

A Benchmark Dataset and Evaluation Framework for Vietnamese Large Language Models in Customer Support

Long S. T. Nguyen[✉], Truong P. Hua, Thanh M. Nguyen, Toan Q. Pham, Nam K. Ngo, An X. Nguyen, Nghi D. M. Pham, Nghia H. Nguyen, and Tho T. Quan^{*}^{1b}

¹ URA Research Group, Faculty of Computer Science and Engineering, Ho Chi Minh City University of Technology (HCMUT), Ho Chi Minh City, Vietnam

² Vietnam National University Ho Chi Minh City, Ho Chi Minh City, Vietnam

Abstract. With the rapid advancement of Artificial Intelligence, Large Language Models (LLMs) have become indispensable in Question Answering (QA) systems, enhancing response efficiency and reducing human workload, particularly in customer service. The rise of Vietnamese LLMs (ViLLMs) has positioned lightweight open-source models as the preferred choice due to their efficiency, accuracy, and privacy advantages. However, systematic evaluations of their performance in domain-specific contexts remain scarce, making it challenging for enterprises to identify the most suitable LLM for customer support applications, especially given the lack of benchmark datasets reflecting real-world customer interactions. To bridge this gap, we introduce Customer Support Conversations Dataset (CSConDa), a high-quality benchmark comprising over 9,000 QA pairs, meticulously curated from customer interactions with human advisors at a large-scale Vietnamese software company. Covering diverse service-related topics, including pricing inquiries, product availability, and technical troubleshooting, CSConDa serves as a representative dataset for evaluating ViLLMs in real-world scenarios. Furthermore, we present a comprehensive evaluation framework, benchmarking 11 lightweight open-source ViLLMs on CSConDa using not only well-suited automatic metrics but also an in-depth syntactic analysis to uncover their strengths, weaknesses, and underlying linguistic patterns. This analysis provides insights into model behavior, explains performance variations, and identifies critical areas for improvement, guiding future advancements in ViLLM development. Thus, by establishing a robust benchmark for LLM-driven customer service applications, our work provides a quantitative evaluation dataset and a comprehensive ViLLM performance comparison, offering key insights into intrinsic model performance, including accuracy, fluency, and consistency, while enabling informed decision-making for next-generation QA systems. Our dataset is publicly available on Hugging Face.

Keywords: Vietnamese LLMs · Customer Support QA Benchmark · Intrinsic Evaluation of LLMs

^{*} Corresponding author

1 Introduction

The rapid advancement of *Artificial Intelligence* (AI) has revolutionized *Question Answering* (QA) systems, enabling *Large Language Models* (LLMs) to automate information retrieval and generate responses with varying complexity [1]. As one of the fastest-growing economies, Vietnam is witnessing increasing demand for AI-driven solutions, particularly in customer service, where efficiency and accuracy are paramount. Consequently, many Vietnamese enterprises are integrating LLMs into their QA systems to streamline automated consultation.

Model size plays a crucial role in LLM performance, with larger models generally exhibiting superior linguistic capabilities [2]. While proprietary closed-source models, such as GPT-4³, dominate large-scale applications, open-source alternatives offer greater adaptability for domain-specific tasks, scalability, and enhanced data privacy. Our analysis of model usage statistics on *Hugging Face*⁴ (HF), a leading platform for hosting and fine-tuning open-source LLMs, reveals that as of January 2025, among the top 3,000 most frequently downloaded and widely adopted models, those within the 7–9 billion parameter range are the most preferred, as illustrated in Figure 1. This trend underscores the growing preference for *lightweight open-source models*, which balance computational efficiency and practical usability [3], making them particularly well-suited for enterprises developing LLM-based QA systems.

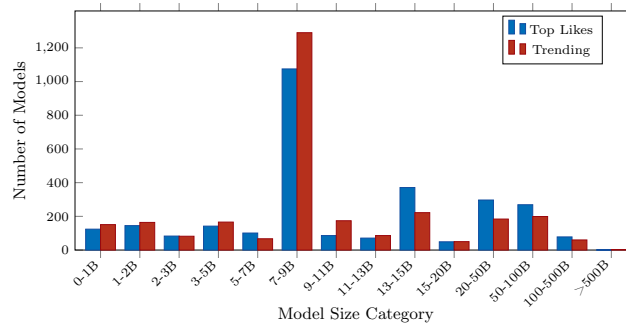


Fig. 1. Distribution of LLM model sizes based on popularity metrics on HF.

Despite the increasing availability of lightweight LLMs in general and *Vietnamese LLMs* (ViLLMs) in particular, enterprises developing LLM-based customer support systems face challenges in selecting the most suitable model for their needs. Therefore, a comprehensive evaluation of LLMs in domain-specific QA is crucial to ensuring their reliability, scalability, and hallucination tendencies before applying external knowledge augmentation or other enhancements. As a result, a benchmark dataset tailored for customer support is essential for

³ <https://openai.com/index/gpt-4/>

⁴ <https://huggingface.co/>

systematically comparing ViLLM performance. However, existing Vietnamese QA datasets [4,5,6,7] primarily focus on factual accuracy within structured contexts, often derived from Wikipedia⁵ articles, covering domains such as news, healthcare, and education. Thus, these datasets fail to reflect the nature of real-world customer interactions, where queries frequently include teencode, abbreviations, and domain-specific jargon. This issue is particularly pronounced in non-English languages like Vietnamese, often leading to misinterpretations, unnatural responses, or even complete failure to generate relevant answers. To the best of our knowledge, no existing dataset has been specifically designed for Vietnamese customer service, nor has any study systematically evaluated ViLLMs in specialized domains such as customer support.

To address this gap, we introduce *Customer Support Conversations Dataset* (CSConDa), a large-scale dataset of nearly 10,000 QA pairs extracted from real-world customer interactions with human advisors at DooPage⁶, a Vietnamese software company serving 30,000 customers and 45,000 advisors through its multi-channel support platform. Through rigorous processing and anonymization, CSConDa serves as a comprehensive benchmark for evaluating ViLLMs in customer support tasks. Our dataset covers diverse topics, including pricing, service inquiries, and technical troubleshooting, and is categorized into three types reflecting different levels of complexity. As a representative dataset, CSConDa enables a systematic assessment of lightweight ViLLMs’ intrinsic ability to generate human-like responses without external knowledge augmentation. Using CSConDa, we systematically evaluate 11 widely adopted lightweight ViLLMs within the 7–9 billion parameter range, employing six carefully selected automatic qualitative and quantitative metrics, along with syntactic analysis, to assess accuracy, hallucination tendencies, and fluency in handling practical customer queries. This evaluation highlights each model’s strengths and weaknesses, identifying key areas for improvement. Therefore, our work equips enterprises with a robust dataset that accurately represents real-world customer interactions, reducing the time required for model assessment while providing valuable insights and comparative analyses to support selecting the most suitable ViLLM for practical deployment. Our key contributions are summarized as follows.

- We conducted a comprehensive survey of publicly available Vietnamese QA datasets, providing an in-depth analysis of their linguistic characteristics.
- We introduced CSConDa, the first large-scale Vietnamese QA dataset derived from real-world customer service interactions, establishing a foundation for evaluating ViLLMs in customer support.
- We conducted a systematic evaluation of 11 widely adopted lightweight open-source ViLLMs, assessing their ability to generate human-like responses in domain-specific QA with selected automatic metrics and syntactic analysis.
- We provided unique and high-quality insights into both the strengths and limitations of existing ViLLMs, identifying key challenges in real-world de-

⁵ <https://www.wikipedia.org/>

⁶ <https://doopage.com/>

ployment and offering actionable recommendations to enhance LLM-based customer support systems.

2 Related Works

2.1 Human-Generated Vietnamese QA Datasets

The availability of text-based Vietnamese QA datasets remains highly limited, as Vietnamese is inherently a low-resource language. Existing datasets, such as [4,5], construct QA pairs from Wikipedia articles, covering various domains. Meanwhile, domain-specific QA datasets, such as [6,7], contain open-ended questions related to news, healthcare, and education, primarily compiled from online news sources and academic materials. Additionally, datasets like [8,9] focus on multiple-choice questions designed for Vietnamese educational assessments. While these datasets provide valuable linguistic resources, they are not well-suited for evaluating ViLLMs in customer support applications, where real-world queries are inherently unstructured and often incorporate teencode, abbreviations, and code-switching between Vietnamese and English. More critically, they lack the conversational spontaneity and domain-specific complexity found in real-world customer support interactions. To bridge this gap, we introduce the first large-scale dataset specifically designed for customer support QA, establishing a representative benchmark for evaluating ViLLMs in practical applications. Additionally, we conduct a comprehensive survey of publicly available text-based human-generated Vietnamese QA datasets, analyzing their statistical properties and linguistic structures in detail, as discussed in Section 3.3.

2.2 Comprehensive Evaluation of ViLLMs

Assessing ViLLMs in QA tasks presents significant challenges due to the scarcity of high-quality datasets, making domain-specific evaluation even more complex. As a result, few studies have systematically analyzed the performance of ViLLMs, particularly in specialized domains. To the best of our knowledge, the only existing study related to our work is [2], presented at the flagship conference *North American Chapter of the Association for Computational Linguistics* (NAACL 2024). This study provides a pioneering evaluation of ViLLMs across various NLP tasks using an extensive set of metrics, including QA, on widely adopted Vietnamese datasets. However, these datasets primarily consist of structured text and lack the spontaneity, informality, and domain-specific nuances inherent to customer support conversations. Furthermore, this study does not focus on domain-specific QA and evaluates a broad range of ViLLMs, including both open-source and closed-source models, without addressing their applicability to real-world customer service scenarios. In contrast, we present the first large-scale evaluation of 11 open-source ViLLMs, specifically targeting the most widely adopted model segment (7–9 billion parameters) and focusing on domain-specific QA in customer support applications. Our findings provide valuable insights into model capabilities, limitations, and areas for improvement, bridging the gap between academic benchmarks and real-world deployment.

3 The CSConDa Benchmark Dataset

3.1 Dataset Creation

CSConDa was constructed through five phases: (i) *worker recruitment*, (ii) *conversation collection*, (iii) *dataset creation*, (iv) *validation and categorization*, and (v) *dataset splitting*, as illustrated in Figure 2.

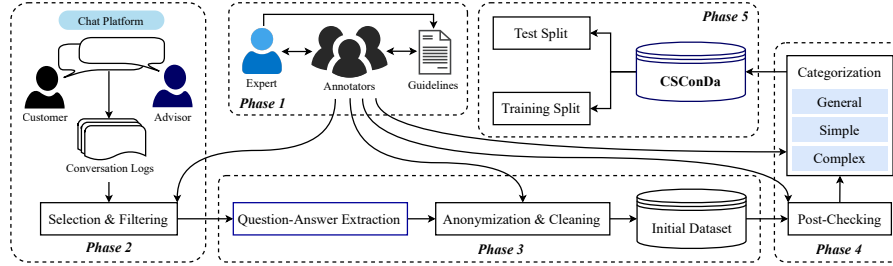


Fig. 2. The five-phase process of CSConDa dataset creation.

Phase 1: Worker Recruitment A team of ten *annotators*, supervised by a domain *expert* in customer service, was recruited. The expert, also a Chief Executive Officer in the business sector, established *guidelines* for conversation selection, anonymization, and dataset categorization. All annotators signed confidentiality agreements, underwent structured training, and maintained continuous communication with the expert.

Phase 2: Conversation Collection Raw *conversation logs* were sourced from a *chat platform*, where *customers* interact with *advisors*. Annotators followed predefined criteria to perform *selection and filtering*, ensuring high-quality exchanges. The filtering process prioritized coherence, topic diversity, and exclusion of sensitive or inappropriate content.

Phase 3: Dataset Creation After filtering, an automated pipeline was developed for *question-answer extraction*, preserving contextual consistency. The extracted QA pairs then underwent *anonymization and cleaning* to remove personally identifiable information and system-generated artifacts (e.g., syntax markers, reaction icons). The output of this phase formed the *initial dataset*.

Phase 4: Validation and Categorization A *post-checking* step ensured compliance with security standards, formatting consistency, and annotation accuracy. The *categorization* process divided the dataset into three types: *General*, *Simple*, and *Complex*, based on conversational complexity, reasoning demands, and required domain knowledge, as defined in the annotation guidelines.

Phase 5: Dataset Splitting After categorizing the dataset, we obtained *CSConDa*, which was then divided into standard *training* and *test* splits. The test split comprises 1,500 representative questions, evenly distributed across the three predefined types, serving as a benchmark for evaluating ViLLMs.

3.2 Dataset Overview

CSConDa is available on Hugging Face, with categorization criteria on the introduction page. Table 1 presents examples of its three types, illustrating linguistic variations in practical customer support conversations, including: (i) **code-switching**, frequently in domain-specific terms; (ii) **abbreviations**, encompassing teencode and informal shorthand; (iii) **acronyms**, formal abbreviations derived from initial letters; and (iv) **typos**, resulting from fast or imprecise typing.

Table 1. Overview of the three predefined types in CSConDa.

| |
|---|
| Type: General |
| Question: âu kê thank kiu e. (English: Okay, thank you.) |
| Answer: dạ vâng ạ. (English: Yes, sure.) |
| Type: Simple |
| Question: Báo giá giúp M nhé. Gọi M lúc 11h số này nef <số điện thoại>. C cần qly 8 page , 2 Zalo , 1 website , 1 YTB , 3-5 người dùng. (English: Please send me the pricing details. Call me at 11 AM at this number <phone number>. I need 8 pages, 2 Zalo accounts, 1 website, and 1 YouTube channel for 3-5 users.) |
| Answer: Em gửi chị chi phí cho nhu cầu của mình gồm 8 page , 2 Zalo OA , 1 web (tổng 13 kênh) và 3 hoặc 5 người dùng ạ. Đối với nền tảng YTB , hiện tại bên em chưa hỗ trợ. Chi phí trên là giá phần mềm trong 01 tháng. (English: I have sent you the pricing details, including 8 pages, 2 Zalo official accounts (OAs), 1 website (total 13 channels), and 3 or 5 users. Currently, we do not support YouTube. The above price covers 1 month of software usage.) |
| Type: Complex |
| Question: à thể đây là teen maps . chứ đâu phải tên business đâu em. tên 1 địa chỉ map đó. Nhưng Isao mà gõ tên map vào đó được. trong khi bên trong cho phép add nhiều locations ? Lỗi file này e ơi, ko down đc . Bên Zalo OA down bt . (English: Oh, so this is the name on Google Maps, not the business name. This is just a map location. But how can I enter a map name when multiple locations are allowed? This file is corrupted and cannot be downloaded. On Zalo OA, the download works fine.) |
| Answer: dạ phần tên này bên Google họ để trong API là tên liên quan Google Business ạ. Phần này chắc để em check lại thêm với dev ạ. File OA để em báo dev kiểm tra lại ạ. Anh ơi, anh có thể cho em xin link file Excel bị lỗi ở trên Zalo OA được không ạ? (English: Yes, Google categorizes this name under Google Business in their application programming interface (API). I will verify this with our development team. Regarding the file issue, I will report it to the developers. Could you send me the link to the corrupted Excel file via Zalo OA?) |

3.3 Dataset Analysis

Overall Statistics Table 2 presents key statistics on the structural composition of CSConDa, detailing the occurrence and frequency of acronyms, abbreviations, and typos within sentences. These features are first averaged at the sentence level, then across all records in each split to obtain the final values. Notably, customer questions exhibit a high prevalence of informal linguistic patterns, such as teencode and shorthand expressions, while answers are consistently longer. These trends reflect common characteristics of customer support interactions.

Table 2. Overall summary of CSConDa.

| | Training | Test | All |
|------------------------|----------|-------|--------|
| Number of samples | 8,349 | 1,500 | 9,849 |
| Question length | 16.31 | 19.69 | 16.82 |
| Answer length | 41.37 | 29.75 | 39.60 |
| Vocabulary size | 4,432 | 2,274 | 4,683 |
| Abbreviation count | 10,164 | 2,179 | 12,343 |
| Acronym count | 5,710 | 1,191 | 6,901 |
| Typos count | 1,537 | 270 | 1,807 |
| Abbreviation frequency | 0.10 | 0.09 | 0.10 |
| Acronym frequency | 0.06 | 0.05 | 0.07 |
| Typos frequency | 0.02 | 0.01 | 0.02 |

Type-based Statistics Table 3 presents a detailed breakdown of CSConDa statistics across different types in each split. A notable observation is that average question length increases progressively with type complexity. The test split, which serves as the benchmark dataset for evaluation, is evenly distributed across the three types, ensuring a diverse representation of linguistic characteristics.

Table 3. Detailed analysis of CSConDa across different types.

| | Training | | | Test | | | All | | |
|------------------------|----------|--------|---------|---------|--------|---------|---------|--------|---------|
| | General | Simple | Complex | General | Simple | Complex | General | Simple | Complex |
| Number of samples | 3,023 | 4,711 | 614 | 500 | 500 | 500 | 3,523 | 5,211 | 1,114 |
| Question length | 9.01 | 18.03 | 38.95 | 10.16 | 20.11 | 28.79 | 9.18 | 18.23 | 34.39 |
| Answer length | 43.24 | 39.21 | 48.70 | 14.41 | 29.34 | 45.50 | 39.16 | 38.27 | 47.26 |
| Vocabulary size | 2,497 | 3,678 | 2,072 | 964 | 1,446 | 1,663 | 2,624 | 3,782 | 2,445 |
| Abbreviation count | 2,208 | 6,758 | 1,198 | 434 | 852 | 893 | 2,642 | 7,610 | 2,091 |
| Acronym count | 1,339 | 3,739 | 632 | 263 | 456 | 472 | 1,602 | 4,195 | 1,104 |
| Typos count | 235 | 1,133 | 169 | 58 | 104 | 108 | 293 | 1,237 | 277 |
| Abbreviation frequency | 0.10 | 0.10 | 0.08 | 0.10 | 0.10 | 0.08 | 0.10 | 0.10 | 0.08 |
| Acronym frequency | 0.08 | 0.05 | 0.04 | 0.06 | 0.06 | 0.04 | 0.06 | 0.05 | 0.04 |
| Typos frequency | 0.01 | 0.02 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.12 | 0.01 |

Comparison with Existing Datasets We conduct a comprehensive comparison of publicly available Vietnamese QA datasets, focusing exclusively on text-based QA. Our analysis examines key characteristics, including the *number of samples* (NoS), *average question length* (AQL), *average answer length* (AAL), and the *presence of abbreviations* (Abb.), *acronyms* (Acr.), and *typos* (Typ.), as summarized in Table 4. Unlike prior datasets, which are primarily derived from structured or semi-structured sources such as Wikipedia, online articles, and educational materials, CSConDa is uniquely constructed from real-world human interactions. As a result, it is the only dataset that authentically captures the informal and context-dependent nature of customer support conversations. Another key distinction lies in the QA paradigm, while most existing Vietnamese QA datasets adopt an extractive approach, where answers are directly retrieved from a predefined context, CSConDa aligns more closely with conversational QA, where responses are inherently open-ended and require contextual reasoning.

4 ViLLMs Comprehensive Evaluation

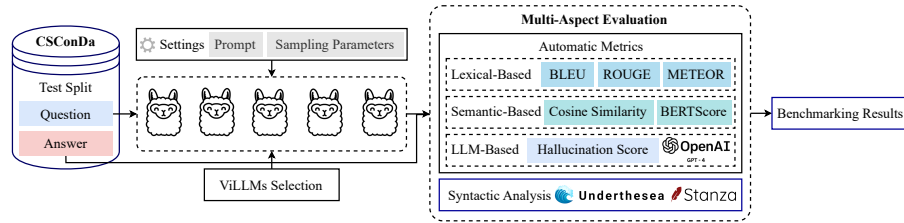
4.1 Proposed Multi-Aspect Evaluation Framework

Overall Framework We present a comprehensive framework for evaluating ViLLMs in domain-specific QA, particularly in customer support, as illustrated in Figure 3. Selected ViLLMs are assessed via inference using questions from the CSConDa test split as input prompts. Our evaluation considers accuracy, fluency, human-likeness, and hallucination tendencies, leveraging automatic metrics benchmarked against human-annotated reference answers. Beyond numerical assessment, we analyze the syntactic structure of ViLLM-generated responses for

Table 4. Comparison of CSConDA with existing Vietnamese QA datasets.

| Dataset | QA Task | Domain | NoS | AQL | AAL | Abb. | Acr. | Typ. | Source |
|-----------------------|------------------------------|-------------------------|--------------|-------------|-------------|------|------|------|------------------------------|
| CSConDA (Ours) | Conversational | Customer Support | 9,862 | 17.2 | 40.7 | ✓ | ✓ | ✓ | Human Conversations |
| UIT-ViQuAD [4] | Extractive | General | 23,074 | 12.2 | 8.2 | ✗ | ✓ | ✗ | Wikipedia |
| UIT-ViCoQA [7] | Conversational | General | 10,000 | 9.4 | 9.7 | ✗ | ✓ | ✗ | Online Healthcare Articles |
| UIT-ViWikiQA [10] | Extractive | General | 23,074 | 14.49 | 39.21 | ✗ | ✓ | ✗ | Wikipedia |
| VIMQA [5] | Extractive, Multi-Hop | General | 9,044 | 15.9 | 2.7 | ✗ | ✓ | ✗ | Wikipedia |
| VnYQA [11] | Extractive | General | 100 | 14.0 | 10.0 | ✗ | ✓ | ✗ | Wikipedia |
| UIT-ViMMRC [8] | Multiple-Choice | General | 2,783 | 12.5 | 7.5 | ✗ | ✓ | ✗ | Primary School Reading Texts |
| ViMedQA [12] | Extractive, Abstractive | Healthcare | 44,313 | 13.49 | 24.16 | ✗ | ✓ | ✗ | Online Medical Documents |
| UIT-ViNewsQA [6] | Extractive | Healthcare | 22,057 | 10.6 | 10.7 | ✗ | ✓ | ✗ | Online Healthcare Articles |
| UIT-ViCoV19QA [13] | Abstractive, Knowledge-Based | Healthcare | 4,500 | 31.79 | 119.76 | ✗ | ✓ | ✗ | Official Health Documents |
| ViHealthQA [14] | Abstractive, Knowledge-Based | Healthcare | 10,013 | 103.87 | 495.33 | ✗ | ✓ | ✗ | Online Healthcare Articles |
| VNHSGE [9] | Multiple-Choice | Education | 7,127 | 64.6 | 1.0 | ✗ | ✓ | ✗ | High School Graduation Tests |
| ViRHE4QA [15] | Extractive, Abstractive | Education | 9,758 | 17.09 | 24.18 | ✗ | ✓ | ✗ | University Regulations |
| [16] | Knowledge-Based | Education | 985 | 74.3 | 415.6 | ✗ | ✓ | ✗ | University Regulations |
| [17] | Knowledge-Based | Legal | 5,992 | 17.3 | 1.58 | ✗ | ✓ | ✗ | Civil Law |
| VlogQA [18] | Extractive | Food & Lifestyle | 8,047 | 10.09 | 3.22 | ✗ | ✓ | ✗ | YouTube Transcriptions |
| ViRe4MRC [19] | Extractive | Food & Lifestyle | 6,637 | 10.72 | 6.37 | ✓ | ✓ | ✓ | Online Food Reviews |

deeper linguistic insights. These findings identify model strengths and limitations, highlight key areas for improvement, and provide businesses with valuable guidance in selecting and refining ViLLMs.

**Fig. 3.** Overview of our multi-aspect evaluation framework for ViLLMs.

ViLLMs Selection We selected 11 *lightweight open-source ViLLMs* with *state-of-the-art* (SOTA) performance in their respective benchmarks, ranging from 7–9 billion parameters, publicly available as of January 2025. Among them, four models are specifically fine-tuned for Vietnamese: URA-LLaMa-2.1 8B [2], GemSura 7B [2], Vistral 7B [20], and VinaLLaMA 7B [21]. The remaining seven are multilingual models trained on multiple Asian languages, including Vietnamese: SeaLLMs 7B [22], Sailor 7B [23], Qwen2 7B [24], Ghost 8B [25], SEA-LION 7B [26], BLOOMZ 7B [27], and Aya-ExpansE 8B [28]. For benchmarking, we used the latest instruction-tuned or chat-oriented versions to improve domain adaptability and maintain consistent prompt formatting.

Settings To ensure fair benchmarking, all ViLLMs were evaluated under identical inference conditions, including *prompting* and *sampling parameters*. Zero-shot prompting was used to assess their intrinsic ability to generate flu-

ent, accurate responses. We fine-tuned sampling parameters to balance output diversity and stability while constraining response length to match the average human-written responses in CSConDa. All experiments were conducted on a single NVIDIA A100 (40GB) GPU to ensure computational uniformity.

Automatic Metrics We employed a *multi-aspect evaluation* framework integrating *lexical-based*, *semantic-based* [29], and modern *LLM-based* metrics to comprehensively assess ViLLM performance in domain-specific QA. For lexical similarity, we used *BLEU-2* to measure bigram precision, *ROUGE-L* for recall-oriented longest common subsequence overlap, and *METEOR*, which incorporates stemming, synonym matching, and word order flexibility to evaluate word choice, phrasing accuracy, and textual coherence. For semantic alignment, we computed *Cosine Similarity* (Cos. Sim.) at the sentence level and applied *BERTScore* to assess contextual token similarity between ViLLM outputs and human references, both leveraging the SOTA *Vietnamese embedding model*⁷ to capture deeper semantic relationships beyond surface-level word overlap. For *hallucination detection* (Hallu. Score), we employed *Kolena’s prompt-based metric*⁸, which utilizes *OpenAI’s GPT-4* to identify misinformation and fabricated content, ensuring response reliability. However, in real-world scenarios, some ViLLMs fail to generate responses or produce nonsensical loops, significantly impacting usability. To address this, we introduce a *penalty factor* ρ in the final metric computation. Let x_i be the individual metric score for the i -th question, A the number of successfully generated answers, and N the total number of test samples. The adjusted score for metric M is computed as Equation 1.

$$\text{Score}_M = \left(\frac{\sum_{i=1}^N x_i}{N} \right) \times \rho, \quad \rho = \left(\frac{A}{N} \right)^{M_c}, \quad (1)$$

where $M_c = 1$ if a higher value of M indicates better performance and $M_c = -1$ otherwise. This ensures that models failing to generate valid responses are fairly penalized in the final evaluation.

Syntactic Analysis Structural analysis is a key component of our novel LLM evaluation framework, revealing syntactic inefficiencies that traditional metrics often overlook. We introduce five distinctive evaluation metrics: (i) *Word Count*, which analyzes meaningful word segments in an answer; (ii) *Part-of-Speech Ratio* (POS Ratio), which measures information richness by computing the ratio of content to function words in the LLM’s response, as illustrated in Figure 4; (iii) *Phrase Ratio*, similar to the POS Ratio but at the phrase level, as shown in Figure 5; (iv) *Named Entity Difference* (NE Diff.), which quantifies discrepancies in named entities (e.g., people, locations) between the answer and question to assess content consistency; and (v) *Dependency Length* (Dep. Length), which reflects sentence complexity by averaging the distance between heads and their dependents over the total dependencies, as illustrated in Figure 6. For implementation, we employ two renowned *Vietnamese Natural Language Processing*

⁷ <https://huggingface.co/dangvantuan/vietnamese-embedding>

⁸ <https://docs.kolena.com/metrics/prompt-based-hallucination-metric/>

toolkits, namely *Stanza* [30] for POS tagging and *underthesea*⁹ for other tasks contributing to our evaluation method.

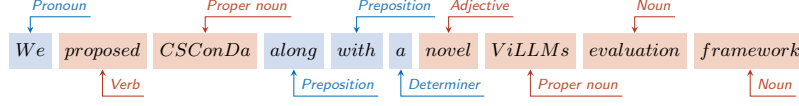


Fig. 4. Illustration of POS Ratio computation. Content words include nouns, verbs, adjectives, adverbs, and proper nouns, whereas function words consist of pronouns, determiners, and prepositions. The POS ratio for this sentence is computed as $POS_{C/F} = \frac{7}{3}$.

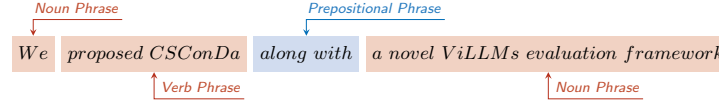


Fig. 5. Illustration of Phrase Ratio computation. Content phrases comprise noun phrases and verb phrases, while function phrases include prepositional phrases. Thus, the phrase ratio for this sentence is determined as $PH_{C/F} = \frac{3}{1}$.

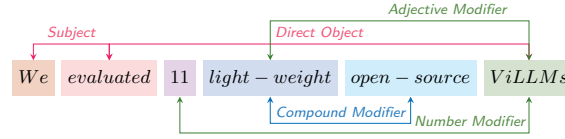


Fig. 6. Illustration of Dependency Length computation with annotated syntactic relations. Arrows represent hierarchical dependencies between words, with colors corresponding to their head words for clarity. Then, the dependency length of the sentence is calculated as $DL = \frac{1+4+3+2+1}{5}$.

Benchmarking Results To rank ViLLMs across different question types in CSConDa and determine the overall ranking, let $r_{i,T}$ be the rank of model X for metric i within type T , and $R_T(X)$ denote its ranking for T . The performance score for each type and the overall ranking are computed as Equation 2.

$$\text{Score}_{X,T} = \frac{1}{\sum_i r_{i,T}}, \quad \text{Score}_{X,\text{overall}} = \frac{1}{\sum_T R_T(X)}. \quad (2)$$

A higher score indicates better performance. This ranking is derived solely from automatic metrics, while syntactic analysis provides complementary insights beyond quantitative evaluation.

4.2 Benchmarking Results

Table 5, Table 6, and Table 7 present benchmark scores of selected ViLLMs across multiple evaluation metrics for the three CSConDa categories, providing a comparative view of model performance. Additionally, Figure 7, Figure 8, and Figure 9 illustrate key syntactic patterns, offering insights into how ViLLMs process different linguistic structures and their impact on performance variations.

⁹ <https://github.com/undertheseanlp/underthesea>

4.3 Analysis and Insights

We analyze the benchmarking results and highlight the following key insights.

Performance of Lightweight Open-Source ViLLMs in Customer Support. Based on the overall ranking, the top five models are Vistral 7B, which achieves the strongest results, followed by SeaLLMs 7B, Sailor 7B, SEA-LION 7B, and GemSUra 7B. However, while these models lead among those evaluated, they are not definitive solutions. Their poor performance on CSConDa benchmarks suggests that current ViLLMs still struggle to generate accurate and fluent responses, particularly when handling queries with linguistic variations.

Vietnamese-Finetuned LLMs vs. Multilingual Models. Our evaluation finds no significant performance gap between LLMs fine-tuned for Vietnamese and multilingual models supporting Vietnamese. This result is expected due to the limited availability of high-quality customer support data in public datasets.

ViLLMs Lack Human-Like Writing Styles. While BERTScore remains acceptable across query types, Cos. Sim. and other lexical-based metrics perform significantly worse. This suggests that ViLLMs produce semantically similar words with different surface forms, leading to responses that deviate from natural human-like writing in this domain.

Table 5. Benchmarking results of ViLLMs on general-type questions.

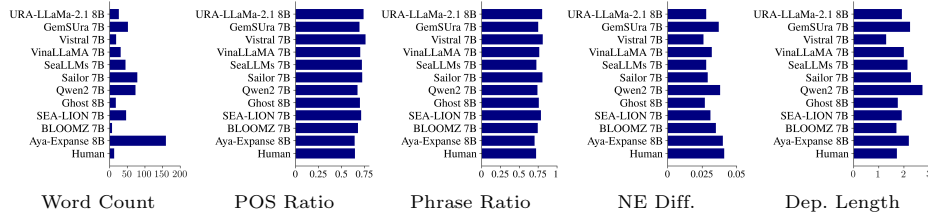
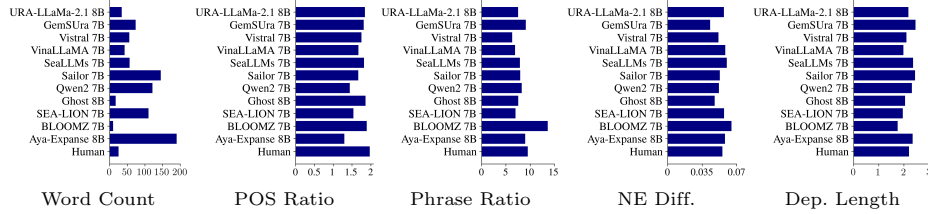
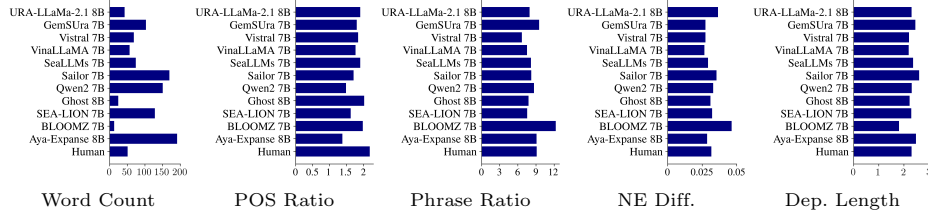
| Model Name | BLEU-2 ↑ | ROUGE-L ↑ | METEOR ↑ | Cos. Sim. ↑ | BERTScore ↑ | Hallu. Score ↓ |
|------------------|--------------|--------------|--------------|--------------|--------------|----------------|
| URA-LLaMa-2.1 8B | 0.004 | 0.161 | 0.072 | 0.219 | 0.654 | 0.782 |
| GemSUra 7B | 0.006 | 0.151 | 0.087 | 0.226 | 0.645 | 0.694 |
| Vistral 7B | 0.005 | 0.171 | 0.089 | 0.269 | 0.659 | 0.522 |
| VinaLLaMA 7B | 0.008 | 0.153 | 0.087 | 0.251 | 0.645 | 0.518 |
| SeaLLMs 7B | 0.001 | 0.166 | 0.047 | 0.247 | 0.661 | 0.598 |
| Sailor 7B | 0.012 | 0.165 | 0.080 | 0.246 | 0.662 | 0.588 |
| Qwen2 7B | 0.006 | 0.133 | 0.086 | 0.261 | 0.634 | 0.612 |
| Ghost 8B | 0.009 | 0.142 | 0.091 | 0.238 | 0.632 | 0.664 |
| SEA-LION 7B | 0.011 | 0.188 | 0.086 | 0.241 | 0.678 | 0.452 |
| BLOOMZ 7B | 0.003 | 0.137 | 0.046 | 0.223 | 0.634 | 0.448 |
| Aya-Expanse 8B | 0.010 | 0.102 | 0.086 | 0.260 | 0.601 | 0.798 |

Table 6. Benchmarking results of ViLLMs on simple-type questions.

| Model Name | BLEU-2 ↑ | ROUGE-L ↑ | METEOR ↑ | Cos. Sim. ↑ | BERTScore ↑ | Hallu. Score ↓ |
|------------------|--------------|--------------|--------------|--------------|--------------|----------------|
| URA-LLaMa-2.1 8B | 0.009 | 0.221 | 0.075 | 0.266 | 0.664 | 0.878 |
| GemSUra 7B | 0.015 | 0.213 | 0.104 | 0.274 | 0.659 | 0.846 |
| Vistral 7B | 0.016 | 0.227 | 0.104 | 0.325 | 0.675 | 0.746 |
| VinaLLaMA 7B | 0.011 | 0.191 | 0.102 | 0.322 | 0.648 | 0.868 |
| SeaLLMs 7B | 0.015 | 0.218 | 0.105 | 0.329 | 0.667 | 0.568 |
| Sailor 7B | 0.012 | 0.213 | 0.088 | 0.304 | 0.670 | 0.774 |
| Qwen2 7B | 0.010 | 0.173 | 0.112 | 0.337 | 0.640 | 0.894 |
| Ghost 8B | 0.012 | 0.187 | 0.110 | 0.323 | 0.645 | 0.900 |
| SEA-LION 7B | 0.013 | 0.206 | 0.073 | 0.253 | 0.674 | 0.616 |
| BLOOMZ 7B | 0.003 | 0.159 | 0.039 | 0.168 | 0.635 | 0.674 |
| Aya-Expanse 8B | 0.017 | 0.171 | 0.115 | 0.347 | 0.631 | 0.984 |

Table 7. Benchmarking results of ViLLMs on complex-type questions.

| Model Name | BLEU-2 \uparrow | ROUGE-L \uparrow | METEOR \uparrow | Cos. Sim. \uparrow | BERTScore \uparrow | Hallu. Score \downarrow |
|------------------|-------------------|--------------------|-------------------|----------------------|----------------------|---------------------------|
| URA-LLaMa-2.1 8B | 0.011 | 0.238 | 0.075 | 0.298 | 0.659 | 0.918 |
| GemSura 7B | 0.022 | 0.242 | 0.107 | 0.335 | 0.658 | 0.902 |
| Vistral 7B | 0.019 | 0.247 | 0.098 | 0.349 | 0.667 | 0.782 |
| VinaLLaMA 7B | 0.019 | 0.232 | 0.107 | 0.356 | 0.652 | 0.912 |
| SeaLLMs 7B | 0.022 | 0.247 | 0.105 | 0.365 | 0.666 | 0.642 |
| Sailor 7B | 0.015 | 0.231 | 0.089 | 0.345 | 0.665 | 0.864 |
| Qwen2 7B | 0.015 | 0.213 | 0.108 | 0.371 | 0.641 | 0.966 |
| Ghost 8B | 0.020 | 0.228 | 0.108 | 0.361 | 0.648 | 0.932 |
| SEA-LION 7B | 0.012 | 0.217 | 0.067 | 0.288 | 0.670 | 0.662 |
| BLOOMZ 7B | 0.004 | 0.146 | 0.034 | 0.168 | 0.606 | 0.712 |
| Aya-Expans 8B | 0.027 | 0.223 | 0.115 | 0.393 | 0.641 | 0.982 |

**Fig. 7.** Syntactic analysis of general-type questions.**Fig. 8.** Syntactic analysis of simple-type questions.**Fig. 9.** Syntactic analysis of complex-type questions.

ViLLMs Are Verbose and Structurally Inflated. Syntactic analysis shows that human responses are consistently more concise across all query types, as indicated by their lower word count. In contrast, ViLLMs over-generate words and sentences, producing longer responses. Furthermore, Hallu. Score is lowest for general-type questions, where ViLLMs provide the shortest answers. These findings suggest that verbosity induces semantic drift, reducing alignment with the original query and increasing the likelihood of errors and hallucinations.

ViLLMs Exhibit Stronger Structural Dependency. The Dep. Length metric shows that human responses consistently have the shortest dependency lengths, indicating minimal syntactic complexity and greater structural flexibility. This suggests a more fluid and unstructured style, characteristic of human-like responses. In contrast, ViLLMs, trained predominantly on structured data, exhibit higher dependency parsing scores, leading to rigid sentence constructions and a writing style that lacks natural variation.

Towards Structurally-Aware Fine-Tuning of ViLLMs. Although POS Ratio, Phrase Ratio, and NE Diff. metrics indicate that ViLLMs maintain internal consistency, their responses remain overly long and structurally rigid, often leading to semantic drift and increased errors. Our comprehensive benchmarking highlights structural constraints as a key barrier to semantic quality, particularly in customer support interactions, where natural and adaptive communication is essential. Addressing these limitations requires a refined fine-tuning approach that enhances structural efficiency, promotes syntactic adaptability, and optimizes linguistic economy.

5 Conclusion, Limitation, and Future Work

In this work, we introduce CSConDa, the first Vietnamese QA dataset for benchmarking customer support interactions, alongside a novel evaluation framework. Our approach extends beyond conventional lexical and semantic assessments by integrating advanced hallucination detection and syntactic analysis. We outline the dataset construction process and provide an in-depth linguistic characterization, highlighting CSConDa’s distinctions from existing resources. Using the test split, we conduct a comprehensive evaluation of SOTA lightweight open-source ViLLMs, following a survey underscoring the need for rigorous benchmarking. Our findings indicate that while LLMs demonstrate strong grammatical accuracy and fluency, they struggle with structural efficiency and adaptability. In contrast, human responses prioritize brevity, clarity, and contextual relevance, reinforcing the importance of structurally-aware fine-tuning to bridge this gap. CSConDa, along with our evaluation results, offers enterprises a valuable tool for identifying suitable ViLLMs and guiding the development of effective QA systems. Despite the robustness of our evaluation framework, its primary limitation is focusing solely on intrinsic model capabilities. Future work could enhance CSConDa by incorporating contextual information, including relevant details essential for answering queries, ensuring such extensions preserve its realism and task-oriented nature for practical applications.

References

1. F. Li, Y. Wang, Y. Xu, S. Wang, J. Liang, Z. Chen, W. Liu, Q. Feng, T. Duan, Y. Huang, Q. Song, and X. Li, “Performance evaluations of large language models for customer service,” *International Journal of Machine Learning and Cybernetics*, 2024.

2. S. Truong, D. Nguyen, T. Nguyen, D. Le, N. Truong, T. Quan, and S. Koyejo, "Crossing Linguistic Horizons: Finetuning and Comprehensive Evaluation of Vietnamese Large Language Models," in *Findings of the Association for Computational Linguistics: NAACL 2024*, pp. 2849–2900, Association for Computational Linguistics, 2024.
3. Z. Wan, X. Wang, C. Liu, S. Alam, Y. Zheng, J. Liu, Z. Qu, S. Yan, Y. Zhu, Q. Zhang, M. Chowdhury, and M. Zhang, "Efficient Large Language Models: A Survey," *Transactions on Machine Learning Research*, 2024.
4. K. Nguyen, V. Nguyen, A. Nguyen, and N. Nguyen, "A Vietnamese Dataset for Evaluating Machine Reading Comprehension," in *Proceedings of the 28th International Conference on Computational Linguistics*, pp. 2595–2605, International Committee on Computational Linguistics, 2020.
5. K. Le, H. Nguyen, T. Le Thanh, and M. Nguyen, "VIMQA: A Vietnamese Dataset for Advanced Reasoning and Explainable Multi-hop Question Answering," in *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pp. 6521–6529, European Language Resources Association, 2022.
6. K. Van Nguyen, T. Van Huynh, D.-V. Nguyen, A. G.-T. Nguyen, and N. L.-T. Nguyen, "New Vietnamese Corpus for Machine Reading Comprehension of Health News Articles," *ACM Transactions on Asian and Low-Resource Language Information Processing*, vol. 21, no. 5, 2022.
7. S. T. Luu, M. N. Bui, L. D. Nguyen, K. V. Tran, K. Van Nguyen, and N. L.-T. Nguyen, "Conversational Machine Reading Comprehension for Vietnamese Healthcare Texts," in *Advances in Computational Collective Intelligence*, pp. 546–558, Springer International Publishing, 2021.
8. K. V. Nguyen, K. V. Tran, S. T. Luu, A. G.-T. Nguyen, and N. L.-T. Nguyen, "Enhancing Lexical-Based Approach With External Knowledge for Vietnamese Multiple-Choice Machine Reading Comprehension," *IEEE Access*, vol. 8, pp. 201404–201417, 2020.
9. X.-Q. Dao, N.-B. Le, T.-D. Vo, X.-D. Phan, B.-B. Ngo, V.-T. Nguyen, T.-M.-T. Nguyen, and H.-P. Nguyen, "VNHSGE: VietNameese High School Graduation Examination Dataset for Large Language Models," 2023.
10. P. N.-T. Do, N. D. Nguyen, T. Van Huynh, K. Van Nguyen, A. G.-T. Nguyen, and N. L.-T. Nguyen, "Sentence Extraction-Based Machine Reading Comprehension for Vietnamese," in *Knowledge Science, Engineering and Management*, pp. 511–523, Springer International Publishing, 2021.
11. C. T. Nguyen and D. T. Nguyen, "Building a Discourse-Argument Hybrid System for Vietnamese Why-Question Answering," *Computational Intelligence and Neuroscience*, vol. 2021, no. 1, p. 6550871, 2021.
12. M.-N. Tran, P.-V. Nguyen, L. Nguyen, and D. Dinh, "ViMedAQA: A Vietnamese Medical Abstractive Question-Answering Dataset and Findings of Large Language Model," in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*, pp. 252–260, Association for Computational Linguistics, 2024.
13. T. Thai, N. C. Thao-Ha, A. Vo, and S. Luu, "UIT-ViCoV19QA: A Dataset for COVID-19 Community-based Question Answering on Vietnamese Language," in *Proceedings of the 36th Pacific Asia Conference on Language, Information and Computation*, pp. 801–810, Association for Computational Linguistics, 2022.
14. N. T.-H. Nguyen, P. P.-D. Ha, L. T. Nguyen, K. Van Nguyen, and N. L.-T. Nguyen, "SPBERTQA: A Two-Stage Question Answering System Based on Sentence Transformers for Medical Texts," in *Knowledge Science, Engineering and Management*, pp. 371–382, Springer International Publishing, 2022.

15. P.-T. P. Do, D.-N. D. Cao, K. Q. Tran, and K. Van Nguyen, “R2GQA: Retriever-Reader-Generator Question Answering System to Support Students Understanding Legal Regulations in Higher Education,” 2024.
16. D. D. Minh, V. N. Van, and T. D. Cong, “Using Large Language Models for education managements in Vietnamese with low resources,” 2025.
17. T.-H.-Y. Vuong, H.-T. Nguyen, Q.-H. Nguyen, L.-M. Nguyen, and X.-H. Phan, “Improving Vietnamese Legal Question–Answering System Based on Automatic Data Enrichment,” in *New Frontiers in Artificial Intelligence*, pp. 49–65, Springer Nature Switzerland, 2024.
18. T. Ngo, K. Dang, S. Luu, K. Nguyen, and N. Nguyen, “VlogQA: Task, Dataset, and Baseline Models for Vietnamese Spoken-Based Machine Reading Comprehension,” in *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics*, pp. 1310–1324, Association for Computational Linguistics, 2024.
19. T. P. P. Do, N. D. D. Cao, N. T. Nguyen, T. V. Huynh, and K. V. Nguyen, “Machine Reading Comprehension for Vietnamese Customer Reviews: Task, Corpus and Baseline Models,” in *Proceedings of the 37th Pacific Asia Conference on Language, Information and Computation*, pp. 24–35, Association for Computational Linguistics, 2023.
20. C. V. Nguyen, T. Nguyen, Q. Nguyen, H. Nguyen, B. Plüster, N. Pham, H. Nguyen, P. Schramowski, and T. Nguyen, “Vistral-7B-Chat - Towards a State-of-the-Art Large Language Model for Vietnamese,” 2023.
21. Q. Nguyen, H. Pham, and D. Dao, “VinaLLaMA: LLaMA-based Vietnamese Foundation Model,” 2023.
22. W. Zhang, H. P. Chan, Y. Zhao, M. Aljunied, J. Wang, *et al.*, “SeaLLMs 3: Open Foundation and Chat Multilingual Large Language Models for Southeast Asian Languages,” 2024.
23. L. Dou, Q. Liu, G. Zeng, J. Guo, J. Zhou, X. Mao, Z. Jin, W. Lu, and M. Lin, “Sailor: Open Language Models for South-East Asia,” in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 424–435, Association for Computational Linguistics, 2024.
24. “Qwen2 Technical Report,” 2024.
25. Ghost X, Hieu Lam, “Ghost 8b beta,” 2024.
26. D. Ong and P. Limkonchotiawat, “SEA-LION (Southeast Asian Languages In One Network): A Family of Southeast Asian Language Models,” in *Proceedings of the 3rd Workshop for Natural Language Processing Open Source Software (NLP-OSS 2023)*, pp. 245–245, Association for Computational Linguistics, 2023.
27. N. Muennighoff, T. Wang, L. Sutawika, A. Roberts, S. Biderman, T. L. Scao, M. S. Bari, S. Shen, Z.-X. Yong, H. Schoelkopf, *et al.*, “Crosslingual Generalization through Multitask Finetuning,” 2022.
28. J. Dang, S. Singh, D. D’souza, A. Ahmadian, A. Salamanca, M. Smith, A. Peppin, S. Hong, M. Govindassamy, T. Zhao, *et al.*, “Aya Expanse: Combining Research Breakthroughs for a New Multilingual Frontier,” 2024.
29. A. Chen, G. Stanovsky, S. Singh, and M. Gardner, “Evaluating Question Answering Evaluation,” in *Proceedings of the 2nd Workshop on Machine Reading for Question Answering*, pp. 119–124, Association for Computational Linguistics, 2019.
30. P. Qi, Y. Zhang, Y. Zhang, J. Bolton, and C. D. Manning, “Stanza: A Python Natural Language Processing Toolkit for Many Human Languages,” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, 2020.