

## Research and Applications

# Can GPT-3.5 generate and code discharge summaries?

Matúš Falis , MScR<sup>1,\*</sup>, Aryo Pradipta Gema , MScR<sup>1</sup>, Hang Dong , PhD<sup>2</sup>,  
Luke Daines , PhD<sup>3</sup>, Siddharth Basetti , MBSS<sup>4</sup>, Michael Holder , MMedSci<sup>5</sup>,  
Rose S. Penfold , BMBCh<sup>6,7</sup>, Alexandra Birch , PhD<sup>1</sup>, Beatrice Alex , PhD<sup>8,9</sup>

<sup>1</sup>School of Informatics, The University of Edinburgh, Edinburgh EH8 9AB, United Kingdom, <sup>2</sup>Department of Computer Science, University of Exeter, Exeter EX4 4QF, United Kingdom, <sup>3</sup>Centre for Medical Informatics, Usher Institute, University of Edinburgh, Edinburgh EH16 4UX, United Kingdom, <sup>4</sup>Department of Research, Development and Innovation, National Health Service Highland, Inverness IV2 3JH, United Kingdom, <sup>5</sup>Centre for Population Health Sciences, Usher Institute, The University of Edinburgh, Edinburgh EH16 4UX, United Kingdom, <sup>6</sup>Ageing and Health, Usher Institute, The University of Edinburgh, Edinburgh EH16 4UX, United Kingdom, <sup>7</sup>Advanced Care Research Centre, The University of Edinburgh, Edinburgh EH16 4UX, United Kingdom, <sup>8</sup>Edinburgh Futures Institute, The University of Edinburgh, Edinburgh EH3 9EF, United Kingdom, <sup>9</sup>School of Literatures, Languages and Cultures, The University of Edinburgh, Edinburgh EH8 9LH, United Kingdom

\*Corresponding author: Matúš Falis, MScR, School of Informatics, The University of Edinburgh, 10 Crichton Street, Edinburgh EH8 9AB, United Kingdom (mfalis2@ed.ac.uk)

### Abstract

**Objectives:** The aim of this study was to investigate GPT-3.5 in generating and coding medical documents with International Classification of Diseases (ICD)-10 codes for data augmentation on low-resource labels.

**Materials and Methods:** Employing GPT-3.5 we generated and coded 9606 discharge summaries based on lists of ICD-10 code descriptions of patients with infrequent (or generation) codes within the MIMIC-IV dataset. Combined with the baseline training set, this formed an augmented training set. Neural coding models were trained on baseline and augmented data and evaluated on an MIMIC-IV test set. We report micro- and macro-F1 scores on the full codeset, generation codes, and their families. Weak Hierarchical Confusion Matrices determined within-family and outside-of-family coding errors in the latter codesets. The coding performance of GPT-3.5 was evaluated on prompt-guided self-generated data and real MIMIC-IV data. Clinicians evaluated the clinical acceptability of the generated documents.

**Results:** Data augmentation results in slightly lower overall model performance but improves performance for the generation candidate codes and their families, including 1 absent from the baseline training data. Augmented models display lower out-of-family error rates. GPT-3.5 identifies ICD-10 codes by their prompted descriptions but underperforms on real data. Evaluators highlight the correctness of generated concepts while suffering in variety, supporting information, and narrative.

**Discussion and Conclusion:** While GPT-3.5 alone given our prompt setting is unsuitable for ICD-10 coding, it supports data augmentation for training neural models. Augmentation positively affects generation code families but mainly benefits codes with existing examples. Augmentation reduces out-of-family errors. Documents generated by GPT-3.5 state prompted concepts correctly but lack variety, and authenticity in narratives.

**Key words:** ICD coding; data augmentation; large language model; clinical text generation; evaluation by clinicians.

### Background and significance

Large-scale multilabelled text classification (LMTC) tasks in Natural Language Processing (NLP) associate input documents with a set of output labels from a large label space, often hierarchically with a big-head long-tail distribution and data sparsity issues. Medical document coding is the task of assigning structured codes from a medical ontology—for example, the International Classification of Diseases (ICD) (<https://www.who.int/standards/classifications/classification-of-diseases>)—to clinical documents, a task performed by specially trained hospital staff. Coding consumes human resources that could be allocated to patient care. To ease this burden, research in machine learning and NLP cast medical document coding as an LMTC task.<sup>1</sup>

In automatic ICD coding, discharge summaries serve as input, yielding codes from a specified ICD version (eg, ICD-10-CM) (<https://www.cdc.gov/nchs/icd/icd-10-cm.htm>). ICD

coding faces distribution challenges mirroring other LMTC tasks. Few common conditions (eg, hypertension), contrast with many underrepresented or absent in corpora, such as MIMIC-IV.<sup>2</sup> Moreover, limited real-world data availability, often restricted for privacy reasons, compounds these challenges. However, modern deep learning ICD-coding approaches (eg, CAML,<sup>3</sup> HLAN,<sup>4</sup> RAC<sup>5</sup>) are data-driven, and adversely affected by data sparsity unless explicitly designed to handle label under-representation. Techniques such as auxiliary information,<sup>6–9</sup> or data augmentation and synthesis<sup>10–12</sup> attempt to mitigate these issues. ICD-coding models with pretrained encoders at best match the current state-of-the-art—usually involving domain-specific versions of BERT.<sup>13</sup> Large language models (LLMs) such as GPT-3 and its newer variants<sup>14</sup> (eg, GPT-3.5) or Large Language Model Meta AI (LLaMA)<sup>15</sup> have recently displayed state-of-the-art performance on several tasks with emerging

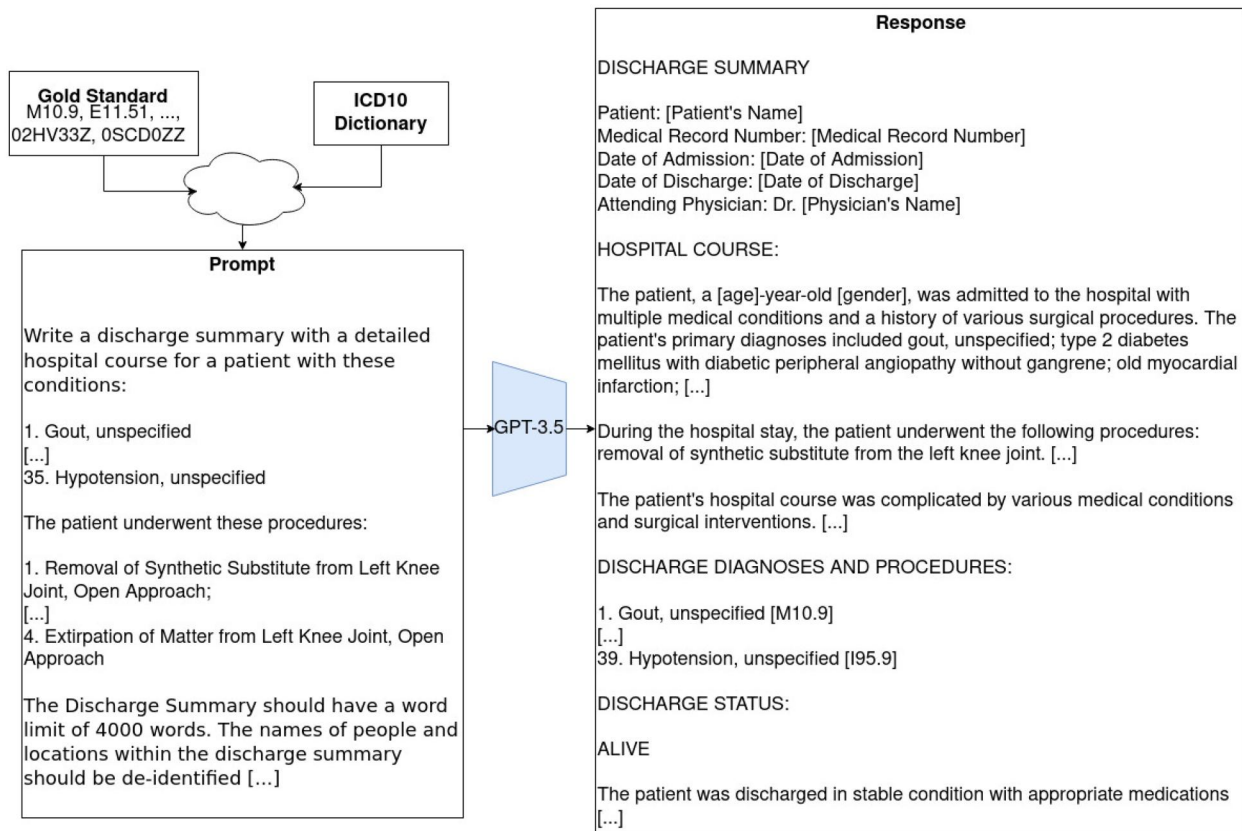
Received: January 2, 2024; Revised: May 12, 2024; Editorial Decision: May 17, 2024; Accepted: May 22, 2024

© The Author(s) 2024. Published by Oxford University Press on behalf of the American Medical Informatics Association.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.







**Figure 1.** An example generation of a synthetic discharge via GPT-3.5.

Micro-averaging assigns equal weight to each prediction, favoring high-population classes (eg, hypertension). Macro-averaging, in contrast, computes the performance for each unique label and averages across the label space, giving each label's average result equal weight regardless of their population. This highlights poor performance in less common classes. Our primary evaluation metrics common with the majority of previous work are micro-F1 and macro-F1 scores. Metrics are further explained in [Supplementary Material 3](#).

#### GPT's coding on real clinical notes

We used the Azure AI Services API (<https://azure.microsoft.com/en-gb/products/ai-services>) to test GPT's ability to assign diagnosis codes based on real clinical notes. The API returns a free text response which we further processed to retrieve the predicted ICD-10 codes. In postprocessing, we verify the response format. For correctly structured arrays of JSON objects we simply extract the predictions. For incorrectly structured outputs, we employ a regular expression pattern to extract all diagnoses and ICD code pairs. The result is a list of predicted diagnoses and corresponding ICD-10 codes for each clinical note. [Figure 3](#) illustrates the API call workflow.

For reproducibility, we specified the model version and the API version as “gpt-3.5-turbo-0613” and “2023-03-15-preview,” respectively. All parameters were set to zero for deterministic responses from *G*, including temperature, top *P*, frequency penalty, and presence penalty. The system prompt directed *G* to act as a clinician assigning ICD-10 diagnosis codes to clinical notes, specifying the expected output format as JSON objects with keys “diagnosis” and “icd\_code.” Refer

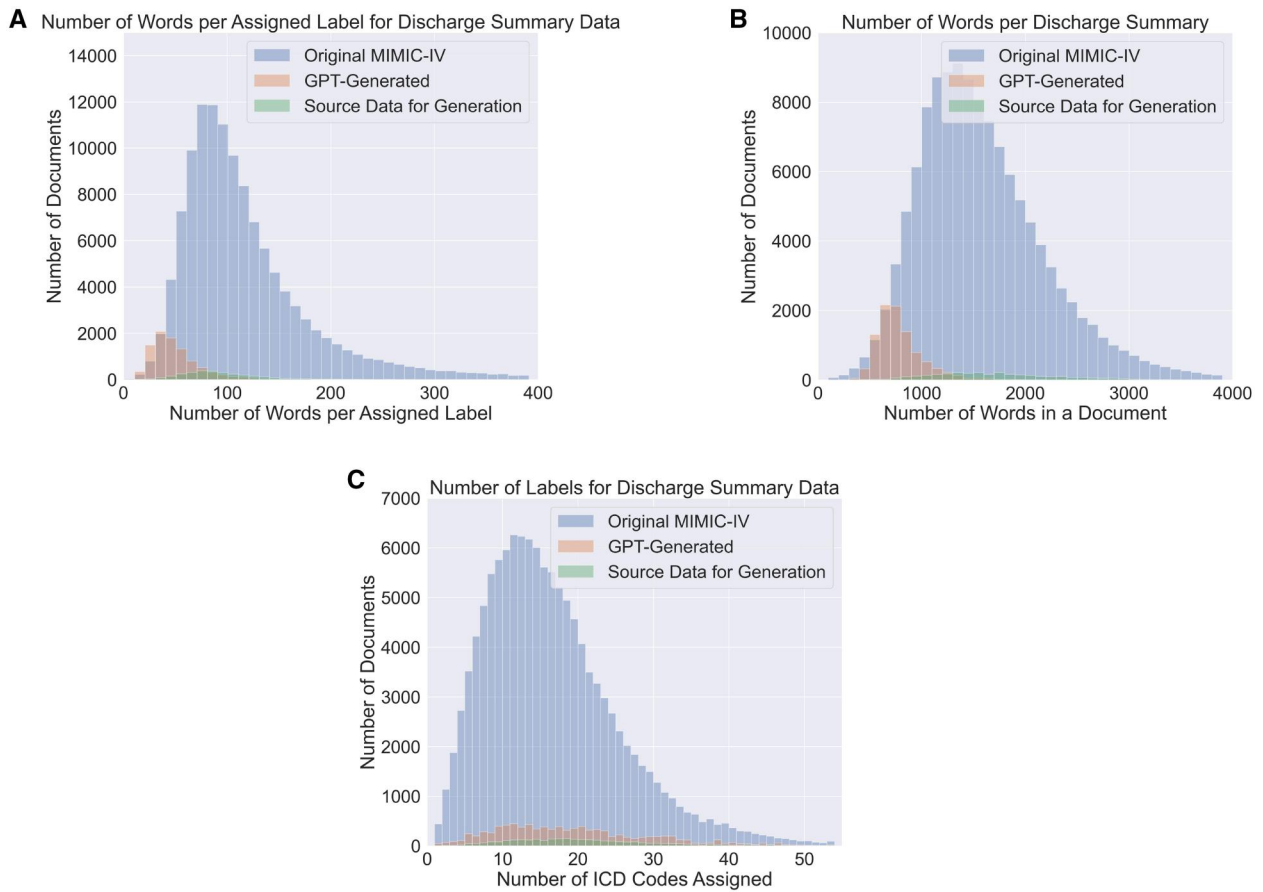
to [Figure 3](#) for all hyperparameter details. Our code implementation can be found in our Github repository ([https://github.com/EdinburghClinicalNLP/chatgpt\\_icd\\_coding](https://github.com/EdinburghClinicalNLP/chatgpt_icd_coding)).

We opted out of human review of the data for 2 reasons. First, the terms of the data use agreement of MIMIC-IV (<https://physionet.org/news/post/415>) did not grant us the authority to permit a third party to process the data for abuse detection. Second, we assessed the likelihood of harmful misuse to be low given the sensitive nature of the clinical notes.

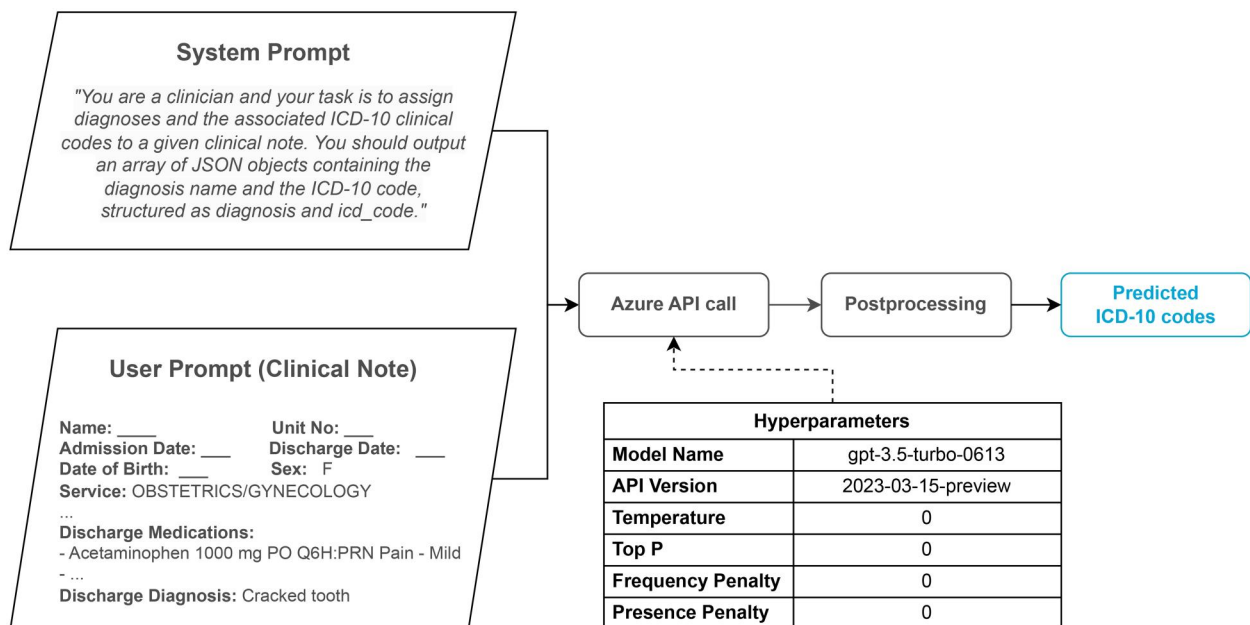
In our evaluation of GPT's performance, we have also employed hierarchical evaluation techniques—set-based hierarchical evaluation<sup>31</sup> and Count-Preserving Hierarchical Evaluation (CoPHE).<sup>32</sup> These metrics award partial credit to mispredicted labels by extending prediction and gold standard sets with their ancestor labels. While set-based evaluation ancestor labels track only the presence of descendants, in CoPHE, ancestor labels link to the count of descendants, penalizing over- and under-predictions within code families.

Comparing set-based and CoPHE results helps to evaluate the model's tendency to over-/under-predict. A lower CoPHE score indicates this phenomenon. See [Supplementary Material 4](#) for details on calculating hierarchical scores.

We utilize macro-averaged metrics from weak hierarchical confusion matrices<sup>10</sup> to summarize in-family versus out-of-family (OOF) prediction errors. These metrics are chosen to explore how expanding the population of codes within code families to a minimum of 100 instances impacts within-family performance. Within-family errors involve false positives that align with false negatives within the same family in the gold standard. On the other hand, an OOF error for a false negative in the gold standard lacks a false positive



**Figure 2.** A comparison between the discharge summary data in MIMIC-IV, seed MIMIC discharge summaries for generation (the source data), and the generated discharge summaries. Subfigures 2A and 2B focus on the number of words in documents, indicating that the GPT-generated data generally contains fewer words overall and per assigned label compared to the real data from MIMIC-IV. Subfigure 2C demonstrates that, although there's a variance in document size, the distribution of the number of labels per document remains relatively similar across the datasets.



**Figure 3.** The workflow of the GPT-3.5 prediction. We used Azure AI Services API to query GPT-3.5 and we employed a postprocessing step to extract the predicted diagnoses and ICD-10 codes for each clinical note.

within the prediction set from the same family to match. Our primary goal in generating synthetic data is to reduce OOF errors, enhancing true positive predictions or ensuring mispredictions occur within family.

### GPT-3.5's coding on synthetic data

The prompt asked *G* to code the conditions and procedures mentioned in the document it generated. This experiment tested *G*'s ability to assign ICD-10 codes to concepts presented in their standard descriptions. Alongside this, the prompt required creating a patient's social and family history, which might have led to the model introducing new conditions like substance abuse and potentially coding them, despite not being part of the initial prompt.

### Acceptability of generated data in clinical practice

Four clinical professionals (coauthors S.B., L.D., M.H., and R.S.P.) assessed the quality of the generated data. As the data was generated based on labels associated with MIMIC-IV discharge summaries, this evaluation included both synthetic discharge summaries generated by *G* and discharge summaries from MIMIC-IV. The clinicians were presented with 20 discharge summaries—10 synthetic and 10 real (based on whose adjusted gold standard the synthetic ones were generated).

Each discharge summary was assessed for:

- Correctness—accuracy in describing patient conditions and procedures;
- Informativeness—clarity and sense in supporting information (eg, test results, medication suggestions);

- Authenticity (patient)—whether such a patient could exist;
- Authenticity (clinical scenario)—whether the hospital course was plausible as reported; and
- Acceptability—suitability of the document for clinical use.

Additionally, they separately evaluated correctness and informativeness for both nonlow-resource and low-resource labels to gauge *G*'s ability to generate low-resource data. Scores from 1 (failure to perform) to 5 (perfect performance) were assigned to each metric, with accompanying comments justifying the score. An example evaluation by a clinician can be seen in [Figure 4](#).

## Results

### Local neural model evaluation

We assessed performance across 3 code subsets: the entire codeset in MIMIC-IV (*overall*), a restricted generation set ( $f_{(\text{gen})}$ ) containing only the 114 low-population candidate generation labels, and a set of codes from the code families present in  $f_{(\text{gen})}$  ( $f$ ). Results are shown in [Table 2](#). Baseline results for CAML and LAAT align closely with prior findings by Nguyen et al<sup>25</sup> for *overall* metrics. While Nguyen et al<sup>25</sup> did not report on Multi-Res CNN, its performance trends were similar to CAML and LAAT in comparison to Edin et al<sup>24</sup> on nonfiltered codesets. LAAT excels in micro-F1, Multi-Res CNN leads in macro-F1, and CAML generally lags behind. Baseline models outperform augmented ones in *overall* micro-F1, a common observation when enhancing

Synthetic Discharge Summary	Non-Low-Resource Labels
DISCHARGE SUMMARY  Patient: [Patient's Name] Medical Record Number: [Medical Record Number] Date of Admission: [Date of Admission] Date of Discharge: [Date of Discharge] Attending Physician: Dr. [Physician's Name]	M10.9: Gout, unspecified; E11.51: Type 2 diabetes mellitus with diabetic peripheral angiopathy without gangrene; [...] 0SCD0ZZ: Extirpation of Matter from Left Knee Joint, Open Approach
HOSPITAL COURSE:  The patient, a [age]-year-old [gender], was admitted to the hospital with multiple medical conditions and a history of various surgical procedures. The patient's primary diagnoses included gout, unspecified; type 2 diabetes mellitus with diabetic peripheral angiopathy without gangrene; old myocardial infarction; [...] mechanical loosening of internal left knee prosthetic joint, initial encounter; [...]	<b>Low-Resource Labels</b> T84.033A: Mechanical loosening of internal left knee prosthetic joint, initial encounter
During the hospital stay, the patient underwent the following procedures: removal of synthetic substitute from the left knee joint, open approach; [...]	<b>Correctness</b> Correctness Score - Non-Low-Resource (1-5): 5 Correctness Score - Low-Resource (1-5): 5  Correctness Comment: <i>All codes are present</i>
The patient's hospital course was complicated by various medical conditions and surgical interventions. The patient received appropriate medical management, [...]	<b>Informativeness</b> Informativeness Score - Non-Low-Resource (1-5): 3 Informativeness Score - Low-Resource (1-5): 2  Informativeness Comment: <i>Within the discharge summary text, the code descriptions are listed in a long paragraph and provided with minimal context or narrative. The description of the low resource code is limited to naming the code only.</i>
The patient's surgical procedures, including the removal of synthetic substitute from the left knee joint and insertion of a spacer, [...]	<b>Authenticity</b> Authenticity Score - Patient (1-5): 1 Authenticity Score - Scenario (1-5): 2  Authenticity Comment: <i>This would be a very complex discharge summary to write given the very large number of codes. Unfortunately the summary has little or no attempt at providing a narrative or context for why the given conditions occurred. There is minimal effort at describing the treatment, or joining up the conditions, or providing a sense of what is most important in this clinical case.</i>
DISCHARGE DIAGNOSES AND PROCEDURES:  1. Gout, unspecified [...]           39. Hypotension, unspecified	<b>Acceptability</b> Acceptability Score (1-5): 1  Acceptability Comment: <i>Too vague, generic phrases, no context or narrative.</i>
DISCHARGE STATUS:  ALIVE  The patient was discharged in stable condition with appropriate medications [...]	

**Figure 4.** An example evaluation of a synthetic discharge summary by a clinical expert.

**Table 2.** A comparison between local neural network models (MRCNN stands for Multi-Res CNN) trained on baseline (*base*) and augmented (*aug*) training sets is evaluated using micro- and macro-averaged F1 scores (*mi* and *ma* respectively) on 3 codesets—*overall* on all codes present in MIMIC-IV; *f* comprising all codes within the families we chose for generation; and  $f_{(\text{gen})}$  corresponding to candidate codes used in generation with a population of at most 100 in the training set.

Experiment	F1 ↑						WHCM error ↓			
	$mi_{\text{overall}}$	$ma_{\text{overall}}$	$mi_f$	$ma_f$	$mi_{f(\text{gen})}$	$ma_{f(\text{gen})}$	$\text{OOF}_f$	$\text{IF}_f$	$\text{OOF}_{f(\text{gen})}$	$\text{IF}_{f(\text{gen})}$
CAML <sub>base</sub>	<b>53.65</b>	3.87	<b>38.43</b>	3.03	17.41	6.64	66.53	25.05	83.81	9.83
CAML <sub>aug</sub>	53.54	<b>3.90</b>	38.41	<b>3.78</b>	<b>20.68</b>	<b>11.86</b>	<b>65.98</b>	<b>23.77</b>	<b>79.79</b>	<b>9.17</b>
LAAT <sub>base</sub>	<b>57.29</b>	<b>6.18</b>	<b>43.59</b>	4.96	<b>26.79</b>	14.48	58.57	<b>28.35</b>	74.03	12.20
LAAT <sub>aug</sub>	57.18	6.09	43.36	<b>5.38</b>	25.70	<b>14.98</b>	<b>55.93</b>	29.78	<b>73.65</b>	<b>11.98</b>
MRCNN <sub>base</sub>	<b>55.66</b>	6.40	40.16	5.04	26.80	13.92	52.72	32.41	<b>69.68</b>	15.24
MRCNN <sub>aug</sub>	54.69	<b>6.46</b>	<b>42.69</b>	<b>5.85</b>	<b>30.39</b>	<b>17.68</b>	<b>49.65</b>	<b>32.36</b>	70.41	<b>10.22</b>

The highest score in each metric for each model pair (baseline versus augmented) is highlighted in bold. Weak hierarchical confusion matrix (WHCM) error rates are produced for codesets *f* and  $f_{(\text{gen})}$ . Performance on the common test set is reported using the macro-averaged proportion of errors that were out-of-family (OOF) and within family (IF). The best (lowest) error rate for each error type for each model pair (baseline versus augmented) is presented in bold.

**Table 3.** Results of GPT-3.5's coding ability on real and self-generated data.

Evaluation set	Leaf-only			Set-based			CoPHE		
	P	R	F1	P	R	F1	P	R	F1
GPT-3.5 real	9.46	33.51	14.76	10.59	44.87	17.13	10.30	44.33	16.72
GPT-3.5 synthetic	59.06	40.72	48.20	66.46	41.32	50.96	67.20	41.55	51.35
Best baseline (LAAT) real	60.42	54.46	57.29	61.28	54.50	57.68	60.84	54.33	57.39

lower-resource label performance. Nonetheless, macro-F1 scores improved for 2 out of 3 models within the *overall* codeset and for all models on *f* and  $f_{(\text{gen})}$ . Multi-Res CNN and CAML macro-F1 scores display sizable relative improvement (26% and 78%, respectively) in  $f_{(\text{gen})}$ .

Augmented models performed on par with or outperformed baseline models in micro-F1 scores for *f* and  $f_{(\text{gen})}$ . Augmented Multi-Res CNN outperforms its baseline in micro-F1 for both *f* and  $f_{(\text{gen})}$ , indicating benefits for the code family from augmenting less-populous members. Augmented LAAT shows improvement in macro-F1 in both *f* and  $f_{(\text{gen})}$  but lags in micro-F1. LAAT's performance may have been biased toward high-population classes and the augmentation's boosting of low-frequency classes (misrepresenting their frequency) may have introduced confusion. Apart from having a recurrent encoder (Bi-LSTM), the LAAT model employed in this experiment is about twice the size of the Multi-Res CNN (21.9M versus 11.9M parameters). This added model complexity may have enabled better performance on already frequent labels, but increased the need for more examples of lower-resource labels.

Comparing within-family and out-of-family errors (Table 2), augmented models generally exhibit fewer out-of-family errors on *f* and  $f_{(\text{gen})}$ . An exception to this is augmented Multi-Res CNN in  $f_{(\text{gen})}$ , whose slight increase in OOF came with a sizable reduction in within-family error.

Unlike baseline models, augmented ones occasionally predicted codes absent from the original data, although incorrectly, except for 1 correctly predicted code (S02.63XA) by a Multi-Res CNN model trained on augmented data. While consistent enhancement in zero-shot code performance was not achieved through augmentation, the potential for improvement is evident.

### GPT's coding ability on real and synthetic data

We examined *G*'s ability to code real MIMIC-IV documents and generate coded documents with explicit code

descriptions in the prompt. The results (Table 3) show that the performance on prompt-guided self-generated (synthetic) data resembles that of local models on the MIMIC-IV test set, not surpassing it. Hierarchical metrics show higher precision, recall, and consequently, F1-score in CoPHE compared to set-based hierarchical evaluation indicating errors coming from within-family misprediction, rather than incorrectly estimating the number of expected labels.

However, the performance on the MIMIC-IV test set is notably low, especially in precision. The improvement in the precision from leaf-only results to hierarchical is minimal. This implies that incorrect predictions were more likely to be out-of-family. Moreover, results on CoPHE are lower than on the set-based hierarchical evaluation indicating a tendency of the model to over-/under-predict within the scope of the family—an issue previously reported in local ICD-coding models<sup>32</sup> and present in the reported hierarchical results for baseline LAAT.

These results demonstrate that *G* can identify ICD-10 codes in self-generated mentions based on provided descriptions if presented within the prompts. Its performance when tasked with standard ICD coding without explicitly identified concepts or nonstandard surface forms of the concepts significantly deteriorates.

### Acceptability of generated data in clinical practice

We calculated the inter-evaluator agreement for the 7 metrics using Fleiss' kappa ( $\kappa$ ).<sup>33</sup> As  $\kappa$  is designed for categorical variables and does not fully capture ordinal scores, also produced the mean scores for each metric ( $\mu$ ). The results are presented in Table 4.

The evaluators' agreement was poor ( $\kappa < 0$ ) in examples from MIMIC-IV for Correctness and Informativeness of nonlow-resource codes, and the acceptability of the discharge summaries. For the other metrics, a  $\kappa > 0$  was reached but never exceeded 0.4 (lower than moderate agreement). All

**Table 4.** Evaluator agreement ( $\kappa$ ) and mean scores ( $\mu$ ) for samples from MIMIC-IV (real), versus GPT-generated (synthetic) data

Metrics	$\kappa_{\text{real}}$	$\mu_{\text{real}}$	$\kappa_{\text{synthetic}}$	$\mu_{\text{synthetic}}$
Correctness—nonlow resource	-0.386	4.175	-0.163	4.375
Correctness—low resource	0.043	4.350	0.206	4.525
Informativeness—nonlow resource	-0.155	4.550	-0.220	2.775
Informativeness—low resource	0.241	4.675	-0.277	3.000
Authenticity—patient	0.340	4.750	-0.078	3.150
Authenticity—scenario	0.373	4.775	-0.333	2.250
Acceptability	-0.056	4.550	0.035	2.225

mean scores are higher than 4. Hence, while the clinicians disagreed on the exact scores, they rated real discharge summaries positively. The disagreement may be due to clinicians being UK-based with significant differences in reporting style within the United Kingdom and the United States (where MIMIC-IV is from).

For GPT-generated summaries, slight agreement was seen in acceptability, and fair agreement in the correctness of low-resource labels. All other metrics had poor agreement. Both correctness metrics scored above 4, with low-resource correctness surpassing 4.5—an encouraging outcome for our primary goal of generating low-resource code data. Mean informativeness in the low-resource scenario and authenticity of scores were at least 3. Once again, performance on the low-resource codes exceeded nonlow-resource codes. Other metrics had  $\mu$  scores above 2. Informativeness and authenticity for nonlow-resource codes had a poor agreement ( $\kappa < 0$ ), while acceptability had some agreement with the lowest mean score of 2.225.

While *G* generally produces correct notes, the clinical evaluators have identified several challenges in the generation of natural-looking clinical notes:

#### GPT-3.5 tends to do verbatim reproductions of the prompted diagnoses list

*G* tends to copy all concepts mentioned in the prompt when generating a clinical note. While instruction following is a desirable behavior, excluding noncrucial details is essential when generating a natural-looking clinical note. Real clinical notes often omit irrelevant and less critical findings for brevity, particularly if the information is inferrable from surrounding contexts such as medications and treatments. For instance, *G* unnecessarily noted a normal BMI.

#### GPT-3.5 may phrase diagnoses in an unnatural manner

*G* tends to use an overly technical and unnatural style when specifying diagnoses. For instance, *G* mentioned “anaemia, which was unspecified,” in the generated clinical note as it was prompted with “D64.9: Anemia, unspecified.” *G* also occasionally introduces vague phrases (eg, “geriatric team provided supportive care, including behavioural interventions and medication management”) without further detail. This contrasts with the more streamlined language of real clinical notes.

#### GPT-3.5 lacks details when introducing supporting information

*G* tends to introduce crucial supporting information without sufficient details. For instance, *G* mentioned “Following a traumatic event” without further specification of the

mentioned traumatic event, which is unacceptable in the clinical setting. This omission limits the overall informativeness of the patient’s medical context, potentially hindering the notes’ usability for a comprehensive view.

#### GPT-3.5 may introduce spurious supporting information

*G* sometimes introduces improbable but possible details. For instance, *G* overemphasized the significance of a patient’s anxiety disorder regarding an episode of syncope and a subsequent facial fracture, which the clinicians consider unlikely.

#### GPT-3.5 failed to present diagnoses as interconnected events

*G* does not effectively present diagnoses as interconnected, resulting in fragmented notes that lack coherence. The clinicians described *G*-generated clinical notes as collections of unrelated facts. For example, *G* presented complications of Type 1 diabetes mellitus separately without illustrating their relation. The lack of coherence between diagnoses may impede the plausibility of the clinical note and undermine the overall acceptability and usefulness of synthetic notes.

#### GPT-3.5 failed to prioritize and emphasize critical diagnoses

*G* struggles to prioritize diagnoses based on clinical significance, which undermines the authenticity of the portrayed scenario. For example, *G* often places critical conditions on the same level as minor issues, such as impacted ear wax, cataracts, and conjunctival hemorrhage. Hence, we concluded that *G* struggles to effectively convey the relative clinical significance of certain diagnoses.

## Discussion, conclusion, and future studies

In this work, we have investigated the capability of GPT-3.5’s potential in augmenting ICD-10 coding for local neural models in low-resource scenarios. While overall performance dipped with synthetic data augmentation, filtered codeset evaluation showed improvements, especially in advanced models like LAAT and Multi-Res CNN. Error analysis indicated augmented models made fewer out-of-family predictions, with some shift to within-family errors (closer to the correct answer). Augmentation showed promise in improving the prediction of generated codes and their siblings. Zero-shot labels did not consistently benefit from the augmentation, emphasizing the need for real data in augmentation success. However, a zero-shot code learned from the synthetic data was predicted correctly. The potential of LLM-generated discharge summaries should further be explored with different (eg, local or specialized) LLMs, prompt engineering, and further supplementing the prompt with external knowledge (eg, from ontologies).

In guided synthetic settings with ICD-10 descriptions, GPT-3.5 showed partial code identification ability displaying lesser over-/under-prediction tendencies than previously reported local models. It, however, struggled in the realistic scenario without in-prompt aid, performing below locally-trained models. Hence, the explored setup of producing a synthetic document based solely on the associated ICD codes is unsuitable for deployment in a clinical setting.

Clinician-evaluated synthetic discharge summaries showed correctness in individual codes, yet lacked naturalness and coherence compared to real data, resulting in lower informativeness, authenticity, and acceptability scores. Synthetic



summaries failed to represent holistic patient narratives or prioritize critical diagnoses.

One potential solution to generating synthetic discharge summaries involves restructuring the prompt to order diagnoses chronologically, providing their corresponding timestamps. This could guide LLMs in creating synthetic notes mirroring the chronological progression of a patient's medical journey, enhancing coherence and prioritization.

Another promising solution is to retrieve real clinical notes as in-context learning examples to help guide the generation process<sup>34</sup> to aid LLMs in generating more realistic and coherent content. As this study focuses on evaluating LLMs' existing capability, we opted to evaluate it in a zero-shot framework. Future work may explore this idea's potential for generating more realistic-looking clinical notes.

## Limitations

In this study, while the annotation experts are involved as coauthors, we ensured that they were independent from the development of the algorithms that involved the synthetic data. While the evaluation utilized few clinical experts ( $n=4$ ), they provided sufficient expertise in evaluating the notes. The study was blinded with respect to the real/synthetic status of documents, but according to the experts, the synthetic data differed from real enough to be distinguishable.

## Author contributions

Matúš Falis proposed exploring data augmentation for ICD coding with GPT-3.5 and developed the main ideas of the publication with Aryo Pradipta Gema, Hang Dong, Alexandra Birch, and Beatrice Alex. Matúš Falis, Aryo Pradipta Gema, Hang Dong, Alexandra Birch, and Beatrice Alex contributed to the experimental design. Matúš Falis designed and performed the document synthesis and the evaluation of GPT's ICD-10 coding performance on GPT-generated data. Aryo Pradipta Gema evaluated GPT's ICD-10 coding performance on MIMIC-IV data. Matúš Falis, Aryo Pradipta Gema, Hang Dong, and Luke Daines designed the evaluation of discharge summaries by clinical staff. Luke Daines, Sidharth Basetti, Michael Holder, and Rose S. Penfold performed the clinical expert evaluation of the data, and Matúš Falis and Aryo Pradipta Gema analyzed its results. Matúš Falis and Aryo Pradipta Gema wrote the manuscript. All authors participated in editing the manuscript. Beatrice Alex and Alexandra Birch supervised the project and provided feedback and input on the experiments and analyses.

## Supplementary material

Supplementary material is available at *Journal of the American Medical Informatics Association* online.

## Funding

This work is supported by the United Kingdom Research and Innovation (grant EP/S02431X/1), UKRI Centre for Doctoral Training in Biomedical AI at the University of Edinburgh, School of Informatics. H.D. was supported by the Engineering and Physical Sciences Research Council (EPSRC, grant EP/V050869/1), Concur: Knowledge Base Construction and

Curation. R.S.P. is a fellow on the Multimorbidity Doctoral Training Programme for Health Professionals, which is supported by the Wellcome Trust (grant 223499/Z/21/Z). B.A. was funded by Legal and General PLC as part of the Advanced Care Research Centre and by the National Institute for Health Research (NIHR) for the Artificial Intelligence and Multimorbidity: Clustering in Individuals, Space and Clinical Context (AIM-CISC, grant NIHR202639). The funders had no role in conduct of the study, interpretation or the decision to submit for publication. The views expressed are those of the authors and not necessarily those of the funders..

## Conflicts of interest

We have identified no competing interests.

## Data availability

The synthetic discharge summary data generated as part of this study will be shared on reasonable request to the corresponding author upon presenting a certificate of completion of the CITI Data or Specimens Only Research course from the Collaborative Institutional Training Initiative program (<https://physionet.org/about/citi-course/>). The data has been accepted for publication and will be made available via PhysioNet (<https://doi.org/10.13026/bnc2-1a81>).

## References

- Dong H, Falis M, Whiteley W, et al. Automated clinical coding: what, why, and where we are? *NPJ Digit Med*. 2022;5(1):159.
- Johnson AE, Bulgarelli L, Shen L, et al. MIMIC-IV, a freely accessible electronic health record dataset. *Sci Data*. 2023;10(1):31.
- Mullenbach J, Wiegrefe S, Duke J, et al. Explainable prediction of medical codes from clinical text. In: *Proceedings of NAACL-HLT*; 2018:1101-1111.
- Dong H, Suárez-Paniagua V, Whiteley Explainable automated coding of clinical notes using hierarchical label-wise attention networks and label embedding initialisation. *J Biomed Inform*. 2021;116:103728.
- Kim BH, Ganapathi V. Read, attend, and code: pushing the limits of medical codes prediction from clinical notes by machines. In: *Machine Learning for Healthcare Conference*. PMLR; 2021:196-208.
- Rios A, Kavuluru R. Few-shot and zero-shot multi-label learning for structured label spaces. In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing. Conference on Empirical Methods in Natural Language Processing*. Vol. 2018. NIH Public Access; 2018:31-32.
- Song C, Zhang S, Sadoughi N, et al. Generalized zero-shot text classification for icd coding. In: *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*; 2021:4018-4024.
- Ren W, Zeng R, Wu T, et al. 2022. Hicu: leveraging hierarchy for curriculum learning in automated icd coding, arXiv, arXiv:2208.02301, preprint: not peer reviewed.
- Wang S, Tang D, Zhang L, et al. Hienet: bidirectional hierarchy framework for automated icd coding. In: *International Conference on Database Systems for Advanced Applications*. Springer; 2022:523-539.
- Falis M, Dong H, Birch A, et al. Horses to zebras: ontology-guided data augmentation and synthesis for icd-9 coding. In: *Proceedings of the 21st Workshop on Biomedical Language Processing. Association for Computational Linguistics*; 2022.
- Kim D, Yoo H, Kim S. 2022. An automatic icd coding network using partition-based label attention, arXiv, arXiv:2211.08429, preprint: not peer reviewed.

12. Barros J, Rojas M, Dunstan J, et al. Divide and conquer: An extreme multi-label classification approach for coding diseases and procedures in Spanish. In: *Proceedings of the 13th International Workshop on Health Text Mining and Information Analysis (LOUHI)*; 2022:138-147.
13. Afkanpour A, Adeel S, Bassani H, et al. 2022. Bert for long documents: a case study of automated icd coding, arXiv, arXiv:2211.02519, preprint: not peer reviewed.
14. Ouyang L, Wu J, Jiang X, et al., et al. Training language models to follow instructions with human feedback. In: Koyejo S, Mohamed S, Agarwal A, eds. *Advances in Neural Information Processing Systems*. Vol 35. Curran Associates, Inc.; 2022:27730-27744. [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/b1efde53-be364a73914f58805a001731-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/b1efde53-be364a73914f58805a001731-Paper-Conference.pdf)
15. Touvron H, Lavril T, Izacard G, et al. 2023. Llama: open and efficient foundation language models, arXiv, arXiv:2302.13971, preprint: not peer reviewed.
16. Zhao WX, Zhou K, Li J, et al. 2023. A survey of large language models, arXiv, arXiv:2303.18223, preprint: not peer reviewed.
17. Singhal K, Azizi S, Tu T, et al. 2022. Large language models encode clinical knowledge, arXiv, arXiv:2212, preprint: not peer reviewed.
18. Ji Z, Lee N, Frieske R, et al. Survey of hallucination in natural language generation. *ACM Comput Surv*. 2023;55(12):1-38. <https://doi.org/10.1145/3571730>
19. Lee P, Bubeck S, Petro J. Benefits, limits, and risks of gpt-4 as an ai chatbot for medicine. *N Engl J Med*. 2023;388(13):1233-1239.
20. Lecler A, Duron L, Soyer P. Revolutionizing radiology with gpt-based models: Current applications, future possibilities and limitations of chatgpt. *Diagn Interv Imaging*. 2023;104(6):269-274.
21. Yeung JA, Kraljevic Z, Luintel A, et al. 2023. Ai chatbots not yet ready for clinical use, medRxiv:2023-03.
22. Kraljevic Z, Bean D, Shek A, et al. 2022. Foresight-deep generative modelling of patient timelines using electronic health records, CoRR.
23. Ghosh S, Evuru CK, Kumar S, et al. 2023. Dale: generative data augmentation for low-resource legal nlp, arXiv, arXiv:2310.15799, preprint: not peer reviewed.
24. Edin J, Junge A, Havtorn JD, et al. 2023. Automated medical coding on mimic-iii and mimic-iv: a critical review and replicability study, arXiv, arXiv:2304.10909, preprint: not peer reviewed.
25. Nguyen TT, Schlegel V, Kashyap A, et al. 2023. Mimic-iv-icd: a new benchmark for extreme multilabel classification, arXiv, arXiv:2304.13998, preprint: not peer reviewed.
26. Vu T, Nguyen DQ, Nguyen A. 2020. A label attention model for icd coding from clinical text, arXiv, arXiv:2007.06351, preprint: not peer reviewed.
27. Li F, Yu H. Icd coding from clinical text using multi-filter residual convolutional neural network. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol 34; 2020:8180-8187.
28. Mikolov T, Sutskever I, Chen K, et al. Distributed representations of words and phrases and their compositionality. *Adv Neural Inf Process Syst*. 2013;26(2):3111-3119.
29. Devlin J, Chang MW, Lee K, et al. 2018. Bert: pre-training of deep bidirectional transformers for language understanding, arXiv, arXiv:1810.04805, preprint: not peer reviewed.
30. Huang CW, Tsai SC, Chen YN. 2022. Plm-icd: automatic icd coding with pretrained language models, arXiv, arXiv:2207.05289, preprint: not peer reviewed.
31. Kosmopoulos A, Partalas I, Gaussier E, et al. Evaluation measures for hierarchical classification: a unified view and novel approaches. *Data Min Knowl Disc*. 2015;29(3):820-865.
32. Falis M, Dong H, Birch A, et al. Cophe: a count-preserving hierarchical evaluation metric in large-scale multi-label text classification. In: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*; 2021:907-912.
33. Fleiss JL. Measuring nominal scale agreement among many raters. *Psychol Bull*. 1971;76(5):378-382.
34. Lewis PSH, Perez E, Piktus A, et al. Retrieval-augmented generation for knowledge-intensive NLP tasks. In: Larochelle H, Ranzato M, Hadsell R, et al., eds. *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, Virtual*; 2020. <https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html>.