

EP 44: How does ChatGPT work?



ALEX XU

FEB 4, 2023



181



1



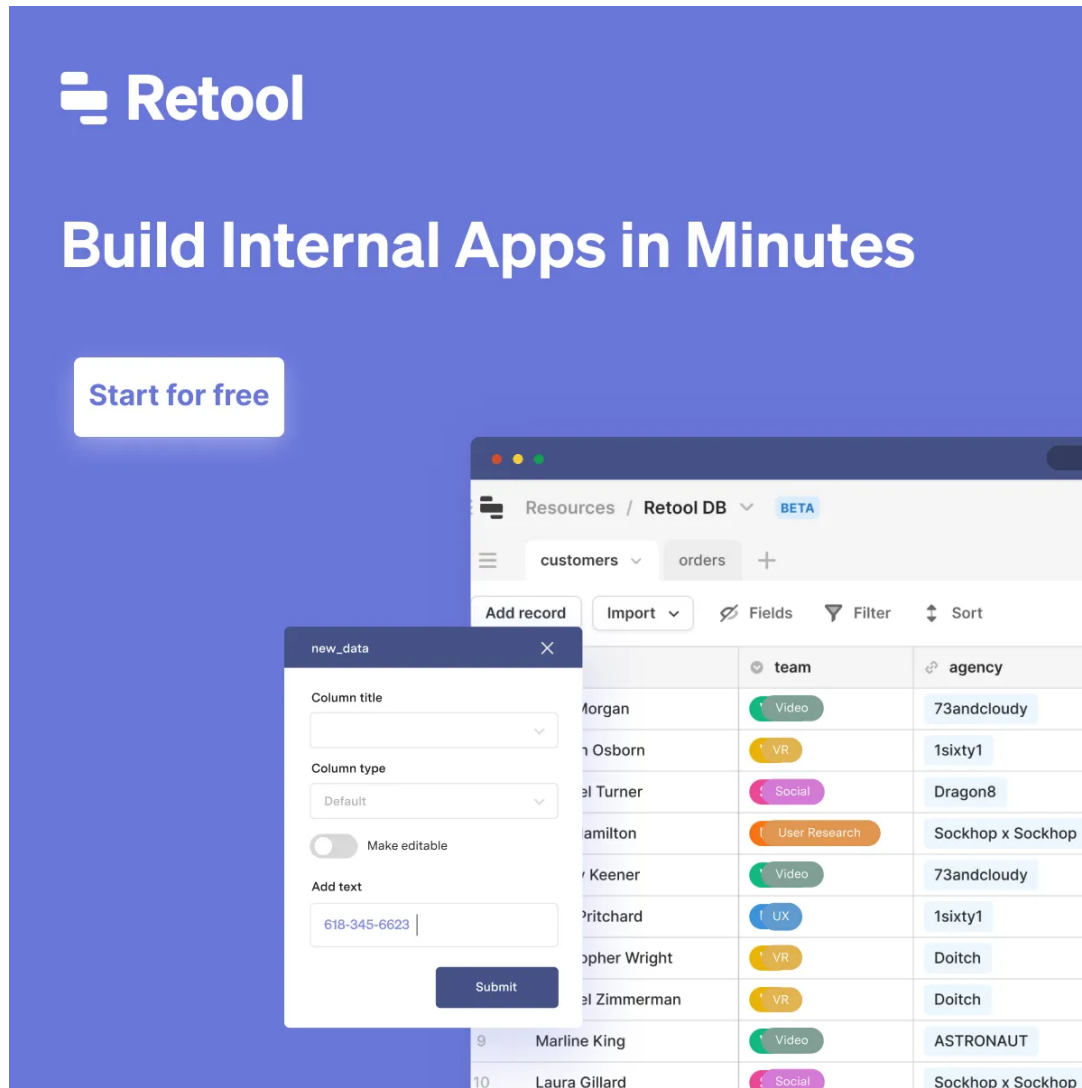
Share



This week's system design refresher:

- 8 Key Data Structures That Power Modern Databases
- How does ChatGPT work?
- Does the cloud really save costs?
- Amazon's system architecture (1998 edition)
- New Machine Learning System Design Interview Book by ByteByteGo

Retool is the fast way to build internal tools (Sponsored)



Building business software today is slow. You often spend more time on boilerplate code and redundant work than you do on actually solving the problem at hand.

Retool is a new approach. We move the starting line with a platform that makes it much faster to connect to any data source, design and develop at the same time, and deploy software securely.

Companies like Amazon and Plaid use Retool to build apps and workflows that help teams work faster. Retool is free for teams of up to 5, and early-stage startups can get \$25,000 in free credits for paid plans.

8 Key Data Structures That Power Modern Databases

8 Key Data Structures That Power Modern Databases



Thanks for reading ByteByteGo Newsletter!
Subscribe for free to receive new posts and
support my work.

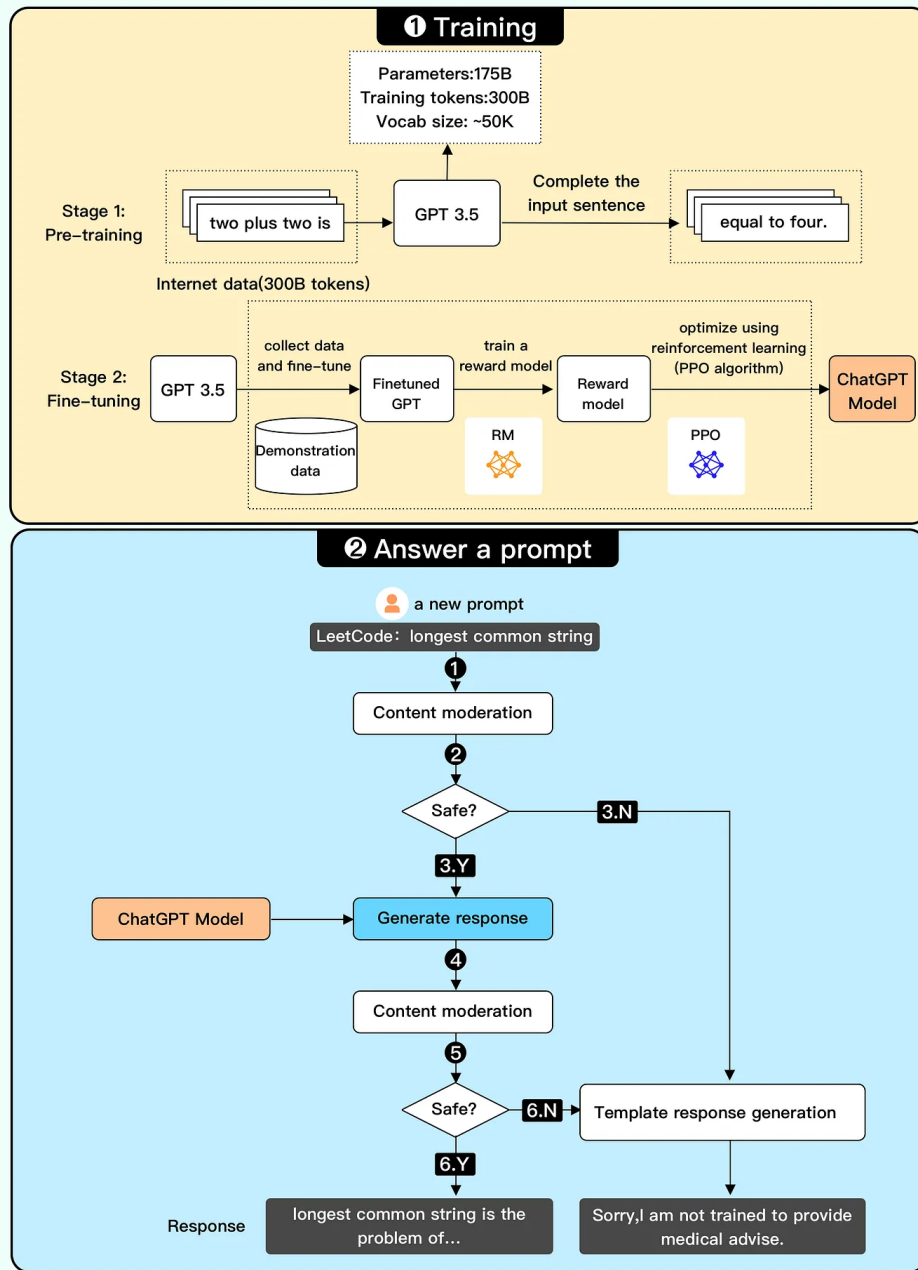
How does ChatGPT work?

Since OpenAI hasn't provided all the details, some parts of the diagram may be inaccurate.

We attempted to explain how it works in the diagram below. The process can be broken down into two parts.

How does ChatGPT-like System Work?

ByteByteGo.com



1. Training. To train a ChatGPT model, there are two stages:

- **Pre-training**: In this stage, we train a GPT model (decoder-only transformer) on a large chunk of internet data. The objective is to train a model that can predict future words given a sentence in a way that is grammatically correct and semantically meaningful similar to the internet data. After the pre-training stage, the model can complete given sentences, but it is not capable of responding to questions.
- **Fine-tuning**: This stage is a 3-step process that turns the pre-trained model into a question-answering ChatGPT model:

- 1). Collect training data (questions and answers), and fine-tune the pre-trained model on this data. The model takes a question as input and learns to generate an answer similar to the training data.
- 2). Collect more data (question, several answers) and train a reward model to rank these answers from most relevant to least relevant.
- 3). Use reinforcement learning (PPO optimization) to fine-tune the model so the model's answers are more accurate.

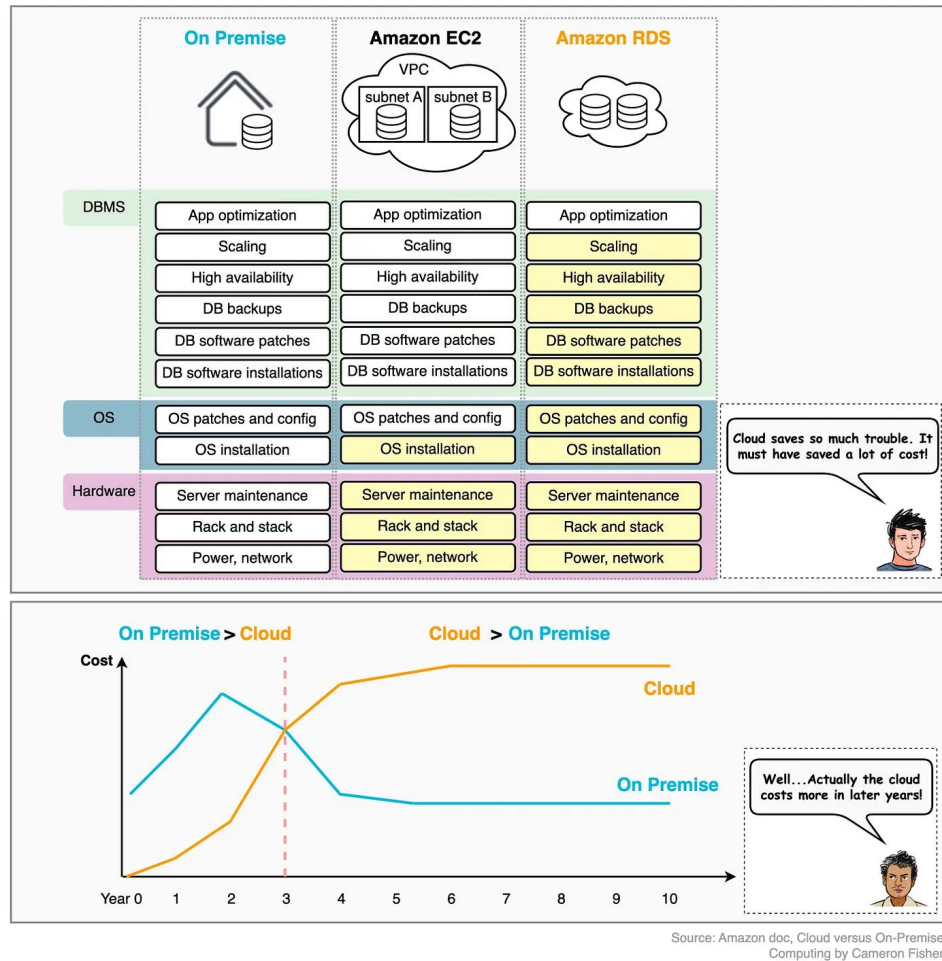
2. Answer a prompt

- Step 1: The user enters the full question, “Explain how a classification algorithm works”.
- Step 2: The question is sent to a content moderation component. This component ensures that the question does not violate safety guidelines and filters inappropriate questions.
- Steps 3-4: If the input passes content moderation, it is sent to the chatGPT model. If the input doesn't pass content moderation, it goes straight to template response generation.
- Step 5-6: Once the model generates the response, it is sent to a content moderation component again. This ensures the generated response is safe, harmless, unbiased, etc.
- Step 7: If the input passes content moderation, it is shown to the user. If the input doesn't pass content moderation, it goes to template response generation and shows a template answer to the user.

Does the cloud really save costs?

Let's look at this question in a longer time range to see what the cloud really brings us.

Does the Cloud Save Costs?



When a company or a business line initially starts, product-market fit (PMF) is key. The cloud enables quick setup to run the system with minimal necessary hardware. The cost is also transparent.

For example, if we run the databases on-premise, we need to take care of hardware setup, operating system installation, DBMS maintenance, etc. But if we use Amazon RDS (Relational Database Service), we just need to take care of application optimization. This saves us the trouble of hiring Linux admins and DB admins.

Later, if the business model doesn't work, we can just stop using the services to save costs without thinking about how to deal with the hardware.

In research conducted by Cameron Fisher, the cloud starts from **almost zero cost**. Over time, the cost starts to accumulate on subscriptions and deployment consulting. Ironically, because it is so easy to allocate services to the cloud for scalability or reliability reasons, an organization tends to **overuse** the cloud after adopting the

cloud. It is essential to set up a monitoring framework for cost transparency.

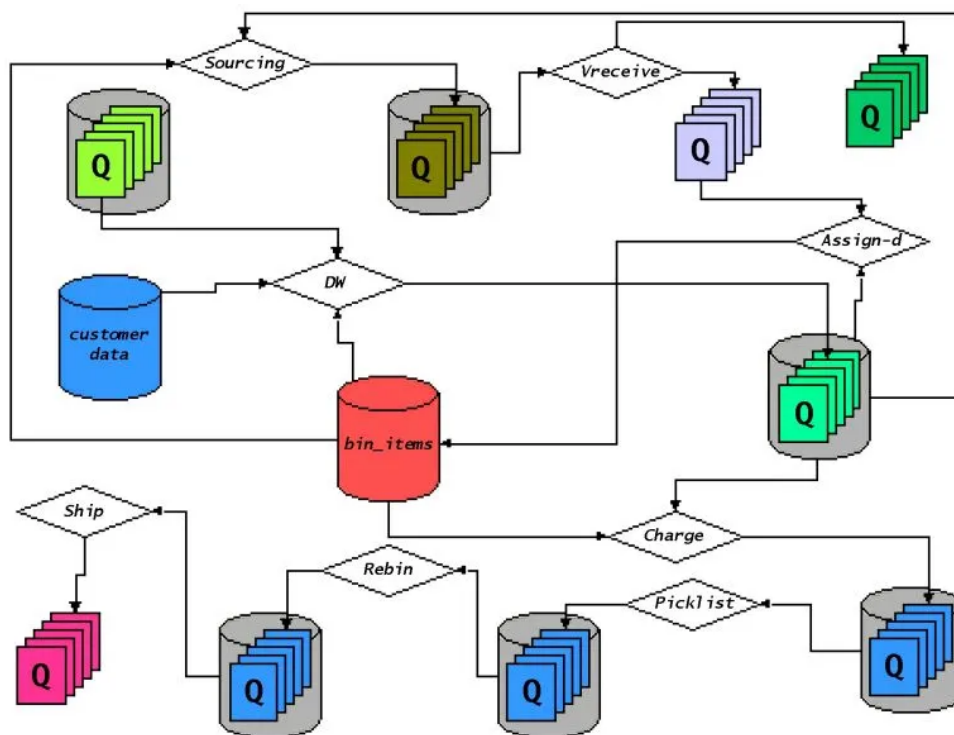
👉 Over to you: Which notable companies use on-premise solutions and why?

Reference:

1. AWS guide: Choosing between Amazon EC2 and Amazon RDS
2. Cloud versus On-Premise Computing by Cameron Fisher, MIT

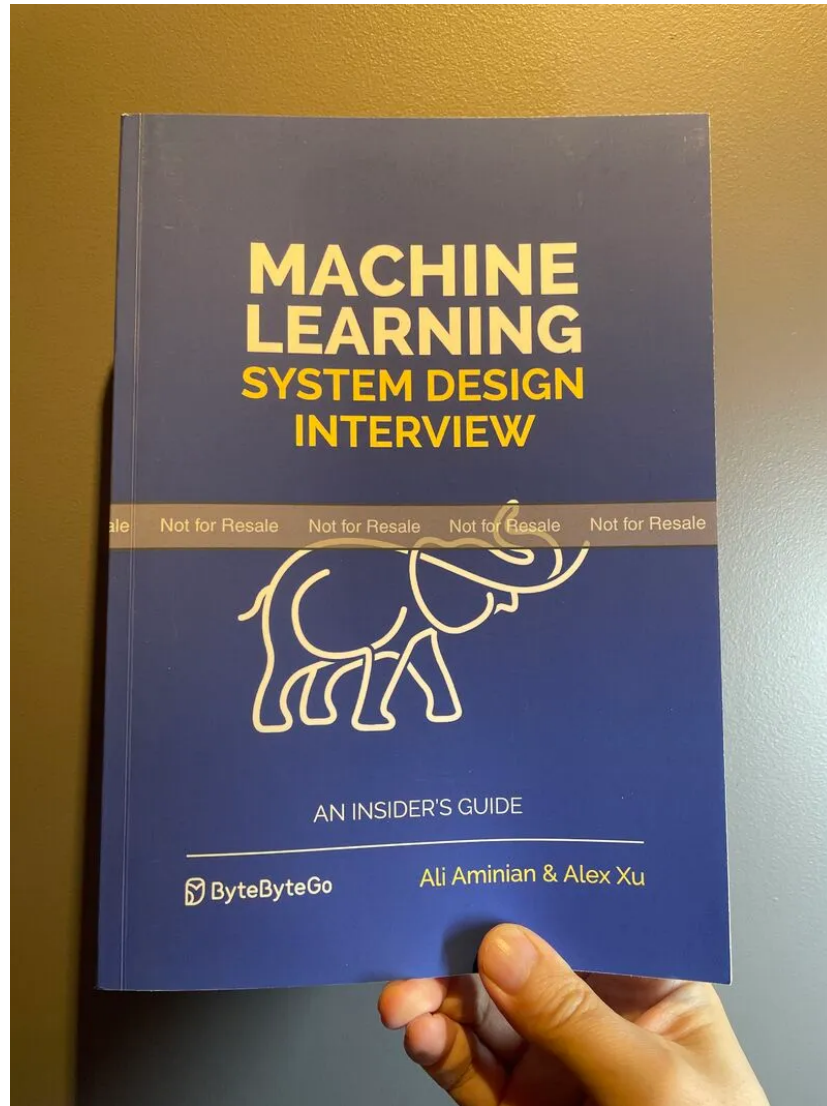
Amazon's system architecture

In 1998, Amazon's system architecture looked like this. The simplicity of the architecture is amazing.



You can read the 25-year-old internal document that changed Amazon's system design and development here: <https://lnkd.in/e5EGHFiu>

New Machine Learning System Design Interview Book



Some stats about the book:

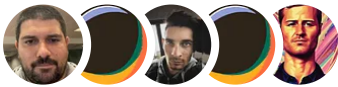
- 10 real machine learning system design interview questions with detailed solutions.
- 211 diagrams to explain how different ML systems work.
- 300+ pages.

Feels so good to hold it in my hand. Thanks to everyone who helped us make this happen.

Paperback version of the book: <https://geni.us/tVsKGey>

Digital version of the book: <https://bytebytego.com>

Thanks for reading ByteByteGo Newsletter!
Subscribe for free to receive new posts and
support my work.



181 Likes

1 Comment



Write a comment...



Boyd Wold Mar 18

I asked chatGPT if the explanation given here was accurate and comprehensive. It explained the corrections and then updated the article. <sigh>

Your explanation provides a good overview of the ChatGPT training process and how it answers a prompt, particularly the emphasis on safety checks and content moderation. However, I would like to suggest a few modifications and clarifications to make it more comprehensive:

Training:

Pre-training: Specify that the model learns to predict the next word in a sequence, given the context of previous words.

Fine-tuning:

For the first step, mention that the training data can include demonstrations and comparisons.

For the second step, clarify that the reward model is trained using human preferences, where human evaluators rank different model-generated responses based on relevance and quality.

For the third step, specify that the model is fine-tuned using Proximal Policy Optimization (PPO).

Answer a prompt:

Consider renumbering the steps for better readability and understanding.

[Expand full comment](#)

 LIKE (3)  REPLY  SHARE

...

© 2023 ByteByteGo · [Privacy](#) · [Terms](#) · [Collection notice](#)
[Substack](#) is the home for great writing