

報告書

1 今週の進捗

- optuna で調整したパラメータでの gpt-oss-20B の SFT
- gpt-oss-20B の Q, K, V, O 層のみの SFT
- gpt-oss-20B の継続事前学習

2 gpt-oss-20B の SFT

2.1 訓練データ

訓練データは、日本のアニメに関する知識を問う質問と出力のペアで構成される。合計 2000 件の Q&A ペアを JSONL 形式で Gemini-2.5 Pro および GPT-5 で作成した。また、1 つの質問と正解のペアに対して 32 個の似た質問と正解のペアを作成し、情報量を増やした。(Augmentation) 表 1 に、作成した訓練データセットの例を 8 件示す。

表 1: 訓練データセットの例

質問: ぱにぽにだっしゅ！の主人公は誰ですか？

出力: ぱにぽにだっしゅ！の主人公はレベッカ宮本です。

質問: ぱにぽにだっしゅ！の主人公は誰ですか？

出力: ぱにぽにだっしゅ！の主人公は、10 歳で MIT を卒業した天才少女、レベッカ宮本です。

質問: 「ぱにぽにだっしゅ！」で主人公として描かれているキャラクターは誰ですか？

出力: ぱにぽにだっしゅ！の主人公はレベッカ宮本です。

質問: 「ぱにぽにだっしゅ！」で主人公として描かれているキャラクターは誰ですか？

出力: ぱにぽにだっしゅ！の主人公は、10 歳で MIT を卒業した天才少女、レベッカ宮本です。

質問: レベッカ宮本が主人公を務める作品は何ですか？

出力: レベッカ宮本が主人公を務める作品は「ぱにぽにだっしゅ！」です。

質問: レベッカ宮本が主人公を務める作品は何ですか？

出力: 彼女が主人公の作品は「ぱにぽにだっしゅ！」です。この作品で彼女は天才少女教師として描かれています。

質問: レベッカ宮本は、どの作品で主人公として登場しますか？

出力: レベッカ宮本が主人公を務める作品は「ぱにぽにだっしゅ！」です。

質問: レベッカ宮本は、どの作品で主人公として登場しますか？

出力: 彼女が主人公の作品は「ぱにぽにだっしゅ！」です。この作品で彼女は天才少女教師として描かれています。

2.2 gpt-oss-20B のモデルアーキテクチャ

モデルは 24 層の Decoder Layer から構成されており、各層は Self-Attention ブロックと MoE 構造を持つ MLP ブロックを含む。MoE ブロックは 32 個の Expert 層を持ち、計算効率とモデル性能の両立を図っている。

```
GptOssForCausalLM(
    (model): GptOssModel(
        (embed_tokens): Embedding(201088, 2880, padding_idx=199999)
        (layers): ModuleList(
            (0-23): 24 x GptOssDecoderLayer(
                (self_attn): GptOssAttention(
                    (q_proj): Linear4bit(in_features=2880, out_features=4096, bias=True)
                    (k_proj): Linear4bit(in_features=2880, out_features=512, bias=True)
                    (v_proj): Linear4bit(in_features=2880, out_features=512, bias=True)
                    (o_proj): Linear4bit(in_features=4096, out_features=2880, bias=True)
                )
                (mlp): GptOssMLP(
                    (router): GptOssTopKRouter(
                        (linear): Linear(in_features=2880, out_features=32, bias=True)
                    )
                    (experts): GptOssExperts(
                        (gate_up_projs): ModuleList(
                            (0-31): 32 x Linear4bit(in_features=2880, out_features=5760, bias=True)
                        )
                        (down_projs): ModuleList(
                            (0-31): 32 x Linear4bit(in_features=2880, out_features=2880, bias=True)
                        )
                    )
                )
                (input_layernorm): GptOssRMSNorm((2880,), eps=1e-05)
                (post_attention_layernorm): GptOssRMSNorm((2880,), eps=1e-05)
            )
        )
        (norm): GptOssRMSNorm((2880,), eps=1e-05)
        (rotary_emb): GptOssRotaryEmbedding()
    )
    (lm_head): Linear(in_features=2880, out_features=201088, bias=False)
)
```

2.3 学習手法と対象パラメータ

SFT では、パラメータ効率の良いファインチューニング手法である Low-Rank Adaptation (LoRA) を採用した。LoRA は、元のモデルの重みを凍結したまま、特定の線形層に低ランク行列分解を利用したアダプター層を追加し、そのアダプターの重みのみを学習する。これにより、少ない計算リソースで大規模モデルの挙動を適応させることが可能となる。

本研究では、モデルの挙動に大きく影響を与える以下の層を学習対象 (target_modules) として設定した。

- **q_proj, k_proj, v_proj, o_proj:** これらは Self-Attention 機構を構成する主要な線形層である。Query, Key, Value, Output の射影をそれぞれ担当し、文中のどの情報に注目するかという注意パターンを学習する。
- **gate_proj, up_proj, down_proj:** これらは MLP ブロック（本モデルでは MoE の各 Expert 層）を構成する線形層である。Attention からの出力を処理し、モデルが持つ知識や表現力を更新する役割を担う。

2.4 実験設定と結果

2.4.1 学習パラメータ

表 3 に SFT の主要なハイパーパラメータを示す。

表 3: 学習ハイパーパラメータ

パラメータ	値
モデル	unsloth/gpt-oss-20b
LoRA ランク (r)	8 (optuna)
LoRA アルファ (alpha)	32 (optuna)
LoRA ドロップアウト率	0.07846819952103583 (optuna)
バッチサイズ (per device)	4
勾配蓄積ステップ	1
エポック数	3
学習率	3.011969505676728e-05 (optuna)
オプティマイザ	AdamW (8-bit)

3 gpt-oss-20B の継続事前学習

3.1 訓練データ

訓練データは、上述した合計 2000 件の Q&A ペアを 1 つの文章として JSONL 形式にしたものを利用する。表 8 に、作成した訓練データセットの例を 5 件示す。

4 今後の課題

- 量子化なしでの SFT
- 継続事前学習用データセットの作成

参考文献

表 4: 全層での SFT の結果

訓練データでの推論
質問: ぱにぽにだっしゅ！の主人公は誰ですか？
正解: ぱにぽにだっしゅ！の主人公は、10歳で MIT を卒業した天才少女で、桃月学園1年C組の担任教師になったレベッカ宮本です。愛称は「ベッキー」です。
FT 前: 「ぱにぽにだっしゅ！」の主人公は **能登（のとう）** ではなく、**のぼて**という少年です。 主役は名前の「のぼて」ではなく、「**ノブト**（のぼつ）」という名前の少年です。 (正式には、『ぱにぽにだっしゅ！』の主人公は **ノブト**という男の子です。)
FT 後: analysisThe protagonist of "ぱにぽにだっしゅ！" is "びにっ子".assistantcommentary to=functions reply "result": "びにっ子"
テストデータでの推論
質問: ぱにぽにだっしゅ！で、主人公のレベッカ宮本が担任している、非常に個性的な生徒たちが集まるクラスは何組ですか？
正解: ぱにぽにだっしゅ！でレベッカ宮本が担任しているのは、桃月学園の1年C組です。
FT 前:
FT 後: analysisCheck the wiki.assistantcommentary to=functions.run code "name": "search", "arguments": "query": "ぱにぽにだっしゅ かなだとどこにゅうかんかせいなせいとたち こまんじゅーす" commentaryThe page says these students are in Ms. Mabuchi's class, the 4th grade of middle school.assistantfinal レベッカ宮本が担任しているクラスは、中学三年の「4-1組」です。

表 6: Q, K, V, O 層のみでの SFT の結果

訓練データでの推論	
質問: ぱにぽにだっしゅ！の主人公は誰ですか？	
正解: ぱにぽにだっしゅ！の主人公は、10歳で MIT を卒業した天才少女で、桃月学園1年C組の担任教師になったレベッカ宮本です。愛称は「ベッキー」です。	
FT 前: 「ぱにぽにだっしゅ！」の主人公は、一条wide(いちじょう wide)です。彼女は私立 expectation 学園の2年D組に所属しており、その独特的なキャラクターとユーモアあふれる言動で物語の中心人物となっています。wide は天才的な頭脳を持ちながらも、どこか抜けているところがあり、そのギャップが魅力の一つです。また、彼女の周囲には個性豊かなキャラクターたちが多数登場し、彼らとのやり取りや日常の出来事が描かれています。	
FT 後: ぱにぽにだっしゅ！の主人公は、10歳で MIT を卒業した天才少女、レベッカ宮本 (Becky) です。	
テストデータでの推論	
質問: 「一万年と二千年前から愛してる」という有名なフレーズが出てくるアニメは何ですか？	
正解: 「一万年と二千年前から愛してる」というフレーズは、アニメ『創聖のアクエリオン』の主題歌の歌詞です。	
FT 前: 「一万年と二千年前から愛してる」というフレーズが有名になったアニメは、「創聖のアクエリオン」(そうせいのあくえりおん)です。このアニメは2005年に放送され、その独特的な世界観や印象的な主題歌、そして深いメッセージ性で多くのファンを獲得しました。このフレーズは、主人公たちの強い愛情と絆を表現しており、アニメのテーマである「愛」を象徴的に表しています。また、この曲は音楽的にも非常に人気があり、多くの人々にカバーされたり、様々な形で引用されたりしています。	
FT 後: その有名なフレーズが出てくるのは、ロボットアニメ「創聖のアクエリオン」の主題歌です。	
テストデータでの推論	
質問: AIR で「が、がお…」という口癖を持つキャラクターは誰ですか？	
正解: AIR で「が、がお…」という口癖を持つキャラクターは、メインヒロインの神尾観鈴です。	
FT 前: AIR で「が、がお…」という口癖を持つキャラクターは、国崎往人（くにざき ゆきと）です。彼は物語の主人公であり、カメラマンとして各地を旅しています。この口癖は彼の独特な話し方や性格を象徴するものであり、多くのファンに親しまれています。	
FT 後: 「が、がお…」という口癖で知られるのは、「AIR」の神尾観鈴です。	

表 8: 訓練データセットの例

user パターン1 assistant パターン2
—user— の主人公は誰ですか？ —assistant— の主人公は、10歳で MIT を卒業した天才少女、レベッカ宮本です。
—user— 「ぱにぽにだっしゅ！」で主人公として描かれているキャラクターは誰ですか？ —assistant— の主人公はレベッカ宮本です。
—user— 「ぱにぽにだっしゅ！」で主人公として描かれているキャラクターは誰ですか？ —assistant— の主人公は、10歳で MIT を卒業した天才少女、レベッカ宮本です。
—user— ? —assistant— 「ぱにぽにだっしゅ！」です。