

Automatic ICD-10 Code Association: A Challenging Task on French Clinical Texts

Yakini Tchouka, Jean-François Couchot, David Laiymani

Université de Franche-Comté
CNRS, Institut FEMTO-ST
F-25000 Besançon, France
firstname.name@femto-st.fr

Philippe Selles, Azzedine Rahmani

Nord Franche-Comté Hospital
F-90400 Trevenans, France
firstname.name@hnfc.fr

Abstract—Automatically associating ICD codes with electronic health data is a well-known NLP task in medical research. NLP has evolved significantly in recent years with the emergence of pre-trained language models based on Transformers architecture, mainly in the English language. This paper adapts these models to automatically associate the ICD codes. Several neural network architectures have been experimented with to address the challenges of dealing with a large set of both input tokens and labels to be guessed. In this paper, we propose a model that combines the latest advances in NLP and multi-label classification for ICD-10 code association. Fair experiments on a Clinical dataset in the French language show that our approach increases the F_1 -score metric by more than 55% compared to state-of-the-art results.

Index Terms—natural language processing, icd-10, clinical document, unstructured data, multi-label classification, supervised learning, health, transformers

I. INTRODUCTION

For a more accurate long term follow-up, patient’s stay in a health center is usually reported in a digital documents which constitute the patient’s medical record. Written by the patient’s physicians, it is composed of operating reports, clinical notes, liaison letters etc. In a large number of countries each patient record is then classified according to the International Classification of Diseases (ICD). ICD is a medical classification system used for coding diseases and other health-related conditions. It is maintained by the World Health Organization (WHO) and is widely used globally as a standard for tracking health statistics, billing for medical services, and conducting research. In its 10th edition [21], ICD is organized into chapters based on different body systems and disease categories. The chapters are further divided into subcategories that provide more specific information about the condition being coded. Each code consists of an alphanumeric string that includes a category code, a subcategory code, and a descriptor code (up to seven characters). The ICD-10 classification system is used by healthcare providers and organizations worldwide to standardize the coding of medical conditions and facilitate the sharing of health information across different systems and platforms. This classification is the common foundation for epidemiology, public health research, and healthcare management. In addition, the reimbursement of medical expenses by public or private insurance companies directly depends on the codes associated with these medical records. This makes it

even more important to associate the right codes with each patient’s record. Finally, it should be noted that a more or less complex patient record can generate several ICD-10 codes

Typically, in a hospital, the responsibility for the ICD-10 classification falls on the medical coders. Staff performing this task is specially trained professionals who use medical documentation to assign the appropriate ICD-10 codes to medical records. Medical coders work closely with healthcare providers, nurses, and other staff to ensure that the medical records are accurately encoded into this classification. In some hospitals the ICD-10 classification is performed by the physicians. However, regardless of how medical coding is managed, accuracy, and attention to detail are crucial to ensure that the data generated is reliable and useful for patient care and management. That is why automatically associating ICD codes to a medical record is a task that has been widely addressed in medical research in recent years [6], [2], [28], [9], [11].

With the recent advances in Natural Language Processing (NLP) and since medical records are unstructured medical documents, it makes sense to apply these theoretical and technological advances in the context of ICD-10 classification. Clearly, the emergence of the Transformers architecture [27], [10] has taken natural language processing to a new precision level. Several works have shown that the representations produced by these models are the most accurate and it is the most used architecture today in a large number of machine learning tasks (from text, to computer vision and time series) [27], [13], [16].

Nevertheless, ICD-10 automatic classification is a multi-label text classification task with tough challenges. For instance, the ICD-10 classification consists of about 140,000 codes (procedure codes and medical codes). Unless one has a huge dataset, extremely important physical resources, and an extremely long period of time, it seems to be unrealistic to believe that one could associate to a patient record one of the 140,000 existing codes with a high degree of accuracy.

This large number of labels clearly stresses existing deep learning models to their limits. Another challenge is the size of the medical notes which far exceeds the usual limit of transformer architectures (typically 512 tokens). Finally, working on non-English data is also challenging since the

vast majority of open-source models available are trained on English corpus.

In this paper, we propose to address the three previous challenges for the ICD-10 classification of French medical records. We developed a deep learning model that combines the latest advances in Natural Language Processing. This approach makes it possible to associate a non-negligible part of the existing ICD-10 codes on French-language patient records with an F_1 -score outperforming with more than 55% latest state of the art approach.

This paper is organized as follows. Section II starts with recalling state of the art of associating ICD codes to medical records. Section III presents the dataset used on the one hand to validate our approach and on the other hand to fairly compare the F_1 -scores obtained by our approach with those obtained with already existing approaches. The architecture of our ICD code association model is presented in Section IV. Results are presented and analyzed in Section V. Concluding remarks and future work are finally given in Section VI.

II. RELATED WORK

A. Natural Language Processing

NLP has significantly evolved in recent years with the joint appearance of the Transformers model [27] and their generalization ability to transfer learning. ELMo [23] and BERT [10] have shown this effectiveness which provides more accurate contextualized representations. Several pre-trained models then appeared such as BERT, RoBERTa [18] These models are pre-trained on a large amount of general domain English text to capture the ability to model text data, and then refined on common classification tasks. In French two main models have been proposed i.e FlauBERT [14], CamemBERT [19]. Note that some multi-lingual models also exist such as XLM-R [7]. Some models are also trained on domain-specific text corpus. For example, ClinicalBERT[1] and BioBERT[15] have been trained on medical data to address medical domain tasks. Unfortunately, there is no such model in the French language, leading to a gap between the usage of machine learning approaches on French documents compared to the same approach in English ones. In general, Transformers models have a limited input size (512 tokens in practice). In the case of clinical documents this limit can become very penalizing since a typical patient document is generally much larger than 512 words or tokens. In [22] the authors proposed some hierarchical methods to tackle this problem. They divided the document into several segments that can be processed by a Transformers. Then the encoding of the segments is aggregated into the next layer (Linear, recurrent neural networks or other layer of Transformers). Recently, the sparse-attention system i.e. the *LongFormer* model has been proposed in [3]. It is composed of a local attention (attention between a window of neighbour tokens) and a global attention that reduces the computational complexity of the model. They can therefore be deployed to process up to 4096 tokens.

B. ICD Code Association

The automatic association of ICD codes is one of the most addressed challenges in medical research. With the emergence of neural networks and the evolution of natural language processing, several authors have tried to tackle this task. [6] and [2] used recurrent neural networks (RNNs) to encode Electronic Health Records (EHR) and predict diagnostic outcomes. On the other hand, [25] and [20] have used the attention mechanism with RNNs and CNNs to implement more accurate models.

The work of [29] and [26] present various ways to consider the hierarchical structure of codes. [29] used a sequence tree LSTM to capture the hierarchical relationship between codes and the semantics of each code. [5] proposed to train the integration of ICD codes in a hyperbolic space to model the code hierarchy. They used a graph neural network to capture code co-occurrences. LAAT [28] integrated a bidirectional LSTM with an attention mechanism that incorporates labels.

EffectiveCAN [17] used a squeeze-and-excitation network and residual connections as well as extraction of representations from all layers of the encoder for label attention. The authors also introduced focal loss to address the problem of long-tail prediction with 58.9% of F_1 -score on MIMIC 3 [12]. ISD [30] used shared representation extraction between high frequency layers and low frequency layers and a self-distillation learning mechanism to mitigate the distribution of long-tailed codes.

Recently [11] proposed the PLM-ICD system that focuses on document encoding with multi-label classification. They used an encoding model based on the Transformers architecture adapted to the medical corpus. Associating ICD-10 codes is finding the codes corresponding to medical documents in a large set of codes. For instance MIMIC 3 [12] contains more than 8,000 codes, and handling such large set of labels in classification is a challenging problem in machine learning. To overcome this problem, the authors used the Label-Aware Attention (LAAT) mechanism proposed in [28] which integrates labels in the encoding of documents. Finally, to solve the problem of long sequences they used the hierarchical method. PLM-ICD is the current state-of-the-art model that achieved 59.8% of F_1 -score on MIMIC 3 [12] and 50.4% on MIMIC 2 [24].

In French, [9] proposed Convolutional Neural Networks (CNN) models with multi-label classification to automatically associate ICD-10 code. The authors used FastText [4] vectors with the skip-Gram algorithm for the encoding of documents. They first considered all the final labels of the dataset, then grouped them into families to reduce the number of classes. This model is trained on a private dataset of 28,000 clinical documents and reached 39% of F_1 -score with 6,116 codes and 52% with 1,549 codes.

III. DATASET

This work is in collaboration with The Hopital Nord Franche-Comté (HNFC), a French public health center that provided us with patient stays. For privacy reasons, all our

experiments were conducted on site and no data was taken out of the hospital.

A patient’s stay is a set of successive visits in possibly different departments of the hospital. Each department produces a clinical document that describes the patient’s stay in that department. These clinical documents are used by the medical coding specialists to associate the corresponding ICD-10 codes. We finally obtain a set of unstructured textual documents corresponding to the global stay of the patient to which a set of codes is associated. As clinical documents, we have for example operating reports, discharge letters, external reports or clinical notes. The obtained dataset, further denoted as ICD-10-HNFC dataset is therefore a database of groups of medical documents with associated codes. This system is well illustrated in Fig. 1.

ICD-10-HNFC dataset is built for supervised deep learning. In supervised learning, to have an accurate model, there are several factors to consider. Is there enough training data? Is the number of classes consistent with the volume of data available? Is the frequency of classes in the dataset balanced? It is always difficult to find the perfect dataset that answers all these questions. In this paper, we worked not only on the main dataset, which consists in associating the raw ICD codes it contains but also on the sub-datasets such as associating the most frequent codes or code families instead of the raw codes.

Class Reduction

As mentioned, the ICD is a classification system that is composed of thousands of codes. Given the large number of labels present in our basic dataset (shown in Table I), it is difficult to approach this classification task by considering all the classes present. By doing so, the results of the constructed model will be far from perfect. The most precise models to date in English for ICD-10 code association is PLM-ICD which reached 59.8% on MIMIC 3 with 8,922 labels [12] and 50.4% on MIMIC 2 with 5,031 labels [24]. This proves the difficulty of this task. The first sub-dataset consists in reducing the codes to the first 3 characters seen as a family. Therefore, instead of considering the raw codes, we will group them into families. This reduces in a consequent way the number of classes to be treated by the model. This dataset is presented in Table I. We can see via the description "line Code with less than 10 examples" in Table I that the reduction of the classes not only allows to have a more reasonable number of classes but also increase the frequency of the codes in the dataset.

Code Frequency

Associating ICD-10 codes is a very frequent task in health centers. As a result, some codes occur more frequently than others. Finding the most frequent codes automatically can only be useful. Our second sub-dataset consists in building models based on the number of codes (K) that we consider more relevant. We evaluate the relevance based on the frequency of the code in the dataset. Thus, a model built on such a dataset will be able to associate the integrated codes with a better classification performance.

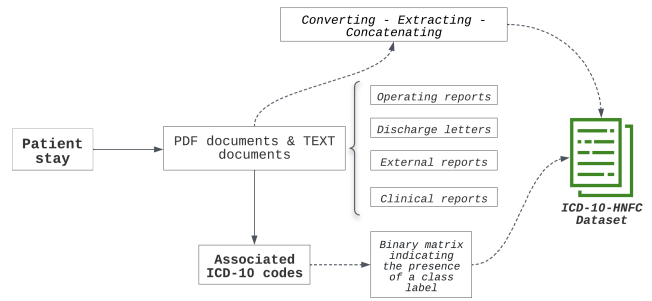


Fig. 1. ICD-10-HNFC Dataset Construction

Additional Label

With the code frequency strategy (K), the dataset is therefore composed of entries whose association belongs only to the K most relevant codes. To keep the coherence of our dataset an additional label is introduced to represent the codes which are not considered as relevant (least frequent codes). So instead of having K classes, the model will have $K + 1$ ones. The additional class mean concretely in the association that there are one or more additional codes to associate.

TABLE I
DESCRIPTIVE STATISTICS OF ICD-10-HNFC DATASET

	Dataset	Dataset with class reduction
Documents	56014	-
Tokens	41868993	-
Average sequence length	747	-
Total ICD codes	416125	415830
Unique ICD codes	6160	1564
Codes with less than 10 examples	3722	523
Codes with 100 examples or more	641	471

IV. MODEL ARCHITECTURE

This section presents the different components of the model architecture we have developed and justifies the choices made to design it. As previously exposed, as we deal with the French language our choice was to fine-tune pre-trained transformer-based French models i.e. CamemBERT [19] and FlauBERT [14] for the implementation of the model architecture.

A. Global Document Representation

As mentioned, Transformers main constraint is the limitation of the number of tokens present in an input sequence. Since the average size of the clinical notes of ICD-10-HNFC dataset exceeds this limit (747 versus 512 as shown in Table I), basic Transformers cannot be used. Recently, [8] summarized the available methods for processing long sequences via Transformers. They can be summarized as hierarchical Transformers and sparse-attention Transformers in which we can find the *Longformer* model of [3] early mentioned. *Longformer* can process up to 4096 tokens per sequence allows to meet this limit. Unfortunately, there is no French pre-trained *Longformer*

model to date. Therefore, in this paper, we will use the hierarchical method to tackle this problem.

Hierarchical Transformers[22], [8] are built on top of Transformers architecture. A document D , is first divided into segments $[t_0, t_1, \dots, t_{|D|}]$, each of which must have less than 512 tokens (the limit of Transformers). These segments are encoded independently using a typically pre-trained Transformers. We then obtain a list of segment representations which must be aggregated to obtain the whole document D representation. There are several ways to do this aggregation. The aggregator can be an average of the representations of all the segments of the document (mean pooling) or the maximum of the representations in each dimension of the segments (max pooling) or stacking the segment representations into a single sequence. The aggregated sequence thus serves as an input to the next layer.

1) *Classification of a Large Number of Labels:* To overcome the problem of a large set of labels since ICD-10-HNFC contains more than 6,000 codes, we used the Label-Aware Attention (LAAT) system as in [11]. LAAT consists in integrating the labels into the document representation. Label-Aware Attention captures important text fragments related to certain labels. Let H be the stacking representation of an input sequence. First, a label-wise attention weight matrix Z is computed as follows:

$$Z = \tanh(VH)$$

$$A = \text{softmax}(WZ)$$

where V and W are linear transforms. The i^{th} row of A represents the weights of the i^{th} label. The softmax function is performed for each label to form a distribution over all tokens. Then, the matrix A is used to perform a weighted-sum of H to compute the label-specific document representation:

$$D = HA^T$$

The i^{th} row of D represents the document representations for the i^{th} label. Finally, D is used to make predictions by computing the inner product between each row of D and the related label vector.

In this paper, several architectures were experimented such as the model without long sequence processing, the model with long sequence processing (max/mean pooling), and the model with LAAT. The global architecture is illustrated in Fig. 2.

V. EXPERIMENTS AND ANALYSIS

In this section, we present the results of the experiments conducted with the previously detailed architectures and dataset. We compare the results of recent works (PLM-ICD[11], CNN[9]) on the association of ICD-10 codes with

the most precise model of this paper. To evaluate our model we use the most used performance measures in classification: Precision, Recall, F_1 -score. The micro average system is used to obtain the aggregation of the performances.

TABLE II
ICD-10 ASSOCIATION RESULTS OF THE DIFFERENT ARCHITECTURES ON THE VALIDATION ICD-10-HNFC DATASET

Models	Labels	Precision	Recall	F_1 -score
FlauBERT (512 tokens)	1564	0.48	0.31	0.38
Hierarchical Mean FlauBERT		0.54	0.39	0.45
Hierarchical Max FlauBERT		0.53	0.40	0.46
FlauBERT + LAAT		0.57	0.51	0.54
CamemBERT + LAAT		0.56	0.53	0.55
FlauBERT + LAAT	6160	0.41	0.43	0.42
CamemBERT + LAAT		0.52	0.4	0.45

A. Paper Models

First, we conducted the experiments on the ICD-10-HNFC dataset with class reduction (1564 labels) as detailed in Table I. This experimentation is performed with all the architectures developed in this paper. They are listed in Table II. Then we trained another model on the global ICD-10-HNFC dataset (6101 labels) with the architectures that obtained the highest F_1 -score in the previous experiment. The results are shown in the Table II. The results confirm the effects of the different components that constitute our architectures. In summary, the LAAT approach outperforms the hierarchical methods which are better than the base truncated model.

B. K-based Models

As detailed in Section III, different models have been trained based on a number (K) of labels (i.e. the most frequent codes). We present here the evaluation of these models with K in [10, 50, 100, 200]. As shown in Table III, models are less and less accurate when we increase the number of labels (classes). This is simply due to the aggregation of performances. The more different codes there are, the fewer instances of each code there are in the dataset, and the less easy the contextualization is.

TABLE III
RESULTS OF THE ICD-10 ASSOCIATION OF K-BASED MODELS

K	Precision	Recall	F_1 -score
10	84	80.5	82.1
50	78.2	65.1	71
100	77.2	58.4	66.5
200	71.9	52.6	60.8

C. Comparison with other Works

The Table IV shows the model with the highest F_1 -score of this paper with the results of previous work on ICD-10 code association. It is difficult to compare the results, since these works do not use the same evaluation dataset and English works can benefit from specialized models such as ClinicalBERT [1]. For French baseline, we implemented and trained the model proposed in [9] on ICD-10-HNFC dataset. The result is shown in parallel with our proposal. Our model clearly outperforms the classification method used in [9]. On the same validation dataset, with class reduction (1564 labels) the F_1 -score goes from 0.35 obtained with the model proposed in [9] to 0.55 with our proposal, i.e. an improvement of 57%.

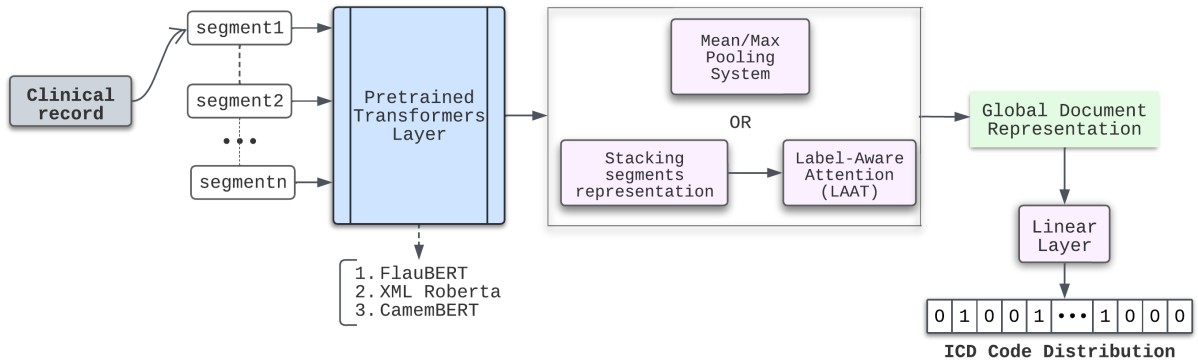


Fig. 2. Global Architecture

With the raw codes (6161 labels), the F_1 -score goes from 0.27 to 0.45, i.e. an improvement of 66.6%. The difference in scores with the results of PLM-ICD can be explained by the use of a context specific (medical) Transformers which has a vocabulary more adapted to the content of the documents.

TABLE IV

RESULTS COMPARISON WITH THE PREVIOUS WORK ON ICD-10 ASSOCIATION. THE STATE OF THE ART WORKS WITH THEIR RESULTS ARE IN *italic*. THE EXPERIMENTS DONE IN THIS PAPER WITH ICD-10-HNFC DATASET ARE PRESENTED IN THE OTHER PART. THE HIGHEST SCORES IN EACH PART IN RELATION TO THE NUMBER OF LABELS ARE MARKED IN **BOLD**

Models	Language	Dataset	Labels	F_1 -score
<i>PLM-ICD[11]</i>	<i>English</i>	<i>MIMIC 2[24]</i> <i>MIMIC 3[12]</i>	5,031 8,922	0.5 0.59
[9]	<i>French</i>	[9]	6,116 1,549	0.39 0.52
PROPOSAL	French	ICD-10-HNFC	6,161 1,564	0.45 0.55
[9]			6,161 1,564	0.27 0.35

VI. CONCLUSION

In this paper, we address the challenges of automatically associating ICD-10 codes to French clinical unstructured data. We have experimented several Transformers architectures to address the challenges of large input tokens and large numbers of labels. We therefore propose an ICD-10 association model that uses the latest advances in natural language processing and achieves the highest results in the French language to date. Our future work will focus on the use of Large Language Models and few-shots learning techniques to the ICD-10 classification.

ACKNOWLEDGMENT

This work is (partially) supported by the EIPHI Graduate School (contract ANR-17-EURE-0002).

REFERENCES

- [1] Emily Alsentzer, John R Murphy, Willie Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, and Matthew McDermott. Publicly available clinical bert embeddings. *arXiv preprint arXiv:1904.03323*, 2019.
- [2] Tal Baumel, Jumana Nassour-Kassis, Raphael Cohen, Michael Elhadad, and Noémie Elhadad. Multi-label classification of patient notes: case study on icd code assignment. In *Workshops at the thirty-second AAAI conference on artificial intelligence*, 2018.
- [3] Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- [4] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomáš Mikolov. Enriching word vectors with subword information. *CoRR*, abs/1607.04606, 2016.
- [5] Pengfei Cao, Yubo Chen, Kang Liu, Jun Zhao, Shengping Liu, and Weifeng Chong. 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, 2020.
- [6] Edward Choi, Mohammad Taha Bahadori, Andy Schuetz, Walter F Stewart, and Jimeng Sun. Doctor ai: Predicting clinical events via recurrent neural networks. In *Machine learning for healthcare conference*, pages 301–318. PMLR, 2016.
- [7] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. *CoRR*, abs/1911.02116, 2019.
- [8] Xiang Dai, Ilias Chalkidis, Sune Darkner, and Desmond Elliott. Revisiting transformer-based models for long document classification. *arXiv preprint arXiv:2204.06683*, 2022.
- [9] Clément Dalloux, Vincent Claveau, Marc Cuggia, Guillaume Bouzillé, and Natalia Grabar. Supervised learning for the icd-10 coding of french clinical narratives. In *MIE 2020-Medical Informatics Europe conference-Digital Personalized Health and Medicine*, pages 1–5, 2020.
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [11] Chao-Wei Huang, Shang-Chi Tsai, and Yun-Nung Chen. Plm-icd: automatic icd coding with pretrained language models. *arXiv preprint arXiv:2207.05289*, 2022.
- [12] 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. MIMIC-III, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016.
- [13] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *ACM computing surveys (CSUR)*, 54(10s):1–41, 2022.
- [14] Hang Le, Loïc Vial, Jibril Frej, Vincent Segonne, Maximin Coavoux, Benjamin Lecouteux, Alexandre Allauzen, Benoit Crabbé, Laurent Besacier, and Didier Schwab. Flaubert: Unsupervised language model pre-training for french. *arXiv preprint arXiv:1912.05372*, 2019.
- [15] Jinhuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240, 2020.
- [16] Tianyang Lin, Yuxin Wang, Xiangyang Liu, and Xipeng Qiu. A survey of transformers. *AI Open*, 2022.
- [17] Yang Liu, Hua Cheng, Russell Klopfer, Matthew R Gormley, and Thomas Schaaf. 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, 2021.
- [18] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [19] Louis Martin, Benjamin Muller, Pedro Javier Ortiz Suárez, Yoann Dupont, Laurent Romary, Éric Villemonte de La Clergerie, Djamel Seddah, and Benoît Sagot. Camembert: a tasty french language model. *arXiv preprint arXiv:1911.03894*, 2019.
- [20] James Mullenbach, Sarah Wiegrefe, Jon Duke, Jimeng Sun, and Jacob Eisenstein. Explainable prediction of medical codes from clinical text. *arXiv preprint arXiv:1802.05695*, 2018.
- [21] World Health Organization et al. Icd-10. international statistical classification of diseases and related health problems: Tenth revision 1992, volume 1= cim-10. classification statistique internationale des maladies et des problèmes de santé connexes: Dixième révision 1992, volume 1. 1992.
- [22] Raghavendra Pappagari, Piotr Zelasko, Jesús Villalba, Yishay Carmiel, and Najim Dehak. Hierarchical transformers for long document classification. In *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 838–844. IEEE, 2019.
- [23] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. 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 2227–2237, New Orleans, Louisiana, June 2018. Association for Computational Linguistics.
- [24] Mohammed Saeed, Mauricio Villarroel, Andrew T Reisner, Gari Clifford, Li-Wei Lehman, George Moody, Thomas Heldt, Tin H Kyaw, Benjamin Moody, and Roger G Mark. Multiparameter intelligent monitoring in intensive care ii (mimic-ii): a public-access intensive care unit database. *Critical care medicine*, 39(5):952, 2011.
- [25] Haoran Shi, Pengtao Xie, Zhiting Hu, Ming Zhang, and Eric P Xing. Towards automated icd coding using deep learning. *arXiv preprint arXiv:1711.04075*, 2017.
- [26] Shang-Chi Tsai, Ting-Yun Chang, and Yun-Nung Chen. Leveraging hierarchical category knowledge for data-imbalanced multi-label diagnostic text understanding. In *Proceedings of the Tenth International Workshop on Health Text Mining and Information Analysis (LOUHI 2019)*, pages 39–43, 2019.
- [27] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [28] Thanh Vu, Dat Quoc Nguyen, and Anthony Nguyen. A label attention model for icd coding from clinical text. *arXiv preprint arXiv:2007.06351*, 2020.
- [29] Pengtao Xie and Eric Xing. 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, 2018.
- [30] Tong Zhou, Pengfei Cao, Yubo Chen, Kang Liu, Jun Zhao, Kun Niu, Weifeng Chong, and Shengping Liu. 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, 2021.