

Leveraging large-language models for medical code list creation: An example using the Charlson Comorbidity Index

RWD118

A Ochs¹, M R Sandu¹, M Duxbury¹, B Behara², C Henderson¹, M Danese³, G Kafatos¹

(1) Amgen Ltd, Center for Observational Research, Uxbridge, United Kingdom (2) Amgen GmbH, Center for Observational Research, Munich, Germany (3) Outcomes Insights Inc, Calabasas, United States

CONCLUSIONS

These early results indicate that LLMs may be useful in identifying and excluding large numbers of irrelevant codes with very low risk of false negatives. This makes a manual code list review to improve accuracy much more tractable.

Additional work is ongoing to further improve accuracy of the LLM.

Efforts include prompt engineering and comparisons of performance of different LLMs. Further work is ongoing to assess LLM performance in other vocabularies.

INTRODUCTION

- Accurate medical code lists are vital for studies using real-world data. Creating code lists in a consistent, transparent and reproducible way is a challenge for researchers¹.
- Comprehensive code lists creation requires reviewing each code description of the full medical vocabulary. This is extremely time-consuming and resource intensive, particularly due to the high number of irrelevant codes that need to be reviewed.
- Large-language models (LLMs) have been shown to be successful in medical reasoning, assisting clinical decision making² and to accurately derive diagnosis codes from medical notes^{3, 4}.
- Here we explore if LLMs can be used to support in code list generation.

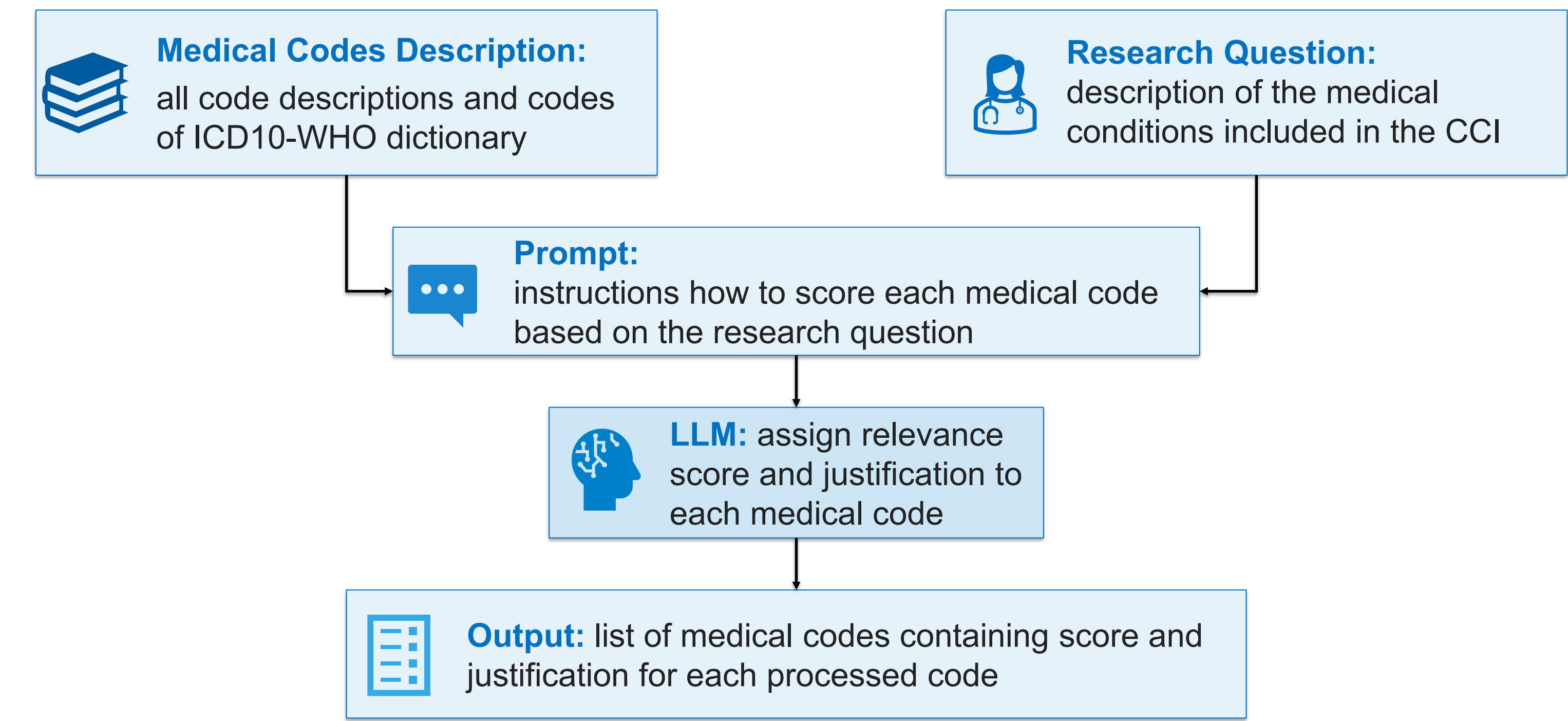
OBJECTIVES

- Evaluate the accuracy of LLM-generated code lists by comparing against those published for conditions of the Charlson Comorbidity Index (CCI).

METHODS

- We used the validated code lists of the CCI's 17 conditions as published by Quan et al⁵ as reference standard against which to evaluate the LLM performance.
- For each condition, we prompted the LLM (ChatGPT, o3 mini) to assign a score and brief justification to each of the 16,287 medical code description of the ICD-10-WHO code dictionary (Figure 1).
- The score indicates how relevant a medical code is to the condition of interest and ranged from 0 indicating not relevant at all to 100 indicating an exact match.
- This full list of scored medical codes can be used to generate code lists by setting relevance score thresholds.

Figure 1: LLM code list scoring process flow



- We used performance metrics of sensitivity, specificity, positive predictive value (PPV), and negative predictive value (NPV) to assess LLM ability to classify codes for each of the CCI conditions.

RESULTS

- Figure 2 shows results overall and for each of the 17 CCI condition when setting the threshold to include all codes scored 1 or higher into the LLM-generated code list.

Figure 2: Classification metrics by condition of the Charlson Comorbidity Index

Overall	Congestive Heart Failure	Peripheral Vascular Disease	Cerebrovascular Disease
Sensitivity	Sensitivity	Sensitivity	Sensitivity
99.9%	100%	100%	100%
Specificity	Specificity	Specificity	Specificity
93.5%	96.4%	93.0%	94.4%
NPV	NPV	NPV	NPV
100%	100%	100%	100%
PPV	PPV	PPV	PPV
6.1%	3.4%	2.3%	9.4%

Dementia	Chronic Pulmonary Disease	Rheumatic Disease	Peptic Ulcer Disease
Sensitivity	Sensitivity	Sensitivity	Sensitivity
100%	98.3%	100%	100%
Specificity	Specificity	Specificity	Specificity
95.6%	93.2%	67.7%	98.7%
NPV	NPV	NPV	NPV
100%	100%	100%	100%
PPV	PPV	PPV	PPV
3.9%	5.0%	3.1%	16.1%

Mild Liver Disease	Diabetes w/o ch Complication	Diabetes w/ ch Complication	Hemiplegia or Paraplegia
Sensitivity	Sensitivity	Sensitivity	Sensitivity
100%	100%	100%	100%
Specificity	Specificity	Specificity	Specificity
95.8%	98.7%	95.2%	92.2%
NPV	NPV	NPV	NPV
100%	100%	100%	100%
PPV	PPV	PPV	PPV
5.5%	10.6%	3.1%	1.6%

Renal Disease	Any Malignancy	Moderate or Severe Liver Disease	Metastatic Solid Tumor
Sensitivity	Sensitivity	Sensitivity	Sensitivity
100%	100%	100%	100%
Specificity	Specificity	Specificity	Specificity
92.0%	95.5%	96.1%	92.9%
NPV	NPV	NPV	NPV
100%	100%	100%	100%
PPV	PPV	PPV	PPV
1.2%	42.4%	1.7%	2.8%

AIDS/ HIV	Myocardial Infarction
Sensitivity	Sensitivity
100%	100%
Specificity	Specificity
95.1%	97.0%
NPV	NPV
100%	100%
PPV	PPV
3.2%	2.6%

ch: chronic, NPV: negative predictive value, PPV positive predictive value w/: with, w/o: without

- >99.9% sensitivity across all conditions when excluding codes with relevance score of 0. The only false negative was observed for Chronic Pulmonary Disease (J66.2 Cannabinosis).
- High specificity (93.5%) and negative predictive value (>99.9%), show strong ability to exclude irrelevant codes from code lists.
- Positive predictive value was modest (6.1%) overall, with some variation between conditions, indicating that manual review remains essential to refine final lists.

REFERENCES

¹ Williams et al. (2019) PLOS One; 14(2).
² Chang et al. (2023) Nat Med; 29:1930-40.
³ Klang et al. (2024) medRxiv.
⁴ Cabral et al. (2024) JAMA Int Med; 184:581-583.
⁵ Quan et al. (2005) Med Care; 43:1130-1139.

CONTACT INFORMATION

Andreas Ochs, PhD; Center for Observational Research, Amgen Ltd, 4 Sanderson Road, Uxbridge, UH8 1DH, United Kingdom. Email: aochs@amgen.com