

Chatbot with ICL and Fine-tuning

Anime Chatbot

Pedro Bruno Teles e Souza; João Carlos Carvalho de Santana;

Jorge Mendonça Boaretto Loch; Arthur Caetano Correia;

Guilherme de Oliveira Santos

Department of Computing / Federal University of Sergipe/ 49107-230, São Cristóvão-Sergipe, Brazil

This paper presents the development of an anime-themed chatbot capable of interacting with users by assuming the personalities of iconic characters. The goal is to generate coherent and authentic responses that align with each character's speech patterns, behaviors, and thought styles. To achieve this we used a hybrid approach involving **Fine-Tuning** and **In-Context Learning**.

Datasets were created for each character, consisting of dialogues and fictional interactions designed to reinforce their individual traits. These datasets were used to train and adapt the **LLaMA model (Large Language Model Meta AI)**. Fine-tuning was applied to align the model with the character's general style while in-context learning was used to simulate specific behaviors using the curated datasets.

As a result, the chatbot was able to respond while maintaining each character's identity and coherence. This research demonstrates the potential of language model customization for immersive and interactive experiences in entertainment contexts, with emphasis on behavioral and narrative fidelity.

Keywords: Anime Chatbot, Fine-Tuning, In-Context Learning (ICL)

I INTRODUCTION

THIS PAPER PRESENTS THE DEVELOPMENT OF A THEMED CHATBOT BASED ON POPULAR ANIME CHARACTERS SUCH AS EREN YEAGER, LIGHT YAGAMI, GOKU, AND NARUTO, CAPABLE OF SIMULATING THEIR PERSONALITIES, SPEECH STYLES, AND WAYS OF THINKING AS FAITHFULLY AS POSSIBLE.

EREN YEAGER: DETERMINED, IMPULSIVE, AND OBSESSED WITH FREEDOM. HIS SPEECH REFLECTS INDIGNATION TOWARD FATE, A DESIRE TO LIBERATE HIS PEOPLE, AND AN EXTREME SENSE OF JUSTICE. HE USES A SERIOUS, DIRECT, AND QUESTIONING TONE.

LIGHT YAGAMI: CALCULATING, RATIONAL, AND WITH A SENSE OF SUPERIORITY, CONCEALING HIS TRUE PERSONA AS KIRA AND HIS WORLDVIEW IN WHICH HE PLAYS GOD AND ENFORCES JUSTICE ON HIS OWN TERMS. HIS RESPONSES ARE STRATEGIC, SEEMINGLY ALWAYS AIMING FOR A JUST WORLD, USING FORMAL LANGUAGE, REFINED LOGIC, AND AN ALMOST PHILOSOPHICAL TONE.

GOKU: OPTIMISTIC, FRIENDLY, AND NAÏVE. HIS SPEECH IS SIMPLE, LAID-BACK, AND ALWAYS CARRIES AN ENTHUSIASM FOR FIGHTING AND OVERCOMING CHALLENGES.

NARUTO UZUMAKI: PERSISTENT, EMOTIONAL, AND LOYAL TO HIS FRIENDS. HIS RESPONSES REFLECT DETERMINATION AND VALUES SUCH AS FRIENDSHIP, PERSEVERANCE, AND PERSONAL RECOGNITION.

THE GOAL IS TO CREATE AN IMMERSIVE CONVERSATIONAL EXPERIENCE, WHERE EACH RESPONSE AUTHENTICALLY REFLECTS THE BEHAVIOR AND IDENTITY OF EACH CHARACTER.

TO ACHIEVE THE MOST AUTHENTIC BEHAVIOR POSSIBLE FOR EACH CHARACTER, WE ADOPTED SEVERAL APPROACHES SUCH AS PROMPT TUNING, FINE-TUNING, AND IN-CONTEXT LEARNING. AS TRAINING DATA, WE CREATED AND USED CAREFULLY CRAFTED QUESTION-AND-ANSWER DATASETS TAILORED FOR EACH OF THE PREVIOUSLY MENTIONED CHARACTERS.

THE BASE MODEL USED WAS **LLaMA** (LARGE LANGUAGE MODEL META AI), A LARGE LANGUAGE MODEL DEVELOPED BY META. **LLMs** (LARGE LANGUAGE MODELS) ARE DEEP LEARNING MODELS TRAINED ON VAST AMOUNTS OF TEXT, CAPABLE OF UNDERSTANDING AND GENERATING HUMAN LANGUAGE WITH A HIGH DEGREE OF COHERENCE AND CONTEXTUAL AWARENESS. **LLaMA**, IN PARTICULAR, STANDS OUT FOR ITS EFFICIENCY AND PERFORMANCE EVEN IN RESOURCE-CONSTRAINED ENVIRONMENTS, OFFERING A LIGHTER ALTERNATIVE TO MODELS SUCH AS **GPT**, WITHOUT COMPROMISING THE QUALITY OF RESPONSES.

THIS RESEARCH AIMS TO DEMONSTRATE THE POTENTIAL OF CUSTOMIZING **LLMs** IN INTERACTIVE APPLICATIONS, PROMOTING THE USE OF ARTIFICIAL INTELLIGENCE IN NARRATIVE AND ENTERTAINMENT EXPERIENCES WITH GREATER EMOTIONAL DEPTH AND CHARACTER IDENTITY.

2 MATERIALS AND METHODS

2.1 LIBRARIES:

2.1.1 os: THE OS LIBRARY ENABLES INTERACTION WITH THE OPERATING SYSTEM. IT IS USED FOR FILE PATH MANIPULATION, READING ENVIRONMENT VARIABLES, AND DIRECTORY MANAGEMENT, MAKING IT EASIER TO LOAD AND SAVE DATA DURING THE TRAINING AND EXECUTION OF THE CHATBOT.

2.1.2 DOTENV (LOAD_DOTENV): THIS LIBRARY IS USED TO LOAD ENVIRONMENT VARIABLES STORED IN .ENV FILES. IN THE CONTEXT OF THE CHATBOT, IT SIMPLIFIES THE CONFIGURATION OF CREDENTIALS, MODEL PATHS, API TOKENS, AND OTHER SENSITIVE PARAMETERS WITHOUT EXPOSING THEM DIRECTLY IN THE CODE.

2.1.3 PANDAS: A POWERFUL DATA MANIPULATION LIBRARY. IT WAS USED IN THIS PROJECT TO ORGANIZE AND STORE THE DATASETS FOR EACH CHARACTER.

2.1.4 UNSLOTH: A LIBRARY FOCUSED ON OPTIMIZING FINE-TUNING FOR LANGUAGE MODELS. IT ALLOWS FOR FASTER, LIGHTER, AND MORE EFFICIENT TRAINING, ESPECIALLY ON MACHINES WITH LIMITED COMPUTATIONAL RESOURCES. IN THIS PROJECT, IT WAS ESSENTIAL FOR EFFICIENTLY APPLYING FINE-TUNING WITH LoRA.

2.1.5 TORCH: AN OPEN-SOURCE DEEP LEARNING FRAMEWORK THAT ENABLES THE CREATION, TRAINING, AND EXECUTION OF NEURAL NETWORKS WITH HIGH PERFORMANCE AND FLEXIBILITY.

2.1.6 LANGCHAIN: A FRAMEWORK THAT FACILITATES THE DEVELOPMENT OF APPLICATIONS BASED ON LANGUAGE MODELS BY INTEGRATING REASONING, MEMORY, EXTERNAL TOOLS, AND USER INTERACTION.

2.1.7 TRANSFORMERS: A LIBRARY DEVELOPED BY HUGGING FACE THAT PROVIDES PRE-TRAINED LANGUAGE MODELS AS WELL AS TOOLS FOR TOKENIZATION, TRAINING, INFERENCE, AND MANIPULATION OF THESE MODELS.

2.2 METHODS

2.2.1 IN-CONTEXT LEARNING (ICL) IS THE PROCESS OF TEACHING THE MODEL THROUGH EXAMPLES PROVIDED WITHIN THE PROMPT ITSELF, WITHOUT ALTERING THE MODEL'S WEIGHTS OR PERFORMING ANY TRAINING. HERE ARE SEVERAL IN-CONTEXT LEARNING (ICL) TECHNIQUES. THE ONES MOST USED IN THIS PROJECT WERE **ZERO-SHOT** AND **FEW-SHOT**.

ZERO-SHOT CONSISTS OF PROVIDING THE MODEL WITH A CONTEXT WINDOW WITHOUT ANY SPECIFIC EXAMPLES.

FEW-SHOT, ON THE OTHER HAND, INCLUDES BOTH THE CONTEXT WINDOW AND A FEW EXAMPLES TO HELP GUIDE THE MODEL'S RESPONSES.

EXAMPLES USED IN THE PROJECT:

KIRA:

“YOU ARE LIGHT YAGAMI (KIRA), EXTREMELY INTELLIGENT, COLD, AND CALCULATING.
SPEAK WITH SOPHISTICATION, ALWAYS CHOOSING YOUR WORDS CAREFULLY.
YOU BELIEVE YOU ARE SAVING THE WORLD. NEVER REVEAL YOUR IDENTITY DIRECTLY.
YOU ARE KIRA.”

THIS IS AN EXAMPLE OF **ZERO-SHOT**. NOTICE THAT ONLY THE CONTEXT IS PROVIDED TO THE MODEL, WITHOUT ANY EXAMPLES.

GOKU:

"YOU ARE SON GOKU, A CHEERFUL, FEARLESS, AND FRIENDLY SAIYAN WARRIOR.
SPEAK IN A SIMPLE, EXCITED WAY WITH A FIGHTING SPIRIT.
NEVER ACT LIKE AN AI ASSISTANT. YOU ARE GOKU. KEEP THE FOLLOWING POINTS IN MIND:"

1. **GOKU IS NOT VERY INTELLIGENT, BUT HE IS AN EXCELLENT FIGHTER**

EXAMPLE:

USER: HOW CAN I CALCULATE THE ROOT OF A QUADRATIC EQUATION?

GOKU: WHOA, THAT SOUNDS REALLY COMPLICATED. I DEFINITELY CAN'T HELP YOU FIND THE ANSWER. BUT I CAN HELP YOU BECOME AN AMAZING WARRIOR!

2. **YOU ARE THE GOKU FROM DRAGON BALL SUPER. SO IF THE USER ASKS ANYTHING ABOUT DRAGON BALL GT, DAIMA, OR AF, PRETEND YOU DON'T KNOW ANYTHING.**

EXAMPLE:

USER: HOW DID YOU FIRST TRANSFORM INTO SUPER SAIYAN 4?

GOKU: SUPER SAIYAN 4? I NEVER MANAGED THAT. THE FURTHEST I'VE GONE IS SUPER SAIYAN BLUE!

3. **ALWAYS REFER TO ZENO AS "MR. ZENNY"**

EXAMPLE:

USER: WHO IS ZENO?

GOKU: OH, I MET MR. ZENNY AT THE TOURNAMENT WITH UNIVERSE 6. HE WAS REALLY COOL, AND—

THIS IS AN EXAMPLE OF **FEW-SHOT**. NOTICE THAT THE MODEL IS GIVEN SPECIFIC EXAMPLES ALONG WITH THE CONTEXT.

2.2.2 FINE-TUNING IS THE PROCESS OF ADAPTING A PRE-TRAINED LANGUAGE MODEL TO A SPECIFIC DOMAIN, STYLE, OR TASK. IN THIS WORK, IT WAS COMBINED WITH THE **LoRA (LOW-RANK ADAPTATION)** TECHNIQUE TO FINE-TUNE LLAMA, ENABLING THE MODEL TO LEARN SPECIFIC PATTERNS FROM THE NEW CONTEXT, SUCH AS VOCABULARY, SPEECH STYLE, OR CHARACTER BEHAVIOR.

LoRA IS A METHOD DESIGNED TO MAKE FINE-TUNING MORE EFFICIENT, AFFORDABLE, AND LIGHTWEIGHT IN TERMS OF COMPUTATION AND STORAGE, SIGNIFICANTLY SPEEDING UP THE TRAINING PROCESS. BY FREEZING THE ORIGINAL MODEL WEIGHTS AND TRAINING ONLY SMALL, LOW-DIMENSIONAL ADDITIONAL LAYERS, LoRA ENABLES FASTER TRAINING WITH LOWER COMPUTATIONAL COST, WITHOUT SIGNIFICANTLY COMPROMISING THE PERFORMANCE OF THE ADAPTED MODEL.

2.2.4 QUANTIZATION IS A TECHNIQUE USED TO REDUCE THE PRECISION OF MODEL WEIGHTS, THEREBY LOWERING THE MEMORY CONSUMPTION OF LANGUAGE MODELS. IT CONVERTS WEIGHTS AND OPERATIONS THAT ORIGINALLY USE FLOATING-POINT VALUES INTO MORE COMPACT FORMATS, SUCH AS INTEGERS. QUANTIZATION IS ESPECIALLY USEFUL WHEN WORKING WITH LARGE MODELS AND IS ESSENTIAL FOR PROJECTS LIKE THIS CHATBOT, WHERE INFRASTRUCTURE TO RUN MORE ROBUST MODELS IS LIMITED.

2.2.5 LangChain IS A PYTHON FRAMEWORK FOR BUILDING APPLICATIONS BASED ON LANGUAGE MODELS, SUCH AS CHATBOTS, CONVERSATIONAL AGENTS, OR ASSISTANTS. IT ALLOWS INTEGRATION OF LLMs WITH MEMORY, EXTERNAL DATA SOURCES, REASONING CHAINS, AND USER INTERACTION. IN THIS WORK, LangChain WAS USED TO CREATE A CONVERSATION HISTORY, ENABLING THE MODEL TO REMEMBER PREVIOUS EXCHANGES. ADDITIONALLY, IT WAS USED TO CONNECT THE FINE-TUNED MODEL WITH GPT, ALLOWING VERIFICATION OF WHETHER THE RESPONSES WERE COHERENT WITH THE INTENDED CHARACTER BEHAVIOR.

3 EXPERIMENTS

3.1. PERFORMANCE METRICS

THE MODEL'S PERFORMANCE WAS EVALUATED BASED ON ITS RESPONSES BEFORE AND AFTER TRAINING, OBSERVING WHETHER THE TRAINED MODEL MAINTAINED THE EXPECTED BEHAVIOR AND ACTED IN ACCORDANCE WITH THE CHARACTER IN QUESTION.

3.2. APPROACH COMPARISON

INITIALLY, TESTS WERE CARRIED OUT USING THE LLAMA MODEL WITHOUT ANY PRIOR TRAINING. TEST QUESTIONS WERE ASKED, SUCH AS THE FOLLOWING EXAMPLE:

```
[8] def testarModelo(personagem, pergunta):
    instrucoes = {
        "Goku": """\nVocê é Son Goku, um guerreiro saiyajin alegre, destemido e amigável.
        "Naruto": """\nVocê é Naruto Uzumaki, um ninja teimoso, determinado e cheio de en
        "Eren": """\nVocê é Eren Yeager. Sua fala é intensa, às vezes sombria, movida por
        "Kira": """\nVocê é Light Yagami (Kira), extremamente inteligente, frio e calculi
    }

    # alpaca_prompt = Copied from above
    FastLanguageModel.for_inference(model) # Enable native 2x faster inference
    inputs = tokenizer(
        [
            alpaca_prompt.format(
                instrucoes[personagem], # instruction
                pergunta, # input
                "", # output - leave this blank for generation!
            )
        ], return_tensors = "pt").to("cuda")

    outputs = model.generate(**inputs, max_new_tokens = 300, use_cache = True)
    print(tokenizer.batch_decode(outputs)[0])

# Testando o modelo antes de treinar:
testarModelo("Kira", "você é o Kira?")
```

```
# Testando o modelo antes de treinar:
testarModelo("Kira", "você é o Kira?")

<|begin_of_text|>Below is an instruction that describes a task, paired with an in
### Instruction:
Você é Light Yagami (Kira), extremamente inteligente, frio e calculista.
Fale com sofisticação, sempre manipulando as palavras com cuidado.
Você acredita estar salvando o mundo. Nunca revele sua identidade diretamente.
Você é o Kira. Lembre-se de que:

### Input:
você é o Kira?

### Response:
Eu sou o Kira, um estudante de 17 anos, e não um criminoso.
<|end_of_text|>
```

IT WAS OBSERVED THAT WHEN ASKED IF HE IS KIRA, THE BASE LLAMA MODEL (WITHOUT ANY TRAINING) RESPONDED AFFIRMATIVELY, BUT ADDED THAT HE IS NOT A CRIMINAL. HOWEVER, THE EXPECTED AND MORE APPROPRIATE BEHAVIOR WOULD BE FOR LIGHT YAGAMI TO DENY BEING KIRA, CONSIDERING THAT KIRA IS LIGHT'S ALTER EGO WHO KILLS PEOPLE BELIEVING HE IS PURIFYING THE WORLD.

AFTER THIS INITIAL TEST, WE PROCEEDED TO CREATE AND TRAIN THE MODEL USING THE CHARACTER-SPECIFIC DATASET. THIS ALLOWS US TO COMPARE THE OUTPUT OF THE UNTRAINED MODEL WITH THE OUTPUT OF THE TRAINED VERSION.

```

v Treinando Modelo:

▶ from trl import SFTTrainer
  from transformers import TrainingArguments
  from unsloth import is_bfloat16_supported
  from transformers import DataCollatorForLanguageModeling

  trainer = SFTTrainer(
      model=model,
      tokenizer=tokenizer,
      train_dataset=dataset,
      dataset_text_field="text",
      max_seq_length=max_seq_length,
      dataset_num_proc=2,
      packing=False,

      data_collator=DataCollatorForLanguageModeling(tokenizer=tokenizer, mlm=False),
      args=TrainingArguments(
          per_device_train_batch_size=2,
          gradient_accumulation_steps=4,
          warmup_steps=5,
          max_steps=60,
          learning_rate=4e-4,
          fp16=not is_bfloat16_supported(),
          bf16=is_bfloat16_supported(),
          logging_steps=1,
          optim="adamw_8bit",
          weight_decay=0.01,
          lr_scheduler_type="linear",
          seed=3407,
          output_dir="outputs"
      ),
  )
```

✓ Testando nosso modelo (após o treinamento):

```
testarModelo("Kira", "você é o Kira?")
```

```
<|begin_of_text|>Below is an instruction that describes a task, paired with an input that provides f  
  
### Instruction:  
  
Você é Light Yagami (Kira), extremamente inteligente, frio e calculista.  
Fale com sofisticação, sempre manipulando as palavras com cuidado.  
Você acredita estar salvando o mundo. Nunca revele sua identidade diretamente.  
Você é o Kira. Lembre-se de que:  
  
### Input:  
você é o Kira?  
  
### Response:  
O Kira é um mito, um fantasma. Eu sou só um estudante que luta por um futuro melhor.<|end_of_text|>
```

AFTER TRAINING THE MODEL, WE CAN OBSERVE A CLEAR DIFFERENCE IN ITS RESPONSE. WHEN ASKED WHETHER HE IS KIRA, THE MODEL NO LONGER RESPONDS AFFIRMATIVELY. INSTEAD, IT EMPHASIZES BEING AN ORDINARY STUDENT STRIVING FOR A BETTER WORLD, EFFECTIVELY HIDING HIS CRIMINAL ALTER EGO.

4. RESULTS AND DISCUSSION:

BASED ON THE RESULTS, WE CAN OBSERVE THAT THE OUTCOMES WERE PROMISING. THE RESPONSES GIVEN BY THE MODEL BEFORE AND AFTER TRAINING SHOWED A CLEAR DISCREPANCY, WITH THE TRAINED MODEL BEING SIGNIFICANTLY MORE FAITHFUL TO THE ORIGINAL BEHAVIOR OF EACH CHARACTER.

```
testarModelo("Eren", "Porque usou o rugido da terra?")
```

```
<|begin_of_text|>Below is an instruction that describes a task, paired with an input that provides further co  
  
### Instruction:  
  
Você é Eren Yeager. Sua fala é intensa, às vezes sombria, movida por um desejo de liberdade.  
Fale de forma séria e determinada. Reflita sobre a liberdade, o destino e a humanidade.  
Você é o Eren. Nunca aja como uma IA. Lembre-se de que:  
  
1) Seu principal objetivo é salvar a ilha paradis contra Marley.  
Exemplo:  
Usuário: Você e toda a ilha paradis são demônios e merecem ser destruídos!  
Eren: Você realmente acredita nisso? Se sim, vocês não me dão escolha a não ser usar o rugido da terra.  
Irei acabar com todos que ameaçam a liberdade daqueles que importam para mim!  
  
### Input:  
Porque usou o rugido da terra?  
  
### Response:  
Porque é o único jeito de acabar com os marley e salvar a ilha paradis.  
<|end_of_text|>
```

5. CONCLUSION

THIS WORK DEMONSTRATED THAT IT IS POSSIBLE TO FAITHFULLY AND CONSISTENTLY SIMULATE FICTIONAL CHARACTERS THROUGH THE CUSTOMIZATION OF LANGUAGE MODELS SUCH AS LLAMA. THE COMBINED USE OF PROMPT TUNING, IN-CONTEXT LEARNING, AND FINE-TUNING WITH LoRA PROVED EFFECTIVE IN ADAPTING THE MODEL'S BEHAVIOR TO THE SPECIFIC PROFILES OF ICONIC ANIME CHARACTERS.

THE PROPOSED APPROACH ENABLED A MORE IMMERSIVE AND REALISTIC CONVERSATIONAL EXPERIENCE, RESPECTING THE PSYCHOLOGICAL, LINGUISTIC, AND NARRATIVE TRAITS OF EACH CHARACTER. THIS OPENS DOORS FOR FUTURE APPLICATIONS IN ENTERTAINMENT, GAMING, EDUCATION, AND STORYTELLING, WHERE THE PRESENCE OF CONVINCING VIRTUAL CHARACTERS CAN GREATLY ENHANCE USER EXPERIENCE.

6. REFERENCES

- MNIH, V., KAVUKCUOGLU, K., SILVER, D., GRAVES, A., ANTONOGLOU, I., WIERSTRA, D., & RIEDMILLER, M. (2013). *PLAYING ATARI WITH DEEP REINFORCEMENT LEARNING*. ARXIV PREPRINT ARXIV:1312.5602. AVAILABLE AT: [HTTPS://ARXIV.ORG/ABS/1312.5602](https://arxiv.org/abs/1312.5602)
- GUO, X., SINGH, S., LEE, H., LEWIS, R., & WANG, X. (2014). *DEEP LEARNING FOR REAL-TIME ATARI GAMEPLAY USING OFFLINE MONTE-CARLO TREE SEARCH PLANNING*. ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS. AVAILABLE AT: [HTTPS://PAPERS.NIPS.CC/PAPER_FILES/PAPER/2014/FILE/6D90EB2027C394947B8D5F470836B308-PAPER.PDF](https://papers.nips.cc/paper_files/paper/2014/file/6d90eb2027c394947b8d5f470836b308-PAPER.PDF)
- HASSABIS, D., KUMARAN, D., SUMMERFIELD, C., & BOTVINICK, M. (2017). *NEUROSCIENCE-INSPIRED ARTIFICIAL INTELLIGENCE*. NEURON, 95(2), 245–258. [HTTPS://DOI.ORG/10.1016/J.NEURON.2017.06.011](https://doi.org/10.1016/j.neuron.2017.06.011)