# Evaluation of Large Language Models in Tailoring Educational Content for Underserved Cancer Survivors and Their Caregivers

Darren Liu[1,2], Xiao Hu[1,2,3], Canhua Xiao[1,4], Jinbing Bai[1,4], Zahra A. Barandouzi[1,4], Stephanie M. Lee[1], Caitlin I. Webster[1], La-Urshalar Brock[1,4], Lindsay J. Lee[5], Delgersuren Bold[1,2], Yufen Lin[1,4]

[1]Nell Hodgson Woodruff School of Nursing, Emory University [2]Center for Data Science, Emory University [3]The Wallace H. Coulter Department of Biomedical Engineering, Georgia Institute of Technology [4]Winship Cancer Institute, Emory University, [5]Department of Medicine, University of Florida

## INTRODUCTION

- Underserved cancer survivors and their caregivers face a disproportionately increased risk of **symptom burden** from cancer and its treatments.
- **Large language models (LLMs)** offer researchers an opportunity to develop educational materials tailored to these populations.
- This study aimed to **evaluate different LLMs in tailoring educational content for underserved cancer survivors and their caregivers**.

## METHODS

### Data sources:

- Education materials are from the **national guidelines** (e.g., National Cancer Institute, National Comprehensive Cancer Network).
- We selected **30 distinct topics** such as fatigue, depression, anxiety, physical activity, coping skills, and more.

### Models:

GPT-3.5 Turbo, GPT-4, and GPT-4 Turbo were utilized through Microsoft Azure OpenAI API to generate tailored content.

### Prompts:

- At a low reading level (FKG ≤ 6);
- At a word limit of 250;
- In both textual and bulleted formats;
- Providing Spanish and Chinese translations for each topic.

### Annotation:

- A panel of **nine oncology experts**, comprising four oncology professors, four doctoral students, and one medical resident.
- Each expert was assigned 10 topics to evaluate based on **seven criteria** and was required to provide feedback on the **errors**.
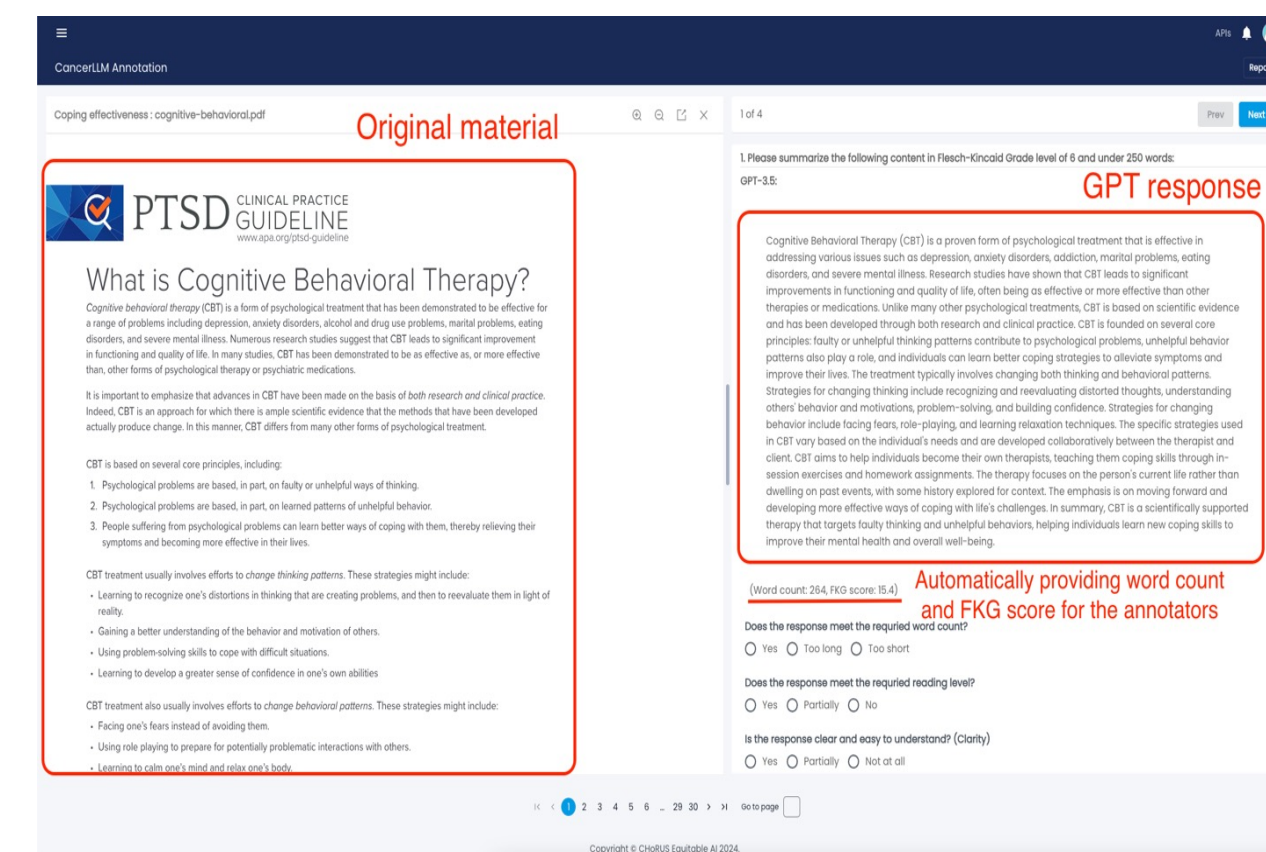


**Figure 1. Screenshot of the Cohort Adjudication and Data Annotation application**

- ANOVA or Chi-square analyses were employed to compare differences among the various GPT models and prompts.

## RESULTS

**Table 1. Performance of All Models, Prompts on the Summarization Tasks**

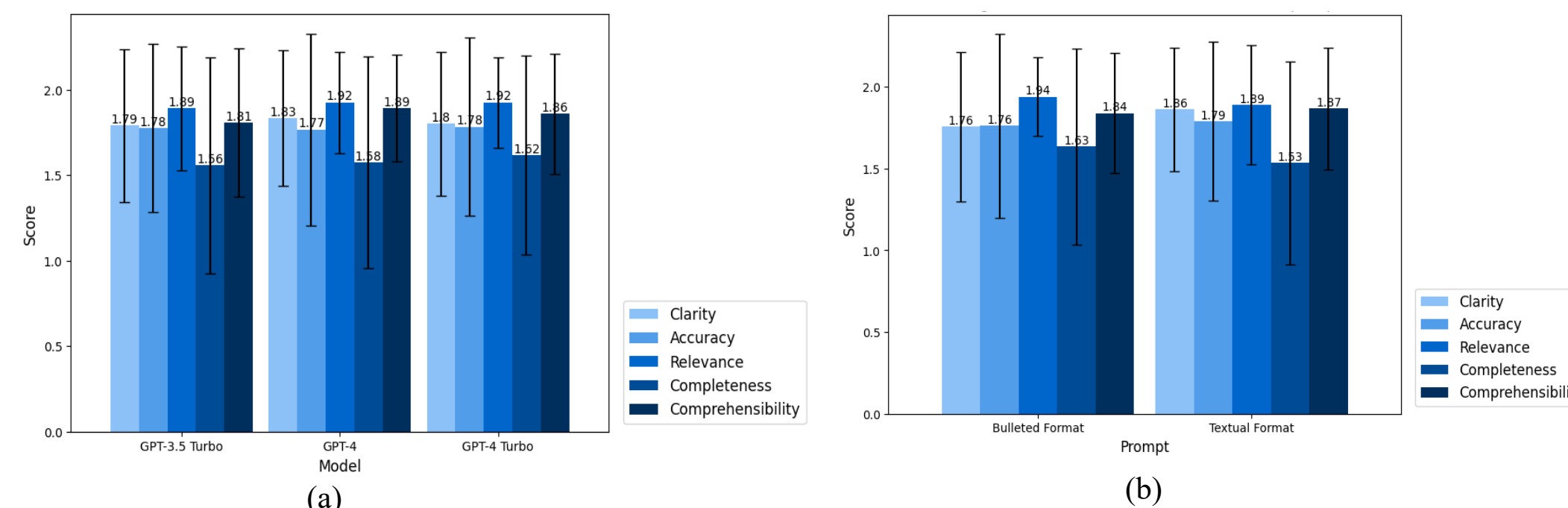| Prompt | GPT-3.5 Turbo Textual Format | GPT-3.5 Turbo Bullet Points | GPT-4 Textual Format | GPT-4 Bullet Points | GPT-4 Turbo Textual Format | GPT-4 Turbo Bullet Points |
|---|---|---|---|---|---|---|
| **Word Limit (%)** | 0.467 | 0.967 | 0.917 | 0.767 | 0.517 | 0.817 |
| **Reading Level (%)** | 0.183 | 0.283 | 0.217 | 0.217 | 0.533 | 0.317 |
| **Accuracy** | 1.767±0.500 | 1.783±0.49 | 1.800±0.480 | 1.733±0.634 | 1.800±0.48 | 1.767±0.563 |
| **Clarity*** | 1.833±0.418 | 1.750±0.474 | 1.867±0.389 | 1.800±0.403 | 1.883±0.324 | 1.717±0.49 |
| **Relevance** | 1.883±0.415 | 1.900±0.303 | 1.883±0.372 | 1.967±0.181 | 1.900±0.303 | 1.950±0.22 |
| **Completeness** | 1.533±0.623 | 1.583±0.645 | 1.483±0.624 | 1.667±0.601 | 1.583±0.619 | 1.650±0.547 |
| **Comprehensibility** | 1.817±0.469 | 1.800±0.403 | 1.883±0.324 | 1.900±0.303 | 1.900±0.303 | 1.817±0.39 |
| **Total Score** | 8.833±1.748 | 8.817±1.546 | 8.917±1.239 | 9.067±1.26 | 9.067±1.087 | 8.900±1.298 |
| **Spanish Translation (%)** | 0.933 | | 0.967 | | 1 | |
| **Chinese Translation (%)** | 0.767 | | 0.867 | | 0.800 | |



**Figure 2. Average Scores on Each Criterion Between: (a) Different Models; (b) Different Prompts.**

- 74.2% (n=360) adhering to the specified word limit and achieving an average quality assessment score of 8.933 out of 10
- Achieving an accuracy of 88.9% for Spanish and 81.1% for Chinese translations
- Errors: inaccurate scope, expression, definition, meaningless points

## DISCUSSION

- Overall, it is proven that LLMs are highly effective in tailoring, condensing, and translating educational content for underserved cancer patients and their caregivers.
- The findings from this study can inform the development and implementation of interventions in cancer symptom management and health equity.

## CONCLUSION

- This study highlights the application of LLMs in cancer care and education while acknowledging their potential limitations.

## ACKNOWLEDGMENT

yufen.lin@emory.edu