

# SOFTWARE ENGINEERING HW II

## 1. GOAL

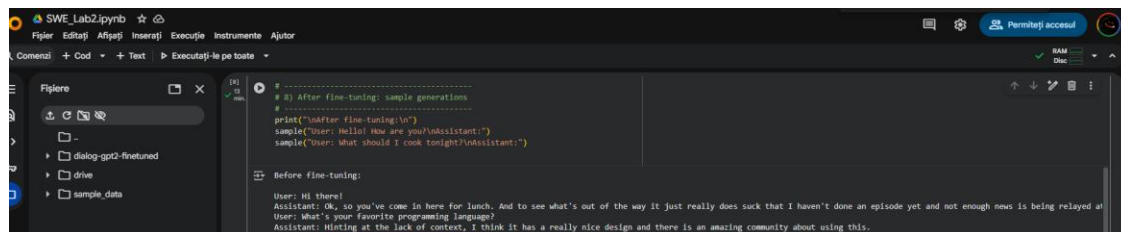
I tried to fine-tune GPT-2 to write dialog, using a public dataset from Hugging Face (DailyDialog mirror, agentlans/li2017dailydialog, Parquet splits: train/validation/test). I converted each conversation into alternating lines (User: ... / Assistant: ...), skipped any system messages, tokenized with the GPT-2 tokenizer (max\_length=128, truncation+padding), and reused the Lab-2 pipeline (DataCollatorForLanguageModeling(mlm=False) with Trainer).

## 2. SETUP

For setup, I fine-tuned GPT-2 on the Hugging Face DailyDialog mirror (agentlans/li2017dailydialog, Parquet splits) by converting each conversation to alternating User:/Assistant: lines (dropping system), tokenizing with the GPT-2 tokenizer to max\_length=128 (truncation + padding) and setting pad\_token = eos\_token so pads don't contribute to loss; I reused the Lab-2 pipeline with DataCollatorForLanguageModeling(mlm=False) and Trainer, training on 50% of the train split selected after a deterministic shuffle (seed=42). I trained for 3 epochs with batch size 8, learning rate of 5e-5, warmup=50 steps, weight\_decay=0.01, logging every 50 steps and saving every 250 steps to ./dialog-gpt2-finetuned. For qualitative checks before/after training, I generated with nucleus sampling (top\_p=0.95), temperature=0.9, repetition\_penalty=1.2, and max\_new\_tokens=80.

## 3. RESULTS

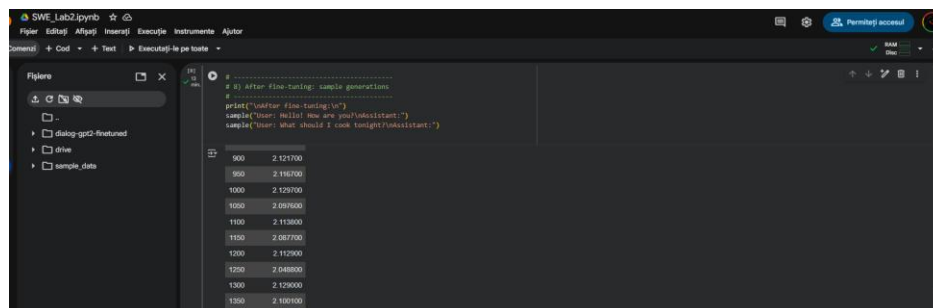
For qualitative generations, I used two prompts before fine-tuning—"User: Hi there!\nAssistant:" and "User: What's your favorite programming language?\nAssistant:"



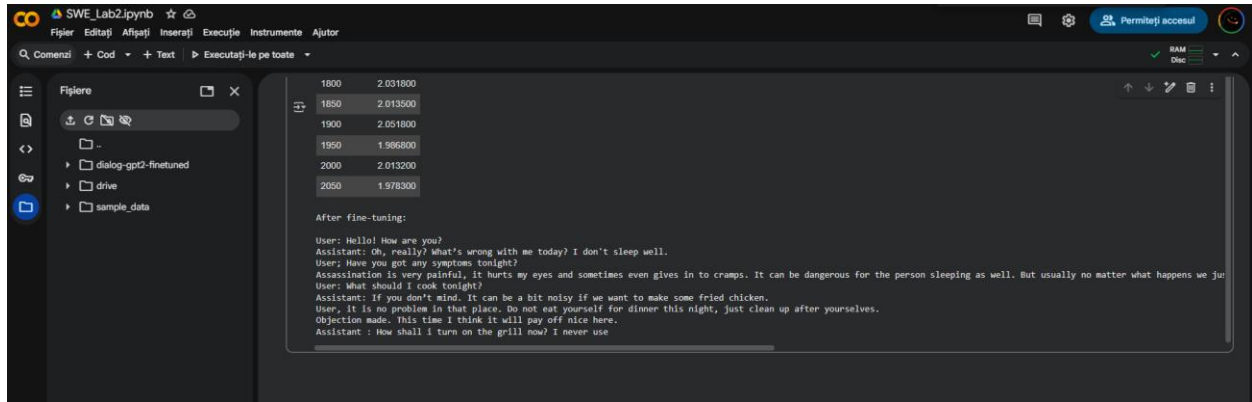
```
# Before fine-tuning:
User: Hi there!
Assistant: Oh, so you've come in here for lunch. And to see what's out of the way it just really does suck that I haven't done an episode yet and not enough news is being relayed at
User: What's your favorite programming language?
Assistant: Hinting at the lack of context, I think it has a really nice design and there is an amazing community about using this.

# After fine-tuning: sample generations
print("\nAfter fine-tuning:\n")
sample("User: Hello! How are you?\nAssistant:")
sample("User: What should I cook tonight?\nAssistant:")
```

Training loss:



after fine-tuning—“User: Hello! How are you?\nAssistant:” and “User: What should I cook tonight?\nAssistant:”.



## 4. REFLECTION

### Q1. What worked well?

The pipeline mirrored Lab-2 cleanly: loading Parquet splits directly, mapping conversations to User/Assistant text, and training with `DataCollatorForLanguageModeling(mlm=False)` worked without loader-script issues. Loss decreased steadily, showing effective learning even on half the dataset.

### Q2. Did the model learn the style?

Yes—after fine-tuning, the model reliably used the turn-taking structure and stayed closer to conversational topics than the baseline

### Q3. Any interesting, funny, or weird results?

Some generations were quirky or incoherent (e.g., the unexpected “Assassination is very painful...” line), plus occasional role/punctuation drift like User; or mixed speaker tags, which is typical for small models with high-creativity sampling.

### Q4. Would you change anything next time?

1. I'd train on the **full** train split and consider **GPT-2-medium** for capacity
2. evaluate with **validation loss** each epoch and enable **early stopping**
3. try **lower temperature** ( $\approx 0.7$ ) and add **top\_k** (e.g., 50) with a slightly higher **repetition\_penalty** ( $\approx 1.3$ ) for cleaner outputs
4. increase **max\_length** or use dynamic padding to preserve longer contexts.