# Esports Data-to-commentary Generation on Large-scale Data-to-text Dataset

**Zihan Wang**[1] and **Naoki Yoshinaga**[2]

[1]Graduate School of Information Science and Technology, The University of Tokyo
[2]Institute of Industrial Science, The University of Tokyo
zwang@tkl.iis.u-tokyo.ac.jp, ynaga@iis.u-tokyo.ac.jp

## Abstract

Esports, a sports competition using video games, has become one of the most important sporting events in recent years. Although the amount of esports data is increasing than ever, only a small fraction of those data accompanies text commentaries for the audience to retrieve and understand the plays. Therefore, in this study, we introduce a task of generating game commentaries from structured data records to address the problem. We first build a large-scale esports data-to-text dataset using structured data and commentaries from a popular esports game, League of Legends. On this dataset, we devise several data preprocessing methods including linearization and data splitting to augment its quality. We then introduce several baseline encoder-decoder models and propose a hierarchical model to generate game commentaries. Considering the characteristics of esports commentaries, we design evaluation metrics including three aspects of the output: correctness, fluency, and strategic depth. Experimental results on our large-scale esports dataset confirmed the advantage of the hierarchical model, and the results revealed several challenges of this novel task.

Screenshot of "`CHAMPION_KILL`" event:



**Input** (one-minute structured data (40 events); an excerpt):

```
{...
 {"type": "CHAMPION_KILL",
  "timestamp": 191394,
  "position": {"x": 4511, "y": 13554},
  "killedId": 1,
  "victimId": 6,
  "assistingParticipantIds": [2, 3]}}
 ...
}
```

**Output** (play-by-play commentary):
. . . lwx gets tagged the death sentence pulled back Mickey the first blood King from Europe but what more can they get one that continues to chase one of the realm . . .

Figure 1: Game commentary based on one-minute play-by-play data records (a series of events in JSON format).

## 1 Introduction

The benefits of video gaming have been brought to the attention of the public in recent years (West et al., 2017). Some video games have matured into having sports competitions called esports, also known as electronic sports or e-sports, where the audience watch players playing video games (Reitman et al., 2020; Hamari and Sjöblom, 2017). Esports contests have become one of the important sporting events in recent years; the esports game "League of Legends (LoL)," as one of the esports games with the most peak views[1] (Ringer et al., 2022), has attracted more than 40M audience all over the world in its 2021 World Championship. LoL is also a demonstration sports event in the 2018 and 2022 Asian Games for its popularity (Hallmann and Giel, 2018; Jenny et al., 2017).

Despite the audience's great enthusiasm, the audience often meet much inconvenience to enjoy watching the games. The whole championship contest includes many individual games, and the average time of the games is long. For example, LoL's 2019 World Championship contest consists of more than 200 individual games with an average time of more than 40 minutes. The audience have trouble understanding the players' important actions in the

---

[1]https://escharts.com/

games, because esports games often have too much focus at once. To fully enjoy watching esports games and to effectively learn playing strategies from skillful players, textual game commentaries are beneficial for the audience to understand the player's actions more easily (Lavelle, 2010). For example, Figure 1 shows a screenshot of LoL and its corresponding commentary, through which the audience can understand the player's actions more easily. However, since it is very costly for human experts to provide play-by-play commentaries for individual games, only a fraction of esports games are provided with paired textual commentaries.

To improve the audience's experience in watching esports games, we propose a method to generate game commentaries from esports structured data records. In the literature, data-to-text generation has been applied in game commentaries generation such as basketball (Wiseman et al., 2017) and chess (Modgil et al., 2013). Compared to basketball data, esports data has much greater lengths containing play-by-play commentaries. Compared to chess data, esports data is not turn-based. There is a lack of research considering the characteristics of esports data in existing studies, and thus this work is meant to fill this gap in research.

In this study, we introduce a task named "game commentary generation from structured data records," using LoL as the target esports game. The goal of the task is set to maximize the correctness, fluency, and strategic depth of the output. Broadly, the overall workflow of addressing this new task includes three main processes: we first build a large-scale data-to-text generation dataset from commentaries obtained from subtitles of YouTube contest videos on esports games and corresponding structured data records obtained using LoL official APIs; we then design a group of baseline encoder-decoder models and a hierarchical model; we also set evaluation metrics of esports data-to-commentary generation. To investigate the performance of the proposed method, we conduct evaluations based on the metrics introduced in the following sections. We also use visual examples of game scenes paired with golden commentaries and generated outputs to represent the performance of the model.

The contributions of this paper is as follows:

- We set up a task of data-to-commentary generation for one of the most popular esports games, League of Legends; we built and will release two versions of large-scale datasets to facilitate research activities on the task.
- We designed evaluation criteria for esports data-to-commentary generation, which reflects the features of strategy esports games.
- We have revealed an effective combination of network structures for encoder-decoder models for the esports data-to-commentary task, and confirmed the impact of a hierarchical encoder that captures the structured data.

## 2 Related Work

In this section, we review data-to-text generation tasks on physical sports, board games, and esports to highlight the characteristics of our esports data-to-commentary task on structured data records.

**Data-to-text Generation for Sports Games.** Many studies focus on transcribing game score data into textual summaries of the entire games (Wiseman et al., 2017; Puduppully et al., 2019; Rebuffel et al., 2020; Dou et al., 2018; van der Lee et al., 2017; Taniguchi et al., 2019; van der Lee et al., 2018), mainly including basketball and soccer games. The essential difference between these sports games and esports is that the basketball data only records certain important values (score, player number, etc.), while esports data provides much more detail about the games. As a result, the average length of commentaries per game is much longer for our LoL dataset than that for the basketball game data. This large gap in length creates even greater challenges for the esports data-to-commentary generation task.

**Data-to-text Generation for Turn-based Board Games.** There are studies focusing on generating play-by-play commentaries with grounded move expressions from chess and shogi (Japanese chess) games (Modgil et al., 2013; Kameko et al., 2015; Jhamtani et al., 2018). For both esports and chess, it is able to recreate a game from the data records. Nevertheless, in chess games, the data is recorded in individual turns, and the duration of the time between two turns is not taken into account. However, in esports games, the data is not turn-based, and thus this leads to the problem that esports games need to interpret multiple actions at the same time, which is also reflected by the esports data.

**Game Video Summarization from Visual Data.** Several researchers aim to produce textual summaries from a game video to provide a quick way to overview the game's full content (Pasunuru and

Bansal, 2018). For example, (Ishigaki et al., 2021) proposes a study concerning racing games, which explores commentary generation methods on game video data. These kinds of studies aim to explore the multimodal processing between visual and textual information, while ours focuses on the transcription from structured data records to natural language text. Additionally, some of these studies produce only brief summaries describing the game (Khan and Pawar, 2015), while ours focuses on more comprehensive commentaries with detailed descriptions and more strategic content.

## 3 Esports Data-to-text Dataset

The first contribution of this work is constructing a large-scale data-to-text dataset for a major esports game, League of Legends (LoL). The large-scale dataset is the core building block to assess the feasibility of data-to-text technology on the task. We build two datasets, called LoL19 and LoL19-21, from all games in the highest-level tournaments of the 2019 and 2019 to 2021 Season World Championship of League of Legends, respectively. Since the 2020 and 2021 Season World Championship of League of Legends are conducted under a special configuration under the COVID-19 pandemic, we refer to the regular-sized LoL19 dataset and the larger-sized LoL19-21 dataset as "core" and "extended" datasets.

In this section, we first introduce the basics of the target game, LoL, and then discuss methods of extracting the structured data records of LoL.[2] We also introduce several data preprocessing methods and the expansion to construct a large-scale dataset.

### 3.1 Basics of League of Legends

In this study, we choose LoL as our research target because of its great popularity and its representative role as an esports game.

In LoL games, each player controls one game character called "champion" with a set of unique abilities that improve over the course of a game and contribute to the team's overall strategy (Cannizzo and Ramírez, 2015). There are two teams competing in one game map, each containing five players. Both teams' goal is to destroy the opponent's base while protecting their own. In the game progress, the champions can kill the enemy's champions, kill non-player controlled characters called "monsters,"

| Event type | Explanation |
|---|---|
| ITEM_PURCHASED | The player purchases an item from the shop |
| ITEM_SOLD | The player sells an item |
| ITEM_UNDO | The player undoes an action of an item |
| ITEM_DESTROYED | The player destroyed an item |
| BUILDING_KILL | The team destroys an enemy building |
| CHAMPION_KILL | The player kills an enemy champion |
| SKILL_LEVEL_UP | The player upgrades a skill |
| WARD_PLACED | The player places a ward |
| WARD_KILL | The player kills an enemy ward |
| ELITE_MONSTER_KILL | The team kills an elite monster |

Table 1: Types of game events.

and destroy buildings to earn resources. They then use the resources to improve their abilities by purchasing items and upgrading abilities.

Different from physical sports games (e.g., basketball) and boards games (e.g., chess), the LoL mechanism is more complicated on diversed gameplay contents, complexity of rules, and scale of data records per game. These factors can be the main obstacles to data-to-text generation.

### 3.2 Data Extraction and Preprocessing

To build the core dataset named LoL19, we target the games of the highest-level tournament in the 2019 Season World Championship of League of Legends. We collect the structured data records of the game plays from the LoL official API site as input, and extract subtitles of YouTube[3] videos on the game plays as output commentaries. These data is strictly paired with the game IDs provided by the game's Match History site.[4]

Since the obtained pairs of structured data records and commentaries are too large, we split them into pieces and format the input structured data so that a common encoder-decoder model can be applied. Later, we will introduce the method to build a large-scale dataset with extensional data.

#### 3.2.1 Textual Commentaries

We collect the YouTube subtitles of LoL contest videos as the output commentaries of this task. The link to every contest game video can also be found on the Match History site. We apply sentence seg-

---

[2]Riot Developer Portal: https://developer.riotgames.com/

[3]https://www.youtube.com/
[4]https://lol.fandom.com/wiki/2019_Season_World_Championship/Match_History

|  | esports (our) | | basketball | chess | racing |
|---|---|---|---|---|---|
|  | **LoL19 (core)** | **LoL19-21 (extended)** | **RotoWire** | **GameKnot** | **Assetto Corsa** |
| The number of games | 220 | 650 | 4,853 | 11,578 | 1,389 |
| The number of examples | 3,490 | 10,590 | 4,853 | 298,008 | 2,473 |
| The average number of tokens (input) | 540.47 | 541.10 | 628.00 | 25.73 | - |
| The average number of tokens (output) | 374.68 | 373.89 | 337.10 | 20.55 | - |
| The number of event types (Table 1) | 10 | 10 | N/A | N/A | N/A |
| The average number of events per examples | 49.13 | 48.58 | N/A | N/A | N/A |

Table 2: Statistics of the esports data-to-text datasets and common data-to-text generation datasets, including RotoWire (Wiseman et al., 2017), GameKnot (Jhamtani et al., 2018), and Assetto Corsa (Ishigaki et al., 2021).

mentation to the obtained subtitles to obtain a sequence of sentential commentaries of each game.

### 3.2.2 Structured Data Records on Game Plays

In LoL esports games, every move made by each player - including mouse movement and key presses - is recorded, and these records are provided by the LoL official API site (RiotGames, 2021). In other words, from LoL data, we can strictly restore the entire game from the structured data records, which is impossible in physical sports games such as basketball and soccer.

However, the complete data has duplicate information, and this large data volume is a heavy burden for storage and the next processing. For this reason, we have chosen another data type provided by the official API, named "event-based data frame." In this data frame, only the key events in each game are recorded as shown in Figure 1. The event types and their corresponding explanations are listed in Table 1. Each event is defined by the official API as an update of certain game status, and includes the key named "type," which presents a general definition of the event. Note that different types of events usually have different sets of keys. For example, the "BUILDING_KILL" event involves information about "killerId" while the "ITEM_PURCHASED" event does not have this key. The event-based data frames of LoL game plays are stored in JSON format. An example of "CHAMPION_KILL" event is shown in Figure 1.

### 3.2.3 Data Splitting

Since the average length of the commentary of an esports game is over 10K words, which is much longer than the outputs of widely used data-to-text datasets, we cannot exploit a neural model with common network architecture such as Transformer (Vaswani et al., 2017) for generation.

We thus decompose the obtained pair of the sequence of structured data and the sequence of sentential commentaries for each game into pieces. We first split the sequence of events on the basis of unit duration of one minute of the game play. We then split the sequence of sentential commentaries by matching the timings of each sentential commentary with the duration of each subsequence of events.

We ultimately obtained the core dataset that contains 3,490 data-commentary pairs (Table 2). An example of input and output for one-minute portion of LoL game plays is excerpted in Figure 1.

### 3.2.4 Data Formatting

We next resolve the format difference between the input data (in JSON format; namely, nested list) and natural language text, to make common encoder-decoder models (Sutskever et al., 2014; Vinyals et al., 2016) applicable to our task. Concretely, we linearize the structured input. For each key-value pair in the top-level list of each event in the JSON format, we recursively place a concatenation of the value and the key with a delimiter "|" while inserting a space between individual key-value pairs. It is worth mentioning that this is done only at the top level and we keep the internal lists as they are in their original JSON format. Following this procedure, "CHAMPION_KILL" event in the JSON format (Figure 1) is linearized into the following sequence:

```
CHAMPION_KILL|type 191394|timestamp
{x:  4511, y:  13554}|position
1|killerId 6|victimId [2,
3]|assistingParticipantIds
```

### 3.2.5 Extended Dataset

We additionally obtain datasets from all games in the highest-level tournament of the 2020 and 2021 Season World Championship of League of Legends

to obtain the extended LoL19-21 dataset by the same measure. Both versions of the esports data-to-commentary generation datasets will both be published to the community.

In conclusion, Table 2 lists the statistics of both versions of our esports data-to-commentary generation dataset. It also includes the statistics of relevant data-to-text generation datasets (Wiseman et al., 2017; Jhamtani et al., 2018; Ishigaki et al., 2021) for comparison.

## 4 Esports Data-to-text Task

In this section, we detail the task settings, notations, and evaluation metrics of our esports data-to-commentary task.

### 4.1 Task Settings and Notations

The input of our esports data-to-commentary task is structured data on one-minute game plays, while the output is a textual commentary that helps understand the game plays. In what follows, we define the input and output of our task:

**Event** An event $e_j$ is a sequence of key-value pairs, $r_j^k$, namely, $e_j = \{r_j^1, r_j^2, \ldots, r_j^{N_j}\}$, where key-value pairs $r_j^k$ consist of key $k_j^k$ and value $v_j^k$ and $N_j$ is the number of key-value pairs at the top-level list in event $e_j$. We handle both key and value as raw strings.

**Input** An input to our data-to-commentary task, $x_i$, is a sequence of linearized events $e_i^k$, namely, $x_i = \{e_j^1, e_j^2, \ldots, e_j^{N_i}\}$, where $N_i$ is the number of events in input $x_i$. The events in input $x_i$ is ordered in an ascending order of the value for key "timestamp" of $e_i^k$.

**Output** The output for input $x_i$, $y_i$ is a textual commentary, which is denoted by sequence of tokens $c_i^k$, namely, $y_i = \{c_i^1, e_i^2, \ldots, c_i^{T_i}\}$, where $T_i$ is the number of tokens in $y_i$. We refer to the $p$-th token to the $q$-th token of a commentary $y_i$ as $c_i^{p:q}$. Therefore, the commentary $y_i$ can also be noted as $y_i = c_i^{1:T_i}$.

### 4.2 Evaluation Metrics

Following the existing data-to-text tasks on sports game summary (Puduppully et al., 2019; Rebuffel et al., 2020) and board game commentary (Jhamtani et al., 2018), we adopt two automatic metrics that evaluate the quality of system outputs in terms of correctness and fluency. In addition, we perform human judgement on the quality of system outputs

| Strategic depth | Score |
|---|---|
| On the basis of the criteria below, the strategic considerations are inspiring, which can help learn from the skillful players | 5 |
| On the basis of the criteria below, the strategic considerations are sufficient and relevant | 4 |
| On the basis of explaining the facts, the commentary also reflects several strategic considerations, such as the player's intention and the team's arrangement | 3 |
| The commentary only reflects the core event of the game moment described by the structured data | 2 |
| The commentary is nonsense or only reflects few facts of the game moment described by the structured data | 1 |

Table 3: Scoring criteria of the strategic depth metric.

in terms of strategic depth, which measures the informativeness as commentaries. We also use human judgement to evaluate the appropriateness of using existing metrics for correctness and fluency.

#### 4.2.1 Correctness

The correctness evaluation measures how accurate the output is for describing the input. We use the word-level accuracy, which measures the percentage of words in the system outputs that appear in the vocabulary of the reference commentary.

We also calculate BLEU (Papineni et al., 2002) scores to evaluate the correctness. Here we use SacreBLEU (Post, 2018) to compute BLEU scores.

Besides, it also matters in the correctness evaluation whether the system output describes the events that occur in the same order as the input. Therefore, we also analyze how well the system orders the structured data records by measuring the normalized Damerau-Levenshtein distance[5] (Brill and Moore, 2000) between $c_i^{1:T_i}$ and $\hat{c}_i^{1:\hat{T}_i}$.

#### 4.2.2 Fluency

The fluency evaluation measures how fluent as natural language text the system output is. Similar to the other data-to-text generation studies, we calculate perplexity scores for the fluency evaluation.

#### 4.2.3 Strategic depth

The strategic depth evaluation measures whether the system output contains useful information on the LoL players' actions and events in the game. Although the above automatic metrics can arguably measure the general qualities of the system output,

---

[5] https://github.com/life4/textdistance

we also want the output to contain strategically relevant commentaries, such as reflecting the players' intentions and the team's arrangement regarding the combat. Since strategic depth is difficult to determine by automatic metrics, we ask human subjects to collect the scores, as shown in Table 3.

For instance, a commentary sentence *"Tiana's been behind enemy lines for so long as g2 come to try to kill the Gallio"* is considered as containing good strategic depth, because it includes the team players' combat arrangements. On the contrary, *"lwx gets tagged the death sentence"* contains relatively bad strategic depth, because it only repeats the scene that just happened in the game play.

## 5   Esports Data-to-text Generation

In this section, we first introduce baseline encoder-decoder models (Sutskever et al., 2014; Vaswani et al., 2017; Chen et al., 2018) for our esports data-to-text generation task. We then present a hierarchical encoder-decoder model that exploits the hierarchical structure of the input sequence of events.

### 5.1   Encoder-Decoder Baseline

We adopt flat encoder-decoder models (Sutskever et al., 2014; Vaswani et al., 2017; Chen et al., 2018) as baselines to generate a game commentary from a given one-minute structured data.

The encoder takes structured data, a sequence of linearized events, as input and converts it into fixed-length vectors. The decoder then transforms the resulting vectors into output sequences. In this study, we explore four encoder-decoder models (Chen et al., 2018) that use either long-short term memory (LSTM) (Hochreiter and Schmidhuber, 1997) or Transformer (Vaswani et al., 2017) as implementations of the encoder and the decoder and discover the best configurations for the different encoder-decoder model structures.

The output of this task should contain specific entities such as player and team names, which do not appear in the training data. For this reason, we incorporate the copy mechanism (Roberti et al., 2019; See et al., 2017) in to the flat decoder to deal with such unknown entities.

### 5.2   Hierarchical Model

The input of this task is composed of events and has an evident 2-level structure, namely, each event is also composed of a sequence of key-value pairs. Since the flat encoders do not distinguish the in-
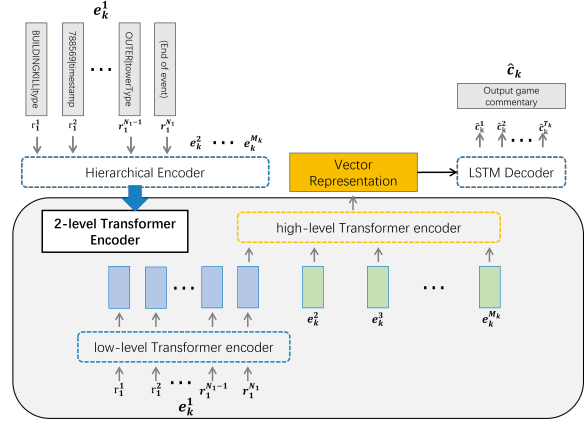


Figure 2: Overall hierarchical model structure of esports data-to-commentary generation.

teraction among events and the interaction among key-value pairs within each event, they will fail to understand the input. For this reason, we apply a hierarchical encoder (Brady and Alvarez, 2011; Rebuffel et al., 2020) to replace the flat encoder.

Figure 2 shows our model which consists of a hierarchical Transformer-based encoder and an LSTM-based decoder. The hierarchical encoder consists of a low-level encoder which feeds key-value pairs in each event $e_i$ to obtain its representation and a high-level encoder which feeds event representations encoded from individual events by the low-level encoder. The hierarchical attention is implemented by learning a first set of attention scores over events $e_i$ and a second set over data records $r_i^j$ belonging to event $e_i$.

## 6   Experiment

This section evaluates the proposed models on esports data-to-commentary generation. We first use the core dataset to evaluate the models with different combinations of encoders and decoders to come up with the best configurations of the encoder and decoder. We then evaluate our hierarchical model based on the best configuration using both core and extended datasets. Next, we provide an analysis of the system outputs through game screenshots to reveal the remaining challenges of our task.

### 6.1   Settings

**Datasets**   To perform training and testing, we split both the core and extended datasets into train, validation, and test sets with a ratio of 8:1:1.

**Models**   We conduct four encoder-decoder models as baseline, with different combination of en-

| Models | Correctness | | | Fluency | Strategic depth |
|---|---|---|---|---|---|
| | Acc (%) | BLEU | Text distance | pplx | Human score |
| Gold | 100 | 100 | 100 | N/A | 3.164 |
| LSTM Enc & Dec | 12.32 | 2.2 | 69.56 | **1.51** | 2.438 |
| LSTM Enc & Transformer Dec | 9.77 | 1.6 | 70.09 | 2.30 | 2.290 |
| Transformer Enc & Dec | 8.17 | 1.4 | **70.22** | 1.98 | 2.312 |
| Transformer Enc & LSTM Dec | 17.00 | 2.8 | 70.19 | 3.04 | 2.548 |
| Hierarchical Transformer Enc & LSTM Dec | **17.06** | **3.3** | 69.47 | 2.10 | **2.772** |

Table 4: Experimental results on the core LoL19 Esports Data-to-text Generation Dataset.

| Models | Correctness | | | | Fluency | |
|---|---|---|---|---|---|---|
| | Acc (%) | BLEU | Text distance | Human score | pplx | Human score |
| Gold | 100 | 100 | 100 | - | N/A | - |
| Transformer Enc & LSTM Dec | 15.91 | 2.9 | 57.84 | 2.754 | 3.40 | 2.930 |
| Hierarchical Transformer Enc & LSTM Dec | **15.94** | **3.1** | **65.10** | **3.014** | **2.34** | **3.130** |

Table 5: Experimental results on the extended LoL19-21 Esports Data-to-text Generation Dataset.

| | |
|---|---|
| Word embedding size | 600 |
| Hidden transformer size | 1024 |
| Decoder dropout | 0.5 |
| Training steps | 20,000 |
| Starting learning rate | 0.001 |
| Learning rate since the 10,000th step | 0.0005 |

Table 6: Experimental settings of the training process.

coder and decoder models. We also conduct a hierarchical model as introduced above to test how a 2-level encoder structure affects the performance of processing the input data of this task.

Specifically, **LSTM Enc & Dec** is a standard LSTM encoder-decoder model (Sutskever et al., 2014). **LSTM Enc & Transformer Dec** is a hybrid encoder-decoder model that uses an LSTM encoder and a Transformer decoder. **Transformer Enc & Dec** is a standard Transformer encoder-decoder model (Vaswani et al., 2017). **Transformer Enc & LSTM Dec** is a hybrid encoder-decoder model that uses a Transformer encoder and an LSTM decoder. **Hierarchical Transformer Enc & LSTM Dec** is a hybrid encoder-decoder system with a 2-level Transformer encoder, which is based on the best configuration of the four baseline models.

**Training**  We use PyTorch[6] and OpenNMT[7] libraries to implement the proposed models. The parameters of the models are listed in Table 6. We re-

duce the vocabulary size of the models to 50,000 to prevent low-frequency words from affecting training. The out-of-vocabulary words are converted into special tokens "⟨UNK⟩"s. We used an Adam optimizer (Kingma and Ba, 2015) for optimization.

**Evaluation procedure**  We use the word-level accuracy, SacreBLEU (Post, 2018), and normalized Damerau-Levenshtein distance to measure the correctness of system outputs, and use perplexity to measure the fluency of system outputs. To confirm the appropriateness of using these metrics, we also evaluate the correctness and fluency by a human subject using a Likert scale (Joshi et al., 2015) of 1-5 for several representative models.

To evaluate the strategic depth, we randomly select 100 data records from the core dataset, along with their gold commentaries and system outputs. We hire 5 human annotators who are capable of understanding the game rules and commentaries of LoL. We collect the annotators' scores of 1-5 scale based on the criteria of Table 3, and then calculate the average number as the strategic depth scores.

### 6.2 Results

Table 4 shows the results on the core dataset. The Transformer encoder & LSTM decoder model achieves the best performance on accuracy, BLEU, and strategic depth score. The LSTM model shows the best performance in terms of perplexity, while the Transformer model shows the best performance on text distance. The results also confirm improved performance brought by the hierarchical structure.

Commentary (reference):
... g2 will use that to try and get a Pavitt turret maybe even the full thing the dredge line doesn't ...

Commentary (system outputs):
... g2 will use that to try and get a [unk] turret maybe even the full thing that *feels like FBX may be trying to get some value* ...

Figure 3: The screenshot shows a "CHAMPION_KILL" event when a player's champion is about to be killed by the enemy team.

We compare the hierarchical model with the best configuration of the four baseline models on the extended dataset. Table 5 shows the results. These results reconfirm that our hierarchical network improves the performance of the generation over all the evaluation metrics. We still perform human judgement for evaluating correctness and fluency. These results are also listed in the table, with 0.940 Spearman correlation with BLEU and 0.871 with perplexity, which proves the appropriateness of using existing metrics for correctness and fluency.

## 6.3 Analysis

To better understand the limit of the current model on the esports data-to-commentary task, we look into the system outputs while referring to the corresponding screenshot for a more direct representation. Figures 3 and 4 show several examples of the commentaries generated by our hierarchical model.

In the first example, the system output correctly describes the core event in which a champion is getting killed. However, it contains an unknown token. It also has hallucination of the game moment, such as "feels like FBX may be trying to get some value," which leads to an erroneous judgement.

In the second example, the system output correctly describes the core event in which a building is getting destroyed. However, it repeats some hallucinated content, such as "it's going to get caught out inside the lane turret." More importantly, the reference output contains a joke about the player's action - "everyone else rolled their characters on



Commentary (reference):
... he's just farming. everyone is fighting everywhere. Jackielove is just farming for himself. everyone else rolled their characters on the PvP server, Jackie love rolled on PvE ...

Commentary (system outputs):
... Jackielove is just farming a tone of these towers *it's going to get caught out inside the lane turret* and two towers on two here's the rift the rift Herald is going to be very important ...

Figure 4: The screenshot shows a "BUILDING_KILL" event when a player is about to destroy an enemy building.

the PvP server, Jackie love rolled on PvE," which means the player's action alone makes him look like playing a different game with the others. The ability to make such jokes can greatly improve the interestingness and intuitiveness of the commentary to help the audience understand it better. Nevertheless, it is still challenging for the data-to-text system to learn to generate such inspiring content.

## 7 Conclusions

This work introduced a novel task of generating game commentaries from structured data of esports games. We built the first data-to-text generation dataset on strategic esports games. We next explored the best combination of the encoder-decoder architectures, and proposed a hierarchical model that can capture 2-level structure of the structured data input. We also discussed evaluation metrics for our task to measure the correctness, fluency, and strategic depth of the system outputs. The hierarchical model was capable of correctly describing the game content in some cases. However, the generated commentaries had limited strategic depth, and they contained hallucination, repetition, and wrong judgement. It was also challenging for the system to generate interesting and intuitive content.

In the future, we plan to explore the domain difference between the commentaries to get a better understanding of the generation.

## Ethical Considerations

In the data collection process, we have strictly followed the policies of RiotGames API and YouTube. The former is the publisher of LoL game records. The later provides subtitles of LoL contest videos, which we used as game commentaries in our work.

## References

Timothy F Brady and George A Alvarez. 2011. Hierarchical encoding in visual working memory: Ensemble statistics bias memory for individual items. *Psychological science*, 22(3):384–392.

Eric Brill and Robert C. Moore. 2000. An improved error model for noisy channel spelling correction. In *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics*, pages 286–293, Hong Kong. Association for Computational Linguistics.

Alejandro Cannizzo and Esmitt Ramírez. 2015. Towards procedural map and character generation for the moba game genre. *Ingeniería y Ciencia*, 11(22):95–119.

Mia Xu Chen, Orhan Firat, Ankur Bapna, Melvin Johnson, Wolfgang Macherey, George Foster, Llion Jones, Mike Schuster, Noam Shazeer, Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Zhifeng Chen, Yonghui Wu, and Macduff Hughes. 2018. The best of both worlds: Combining recent advances in neural machine translation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 76–86, Melbourne, Australia. Association for Computational Linguistics.

Longxu Dou, Guanghui Qin, Jinpeng Wang, Jin-Ge Yao, and Chin-Yew Lin. 2018. Data2text studio: Automated text generation from structured data. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 13–18.

Kirstin Hallmann and Thomas Giel. 2018. esports–competitive sports or recreational activity? *Sport management review*, 21(1):14–20.

Juho Hamari and Max Sjöblom. 2017. What is esports and why do people watch it? *Internet research*.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.

Tatsuya Ishigaki, Goran Topic, Yumi Hamazono, Hiroshi Noji, Ichiro Kobayashi, Yusuke Miyao, and Hiroya Takamura. 2021. Generating racing game commentary from vision, language, and structured data. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 103–113, Aberdeen, Scotland, UK. Association for Computational Linguistics.

Seth E Jenny, R Douglas Manning, Margaret C Keiper, and Tracy W Olrich. 2017. Virtual (ly) athletes: where esports fit within the definition of "sport". *Quest*, 69(1):1–18.

Harsh Jhamtani, Varun Gangal, Eduard Hovy, Graham Neubig, and Taylor Berg-Kirkpatrick. 2018. Learning to generate move-by-move commentary for chess games from large-scale social forum data. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1661–1671, Melbourne, Australia. Association for Computational Linguistics.

Ankur Joshi, Saket Kale, Satish Chandel, and D Kumar Pal. 2015. Likert scale: Explored and explained. *British journal of applied science & technology*, 7(4):396.

Hirotaka Kameko, Shinsuke Mori, and Yoshimasa Tsuruoka. 2015. Learning a game commentary generator with grounded move expressions. In *2015 IEEE Conference on Computational Intelligence and Games (CIG)*, pages 177–184. IEEE.

Yasmin S Khan and Soudamini Pawar. 2015. Video summarization: survey on event detection and summarization in soccer videos. *International Journal of Advanced Computer Science and Applications*, 6(11).

Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Katherine L Lavelle. 2010. A critical discourse analysis of black masculinity in nba game commentary. *The Howard Journal of Communications*, 21(3):294–314.

Sanjay Modgil, Francesca Toni, Floris Bex, Ivan Bratko, Carlos I Chesnevar, Wolfgang Dvořák, Marcelo A Falappa, Xiuyi Fan, Sarah Alice Gaggl, Alejandro J García, et al. 2013. The added value of argumentation. In *Agreement technologies*, pages 357–403. Springer.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Ramakanth Pasunuru and Mohit Bansal. 2018. Game-based video-context dialogue. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 125–136.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on*

*Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Ratish Puduppully, Li Dong, and Mirella Lapata. 2019. Data-to-text generation with content selection and planning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 6908–6915.

Clément Rebuffel, Laure Soulier, Geoffrey Scoutheeten, and Patrick Gallinari. 2020. A hierarchical model for data-to-text generation. In *European Conference on Information Retrieval*, pages 65–80. Springer.

Jason G Reitman, Maria J Anderson-Coto, Minerva Wu, Je Seok Lee, and Constance Steinkuehler. 2020. Esports research: A literature review. *Games and Culture*, 15(1):32–50.

Charles Ringer, Mihalis A Nicolaou, and James Alfred Walker. 2022. Autohighlight: Highlight detection in league of legends esports broadcasts via crowd-sourced data. *Machine Learning with Applications*, page 100338.

RiotGames. 2021. Riot games api for developer. https://developer.riotgames.com/.

Marco Roberti, Giovanni Bonetta, Rossella Cancelliere, and Patrick Gallinari. 2019. Copy mechanism and tailored training for character-based data-to-text generation. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 648–664. Springer.

Abigail See, Peter J Liu, and Christopher D Manning. 2017. Get to the point: Summarization with pointer-generator networks. *arXiv preprint arXiv:1704.04368*.

Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27:3104–3112.

Yasufumi Taniguchi, Yukun Feng, Hiroya Takamura, and Manabu Okumura. 2019. Generating live soccer-match commentary from play data. In *Proceedings of the thirty-third AAAI Conference on Artificial Intelligence*, pages 7096–7103.

Chris van der Lee, Emiel Krahmer, and Sander Wubben. 2017. Pass: A dutch data-to-text system for soccer, targeted towards specific audiences. In *Proceedings of the 10th International Conference on Natural Language Generation*, pages 95–104.

Chris van der Lee, Bart Verduijn, Emiel Krahmer, and Sander Wubben. 2018. Evaluating the text quality, human likeness and tailoring component of pass: A dutch data-to-text system for soccer. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 962–972.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. 2016. Order matters: Sequence to sequence for sets.

Greg L West, Benjamin Rich Zendel, Kyoko Konishi, Jessica Benady-Chorney, Veronique D Bohbot, Isabelle Peretz, and Sylvie Belleville. 2017. Playing super mario 64 increases hippocampal grey matter in older adults. *PLoS One*, 12(12):e0187779.

Sam Wiseman, Stuart M Shieber, and Alexander M Rush. 2017. Challenges in data-to-document generation. In *EMNLP*.