

1 **Title:** Assessing GPT-3.5 and GPT-4 in Generating International Classification of Diseases  
2 Billing Codes

3  
4 **Authors:** Ali Soroush<sup>1,2</sup>, Benjamin S. Glicksberg<sup>1,2</sup>, Eyal Zimlichman<sup>3,4</sup>, Yiftach Barash<sup>4,5</sup>,  
5 Robert Freeman<sup>6</sup>, Alexander W. Charney<sup>1,2</sup>, Girish N Nadkarni\*<sup>1,2</sup>, Eyal Klang\*<sup>1,2,4,7</sup>

6  
7 **Affiliations:**

- 8 1. Division of Data-Driven and Digital Medicine (D3M), Icahn School of Medicine at Mount  
9 Sinai, New York, New York, USA
- 10 2. The Charles Bronfman Institute of Personalized Medicine, Icahn School of Medicine at  
11 Mount Sinai, New York, New York, USA.
- 12 3. Central Management, Sheba Medical Centre, Ramat-Gan, Israel
- 13 4. Faculty of Medicine, Tel Aviv University, Tel Aviv, Israel.
- 14 5. Department of Diagnostic Imaging, Chaim Sheba Medical Center, Ramat Gan, Israel.
- 15 6. Mount Sinai Health System, New York, New York, USA
- 16 7. ARC Center for Digital Innovation, Chaim Sheba Medical Center, Ramat Gan, Israel.

17  
18  
19 **Author Contributions:** **Ali Soroush:** Methodology, Software, Formal Analysis, Investigation,  
20 Data Curation, Writing - Review and Editing, Visualization. **Benjamin S. Glicksberg:** Writing -  
21 Review and Editing. **Eyal Zimlichman:** Writing - Review and Editing. **Yiftach Barash:** Writing  
22 - Review and Editing. **Robert Freeman:** Writing - Review and Editing. Alexander W. Charney:  
23 Writing - Review and Editing. **Girish N Nadkarni\*:** Supervision, Writing - Review and Editing.  
24 **Eyal Klang\*:** Conceptualization, Methodology, Formal Analysis, Investigation, Data Curation,  
25 Writing - Original Draft, Writing - Review and Editing, Supervision.

26  
27 \* Equal Contribution

28  
29  
30 **Corresponding author:**

31 Ali Soroush, MD, MS  
32 Assistant Professor Of Medicine  
33 Henry D. Janowitz Division of Gastroenterology  
34 Division of Data-Driven and Digital Medicine  
35 The Charles Bronfman Institute of Personalized Medicine  
36 Icahn School of Medicine at Mount Sinai  
37 [Ali.soroush@mountsinai.org](mailto:Ali.soroush@mountsinai.org)

1 **Abstract:**

2 Background: Large Language Models (LLMs) like GPT-3.5 and GPT-4 are increasingly entering  
3 the healthcare domain as a proposed means to assist with administrative tasks. To ensure safe  
4 and effective use with billing coding tasks, it is crucial to assess these models' ability to generate  
5 the correct International Classification of Diseases (ICD) codes from text descriptions.

6 Objectives: We aimed to evaluate GPT-3.5 and GPT-4's capability to generate correct ICD  
7 billing codes, using the ICD-9-CM (2014) and ICD-10-CM and PCS (2023) systems.

8 Methods: We randomly selected 100 unique codes from each of the most recent versions of the  
9 ICD-9-CM, ICD-10-CM, and ICD-10-PCS billing code sets published by the Centers for  
10 Medicare and Medicaid Services. Using the ChatGPT interface (GPT-3.5 and GPT-4), we  
11 prompted for the ICD codes that corresponding to each provided code description. Outputs were  
12 compared with the actual billing codes across several performance measures. Errors were  
13 qualitatively and quantitatively assessed for any underlying patterns.

14 Results: GPT-4 and GPT-3.5 demonstrated varied performance across each ICD system. In  
15 ICD-9-CM, GPT-4 and GPT-3.5 achieved an exact match rate of 22% and 10%, respectively.  
16 13% (GPT-4) and 10% (GPT-3.5) of generated ICD-10-CM codes were exact matches. Notably,  
17 both models struggled considerably with the procedurally focused ICD-10-PCS, with neither  
18 GPT-4 or GPT-3.5 producing any exactly matched codes. A substantial number of incorrect  
19 codes had semantic similarity with the actual codes for ICD-9-CM (GPT-4: 60.3%, GPT-3.5:  
20 51.1%) and ICD-10-CM (GPT-4: 70.1%, GPT-3.5: 61.1%), in contrast to ICD-10-PCS (GPT-4:  
21 30.0%, GPT-3.5: 16.0%).

22 Conclusion: Our evaluation of GPT-3.5 and GPT-4's proficiency in generating ICD billing  
23 codes from ICD-9-CM, ICD-10-CM and ICD-10-PCS code descriptions reveals an inadequate  
24 level of performance. While the models appear to exhibit a general conceptual understanding of  
25 the codes and their descriptions, they have a propensity for hallucinating key details, suggesting  
26 underlying technological limitations of the base LLMs. This suggests a need for more rigorous  
27 LLM augmentation strategies and validation prior to their implementation in healthcare contexts,  
28 particularly in tasks such as ICD coding which require significant digit-level precision.

29

30 **Keywords:** Large Language Models, GPT-3.5, GPT-4, International Classification of Diseases,  
31 ICD Coding, Medical Terminology, Healthcare Applications

## 1 **Introduction:**

2 The International Classification of Diseases (ICD) terminology is the most widely used  
3 administrative coding system in the world (1-3). It provides a standardized representation of  
4 medical conditions and procedures, playing a critical role in clinical recordkeeping, public health  
5 surveillance, research, and billing (2).

6 Large Language Models (LLMs), such as GPT-3.5 and its more advanced successor GPT-4, have  
7 shown remarkable and varied capabilities that have potential to impact many domains (4-6).

8 LLMs leverage the power of state-of-the-art deep learning algorithms, trained on extensive  
9 textual data (6). These models are capable of impressive feats, ranging from correctly answering  
10 medical board exam questions to writing poetry (4, 7). In the realm of healthcare, they could  
11 potentially support clinicians by automating administrative tasks, offering decision support, or  
12 communicating with patients (8-10).

13 These models have an ability to create text that appears human-generated, offering hope that this  
14 technological advancement will unlock important downstream natural language processing tasks  
15 like assigning ICD codes based on unstructured medical text descriptions. However, the accuracy  
16 of LLMs for administrative tasks in medicine has not been thoroughly assessed yet. A prominent  
17 concern is the models' propensity to 'hallucinate,' i.e., generate plausible sounding, but factually  
18 incorrect information (11-14). Before LLMs can be used to automate burdensome administrative  
19 tasks like assigning ICD billing codes based on clinical documentation, the propensity for  
20 hallucination must be delineated and, ultimately, mitigated.

21 Our study analyzes how GPT-3.5 and GPT-4 fare in the context of ICD billing code generation.  
22 We systematically assess their performance in matching ICD billing codes to their descriptions  
23 across multiple versions of the classification system.

24

## 25 **Methods:**

### 26 *ICD Code Selection:*

27 We obtained most recent lists of ICD-9-CM-CM (2014), ICD-10-CM (2023), and ICD-10-PCS  
28 (2023) billing codes from the Centers for Medicare & Medicaid Services (CMS)(15). From each  
29 list, we randomly selected 100 unique codes, leading to a total of 300 ICD billing codes for our  
30 dataset.

### 31 *LLMs Used:*

32 For this study, we deployed two commercially available Large Language Models (LLMs): GPT-  
33 3.5 and GPT-4, both of which were developed by OpenAI (5, 16). The underlying data used to  
34 train these models has not been publicly released, but presumably contains a combination of all  
35 publicly available online data as well as private datasets. Their data is inclusive to September  
36 2021.

### 37 *ICD Code Generation:*

38 We used the public ChatGPT interface to produce the corresponding ICD billing code for each of  
39 the 300 code descriptions. We used the following natural language prompt, where <ICD coding  
40 system> refers to ICD-9-CM-CM, ICD-10-CM, or ICD-10-PCS: “*Dear GPT, we will feed you a*  
41 *list of <ICD coding system> descriptions. Please write the matching <ICD coding system>*”

42 *codes. Please, return in a table format.*” The descriptions were provided to the LLMs in batches  
43 of 10 per prompt to improve processing efficiency.

#### 44 *Performance Evaluation:*

45 To assess LLM billing code generation performance, we determined the number of exact code  
46 matches, billable codes, nonbillable codes, and nonexistent codes for the output codes.  
47 Nonbillable codes are generally more non-specific than billing codes. We labeled exact matches  
48 by comparing each output code with the original code for its originating description. We labeled  
49 billable codes and obtained their descriptions by matching output codes with any code from the  
50 same coding system present in the original CMS billing code lists. To label the nonbillable codes  
51 and obtain their descriptions, we applied the UMLS Metathesaurus (17) to the remaining  
52 unmatched codes. Any remaining codes without an assigned description were considered to be  
53 nonexistent codes. For all generated codes, we assessed *semantic* and *syntactic* similarity  
54 measures to interrogate the LLMs' broader understanding of the ICD coding systems. To assess  
55 semantic similarity, two physicians (EK, AS) assessed for any meaningful conceptual similarity  
56 between a generated billable code's description and its original code description. In terms of  
57 syntactic similarity, we ascertained if a generated code differed from its original code by a single  
58 character or less, including differences in code length. Syntactic and semantic similarity  
59 percentages were calculated by dividing each similarity count by the total number of generated  
60 codes.

#### 61 *Error Analysis:*

62 We conducted error analyses for each coding system to identify with which specific contexts the  
63 models excelled or struggled. We first measured the semantic and syntactic similarities of valid

64 billing codes that did not match with their original codes. To assess the effect of code and  
65 description complexity on model performance, we assessed the relationships between  
66 performance and code length and description length respectively. We could not assess this for  
67 ICD-10-PCS due to low number of correctly matched codes (3 for GPT-4 and 0 for GPT-3.5).  
68 We additionally performed a qualitative analysis of the nature and context of errors, focusing on  
69 how well the models handled unspecified conditions, complex conditions, and semantic  
70 variations in descriptions. Two physicians performed this analysis using a consensus approach.

#### 71 *Statistical Analysis:*

72 We summarized the ICD code generating performance for each model using descriptive  
73 statistics. We calculated counts and percentages for exact code matches, billable codes, and  
74 nonbillable codes exact match counts and percentages for each score across the ICD versions.  
75 Counts and percentages were similarly calculated for semantic and syntactic similarity. We used  
76 Fisher's exact test to compare each performance metric between GPT-4 with GPT-3.5. To assess  
77 the relationships between exact matches and code and description length respectively, we applied  
78 the Mann-Whitney U test. We reported the median and interquartile range (IQR) values for the  
79 lengths for each case. We considered a p-value less than 0.05 statistically significant, indicating a  
80 meaningful performance difference between the two models on a given ICD code dataset. We  
81 coded all analyses in Python (Version 3.9.16).

82

## 83 **Results**

### 84 *Performance Evaluation:*

85 When assessing the ICD code generation performance of GPT-3.5 and GPT-4 from our  
86 predefined list of code descriptions (**Supplementary tables 1-3**), we found variable results  
87 across the different ICD coding systems (**Table 1**). For ICD-9-CM, GPT-4 generated exactly  
88 matched codes for 22% of cases, billable codes for 72% of cases, nonbillable codes for 26% of  
89 cases, and nonexistent codes for 2% of cases. GPT-3.5 generated a significantly lower rate of  
90 exactly matched codes (10%,  $p=0.033$ ), but had similar rates of billable (76%,  $p=0.629$ ),  
91 nonbillable (20%,  $p=0.401$ ), and nonexistent codes (4%,  $p=0.683$ ). For the ICD-10-CM system,  
92 both models produced fewer exact match and nonbillable codes and more nonexistent codes.  
93 GPT-4 generated exactly matched codes for 13% of cases, billable codes for 77% of cases,  
94 nonbillable codes for 3% of cases, and nonexistent codes for 20% of cases. In comparison to  
95 GPT-4, GPT-3.5 generated similar rates of exactly matched (5%,  $p=0.081$ ), billable (67%,  
96  $p=0.156$ ), nonbillable (4%,  $p=1.000$ ), and nonexistent codes (28%,  $p=0.246$ ). GPT-4 and GPT-  
97 3.5 both had the lower ICD code generation performance for the ICD-10-PCS system, with both  
98 producing no exact match codes or nonbillable codes. GPT-4 was able to generate billable codes  
99 for 39% of cases, with the remainder (61%) of the codes being nonexistent GPT-3.5 had similar  
100 rates for billable (30%,  $p=0.234$ ) and nonexistent (70%,  $p=0.234$ ) codes.

101  
102 **Table 1:** Validity and Accuracy of GPT-generated ICD codes; abbreviations: ICD International  
103 Classification of Diseases

<b>ICD-9-CM System</b>	<b>GPT-4</b>	<b>GPT-3.5</b>	<b>P-value</b>
Exactly Matched Codes (n, % of total)	22 (22%)	10 (10%)	0.033
Billable Codes (n, % of total)	72 (72%)	76 (76%)	0.629
Nonbillable Codes (n, % of total)	26 (26%)	20 (20%)	0.401
Nonexistent Codes (n, % of total)	2 (2%)	4 (4%)	0.683



<b>ICD-10-CM System</b>	<b>GPT-4</b>	<b>GPT-3.5</b>	<b>P-value</b>
Exactly Matched Codes (n, % of total)	13 (13%)	5 (5%)	0.081
Billable Codes (n, % of total)	77 (77%)	67 (67%)	0.156
Nonbillable Codes (n, % of total)	3 (3%)	5 (5%)	1.000
Nonexistent Codes (n, % of total)	20 (20%)	28 (28%)	0.246
<b>ICD-10-PCS System</b>			
Exactly Matched Codes (n, % of total)	0 (0%)	0 (0%)	1.000
Billable Codes (n, % of total)	39 (39%)	30 (30%)	0.234
Nonbillable Codes (n, % of total)	0 (0%)	0 (0%)	1.000
Nonexistent Codes (n, % of total)	61 (61%)	70 (70%)	0.234

104

105 *Semantic and syntactic similarity:*

106 **Table 2** reveals that both models demonstrated a high level of semantic, or conceptual, similarity  
107 for all three ICD systems, though GPT-4 outperformed GPT-3.5. When assessing ICD-9-CM,  
108 GPT-4 achieved a semantic similarity of 60% versus 43% for GPT 3.5 (p=0.003). In the case of  
109 ICD-10-CM, GPT-4 and GPT-3.5 both attained higher semantic similarities (74% and 63%  
110 respectively) than they did for ICD-9-CM. However, both models struggled with the ICD-10-  
111 PCS system, showing semantic similarity for only 30% (GPT-4) and 16% (GPT-3.5) of the  
112 codes.

113 Rates of syntactic, or character-level, code similarity were more variable. GPT-4 displayed a  
114 high degree of syntactic similarity when generating ICD-9-CM codes (60%), but this measure  
115 dropped greatly for ICD-10-CM (36%) and even more so for ICD-10-PCS (3%). GPT-3.5  
116 demonstrated a similar, but consistently lower scoring pattern of syntactic similarity scores  
117 across all three ICD terminologies, with rates of 43% for ICD-9-CM, 19% for ICD-10-CM, and  
118 0% for ICD-10-PCS.

119

120 **Table 2:** Semantic and syntactic similarity of billable GPT-generated ICD codes for all  
 121 generated codes; abbreviations: ICD International Classification of Diseases

<b>ICD-9-CM</b>	<b>GPT-4</b>	<b>GPT-3.5</b>	<b>P-value</b>
Syntactically Similar Codes (n, %)	60 (60%)	43 (43%)	0.003
Semantically Similar Codes (n, %)	69 (69%)	56 (56%)	0.002
<b>ICD-10-CM</b>			
Syntactically Similar Codes (n, %)	36 (36%)	19 (19%)	0.007
Semantically Similar Codes (n, %)	74 (74%)	63 (63%)	0.008
<b>ICD-10-PCS</b>			
Syntactically Similar Codes (n, %)	3 (3%)	0 (0.0%)	0.240
Semantically Similar Codes (n, %)	30 (30%)	16 (16%)	0.003

122

123 *Error analysis:*

124 When repeating the similarity analysis on only the incorrect codes, we found that similar rates of  
 125 semantic similarity as in the overall analysis (**Table 3**). For the codes produced by GPT-4, 60.3%  
 126 of incorrect ICD-9-CM codes, 70.1% of incorrect ICD-10-CM codes, and 30% of incorrect ICD-  
 127 10-PCS codes were semantically similar. GPT-3.5 produced a similar semantic similarity pattern,  
 128 but with consistently lower scores: 51.1% for ICD-9-CM, 61.1% for ICD-10-CM, and 16.0% for  
 129 ICD-10-PCS.

130

131 **Table 3:** Semantic and syntactic similarity of GPT-generated billable GPT-generated ICD codes  
 132 for all incorrectly generated codes; abbreviations: ICD International Classification of Diseases

<b>ICD-9-CM</b>	<b>GPT-4</b>	<b>GPT-3.5</b>	<b>P-value</b>
Syntactically Similar Codes (n, %)	38 (48.7%)	33 (36.7%)	0.121
Semantically Similar Codes (n, %)	47 (60.3%)	46 (51.1%)	0.277
<b>ICD-10-CM</b>			
Syntactically Similar Codes (n, %)	36 (36%)	19 (19%)	0.007
Semantically Similar Codes (n, %)	74 (74%)	63 (63%)	0.008
<b>ICD-10-PCS</b>			
Syntactically Similar Codes (n, %)	3 (3%)	0 (0.0%)	0.240
Semantically Similar Codes (n, %)	30 (30%)	16 (16%)	0.003

Syntactically Similar Codes (n, %)	23 (26.4%)	14 (14.7%)	0.065
Semantically Similar Codes (n, %)	60 (70.1%)	58 (61.1%)	0.215
<b>ICD-10-PCS</b>			
Syntactically Similar Codes (n, %)	3 (3.0%)	0 (0.0%)	0.246
Semantically Similar Codes (n, %)	30 (30.0%)	16 (16.0%)	0.028

133

134 Our error analysis unveiled several findings. With all code systems, the models tended toward  
 135 generating more general nonbillable codes for complex and lengthy code descriptions such as  
 136 case 1 in **Supplementary Table 1** ("Other open skull fracture with cerebral laceration and  
 137 contusion, with moderate [1-24 hours] loss of consciousness", ICD-9-CM code 80363) other  
 138 similar cases 11, 14, 20. **Table 4** presents quantitative analyses of the associations between  
 139 model exact match performance and code length and description length respectively. Overall,  
 140 code and description length were inversely associated with exact matches. Significant  
 141 relationships were present for ICD-9-CM descriptions lengths and GPT-4 and GPT-3.5 exact  
 142 match performance, ICD-10-CM code lengths and GPT-4 and GPT-3.5 exact match  
 143 performance, and ICD-10-CM description length and GPT-4 exact match performance.

144

145 **Table 4:** ICD-9-CM and ICD-10-CM length of codes and descriptions association with exact  
 146 matches (Mann Whitney U test); abbreviations: ICD International Classification of Diseases,  
 147 IQR interquartile range

ICD-9-CM	GPT-4			GPT-3.5		
	Exact Matches	Not Exact Matches	P-value	Exact Matches	Not Exact Matches	P-value
Code length (median, IQR)	4.0 (4.0 - 5.0)	4.0 (4.0 - 5.0)	0.917	5.0 (4.0 - 5.0)	4.0 (4.0 - 5.0)	0.410
Description length (median, IQR)	38.0 (26.5 - 54.0)	46.0 (35.0 - 75.0)	0.081	27.5 (23.5 - 52.5)	46.5 (35.0 - 74.5)	0.025
<b>ICD-10-CM</b>						
	<b>GPT-4</b>			<b>GPT-3.5</b>		

	<b>Exact Matches</b>	<b>Not Exact Matches</b>	<b>P-value</b>	<b>Exact Matches</b>	<b>Not Exact Matches</b>	<b>P-value</b>
Code length (median, IQR)	5.0 (4.0 - 7.0)	7.0 (7.0 - 7.0)	< 0.001	5.0 (4.0 - 5.0)	7.0 (6.0 - 7.0)	0.001
Description length (median, IQR)	49.0 (38.0 - 63.0)	84.0 (62.0 - 112.5)	<0.001	63.0 (38.0 - 77.0)	80.0 (59.5 - 111.0)	0.131

148

149 Certain condition categories consistently had worse exact match accuracy and semantic  
150 similarity scores, driven again by a tendency for the LLMs to rely on non-billing codes for  
151 complex scenarios. For example, generic non-billable codes were generated for pregnancy-  
152 related disorders such as cases 9 ("Major puerperal infection", ICD-9-CM code 67082) and 20  
153 ("Mild or unspecified pre-eclampsia, unspecified as to episode of care", ICD-9-CM code 64224).  
154 Even when the models generated billable codes, the codes often differed in severity or  
155 specificity, as observed in ICD-9-CM case 5 ("Injury to hypoglossal nerve", ICD-9-CM code  
156 9517) and others (7, 12, 32, 33, 50, 62).

157 In the context of ICD-10-CM code generation (**Supplemental Table 2**), GPT-4 had trouble  
158 achieving the high level of specificity required by this coding system, as only 14 generated codes  
159 were exact matches. Nevertheless, the overall level of specificity was generally higher than with  
160 ICD-9-CM, as demonstrated by case 7 ("Acute embolism and thrombosis of unspecified deep  
161 veins of right lower extremity", ICD-10-CM code I82401), case 15 ("Benign carcinoid tumor of  
162 the transverse colon", ICD-10-CM code D3A023), and others (14, 27, 34, 46). The low number  
163 of non-billable codes (3%) generated and high rate of semantic similarity for generated codes  
164 supports the overall trend toward specificity for generated ICD-10-CM code. Incorrectly  
165 assigned billable codes had errors related to highly specific features of the description such as  
166 injury type, laterality, complexity, and complications (e.g. cases 1, 14, 41, 75). Nearly all valid  
167 codes were aligned with the temporal features of the code (initial vs subsequent encounter),

168 which is represented by the last letter in the ICD-10-CM code. Nonexistent codes were most  
169 frequent for codes related to maternal care (0/2, 0.0%), vehicle-related injuries (4/8, 50.0%), and  
170 joint-related conditions (4/7, 57.1%).

171 For ICD-10-PCS, GPT-4 was unable to correctly match any codes. However, the model was able  
172 to achieve some semantic similarity for 30% of generated codes. These codes differed in much  
173 more in terms of anatomical locations, laterality, maneuvers, and procedures. There was a high  
174 degree of concordance for procedure approach, which is represented by the last letter in the ICD-  
175 10-PCS code.

176

## 177 **Discussion**

178 Our study is the first to evaluate the ICD billing code mapping performance of both GPT-3.5 and  
179 GPT-4. Their underlying LLM technology holds the promise of automating the mapping of these  
180 core administrative terminologies in healthcare, with significant implications for billing, clinical  
181 decision making, quality improvement, research, and health policy.

182 The findings revealed that GPT-4 generally outperformed GPT-3.5 in generating exact match  
183 ICD billing codes, billing codes, and non-billing codes though overall performance was not  
184 reliable across both. Performance varied considerably across different versions of the ICD  
185 system. Both models showed some success with ICD-9-CM, but struggled more with the newer  
186 ICD-10-CM and ICD-10-PCS systems. This reflects the increased complexity, granularity, and  
187 comprehensiveness of these more recent versions (18) and highlights the difficulty faced by  
188 LLMs in handling a critical healthcare-related task.

189 Interestingly, despite their struggle with exact code generation, the models often generated codes  
190 that were semantically or syntactically similar to the actual codes. Both models demonstrated a  
191 better grasp of the more commonly used and queried ICD-9-CM and ICD-10-CM diagnosis code  
192 systems compared to the ICD-10-PCS procedure code system. This suggests that the base  
193 versions of these LLMs can parse the general descriptive nature of ICD codes, but could not  
194 attain precision, potentially due to a tendency to generalize or, worse, hallucinate data when the  
195 input is sparse or ambiguous (6, 12). The performance differences across systems may be due to  
196 the broader availability of ICD-9-CM and ICD-10-CM coding references in the LLM training  
197 data (19), which is drawn from publicly available web and print data. Alternatively, this behavior  
198 may be due to the underlying architecture of the LLMs we tested, as they have been trained to  
199 produce generalized responses for a public audience, rather than exact and technical responses,  
200 unless specified otherwise (5, 16, 20).

201 Our formal error analysis highlighted additional patterns of poor ICD billing code generation  
202 performance. Complex and lengthy descriptions, as well as certain condition categories, proved  
203 challenging for the model to parse. It should be noted that nonbillable codes, shorter codes, and  
204 shorter descriptions usually reflect more generalized conditions. These patterns of poor  
205 performance must be addressed before GPT-based LLMs can be used to interface with ICD  
206 terminologies for billing purposes. Additional prompt engineering to encourage exact matching  
207 of ICD billing codes may be able to improve accuracy (21).

208 The inability of LLMs to parse ICD codes consistently has parallels with other previously  
209 identified general technical limitations of these models. Due to the nature of how LLMs  
210 algorithmically break up text into “tokens”, some have difficulty with tasks that require an  
211 understanding the structure of information contained within each token (22-24). Similarly, LLMs

212 are unable to consistently reverse the spelling of words, perform math problems, or complete  
213 complex coding tasks without the use of external software or model fine-tuning. When we  
214 evaluated ICD codes in OpenAI's tokenization tool, we observed that tokenization behavior was  
215 not consistent with the underlying ICD terminology structures, leaving the models without key  
216 digit-level hierarchical information during the model training process (25, 26). To address this  
217 technical issue, an GPT-3.5 or GPT-4 ICD tool will need to be linked to additional software  
218 layers that can facilitate an interface between LLMs and ICD billing codes (11, 27-32).

219 This study's results reaffirm the tendency for LLMs to produce realistic-appearing, but incorrect  
220 information—a significant obstacle to their use in the healthcare setting. While this study showed  
221 that such “hallucination” behavior did not significantly impact broad semantic understanding, it  
222 often compromised precision to a degree that is not acceptable for medical coding purposes (33).

223 As LLM technology is increasingly embedded into healthcare, understanding these limitations is  
224 crucial. LLM accuracy while interfacing with ICD terminologies must be improved before they  
225 can help realize their potential in streamlining administrative tasks, supporting clinicians, and  
226 improving patient care. Additional LLM performance enhancement strategies such as prompt  
227 engineering, database linkage, model fine-tuning, and others can potentially address this in the  
228 near future.

#### 229 *Limitations:*

230 Our study has a few limitations. First, we used a sample of conditions and procedures for testing,  
231 which may not represent the codes most frequently used in real-world coding scenarios. Second,  
232 we did not evaluate advanced LLM performance enhancement strategies such as prompt  
233 engineering, database-linkage, or model fine-tuning. Lastly, we did not evaluate the models'

234 performance in the context of real-world clinical narratives, which often involves complex and  
235 ambiguous language that does not always have a clear ground-truth mapping.

236 *Conclusion:*

237 Our evaluation of baseline GPT-3.5 and GPT-4's proficiency in generating ICD billing codes  
238 from the ICD-9-CM, ICD-10-CM and ICD-10-PCS coding systems reveals an inadequate level  
239 of performance for downstream use cases. While the models do exhibit an understanding of the  
240 conditions, as demonstrated by the generation of codes semantically related to the actual ones,  
241 they tend to exhibit a propensity for data hallucination, producing codes that are not entirely  
242 accurate. This suggests a need for rigorous validation and refinement processes prior to their  
243 implementation in healthcare contexts, particularly in tasks such as ICD billing coding which  
244 require precision. Further, the integration with external tools may be necessary to enhance their  
245 performance.

246

247 **Data Availability Statement:** The data that support the findings of this study are available from  
248 the corresponding author, AS, upon reasonable request.

249

250 **Funding/Support:** This research did not receive any specific grant from funding agencies in the  
251 public, commercial or not-for-profit sectors.

252

253 **Conflict of Interest Statement:** BSG: no relevant conflicts of interest but is an employee of  
254 Character Biosciences. GN: Consultancy agreements with AstraZeneca, BioVie, GLG  
255 Consulting, Pensieve Health, Reata, Renalytix, Siemens Healthineers, and Variant Bio; research  
256 funding from Goldfinch Bio and Renalytix; honoraria from AstraZeneca, BioVie, Lexicon,



257 Daiichi Sankyo, Meanrini Health and Reata; patents or royalties with Renalytix; owns equity and  
258 stock options in Pensieve Health and Renalytix as a scientific cofounder; owns equity in Verici  
259 Dx; has received financial compensation as a scientific board member and advisor to Renalytix;  
260 serves on the advisory board of Neurona Health; and serves in an advisory or leadership role for  
261 Pensieve Health and Renalytix. All other authors: no conflicts of interest to declare.

262

263 **Ethics Approval:** Since all data and responses were publicly available, approval from the  
264 institutional review board was not sought.

1

## References

1. Organization WH. History of the development of the ICD [June 19, 2023]. Available from: <https://cdn.who.int/media/docs/default-source/classification/icd/historyoficd.pdf>.
2. Organization WH. Importance of ICD. Available from: <https://www.who.int/standards/classifications/frequently-asked-questions/importance-of-icd>.
3. Wood PH. Applications of the International Classification of Diseases. *World Health Stat Q*. 1990;43(4):263-8. PubMed PMID: 2293495.
4. Bubeck S, Chandrasekaran V, Eldan R, Gehrke J, Horvitz E, Kamar E, et al. Sparks of Artificial General Intelligence: Early experiments with GPT-4 2023 March 01, 2023:[arXiv:2303.12712 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2023arXiv230312712B>.
5. OpenAI. GPT-4 Technical Report 2023 March 01, 2023:[arXiv:2303.08774 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2023arXiv230308774O>.
6. Zhao WX, Zhou K, Li J, Tang T, Wang X, Hou Y, et al. A Survey of Large Language Models 2023 March 01, 2023:[arXiv:2303.18223 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2023arXiv230318223Z>.
7. Nori H, King N, McKinney SM, Carignan D, Horvitz E. Capabilities of GPT-4 on Medical Challenge Problems 2023 March 01, 2023:[arXiv:2303.13375 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2023arXiv230313375N>.
8. 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):1233-9. doi: 10.1056/NEJMs2214184. PubMed PMID: 36988602.
9. Patel SB, Lam K. ChatGPT: the future of discharge summaries? *Lancet Digit Health*. 2023;5(3):e107-e8. Epub 20230206. doi: 10.1016/S2589-7500(23)00021-3. PubMed PMID: 36754724.
10. Moor M, Banerjee O, Abad ZSH, Krumholz HM, Leskovec J, Topol EJ, et al. Foundation models for generalist medical artificial intelligence. *Nature*. 2023;616(7956):259-65. Epub 20230412. doi: 10.1038/s41586-023-05881-4. PubMed PMID: 37045921.
11. Manakul P, Liusie A, Gales MJF. SelfCheckGPT: Zero-Resource Black-Box Hallucination Detection for Generative Large Language Models 2023 March 01, 2023:[arXiv:2303.08896 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2023arXiv230308896M>.
12. McKenna N, Li T, Cheng L, Hosseini MJ, Johnson M, Steedman M. Sources of Hallucination by Large Language Models on Inference Tasks 2023 May 01, 2023:[arXiv:2305.14552 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2023arXiv230514552M>.
13. Li J, Cheng X, Zhao WX, Nie J-Y, Wen J-R. HaluEval: A Large-Scale Hallucination Evaluation Benchmark for Large Language Models 2023 May 01, 2023:[arXiv:2305.11747 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2023arXiv230511747L>.

14. Ji Z, Lee N, Frieske R, Yu T, Su D, Xu Y, et al. Survey of Hallucination in Natural Language Generation 2022 February 01, 2022:[arXiv:2202.03629 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2022arXiv220203629J>.
15. Services CfMM. ICD Code Lists 2022 [updated September 28, 2022]. Available from: <https://www.cms.gov/medicare/coordination-benefits-recovery-overview/icd-code-lists>.
16. Ouyang L, Wu J, Jiang X, Almeida D, Wainwright CL, Mishkin P, et al. Training language models to follow instructions with human feedback 2022 March 01, 2022:[arXiv:2203.02155 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2022arXiv220302155O>.
17. Bodenreider O. The Unified Medical Language System (UMLS): integrating biomedical terminology. *Nucleic Acids Res.* 2004;32(Database issue):D267-70. doi: 10.1093/nar/gkh061. PubMed PMID: 14681409; PubMed Central PMCID: PMC308795.
18. Topaz M, Shafran-Topaz L, Bowles KH. ICD-9 to ICD-10: evolution, revolution, and current debates in the United States. *Perspect Health Inf Manag.* 2013;10(Spring):1d. Epub 20130401. PubMed PMID: 23805064; PubMed Central PMCID: PMC3692324.
19. Razeghi Y, Logan RL, IV, Gardner M, Singh S. Impact of Pretraining Term Frequencies on Few-Shot Reasoning 2022 February 01, 2022:[arXiv:2202.07206 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2022arXiv220207206R>.
20. Kenton Z, Everitt T, Weidinger L, Gabriel I, Mikulik V, Irving G. Alignment of Language Agents 2021 March 01, 2021:[arXiv:2103.14659 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2021arXiv210314659K>.
21. White J, Fu Q, Hays S, Sandborn M, Olea C, Gilbert H, et al. A Prompt Pattern Catalog to Enhance Prompt Engineering with ChatGPT 2023 February 01, 2023:[arXiv:2302.11382 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2023arXiv230211382W>.
22. Yuan Z, Yuan H, Tan C, Wang W, Huang S. How well do Large Language Models perform in Arithmetic tasks? 2023 March 01, 2023:[arXiv:2304.02015 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2023arXiv230402015Y>.
23. Kim J, Hong G, Kim K-m, Kang J, Myaeng S-H, editors. Have You Seen That Number? Investigating Extrapolation in Question Answering Models 2021 November; Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.
24. Nogueira R, Jiang Z, Lin J. Investigating the Limitations of Transformers with Simple Arithmetic Tasks 2021 February 01, 2021:[arXiv:2102.13019 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2021arXiv210213019N>.
25. OpenAI. Tokenizer. Available from: <https://platform.openai.com/tokenizer>.
26. Brown TB, Mann B, Ryder N, Subbiah M, Kaplan J, Dhariwal P, et al. Language Models are Few-Shot Learners 2020 May 01, 2020:[arXiv:2005.14165 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2020arXiv200514165B>.
27. Peng B, Galley M, He P, Cheng H, Xie Y, Hu Y, et al. Check Your Facts and Try Again: Improving Large Language Models with External Knowledge and Automated Feedback 2023 February 01, 2023:[arXiv:2302.12813 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2023arXiv230212813P>.

28. Huang J, Gu SS, Hou L, Wu Y, Wang X, Yu H, et al. Large Language Models Can Self-Improve2022 October 01, 2022:[arXiv:2210.11610 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2022arXiv221011610H>.
29. Azaria A, Mitchell T. The Internal State of an LLM Knows When its Lying2023 April 01, 2023:[arXiv:2304.13734 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2023arXiv230413734A>.
30. Lin S, Hilton J, Evans O. Teaching Models to Express Their Uncertainty in Words2022 May 01, 2022:[arXiv:2205.14334 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2022arXiv220514334L>.
31. Lu P, Peng B, Cheng H, Galley M, Chang K-W, Nian Wu Y, et al. Chameleon: Plug-and-Play Compositional Reasoning with Large Language Models2023 April 01, 2023:[arXiv:2304.09842 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2023arXiv230409842L>.
32. Dai H, Liu Z, Liao W, Huang X, Cao Y, Wu Z, et al. AugGPT: Leveraging ChatGPT for Text Data Augmentation2023 February 01, 2023:[arXiv:2302.13007 p.]. Available from: <https://ui.adsabs.harvard.edu/abs/2023arXiv230213007D>.
33. The Lancet Digital H. ChatGPT: friend or foe? Lancet Digit Health. 2023;5(3):e102. Epub 20230206. doi: 10.1016/S2589-7500(23)00023-7. PubMed PMID: 36754723.