

# Supplementary material for “Character-aware audio-visual subtitling in context”

Jaesung Huh<sup>✉</sup> and Andrew Zisserman<sup>✉</sup>

Visual Geometry Group, Department of Engineering Science, University of Oxford  
{jaesung,az}@robots.ox.ac.uk

## 1 Character recognition accuracy on *Bazinga!*-gold-TV

Tab. 1 shows the character recognition performance on *Bazinga!*-gold-TV, comparing with that from LLR [1]. The reported performance of LLR corresponds to its highest character recognition accuracy (**Acc.**). This peak accuracy is achieved by adjusting the threshold during the nearest centroid classification process. We showcase that our method achieves higher accuracy than LLR in all TV shows. It also achieves higher recall for both ‘all characters’ and ‘main characters’ compared to LLR. LLR exhibits higher precision for ‘all characters’ and ‘main characters’ in a few TV shows. This difference in performance is due to how each method handles short segments. LLR does not assign speakers to some of short segments, instead marking the character as **unknown**. In contrast, our LLM tends to predict speakers for these **unknown** segments after local embedding classification. It often assigns one of the characters appearing in the dialogue, even when we explicitly instruct the LLM that it doesn’t have to make an assignment if it’s unsure. Thus, the precision decreases but the recall increases.

## 2 Effects of utilising spatial regions and speech enhancement on audio exemplar yield and performance

Tab. 2 demonstrates how cropping lip areas and speech enhancement improve exemplar yield. Unlike LLR, which uses whole frames without speech enhancement, our method identifies characters by focusing on the spatial region of speakers and reducing background noise to decrease false positive peaks, increasing exemplar yield by 8.5% and 0.9% respectively. This showcases the effectiveness of our cropping-based approach in terms of exemplar yield.

## 3 List of main characters per series

Tab. 3 lists the main characters per series, both in LLR-TV and *Bazinga!*-gold-TV. Note that we report the names of main characters in *Bazinga!*-gold-TV as they are in the dataset annotation.

**Table 1:** Character recognition performance on *Bazinga!*-gold-TV test set. **Acc.** is the character recognition accuracy for those which overlap with one of the groundtruth timestamps. **Prec.** and **Rec.** indicate the precision and recall of overall audio segments respectively, while **Prec.(M)** and **Rec.(M)** are those of main characters in TV shows.

	LLR [1]					Ours				
	Acc.	Prec.	Rec.	Prec.(M)	Rec.(M)	Acc.	Prec.	Rec.	Prec.(M)	Rec.(M)
<b>B.G.</b>	62.7	67.7	64.1	<b>72.7</b>	71.1	<b>68.7</b>	<b>69.1</b>	<b>70.1</b>	69.1	<b>75.9</b>
<b>B.B.</b>	60.5	69.4	61.4	<b>83.6</b>	64.9	<b>67.5</b>	<b>72.6</b>	<b>68.6</b>	76.3	<b>69.4</b>
<b>Buffy</b>	57.6	62.4	58.1	71.7	57.5	<b>61.6</b>	<b>62.6</b>	<b>62.1</b>	<b>75.0</b>	<b>65.0</b>
<b>Friends</b>	70.9	<b>82.1</b>	72.9	<b>85.8</b>	79.7	<b>74.2</b>	75.6	<b>76.4</b>	76.0	<b>83.2</b>
<b>GoT</b>	52.6	55.4	55.7	<b>57.3</b>	59.7	<b>61.6</b>	<b>62.9</b>	<b>65.3</b>	54.1	<b>63.7</b>
<b>Lost</b>	51.4	61.2	53.4	62.6	56.6	<b>60.4</b>	<b>63.9</b>	<b>62.6</b>	<b>65.4</b>	<b>67.3</b>
<b>TBBT</b>	80.2	<b>88.8</b>	80.4	<b>90.9</b>	85.4	<b>81.5</b>	82.7	<b>81.8</b>	83.2	<b>86.2</b>
<b>Office</b>	68.7	<b>77.5</b>	68.9	<b>84.8</b>	74.1	<b>70.8</b>	74.4	<b>71.2</b>	75.4	<b>75.0</b>
<b>W.D.</b>	49.7	<b>60.3</b>	51.9	<b>57.0</b>	52.9	<b>54.7</b>	57.6	<b>57.1</b>	53.5	<b>64.0</b>

**Table 2:** Effect of using local visual predictions (vis) and speech enhancement (SE) on exemplar yield on the LLR-TV.

	LLR		LLR + vis		LLR + vis + SE (Ours)	
	# of exemplars	% of total	# of exemplars	% of total	# of exemplars	% of total
1-1. VAD + ASR	5554	100	5554	100	5554	100
1-2. Audio-visual speaker recognition	2192	39.5	3332	60.0	3409	<b>61.4</b>
1-3. Audio filtering	1213	21.8	1681	30.3	1734	<b>31.2</b>

## 4 LLM prompt

Tab. 4 shows the LLM prompt we’ve used to determine the **unknown** character in the dialogue. We adopt the strategy introduced in [2]. Given the query dialogue with the **unknown** character we want to classify, we first ask the LLM to summarize the dialogue. Then, we ask the model who would be the **unknown** in the dialogue based on the generated summary. A list of characters with their corresponding indices is provided in the prompt. We compute the softmaxed logits of the tokens corresponding to each index and choose the largest one to assign the speaker. We use this same prompt for experiments in LLR-TV and *Bazinga!*-gold-TV.

There can be multiple **unknowns** in the query dialogue. We mark the sentence for which we want to identify the speaker by placing asterisks (\*\*) before and after it. This explicitly instructs the LLM to predict the speaker only for that specific sentence.

## References

1. Korbar, B., Huh, J., Zisserman, A.: Look, listen and recognise: character-aware audio-visual subtitling. In: International Conference on Acoustics, Speech, and Sig-

**Table 3:** List of main characters per show.

Show	Main characters
<b>LLR-TV</b>	
<b>Frasier</b>	Frasier, Martin, Niles, Roz, Daphne
<b>Seinfeld</b>	Jerry, Elaine, George, Kramer
<b>Scrubs</b>	J.D., Dr.Cox, Dr.Kelso, Carla, Turk, Elliot
<b>Bazingal-gold-TV</b>	
<b>Battlestar Galactica (B.G.)</b>	Admiral William Adama, President Laura Roslin, Captain Kara Thrace, Cpt. Lee Adama, Number six, Dr. Gaius Baltar, Lt. Sharon Valerii
<b>Breaking Bad (B.B.)</b>	Walter White, Skyler White, Jesse Pinkman, Hank Schrader, Marie Schrader
<b>Buffy the Vampire Slayer (Buffy)</b>	Buffy Summers, Willow Rosenberg, Xander Harris, Angel, Rubert Giles
<b>Friends</b>	Rachel Green, Monica Geller, Phoebe Buffay, Joey Tribbiani, Chandler Bing, Dr. Ross Geller
<b>Game of Thrones (GoT)</b>	Danerys Targaryen, Jon Snow, Jorah Mormont, Tyrion Lannister, Catelyn Stark, Sansa Stark, Arya Stark, Cersei Lannister, Eddard Stark, Robert Baratheon
<b>Lost</b>	Dr. Jack Sheperd, Kate Austen, Sayid Jarrah, Hugo Reyes, Sunhwa Kwon, Charlie Pace, Clarie Littleton, Michael Dawson, John Locke, Shannon Rutherford, James Ford
<b>The Big Bang Theory (TBBT)</b>	Sheldon Cooper, Penny, Howard Wolowitz, Raj Koothrappali, Leonard Hofstadter
<b>The Office (Office)</b>	Jim Halpert, Michael Scott, Ryan Howard, Pam Beesly, Dwight Schrute, Stanley Hudson, Phyllis Vance, Angela Martin
<b>Lost</b>	Rick Grimes, Lori Grimes, Carl Grimes, Carol Peletier, Shane Walsh, Andrea Harrison, Dale Horvath, Glenn Rhee

- nal Processing (2024)
2. Park, R.Y., Windsor, R., Jamaludin, A., Zisserman, A.: Automated spinal mri labelling from reports using a large language model. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 101–111. Springer (2024)

**Table 4:** LLM prompt to determine unknown.

---

**PROMPT FOR DETERMINING [UNKNOWN]**

**system** : You are a AI assistant to analyze the transcript of TV shows. Your job is to figure out who are [UNKNOWN]s in a dialogue in TV shows. Tell the truth and answer as precisely as possible.

*(provide the dialogue here)*

**user** : Write a summary for the above conversation.

**assistant** : *(model generates the summary)*

**user** : Based on the summary, your job is to identify the name of the speaker of '[UNKNOWN]' when the line starts and ends with '\*\*'. You must use the context and the flow of the dialogue, using the speakers' names and what they speak. The list of speakers with their corresponding tokens are provided below. Choose [UNKNOWN] if his or her name is not in the dialogue, or when you are not sure.

*(provide the list of speakers here)*

Only output one number after ANSWER:

**EXAMPLE OF THE DIALOGUE**

Dr.Cox : You've got an opposable thumb.

Dr.Cox : You can use it.

Dr.Cox : God, I hate Halloween.

Carla : Somebody needs to adjust their attitude if they want some candy.

Dr.Cox : You mean, the popcorn balls and the deformed lollipops.

Dr.Cox : I mean, honestly, where do you get this crap anyway?

\*\*[UNKNOWN] : I made it.\*\*

NurseRoberts : If you want name brand candy, my fish is packed with peanuts.

Dr.Cox : Of course it is.

Carla : Oh, what's the matter?

Carla : Did Raggedy Ann scare you?

Dr.Cox : What are you, a rat?

**EXAMPLE OF THE LIST OF SPEAKERS**

1: Dr.Cox, 2: Carla, 3: NurseRoberts, 4: [UNKNOWN]

---