# NTU ADL Homework 3, 2023 Fall.
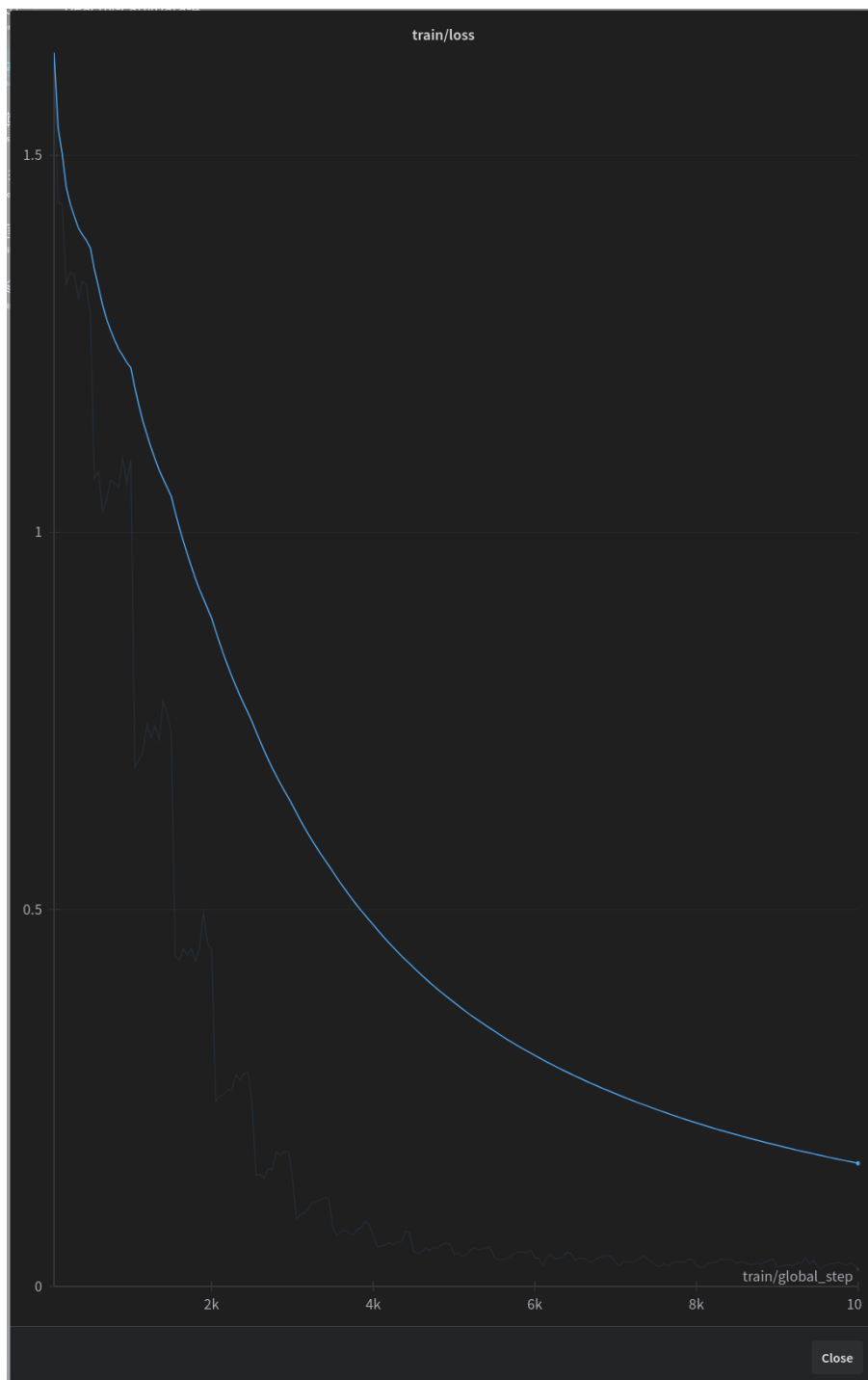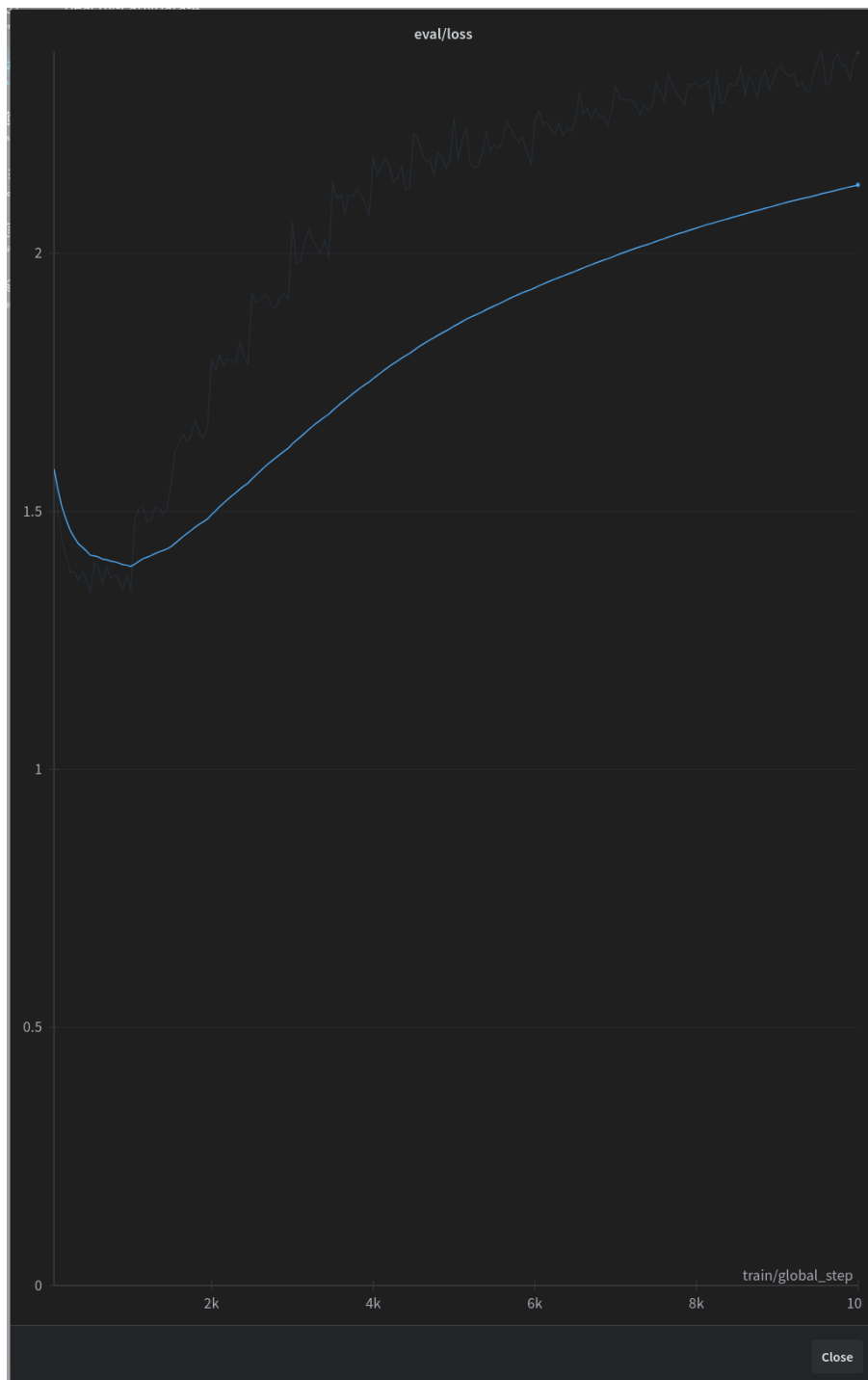
## LLM Tuning

### Describe & Performance:

- I used all the training data while training, at first, I didn't think too much, just normally trained it with all training data, it took 10k training steps, at that time, I decided to evaluate every 500 steps, followed by this thought, something happened…

- **Above picture is my training loss curve, it looks pretty well, right? Let me show you another figure...**



- **This picture is my evaluation loss curve, it looks like a NIKE logo, however, it's not a good symbol...**

- **Combined these two figures, it shows that my adapter was overfitting after about 1.1k steps, I've found this situation at that time, however, I just decided to finishing this training period, and it took me about 110 hours.**

- **Following picture is my training hyper-parameters.**

```
python train.py \
    --model_name_or_path /home/cjweibb/workspace/2023_NTUADL/testSpace/Taiwan-LLM-7B-v2.0-chat \
    --output_dir ./output/test\
    --logging_steps 50 \
    --save_strategy steps \
    --data_seed 42 \
    --save_steps 500 \
    --save_total_limit 40 \
    --evaluation_strategy steps \
    --eval_dataset_size 1024 \
    --max_eval_samples 1000 \
    --per_device_eval_batch_size 1 \
    --max_new_tokens 32 \
    --dataloader_num_workers 1 \
    --group_by_length \
    --logging_strategy steps \
    --remove_unused_columns False \
    --do_train \
    --do_eval \
    --lora_r 64 \
    --lora_alpha 16 \
    --lora_modules all \
    --double_quant \
    --quant_type nf4 \
    --bits 4 \
    --warmup_ratio 0.03 \
    --lr_scheduler_type constant \
    --gradient_checkpointing \
    --dataset /home/cjweibb/workspace/2023_NTUADL/testSpace/data/train.json \
    --source_max_len 128 \
    --target_max_len 128 \
    --per_device_train_batch_size 1 \
    --gradient_accumulation_steps 16 \
    --learning_rate 0.0003 \
    --adam_beta2 0.999 \
    --max_grad_norm 0.3 \
    --lora_dropout 0.01 \
    --weight_decay 0.0 \
    --seed 42 \
    --dataset_format alpaca \
    --report_to wandb
```

- **Following picture is my model's performances on public_test set:**

```
Loading checkpoint shards: 100%|                                                          [7/12]
                                   | 2/2 [00:05<00:00,  2.72s/it]
100%|
                                   | 250/250 [03:24<00:00,  1.22it/s]
Qlora-finetuned Inference results:
Mean perplexity: 3.4397967281341555
(ntu_adl) cjweibb@mvclab-cjw:~/workspace/2023_NTUADL/testSpace$ bash eval.sh
Loading checkpoint shards: 100%|
                                   | 2/2 [00:05<00:00,  2.99s/it]
100%|
                                   | 250/250 [03:29<00:00,  1.19it/s]
Qlora-finetuned Inference results:
Mean perplexity: 3.6434353761672975
(ntu_adl) cjweibb@mvclab-cjw:~/workspace/2023_NTUADL/testSpace$ bash eval.sh
Loading checkpoint shards: 100%|
                                   | 2/2 [00:05<00:00,  2.73s/it]
100%|
                                   | 250/250 [03:33<00:00,  1.17it/s]
Qlora-finetuned Inference results:
Mean perplexity: 5.025580849647522
(ntu_adl) cjweibb@mvclab-cjw:~/workspace/2023_NTUADL/testSpace$ bash eval.sh
Loading checkpoint shards: 100%|
                                   | 2/2 [00:05<00:00,  2.67s/it]
100%|
                                   | 250/250 [03:32<00:00,  1.17it/s]
Qlora-finetuned Inference results:
Mean perplexity: 17.62995908117294
```

- **As you can see, there are four perplexity scores listed from top to bottom, corresponding to ordered checkpoints: "500 steps, 1500 steps, 2000 steps, 10000 steps." The results confirm what I mentioned earlier about "overfitting," and my best performance is at 500 steps with a perplexity score of 3.439.**

# LLM Inference Strategies

## Zero-Shot, Few-Shot:

### Zero-Shot inference :

- **For Zero-Shot inference, I just directly inply "Taiwan-LLM-7B-v2.0-chat" model on public_test dataset with max_length=2048, below picture shows my custom prompt, there is a little different with original one given by TAs.**

```
def get_prompt(instruction: str) -> str:
    '''Format the instruction as a prompt for LLM.'''
    return f"你是人工智慧助理，以下是用戶和人工智能助理之間的對話。你要對用戶的問題提供有用、安全、詳細和禮貌的回答。\n請回答下面的問題。\nUSER: {instruction} ASSISTANT:"
```

- **Following screenshot is my mean perplexity score of Zero-Shot Inferencing:**

Mean perplexity: 5.452533513069153

### Few-Shot inference :

- **For Few-Shot inference, I just directly inference "Taiwan-LLM-7B-v2.0-chat" model on public_test dataset with max_length=2048, and I have done an experiment, following are my experiment steps:**

  - **1. Randomly Select 10 different datas from public_test dataset as Few-Shot prompts' instructions and outputs.**

  - **2. Remove datas that have been selected as Few-Shot prompts.**

  - **3. Modify the prompt as shown below.**

```
你是人工智慧助理，以下是用戶和人工智能助理之間的對話。你要對用戶的問題提供有用、安全、詳細和禮貌的回答。
以下為幾個範例：
用戶：左衛大將軍宇文述勸說煬帝選挑選幾驍名精銳騎兵突圍齣去，納言蘇威說： 據守城池我方還有餘力，而輕騎則是對方的長處，陛下是萬乘之主，怎麼可以輕易行動？
幫我把這句話翻譯成文言文 助理：左衛大將軍宇文述勸帝簡精銳數驍騎潰圍而齣，納言蘇威曰： 城守則我有馀力，輕騎乃彼之所長，陛下萬乘之主，豈宜輕動！
用戶：盧延貴非常害怕，急急忙忙逃走，他卻連連呼籲： 不要害怕，我是個人！
這句話在中國古代怎麼說： 助理：見人即行起而來，延貴懼而卻走。此物連呼： 無懼，吾乃人也。
用戶：秦領是都人，從江夏郡尉升任南陽太守，上任時經過宜城城內，看見一傢朝東的房子，他停車觀看，說：這個住人的地方可以築墓。
翻譯成文言文： 助理：韻，郡人也，以江夏郡尉韵為南陽太守。徑宜城中，見一傢東嚮，韻住車視之，曰：此居處可作塚。
用戶：不用兵而巧設腦局取勝，他就是這樣做的。南蠻
翻譯成古文： 助理：而驟信譬懍，不假兵以詐勝，斯之謂也。南蠻
用戶：這是二不可。
這句話在古代怎麼說： 助理：不可二。
用戶：有的說不吃葷腥，不近女色，遇上這樣的好藥就能成為地仙。
把這句話翻譯成文言文： 助理：或不葷血，不色欲，遇之必能降真為地仙矣。
用戶：四人祝壽完畢，趨步離開。
這句話在古代怎麼說： 助理：四人為壽已畢，趨去。
用戶：六月，西瓜、甜瓜、蓮子、鑾瓜。
翻譯成白話文： 助理：六月，獻西瓜、甜瓜、蓮子、鑾瓜。
用戶：翻譯成現代文。
上遣使者諭釋，冠蓋相望。 助理：唐僖宗接連不斷地派遣使臣前往李剋用處進行規勸解釋。
用戶：孝文貞皇後林氏，是平原人。
幫我把這句話翻譯成文言文 助理：孝文貞皇後林氏，平原人也。
請根據以上的範例進行回答。
用戶：翻譯成文言文：
於是，廢帝讓潘麘之的堂侄、直將軍潘攸之賜潘麘之毒藥，命潘麘之自殺。 助理：
```

  - **4. Repeat above steps for 5 times.**

- **Following figures are my Few-Shot inference perplexity scores:**

Few-Shot Inference #0 results:
Samples: [{'id': '7b7ead70-1353-433f-a59f-7704594cce59', 'instruction': '希望您以後留意，不要再齣這樣的事，你的小女兒病就會好。\n這句話在古代怎麼說：', 'output': '以後幸長官留意，勿令如此。'}, {'id': 'b4c07b51-3bd7-4512-a196-f2db6e90c371', 'instruction': '十二月，甲辰，憲宗加封山南東道節度使於為同平章事。\n幫我把這句話翻譯成文言文', 'output': '十二月，甲辰，加山南東道節度使於國同平章事。'}, {'id': '6f0e2f00-beaa-45ac-baad-9c8072fc4edf', 'instruction': '將下麵句子翻譯成古言文：\n免除冬年租賦、丁口力役一半。', 'output': '免今載租、庸半。'}, {'id': '0020ed81-89b9-4ad4-9860-eb9ebaf2b139', 'instruction': '將下麵句子翻譯成現代文：\n鄧池沼繼紳，波妾後又。', 'output': '愛卿乃是池沼中的曾紳，岸渠褻的俊倓。'}, {'id': '7c700d23-3856-4f3d-87c9-f6a52e207cd2', 'instruction': '將下麵句子翻譯成文言文：\n唯懷英即康慎典，因為避諱後嫌帝司王殊友真的名字而改為康懷英。', 'output': '懷英即懷貞也，避帝名改甚。'}, {'id': '72335b5a-ad07-4255-8a9a-0469ee1dd9f7', 'instruction': '文言文翻譯：\n鄧丹尼以陽城之戰為彥郵所敦，結之。彥郵曰：臣當時惟知為蔡主竭力，今日死生惟命。', 'output': '答案：契丹主因降城之戰致符彥郵扣敗，追附司彥郵。彥郵說：臣當時雙知為蔡主蠣盡全力，今日死生聽所決定。'}, {'id': '513d08b5-40a6-4896-98fb-3ee9adbc13c5', 'instruction': '翻譯成文言文：\n癸未，有黑氣如龍，齣現在宮殿裏。\n答案：', 'output': '癸未，有黑氣如龍，見於殿內。'}, {'id': '24a56c5d-2ed6-4f85-a164-b8e21d31d653', 'instruction': '文言文翻譯：\n既練而齣，則已未練而反，則斯；既練而反，則遂之·再期之喪·三年也；期之喪，二年也。', 'output': '答案：作為嬸姆，在為公婆服喪期間被丈夫休棄，因為恩斷義絕，所以也就不再繼續為公婆服喪。'}, {'id': '6944b5b5-3dcc-49b3-9a1a-13914419f9f4', 'instruction': '六月，太祖上籤推薦丁會為瀘州節度使，是因為李罕之病重。\n翻譯成古文：', 'output': '六月，帝籤丁會為瀘州節度使，以李罕之疾亟故也。'}, {'id': '5cec4a4c-a77d-42e4-ab26-74fe9bec1a05', 'instruction': '文言文翻譯：\n何則？', 'output': '答案：什麼緣故呢？'}]
Mean perplexity: 5.778133679926396

Few-Shot Inference #1 results:
Samples: [{'id': '920cdace-4403-4982-be69-d9f32eaa186a', 'instruction': '至順元年，拜西行動監察禦史，建明八事：一曰正君道，二曰結人心，三曰崇禮讓，四曰正綱紀，五曰審銓衡，六曰勵孝行，七曰紓民力，八曰修軍政。\n翻譯成現代文。', 'output': '至順元年，拜西行動監察禦史，他提齣瞭八條建議：一，端正君道；二，團結人心；三，崇尚禮讓；四，整頓綱紀；五，審查官吏推舉製度；六，鼓勵孝行；七，減輕百姓負擔；八、修整軍政。'}, {'id': '0e79a40c-b55f-40ae-a40a-b110273a68cf', 'instruction': '翻譯成文言文：\n侯是迎拜車下，王丹下車答謝。\n答案：', 'output': '侯是迎拜車下，王丹下車答謝。'}, {'id': 'a11c6867-8e71-406f-8c14-94b41d473fbe', 'instruction': '能廕侍政，此謂正紀。\n翻譯成現代文：', 'output': '能守信於民，這叫作端正綱紀。'}, {'id': '05777600-2826-4681-a41b-3590c83bad2c', 'instruction': '取吾犯者，讀以故違敕議。\n翻譯成白話文：', 'output': '激敢有違犯者，讀以故意違抗罪處理。'}, {'id': 'cd27e23d-b707-412c-96f2-e59529ac7b4c', 'instruction': '文言文翻譯：\n通、容歸漢。\n output': '答案：夏侯淵、韓容迴到漢朝。'}, {'id': '7d74d352-a907-4357-b7a2-5c6c451186e9', 'instruction': '皇上對高潁說：晉王妃有種置附身，說曾王必欲不天下。你說怎麼辦？\n翻譯成古文：', 'output': '謂頭曰：晉王妃有神憑之，言王必有天下，若之何？'}, {'id': '882d0d87-7e22-462e-b5d6-0bfc7fe452af', 'instruction': '又以為黃門侍郎，令長史劉孝勝深布心腹，便工畫《陳平渡河歸漢圖》以邊之。\n把這句翻譯成現代文。', 'output': '蕭紀又任命劉璋為黃門侍郎，令長史劉孝勝安排心腹之人，讓畫師畫《陳平渡河歸漢圖》送給劉璋。'}, {'id': 'a4bd0ee3-598b-4914-a70b-2f1061139cff', 'instruction': '將下麵句子翻譯成現代文：\n農乃建之，以定備注，朝廷是之。', 'output': '常景於是建正例，以定朝偏，朝廷肯定瞭他的做法。'}, {'id': '7192faf3-bd87-4489-9df7-593f410fca90', 'instruction': '地勢高而乾燥，他們築土為城牆，架木建房子，用土覆蓋在房屋上麵。\n翻譯成文言文：', 'output': '其地高燥，築土為城，架木為屋，土覆其上。'}, {'id': '513d08b5-40a6-4896-98fb-3ee9adbc13c5', 'instruction': '翻譯成文言文：\n癸未，有黑氣如龍，齣現在宮殿裏。\n答案：', 'output': '癸未，有黑氣如龍，見於殿內。'}]
Mean perplexity: 5.704349031051

Few-Shot Inference #2 results:
Samples: [{'id': '7f9a3c0a-83f9-4504-99ce-406a7f755288', 'instruction': '己巳，詔廕建官招集衆兀氏。\n翻譯成現代文：', 'output': '初六，下詔停止派遣官員招聚衆兀兒人。'}, {'id': '84184220-3543-4109-96a2-07a08731358d', 'instruction': '將下麵句子翻譯成文言文：\n四年，隆金符，任命也為蘭、蓬等處都元帥府參謀。', 'output': '四年，隆金符，為蘭蓬等處都元帥府參議。'}, {'id': '20bc4bb5-5f51-41fb-83cc-4f5c36dd14d1', 'instruction': '翻譯成現代文：\n上宮中大妻，裸婦人而親之，主後於戶庭舞蠣·上殿日：外接案臣。', 'output': '朝夜在宮中大擺宴席，命婦女脫光衣服，讓大傢欣賞，皇後用扇子搧往調戲。明帝大怒說：真是褻傢的褻麵相！'}, {'id': 'a11c6867-8e71-406f-8c14-94b41d473fbe', 'instruction': '能廕侍政，此謂正紀。\n翻譯成現代文：', 'output': '能守信於民，這叫作端正綱紀。'}, {'id': '1677d468-f133-439c-8c85-59ccd2a71ae1', 'instruction': '宋敝公罕之喪，字彌俊突。\n這句話在古代怎麼說：', 'output': '宋厭公罕之喪，字彌俊突。'}, {'id': '4a16bd5b-8bde-4c02-a601-f764854eba0c', 'instruction': '將下麵句子翻譯成文言文：\n隴西太守張紆以權宜之計將號吾故興退去，羌人立刻就解散繼，各自迴到原先生活的地方，迷吾繼退到黃河北邊的歸義城居住。', 'output': '隴西太守張紆權宜放遣，而羌衆驚歎，各歸故地·迷吾退居河北歸義城。'}, {'id': 'd48e04f4-6fab-47ec-b11c-61062810981a', 'instruction': '世祖說：你甚認為夷炎比莫孝更賢嗎？\n這句話在中國古代怎麼說：', 'output': '帝曰：汝以夢炎賢於李耶？'}, {'id': '8fbcd1f8-6304-47ea-a5d2-2961d18bf640', 'instruction': '微自艺老病求解，凡又與書答曰：曹玉蕃鼠，自立為帝，是猶士羅劉狗之有任·\n翻譯成現代文：', 'output': '杜貴自稱年老多病，請求辭官，誡又給他寫信說：曹玉蕃起兵，自立為帝，是猶士羅劉狗之有任。'}, {'id': 'bfce72e4-39a4-419d-aba0-9b5292cd7fb5', 'instruction': '翻譯成現代文：\n乃遣知在進奉詔書，諭以深旨，乃承昭副茲麾下，擁彼舊封。\n答案：', 'output': '於是派遣操知古進道地奉著詔書，以懇切的旨意曉諭他，而且任命李承昭在他的部下擁任副職，鎮擁他的舊封之地。'}]
Mean perplexity: 5.576618615289529

Few-Shot Inference #3 results:
Samples: [{'id': '6b39c935-4ffe-4cec-b813-926904a8863c', 'instruction': '翻譯成現代文：\n齣後彼嚮憂懼，與母常氏求托女巫，禱厭孝文疾不起：一旦得如文明太後輔少主稱命者，賞願不費。\n答案：', 'output': '此後，皇後逐漸感到憂慮和恐懼，於是便和母親常氏一起求托女巫，祈禱用巫術讓孝文帝一病不起，如果有朝一日能像文明太後那樣輔佐少主臨朝受命的話，實願無既。'}, {'id': 'db35b753-351d-4314-9323-f907808aa7b7', 'instruction': '吾兵以賽轉，往無不剋，鳥 居此物，使暴除百姓故！\n把這句話翻譯成現代文。', 'output': '我軍以賽超事責任直前，攻無不剋，何必勤用這像夥，讓他暴虐殘害百姓呢！'}, {'id': 'f23e2444-0ec8-413b-9f3e-87a8a431cb96', 'instruction': '隆下諒暗始闢，大臣皆有願願。\n翻譯成現代文：', 'output': '隆下賜居喪初，大臣們都懷有非份的企圖。'}, {'id': 'e0614eb6-e140-48da-a7f3-699a3eb77b66', 'instruction': '但是，祭祀中的登堂分食品、顯禮獻酒、飲酒解，這些事一定要由國君的嫡長子來做。\n把這句話翻譯成文言文：', 'output': '其在宗廟之中，則如外朝之位·宗人授事，以嚮以官。'}, {'id': 'a94dd344-261c-48ba-bf6c-d915043f2332', 'instruction': '翻譯成現代文：\n坤日：新天子將兵二十年，位至大總督，所領精兵三十萬。天師人事，其可得達？\n答案：', 'output': '姚坤說：新天子手兵二十年，官做到大總督，率領的精兵三十萬·天師人和，怎達得嗎？'}, {'id': 'dd3f60d2-a734-41cd-ac81-d2cad80835bd', 'instruction': '將下麵句子翻譯成現代文：\n李柏慮愚恐己鴝解聯，推構委節，歸之有在，自古偵逆未有若斯之甚者。', 'output': '零柏本應愚處恐愚，他卻擁自為己解脫，推脫罪責，歸罪於人，自古以來像越錯埒之人沒有像他這樣嚴重的。'}, {'id': 'be746664-fce7-4baa-ac15-3dac0db41fa5', 'instruction': '將下麵句子翻譯成現代文：\n尋除兗州剌史。', 'output': '不久又授為兗州剌史。'}, {'id': '4a16bd5b-8bde-4c02-a601-f764854eba0c', 'instruction': '將下麵句子翻譯成文言文：\n隴西太守張紆以權宜之計將號吾故興退去，羌人立刻就解散繼，各自迴到原先生活的地方，迷吾繼退到黃河北邊的歸義城居住。', 'output': '隴西太守張紆權宜放遣，而羌衆驚歎，各歸故地·迷吾退居河北歸義城。'}, {'id': '3cc156cd-bf21-4774-bf93-a7fb55e644cf', 'instruction': '將下麵句子翻譯成現代文：\n建武帝蕭鸞從從他做皇帝，沿襲陳舊的一套做法，當時流行風氣不好文學，輔校羊祖沒有瞭，學校設施簡陋，先前行事規範再也難找，劉跋承繼馬融、鄭玄之後，一個時期學生門徒把他看作是師長楷模·\n翻譯成現代文：', 'output': '建武鹽立，因循舊緒，時不好文，輔相無術，學校雖設，前軌難追·劉輩孤高，鄭之後，一時學徒立為師範。'}]
Mean perplexity: 5.3158752972881

Few-Shot Inference #4 results:
Samples: [{'id': '6af2a622-d9ae-49db-89b5-d56f588df4ec', 'instruction': '翻譯成現代文：\n江陵陷，周武帝尤重之，授太學博士。\n答案：', 'output': '江陵淪陷後，周武帝特尊器重他，命他為太學博士。'}, {'id': '1dc6a04e-ba16-4f2a-844e-0dce2176462d', 'instruction': '五月甲申，京師雨雹·地震·\n把這句話翻譯成現代文。', 'output': '五月初五，京城降冰雹，發生地震。'}, {'id': 'be746664-fce7-4baa-ac15-3dac0db41fa5', 'instruction': '將下麵句子翻譯成現代文：\n尋除兗州剌史。', 'output': '不久又授為兗州剌史。'}, {'id': '9665ec3a-aea0-4a07-b230-d0b8959b2bc5', 'instruction': '翻譯成現代文：\n故日：平原之隰，衍有徒戎。\n答案：', 'output': '所以說：平原之隰，衍有徒戎。'}, {'id': 'f94d6b2d-abba-4780-9b74-332cdf78843d', 'instruction': '丙辰及攻喪，費病死。\n把製錄這句話翻譯成文言文：', 'output': '丙辰時攻喪，費病死。'}, {'id': 'a94dd344-261c-48ba-bf6c-d915043f2332', 'instruction': '翻譯成現代文：\n坤日：新天子將兵二十年，位至大總督，所領精兵三十萬。天師人事，其可得達？\n答案：', 'output': '姚坤說：新天子手兵二十年，官做到大總督，率領的精兵三十萬·天師人和，怎達得嗎？'}, {'id': 'e1b44f71-141f-46c9-b2cc-d60d072de77d', 'instruction': '翻譯成文言文：\n劉澄之著《永初山川記》說，夏水，古文以為滄浪·漁父所歌也。', 'output': '劉澄之著《永初山川記》雲：夏水，古文以為滄浪·漁父所歌也。'}, {'id': '16e5de98-dcf1-45df-9d6d-d4e9b62b5e60', 'instruction': '辭今，漢代的舊賺，屬於徐·\n這句話在中國古代怎麼說：', 'output': '辭令，漢舊賺，屬徐。'}, {'id': 'f472f1de-ef5e-4097-9 8b0-860ccd1440c9', 'instruction': '三將至，繆公素服郊迎，嚮三人哭日：孤以不用百裏俟、蹇叔言以辱三子，三子何罪乎？\n翻譯成現代文：', 'output': '三將軍迴到秦國，繆公穿戴凶服親迎於郊，體靡三人大哭說：都是我不聽百裏俟、蹇叔的話讓讓三位受職要辱，三位有什麼罪呢？'}]
Mean perplexity: 5.616713916261991

- **Upon observing the results and the samples I selected as prompts, I've noticed that if the samples we choose are too imbalanced, the performance tends to be poorer.**

- **After examining the dataset, it can be broadly categorized into two main branches: "Classical Chinese translated into Modern Chinese" and "Modern Chinese translated into Classical Chinese." The imbalance refers to a situation where the samples are skewed towards a particular direction,this will affecting the model's understanding of the prompt.**

- **On another note, I have also observed that the accuracy and level of detail in the text content impact the output in some aspects. Accuracy here pertains to the correctness of the text itself or the accuracy of the meaning of words.**

# Comparison:

- **After observing the results of Zero-Shot and Few-Shot approaches, we can see that Few-Shot does indeed achieve better results than Zero-Shot based on the provided samples.**

- **However, we also notice from the perplexity results that the impact of samples on Few-Shot Inference is significant. Therefore, if one plans to use Few-Shot Inference in the future, it may be crucial to ensure the correctness of the sample references.**

- **Lastly, incorporating the results after fine-tuning with Qlora, we can see that, compared to training a new language model directly, fine-tuning through adapters for specific application contexts is a more user-friendly and feasible approach, especially for smaller user groups.**

- **Training a language model often requires high computational resources, typically several A100 GPUs, which is not feasible for the average person. On the other hand, using adapters allows us to achieve good fine-tuned model performance at significantly lower costs. Therefore, I believe adapters offer a very promising framework.**