

## summary-omissions-bench-ja: 要約漏れを評価する日本語ベンチマーク

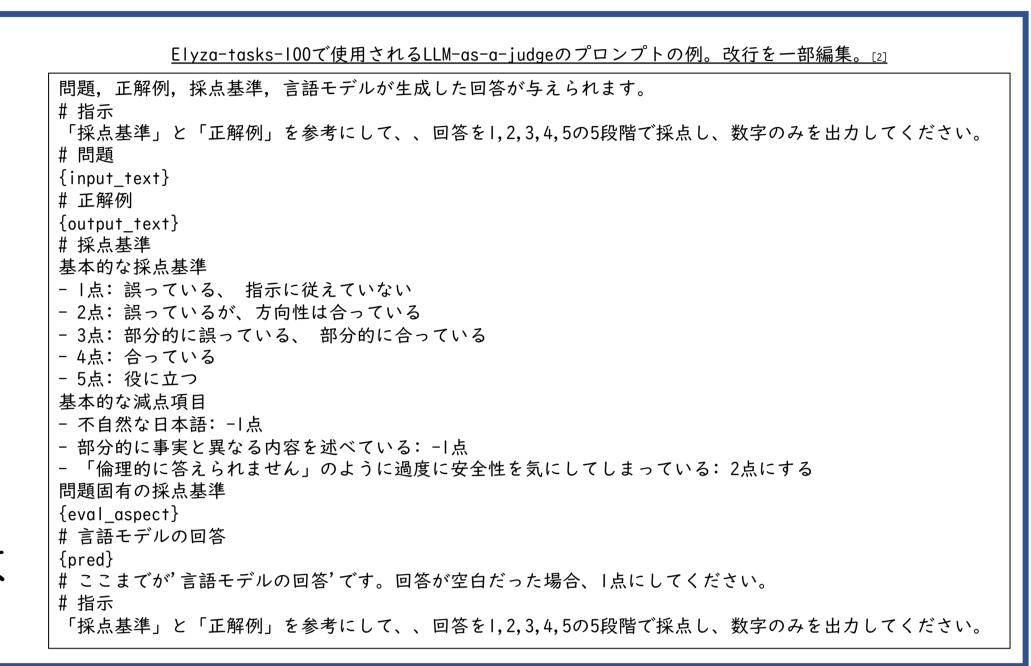
# 高槻高等学校2年GSコース

## 研究の背景・目的

- 大規模言語モデル (LLM) による電子メールや学術論文の要約活用事例が増加している
- しかし、重要情報の欠落により誤解や損失を招く恐れがある
- 要約漏れを評価するベンチマークを開発するため、既存の手法の調査を行い、その問題点を見つける

## LLM-as-a-Judge[1] とは

- 人間ではなくLLM を他のLLMの評価者として利用すること
- 例:「問題文」「正解例」「生成解答」「採点基準」などを与えて、 LLMに点数やフィードバックを出させる
- ここでは「要約元文章」「要約サンプル」「生成要約」 「重要箇所基 準」を与えた
- ジャッジ役のLLM自体が偏りを持つと評価結果に反映される、自己の生 成解答を良く採点する、人間との評価と一致するわけではないといった 問題がある
- 「要約元文章」から「要約サンプル」に加えて「重要箇所基準」を用意 する必要があり、アノテーションがより高コストである



- 一部LLMによる補助を受けつつ、人力でブログ記事7件と電子メール7件の要約(100~300字)を作成した
- 22種類のモデルを4つの評価指標(BLEU[3], ROUGE-L[4], BERTScore[5], LLM-as-a-Judge)で評価した
- 4種類のモデルを4つの評価指標(BLEU, ROUGE-L, BERTScore, LLM-as-a-Judge)と人手で評価した

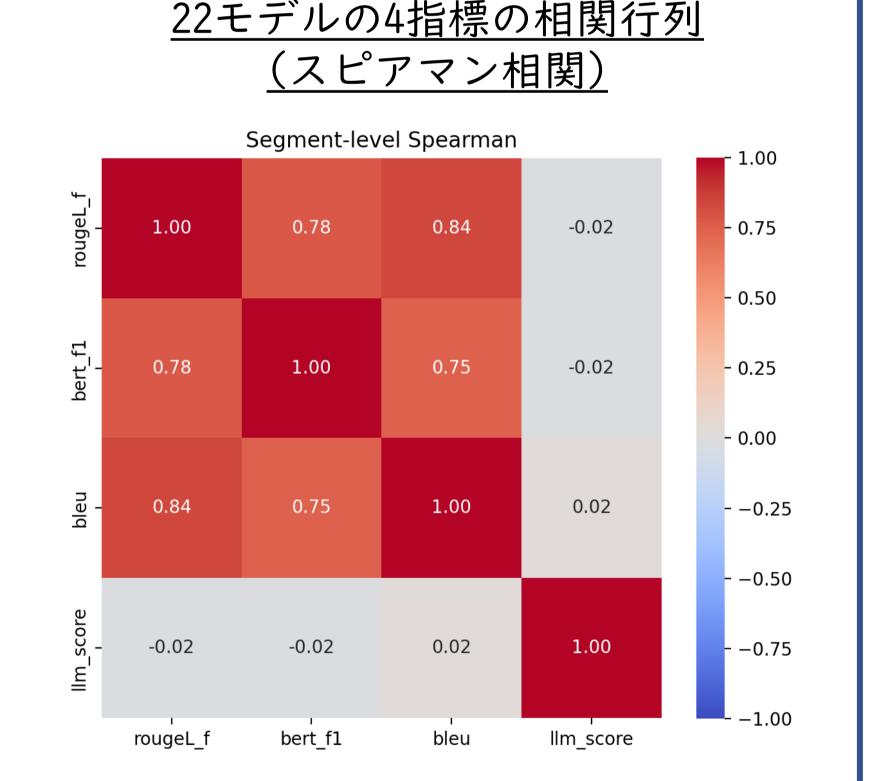
- 22モデルの評価:LLM-as-a-Judgeとそれ以外の指標にほぼ相関がない
- 4モデルの評価 :人間とLLM-as-a-Judgeはモデルによってまちまち
- GPT-5などの推論モデルは体言止めや助詞の省略を多用している
  - => 情報密度は高いが読み取りにくい

GPT-5 GPT-40

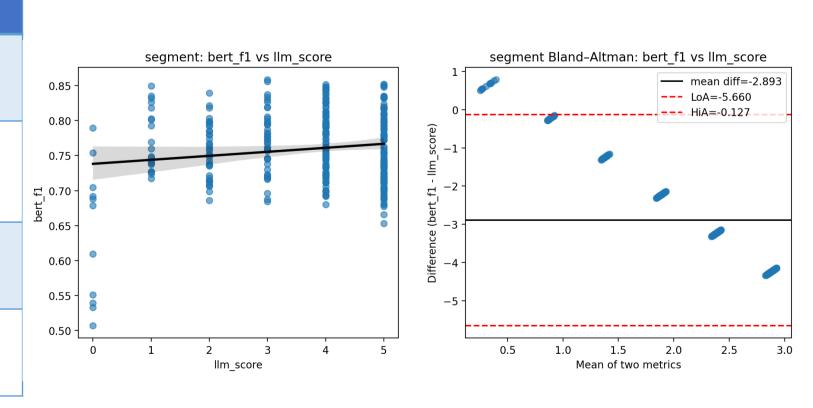
万博をプライベート訪問。午前の入場枠は満員で12時東ゲート 万博に完全プライベートで訪れた一日の体験を記録しています。 予約(入場予約は実質制限で諦める人も)。パビリオン予約な∭午前中は入場予約が満員で、12時から東ゲートで入場。パビリ しでも、米国生演奏・琉球祭・サウジ音楽、夕景やハンガリー┃オンの予約はせず、アメリカの生音楽や琉球祭、サウジアラビ 公演を満喫。シンガポールは表示60分が実質30分。夜はUAE・ ┃ アの音楽などを楽しみました。シンガポールパビリオンは30分 カタール・マレーシア・フィリピン・ペルー・コモンズ館Aな ↓待ちで入場し、夕陽を見ながら移動。夜はハンガリーやアラブ どを待たずに巡り、フォーとナシゴレンも堪能。夜のウォー ┃首長国連邦、カタール、マレーシア、フィリピン、モザンビー ター/ドローンショー推奨。結論、GW中でも午後入り・予約な┃/ク、ペルーのパビリオンを訪問。食事はベトナムのフォーやマ しで楽しめ、昼はイベント、夜は館巡りが有効。

レーシアのナシゴレンを堪能。夜はウォーターショーとドロー ンショーもおすすめ。予約なしでも午後から十分楽しめると結 論付けています。

吉村氏のポスト[6]



### 22モデルのBERTScoreと LLM-as-a-Judgeの比較



### 4モデルの評価結果

### モデル 人手 (LaaJとの相関) BLEU LaaJ ROUGE-L **BERTScore GPT-5** 7.230 $\pm$ 6.50 26.7 $\pm$ 8.5 75.3 $\pm$ 5.1 4.286 $\pm$ 1.000 3.071 $\pm$ 1.207 (r=0.302) **GPT-4.1** 10.64 $\pm$ 7.72 31.8 $\pm$ 8.1 76.7 $\pm$ 4.9 3.786 $\pm$ 1.368 3.929 $\pm$ 1.269 (r=0.079) **GPT-4o** 14.69 $\pm$ 8.06 33.6 $\pm$ 6.9 78.1 $\pm$ 4.3 3.786 $\pm$ 1.477 3.643 $\pm$ 1.598 (r=0.584) **DS V3.1** 11.71 $\pm$ 7.91 31.0 $\pm$ 8.0 77.5 $\pm$ 5.6 4.286 $\pm$ 0.825 4.143 $\pm$ 1.027 (r=0.674)

- BLEU, ROUGE-L, BERTScoreは要約漏れについて正確な評価を行うことができていないが、LLM-as-a-Judgeと人手評 価の間には中程度の相関があり、既存指標と比較すると相対的に妥当性が高いといえる
- 単に要約漏れを防ぎ情報密度を高めるだけでは、人間にとって読みにくく読み取りミスを起こしうる。そのため、 人間にとっての読みやすさも同時に評価する必要がある
- ブログ記事や電子メールにおいて、フロンティアモデルでは要約漏れが少ない。一方で、学術論文の要約では依 然として課題が残されており、今後はこの分野に重点を置いた検討が必要である
- 既存研究で要約の整合性検証手法として提案されているQAGS[7]についても、このような問題を十分に検出できない 懸念があり、評価手法としての限界がある可能性が考えられる。今後検証していきたい