

Mind the Gap: A Decade-Scale Empirical Study of Multi-Stakeholder Dynamics in VR Ecosystem

YIJUN LU, Waseda University, Japan

HIRONORI WASHIZAKI*, Waseda University, Japan

NAOYASU UBAYASHI, Waseda University, Japan

NOBUKAZU YOSHIOKA, Waseda University, Japan

CHENHAO WU, Waseda University, Japan

MASANARI KONDO, Kyushu University, Japan

YUYIN MA, Xinjiang University, China

JIONG DONG, Xuchang University, China

JIANJIN ZHAO, Beijing University of Posts and Telecommunications, China

DONGQI HAN, Beijing University of Posts and Telecommunications, China

In the development and evolution of VR ecosystem, platform stakeholders continuously adapt their products in response to user and technical feedback, often reflected in subtle shifts in discussion topics or system updates. A comprehensive understanding of these changes is essential for identifying gaps between user expectations and developer actions, which can guide more effective quality assurance and user-centered innovation. While previous studies have analyzed either user reviews or developer discussions in isolation, such approaches typically fail to reveal how specific user concerns are (or are not) addressed by corresponding technical activities. To address this limitation, our study introduces a multi-view empirical framework that systematically compares and aligns stakeholder perspectives. By applying topic modeling and quantitative impact analysis to 944,320 user reviews and 389,477 developer posts, we identify not only the overlap in concerns (e.g., performance, input methods), but also clear gaps in areas like inclusivity and community safety (e.g., LGBTQ+ representation, child-friendly content). Our findings show that while users repeatedly raise such issues, they are rarely discussed in developer forums. These insights enable data-driven recommendations for closing the user-developer gap in VR ecosystems, offering practical implications for platform governance and the design of next-generation VR systems.

CCS Concepts: • **Software and its engineering** → **Software design engineering**; • **Human-centered computing** → **Empirical studies in HCI**.

Additional Key Words and Phrases: Empirical Study, User Experience, Virtual Reality, Topic Modeling, Large Language Models

Authors' Contact Information: Yijun Lu, Waseda University, Tokyo, Japan, yijun@ruri.waseda.jp; Hironori Washizaki, Waseda University, Tokyo, Japan, washizaki@waseda.jp; Naoyasu Ubayashi, Waseda University, Tokyo, Japan, ubayashi@aoni.waseda.jp; Nobukazu Yoshioka, Waseda University, Tokyo, Japan, nobukazu@engineerale.ai; Chenhao Wu, Waseda University, Tokyo, Japan, wuchenhao@toki.waseda.jp; Masanari Kondo, Kyushu University, Kyushu, Japan, kondo@ait.kyushu-u.ac.jp; Yuyin Ma, Xinjiang University, Urumqi, China, mayuyin@xju.edu.cn; Jiong Dong, Xuchang University, Xuchang, China, jiongdong@xcu.edu.cn; Jianjin Zhao, Beijing University of Posts and Telecommunications, Beijing, China, jianjinzhao@bupt.edu.cn; Dongqi Han, Beijing University of Posts and Telecommunications, Beijing, China, handongqi@bupt.edu.cn.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

Manuscript submitted to ACM

1

ACM Reference Format:

Yijun Lu, Hironori Washizaki, Naoyasu Ubayashi, Nobukazu Yoshioka, Chenhao Wu, Masanari Kondo, Yuyin Ma, Jiong Dong, Jianjin Zhao, and Dongqi Han. 2025. Mind the Gap: A Decade-Scale Empirical Study of Multi-Stakeholder Dynamics in VR Ecosystem. 1, 1 (August 2025), 30 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

Virtual Reality (VR) has rapidly evolved from a niche technology into a widely adopted medium, fundamentally transforming both the ways users interact with software and how developers approach application design. The global VR games market was valued at \$47.1 billion in 2024 and is projected to surge to \$346.0 billion by 2033 [45]. This remarkable expansion is driven by increasingly immersive user experiences, which generate rich streams of user feedback that illuminate VR-specific needs and usability challenges [13]. Understanding these nuanced user perspectives is essential, as VR players frequently express concerns (e.g., motion sickness, inclusivity, and hardware limitations) that are less prevalent in traditional software [30]. Addressing these unique issues is critical for developers seeking to align products with evolving user expectations and to fully realize the potential of VR platforms [30]. At the same time, the VR ecosystem involves a diverse set of stakeholders, from independent developers to global studios, whose priorities, constraints, and interpretations often diverge [39]. Bridging the persistent gap between user experiences and developer practices is therefore central to advancing VR technology.

However, the current research landscape is dominated by single-stakeholder analytics, with most prior studies focusing exclusively on either user feedback or developer discourse. On the user side, analyses of app store reviews (including those for VR applications) have revealed user preferences and criticisms, most notably frequent requests for enhanced immersion, richer content, and better accessibility [10, 24]. For example, Dong et al.[10] found that while users value aspects like music and gameplay, their most persistent complaints concern bugs, insufficient content, and high costs. On the developer side, empirical analyses of Q&A forums, code repositories, and technical discussions have revealed the practical challenges that developers encounter, including toolchain adaptation, cross-platform support, and performance optimization[4]. Although each of these research streams offers important perspectives, their separation makes it challenging to trace the translation of user concerns into developer priorities and technical solutions.

This separation has led to a well-documented misalignment between user needs and developer focus, especially in rapidly evolving domains like VR. As recent work shows, user feedback in VR shifts quickly from foundational concerns such as physical comfort to higher-level expectations involving social interaction, content diversity, and community safety [13, 30]. However, developers may not adapt with these changes, resulting in critical issues such as inclusivity, accessibility, and safety receiving insufficient attention in development processes. For example, our empirical results show that inclusivity topics, e.g., the representation of LGBTQ+ users or child-friendly content, are discussed by users but receive minimal attention in developer forums. Conversely, developers devote considerable discussion to technical challenges that may not directly address user frustrations or desires. This persistent gap underscores the necessity of integrated, multi-stakeholder empirical analysis in VR ecosystems.

To address this critical gap, we conduct the first large-scale empirical study that systematically integrates and compares the perspectives of both VR users and developers. Leveraging a dataset of 944,320 user reviews and 389,477 developer posts from three major VR platforms (Mate, SteamVR, and VivePort), we provide a comprehensive view of stakeholder dialogue and its evolution over time. For efficient analysis of this extensive corpus, we introduce a hybrid approach, HTModel (See Section 3.3), which decouples topic discovery from semantic interpretation by combining scalable clustering with LLM-based labeling to produce interpretable topic maps. We perform topic modeling separately

on user and developer datasets, then merge and categorize the resulting topics into stakeholder-relevant themes. This enables absolute and relative impact analysis for each topic across both groups. Our unified framework identifies overlooked user concerns, highlights stakeholder alignment, and provides actionable insights for research and practice. In summary, this work fills a key gap in VR analytics by offering a reproducible, multi-view framework that links user experience with developer challenges. The workflow is summarized in Fig. 1.

Specifically, we address the following research questions (RQs), along with a summary of the key findings for each. Detailed results are presented in Section 4.

- **RQ1: What topics are most frequently discussed by VR developers and users?**

Motivation: This research question aims to reveal the most prominent topics discussed by each stakeholder group and highlight key mismatches in their priorities.

Results: Technical performance issues such as hardware support (users 2.6%, developers 2.1%) and input methods (users 1.8%, developers 2.8%) are common to both groups. In contrast, inclusivity and community-related topics including LGBTQ+ representation (users 1.9%), child-friendly content (users 2.0%), and female character inclusivity (users 2.4%) are frequently raised by users but are virtually absent in developer forums. Users also mention language support (1.7%) and pricing (3.7%) more often, whereas developers discuss toolchains (4.0%), cross-platform issues (3.8%), and server management (2.8%) more frequently. These findings quantify a clear misalignment, especially regarding social, inclusivity, and accessibility issues.

- **RQ2: How do these topics evolve over time from different perspectives?**

Motivation: This research question investigates whether the focus of each stakeholder group shifts over time and how their evolving concerns converge or diverge.

Results: Software-related topics accounted for up to 60% of user discussions in earlier years. Since 2021, mentions of content quality, replay value, and inclusivity have increased by more than 30%. Community experience, child-friendly design, and diverse avatars have become notably more prominent in recent user reviews. Developer focus on user experience increased to about 20% by 2024, but inclusivity, safety, and accessibility topics remain almost absent in developer posts, even as these themes appear in up to 6% of user comments during the same period. Developer attention to new user concerns often lags behind shifts in user discussions.

The primary contributions of this study are summarized as follows:

- We conduct the first large-scale empirical study integrating and systematically comparing the perspectives of both users and developers in VR ecosystems, analyzing 944,320 user reviews and 389,477 developer posts collected from three major VR platforms.
- We identify and quantify significant gaps between user-centric concerns (e.g., inclusivity, accessibility, pricing) and developer-focused priorities (e.g., toolchains, cross-platform compatibility), revealing critical gaps between user expectations and developer practices.
- We track the temporal evolution of discussion topics within each stakeholder group, highlighting areas where developer attention lags behind emerging user needs, such as community experience, diverse avatars, and child-friendly content.
- To facilitate the analysis, we propose HTModel, a hybrid topic-modeling pipeline combining traditional clustering methods with LLM-based semantic labeling. HTModel significantly reduces manual labeling effort compared to traditional topic models and achieves over 99% cost savings in LLM API usage compared to purely LLM-driven approaches.

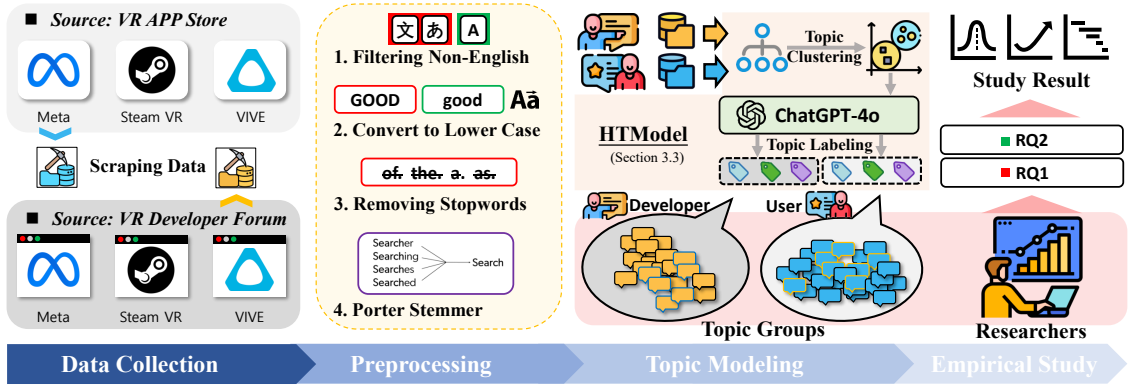


Fig. 1. Overview of the research workflow, encompassing data collection from major VR platforms, comprehensive preprocessing, a hybrid topic modeling pipeline, and empirical analysis. This workflow enables systematic investigation of user–developer dynamics in the VR ecosystem by integrating large-scale data acquisition, standardized text processing, advanced topic extraction, and multi-view analysis.

- We provide the first publicly available multi-stakeholder textual dataset from major VR platforms, supporting future reproducibility and empirical research in the VR ecosystem. The dataset is available here.

The remainder of this paper is structured as follows. Section 2 reviews the relevant literature. Section 3 details the data collection, processing, and topic extraction methodology. Section 4 presents the empirical study of user reviews and developer posts. Section 5 discusses the broader implications of our findings. Section 6 outlines threats to validity. Finally, Section 7 concludes the paper. Additional methodological details are provided in the Appendix.

2 Related Work

Building on this foundation, our work integrates three key research streams: (1) empirical multi-stakeholder analysis in software engineering, (2) empirical studies focused on VR systems, and (3) advances in empirical methods for topic modeling. By synthesizing these perspectives, we address existing limitations and provide a unified framework for comprehensive multi-view analysis in the VR domain.

2.1 Multi-stakeholder Empirical Studies in Software Engineering

Research in software engineering has long emphasized the need to align perspectives across stakeholder groups. For instance, Hassan et al.[19] showed that many user-reported issues in Google Play app reviews are never addressed in developer discussions, leading to unresolved bugs and unmet user expectations. Buchan et al.[6] used repertory grid techniques in agile teams to uncover subtle misalignments between user representatives and developers, illustrating how differing expectations can hinder collaboration. Similarly, Hasan et al.[18] surveyed 181 professional developers and found that demands from different stakeholders, such as peers, managers, and users, substantially impact developer satisfaction and productivity. Lenberg et al.[26] further demonstrated that organizational value misalignment is significantly associated with reduced team performance and increased conflict. Additionally, Mauerer et al. [31] conducted a longitudinal study of 25 open-source projects, revealing that misalignment between social structures and technical dependencies does not always result in poor code quality, highlighting the complexity of stakeholder

alignment in software ecosystems. Despite this growing body of work, few studies have directly examined the distinct challenges of stakeholder alignment in VR systems.

2.2 VR Empirical Studies with a Single Stakeholder Focus

Empirical research in VR has predominantly relied on single-stakeholder perspectives, which provide useful but incomplete insights. On the user side, text-mining studies of app reviews are common. Lu et al. [30] analyzed over 176,000 VR game reviews, identifying user concerns such as performance, tracking, and compatibility. Dong et al. [10] examined more than 105,000 social VR reviews and highlighted the demand for avatar customization. Li et al. [28] further analyzed over one million user reviews to develop a taxonomy of VR-specific quality attributes. Additional studies focus on immersive usability: Zhang et al. [50] conducted a meta-analysis of Cinematic VR, revealing inconsistent measurement of presence and immersion, while Rzig et al. [40] found that 74% of open-source VR projects lacked adequate test coverage. Singh and O'Hagan [41] used topic modeling on 40,000 Rec Room reviews, identifying emergent issues such as harassment through sentiment analysis.

On the developer side, separate studies have addressed VR-specific tool and testing challenges. Karre et al. [23] surveyed VR software teams across several countries and found that hybrid engineering practices are needed but understudied. Survey-based work by Ashtari et al. [2] and Dhia et al. [40] report ongoing deficiencies in testing workflows, tool support, and quality assurance. While these studies enhance our understanding of either user or developer needs in VR, there is still a lack of large-scale, integrated analysis of both perspectives. It remains uncertain how developer priorities correspond with user concerns or how quickly developers respond to shifting user expectations. Addressing these gaps, our work is the first to systematically compare user reviews and developer discussions in the VR domain through a multi-stakeholder lens.

2.3 Advances in Empirical Methods for Topic Modeling

Automated methods for topic analysis have evolved significantly, ranging from early semantic approaches like Latent Semantic Analysis and Non-negative Matrix Factorization (LSA, NMF) to probabilistic techniques such as Latent Dirichlet Allocation (LDA) [5, 8]. While LDA improves topic discovery, it still suffers from label ambiguity and requires expert interpretation [13, 37]. Embedding-based models such as BERTopic and Top2Vec [1, 15] further enhance coherence and handle short texts more effectively. However, these models retain the need for manual labeling of topic clusters [37].

The rise of large language models (LLMs) has enabled the automation of topic labeling. Approaches like TopicGPT [36], generative semantic labeling techniques [25], and LLM-based frameworks such as LimTopic [3] and Kapoor's method for BERTopic labels [22] demonstrate improved interpretability, often matching or surpassing human judgment [3, 22]. At the same time, iterative methods like LITA [7] show that combining LLMs with traditional clustering can improve topic coherence and reduce labeling effort. Despite these advancements, purely LLM-driven labeling remains expensive and difficult to scale to millions of documents, as costs grow linearly with data volume [25, 36]. Our HTModel builds on these foundations by first applying efficient traditional clustering techniques (i.e., LDA, BERTopic, NMF, LSA, and Top2Vec.) at scale, then selectively using LLMs to assign semantic labels to clusters. This hybrid design achieves interpretable, semantically rich topic extraction across large corpora while drastically reducing manual annotation and API costs.

Table 1. Dataset Characteristics. The $D_{developer}$ dataset covers the period from 2013 to 2024, and the D_{user} dataset spans from 2015 to 2024. All data were collected and finalized between July 1–15, 2024.

Dataset	Platforms	Number of Collected Data	Number of Processed Data
$D_{developer}$	VivePort	43,053	41,782
	Meta	352,463	342,104
	SteamVR	5,932	5,591
	Subtotal	401,448	389,477
D_{user}	VivePort	17,844	15,964
	Meta	393,355	351,480
	SteamVR	646,642	576,876
	Subtotal	1,057,841	944,320
Total		1,459,289	1,333,797

3 Methodology

To support the subsequent VR empirical study and automated analysis, we collect VR-related textual data, apply noise reduction, and introduce a hybrid pipeline that combines traditional topic models with GPT-4o to realize enhanced topic extraction.

3.1 Data Collection

To enable a comprehensive multi-stakeholder analysis, we collected large-scale textual data from three leading VR gaming platforms: Meta [33], SteamVR [44], and Viveport [48]. These platforms are widely recognized in the literature for their extensive VR game catalogs and vibrant user communities, making them representative sources for studying user experience in the VR domain [11, 13, 30]. By focusing on VR_Only games, we ensured that the collected user reviews directly reflect authentic VR experiences, in line with established empirical research practices [10, 30].

In addition to user reviews, our work uniquely incorporates large-scale developer forum data from the Meta Forum [34], SteamVR Forum [43], and HTC Vive Forum [20]. While forum-based analysis has proven valuable for exploring collaboration and technical knowledge in software engineering and open source communities [21, 46], the systematic study of VR developer forums remains rare. By combining both user and developer sources, our dataset enables a holistic investigation of the dialogue and gaps between end-user expectations and developer practices in the VR ecosystem.

We developed customized web crawlers to systematically collect and synchronize user reviews and developer discussions from these platforms. The resulting datasets, D_{user} and $D_{developer}$, support robust, multi-perspective topic modeling and empirical analysis. See Table 1 for detailed statistics.

3.2 Data Processing

To ensure the reliability and consistency of subsequent analyses, both D_{user} and $D_{developer}$ datasets underwent rigorous preprocessing, following best practices in large-scale empirical software analytics [17, 32, 42]. The main steps are summarized as follows:

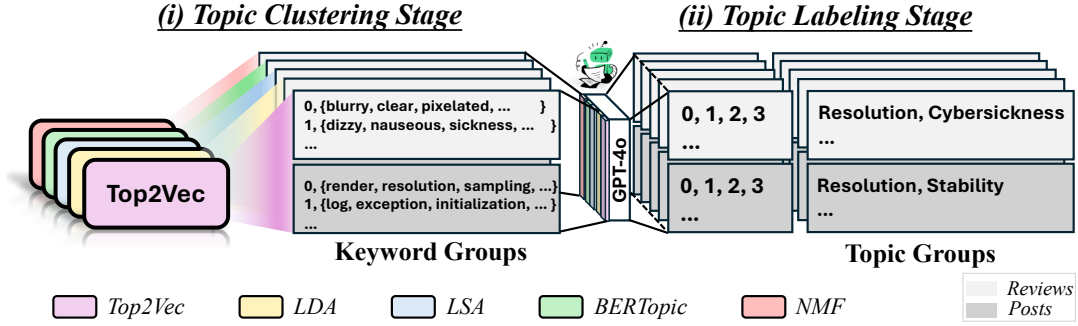


Fig. 2. Implementation process of the proposed Hybrid Topic Modeling (HTModel). The pipeline comprises two main stages: (i) parallel topic clustering using multiple algorithms (Top2Vec, LDA, LSA, BERTopic, and NMF) to generate diverse keyword groups, and (ii) automated topic labeling with ChatGPT-4o to produce unified, interpretable topic groups for downstream analysis.

- **Filtering Non-English:** Only English-language content was retained using the *lingua-language-detector*[42], as many state-of-the-art natural language processing models are primarily trained and evaluated on English corpora, and thus exhibit optimal performance or applicability for English text[29].
- **Convert to Lowercase:** All textual data were converted to lowercase using Pandas [32], ensuring case-insensitive processing and vocabulary normalization, thereby reducing data sparsity and potential analytical artifacts in downstream modeling.
- **Removing Stopwords:** Standard stopwords removal, including the application of a domain-specific (VR-related) stopwords list, was conducted with NLTK [17]. This step serves to eliminate high-frequency, semantically uninformative terms, thus enhancing the interpretability and coherence of extracted topics.
- **Porter Stemmer:** The Porter stemming algorithm [35] was employed to standardize inflected word forms to their roots, reduce feature dimensionality, and improve the effectiveness of topic modeling and other downstream natural language processing tasks [29, 37].

These steps collectively ensure high-quality input for subsequent topic extraction. The final dataset statistics after processing are shown in Table 1.

3.3 Topic Extraction

To achieve robust and interpretable topic analysis for multi-stakeholder VR textual data, we propose the HTModel pipeline to extract topics. First, we employ five topic modeling approaches, namely LDA, LSA, NMF, BERTopic, and Top2Vec, to cluster large-scale textual data effectively, thereby reducing reliance on computationally intensive LLM APIs. We then leverage GPT-4o to enhance topic interpretability by generating semantically refined labels for the extracted keyword clusters, thus eliminating the need for manual annotation. The overall workflow of the proposed HTModel, encompassing both topic clustering and topic labeling, is illustrated in Figure 2.

3.3.1 Topic Clustering. Building on the foundational work of Roman et al. [12], we select the LDA, BERTopic, Top2Vec, LSA, and NMF topic models to perform a comparative analysis across datasets $D_{developer}$ and D_{user} . LDA represents generative models, and the LSA and NMF models are based on matrix decomposition. In contrast, the BERTopic and Top2Vec models employ embedding-based clustering techniques.

These five models are chosen because they collectively represent the three dominant families of topic modeling approaches in recent research: (1) probabilistic models (LDA), (2) matrix factorization models (LSA, NMF), and (3) embedding-based clustering models (BERTopic, Top2Vec). This selection ensures a comprehensive and balanced comparison, capturing the methodological diversity and current state-of-the-art in the topic modeling literature [1, 5, 8, 12, 15].

To ensure a consistent and fair evaluation, we implement all five topic models with a fixed number of topics ($T = 50$) and maintain the number of keywords per topic at $Word = 10$, aligning with parameters commonly used in manual annotation [11, 13, 30]. This consistency allows for a direct comparison between traditional manual labeling and our GPT-based approach. Note that the topics generated by these models are essentially clusters of keywords that serve as the foundation for subsequent semantic refinement and labeling.

3.3.2 Topic Labeling. After applying the five conventional topic models for topic clustering, the large-scale datasets yield 250 representative topics, each represented by a group of 10 related keywords. These keyword groups are then input to GPT-4o to generate the topic labels for each group, a process that replaces traditional manual annotation tasks. Drawing inspiration from recent advances in LLM-based topic modeling [3, 7, 22, 25, 36], we further refine and adapt the prompts to suit the characteristics of our datasets and research objectives.

All experiments utilize OpenAI GPT-4o (API access dates: November 1, 2024, to January 1, 2025), during which the GPT-4o model parameters remain stable, ensuring consistency and reliability across experimental trials. Moreover, we conduct systematic prompt engineering and confirm that minor variations in prompt design have minimal impact on the stability and accuracy of the semantic labeling process.

Prompt Design for Semantic Labeling of D_{user} and $D_{developer}$	
Prompt 1: User Review Topic Labeling This analysis categorizes VR user reviews into 50 topics. One topic is defined by the following keywords: {keywords}. Generate a concise, semantically meaningful label that best represents this topic. Use the format: Topic: <topic label>.	Prompt 2: Developer Post Topic Labeling This analysis categorizes VR developer discussions into 50 topics. One topic is defined by the following keywords: {keywords}. Generate a concise, semantically meaningful label that best represents this topic. Use the format: Topic: <topic label>.

This two-stage pipeline enables interpretable and scalable topic extraction for millions of documents across both user and developer datasets. We conducted both qualitative and quantitative evaluations of the hybrid model and ultimately selected LDA_GPT as the analytical framework. Through multiple rounds of sensitivity analysis, we determined the optimal number of topics to be $T_{D_{developer}} = 40$ for developer discussions and $T_{D_{user}} = 50$ for user discussions, with the number of keywords set to $K = 20$. Further methodological details can be found in Section A.

4 Empirical study

As mentioned previously, a multiview analysis is conducted in this study to identify the key concerns of VR users and the primary priorities of developers. Unlike previous studies that focus primarily on a single perspective (e.g., users or developers), the current study integrates both user and developer perspectives to systematically investigate similarities and differences and analyze how these evolve over time in the VR domain. Thus, this study is to explore a more comprehensive understanding of these dynamics by leveraging a dataset of user reviews (denoted with **blue** R-prefixed IDs) and developer posts (denoted with **orange** P-prefixed IDs).

After topic modeling and automated labeling, each comment and post is assigned to a specific topic. As shown in Figures 4 and 5, the average number of reviews or posts per topic is 21,582 for the user dataset and 10,246 for the developer dataset. According to standard statistical formulas, the minimum required sample sizes to achieve a 95% confidence level with a $\pm 5\%$ margin of error are 378 and 370, respectively. To further enhance the robustness and reliability of our qualitative analysis, we conservatively select 500 reviews or posts per topic. The qualitative analysis is independently conducted by two authors per dataset, each possessing more than five years of development experience and substantial expertise in VR.

4.1 (RQ1) Which topics are most frequently discussed by VR developers and users

4.1.1 Motivation. User reviews and developer forum discussions serve as critical resources for understanding the VR domain, offering dual perspectives on players' perceptions of game quality and the technical challenges faced by developers.

In this study, we explore overlaps and divergences between perspectives of the users and developers by conducting a multiview textual data analysis, aiming to bridge the experience-implementation gap between user feedback and developer challenges. The results of this analysis provide actionable insights for designing VR systems that better align with the user's expectations while also addressing the real-world technical and practical constraints of developers.

4.1.2 Approach. In Section A.2, we evaluate the effectiveness of the proposed HTModel and select the LDA_GPT model for topic extraction. Based on the scale and content of the datasets, we determine the optimal number of topics using the coherence score, resulting in $T_{D_{developer}} = 40$ topics and $T_{D_{user}} = 50$ topics for the developer and user datasets, respectively.

To further analyze the similarities and differences in focus between developers and users, a two-step method is proposed to calculate the intra-dataset and the inter-dataset topic similarity. Specifically, we employ the Sentence-BERT (S-BERT) model [38] and cosine similarity to address the potential redundancy in topic labels in the first step. We then calculate the inter-dataset topic similarity, the cosine similarity is shown as follows:

$$CoS_{A_i, B_i} = \frac{\sum_{i=1}^n A_i \cdot B_i}{\sqrt{\sum_{i=1}^n A_i^2} \cdot \sqrt{\sum_{i=1}^n B_i^2}} \quad (1)$$

where A_i and B_i denote the feature vectors of the two topic labels. We conduct a sensitivity analysis to determine the optimal merging threshold. When the threshold is set too low (e.g., 0.7), semantically distinct topics such as *Bug Reporting* and *Troubleshooting* are erroneously merged, resulting in excessive generalization. Conversely, a high threshold (e.g., 0.9) prevents the merging of highly similar topics, such as *Graphics Optimization* and *Rendering Techniques*, leading to fragmented and redundant labels. Through empirical validation, we set the merging threshold to 0.8, which achieves the best trade-off between topic coherence and granularity. After the first step, we extract $T_{D_{user}} = 43$ and $T_{D_{developer}} = 37$ topics from user and developer datasets, respectively. That is, $T_{D_{user}} = 43$ and $T_{D_{developer}} = 37$ topics are assigned to associated documents, aiming to quantify the document count for each discussion topic. We then obtain the similarities and differences between developers and users, as shown in Figure 3.

4.1.3 Results. Figure 4 shows the 43 topics extracted from the user reviews ranked by review volume. On average, each topic includes 21,582 reviews. Notably, 13 topics exceed this average, indicating key areas of user interest and concern. In addition, Figure 5 shows the distribution of the 37 topics extracted from developer forum discussions ranked by post



Fig. 3. Topic distribution after automated merging using S-BERT semantic analysis and cosine similarity. Blue nodes indicate user-focused topics, yellow indicate developer-focused topics, and green indicate topics of mutual interest. The horizontal axis reflects relative focus (left: user-oriented, right: developer-oriented), and the vertical axis shows variance in discussion prominence between the two groups. Positions are derived from S-BERT embeddings projected via t-SNE.

volume. On average, each topic includes 10,246 discussions. Notably, 12 topics exceed this average, highlighting key areas of developer focus and technical challenges.

Figure 3 displays the shared and distinct focus areas of users and developers, where green circles indicate overlapping topics between game reviews and forum discussions, blue circles represent topics unique to reviews, and orange circles denote those exclusive to forums.

Overlapping Focus. The analysis revealed distinct yet overlapping priorities between the users and developers. For example, in terms of **input methods** (1.8% & 2.8%), users appreciate support for diverse peripherals, e.g., hands on throttle-and-stick HOTAS devices, mouse and keyboard, controllers, and Vive wands (R580062), while developers focus on addressing technical challenges and proposing workarounds for input-related issues (P1182045). In terms of **multiplayer** considerations (2.1% & 2.1%), users advocate for more cooperative experiences (R139393639), and developers focus on troubleshooting implementation errors (P750531).

For **hardware support** (2.6% & 2.1%), users emphasize the need for broader compatibility, including specialized peripherals, e.g., flight hardware (R670366), and developers highlight the impact of hardware compatibility on development efficiency (P96008). In addition, in terms of **target tracking** (2.1% & 2.4%), users report issues associated with tracking specific body parts, e.g., hands and feet (R705224), whereas developers analyze factors like LED positioning to improve accuracy (P140481). Regarding **customization** (3.8% & 6.7%), users express dissatisfaction with limited options

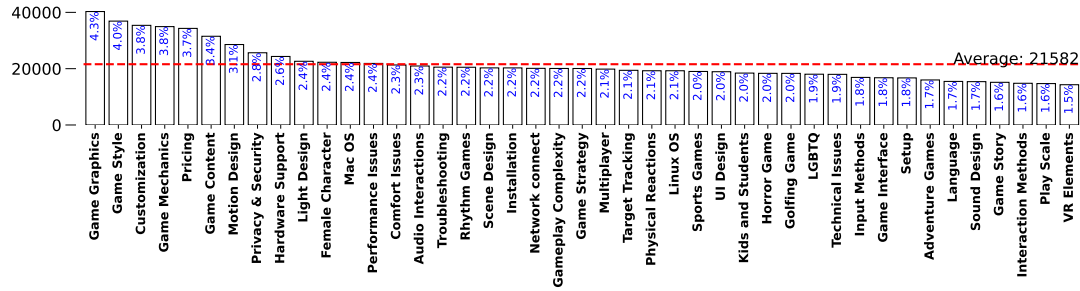


Fig. 4. Topic distribution in D_{user} (43 extracted topics). The X-axis represents the topics, and the Y-axis indicates the corresponding number of user reviews. The red line denotes the average number of reviews across all topics, while the blue percentages indicate the proportion of reviews for each topic.

(R581554), and developers investigate strategies to enhance personalization, e.g., defining hairstyles through arrays (P153285).

Performance (2.4% & 2.8%) was identified as a recurring concern. Users acknowledge that initial issues have been resolved through updates (R157371048), and developers propose optimizations, e.g., tailoring speaker placement based on head profiles (P68101). In addition, **game elements** (2.4% & 2.5%) are critical in terms of player retention, with users praising immersive narratives (R544047), and developers focusing on interaction-driven storytelling over more conventional mechanics (P109287).

Privacy and security (2.8% & 3.3%) are also significant concerns. Users request private spaces in games like *Star Trek VR* (R204994) and raise alarms about excessive data collection (R119697824). Developers address these concerns by restricting raw camera data access and providing alternative methods, e.g., the Wave 5.1.0 SDK, for application development (P50679). In addition, **audio and visual elements** (8.1% & 2.7%) are pivotal for user immersion (R65580), while developers explore integrating modern music into VR environments (P557552). In terms of **interaction elements** (3.1% & 4.4%), well-designed mechanics in games like *Half-Life: Alyx*, *Lone Echo 1*, and *S.T.A.L.K.E.R* enhance engagement considerably (R516719, R337390, R364171), with developers striving to improve the quality of interactions through innovations, e.g., gesture-based control schemes (P55076).

Finding 1: The analysis reveals substantial overlap between user priorities and developer efforts, particularly in areas such as input support, hardware compatibility, tracking, privacy, and immersive features. This alignment indicates a continuous feedback loop, where user experiences and expectations inform technical solutions, and developer innovations in turn shape user perceptions and adoption. Strengthening the interplay between users and developers through transparent communication, responsive design iterations, and a shared focus on emerging challenges can accelerate progress toward more effective and user-aligned VR systems.

User Focus. In addition to the overlapping topics, several user-specific concerns do not appear in the developer discussions. A prominent issue is **LGBTQ+ representation** (1.9%). Some users reported experiencing LGBTQ+-related discrimination, criticizing the VR community management as unfair (R138371, R183044302), while others expressed discomfort with what they perceived as excessive LGBTQ+ promotion in games (R135166). Discussions also address **children, students, and females** (4.4%). Users tend to highlight the lack of VR games tailored to children (R87842) and recognize VR's potential to enhance realistic learning experiences for students (R57457). This aligns with developer

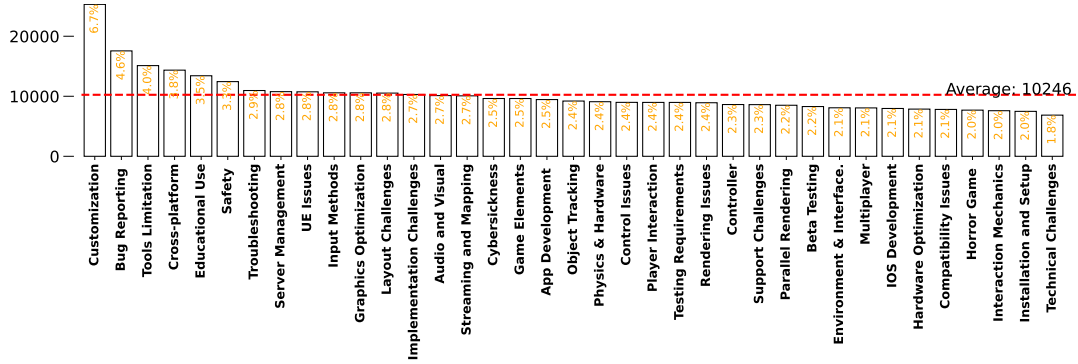


Fig. 5. Topic distribution in $D_{developer}$ (37 extracted topics). The X-axis represents the topics, and the Y-axis indicates the corresponding number of developer posts. The red line denotes the average number of posts across all topics, while the orange percentages indicate the proportion of posts for each topic.

efforts to explore VR applications in educational training (P2478). In addition, users suggest that increasing female character representation and expanding avatar customization could make VR gaming more appealing to female players (R194495, R248678).

Another key concern is **language support** (1.7%). Users criticize the limited language options in VR games, e.g., the absence of Dutch language support (R234861). In addition, pricing issues are frequently debated, with users evaluating game quality relative to cost (R140819314). Finally, discussions focus on **game-specific topics** (4.2%). For example, in golf-related games, users request more diverse course designs (R187950), and in dynamic action games, users point out that movement restrictions, e.g., prolonged standing requirements, create accessibility barriers for players with mobility limitations. To address such issues and improve accessibility, users propose alternative control methods, e.g., alternative button inputs or specific actions (R213644).

Finding 2: User discussions highlight a growing demand for greater inclusivity and accessibility in VR, spanning fair LGBTQ+ representation, child-friendly content, expanded language support, and accommodations for players with diverse physical abilities. Addressing these concerns is essential for fostering a more equitable and engaging VR environment. Future research and development should prioritize inclusive design, diverse content, and accessible interaction mechanisms to broaden participation and ensure that VR experiences are welcoming to all user groups.

Developer Focus. Developer forums highlight several technical challenges encountered during the development process. A key area of discussion is **development tools** (4.0%), as efficient tools significantly enhance productivity (P1073653, P650628). **UE development** (2.8%) is another prominent topic, with developers addressing cross-platform development, environment setup, and interface design. Issues with tools like Unreal Engine and Unity, particularly compatibility problems arising from version updates, are frequently reported (P10239). **Device compatibility** (2.1%) remains a critical concern, with developers proposing solutions such as delaying head pose updates until the camera latches the reference frame to address rendering challenges (P316431). **Server management** (2.8%) is also widely discussed, as larger projects require robust server infrastructure to handle increasing resource demands (P332852).

Developers also face **layout challenges** (2.8%), e.g., design constraints related to Hydra (P70745), and **support challenges** (2.3%), including plugin compatibility issues with hardware, as demonstrated in inquiries about DK2 support

(P235055). Another significant topic is **cybersickness** (2.5%), where developers analyze the factors that contribute to motion discomfort and propose design solutions. Some discussions reference Gavgani's [14] work on VR-induced cybersickness (P585178), and these technical discussions align with user concerns regarding **comfort issues**, e.g., reports of discomfort during VR use (R722882).

Finding 3: Developer discussions reveal a strong focus on overcoming technical barriers in VR, including the adoption of efficient development tools, ensuring device and platform compatibility, managing server resources, and addressing layout and plugin support challenges. Notably, developers are actively engaged in mitigating cybersickness and enhancing user comfort, aligning technical solutions with end-user needs. These insights highlight the necessity for continued innovation in toolchains, robust compatibility frameworks, and user-centered design strategies. Future research and development should prioritize integrated approaches that streamline development workflows, facilitate cross-platform support, and proactively address factors affecting user comfort and health, ultimately promoting the creation of seamless and enjoyable VR experiences.

4.2 (RQ2) How topics evolve over time from different perspectives ?

4.2.1 Motivation. Discussions in the VR domain exhibit unique temporal dynamics because user concerns and developer challenges evolve with emerging technological advancements. For example, user feedback regarding hardware devices may shift with iterations in HMD technology, and developer discussions on technical obstacles may decline as solutions mature. By addressing RQ2, our goal is to investigate how user needs and topics evolve, resulting in a continuously expanding VR community.

4.2.2 Approach. We employ absolute and relative impact metrics in different time periods, which is inspired by Han et al. [16] and Lee et al. [46] to examine dynamic changes in topics within specific fields (e.g., IoT).

We adopt the widely recognized open-coding procedure [9] for the manual inspection and categorization of topics into higher-level categories, as suggested in [46]. For example, topics such as *game graphics* and *UI design*, both related to game design, are categorized under the broader "software" category. This iterative process preserves the distinct characteristics of the dataset while consolidating similar topics into more representative and comprehensive categories. The classification is conducted independently by two authors to ensure objectivity and consistency. The iteration continues until consensus is reached among the evaluators and no further optimization can be made. As a result, all topics are organized into four major categories, as shown in Figure 6.

Calculation of Absolute Impact. First, we apply topic popularity metrics to compute the popularity of a topic z_t in corpus c_j for a post or review d_i , where i represents any topic in corpus c_j . Formally, the popularity of each topic is defined as follows.

$$\text{popularity}(z_t, c_j) = \frac{|d_i|}{|c_j|} : \quad \text{dominant}(d_i) = z_t, \quad 1 \leq i \leq c_j, \quad 1 \leq j \leq T \quad (2)$$

We apply LDA to corpus c_j to obtain a set of T topics (z_1, \dots, z_t) . We express the probability for a specific topic z_t in a post or review d_i as $\theta(d_i, z_t)$ to define the absolute impact metric of a topic z_t in a month m as follows:

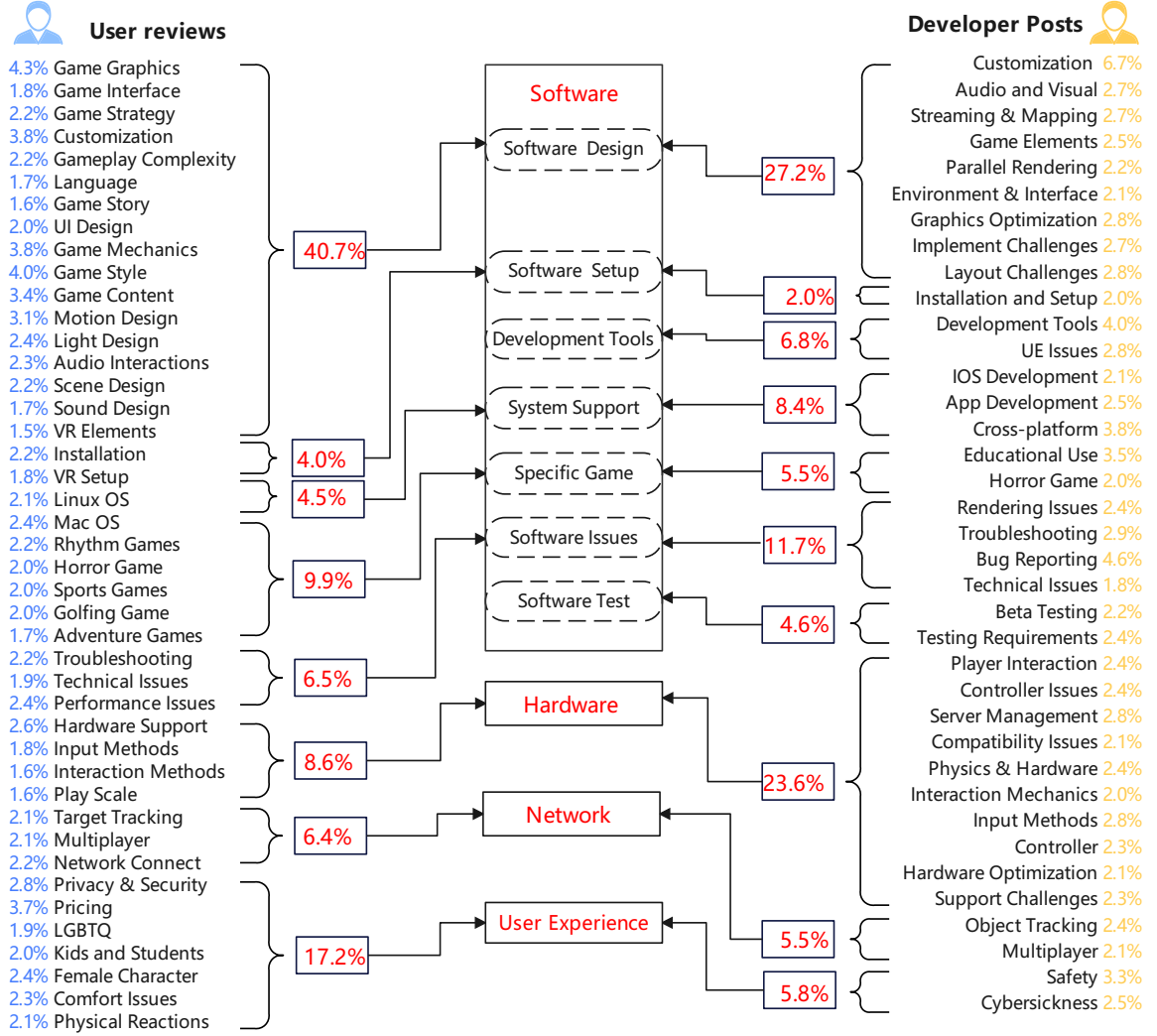


Fig. 6. VR Topics with Categories and Subcategories

$$\text{impact}_{\text{absolute}}(z_t; m) = \sum_{d_i=1}^{D(m)} \theta(d_i; z_t) \quad (3)$$

where $D(m)$ denotes the total number of reviews or posts in month m , and $\theta(d_i; z_t)$ represents the probability that the post or review d_i belongs to topic z_t .

We further extend this definition to calculate the absolute impact of a category C as follows:

$$\text{impact}_{\text{absolute}}(C; m) = \sum_{z_t \in C} \text{impact}_{\text{absolute}}(z_t; m) \quad (4)$$

The category C belongs to four major category of VR topic i.e., Hardware, Software, Network and User Experience.

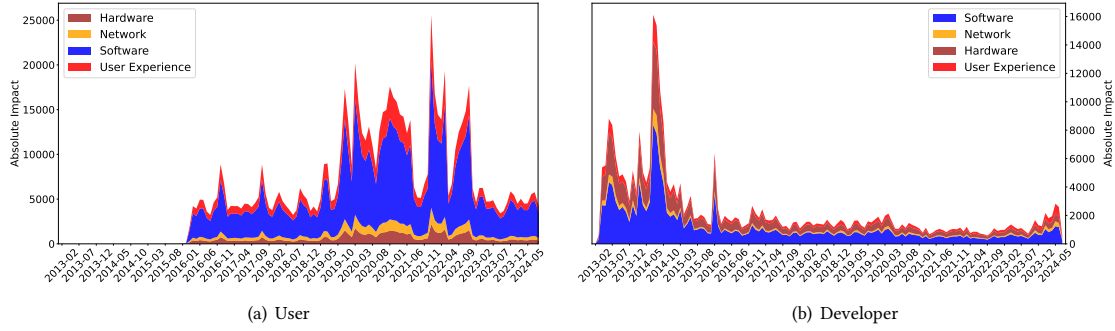


Fig. 7. Temporal trends in absolute impact scores of topic popularity for D_{user} and $D_{developer}$ (2013–2024). The X-axis is divided by month, with labels shown at five-month intervals, and the Y-axis represents the absolute impact scores of topics over time.

Calculation of Relative Impact. The relative impact measures the proportion of a topic z_t relative to all discussions during a specific time. It is defined as:

$$\text{impact}_{\text{relative}}(z_t; m) = \frac{1}{|D(t)|} \sum_{d_i=1}^{D(m)} \theta(d_i; z_t) \quad (5)$$

where $|D(t)|$ represents the total number of posts in a month m . The θ shows the probability of a particular topic z_t for a post or review d_i . The relative impact metric estimates the proportion of posts for a specific topic z_t relative to all posts or reviews in a particular month m .

Similarly, the relative impact of a category C can be calculated as:

$$\text{impact}_{\text{relative}}(C; m) = \sum_{z_t \in C} \text{impact}_{\text{relative}}(z_t; t) \quad (6)$$

where C is the set of posts related to one of the four major categories of posts or reviews.

4.2.3 Results. The impact of VR topics and categories was quantified using Eqs.(2)–(6), enabling dynamic analysis across four categories, i.e., software, hardware, network, and user experience, from January 2013 to July 2024, as shown in Figure 7 and Figure 8. In the following descriptions, we examine trends in absolute and relative impacts for these categories in both user and developer datasets. For absolute impact, we identify key inflection points and major transitions within each category. For relative impact, we highlight interaction patterns and fluctuations among categories, focusing on periods of overlap, divergence, and sudden shifts.

Topic Absolute Impact from User. Figure 7 (a) illustrates the temporal variation of absolute impact across the four major categories from 2015 to 2024. Software-related topics consistently dominate user discussions, maintaining the highest absolute impact throughout the period. While the overall pattern remains stable, a pronounced surge in absolute impact is observed from late 2019 to early 2022, followed by a return to more stable levels.

Software. In **August 2016**, users expressed concerns about visual clarity, with blurry graphics negatively affecting immersion (R502859). By **December 2017**, customization emerged as a central topic, with players requesting greater control over movement mechanics (R674926). In **December 2019**, users shifted focus to game design, calling for more intelligent NPCs and an increased number of levels to enhance engagement (R735575). In **March 2020**, the emphasis turned to content depth, as users demanded expanded gameplay functionalities (R387783). By **January 2021**, users

reported frustrations with the complexity of VR setup processes, citing them as barriers to accessibility (R385168). Customizable shortcuts were emphasized again in **January 2022**, with users highlighting their potential to improve usability and efficiency (R99936). Most recently, in **December 2023**, software update failures became a major concern, as players reported that games became unplayable due to malfunctioning updates (R152968406).

Hardware. In **August 2016**, controller compatibility issues were widely reported, particularly mismatches between hardware and VR software environments (R48028). By **December 2017**, users began requesting broader support for diverse input devices such as steering wheels and flight controllers (R681056). In **December 2019**, platform compatibility concerns emerged, with users criticizing Oculus for lacking cross-platform support, which limited accessibility (R76071). In **March 2020**, high system requirements frustrated users, who found VR inaccessible due to hardware demands (R557044). By **January 2021**, ergonomic design became a focal point, with complaints about hand controllers like the Vive wand being misaligned with natural hand positions (R386869). In **January 2022**, users emphasized limitations in storage capacity on standalone headsets, which affected performance and load times (R752039). Most recently, in **December 2023**, HMD–software compatibility issues resurfaced, as users reported failures when using certain headsets with specific applications (R152859220).

Network. Users initially reported head tracking deficiencies in **August 2016**, highlighting precision and responsiveness issues that significantly impacted immersion (R506471). By **December 2017**, discussions shifted towards multiplayer modes and cross-brand controller connectivity problems (R12218). Network instability, including frequent latency spikes and disconnections, became a critical topic by **December 2019** (R05846). User concerns in **March 2020** emphasized sluggish target-tracking responsiveness, noting that delayed inputs severely disrupted immersive experiences (R562457). Bluetooth convenience and reliability emerged as key discussion points by **January 2021** (R237800). By **January 2022**, users increasingly reported issues related to hand tracking accuracy, expressing frustration over frequent recognition errors (R277222). Most recently, in **December 2023**, users highlighted limitations in camera tracking compatibility with HOTAS devices, affecting gameplay in specialized VR applications (R54298220).

User Experience. Early VR users frequently discussed physiological discomfort, notably VR-induced nausea and dizziness in **August 2016** (R320532). Economic considerations regarding the cost-effectiveness of VR hardware and software became prevalent in **December 2017**, indicating users' concerns about long-term investment value (R505014). By **December 2019**, the increasing participation of younger users led to significant discussions about community management, especially regarding behavior moderation to ensure immersive experiences for all (R324096). By **March 2020**, the inclusivity of environmental settings, particularly for child players, was critically reviewed, with suggestions for more adaptive design (R270975). Social inclusivity became increasingly prominent in **January 2021**, with specific calls to address discriminatory behavior towards LGBTQ+ users in social VR platforms (R140921). Data privacy became a dominant concern by **January 2022**, highlighting debates about ethical data handling and user privacy protections in VR platforms like Facebook (R142345). In the latest discussions of **December 2023**, users advocated for improved physics engines to enhance immersive realism through more precise physical feedback (R154626566).

Finding 4: This study shows a clear shift in user concerns from basic functionality (“Can it work?”) to quality and experience (“Does it work well?”), reflecting higher expectations as VR technology matures. User discussions have moved from feasibility and usability to stability, smooth interaction, social inclusivity, and immersion. For future research and development, it is essential to go beyond basic functionality and create VR systems that are reliable, user-friendly, and inclusive. Developers and researchers should focus on improving software reliability, ergonomic hardware design, and tailored user experiences to support broader adoption and long-term satisfaction.

Topic Absolute Impact from Developer. Figure 7 (b) presents the absolute impact trends across four categories from 2013 to 2024. Software topics overwhelmingly dominate developer discussions, especially during the early years, with a pronounced peak around 2014–2015. Following this surge, the absolute impact across all categories declines sharply and stabilizes at a much lower level, with only minor fluctuations observed in recent years.

Software. Software discussions showed evolving priorities, beginning in **June 2013** with compatibility concerns, exemplified by Descope supporting Rift with Windows (P105980). By **August 2014**, attention shifted towards video playback, particularly DK2’s Video Player functionality (P193258). Rendering issues, notably SpaceEngine in direct mode, became central in **December 2014** (P284695). Audio configuration discussions emerged in **January 2016**, emphasizing stereo versus surround sound preferences (P342270). By **May 2024**, developers tackled compatibility issues like Quest 3 and Unreal Engine 5.3.2 passthrough displays showing black screens (P1193769).

Hardware. Initially, in **June 2013**, high-resolution devices were prioritized (P80843). Compatibility discussions, particularly regarding the Digital Combat Simulator with DK2, gained prominence in **August 2014** (P181179). Attention moved to accessory improvements such as Rift lens replacements in **December 2014** (P284366). GPU optimization discussions were significant by **January 2016** (P377721). Recent hardware discussions in **May 2024** included multimodal interaction challenges with Oculus controllers (P1199093). **Network.** Initial conversations in **June 2013** explored Wi-Fi-based tracking feasibility (P75710). By **August 2014**, developers discussed controller connectivity and latency issues (P217100). Optimization of DK2 connectivity became prominent by **December 2014** (P332607). Tracking accuracy improvements, such as increased LED use in Oculus DK2, dominated in **January 2016** (P324142). Renewed attention to motion capture and body tracking characterized discussions in **May 2024** (P1196371).

User Experience. UX discussions evolved from basic comfort to advanced personalization. Cybersickness mitigation, such as Tai Chi applications, was a focus in **June 2013** (P61346). By **August 2014**, attention turned to screen door effects and visual discomfort (P226422). Interpupillary distance adjustments for motion sickness were central in **December 2014** (P330891). Rotational comfort strategies gained prominence in **January 2016** (P346572). By **May 2024**, discussions shifted towards user-customized VR solutions (P1194513).

Finding 5: VR developer discussions reflect a clear evolution from basic functional concerns ("Will it work?") to sophisticated user-centric enhancements ("How can we make it awesome?"). Trends indicate a significant shift towards software integration, hardware optimization, robust networking, and personalized UX solutions, aligning with technological advances and user expectations. Future research and development should prioritize integrated system approaches, improved connectivity, and highly customized user experiences to foster sustained adoption and satisfaction.

Topic Relative Impact from User. Figure 8(a) presents the relative changes in topic popularity, illustrating the distribution and temporal evolution of each topic within these categories. Since **early 2015**, user interests on VR platforms have experienced rapid differentiation, eventually stabilizing into a structure dominated by **Software** topics, with **User Experience** following closely behind. The **Software** category has consistently maintained approximately **60%** of the relative impact, serving as the core of community discussions, while **User Experience** has remained stable at around **20%–25%**. In contrast, both **Hardware** and **Network** topics have consistently accounted for less than **10%** each, indicating their marginal status within the overall discourse.

During **April 2017**, **October 2017**, **September 2018**, and **August 2021**, the topics of **Hardware** and **Network** repeatedly exhibited convergence and simultaneous surges in user attention. This phenomenon reflects users’ dual focus

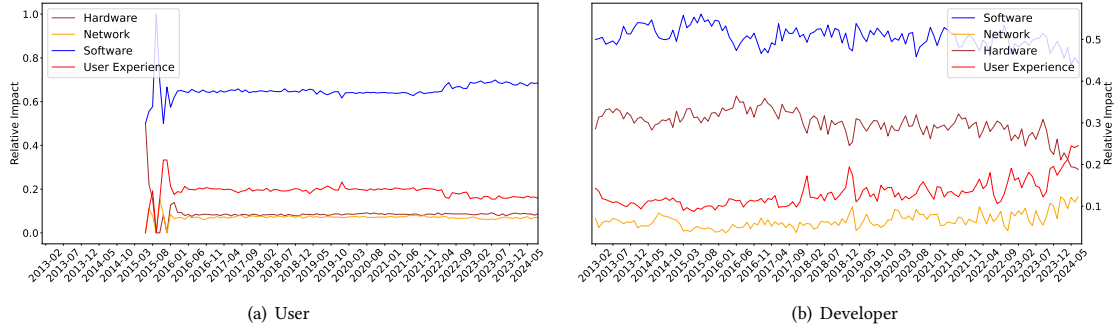


Fig. 8. Temporal trends in relative impact scores of topic popularity for D_{user} and $D_{developer}$ (2013–2024). The X-axis is divided by month, with labels shown at five-month intervals, and the Y-axis represents the relative impact scores, indicating the proportion of attention each topic received over time.

on VR device performance upgrades and the optimization of online experiences. For example, users anticipated new gameplay enabled by enhanced hardware capabilities (e.g., GTX 1060 support, [R465148](#)) and simultaneously showed increased sensitivity and controversy regarding network experiences following platform updates ([R54447](#)). Over time, users placed greater emphasis on system-level synergy among content depth, device capability, and network services, as seen in discussions on desktop mapping and highly immersive interaction ([R332041](#), [R615691](#)). Users frequently called for holistic, end-to-end optimization in response to the tension between hardware advancements and network service limitations. By **August 2021**, a diversified experiential system characterized by community activity, content innovation, and platform collaboration had become the standard for high-quality VR experiences. The integration of **Hardware**, **Network**, and community engagement continued to drive the evolution of the VR ecosystem, as illustrated by discussions on MOD experiences ([R621037](#)) and the sense of social belonging ([R616620](#)).

In **September 2019**, **March 2022**, and **August 2022**, the topics of **User Experience** and **Software** exhibited alternating rises and differentiation. When **User Experience** increased, users tended to raise higher expectations regarding onboarding, operational details, and fairness (e.g., lack of tutorials, [R135539](#)), while the **Software** topic focused on emotional resonance, aesthetics, and enjoyment (e.g., [R736720](#)). With improvements in platform functionality and compatibility, the **Software** topic received positive feedback for technological advancements (such as interface design, visual quality, and functional innovation, [R120065517](#), [R300787](#)); however, when content depth and diversity were lacking, the **User Experience** topic declined due to issues like short storylines or repetitive tasks ([R121041547](#), [R268013](#)).

Entering **2024**, the focal points of **User Experience** and **Software** topics have demonstrated both deep integration and further differentiation. On the **User Experience** front, extremely positive feedback and diverse demands have become mainstream (e.g., “best VR game,” [R167402269](#); “peaceful mode,” [R166799844](#)), alongside a growing user demand for emotion regulation and extreme gameplay experiences. Meanwhile, the **Software** topic features a remarkably high degree of interaction freedom and emotional release (e.g., “sadistic rage,” [R168182744](#)), as well as reflections on moral and social issues (e.g., [R168235987](#)).

Finding 6: The long-term evolution of user interests on VR platforms reveals a persistent emphasis on **software functionality and innovation**, alongside increasing attention to **user experience**, while **hardware** and **network** periodically become prominent during ecosystem upgrades. As user expectations expand from isolated features to integrated “content-hardware-network-community” experiences and emotional value, achieving high-quality VR will require balancing core usability, content richness, emotional engagement, and systemic synergy.

Topic Relative Impact from Developer. Figure 8 (b) illustrates the distribution and temporal evolution of topic popularity across categories. From **2013** to **2024**, developer forum discussions exhibited pronounced cyclical trends and distinct evolutionary phases. **Software** consistently dominated the discourse, accounting for **45%–55%** of attention. **Hardware** initially alternated with software as the primary focus, but after **2015** its prominence stabilized at **25%–35%**. **User Experience** demonstrated continuous growth, rising to approximately **20%** by **2024** and becoming the second most prominent topic. In contrast, **Network** maintained a relatively low profile, generally below **10%**, except for occasional surges during major platform events.

Between **April 2013** and **December 2013**, **Software** topics increased markedly, with developers focusing on the evaluation and implementation of emerging platform tools (e.g., VorpX VR middleware), API integration, and content ecosystems (e.g., Mate), which stimulated vibrant community discussion (P45260, P94111, P132473). In contrast, interest in **Hardware** declined, centering on the performance, peripheral integration, and compatibility of headsets such as the Rift (P50428, P110594, P128752), indicative of a stabilization in hardware development. **User Experience** topics also decreased, with discussions on gesture tracking, embodied interaction, and environmental adaptation (e.g., Leap Motion, P109799, P34733, P27161) becoming more diffuse as interaction paradigms matured. **Network** remained the least discussed, with only brief surges during device adaptation and low-level debugging (e.g., Rift network configuration across different computers and resolutions, P35154, P40218, P114457), further subsiding as technical challenges were addressed.

From **June 2014** to **December 2015**, the release of next-generation hardware platforms (such as Oculus CV1) and advancements in core parameters like FoV and resolution became the main drivers of heightened developer interest. The prominence of **Hardware** topics increased sharply, with community discussions focusing on device compatibility, peripheral integration, and proof-of-concept validation (P160496, P186772, P193366). Meanwhile, **Software** topics gradually stabilized, as developers concentrated on content innovation and ecosystem compatibility (e.g., new versions of the Mate platform, P255485, P227487, P145477); the toolchain matured, and the pace of innovation slowed slightly. **User Experience** discussions centered on flagship titles (such as Oculus and Half-Life 2), with a focus on immersive content and product feedback, though overall attention declined (P197393, P332755, P320324), marking a community shift towards hardware-centric concerns. **Network** topics primarily addressed community management, authentication, and account system development, with relatively limited engagement (P224722, P242079, P148747).

In **October 2018**, ecosystem synergy and diversified content innovation emerged as new focal points within the community. The prominence of **Network** and **User Experience** topics increased significantly, with developers actively discussing multi-location collaboration, content synchronization, and the launch of new services on platforms such as Oculus and SteamVR, thereby advancing cross-platform interoperability (P731934, P676749, P696682). On the **User Experience** front, content uploading, device compatibility, and community events became central concerns, and there was a clear rise in demands for ecosystem openness and interactive innovation (e.g., Google Drive integration, device list upgrades, P21382, P692004, P676755). In contrast, attention to **Software** and **Hardware** topics waned, with community

discussion focusing more on routine issues such as APIs, content distribution channels, and device crash resolution (P21317, P21335, P21147, P21623, P698394).

In **April 2022**, focus on **User Experience** topics declined, as developers concentrated on controller comfort, interaction commands, and scenario details for devices like VIVE (e.g., Vive palm pinch, peripheral evaluations, P47320, P47150), reflecting mature basic interactions and device adaptation. The lack of major breakthroughs led to reduced discussion intensity. By **November 2022**, **User Experience** surpassed hardware as the second most discussed topic. **Hardware** discussions addressed bug fixes, display issues, and compatibility (P50113, P50034, P998953), while **User Experience** centered on detail optimization, third-party integration, and community Q&A (e.g., Pimax compatibility, soundscape commands, P995868).

From **September 2022** to **June 2024**, community discussions developed a dual-core structure centered on **Software** innovation and **User Experience**. **Software** topics focused on custom development and multi-platform integration for emerging platforms such as OpenXR, Unity, and SteamVR APIs (P170249942, P170126759), with engine features and content distribution as mainstream concerns. **User Experience** continued to rise, encompassing multi-scenario adaptation, interaction optimization, and community support for devices like Meta Quest and Pimax (P999795, P994961, P995868). **Hardware** topics exhibited periodic peaks, mainly related to the release of new HMD models, driver compatibility, and module updates (P169583019, P169755139, P170004160). While **Network** remained less prominent, it showed notable fluctuations during cross-platform synchronization, cloud services, and collaborative node discussions (P999199, P170054170, P995949).

Finding 7: VR developers' focus has shifted from **hardware** and **network** adaptation to **software innovation** and **user experience**. These have become the main drivers of ecosystem growth, especially with increased attention to **multi-platform integration** and **personalized experiences**. While **hardware** and **network** receive less ongoing focus, they regain prominence with major device launches. Sustained community vitality relies on **ecosystem openness, content innovation, user feedback, and multi-platform connectivity**.

5 Implication

This section is organized into two parts: the first highlights methodological innovations and their broad applicability, while the second offers actionable recommendations for developers and researchers to address key gaps identified within the evolving VR ecosystem.

5.1 Methodological Implications: Scalable and Interpretable Text Analytics

The results verify that the proposed HTModel is an effective and scalable approach for large-scale textual analysis in software engineering. First, HTModel achieves higher label precision, greater interpretability, and lower annotation costs than LLM-only pipelines by decoupling clustering from LLM-based semantic refinement. Second, its multidimensional evaluation framework, integrating coherence, coverage, diversity, label accuracy, and label usefulness, ensures robust topic extraction and strong domain adaptability across heterogeneous real-world datasets. Third, HTModel can be extended with neural, hierarchical, or domain-adaptive models to support multilingual and cross-platform contexts, effectively capturing the evolving complexity of software ecosystems.

5.2 Practical Implications: Enhancing User Experience through Multi-Perspective Insights

Our empirical findings highlight critical areas where targeted efforts from both developers and researchers can significantly improve the VR ecosystem, enhancing user satisfaction, inclusivity, and long-term platform sustainability. Integrating user and developer perspectives uniquely reveals previously unrecognized challenges, providing actionable guidance extending beyond conventional single-stakeholder studies.

Ensure inclusive and adaptive design. User discussions consistently emphasize the necessity for adaptive interfaces accommodating diverse physical and cognitive abilities. Developers should:

- Implement adaptive interfaces that support varying accessibility requirements, reducing barriers for users with physical or cognitive impairments. *For example, in VR fitness games, offer configurable control schemes and gesture sensitivity settings so that both wheelchair users and standing players can fully participate.*
- Enable extensive avatar personalization options to authentically represent diverse user identities across cultural, gender, and personal expression dimensions. *For example, allow players to customize avatars with cultural attire, gender-neutral features, or assistive devices such as hearing aids or prosthetic limbs.*
- Develop environments specifically designed for vulnerable or underrepresented groups such as children, elderly individuals, and persons with disabilities. *For example, design child-friendly VR classrooms with simplified navigation and content filtering for age-appropriate learning.*

Enhance moderation and safety frameworks. Persistent user concerns about inadequate moderation necessitate proactive safety measures. Developers are advised to:

- Embed intuitive reporting mechanisms directly into VR platforms to facilitate user-driven moderation. *For example, integrate a “Report” button into the user interface that can be activated via a quick gesture or voice command during interactions.*
- Integrate automated moderation technologies capable of real-time detection and intervention against harassment or inappropriate content. *For example, use AI-based voice chat monitoring to detect hate speech instantly and mute the offending player until reviewed by a moderator.*
- Clearly articulate and rigorously enforce community guidelines to sustain respectful, supportive interactions among users. *For example, display a brief but visible code of conduct whenever users join a multiplayer session.*

Optimize performance equitably. Performance inequities stemming from hardware disparities create accessibility barriers. Developers should:

- Prioritize performance optimization across a diverse spectrum of devices, including mid- to low-tier VR headsets. *For example, include a “performance mode” that automatically reduces rendering load while maintaining gameplay mechanics for budget devices.*
- Conduct regular assessments of feature impacts on rendering performance, latency, and overall usability. *For example, run automated benchmark tests for every major update to ensure frame rates remain stable across supported hardware.*
- Refrain from over-reliance on high-end hardware specifications when developing new functionalities. *For example, ensure that newly introduced visual effects, such as advanced particle simulations, have simplified fallback versions for less powerful systems.*

5.2.1 Implications For Researchers. Examine user–developer alignment longitudinally. Researchers must investigate the dynamic interactions between user expectations and developer initiatives. Specifically, researchers are encouraged to:

- Conduct longitudinal studies incorporating multi-stakeholder analyses to expose subtle divergences that single-perspective approaches might miss. *For example, track how requests for accessibility features from users are addressed (or ignored) across multiple software release cycles.*
- Facilitate early detection of user–developer mismatches, thereby informing timely, evidence-driven interventions. *For example, apply topic modeling quarterly to compare developer forum priorities with user reviews, flagging mismatches in content focus.*
- Leverage comprehensive datasets to systematically track gaps between inclusive design demands and developer technical priorities. *For example, map keyword clusters from both sides to quantify underrepresented user demands.*

Explore the role of inclusivity in user retention. Sustained user engagement is closely linked to inclusivity in design. Researchers are encouraged to:

- Investigate how inclusive design elements such as diverse avatars, accessible interactions, and moderated environments affect long-term user engagement and loyalty. *For example, compare retention metrics between platforms that offer inclusive avatar customization and those that do not.*
- Utilize mixed-method and longitudinal research designs to uncover causal relationships between inclusive practices and sustained user retention. *For example, combine large-scale user log data with follow-up surveys to understand why users remain active on certain inclusive platforms.*

Assess intervention efficacy rigorously. To bridge identified user–developer gaps, we highlight the need for:

- Implement controlled empirical evaluations of strategies designed to address gaps, using structured feedback mechanisms, inclusive quality benchmarks, and participatory design methodologies. *For example, run an A/B test comparing standard onboarding with an adaptive onboarding flow for users with different interaction needs.*
- Identify and validate best practices for adaptability and inclusivity in VR through systematic, evidence-based testing. *For example, develop a public benchmark dataset of inclusive design implementations and their observed impact on user satisfaction.*

Moreover, researchers should explore emerging topics such as the ethical implications of data privacy in immersive environments and the psychological impacts of prolonged VR exposure, areas that surfaced prominently in user discussions but remain underexplored academically. By addressing these topics, future research can contribute substantially to the ethical and sustainable advancement of VR technology.

6 Threats to Validity

6.1 Internal validity

The main internal threat relates to the reliability of the proposed HTModel, i.e., the effectiveness of topic assignment and label consistency. HTModel integrates efficient clustering algorithms with LLM-based semantic understanding, reducing manual annotation compared to traditional clustering models and lowering annotation costs compared to pure LLM-based approaches. HTModel provides high precision, interpretability, and efficiency. However, the number of evaluators and the topics that occur less frequently are easily ignored. For example, although a comprehensive evaluation framework is applied, integrating quantitative metrics such as coherence, coverage, and diversity, and

expert-based qualitative assessments such as label accuracy and usefulness, the qualitative evaluation involves only two evaluators.

To address these limitations, future work should consider integrating domain-adaptive or hierarchical models and semi-supervised learning approaches to improve topic boundary definition and adapt to evolving vocabularies.

6.2 External validity

The main threat to external validity concerns the generalizability of our findings beyond the current datasets and domain. Our analysis is based on textual data from three major VR gaming platforms and their associated developer forums. While these platforms cover a substantial portion of the VR ecosystem, they do not capture the full range of VR applications, such as those in education, healthcare, or industrial settings. Additionally, discussions from alternative channels, for example social media platforms or private repositories, are not included and may present different topical distributions or community concerns. The study is also limited to English-language data, with more than 100,000 non-English entries excluded during preprocessing. This language restriction may lead to the omission of region-specific issues and culturally distinct perspectives, as highlighted in previous work [11, 13, 30].

To address these limitations, future research should include a wider range of platforms, application domains, and languages to better assess the generalizability and applicability of the proposed method and findings.

7 Conclusion

This decade-spanning multi-stakeholder analysis of the VR ecosystem reveals both convergences and enduring disparities in priorities between users and developers. By examining 944,320 user reviews alongside 389,477 developer forum posts from Meta, SteamVR, and Viveport, the study finds that certain technical concerns—most notably performance optimization and input methods—consistently emerge as shared focal points. At the same time, critical issues surrounding inclusivity, user safety, and community experience remain markedly under-addressed by developers despite persistent and vocal user concern. Temporal trend analysis further exposes an asynchronous evolution of stakeholder interests: over the years, user discourse has increasingly gravitated toward emotional, inclusive, and social dimensions of VR engagement, whereas developer discussions have remained largely anchored in backend technical improvements and platform integration efforts. This misalignment in focus is not merely anecdotal but empirically evident, underscoring a substantive gap between what users value in emerging VR experiences and where developers are investing their effort.

Beyond these empirical insights, the study makes a broader scholarly contribution through its introduction of a multi-view analytical framework grounded in a novel hybrid topic modeling method (HTModel). By integrating classical topic modeling with a state-of-the-art language model, HTModel achieved high topic-labeling accuracy (exceeding 80%) while significantly reducing the need for manual annotation, thus demonstrating a cost-effective approach to large-scale text analytics. Notably, compared to fully LLM-based topic modeling pipelines, HTModel reduces computational cost by over 99% while maintaining high interpretability and scalability. The resulting framework and the curated cross-platform dataset are made available for reuse, providing valuable tools for further research on software ecosystems.

Overall, The findings and methodological advances presented here not only deepen our understanding of user–developer dynamics in VR but also offer a generalizable approach for examining multi-perspective data in other domains. Looking forward, this work opens important avenues for future investigation, such as extending analyses to non-English and region-specific content, incorporating additional stakeholder groups (e.g., platform designers or community moderators), and developing intervention strategies to better synchronize developer priorities with the inclusive, community-oriented expectations of users.

References

- [1] Dimo Angelov. 2020. Top2vec: Distributed representations of topics. *arXiv preprint arXiv:2008.09470* (2020).
- [2] Amirhossein Ashtari, Michael Nebeling, and Moritz Speicher. 2023. An Empirical Study on Current Practices and Challenges of Core AR/VR Developers. *ACM CHI Workshop on AR/VR Development* (2023).
- [3] Ibrahim Al Azher, Venkata Devesh Reddy Seethi, Akhil Pandey Akella, and Hamed Alhoori. 2024. Limtopic: Llm-based topic modeling and text summarization for analyzing scientific articles limitations. In *Proceedings of the 24th ACM/IEEE Joint Conference on Digital Libraries*. 1–12.
- [4] Anton Barua, Stephen W Thomas, and Ahmed E Hassan. 2014. What are developers talking about? An analysis of topics and trends in Stack Overflow. *Empirical Software Engineering* 19, 3 (2014), 619–654.
- [5] David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent dirichlet allocation. *Journal of machine Learning research* 3, Jan (2003), 993–1022.
- [6] Jim Buchan, Muneera Bano, Didar Zowghi, Stephen MacDonell, and Amrita Shinde. 2021. Alignment of Stakeholder Expectations about User Involvement in Agile Software Development. *Empirical Software Engineering* 26, 3 (2021), 1275–1312.
- [7] Chia-Hsuan Chang, Jui-Tse Tsai, Yi-Hang Tsai, and San-Yih Hwang. 2025. LITA: An Efficient LLM-Assisted Iterative Topic Augmentation Framework. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 449–460.
- [8] Jonathan Chang, Sean Gerrish, Chong Wang, Jordan Boyd-Graber, and David Blei. 2009. Reading tea leaves: How humans interpret topic models. *Advances in neural information processing systems* 22 (2009).
- [9] John W Creswell and Cheryl N Poth. 2016. *Qualitative inquiry and research design: Choosing among five approaches*. Sage publications.
- [10] Jiong Dong, Kaoru Ota, and Mianxiong Dong. 2024. User Experience of Different Groups in Social VR Applications: An Empirical Study Based on User Reviews. *IEEE Transactions on Computational Social Systems* (2024).
- [11] Jiong Dong, Kaoru Ota, and Mianxiong Dong. 2024. What Are the Points of Concern for Players about VR Games: An Empirical Study based on User Reviews in Different Languages. *ACM Games: Research and Practice* 2, 4 (2024), 1–18.
- [12] Roman Egger and Joanne Yu. 2022. A topic modeling comparison between lda, nmf, top2vec, and bertopic to demystify twitter posts. *Frontiers in sociology* 7 (2022), 886498.
- [13] Rain Epp, Dayi Lin, and Cor-Paul Bezemer. 2021. An empirical study of trends of popular virtual reality games and their complaints. *IEEE Transactions on Games* 13, 3 (2021), 275–286.
- [14] Alireza Mazloumi Gavani, Keith V Nesbitt, Karen L Blackmore, and Eugene Nalivaiko. 2017. Profiling subjective symptoms and autonomic changes associated with cybersickness. *Autonomic Neuroscience* 203 (2017), 41–50.
- [15] Maarten Grootendorst. 2022. BERTopic: Neural topic modeling with a class-based TF-IDF procedure. *arXiv preprint arXiv:2203.05794* (2022).
- [16] Junxiao Han, Emad Shihab, Zhiyuan Wan, Shuiguang Deng, and Xin Xia. 2020. What do programmers discuss about deep learning frameworks. *Empirical Software Engineering* 25 (2020), 2694–2747.
- [17] Nitin Hardeniya, Jacob Perkins, Deepti Chopra, Nisheeth Joshi, and Iti Mathur. 2016. *Natural language processing: python and NLTK*. Packt Publishing Ltd.
- [18] Khalid Hasan, Partho Chakraborty, Rifat Shahriyar, Anindya Iqbal, and Gias Uddin. 2021. A Survey-Based Qualitative Study to Characterize Expectations of Software Developers from Five Stakeholders. *ACM Symposium on Application Performance Engineering* (2021).
- [19] Safwat Hassan, Wei Wang Shang, Chakkrit Tantithamthavorn, Cor-Paul Bezemer, and Ahmed E. Hassan. 2018. Studying the dialogue between users and developers of free apps in the Google Play store. *Empirical Software Engineering* 23, 3 (2018), 1275–1312.
- [20] HTC Vive Developer Forums. 2024. HTC Vive Developer Forums. Retrieved Jul 2024 from <https://forum.htc.com/forum/24-vive-developer-forums/>.
- [21] Arthur Kamienski and Cor-Paul Bezemer. 2021. An empirical study of Q&A websites for game developers. *Empirical Software Engineering* 26, 6 (2021), 115.
- [22] Satya Kapoor, Alex Gil, Sreyoshi Bhaduri, Anshul Mittal, and Rutu Mulkar. 2024. Qualitative insights tool (qualit): Llm enhanced topic modeling. *arXiv preprint arXiv:2409.15626* (2024).
- [23] Sai Anirudh Karre, Neeraj Mathur, and Y. Raghu Reddy. 2019. Is Virtual Reality Product Development Different? An Empirical Study on VR Product Development Practices. In *ISEC '19*. 1–11.
- [24] Hammad Khalid, Emad Shihab, Meiyappan Nagappan, and Ahmed E. Hassan. 2015. What Do Mobile App Users Complain About? *IEEE Software* 32, 3 (2015), 70–77. doi:10.1109/MS.2014.50
- [25] Diego Kozłowski, Carolina Pradier, and Pierre Benz. 2024. Generative AI for automatic topic labelling. *arXiv preprint arXiv:2408.07003* (2024).
- [26] Per Lenberg, Robert Feldt, and Lars Göran Wallgren Tengberg. 2018. Misaligned Values in Software Engineering Organizations. *arXiv preprint arXiv:1810.06104* (2018).
- [27] Chen Li, Yixiao Ge, Jiayong Mao, Dian Li, and Ying Shan. 2023. Taggpt: Large language models are zero-shot multimodal taggers. *arXiv preprint arXiv:2304.03022* (2023).
- [28] Shuqing Li, Lili Wei, Yepang Liu, Cuiyun Gao, Shing-Chi Cheung, and Michael R. Lyu. 2023. Towards Modeling Software Quality of Virtual Reality Applications from Users’ Perspectives. *arXiv preprint arXiv:2308.06783* (2023).
- [29] Dayi Lin, Cor-Paul Bezemer, Ying Zou, and Ahmed E Hassan. 2019. An empirical study of game reviews on the Steam platform. *Empirical Software Engineering* 24 (2019), 170–207.
- [30] Yijun Lu, Kaoru Ota, and Mianxiong Dong. 2024. An Empirical Study of VR Head-Mounted Displays Based on VR Games Reviews. *ACM Games: Research and Practice* 2, 3 (2024), 1–20.

- [31] Wolfgang Mauerer, Mitchell Joblin, Damian Tamburri, Carlos Paradis, Rick Kazman, and Sven Apel. 2021. In Search of Socio-Technical Congruence: A Large-Scale Longitudinal Study. *arXiv preprint arXiv:2105.08198* (2021).
- [32] Wes McKinney et al. 2011. pandas: a foundational Python library for data analysis and statistics. *Python for high performance and scientific computing* 14, 9 (2011), 1–9.
- [33] Meta. 2024. Meta Official Website. Retrieved Jul 2024 from <https://www.meta.com/>.
- [34] Meta Community Forums. 2024. Meta Community Forums. Retrieved Jul 2024 from <https://communityforums.atmeta.com/>.
- [35] Yudi Permana, Arvita Emarilis, et al. 2021. Stemming analysis indonesian language news text with Porter algorithm. In *Journal of Physics: Conference Series*, Vol. 1845. IOP Publishing, IOP Publishing, 012019.
- [36] Chau Minh Pham, Alexander Hoyle, Simeng Sun, Philip Resnik, and Mohit Iyyer. 2023. Topicgpt: A prompt-based topic modeling framework. *arXiv preprint arXiv:2311.01449* (2023).
- [37] Daniel Ramage, Christopher D Manning, and Susan Dumais. 2011. Partially labeled topic models for interpretable text mining. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*. 457–465.
- [38] N Reimers. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. *arXiv preprint arXiv:1908.10084* (2019).
- [39] Gema Rodríguez-Pérez, Reza Nadri, and Meiyappan Nagappan. 2021. Perceived diversity in software engineering: a systematic literature review. *Empirical Software Engineering* 26, 5 (2021), 102.
- [40] Dhia Elhaq Rzig, Nafees Iqbal, Isabella Attisano, Xue Qin, and Foyzul Hassan. 2022. Characterizing Virtual Reality Software Testing. *arXiv preprint arXiv:2211.01992* (2022).
- [41] Angelo Singh and Joseph O’Hagan. 2024. Exploring Topic Modelling of User Reviews as a Monitoring Mechanism for Emergent Issues Within Social VR Communities. *arXiv preprint arXiv:2406.03989* (2024).
- [42] Peter M. Stahl. 2020. Language Detection. Retrieved Jul 2, 2022 from <https://github.com/pemistahl/lingua>.
- [43] Steam Community - App 250820. 2024. Steam Community Discussions for VR. Retrieved Jul 2024 from <https://steamcommunity.com/app/250820>.
- [44] Steam VR Store. 2024. Steam VR Store Page. Retrieved Jul 2024 from <https://store.steampowered.com/vr/>.
- [45] VIVE Team. 2024. Virtual Reality Gaming Market Size, Share, Trends and Forecast by Segment, Device, Age Group, Type of Games, and Region, 2025-2033. <https://www.imarcgroup.com/virtual-reality-gaming-market/>. Accessed: 2025-07-23.
- [46] Gias Uddin, Fatima Sabir, Yann-Gaël Guéhéneuc, Omar Alam, and Foutse Khomh. 2021. An empirical study of iot topics in iot developer discussions on stack overflow. *Empirical Software Engineering* 26 (2021), 1–45.
- [47] L Vergni, F Todisco, and B Di Lena. 2021. Evaluation of the similarity between drought indices by correlation analysis and Cohen’s Kappa test in a Mediterranean area. *Natural Hazards* 108, 2 (2021), 2187–2209.
- [48] Viveport. 2024. Viveport Official Website. Retrieved Jul 2024 from <https://www.viveport.com/>.
- [49] Mengting Wan, Tara Safavi, Sujay Kumar Jauhar, Yujin Kim, Scott Counts, Jennifer Neville, Siddharth Suri, Chirag Shah, Ryen W White, Longqi Yang, et al. 2024. Tnt-llm: Text mining at scale with large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 5836–5847.
- [50] Yawen Zhang, Han Zhou, Zhoumingju Jiang, and Qinyuan Lei. 2025. Exploring Viewing Modalities in Cinematic Virtual Reality: A Systematic Review and Meta-Analysis of User Experience. *Proceedings of the ACM on Human-Computer Interaction* (2025).

A Appendix

This appendix provides a detailed description of the evaluation procedures and parameter selection process for the proposed HTModel.

A.1 Model Evaluation

To evaluate the performance of the proposed HTModel for topic extraction in VR textual data, we design a comprehensive evaluation framework that combines both quantitative and qualitative evaluations. Here, the quantitative evaluation assesses the clustering effectiveness, and the qualitative evaluation examines the labeling quality.

A.1.1 Quantitative Evaluation. We quantitatively evaluate the efficiency of conventional topic models in topic clustering using three key metrics, i.e., topic coherence, coverage, and diversity. The results are presented in Figure 9 (a).

Topic coherence is a critical metric used to assess the semantic coherence and meaningfulness of the topics generated by a topic model. We calculate topic coherence using the CoherenceModel class in Gensim, a widely used open-source Python library for topic modeling and document similarity analysis. The formula used to calculate topic coherence is

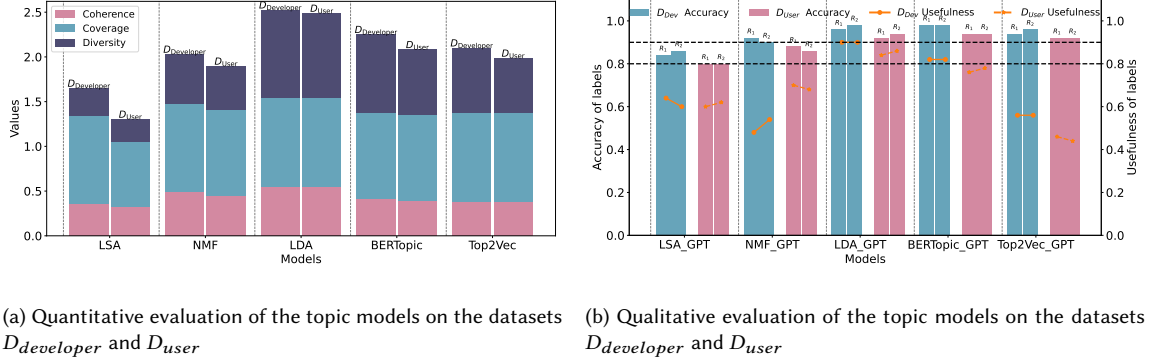


Fig. 9. Quantitative (a) and qualitative (b) evaluation of the topic models on the datasets $D_{developer}$ and D_{user} .

expressed as follows:

$$Coh_m = \frac{1}{T} \sum_{t=1}^T \sum_{i=1}^{|W_t|} \sum_{j=i+1}^{|W_t|} s(w_{ti}, w_{tj}) \quad (7)$$

where $m \in [1, 5]$ represents the model index, $t \in [1, 10]$ denotes the topic index, and W_t denotes the set of words in topic t . Here, w_{ti} and w_{tj} are the i -th and j -th words in topic W_t , respectively, and $s(w_{ti}, w_{tj})$ measures the similarity between these words, which is typically computed using word co-occurrence statistics.

Topic coverage quantifies the proportion of documents covered by the generated topics. Here, a higher coverage rate indicates that the topics are more representative of the dataset. The formula used to calculate topic coverage is expressed as follows:

$$Cov_m = \frac{1}{D} \sum_{d=1}^D \mathbb{I} \left(\sum_{t=1}^T \mathbb{I}(P_{dt} > 0) > 0 \right) \quad (8)$$

where D denotes the total number of documents in the dataset. The topic coverage is determined by checking whether each document d is associated with at least one topic, as indicated by the probability P_{dt} of document d belonging to topic t being greater than 0. The indicator function $\mathbb{I}(P_{dt} > 0)$ evaluates this condition.

Topic diversity measures the distinctiveness of the topics generated by a model, where higher diversity values indicate that the topics cover a broader range of content and aspects in the dataset. The formula used to calculate topic diversity is expressed as follows:

$$Div_m = \frac{1}{T} \sum_{t=1}^T \frac{|Unique(W_t)|}{|W_t|} \quad (9)$$

where $|Unique(W_t)|$ represents the number of unique words in topic t , and $|W_t|$ is the total number of words in topic t . Here, a higher ratio of unique words to total words indicates greater topic diversity.

A.1.2 Qualitative Evaluation. We manually evaluate the labeling quality generated by GPT-4o using two key metrics: *label accuracy* and *label usefulness*. While our evaluation procedure is inspired by methodologies from previous studies [25, 27, 49], both metrics are uniquely defined in this work to better reflect the specific requirements of our research.

This section provides a comprehensive assessment of the labels produced by the five models, with the results visualized in Figure 9 (b). The details of our evaluation process and metric definitions are provided below.

Label accuracy measures how well the assigned label can be directly inferred from the keyword group and is categorized into three levels:

$$Acc_t = \begin{cases} 1, & \text{Bad (Unrelated or misleading relative to the keyword group.)}, \\ 2, & \text{Acceptable (Captures part of the meaning but lacks clarity or completeness.)}, \\ 3, & \text{Perfect (Clearly and comprehensively represents the keyword group.)}. \end{cases} \quad (10)$$

Label usefulness evaluates the research significance of the label and is classified into the following three levels.

$$Use_t = \begin{cases} 1, & \text{Low (Does not contribute meaningful value or insight)}, \\ 2, & \text{Good (Provides some value or partial contribution to the research)}, \\ 3, & \text{High (Offers significant value and demonstrates clear innovation or research impact)}. \end{cases} \quad (11)$$

A.1.3 Evaluation Process. The qualitative evaluation was independently conducted by two experts, both of whom hold PhDs in Software Engineering and currently serve as senior UX designers at leading VR companies. Prior to assigning scores, the experts discussed and standardized the evaluation criteria, clarifying VR-related terminology and the interpretation of each rating level to ensure consistent and objective assessments.

For each topic label generated by the LDA_GPT, LSA_GPT, NMF_GPT, BERTopic_GPT, and Top2Vec_GPT models, both experts independently rated both *accuracy* and *usefulness*. For instance, in terms of accuracy, a label such as “Gameplay Complexity” assigned to a cluster containing keywords like “challenge, difficulty, progression, level, strategy” was rated as *Perfect* (score 3), since it precisely captures the semantic content of the cluster. In contrast, a label like “Miscellaneous” applied to a diverse set of unrelated keywords was rated as *Bad* (score 1), due to its lack of specificity and interpretability. Labels such as “Customization” that sufficiently reflected most but not all of the cluster’s meaning were scored as *Acceptable* (score 2).

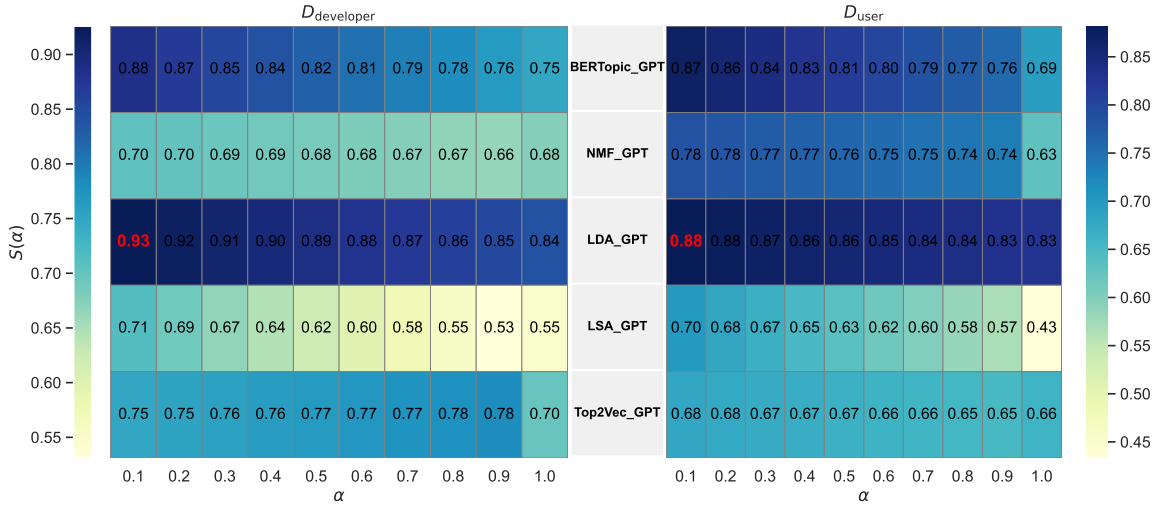
Similarly, for usefulness, a label such as “Technical Issues” was considered *Low* (score 1) as it is too broad for actionable development, while “Audio and Visual” was considered *Good* (score 2), offering partial guidance but lacking specificity. In contrast, labels such as “Privacy & Security” were assigned *High* (score 3), as they are immediately actionable for design or engineering decisions.

To assess the consistency of the evaluation results, we employ *Cohen’s Kappa coefficient* [47], which measures inter-rater reliability on a scale of [0, 1], where higher values indicate stronger agreement. The results for dataset $D_{developer}$ are $K_{LDA} = 0.83$, $K_{BERTopic} = 0.80$, $K_{Top2Vec} = 0.85$, $K_{LSA} = 0.84$, and $K_{NMF} = 0.87$. In addition, the results for dataset D_{user} are $K_{LDA} = 0.81$, $K_{BERTopic} = 0.82$, $K_{Top2Vec} = 0.84$, $K_{LSA} = 0.83$, and $K_{NMF} = 0.85$. These values indicate strong agreement between the raters.

To quantify the metrics, we calculate the proportion of labels with scores ≥ 2 for both accuracy and usefulness, transforming them into a normalized scale of 0–1. Here, values closer to 1 indicate higher label quality, where 1 represents perfect accuracy or usefulness. This approach enables an objective comparison of the models’ performance in generating meaningful and precise labels.

Table 2. Comparison of Cost, Human Involvement, and Interpretability for Major Topic Modeling Pipelines

Method	LLM Usage	Cost (1M docs)	Human Effort	Interpretability
LDA/LSA/NMF	×	-	High (manual labeling)	Low (keywords only)
BERTopic/Top2Vec	×	-	High (manual labeling)	Low (keywords only)
TopicGPT	✓	High (\$2,700–\$10,800)	Low (validation)	High
HTModel	✓	Low (\$1–\$10)	Low (validation)	High

Fig. 10. Distribution of Comprehensive Scores (S) for $D_{\text{developer}}$ and D_{user} Across Topic Models and Weight Parameters (α)

A.1.4 Baseline Comparison. Recent advances such as TopicGPT [36] demonstrate the potential of large language models (LLMs) for fully automated, highly interpretable topic modeling. Here, we compare the computational cost, human effort, and interpretability among traditional topic models (LDA, LSA, NMF, BERTopic, Top2Vec), end-to-end LLM-based frameworks (e.g., TopicGPT) and hybrid pipelines (traditional topic model clustering + LLM-based labeling, as in our HTModel).

Table 2 summarizes key differences. Traditional models are efficient and require almost no computational cost, but manual intervention is needed to interpret and label topics, which is labor-intensive and subjective. Hybrid approaches like HTModel automate topic extraction and use LLMs to label a limited set of topics, resulting in very low cost (e.g., labeling 50–100 topics via GPT-4o costs less than \$ 10 for over one million documents) and much greater interpretability, with minimal human involvement required for validation. In contrast, end-to-end LLM frameworks such as TopicGPT offer the highest interpretability, automatically generating natural language topics and assignments, but at much higher computational cost (e.g., \$88–\$155 for tens of thousands of documents), since LLMs are called repeatedly for generation, assignment, and refinement.

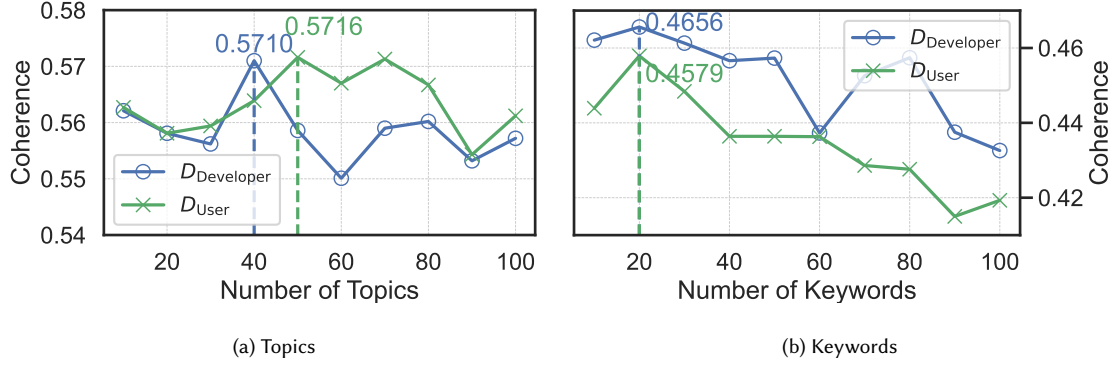


Fig. 11. Coherence scores versus the number of topics (a) and keywords (b) for $D_{developer}$ and D_{user} .

A.2 Model Selection

To support the subsequent multi-perspective empirical analysis, we design the following weighted scoring formula to select the optimal pipeline for topic modeling and label generation:

$$S_m = \alpha \cdot \frac{Coh_m + Cov_m + Div_m}{3} + (1 - \alpha) \cdot \frac{Q_m}{2} \quad (12)$$

where S_m denotes the comprehensive model scores of the m -th model, and $\alpha \in [0, 1]$ is a weight parameter that balances the importance of quantitative and qualitative evaluation metrics. $Q_m = \frac{Q_{m1} + Q_{m2}}{2}$, where $Q_{m1} = \frac{1}{T} \sum_{t=1}^T (Acc_t + Use_t)$ denotes the scores assigned by the evaluators. A larger value of α emphasizes quantitative evaluation (i.e., topic coherence, diversity, and coverage), while a smaller α gives greater weight to qualitative evaluation (i.e., label accuracy and usefulness). This parameter can be flexibly adjusted according to the specific requirements of the task and application scenario.

We conduct a systematic sensitivity analysis of α in the range from 0.1 to 1.0. As shown in Figure 10, the LDA-GPT model achieves its highest comprehensive score S when $\alpha = 0.1$ ($S_{LDA-GPT} = 0.93$ on $D_{developer}$ and $S_{LDA-GPT} = 0.88$ on D_{user}). Therefore, we adopt $\alpha = 0.1$ as the default weight setting for subsequent experiments.

After determining the optimal model and weight parameter, we further tune two key hyperparameters of the LDA model: the number of topics (T) and the number of keywords per topic (K). Specifically, for both the developer and user datasets, we perform a sensitivity analysis by varying the number of topics T over $\{10, 20, 30, \dots, 100\}$ and compute the topic coherence score for each configuration. The results indicate that the highest coherence is achieved at $T = 40$ for the developer dataset and $T = 50$ for the user dataset. As shown in Figure 11(a). Thus, we select $T_{D_{Developer}} = 40$ and $T_{D_{User}} = 50$ as the optimal numbers of topics for subsequent analysis.

Building upon the optimal topic numbers, we further perform a sensitivity analysis on the number of keywords per topic K (i.e., $K = 10, 20, 30, \dots, 100$). The results show that the topic coherence score reaches its maximum when $K = 20$ for both datasets (see Figure 11(b)). Therefore, we set the number of keywords per topic to 20, striking a balance between the richness of topic description and the semantic cohesion of the resulting topics.

Finding 8: The proposed HTModel exhibits strong and consistent performance across different model combinations. All hybrid models achieve over 80% labeling accuracy, with LDA_GPT, BERTopic_GPT, and Top2Vec_GPT exceeding 90%. Notably, only LDA_GPT and BERTopic_GPT surpass 80% in labeling usefulness, with LDA_GPT achieving the best overall performance on our multi-stakeholder dataset. Importantly, hybrid pipelines offer an effective trade-off between cost, automation, and interpretability, making them particularly well-suited for large-scale empirical studies.