



北京大学  
PEKING UNIVERSITY

# LiveBot: Generating Live Video Comments Based on Visual and Textual Contexts

Shuming Ma<sup>1</sup> Lei Cui<sup>2</sup> Damai Dai<sup>1</sup> Furu Wei<sup>2</sup> Xu Sun<sup>1</sup>

<sup>1</sup>MOE Key Lab of Computational Linguistics, School of EECS, Peking University

<sup>2</sup>Microsoft Research Asia

{shumingma, daidamai, xusun}@pku.edu.cn

{lecu, fuwei}@pku.edu.cn

Microsoft

Research

微软亚洲研究院

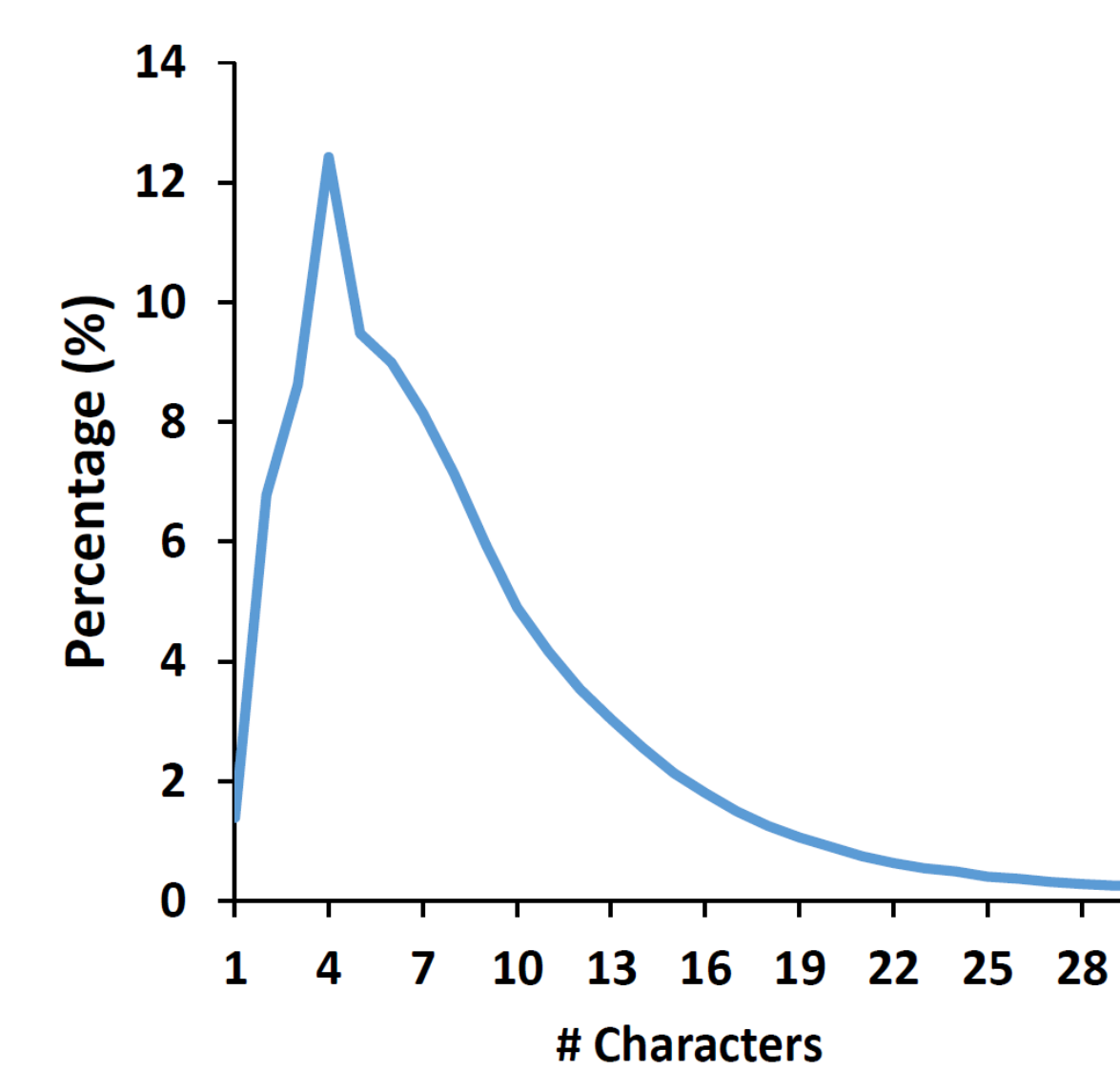
## The Live Video Comment Task and Dataset



时间	弹幕内容 (1000)	发送时间
03:03	二筒, 再见	06-09 14:32
03:13	我来看看猫为啥要虐我.....	06-09 20:12
02:57	你在喵星一定要好好的啊! 二筒	06-10 10:45
00:11	...	06-10 10:50
00:29	喵喵喵	06-10 10:50
00:04	二筒我来看看你啦	06-10 11:00
02:07	寄花姐姐也要在喵星快乐哦	06-10 11:02
02:45	猫传腹: 猫传染性腹膜炎, 一旦...	06-10 12:02
00:33	在泉水等复活的路过	06-10 14:01
02:53	二筒等我, 我来探望一下	06-10 14:43
01:55	二筒在喵星要开开心心的	06-10 14:47
00:12	二筒一路走好	06-10 15:06
01:33	猫弹真好玩	06-10 16:17
00:36	二筒一路走好	06-10 17:12
03:02	走好	06-10 17:14
00:10	想二筒了TAT	06-10 18:33

Live video commenting, which is also called "**video barrage**" ("**弹幕**" in Chinese or "**Danmaku**" in Japanese), is an emerging feature on online video sites that allows real-time comments from viewers to **fly across the screen like bullets** or **roll at the right side of the screen**.

Statistic	Train	Test	Dev	Total
#Video	2,161	100	100	2,361
#Comment	818,905	42,405	34,609	895,929
#Word	4,418,601	248,399	193,246	4,860,246
Avg. Words	5.39	5.85	5.58	5.42
Duration (hrs)	103.81	5.02	5.01	113.84



Interval	Edit Distance	TF-IDF	Human
0-1s	11.74	0.048	4.3
1-3s	11.79	0.033	4.1
3-5s	12.05	0.028	3.9
5-10s	12.42	0.025	3.1
>10s	12.26	0.015	2.2

The dataset is crawled from a popular **video streaming** website called **Bilibili.com**.

## An Example from the Live Commenting Dataset

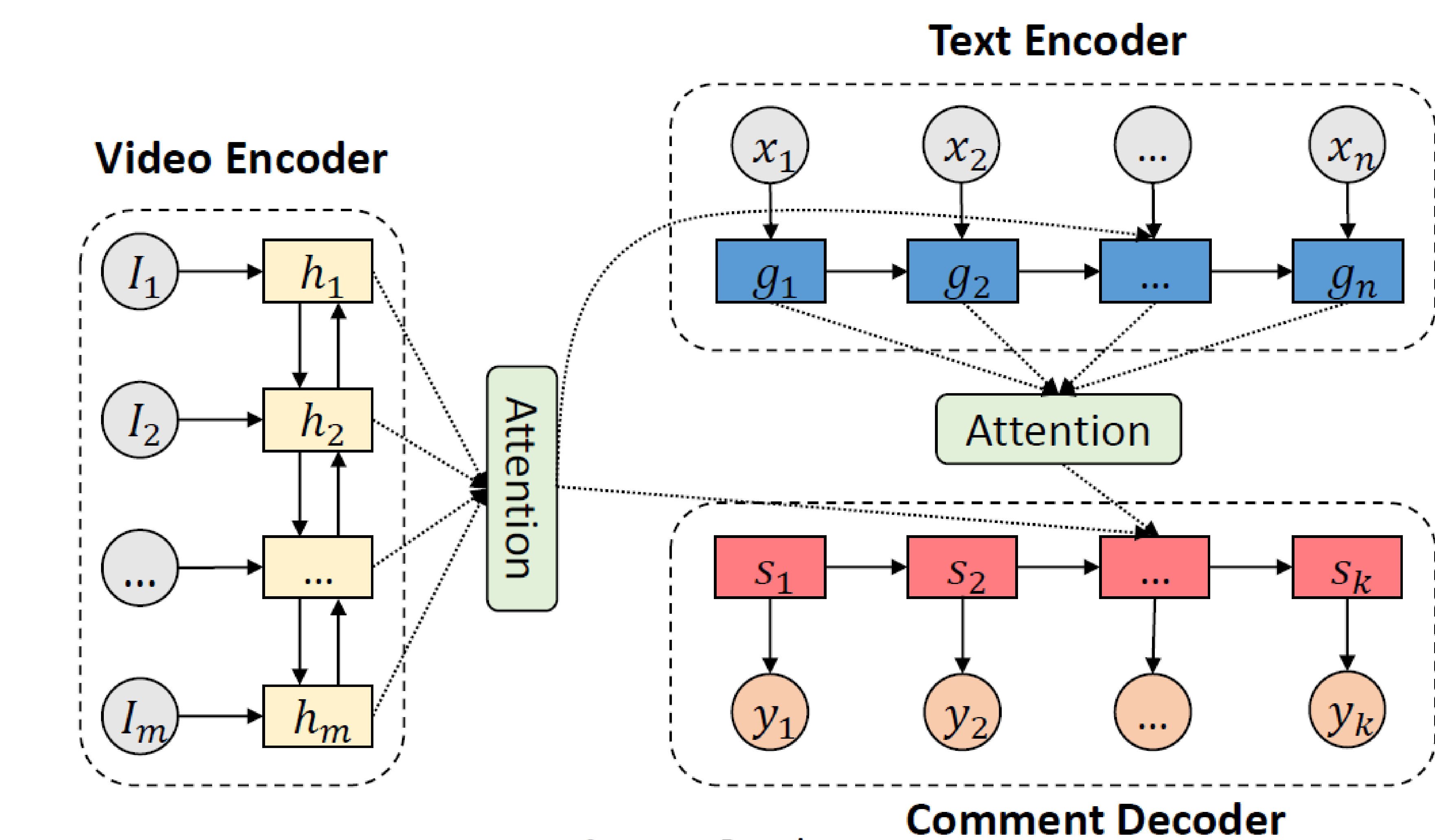


Three selected **frames** from the videos to demonstrate the content

Time Stamp	Comments
0:48	橙猫是短腿吗 (Is the orange cat short leg?)
1:06	根本停不下来 (Simply can't stop)
1:09	哎呀好可爱啊 (Oh so cute)
1:52	天哪这么多, 天堂啊 (OMG, so many kittens, what a paradise!)
1:56	这么多只猫 (So many kittens!)
2:39	我在想猫薄荷对老虎也有用吗 (I am wondering whether the catmint works for the tiger.)
2:41	猫薄荷对老虎也有用 (Catmint also works for the tiger.)
3:41	活得不如猫 (The cat lives even better than me)
3:43	两个猫头挤在一起超可爱 (It's so cute that two heads are together)

Several selected **live comments** paired with the **time stamps** when the comments appear on the screen.

## Live Commenting Models



Two approaches to generate the comments based on the **visual contexts** and the **textual contexts**.

The two approaches are based on two popular architectures for text generation: **recurrent neural network (RNN)** and **transformer**.

Both two models consist of a **video encoder**, a **text encoder**, and a **comment decoder**. We denote two approaches as **Fusional RNN Model** and **Unified Transformer Model**, respectively.

## Experiments

	Model	#I	#C	Recall@1	Recall@5	Recall@10	MR	MRR
Video Only	S2S-I	5	0	4.69	19.93	36.46	21.60	0.1451
	S2S-IC	5	0	5.49	20.71	38.35	20.15	0.1556
	Fusional RNN	5	0	10.05	31.15	48.12	19.53	0.2217
	Unified Transformer	5	0	<b>11.40</b>	<b>32.62</b>	<b>50.47</b>	<b>18.12</b>	<b>0.2311</b>
Comment Only	S2S-C	0	5	9.12	28.05	44.26	19.76	0.2013
	S2S-IC	0	5	10.45	30.91	46.84	18.06	0.2194
	Fusional RNN	0	5	13.15	<b>34.71</b>	<b>52.10</b>	17.51	0.2487
	Unified Transformer	0	5	<b>13.95</b>	34.57	51.57	<b>17.01</b>	<b>0.2513</b>
Both	S2S-IC	5	5	12.89	33.78	50.29	17.05	0.2454
	Fusional RNN	5	5	17.25	37.96	<b>56.10</b>	16.14	0.2710
	Unified Transformer	5	5	<b>18.01</b>	<b>38.12</b>	55.78	<b>16.01</b>	<b>0.2753</b>

Model	Fluency	Relevance	Correctness
S2S-IC	4.07	2.23	2.91
Fusion	<b>4.45</b>	2.95	3.34
Transformer	4.31	<b>3.07</b>	<b>3.45</b>
Human	4.82	3.31	4.11

Our models outperform the existing video-to-text models in terms of both **automatic evaluation** and **human evaluation**.

## Conclusions and Future Work

- ▣ **A new task:** live video commenting
- ▣ **A new dataset:** the live comment dataset
- ▣ **Two new baselines:** Fusional RNN and Unified Transformer
- ▣ **An evaluation protocol**

### Future

- ▣ **A better video/image representation**
- ▣ Generating **low-frequent/surprising** comments
- ▣ Integrating **audio** information
- ▣ **Other tasks:** video MT, video summarization, video QA



Scan the QR code for the codes and dataset