







StreamVis: An Analysis Platform for YouTube Live Chat Audience Interaction, Trends and Controversial Topics

Gabriela B. Kurtz¹^a, Stéfano de P. Carraro², Carlos R. G. Teixeira²^b, Leonardo D. Bandeira³^c,
Bernardo L. Müller², Roberto Tietzmann²^d, Milene S. Silveira²^e and Isabel H. Manssour²^f

¹University Canada West, Vancouver, Canada

²Pontifical Catholic University of Rio Grande do Sul (PUCRS), Porto Alegre, Brazil

³Federal University of Rio Grande do Sul (UFRGS), Porto Alegre, Brazil

{stefano.c, b.lorenzi}@edu.pucrs.br, {rtietz@pucrs.br, milene.silveira, isabel.manssour}@pucrs.br

Keywords: YouTube, Social Media Analysis, Visualization, Streaming, Chat Interaction.

Abstract: This paper presents StreamVis, an easy-to-use platform that provides stats and visual representations to analyze live chat data from YouTube. StreamVis uses Python and Google's YouTube Data API for data gathering, combined with libraries such as NLTK for natural language processing, Pandas for data analysis, and Matplotlib for visualization. Its interactive dashboard facilitates real-time data visualization through frequency charts, word clouds, and sentiment analysis, providing deep insights into audience engagement patterns. A case study analyzing the NFL's first game in Brazil broadcast on Cazé TV demonstrates how StreamVis reveals trends in audience interactions during critical moments, like game highlights and performances. StreamVis is different from previous tools because it has a user-friendly interface, enabling non-technical users (such as journalists and other media professionals) to perform complex data analysis with a large volume of content, helping them to understand how live chat dynamics influence media consumption.

1 INTRODUCTION


Television producers have always used various methods to measure audience response and guide production efforts (Barnouw, 1990), (Sterling and Kittross, 2001), (Bourdon and Méadel, 2014). The most traditional, inherited from the radio, included letters and phone calls to the broadcaster regarding aired programs. This was followed by audience measurement techniques, which became more automated over the years. In parallel, qualitative studies have also been regularly used to assess audience appreciation beyond mere quantitative measurements (Marc, 1996). More recently, online data collection has played a crucial role in shaping strategies and narratives, providing immediate and detailed insights into audience behavior. With the restrictions on data collection on X (for-


merly Twitter) after the acquisition by Elon Musk, the attention flowed to other platforms.


Building on the growing popularity of live broadcasts on YouTube (YT) and similar platforms, the analysis of audience comments becomes increasingly significant. These comments not only reflect the viewers' immediate reactions but also reveal broader social dynamics (Burgess and Green, 2018), such as public engagement, sentiment, and community building around the content being shared.


The massive engagement found on such platforms has shaped sports institutions and promoters such as the International Olympic Committee (IOC) and the National Football League (NFL) to adopt an innovative broadcasting strategy. During the Paris 2024 Olympic Games, Brazilian audiences had multiple viewing options: they could watch the broadcasts on TV Globo, the official broadcasting partner of the event, on both free-to-air and pay TV, through the IOC's own channels on its applications, or via the streaming partner and YT channel of Cazé TV¹, also available for free. On September 6, 2024, the first


¹<https://www.youtube.com/cazetv>


^a <https://orcid.org/0000-0002-8730-3383>

^b <https://orcid.org/0000-0001-6829-1682>

^c <https://orcid.org/0009-0006-0043-912X>

^d <https://orcid.org/0000-0002-8270-0865>

^e <https://orcid.org/0000-0003-2159-551X>

^f <https://orcid.org/0000-0001-9446-6757>

NFL game played in Brazil showcased the Green Bay Packers facing the Philadelphia Eagles on a Friday night, featuring a similar media arrangement. This multi-channel approach underscored the growing significance of digital platforms in sports broadcasting, reflecting a shift towards more flexible and audience-driven viewing experiences.

The main goal of this work is to provide a deeper understanding of how interactive platforms influence audience engagement and media consumption by analyzing user-generated comments. We explore how the comment streams accompanying live broadcasts, particularly in informal and interactive YouTube channels like Cazé TV, contribute to the viewer experience. By analyzing the content, frequency, and sentiment of these comments, we seek to uncover patterns that highlight the role of live-streaming platforms in shaping contemporary media consumption. An interdisciplinary team featuring Computer Science and Communication researchers defined the scope, tool feature set, and prototype case study (NFL in Brazil).

To achieve this goal, the authors developed a platform called StreamVis that collects comments accompanying live broadcasts on YouTube and processes the collected data, generating a dashboard with frequency charts, word clouds, sentiment analysis, and other visual analyses described further in this study. These findings are presented as a report, allowing users to capture and interpret the dynamics of audience interactions in a short time, offering valuable insights into how viewers engage with live content.

Our main contributions are outlined as follows:

- An innovative platform for collecting and analyzing live broadcast comments on YouTube, which enables interactive visualizations of audience engagement through frequency charts, word clouds, and sentiment analysis.
- A step-by-step approach that combines data collection and visual analysis, providing new insights into audience behavior during live events on interactive platforms.
- A case study demonstrating the application of the proposed platform in analyzing comments from Cazé TV's NFL in Brazil broadcast, illustrating how informal and interactive broadcasting formats impact viewer experiences.

The remainder of the paper is organized as follows. Sections 2 and 3 present the background and related work. The proposed approach is described in Section 4. Section 5 details the obtained results by presenting one use case. A discussion on contributions, limitations, and future directions is presented in Section 6. The last section outlines our conclusions.

2 BACKGROUND

This section presents an introduction related to YouTube Lives and the Cazé TV channel in the context of sports broadcasts.

2.1 YouTube Lives

YouTube, owned by Alphabet, is one of the largest social media and streaming platforms in the world. According to the Digital 2024 Global Overview Report, YouTube is the second-ranked social media channel when it comes to time spent on platforms, with an average of slightly over 28 hours per month per user—it only loses to TikTok, with 34 hours². As of July 2024, India had the largest YouTube audience in the world, with roughly 476 million users actively engaging with the video-sharing platform, a little over one-third of its population. Both the United States, ranked second, and Brazil, ranked third, had over two-thirds of their populations watching YouTube. The US had approximately 238 million viewers, while Brazil had about 147 million people using YouTube³.

Live streaming, once a niche technology used by early internet pioneers, has evolved into one of the most popular and influential forms of broadcasting. Initially introduced in the 1990s, live streaming has since become a cornerstone of digital communication, utilized by tech giants such as Google, Microsoft, and Apple. Its growth was further accelerated by the global shift to remote work and study during the COVID-19 pandemic. The platform that truly revolutionized the accessibility and popularity of live streaming, however, was YouTube. Beginning with its first live event in 2008, YouTube's foray into live streaming signaled a pivotal moment in the technology's mainstream adoption, paving the way for a broader cultural and technological shift that redefined media consumption⁴.

YouTube Live offers a variety of features that make live streaming accessible and flexible for creators. It allows users to stream directly from a webcam, mobile device, or through more advanced setups with encoders. Creators can use real-time interactions like polling and visual overlays to engage their audience, while options like monetization and scheduling enhance the overall live experience. The live chat feature is a key component of YouTube Live, enabling real-time communication between creators and viewers. Creators can moderate the chat through tools like slow mode, where messages are spaced out to control

²<https://tinyurl.com/ydvmec8t>

³<https://tinyurl.com/39bshxj7>

⁴<https://tinyurl.com/2kv9w2we>

the flow, and can even filter or block inappropriate comments. This interaction fosters stronger community engagement, encouraging real-time feedback and participation. YouTube Live also provides creators with detailed insights about their streams. Metrics such as concurrent viewers, chat activity, and overall audience demographics help creators evaluate their content's performance. These analytics allow streamers to better understand their audience, refine their content, and optimize future live broadcasts, making YouTube Live not only a platform for real-time interaction but also a tool for long-term content strategy development⁵.

2.2 Cazé TV and Sports Streaming

The media strategy used by the league deserves attention: the event was broadcasted across three different platforms (RedeTV on free-to-air TV, ESPN Brazil on cable TV, and CazéTV via online streaming), and scheduled for the NFL's first Friday night primetime slot. This approach contributed to the success of the event, not just in the host state but throughout Brazil⁶.

Casimiro Miguel, popularly known as Cazé, is a reference in sports broadcasting on YouTube, with over 16 million subscribers to his channel, more than 7,000 videos published, and over 1.8 billion views. The streamer has gained worldwide recognition for his coverage of major events such as the World Cup and the Olympics, ranking 300th globally in terms of subscribers⁷. During the 2022 World Cup, Casimiro made history by becoming the first Brazilian streamer to broadcast tournament matches on YouTube. Thanks to an agreement with FIFA and LiveMode, he streamed 22 matches, including those of the Brazilian national team, the semifinals, and the final. His channel, CazéTV, peaked at 6 million simultaneous viewers, setting a live viewership record on YouTube, ranking among the top 4 of the 5 most-watched live streams⁸.

This success marked a turning point for sports streaming, demonstrating that digital platforms could compete with traditional television broadcasters⁹.

Another important milestone was his role in broadcasting the first NFL game in Brazil. The event offered Casimiro an opportunity to expand his reach beyond soccer. Although American football does not have the same penetration as traditional soccer in

Brazil, the broadcast lasted over 7 hours and garnered nearly 3 million views, as shown in Figure 1¹⁰.

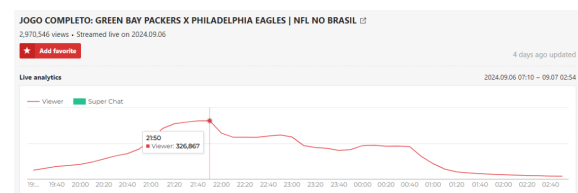


Figure 1: Total views of the NFL broadcast. Source: <https://playboard.co/en/video/W049t2oII4Q>.

3 RELATED WORK

Research about YouTube has been quite abundant, since it is a well-established platform. In our work, we aim to expand on those findings, acknowledging the previous works. When looking for those works, we filtered investigations that were related to YouTube comments, both in regular uploaded videos and Lives. We noticed a lower volume of comment analysis in YouTube livestreams, suggesting room for further analysis.

When it comes to sentiment analysis on YouTube, recent studies have made progress in using machine learning and deep learning to understand user opinions from the platform's large amount of comments. For instance, Chalkias et al. (Chalkias et al., 2023) used tools like VADER and TextBlob to study comments on educational YouTube videos, finding that most comments were neutral, but many were positive when videos used engaging features like animations. Mehta and Deshmukh (Mehta and Deshmukh, 2022) focused on YouTube ads, using machine learning models like Decision Trees, which showed strong results in predicting how viewers would react. These studies highlight how automated systems are improving in handling large datasets and identifying patterns in user sentiment across different types of YouTube content.

Earlier research laid the groundwork for these advancements. Singh and Tiwari (Singh and Tiwari, 2021) used several machine learning models, finding that Support Vector Machine (SVM) worked best for analyzing comment sentiment. Sentiment analysis is also a topic for multi-language adaptations other than English. Yasin et al. (Yasin et al., 2020) tackled the challenge of bilingual sentiment analysis in English and Roman Urdu, successfully using logistic regression to rank videos based on comments. This is just one example of this analysis being done in multiple languages. Uryupina et al. (Uryupina et al.,

⁵<https://tinyurl.com/2nwj4eb8>

⁶<https://tinyurl.com/mwznhf66>

⁷<https://tinyurl.com/p6zm5wcy>

⁸<https://tinyurl.com/yau4h6vs>

⁹<https://tinyurl.com/djh43p7t>

¹⁰<https://tinyurl.com/3hcejrm9>

2014) introduced the SenTube corpus, a dataset designed for sentiment analysis in multiple languages. Together, these studies show how sentiment analysis on YouTube has grown from basic dataset creation to more advanced models that can handle a wider range of content and languages.

Sentiment analysis has been widely used by researchers on YouTube, but there were also other approaches regarding user participation in the comments. It is important to note that the following researches were done in pre-recorded posted videos on YouTube, not live streams. Recent studies have explored various aspects of YouTube comments, from informal learning to public discourse and content moderation. Dubovi and Tabak (Dubovi and Tabak, 2020) focused on how YouTube comments foster knowledge co-construction, particularly in science videos. Their research demonstrated that user comments often go beyond simple information sharing, engaging in discussions and debates that lead to deeper learning.

Other researchers have examined the broader dynamics of YouTube comments, particularly their impact on engagement and content quality. Siersdorfer et al. (Siersdorfer et al., 2010) analyzed over six million comments to understand the factors that influence comment ratings and usefulness. They found a clear correlation between positive sentiment and higher community ratings, while offensive language led to negative ratings. Their work also introduced machine learning classifiers to predict which comments would be rated positively or negatively by the community, providing insights into how sentiment and language shape user interactions on the platform.

Comment analysis on YouTube has also been applied to specific issues such as hate speech and spam. Latorre and Amores (Latorre and Amores, 2021) used topic modeling to identify racist and xenophobic comments targeting migrants and refugees, showing that 19.35 percent of the analyzed comments contained hate speech. This highlights the darker side of YouTube comments, where far-right rhetoric and nationalist views are prevalent in certain channels. Abdullah et al. (Abdullah et al., 2018) took a different approach, focusing on spam detection in YouTube comments. They compared various spam filtering techniques, finding that low-complexity algorithms can achieve high accuracy in identifying spam content, suggesting that YouTube's built-in tools could be enhanced to combat the spread of malicious content. These studies show that while YouTube comments offer valuable insights into user sentiment and interaction, they also present challenges related to content moderation and the spread of harmful speech.

More specifically about YouTube comments analysis in live streams, we have some relevant findings. Recent research focuses on understanding user behavior, emotional intensity, and content moderation during live events. Sentiment analysis was one of the tools used to analyze live stream comments. Tirpude et al. (Tirpude et al., 2024) developed a system to analyze sentiments in live chat through Natural Language Processing (NLP) techniques. By using Fast-Text for sentiment scoring and emoji analysis, they provided real-time insights into audience reactions, enabling content creators to adjust their approach dynamically during live streams. Similarly, Liebeskind et al. (Liebeskind et al., 2021) investigated engagement patterns in YouTube live chats, specifically during political events such as Donald Trump's 2020 campaign. Their study revealed that live comments were highly emotional, with a significant portion being abusive, but frequent commenters tended to use less offensive language, emphasizing how emotional involvement plays a role in live chat dynamics.

Emotional intensity in live streams has been another focus, as evidenced by the works of Luo et al. (Luo et al., 2020) and Guo and Fussell (Guo and Fussell, 2020). Luo et al. explored how emotions are amplified in live chat compared to comments posted after events, finding that shared real-time experiences intensify both positive and negative sentiments. This emotional amplification has implications for content moderation, as heightened emotions can lead to an increase in abusive or toxic comments. Guo and Fussell took this further by examining emotional contagion in live-streaming environments, showing that the sentiments in chat messages have a stronger influence on subsequent messages than the content of the video itself. Their findings suggest that audience interactions can significantly shape the overall sentiment of live chat, often more than the live content itself.

In addition to understanding emotional dynamics, other researchers have focused on combating challenges such as spam and toxicity in live chats. Yousukkee and Wisitpongphan (Yousukkee and Wisitpongphan, 2021) analyzed spammers' behavior in live streams, using machine learning models to differentiate between spam and legitimate user engagement. Their decision tree classifier achieved high accuracy in detecting repetitive, irrelevant content. Complementing this, Tarafder (Tarafder et al., 2023) developed an automated tool to identify and flag toxic comments in real-time, addressing the increasing need for content moderation during live streams, especially as platforms like YouTube have seen a surge in usage during the pandemic.

An interesting recent tool that is worth mentioning

is the “CatchLive”, developed by Yang et al. (Yang et al., 2022). The system provides a real-time summary by segmenting the stream into high-level sections and highlighting key moments based on both stream content and user interaction data, such as chat messages and likes. The paper discusses the development of two core algorithms—one for stream segmentation and another for identifying highlight moments. Even though this tool was not developed for analyzing comments on YouTube Livestreams, it was an interesting inspiration for our team when developing our solution, since the main goal is to provide a comprehensive summarization of user comments and engagement in Live Streamings.

4 StreamVis PLATFORM

To facilitate the analysis of YT live chat, we developed an online platform for gathering and analyzing data through an interactive dashboard. The dashboard enables temporal analysis of chat interactions, allowing users to identify influential viewers, trending topics, and controversial discussions. Its easy-to-use interface allows users to gain insights from live chat data without the need for technical expertise or programming skills.

Figure 2 presents the main components of StreamVis. The following sections describe these components and the platform’s main functionalities.

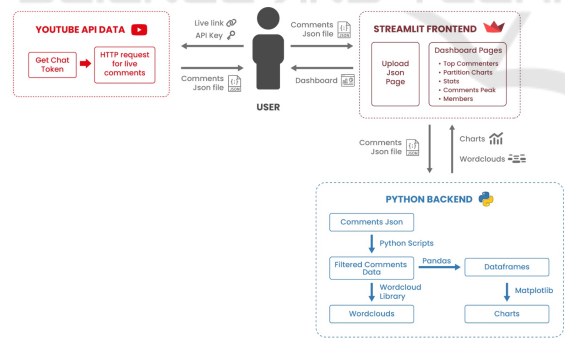


Figure 2: StreamVis main components.

4.1 Implementation Environment

We used the Python programming language and Google’s official YouTube Data API for data gathering. This API returns a payload with relevant data, such as comment message, author, and timestamp. The decision to use Python as the primary programming language for this project was based on its library support and performance in data analysis and natural language processing.

In addition to the Python language, we also used Streamlit¹¹ to develop the dashboard interface.

In Table 1, we outline the technologies utilized and their essential roles in the project. Integrating these libraries with Python provided a robust framework for conducting large-scale analysis of YouTube comments, enabling seamless interaction with APIs, efficient data processing, natural language analysis, and the creation of insightful visualizations.

Table 1: Libraries used and their respective functions.

Library	Description
requests	Used to retrieve YouTube comments in real-time by interacting with the YouTube Data API, enabling dynamic data collection directly from the platform.
python-dateutil	Used for parsing and manipulating dates, ensuring consistent processing and visualization by normalizing date and time data.
pytz	Used for managing time zones and converting timestamps to a uniform time reference, facilitating accurate time-based analysis across different geographic locations.
python-dotenv	Used for securely managing environment variables and safeguarding sensitive information in a .env file.
nltk	Used for preprocessing comments, particularly for removing common stopwords in Portuguese, ensuring that only meaningful words remained in text-based analyses.
pandas	Used for organizing comments into DataFrames and conducting detailed analysis, such as identifying peaks in commenting activity.
WordCloud	Used for generating visual representations of the most frequently occurring words in the comments, providing intuitive insights into prominent themes and topics.
matplotlib	Used for generating charts and plots, aiding in the interpretation of resampled data and identification of trends in user behavior.

4.2 Data Gathering

As a tool for data gathering, we developed a small and reusable script in Python that uses the YouTube Data API to make requests. This script interacts with YouTube’s API to retrieve a payload con-

¹¹<https://streamlit.io/>

taining relevant data, such as the comment message, author details, and timestamp. Key functions within the script, `get_live_details()` and `get_chat_messages()`, respectively, are responsible for retrieving live stream details and the live chat ID, and gathering chat messages, calculating elapsed time from the stream's start, and processing each comment, author, and time data into a structured format.

The script continuously runs in a loop, collecting data at specified intervals and automatically saving the gathered data into a JSON file. This data storage ensures persistence and can be later used as input to the developed dashboard, enabling further visualizations like frequency charts, word clouds, and sentiment analysis. Additional functions manage the reading and writing of data and filter out duplicate entries, ensuring only new, relevant comments are added to the dataset.

This script's main benefit is its ease of use¹²; users only need to provide their Google API Key and the YouTube video ID, which can be easily obtained from the video URL. This simplicity allows even non-technical users to deploy the script effortlessly. Furthermore, the script's automated real-time data collection with subsequent visual analysis through the dashboard described in the next section enhances the understanding of audience behavior during live YouTube events. The proposed platform is particularly advantageous for informal and interactive broadcasting formats, as it captures the dynamic flow of viewer responses, which traditional data-gathering methods might overlook.

4.3 Dashboard Description

StreamVis is designed to provide insights into user engagement through YouTube comments and live chat data. Its dashboard allows users to examine critical metrics, such as chat volume, trending topics, and influential users.

The dashboard offers a simple interface, incorporating line and bar charts and word clouds. It enables users to filter data, zoom, and access detailed insights, allowing for an in-depth exploration of audience behavior without requiring programming expertise. Additionally, customizable views enhance the tool's adaptability to diverse use cases, supporting various analytical needs.

A sidebar menu allows users to navigate between various features of the dashboard. The available options included are shown in Figure 3 and explained in detail below.

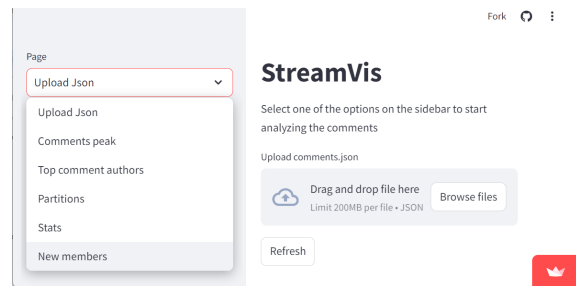


Figure 3: StreamVis main page presenting menu options.

- **Upload JSON** - Allows you to load a JSON file with the data collected by the script described in Section 4.2.
- **Comment Peaks** - One of the primary features of StreamVis is the ability to identify and analyze peaks in comment activity. The Comment Peaks page exemplified in Figure 4 generates an interactive time-series chart that visualizes moments of high comment volume. These peaks often correspond to key events in the video, such as significant announcements or moments of heightened viewer interest. This chart allows users to zoom in on specific time intervals and explore trends inside the collected data. In addition to visualizing the peaks, StreamVis generates a corresponding word cloud for each peak based on the comments made during that time frame.
- **Top Commenters** - Highlights the most active participants and offers a detailed breakdown of all their comments, including timestamps and content, allowing analysts to explore the volume of engagement and the nature of the contributions made by key users. It is also possible to configure the number to the top commenters you want to analyze, as shown in Figure 5.
- **Partitions** - For further analysis, allow the user to divide the dataset into custom time intervals for further analysis. Sliders allow users to adjust the granularity of the data displayed, such as resampling comment data by minutes or hours for more refined insights.
- **Stats** - Provides general statistics about the dataset, such as the total number of comments, unique commenters, and average comments per user.
- **New Members** - As illustrated in Figure 6, displays new users who have recently joined a paid membership plan, helping to sponsor the streaming channel and its projects. This function shows the timing and volume of memberships throughout the live stream.

¹²The script is available at <https://github.com/DAVINTLAB/StreamVis>

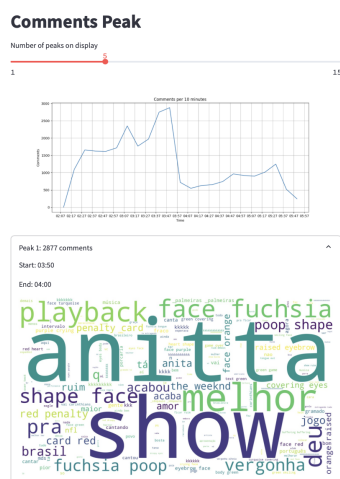


Figure 4: StreamVis comments peak page.



Figure 5: StreamVis top commenters page.

StreamVis provides an easy-to-use platform for analyzing YouTube comments and live chat data. Through its interactive components, including time-series charts, word clouds, and user activity tracking, the tool facilitates in-depth exploration of audience engagement.

The flexibility afforded by customizable time intervals, filtering capabilities, and visual representations allows the platform to adapt to diverse analytical contexts. Whether analyzing overall chat dynamics to understand general audience engagement, zooming in on key moments of significance, or identifying stand-out trends and commenters, StreamVis provides the tools necessary to explore the data from various perspectives.

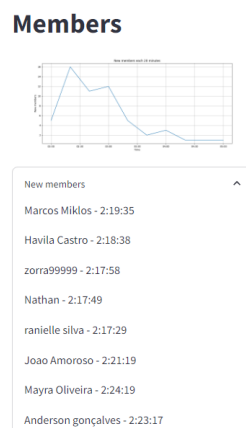


Figure 6: StreamVis new members page.

5 CASE STUDY

In this section, we present a case study to exemplify the use and potential of the developed approach.

5.1 NFL and Its First Game in Brazil

The NFL has been heavily investing in its internationalization strategy since 2007, when it held its first game outside the United States, between the New York Giants and the Miami Dolphins at Wembley Stadium in London. This event marked the beginning of the NFL International Games, which have since expanded, with London becoming the main hub for the league's international games. Between 2007 and 2023, more than 30 regular season games were played in the English capital, at venues such as Wembley and Tottenham Hotspur Stadium¹³.

The peak of this expansion reached Brazil in 2024, with the first NFL game on Brazilian soil, held on September 6th at Neo Química Arena in São Paulo, the stadium of the Corinthians football team. The event attracted more than 47,000 people to the stadium and was a milestone for the NFL's presence in South America, breaking audience records both in the American and Brazilian broadcasts¹⁴.

5.2 Game Analysis Through the Proposed Platform

The Eagles and Packers game started at 9:15 PM Brasilia time (Brazil); however, the official live broadcast on Cazé TV's YouTube channel began two hours

¹³<https://tinyurl.com/hj6622mk>

¹⁴<https://tinyurl.com/ny6y2rr4>

earlier. The analysis presented below refers exclusively to the game period, from the kickoff (9:15 PM) until its conclusion (past midnight). In total, the match lasted approximately 3 hours and 50 minutes, during which data collection occurred, resulting in 28,524 comments during the live stream, representing an average of 2 comments per second. This figure is noteworthy, as a higher number of posts was initially expected. The authors feared a bottleneck in data collection due to limitations of the YouTube API, which did not occur, making it possible to capture nearly all the comments made.

A total of 9,500 users made comments during the live stream, with an average of 3 comments per user. This number can also be considered low, as the number of views approached 3 million, with comments representing less than 0.5 of the views. The most active user posted 91 comments, averaging approximately 17 comments during the 30-minute periods in which they were most active. The 10 most active users averaged around 80 comments, showing varied frequency patterns throughout the game, as illustrated in Figure 7.

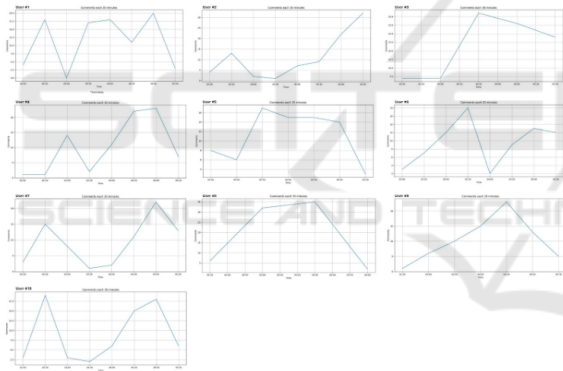


Figure 7: Comments throughout the live stream by the top 10 users.

Figure 7 shows the distribution of all comments over the course of the game, allowing for a deeper analysis of user behavior. What one can observe in the ten charts shown in Figure 7 is that there are distinct behaviors among users. Primarily, users' posting behavior does not follow a consistent curve, as one might expect from a professional commentator or a game narrator. Driven by fan support for the team or attention to specific plays, the curves suggest different motivations for engagement, which becomes evident when examining individual messages more closely.

Considering the overall numbers of the comments, there were 135,498 words with 23,172 unique words, and an average of 4.75 words per comment. When all comments made during the game broadcast are aggregated into a line chart based on the comment fre-

quency one can observe an increasing interest toward the middle of the match and a decreasing interest in the second half, represented in Figure 8.

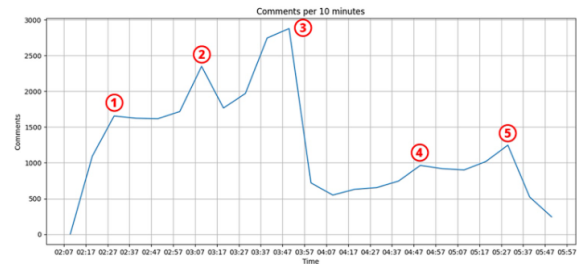


Figure 8: Comments throughout the live stream.

The three highest peaks of interest in terms of comments are clustered in the first half. Typically, fan enthusiasm and key plays prompt increased commenting. Our platform can analyze each peak in detail using segmented word clouds. The first peak, comprising 1,655 comments (5.8% of the total), occurred during one of the initial commercial breaks (Figure 9), highlighting a common second-screen behavior where viewers shift their attention to chat and comments during pauses in the main broadcast. This break followed a field goal by the Green Bay Packers, a moment that sparked a trend seen throughout the broadcast: references and comparisons between American football (NFL) and soccer (the dominant sport in Brazil and the location of the event, held at the Corinthians football club stadium). In this context, the recurrence of the word 'Goiás' stands out, referring to Goiás Futebol Clube, which shares a similar visual identity with the Green Bay Packers. Viewers played on this similarity with comments like 'VAMO GOIÁS, 3X0 JÁ' (LET'S GO GOIÁS, 3-0 NOW), 'Goiás 3 X 0 Corinthians?' and 'CARACA KKK GOIÁS TÁ JOOOGANDO!!!!' (WOW LOL GOIÁS IS PLAYIIING!!!!). Moreover, much of the discussion centered on 'Anitta,' the Brazilian singer performing at the halftime show, with comments such as 'we want Anitta,' 'When is Anitta's show???' and 'Come on, Anitta.'

The second peak also occurred during the break following another score, with 2,349 comments (8.2%), already indicating a higher frequency of viewers as the main break and the show approached. At this point, mentions of the singer 'Anitta' were already considerably higher than any other topic. During the third peak, Anitta's show finally took place, directly mentioned in the comments, with 2,877 comments (10% from the total). The comments reveals a high volume of criticism directed at the show, with words like 'playback' and 'shame' standing out. Throughout Anitta's performance, the comments



Figure 9: Wordcloud of peak one comments.

were divided between criticism of her use of playback and her decision to sing only her international songs rather than her native Portuguese repertoire of songs. It is also worth noting the word ‘better,’ which, when analyzed in isolation, might seem positive, but upon reviewing the comments, its tone aligned with the critical sentiment of the audience: ‘wasn’t there a better singer??’, ‘the game deserved something better,’ ‘so many better and quality artists to choose from,’ among others. Few comments used ‘better’ in a complimentary manner. The presence of emojis in the comments was also noteworthy, with a significant volume of expressions representing shame and criticism. ‘shame’ / Fuchsia Poop / poop shape (poop emoji?) / face fuchsia (shame?)”

The fourth peak occurred during a rare game moment: a field goal hitting the post, which allows for a direct comparison to soccer. This peak had 964 comments (3.3%), with most of them once again focusing on the game, consistently making references to soccer. The pitch and its quality were frequently mentioned. Athlete LeBron James criticized the pitch on his social media, which reignited discussions in the live chat about its quality and concerns that the sport might be ‘ruining’ the stadium’s pitch. Some of the comments included: ‘LeBron criticized the pitch, oh boy,’ ‘Is it just me, or is the pitch deteriorating?’, ‘EVEN LEBRON COMPLAINED ABOUT THE PITCH,’ ‘COME ON CORINTHIANS, WHAT ARE THEY DOING TO OUR PITCH, IT’S A DISGRACE.’”

Finally, peak 5 had 1,247 comments (4.3%), around the time of an error in the official broadcast that incorrectly announced a touchdown and mistakenly changed the score, a fact that was highlighted by

the commentators and caught viewers’ attention. In spite of that, this incident itself did not stand out in the comments, once more focused on the connection to soccer. The presence of the words ‘Corinthians,’ ‘Timão’ (a popular nickname for the team in Brazil), and ‘Vai Corinthians’ (a characteristic chant among the fans) versus ‘Goias’ brought soccer back into the spotlight. Another analysis enabled by the tool relates to the channel’s members. The channel has a membership area that offers benefits in exchange for a monthly subscription. These benefits range from a badge in the chat to the opportunity to make more frequent comments. Regarding memberships, there were 57 new members during the game. It is possible to observe a trend that the beginning of the games has a higher potential for attracting new members, with this number decreasing significantly as the match progresses.

6 DISCUSSION

While many studies, such as those by Mehta and Deshmukh (Mehta and Deshmukh, 2022), focused on pre-recorded videos and YouTube ads, this work specifically targets live streaming environments, which present unique challenges and opportunities in real-time data gathering and its subsequent analysis.

The idea of collecting and analyzing live comments and chats during YT broadcasts has proven rich in possibilities and relatively simple to implement. Firstly, StreamVis does not automatically provide a way to receive live video or audio feeds, as typical in TV broadcasts or platforms like StreamYard. Thus, the primary feedback channel chosen by YouTube is text, which includes using slang, hashtags, and emojis, adding different inflections to the interaction. Previous research on comment analysis has mainly addressed content moderation and spam detection in pre-recorded content, as seen in the works of Latorre and Amores (Latorre and Amores, 2021) or Abdullah et al. (Abdullah et al., 2018). In contrast, our approach focuses on providing insights into audience engagement and interaction dynamics during live streams, making it more relevant to dynamic broadcasting scenarios.

Besides, we are not directly analyzing video or audio content, which would require significantly more computational power. This allows our platform to be used on various machines without needing high-end hardware. This accessibility is advantageous, as it opens opportunities for analysts, consultants, and professionals interested in understanding and evaluating engagement with YT live content. Our platform im-

plementation emphasizes visual analytics using, e.g., frequency charts and word clouds, allowing users to interpret trends in live chat data interactively. This visual approach sets this tool apart from traditional text-based sentiment analysis methods used in earlier studies.

The ability to perform a visual analysis of the data allows the StreamVis users to develop a macro understanding of the dynamics during the broadcast and audience engagement, and a granular view that enables them to examine conversations in detail and identify situations or moments that acted as turning points. This dual perspective helps users capture broader trends while also pinpointing specific interactions or events that significantly influenced the flow of the live broadcast. When aligned with an understanding of what was happening in the broadcast during those specific minutes, the visual nature of peaks allows us to build an insight into how people collectively react to what is being shown. We recognize that different events can elicit varying behaviors, and even the behavior of the most active users is not necessarily convergent.

The choice of an NFL game as the proof-of-concept of our platform was intentional for several reasons. First, the novelty of the sporting event guaranteed attention and was likely to drive engagement. Additionally, the growing interest in American football in Brazil is a proxy to the richer strata of the population, one with widely available internet services and devices, consequently abundant participation on YT and other platforms.

7 FINAL REMARKS

The work done in this study contributes to the existing area research by introducing a platform specifically designed for YouTube live chat analysis, which allows real-time data gathering, visualization of user engagement, trending topics, and emotional intensity during live events. We provide all the codes of our data gathering and visual analytics approach at the GitHub¹⁵, and the StreamVis app is available online¹⁶ for anyone who wants to use it. Thus, it is possible to gather and analyze data from different contexts.

As future possibilities for this research, we believe that integrating other natural language processing models and AI infrastructures will enable a more detailed sentiment analysis, potentially allowing for seamless operation in multiple languages. This would

make it possible to observe the reactions of international audiences to the same event that attracts viewers from various countries, such as the Olympic Games, the UEFA Champions League, or other similar sporting events.

ACKNOWLEDGEMENTS

Carraro and Müller would like to thank the Tutorial Education Program (PET). Manssour would like to thank the financial support of the CNPq Scholarship - Brazil (303208/2023-6).

While preparing and revising this manuscript, we used ChatGPT, Grammarly, and Google Translate to enhance clarity and grammatical precision, as English is not our first language. The authors take full responsibility for the content and its technical accuracy.

REFERENCES

- Abdullah, A. O., Ali, M. A., Karabatak, M., and Sengur, A. (2018). A comparative analysis of common youtube comment spam filtering techniques. In *2018 6th International Symposium on Digital Forensic and Security (ISDFS)*, pages 1–5.
- Barnouw, E. (1990). *Tube of Plenty: The Evolution of American Television*. Oxford University Press.
- Bourdon, J. and Méadel, C. (2014). *Television Audiences Across the World: Deconstructing the Ratings Machine*. Palgrave Macmillan UK.
- Burgess, J. and Green, J. (2018). *YouTube: Online Video and Participatory Culture*. Digital Media and Society. Polity Press.
- Chalkias, I., Tzafilkou, K., Karapiperis, D., and Tjortjis, C. (2023). Learning analytics on youtube educational videos: Exploring sentiment analysis methods and topic clustering. *Electronics*, 12(18).
- Dubovi, I. and Tabak, I. (2020). An empirical analysis of knowledge co-construction in youtube comments. *Computers & Education*, 156:103939.
- Guo, J. and Fussell, S. R. (2020). A preliminary study of emotional contagion in live streaming. In *Companion Publication of the 2020 Conference on Computer Supported Cooperative Work and Social Computing, CSCW '20 Companion*, page 263–268, New York, NY, USA. Association for Computing Machinery.
- Latorre, J. P. and Amores, J. J. (2021). Topic modelling of racist and xenophobic youtube comments. analyzing hate speech against migrants and refugees spread through youtube in spanish. In *Ninth International Conference on Technological Ecosystems for Enhancing Multiculturality (TEEM'21)*, TEEM'21, page 456–460, New York, NY, USA. Association for Computing Machinery.

¹⁵<https://github.com/DAVINTLAB/StreamVis>

¹⁶<https://davint-live-comments.streamlit.app>

- Liebeskind, C., Liebeskind, S., and Yechezkely, S. (2021). An analysis of interaction and engagement in youtube live streaming chat. In *2021 IEEE SmartWorld, Ubiquitous Intelligence & Computing, Advanced & Trusted Computing, Scalable Computing & Communications, Internet of People and Smart City Innovation (SmartWorld/SCALCOM/UIC/ATC/IOP/SCI)*, pages 272–279.
- Luo, M., Hsu, T. W., Park, J. S., and Hancock, J. T. (2020). Emotional amplification during live-streaming: Evidence from comments during and after news events. *Proc. ACM Hum.-Comput. Interact.*, 4(CSCW1).
- Marc, D. (1996). *Demographic Vistas: Television in American Culture*. University of Pennsylvania Press.
- Mehta, T. and Deshmukh, G. (2022). Youtube ad view sentiment analysis using deep learning and machine learning. *International Journal of Computer Applications*, 184(11):10–14.
- Siersdorfer, S., Chelaru, S., Nejd, W., and San Pedro, J. (2010). How useful are your comments? analyzing and predicting youtube comments and comment ratings. In *Proceedings of the 19th International Conference on World Wide Web, WWW '10*, page 891–900, New York, NY, USA. Association for Computing Machinery.
- Singh, R. and Tiwari, A. (2021). Youtube comments sentiment analysis. *International Journal of Scientific Research in Engineering and Management (IJSREM)*, 5(5):1–11.
- Sterling, C. and Kittross, J. M. (2001). *Stay Tuned: A History of American Broadcasting*. Routledge Communication Series. Taylor & Francis.
- Tarafder, T., Vashisth, H. K., and Arora, M. (2023). Automated tool for toxic comments identification on live streaming youtube. In *International Conference on MAchine inTelligence for Research & Innovations*, pages 47–56. Springer.
- Tirpude, S., Thakre, Y., Sudan, S., Agrawal, S., and Ganorkar, A. (2024). Mining comments and sentiments in youtube live chat data. In *2024 4th International Conference on Intelligent Technologies (CONIT)*, pages 1–6.
- Uryupina, O., Plank, B., Severyn, A., Rotondi, A., Moschitti, A., et al. (2014). Sentube: A corpus for sentiment analysis on youtube social media. In *LREC*, pages 4244–4249.
- Yang, S., Yim, J., Kim, J., and Shin, H. V. (2022). Catchlive: Real-time summarization of live streams with stream content and interaction data. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems, CHI '22*, New York, NY, USA. Association for Computing Machinery.
- Yasin, S., Ullah, K., Nawaz, S., Rizwan, M., and Aslam, Z. (2020). Dual language sentiment analysis model for youtube videos ranking based on machine learning techniques. *Pakistan Journal of Engineering and Technology*, 3(2):213–218.
- Yousukkee, S. and Wisitpongphan, N. (2021). Analysis of spammers' behavior on a live streaming chat. *IAES International Journal of Artificial Intelligence*, 10(1):139.