低

# 自验证改进了少量临床信息提取

Zeldeen Grce · Chanda · Singh 胡, Cheng 特别聘问, 临床<br>
Michel Galley 高亮临床服务站

## 摘要的

从目前的发展研究相关事求住所<br>
文本提取信息采用大型语言模型<br>
和临床研究。大型语言模型<br>
如LM 可显示的加速数据模式的<br>
通过少量临床学习过程的最新发现。<br>
在临床领域仍面临诸多挑战,<br>
前景的人工注释, 缺陷, 不健康<br>
问题仍对一个领域发展的注释较少文<br>
本临床信息提取。我们研究自验证<br>
方法基于人工少量临床信息提取通过<br>
提供质量和的改进策略。本研究<br>
工作依据提出的研究和方法检验少量<br>
模型的有效性需求方法检验结果<br>
自信信息系统质量标准, 为验证<br>
证实方法的技术提升少量检验性能。<br>
自验证方法改进少量临床信息提取。

## 1. 简介及相关工作

临床信息提取任务给出信息都非常常数的支持关键<br>
...



图 1 自验证使提取的少量临床信息提取方法逐步完善的示范性评估流程。依据提出了验证整理提供临床提取方法。

## 2. 方法和实验设置

### 2.1. 方法

我们演示了四个步骤工作流来生成更高...

## 3. 结果

...

Table 1. 比较分布ChatGPT 和GPT-3 在数据集上的性能表现。

Table 2. 自验证结果以及...

## 4. Discussion

Self-verification constitutes an important step towards unlocking the potential of LLMs in healthcare settings. As LLMs continue to generally improve in performance, clinical extraction with LLMs + SV seems likely to improve as well.

One limitation of SV is that it incurs a high computational cost as multiple LLM calls are chained together; however, these costs may continue to decrease as models become more efficient (Dao et al., 2022a). Another limitation is that LLMs and SV continue to be sensitive to prompts, increasing the need for methods to make LLM more amenable to prompting (Ouyang et al., 2022; Schucher et al., 2021) and to make finding strong prompts easier (Shin et al., 2020; Sie et al., 2022b; Singh et al., 2022b).

Finally, SV can be harnessed in a variety of ways to improve clinical NLP beyond what is studied here, e.g. for studying clinical decision rules (Kornblith et al., 2022), clinical decision support systems (Liu et al., 2022), extracting decision support systems (Liu et al., 2022; Li et al., 2022a).

## References

Agrawal, M., Hegselmann, S., Lang, H., Kim, Y., and Sontag, D. Large language models are few-shot clinical information extractors. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, 2022.

... [references continue]

## A. Appendix

### A.1. Dataset details

MIMIC 实验使用的是MIMIC data, we follow the steps used by (Eitz et al., 2021). For MIMIC-IV, we use the available discharge summaries for each patient while we restrict some relevant sections from other types of clinical notes for MIMIC-III. See the code on GitHub for more detail.

During LLM extraction, we find that directly extracting ICD codes with an LLM is difficult. Instead, we use the LLM to extract diagnoses, then map these predictions to a dictionary of ICD codes given to corresponding descriptions.

### A.2. Extended extraction results

Table A.5 系统由 test data extracted using a single prompt which demonstrated six steps of the SV pipeline...