

# MBAi 448 | Winter 2026

## Transformers

# Today

- Assignment 4 / week 5 takeaways *10 minutes*
- Transformers, so what? *10 minutes*
- What is the right AI to use? *20 minutes*
- Quick break *5 minutes*
- Transformers case studies *25 minutes*
- Assignment walkthrough *20 minutes*

# Assignment 4 / week 5 takeaways

Transformers

# Week 5 takeaways

- **Transformers fundamentally changed how machines understand language—and beyond.**  
Unlike earlier models that processed text sequentially, transformers use attention to evaluate all parts of an input at once, enabling far better handling of context, ambiguity, and long-range dependencies. This architectural shift underpins today's large language models and is now extending to images, audio, and video.
- **“Attention” is the business-relevant breakthrough.**  
Attention mechanisms allow models to focus on what matters most in an input (e.g., key words in a customer request or critical clauses in a document). For businesses, this enables more accurate classification, summarization, routing, and generation—core capabilities behind copilots, chatbots, and automation tools.
- **Different transformer architectures map to different business tasks.**  
Encoder models (e.g., BERT) excel at understanding and classification, decoders (e.g., GPT) at generation, and encoder-decoder models (e.g., T5, BART) at transformation tasks like translation or summarization. Choosing the right architecture is a strategic product decision, not just a technical one.

# Week 5 takeaways

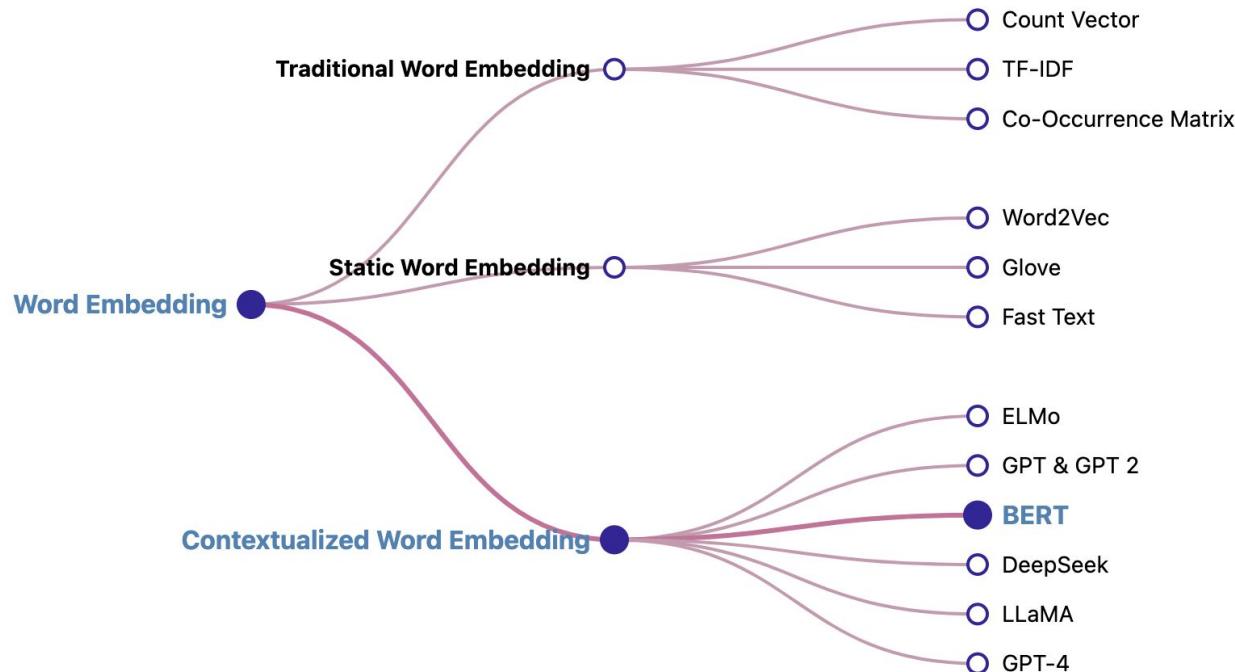
- **Self-supervised learning ~~unlocked scale and lowered dependency on~~ enables deep learning at scale in the absence of explicitly labeled data.**

Transformers learn from massive amounts of raw data by predicting missing or reordered information, eliminating the need for expensive human labeling. This is why AI capabilities scaled so rapidly—and why data access, not algorithms, is often the key competitive advantage.
- **Benchmark performance doesn't equal business readiness.**

While models can achieve “gold-medal” scores on standardized tests, real-world deployment ~~introduces risks~~ tasks rarely resemble standardized tests. Further, the models themselves are but statistical learners and so implementations must account for ~~introduces risks~~: hallucinations, bias, ambiguity, and data drift. Effective applications require guardrails, evaluation frameworks, and thoughtful human-in-the-loop design.
- **The strategic value lies in applications, not models.**

~~Transformers are best understood as general purpose reasoning engines.~~ Competitive advantage comes from how firms integrate them into workflows—routing decisions, drafting responses, summarizing knowledge, or powering agentic systems—not from building the largest model.

# Transformers did not invent word embeddings



# Transformers represent words as tokens

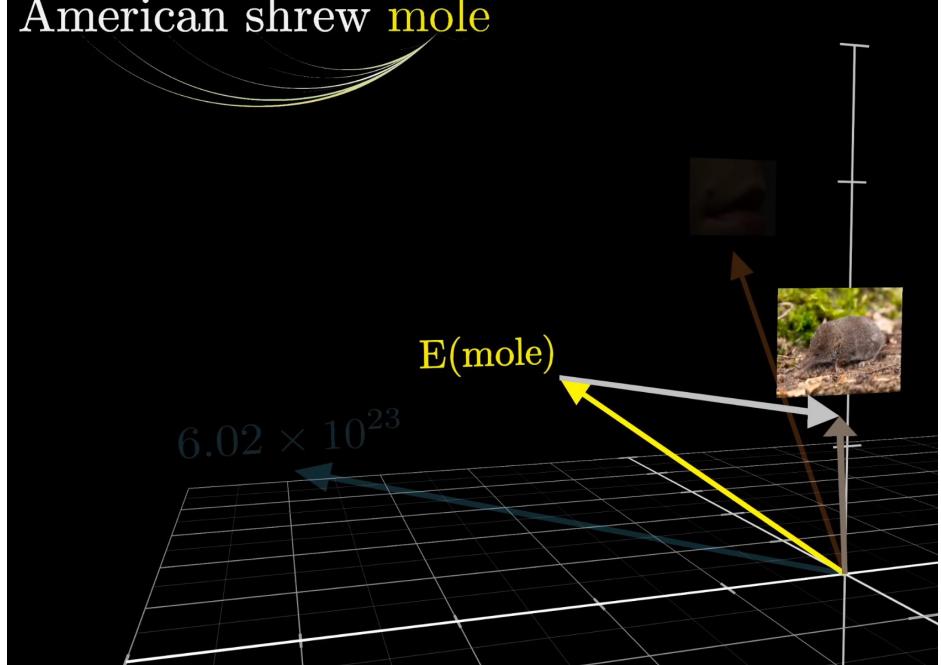
Edit `mole` 1/6 ⌂ ⌃ ⌁ ⌄ The highlighted chips match the tokens the model receives.

30...	5765	6...	1...	17463	10...	34...	10...	17463	3...	15883	70513	1753			
"The	American	shrew	mole	exhal	es	one	mole	of	carbon	dioxide	every				
2944	6...	34248	290	41801	134212	1	3164	12183	3...	1803	6...	2...	161485	1...	1
month",	advised	the	exhibit	curator.	They	moved	to	hand	me	a pamphlet,					
3...	3...	17829	1023	14...	2..	3291	17463	5862	290	3611	3...	1043	1803	1	
and	I	noticed	they	had	a small	mole	near	the base	of	their hand.					

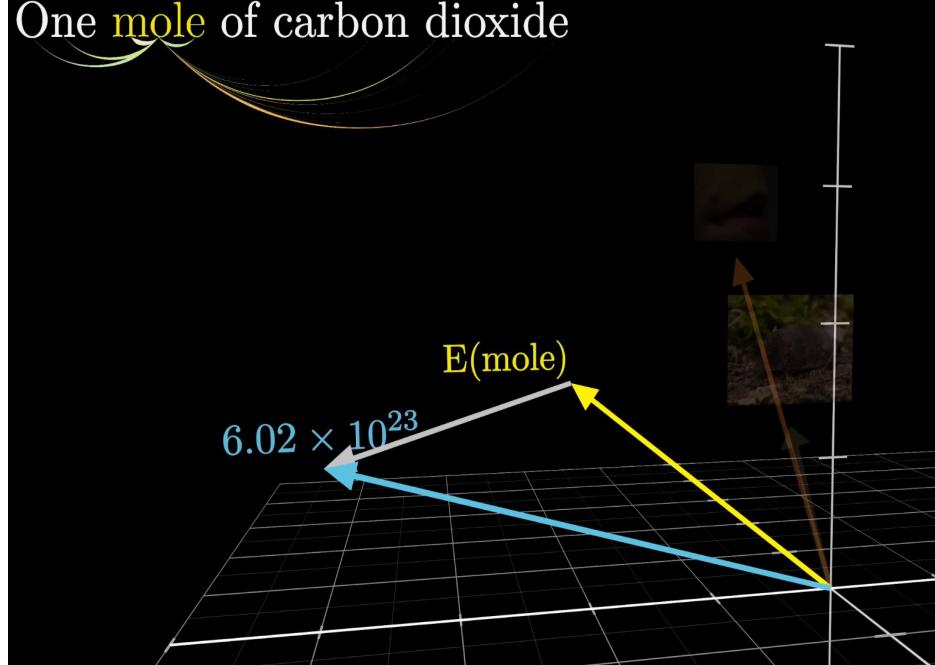
Token count: **44** Live tokenizer view

# Different word senses use the same token?

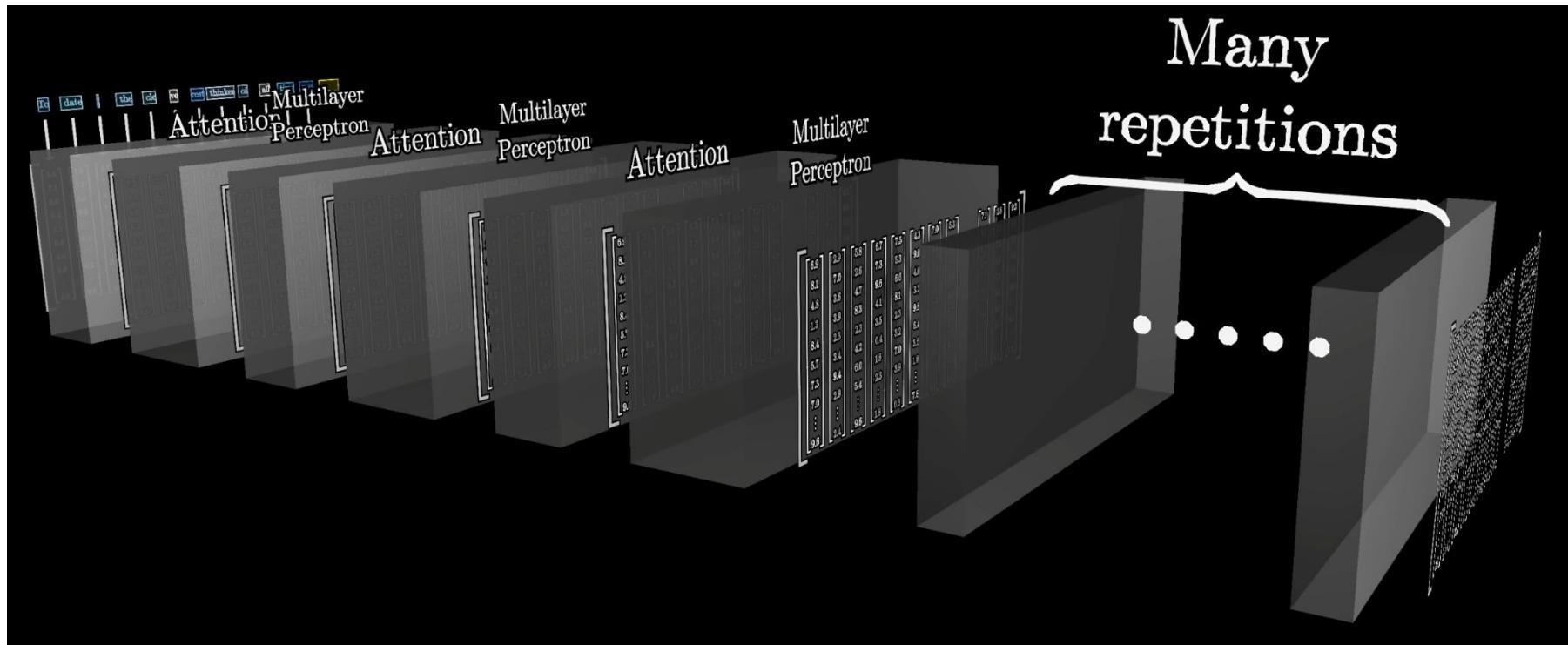
American shrew **mole**

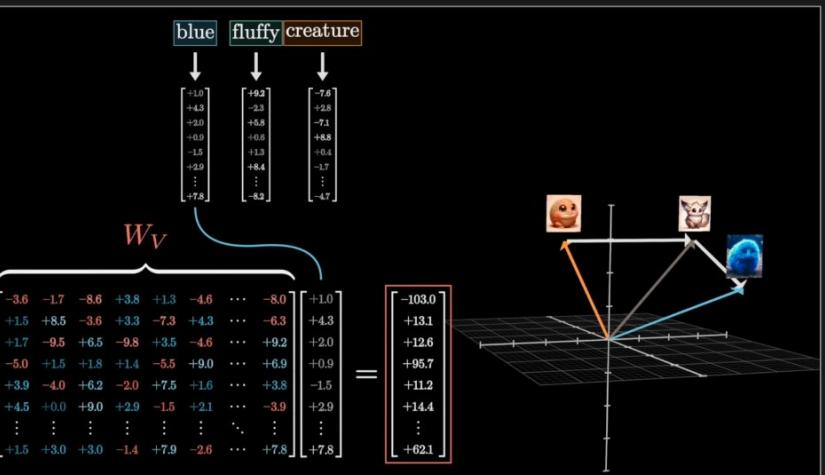
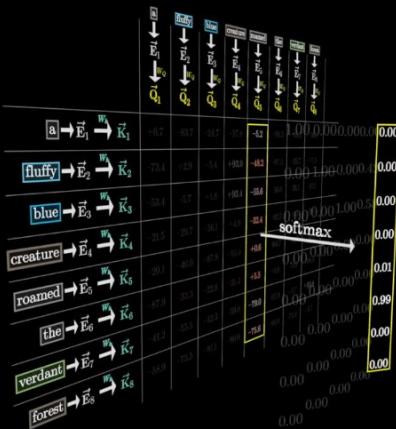
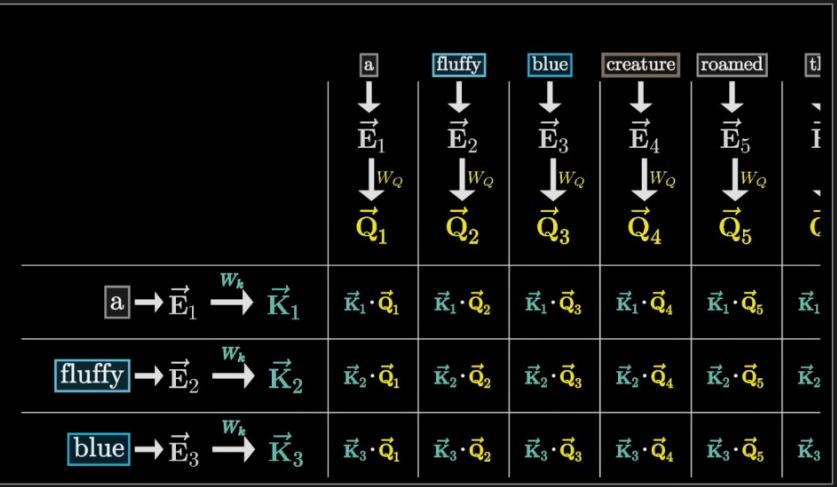


One **mole** of carbon dioxide



# Large models have capacity to learn a lot



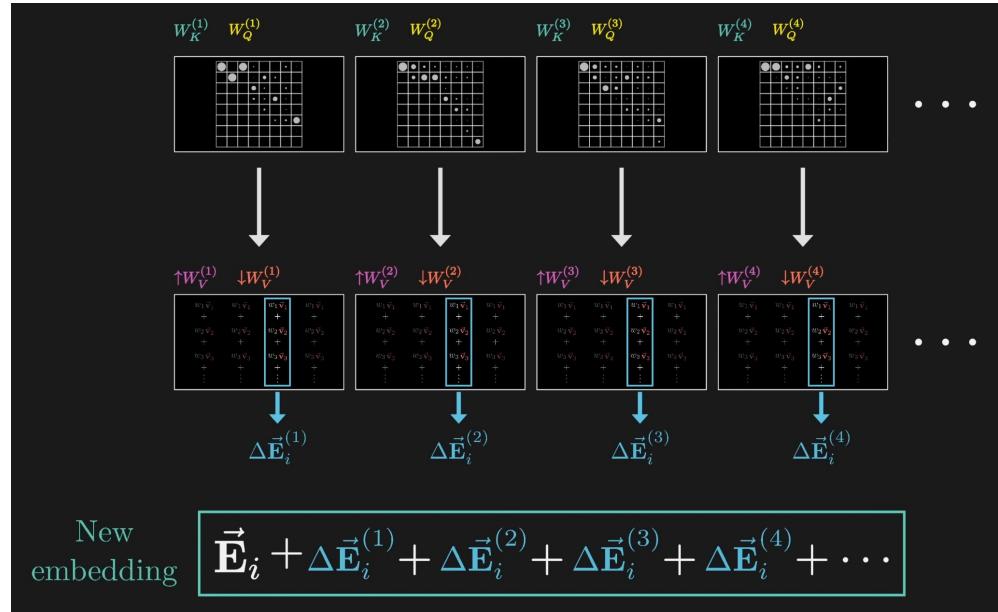
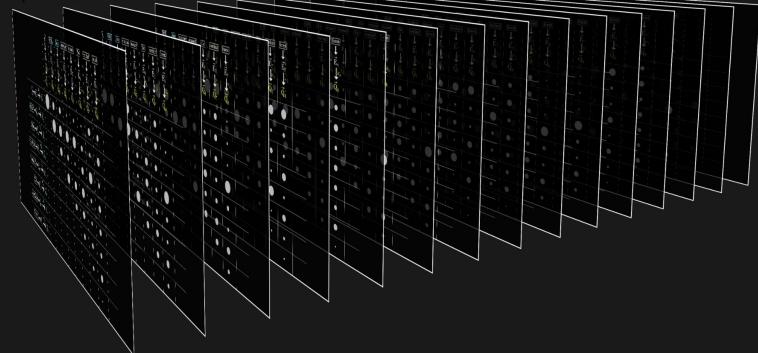


	$\vec{E}_1$	$\vec{E}_2$	$\vec{E}_3$	$\vec{E}_4$	$\vec{E}_5$	$\vec{E}_6$	$\vec{E}_7$	$\vec{E}_8$
$\vec{E}_1 \xrightarrow{W_V} \vec{V}_1$	1.00 $\vec{V}_1$	0.00 $\vec{V}_1$						
$\vec{E}_2 \xrightarrow{W_V} \vec{V}_2$	0.00 $\vec{V}_2$	1.00 $\vec{V}_2$	0.00 $\vec{V}_2$	0.42 $\vec{V}_2$	0.00 $\vec{V}_2$	0.00 $\vec{V}_2$	0.00 $\vec{V}_2$	0.00 $\vec{V}_2$
$\vec{E}_3 \xrightarrow{W_V} \vec{V}_3$	0.00 $\vec{V}_3$	0.00 $\vec{V}_3$	1.00 $\vec{V}_3$	0.58 $\vec{V}_3$	0.00 $\vec{V}_3$	0.00 $\vec{V}_3$	0.00 $\vec{V}_3$	0.00 $\vec{V}_3$
$\vec{E}_4 \xrightarrow{W_V} \vec{V}_4$	0.00 $\vec{V}_4$							
$\vec{E}_5 \xrightarrow{W_V} \vec{V}_5$	0.00 $\vec{V}_5$	0.00 $\vec{V}_5$	0.00 $\vec{V}_5$	0.00 $\vec{V}_5$	0.01 $\vec{V}_5$	0.00 $\vec{V}_5$	0.00 $\vec{V}_5$	0.00 $\vec{V}_5$
$\vec{E}_6 \xrightarrow{W_V} \vec{V}_6$	0.00 $\vec{V}_6$	0.00 $\vec{V}_6$	0.00 $\vec{V}_6$	0.00 $\vec{V}_6$	0.99 $\vec{V}_6$	1.00 $\vec{V}_6$	0.00 $\vec{V}_6$	0.00 $\vec{V}_6$
$\vec{E}_7 \xrightarrow{W_V} \vec{V}_7$	0.00 $\vec{V}_7$	1.00 $\vec{V}_7$	1.00 $\vec{V}_7$					
$\vec{E}_8 \xrightarrow{W_V} \vec{V}_8$	0.00 $\vec{V}_8$							
	$\Delta\vec{E}_1$	$\Delta\vec{E}_2$	$\Delta\vec{E}_3$	$\Delta\vec{E}_4$	$\Delta\vec{E}_5$	$\Delta\vec{E}_6$	$\Delta\vec{E}_7$	$\Delta\vec{E}_8$

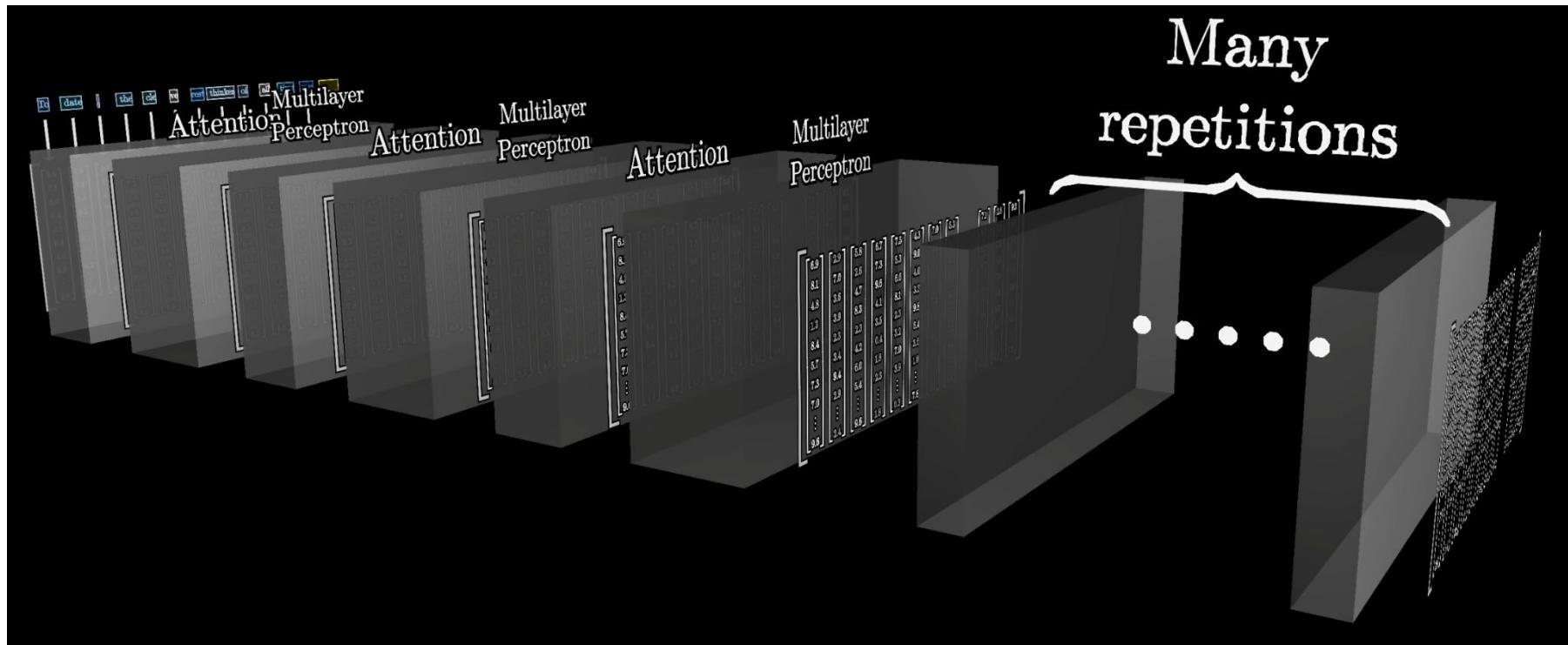
# Multi-head attention learns different relations

## Multi-headed attention

$$\begin{matrix} W_Q^{(1)} & W_Q^{(2)} & W_Q^{(3)} & W_Q^{(4)} & W_Q^{(5)} & W_Q^{(6)} & W_Q^{(7)} & W_Q^{(8)} & W_Q^{(9)} & \dots \\ W_K^{(1)} & W_K^{(2)} & W_K^{(3)} & W_K^{(4)} & W_K^{(5)} & W_K^{(6)} & W_K^{(7)} & W_K^{(8)} & W_K^{(9)} & \dots \\ \downarrow W_V^{(1)} & \downarrow W_V^{(2)} & \downarrow W_V^{(3)} & \downarrow W_V^{(4)} & \downarrow W_V^{(5)} & \downarrow W_V^{(6)} & \downarrow W_V^{(7)} & \downarrow W_V^{(8)} & \downarrow W_V^{(9)} & \dots \\ \uparrow W_V^{(1)} & \uparrow W_V^{(2)} & \uparrow W_V^{(3)} & \uparrow W_V^{(4)} & \uparrow W_V^{(5)} & \uparrow W_V^{(6)} & \uparrow W_V^{(7)} & \uparrow W_V^{(8)} & \uparrow W_V^{(9)} & \dots \end{matrix}$$



# Attention is not all you need



# Transformers, so what?

Transformers

# AI enables a variety of discrete capabilities

Capability	Task	Examples
Classify / Label	Assign inputs to predefined categories	Type, triage, diagnose, etc.
Extract	Identify a specific part of an unstructured input	Segmentation of image, audio, etc.
Group	Divide the input into distinct subsets	Cluster, detect anomalies, etc.
Score / Predict	Estimate a continuous numeric value	Forecast, quantify, rank, etc.
Search	Explore to find the best options for the input	Retrieve, route, match, etc.

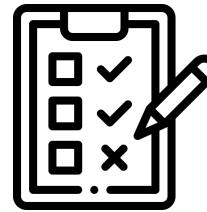
# Different types of AI work with different data



## Type

---

- Structured (tables, defined fields)
- Text (documents, messages, descriptions)
- Images, audio, video, sensor data
- Multimodal



## Quality

---

- Volume / size
- Representativeness
- Labels
- Accessibility

# Real world constraints inform the tech choice

Constraint	Consideration
Explainability	Must decisions be justified to regulators, customers, or internal stakeholders?
Latency	Real-time (milliseconds)? Interactive (seconds)? Batch (minutes/hours)?
Cost	What can you spend per prediction? On infrastructure? On development?
Availability	Do you have enough examples? Can you acquire or create more?
Resources	What does your team have the expertise and capacity to build and maintain?

# Deciding what AI to use when requires knowing your task and data circumstances

1. Can you specify the rules?
2. Can you identify the features?
3. Is there a pretrained model for your data type?
4. Is there data you can use for fine-tuning?
5. Do you have the resources to train a model from scratch?

# Transformers enabled the training of larger, more general pretrained models

1. Can you specify the rules?
2. Can you identify the features?
3. Is there a pretrained model for your data type?
4. Is there data you can use for fine-tuning?
5. Do you have the resources to train a model from scratch?

# Transformers enabled the training of larger, more general pretrained models

1. Can you specify the rules?
2. Can you identify the features?
3. Is there a pretrained model for your data type?
4. Is there data you can use for fine-tuning?
5. Do you have the resources to train a model from scratch?

# Transformers use cases

Transformers

# Case Study: Google Search

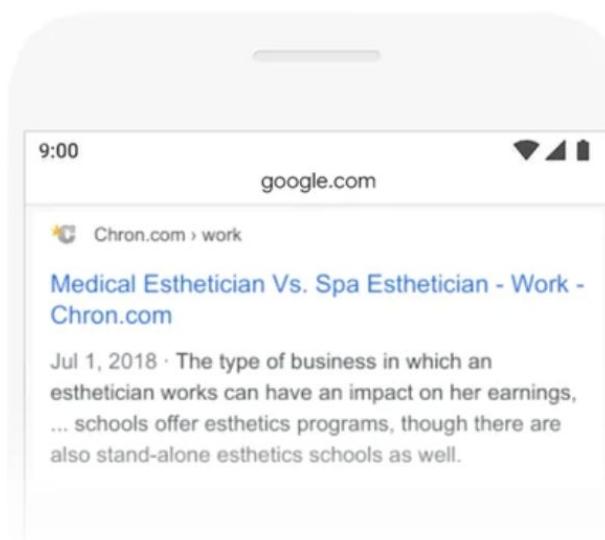
2019 brazil traveler to usa need a visa

BEFORE



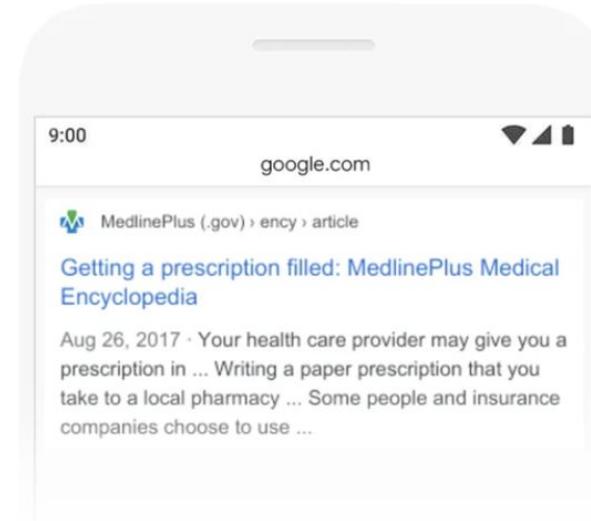
do estheticians stand a lot at work

BEFORE



Can you get medicine for someone pharmacy

BEFORE



# Case Study: Google Search

## Cracking your queries

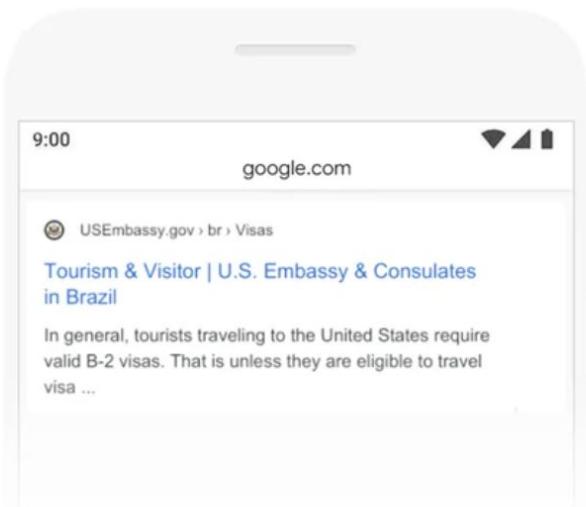
So that's a lot of technical details, but what does it all mean for you? Well, by applying BERT models to both ranking and featured snippets in Search, we're able to do a much better job helping you find useful information. In fact, when it comes to ranking results, BERT will help Search better understand one in 10 searches in the U.S. in English, and we'll bring this to more languages and locales over time.

Particularly for longer, more conversational queries, or searches where prepositions like "for" and "to" matter a lot to the meaning, Search will be able to understand the context of the words in your query. You can search in a way that feels natural for you.

# Case Study: Google Search

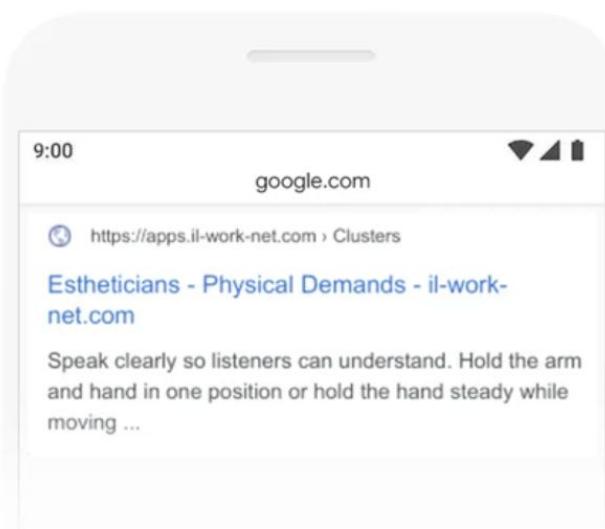
2019 brazil traveler to usa need a visa

AFTER



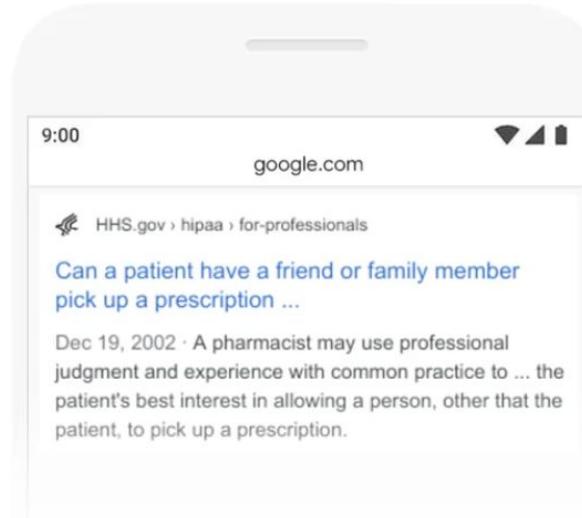
do estheticians stand a lot at work

AFTER



Can you get medicine for someone pharmacy

AFTER



# Case Study: AirBnb

Inspiration for your next trip



**South Lake Tahoe**

160 miles away



**Paso Robles**

154 miles away



**Arnold**

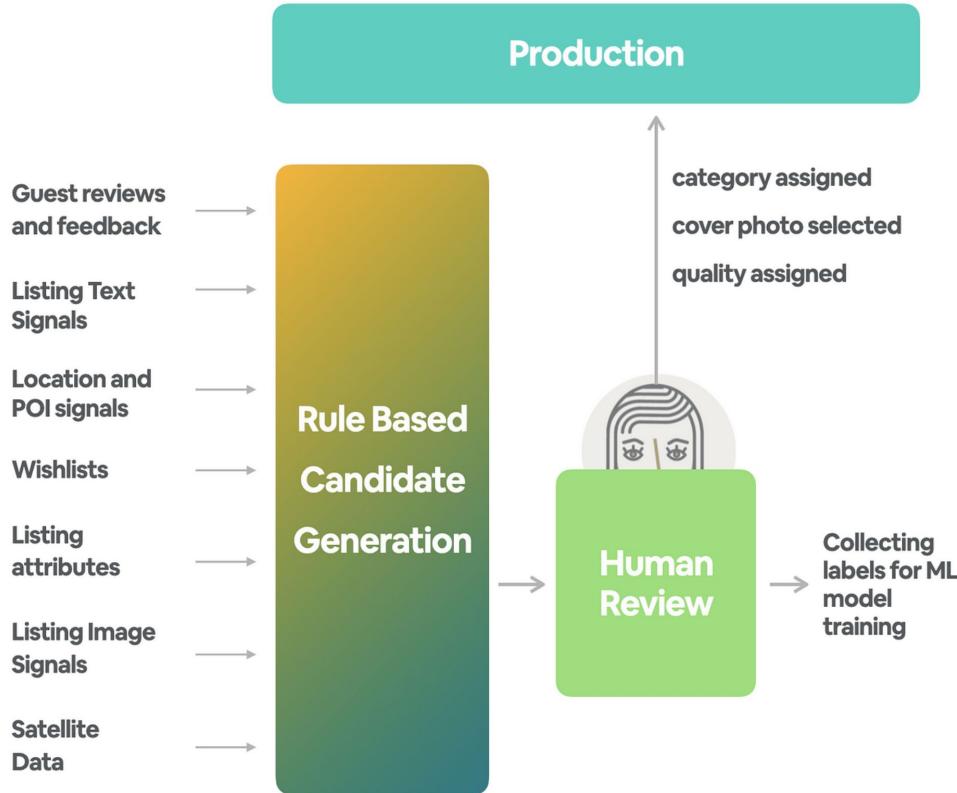
117 miles away



**Carmel-by-the-Sea**

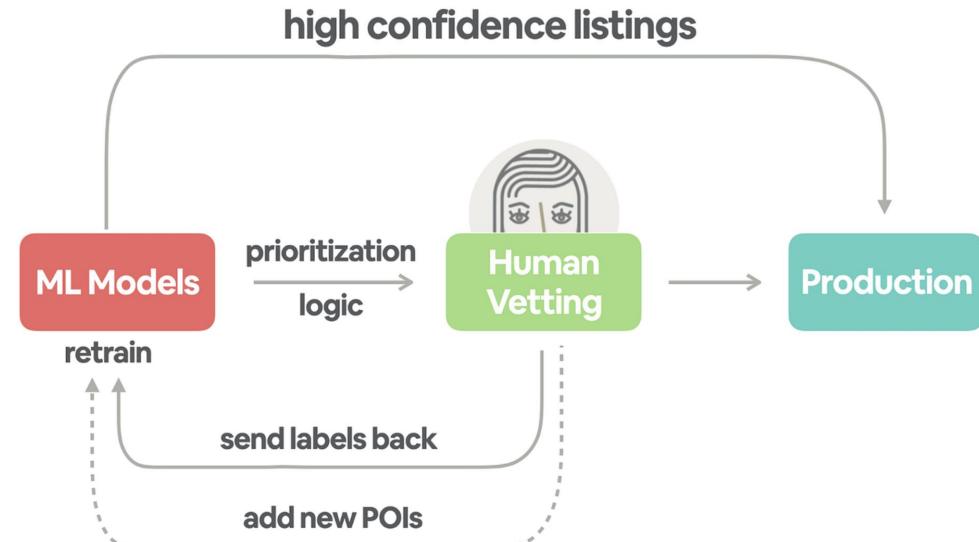
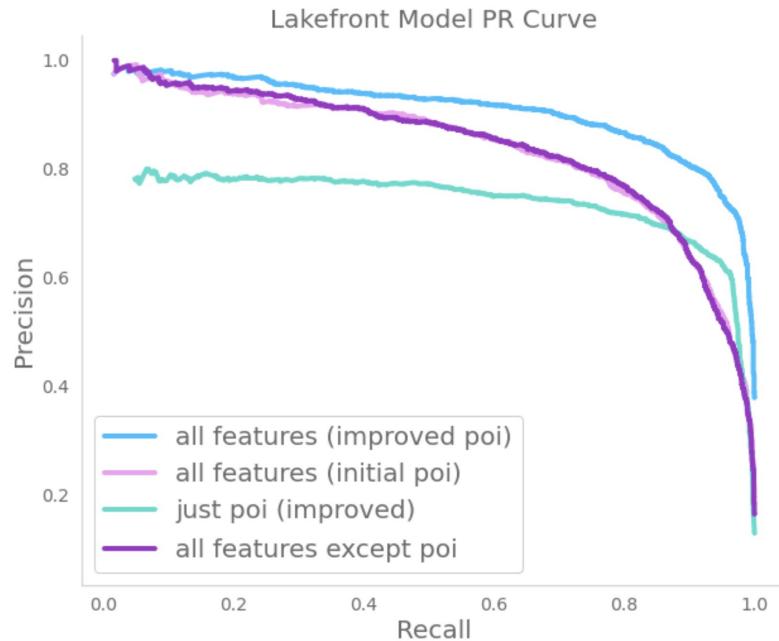
65 miles away

# Case Study: Airbnb

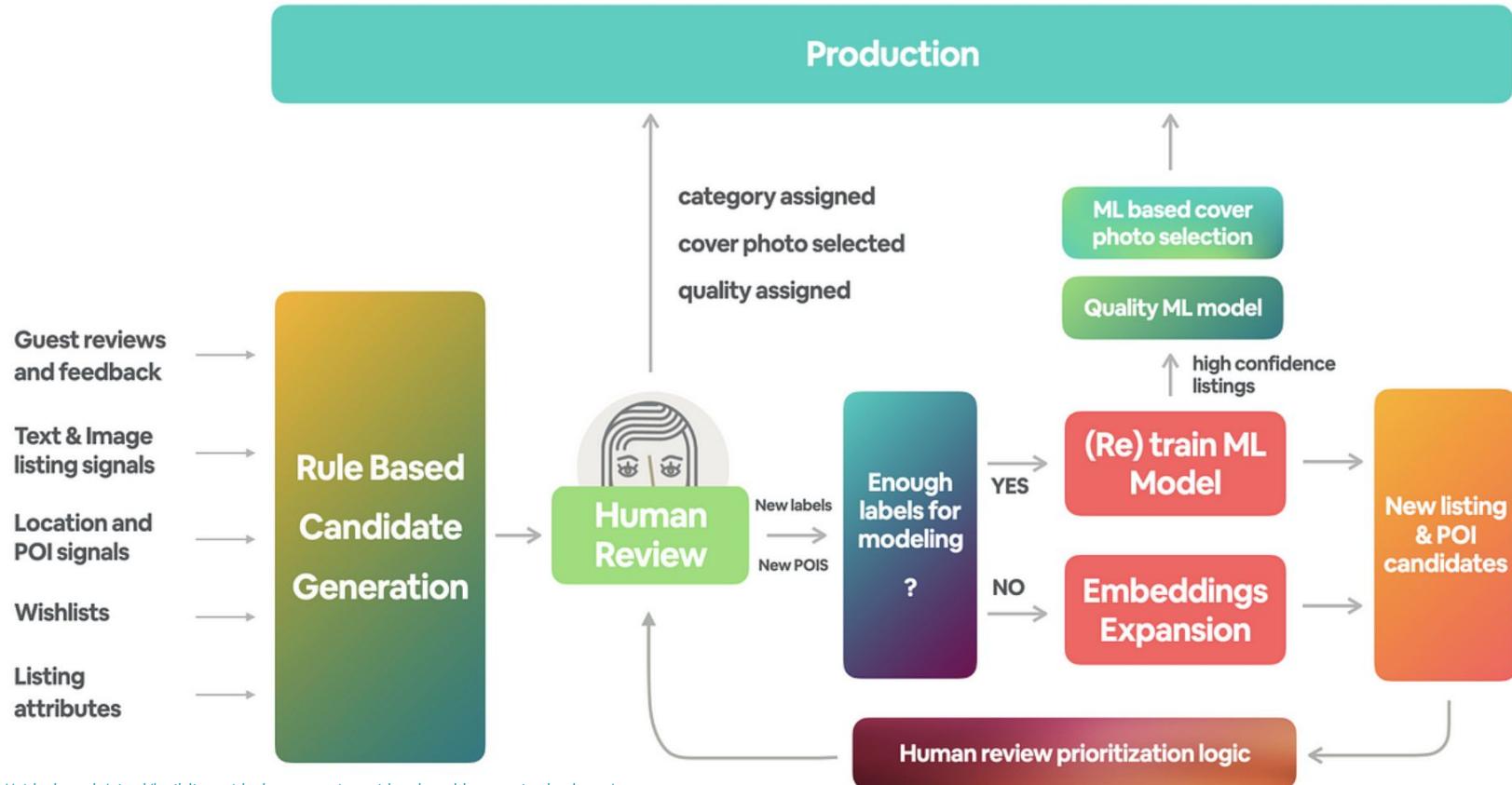


*"With the Summer launch just a few months away, we knew that we could not manually curate all the categories, as it would be very time consuming and costly. We also knew that we could not generate all the categories in a rule-based manner, as this approach would not be accurate enough. Finally, we knew we could not produce an accurate ML categorization model without a training set of human-generated labels. Given all of these limitations, we decided to combine the accuracy of human review with the scale of ML models to create a human-in-the-loop system for listing categorization and display."*

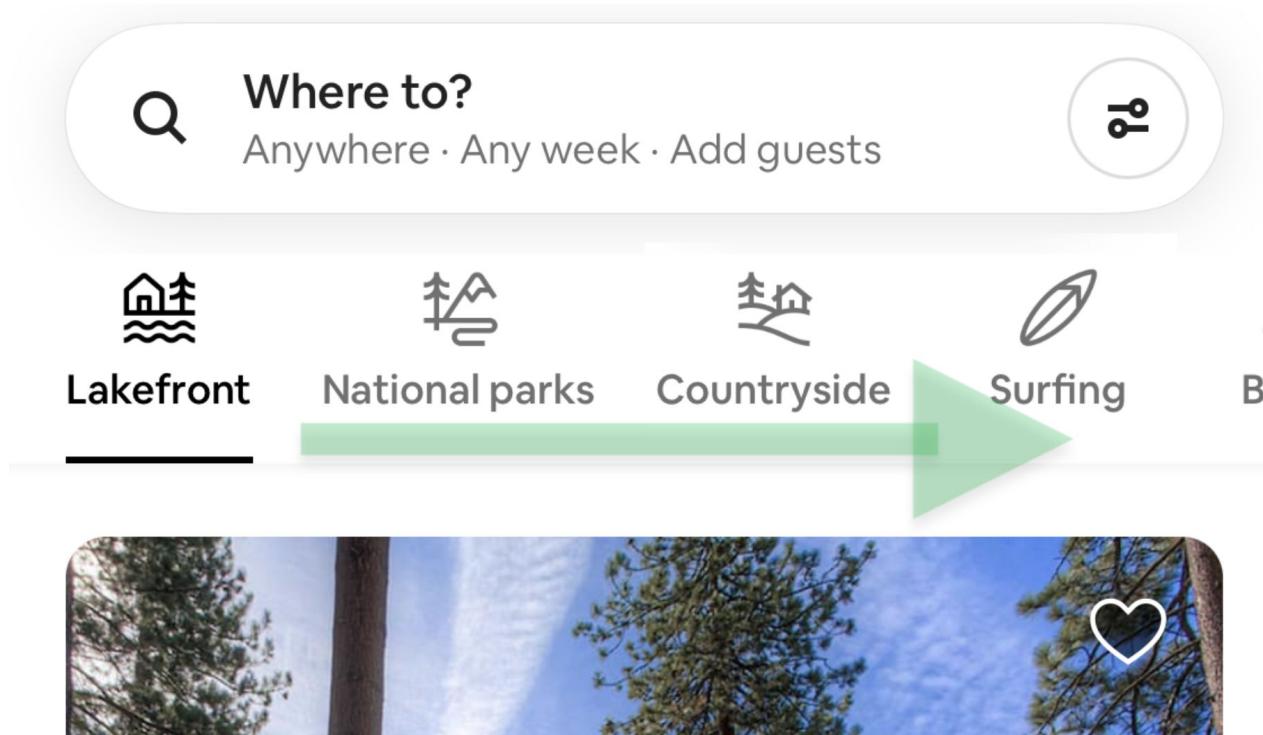
# Case Study: Airbnb



# Case Study: AirBnb



# Case Study: Airbnb

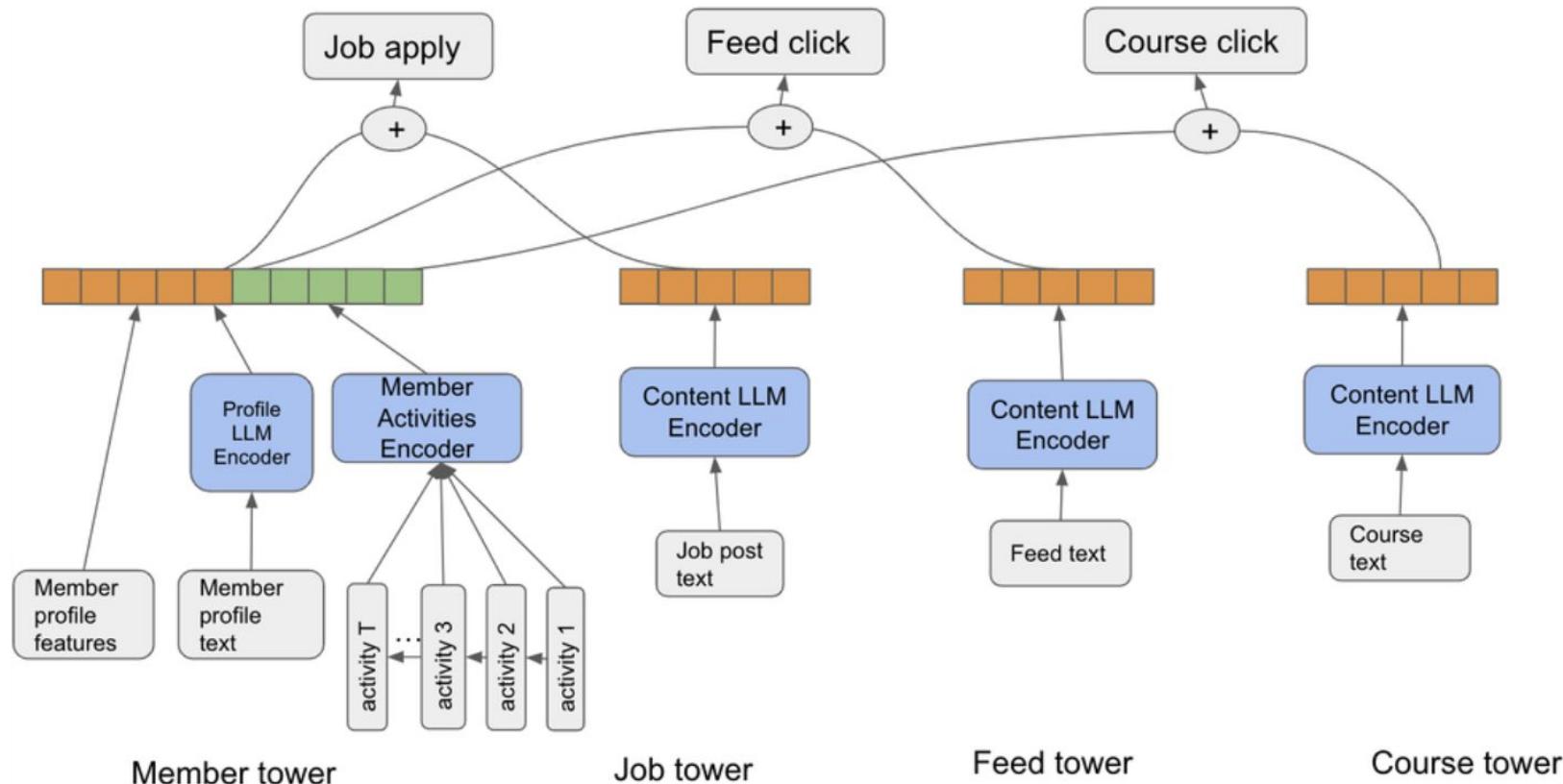


# Case Study: LinkedIn

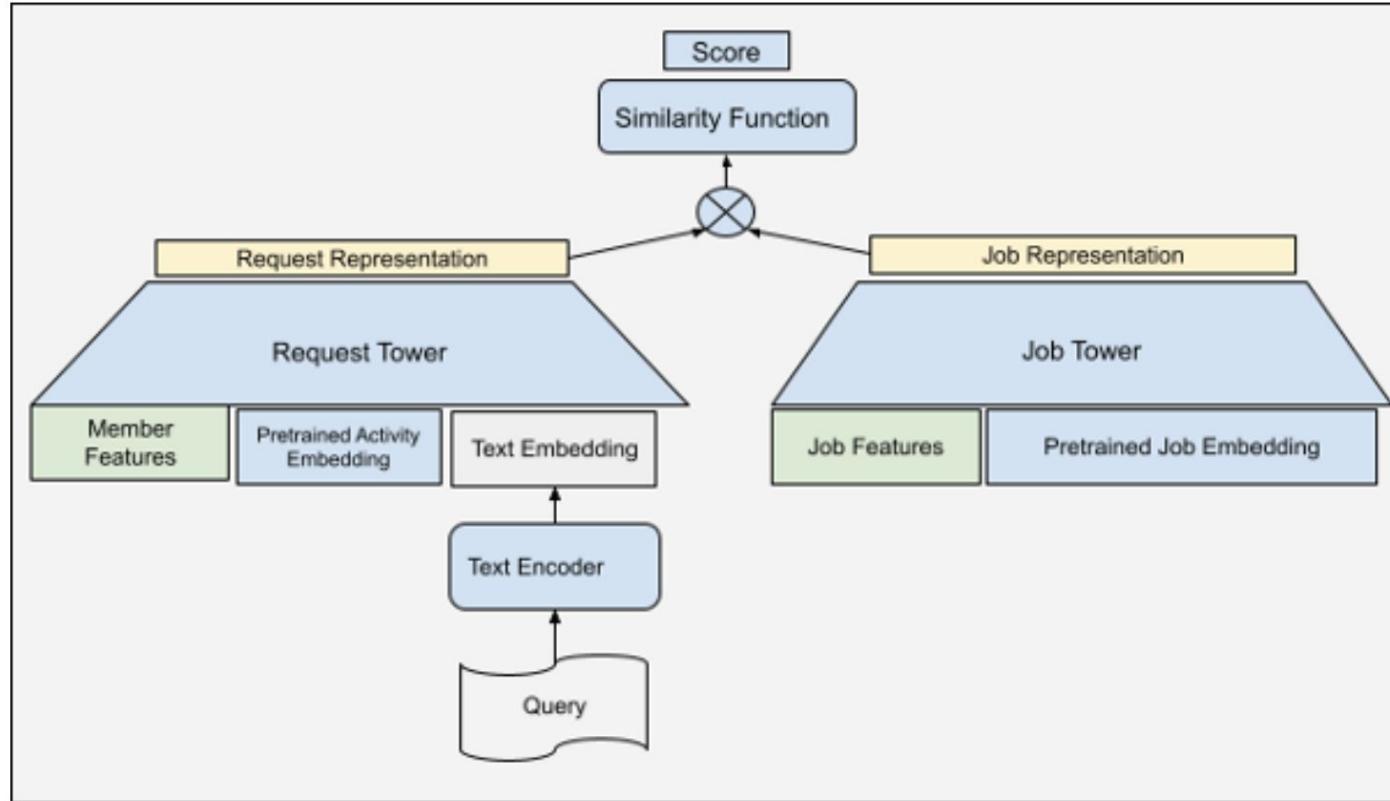
## Problem

- LinkedIn wants to recommend Jobs You Might Be Interested In (JYMBII), but discerning which these are is a challenge

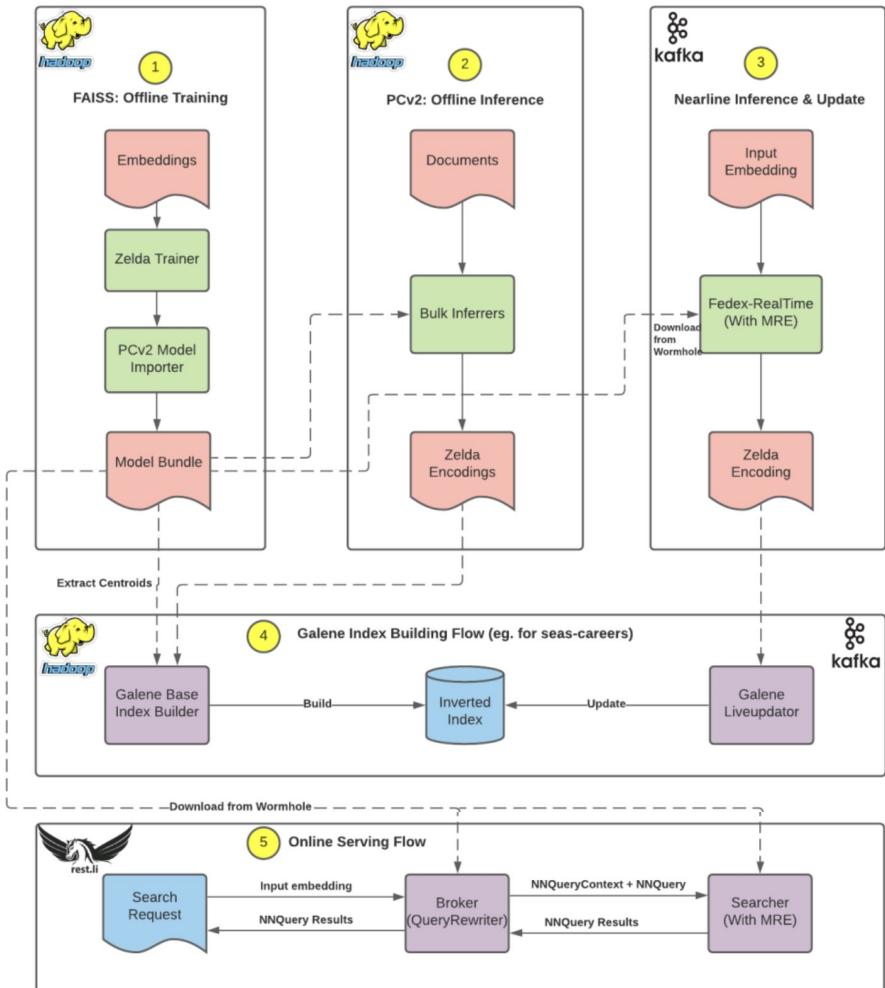
# Case Study: LinkedIn



# Case Study: LinkedIn



# Case Study: LinkedIn



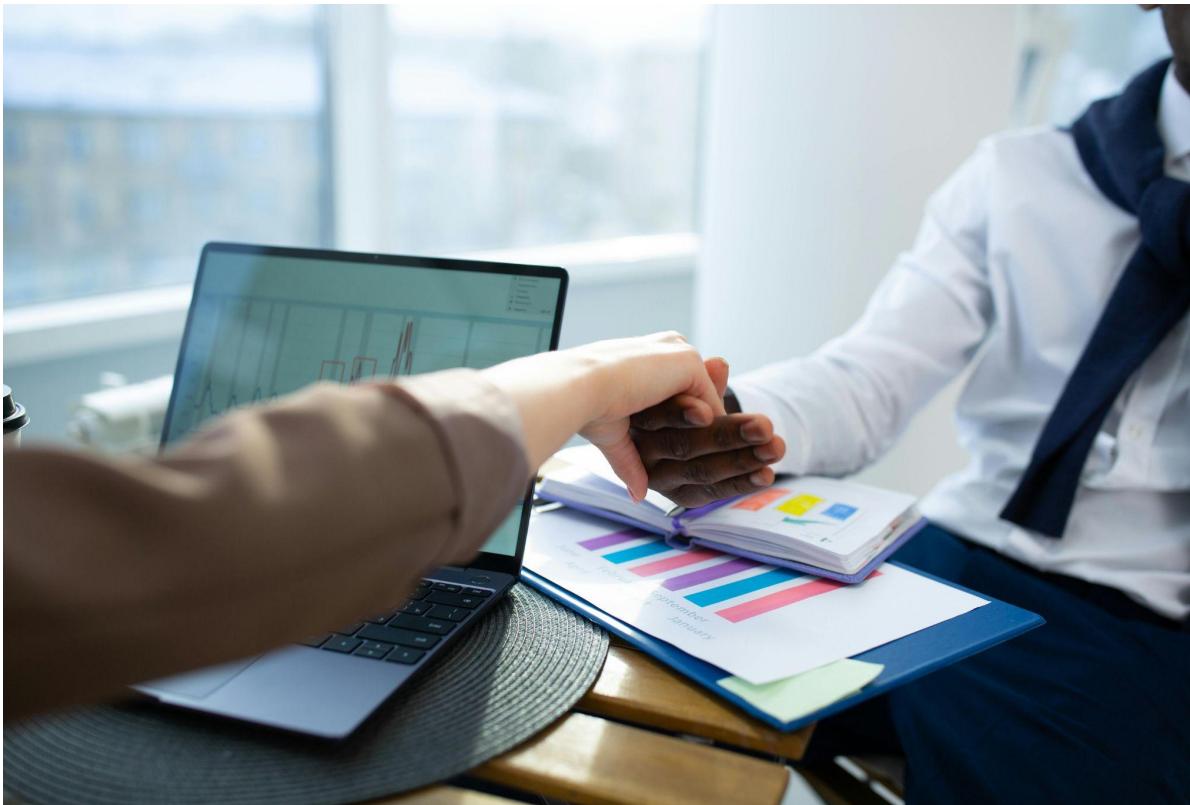
# What is the right AI to use?

Transformers

# Loan origination



# Grant review



# Ask an AI model

Transformers

# Getting AI to think with you, not for you



**Stay mindful of bias**



**Question solutions**



**Expand your perspectives**



**Spot weaknesses or  
vulnerabilities**



**Start with your own  
ideas**

# Assignment walkthrough

Transformers