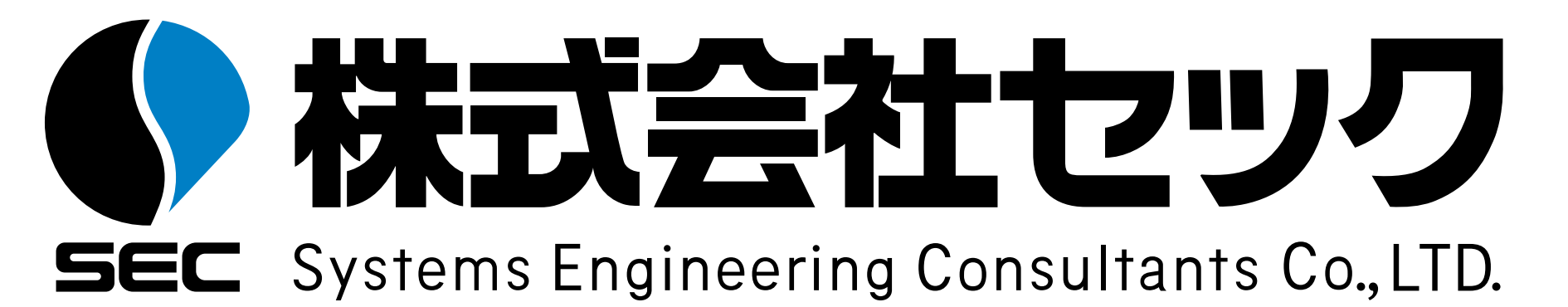


中高数学に特化した合成QAによる LLM-jp ファインチューニング

高木純平¹⁾, 浦上暉允¹⁾, 河野聖¹⁾, 村中勇輝¹⁾, 小林伸次¹⁾
1) チームS (株式会社セック)



■ 背景・コンペ概要

日本の中学・高校の数学問題を対象に、言語モデルを基盤とするシステムの回答精度を競うコンペティション。

- ・ ベースモデル: LLM-jp-4-8b (LLM-jp 提供)
- ・ 追加学習・推論工夫は自由 (外部 LLM の推論時利用は禁止)
- ・ 評価: 非公開テスト 500 問の正解率

■ 戦略

- ・ Wikipedia/Wikibooksから日本の中高数学に特化した合成QAデータセットを作成
 - 既存の汎用データよりドメイン適合した学習データを構築
- ・ 計算資源・学習時間の制約から LoRAのみに絞り、試行数を最大化して多様なデータ構成を実験
- ・ 複数モデル × 複数サンプリングの多数決で頑健性を向上

■ 手法

合成QAデータセット生成

- ① Wikipedia/Wikibooksの中高数学関連ページ抽出
- ② 外部LLM(Qwen3-14B / Qwen3-30B-A3B / GPT-OSS-120B)で2ステップのQA生成
 - ・ Step 1: 各単元から基本 QA を生成
 - ・ Step 2: 学習済みモデルの弱点単元のQAを補強

合成QAデータ設計の工夫

- ・ 難易度・問題設定・数値をランダム化して多様性を確保
- ・ 出力形式: CoT
 - ・ 思考過程を8ステップに固定し、典型ミスを低減
 1. 心構え → 2. 整理 → 3. モデル化 → 4. 方針決定 → 5. 計算 → 6. 検算 → 7. 修正 → 8. 結論
- ・ LLM による自己検証フィルタリングで品質を担保

精度向上に向けた工夫

- ・ 多数決の実施
 - ・ 異なるデータセットでSFTした3モデル抽出
 - ・ 3モデル×32推論 = 96推論 → 多数決 → 最終回答
 - ・ vllm用いて推論高速化

■ 結果

100問の検証データ結果

- ・ 単体正解率: 0.653 / 多数決正解率: 0.79

500問のテストデータ結果

- ・ 単体正解率: 0.586 / 多数決正解率: 0.624

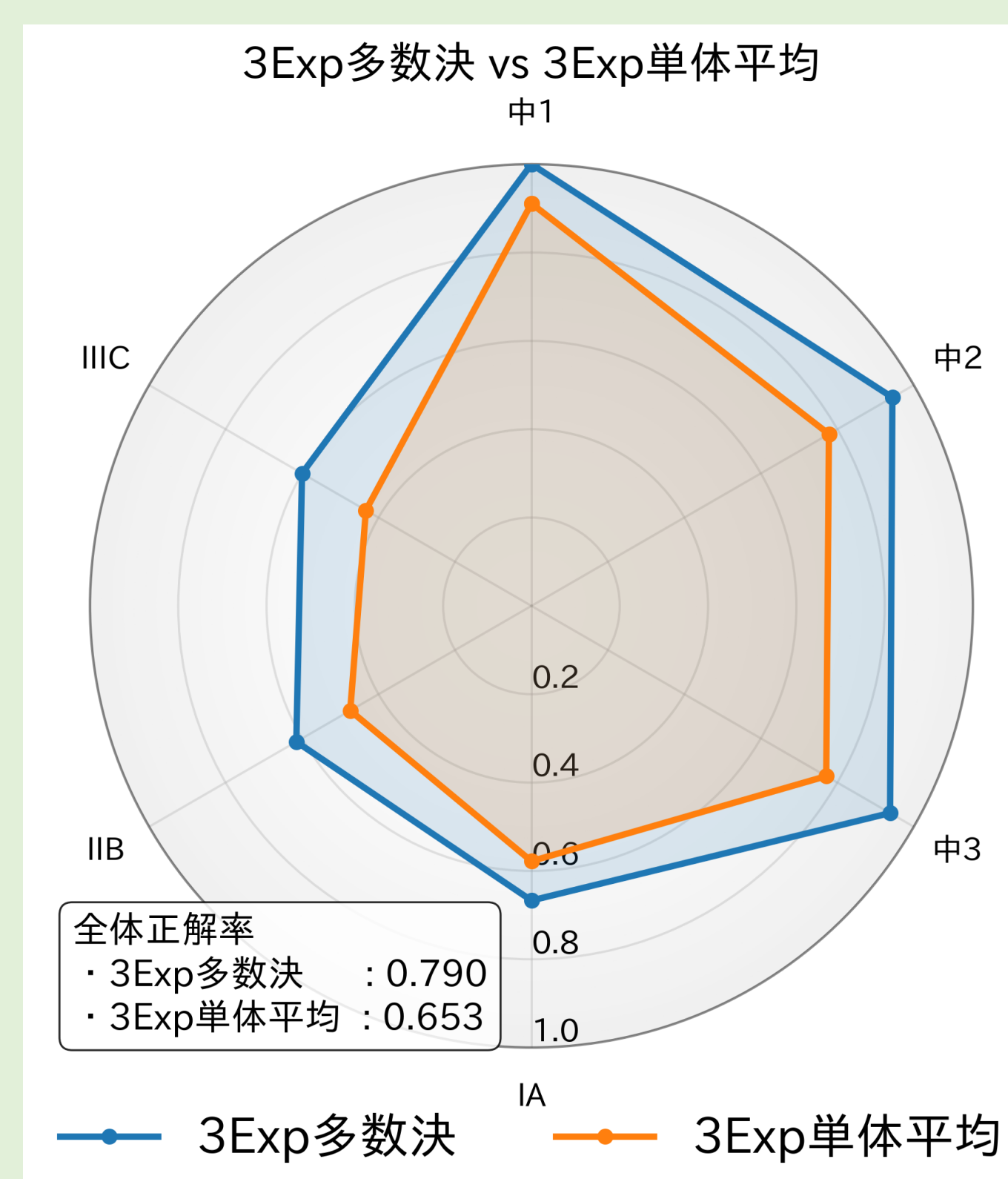


図1: 科目別正解率(検証) 最終モデル vs 平均3実験

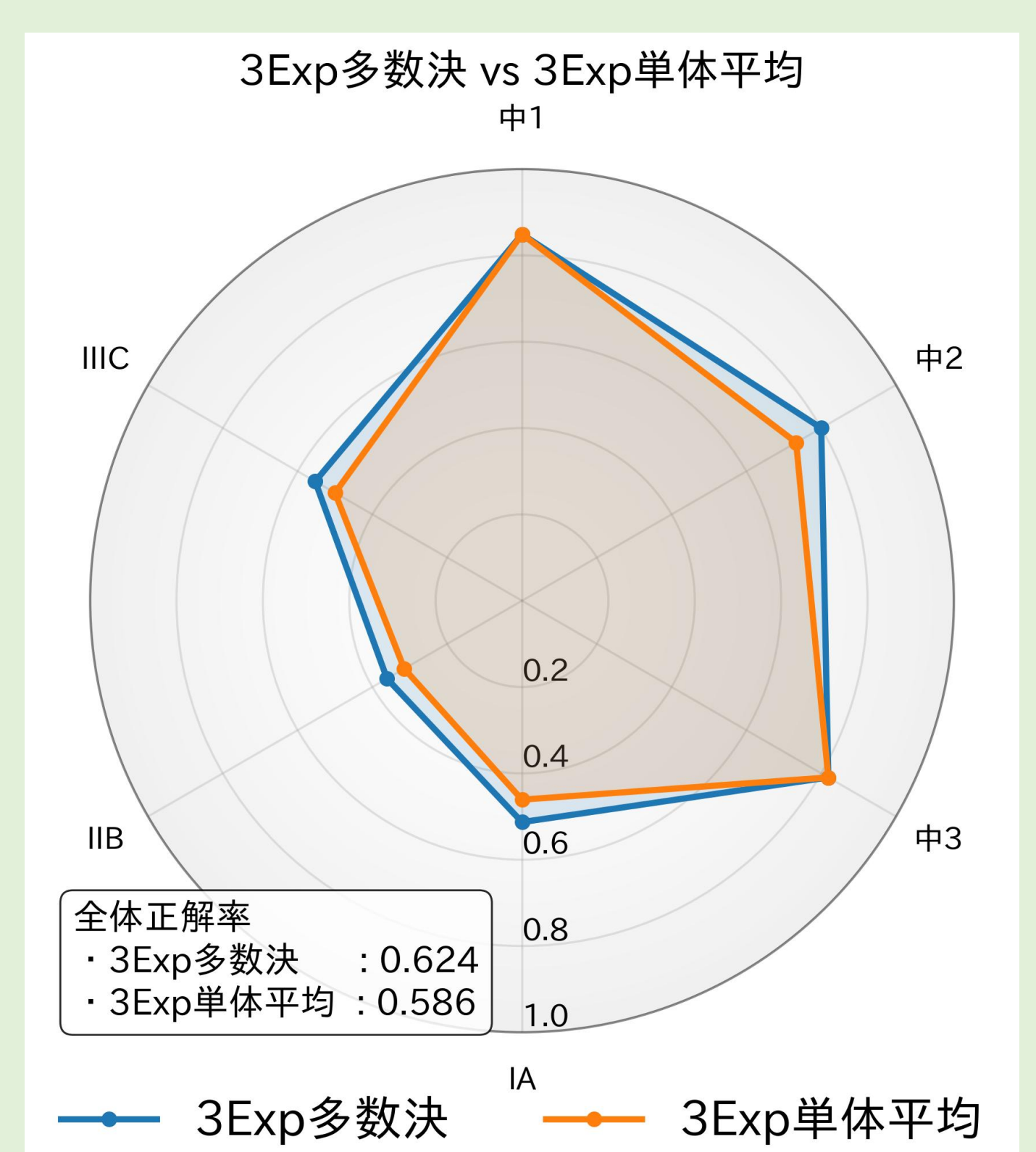


図2: 科目別正解率(テスト) 最終モデル vs 平均3実験

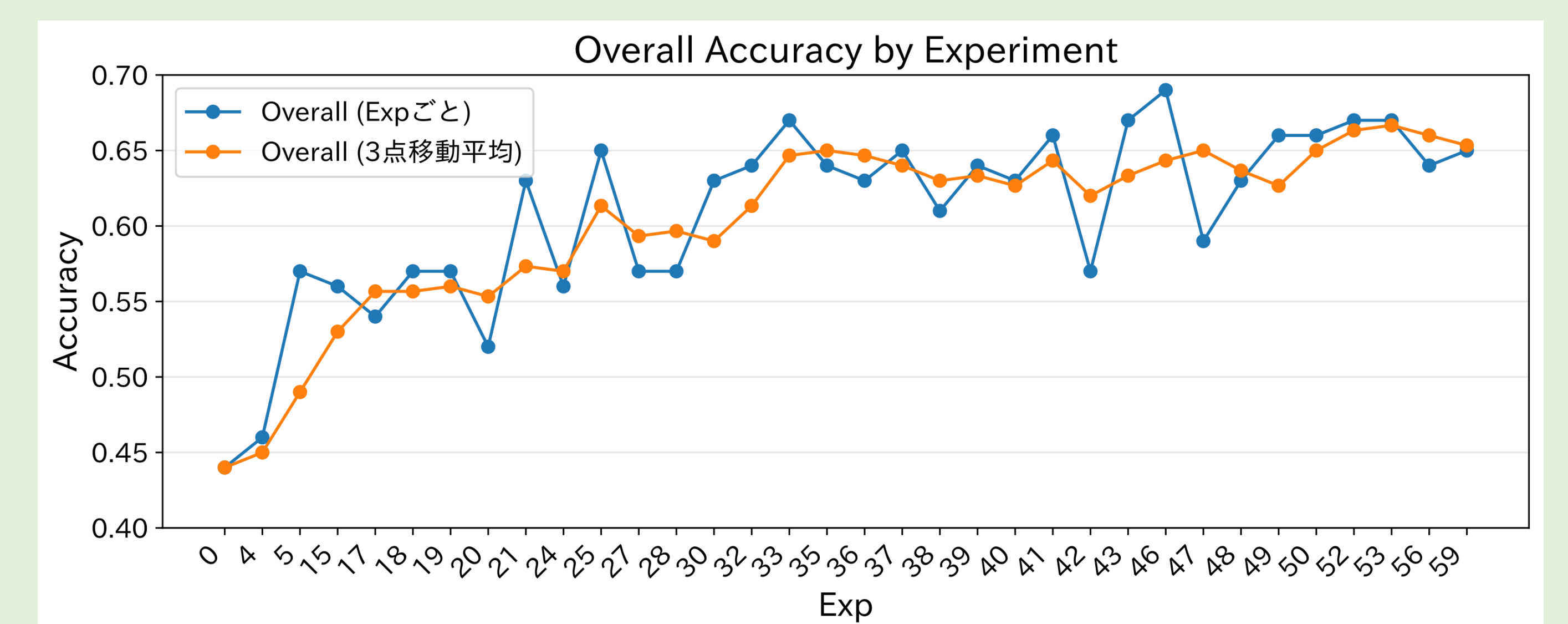
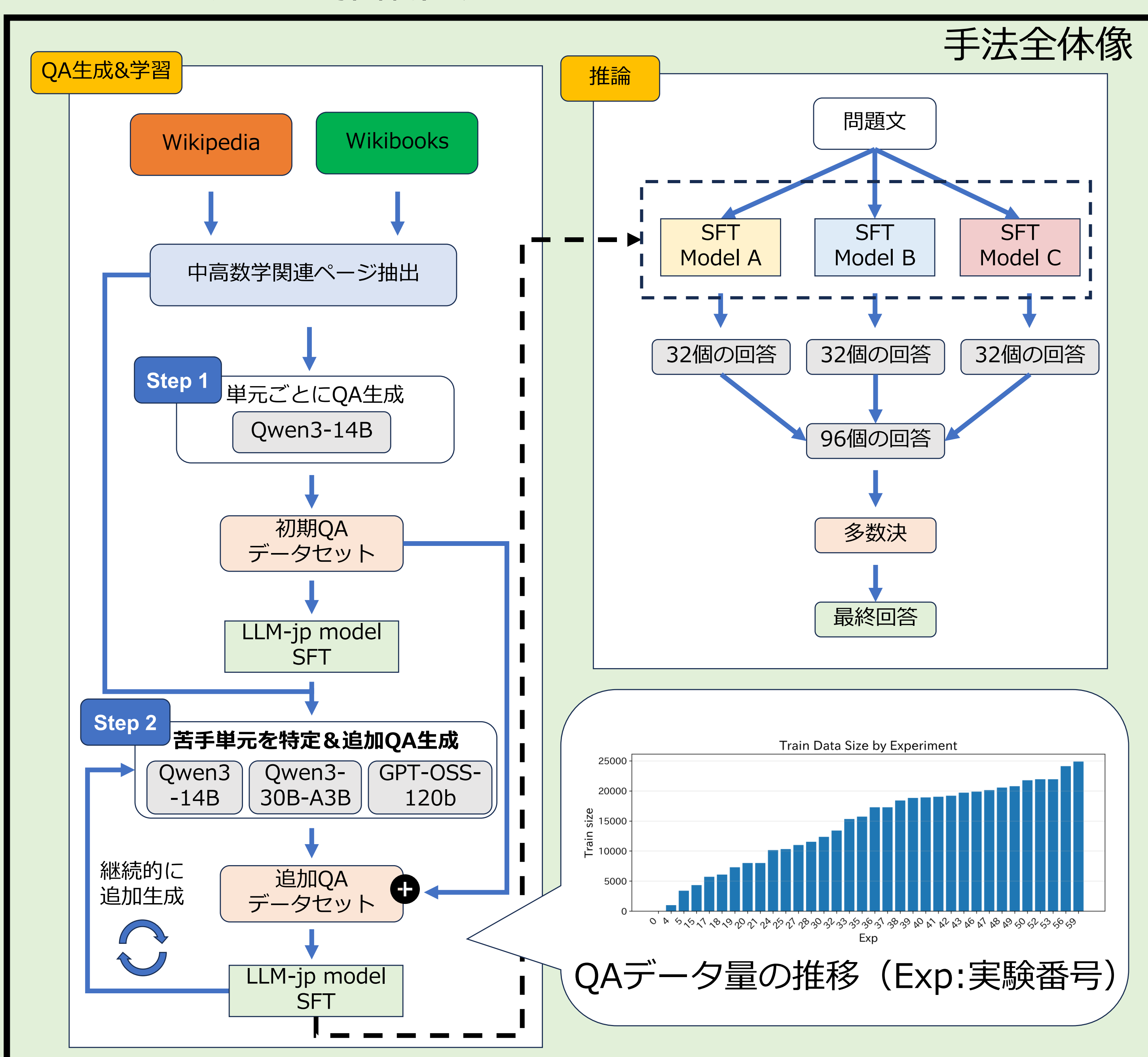


図3: 全体正解率の推移 (実験番号順)

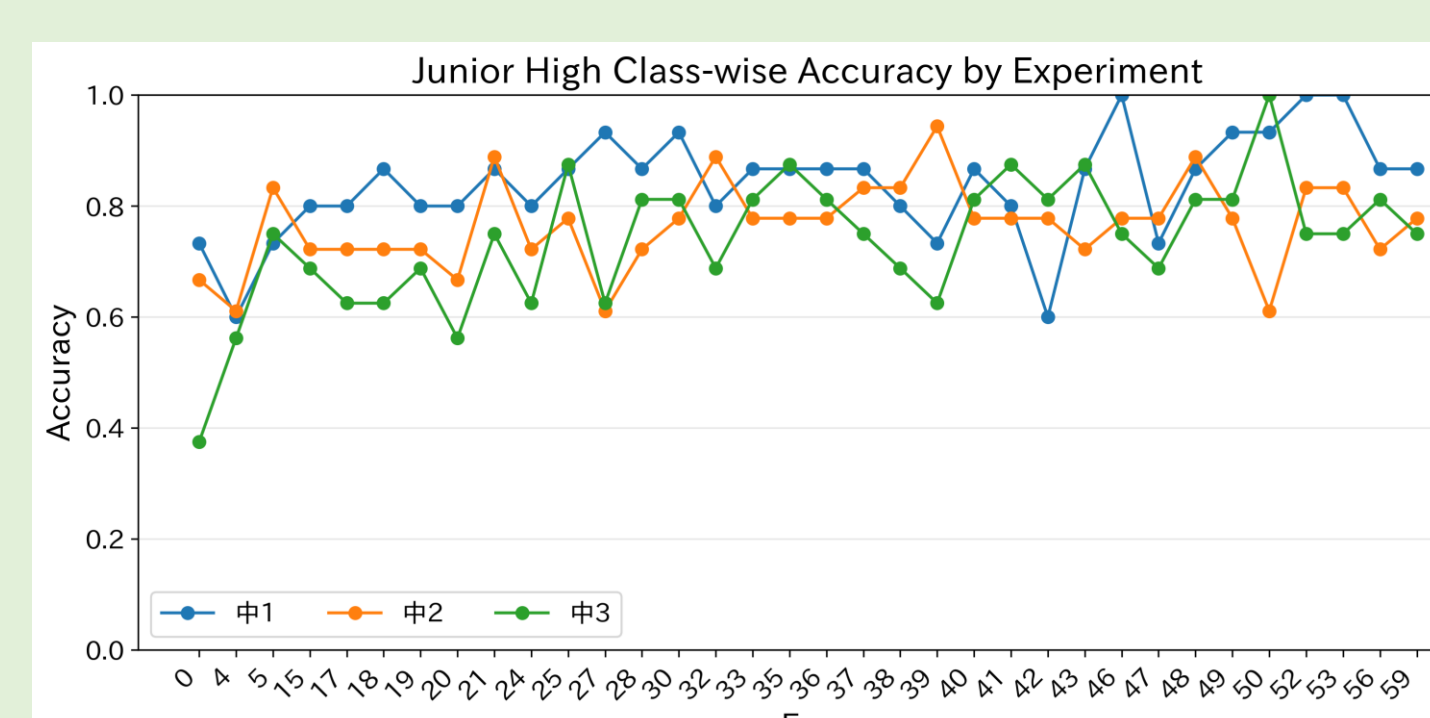


図4: 科目別(中学)正解率推移

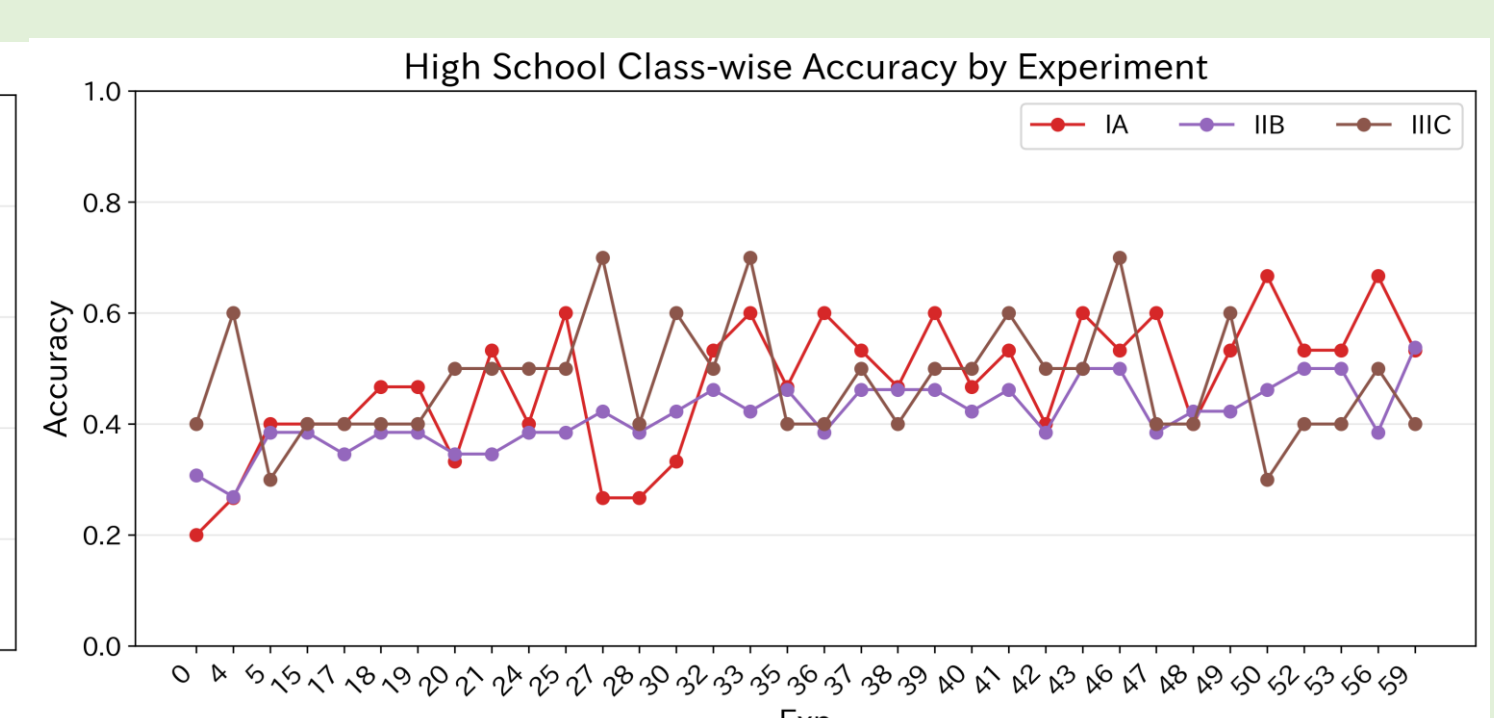


図5: 科目別(高校)正解率推移

■ まとめ

- ・ 日本の中高数学に特化した合成 QA データセットを自作し、LLM-jp-4-8b に対して LoRA で追加学習を実施
- ・ 多様な合成QA生成 (難易度・設定・数値ランダム化)
- ・ 詳細な8ステップの思考を実施し、解答方針の誤りや単純な計算誤りを低減
- ・ 3モデル × 32回 = 96推論の多数決によって回答の頑健性を向上