Don't Just Pay Attention, PLANT It: Transfer L2R Models to Fine-tune Attention in Extreme Multi-Label Text Classification For ICD Coding

Anonymous ACL submission

Abstract

The keystone of state-of-the-art Extreme Multi-Label Text Classification (XMTC) models is the multi-label attention layer within the decoder, which deftly directs label-specific focus to salient tokens in input text. Nonetheless, the process of acquiring these optimal attention weights is onerous and resource-intensive. To alleviate this strain, we introduce PLANT -Pretrained and Leveraged AtteNTion - an innovative transfer learning strategy to fine-tune XMTC decoders. The central notion involves transferring a pretrained learning-to-rank (L2R) model, utilizing its activations as attention weights, thereby serving as the 'planted' attention layer in the decoder. On the full MIMIC-III dataset, PLANT excels in four out of seven metrics and surpasses in five for the top-50 code set, demonstrating its effectiveness. Remarkably, for the rare-50 code set, PLANT achieves a significant 12.7 - 52.2% improvement in four metrics. On MIMIC-IV, it leads in three metrics. Notably, in low-shot scenarios, PLANT matches traditional attention models' precision despite using significantly less data ($\frac{1}{10}$ for precision at 5, $\frac{1}{5}$ for precision at 15), highlighting its efficiency with skewed label distributions.

1 Introduction

004

007

009

013

015

017

019

021

022

028

Extreme Multi-Label Text Classification (XMTC) addresses the problem of automatically assigning each data point with most relevant subset of labels from an extremely large label set, often containing hundreds of thousands, even millions of labels and samples in various real-world XMTC applications. One major application of XMTC is in the global healthcare system, specifically in the context of the International Classification of Diseases (ICD)¹. ICD coding is the process of assigning codes representing diagnoses and procedures performed during a patient visit using clinical notes documented by health professionals (Table 1). ICD codes are

¹https://www.who.int/standards/ classifications/classification-of-diseases

998.32 : Disruption of external operation wound
··· wound infection, and wound breakdown ···
428.0 : Congestive heart failure
··· DIAGNOSES: 1. Acute congestive heart failure
2. Diabetes mellitus 3. Pulmonary edema · · ·
202.8 : Other malignant lymphomas
··· a 55 year-old female with non Hodgkin's lymphoma
and acquired C1 esterase inhibitor deficiency · · ·
770.6 : Transitory tachypnea of newborn
··· Chest x-ray was consistent with transient tachypnea
of the newborn · · ·
424.1 : Aortic valve disorders
\cdots mild aortic stenosis with an aortic valve area of
1.9 cm squared and $2+$ aortic insuffiency \cdots

Table 1: Examples of clinical text fragments and their corresponding ICD codes (Li and Yu, 2020).

used for both epidemiological studies and billing of services (Bottle and Aylin, 2008). XMTC has been utilized to automate the manual ICD coding performed by clinical coders which is time intensive and prone to human errors (O'malley et al., 2005; Nguyen et al., 2018). 041

042

044

045

046

047

054

056

060

061

062

063

Main Challenge: Building XMTC models is challenging because datasets often consist of texts with multiple lengthy narratives - more than 1500 tokens (i.e., words) on average. However, only a small fraction of tokens are most informative with regard to assigning relevant labels. Automatically assigning labels become even more challenging when, (1) the label space is extremely high dimensional, and, (2) the label distribution is heavily skewed. For example, in automatic ICD coding, there are over 18000 and 170000 codes in ICD-9-CM and ICD-10-CM/PCS², respectively. The skewness of ICD-9-CM label distribution in the MIMIC-III dataset (Johnson et al., 2016) is evident from the fact that approximately 5411 out of all the 8929 codes appear less than 10 times (refer to Appendix A.1, Figure 5 for a visual).

²https://www.cdc.gov/nchs/icd/icd10cm_pcs_ background.htm

How SOTA models address the main challenge in XMTC? (Red Box) In state-of-the-art (SOTA) 065 NLP models, the inclusion of attention mecha-066 nisms is crucial, benefiting various applications like Machine Translation, Summarization, Text Representation, Sentiment Analysis, and Question Answering (Vaswani et al., 2017; Tang et al., 2018; Xu et al., 2020; Kiela et al., 2018; Wang et al., 071 2020; Dehghani et al., 2018). In XMTC, these attention mechanisms play a vital role in addressing 073 the challenges of high-dimensional label spaces and skewed label distributions. XMTC models (Mullenbach et al., 2018; Xie et al., 2019; Li and Yu, 2020; Cao et al., 2020; Vu et al., 2021; Zhou 077 et al., 2021; Liu et al., 2021; Yuan et al., 2022; Zhang et al., 2022; Yang et al., 2022) consistently feature a multi-label attention layer, dynamically allocating label-specific attention weights to the most informative tokens in input text. Refer to the components highlighted in red in Figure 1, which illustrate this critical attention layer in action. Regardless of the specific encoder architecture, removing this attention layer leads to a significant drop in performance. 087

Main Shortfall in Red Box: Current SOTA XMTC models often begin with random attention 089 weights, requiring them to rank all tokens for each label from scratch, a computationally intensive process, which, given the high-dimensional label space characteristic of XMTC datasets, leads to extensive 093 computational requirements and prolonged training 094 times. Moreover, the presence of heavily skewed la-095 bel distributions further exacerbates this challenge, as rare labels necessitate even longer training durations and increase the risk of overfitting (Figure 4). Corroborating the issue of rare codes, the study in (Edin et al., 2023) reveals that SOTA models 100 exhibit considerable difficulties when predicting 101 rare ICD diagnosis codes (Figure 2). Models tend 102 to perform similarly across codes with compara-103 ble frequencies, implicating the higher proportion of rare codes in ICD as a significant factor in per-105 formance disparities. Correlations between code 106 frequency and F1 score are moderately high, indi-107 cating that rare codes are predicted with less accuracy than common ones. This inherent complexity 109 underscores the need for efficient mechanisms to 110 learn optimal attention configurations in XMTC 111 models, as starting with random weights may not 112 suffice. 113



Figure 1: Architecture of PLANT showcasing the integration of contemporary SOTA components (grey box), multi-label attention (red box), planted attention (green box), and mutual information gain (yellow box) to enhance label prediction efficacy.

Our Contributions (Green Box):

1. To address this shortfall, we propose PLANT – 115 Pretrained and Leveraged AtteNTion, a novel 116 transfer learning mechanism to fine-tune atten-117 tion in XMTC. The core idea is to bootstrap 118 using mutual information gain (refer to the 119 yellow components in Figure 1) a standalone 120 model that learns to rank (L2R) tokens based 121 on their relevance to labels. The pretrained 122 L2R model (refer to the green components in 123 Figure 1) that leverages its activations as atten-124 tion weights serves as the 'planted' attention 125 layer in the XMTC decoder. This tranferring 126 of the L2R model ensures the decoder starts 127



Figure 2: Comparative analysis of model performance from (Edin et al., 2023) on rare versus common ICD diagnosis codes, highlighting that rare codes have near zero macro-F1scores.

129 130

131

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

160

161

162

163

164

with well-informed attention weights rather than training from scratch with randomly initialized weights. Subsequent fine-tuning enables not only efficient convergence toward optimal attention weight configurations but also enhances the model's ability to prioritize salient features of the input texts while minimizing the risk of overfitting by adapting them to the specific characteristics of the target dataset. Notably, we compared PLANT with a SOTA model LAAT (Vu et al., 2021) on MIMIC-IV-full, showing that PLANT avoids overfitting during training (Figure 4).

- 2. Addressing the shortfall in rare code prediction (Edin et al., 2023), PLANT is particularly effective in dealing with high dimenstional skewed label distributions in a low shot setting. It demonstrates comparable precision to traditional attention models, even with substantially less data – $\frac{1}{10}$ for precision at 5, $\frac{1}{5}$ for precision at 15 (Figure 3).
- 3. We introduce the *inattention* technique, which strategically filters out less relevant tokens, enhancing the significance of attention weights and enabling a sharper focus on critical elements within a token sequence. Additionally, inspired by Backpropagation-Through-Timefor-Text-Classification (Howard and Ruder, 2018), we propose a *stateful decoder* that accumulates information across segments, enabling cumulative predictions. This mechanism utilizes batch-level states, improving adaptability to large documents and model convergence, eliminating text truncation needs, and ensuring stable GPU memory usage, thereby enhancing both performance and efficiency (Table 7).

4. We extensively evaluated PLANT on benchmark MIMIC-III and newly available MIMIC-166 IV datasets, widely used in automatic ICD 167 coding research. Compared to 10 existing 168 SOTA models (Section 3.2), PLANT outper-169 formed them across 7 different evaluation 170 metrics. Specifically, in MIMIC-III-full, 171 MIMIC-III-top50, MIMIC-III-rare50, and MIMIC-IV-full datasets, PLANT exhibited 173 significant performance improvements (Ta-174 ble 3, Table 4, Table 5, Table 6). We 175 also conducted rigorous ablation analysis 176 (Section 4) and made our trained models and code available at https://anonymous. 178 4open.science/r/brainsplant/.

165

172

177

179

181

182

183

184

185

186

188

189

190

191

192

193

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

2 Approach

XMTC: The input is a set of documents and their corresponding labels, $\mathcal{D} = \{(\boldsymbol{x}_i, \boldsymbol{y}_i) \mid \boldsymbol{y}_i \in$ $\{0,1\}^{|\mathcal{L}|}, i = 1, ..., d\}$, where \mathcal{L} is the set of labels. The goal of XMTC is to learn a prediction function $\hat{y}(\boldsymbol{x}_i) \in \mathbb{R}^{|\mathcal{L}|}$. The function \hat{y} should be optimized such that the $\hat{y}(x_{il})$ is high when $y_{il} = 1$ (i.e., label l is relevant to x_i), and is low when $y_{il} = 0.$

Intuition behind our XMTC model - PLANT (Figure 1): The intuitive flow starts with document tokenization into embeddings processed by a pretrained AWD-LSTM to grasp textual contexts. The decoder introduces planted attention (green box), leveraging a L2R model's ability to rank token significance by label relevance, enriching the model with a pre-understanding of token-label dynamics. This is adeptly paired with multi-label attention (red box), merging learned and pretrained insights for feature prominence. Additionally, mutual information gain (yellow box) is utilized to enhance the decision-making process by calculating the relevance of each token to the potential labels, providing an informed basis for further attention refinement. A subsequent boost attention phase fine-tunes this for label-specific discernment, culminating in a sigmoid-derived label probability prediction. Section 2.1 provides a detailed description of the L2R model components, while Section 2.2 explains how we utilize the pretrained L2R model for planted attention, illustrating the integration of the green boxes in Figure 1.

214

215

216

217

218

219

222

223

234

239

241

242

244

245

246

247

248

252

253

2.1 L2R Model

We use superscript to denote the id of a label and subscript to denote the id of a token. The training set of the L2R model contains a set of labels $\mathcal{L} = \{l^{(1)}, l^{(2)}, \dots, l^{(m)}\}$, and a set of tokens $\mathcal{T} = \{t_1, t_2, \dots, t_n\}$. Furthermore, $G = [g^{(1)}, g^{(2)}, \dots, g^{(m)}] \in \mathbb{R}^{n \times m}$, and $g^{(i)} =$ $[g_1^{(i)}, g_2^{(i)}, \dots, g_n^{(i)}]^T \in \mathbb{R}^n$, where $g_j^{(i)}$ denotes the relevance of the token t_j with respect to label $l^{(i)}$. We represent each label $l^{(i)}$ and token t_j with word embeddings $e_{l^{(i)}}$ and e_{t_j} , respectively. A feature vector

$$\boldsymbol{x}_{j}^{(i)} = \Psi\left(\boldsymbol{e}_{l^{(i)}}, \boldsymbol{e}_{t_{j}}\right)$$
(1)

is created from each label-token pair $(l^{(i)}, t_j)$, $i = 1, 2, \dots, m; j = 1, 2, \dots, n$, by concatenating the corresponding word embeddings $e_{l^{(i)}}$ and e_{t_j} . The feature matrix, $X^{(i)} = [x_1^{(i)}, \dots, x_n^{(i)}]$ and the corresponding scores, $g^{(i)} = [g_1^{(i)}, g_2^{(i)}, \dots, g_n^{(i)}]^T$ then form an 'instance'. The training set can be denoted as $\{(X^{(i)}, g^{(i)})\}_{i=1}^m$. The L2R model is associated with a ranking function, $f : x_j^{(i)} \mapsto \mathbb{R}$. At any point in the training, the model outputs the score $z^{(i)} = [f(x_1^{(i)}), \dots, f(x_n^{(i)})]^T \in \mathbb{R}^n$. We direct readers to Appendix A.2 for detailed specifics about the L2R model, including our methods for bootstrapping it with mutual information gain and subsequent training procedures.

2.2 Leveraging L2R as Pretrained Attention

Pretrained and Fine-tuned AWD-LSTM: We use the AWD-LSTM architecture (Merity et al., 2017) as LM in our experiments³. That means, AWD-LSTM model learns hidden features from a sequence of *n* tokens $\langle t_1, t_2, \dots, t_n \rangle$, where each token is represented by word embedding $e_{t_j} \in \mathbb{R}^{s_e}$. The hidden feature learned by AWD-LSTM corresponding to the j^{th} token is represented as:

$$h_j = \mathsf{AWD}\text{-}\mathsf{LSTM}(\langle e_{t_1}, \cdots, e_{t_j} \rangle), h_j \in \mathbb{R}^{s_e}$$
(2)

Note that all the pretrained word embeddings e_{t_j} and the parameters of the AWD-LSTM model are finetuned on the target task using the mechanisms proposed in Howard and Ruder (2018). **Decoder – PLANT L2R as Attention**: To allocate label-specific attention weights to the most informative tokens in the sequence $\langle t_1, t_2, \cdots, t_n \rangle$ we take the following three steps.

254

255

257

258

259

260

261

262

263

264

266

267

270

271

272

273

274

275

276

277

278

279

283 284

287

291

294

296

297

298

300

First, the hidden features h_1, h_2, \dots, h_n of the sequence $\langle t_1, t_2, \dots, t_n \rangle$ are concatenated to formulate the matrix $\boldsymbol{H} = [\boldsymbol{h}_1, \boldsymbol{h}_2, \dots, \boldsymbol{h}_n]^T \in \mathbb{R}^{n \times s_e}$. To transform \boldsymbol{H} into label-specific vectors, we compute label-specific attention weights as:

$$A = \operatorname{softmax}(HU^T), A \in \mathbb{R}^{n \times |\mathcal{L}|}$$
 (3)

where $U \in \mathbb{R}^{|\mathcal{L}| \times s_e}$ is the label embedding matrix. The i^{th} column in A represents the attention weights corresponding to the i^{th} label in \mathcal{L} for each of the n tokens. To ensure the bulk of the weight is placed on the most informative tokens, the softmax is applied at the column level. Here A denotes the *learned* attention weights.

Second, we perform attention planting by utilizing two types of attention weights: staticplanted (S) and differentiable-planted (P). The static-planted attention (S) remains constant and is based on mutual information gain, while the differentiable-planted attention (P) comprises trainable parameters. These mechanisms enhance the model's ability to prioritize relevant tokens. We determine the static-planted attention as $\boldsymbol{S} = \left[\boldsymbol{g}^{(1)}, \boldsymbol{g}^{(2)}, \cdots, \boldsymbol{g}^{(|\mathcal{L}|)} \right] \in \mathbb{R}^{n \times |\mathcal{L}|},$ is comprised of individual vectors $g^{(i)} = [g_1^{(i)}, g_2^{(i)}, \cdots, g_n^{(i)}]^T \in \mathbb{R}^n$. Each element $g_j^{(i)}$ of these vectors represents the relevance of token t_i with respect to label $l^{(i)}$, as precisely defined in section 2.1. We determine the differentiableplanted attention by computing feature vectors $m{x}_{j}^{(i)} = \Psi\left(m{e}_{l^{(i)}},m{e}_{t_{j}}
ight)$ for each label-token pair $(l^{(i)}, t_j), i = 1, 2, \cdots, |\mathcal{L}|; j = 1, 2, \cdots, n$ as per equation 1. Then utilizing pretrained embeddings $e_{i(i)}$ and e_{t_i} from the L2R model in section 2.1, the pretrained L2R model computes scores $m{P} = \left[m{p}^{(1)}, m{p}^{(2)}, \cdots, m{p}^{(|\mathcal{L}|)}
ight] \in \mathbb{R}^{n imes |\mathcal{L}|},$ where $p^{(i)} = \left[f\left(x_1^{(i)} \right), \cdots, f\left(x_n^{(i)} \right) \right]^T \in \mathbb{R}^n$, and fis the ranking function from equation 7. In a departure from the standard attention approach, we introduce inattention, a pre-softmax thresholding technique that strategically elevates the significance of attention weights. By effectively zeroing out less relevant tokens, this method ensures maximal focus

³We used the pretrained LM from https://docs.fast. ai/text.models.awdlstm.html

304

311

314

315

318

319

321

322

323

324

325

on pivotal tokens:

$$\boldsymbol{P} = \mathsf{softmax}(\mathsf{threshold}(\boldsymbol{P}, k))$$
 (4)

where both threshold (Appendix A.3) and softmax are applied at the column level.

Third, to compute the label-specific vectors, we perform linear combinations of the hidden features 306 h_1, h_2, \cdots, h_n using the attention weights from three sources: the *learned* attention weights in each column of A, the *static-planted* attention weights in each column of S, and the *differentiable-planted* 310 attention weights in each column of P. This is followed by element-wise matrix multiplication 312 with a weight matrix $\boldsymbol{W} \in \mathbb{R}^{|\mathcal{L}| \times s_e}$:

$$\boldsymbol{V} = (\boldsymbol{A}^T \boldsymbol{H} + \boldsymbol{S}^T \boldsymbol{H} + \boldsymbol{P}^T \boldsymbol{H}) \odot \boldsymbol{W}, \boldsymbol{V} \in \mathbb{R}^{|\mathcal{L}| \times s_c}$$
(5)

The purpose of W is to boost attention. The i^{th} row v_i of V, can be thought of as the information regarding the i^{th} label captured by *attention* from the token sequence $\langle t_1, t_2, \cdots, t_n \rangle$. Finally, this labelspecific information is summed and added with a label-specific bias followed by sigmoid activation to produce predictions:

$$\hat{\boldsymbol{y}} = \mathsf{sigmoid}(\boldsymbol{1}\boldsymbol{V}^T + \boldsymbol{b}); \boldsymbol{1} \in \mathbb{R}^{s_e}; \boldsymbol{b}, \hat{\boldsymbol{y}} \in \mathbb{R}^{|\mathcal{L}|}$$
(6)

The training objective is to mimimize the binary cross-entropy loss between \hat{y} and the target y as:

$$\operatorname{Loss}(\boldsymbol{y}, \hat{\boldsymbol{y}}, \theta) = \sum_{i=1}^{|\mathcal{L}|} y_i \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i)$$

where θ denotes all trainable model parameters.

Inattention: In contrast to traditional attention mechanisms, we introduce *inattention* (Equation 4) a novel technique that strategically enhances attention weights' significance by filtering less relevant 331 tokens, ensuring focus on critical elements within a token sequence. Our ablation analysis consistently identifies the optimal threshold parameter k in the 333 range [1, 10k'], where k' is from the nDCG@k loss function (Equation 8) for L2R model pretraining (Section 2.1). This aligns with our motivation to use an L2R model to learn token ranks, concen-337 trating attention on informative tokens with higher ranks while reducing attention to less relevant to-339 kens.

Stateful Decoder: Our decoder innovates with a 341 stateful mechanism inspired by backpropagation through time (BPTT) (Howard and Ruder, 2018). 343

Segmentation into fixed-size batches preserves the *state*, consisting of the last hidden feature h_n and prediction \hat{y}_{b} for each batch. This state guides subsequent batches, allowing cumulative predictions through initializing the AWD-LSTM encoder with h_n and continuously adding predictions. Gradients propagate back across batches, improving adaptability to large documents and model convergence. Our stateful decoder eliminates the need for text truncation (Li and Yu, 2020; Xie et al., 2019), preventing performance loss, and ensures stable GPU memory usage by processing long texts in manageable batches.

344

345

346

348

349

350

351

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

381

382

384

385

388

389

390

391

392

Discriminative Fine-tuning and Gradual Unfreezing: To fine-tune our pretrained model effectively for attention planting, we employ two essential strategies. First, we leverage discriminative fine-tuning (Howard and Ruder, 2018). This technique assigns distinct learning rates (LR) to different parameter groups $\theta^l \in \left\{ \theta^e, \theta^p, \theta^d \right\}$ corresponding to AWD-LSTM encoder, planted decoder, and the remaining model components. This approach optimizes the pretrained model by focusing on areas that need the most adjustment. The update rule for discriminative fine-tuning is as follows:

$$\theta_t^l = \theta_{t-1}^l - \eta^l \cdot \nabla_{\theta^l} J(\theta)$$

where $\nabla_{\theta^l} J(\theta)$ is the gradient with respect to the model's loss function. For experiments we applied half or a third of the LR for already proficient L2R model parameters compared to others. Second, we embrace gradual unfreezing (Howard and Ruder, 2018). This method fine-tunes the model in a layerwise sequence, starting from the last layer and moving gradually towards the initial layers.

Bidirectional Language Model: For the MIMIC-III-full and MIMIC-IV-full (Table 2), we pretrain both a forward and backward LM. We fine-tune an XMTC model for each LM independently and average the classifier predictions. On MIMIC-III-full P@15 increased from 60.61 to 61.67, and on MIMIC-IV-full, from 54.5 to 55.6.

3 **Experiments**

Experimental Setup 3.1

Datasets: We compare PLANT to SOTA ICD coding models (Yang et al., 2022; Zhang et al., 2022; Yuan et al., 2022; Liu et al., 2021; Vu et al., 2021; Li and Yu, 2020; Cao et al., 2020; Xie et al., 2019; Mullenbach et al., 2018). Our primary datasets are MIMIC-III (Johnson et al., 2016)

and the newly available MIMIC-IV (Johnson et al., 2023). These datasets contain rich textual and structured records from ICU settings, with a focus on discharge summaries. These summaries are meticulously annotated with ICD-9 codes (MIMIC-III) and ICD-10 codes (MIMIC-IV) to represent diagnoses and procedures. MIMIC-III comprises 52,722 discharge summaries and 8,929 unique ICD-400 9 codes. We follow the methodology in (Mul-401 lenbach et al., 2018), including patient ID-based 402 splits for full-code experiments and a subset of 403 50 frequent codes. We also evaluate our model 404 on the few-shot MIMIC-III-rare50 dataset (Yang 405 et al., 2022), featuring 50 rare ICD codes. Addi-406 tionally, we explore MIMIC-IV, with 122,279 dis-407 charge summaries and 7,942 unique ICD-10 codes, 408 following Edin et al. (2023). We denote these 409 datasets as MIMIC-III-full, MIMIC-III-top50, 410 MIMIC-III-rare50, and MIMIC-IV-full. Refer 411 to Table 2 for dataset statistics. 412

MIMIC-III-full	MIMIC-IV-full
52,723	122,279
41,126	65,659
8,929	7,942
14(10 - 20)	14(9 - 20)
1,375(965 - 1,900)	1,492(1,147-1,931)
90.5/3.1/6.4	72.9/10.9/16.2
	MIMIC-III-full 52,723 41,126 8,929 14(10 - 20) 1,375(965 - 1,900) 90.5/3.1/6.4

Table 2: Descriptive statistics for MIMIC-III-full and MIMIC-IV-full discharge summary training sets.

413 Preprocessing, Implementation and Hyperpa414 rameters: We direct readers to Appendix A.4 and 415 Appendix A.5 for specifications.

Evaluation metrics: To compare with prior ICD 416 studies, we use various metrics, focusing on mi-417 cro and macro F1 scores, AUC, and P@k. Micro-418 419 averaging treats each (text, code) pair individually, while macro-averaging computes metrics per label. 420 Micro-R reflects the ratio of true positives to the 421 sum of true positives and false negatives for each 422 label, while Macro-R represents the average recall 423 across all labels. Precision follows a similar calcu-424 lation pattern. Macro-averaged metrics prioritize 425 infrequent labels. P@k denotes the proportion of 426 the k top-scored labels that match the ground truth. 427 Baselines: These included models such as CAML 428 (Mullenbach et al., 2018), MSATT-KG (Xie et al., 429 2019), MUltiResCNN (Li and Yu, 2020), Hyper-430 Core (Cao et al., 2020), LAAT/JointLAAT (Vu 431 et al., 2021), ISD (Zhou et al., 2021), Effective-432 CAN (Liu et al., 2021), MSMN (Yuan et al., 2022), 433 DiscNet (Zhang et al., 2022), and KEPTLong-434 former (Yang et al., 2022). 435

3.2 Main Results

MIMIC-III-full (Table 3): PLANT demonstrated 437 notable enhancements over existing SOTA mod-438 els. Specifically, when compared with Effective-439 CAN, LAAT, and DiscNet, PLANT yielded su-440 perior performance in terms of micro-F1, P@5, 441 P@8, and P@15 metrics, with improvements of 442 0.5%, 2.7%, 0.6%, and 0.3%, respectively. Signifi-443 cantly, PLANT achieved a remarkable P@5 score 444 of 84%, indicative of an average of 4.2 correct pre-445 dictions among the top 5; while demonstrating only 446 a slightly lower micro-AUCthan DiscNet.

Model	AUC		F1		P@k		
Model	Macro	Micro	Macro	Micro	P@5	P@8	P@15
CAML/DR-CAML	89.7	98.6	8.8	53.9	-	70.9	56.1
MSATT-KG	91.0	99.2	9.0	55.3	-	72.8	58.1
MultiResCNN	91.0	98.6	8.5	55.2	-	73.4	58.4
HyperCore	93.0	98.9	9.0	55.1	-	72.2	57.9
LAAT/JointLAAT	92.1	98.8	10.7	57.5	81.3	73.8	59.1
ISD	93.8	99.0	11.9	55.9	-	74.5	-
Effective-CAN	92.1	98.9	10.6	58.9	-	75.8	60.6
MSMN	95.0	99.2	10.3	58.4	-	75.2	59.9
DiscNet	95.6	99.3	14.0	58.8	-	76.5	61.4
PLANT (Ours)	90.4	98.9	10.1	59.4^{*}	84.0*	77.1^{*}	61.7^{*}

Table 3: Results (in %) on the MIMIC-III-full test set. We ran our model 5 times each with different random seeds for initialization and report mean scores. * indicates that the performance difference between PLANT and the next best is significant (p < 0.01, using the Approximate Randomization test). All scores in tables 3, 4, 5 and 6 are reported under the same experimental setup.

MIMIC-III-top50 (Table 4): PLANT outperformed the previous SOTA baseline models of MSMN and LAAT with regard to macro-F1, micro-F1, P@8 and P@15, respectively; while matching micro-AUCwith ISD and achieving a slightly lower P@5 as compared to MSMN. PLANT produced improvements of 0.4%, 0.3%, 0.3% and 1.4% for macro-F1, micro-F1, P@8 and P@15, respectively.

Madal	AUC		F1		P@k		
Model	Macro	Micro	Macro	Micro	P@5	P@8	P@15
CAML/DR-CAML	88.4	91.6	57.6	63.3	61.8	-	-
MSATT-KG	91.4	93.6	63.8	68.4	64.4	-	-
MultiResCNN	89.9	92.8	60.6	67.0	64.1	-	-
HyperCore	89.5	92.9	60.9	66.3	63.2	-	-
LAAT/JointLAAT	92.5	94.6	66.6	71.6	67.5	54.7	35.7
ISD	93.5	94.9	67.9	71.7	68.2	-	-
Effective-CAN	92.0	94.5	66.8	71.7	66.4	-	-
MSMN	92.8	94.7	68.3	72.5	68.0	-	-
PLANT (Ours)	93.1	94.9	68.7	72.8	67.2	55.0^{*}	36.3^{*}

Table 4: Results on the MIMIC-III-top50 test set.

MIMIC-III-rare50 (Table 5): PLANT surpassed the prior SOTA baseline, KEPTLongformer, by astounding margins. Specifically, by 12.9% in macro-AUC, 12.7% in micro-AUC, 52.2% in

455

456

436

Model	AL	JC	F1		
WIGGET	Macro	Micro	Macro	Micro	
MSMN	75.3	76.2	17.1	17.2	
KEPTLongformer	82.7	83.3	30.4	32.6	
PLANT (Ours)	95.6^{*}	96 .0*	82.6^{*}	84.2^{*}	

Table 5: Results on the MIMIC-III-rare50 test set.

macro-F1, and 51.6% in micro-F1. Intriguingly, it's worth noting that these remarkable results were achieved by training with only unfrozen PLANT layers, without even utilizing the entire model's capacity. This underscores the extraordinary potential of PLANT in delivering outstanding performance with efficient training strategies in low-shot settings.

> MIMIC-IV-full (Table 6): PLANT outperformed previous SOTA baseline model of LAAT with regard to P@8 and P@15 while matching micro-AUC with LAAT. PLANT produced improvements of 1.7%, 1.3% for P@8 and P@15, respectively.

Model	AUC		F1		P@k		
Woder	Macro	Micro	Macro	Micro	P@5	P@8	P@15
CAML/DR-CAML	91.1	98.5	16.0	55.4	-	66.8	52.2
MultiResCNN	94.5	99.0	21.1	56.9	-	67.8	53.5
LAAT/JointLAAT	95.4	99.0	20.3	57.9	-	68.9	54.3
PLANT (Ours)	94.8	99.0	19.6	57.1	78.1^{*}	70.6^{*}	55.6^{*}

Table 6: Results on the MIMIC-IV-full test set. The comparitive results are reported from Edin et al. (2023).

4 Analysis

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

Firstly, except for the Gradual Unfreezing and Bidirectionality, we selectively unfreeze the layers in decoder, keeping the encoder frozen—meaning no backpropagation was performed on their weights during training. This ensures that performance improvements are attributed directly to the decoder, our primary focus. Secondly, all reported performance metrics stem from the full test sets of both MIMIC-III-full and MIMIC-IV-full datasets. Thirdly, reported enhancements were statistically significant (p < 0.01, using the Approximate Randomization test).

Impact of PLANT (Figure 3,4): We evaluate 487 PLANT and LAAT (Vu et al., 2021) in contexts 488 with skewed label distributions. PLANT uses pre-489 trained L2R activations P and mutual information gain S, initializing the decoder's attention weights. 491 While LAAT relies solely on learned attention A, 492 initialized randomly and learned from scratch. That 493 is LANT omits P and S form Equation 5. Our 494 analysis involves training both PLANT and LAAT 495

models across varying fractions of a balanced train-496 ing dataset, with both models trained for up to five 497 epochs. The test set remains constant, and we mea-498 sure P@5 and P@15 as the performance metric for 499 both models. The results were notable: the PLANT 500 model consistently matched or surpassed the LAAT 501 model's performance across all training sizes, even 502 with significantly less data. For instance, in the case of MIMIC-IV-full, PLANT achieved a P@5 of 504 0.50 and P@15 of 0.37 with a smaller training split 505 of 1090 and 2743 instances, respectively, matching the performance of the LAAT model trained on a 507 significantly larger split of 10, 337 and 12, 902 in-508 stances. Similarly, in the case of MIMIC-III-full, PLANT achieved a P@5 of 0.47 and P@15 of 0.30, 510 trained with only 136 and 235 instances, respec-511 tively. This performance equates to that of the 512 LAAT model trained on a dataset comprising 1342 513 and 1578 instances. These findings are visually 514 represented in Figure 3 through vertical and hori-515 zontal lines, illustrating the substantial efficiency 516 gains of PLANT in terms of training data require-517 ments while maintaining or improving model per-518 formance. Since PLANT achieves comparable 519 performance to LAAT with significantly less data, 520 which also implies a lower number of instances per label (aka skewed label distribution), this out-522 come underscores the inefficiencies of the LAAT 523 approach in such scenarios. To examine overfit-524 ting (Figure 4), we trained both PLANT and LAAT 525 on MIMIC-IV-full for 60 epochs. While PLANT 526 remained stable, LAAT began overfitting after 40 527 epochs, diverging train and test loss, leading to a 528 decline in P@15.



Figure 3: P@15 for PLANT vs. LAAT (Vu et al., 2021) with different number of training examples on MIMIC-III-full and MIMIC-IV-full.



Figure 4: PLANT does not overfit on MIMIC-IV-full, LAAT (Vu et al., 2021) does.

Ablation	MIMIC-III-full	MIMIC-IV-full
Without Inattention	50.95	42.40
With Inattention	51.05	42.51
Stateless	52.80	43.38
Stateful	52.90	44.22
- disc	51.40	43.29
+ disc	52.21	44.34
full unfreezing	57.78	49.78
gradual unfreezing	58.31	50.97

P@15 for MIMIC-III-full Table 7: and MIMIC-IV-full (train split 49, 579) test set.

531

532

533

534

536

537

538

541

543

545

546

547

549

550

551

552

554

Impact of Inattention (Table 7): We investigated the impact of the inattention threshold k (Equation 4) within PLANT on MIMIC-III-full and MIMIC-IV-full. The training splits comprised 22,525 instances (average 49 instances per label) and 49,579 instances (average 97 instances per label) for the respective datasets. We trained each model for 5 epoch and measured P@15. For MIMIC-III-full, the model without inattention (k = 72) achieved a P@15 of 50.95, while the model with inattention (k = 56) achieved a slightly higher P@15 of 51.05. In the case of MIMIC-IV-full, the model without inattention attained a P@15 of 42.4, which improved to 42.51 with inattention (k = 8).

Impact of Sateful Decoder (Table 7): On the MIMIC-III-full training dataset, using the stateful decoder for three epochs yielded a P@15 of 52.9, a slight improvement over 52.8 without it. 548 Similarly, on the MIMIC-IV-full (training split of 49, 579), employing the stateful decoder for seven epochs significantly boosted P@15, from 43.28 to 44.22. These improvements highlight the stateful decoder's role in enhancing PLANT's performance with extensive text data.

Impact of Discriminative Fine-tuning and 555 Gradual Unfreezing (GU) (Table 7): On the MIMIC-III-full, training PLANT for one epoch with discriminative fine-tuning, applying half the learning rate to L2R parameters, improved P@15 from 51.40 to 52.21 on the test set. Similarly, on 560 MIMIC-IV-full (training split of 49, 579), training 561 PLANT for seven epochs with a third of the learning rate for L2R parameters increased P@15 from 563 43.29 to 44.34. For GU we explored two scenarios: one gradually unfreezing the model layer by layer, 565 and the other unfreezing the entire model simultane-566 ously. Both models were trained for 10 epochs. On 567 the MIMIC-III-full, GU increased P@15 from 57.78 to 58.31; and on MIMIC-IV-full from 49.78

518.81: Acute respiratory failure
PLANT:patient had a gcs3t and required intubationfio2 ··· temp po2 pco2 ph
MSATT-KG: left hemothorax, ETOH, depression, stable discharge condition
CAML:small apical pneumothorax remained unchanged now tolerating a
Text-CNN:revealed a persistent left pleural effusion and due to concern for loculated hemothorax
530.81: Esophageal reflux
PLANT: gastroesophageal reflux home o2 gerd osteoporosis one puff hospital1 prilosec 20mg
MSATT-KG: tracheostomy & feeding gastrostomy GERD, anxiety
CAML: rib fx requiring tracheostomy & feeding gastrostomy,, GERD, anxiety, cataracts
Text-CNN: right thoracotomy, decortication of lung, mobilization of liver off of chest wall

Table 8: Interpretability evaluation results for different models.

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

598

599

600

601

602

603

604

605

606

607

608

609

610

to 50.97.

Interpretability Case Study (Table 8): We compare PLANT's interpretability against three baselines: MSATT-KG, CAML, and Text-CNN(Kim, 2014). While PLANT selects top 5 tokens per label based on attention values, baseline methods extract informative *n*-grams. MSATT-KG employs multiscale and label-dependent attention, while CAML and Text-CNN use label-dependent attention and different phrase selection strategies. CAML uses a receptive field, and Text-CNN selects positions based on maximum channel values. In the interpretability case study, PLANT attends to tokens like 'intubation', 'fio2', and 'pc02'. 'fio2' represents Fraction of Inspired Oxygen, critical in determining oxygen concentration delivered to a patient. 'PCO2' signifies partial pressure of carbon dioxide, indicative of conditions like respiratory acidosis or alkalosis. In another example, informative tokens include 'gastrophageal', 'reflux', 'gerd', and 'prilosec', where 'gerd' denotes Gastroesophageal Reflux Disease and 'prilosec' is a proton pump inhibitor.

Related Work: Automatic ICD Coding 5

Xie and Xing (2018) introduced LSTM with tree structures and adversarial learning, Prakash et al. (2017) utilized condensed memory neural networks on MIMIC-III (Johnson et al., 2016). Baumel et al. (2017) proposed a hierarchical GRU network. Further enhancements include Xie et al. (2019)'s convolutions and multi-scale feature attention, Li and Yu (2020)'s convolutional layers, and Cao et al. (2020)'s graph convolution and hyperbolic representation. Vu et al. (2021) introduced LSTMbased attention models, Zhou et al. (2021) proposed shared representation networks, Liu et al. (2021) improved convolutional networks, and Yuan et al. (2022) introduced multi-synonyms attention networks. Zhang et al. (2022) addressed discourse structure and code-description reconciliation, including physician informal abbreviations.

611 Limitations

The PLANT method, while effective, presents a notable trade-off in terms of computational resources. 613 The necessity to pretrain and load the L2R model 614 imposes a substantial memory overhead compared 615 to traditional attention mechanisms. Consequently, 616 our memory constraints limited the number of 617 epochs for which PLANT could be trained. This as-618 pect of PLANT, particularly its scalability to larger XMTC datasets, warrants further investigation. Future work will explore strategies to optimize mem-621 ory usage, ensuring that the benefits of PLANT can be harnessed more broadly without the current 623 limitations on training duration and dataset size.

5 Broader Impacts and Ethical 6 Considerations

Our research contributes to the broader field of nat-627 ural language processing (NLP) and machine learn-628 ing (ML), advancing the SOTA in XMTC. In the context of XMTC, our research has the potential to significantly impact various sectors, including 631 healthcare, finance, and e-commerce. By automating labor-intensive tasks such as medical coding and diagnosis, these models can enhance healthcare accessibility, particularly in underserved communities. This can lead to improved patient outcomes and reduced disparities in healthcare access. Ad-637 ditionally, in education, XMTC models can support personalized learning experiences by categorizing educational resources and recommending tailored learning materials to students. Further-641 more, XMTC can contribute to policy development 642 by analyzing public opinion and sentiment from social media and news sources, providing valuable insights for policymakers to develop evidencebased policies and interventions. These applications demonstrate the diverse and far-reaching societal implications of XMTC technology. How-648 ever, we acknowledge the importance of ensuring that automated systems do not perpetuate biases or discrimination present in the data. Therefore, we prioritize fairness, transparency, and accountability in our model development process. In summary, while our research presents exciting opportunities for automation and efficiency gains, we recognize the importance of ethical considerations and broader societal impacts. By upholding ethical 657 principles and promoting responsible AI develop-658 ment, we aim to maximize the positive impact of our work while mitigating potential risks.

References

Tal Baumel, Jumana Nassour-Kassis, Raphael Cohen, 662 Michael Elhadad, and Noémie Elhadad. 2017. Multi-663 label classification of patient notes a case study on icd 664 code assignment. arXiv preprint arXiv:1709.09587. 665 Alex Bottle and Paul Aylin. 2008. Intelligent infor-666 mation: a national system for monitoring clinical 667 performance. Health services research, 43(1p1):10-668 31. 669 Christopher JC Burges. 2010. From ranknet to lamb-670 darank to lambdamart: An overview. Learning, 671 11(23-581):81. 672

661

673

674

675

676

677

678

679

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

- Pengfei Cao, Yubo Chen, Kang Liu, Jun Zhao, Shengping Liu, and Weifeng Chong. 2020. HyperCore: Hyperbolic and co-graph representation for automatic ICD coding. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3105–3114, Online. Association for Computational Linguistics.
- Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. 2018. Universal transformers. *arXiv preprint arXiv:1807.03819*.
- Joakim Edin, Alexander Junge, Jakob D Havtorn, Lasse Borgholt, Maria Maistro, Tuukka Ruotsalo, and Lars Maaløe. 2023. Automated medical coding on mimiciii and mimic-iv: A critical review and replicability study. *arXiv preprint arXiv:2304.10909*.
- Jeremy Howard and Sebastian Ruder. 2018. Universal language model fine-tuning for text classification. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 328–339, Melbourne, Australia. Association for Computational Linguistics.
- Alistair EW Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, Tom J Pollard, Sicheng Hao, Benjamin Moody, Brian Gow, et al. 2023. Mimic-iv, a freely accessible electronic health record dataset. *Scientific data*, 10(1):1.
- Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9.
- Douwe Kiela, Changhan Wang, and Kyunghyun Cho. 2018. Dynamic meta-embeddings for improved sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1466–1477, Brussels, Belgium. Association for Computational Linguistics.
- Yoon Kim. 2014. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*.
- Fei Li and Hong Yu. 2020. Icd coding from clinical text using multi-filter residual convolutional neural

822

823

824

825

826

827

828

715

716

717

- 739 740 741 742 743 744
- 745 747 750 751 754 755 756 757 758
- 761

765 766 767

771

network. In proceedings of the AAAI conference on artificial intelligence, volume 34, pages 8180-8187.

- Yang Liu, Hua Cheng, Russell Klopfer, Matthew R. Gormley, and Thomas Schaaf. 2021. Effective convolutional attention network for multi-label clinical document classification. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 5941-5953, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Stephen Merity, Nitish Shirish Keskar, and Richard Socher. 2017. Regularizing and optimizing lstm language models. arXiv preprint arXiv:1708.02182.
- James Mullenbach, Sarah Wiegreffe, Jon Duke, Jimeng Sun, and Jacob Eisenstein. 2018. Explainable prediction of medical codes from clinical text. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 1101–1111, New Orleans, Louisiana. Association for Computational Linguistics.
- Anthony N Nguyen, Donna Truran, Madonna Kemp, Bevan Koopman, David Conlan, John O'Dwyer, Ming Zhang, Sarvnaz Karimi, Hamed Hassanzadeh, Michael J Lawley, et al. 2018. Computer-assisted diagnostic coding: effectiveness of an nlp-based approach using snomed ct to icd-10 mappings. In AMIA Annual Symposium Proceedings, volume 2018, page 807. American Medical Informatics Association.
- Kimberly J O'malley, Karon F Cook, Matt D Price, Kimberly Raiford Wildes, John F Hurdle, and Carol M Ashton. 2005. Measuring diagnoses: Icd code accuracy. Health services research, 40(5p2):1620–1639.
- Aaditya Prakash, Siyuan Zhao, Sadid Hasan, Vivek Datla, Kathy Lee, Ashequl Qadir, Joey Liu, and Oladimeji Farri. 2017. Condensed memory networks for clinical diagnostic inferencing. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 31.
- Gongbo Tang, Mathias Müller, Annette Rios, and Rico Sennrich. 2018. Why self-attention? a targeted evaluation of neural machine translation architectures. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 4263-4272, Brussels, Belgium. Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. Advances in neural information processing systems, 30.
- Thanh Vu, Dat Quoc Nguyen, and Anthony Nguyen. 2021. A label attention model for icd coding from clinical text. In Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI'20.

- Kai Wang, Weizhou Shen, Yunyi Yang, Xiaojun Quan, and Rui Wang. 2020. Relational graph attention network for aspect-based sentiment analysis. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 3229– 3238, Online. Association for Computational Linguistics.
- Pengtao Xie and Eric Xing. 2018. A neural architecture for automated icd coding. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1066-1076.
- Xiancheng Xie, Yun Xiong, Philip S. Yu, and Yangyong Zhu. 2019. Ehr coding with multi-scale feature attention and structured knowledge graph propagation. In Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM '19, page 649-658, New York, NY, USA. Association for Computing Machinery.
- Song Xu, Haoran Li, Peng Yuan, Youzheng Wu, Xiaodong He, and Bowen Zhou. 2020. Self-attention guided copy mechanism for abstractive summarization. In Proceedings of the 58th annual meeting of the association for computational linguistics, pages 1355-1362.
- Zhichao Yang, Shufan Wang, Bhanu Pratap Singh Rawat, Avijit Mitra, and Hong Yu. 2022. Knowledge injected prompt based fine-tuning for multi-label fewshot icd coding. In Proceedings of the Conference on Empirical Methods in Natural Language Processing. Conference on Empirical Methods in Natural Language Processing, volume 2022, page 1767. NIH Public Access.
- Zheng Yuan, Chuanqi Tan, and Songfang Huang. 2022. Code synonyms do matter: Multiple synonyms matching network for automatic ICD coding. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 808-814, Dublin, Ireland. Association for Computational Linguistics.
- Shurui Zhang, Bozheng Zhang, Fuxin Zhang, Bo Sang, and Wanchun Yang. 2022. Automatic ICD coding exploiting discourse structure and reconciled code embeddings. In Proceedings of the 29th International Conference on Computational Linguistics, pages 2883-2891, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Tong Zhou, Pengfei Cao, Yubo Chen, Kang Liu, Jun Zhao, Kun Niu, Weifeng Chong, and Shengping Liu. 2021. Automatic ICD coding via interactive shared representation networks with self-distillation mechanism. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 5948-5957, Online. Association for Computational Linguistics.

A Appendix

A.1 Skewness of Codes



Figure 5: The skewness of ICD-9-CM code distribution for MIMIC-III (Johnson et al., 2016).

A.2 L2R Model (continued from Section 2.1)

The ranking function, $f : x_j^{(i)} \mapsto \mathbb{R}$, of the L2R model is an *L* layered feed forward network,

$$f(\boldsymbol{x_j^{(i)}}) = y^L, y^{(l)} = a(W^{(l)} \cdot y^{(l-1)} + b^{(l)}),$$
(7)

where $y^{(l)}$ is layer l output, $y^{(0)} = x$ is input, $W^{(l)}$ is layer l weight matrix, $b^{(l)}$ is layer l bias vector, and $a(\cdot)$ is the activation function. In our experiments we trained the L2R model with L = 2.

At any point in the training, the model outputs the score $\boldsymbol{z}^{(i)} = \left[f\left(\boldsymbol{x}_{1}^{(i)}\right), \cdots, f\left(\boldsymbol{x}_{n}^{(i)}\right) \right]^{T} \in \mathbb{R}^{n}$. The objective of the L2R model is to minimize the total loss,

$$\sum_{i=1}^{m} \mathsf{nDCG@k}\left(\boldsymbol{z}^{(i)}, \boldsymbol{g}^{(i)}\right), \tag{8}$$

where nDCG@k is the maximum allowable DCG@k, which is defined as:

$$\mathsf{DCG}@\mathsf{k}\left(\boldsymbol{z^{(i)}},\boldsymbol{g^{(i)}}\right) := \sum_{l \in \mathsf{rank}_k(\boldsymbol{z^{(i)}})} \frac{2^{\boldsymbol{g_l^{(i)}}}}{\log(l+1)}.^4$$

Bootstrapping L2R Model: Let (I, J) be a pair of random variables for the label $l^{(i)}$ and token t_j over the space $\mathcal{I} \times \mathcal{J}$, where $\mathcal{I} = \{\text{label } i \text{ present}, \text{label } i \text{ not present}\}$ and $\mathcal{J} = \{\text{token } j \text{ present}, \text{token } j \text{ not present}\}$. Then, g_j^i is defined as the mutual information gain of I and J:

$$g_j^{(i)} = \sum_{x \in \mathcal{I}, y \in \mathcal{J}} P_{(I,J)}(x,y) \log\left(\frac{P_{(I,J)}(x,y)}{P_I(x)P_J(y)}\right),$$

where $P_{(I,J)}$ is the joint, and P_I , P_J are the marginal probability mass function of I and J, respectively.

854

855

856

857

858

859

860

861

862

863

868

869

870

871

872

873

874

875

876

877

878

879

880

882

883

884

885

887

888

889

890

891

892

893

894

Training L2R Model: Gradient update rule to train the L2R model on $\left\{ \left(\boldsymbol{X}^{(i)}, \boldsymbol{g}^{(i)} \right) \right\}_{i=1}^{m}$ are defined as follows. Let $I^{(i)}$ denote the set of pairs of token indices $\{j, k\}$, such that $g_{j}^{(i)} > g_{k}^{(i)}$. Also, let $z_{j}^{(i)} = f\left(\boldsymbol{x}_{j}^{(i)} \right)$ and $z_{k}^{(i)} = f\left(\boldsymbol{x}_{k}^{(i)} \right)$. The parameters of L2R model, $w_{p} \in \mathbb{R}$, are updated as (Burges, 2010):

(:)

$$\delta w_p = -\eta \sum_j \lambda_j \frac{\partial z_j^{(i)}}{\partial w_k},$$
86

$$\lambda_j = \sum_{k:\{j,k\}\in I^{(i)}} \lambda_{jk} - \sum_{k:\{k,j\}\in I^{(i)}} \lambda_{kj},$$
869

$$\lambda_{jk} = -\frac{\sigma}{1 + e^{\sigma\left(z_j^{(i)} - z_k^{(i)}\right)}} |\Delta \mathsf{nDCG@k}|_{jk}, \tag{86}$$

where $|\Delta nDCG@k|_{jk}$ denotes the change in nDCG@k by swapping j and k in $rank(z^{(i)})$.

A.3 Threshold

threshold(
$$p^i, k$$
) =

$$\begin{cases} p_j, & \text{if } p_j > k^{th} \text{ largest } p \\ 0 & \text{otherwise.} \end{cases}$$

A.4 Preprocessing

Following prior research (Mullenbach et al., 2018; Xie et al., 2019; Li and Yu, 2020), we tokenize and lowercase all text while eliminating non-alphabetic tokens containing numbers or punctuation. A distinctive feature of our approach is the absence of preprocessed word embeddings. Instead, we finetune a pretrained AWD-LSTM model on our target dataset, allowing for parameter refinement, including word embeddings, and the generation of context-specific embeddings for new words in the dataset. While the concept of fine-tuning pretrained models is not new (Howard and Ruder, 2018), our innovation lies in its application to the XMTC domain. Contrary to previous practices (Li and Yu, 2020), we refrain from truncating text, as our experiments and findings align with those of Zhang et al. (2022), which demonstates substantial performance variation due to truncation. To handle longer texts, we employ our stateful decoder (refer to Section 2.2).

A.5 Implementation and Hyperparameters

We ensure robustness across diverse XMTC datasets by fine-tuning hyperparameters on the

- 832 833
- 83
- 00
- 840
- 841
- 842
- .

84

847

849

⁴here rank_k($z^{(i)}$) returns the k largest indices of $g^{(i)}$ ranked in descending order.

895	MIMIC-III-full and MIMIC-IV-full validation
896	sets. Experiments are conducted on an NVIDIA
897	QUADRO RTX 8000 GPU with 48 GB VRAM.
898	We utilize the AWD-LSTM LM with an embedding
899	size of 400, 3 LSTM layers with 1152 hidden ac-
900	tivations, and the Adam Optimizer with $\beta_1 = 0.9$,
901	$\beta_2 = 0.99$, and weight decay of 0.01. During
902	fine-tuning, we apply dropout rates and weight
903	dropout, with a batch size of 384, BPTT of 80,
904	20 epochs, and a learning rate of $1e - 5$. Classi-
905	fier training also includes dropout rates and weight
906	dropout, with a batch size of 16, BPTT of 72, and
907	discriminative fine-tuning with gradual unfreezing
908	over $115~{\rm epochs}$ (on MIMIC-III-full), alongside
909	scheduled weight decay and learning rate ranges.