# *The 2 Cutting Edges of GenAI:*
## *Coding Liberator or Job Terminator?*

Michael Wu, PhD (@mich8elwu)
chief AI strategist @ PROS

2024.04.19

**PROS.**
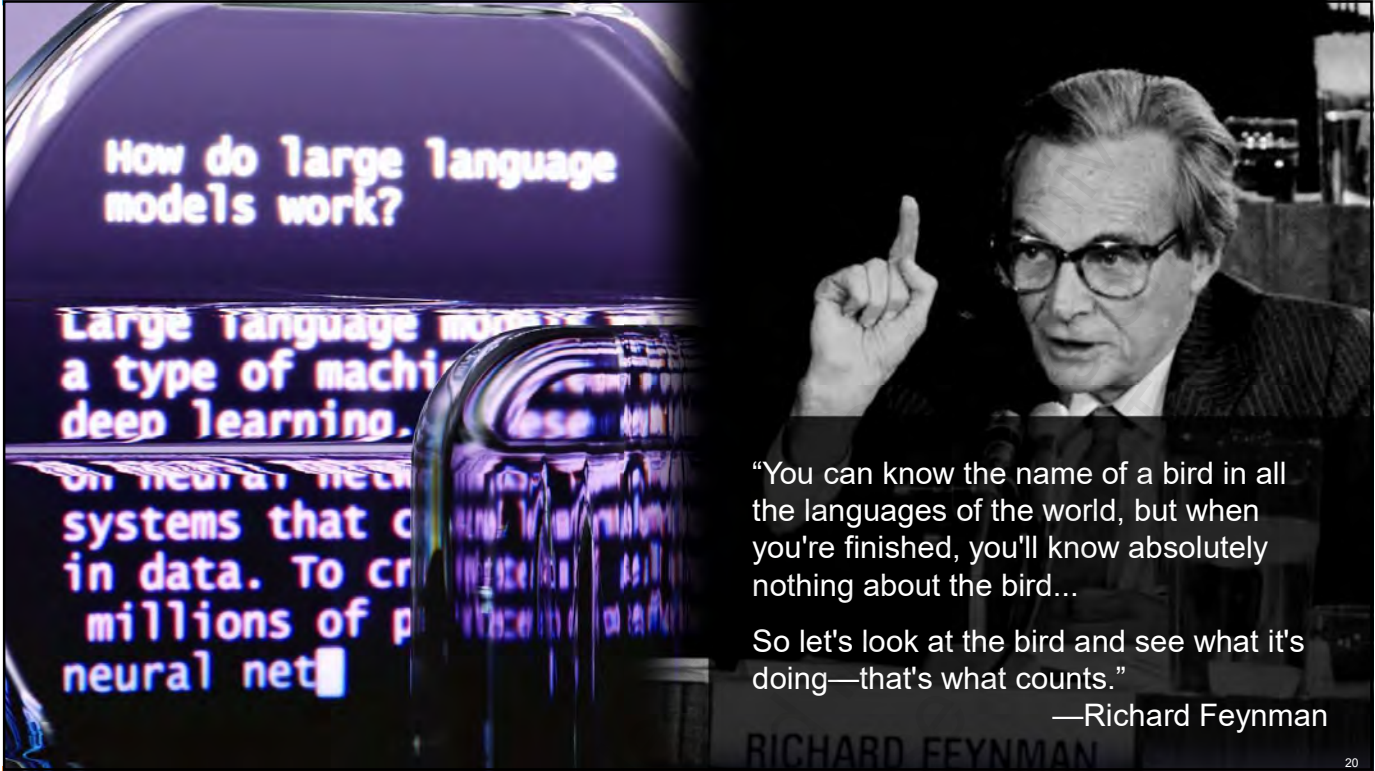
1

---

Michael Wu, PhD (@mich8elwu)
chief AI strategist @ PROS
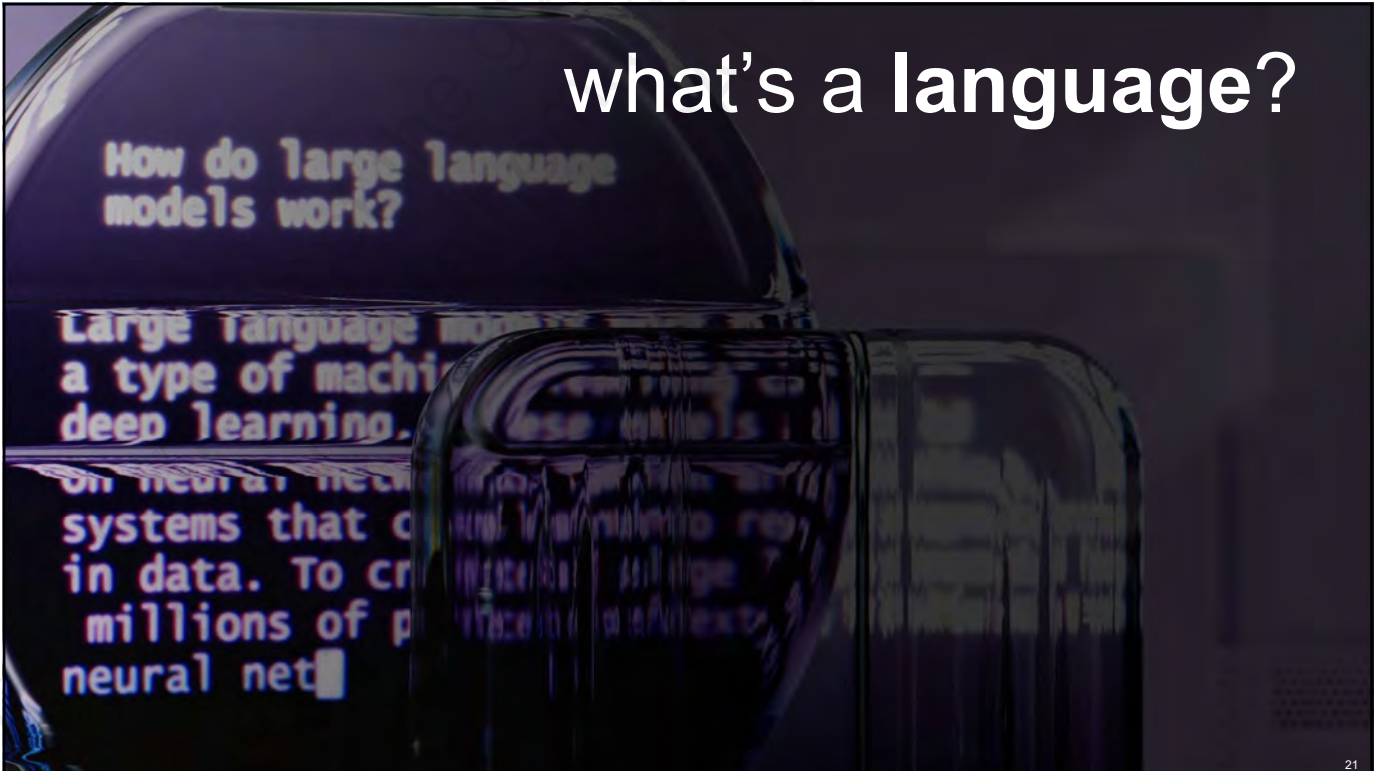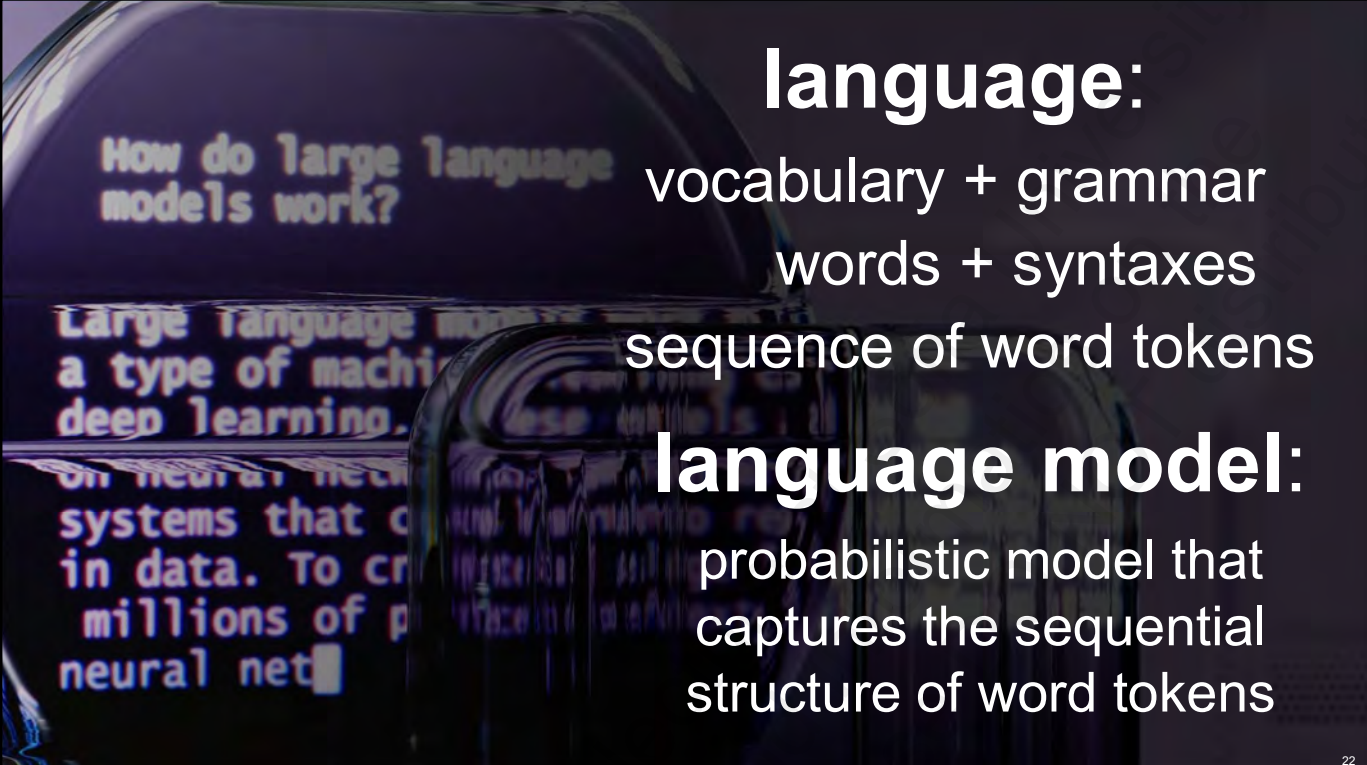
2024.04.19

**PROS.ai**

12

12

"You can know the name of a bird in all the languages of the world, but when you're finished, you'll know absolutely nothing about the bird...

So let's look at the bird and see what it's doing—that's what counts."
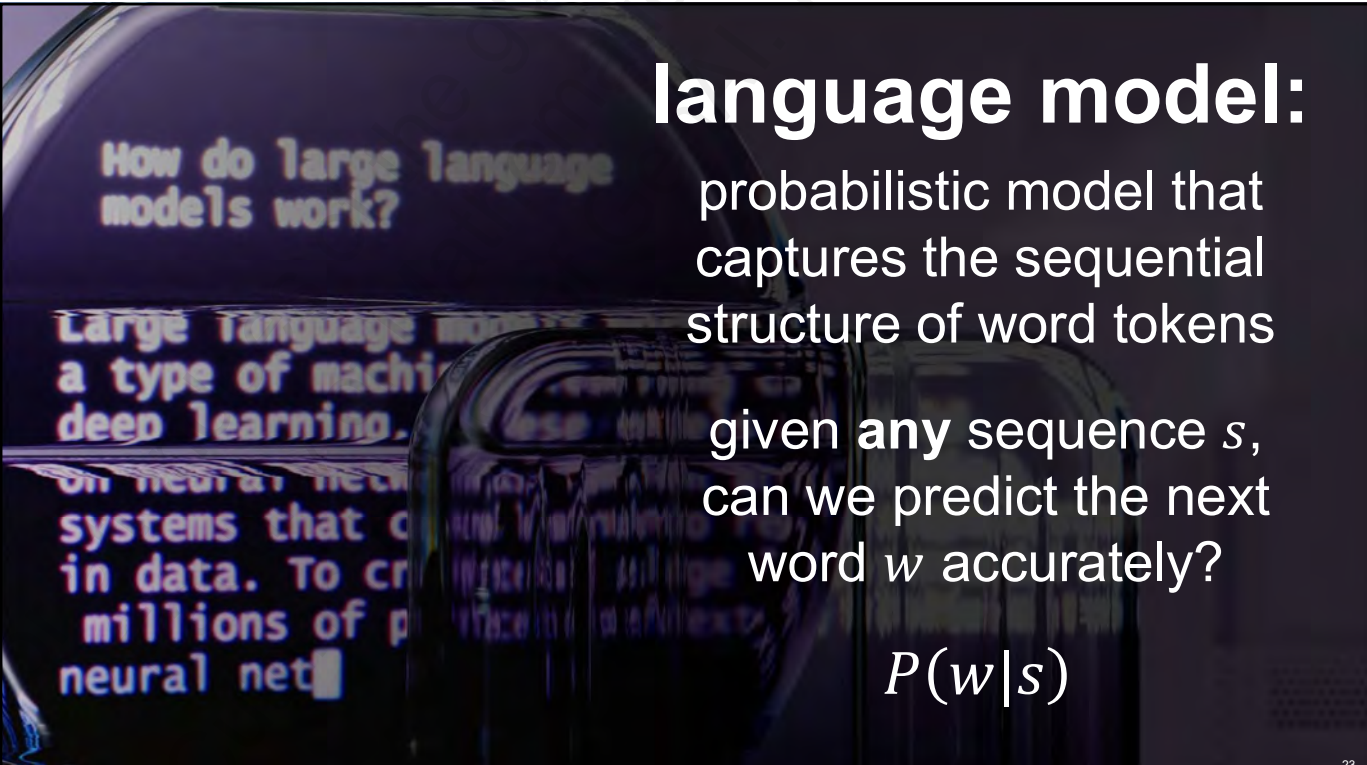—Richard Feynman

20



what's a **language**?

21

**language**:

vocabulary + grammar
words + syntaxes
sequence of word tokens

**language model**:

probabilistic model that
captures the sequential
structure of word tokens

22

22

---

**language model:**

probabilistic model that
captures the sequential
structure of word tokens

given **any** sequence $s$,
can we predict the next
word $w$ accurately?

$$P(w|s)$$
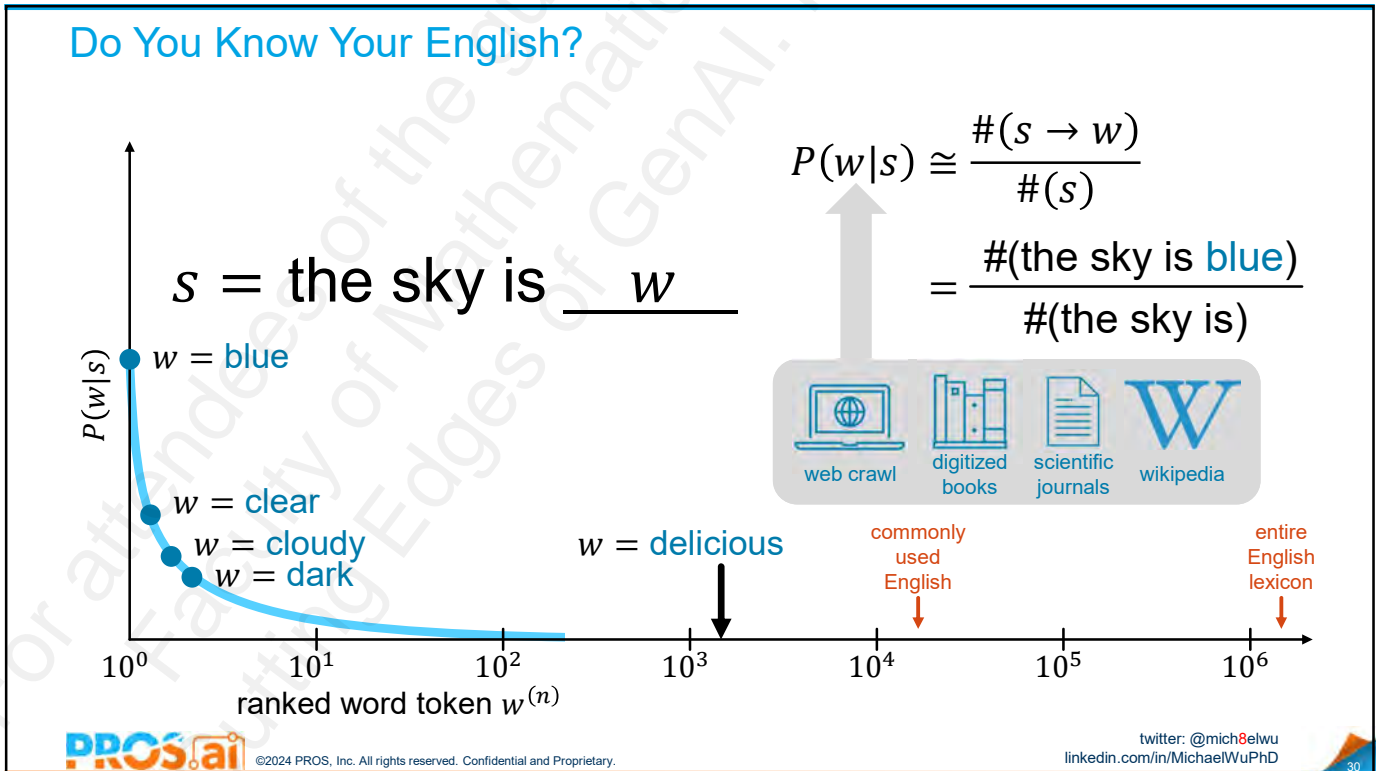
23

23

## language model:

given ***any*** sequence $s$, can we predict the next word $w$ accurately?

$$P(w|s)$$

$s = $ *the water is filled with ribulose-bisphosphate-carboxylase-oxygenase, it's very* _____

How do large language models work?

Large language a type of machi deep learning. systems that in data. To c millions of neural net

24

24

---

### Do You Know Your English?

$$P(w|s) \cong \frac{\#(s \to w)}{\#(s)}$$

$$= \frac{\#(\text{the sky is blue})}{\#(\text{the sky is})}$$

$s = $ the sky is _____ $w$

web crawl · digitized books · scientific journals · wikipedia

$P(w|s)$

$w = $ blue

$w = $ clear
$w = $ cloudy
$w = $ dark

$w = $ delicious

commonly used English

entire English lexicon

$10^0$  $10^1$  $10^2$  $10^3$  $10^4$  $10^5$  $10^6$

ranked word token $w^{(n)}$

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

30

30

## Do You Know Your English?

$$s = \text{the sun just set, the sky is } \underline{\quad w \quad}$$

$$P(w|s) \cong \frac{\#(s \to w)}{\#(s)}$$

$$= \frac{\#(\text{the sun just set, the sky is } w)}{\#(\text{the sun just set, the sky is})}$$

$P(w|s)$

w = blazing
w = golden
w = spectacular

**attention mechanism = spotlight that boosts the word probability**

w = blue
w = clear
w = cloudy
w = dark

web crawl  digitized books  scientific journals  wikipedia  w = spectacular

commonly used English

entire English lexicon

$10^0 \quad 10^1 \quad 10^2 \quad 10^3 \quad 10^4 \quad 10^5 \quad 10^6$

ranked word token $w^{(n)}$

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

33

---

## Longer Sequence → More Concentrated Word Probability

$$s = \text{the sun just set, the sky is } \underline{\quad w \quad}$$

$$P(w|s) \cong \frac{\#(s \to w)}{\#(s)}$$

$$= \frac{\#(\text{the sun just set, the sky is } w)}{\#(\text{the sun just set, the sky is})}$$

$P(w|s)$

**the word probability distribution sharpens**

web crawl  digitized books  scientific journals  wikipedia

commonly used English

entire English lexicon

$10^0 \quad 10^1 \quad 10^2 \quad 10^3 \quad 10^4 \quad 10^5 \quad 10^6$

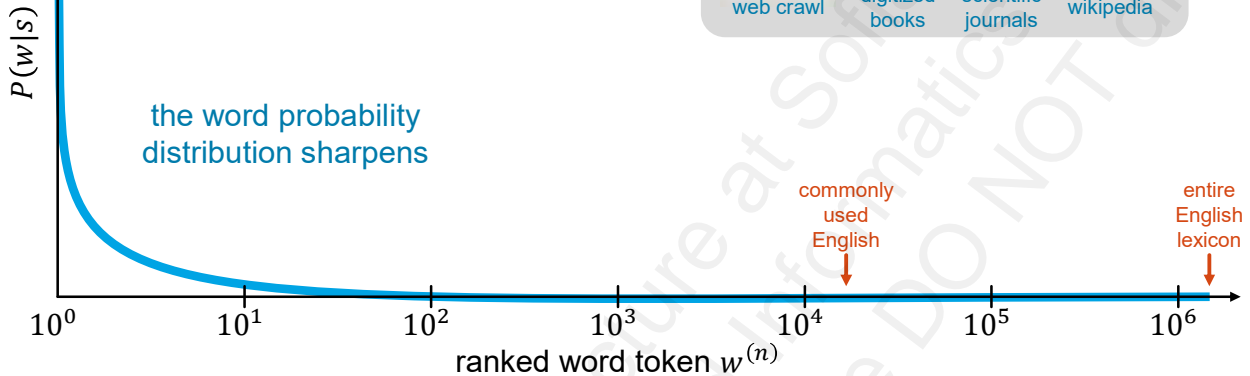ranked word token $w^{(n)}$

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

34

## Longer Sequence → More Concentrated Word Probability

$$s = \text{it's raining hard, the sun just set, the sky is } \underline{\quad w \quad}$$

$$P(w|s) \cong \frac{\#(s \to w)}{\#(s)}$$

web crawl · digitized books · scientific journals · wikipedia

$P(w|s)$

the word probability distribution sharpens

commonly used English

entire English lexicon

ranked word token $w^{(n)}$

$10^0 \quad 10^1 \quad 10^2 \quad 10^3 \quad 10^4 \quad 10^5 \quad 10^6$

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

35

## Longer Sequence → More Concentrated Word Probability

$$s = \text{it's raining hard, the sun just set, the sky is } \underline{\quad w \quad}$$

$$P(w|s) \cong \frac{\#(s \to w)}{\#(s)}$$

web crawl · digitized books · scientific journals · wikipedia

~500B tokens $\ll \infty$

$P(w|s)$

the word probability distribution sharpens even more → able to stay on-topic well

commonly used English

entire English lexicon

ranked word token $w^{(n)}$

$10^0 \quad 10^1 \quad 10^2 \quad 10^3 \quad 10^4 \quad 10^5 \quad 10^6$
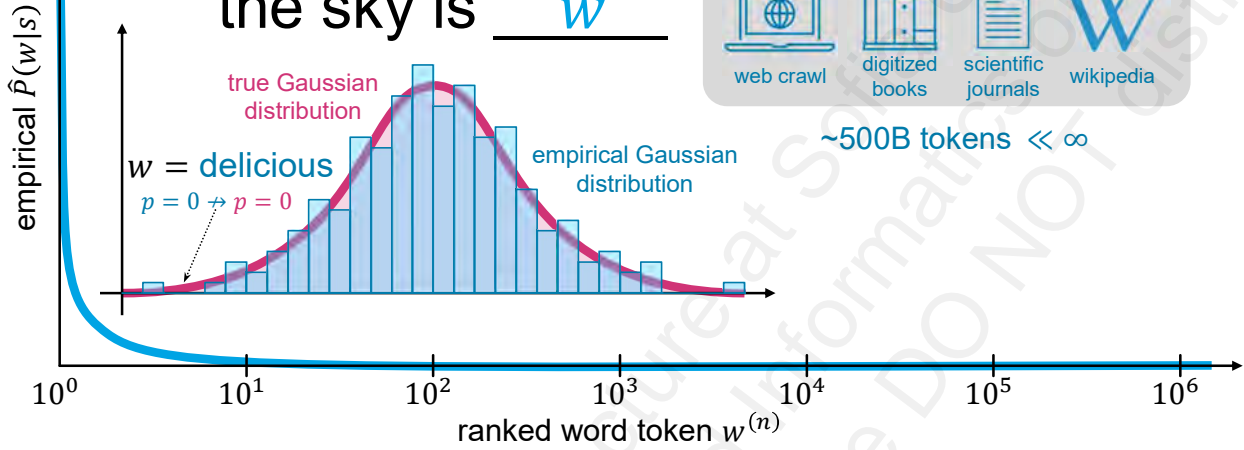
twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

36

## Empirical Word Probability to Language Model

$s =$ it's raining hard, the sun just set, the sky is _____ $w$

$$P(w|s) \cong \frac{\#(s \to w)}{\#(s)}$$

web crawl · digitized books · scientific journals · wikipedia

~500B tokens $\ll \infty$

empirical $\hat{P}(w|s)$

true Gaussian distribution

empirical Gaussian distribution

$w =$ delicious
$p = 0 \nrightarrow p = 0$

ranked word token $w^{(n)}$

$10^0$   $10^1$   $10^2$   $10^3$   $10^4$   $10^5$   $10^6$

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

37

37



38

38

## Empirical Word Probability to Language Model

$s =$ it's raining hard, the sun just set, the sky is ___ $w$ ___

$$P(w|s) \cong \frac{\#(s \to w)}{\#(s)}$$

web crawl · digitized books · scientific journals · wikipedia

~500B tokens $\ll \infty$

true Gaussian distribution

empirical Gaussian distribution

$w =$ delicious
$p = 0 \not\to p = 0$

$p \neq 0 \longrightarrow$ anything is **possible** under the true distribution

**the ability to create novel content = "generative"**

empirical $\hat{P}(w|s)$

ranked word token $w^{(n)}$

$10^0 \quad 10^1 \quad 10^2 \quad 10^3 \quad 10^4 \quad 10^5$

twitter: @
linkedin.com/in/Mich

40

---

## Empirical Word Probability to Language Model

web crawl · digitized books · scientific journals · wikipedia

$$\hat{P}(w|s) = \frac{\#(s \to w)}{\#(s)}$$

$P(w|s) =$ language model

empirical $\hat{P}(w|s)$

ranked word token $w^{(n)}$

$10^0 \quad 10^1 \quad 10^2 \quad 10^3 \quad 10^4 \quad 10^5 \quad 10^6$

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

41

## Empirical Word Probability to Language Model



web crawl · digitized books · scientific journals · wikipedia

$$\hat{P}(w|s) = \frac{\#(s \rightarrow w)}{\#(s)}$$

**TRANSFORMER**

$$P(w|s) = \text{language model}$$

empirical $\hat{P}(w|s)$

ranked word token $w^{(n)}$

$10^0 \quad 10^1 \quad 10^2 \quad 10^3 \quad 10^4 \quad 10^5 \quad 10^6$

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

42

42

## Empirical Word Probability to Language Model



web crawl · digitized books · scientific journals · wikipedia

**TRANSFORMER**

$$\hat{P}(w|s) = \frac{\#(s \rightarrow w)}{\#(s)}$$

**TRANSFORMER**

$$P(w|s) = \text{language model}$$

empirical $\hat{P}(w|s)$

ranked word token $w^{(n)}$

$10^0 \quad 10^1 \quad 10^2 \quad 10^3 \quad 10^4 \quad 10^5 \quad 10^6$

ENCODER
Add & Norm
Feed forward
Add & Norm
Multi-head attention
Positional Encoding
Embedding
INPUT

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

43

43

## Empirical Word Probability to Language Model

web crawl · digitized books · scientific journals · wikipedia

empirical $\hat{P}(w|s)$

$$\hat{P}(w|s) = \frac{\#(s \to w)}{\#(s)}$$

TRANSFORMER

$P(w|s) =$ language model

**Attention Is All You Need**

= spotlight that boosts the word probability

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

Łukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

**Abstract**

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

*NIPS(2017)*

ranked word token $w^{(n)}$

($10^0$, $10^1$, $10^2$, $10^3$, $10^4$, $10^5$, $10^6$)

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

PROS.ai

44

---

## Empirical Word Probability to Language Model

web crawl · digitized books · scientific journals · wikipedia

$P(w|s)$

$$\hat{P}(w|s) = \frac{\#(s \to w)}{\#(s)}$$

TRANSFORMER

$P(w|s) =$ language model

TRANSFORMER

ENCODER — DECODER — GPT

Linear · Softmax · OUTPUT PROBABILITIES

Add & Norm · Feed forward · Multi-head attention · Positional Encoding · Embedding · INPUT

Add & Norm · Multi-head attention · OUTPUTS (shifted right)

$P(w|s)$

sample a word $w$ from $P(w|s)$

$s = \{s + w\}$

ranked word token $w^{(n)}$

($10^0$, $10^1$, $10^2$, $10^3$, $10^4$, $10^5$, $10^6$)

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

PROS.ai

45

Created with the help of Beautiful.ai

**GPT: Generative Pre-trained Transformer
is a large language model (LLM)**

random number generator
from a distribution over *all* words
given *any* word sequence
trained with human written text
using transformer architecture

46

46

doesn't sound
very intelligent

47

47

## From GPT to ChatGPT

**pre-trained via supervised/self-supervised**

web crawl | digitized books | scientific journals | wikipedia

$$\hat{P}(w|s) = \frac{\#(s \to w)}{\#(s)}$$

**GPT**

$P(w|s) =$ language model

**plausible text continuation ≠ good responses**

- supervised *transfer learning* to finetune the model to follow instructions + provide answers

**good responses ≠ good dialog**

- *reinforcement learning* with human-in-the-loop ranking of good dialogue responses

**transfer learning** → **reinforcement learning** →

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

53

53

ChatGPT has no knowledge of the real world

54

54

"*generative*" AI are, by definition and construction, *hallucinatory*

55



very convincing legal briefings citing non-existent court cases

56

Steven Schwartz
+ Peter LoDuca fined $5K

57



hallucination:
a feature
or a bug?

59

feature: for *design* + *creative* use cases

60

bug: for *fact-based* applications

61

# fact-based applications require *grounding*



**empirical Gaussian distribution**

**true Gaussian distribution**

**when there is no data empirically, don't make things up**

conventional info. retrieval (e.g. search) → chatGPT to summarize search result

*retrieval augmented generation*

I don't know

62

# 2 ways to work with LLM



## model fine-tuning

**pros**
- knowledge encoded into the model parameters
- can teach it anything

**cons**
- costly: 25,000 × nvidia A100 for ~100 days ~$63M → GPT4
- must be retrained when there's new data or new LLMs
- hard to iterate, slow time to market

## RAG: prompting

**pros**
- no upfront cost
- no retraining on new data
- easily swap in/out different LLM
- easy to iterate, fast time to market

**cons**
- limited context length (GPT4: ~128k tokens)
- knowledge accuracy depends on retrieval mechanism (search)

65

## A Language Guru with Broad General Knowledge

**think of ChatGPT as a colleague**
- reads lightning fast
- understands any language
- forgetful: small working memory (limited context length)
  - GPT3.5: ~4K tokens
  - GPT4: ~128K tokens
- has broad (non-specific) knowledge
- very imaginative, but overconfident

**how could you leverage and work with someone with such skill?**



LANGUAGES of the WORLD

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

67



68

## Beyond ChatGPT

| data | generic generative AI | | | | | | specialized generative AI | | |
|---|---|---|---|---|---|---|---|---|---|
| | textual | | visual | | audio | | game | specialized design | |
| | text | code | image | video | speech | music | 3D model | biotech | other |
| model | BERT<br>GPT<br>Mistral<br>Claude<br>LaMDA<br>Gemini<br>Perplexity<br>LLaMA | Codex/GPT4<br>Github copilot<br>tabnine<br>stability.ai<br>CodeWhisperer | Dall-E2<br>Make-a-Scene<br>Craiyon<br>Midjourney<br>stable diffusion<br>Imagen<br>nvidia eDiff-I | X-Clip<br>Make-a-Video<br>Imagen Video<br>Sora | Whisper<br>voicebox | Jukebox<br>Riffusion<br>dance diffusion<br>musicLM | DreamFusion<br>nvidia Get3D<br>human MDM | AlphaFold<br>RoseTTAFold | |
| | | | | | **more models to come** | | | | |
| application | general writing<br>summarize +<br>note taking<br>compare/contrast<br>content creation<br>question/answer<br>realtime translation | code generation<br>documentation<br>text to SQL<br>web app builder | image generation<br>media/advertising<br>2D design<br>social media | video generation<br>video edit/modify | voice synthesis<br>voice cloning | song/music creation | | | |
| | | | | | **more use cases to come** | | | | |
| | **many many more start-ups** | | | | | | | | |

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

69

---



" *If all you have is a hammer, everything looks like a nail* "

—Abraham Maslow

70

**generative nature**

**not pre-scriptable
not pre-programmable
created at the moment**

75

75



76

77



78

**Kristina Kashtanova**

79



## Good Use Cases

### Github copilot/Codex

| use case | ask *any* natural language question → SQL generation for a known DB schema |
|---|---|

**prompt** — DB schema, data dictionary, column definitions, etc.

**example** — SQL

"what is the total margin lift for my French customers last quarter?"

€5.78M

**guardrail**
- read-only
- respect access permission

**POLICE LINE • DO NOT CROSS • POLICE LINE • DO NOT CROSS • POLICE LINE**

should **ALWAYS** have some guardrails b/c so much is unknown with GenAI

**value/adoption**
- explain aggregated results
- human languages are imprecise
- step through calculations
- without trust there is no value

80

**generative nature**

**not pre-scriptable
not pre-programmable
created at the moment**

81

81



**why does it work so well?**

96

96

high intra-dimensional correlation

97

97



80% masked

98

98

simple musical tunes: Beethoven – Fur Elise

99



orchestral music: Beethoven – 5th Symphony

100

speech: Steve Job – Love what you do

101



30% masked

102

pop song: PSY – Gangnam Style

103



simple musical tunes

orchestral music

speech

pop song

104

## ChatGPT and Beyond

### 50th b-day poem from ChatGPT

Happy birthday, Dr. Michael Wu
Today we celebrate all that you do
Your wisdom and insight we hold so dear
Your guidance and leadership, always so clear

In the field of revenue management
You are a true expert and phenomenon
Your knowledge and skills have paved the way
For success and growth, day by day

Your dedication and hard work
Are an inspiration, no need to lurk
Your passion for learning and innovation
Have brought us all great motivation

So on this special day, we honor you
For all that you have accomplished, it's true
We wish you joy, happiness, and health
May your dreams and goals be realized with wealth

Happy birthday, Dr. Michael Wu
We are grateful for all that you do
May your year ahead be filled with delight
And your future be forever bright

**Boris Eldagsen wins Sony World Photography Award**

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

105

**Jason Allen wins Colorado State Fair Digital Art Competition**

106

107



108

correlated data attributes are highly predictable

109



where does that leave software engineers and data scientists?

116

**enterprise codegen use cases:**

- inline auto-complete
- commenting
- test cases
- documentation

*coding liberator*

- syntax debugger
- error-log interpreter
- code translator
  (R → python, C++ → Java)
- refactoring suggestion

*could be a job terminator*

117

117

**" engineering is a creative practice "**

Cory Arcangel

118

119



## Empirical Word Probability to Language Model

$$\hat{P}(w|s) = \frac{\#(s \to w)}{\#(s)}$$

TRANSFORMER

$$P(w|s) = \text{language model}$$

TRANSFORMER

GPT

$$P(w|s)$$

sample a word $w$ from $P(w|s)$

$$s = \{s + w\}$$

web crawl • digitized books • scientific journals • wikipedia

ranked word token $w^{(n)}$

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

120

120

most data scientists in the industry work with some sort of structured tabular data

127

127

## Properties of Tabular Data

*tabular data*:
lower dimensionality
lack of internal correlation

twitter: @mich8elwu
linkedin.com/in/MichaelWuPhD

128

128

129



131

the future of software development is one where engineers collaborate with different kinds of AI tools

**users must have the relevant expertise in discerning and mitigating risk of hallucinations**

132

132



AI *will* ***not*** replace you,
but someone using
AI ***will***

**users must have the relevant expertise in discerning and mitigating risk of hallucinations**

133

133

145



149