LLMs in the wilderness: Leveraging the power of LLMs in our business





VBYO 2025

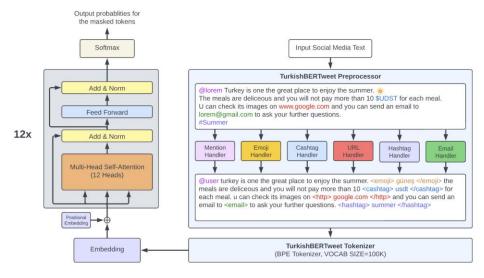




Ali Najafi

PhD in Computer Science, Sabancı University, 2024-Present MSc in Computer Science, Sabancı University, 2022-2024

o najafi-ali.com



Published Works:



- Turkishbertweet: Fast and reliable large language model for social media analysis
- First public dataset to study 2023
 Turkish general election



On going Projects:



- Social Media Data Adaptors: Standardization and Anonymization of Social Network Datasets
- Cross-Platform Bot Detection via Modeling of Social Media Account Behavior
- Analyzing the Impact of Voting Agency Discrepancies on Social Media Discourse During Turkey's 2023 Election Nights

By the end of this session, You will know

- How to approach an NLP project
- What resources exist
- What the pros and cons of LLMs are
- How to apply models to a real project
- How realistic the LLM-Hype is
- How to finetune generative Language models

Hands on Project

- Amazon company after witnessing the advances in the AI field, decides to develop a Chatbot for its customers to help them purchase shoes!
- They recruited a Machine Learning Engineer to develop it!
- They provided the details of their shoe products and their prices!
- Now, let's say you are the ML guy, how can you build such a system?



Let's check out the data!

- The dataset is available at <u>amazon-uk-shoes-dataset</u>
- The downloaded data is in JSON format!
- Available information:
 - url: Link to the product
 - title: title of the product
 - asin: Product Unique ID
 - price: Price of the product
 - brand: Brand of the product
 - product_details: Details of product
 - breadcrumbs: a navigational aid that allows users to keep track of their current location on a website or interface
 - images_list: links of product images
 - features: features of product



amazon_uk_shoes_dataset.json

	<i>હ</i>	url	~	T title	~	T a	sin	~	T	price	~	□ brand ✓	T	product_d	letails	Â	ABOU
1	https://www.an	mazon.co.uk/dp/B08	BLP231I	Geox Jr Sandal Strada B Fisherman,	Br	B08BLP2	231K		£50.00			Visit the Geox Store	Package D	imensions	: 31.	2 x	
2	https://www.am	mazon.co.uk/dp/B08	N587YZ	Fila Women's Oakmont Tr Sneaker		B08N587	YZ9		£49.57 -	£234.95		Fila	Product D	imensions	: 32.	51 x	Last
3	https://www.an	mazon.co.uk/dp/B09	18K4H1	Gabor Rollingsoft Trainers in Plus	Si	B0918K4	H1W		No data.			Gabor	Is Discon	tinued By M	1anufact	urer	Own
4	https://www.an	mazon.co.uk/dp/B07	KMB98C	Merrell Women'S Bare Access Xtr Tr	ail	B07KMB9	8CG		£67.00 -	£182.44		Visit the Merrell Store	Package D	imensions	: 28.	96 x	Crea
5	https://www.an	mazon.co.uk/dp/B08	CN3S1ZI	Desigual Women's Shoes_Runner_cmot	l S	B08CN3S	1ZK		£38.96 -	£81.10		Desigual	Package D	imensions	: 34.	6 x	Size
6	https://www.an	mazon.co.uk/dp/B07	KJX15D	Aquatalia Men's Adamo Dress Calt P	enn:	B07KJX1	5DV		£161.63	- £202.04		Aquatalia	Product D	imensions	: 24.	13 x	Displa
7	https://www.an	mazon.co.uk/dp/B08	9GYT1H	Womens Convertible Slip-on, Tan, 1	2	B089GYT	1HP		£52.82			Unknown	Package D	imensions	: 32	x 18	ama
8	https://www.an	mazon.co.uk/dp/B01	614251	Kappa Delhi Footwear Unisex, Mesh,	Woi	B016142	251K		£46.18			Visit the Kappa Store	Date Firs	t Available	2 : 3/	0 De	TABL
9	https://www.an	mazon.co.uk/dp/B07	MJ4DRR:	Naturalizer Women's Cairo 4 Sneake	r, I	B07MJ4D	RRZ		£41.66 -	£73.51		Naturalizer	Package D	imensions	: 32.	76 x	
10	https://www.an	mazon.co.uk/dp/B07	GX8GPJI	Rockport - Womens Cobb Hill Hattie	Ηi	B07GX8G	SPJN		£120.35	- £130.73		Visit the Rockport Store	Package D	imensions	: 30.	48 x	& ur
11	https://www.an	mazon.co.uk/dp/B01	2DS5U2	Reebok Men's NPC Ii Fashion Sneake	r	B012DS5	SU2Y		£47.00 -	£156.93		Visit the Reebok Store	Is Discon	tinued By M	Manufact	urer	T ti
12	https://www.an	mazon.co.uk/dp/B08	D3F9JQ	adidas Women's Superstar Bold W Gy	mna	B08D3F9	JQY		£85.00 -	£93.61		Visit the adidas Store	Package D	imensions	: 28.	8 x	T as
13	https://www.an	mazon.co.uk/dp/B01	GRUZWA	Geox Women's D Stinge a Fashion Sn	eak	B01GRUZ	OAW.		£64.95 -	£117.00		Visit the Geox Store	Is Discon	tinued By M	Manufact	urer	
14	https://www.an	mazon.co.uk/dp/B07	SQ5MLX	K-Swiss Men's Clean Court Ii CMF L	OW-	B07SQ5M	ILXZ		£20.99 -	£67.57		K-Swiss	Package D	imensions	: 30.	48 x	T pr
15	https://www.an	mazon.co.uk/dp/B00	LIHRNQ	Volcom Mens Sub Zero Boot Winter		B00LIHR	RNQI		£65.11 -	£238.10		Visit the Volcom Store	Package D	imensions	: 36.	3 x	T br
16	https://www.an	mazon.co.uk/dp/B07	H3XFPR	ASICS Mens Gel-Ziruss 2 Running Sh	oes	B07H3XF	PR9		£99.79 -	£124.99		ASICS	Product D	imensions	: 25.	25 x	T pr
17	https://www.an	mazon.co.uk/dp/B07	KQKYZY	bugatti Women's 432636065969 Low-T	ор	B07KQKY	ZYR		£28.77 -	£62.87		bugatti	Is Discon	tinued By M	Manufact	urer	
18	https://www.an	mazon.co.uk/dp/B08	HJJ5N7	Pablosky Girl's 495000 Sandal		В08НЈЈ5	N7R		£25.76 -	£36.69		Pablosky	Package D	imensions	: 23	x 16	T bi
19	https://www.an	mazon.co.uk/dp/B08	635QDS	HROYL Girls Dance Shoes Satin Lati	n/C	B08635Q	DSF		£25.99 -	£28.99		Visit the HROYL Store	Is Discon	tinued By M	Manufact	urer	T im
20	https://www.an	mazon.co.uk/dp/B08	4GB1DS	MILLET Unisex's Hike Up Mid GTX SI	ouc	B084GB1	DS3		£108.49	- £143.40		Visit the MILLET Store	Package D	imensions	: 35	x 25	T fe
21	https://www.am	mazon.co.uk/dp/B07	T26WQF	adidas Originals Men's Superstar F	oun	B07T26W	/QF5		£27.50 -	£214.11		adidas Originals	Product D	imensions	: 30	x 35	Δ, 16
22	https://www.an	mazon.co.uk/dp/B08	24F5JRI	victoria Women's Tenis Piel/Viruta	s G.	B0824F5	JRM		£42.49 -	£75.49		VICTORIA	Package D	imensions	: 29	x 29	
23	https://www.an	mazon.co.uk/dp/B08	XJG46Q	ZYLDK Garden Shoes Children's Unis	ex	B08XJG4	16Q4		£8.06 -	£9.10		ZYLDK	Package D	imensions	: 23	x 19	

	ABOUT THIS FILE							
X	Last Updated	3 years ago						
r	Owner	Crawl Feeds						
X	Created	3 years ago						
	Size	13.16 MB						
x 8	Displaying 9 columns, 11,605 rows in table amazon_uk_shoes_dataset							
e x	TABLE COLUMN	ABLE COLUMNS						
Х	€ url (i)							
r	T title (i) T asin (i) T price (i)							
er								
Х								
	T brand (i)	T brand (i)						
x	T product_details (i) T breadcrumbs (i)							
6								
r	T images_list (∏ images_list ① ∏ features ①						
5	T features (i)							
9								

NLP Tasks

Language Translation

• EN -> TR, EN -> FA

Text Classification/Regression

- Sentiment Analysis
- Personality Detection
- Topic Detection
- Hate Speech Detection
- Stance Prediction

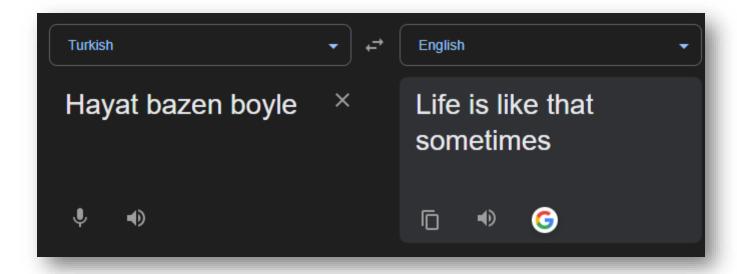
Token Classification

• Name Entity Recognition (NER)

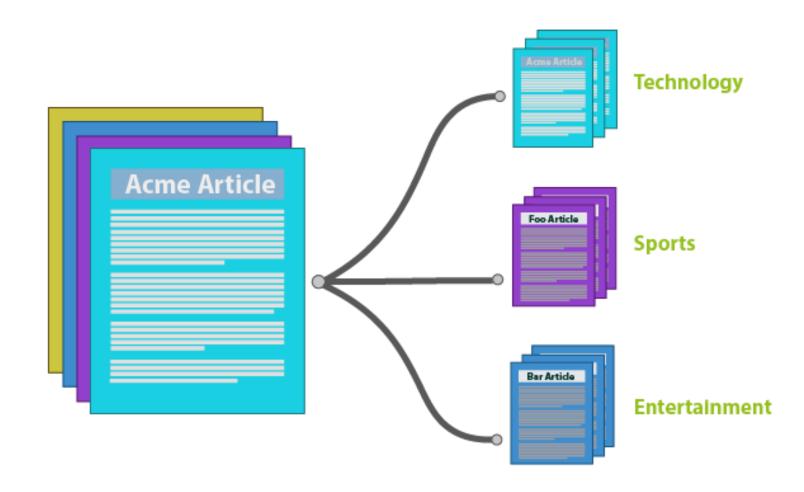
Text Generation

- Question Answering
- Text Summarization

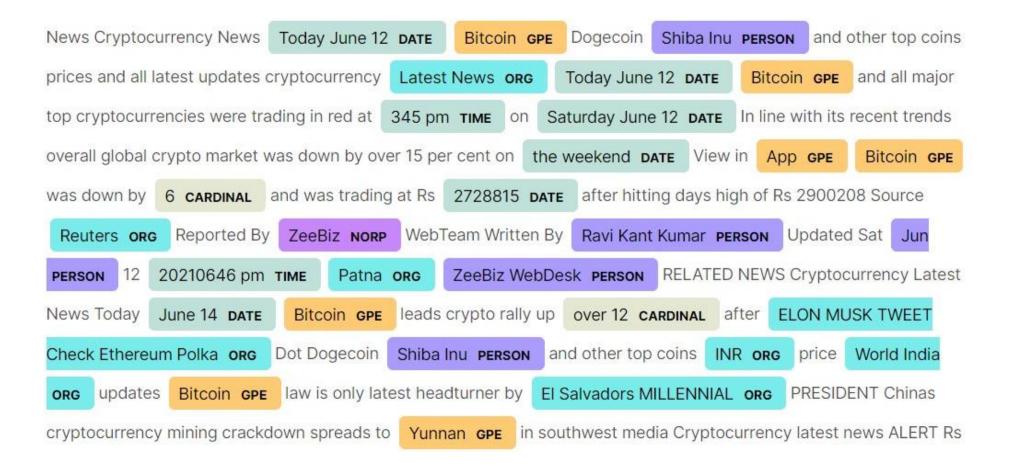
Language Translation



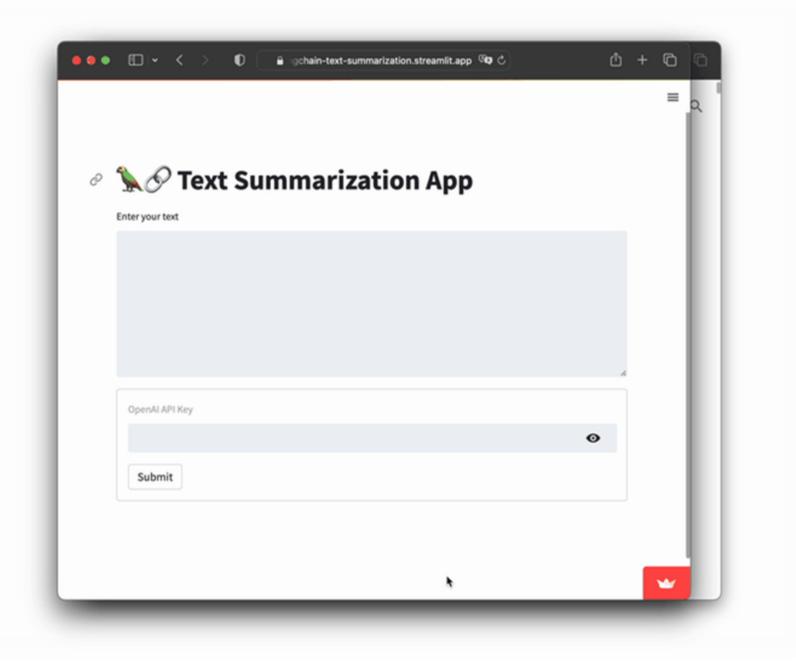
Text Classification



Token Classification



Text Summarization



Question Answering

rammin... I wonder now 5 your website doing:

give me my daily briefin

In the last 24 hours, you're traffic was 8 visitors. Out of the 8 visitors 50.0% were new visitors to your site. Your top 3 sources for new traffic came from t.co ,google , and (direct) , t.co represented the biggest chunk at 60.0%. Your most engaged users were from google They spend an average of 130.4 min on your site and visited 2.5 pages per session. I predict you will have 7 visitors to your site in the next 24 hours.

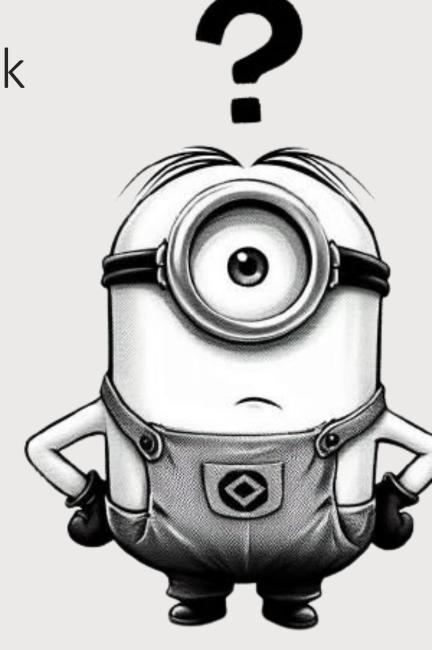
Hello, friend. I can help you with your data.

Ask me anything...



We should ask ourselves

- Do we need LLMs?
- Can we afford them?
- Do we have access to skilled people?
- Is the outcome significant enough to invest in them?





Which model?

There are many already pre-trained models available on Hugging Face

BUT Which model should I choose?

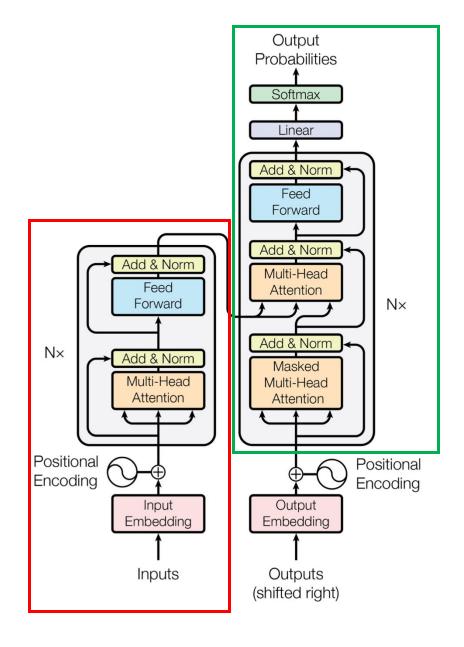
Several factors help us answer this question!

- 1. NLP Task
- 2. Computational Resources
- 3. Domain Requirement
- 4. LICENSE

NLP Task

• If your task is other than text generation, Encoder models are available!

If your task contains generation, Encoder-Decoder OR
 Decoder models are available. Need to mention that you
 can still use these types of models for non-generative
 tasks as well!



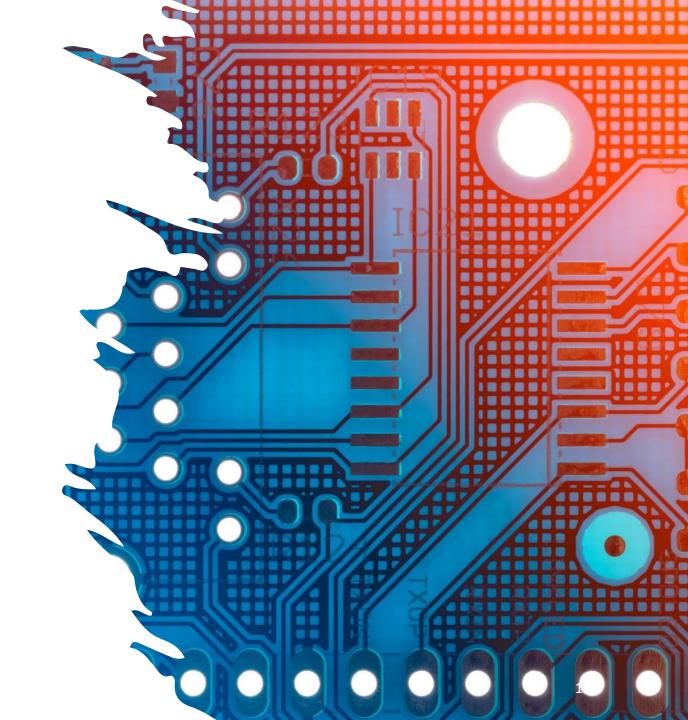
Computational Resources

Computation power plays a great role in model selection.

Computation power means the GPU infrastructure of the company!

There are two main things about GPUs that we are interested in:

- RAM
- Speed



Domain Requirements

- Each model has been trained on specific data for instance:
 - Internet Data
 - Code Data
 - ...
- In our project, models may have seen data about shoes, but they have no clue about the available products of Amazon company.
- Language matters: If you are targeting customers of the UK, then the model needs to interact in English. What if you want a chatbot that can handle customers with different languages?





Terminologies in LLM World

- Pre-Training: Training a language model from scratch
- Finetuning: Training the pre-trained model on a domain-specific dataset(s)
- Instruction Finetuned: The model is finetuned on a dataset that contains Question-Answer pairs.
- Zero Shot: Using model without finetuning
- Tokenizer Special tokens: Tokens that the tokenizer maps into a single ID.

Model Architectures and their objectives

models like _____ models can be trained on _____ objectives.

- 1. Encoder; BERT-based; classification or masked language modeling
- 2. Encoder-Decoder; T5-based; sequence-to-sequence or masked language modeling
- 3. Decoder; GPT-based; causal language modeling or masked language modeling

Language Modeling

Causal Language Modeling: predict the next token

• Input: Hard work always pays _____.

• Output: off

Masked Language Modeling: predict the masked token

• Input: Hard work ____ pays off.

Output: always

Sequence to sequence: predict the corrupted tokens

• Input: Hard _____ pays off.

• Output: work always

Finetuning Approaches

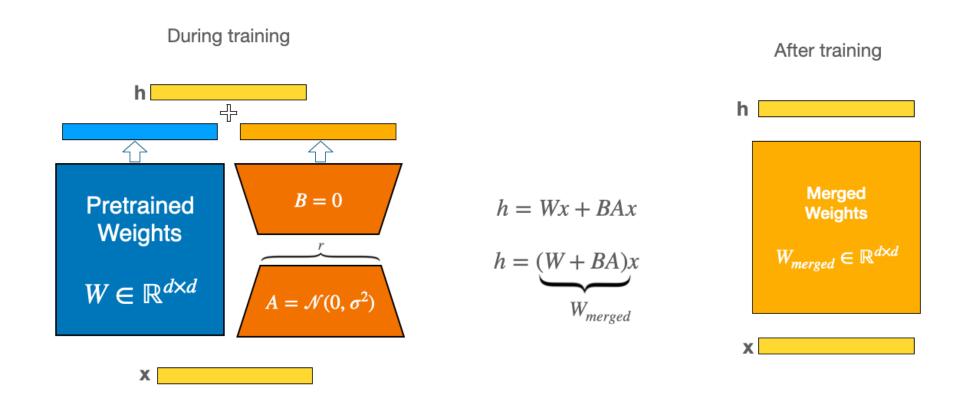
Full-Finetuning: Updating directly the wights of the models

LoRA Finetuning: Just updating the LoRA adapter's weights and keeping the model's weights intact

RAG Finetuning: Finetuning with improved prompts using our knowledge database.

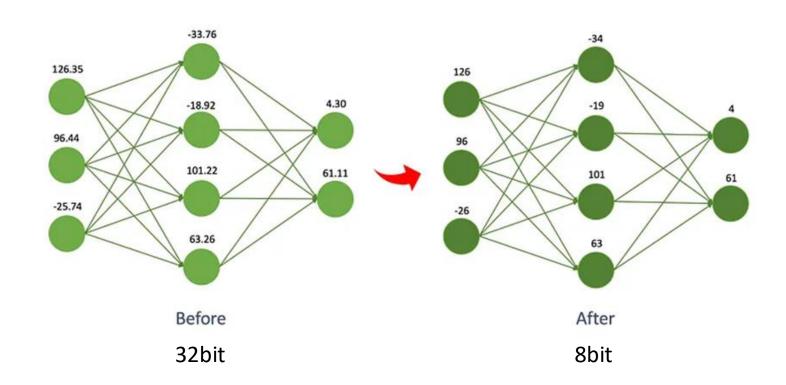
RLHF finetuning: Finetuning the model with a dataset based on human feedbacks using reinforcement learning

LoRA Finetuning



Quantization

- It is not simply rounding the numbers!
- It enables us to load huge models on our devices!
- 2bit
- 4bit
- 8bit



Libraries available for quantization

- Bitsandbytes library of Huggingface
 - https://github.com/bitsandbytes-foundation/bitsandbytes
- AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration
 - https://github.com/mit-han-lab/llm-awq
- SeeDOT:
 - normally used in the IoT field
 - Developed by Microsoft
 - https://github.com/microsoft/EdgeML/tree/master/tools/SeeDot



Now we are ready to go!

What is Amazon expecting?

A Chatbot that can interact with the customer and provide answers based on the provided data!

You may find an open-source chatbot that handles conversations very well but provides irrelevant information!

As a result, we need to customize our model!

Or in other words, finetune our model!



Let's start!

• Resources:

https://huggingface.co/spaces/Vokturz/can-it-run-llm

Model Selection:

https://huggingface.co/models

- Data Exploration
- Data Preparation
- Finetuning
- Evaluation



project

```
questions = [
f"Can you recommend some {category} shoes from {brand_name}?",
f"What are the best {category} shoes from {brand_name}?",
f"Can you show me some {category} shoes from {brand_name}?",
f"Recommend me some {category} shoes from {brand_name}.",
f"What {category} shoes do you have from {brand_name}?",
f"Show me some {category} shoes from {brand_name}?",
f"{category} shoes from {brand_name}?",
```

```
answers = [
  f"Here are some {category} shoes from {brand_name}:",
  f"Here are some {category} shoes from {brand_name} that you might like:",
  f"Here are some {category} shoes from {brand_name} that you might find interesting:",
  f"Here are some {category} shoes from {brand_name} that you might find appealing:",
  f"Here are some {category} shoes from {brand_name} that you might like:",
]
```

```
rejection_answers = [
    f"Unfortunately, we do not have any {category} shoes from {brand_name} at the moment.",
    f"Sorry, we do not have any {category} shoes from {brand_name}.",
    f"No {category} shoes found for {brand_name}.",
    f"Currently, we do not have any {category} shoes from {brand_name}.",
    f"Sorry, we do not have any {category} shoes from {brand_name} at the moment.",
    f"Unfortunately, we do not have any {category} shoes from {brand_name}.",
    f"Currently there are no {category} shoes available from {brand_name}.",
]
```

Finetuning Data Preparation

So far, we have tabular data. To give a conversational taste to the communication between customers and the chatbot, we should transform our dataset into a textual conversational dataset.

Deploying the LLMs

- Transformers library is unfortunately slow, but it has a great programming interface and is great for finetuning.
- As a result, for deploying LLMs, it is highly recommended to use libraries like **vLLM** or **TGI**.
- The following web pages contain rich content about these libraries.
 - vLLM:
 - How does vLLM optimize the LLM serving system?
 - vLLM Doc
 - <u>TGI</u>:
 - Deploying Large Language Models With HuggingFace TGI

