

Pretraining

- a. Using distilgpt for textgeneration
- b. Chat generation using mistralai/Mistral-7B-Instruct-v0.2
- c. Summarize text using T5-small

```
1 # =====
2 # a.Using distilgpt2 for Text Generation
3 # =====
4 !pip install -q transformers datasets evaluate accelerate sentencepiece
5 !pip install -q torch
6
7 import torch
8 from transformers import pipeline, AutoModelForCausalLM, AutoTokenizer
9 print("✅ Environment ready. Torch version:", torch.__version__)
10 from transformers import AutoModelForCausalLM, AutoTokenizer
11 import torch
12
13 model_name = "distilgpt2"
14 tokenizer = AutoTokenizer.from_pretrained(model_name)
15 model = AutoModelForCausalLM.from_pretrained(model_name)
16
17 prompt = "Artificial intelligence is transforming the world because"
18 inputs = tokenizer(prompt, return_tensors="pt")
19
20 outputs = model.generate(
21     **inputs,
22     max_length=60,
23     num_return_sequences=3,
24     temperature=0.8,
25     top_k=50,
26     top_p=0.95,
27     do_sample=True
28 )
29
```

30

31

[Show hidden output](#)

```
1 for i, output in enumerate(outputs):
2     print(f"\n--- Generated Text {i+1} ---\n")
3     print(tokenizer.decode(output, skip_special_tokens=True))
```

--- Generated Text 1 ---

Artificial intelligence is transforming the world because the technology is moving us closer to the natural world, not closer to the natural world

--- Generated Text 2 ---

Artificial intelligence is transforming the world because of advances in the development of artificial intelligence, but it's still not new. A new

--- Generated Text 3 ---

Artificial intelligence is transforming the world because it's the kind of place where humans can communicate with each other.

```
1 # =====
2 # b. Chat Generation using mistralai/Mistral-7B-Instruct-v0.2
3 # =====
4
5 from transformers import AutoModelForCausalLM, AutoTokenizer, pipeline
6
7 model_id = "mistralai/Mistral-7B-Instruct-v0.2"
8 pipe = pipeline("text-generation", model=model_id, torch_dtype=torch.bfloat16, device_map="auto")
9
10 prompt = "You are a helpful assistant. User: What is quantum computing? Assistant:"
11
12 result = pipe(prompt, max_new_tokens=120, temperature=0.7)
13
14
```

Show hidden output

```
1 print(result[0]['generated_text'])
```

You are a helpful assistant. User: What is quantum computing? Assistant: Quantum computing is a type of computing technology that utilizes quantum

```
1 # =====
2 # c. Summarize text using T5-small
3 # =====
4
5 from transformers import pipeline
6
7 summarizer = pipeline("summarization", model="t5-small")
8
9 text = """
10 Artificial Intelligence (AI) is a field of computer science that aims to create machines
11 that can perform tasks that would normally require human intelligence. This includes
12 learning, reasoning, problem-solving, perception, and language understanding.
13 """
14
15 summary = summarizer(text, max_length=50, min_length=20, do_sample=False)
16
```

Show hidden output

```
1 print(" Summary:", summary[0]['summary_text'])
2
```

Summary: artificial intelligence is a field of computer science that aims to create machines that can perform tasks that would normally require h

Fine-tuning

- a. Using BERT(BERT-base) for sentiment analysis using IMDb sentiment dataset
 - i. Using bert-base-uncased
 - ii. Using Hugging Face Trainer API

```
1 # Install required libraries
2 !pip install -q transformers datasets accelerate evaluate
3
4 # Import necessary modules
5 from datasets import load_dataset
6 from transformers import AutoTokenizer, AutoModelForSequenceClassification, TrainingArguments, Trainer
```

```
7 import numpy as np
8 import evaluate
9
10 # Load the IMDb dataset
11 dataset = load_dataset("imdb")
12
13 # Load the BERT tokenizer
14 tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
15
16 # Tokenization function
17 def tokenize_function(examples):
18     return tokenizer(examples["text"], padding="max_length", truncation=True, max_length=512)
19
20 # Tokenize datasets
21 tokenized_datasets = dataset.map(tokenize_function, batched=True)
22
23 # Use small subsets for fast training
24 small_train_dataset = tokenized_datasets["train"].shuffle(seed=42).select(range(1000))
25 small_eval_dataset = tokenized_datasets["test"].shuffle(seed=42).select(range(200))
26
27 # Load model
28 model = AutoModelForSequenceClassification.from_pretrained("bert-base-uncased", num_labels=2)
29
30 # Training arguments (fixed: use 'eval_strategy' instead of 'evaluation_strategy')
31 training_args = TrainingArguments(
32     output_dir="bert_imdb_trainer",
33     eval_strategy="epoch",           # Fixed: was 'evaluation_strategy'
34     per_device_train_batch_size=8,
35     per_device_eval_batch_size=8,
36     num_train_epochs=1,
37     weight_decay=0.01,
38     logging_dir="logs",
39     logging_steps=10,
40     save_strategy="epoch",
41     load_best_model_at_end=True,
42     metric_for_best_model="accuracy",
43 )
```

```
44  
45 # Load accuracy metric  
46 metric = evaluate.load("accuracy")  
47  
48 # Compute metrics  
49 def compute_metrics(eval_pred):  
50     logits, labels = eval_pred  
51     predictions = np.argmax(logits, axis=-1)  
52     return metric.compute(predictions=predictions, references=labels)  
53  
54 # Initialize Trainer  
55 trainer = Trainer(  
56     model=model,  
57     args=training_args,  
58     train_dataset=small_train_dataset,  
59     eval_dataset=small_eval_dataset,  
60     compute_metrics=compute_metrics,  
61 )  
62  
63 # Train  
...  
...
```

Show hidden output

```
1 import torch  
2  
3 # Sample inference with proper device handling  
4 text = "This is a fantastic movie with great acting!"  
5 inputs = tokenizer(text, return_tensors="pt", padding=True, truncation=True, max_length=512)  
6  
7 # Move inputs to the same device as the model  
8 device = model.device # Automatically detects GPU if available  
9 inputs = {k: v.to(device) for k, v in inputs.items()}  
10  
11 # Forward pass  
12 with torch.no_grad():  
13     outputs = model(**inputs)  
14
```

```
15 # Get prediction
16 logits = outputs.logits
17 prediction = torch.argmax(logits, dim=-1).item()
18
19 print("Text:", text)
20 print("Sentiment:", "Positive" if prediction == 1 else "Negative")
```

```
Text: This is a fantastic movie with great acting!
Sentiment: Positive
```

b. Fine-tuning GPT for summarization.

i. Using gpt2 for text summarization using DailyMail/XSum dataset

ii. Using Hugging face Pipeline

```
1 !pip install -U transformers datasets accelerate sentencepiece -q
2
3 import os
4 os.environ["WANDB_DISABLED"] = "true"
5
6 from datasets import load_dataset
7 from transformers import AutoTokenizer, GPT2LMHeadModel
8 from transformers import DataCollatorForLanguageModeling, TrainingArguments, Trainer
9 from transformers import pipeline
10
11 dataset = load_dataset("cnn_dailymail", "3.0.0")["train"].select(range(500))
12
13 tokenizer = AutoTokenizer.from_pretrained("gpt2")
14 tokenizer.pad_token = tokenizer.eos_token
15 model = GPT2LMHeadModel.from_pretrained("gpt2")
16
17 def preprocess(b):
18     t = ["summarize: " + x for x in b["article"]]
19     return tokenizer(t, truncation=True, padding="max_length", max_length=256)
20
21 tokenized = dataset.map(preprocess, batched=True, remove_columns=["article", "highlights"])
22 collator = DataCollatorForLanguageModeling(tokenizer=tokenizer, mlm=False)
```

```
23  
24 args = TrainingArguments(  
25     output_dir=".gpt2-sum",  
26     per_device_train_batch_size=2,  
27     gradient_accumulation_steps=4,  
28     num_train_epochs=1,  
29     learning_rate=5e-5,  
30     logging_steps=50,  
31     fp16=True  
32 )  
33  
34 trainer = Trainer(model=model, args=args, train_dataset=tokenized, data_collator=collator)  
35 trainer.train()  
36  
37
```

Show hidden output

```
1 summarizer = pipeline("text-generation", model=model, tokenizer=tokenizer, max_new_tokens=60)  
2 print(summarizer("summarize: The Eiffel Tower is visited by millions every year.")[0]["generated_text"])
```

```
Device set to use cuda:0  
summarize: The Eiffel Tower is visited by millions every year. But for the second year running, the hotel's owners are making it back to normal.
```

```
1 # Install required libraries  
2 !pip install -q transformers datasets accelerate evaluate  
3  
4 # Import modules  
5 from datasets import load_dataset  
6 from transformers import AutoTokenizer, AutoModelForSequenceClassification, TrainingArguments, Trainer  
7 import numpy as np  
8 import evaluate  
9 import torch  
10  
11 # Load IMDb dataset  
12 dataset = load_dataset("imdb")  
13
```

```
14 # Load DistilBERT tokenizer
15 tokenizer = AutoTokenizer.from_pretrained("distilbert-base-uncased")
16
17 # Tokenization function
18 def tokenize_function(examples):
19     return tokenizer(examples["text"], truncation=True, padding="max_length", max_length=512)
20
21 # Tokenize datasets
22 tokenized_datasets = dataset.map(tokenize_function, batched=True)
23
24 # Small subsets for quick training
25 small_train_dataset = tokenized_datasets["train"].shuffle(seed=42).select(range(1000))
26 small_eval_dataset = tokenized_datasets["test"].shuffle(seed=42).select(range(200))
27
28 # Load DistilBERT model for classification
29 model = AutoModelForSequenceClassification.from_pretrained("distilbert-base-uncased", num_labels=2)
30
31 # Training arguments
32 training_args = TrainingArguments(
33     output_dir="distilbert_imdb",
34     eval_strategy="epoch",
35     save_strategy="epoch",
36     per_device_train_batch_size=16,    # DistilBERT is lighter → can use larger batch
37     per_device_eval_batch_size=16,
38     num_train_epochs=1,
39     weight_decay=0.01,
40     logging_dir="logs",
41     logging_steps=10,
42     load_best_model_at_end=True,
43     metric_for_best_model="accuracy",
44     report_to="none",   # Disable wandb
45 )
46
47 # Accuracy metric
48 metric = evaluate.load("accuracy")
49
50 def compute_metrics(eval_pred):
```

```
51     logits, labels = eval_pred
52     predictions = np.argmax(logits, axis=-1)
53     return metric.compute(predictions=predictions, references=labels)
54
55 # Trainer
56 trainer = Trainer(
57     model=model,
58     args=training_args,
59     train_dataset=small_train_dataset,
60     eval_dataset=small_eval_dataset,
61     compute_metrics=compute_metrics,
62 )
63
64 # Train
65 trainer.train()
66
67 # Move model to correct device for inference
```

Show hidden output

```
1 text = "This is a fantastic movie with great acting!"
2 inputs = tokenizer(text, return_tensors="pt", truncation=True, padding=True, max_length=512)
3 inputs = {k: v.to(model.device) for k, v in inputs.items()}
4
5 with torch.no_grad():
6     outputs = model(**inputs)
7     pred = torch.argmax(outputs.logits, dim=-1).item()
8
9 print("Text:", text)
10 print("Sentiment:", "Positive" if pred == 1 else "Negative")
```

```
Text: This is a fantastic movie with great acting!
Sentiment: Positive
```

For text summarization Swap the models (bart-large-cnn, t5-small, pegasus-xsum) and change parameters and use trainer API

```
1 # Install required libraries
2 !pip install -q transformers datasets evaluate rouge_score accelerate torch
3
4 import torch
5 from datasets import load_dataset
6 from transformers import AutoTokenizer, AutoModelForSeq2SeqLM, Seq2SeqTrainingArguments, Seq2SeqTrainer, DataCollatorForSeq2Seq
7 import numpy as np
8 import evaluate
9
10 # === CONFIG: 3 MODELS ===
11 models_to_train = [
12     {"name": "facebook/bart-large-cnn", "batch_size": 4, "epochs": 1, "prefix": ""},
13     {"name": "t5-small", "batch_size": 8, "epochs": 1, "prefix": "summarize: "},
14     {"name": "google/pegasus-xsum", "batch_size": 4, "epochs": 1, "prefix": ""}
15 ]
16
17 # Load CNN/DailyMail dataset
18 dataset = load_dataset("cnn_dailymail", "3.0.0")
19 train_dataset = dataset["train"].shuffle(seed=42).select(range(1000))
20 eval_dataset = dataset["validation"].shuffle(seed=42).select(range(200))
21
22 # Load ROUGE
23 rouge = evaluate.load("rouge")
24
25 def compute_metrics(eval_pred):
26     preds, labels = eval_pred
27     decoded_preds = tokenizer.batch_decode(preds, skip_special_tokens=True)
28     labels = np.where(labels != -100, labels, tokenizer.pad_token_id)
29     decoded_labels = tokenizer.batch_decode(labels, skip_special_tokens=True)
30     result = rouge.compute(predictions=decoded_preds, references=decoded_labels, use_stemmer=True)
31     return {k: round(v, 4) for k, v in result.items()}
32
33 # === TRAINING LOOP ===
34 for idx, config in enumerate(models_to_train):
35     print(f"\n{'='*70}")
36     print(f"TRAINING MODEL {idx+1}/3: {config['name']}")
37     print(f"{'='*70}\n")
```

```
38
39 MODEL_NAME = config["name"]
40 BATCH_SIZE = config["batch_size"]
41 EPOCHS = config["epochs"]
42 PREFIX = config["prefix"]
43
44 # Load tokenizer & model
45 tokenizer = AutoTokenizer.from_pretrained(MODEL_NAME)
46 model = AutoModelForSeq2SeqLM.from_pretrained(MODEL_NAME)
47
48 # Preprocess function
49 def preprocess(examples):
50     inputs = [PREFIX + doc for doc in examples["article"]]
51     model_inputs = tokenizer(inputs, max_length=512, truncation=True, padding="max_length")
52     with tokenizer.as_target_tokenizer():
53         labels = tokenizer(examples["highlights"], max_length=128, truncation=True, padding="max_length")
54     model_inputs["labels"] = labels["input_ids"]
55     return model_inputs
56
57 # Tokenize
58 tokenized_train = train_dataset.map(preprocess, batched=True, remove_columns=["article", "highlights", "id"])
59 tokenized_eval = eval_dataset.map(preprocess, batched=True, remove_columns=["article", "highlights", "id"])
60
61 data_collator = DataCollatorForSeq2Seq(tokenizer, model=model)
62
63 # Training args
64 training_args = Seq2SeqTrainingArguments(
65     output_dir=f"./cnn_model_{idx}_{MODEL_NAME.split('/')[-1]}",
66     eval_strategy="epoch",
67     save_strategy="epoch",
68     per_device_train_batch_size=BATCH_SIZE,
69     per_device_eval_batch_size=BATCH_SIZE,
70     num_train_epochs=EPOCHS,
71     weight_decay=0.01,
72     predict_with_generate=True,
73     fp16=torch.cuda.is_available(),
74     logging_steps=10,
```

```
75     load_best_model_at_end=True,
76     metric_for_best_model="rouge1",
77     report_to="none",
78     disable_tqdm=True,
79   )
80
81 # Trainer
82 trainer = Seq2SeqTrainer(
83     model=model,
84     args=training_args,
85     train_dataset=tokenized_train,
86     eval_dataset=tokenized_eval,
87     tokenizer=tokenizer,
88     data_collator=data_collator,
89     compute_metrics=compute_metrics,
90   )
91
92 # Train
93 trainer.train()
94
95 # Save
96 save_dir = f"./trained_cnn_{MODEL_NAME.split('/')[-1]}[-1]}"
97 model.save_pretrained(save_dir)
98 tokenizer.save_pretrained(save_dir)
99 print(f"Model saved to: {save_dir}\n")
```

Show hidden output

```
1 from transformers import pipeline
2 import torch, os, warnings, textwrap
3
4 # Hide all warnings
5 warnings.filterwarnings("ignore")
6 os.environ["TOKENIZERS_PARALLELISM"] = "false"
7
8 # Test article
9 text = """London (CNN) -- A British man has been charged with murder after his wife and daughter were found
```

```
10 The bodies of 43-year-old Samantha Baldwin and her 7-year-old daughter, Lily, were discovered Monday at thei
11 The suspect, 45-year-old David Baldwin, was arrested at the scene and remains in custody, police said. He is
12 A post-mortem examination is being carried out to establish the cause of death, police said. Neighbors descr
13
14 # Your actual saved model folders
15 models = [
16     ("BART-Large-CNN",    "/content/trained_cnn_bart-large-cnn"),
17     ("T5-Small",          "/content/trained_cnn_t5-small"),
18     ("Pegasus-XSum",     "/content/trained_cnn_pegasus-xsum")
19 ]
20
21 print("=" * 90)
22 print("CNN/DAILYMAIL SUMMARIZATION RESULTS".center(90))
23 print("=" * 90 + "\n")
24
25 for name, folder in models:
26     if not os.path.exists(folder):
27         print(f"ERROR: Folder not found → {folder}")
28         print("  Make sure the model was saved correctly.\n")
29         continue
30
31     summarizer = pipeline(
32         "summarization",
33         model=folder,
34         tokenizer=folder,
35         device=0 if torch.cuda.is_available() else -1,
36         max_length=80,
37         min_length=30,
38         do_sample=False,
39         truncation=True,
40         clean_up_tokenization_spaces=True
41     )
42
43     summary = summarizer(text)[0]["summary_text"]
44
45     print(f"{name.upper()}")
46     print("-" * 70)
```

```
47     print(textwrap.fill(summary, width=88))

=====
CNN/DAILYMAIL SUMMARIZATION RESULTS
=====

Device set to use cuda:0
Asking to truncate to max_length but no maximum length is provided and the model has no predefined maximum length. Default to no truncation.
BART-LARGE-CNN

Samantha Baldwin, 43, and her 7-year-old daughter, Lily, were found dead at their home.
David Baldwin, 45, was arrested at the scene and remains in custody, police say. A post-mortem examination is being carried out to establish the cause of death.

Device set to use cuda:0
Your max_length is set to 200, but your input_length is only 148. Since this is a summarization task, where outputs shorter than the input are typical.
T5-SMALL

the bodies of Samantha Baldwin, 43, and her 7-year-old daughter, Lily, were found Monday
at their home in southern England. David Baldwin, 45, was arrested at the scene and
remains in custody, police say. Neighbors described the family as "quiet and private"

Device set to use cuda:0
PEGASUS-XSUM

The bodies of 43-year-old Samantha Baldwin and her 7-year-old daughter, Lily, were
discovered Monday at their home in Basingstoke, Hampshire, police said.
```

For Question Answering Use RoBERTa variants, GPT-3/4 fine-tuned on SQuAD dataset

```
1 !pip install -q transformers datasets accelerate sentencepiece
2
3 from datasets import load_dataset
4 from transformers import AutoTokenizer, AutoModelForQuestionAnswering
5 from transformers import TrainingArguments, Trainer, DefaultDataCollator
6 import torch
7 import os
8
9 os.environ["WANDB_DISABLED"] = "true"
10
11 dataset = load_dataset("squad")
```

```
12
13 tokenizer = AutoTokenizer.from_pretrained("roberta-base")
14 model = AutoModelForQuestionAnswering.from_pretrained("roberta-base")
15
16 def preprocess(example):
17     inputs = tokenizer(
18         example["question"],
19         example["context"],
20         max_length=384,
21         truncation=True,
22         padding="max_length",
23         return_offsets_mapping=True
24     )
25     offset = inputs["offset_mapping"]
26     answers = example["answers"]["text"][0]
27     start_char = example["answers"]["answer_start"][0]
28     end_char = start_char + len(answers)
29
30     start_pos = 0
31     end_pos = 0
32     for i, (s, e) in enumerate(offset):
33         if s <= start_char < e:
34             start_pos = i
35         if s < end_char <= e:
36             end_pos = i
37
38     inputs["start_positions"] = start_pos
39     inputs["end_positions"] = end_pos
40     inputs.pop("offset_mapping")
41     return inputs
42
43 tokenized = dataset.map(preprocess, batched=False)
44
45 collator = DefaultDataCollator()
46
47 args = TrainingArguments(
48     output_dir=". ./roberta-qa",
```

```
49     per_device_train_batch_size=8,  
50     num_train_epochs=1,  
51     logging_steps=100,  
52     report_to=[]  
53 )  
54  
55 trainer = Trainer(  
56     model=model,  
57     args=args,  
58     train_dataset=tokenized["train"],  
59     eval_dataset=tokenized["validation"],  
60     data_collator=collator,  
61     tokenizer=tokenizer  
62 )  
63  
64 trainer.train()  
65  
66 from transformers import pipeline  
67 qa = pipeline("question-answering", model=model, tokenizer=tokenizer)
```

Show hidden output

```
1 print(qa({"question":"Who founded Microsoft?", "context":"Microsoft was founded by Bill Gates and Paul Allen.  
2  
{'score': 0.8770565986633301, 'start': 25, 'end': 50, 'answer': 'Bill Gates and Paul Allen'}
```

Build a simple rule-based chatbot for FAQs.

```
1 def chatbot(message):  
2     msg = message.lower()  
3  
4     # Greetings  
5     if "hello" in msg or "hi" in msg:  
6         return "Hello! How can I help you today?"  
7  
8     # About company
```

```
9     elif "what is your name" in msg:
10         return "I'm a simple FAQ chatbot."
11
12     elif "what do you do" in msg:
13         return "I answer frequently asked questions."
14
15     # Timings
16     elif "working hours" in msg or "open" in msg:
17         return "We are open from 9 AM to 6 PM, Monday to Saturday."
18
19     # Pricing
20     elif "price" in msg or "cost" in msg:
21         return "Pricing depends on the service. Basic plans start at $10/month."
22
23     # Support
24     elif "support" in msg or "help" in msg:
25         return "You can contact support at support@example.com."
26
27     # Location
28     elif "location" in msg or "where are you" in msg:
29         return "We are located in Chennai, India."
30
31     # Goodbye
32     elif "bye" in msg or "goodbye" in msg:
33         return "Goodbye! Have a great day!"
34
35     # Default fallback
36     else:
37         return "I'm not sure about that. Could you ask something else?"
38
39
40 # Chat loop
41 print("Type 'exit' to stop.\n")
42 while True:
43     user = input("You: ")
44     if user.lower() == "exit":
45         print("Chatbot: Goodbye!")
```

```
46         break
47     print("Chatbot:", chatbot(user))
--
```

Type 'exit' to stop.

```
You: support
Chatbot: You can contact support at support@example.com.
You: bye
Chatbot: Goodbye! Have a great day!
You: exit
Chatbot: Goodbye!
```

Apply various prompting techniques and view the difference

```
1 !pip install transformers accelerate sentencepiece -q
2
3 from transformers import AutoTokenizer, AutoModelForSeq2SeqLM, pipeline
4
5 model = "google/flan-t5-large" # fits on Colab T4 with 16GB
6 pipe = pipeline("text2text-generation", model=model, tokenizer=model, max_new_tokens=150)
7
8 def run(prompt):
9     print("\n--- PROMPT ---")
10    print(prompt)
11    print("--- OUTPUT ---")
12    print(pipe(prompt)[0]["generated_text"])
13    print("\n")
```

Show hidden output

```
1
2
3 # 1. Zero-shot
4 run("Explain blockchain in simple terms.")
5
6 # 2. One-shot
7 run("""
8 Q: What is AI?
```

```
9 A: AI means machines that can think and learn like humans.  
10  
11 Q: What is blockchain?  
12 A:  
13 """)  
14  
15 # 3. Few-shot  
16 run("")  
17 Q: What is AI?  
18 A: AI refers to machines performing tasks that require human intelligence.  
19  
20 Q: What is Machine Learning?  
21 A: ML is a subset of AI that learns patterns from data.  
22  
23 Q: What is blockchain?  
24 A:  
25 """)  
26  
27 # 4. Chain-of-thought  
28 run("Why do objects fall to the ground? Let's think step by step.")  
29  
30 # 5. Self-consistency style  
31 run("")  
32 Why does ice float on water?  
33  
34 Let's think step by step:  
35 """)  
36  
37 # 6. Instruction prompting  
38 run("")  
39 Summarize the following in 3 lines:  
40 Artificial Intelligence enables machines to learn from data and make decisions.  
41 """)  
42  
43 # 7. Role prompting  
44 run("")  
45 You are a polite science teacher.
```

```
46 Explain how the Internet works.  
47 """)  
48  
49 # 8. Delimiter prompting  
50 run("")  
51 Summarize this text: <text>Machines learn patterns from data in ML.</text>  
52 """)  
53  
54 # 9. Contrastive prompt  
55 run("")  
56 Bad summary: It's about machines.  
57 Good summary: Give a precise, clear summary.  
58  
59 Text: Deep learning uses neural networks with layers to learn patterns from data.  
60  
61 Summary:  
62 "")
```

--- PROMPT ---

You are a polite science teacher.
Explain how the Internet works.

--- OUTPUT ---

The Internet is a network of computers that communicate with each other. The computers are called nodes. The nodes are called nodes. The nodes a

--- PROMPT ---

Summarize this text: <text>Machines learn patterns from data in ML.</text>

--- OUTPUT ---

Machine learning (ML) is a branch of computer science that focuses on learning patterns from data.

--- PROMPT ---

Bad summary: It's about machines.
Good summary: Give a precise, clear summary.

Text: Deep learning uses neural networks with layers to learn patterns from data.

Summary:

--- OUTPUT ---

Deep learning uses neural networks with layers to learn patterns from data.

Compare extractive and abstractive summarization using Hugging Face and TextRank

```
1 !pip install -U transformers datasets summa -q
2
3 from summa.summarizer import summarize
4 from transformers import pipeline
5
6 text = """Artificial Intelligence is transforming industries across the globe.
7 Companies are rapidly adopting AI solutions to automate tasks, increase efficiency,
8 enhance customer experience, and make better decisions. The rise of AI-driven
9 innovation has created new opportunities in healthcare, finance, education, and transportation."""
10
```

```
--  
11 extractive = summarize(text, ratio=0.4)  
12  
13 abstractive_model = pipeline("summarization", model="t5-base", tokenizer="t5-base")  
14 abstractive = abstractive_model("summarize: " + text, max_length=60, min_length=20)[0]["summary_text"]  
15  
16
```

Show hidden output

```
1 print("Extractive Summary (TextRank):\n", extractive, "\n")  
2 print("Abstractive Summary (HF):\n", abstractive)  
3
```

Extractive Summary (TextRank):

Companies are rapidly adopting AI solutions to automate tasks, increase efficiency,
The rise of AI-driven

Abstractive Summary (HF):

companies are rapidly adopting AI solutions to automate tasks . the rise of AI-driven innovation has created new opportunities .

Compare retrieval-based vs generative chatbot responses.

```
1 !pip install -q transformers scikit-learn  
2  
3 from sklearn.feature_extraction.text import TfidfVectorizer  
4 from sklearn.metrics.pairwise import cosine_similarity  
5 from transformers import AutoModelForSeq2SeqLM, AutoTokenizer  
6  
7 corpus = [  
8     "Hello, how can I help you?",  
9     "AI stands for Artificial Intelligence.",  
10    "Machine learning is a subset of AI.",  
11    "The weather is nice today.",  
12    "Chatbots can be retrieval-based or generative."  
13 ]  
14  
15 vectorizer = TfidfVectorizer()  
16 X = vectorizer.fit_transform(corpus)
```

```
17
18 def retrieval_chatbot(query):
19     q_vec = vectorizer.transform([query])
20     sims = cosine_similarity(q_vec, X).flatten()
21     return corpus[sims.argmax()]
22
23 tokenizer = AutoTokenizer.from_pretrained("google/flan-t5-small")
24 model = AutoModelForSeq2SeqLM.from_pretrained("google/flan-t5-small")
25
26 def generative_chatbot(query):
27     inputs = tokenizer(query, return_tensors="pt")
28     output = model.generate(**inputs, max_new_tokens=50)
29     return tokenizer.decode(output[0], skip_special_tokens=True)
30
31 query = "What is AI?"
32
```

Show hidden output

```
1 print("User:", query)
2 print("\nRetrieval-based Response:\n", retrieval_chatbot(query))
3 print("\nGenerative Response:\n", generative_chatbot(query))
```

User: What is AI?

Retrieval-based Response:

Machine learning is a subset of AI.

Generative Response:

AI (disambiguation) AI is a term used to describe a person's behavior.

