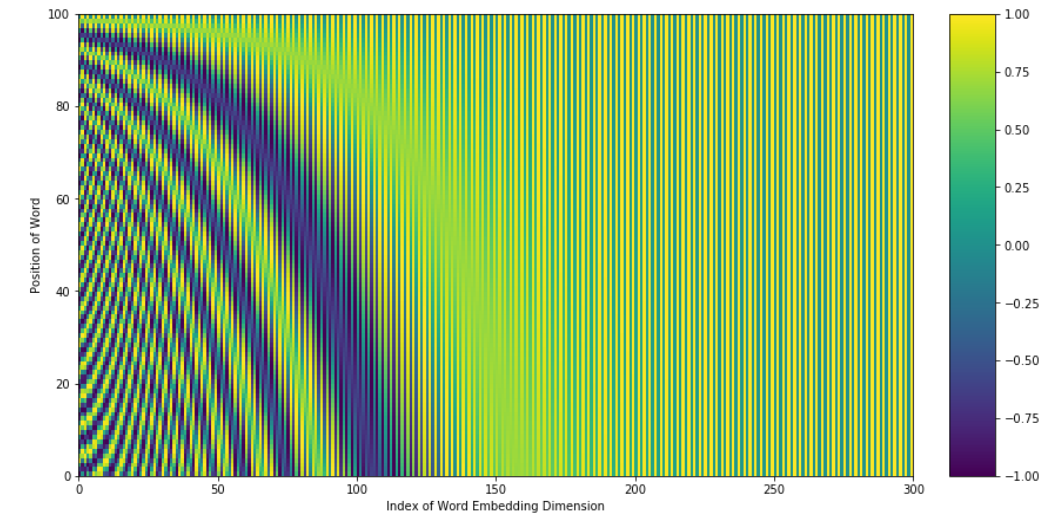
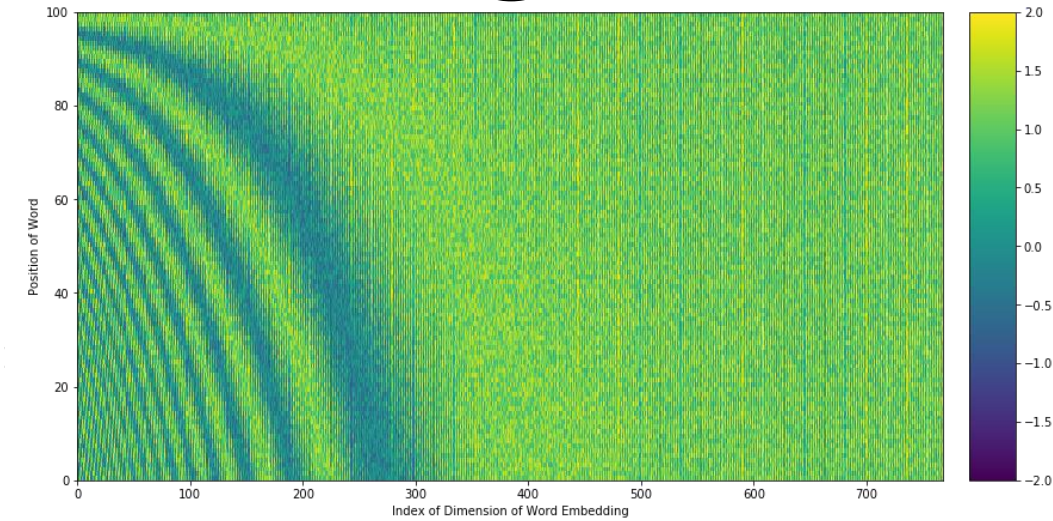
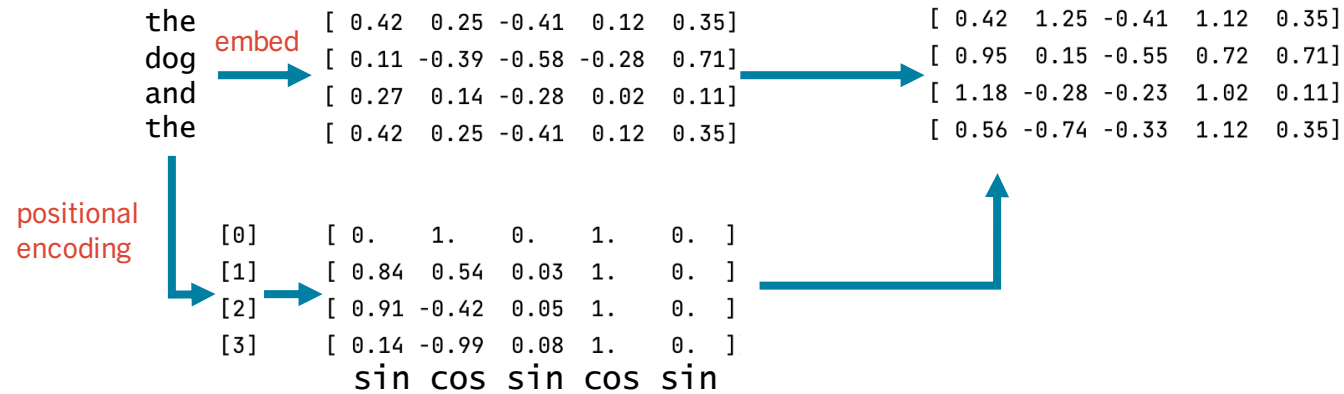
Building GPT



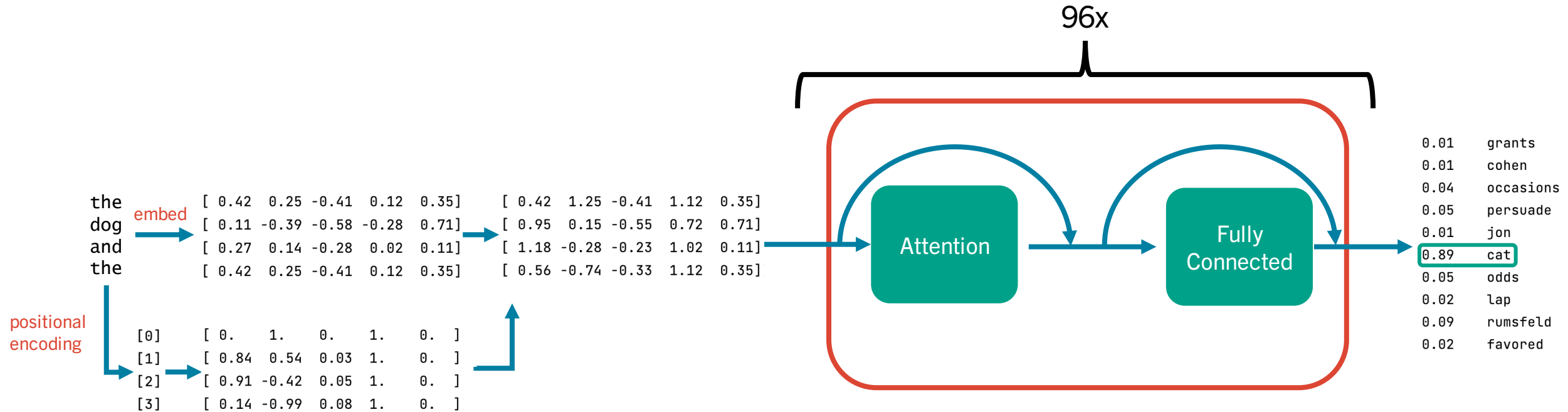
Building GPT: The Transformer



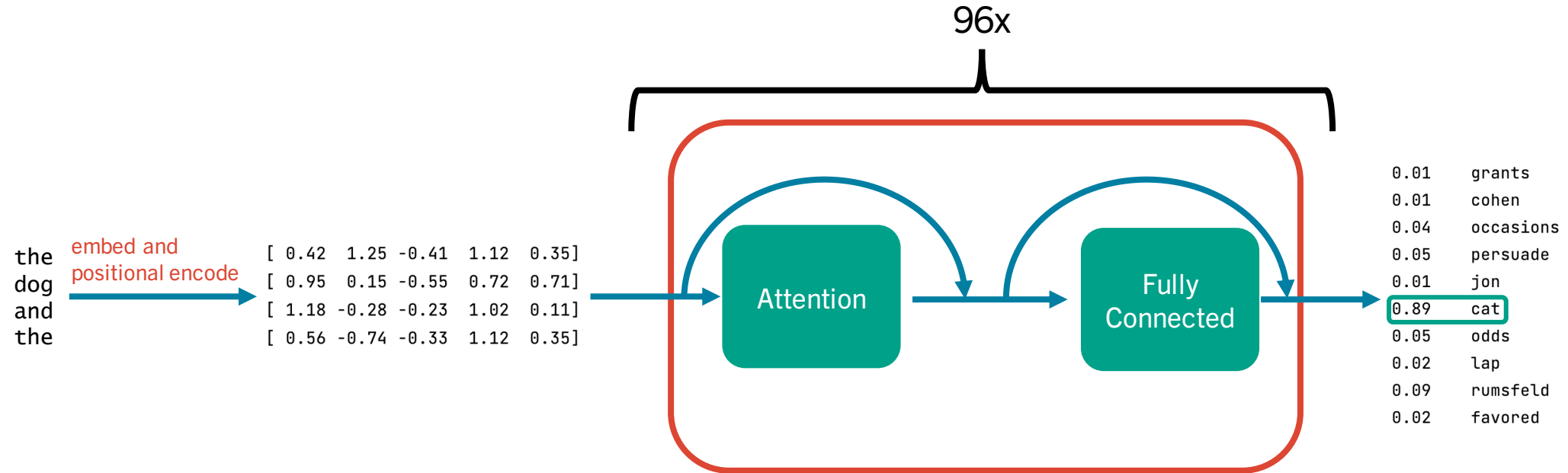
Building GPT: Positional Embedding



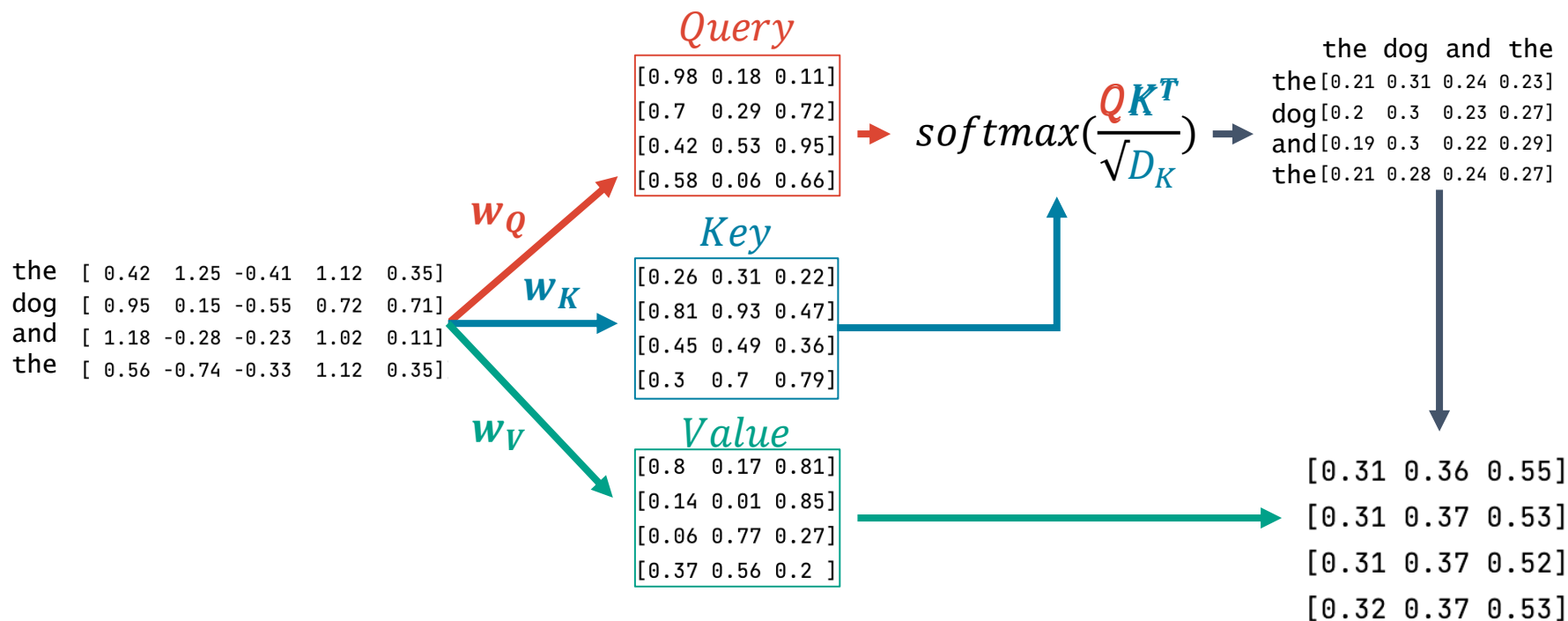
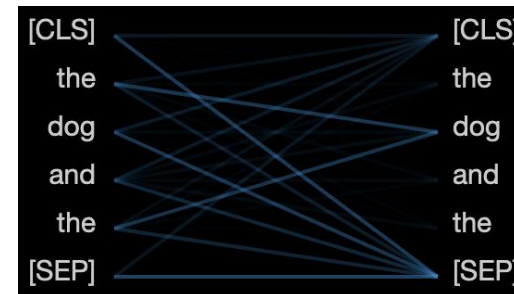
Building GPT



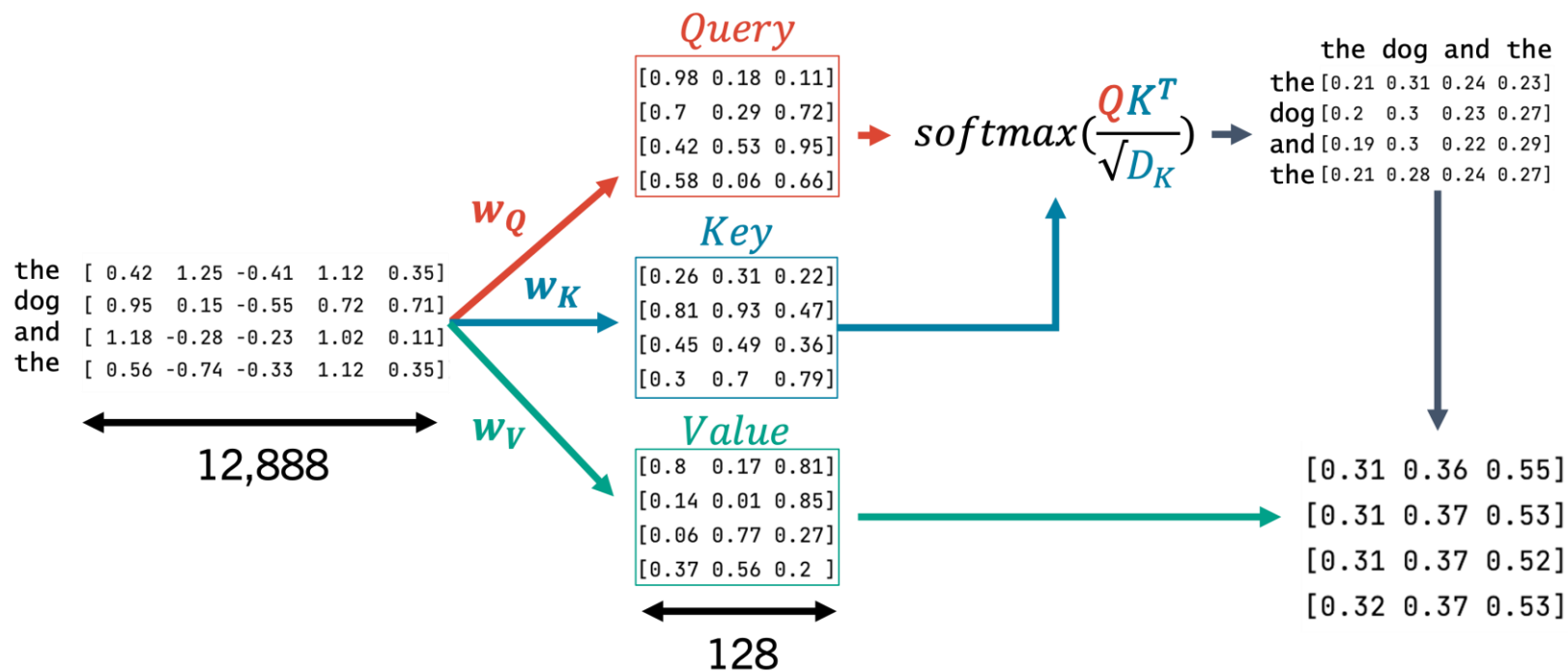
Building GPT



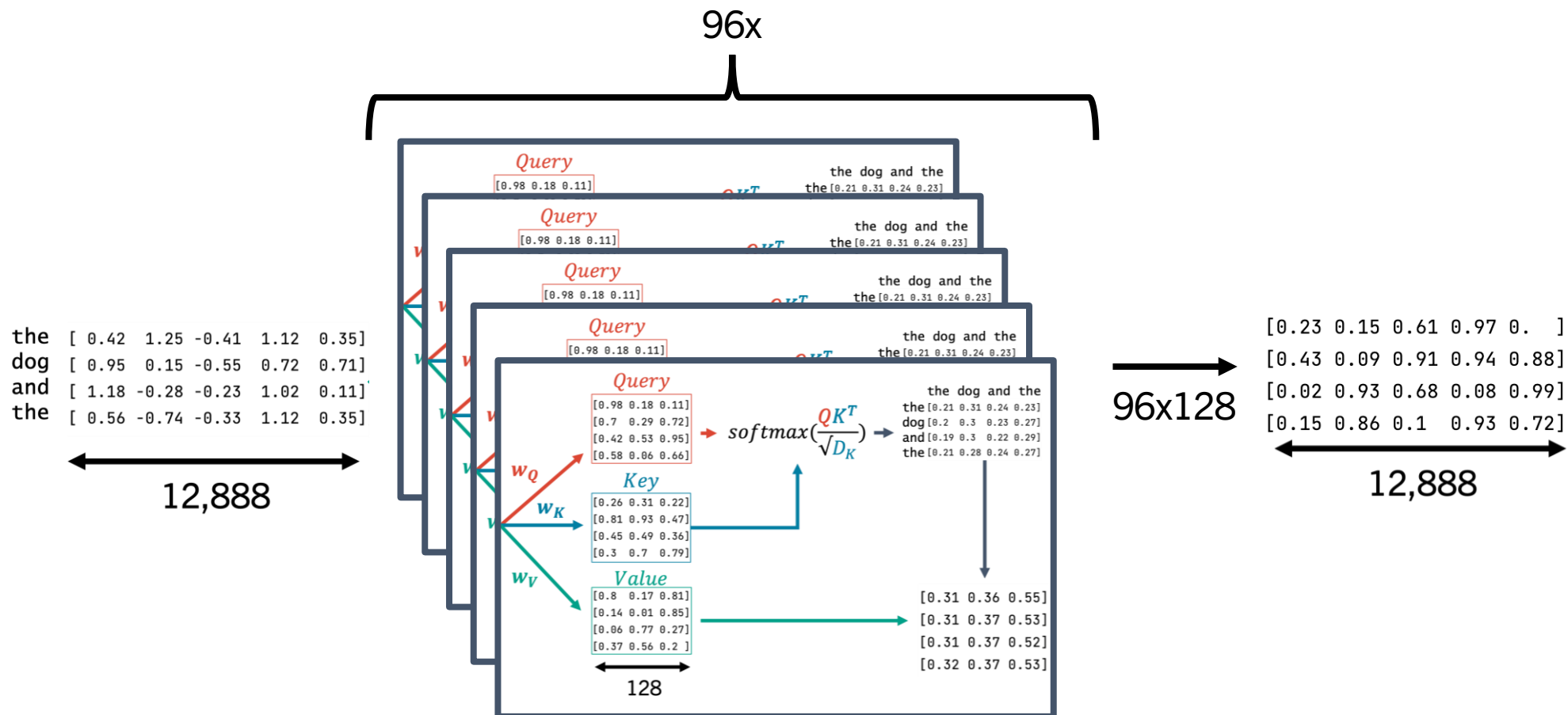
Building GPT: Attention



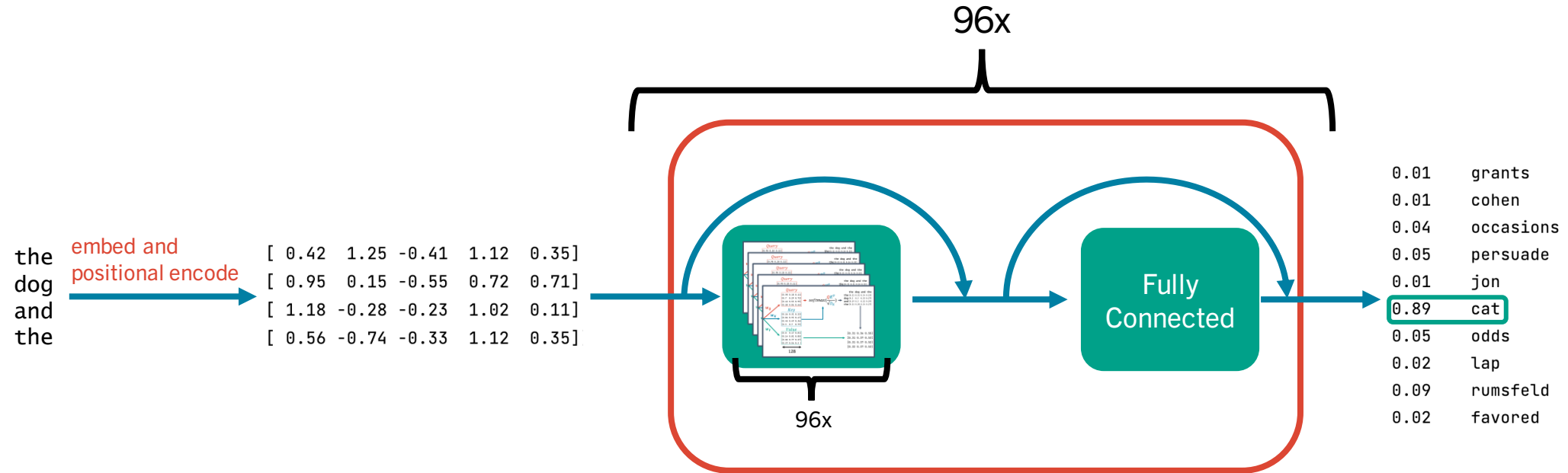
Building GPT: Attention



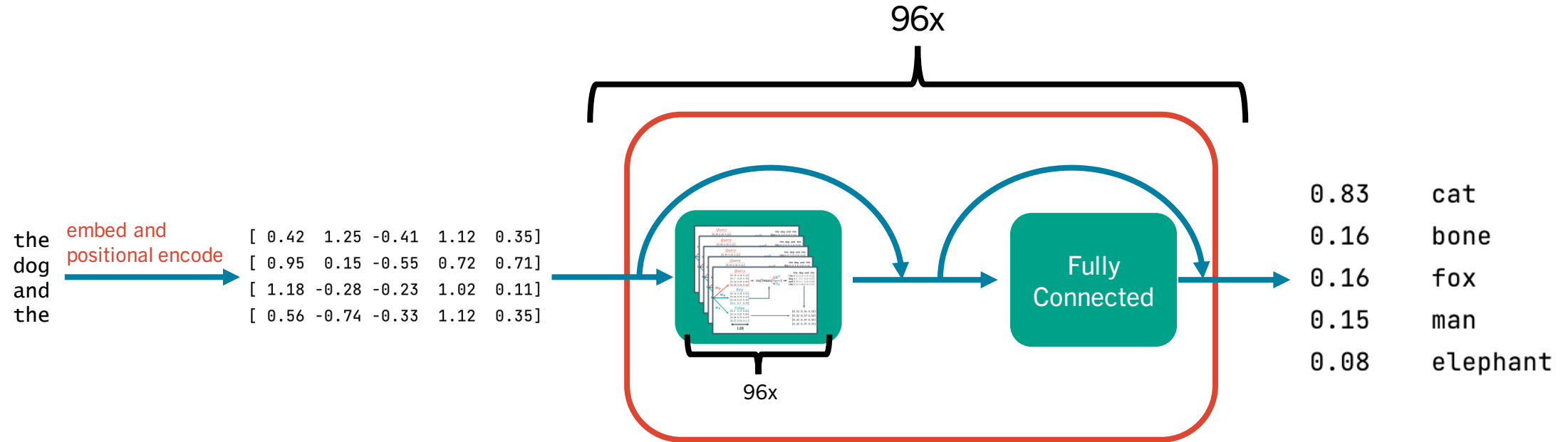
Building GPT: Attention



Building GPT



Building GPT: Top-P



Building GPT: Top-P

Top 10 documentaries about artificial intelligence:

1. AlphaGo (2017)

2017 = 96.15%

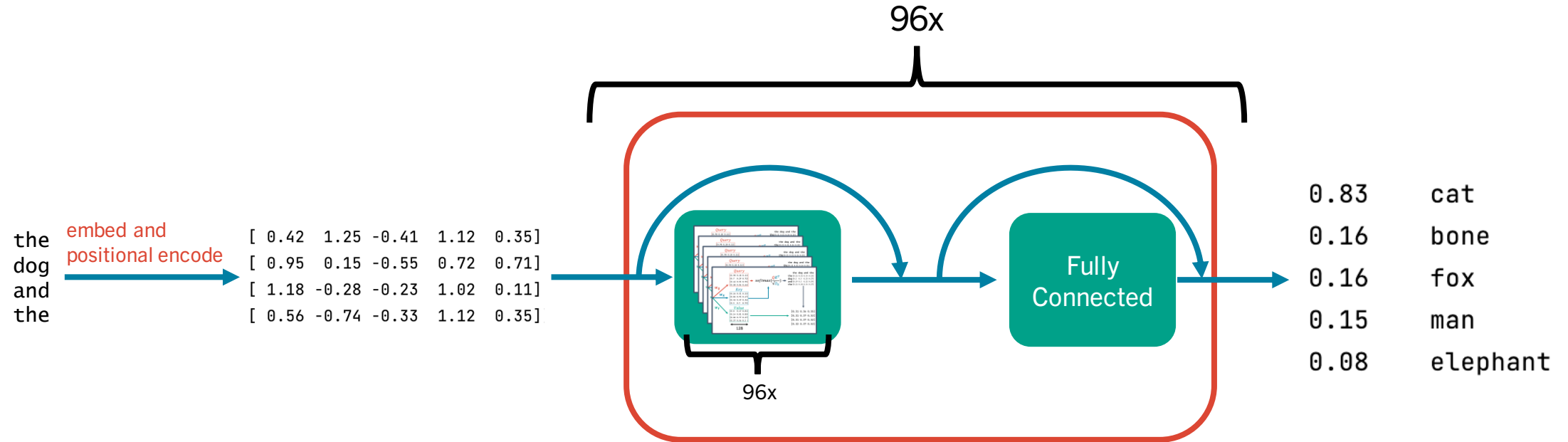
2016 = 2.79%

2018 = 0.88%

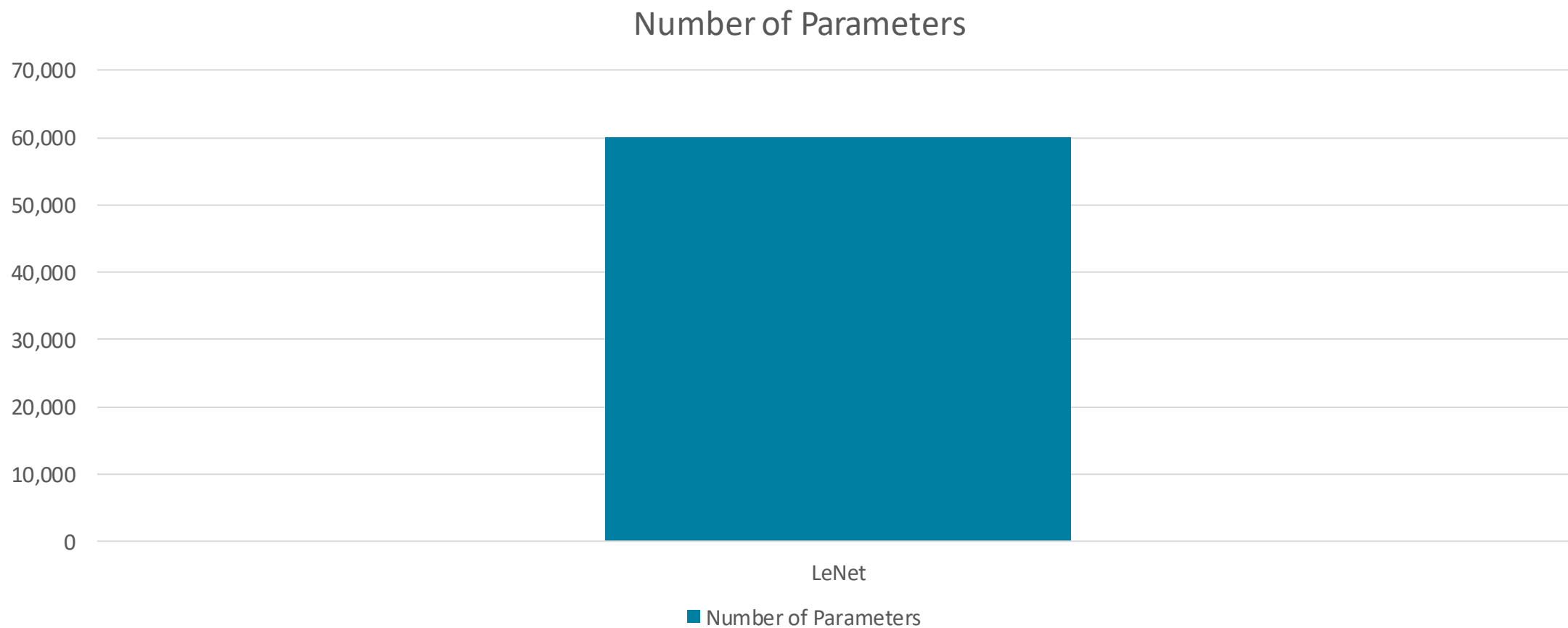
2015 = 0.07%

2019 = 0.03%

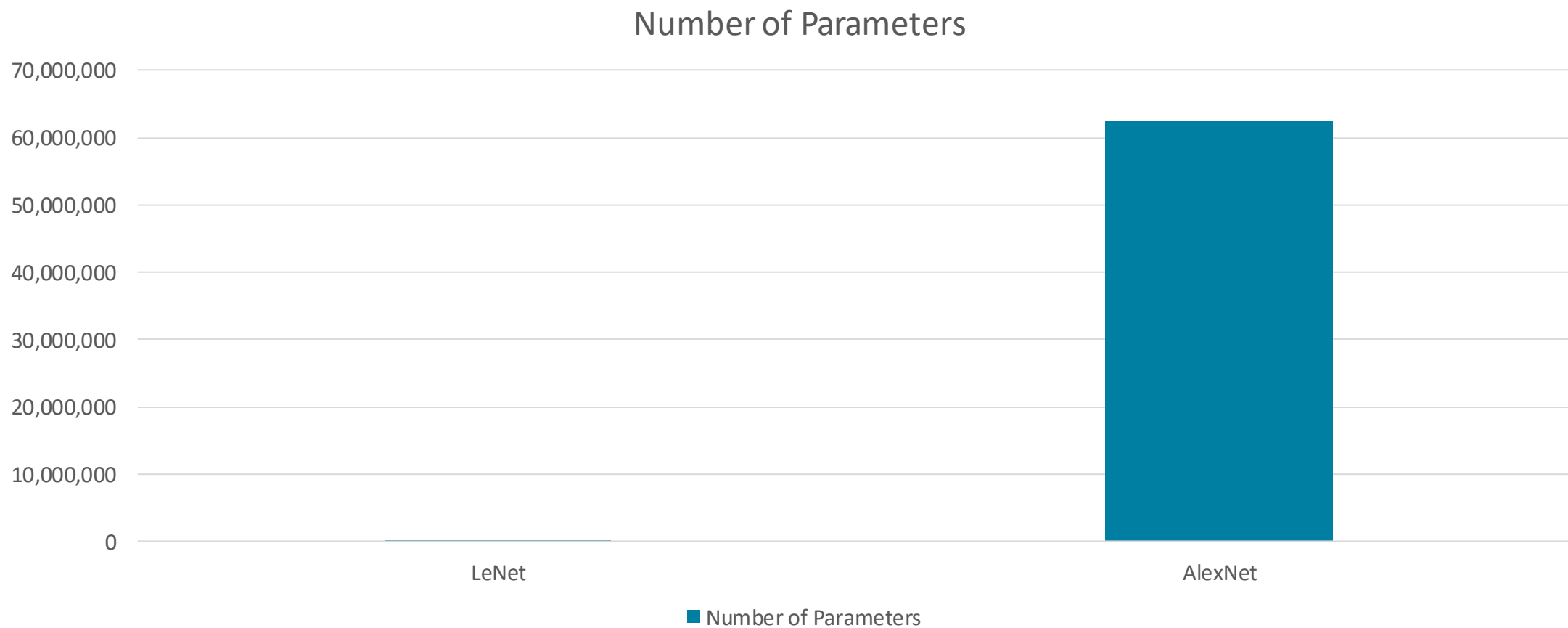
Building GPT



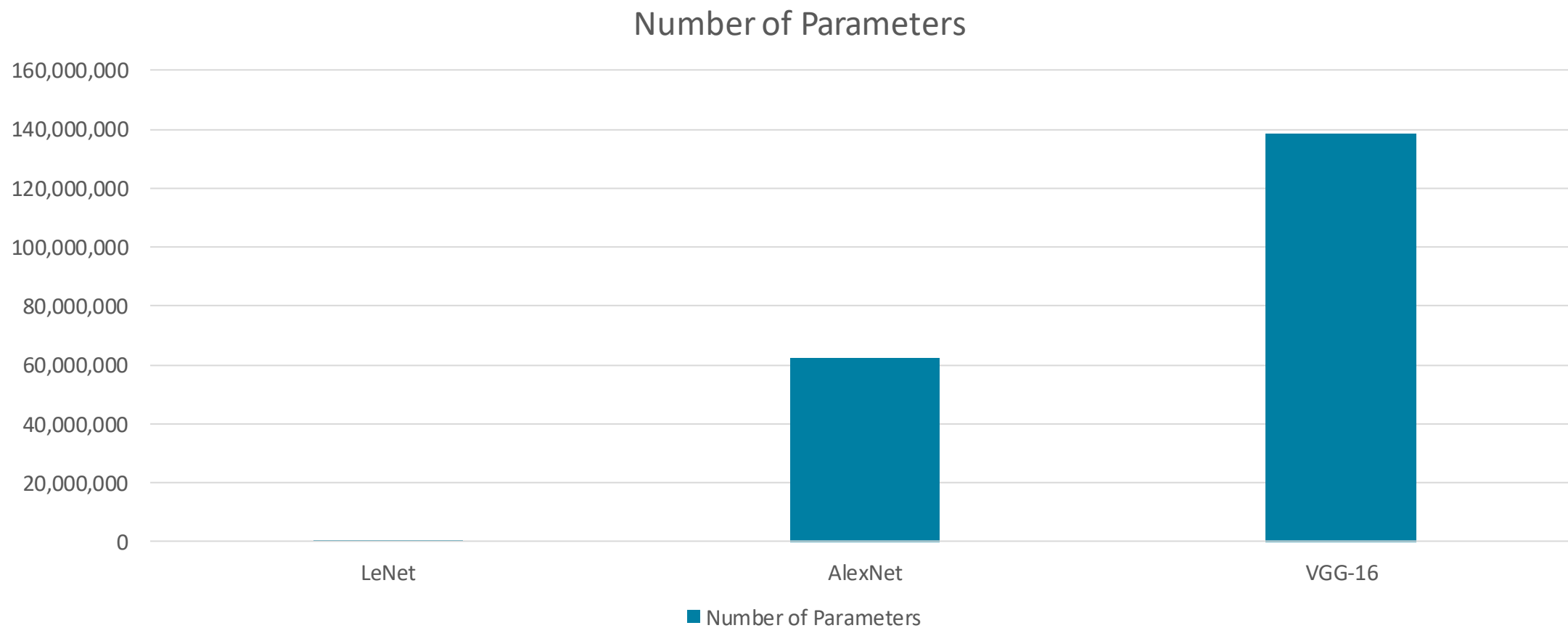
Scale of GPT



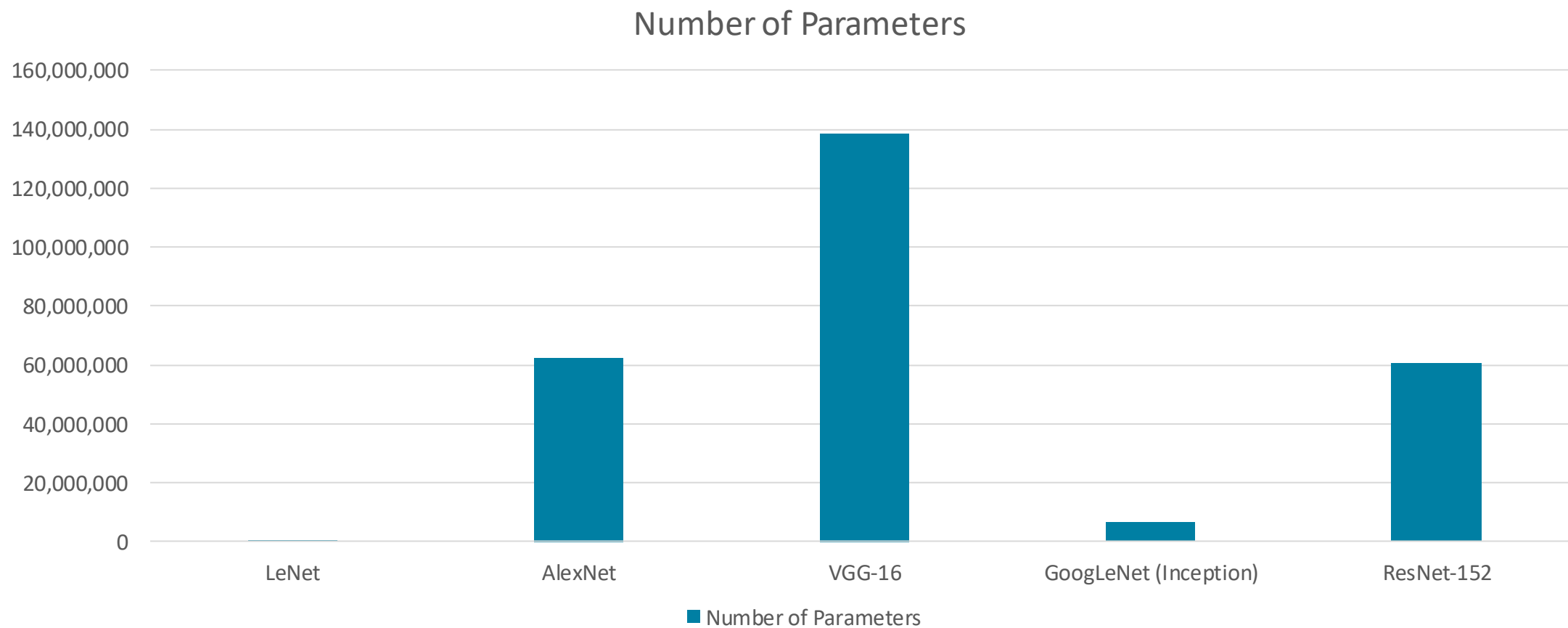
Scale of GPT



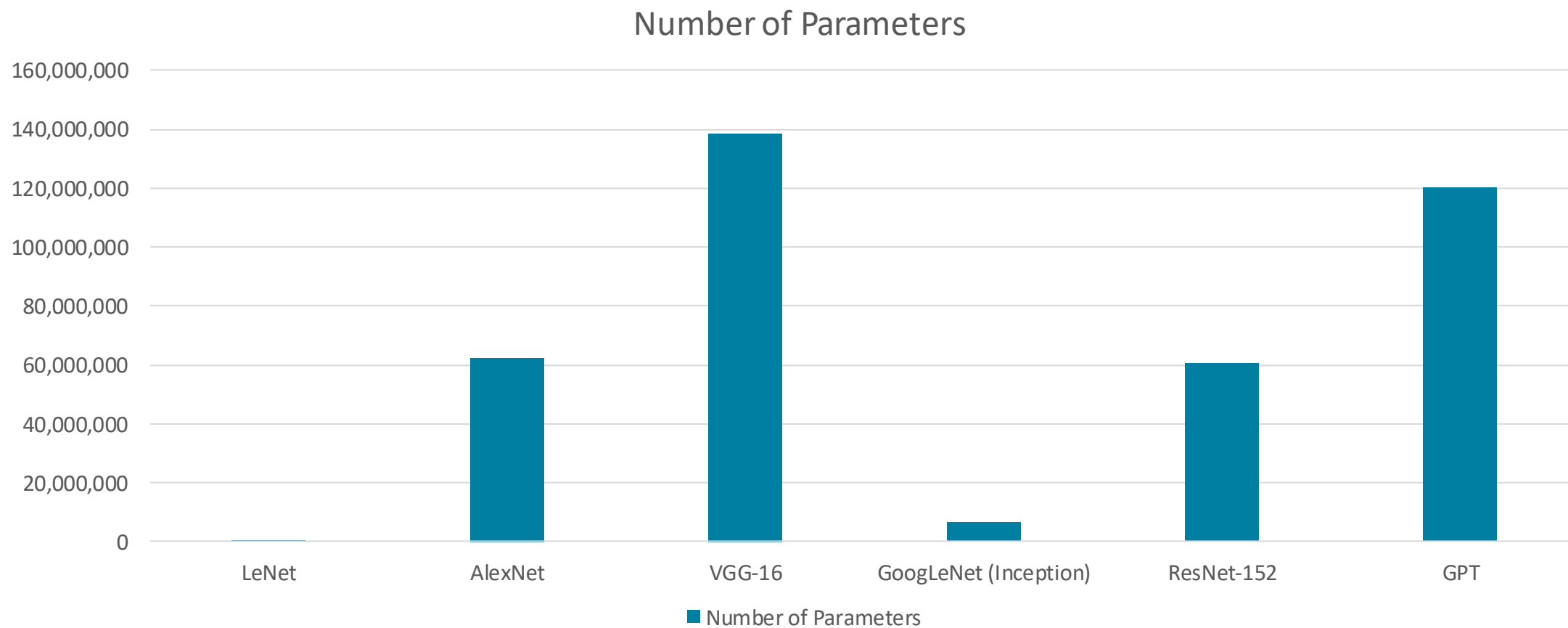
Scale of GPT



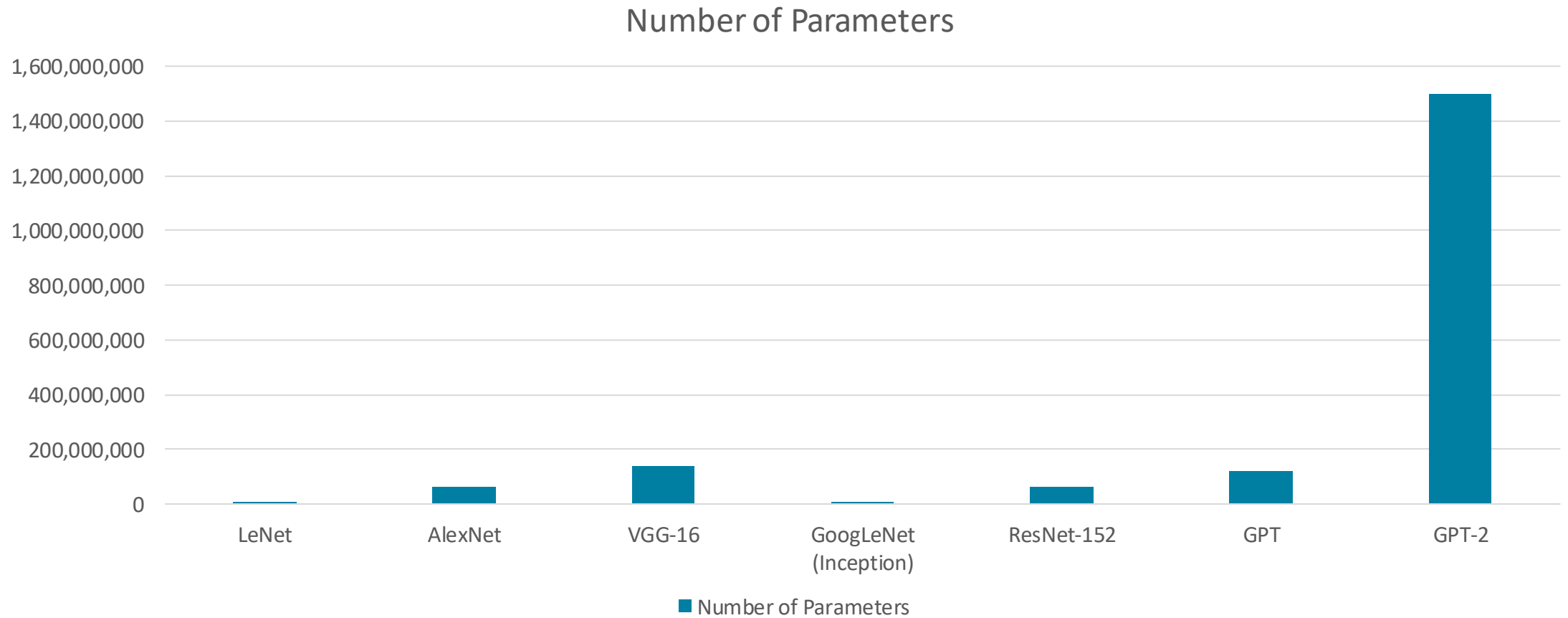
Scale of GPT



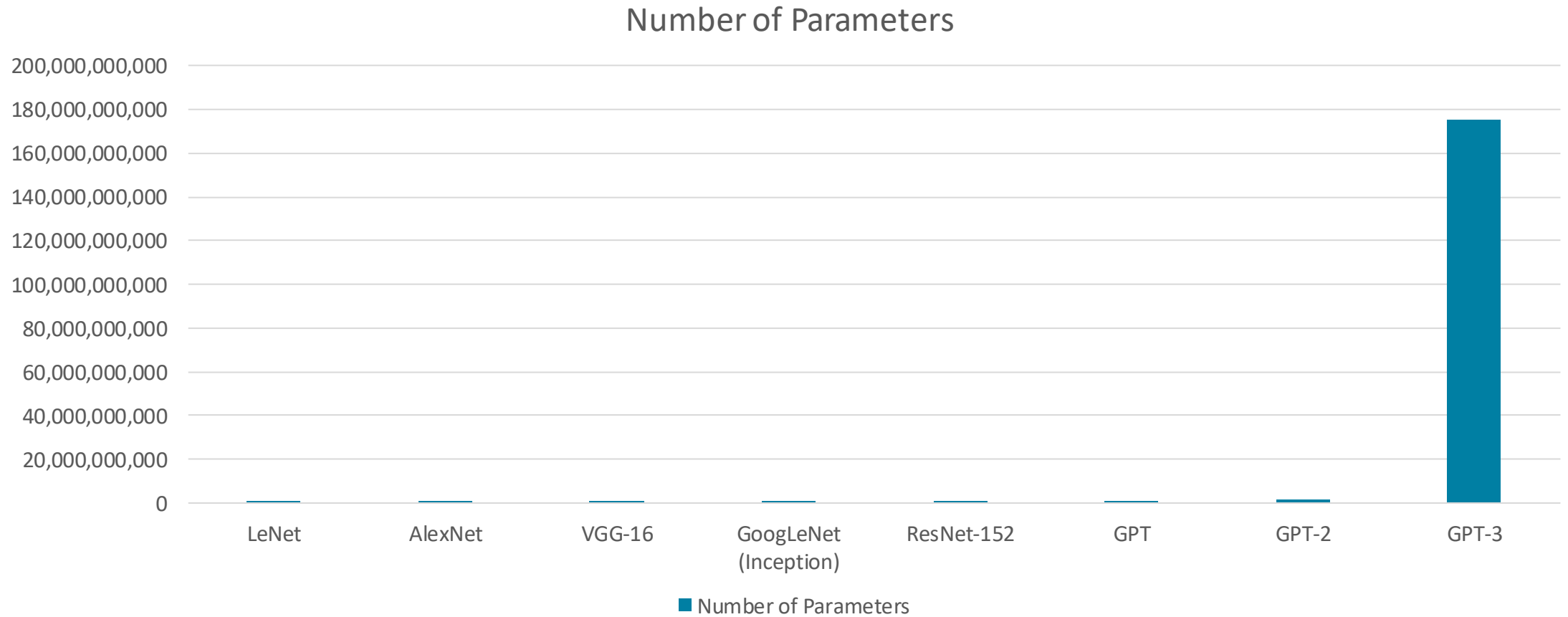
Scale of GPT



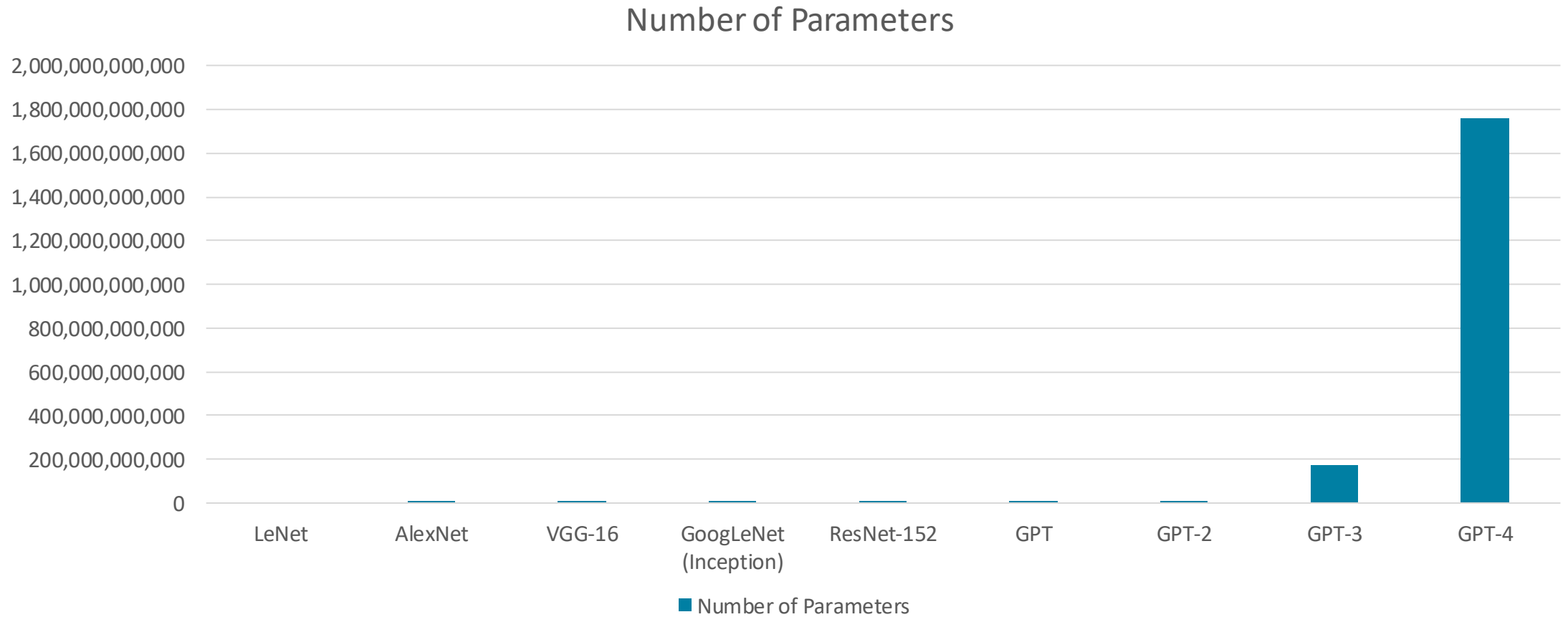
Scale of GPT



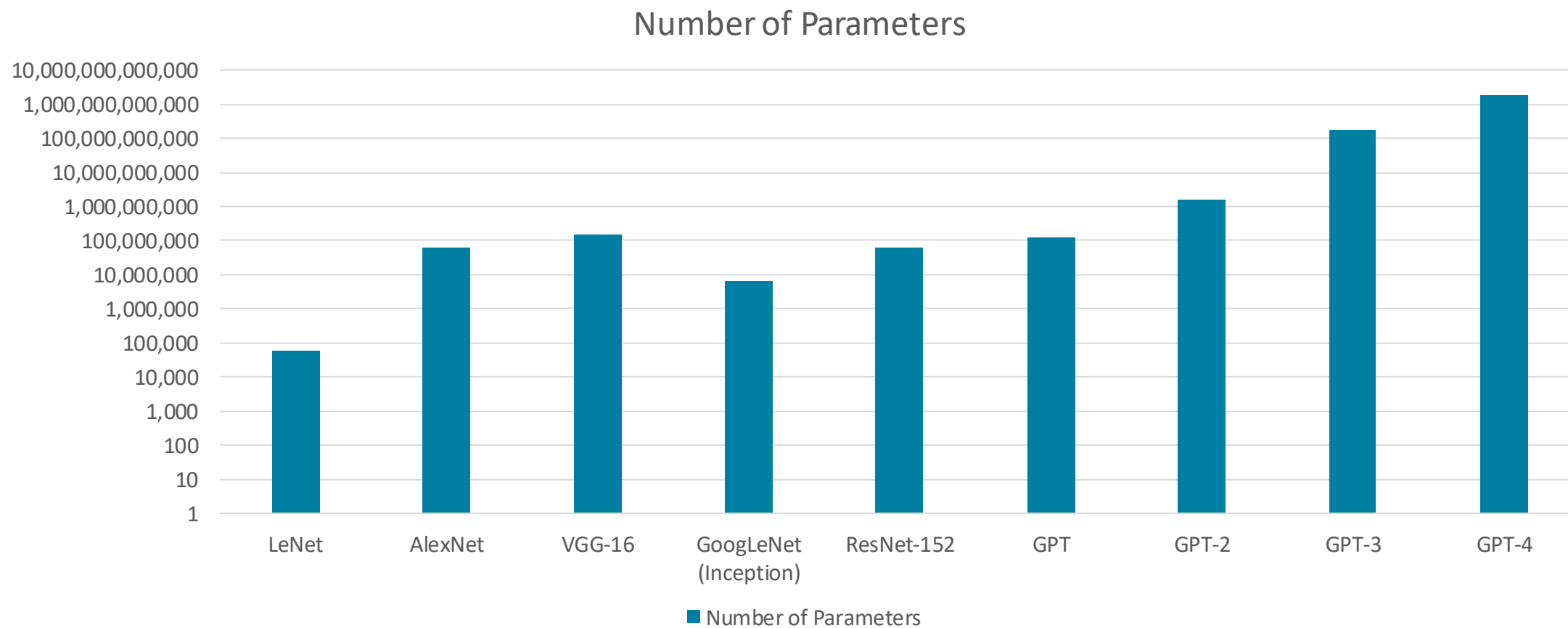
Scale of GPT



Scale of GPT



Scale of GPT



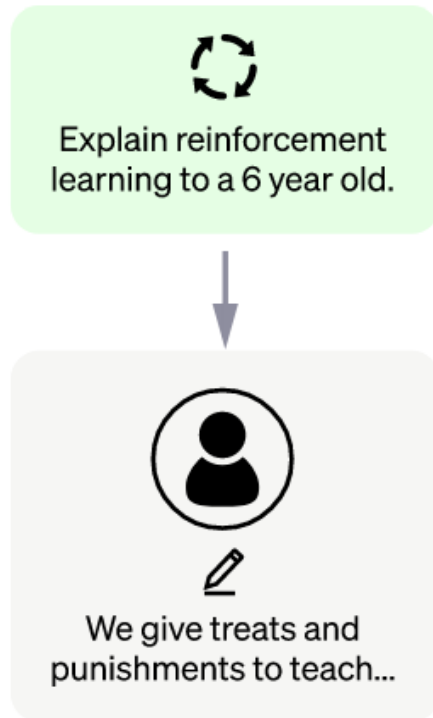
GPT's Training Data

- 1 token $\approx \frac{3}{4}$ word
- Some datasets are sampled more times than others
- Common Crawl: billions of webpages collected over 7 years
- Webtext2: Dataset of webpages that have been shared on Reddit
- Books1: Free ebooks (?)
- Books2: Secret!
- English Wikipedia

Dataset	Quantity (tokens)	Weight in training mix
---------	----------------------	---------------------------

The training innovation of ChatGPT

Human annotators write answers to questions



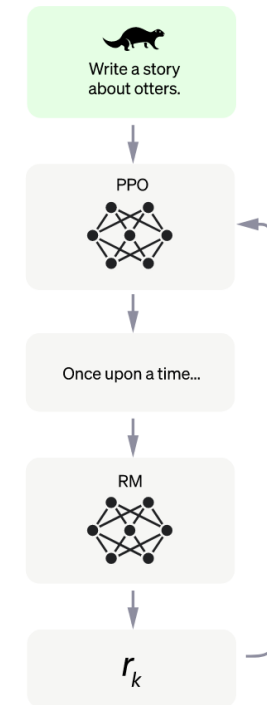
The generalist GPT model is taught from these Q&A pairs

Human annotators write more answers, and someone else ranks them



A separate model learns to rate the quality of an answer

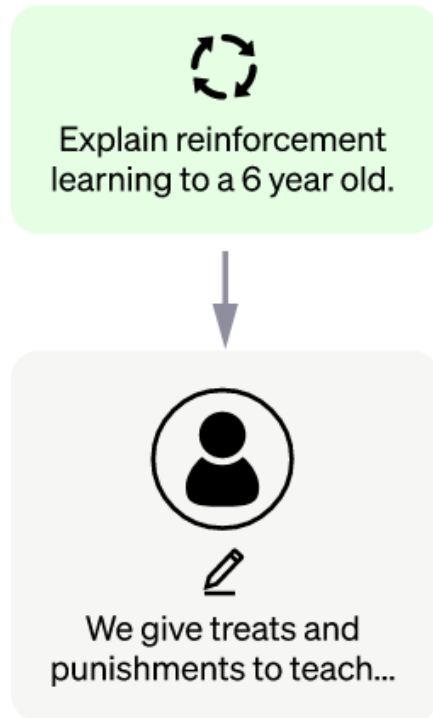
GPT writes answers to sampled questions



The reward model rates each answer, allowing GPT to keep learning

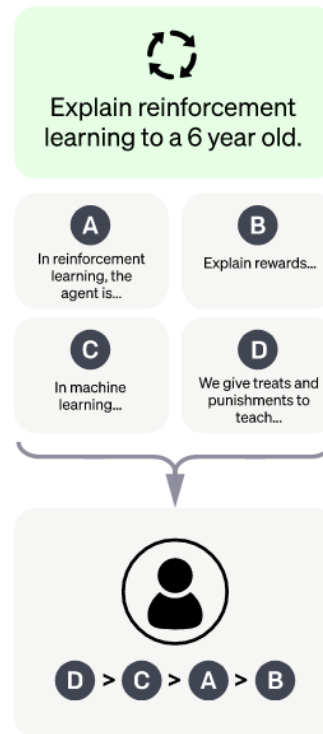
The training innovation of ChatGPT

Human annotators write answers to questions



The generalist GPT model is taught from these Q&A pairs

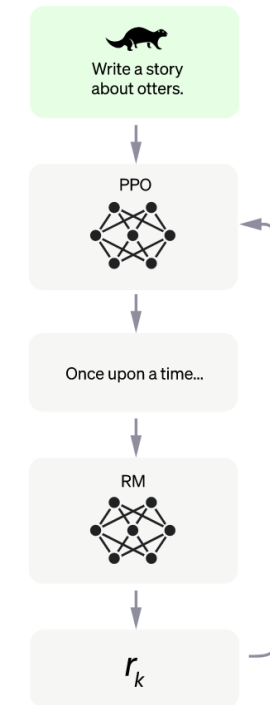
Human annotators write more answers, and someone else ranks them



A separate model learns to rate the quality of an answer

No more humans involved!

GPT writes answers to sampled questions



The reward model rates each answer, allowing GPT to keep learning