

# Game-MUG: Multimodal Oriented Game Situation Understanding and Commentary Generation Dataset

Anonymous ACL submission

## Abstract

The dynamic nature of esports makes the situation relatively complicated for average viewers. Esports broadcasting involves game expert casters, but the caster-dependent game commentary is not enough to fully understand the game situation. It will be richer by including diverse multimodal esports information, including audiences' talks/emotions, game audio, and game match event information. This paper introduces GAME-MUG, a new multimodal game situation understanding and audience-engaged commentary generation dataset and its strong baseline. Our dataset is collected from 2020-2022 LOL game live streams from YouTube and Twitch, and includes multimodal esports game information, including text, audio, and time-series event logs, for detecting the game situation. In addition, we also propose a new audience conversation augmented commentary dataset by covering the game situation and audience conversation understanding, and introducing a robust joint multimodal dual learning model as a baseline. We examine the model's game situation/event understanding ability and commentary generation capability to show the effectiveness of the multimodal aspects coverage and the joint integration learning approach.

## 1 Introduction

The recent advent of esports has led to a trendy and rapidly growing industry, capturing the attention of a large and continuously expanding global audience. Within a few seconds of a game event, numerous aspects demand attention, such as player action, skills demonstrations, team cooperation, gain and loss, and the key items contributing to the specific game events. This requires the audience to quickly digest complicated information whenever something significant happens in the game. Unlike conventional sports broadcasting like NBA games (Yu et al., 2018), where the fundamental sport's concepts are easily comprehensible, this dynamic nature of esports introduces complexity,

making it challenging for the average audience to grasp the game situation fully. Therefore, we need to find a way to assist the audience in understanding the game situation better. Esports competition organisers address this issue by involving one or two casters to explain the game situation during live streaming. However, this heavily relies on the specific casters, making it difficult for them to provide more diverse information, including audience opinions, feelings, and detailed game match information. In addition, different casters may prioritise various game aspects, leaving many online esports game resources unexplained. Therefore, it is essential to explore methods for automatically generating game-related commentary that comprehensively understand the game situation, incorporating multiple aspects, such as audience discussion, emotions, and domain-specific information.

Existing esports game commentary datasets (Tanaka and Simo-Serra, 2021; Wang and Yoshinaga, 2022; Zhang et al., 2022) only utilise single-modal information as input to generate textual commentary, disregarding the potential richness of multiple aspects that can provide valuable information about the game. The lack of multimodal resources hinders researchers interested in commentary generation for Multiplayer Online Battle Arena (MOBA) games from determining the best approach to leverage information from various sources to address the game commentary task. Moreover, previous works primarily focus on providing accurate game-related facts (Wang and Yoshinaga, 2022; Zhang et al., 2022) in the generated commentary for the audience, neglecting the importance of infusing human-like qualities and emotions to engage the audience better. Due to the lack of resources, existing game commentary generation models (Tanaka and Simo-Serra, 2021; Zhang et al., 2022; Wang and Yoshinaga, 2022) simply employ an encoder-decoder to process raw game

information and generate human-like commentary without fully understanding the game situations.

We introduce GAME-MUG, a multimodal game situation understanding and commentary generation dataset, and its strong baseline. Our dataset incorporates publicly available League of Legends (LOL) resources with professional caster comments from popular live streaming platforms, YouTube and Twitch, with multimodal information, including game event logs, caster speech audio, and game-related natural language discussions encompassing both human casters' commentaries and audience chats and emotions. Inspired by the joint integration of natural language understanding and generation tasks, we propose a strong baseline model that employs joint integration framework to comprehend game situations from multimodal information and generate game commentary based on this understanding of game situations and emotions. To conduct the game commentary generation, we summarise the game situation and audience conversation via multi-modality sources. Our contribution can be summarised as follows:

- Introducing a multimodal game understanding and commentary generation dataset to provide a full understanding of the game situations with caster comments and diverse information, including audience conversation, caster speech audio, and game event logs.
- Proposing a joint integration framework to generate more human-like commentary with the help of game situation understanding
- Conducting extensive experiments to show the effectiveness of multimodality in game understanding and commentary generation.

## 2 Related Work

### 2.1 Game-related Datasets

Most datasets in the game domain are proposed for commentary generation across different games, such as live-streamed MOBA games (Tanaka and Simo-Serra, 2021; Wang and Yoshinaga, 2022; Zhang et al., 2022) as well as pre-recorded esports games (Ishigaki et al., 2021; Li et al., 2019; Shah et al., 2019) or traditional sports (Yu et al., 2018), while several datasets also focus on classification tasks related to scene understanding as shown in Table 1. CS-lol (Xu et al., 2023) proposed

a task of viewer comment retrieval, while MOBA-LoL (Ringer et al., 2019) proposed two classification tasks on their dataset. On top of predicting game event types, they also provide multi-view to understand the game context, by indicating the streamer's emotional state. Among all the datasets proposed for game commentary generation, most datasets allow only a single modality as the input, video only, or game information only. Some datasets allow multimodal input, but it was not for MOBA games. So far, no previous work utilises audience emotion when they build datasets to generate more human-like commentary for MOBA games. Our dataset provides both audience emotion and rich multimodal input, including audio, audience chat, and game information.

### 2.2 Visual-Linguistic Generation

Most works in video captioning or commentary generation for games used encoder-decoder structure (Yu et al., 2018; Li et al., 2019; Shah et al., 2019; Ishigaki et al., 2021; Tanaka and Simo-Serra, 2021; Zhang et al., 2022; Wang and Yoshinaga, 2022), and some (Tanaka and Simo-Serra, 2021; Zhang et al., 2022; Wang and Yoshinaga, 2022) experimented with several types of structures like unified encoder-decoder, pretraining method, rule-based models, and hybrid models. Some works (Li et al., 2019; Wang and Yoshinaga, 2022; Zhang et al., 2022; Ishigaki et al., 2021; Yu et al., 2018) applied recurrent seq2seq models like LSTM/GRU structures for both encoding the input and decoding for commentary, some (Tanaka and Simo-Serra, 2021; Wang and Yoshinaga, 2022; Zhang et al., 2022) used transformer-based models for generating commentary. However, no model used dense interaction/fusion among different input modalities. Previous models either lack multimodal input or concatenate different modality features as one feature vector or via simple tensor operation. The semantic gap between different modalities is ignored. In addition, no work tried dual learning of understanding game scenes and generating commentary due to limited information provided by datasets. Our method uses the audience's chats and opinions to understand the game context to facilitate the automatic generation of commentary.

### 3 Game-MUG

We introduce a new game commentary dataset using multimodal game situational information,

Dataset	# Matches	Modality sources	Core Task
FSN (Yu et al., 2018)	50	video, transcript	Game commentary generation
Getting Over It (Li et al., 2019)	8	video, audio, transcript	Game commentary generation
Minecraft (Shah et al., 2019)	3	video, transcript	Game commentary generation
MOBA LoL (Ringer et al., 2019)	-	video, audio, streamer’s image	Streamer emotion prediction, game event type prediction
Car Racing (Ishigaki et al., 2021)	1,389	video, game info, transcript	Game commentary generation
LoL-V2T (Tanaka and Simo-Serra, 2021)	157	video, transcript	Game commentary generation
eSports Data-to-Text (Wang and Yoshinaga, 2022)	-	game info, transcript	Game commentary generation
Dota2-Commentary (Zhang et al., 2022)	234	game info, transcript	Game commentary generation
CS-lol (Xu et al., 2023)	20	transcript, chat	Viewer comment retrieval
Game-MUG (ours)	216	audio, chat, game info, transcript	Game commentary generation, game event type prediction

Table 1: Summary of existing game datasets

called Game-MUG. It features three modalities: game match event logs, audio features derived from signal data and textual discussions, such as caster comment transcript and audience chat. It comprises 70k clips with transcripts and 3.7M audience chats collected from 45 LOL competition live streams. Each live stream has an average of 4.8 individual matches, leading to 216 game matches and 15k game events. Game matches are sourced from 3 distinct leagues between 2020 and 2022, including [Tencent League of Legends Pro League](#), [League of Legends Champions Korea](#) and [World Championships](#). These top-tier league matches in various regions attract many views (from 507K to 7.2M), which derives abundant audience chats in multiple languages. We collect caster commentaries and audience live chats from two different livestream platforms: [Twitch](#), which contributes 150 matches, and [YouTube](#), which contributes 66 matches. In addition to this, we crawl game events from the [League of Legends Competitive Statistics Website](#)<sup>1</sup>.

### 3.1 Data Collection

**Gaming Human Commentary Transcription.** We collect human commentaries by transcribing the raw live stream files<sup>2</sup>. Due to the substantial size of live-stream videos, we use [YT-DLP](#) and [Twitch-DL](#) only to download their high-definition (44.1kHz) audio and utilise a speech recognition model named Whisper (Radford et al., 2022) for speech-to-text conversion. Whisper is a large supervised model that implies the encoder-decoder architecture from Transformer (Vaswani et al., 2017). We use Whisper [medium English model](#) and set the compression ratio to 1.7 without previous text conditions for speech-to-text recognition, which slightly trades off the transcript accuracy but maximises its robustness. Each transcribed text is

<sup>1</sup><https://gol.gg/esports/home/>

<sup>2</sup>YouTube and Twitch disable their Automatic Speech Recognition tools on game live streams

paired with its start and end timestamps in seconds.

**Audience Live Chats Collection.** Audience live chats are scrapped from the live stream platforms. We employ a multiplatform software named [Chat Downloader](#) to scrap the chat content from YouTube and Twitch. Because of the multilingual nature of live chats, we use [Lingua](#) to identify different languages and apply a special label called “emo” for chat instances that only include emotes or emojis. Live stream platforms automatically filters out hateful and toxic contents and we further filter out the live chats without any content<sup>3</sup> and associate reminders with their respective timestamps in seconds.

**Game Events Collection.** Game events are collected from the [League of Legends Competitive Statistics Website](#) by a scrapper; it first finds the game-related HTML tags and extracts the contents from the selected tags. It is worth noticing that sometimes the contents of the tags can be empty, which means a minion or a non-epic monster triggers this event. Our scrapper automatically populates missing contents in the tags and links them to game timestamps, constructing complete game event instances. We categorise collected game events into the following 6 different classes in our dataset: **1) Kill:** A game character is defeated; **2) Non-Epic Monster:** A jungle monster is eliminated; **3) Tower:** A turret/inhibitor is destroyed; **4) Dragon:** A dragon is eliminated; **5) Plate:** A turret’s defensive barrier is shattered; **6) Nexus:** An nexus is destroyed, leading to the end of the game.

**Audio Feature Extraction.** It is known that human speech tone fluctuates based on emotions (Kienast and Sendlmeier, 2000) and audio modality demonstrates a notable advantage over video in capturing emotional fluctuations (Wu et al., 2021). Therefore, we extract audio features from the caster speech audio to enrich emo-

<sup>3</sup>We de-identified all chats by masking their usernames. Details can be seen in Appendix A.

Categories	GPT-3.5	GPT-4	Tie
<b>Kill</b>	25.78%	51.56%	22.66%
<b>Tower</b>	14.20%	59.66%	26.14%
<b>Dragon</b>	17.71%	66.67%	15.63%
<b>Overall</b>	18.75%	58.75%	22.50%

Table 2: Pairwise comparison between GPT-3.5 and GPT-4 summaries, the overall agreement coefficient (Krippendorff, 2011) is 0.64 from nine human annotators. In most cases, annotators choose GPT-4 summaries over GPT-3.5 or think they are similar.

Event	# of events	Avg per match	Percentage
<b>Kill</b>	5,548	25.69	36.45%
<b>Tower</b>	2,889	13.38	18.98%
<b>Dragon</b>	1,646	7.62	10.81%
<b>Other</b>	5,138	23.79	33.76%
<b>Total</b>	15,221	70.47	100%

Table 3: Distributions of the more important game events in our collected dataset, where the less important ones, **Non-Epic Monster**, **Plate** and **Nexus**, are categorised into **Other** as an initial step for analysis.

tional representation within diverse domain data. The Geneva Minimalistic Acoustic Parameter Set (GeMAPS) (Eyben et al., 2016) is commonly used for voice research and it encompasses 18 Low-Level Descriptors, which covers features related to frequency, amplitude and spectral parameters. We utilise `audiofile` to convert raw audio files into audio waveforms, and then extract audio features with a sampling rate of 50Hz using openSMILE (Eyben et al., 2010), a tool commonly used for vocal emotion recognition (Doğdu et al., 2022).

### 3.2 Data Annotation

**Game Situation Summary Annotation.** Inspired by the success of Standford Alpaca (Taori et al., 2023), we make use of GPT-3.5 (Ouyang et al., 2022) and GPT-4 (OpenAI, 2023) to condense all 70,711 human commentaries into concise summaries with emotional clues from audience chats as detailed in Appendix Algorithm 1. To ensure the annotation quality, we conduct a pairwise human evaluation between the summaries from GPT-3.5 and GPT-4. As shown in Table 2, GPT-4 excels GPT-3.5 in all three categories, indicating GPT-4’s summaries are better aligned with human understanding. Therefore, we choose GPT-4’s summaries as ground truth annotations in our dataset.

### 3.3 Data Processing

Considering each live stream can be treated as a chronological sequence comprised of game events, human commentaries and live chats, we match them via their timestamps. As game events’ timestamps are reset after each match, we manually adjust them to align with live stream seconds prior to the matching process. Additionally, background music before the commencement of each live stream is also removed manually, since there is no game-related factual information to help with game situation understanding.

## 4 Data Analysis

Our dataset includes 70,711 transcript sentences with an average duration of 12.2 secs and 3,657,611 chat instances. 15,221 game events are collected from 216 game matches. Not all events are equally important for the human caster and audience; **Kill**, **Tower**, and **Dragon** events usually attract more interest than other events. Therefore, we categorise all other events into **Other** as an initial input processing step for our following analysis in Section 4 and experiments in Section 6.3. We present the statistics of each event category in Table 3.

### 4.1 Game Keyword Analysis

Different from other domains, game-related data contains numerous keywords that rarely appear in everyday conversations. We manually extract 2,003 unique keywords from the caster speech transcript in our dataset and clean the typos and misspells while retaining essential abbreviations, such as character’s skills denoted by Q, W, E, and R. As shown in Figure 1, extracted keywords can be categorised into 5 different classes, including skill, player, team, character and item. To better address the importance of each keyword, we compute their Term Frequency - Inverse Document Frequency (TF-IDF) based on the game events with different time windows, specifically 15 seconds and 30 seconds. The transcripts encompassed within these windows are treated as a singular document to compute TF-IDF values. This allows us to identify key terms closely associated with game events. Depending on the precise timing of the event, such a window might encapsulate one or several transcripts. This calculation is performed using the Scikit-learn library (Pedregosa et al., 2011) with normalisation. Figure 1 shows a sample visualisation of the keywords’ characteristics when the

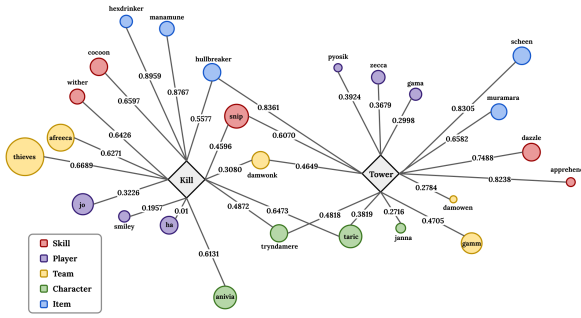


Figure 1: The visualisation for keyword analysis with top 15 words from **Kill** and **Tower** Event, where the time window is 30s. Entities related to each event, Kill and Tower, are remarkably different, such as skill ‘cocoon’ for **Kill** and ‘apprehend’ for **Tower**.

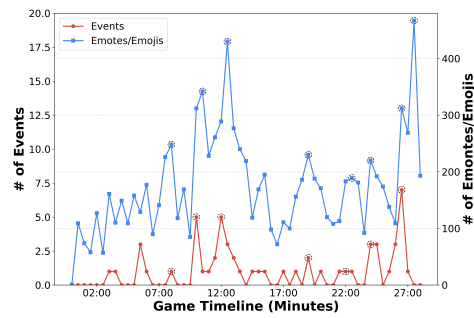


Figure 2: The concurrent plot for audience chat analysis with the numbers of emotes, emojis, and game events along the same timeline. A positive correlation can be observed between the number of audience chat emojis and the number of game events happening within the same time window.

331 window of time equals 30 seconds. We select the  
 332 top 15 keywords for **Kill** and **Tower** events and dif-  
 333 ferentiate their types by distinct colours. The size  
 334 of each keyword’s node depends on the normalised  
 335 occurrence of the keyword, whereas the distance be-  
 336 tween the event and keyword nodes is determined  
 337 by the normalised TF-IDF values. From Figure 1  
 338 we can see that **Kill** and **Tower** are more related  
 339 to items to attack, skills that either increase the  
 340 damage for attacking enemies or limit the ability of  
 341 enemies moving to avoid damage or fighting back.  
 342 This reflects the typical player’s actions in games,  
 343 which often involve attacking opponents, indicat-  
 344 ing that the text in our dataset effectively describes  
 345 the game scene and offers a robust understanding of  
 346 the situation. Moreover, we can see that team, play-  
 347 ers, and character names are frequently mentioned  
 348 or discussed by commentators when these cases  
 349 happened; though the names might depend on spe-  
 350 cific games, it demonstrates the multiple aspects  
 351 people could focus on about the game situation.

## 352 4.2 Audience Chat Analysis

353 The audience tends to send many emotes and emo-  
 354 jis in chat to express their sentiments. We re-  
 355 trieve emotes and emojis based on their distinct  
 356 formats found in publicly available sources<sup>45</sup> and  
 357 then count the number of emotes and emojis per  
 358 30-second window in each match. The counts of  
 359 emotes, emojis, and game events are plotted con-  
 360 currently on the same timeline, shown in Figure 2.  
 361 It is not hard to discover that the number of emotes  
 362 correlates with the game situation, since audiences  
 363 tend to send more emotional expressions in chats to

share their feelings when a dramatic turning point  
 or a series of events happens.

## 366 5 Proposed Baseline

367 Based on Game-MUG, we proposed a joint inte-  
 368 gration framework that generates summaries of the  
 369 game commentaries based on understanding the  
 370 game situation through multimodal data. For game  
 371 situation understanding, we implemented and fine-  
 372 tuned a multimodal transformer encoder that en-  
 373 codes text and audio data. For game commentary  
 374 generation, we employ a pre-trained decoder and  
 375 encoded game information to generate new sum-  
 376 maries. The quality of generated summaries is eval-  
 377 uated by both automatic metrics and humans. We  
 378 partition our dataset into 206 matches for training  
 379 and 10 matches for testing.

### 380 5.1 Input Processing

381 Given an  $i$ -th event  $E_i$  happening at  $t_{ei}$  of a game,  
 382 we try to predict its event type via the multimodal  
 383 information provided in our dataset and the game  
 384 situation understanding module, and generate a  
 385 commentary summary via the game commentary  
 386 summarisation module. Taking  $m$  most recent  
 387 game events which happened before  $E_i$  as a histor-  
 388 ical reference, we extract the time-series event se-  
 389 quence as  $\mathbb{E} = \{E_{i-m}, \dots, E_{i-2}, E_{i-1}\}$ . Assum-  
 390 ing that the input window size for transcript and  
 391 chat is  $w$ , we extract a time-series sequence consist-  
 392 ing of  $x$  transcript clips  $\mathbb{T} = \{T_{s-x}, \dots, T_{s-1}, T_s\}$ ,  
 393 where  $T_s$  refers to the  $s$ -th transcript clip in the cur-  
 394 rent game. These clips fully cover the time period  
 395 from  $(t_{ei} - w)$  to  $t_{ei}$ , meaning that the timestamp  
 396  $(t_{ei} - w)$  falls within the time frame covered by

<sup>4</sup><https://www.frankerfacez.com/emoticons/>

<sup>5</sup><https://github.com/carpedm20/emoji/>

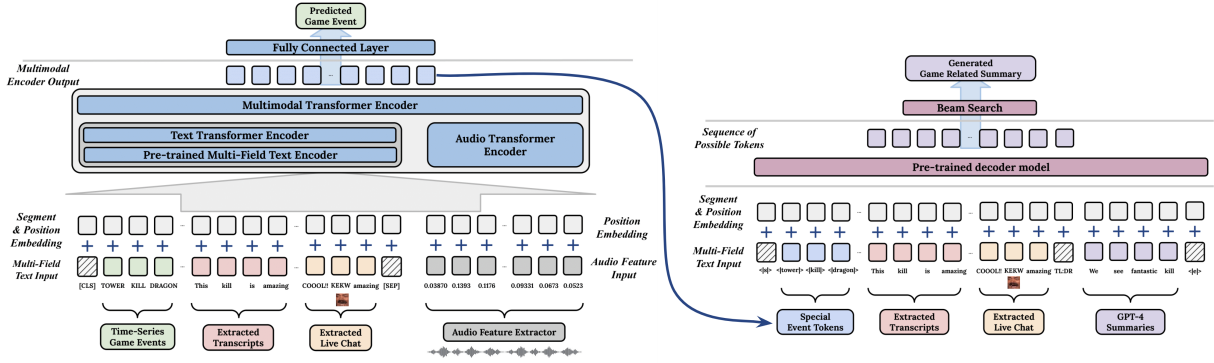


Figure 3: Joint integration framework of Game Situation Understanding and Game Commentary Summarisation

397  $T_{s-x}$ , and  $t_{ei}$  falls within the time frame covered 398 by  $T_s$ . The time-series sequence of chats  $\mathbb{C}$  is 399 extracted based on their specific timestamps between 400  $(t_{ei} - w)$  and  $t_{ei}$ . For the audio component, given 401 the window size  $w_a$ , the audio feature sequence 402 is extracted as  $\mathbb{A}$  within the time period between 403  $(t_{ei} - w_a)$  and  $t_{ei}$ . This results in a vector consist- 404 ing of  $w_a * 50$  values that serve as the input for the 405 audio transformer, given that the audio features are 406 sampled at a rate of 50Hz.

## 407 5.2 Game Situation Understanding

408 The model architecture is shown on the left of Fig- 409 ure 3. On the text side, the input is a combination 410 of multi-field sequential time-series data from pre- 411 vious event  $\mathbb{E}$ , caster transcript  $\mathbb{T}$  and audience 412 chat  $\mathbb{C}$ , with graphical emotional expressions in 413 chats being converted into their text representa- 414 tion. Since chats tend to contain many repetitions 415 in phrases and emotions, we truncate the input se- 416 quence up to 256 tokens. Following the approaches 417 in BERT (Devlin et al., 2019), we insert a [CLS] 418 token at the beginning and a [SEP] token at the 419 end of the input sequence, creating the input em- 420 beddings by summing the token, segment, and po- 421 sition embeddings. These input embeddings are 422 initially passed into a pre-trained multi-field text 423 encoder. The [CLS] token output from this pre- 424 trained multi-field text encoder is then forwarded 425 to the text transformer encoder to project the text 426 representation into a common space. On the audio 427 side, the combination of audio feature  $\mathbb{A}$  and po- 428 sition embedding are fed into an audio transformer, 429 which maps the audio into the same common space 430 as the text. The text and audio representations are 431 then concatenated to form a single vector, which 432 serves as the input for the multimodal transformer 433 encoder followed by a fully connected layer to pre-

dict the subsequent game event. We take advantage 394 of existing pre-trained models in our multi-field 395 text encoder including BERT (Devlin et al., 2019), 396 RoBERTa (Liu et al., 2019), DeBERTa (He et al., 397 2021), and XLNet (Yang et al., 2019). More details 398 can be found in Section 6.1. 399

## 440 5.3 Game Commentary Summarisation

441 After fine-tuning the game situation understand- 442 ing model, we obtain the event representations from 443 it before the fully connected layer and incorpo- 444 rate these representations along with transcrip- 445 ts and chats into the pre-trained generative 446 model for summarisation. We calculate the mean of 447 each event representation by inference the trained 448 game situation understanding model with all the 449 matches in our dataset to get the special event em- 450 beddings. These embeddings are then added to the 451 decoder models' vocabulary as  $\langle \text{kill} \rangle$ ,  $\langle \text{tower} \rangle$  452 and  $\langle \text{dragon} \rangle$  to enhance efficiency during sum- 453 mary generation. Similar to the encoder model, 454 we truncate the chat sequence up to 256 tokens for 455 emotion extraction before combination them with 456 special event tokens and transcripts. As shown 457 in the right of Figure 3, a special [TL;DR] token 458 and GPT-4 summary are concatenated to the se- 459 quence as a reference during fine-tuning. We utilise 460 two different pre-trained decoders, including GPT- 461 2 (Radford et al., 2019) and Pythia (Biderman et al., 462 2023). More details can be found in Section 6.1. 463

## 463 6 Experiments and Results

### 464 6.1 Experiment Setup

465 **Game Situation Understanding** We test four 466 pre-trained encoder models with their large 467 settings as the baseline multi-field text en- 468 coders: BERTLARGE, RoBERTaLARGE, DeBER- 469 TaV3LARGE, and XLNetLARGE. The text and au-

Chat	Audio	Game Events	BERT				DeBERTaV3				RoBERTa				XLNet			
			Kill	Tower	Dragon	All	Kill	Tower	Dragon	All	Kill	Tower	Dragon	All	Kill	Tower	Dragon	All
✗	✗	✗	77.98	47.75	8.45	61.97	79.46	62.16	1.41	65.06	79.17	62.16	8.45	65.83	93.75	10.81	4.23	63.71
✓	✗	✗	86.01	20.72	9.86	61.58	81.55	62.16	0.00	66.22	79.46	59.46	7.04	65.25	96.43	0.90	5.63	63.51
✗	✓	✗	83.63	37.84	14.08	64.29	77.08	36.94	49.30	64.67	78.57	62.16	25.35	67.76	72.02	55.86	22.54	61.78
✗	✗	✓	80.55	51.35	17.19	64.96	72.35	61.26	17.19	62.18	78.50	58.56	35.94	67.95	95.22	15.32	0.00	63.25
✓	✓	✗	75.00	48.65	11.27	60.62	82.44	43.24	43.66	68.73	77.38	63.06	15.49	65.83	67.86	53.15	14.08	57.34
✓	✗	✓	83.22	51.35	18.03	66.81	81.82	58.56	40.98	70.74	80.07	55.86	36.07	68.34	80.07	62.16	21.31	67.90
✗	✓	✓	84.97	32.43	42.62	66.59	79.72	49.55	34.43	66.38	84.62	51.35	26.23	68.78	76.22	60.36	21.31	65.07
✓	✓	✓	83.57	43.24	18.03	65.07	86.71	31.53	59.02	69.65	80.42	52.25	31.15	67.03	83.92	53.15	18.03	67.69

Table 4: The effect of Chat, Audio and previous Game Events on 2 different Game Situation Understanding Models.

470 dio transformer encoder and the multimodal trans-  
471 former encoder are all 8-head and 6-layer encoder  
472 structures and 1024 embedding dimension. The  
473 entire model is trained using AdamW (Loshchilov  
474 and Hutter, 2019) with 2 epochs for each instance,  
475 with a dropout value of 0.1 (Srivastava et al., 2014),  
476 a learning rate of 1e-6, and a learning rate decay  
477 rate of 0.95 for every 2 epochs. **Game Comment-**  
478 **ary Summarisation** We adopt two pre-trained de-  
479 coder models as the baseline commentary summari-  
480 sation models: 762M GPT2 with 1280 dimension  
481 size and 410M Pythia with 1024 embedding size.  
482 We apply Principal Component Analysis (Wold  
483 et al., 1987) to the game event embeddings when  
484 their dimensions are larger than the embeddings  
485 of pre-trained models for fine-tuning consistency.  
486 All models are trained using AdamW for 3 epochs,  
487 with a learning rate of 1e-5, and a warmup step of 5.  
488 Our implementations are based on PyTorch (Paszke  
489 et al., 2019) and HuggingFace Transformers (Wolf  
490 et al., 2020), with the help of Scikit-learn (Buit-  
491 inck et al., 2013). All experiments are run on a test  
492 bench with 24GB NVIDIA RTX 3090 GPU.

## 493 6.2 Evaluation Metrics

494 We evaluate the game situation understanding  
495 model with a multi-class accuracy metric, directly  
496 comparing the predicted game event with the  
497 ground truth for each event class. Generated sum-  
498 maries are evaluated with ROUGE (Lin, 2004) and  
499 BERTScore (Zhang et al., 2020), common auto-  
500 matic evaluation metrics. To have the best correla-  
501 tion with humans, we choose a RoBERTaLARGE  
502 version of BERTScore, which deploys a RoBERTa  
503 model to compare the similarity between the model  
504 generations and references. All results are reported  
505 for a single run of the experiments.

## 506 6.3 Results

507 **Overall Performance** As illustrated in Table 4,  
508 when all input features are utilised, DeBERTaV3  
509 notably outperforms the others in overall accuracy  
510 as well as **Kill** and **Dragon** categories by trading

Special Event Token	GPT2		Pythia	
	BertScore	ROUGE-L	BertScore	ROUGE-L
✗	76.15	18.52	74.45	13.24
✓	76.38	17.10	75.37	15.98

Table 5: The effect of special event tokens on 2 different Game Commentary Summarisation Models.

511 off the performance on **Tower**. Trailing behind De-  
512 BERTaV3, the overall performance of RoBERTa  
513 and XLNet is similar, with a margin difference of  
514 less than 1%. It is worth noting that RoBERTa  
515 excels in the **Dragon** category, while XLNet ex-  
516 cels in the **Kill** and **Tower** categories. Although  
517 BERT achieves an overall accuracy of 65.07%, it  
518 ranks last among the four encoder variants. This is  
519 likely attributable to the other models’ more robust  
520 optimisation built upon BERT’s architecture. In  
521 addition, all models produce better prediction ac-  
522 curacy for **Kill** than for **Tower** and **Dragon**. This  
523 trend is primarily due to the imbalanced event data  
524 since the average number of **Kill** instances per  
525 match is 25.69, which is double the average number  
526 of **Tower** instances (11.62) and triple the average  
527 number of **Dragon** instances (7.62). Regarding  
528 the game commentary summarisation results pre-  
529 sented in Table 5, we note that GPT2 consistently  
530 outperforms Pythia across both evaluation metrics,  
531 irrespective of special event tokens.

532 **Ablation Studies** To further analyse the effective-  
533 ness of our data, we conduct ablation studies to  
534 compare 3 different input combinations with tran-  
535 scriptions for the game situation understanding model:  
536 **1) Audio:** with and without audio features as part of  
537 the sequence input; **2) Chat:** with and without chat  
538 as part of the sequence input; **3) Game Events:**  
539 with and without game events as part of the se-  
540 quence input. The results are presented in Table 4.  
541 We observed that supplementing the model with ad-  
542 ditional input data improves its capability for under-  
543 standing game situations. This results in a notice-  
544 able performance increase across all three models,  
545 particularly for the rare **Dragon** event, albeit with

Audio	BERT				DeBERTaV3			
	Kill	Tower	Dragon	All	Kill	Tower	Dragon	All
5s	80.55	42.34	25.00	63.89	83.62	35.14	57.81	68.59
10s	83.62	37.84	23.44	64.53	83.28	35.14	59.38	68.59
15s	84.30	40.54	18.75	<b>64.96</b>	86.35	28.83	57.81	<b>68.80</b>

Table 6: Hyperparameter testing on the Game Situation Understanding Models for different audio time windows (rounded to the nearest integer in order to obtain enough data to match the audio transformer embedding size which should be a multiple of 8), where input transcript and chat time windows are 30s, and the number of previous game events is 5. A larger audio time window may lead to higher performance with a small margin.

a slight trade-off in performance for other events. Specifically, individually incorporating audio or previous game events into the transcript yields a greater improvement than adding chat data alone. Furthermore, combining two types of additional inputs surpasses the performance achieved with just a single extra input. We also conduct experiments both in the presence and absence of the **Special Event Token**, defined as the intermediate embedding before the fully connected layer within the game situation understanding model, as illustrated in Figure 3. Other inputs, such as transcripts, chats, and GPT-4 summaries, are essential for fine-tuning since omitting any of these causes a significant drop in generation performance. The results of these experiments are shown in Table 5. We observed the addition of a special event token can guide model generation, leading to improvements in BertScore for both GPT2 and Pythia.

**Hyperparameter Testing** The audio hyperparameter testing for the three different variations of the Game Situation Understanding Model is in Table 6, where input transcript and chat time windows are set to 30 seconds, and the number of previous game events are set to 5. We observe that the performance of each model is barely influenced by the input length of the audio features, as the difference is within a 1% margin. We also explore the effectiveness of different numbers of previous game events and results are shown in Table 7, where input transcript and chat time windows are set to 30 seconds, and the audio time window is set to 15 seconds. Increasing the number of previous game events improves the models’ aggregate performance up until a specific threshold. However, it is observed that when this threshold is surpassed, there is a discernible decrement in performance. We hypothesise that the performance decline is due

Game Events	BERT				DeBERTaV3			
	Kill	Tower	Dragon	All	Kill	Tower	Dragon	All
3	85.39	29.73	18.84	63.32	88.64	25.23	53.62	69.26
5	84.30	40.54	18.75	<b>64.96</b>	86.35	28.83	57.81	68.80
7	80.58	40.91	21.67	62.95	85.25	42.73	56.67	<b>70.98</b>
9	79.32	50.00	12.50	63.32	71.80	64.15	0.00	60.51

Table 7: Hyperparameter testing on the Game Situation Understanding Model for different numbers of previous game events, where input transcript and chat time windows are 30s and the audio time window is 15s. A large number of previous game events may include less relevant histories and lead to a worse performance.

Category	GPT2			Pythia		
	Event	Coherence	Overall	Event	Coherence	Overall
<b>Kill</b>	75.31%	75.31%	66.67%	24.69%	24.69%	33.33%
<b>Tower</b>	60.74%	59.26%	59.26%	39.26%	40.74%	40.74%
<b>Dragon</b>	61.62%	66.67%	59.60%	38.38%	33.33%	40.40%
<b>All</b>	64.76%	65.71%	61.27%	35.24%	34.29%	38.73%

Table 8: Human evaluation comparison between GPT2 and Pythia summaries. Appendix D and Figure 4 show the details for event inclusion, coherence and overall quality. GPT2 gains better support from human annotators across all 3 aspects compared to Pythia.

to the extended length of the previous events, which have less correlation with the target event.

**Human Evaluation** Automatic metrics may not correlate well with human judgments in different aspects (Durmus et al., 2020), therefore we conduct the human evaluation to enrich the comprehensiveness of the results. We randomly collected testing samples for evaluating the summaries from GPT2 and Pythia and recruited nine workers, all with general background knowledge of League of Legends for evaluation, resulting in 1,890 instances of human feedback. As shown in Table 8, summarisations of GPT2 are more preferred by humans in all categories which aligns with the results from automatic evaluation metrics.

## 7 Conclusions

We introduce GAME-MUG, a multimodal dataset for game situation understanding and game commentary generation, and propose a joint integration baseline model. It contains diverse game-related information from game event logs, caster comments, audience conversations and caster speech audio. The combination of multimodal data improves the model’s understanding of the game situations while providing the game situation information leads to more human-like game commentary generation. Finally, we will make our dataset publicly available, hoping it will lead to novel applications.



## 612 Limitations

613 In this work, we only consider League of Legends  
614 as the representation of the MOBA game due to  
615 its popularity (Duan et al., 2023). This constrains  
616 the range of game scenarios covered by GAME-  
617 MUG and consequently limits the scope of poten-  
618 tial applications built upon it. We encourage future  
619 studies to incorporate a variety of MOBA games  
620 to further enrich the diversity of game situations.  
621 Furthermore, we used GPT-generated summaries  
622 for annotation in our dataset. We may apply other  
623 generative AI if new models emerge later on.

## 624 Ethics Statement

625 All the experiments strictly follow the Code of  
626 Ethics. In Section 6.3 human evaluation and Ap-  
627 pendix, we include the instructions and screenshots  
628 of the interface in the human evaluation and report  
629 the background of human judges. More detailed  
630 information about the recruitment process will be  
631 shared after the paper acceptance. We inform the  
632 human evaluators what the task is about and tell  
633 them that their responses will be used to assess the  
634 ability of language generation models.

## 635 References

636 Stella Biderman, Hailey Schoelkopf, Quentin Anthony,  
637 Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mo-  
638 hammad Aflah Khan, Shivanshu Purohit, USVSN Sai  
639 Prashanth, Edward Raff, Aviya Skowron, Lintang  
640 Sutawika, and Oskar van der Wal. 2023. *Pythia:  
641 A suite for analyzing large language models across  
642 training and scaling.*

643 Lars Buitinck, Gilles Louppe, Mathieu Blondel, Fabian  
644 Pedregosa, Andreas Mueller, Olivier Grisel, Vlad  
645 Niculae, Peter Prettenhofer, Alexandre Gramfort,  
646 Jaques Grobler, Robert Layton, Jake VanderPlas, Ar-  
647 naud Joly, Brian Holt, and Gaël Varoquaux. 2013.  
648 API design for machine learning software: experi-  
649 ences from the scikit-learn project. In *ECML PKDD  
650 Workshop: Languages for Data Mining and Machine  
651 Learning*, pages 108–122.

652 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
653 Kristina Toutanova. 2019. *BERT: Pre-training of  
654 deep bidirectional transformers for language under-  
655 standing.* In *Proceedings of the 2019 Conference of  
656 the North American Chapter of the Association for  
657 Computational Linguistics: Human Language Tech-  
658 nologies, Volume 1 (Long and Short Papers)*, pages  
659 4171–4186, Minneapolis, Minnesota. Association for  
660 Computational Linguistics.

661 Cem Dođdu, Thomas Kessler, Dana Schneider, Maha  
662 Shadaydeh, and Stefan R. Schweinberger. 2022. A

comparison of machine learning algorithms and fea-  
663 ture sets for automatic vocal emotion recognition in  
664 speech. *Sensors*, 22(19). 665

Peng Duan, Xiaohui Wang, Allan Yijia Zhang, and Bin  
666 Ji. 2023. Case analysis on world’s top e-sports events.  
667 In *Electronic Sports Industry in China: An Overview*,  
668 pages 119–133. Springer. 669

Esin Durmus, He He, and Mona Diab. 2020. *FEQA: A  
670 question answering evaluation framework for faith-  
671 fulness assessment in abstractive summarization.* In  
672 *Proceedings of the 58th Annual Meeting of the Asso-  
673 ciation for Computational Linguistics*, pages 5055–  
674 5070, Online. Association for Computational Lin-  
675 guistics. 676

Florian Eyben, Klaus R. Scherer, Björn W. Schuller,  
677 Johan Sundberg, Elisabeth André, Carlos Busso,  
678 Laurence Y. Devillers, Julien Epps, Petri Laukka,  
679 Shrikanth S. Narayanan, and Khiet P. Truong. 2016.  
680 *The geneva minimalistic acoustic parameter set  
681 (gemaps) for voice research and affective computing.*  
682 *IEEE Transactions on Affective Computing*, 7(2):190–  
683 202. 684

Florian Eyben, Martin Wöllmer, and Björn Schuller.  
685 2010. *Opensmile: The munich versatile and fast  
686 open-source audio feature extractor.* In *Proceedings  
687 of the 18th ACM International Conference on Mul-  
688 timedia, MM ’10*, page 1459–1462, New York, NY,  
689 USA. Association for Computing Machinery. 690

Pengcheng He, Xiaodong Liu, Jianfeng Gao, and  
691 Weizhu Chen. 2021. *Deberta: Decoding-enhanced  
692 bert with disentangled attention.* 693

Tatsuya Ishigaki, Goran Topic, Yumi Hamazono, Hi-  
694 roshi Noji, Ichiro Kobayashi, Yusuke Miyao, and Hi-  
695 roya Takamura. 2021. *Generating racing game com-  
696 mentary from vision, language, and structured data.*  
697 In *Proceedings of the 14th International Conference  
698 on Natural Language Generation*, pages 103–113,  
699 Aberdeen, Scotland, UK. Association for Computa-  
700 tional Linguistics. 701

Miriam Kienast and Walter F. Sendlmeier. 2000. Acous-  
702 tical analysis of spectral and temporal changes in  
703 emotional speech. In *Proc. ITRW on Speech and  
704 Emotion*, pages 92–97. 705

Klaus Krippendorff. 2011. Computing krippendorff’s  
706 alpha-reliability. 707

Chengxi Li, Sagar Gandhi, and Brent Harrison. 2019.  
708 *End-to-end let’s play commentary generation using  
709 multi-modal video representations.* In *Proceedings of  
710 the 14th International Conference on the Foundations  
711 of Digital Games, FDG ’19*, New York, NY, USA.  
712 Association for Computing Machinery. 713

Chin-Yew Lin. 2004. *ROUGE: A package for auto-  
714 matic evaluation of summaries.* In *Text Summariza-  
715 tion Branches Out*, pages 74–81, Barcelona, Spain.  
716 Association for Computational Linguistics. 717

718	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	Tsunehiko Tanaka and Edgar Simo-Serra. 2021. <a href="#">Lol-</a>	775
719	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	v2t: <a href="#">Large-scale esports video description dataset</a> .	776
720	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	In <i>2021 IEEE/CVF Conference on Computer Vision</i>	777
721	<a href="#">Roberta: A robustly optimized bert pretraining ap-</a>	<i>and Pattern Recognition Workshops (CVPRW)</i> , pages	778
722	<a href="#">proach</a> .	4552–4561.	779
723	Ilya Loshchilov and Frank Hutter. 2019. <a href="#">Decoupled</a>	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann	780
724	<a href="#">weight decay regularization</a> . In <i>7th International</i>	Dubois, Xuechen Li, Carlos Guestrin, Percy Liang,	781
725	<i>Conference on Learning Representations, ICLR 2019,</i>	and Tatsunori B. Hashimoto. 2023. <a href="#">Stanford alpaca:</a>	782
726	<i>New Orleans, LA, USA, May 6-9, 2019</i> . <a href="#">OpenRe-</a>	An instruction-following llama model. <a href="https://github.com/tatsu-lab/stanford_alpaca">https://</a>	783
727	<a href="#">view.net</a> .	<a href="https://github.com/tatsu-lab/stanford_alpaca">github.com/tatsu-lab/stanford_alpaca</a> .	784
728	OpenAI. 2023. <a href="#">Gpt-4 technical report</a> .	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	785
729	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz	786
730	Carroll L. Wainwright, Pamela Mishkin, Chong	Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all</a>	787
731	Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray,	<a href="#">you need</a> .	788
732	John Schulman, Jacob Hilton, Fraser Kelton, Luke	Zihan Wang and Naoki Yoshinaga. 2022. <a href="#">Esports data-</a>	789
733	Miller, Maddie Simens, Amanda Askell, Peter Welin-	<a href="#">to-commentary generation on large-scale data-to-text</a>	790
734	der, Paul F. Christiano, Jan Leike, and Ryan Lowe.	<a href="#">dataset</a> .	791
735	2022. <a href="#">Training language models to follow instruc-</a>	Svante Wold, Kim Esbensen, and Paul Geladi. 1987.	792
736	<a href="#">tions with human feedback</a> . In <i>NeurIPS</i> .	<a href="#">Principal component analysis</a> . <i>Chemometrics and</i>	793
737	Adam Paszke, Sam Gross, Francisco Massa, Adam	<i>Intelligent Laboratory Systems</i> , 2(1):37–52. Proceed-	794
738	Lerer, James Bradbury, Gregory Chanan, Trevor	ings of the Multivariate Statistical Workshop for Ge-	795
739	Killeen, Zeming Lin, Natalia Gimelshein, Luca	ologists and Geochemists.	796
740	Antiga, Alban Desmaison, Andreas Kopf, Edward	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	797
741	Yang, Zachary DeVito, Martin Raison, Alykhan Te-	Chaumond, Clement Delangue, Anthony Moi, Pier-	798
742	jani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang,	ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,	799
743	Junjie Bai, and Soumith Chintala. 2019. <a href="#">Pytorch:</a>	Joe Davison, Sam Shleifer, Patrick von Platen, Clara	800
744	<a href="#">An imperative style, high-performance deep learning</a>	Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le	801
745	<a href="#">library</a> . In <i>Advances in Neural Information Process-</i>	Scao, Sylvain Gugger, Mariama Drame, Quentin	802
746	<i>ing Systems 32</i> , pages 8024–8035. Curran Associates,	Lhoest, and Alexander M. Rush. 2020. <a href="#">Transform-</a>	803
747	Inc.	<a href="#">ers: State-of-the-art natural language processing</a> . In	804
748	F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel,	<i>Proceedings of the 2020 Conference on Empirical</i>	805
749	B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer,	<i>Methods in Natural Language Processing: System</i>	806
750	R. Weiss, V. Dubourg, J. Vanderplas, A. Passos,	<i>Demonstrations</i> , pages 38–45, Online. Association	807
751	D. Cournapeau, M. Brucher, M. Perrot, and E. Duch-	for Computational Linguistics.	808
752	esnay. 2011. <a href="#">Scikit-learn: Machine learning in</a>	Jingyao Wu, Ting Dang, Vidhyasaharan Sethu, and	809
753	<a href="#">Python</a> . <i>Journal of Machine Learning Research</i> ,	Eliathamby Ambikairajah. 2021. <a href="#">Multimodal affect</a>	810
754	12:2825–2830.	<a href="#">models: An investigation of relative salience of audio</a>	811
755	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brock-	<a href="#">and visual cues for emotion prediction</a> . <i>Frontiers in</i>	812
756	man, Christine McLeavey, and Ilya Sutskever. 2022.	<i>Computer Science</i> , 3.	813
757	<a href="#">Robust speech recognition via large-scale weak su-</a>	Junjie H. Xu, Yu Nakano, Lingrong Kong, and Kojiro	814
758	<a href="#">pervision</a> .	Iizuka. 2023. <a href="#">Cs-lol: A dataset of viewer comment</a>	815
759	Alec Radford, Jeff Wu, Rewon Child, David Luan,	<a href="#">with scene in e-sports live-streaming</a> . In <i>Proceedings</i>	816
760	Dario Amodei, and Ilya Sutskever. 2019. <a href="#">Language</a>	<i>of the 2023 Conference on Human Information In-</i>	817
761	<a href="#">models are unsupervised multitask learners</a> .	<i>teraction and Retrieval, CHIIR '23</i> , page 422–426,	818
762	Charles Ringer, James Alfred Walker, and Mihalis A.	New York, NY, USA. Association for Computing	819
763	Nicolaou. 2019. <a href="#">Multimodal joint emotion and game</a>	<i>Machinery</i> .	820
764	<a href="#">context recognition in league of legends livestreams</a> .	Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Car-	821
765	In <i>2019 IEEE Conference on Games (CoG)</i> , pages	bonell, Russ R Salakhutdinov, and Quoc V Le. 2019.	822
766	1–8.	<a href="#">Xlnet: Generalized autoregressive pretraining for lan-</a>	823
767	Shukan Shah, Matthew Guzdial, and Mark O Riedl.	<a href="#">guage understanding</a> . In <i>Advances in Neural Infor-</i>	824
768	2019. <a href="#">Automated let’s play commentary</a> . <i>arXiv</i>	<i>mation Processing Systems</i> , volume 32. Curran Asso-	825
769	<a href="#">preprint arXiv:1909.02195</a> .	ciates, Inc.	826
770	Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky,	Huanyu Yu, Shuo Cheng, Bingbing Ni, Minsi Wang,	827
771	Ilya Sutskever, and Ruslan Salakhutdinov. 2014.	Jian Zhang, and Xiaokang Yang. 2018. <a href="#">Fine-grained</a>	828
772	<a href="#">Dropout: A simple way to prevent neural networks</a>	<a href="#">video captioning for sports narrative</a> . In <i>2018</i>	829
773	<a href="#">from overfitting</a> . <i>Journal of Machine Learning Re-</i>	<i>IEEE/CVF Conference on Computer Vision and Pat-</i>	830
774	<i>search</i> , 15(56):1929–1958.	<i>tern Recognition</i> , pages 6006–6015.	831

832 Dawei Zhang, Sixing Wu, Yao Guo, and Xiangqun Chen.  
833 2022. [MOBA-E2C: Generating MOBA game com-](#)  
834 [mentaries via capturing highlight events from the](#)  
835 [meta-data](#). In *Findings of the Association for Com-*  
836 *putational Linguistics: EMNLP 2022*, pages 4545–  
837 4556, Abu Dhabi, United Arab Emirates. Association  
838 for Computational Linguistics.

839 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.  
840 Weinberger, and Yoav Artzi. 2020. [Bertscore: Evalu-](#)  
841 [ating text generation with BERT](#). In *8th International*  
842 *Conference on Learning Representations, ICLR 2020,*  
843 *Addis Ababa, Ethiopia, April 26-30, 2020*. OpenRe-  
844 view.net.

## A De-Identification of User Chats

To ensure the anonymity and privacy of individuals involved in the live chats, we implemented a de-identification protocol. The primary objective of this protocol is to mask any information that could potentially reveal the identity of a chat participant. We directly remove all original usernames associated with the chats, ensuring it is infeasible to reverse engineer the original usernames. All de-identified chats are stored in plain text format, without any identifying information. The original raw data are permanently deleted after the de-identification process. By taking these steps, we ensure that our data collection and analysis processes align with ethical guidelines and data protection regulations.

## B Prompt Design

Algorithm 1 illustrates the approach for querying the GPT-4 API. We set the background information as watching a live game streaming via a system prompt. Whenever a game event occurs, we forward the commentary and live chat content to the GPT-4 API through the summary prompts. We design several prompt parameters to guide the GPT-4 generation: `<game streaming platform>` indicates different live stream platforms, `<number of summary words>` control the number of generated words, and `<game-related topics>` adjusts the generated summary to focus on different aspects, such as on player, character, event or overall situation.

## C Full Hyperparameter Testing Results

The complete hyperparameter results are displayed in Table 9. We conducted experiments using 15s and 30s time windows for transcripts and chats, and 5s, 10s, and 15s time windows for audio. Additionally, we experimented with time-series events ranging from 3 to 10.

## D Human Evaluation

We recruited nine volunteers aged between 25 and 30, all holding at least a Bachelor’s degree, to participate in the human evaluation. The group was composed of three females and six males, each with a general understanding of League of Legends. While one participant was a native English speaker, the other eight were proficient in English. For the human evaluation survey, participants were presented with the original transcript, the truncated

chat, and the generated summaries from the baseline models. They are then asked to rank the summaries based on the following four criteria:

- **Game Event Information:** The quality of summaries in terms of the game event-related expressions.
- **Coherence:** The quality of summaries in terms of fluency and logic.
- **Overall:** The overall quality of summaries regarding the above criteria and any other game-related criteria.

The sample evaluation questions are shown in Figure 4.



**Evaluation Sample**

**Original commentary:**  
 His kindred as Viego doing fantastic so far seven out of eight kills dragon spawning in ten. Zekka is going to take the base barrels coming out of it with wards. Gen-G actually just took a base as well. Top and bot lane in the nexus towers area just running out of base. I think they'll be way too late to contest this dragon so TRX should be able to pick this one up pretty easily. Will they

**Audience Chat:**

1. DEFT GIGACHAD oneandonlyNasusWow oneandonlyNasusWow CHOVIY CS KEKW ICANT
2. ZEKAA IS 19 YEARS OLD I THINK SEND Prayge THIS Prayge BLESS Prayge TO Prayge SAVE Prayge CHOVIY Prayge CS DEFT GIGACHAD YUHAN KEKW emily rand ITEM ??? ???? chovy cs NO FLASH DEFT GIGACHAD YUHAN KEKW BigBrother COME TO KANSAS BigBrother IM A PROBLEM BigBrother
3. shurelylias? YOOHAN chovy cs xdd CHOVIY went ludens L0000L KEKWWait CHOVIY CS monkaS deft
4. Pyoshik is rolling 20s every game Pog

**Summary 1:**  
 Viega leapsfrogs GenG securing first dragon of the game while onlookers cheer on Pog EZ

**Summary 2:**  
 entschied for a thrilling fight as DRX secures dragon and tower audience goes wild

**Evaluation Questions**

How well you think those summaries in terms of containing game event information about **'Dragon'**?

Please provide ranking for these summaries above from 1 to 2, where 1 is the **better** and 2 is the **worse**.

	1	2
Summary 1	<input type="radio"/>	<input type="radio"/>
Summary 2	<input type="radio"/>	<input type="radio"/>

How well you think those summaries in terms of **fluency**?

Please provide ranking for these summaries above from 1 to 2, where 1 is the **better** and 2 is the **worse**.

	1	2
Summary 1	<input type="radio"/>	<input type="radio"/>
Summary 2	<input type="radio"/>	<input type="radio"/>

Please rank these summaries **overall** qualities above from 1 to 2, where 1 is the **better** and 2 is the **worse**.

	1	2
Summary 1	<input type="radio"/>	<input type="radio"/>
Summary 2	<input type="radio"/>	<input type="radio"/>

Figure 4: Screenshot of a human evaluation sample. Workers are shown the original commentary with truncated audience chats on the top left. We provide the generated summarisations on the bottom left. The worker ranks these two summarisations in terms of the inclusion of the game event, coherence and overall quality.