

From Documents to Spans: Code-Centric Learning for LLM-based ICD Coding

Anonymous ACL submission

Abstract

ICD coding is a critical yet challenging task in healthcare. Recently, LLM-based methods demonstrate stronger generalization than discriminative methods in ICD coding. However, fine-tuning LLMs for ICD coding faces three major challenges. First, existing public ICD coding datasets provide limited coverage of the ICD code space, restricting a model’s ability to generalize to unseen codes. Second, naive fine-tuning diminishes the interpretability of LLMs, as few public datasets contain explicit supporting evidence for assigned codes. Third, ICD coding typically involves long clinical documents, making fine-tuning LLMs computationally expensive. To address these issues, we propose Code-Centric Learning, a training framework that shifts supervision from full clinical documents to scalable, short evidence spans. The key idea of this framework is that span-level classification improves LLMs’ ability to perform document-level ICD coding. Our proposed framework consists of a mixed training strategy and code-centric data expansion, which effectively improve coding performance on ICD codes out of domain while preserving interpretability. With only the open-source Llama-3.1-8B model, our method outperforms or matches strong discriminative baselines and GPT-4.1-based generative methods, demonstrating its effectiveness and potential for fully automated ICD coding. Code is available at <https://anonymous.4open.science/r/CCL-ICD>.

1 Introduction

ICD codes play a foundational role in medical insurance reimbursement and health data analysis. ICD coding, the process of assigning ICD codes to each patient encounter, is a critically important task in modern healthcare systems. However, ICD coding remains an extremely challenging task. Even expert human coders frequently make errors (Burns et al., 2012; Horsky et al., 2018; Gan et al., 2025).

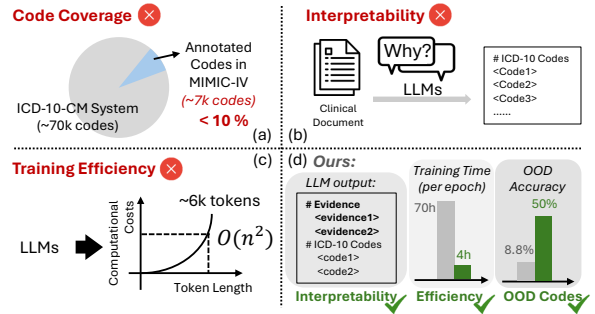


Figure 1: (a) Existing public datasets covers limited ICD codes; (b) Direct fine-tuning lacks interpretability; (c) Training cost of LLM on long documents is expensive; (d) Our solution hugely improves OOD code accuracy via efficient training, and preserves interpretability.

First, multiple codes often apply to a single encounter, leading to pervasive undercoding. In addition, the large and fine-grained ICD vocabulary increases the likelihood of code confusion and misassignment. These challenges motivate extensive research on automated ICD coding.

Early ICD coding approaches adopt discriminative models with label attention mechanisms (Mullenbach et al., 2018; Huang et al., 2022; Edin et al., 2024). Recently, LLM-based methods gain attention due to their stronger generalization capabilities (Motzfeldt et al., 2025; Yuan et al., 2025). Among them, training-free approaches (Li et al., 2024; Motzfeldt et al., 2025) rely on carefully designed prompts or workflows, and place strong demands on the capability of the backbone LLM. However, such high-performance LLMs are often closed-source and require online access, making them incompatible with the strict data privacy and deployment constraints in clinical settings. In contrast, fine-tuning methods (Yuan et al., 2025) enable smaller, locally deployable LLMs to acquire ICD coding capabilities, making them more practical for real-world healthcare applications. Given clinical notes and a simple task prompt, the LLM is

069 optimized to output ICD codes along with their textual descriptions and outperforms discriminative models on out-of-distribution data.

070
071
072 However, this fine-tuning paradigm entails several issues: **limited code coverage, poor interpretability and low training efficiency** (shown in Figure 1). First, existing ICD coding datasets cover only a small subset of ICD codes. For example, even MIMIC-IV (Johnson et al., 2023) includes only about 10% of the 70K ICD-10-CM codes, which severely restricts generalization to unseen codes. Second, most ICD coding datasets provide ICD codes without supporting evidence; as a result, fine-tuning LLMs on such data encourages direct code prediction without explicit evidence grounding. This not only diminishes interpretability, but also precludes human review and correction of the underlying evidence, thereby limiting effective human-AI collaboration. Third, fine-tuning LLMs on long clinical documents is computationally expensive, considering the quadratic complexity with respect to input length. Collectively, these challenges significantly limit the applicability of fine-tuning LLMs for ICD coding.

093 To address the above issues, we propose a novel training framework, Code-Centric Learning (CCL). Unlike traditional paradigms that operate on entire clinical notes, our method centers on code-specific evidence spans. Intuitively, ICD coding can be decomposed into two sub-tasks: locating evidence and assigning codes. If an LLM is trained to assign the correct ICD code to an evidence span, it implicitly learns to recognize such evidence in a long clinical document. Consequently, **strengthening the evidence-level code assignment capability also enhances the model’s ability to locate relevant evidence and assign codes in a full note.**

106 Based on this intuition, CCL consists of two key strategies. First, we adopt a **mixed training** strategy that leverages a limited number of samples with annotated evidence spans together with a large collection of code-related evidence spans. Compared to direct fine-tuning on full clinical notes, this strategy provides explicit interpretability and reduces computational cost by focusing on short evidence spans instead of long clinical documents. Second, we propose a **code-centric data expansion** strategy. We extract code-specific evidence spans from public datasets based on annotated codes, and supplement them using official ICD coding resources. For codes unseen in both official knowledge bases and public datasets, we retrieve

121 the closest codes and evidence to synthesize plausible evidence spans, ensuring full ICD code coverage. Our proposed training framework enables an 8B-scale LLM to outperform both discriminative models and strong GPT-4-based approaches on both in-domain and out-of-domain datasets, while additionally providing explicit interpretability and intervention capability.

128 Our main contributions are summarized as:

- 129
130 • We propose a fine-tuning framework for LLM-based ICD coding that simultaneously addresses three key limitations: (i) limited code coverage in public datasets, (ii) the lack of interpretability in standard fine-tuning, and (iii) the low training efficiency caused by long clinical documents. 131 132 133 134 135 136
- 137 • We introduce a novel training framework consisting of a mixed training strategy and a code-centric data expansion strategy, motivated by the central insight that **span-level classification yields transferable gains for document-level evidence extraction.** 138 139 140 141 142
- 143 • Our proposed framework enables an 8B-scale LLM to achieve competitive performance against both discriminative models and strong GPT-4-based approaches on both in-domain and out-of-domain datasets, while attaining state-of-the-art results on several key metrics. In addition, it provides explicit interpretability and supports human-AI collaboration. 144 145 146 147 148 149 150

151 2 Related Work

152 2.1 Discriminative methods

153 **Label attention.** Discriminative models have long dominated the ICD coding task, whose strong performance can be attributed to label attention mechanisms (Mullenbach et al., 2018), learning an independent query vector for each ICD code. The widely adopted state-of-the-art method is PLM-ICD (Huang et al., 2022) and its variant PLM-CA (Edin et al., 2024). 154 155 156 157 158 159 160

161 **Knowledge injection.** Building on label attention, several studies explore incorporating external ICD-related knowledge. DKEC (Ge et al., 2024) apply a graph network to encode knowledge from multiple sources. AKIL and MRR (Wang et al., 2024b,a) introduce auxiliary coding systems such as DRG, CPT, and procedure codes. Correlation (Luo et al., 2024) models relationships among ICD 162 163 164 165 166 167 168

codes. MSMN (Yuan et al., 2022) and MSAM (Gomes et al., 2024) utilize synonyms to learn code representations. GKI-ICD (Zhang et al., 2025) injects code descriptions, synonyms, and hierarchical relations by synthesizing guidelines. ACE-ICD (Le et al., 2025) improves medical terminology understanding via acronym expansion.

Interpretability. Traditional discriminative models operate as black-box systems, with limited interpretability. To address this challenge, MDACE (Cheng et al., 2023) re-annotated a subset of MIMIC-III (Johnson et al., 2016), in which expert coders performed multi-round cross-validation to relabel ICD codes and provided gold-standard evidence spans for each assignment. Building on MDACE, Edin et al. (2024) apply AttInGrad to map model predictions to evidence spans in the original text. AutoCodeDL (Wu et al., 2024) incorporates dictionary learning to decode dense embeddings into medical concepts. While discriminative models require additional mechanisms to produce evidence, generative models naturally support evidence-based coding.

2.2 Generative Methods

Training-free methods. Early works explored the use of off-the-shelf LLMs for ICD coding, designing prompts and workflows. Boyle et al. (2023) prompts the LLM to predict ICD codes in a hierarchical manner, from chapters and sections to specific codes. MAC (Li et al., 2024) prompts the LLM to perform as different roles, such as coder and physician, reducing coding errors via cross-role verification. MedCodER (Baksi et al., 2025) combines evidence extraction, candidate retrieval, and re-ranking, to overcome the large label space. Similarly, CLH (Motzfeldt et al., 2025) extracts evidence, retrieves candidate codes via the Alphabetic Index and Tabular List, and verifies them.

Fine-tuning methods. Recently, (Yuan et al., 2025) demonstrates that fine-tuning LLMs is more suitable for this task than training-free paradigms and proposes a verification module to fix mistakes. Nesterov et al. (2025) also validates this conclusion on their proposed Russian ICD coding benchmark. However, these methods typically perform supervised fine-tuning (SFT) directly on ICD coding datasets, which leads to low training efficiency on long documents, covers limited ICD code space, and lacks interpretability. Therefore, developing a more suitable fine-tuning paradigm for LLM-based ICD coding is particularly important.

3 Methodology

3.1 Overview

Fine-tuning LLMs for ICD coding is commonly performed at the document level. Given a clinical document x and its associated ICD code set C (followed by code description), conventional document-level training can be cast as supervised fine-tuning (SFT):

$$\min_{\theta} \mathcal{L}_{\text{SFT}}(f_{\theta}(x), C), \quad (1)$$

where the LLM $f_{\theta}(\cdot)$ maps a long clinical document to a set of ICD codes under standard next-token prediction. However, this paradigm is inherently limited by three critical challenges: low training efficiency, high annotation noise, and poor interpretability (See Section 1).

To overcome these challenges, we propose a code-centric learning framework. We perform mixed SFT on two types of instances: (i) a full clinical document, annotated with evidence and ICD codes; (ii) an evidence span and a single ICD code, which can be formulated as:

$$\min_{\theta} \left[\mathcal{L}_{\text{SFT}}(f_{\theta}(x), (E, C)) + \mathcal{L}_{\text{SFT}}(f_{\theta}(e), c) \right], \quad (2)$$

where E denotes evidence spans supporting annotated ICD codes C , e denotes an evidence span with corresponding ICD code c . **The former forces the LLM to extract evidence before assigning codes, enhancing interpretability, while the latter consists of short evidence spans that are easy to acquire, which naturally address issues of limited code coverage and training efficiency.**

Defining (e, c) as an evidence-code pair,

$$(e, c) \sim \mathcal{D}, \quad (3)$$

where \mathcal{D} is a multi-source knowledge base. To ensure full ICD code coverage, we construct \mathcal{D} by integrating three tiers of evidence-code pairs: *gold*, *silver*, and *synthetic*, as:

$$\mathcal{D} = \mathcal{D}_{\text{gold}} \cup \mathcal{D}_{\text{silver}} \cup \mathcal{D}_{\text{syn}}, \quad (4)$$

in which $\mathcal{D}_{\text{gold}}$ consists of scarce, authoritative evidence-code pairs from official ICD resources. $\mathcal{D}_{\text{silver}}$ is mined from public datasets, yielding a larger collection of evidence-code pairs with broader coverage. For the remaining unseen codes, we construct \mathcal{D}_{syn} by synthesizing evidence via

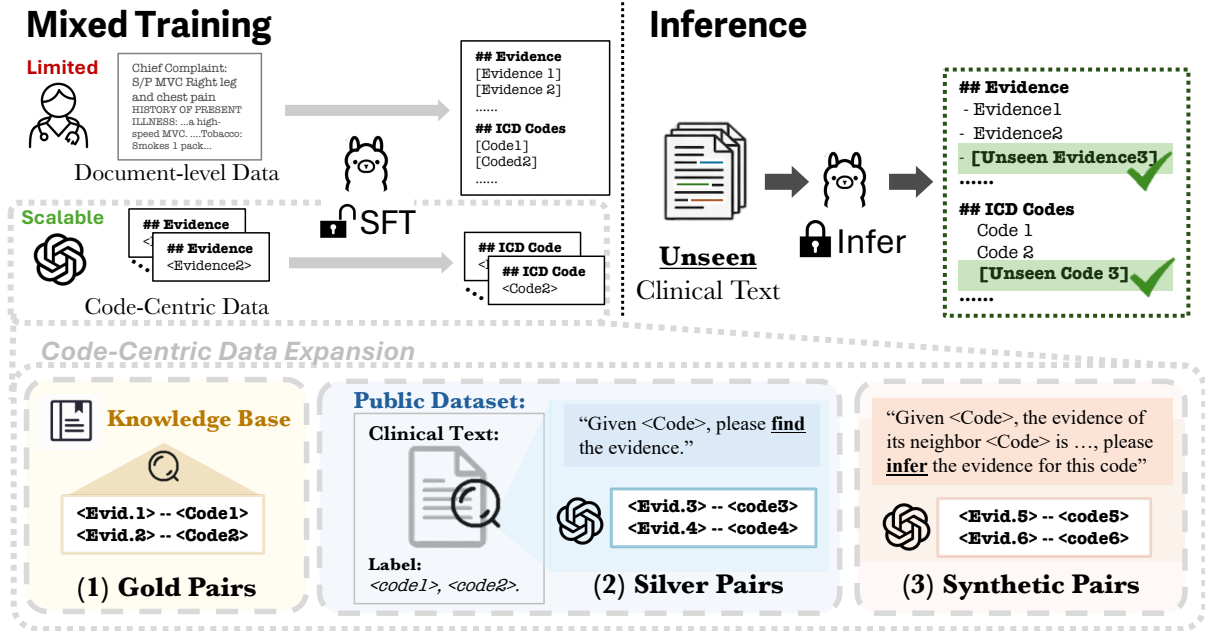


Figure 2: Overview of code-centric learning framework. Mixed Training unifies document-level evidence-based ICD coding and span-level code knowledge injection within a single framework. Code-centric Data Expansion leverages large language models to extract and infer evidence spans for each code from diverse knowledge sources, addressing codes that lack evidence-annotated documents.

LLM, thereby completing coverage over the entire ICD code set. Below, we will describe the details of the mixed training strategy for Eq. 2, and the code-centric data expansion strategy for Eq. 4.

3.2 Mixed Training

Mixed training is designed to bridge the gap between scarce, high-quality document-level ICD coding data and large-scale but weakly supervised code-centric data. By jointly training on these two complementary supervision signals, the model learns to follow human-aligned evidence-based coding workflows while acquiring broad ICD knowledge that cannot be covered by limited expert annotations alone.

Document-level evidence-based ICD coding data refers to medical documents annotated with both ICD codes and supporting evidence. This type of data is very hard to obtain, and therefore extremely scarce and valuable. To our knowledge, MDACE (Cheng et al., 2023) is the only available public dataset that contains such kind of data.

For each clinical document, we first extract the human-annotated evidence spans, preserving their original order in the document. We then order the ICD codes accordingly, and augment them with their textual descriptions from the ICD-10 Tabular List. Finally, we convert text, evidence

and codes into instruction-tuning samples using a unified prompt template (Appendix B). Note that evidence extraction and code assignment are performed jointly within a single generation process, rather than through staged or multi-step pipelines.

Despite MDACE’s small scale, LLMs can learn robust evidence-based ICD coding behavior, as demonstrated in Section 5.2.

Code-centric learning data refers to evidence–code pairs, where the model takes an evidence span as input, instead of a full clinical document, and predicts the corresponding ICD code. The advantage of code-centric learning data is obvious: since each sample contains only a short evidence span and a single ICD code, it enables efficient training, and is scalable and substantially easier to denoise than document-level labels.

Such evidence–code pairs can be obtained from diverse sources. They may originate from high-quality, human-curated resources (e.g., official coding references) or be automatically extracted by LLMs from public ICD coding dataset. Section 3.3 describes how we systematically expand these pairs to increase code coverage.

Training and inference. We fine-tune the LLM using a mixture of document-level, evidence-based ICD coding data and large-scale code-centric learning data under a standard autoregressive super-

vised fine-tuning objective. The evidence-based data teaches the model to follow the human coding workflow and produce structured, interpretable outputs, while the code-centric data injects scalable code-specific knowledge and enables generalization beyond the limited code set covered by dense annotations.

$$\min_{\theta} \mathcal{L}_{\text{SFT}}(\theta) = \sum_{i=1}^N \sum_{t=1}^{T_i} -\log p_{\theta}(y_t^{(i)} | x^{(i)}, y_{<t}^{(i)}). \quad (5)$$

At inference time, we apply the same prompt schema used for evidence-based training (Appendix B). This unified formulation encourages the model to first identify relevant evidence spans and then assign ICD codes, while implicitly leveraging the code-level knowledge acquired during code-centric learning. As a result, the model produces ICD predictions that are both interpretable and accurate, without requiring evidence annotations at test time.

3.3 Code-centric Data Expansion

Under the code-centric learning paradigm, we organize training supervision around individual ICD codes rather than full clinical documents. Accordingly, we construct a multi-tier evidence-code knowledge base composed of gold pairs from human-curated ICD resources, silver pairs mined from noisy labeled datasets, and synthetic pairs inferred by LLMs to ensure full code coverage.

Gold pairs from human knowledge bases. In clinical practice, human coders routinely consult the Alphabetic Index and the Tabular List when assigning ICD codes (See Appendix A.3). These resources naturally form high-quality code-centric learning corpus, but had been overlooked by previous works. CLH (Motzfeldt et al., 2025) first incorporated these resources, treating them as an external corpus for retrieval-augmented generation. In our framework, we treat Alphabetic Index terms paired with their default ICD codes as gold evidence-code pairs.

Silver pairs from noisy labeled datasets. Most large-scale ICD datasets provide only document-level code labels without explicit evidence annotations. To align such data with evidence-based ICD coding, we construct silver evidence-code pairs via a two-stage LLM pipeline: document-level evidence extraction followed by code-level evidence consolidation.

In the first stage, given a clinical note and one of

its assigned ICD codes, the LLM extracts a textual span that plausibly supports the code. We aggregate extracted spans across the dataset for each ICD code, retain unique evidence phrases, and record their frequencies.

$$\mathcal{E}_c = \{e | e = f_{\text{LLM}}(x, c), x \in \mathcal{X}_c\}, \quad (6)$$

where \mathcal{X}_c denotes the set of clinical documents labeled with code c , $f_{\text{LLM}}(x, c)$ extracts supporting evidence spans from document x for code c , yielding a large evidence set \mathcal{E}_c . g_{LLM} then summarizes \mathcal{E}_c into a small set of representative (typical) evidence expressions $\tilde{\mathcal{E}}_c$.

In the second stage, given an ICD code and its frequency-ranked evidence candidates, the LLM infers a small set of representative evidence expressions for that code, forming the silver dataset $\mathcal{D}_{\text{silver}}$.

$$\tilde{\mathcal{E}}_c = f_{\text{LLM}}(c, \mathcal{E}_c), \quad (7)$$

where f_{LLM} denotes the LLM, c denotes the target ICD code, and \mathcal{E} denotes the evidence candidates.

Synthetic pairs inferred by LLMs. Despite combining gold and silver pairs, many ICD codes remain uncovered in public datasets. To achieve full code coverage, we synthesize evidence-code pairs using an LLM guided by ICD knowledge.

For each uncovered target ICD code, we retrieve its nearest neighbor code in the ICD-10-CM hierarchy and related information of this nearest code. Conditioned on this information, the LLM generates concise clinical evidence expressions that plausibly support the target code, forming the synthetic dataset \mathcal{D}_{syn} :

$$e = f_{\text{LLM}}(c, c^*, \mathcal{K}(c^*)), \quad (8)$$

where c denotes the target ICD code, c^* denotes its nearest ICD code, and $\mathcal{K}(c^*)$ represents the associated ICD knowledge of c^* , i.e. potential evidence from gold pairs and silver pairs. These synthetic pairs complement gold and silver data, resulting in a code-centric knowledge base with complete ICD coverage.

Finally, we mix gold, silver and synthetic evidence-code pairs, and convert them into the same prompt format as evidence-based ICD coding, to form a large, high-quality instruction-tuning dataset.

4 Experiments

4.1 Dataset

MIMIC-IV (Johnson et al., 2023) is currently the largest publicly available dataset annotated with ICD-10 codes. However, it contains only ICU clinical notes, while the ICD codes span the entire patient encounter. As a result, many codes are not supported by the available text (Cheng et al., 2023; Edin et al., 2023; Yuan et al., 2025), making MIMIC-IV unsuitable as a reliable benchmark. We therefore treat MIMIC as a large but noise-prone training dataset, and rely on high-quality external benchmarks to measure true ICD coding performance. Since some of the baselines are trained on its training set, our method uses the same training set to extract code-centric data, following Edin et al. (2024)’s split, to ensure fair comparison.

MDACE (Cheng et al., 2023) is an expert-annotated subset of MIMIC-III (Johnson et al., 2016), containing gold-standard evidence span annotations. We follow its official data split, using the training set for document-level evidence-based ICD coding, and using its test set for evaluation.

ACI-Bench (Yim et al., 2023) is a synthetic dataset of clinical notes, based on which Yuan et al. (2025) constructs a new double expert-annotated ICD-10-CM coding benchmark.

In summary, our training primarily relies on MIMIC-IV, supplemented with MDACE. Evaluation is conducted on MDACE and ACI-Bench, with the latter providing a more out-of-distribution benchmark to assess generalization.

4.2 Evaluation

For discriminative models, we use the validated-optimal threshold. For generative models, we extract the alphanumeric code component (letter followed by digits) from LLM’s text output. We filter out codes not present in the test-set label space. Since generative models do not provide probability scores, AUC cannot be computed. We therefore report micro-F1, macro-F1, Precision, and Recall as our primary evaluation metrics.

4.3 Baselines

Discriminative Methods. PLM-ICD (Huang et al., 2022) is a BERT-based model combined with label attention. PLM-CA (Edin et al., 2024) is its variant, replacing label attention with cross attention. GKI-ICD (Zhang et al., 2025) is a knowledge injected version of PLM-CA. We finetune these models on

MIMIC-IV-ICD10 dataset, and find the optimal threshold.

Generative methods. For closed-source LLMs, we include Chain-of-Thought (CoT) (Wei et al., 2022) and its self-consistency variant (CoT-SC) (Wang et al., 2022), as well as task-specific methods, MAC (Li et al., 2024) and CLH (Motzfeldt et al., 2025). We additionally include CoT and SFT of the same backbone LLM. For SFT, we adopt the optimal configuration from Yuan et al. (2025), using code and description as output format.

5 Results

5.1 Comparison with SOTA models

Table 1 shows that our method (CCL) achieves strong performance using only Llama3.1-8B. **Compared with methods of similar scale, CCL substantially outperforms** prompt-based CoT and improves over standard SFT across both in-domain (MDACE) and out-of-domain (ACI-Bench) benchmarks. Despite its small model size, CCL achieves performance **comparable to GPT-4.1-based methods, and even surpasses** them on several key metrics (e.g., Micro-F1 on MDACE and competitive Micro-/Macro-F1 on ACI-Bench), while remaining competitive on the others.

Importantly, CCL also provides explicit evidence for its predicted ICD codes, enabling human coders to review and revise the evidence, and our experiments in 2 show that such revisions lead to measurable performance improvements. Overall, CCL strikes an effective balance between performance, efficiency, and interpretability.

5.2 Ablation Study

We conduct ablation studies on the MDACE dataset using Llama-3.1-8B as the backbone model.

Evidence-based fine-tuning is more effective than code-only fine-tuning. First, even with a small amount of fine-tuning data, i.e. MDACE training set, LLMs achieve clear performance gains over CoT without fine-tuning, demonstrating that fine-tuning is necessary and effective for ICD coding. Second, on the same amount of training data, incorporating human-annotated evidence spans into ground truth (i.e. train LLMs to extract evidence before assigning codes) leads to larger improvements than code-only fine-tuning. This behavior contrasts with prior findings on discriminative models, where evidence-span supervision does not necessarily improve coding accuracy (Cheng et al., 2023).

Method	Backbone	MDACE (In Domain)				ACI-Bench (Out of Domain)			
		Micro-F1	Macro-F1	Recall	Precision	Micro-F1	Macro-F1	Recall	Precision
PLM-ICD (Huang et al., 2022)	RoBERTa (120M)	50.6	26.5	66.6	40.8	39.3	18.3	58.1	39.6
PLM-CA (Edin et al., 2024)	RoBERTa (120M)	50.0	25.8	65.3	40.4	25.1	11.1	64.8	15.5
GKI-ICD (Zhang et al., 2025)	RoBERTa (120M)	50.4	26.3	64.3	41.0	47.8	25.8	62.0	39.0
CoT (Wei et al., 2022)	GPT-4.1	57.0	35.8	59.7	54.4	61.9	48.1	67.6	57.0
CoT-SC (Wang et al., 2022)	GPT-4.1	57.1	36.2	56.1	58.2	61.4	47.4	67.6	56.2
MAC (Li et al., 2024)	GPT-4.1	48.1	31.0	61.2	39.7	51.5	45.6	69.6	40.9
CLH (Motzfeldt et al., 2025)	Qwen3-235BA22B	56.3	43.7	55.2	57.6	67.1	47.7	66.1	68.1
CoT (Wei et al., 2022)	Llama3.1-8B	26.5	9.1	20.2	38.9	43.4	12.4	36.4	53.6
SFT (Yuan et al., 2025)	Llama3.1-8B	57.4	34.5	60.1	55.0	63.4	38.9	57.9	70.1
CCL (Ours)	Llama3.1-8B	59.3	35.2	56.4	62.5	65.8	47.1	64.8	67.0

Table 1: Comparison with previous SOTA methods on in domain and out of domain benchmarks.

Methods	Micro-F1	Macro-F1	Recall	Precision
w/o fine-tuning	26.5	9.1	20.2	38.9
Code-only	42.3	17.0	37.8	48.2
Evidence & Code	49.7	21.4	47.0	52.8
CCL (Gold-only)	53.4	25.2	57.2	50.1
CCL (Ours)	59.3	35.2	56.4	62.5
w/ human evidence	78.0	54.8	73.5	83.0

Table 2: Ablation Study on MDACE Dataset. Code-only and Code & Evidence use MDACE training set, whereas CCL additionally incorporates evidence spans. w/ Human Evidence indicates that, during testing, human-annotated evidence is added to the input, allowing the LLM to continue assigning codes.

Code-centric learning effectively improves coding performance. We first conduct a simplified experiment, using only gold pairs from Alphabetic Index. As shown in Table 2, CCL on gold pairs already yields a clear improvement in coding performance. Furthermore, we extract or synthesize code-specific evidence spans, i.e. silver and synthetic pairs, to achieve full code coverage, which leads to further gains. These results highlight the effectiveness of code-centric learning and code-centric data expansion.

Human-AI Collaboration. Our paradigm explicitly extracts evidence before code assignment, enabling human-in-the-loop ICD coding by allowing clinicians to review and revise LLM-generated evidence. As shown in Table 5.2, replacing model-generated evidence with human-annotated evidence yields substantial performance improvements, highlighting the interpretability, controllability, and practical applicability of our approach.

5.3 Generality and Scalability

We conduct extensive experiments across LLMs with different architectures and parameter scales. Due to hardware constraints, we are limited to fine-tuning models with up to 8B parameters.

Methods	Micro-F1	Macro-F1	Recall	Precision
Llama3.1-8B	59.3	35.2	56.4	62.5
Llama-3.2-1B	40.8	16.0	35.2	48.6
Llama-3.2-3B	53.2	26.0	47.4	60.6
Qwen2.5-0.5B	31.8	10.6	24.9	44.2
Qwen2.5-1.5B	50.6	24.1	45.0	57.8
Qwen2.5-7B	54.2	27.4	48.6	61.2

Table 3: Performance of different LLM under CCL Framework on MDACE Dataset.

Generality across architectures. As shown in Table 3, the proposed framework performs consistently across both Llama and Qwen backbones, demonstrating robust cross-architecture generalization. Across comparable model scales, Llama models consistently outperform Qwen models, which may be due to Llama being primarily pretrained on English data.

Scalability with model size. Within each model family, performance improves steadily as model size increases. Both Llama-3.2 and Qwen2.5 models exhibit clear gains in Micro-F1 and Macro-F1 when scaling from smaller to larger variants. The consistent scaling behavior across architectures further supports the scalability of the approach.

5.4 Analysis of Code-Centric Learning

Code-centric Learning enables the model to acquire knowledge of codes beyond the coverage of the document-level data. As illustrated in Figure 3 (a), we categorize the codes in the test set into two groups: (1) codes covered by the document-level data, for which the model explicitly learns to extract supporting evidence and assign ICD codes, and (2) codes learned exclusively through code-centric learning, for which the model has only been trained to assign codes without explicit evidence extraction. The results in Figure 3 (c) demonstrate that, under the mixed training strategy, the code

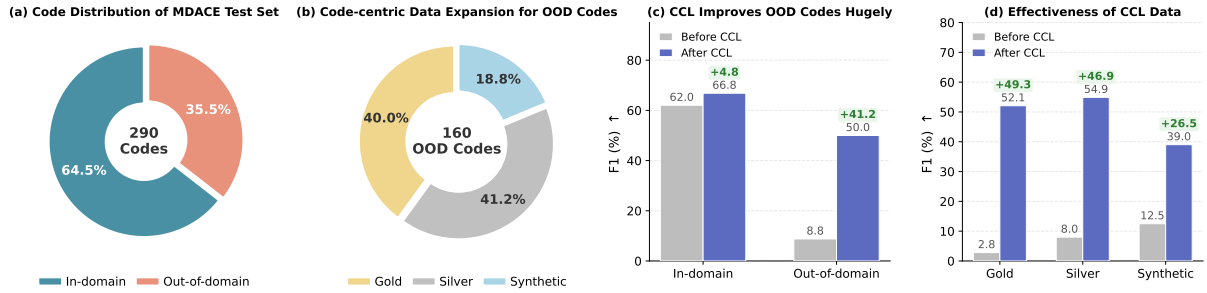


Figure 3: (a) Test-set codes are partitioned into in-domain and out-of-domain (unseen in the document-level training set). (b) Sources of spans constructed for OOD codes, with proportions of three strategies. (c) Our method improves OOD performance using spans only, without documents. (d) Each strategy contributes to improved coding accuracy.

	Ground Truth	Code-only Methods	Evidence-based Methods	CCL (Ours)
Evidence	Residual Pituitary Adenoma Diabetes Hypokalemia glaucoma cataract hypertension GERD	N/A No Evidence Provided 🤖 Hard to check	Hypertension GERD Glaucoma Cataract Hypokalemia pituitary adenoma OA ectopic pregnancy tubal ligation appendectomy Parathyroidectomy	Residual Pituitary Adenoma Diabetes Hypokalemia Glaucoma Cataract Hypertension GERD MS Evidence Provided 👩 I can check
ICD Codes	D35.2: Benign neoplasm of pituitary gland E11.9: Type 2 diabetes mellitus without complications E87.6: Hypokalemia H40.9: Unspecified glaucoma H26.9: Unspecified cataract I10: Essential (primary) hypertension K21.9: Gastro-esophageal reflux disease without esophagitis	D35.2 - Benign neoplasm of pituitary gland E11.9 - Type 2 diabetes mellitus without complications H40.9 - Unspecified glaucoma I10 - Essential (primary) hypertension K21.9 - Gastro-esophageal reflux disease without esophagitis	I10 - Essential (primary) hypertension K21.9 - Gastro-esophageal reflux disease without esophagitis H40.9 - Unspecified glaucoma E87.6 - Hypokalemia E22.0 - Pituitary adenoma H26.8 - Other specified cataract M19.90 - Unspecified osteoarthritis, unspecified site Z33.1 - History of ectopic pregnancy Z98.1 - Tubal ligation status Z90.89 - Acquired absence of other organs Z90.49 - Acquired absence of thyroid gland Correct Wrong codes Less Superfluous codes	D35.2: Benign neoplasm of pituitary gland E11.9: Type 2 diabetes mellitus without complications E87.6: Hypokalemia H40.9: Unspecified glaucoma H26.9: Unspecified cataract I10: Essential (primary) hypertension K21.9: Gastro-esophageal reflux disease without esophagitis G35 - Multiple sclerosis

Figure 4: An example from the test set. Code-only methods cannot generate evidence, making the results difficult for humans to evaluate and revise. Evidence-based methods suffer from limited evidence-annotated data for fine-tuning, and therefore achieve lower accuracy. Our method balances interpretability and accuracy, producing evidence and ICD codes that are highly consistent with human annotations.

562 knowledge learned during code-centric learning
563 can be transferred to document-level ICD coding,
564 which demonstrates the effectiveness of our proposed
565 hybrid training strategy.

566 **Code-Centric Data Expansion effectively**
567 **broadens the code coverage of document-level**
568 **ICD coding.** Overall, CCL data derived from
569 all three knowledge sources consistently improve
570 OOD coding performance, demonstrating that
571 document-level ICD coding can be effectively en-
572 hanced using span-level supervision alone. Further
573 analysis across knowledge sources shows that evi-
574 dence-code pairs mined from official guidelines
575 and public datasets yield larger performance gains,
576 while synthetic data also contributes to OOD code
577 recognition, albeit to a relatively smaller extent.

578 5.5 Case Study

579 Figure 4 shows an example from the test set. Un-
580 like **code-only methods** that achieve high accuracy
581 but only output final codes without supporting evi-
582 dence, and **evidence-based methods** whose accu-
583 racy is constrained by limited evidence-annotated
584 data, our approach is able to generate interpretable

585 evidence while maintaining accurate ICD coding.
586 In this case, the predicted ICD codes and support-
587 ing evidence are highly consistent with the ground
588 truth, which demonstrates that our method effec-
589 tively **balances interpretability and accuracy.**

590 6 Conclusion

591 In this paper, we identify a key finding: span-level
592 learning can effectively boost document-level ICD
593 coding. Based on this observation, we propose
594 a new LLM fine-tuning paradigm for ICD cod-
595 ing, termed code-centric learning. This training
596 paradigm incorporates mixed training and code-
597 centric data expansion, and simultaneously ad-
598 dresses three long-standing challenges in this do-
599 main: limited code coverage, lack of interpretabil-
600 ity and inefficient training on long documents. Our
601 paradigm is efficient, scalable and better suited for
602 real-world deployment in collaboration with human
603 coders. We hope that the community will adopt this
604 new paradigm and advance the long-standing goal
605 of fully automated ICD coding.

606 Limitations

607 **Data quality.** The scarcity of high-quality ICD
608 coding benchmarks remains a fundamental chal-
609 lenge for the research community. Although
610 MIMIC-III and MIMIC-IV are widely used, many
611 studies have repeatedly highlighted their annota-
612 tion inconsistencies and data quality issues (Cheng
613 et al., 2023; Edin et al., 2023; Yuan et al., 2025).
614 Even comparatively higher-quality benchmarks
615 such as MDACE still contain labeling errors
616 (Khadka et al., 2025). We therefore call for the
617 development of more reliable, rigorously curated
618 ICD coding benchmarks.

619 **Simplified setting.** First, existing datasets cover
620 only a limited subset of the ICD ontology, typically
621 a few thousand codes, which fall far short of the
622 tens of thousands of codes present in the full ICD-
623 10 system (Motzfeldt et al., 2025). Furthermore,
624 practical ICD coding is substantially more complex
625 than the simplified one-to-one mapping between
626 evidence spans and codes. ICD codes exhibit hi-
627 erarchical, inclusive, and mutually exclusive rela-
628 tionships, and these interactions strongly influence
629 how evidence should be identified and grouped
630 (CMS and NCHS, 2025). Modeling such relational
631 structures poses significant technical challenges.
632 Currently, the community still lacks such compre-
633 hensive benchmarks. Addressing these challenges
634 requires collective effort of the community.

635 Ethics Statement

636 We use the publicly available clinical dataset. We
637 do not see any ethics issues.

638 References

639 Krishanu Das Baksi, Elijah Soba, John J Higgins, Ravi
640 Saini, Jaden Wood, Jane Cook, Jack I Scott, Nirmala
641 Pudota, Tim Weninger, Edward Bowen, and Sanmitra
642 Bhattacharya. 2025. [MedCodER: A generative
643 AI assistant for medical coding](#). In *Proceedings of
644 the 2025 Conference of the Nations of the Americas
645 Chapter of the Association for Computational Lin-
646 guistics: Human Language Technologies (Volume 3:
647 Industry Track)*, pages 449–459, Albuquerque, New
648 Mexico. Association for Computational Linguistics.

649 Joseph Boyle, Antanas Kascenas, Pat Lok, Maria Li-
650 akata, and Alison O’Neil. 2023. [Automated clinical
651 coding using off-the-shelf large language models](#).
652 In *Deep Generative Models for Health Workshop
653 NeurIPS 2023*.

654 Elaine M Burns, E Rigby, R Mamidanna, A Bottle,
655 P Aylin, P Ziprin, and OD Faiz. 2012. Systematic

review of discharge coding accuracy. *Journal of
656 public health*, 34(1):138–148. 657

Hua Cheng, Rana Jafari, April Russell, Russell Klopfer,
658 Edmond Lu, Benjamin Striner, and Matthew R Gorm-
659 ley. 2023. [Mdace: Mimic documents annotated with
660 code evidence](#). In *Proceedings of the 61st Annual
661 Meeting of the Association for Computational Lin-
662 guistics (Volume 1: Long Papers)*, pages 7534–7550. 663

CMS and NCHS. 2025. [ICD-10-CM Official Guide-
664 lines for Coding and Reporting \(FY 2025\)](#). Technical
665 Report 10, The Centers for Medicare and Medicaid
666 Services (CMS) and the National Center for Health
667 Statistics (NCHS). 668

Joakim Edin, Alexander Junge, Jakob D Havtorn, Lasse
669 Borgholt, Maria Maistro, Tuukka Ruotsalo, and Lars
670 Maaløe. 2023. [Automated medical coding on mimic-
671 iii and mimic-iv: a critical review and replicability
672 study](#). In *Proceedings of the 46th International ACM
673 SIGIR Conference on Research and Development in
674 Information Retrieval*, pages 2572–2582. 675

Joakim Edin, Maria Maistro, Lars Maaløe, Lasse
676 Borgholt, Jakob Drachmann Havtorn, and Tuukka
677 Ruotsalo. 2024. [An unsupervised approach to
678 achieve supervised-level explainability in healthcare
679 records](#). In *Proceedings of the 2024 Conference on
680 Empirical Methods in Natural Language Processing*,
681 pages 4869–4890, Miami, Florida, USA. Association
682 for Computational Linguistics. 683

Yidong Gan, Maciej Rybinski, Ben Hachey, and
684 Jonathan K. Kummerfeld. 2025. [Aligning AI re-
685 search with the needs of clinical coding workflows:
686 Eight recommendations based on US data analysis
687 and critical review](#). In *Proceedings of the 63rd An-
688 nual Meeting of the Association for Computational
689 Linguistics (Volume 1: Long Papers)*, pages 909–922,
690 Vienna, Austria. Association for Computational Lin-
691 guistics. 692

Xueren Ge, Abhishek Satpathy, Ronald Dean Williams,
693 John Stankovic, and Homa Alemzadeh. 2024. [DKEC: Domain knowledge enhanced multi-label
694 classification for diagnosis prediction](#). In *Proceed-
695 ings of the 2024 Conference on Empirical Methods in
696 Natural Language Processing*, pages 12798–12813,
697 Miami, Florida, USA. Association for Computational
698 Linguistics. 699

Goncalo Gomes, Isabel Coutinho, and Bruno Martins.
701 2024. [Accurate and well-calibrated ICD code as-
702 signment through attention over diverse label embed-
703 dings](#). In *Proceedings of the 18th Conference of the
704 European Chapter of the Association for Computa-
705 tional Linguistics (Volume 1: Long Papers)*, pages
706 2302–2315, St. Julian’s, Malta. Association for Com-
707 putational Linguistics. 708

Jan Horsky, Elizabeth A Drucker, and Harley Z Ramel-
709 son. 2018. [Accuracy and completeness of clinical
710 coding using icd-10 for ambulatory visits](#). In *AMIA
711 annual symposium proceedings*, volume 2017, page
712 912. 713

714	Chao-Wei Huang, Shang-Chi Tsai, and Yun-Nung Chen.	Alexandr Nesterov, Andrey Sakhovskiy, Ivan Sviri-	771
715	2022. PLM-ICD: Automatic ICD coding with pre-	dov, Airat Valiev, Vladimir Makharev, Petr Anokhin,	772
716	trained language models . In <i>Proceedings of the 4th</i>	Galina Zubkova, and Elena Tutubalina. 2025. RuC-	773
717	<i>Clinical Natural Language Processing Workshop</i> ,	CoD: Towards automated ICD coding in Russian .	774
718	pages 10–20, Seattle, WA. Association for Computa-	In <i>Proceedings of the 2025 Conference on Empiri-</i>	775
719	tional Linguistics.	<i>cal Methods in Natural Language Processing</i> , pages	776
720	Alistair EW Johnson, Lucas Bulgarelli, Lu Shen, Alvin	2558–2585, Suzhou, China. Association for Computa-	777
721	Gayles, Ayad Shammout, Steven Horng, Tom J Pol-	tional Linguistics.	778
722	lard, Sicheng Hao, Benjamin Moody, Brian Gow,	Xindi Wang, Robert Mercer, and Frank Rudzicz. 2024a.	779
723	et al. 2023. Mimic-iv, a freely accessible electronic	Multi-stage retrieve and re-rank model for automatic	780
724	health record dataset. <i>Scientific data</i> , 10(1):1.	medical coding recommendation . In <i>Proceedings of</i>	781
725	Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H	<i>the 2024 Conference of the North American Chap-</i>	782
726	Lehman, Mengling Feng, Mohammad Ghassemi,	<i>ter of the Association for Computational Linguistics:</i>	783
727	Benjamin Moody, Peter Szolovits, Leo Anthony Celi,	<i>Human Language Technologies (Volume 1: Long</i>	784
728	and Roger G Mark. 2016. Mimic-iii, a freely accessi-	<i>Papers)</i> , pages 4881–4891, Mexico City, Mexico. As-	785
729	ble critical care database. <i>Scientific data</i> , 3(1):1–9.	sociation for Computational Linguistics.	786
730	Supriya Khadka, Xiaorui Jiang, and Vasile Palade. 2025.	Xindi Wang, Robert E. Mercer, and Frank Rudzicz.	787
731	Data quality in clinical coding: A critical analysis	2024b. Auxiliary knowledge-induced learning for	788
732	and preliminary study. <i>medRxiv</i> , pages 2025–08.	automatic multi-label medical document classifica-	789
733	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	tion . In <i>Proceedings of the 2024 Joint International</i>	790
734	Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.	<i>Conference on Computational Linguistics, Language</i>	791
735	Gonzalez, Hao Zhang, and Ion Stoica. 2023. Effi-	<i>Resources and Evaluation (LREC-COLING 2024)</i> ,	792
736	cient memory management for large language model	pages 2006–2016, Torino, Italia. ELRA and ICCL.	793
737	erving with pagedattention. In <i>Proceedings of the</i>	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le,	794
738	<i>ACM SIGOPS 29th Symposium on Operating Systems</i>	Ed Chi, Sharan Narang, Aakanksha Chowdhery, and	795
739	<i>Principles</i> .	Denny Zhou. 2022. Self-consistency improves chain	796
740	Tuan-Dung Le, Shohreh Haddadan, and Thanh Q.	of thought reasoning in language models. <i>arXiv</i>	797
741	Thieu. 2025. Ace-icd: Acronym expansion as data	<i>preprint arXiv:2203.11171</i> .	798
742	augmentation for automated icd coding . <i>Preprint</i> ,	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	799
743	arXiv:2511.07311.	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	800
744	Rumeng Li, Xun Wang, and Hong Yu. 2024. Explor-	et al. 2022. Chain-of-thought prompting elicits rea-	801
745	ing llm multi-agents for icd coding. <i>arXiv preprint</i>	soning in large language models. <i>Advances in neural</i>	802
746	<i>arXiv:2406.15363</i> .	<i>information processing systems</i> , 35:24824–24837.	803
747	Junyu Luo, Xiaochen Wang, Jiaqi Wang, Aoife Chang,	John Wu, David Wu, and Jimeng Sun. 2024. Beyond	804
748	Yaqing Wang, and Fenglong Ma. 2024. CoRelation:	label attention: Transparency in language models for	805
749	Boosting automatic ICD coding through contextual-	automated medical coding via dictionary learning .	806
750	ized code relation learning . In <i>Proceedings of the</i>	In <i>Proceedings of the 2024 Conference on Empiri-</i>	807
751	<i>2024 Joint International Conference on Computa-</i>	<i>cal Methods in Natural Language Processing</i> , pages	808
752	<i>tional Linguistics, Language Resources and Eval-</i>	8848–8871, Miami, Florida, USA. Association for	809
753	<i>uation (LREC-COLING 2024)</i> , pages 3997–4007,	Computational Linguistics.	810
754	Torino, Italia. ELRA and ICCL.	Wen-wai Yim, Yujuan Fu, Asma Ben Abacha, Neal	811
755	Andreas Geert Motzfeldt, Joakim Edin, Casper L. Chris-	Snider, Thomas Lin, and Meliha Yetisgen. 2023. Ac-	812
756	tensen, Christian Hardmeier, Lars Maaløe, and Anna	bench: a novel ambient clinical intelligence dataset	813
757	Rogers. 2025. Code like humans: A multi-agent	for benchmarking automatic visit note generation .	814
758	solution for medical coding . In <i>Findings of the Asso-</i>	<i>Scientific data</i> , 10(1):586.	815
759	<i>ciation for Computational Linguistics: EMNLP 2025</i> ,	Moy Yuan, Han-Chin Shing, Mitch Strong, and Chai-	816
760	pages 22612–22627, Suzhou, China. Association for	tanya Shivade. 2025. Toward reliable clinical coding	817
761	Computational Linguistics.	with language models: Verification and lightweight	818
762	James Mullenbach, Sarah Wiegrefe, Jon Duke, Jimeng	adaptation . In <i>Proceedings of the 2025 Conference</i>	819
763	Sun, and Jacob Eisenstein. 2018. Explainable predic-	<i>on Empirical Methods in Natural Language Process-</i>	820
764	tion of medical codes from clinical text . In <i>Proceed-</i>	<i>ing: Industry Track</i> , pages 173–184, Suzhou (China).	821
765	<i>ings of the 2018 Conference of the North American</i>	Association for Computational Linguistics.	822
766	<i>Chapter of the Association for Computational Lin-</i>	Zheng Yuan, Chuanqi Tan, and Songfang Huang. 2022.	823
767	<i>guistics: Human Language Technologies, Volume</i>	Code synonyms do matter: Multiple synonyms	824
768	<i>1 (Long Papers)</i> , pages 1101–1111, New Orleans,	matching network for automatic ICD coding . In	825
769	Louisiana. Association for Computational Linguis-	<i>Proceedings of the 60th Annual Meeting of the As-</i>	826
770	tics.	<i>sociation for Computational Linguistics (Volume 2:</i>	827

828 *Short Papers*), pages 808–814, Dublin, Ireland. As-
829 sociation for Computational Linguistics.

830 Xu Zhang, Kun Zhang, Wenxin Ma, Rongsheng Wang,
831 Chenxu Wu, Yingtai Li, and S Kevin Zhou. 2025.
832 [A general knowledge injection framework for ICD](#)
833 [coding](#). In *Findings of the Association for Computa-*
834 *tational Linguistics: ACL 2025*, pages 7180–7189,
835 Vienna, Austria. Association for Computational Lin-
836 guistics.

837 Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan
838 Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma.
839 2024. [Llamafactory: Unified efficient fine-tuning](#)
840 [of 100+ language models](#). In *Proceedings of the*
841 *62nd Annual Meeting of the Association for Computa-*
842 *tional Linguistics (Volume 3: System Demonstra-*
843 *tions)*, Bangkok, Thailand. Association for Computa-
844 tional Linguistics.

845	Appendix		
846	A ICD Coding Background		
847	A.1 ICD coding task.		
848	The ICD Coding task recognizes diseases, symptoms, conditions and procedures in a medical document, including discharge summaries, progress notes, and operative reports, and assign standardized ICD (International Classification of Diseases) codes to them. This task plays a critical role in healthcare administration, clinical statistics, reimbursement systems, and medical research.		
856	From a computational perspective, ICD Coding is commonly formulated as a text-to-code prediction problem. Given a patient-level clinical document, the model is required to output a set of ICD codes. The task is characterized by a large label space, hierarchical code structures, and severe label imbalance, which together make ICD Coding a challenging and distinctive problem in clinical natural language processing.		
865	A.2 Distinction from Related Tasks		
866	Although ICD Coding seems to share similarities with several well-studied tasks, its objectives and constraints differ substantially from those tasks, which may lead to confusion.		
870	Difference from Diagnosis. Diagnosis aims to infer or determine what diseases a patient has, often involving clinical reasoning, uncertainty management, and causal inference. In contrast, ICD Coding does not seek to generate new diagnostic conclusions. Instead, it focuses on assigning standardized codes based solely on diagnoses and clinical facts that have already been documented by healthcare professionals. Therefore, ICD coding should be viewed as an information standardization task rather than a diagnostic or decision-making task.		
882	Difference from Multi-label Text Classification. In conventional multi-label settings, labels are typically assumed to be conditionally independent given the input text, and prediction is treated as a parallel, order-agnostic decision process. In contrast, ICD Coding involves strong dependencies among codes. In real-world coding practice, human coders do not assign all codes independently or simultaneously. Instead, they follow official coding guidelines to sequentially identify the principal diagnosis, secondary diagnoses, supplementary conditions, and procedures, with each step		
894	constraining and informing subsequent coding decisions. This process is inherently procedural and generative, rather than purely discriminative.	894	
895		895	
896		896	
897	Difference from Information Extraction. Information extraction tasks typically focus on identifying entities, relations, or events explicitly mentioned in text. By contrast, ICD Coding requires the prediction of standardized code identifiers rather than text spans, and often relies on a holistic understanding of the entire document rather than localized entity mentions.	897	
898		898	
899		899	
900		900	
901		901	
902		902	
903		903	
904		904	
905	A.3 Authoritative Resources in ICD Coding		
906	Alphabetic Index. The Alphabetic Index maps various synonyms, abbreviations, and lexical variants to candidate ICD codes, thereby bridging the gap between natural language expressions and standardized code identifiers. Importantly, the codes suggested by the Alphabetic Index are not definitive; rather, they represent preliminary references that must be further validated.	906	
907		907	
908		908	
909		909	
910		910	
911		911	
912		912	
913		913	
914	Tabular List. The Tabular List is the authoritative, structured listing of all valid ICD codes, organized by chapters, categories, subcategories, and extensions. Each code entry in the Tabular List is accompanied by a formal definition and may include additional annotations such as inclusion terms, exclusion notes, code-first instructions, and combination code indicators. Coders are required to confirm all codes suggested by the Alphabetic Index against the Tabular List before assignment.	914	
915		915	
916		916	
917		917	
918		918	
919		919	
920		920	
921		921	
922		922	
923		923	
924	Coding Guidelines. The Coding Guidelines provide a comprehensive set of rules and conventions that govern how ICD codes should be applied in practice. Guidelines often specify conditional logic (e.g., “code first,” “use additional code,” or “do not code separately”) and clarify how multiple diagnoses or clinical conditions should be represented in a single episode.	924	
925		925	
926		926	
927		927	
928		928	
929		929	
930		930	
931		931	
932	In practical ICD coding workflows, these resources are used in a complementary and sequential manner. The Alphabetic Index supports initial term-to-code lookup, the Tabular List determines valid and precise code selection, and the Coding Guidelines regulate how codes are combined, ordered, and reported.	932	
933		933	
934		934	
935		935	
936		936	
937		937	
938		938	

939 **B Prompts**

940 **B.1 Prompts for Code-centric Learning** 941 **Framework**

942 In this section, we present all the prompts used
943 in our approach. Specifically, we describe the
944 prompts for Mixed Training and Code-centric Data
945 Expansion.

946 Mixed Training relies on two types of data for-
947 mats: (1) document-level evidence-based ICD cod-
948 ing data, and (2) span-level data designed for code-
949 centric learning. We show the prompts of these two
950 different tasks in Table 4.

951 For Code-centric Data Expansion, we show the
952 prompts used to construct Silver Pairs and Syn-
953 thetic Pairs, as Gold Pairs are primarily obtained
954 from the Official Alphabetic Index.

955 To construct Silver Pairs, we employ LLaMA
956 3.1-70B to mine all supporting evidence from each
957 MIMIC-IV sample, followed by deduplication and
958 refinement of the evidence associated with each
959 ICD code. We show the used prompts in Table 6.

960 For Synthetic Pairs, we use GPT-5.1 to infer un-
961 seen ICD codes based on existing Gold and Silver
962 Pairs. We show the prompts in Table 5.

963 **B.2 Prompts for Baselines**

964 In this section, we present all the prompts used by
965 the generative baselines, as shown in Table 7.

966 For Chain-of-Thought (CoT), we adopt the stan-
967 dard CoT prompting strategy.

968 For CoT-SC, we use the same prompt as CoT,
969 but retain only those ICD codes that appear in at
970 least three out of five reasoning runs.

971 For MAC, we make minor modifications to the
972 original prompt to adapt it from ICD-9 to ICD-10,
973 as the original method was evaluated primarily on
974 ICD-9.

975 For CLH, we use the official open-source imple-
976 mentation and apply it directly to our dataset.

977 For the SFT on MIMICIV and Code-only ICD
978 Coding setting on MDACE in the ablation study,
979 we use the same prompt, following Yuan et al.
980 (2025) to add descriptions after ICD codes.

981 **C Implementation Details**

982 For LLM inference, we use vLLM (Kwon et al.,
983 2023). For SFT, we use LLaMa Factory (Zheng
984 et al., 2024). All the experiments are implemen-
985 tated on a single H20 GPU with 96GB of VRAM.

Table 4: Prompt templates used for Mixed Training

Data Type	Prompt Template
Document-level ICD coding Data	<p>Task:</p> <p>You are a clinical coding assistant.</p> <p>Your task is to analyze the provided clinical note, first extract all relevant clinical evidences that supports diagnostic coding, and then output the corresponding ICD-10-CM codes.</p> <p>Example</p> <p>### Clinical Note: ...</p> <p>### Evidence</p> <p>CAD COPD Anemia</p> <p>### ICD-10-CM Codes</p> <p>I25.10 – Atherosclerotic heart disease of native coronary artery without angina pectoris J44.9 – Chronic obstructive pulmonary disease, unspecified D62 – Acute posthemorrhagic anemia</p> <p>---</p> <p>### Clinical Note: \{text\}</p>
Span-level Code-Centric Learning Data	<p>### Evidence: \{evidence\}</p> <p>### ICD-10-CM Codes:</p>

Table 5: Prompt templates used for Code-Centric Data Expansion (Silver Pairs)

Task	Prompt Template
Evidence Extraion	<p>You are a professional ICD–10–CM coder.</p> <p>Your task is to extract the <i>*verbatim minimal text spans*</i> that supports each ICD–10–CM code. If no explicit evidence exists in the note, output: "No evidence found".</p> <p>---</p> <p>Example</p> <p>### Clinical Note: ...</p> <p>### ICD–10–CM Codes ...</p> <p>### Evidence</p> <p>I25.10 – Atherosclerotic heart disease of native coronary artery without angina pectoris > CAD</p> <p>J44.9 – Chronic obstructive pulmonary disease, unspecified > COPD</p> <p>D62 – Acute posthemorrhagic anemia > Anemia</p> <p>---</p> <p>### Clinical Note {text}</p> <p>### ICD–10–CM Codes {diagnosis_codes}</p> <p>### Evidence</p>
Evidence Refinement	<p>You are a professional ICD–10–CM coder.</p> <p>Your task is to update and refine the Evidence Set for the ICD–10–CM code below. Follow these rules:</p> <ol style="list-style-type: none"> 1. Only keep the <i>**most essential**</i> evidence that clearly supports this code. 2. You may reference the Alphabetic Index terms, but you do not need to match them exactly. 3. Use the <i>**Original Evidence Set**</i> as the base. <ul style="list-style-type: none"> – If the MIMIC–IV evidence contains new, meaningful, or more specific expressions, add them. – If not, keep the existing evidence unchanged. 4. Remove duplicates and unify phrasing into <i>**clear, concise, canonical**</i> clinical expressions. 5. Output the <i>**updated Evidence Set only**</i>, as a bullet list. No explanation. <p>---</p> <p>### ICD–10–CM Code {code}</p> <p>### Alphabetic Index Term {alphabetic_index_term}</p> <p>### Original Evidence Set {evidence_set}</p> <p>### New Evidence from MIMIC–IV {mimiciv_evidence}</p> <p>### Updated Evidence Set -</p>

Table 6: Prompt templates used for Code-Centric Data Expansion (Synthetic Pairs)

Task	Prompt Template
Synthesize Evidence	<p>You are a professional ICD–10–CM coding and clinical documentation expert.</p> <p>Your task is to synthesize a focused, audit–defensible list of clinical evidence terms that directly support assignment of the ICD–10–CM code: {code}.</p> <p>Definition of evidence: Evidence refers only to clinical findings or documentation elements that materially support the diagnosis represented by the code.</p> <p>Available references: {reference}</p> <p>Instructions:</p> <ul style="list-style-type: none"> – Use the parent and sibling codes to understand diagnostic scope. – Infer conservatively based on ICD–10–CM conventions and real–world clinical documentation patterns. – Prioritize diagnostic–confirmatory evidence (e.g., imaging findings, explicit diagnoses, anatomical localization). <p>Do NOT include:</p> <ul style="list-style-type: none"> – Mechanism of injury or accident descriptions – General symptoms or nonspecific complaints – Treatment, procedures, immobilization, or care plans – Encounter setting or workflow details – Redundant negative statements unless required to distinguish code type <p>Unspecified code rule:</p> <ul style="list-style-type: none"> – If the code is unspecified, do NOT introduce inferred specificity (e.g., displacement, fracture pattern, severity). <p>Output constraints:</p> <ul style="list-style-type: none"> – Consolidate overlapping or synonymous terms. – Stop generating new items once additional terms no longer add distinct coding value. <p>Output format:</p> <ul style="list-style-type: none"> – <evidence term> – <evidence term> ...

Table 7: Prompt templates used for Generative Baselines

Baselines	Prompt Template
CoT	<p>You are a clinical coding assistant.</p> <p>Your task is to analyze the provided clinical note, and then output the corresponding ICD-10-CM codes.</p> <p>### Clinical Note: {text}</p> <p>Let's think step by step.</p>
MAC-coder	<p>You are an ICD-10 coder.</p> <p>You assign ICD-10 codes to the discharge summary based on the clinical care that the patients received.</p> <p>You cite the discharge summary as evidence when needed.</p> <p>You assign as many as possible ICD-10 codes and explain the reasons for each code.</p> <p>The discharge summary is: {text}</p>
MAC-reviewer	<p>You are a reviewer.</p> <p>You will check the ICD-10 codes assigned by the coder.</p> <p>You can use the ICD-10 dictionary for guidance.</p> <p>Your role is to ensure that the assigned ICD-10 codes are correct.</p> <p>You assign all possible ICD-10 codes and explain the reasons for each code.</p> <p>The discharge summary is: {text}</p> <p>The ICD-10 codes assigned by the coder are: {coder_pred}</p>
MAC-physician	<p>You are a physician who treats patients.</p> <p>You strive to provide the best service to each patient.</p> <p>You document your findings, interventions and results in the discharge summary note.</p> <p>You check all assigned ICD-10 codes and explain the reasons for each code.</p> <p>The discharge summary is: {text}</p> <p>The ICD-10 codes assigned by the coder are: {reviewer_pred}</p>
MAC-patient	<p>You are a patient who received treatment at the hospital.</p> <p>You cooperate fully with the health care system to receive the best service possible.</p> <p>You also check the ICD-10 codes to avoid being overbilled.</p> <p>You check all assigned ICD-10 codes and explain the reasons for each code.</p> <p>The discharge summary is: {text}</p> <p>The ICD-10 codes assigned by the coder are: {reviewer_pred}</p>
MAC-adjustor	<p>When a patient or a physician has different thoughts about the ICD-10 codes, you will review the discharge summary and the ICD codes assigned by the coder and checked by the reviewer.</p> <p>You can add, remove the assigned codes to make them accurate.</p> <p>You can consult the ICD-10 dictionary for assistance.</p> <p>Your duty is to ensure that the assigned ICD-10 codes are valid and exact.</p> <p>You assign all possible ICD-10 codes and explain the reasons for each code.</p> <p>The discharge summary is {text}</p> <p>The ICD-10 codes assigned by the physician are {physician_pred}</p> <p>The ICD-10 codes assigned by the patient are {patient_pred}</p> <p>The ICD-10 codes assigned by the coder are {coder_pred}</p> <p>The ICD-10 codes checked by the reviewer are {reviewer_pred}</p>
SFT / Code-only	<p>You are a clinical coding assistant.</p> <p>Your task is to analyze the provided clinical note, and then output the corresponding ICD-10-CM codes.</p> <p>### Clinical Note: {text}</p>