

Language Models

Guess the next word: I enjoy drinking coffee, because I like...

- Language or code has low information density
- Plenty of equivalent formulations / phrasings
- Simplistic model: Markov chain n-grams

1-gram:

$P(\text{next word} \mid \text{"like"})$

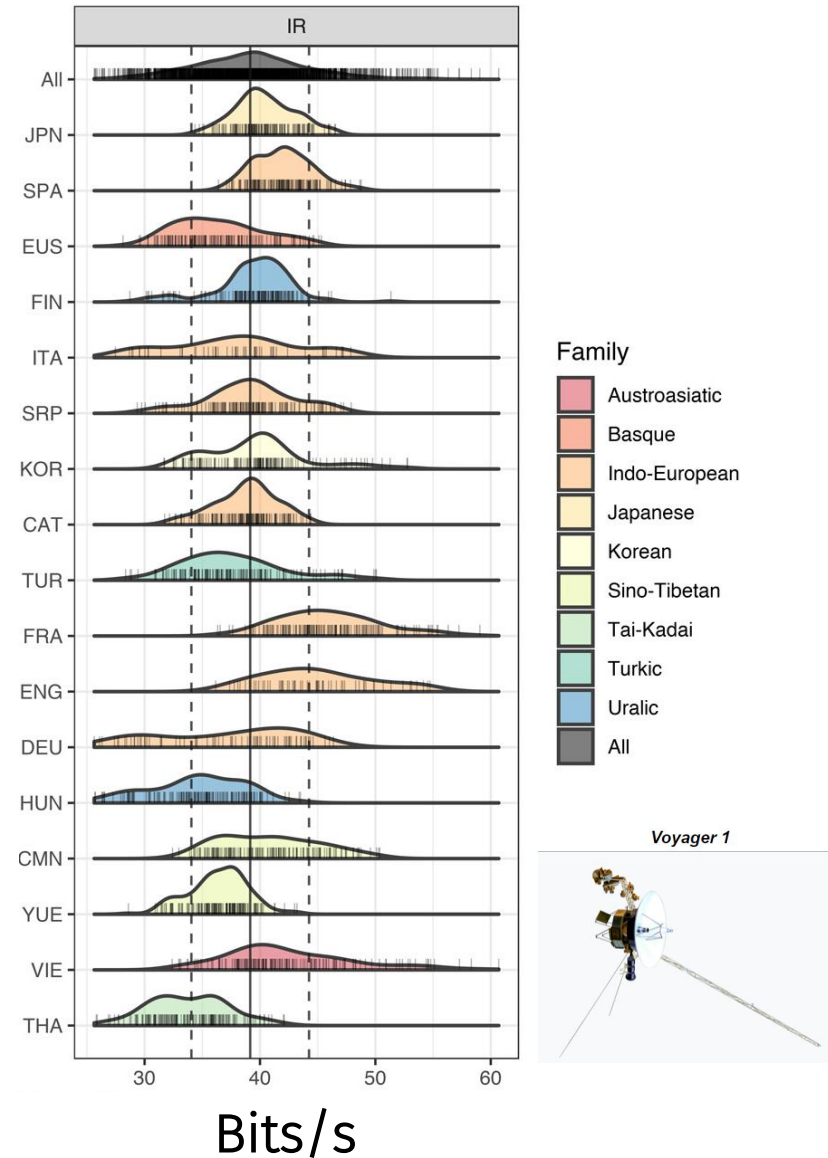
2-gram:

$P(\text{next word} \mid \text{"I like"})$

3-gram:

$P(\text{next word} \mid \text{"because I like"})$

... no coffee!



Tokenisation

- Split text into a finite set of substrings
- May already include some context

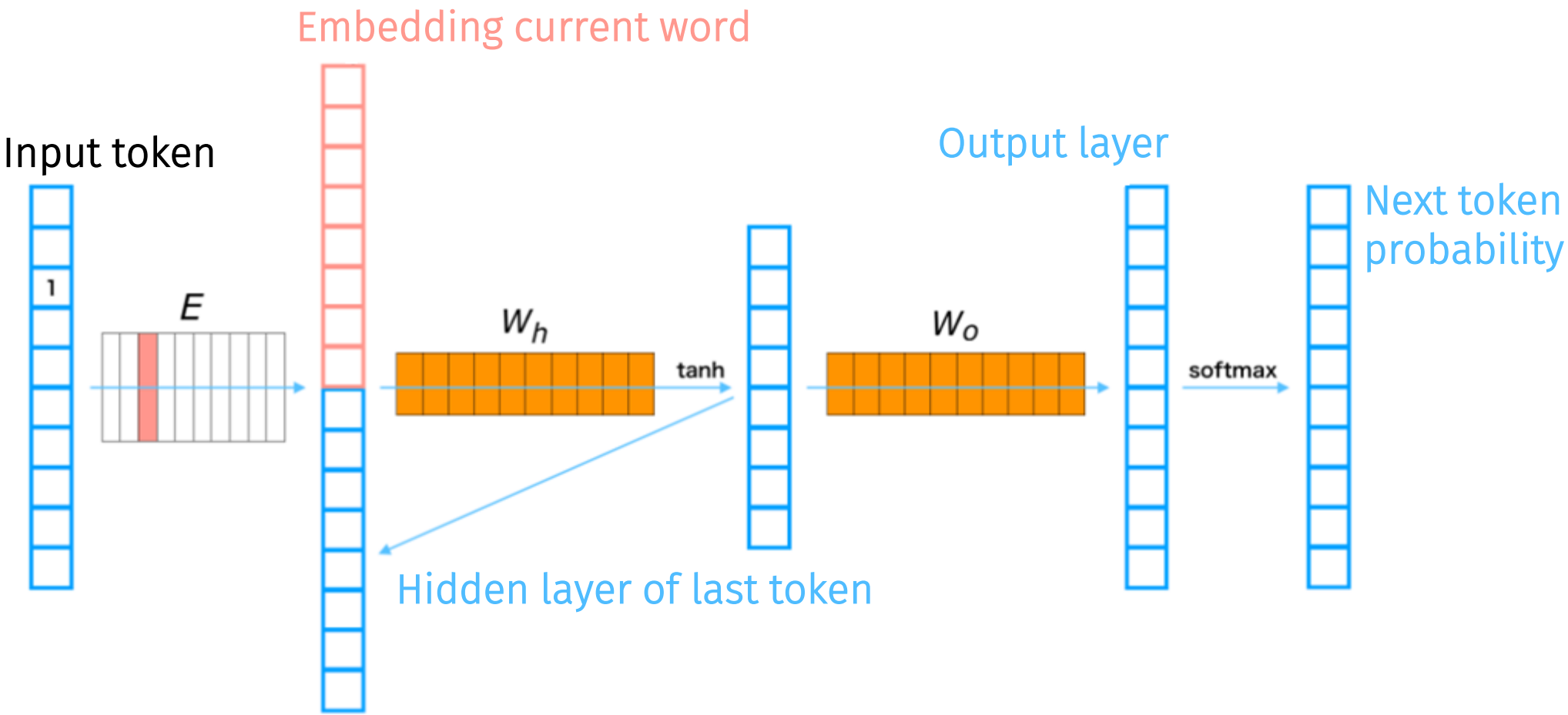


Tokens	Characters
10	32

Caffeine
Club sandwich
Golf Club

Embedding

- Starts from one vector for each token
- Typically fed through neural network to include context



Self-attention

- Output y is weighted sum of all embeddings x so far
- Weights calculated, not trained
- Parameters introduced to tailor to three roles
 - Query, Key, Value

$$\mathbf{q}_i = \mathbf{W}_q \mathbf{x}_i \quad \mathbf{k}_i = \mathbf{W}_k \mathbf{x}_i \quad \mathbf{v}_i = \mathbf{W}_v \mathbf{x}_i$$

$$w'_{ij} = \mathbf{q}_i^T \mathbf{k}_j$$

$$w_{ij} = \text{softmax}(w'_{ij})$$

$$\mathbf{y}_i = \sum_j w_{ij} \mathbf{v}_j .$$

$$\mathbf{y}_i = \sum_j w_{ij} \mathbf{x}_j$$
$$w'_{ij} = \mathbf{x}_i^T \mathbf{x}_j \quad w_{ij} = \frac{\exp w'_{ij}}{\sum_j \exp w'_{ij}}$$

Attention Is All You Need

Architecture: free to choose, as long as it includes self-attention

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

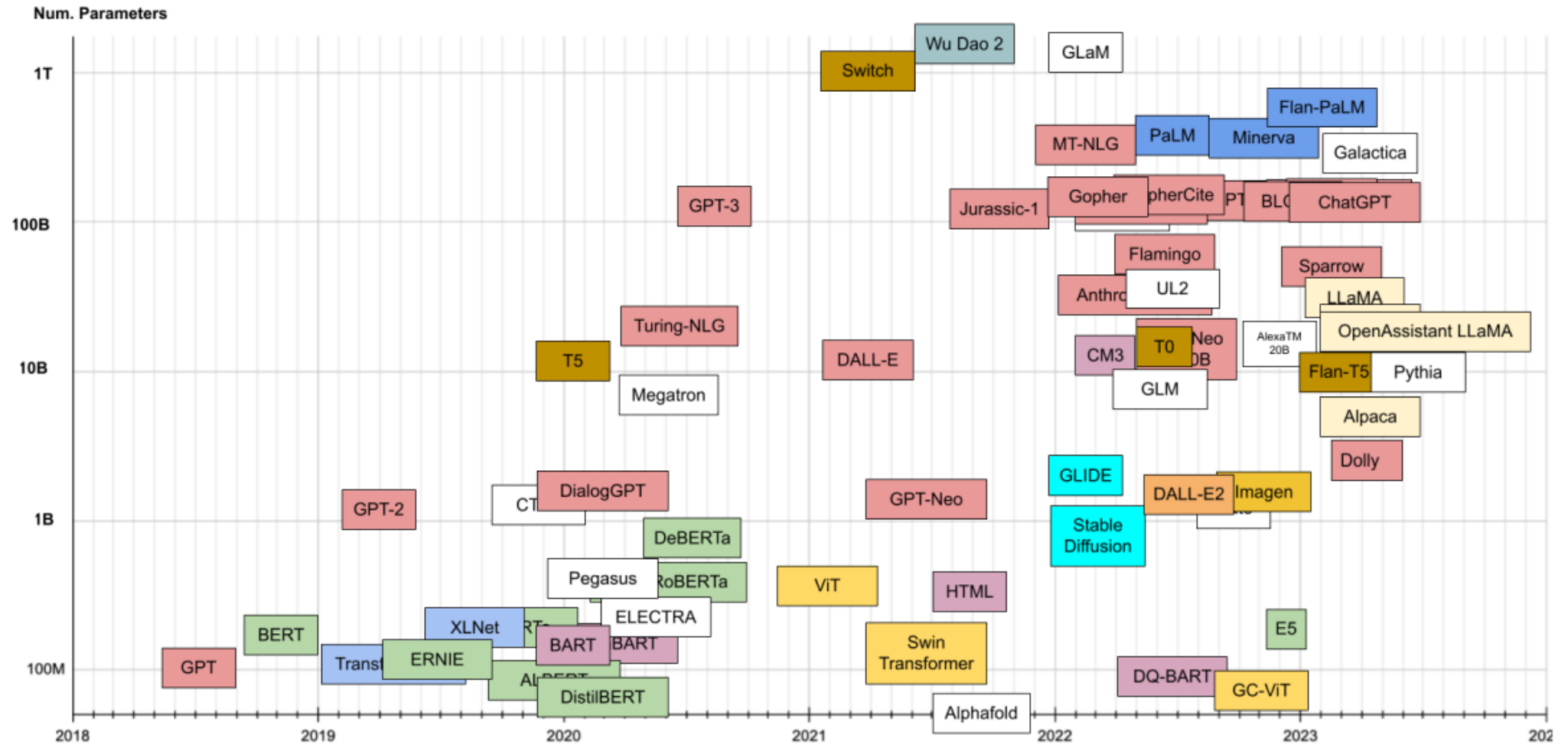
Jakob Uszkoreit*
Google Research
usz@google.com

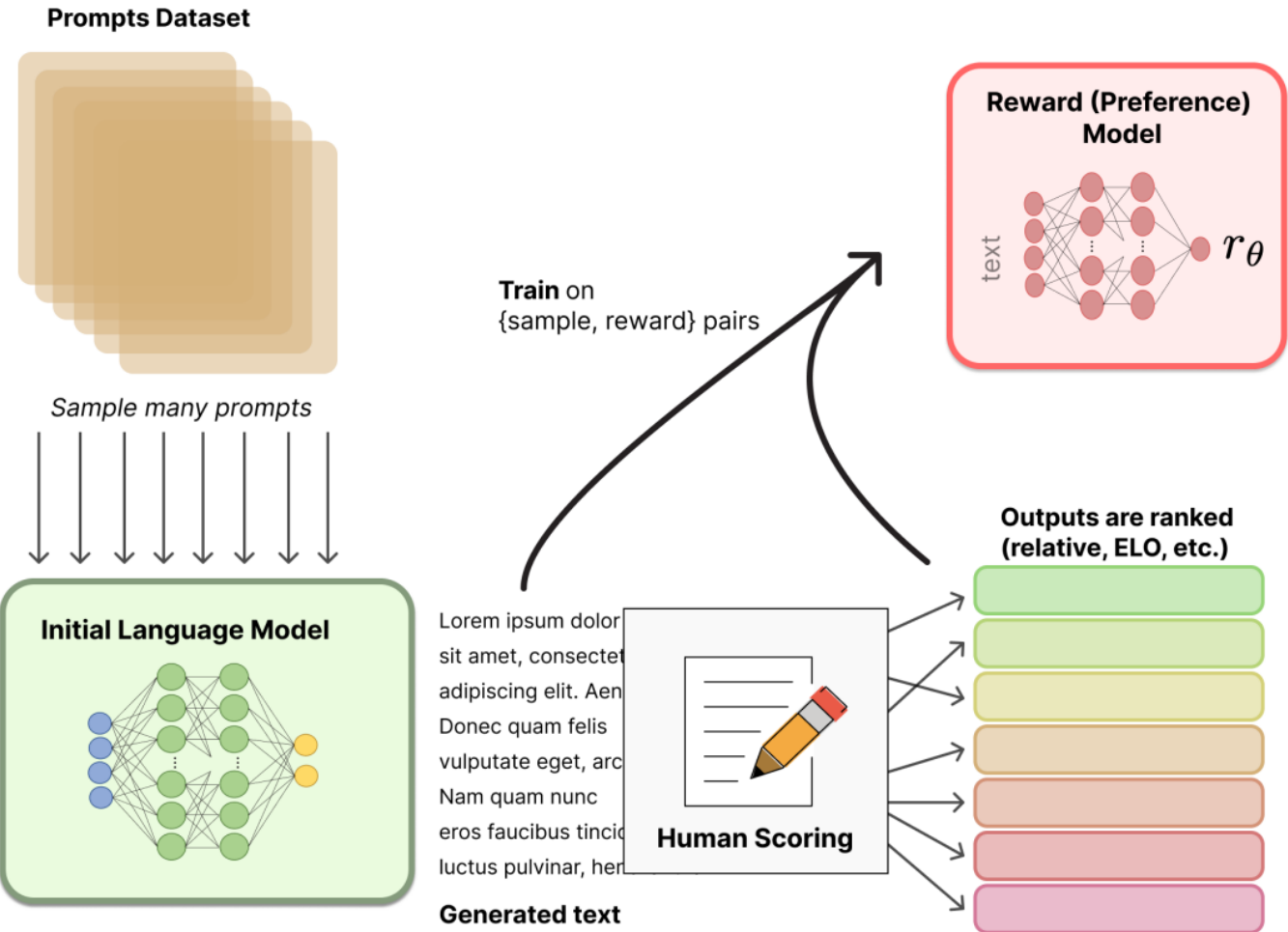
Llion Jones*
Google Research
llion@google.com

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

Lukasz Kaiser*
Google Brain
lukaszkaizer@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com





Reasoning

- Statistical model, not a factual one
- Sounds trustworthy, since human-like

Training data

- (Data) privacy
- Biases

Transparency and Governance

- Intransparent model
- Intransparent access

Reliability

- Code generation: edge cases, API changes
- Statistical output

Summary Large Language Models

- Trying to build long-range context in sequential data
- Expensive in training
- Context-aware predictions
- Helps in data mining, code generation, interpretation, summaries