



DeepLearning.AI

Software Applications

Using generative AI in
software applications

Examples of software applications

Writing

Is there parking for staff?

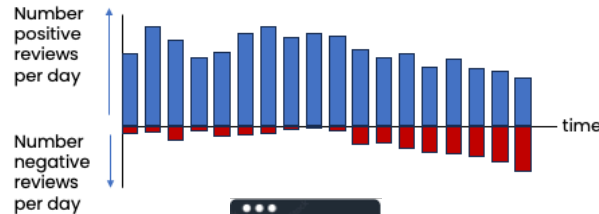
Yes, employees can park on levels 1 and 2 of the office lot. You can get a sticker at [...]



Reading

Classify the sentiment of the following review as either positive or negative:

The mochi is excellent!
Positive



Chatting

Welcome to BettaBurgers!

I'd like a cheeseburger for delivery!

Awesome. Anything else?

That's it thanks!

Ok, it'll be there in 20 minutes. Enjoy!



```

1 import torch
2 from torch.utils.data import DataLoader, TensorDataset
3 from torch import nn
4
5 class SentimentLSTM(nn.Module):
6     def __init__(self, vocab_size, output_size, embedding_dim, hidden_dim, n_layers):
7         super(SentimentLSTM, self).__init__()
8         self.output_size = output_size
9         self.n_layers = n_layers
10        self.hidden_dim = hidden_dim
11
12        self.embedding = nn.Embedding(vocab_size, embedding_dim)
13        self.lstm = nn.LSTM(embedding_dim, hidden_dim, n_layers, batch_first=True)
14        self.fc = nn.Linear(hidden_dim, output_size)
15        self.sigmoid = nn.Sigmoid()
16
17    def forward(self, x, hidden):
18        batch_size = x.size(0)
19        x = self.embedding(x)
20        lstm_out, hidden = self.lstm(x, hidden)
21        lstm_out = lstm_out.contiguous().view(-1, self.hidden_dim)
22        out = self.fc(lstm_out)
23        out = self.sigmoid(out)
24        out = out.view(batch_size, -1)
25        out = out[:, -1]
26        return out, hidden
27
28    def init_hidden(self, batch_size):
29        weight = next(self.parameters()).data
30        hidden = (weight.new(self.n_layers, batch_size, self.hidden_dim).zero_(),

```

Supervised learning for restaurant reputation monitoring

Input A → Output B

Input (A)	Output (B)
Best soup dumplings I've ever eaten.	Positive
The colorful tablecloths made me smile!	Positive
Not worth the 3 month wait for a reservation	Negative



Input: "Best bubble tea I've ever had"

Output: Positive

Prompt-based development

```
prompt = """  
    Classify the following review  
    as having either a positive or  
    negative sentiment:  
  
    The banana pudding was really  
    tasty!  
    """>  
  
response = llm_response(prompt)  
print(response)
```

Instruction text

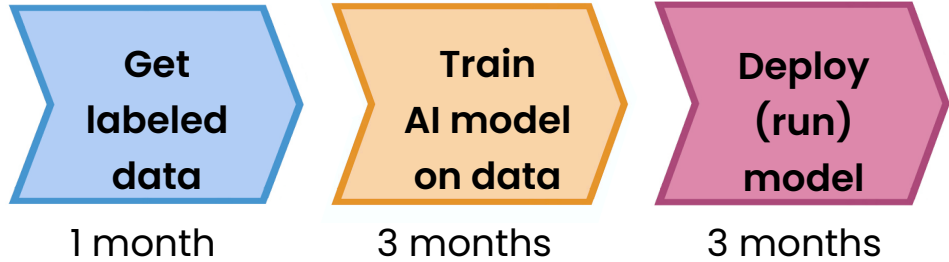
Review text

Code to call LLM

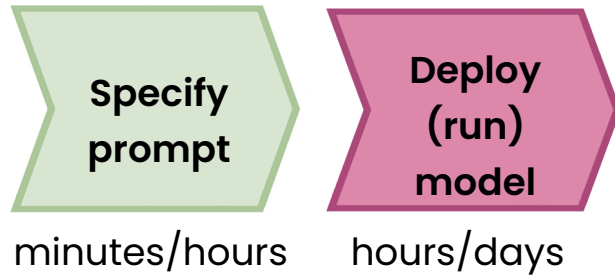
Code to print output

Workflow using Generative AI

Supervised learning



Prompt-based AI





DeepLearning.AI

Software Applications

Trying generative AI code
yourself (optional)

Coding platform (Jupyter notebook)

jupyter l2-guidelines Python 3 (ipykernel) Logout

Not Trusted Python 3 (ipykernel)

File Edit View Insert Cell Kernel Help

Run

Guidelines for Prompting

In this lesson, you'll practice two prompting principles and their related tactics in order to write effective prompts for large language models.

Setup

Load the API key and relevant Python libraries.

In this course, we've provided some code that loads the OpenAI API key for you.

```
In [ ]: import openai
import os

from dotenv import load_dotenv, find_dotenv
_ = load_dotenv(find_dotenv())

openai.api_key = os.getenv('OPENAI_API_KEY')
```

helper function

Throughout this course, we will use OpenAI's gpt-3.5-turbo model

The one
command
to know:

shift

enter

ChatGPT
Prompt Engineering
for Developers

Guidelines
for Prompting

OpenAI DeepLearning.AI

0:00 / 17:19

TRANSCRIPT

NEXT LESSON

Play button

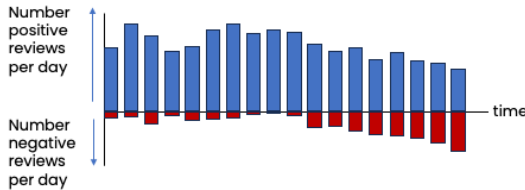
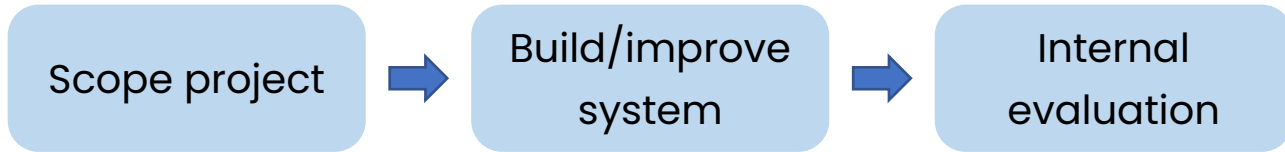


DeepLearning.AI

Software Applications

Lifecycle of a generative
AI project

Lifecycle of a generative AI project



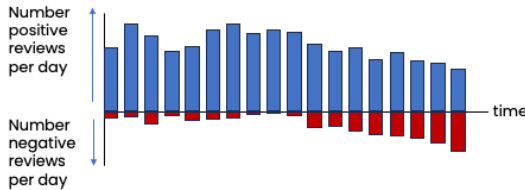
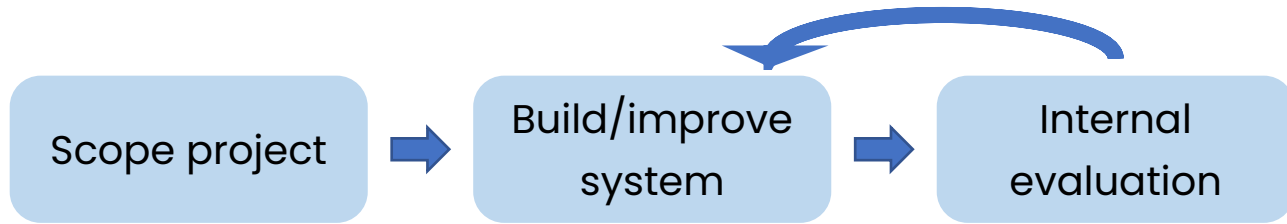
Initially a prototype, that we will improve over time



Classify the sentiment of the following review as either positive or negative:

The custard tart was
Positive

Lifecycle of a generative AI project



Initially a prototype, that we will improve over time

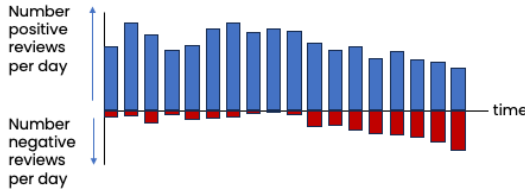
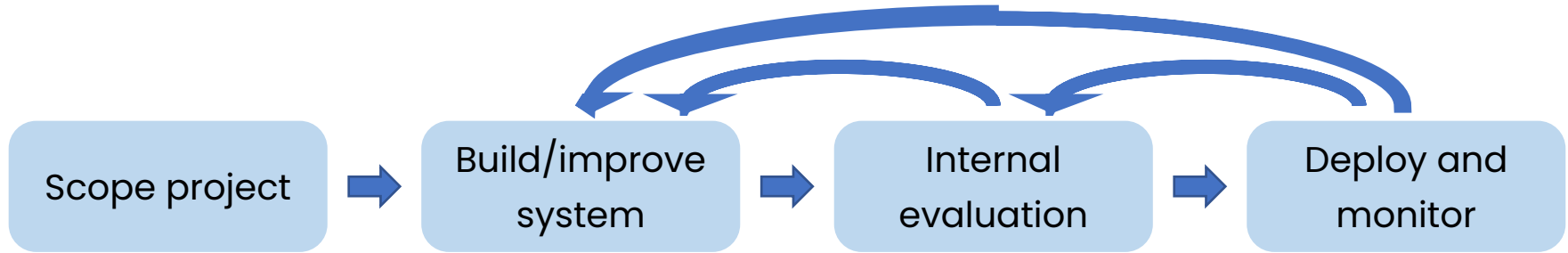


Classify the sentiment of the following review as either positive or negative:

Positive

✘ Incorrect response

Lifecycle of a generative AI project



Initially a prototype, that we will improve over time



Classify the sentiment of the following review as either positive or negative:

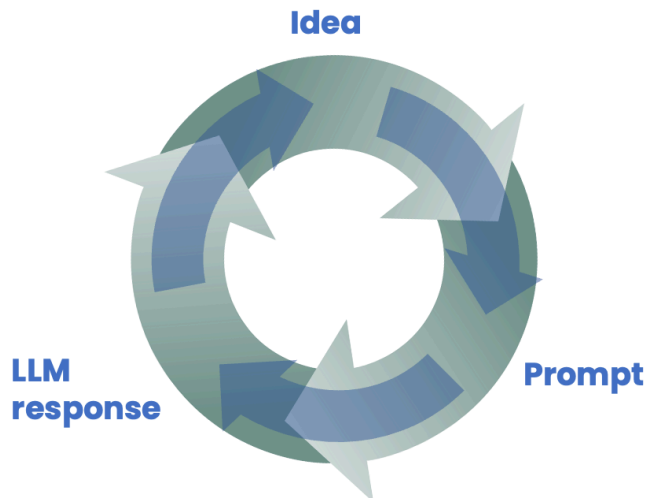
My miso ramen tasted

Positive

Tools to improve performance

Building Generative AI is a highly empirical (experimental) process – we repeatedly find and fix mistakes.

- Prompting



Tools to improve performance

Building Generative AI is a highly empirical (experimental) process – we repeatedly find and fix mistakes.

- Prompting
- Retrieval augmented generation (RAG)



Tools to improve performance

Building Generative AI is a highly empirical (experimental) process – we repeatedly find and fix mistakes.

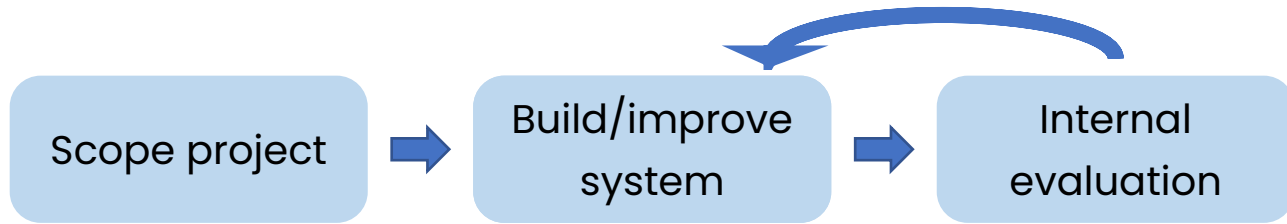
- Prompting
- Retrieval augmented generation (RAG)
 - Give LLM access to external data sources
- Adapt LLM to your task Fine-tune models

Tools to improve performance

Building Generative AI is a highly empirical (experimental) process – we repeatedly find and fix mistakes.

- Prompting
- Retrieval augmented generation (RAG)
 - Give LLM access to external data sources
- Fine-tune models
 - Adapt LLM to your task
- Train LLM from scratch
- Pretrain models

Lifecycle of a generative AI project



Take food
order



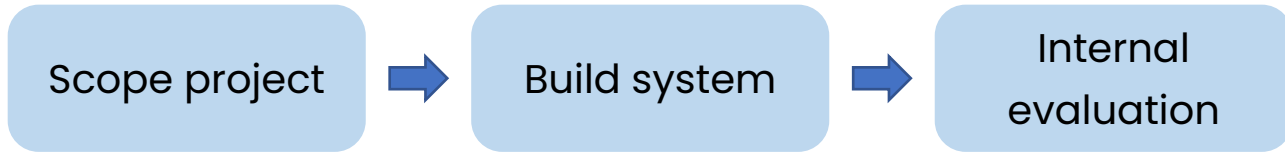
Let team
order...

I'd like a mushroom
and swiss burger

I'm sorry, we don't have
mushrooms, how about
onions?

✘ Incorrect
response

Lifecycle of a generative AI project



Take food
order

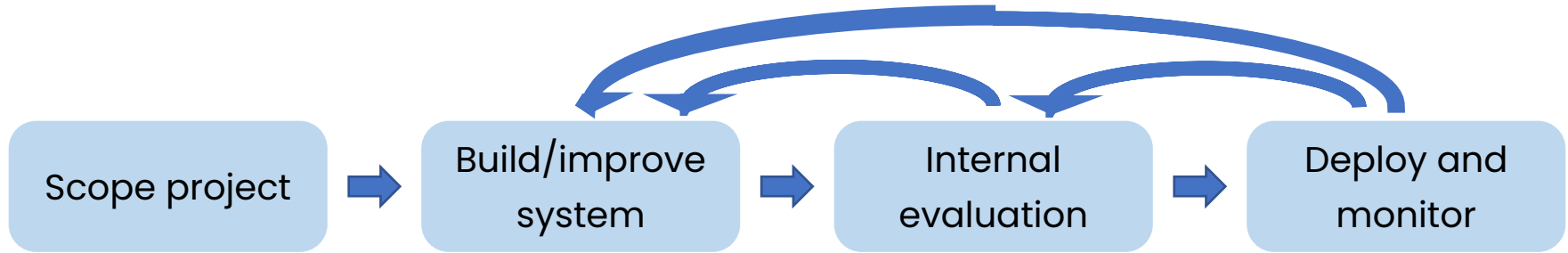


Let team
order...

Do you have pickles?

Yes – would you like to
add some?

Lifecycle of a generative AI project

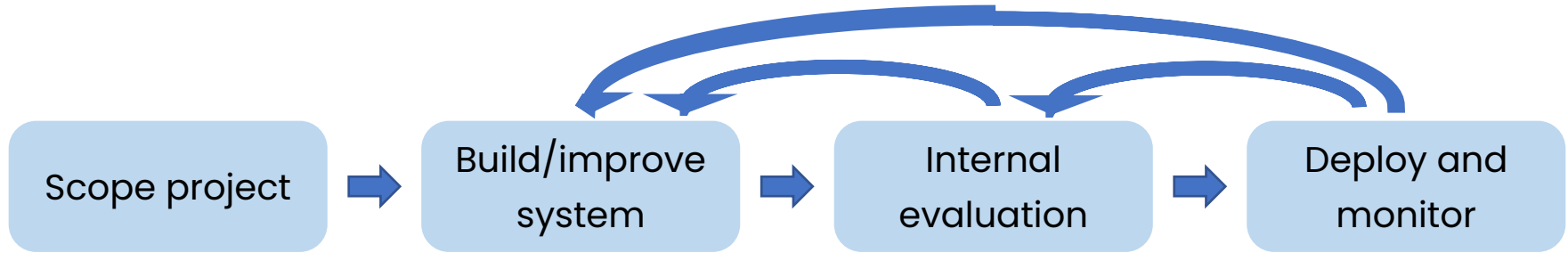


Take food
order



Let customers
order, monitor
LLM responses
How many calories in
the cheeseburger?
I'm sorry, I don't know

Lifecycle of a generative AI project



Take food
order



Let customers
order, monitor
LLM responses
How many calories in
the cheeseburger?
About 375 calories



DeepLearning.AI

Software Applications

Cost intuition

How much does it cost?

Example prices

	OpenAI/GPT3.5	OpenAI/GPT4	Google/PaLM 2	Amazon/Titan Lite
Input tokens	\$0.0015/1K tokens	\$0.03/1K tokens	\$0.0005/1K tokens	\$0.0003/1K tokens
Output tokens	<u>\$0.002/1K tokens</u>	\$0.06/1K tokens	\$0.0005/1K tokens	\$0.0004/1K tokens

What is a token?

the example Andrew 1 token

translate programming 2 tokens

tonkotsu 4 tokens

300 words
400 tokens

Roughly, 1 token = 3/4 words

Estimating Cost

Typical adult reading speed: 250 words/minute.

How much would it cost to keep someone occupied for 1 hour?

$$60 \times 250 = 15,000 \leftarrow \text{output}$$

$$+ 15,000 \leftarrow \text{prompt}$$

$$\hline 30,000$$

↳ 40,000 tokens

$$1 \text{ token} \approx \frac{3}{4} \text{ word}$$

Cost: \$0.002 cents / 1K tokens.

$$\$0.002 \times 40 = \$0.08$$



DeepLearning.AI

Advanced Technologies: Beyond Prompting

Retrieval Augmented Generation (RAG)

Retrieval Augmented Generation (RAG) example

General Chatbot

Is there parking for employees?

I need more specific information about your workplace to answer that question.

Chatbot with RAG

Is there parking for employees?

Yes, employees can park on levels 1 and 2 of the office lot. You can get a sticker at [...]

Retrieval Augmented Generation (RAG) example

1. Given question, search relevant documents for answer

Is there parking for employees?

Company documents



Benefits



Leave
policy



Facilities



Payroll

2. Incorporate retrieved text into an updated prompt

Use the following pieces of context to answer the question at the end.

Parking Policy: All employees may park on levels 1 and 2 of the office lot. Use the entrance on Front St [...]

Is there parking for employees?

Retrieval Augmented Generation (RAG) example

3. Generate answer from the new prompt with additional context

Use the following pieces of context to answer the question at the end.

Parking Policy: All employees may park on levels 1 and 2 of the office lot. Use the entrance on Front St [...]

Yes, employees can park on [...] ^①

→ Link to full onboarding policy

Examples of RAG applications

Chat with PDF files



Accenture-Su... 2 / 4

New tech, new jobs, new value: How generative AI will reinvent supply chains

95%⁽¹⁾ of global executives agree advances in generative AI signify a new era of enterprise intelligence.

The message for CISOs and CPOs? It's not a question of whether generative AI will impact your supply chain, but by how much.

Showcased in the [Accenture Technology Vision 2023](#), four technology trends are driving the next phase of business transformation. One of them, generative AI, is among the biggest leaps forward in AI history. Generative AI's unique power? Rather than analyzing existing information, which AI has been used for until now, it creates new content - from writing code to making movies.

As humans working with generative AI "virtual colleagues" become the norm, every role in every supply chain has the potential to be transformed. Companies that adopt a strong digital core in the way you drive innovation and help accelerate a new performance frontier as part of a broader strategy of Total Enterprise Transformation.

"81% of C-suite executives are planning to explore generative AI technologies in 2023 as part of a broader strategy"

Accenture Technology Vision 2023 research

Supply chain roles will change dramatically...

Generative AI will have a profound impact, supporting people right across the supply chain.

Advising: By putting new kinds of hyper-personalized intelligence into human hands, Generative AI will support many different supply-chain tasks. These range from advising on vendor selection with specific field parameters, to providing step-by-step guidance for asset maintenance processes, or even assisting with real-time queries on supply chain KPIs and alerts.

Supply chain roles potentially impacted: Category managers, supply chain planners, field engineers, plant workers.

Creating: Generative AI will become an essential creative partner for people, revealing new ways to reach and appeal to target audiences and bringing unprecedented speed and innovation in supply chain areas like product design - from solving for sustainable product packaging to creating new materials/parts and personalized products right at the time of sale.

Supply chain roles potentially impacted: Product engineers, retailers, materials planners.

Coding: By enabling rapid conversion of one programming language to another, generative AI will dramatically boost the productivity and speed of professionals working to optimize supply chain processes and onboard new supply chain players; generative AI will also democratize access to AI and analytics for professionals without niche expertise.

Supply chain roles potentially impacted: Data scientists working in supply chain operations.

Automating: Generative AI will transform business process automation. Hyper-efficiency and hyper-personalization will be the norm in back- and front-office operations. For supply chains, this means new customer service interactions with next-gen chatbots advising on complicated resolution issues, drafting new contracts - and even automating alerts and responses to suppliers.

Supply chain roles potentially impacted: Call-center staff, planners, buyers.

Protecting: In time, Generative AI will support enterprise governance and information security, protecting against fraud, improving regulatory compliance for sustainability and responsibility, and proactively identifying risk - all mission-critical supply chain responsibilities.

Supply chain roles potentially impacted: Sustainability and compliance officers, risk managers.

(1) Technology Vision 2023: www.accenture.com/TechVision2023

Chat



What are some roles in supply chain that will be changed by AI?

Generative AI will have a profound impact on supply chain roles. According to 2, some of the roles that will be potentially impacted include data scientists, call-center staff, planners, buyers, sustainability and compliance officers, risk managers, category managers, supply chain planners, field engineers, plant workers, product engineers, retailers, and materials planners. 2

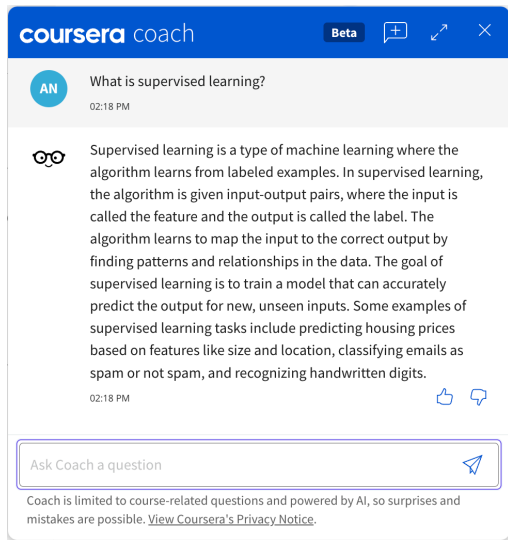
Ask any question...



Chatting with an Accenture report using ChatPDF

Examples of RAG applications

Answer questions based on a website's articles



The screenshot shows the Coursera Coach chat interface. At the top, it says "coursera coach" and "Beta". A user asks, "What is supervised learning?". The response explains supervised learning as a type of machine learning where an algorithm learns from labeled examples. It details how the algorithm maps input-output pairs to find patterns and relationships in data. The goal is to train a model that can accurately predict the output for new, unseen inputs. Examples of supervised learning tasks include predicting housing prices, classifying emails as spam or not spam, and recognizing handwritten digits. There is a "Verified Sources" section at the bottom with three links.

AN What is supervised learning?
02:18 PM

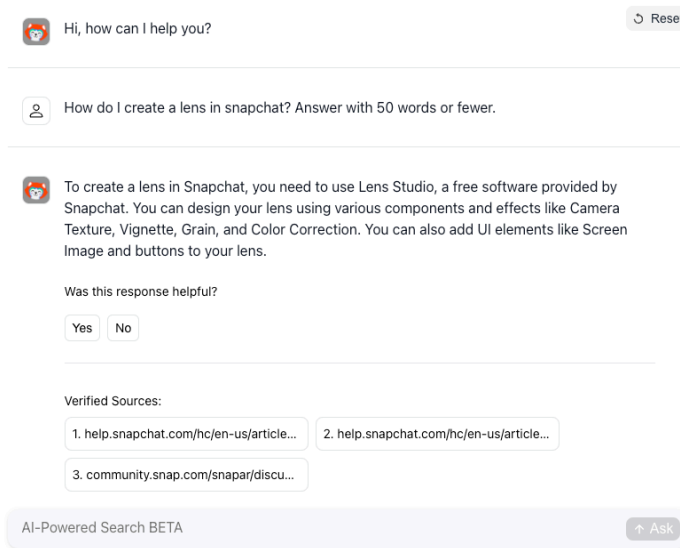
Supervised learning is a type of machine learning where the algorithm learns from labeled examples. In supervised learning, the algorithm is given input-output pairs, where the input is called the feature and the output is called the label. The algorithm learns to map the input to the correct output by finding patterns and relationships in the data. The goal of supervised learning is to train a model that can accurately predict the output for new, unseen inputs. Some examples of supervised learning tasks include predicting housing prices based on features like size and location, classifying emails as spam or not spam, and recognizing handwritten digits.

02:18 PM

Ask Coach a question

Coach is limited to course-related questions and powered by AI, so surprises and mistakes are possible. [View Coursera's Privacy Notice.](#)

Coursera Coach



The screenshot shows the Snapchat AI chat interface. A user asks, "Hi, how can I help you?". The user then asks, "How do I create a lens in snapchat? Answer with 50 words or fewer.". The response explains that to create a lens in Snapchat, you need to use Lens Studio, a free software provided by Snapchat. It lists various components and effects like Camera Texture, Vignette, Grain, and Color Correction, and mentions that you can also add UI elements like Screen Image and buttons to your lens. There is a "Verified Sources" section with three links.

Hi, how can I help you?

How do I create a lens in snapchat? Answer with 50 words or fewer.

To create a lens in Snapchat, you need to use Lens Studio, a free software provided by Snapchat. You can design your lens using various components and effects like Camera Texture, Vignette, Grain, and Color Correction. You can also add UI elements like Screen Image and buttons to your lens.

Was this response helpful?

Yes No

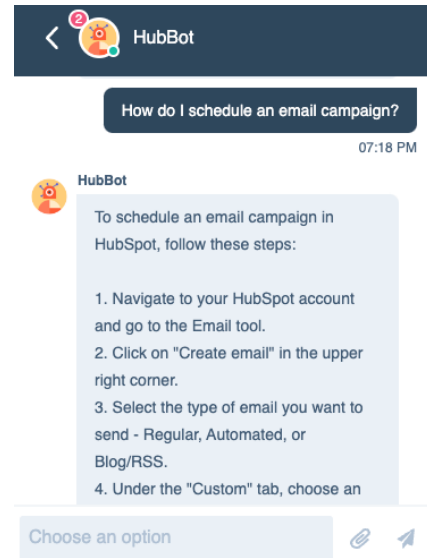
Verified Sources:

- help.snapchat.com/hc/en-us/article...
- help.snapchat.com/hc/en-us/article...
- community.snap.com/snapar/discu...

AI-Powered Search BETA

Ask

Snapchat



The screenshot shows the HubSpot chat interface. A user asks, "How do I schedule an email campaign?". The response provides a list of steps to schedule an email campaign in HubSpot. There is a "Choose an option" button at the bottom.

HubBot

How do I schedule an email campaign?
07:18 PM

HubBot

To schedule an email campaign in HubSpot, follow these steps:

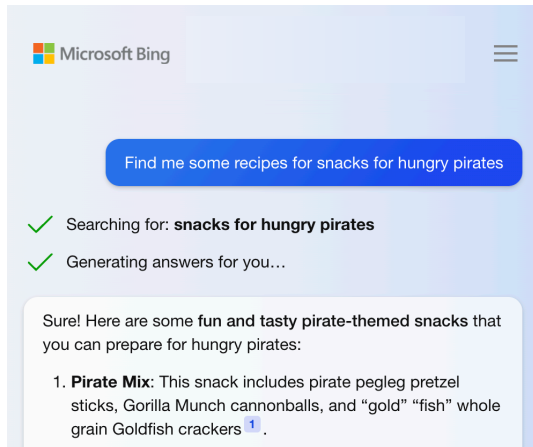
- Navigate to your HubSpot account and go to the Email tool.
- Click on "Create email" in the upper right corner.
- Select the type of email you want to send - Regular, Automated, or Blog/RSS.
- Under the "Custom" tab, choose an

Choose an option

Hubspot

Examples of RAG applications

New form of web search



Microsoft Bing

Find me some recipes for snacks for hungry pirates

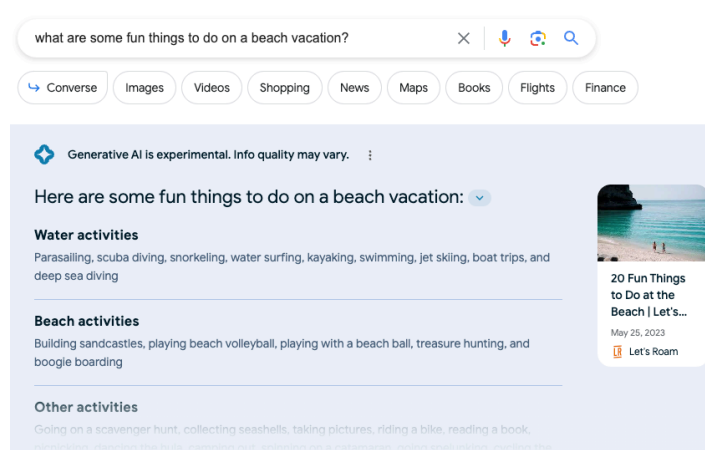
✓ Searching for: **snacks for hungry pirates**

✓ Generating answers for you...

Sure! Here are some fun and tasty pirate-themed snacks that you can prepare for hungry pirates:

1. **Pirate Mix:** This snack includes pirate pegleg pretzel sticks, Gorilla Munch cannonballs, and "gold" "fish" whole grain Goldfish crackers ¹.

Microsoft/Bing Chat



what are some fun things to do on a beach vacation?

Converse Images Videos Shopping News Maps Books Flights Finance

Generative AI is experimental. Info quality may vary.

Here are some fun things to do on a beach vacation:

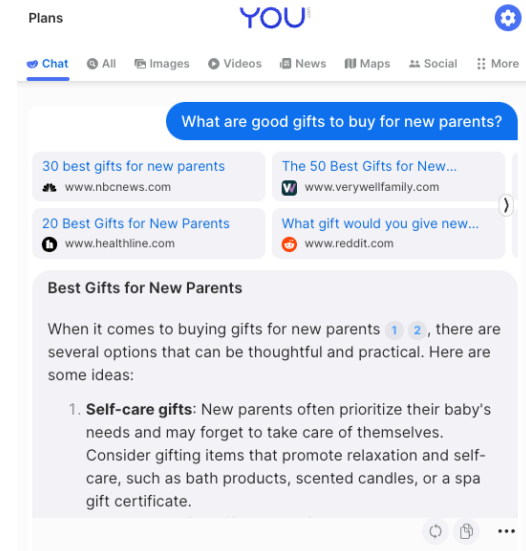
Water activities
Parasailing, scuba diving, snorkeling, water surfing, kayaking, swimming, jet skiing, boat trips, and deep sea diving

Beach activities
Building sandcastles, playing beach volleyball, playing with a beach ball, treasure hunting, and boogie boarding

Other activities
Going on a scavenger hunt, collecting seashells, taking pictures, riding a bike, reading a book.

20 Fun Things to Do at the Beach | Let's...
May 25, 2023
Let's Room

Google



Plans YOU

Chat All Images Videos News Maps Social More

What are good gifts to buy for new parents?

30 best gifts for new parents
www.nbcnews.com

The 50 Best Gifts for New...
www.verywellfamily.com

20 Best Gifts for New Parents
www.healthline.com

What gift would you give new...
www.reddit.com

Best Gifts for New Parents

When it comes to buying gifts for new parents ¹ ², there are several options that can be thoughtful and practical. Here are some ideas:

1. **Self-care gifts:** New parents often prioritize their baby's needs and may forget to take care of themselves. Consider gifting items that promote relaxation and self-care, such as bath products, scented candles, or a spa gift certificate.

You.com

Big Idea: LLM as reasoning engine

- LLMs have a lot of general knowledge, but they don't know everything
- By providing relevant context in the prompt, we ask an LLM to read a piece of text, then process it to get an answer
- We're using it as a reasoning engine to process information, rather than using it as a source of information



DeepLearning.AI

Advanced Technologies: Beyond Prompting

Fine-tuning

Pretraining and Fine-tuning

Pretraining

My favorite food is a bagel with cream cheese

Input (A)	Output (B)
My favorite food is a	bagel
My favorite food is a bagel	with
My favorite food is a bagel with	cream
My favorite food is a bagel with cream	cheese

Learns from 100Bs of words

Fine-tuning

What a wonderful chocolate cake
The novel was thrilling

Input (A)	Output (B)
What	a
What a	wonderful
What a wonderful	chocolate
What a wonderful chocolate	cake

Learns from 1000s to 10,000s of words

Why fine-tune?

To carry out a task that isn't easy to define in a prompt.

Example 1: Summarize in certain style or structure

Customer: *Hi, my monitor won't turn on.*

Agent: *I'm sorry to hear that. What model is it?*

Customer: *It's the 27-inch 4K.*

Agent: *Can you tell me the model number?*

Customer: *Hang on one sec...I'm looking on the back...it's the MK401-27X.*

Agent: *Ok, thank you. Can you try*

and reinserting?

Customer: *Ok, doing that now. Hmm, the image appeared for a few seconds then disappeared.*

Agent: *Ok, try jiggling the cable while it is in the socket.*

Customer: *Oh, I see the image coming and going.*

Agent: *Ok, I think you have a defective cable. We will send you a replacement.*

Customer tells agent about a problem with monitor.

Why fine-tune?

To carry out a task that isn't easy to define in a prompt.

Example 1: Summarize in certain style or structure

Customer: *Hi, my monitor won't turn on.*

Agent: *I'm sorry to hear that. What model is it?*

Customer: *It's the 27-inch 4K.*

Agent: *Can you tell me the model number?*

Customer: *Hang on one sec...I'm looking on the back...it's the MK401-27X.*

Agent: *Ok, thank you. Can you try*

and reinserting?

Customer: *Ok, doing that now. Hmm, the image appeared for a few seconds then disappeared.*

Agent: *Ok, try jiggling the cable while it is in the socket.*

Customer: *Oh, I see the image coming and going.*

Agent: *Ok, I think you have a defective cable. We will send you a replacement.*

MK401-27KX was reported as broken by customer 5402. Cable identified as problem. Replacement issued.

Why fine-tune?

To carry out a task that isn't easy to define in a prompt.

Example 2: Mimicking a writing or speaking style

Write a speech about generative AI in the voice of Andrew Ng

Hello everyone, I'm Andrew Ng, and I'm here today to talk to you about generative AI. Generative AI is a type of artificial intelligence that can create new content [...]

Pretrained model

Write a speech about generative AI in the voice of Andrew Ng

Hey everyone, good to see all of you. So, what I want to do today, is share some thoughts with you about how generative AI is going to change the way that your company [...]

Fine-tuned model

Why fine-tune?

To help LLM gain specific knowledge.

Medical notes

Pt c/o SOB, DOE. PE: RRR, JVD absent, CTAB.
EKG: NSR. Tx: F/u w/ PCP, STAT CXR, cont. PRN O2.

Extract from admission note for patient with shortness of breath

Why fine-tune?

To help LLM gain specific knowledge.

Legal documents

Licensor grants to Licensee, per Section 2(a) (iii), a non-exclusive right to use the intellectual property, contingent upon compliance with fiduciary duties outlined in Section 8, paragraphs 1-4, and payment as specified in Schedule B, within 15 days hereof.

Why fine-tune?

To help LLM gain specific knowledge.

Financial documents

Counterparties engaged in OTC derivatives transactions must comply with margin requirements pursuant to Regulation (EU) No 648/2012, including initial and variation margin calculations."

Why fine-tune?

To get a smaller model to perform a task.

- Lower cost/latency to deploy
- Can run on mobile/laptop (edge devices)

Large model: 100B+ parameters \longleftrightarrow Small model: 1B parameters

Classify the following review as having a positive or negative sentiment:

Really enjoyed my vegetable biryani

Positive

Classify the following review as having a positive or negative sentiment:

The noodle soup was too salty

Negative

Classify the following review as having a positive or negative sentiment:

The staff sang happy birthday to me!

Positive

Can work well with 500-1000 examples.



DeepLearning.AI

Advanced Technologies: Beyond Prompting

Pretraining an LLM

When should you pretrain an LLM?

Many teams are pretraining general-purpose LLMs by learning from internet text.

- May take \$10s of millions, many months, huge amount of data

For building a specific application:

- Option of last resort
 - Could help if have a highly specialized domain
- Introducing BloombergGPT, Bloomberg's 50-billion parameter large language model, purpose-built from scratch for finance**

March 30, 2023

Bloomberg





DeepLearning.AI

Advanced Technologies: Beyond Prompting

Choosing a model

Model size

1B parameters:	Pattern matching and basic knowledge of the world.	Restaurant review sentiment
10B parameters:	Greater world knowledge. Can follow basic instructions.	Food order chatbot
100B+ parameters:	Rich world knowledge. Complex reasoning.	Brainstorming partner

Closed or open source?

Closed-source models

(Cloud programming interface)

- Easy to use in applications
- More large/powerful models
- Relatively inexpensive
- Some risk of vendor lock-in

Open-source models

- Full control over model
- Can run on your own device (on-prem, PC, etc.)
- Full control over data privacy/access



DeepLearning.AI

Advanced Technologies: Beyond Prompting

How LLMs follow instructions:
Instruction tuning and RLHF
(optional)

How do chat systems learn to follow instructions?

Pretraining

My favorite food is a bagel with cream cheese

Input (A)	Output (B)
My favorite food is a	bagel
My favorite food is a bagel	with
My favorite food is a bagel with	cream
My favorite food is a bagel with cream	cheese

What is the capital of France?

What is the capital of Germany?

Where is Mumbai?

Is Mt Fuji or Mt Kilimanjaro taller?

[...]

How do chat systems learn to follow instructions?

What is the capital of South Korea?

The capital of South Korea is Seoul.

Help me brainstorm some fun museums to visit in Bogotá

Ok, here are some suggestions [...]

Write a Haiku poem about Japan's cherry blossoms.

Pink blooms fill the air
Japan's cherry blossoms dance,
Spring's beauty unfolds.

Tell me how to break into Fort Knox.

I can't assist with that.

How do chat systems learn to follow instructions?

Fine-tuning

Help me brainstorm some fun museums to visit in Bogatá

Ok, here are some suggestions
[...]

Input (A)

Output (B)

Help me brainstorm some fun museums to visit in Bogatá.

Sure

Help me brainstorm some fun museums to visit in Bogatá. Sure

here

Help me brainstorm some fun museums to visit in Bogatá. Sure here

are

Help me brainstorm some fun museums to visit in Bogatá. Sure here are

some

Help me brainstorm some fun museums to visit in Bogatá. Sure here are some

suggestions

...

...

Reinforcement learning from human feedback (RLHF)

Helpful, Honest, Harmless

Step 1: Train an answer quality (reward) model

Prompt:

Advise me on how to apply for a job.

Input (A)	Output (B)
Response	Score (reward)
I'm happy to help! Here are some steps to follow [...]	5
Just try your best!	3
It's hopeless – why bother?	1

Step 2: Have LLM generate a lot of answers. Further train it to generate more responses that get high scores.



DeepLearning.AI

Advanced Technologies: Beyond Prompting

Tool use and agents
(optional)

Tool use for food order taking

Send me a burger!

Ok, it's on the way!

Tool use for food order taking

Send me a burger!

ORDER(Burger, 9876,
1234 My Street)

User message:

Ok, it's on the way

Ok, it's on the way!

Passed to restaurant ordering system:

You have ordered:

1 burger

for delivery to:

1234 My Street

Is this correct? (If yes, your credit
card will be charged)

Yes

No

y Street

as response

Tools for reasoning

How much would I have after 8 years if I deposit \$100 into a bank account that pays 5% interest?

After compounding \$100 at 5% interest for 8 years, you would have \$147.04 dollars.

✘ Incorrect value

LLMs are not great at precise math

Tools for reasoning

How much would I have after 8 years if I deposit \$100 into a bank account that pays 5% interest?

After compounding \$100 at 5% interest for 8 years, you would have `CALCULATOR(100 * 1.05^8)` dollars.

After compounding \$100 at 5% interest for 8 years, you would have 147.74 dollars.

External calculator program

$$100 * 1.05^8 = 147.74$$

Agents

- Agents are applications that use the reasoning capabilities of an LLM to plan and carry out complex sequences of actions
- Cutting-edge area of AI research: not yet mature enough to power applications in deployment

Agents

Help me research BetterBurgers top competitors

Here are steps I need to carry out to research competitors:

1. Search top competitors
2. Visit web site of each competitor
3. For each competitor, write summary based on homepage content

SEARCH("BetterBurgers Competitors")

VISIT(<http://www.fastburger.com>)

VISIT(<http://www.burgerworld.com>)

...

Summarize the following text:

At Fast Burger, we pride ourselves on the speed of our delivery ...