Zero-Shot Character Identification and Speaker Prediction in Comics via Iterative Multimodal Fusion

Yingxuan Li¹, Ryota Hinami², Kiyoharu Aizawa¹, Yusuke Matsui¹ ¹The University of Tokyo, ²Mantra Inc.



ACM Multimedia 2024 Melbourne, Australia

Topic

• Zero-Shot character identification and speaker prediction in comics

• Identify characters and predict speakers of unseen comics only from images



Proposal

Research focus

- Predict character names for both text and character regions
- Tackle *zero-shot* tasks without requiring any annotations
- Enhance real-world applicability

Proposed method

- Multimodal fusion approach: Merge text-based large language model (LLM) predictions with image-based classifiers
- Iterative process: Alternately refine speaker prediction using character labels, and character identification using text labels



Overall framework



Data preprocessing

- Object detection: Get character regions and text regions
- Relationship prediction: Get initial relationship scores
- OCR: Get texts
- Character name extraction: Get target labels

Overall framework



- Main pipeline: Three modules
 - Speaker prediction (F):
 - Character identification (G):
 - Label propagation $(H_{t \to c}, H_{c \to t})$: Text labels \leftrightarrow Character labels





- Approach
- Experiments
- Discussion and Conclusion

Iterative character identification



• Step 1: Pseudo label generation

- Select the character-text pair with the highest relationship score
- Take the text region's label as the pseudo label

• Step 2: Character identification

- Fine-tune a classifier with pseudo labels
- Apply the classifier to character regions

Iterative speaker prediction



• Treat pseudo labels as candidate speakers

Iterative process



Iterative performance enhancement

- Predictions from the previous step
- Generated pseudo labels
- Updated prediction results

- \rightarrow Close to true labels
 - \rightarrow Reliable training data
 - \rightarrow High accuracy



- Approach
- Experiments
- Discussion and Conclusion

Quantitative results

- Main results: Simplified the tasks to classifying the character labels for the character and text regions
- **Baselines:** Clustering + mapping clusters to ground truth*
- **Data division:** Divided the test set by the difficulty of relationship prediction

				Speaker pred.			Character id.			
	iter	text	img	Easy	Hard	Total	Easy	Hard	Total	
Baseline K-means+Distance K-means+SGG	-		√ √	34.5* 36.7*	31.8* 34.8*	33.1* 35.7*	37.0* 37.0*	36.7* 36.7*	36.8* 36.8*	Our multimodal method shows a significant improvement over
Proposed							a server	-		unimodal methods
LLM only	0	\checkmark		41.8	45.1	43.6	na n	-	-	
Multimodal	1	\checkmark	\checkmark	51.0	51.2	51.1	45.8	39.6	42.4	
	2	\checkmark	\checkmark	52.4	51.3	51.8	48.5	40.3	44.0	
	3	\checkmark	\checkmark	53.5	49.8	51.6	48.9	37.7	42.8	

Quantitative results

- Main results: Simplified the tasks to classifying the character labels for the character and text regions
- **Baselines:** Clustering + mapping clusters to ground truth*
- Data division: Divided the test set by the difficulty of relationship prediction

				Speaker pred.			Character id.			
	iter	text	img	Easy	Hard	Total	Easy	Hard	Total	
Baseline										Our iterative process is effective
K-means+Distance	-		\checkmark	34.5*	31.8*	33.1*	37.0*	36.7*	36.8*	on both subsets particularly on
K-means+SGG	-		\checkmark	36.7*	34.8*	35.7*	37.0*	36.7*	36.8*	Face
Proposed								1		LUSY
LLM only	0	\checkmark		41.8	45.1	43.6	-	_	-	
Multimodal	1	\checkmark	\checkmark	51.0	51.2	51.1	45.8	39.6	42.4	
	2	\checkmark	\checkmark	52.4	51.3	51.8	48.5	40.3	44.0	
	3	\checkmark	\checkmark	53.5	49.8	51.6	48.9	37.7	42.8	

Qualitative results

• Unimodal vs. Multimodal



Mislabeled predictions (ground truth)



Multimodal



Multimodal

Qualitative results

• One-step vs. Iterative (Accuracy of speaker pred. & character id.)



Iteration 1 (0/3 & 0/3)



Iteration 2 (1/3 & 1/3)



Iteration 3 (3/3 & 2/3)

Zero-shot results

• Results under entirely zero-shot settings



Introduction

- Approach
- Experiments
- Discussion and Conclusion

Discussion

Limitations of proposed method

- The overall accuracy is limited
- More iterations do not necessarily lead to higher accuracy



Iteration 1



Iteration 2

Conclusion

New tasks

- First to integrate the tasks of character identification and speaker prediction in comics
- First to tackle zero-shot tasks with direct applications in real-world scenarios

Iterative multimodal fusion

- Revealing the significant potential of LLMs for comics analysis
- First approach to use multimodal information for character identification and speaker prediction

• Future work

• Refine the model and enhance accuracy

Explore more

- Visit our poster presentation!
 - Poster session: Poster Session 3
 - Posterboard: P175
- Explore our project page for more details



