# ViMQ: A Vietnamese Medical Question Dataset for Healthcare Dialogue System Development

Ta Duc Huy[1], Nguyen Anh Tu[1], Tran Hoang Vu[1], Nguyen Phuc Minh[1], Nguyen Phan[1], Trung H. Bui, and Steven Q. H. Truong[1]

VinBrain
{v.huyta, v.tunguyen, v.vutran, v.minhng, v.nguyenphan,
v.brain01}@vinbrain.net
bhtrung@gmail.com

**Abstract.** Existing medical text datasets usually take the form of question and answer pairs that support the task of natural language generation, but lacking the composite annotations of the medical terms. In this study, we publish a Vietnamese dataset of medical questions from patients with sentence-level and entity-level annotations for the Intent Classification and Named Entity Recognition tasks. The tag sets for two tasks are in medical domain and can facilitate the development of task-oriented healthcare chatbots with better comprehension of queries from patients. We train baseline models for the two tasks and propose a simple self-supervised training strategy with span-noise modelling that substantially improves the performance. Dataset and code will be published at https://github.com/tadeephuy/ViMQ.

**Keywords:** NER · intent classification · medical question dataset · self-supervised · learning with noise.

## 1 Introduction

Named Entity Recognition (NER) involves extracting certain entities in a sentence and classifying them into a defined set of tags such as organizations, locations, dates, time, or person names. In the medical domain, NER tag sets usually contain patient information and clinical terms. Intent Classification (IC) is the task of categorizing a sequence into a set of defined intentions. In the medical domain, IC classes could include the conversation objectives in interactions between patients and clinical experts. NER and IC are two major components of a task-oriented dialogue system. There exist several medical conversation datasets in English [2], German [6], and Chinese [2,9], while Vietnamese medical conversation datasets are not abundant. COVID-19 NER for Vietnamese [7] is a dataset for the named entity recognition task with generalized entity types that can extend its application to other future epidemics. The dataset is crawled from Vietnamese online news websites with the keyword "COVID" and filtered out irrelevant sentences, resulting in 34,984 entities and 10,027 sentences.

To our knowledge, we are the first to publish a Vietnamese Medical Question dataset (ViMQ) that contains medical NER and IC tags set where the applications could be generalized to developing the Natural Language Understanding (NLU) module for healthcare chatbots.

NER annotation includes highlighting the span of a given entity and assigning an appropriate tag. Specifying the spans of the entities can raise consensus issues stemming from the subjectivity of different annotators while classifying entities into tags are far more compliant due to the distinctiveness of the tags. In this paper, we developed an annotation methodology that aims to minimize this effect. In addition, NER annotation tools with poor ergonomic design can contribute to the span-noise, where the start and end indexes of a given span are shifted from their correct ground truth. We propose a training strategy to learn the model with such noise. To summarize, our contributions are:

- We published a Vietnamese Medical Question dataset for healthcare chatbot development.
- We proposed a training strategy to learn the model with span-noise for the NER task.

## 2   ViMQ dataset

### 2.1   Intent and entity types

The ViMQ dataset contains Vietnamese medical questions crawled from the consultation section online between patients and doctors from `www.vinmec.com`, a website of a Vietnamese general hospital. Each consultation consists of a question regarding a specific health issue of a patient and a detailed respond provided by a clinical expert. The dataset contains health issues that fall into a wide range of categories including common illness, cardiology, hematology, cancer, pediatrics, etc. We removed sections where users ask about information of the hospital and selected 9,000 questions for the dataset. We annotated the questions for the NER and IC tasks. The tag sets for two tasks could be applied to a dialogue system for medical consultation scenario, where the chatbot acts as a medical assistant that inquires queries from users for their health issues. Each question is annotated for both tasks. The statistics of the labels are shown in Table 1.

*Entity definition:* The tag set for NER includes SYMPTOM&DISEASE, MEDICAL_PROCEDURES and MEDICINE:

- SYMPTOM&DISEASE: any symptom or disease that appears in the sentence, including disease that are mentioned as part of a medicine such as "rabies" in "rabies vaccine".
- MEDICAL_PROCEDURE: actions provided by clinical experts to address the health issue of the patient, including diagnosing, measuring, therapeutic or surgical procedures.
- MEDICINE: any medicine name that appears in the sentence.

Table 1: Statistics of ViMQ dataset.

| Entity Type | Train | Valid | Test | All |
|---|---|---|---|---|
| SYMPTOM&DISEASE | 10,599 | 1,300 | 1,354 | 13,253 |
| MEDICAL_PROCEDURE | 1,583 | 204 | 213 | 2,000 |
| MEDICINE | 781 | 90 | 108 | 979 |
| **Intent Type** | **Train** | **Valid** | **Test** | **All** |
| Diagnosis | 3,444 | 498 | 484 | 4,426 |
| Severity | 1,070 | 150 | 162 | 1,382 |
| Treatment | 1,998 | 264 | 265 | 2,527 |
| Cause | 488 | 88 | 89 | 625 |

*Intent definition:* The tag sets for IC includes Diagnosis, Severity, Treatment and Cause:

- Diagnosis: questions relating to identification of symptoms or diseases.
- Severity: questions relating to the conditions or grade of an illness.
- Treatment: questions relating to the medical procedures for an illness.
- Cause: questions relating to factors of a symptom or disease.



Fig. 1: Examples from ViMQ dataset. Each block includes the IC tag (right side) and the NER tags of an example from the dataset (above the horizontal line) and its English translation (below the horizontal line).

## 2.2   Annotation process

We propose a novel method to manage the annotation process for the two tasks of NER and IC called Hierarchical Supervisors Seeding (HSS). The method circumvents around building a solid and consensual supervisor team which grows in size as the annotation process progresses. Leveraging a solid supervisor team mitigates the subjectivity and improves the task-oriented competency of individual annotators. Each supervisor, upon designated, does not only make mutual concessions for label disagreements of individual annotators but also acts as a seed to make way for them to become the next generation supervisors.

We first ran a pilot phase where 1,000 samples are selected for annotation, called the pilot set. The author of the annotation guideline is designated as the first supervisor, where two other NLP engineers are individual annotators. The two individual annotators follow the annotation guideline and annotate two overlapping sets of 600 sentences in the pilot set. The intersection of the two sets consisting of 100 sentences is also annotated by the supervisor and is used to measure the annotation quality of the two NLP engineers using F1-score. The intersection set is selected such that it covers a broad range of difficulties. Hard sentences contain rare diseases such as anxiety disorders, thalassemia. Easy sentences contain common diseases or symptoms such as flu, fever, and cold. Both individual annotators are required to achieve an F1-score of at least 0.9 on both tasks on the intersection set before discussing with the supervisor on the disagreement cases. The supervisor has to point out which parts of the annotation guideline could solve these conflicts and update the guideline until they reach a consensus. This session encourages the individual annotators to gain better comprehension and help the supervisor to improve the guideline.

Two individual annotators in the pilot phase are designated as supervisors in the following phases. In the main phase, each supervisor monitors two other individual annotators for a set of the next 1,000 sentences. They follow the same procedure in the pilot phase where each of the individual annotators works with 600 sentences, and the conflicts in the intersection set of 100 sentences are solved by the supervisor of each group in the discussion sessions. All updates to the annotation guideline made by the supervisors are reviewed by the guideline author. We continue to designate the two individual annotators to be supervisors in the next phase in a hierarchical manner. We repeat the process until all sentences are annotated. The intersection sets in each group are aggregated and used as the golden test set. The individual annotators after the pilot phase are medical under-graduates, which are hired at a rate of 0.1 USD per annotated sentence and 0.05 USD per resolved conflict.

## 3   Baseline models

### 3.1   Intent classification

For task intent classification, we follow a common strategy when fine tuning pretrained BERT [1] for sequence classification task, we use PhoBERT [4] to

extract contextual features of sentences then take the hidden state of the first special token ([CLS]) denoted $h_1$, the intent is predicted as:

$$y^i = softmax(W^i h_1 + b^i) \qquad (1)$$

### 3.2 Named entity recognition

A standard approach for NER task is to formulate it as a sequence labeling problem with the Begin-Inside-Outside (BIO) tagging scheme. Similar to the intent classification task, we use PhoBERT to extract contextual embedding with a conditional random fields (CRF) [3] inference layer on top of our model to improve sequence labeling performance.

**Sub-word Alignment**: In our approach, because PhoBERT uses a Byte Pair Encoder (BPE) to segment the input sentence with sub-word units. A word is split into k sub-words. The original word label is assigned to the first sub-word and k-1 sub-words is marked with a special label (X). Final, the predicted label of original word based on the first label of sub-word, we keep a binary vector for active positions.

## 4    Method

As the annotators only agree on the intersection sets, the remaining sets, which were only labeled by single annotators, could be polluted with noise. We develop a training strategy to minimize the effect of remaining noisy labels in the training set. We empirically show that the method makes use of the potentially noisy samples and improves the performance substantially. We apply our training strategy on the ViMQ dataset and the COVID-19 NER Vietnamese dataset and achieve better performance using the standard settings.

Given a set of entities $\{E_i\}$ in a sentence, where $E_i$ is a tuple of $(s_i, e_i, c_i)$ where $s_i$ and $e_i$ is the start and end index of the named entity $i^{th}$ and $c_i$ is its category.

**Span-noise modelling**. We model the span-noise by adding $\delta$ to $s_i$ and $e_i$ with a probability of $p$ during training. Span-noise modeling acts as a regularization method and possibly corrects the noisy span indexes during training model.

**Online self-supervised training**. The training progresses through $N$ iterations, each consisting of $T$ epochs. After training for the first iteration, we start to aggregate the predictions made by the model in each epoch in the $j^{th}$ iteration. Entities with correct predictions for the entity category $c_i$ and have an $IoU > 0.4$ with the span ground truth $(s_i, e_i)$ are saved. We then employ major voting for the start/end indices of the span of each entity to combine the aggregated predictions of $T$ epochs in iteration $j^{th}$ and use the result as labels for training the model in the next iteration.

Table 2: Confusion matrix of the IC task on the ViMQ dataset.

|      | Sev. | Cau. | Tre. | Dia. |
|------|------|------|------|------|
| Sev. | 152  | 0    | 0    | 10   |
| Cau. | 0    | 85   | 0    | 4    |
| Tre. | 1    | 1    | 233  | 30   |
| Dia. | 19   | 2    | 27   | 436  |

Table 3: F1-score of IC task on ViMQ dataset.

|            | F1-Score(%) | |
|------------|-------------|-------------|
| **Baseline**  | ✓ | ✓ |
| **+ self-sup.** | - | ✓ |
| Dia. | $90.22 \pm 0.05$ | $\mathbf{90.55 \pm 0{,}09}$ |
| Sev. | $90.22 \pm 0.38$ | $\mathbf{91.02 \pm 0{,}00}$ |
| Tre. | $\mathbf{89.04 \pm 0{,}23}$ | $88.97 \pm 0.21$ |
| Cau. | $95.15 \pm 0.9$ | $\mathbf{96.05 \pm 0.00}$ |
| Mic.F1 | $90.36 \pm 0.05$ | $\mathbf{90.65 \pm 0.15}$ |
| Mac.F1 | $91.17 \pm 0.23$ | $\mathbf{91.65 \pm 0.08}$ |

## 5  Experiments

### 5.1  Experiments setups

We conduct experiments on our dataset and the COVID-19 NER dataset to compare the performances of the baseline models and the online self-supervised training strategy for the IC and NER tasks which are presented in Sections 3 and 4. It should be noted that we do not add noise in the IC task because span-noise is inapplicable. Our experimental models were implemented PyTorch [5] using Huggingface's Transformers [8], a huge library for pretrained Transformer-based models. We train both of the models for 5 iterations, each with 10 epochs. In the online self-supervised approach, the model starts using pseudo-labels from the second iteration. We set the noise injection probability $p$ to 0.1 and the noise shifting offset $\delta$ to 1 in all of our experiments. We employ AdamW optimizer with learning rate at $5e^{-5}$. In this study, we ran each experiments 5 times with different random seeds and report the average values with their standard deviations.

### 5.2  Results

**Intent classification.** Table 3 shows that the online self-supervised method improves baseline model from 90.36% to 90.65% Micro-F1 and from 91.17% to 91.65% Macro-F1. The confusion matrix in Table 2 shows that the IC model makes mistakes some of the sample having intent Diagnosis of intent Treatment and vice versa.

**Named entity recognition** Tables 4 and 5 show the experiments results on the COVID-19 NER and ViMQ datasets for the NER task. Online self-supervised training consistently improves the F1-score compares to the baseline model. The injection of span-noise adds significant gain to the performance on both datasets with a margin of 1% on COVID-19 NER and 3% on ViMQ In ViMQ dataset, there are 389 sentences with wrong predictions. Boundary errors exists

Table 4: F1-score of NER task on COVID-19 Vietnamese dataset.

|  | F1-Score(%) | | |
|---|---|---|---|
| **Baseline** | ✓ | ✓ | ✓ |
| **+ self-supervised** | - | ✓ | ✓ |
| **+ span-noise** | - | - | ✓ |
| Mic.F1 | 93.34 ± 0.04 | 94.34 ± 0.52 | **94.81 ± 0.33** |
| Mac.F1 | 91.15 ± 0.81 | 91.55 ± 0.81 | **92.49 ± 0.22** |

Table 5: F1-score of the NER task on the ViMQ dataset.

|  | F1-Score(%) | | |
|---|---|---|---|
| **Baseline** | ✓ | ✓ | ✓ |
| **+ self-supervised** | - | ✓ | ✓ |
| **+ span-noise** | - | - | ✓ |
| SYMP.&DIS. | 73.44 ± 0.05 | 73.60 ± 0.29 | **77,21 ± 0.09** |
| MEDICINE | 58.14 ± 2.47 | 63.22 ± 0.47 | **67,34 ± 0.51** |
| MED.PRO. | 58.51 ± 0.41 | 59.01 ± 0.78 | **61,73 ± 0.08** |
| Mic.F1 | 70.78 ± 0.20 | 71.22 ± 0.30 | **74,78 ± 0.05** |
| Mac.F1 | 63.36 ± 0.70 | 65.28 ± 0.51 | **68,76 ± 0.23** |

in 295/389 sentences that contains 362 entities with correct categories but incorrect span indices. Most of them are from class SYMPTOM&DISEASE as labels sometimes contain the inflicted body parts of the patient.

## 6 Discussion

ViMQ dataset can be utilized to develop a NLU module for healthcare chatbots. We propose a standard use case where it is applicable. A model trained for NER and IC tasks on the dataset decomposes a question from user to the system into named entities and intent components. These components are used to retrieve a corresponding medical answer from a pre-installed database. If no respond can be retrieved, the system routes the question to a recommended doctor.

## 7 Conclusion

In this work, we published a Vietnamese dataset of medical questions for the two tasks of intent classification and named entity recognition. Additionally, we proposed a training strategy to learn the model with span-noise modelling. The training strategy demonstrates positive gains on our dataset and COVID-19 Vietnamese NER dataset. Our dataset can be leveraged to develop a NLU module for healthcare chatbots to reduce workload of Telehealth doctors.

# References

1. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers). pp. 4171–4186. Association for Computational Linguistics, Minneapolis, Minnesota (Jun 2019). https://doi.org/10.18653/v1/N19-1423, `https://www.aclweb.org/anthology/N19-1423`
2. He, X., Chen, S., Ju, Z., Dong, X., Fang, H., Wang, S., Yang, Y., Zeng, J., Zhang, R., Zhang, R., Zhou, M., Zhu, P., Xie, P.: Meddialog: Two large-scale medical dialogue datasets (2020)
3. Lafferty, J.D., McCallum, A., Pereira, F.C.N.: Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In: Proceedings of the Eighteenth International Conference on Machine Learning. p. 282–289. ICML '01, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA (2001)
4. Nguyen, D.Q., Nguyen, A.T.: PhoBERT: Pre-trained language models for Vietnamese. In: Findings of the Association for Computational Linguistics: EMNLP 2020. pp. 1037–1042 (2020)
5. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., Garnett, R. (eds.) Advances in Neural Information Processing Systems. vol. 32. Curran Associates, Inc. (2019), `https://proceedings.neurips.cc/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf`
6. Rojowiec, R., Roth, B., Fink, M.: Intent recognition in doctor-patient interviews. In: Proceedings of the 12th Language Resources and Evaluation Conference. pp. 702–709. European Language Resources Association, Marseille, France (2020), `https://www.aclweb.org/anthology/2020.lrec-1.88`
7. Truong, T.H., Dao, M.H., Nguyen, D.Q.: COVID-19 Named Entity Recognition for Vietnamese. In: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (2021)
8. Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Scao, T.L., Gugger, S., Drame, M., Lhoest, Q., Rush, A.M.: Transformers: State-of-the-art natural language processing. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations. pp. 38–45. Association for Computational Linguistics, Online (2020), `https://www.aclweb.org/anthology/2020.emnlp-demos.6`
9. Zeng, G., Yang, W., Ju, Z., Yang, Y., Wang, S., Zhang, R., Zhou, M., Zeng, J., Dong, X., Zhang, R., Fang, H., Zhu, P., Chen, S., Xie, P.: MedDialog: Large-scale medical dialogue datasets. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing. pp. 9241–9250. Association for Computational Linguistics, Online (2020). https://doi.org/10.18653/v1/2020.emnlp-main.743, `https://www.aclweb.org/anthology/2020.emnlp-main.743`