

# MICap: A Unified Model for Identity-aware Movie Descriptions

## SUPPLEMENTARY MATERIAL

Haran Raajesh\*<sup>1</sup>   Naveen Reddy Desanur\*<sup>1</sup>   Zeeshan Khan<sup>2</sup>   Makarand Tapaswi<sup>1</sup>  
<sup>1</sup>CVIT, IIIT Hyderabad, India

<sup>2</sup>Inria Paris and Département d’informatique de l’ENS, CNRS, PSL Research University

<https://katha-ai.github.io/projects/micap/>

\* denotes equal contribution

We present additional insights and results in the supplementary material. In Appendix A, we highlight how our auto-regressive Transformer decoder attends to various memory features. For the id-aware captioning task, we show the relative importance of the 3 visual features, while for the Fill-in-the-blanks (FITB) task, we highlight how our model attends to correct face clusters. Next, in Appendix B, we show qualitative results for both tasks, FITB and id-aware captioning. We also illustrate how our new identity-aware metric, iSPICE, is calculated on some examples. Finally, we end with discussion of some limitations in Appendix C.

### A. Analyzing Model Attention

In this section, we visualize and discuss the attention scores from MICap’s auto-regressive Transformer decoder. In particular, we focus on the cross-attention scores of the last layer as they reveal interesting insights about the features that the captioning model uses. Throughout this section, we analyze MICap trained jointly on id-aware captioning and FITB. All attention scores are obtained in inference mode.

#### A.1. Attention Patterns in Id-aware Captioning

In id-aware full captioning, for a particular videoset  $\mathcal{N} = \{V_i\}_{i=1}^N$ , we first encode the videos to obtain memory tokens  $M$  and pass them through a Transformer decoder auto-regressively to generate one token (word) at a time. If we consider that the number of tokens in the predicted captionset is  $L$ , we can compute a matrix of cross-attention scores  $\alpha = L \times |M|$ , where  $|M|$  is the number of tokens in the decoder memory. Note, while we use multi-head attention, scores over the heads are averaged obtain  $\alpha$ .

We split the  $L$  tokens into 2 groups: (i) one group consists of person id label predictions or *person tokens* (PT); and (ii) the other group consists of all other tokens referred to as *caption tokens* (CT). For visualization, we sum over the attention scores for each of the token types (id labels and

text) and convert our attention map to a matrix of  $2 \times |M|$ .

Next, we also group the memory tokens into 3 types of visual features used in our work: action (I3D), face (Arc-face), and semantic features (CLIP). Thus, we obtain a  $2 \times 3$  matrix of cross-attention scores for each sample.

**Results.** We compute attention scores over all samples of the validation set and plot them as a probability density function in Fig. 1. PT (red) and CT (green) represent the person and caption tokens respectively. We observe that: (i) The model relies on CLIP features to predict captions (depicted by the overall high attention scores from 0.5-0.7). (ii) When predicting person tokens (PT) of the identity-aware captions, the model tends to look at face features (0.1-0.6) more than when predicting caption tokens (0-0.4). (iii) Finally, while action features are useful for captioning, they are less useful for predicting person-id labels. This is expected as action recognition is an identity-agnostic task.

#### A.2. Attention Patterns in FITB

For the FITB task, we analyze how the person id predictions attend to *face features* from the decoder memory. For a videoset  $\mathcal{N} = \{V_i\}_{i=1}^N$  and its corresponding captionset with blanks  $\hat{\mathcal{C}}$  we obtain a cross-attention map of  $\alpha = |\mathcal{B}| \times F$ , where  $|\mathcal{B}|$  is the number of blanks in the captionset, and  $F$  is the number of face detections across the videoset. Each row of this matrix is normalized to sum to 1.

The attention scores and captionsets with blanks are presented in Fig. 2. In the next paragraphs, we will analyze the 3 types (columns) of the presented scores.

**Cross-attention scores for face detections.** In the left column of Fig. 2, we visualize the attention scores directly for each face detection. In the plot, x-axis spans time across different videos. Our model tends to show a diagonal pattern indicating that person id label predictions tend to look at faces in the same video (facilitated through the  $\mathbf{E}^{\text{vid}}$  embeddings). However, as seen in captionset 5, left, row 1, the

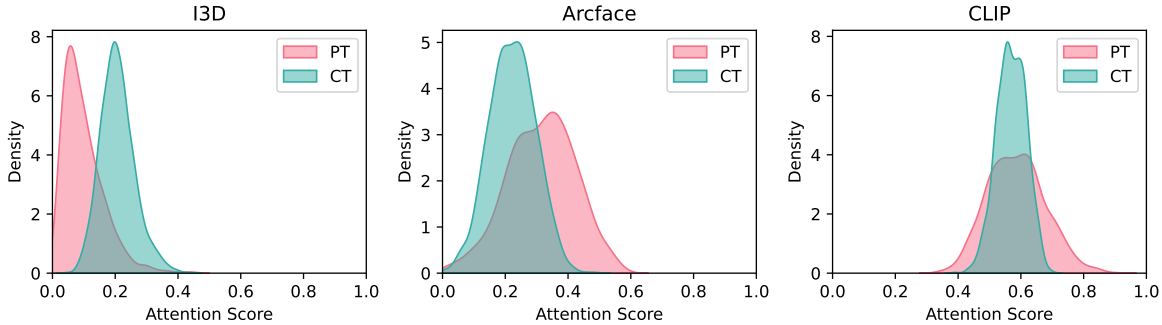


Figure 1. Cross-attention scores density plots for the id-aware captioning task. We group decoder output tokens into two types: person id label tokens (PT), and caption tokens that represent other words (CT). Attention scores are grouped across the three input visual features capturing actions (I3D, left), faces (Arcface, middle), and semantic content (CLIP, right). Please refer to Appendix A.1 for a discussion.

model may also attend to other face detections of the same person across videos. This highlights that being able to attend to faces across videos is useful (compared to [1] that only looks at faces within the same video).

**Cross-attention scores for face clusters grouped by video index.** Shown in the middle column of Fig. 2, we group the  $F$  face detections into clusters, but split them based on video index in the videoset. For example, in captionset 1, we see that faces in cluster 1 appears across videos 1, 2, 4 (C1/V1, C1/V2, C1/V4). This allows us to explain some of the predictions made by our model.

Please note that the face cluster index and person id labels need not match numerically. That is, cluster 2 could be assigned the label P1 and cluster 1 the label P2. These changes are acceptable as we only consider person id labels in a local videoset.

In captionset 3, we see that cluster 2 corresponds to the prediction P1 (first two rows) and cluster 4 (C4/V3) corresponds to person id label P2 (bottom two rows). In the last row of captionset 3, we see that our model predicts P2 for the video id 4 correctly, while looking at cluster 4 in video 3 (C4/V3). Previous work [1] is unable to use such cross-video information.

**Cross-attention scores for clusters.** In the right columns of Fig. 2, we show attention scores directly grouped by cluster ids. Here, the original attention map of  $|\mathcal{B}| \times F$  is grouped to  $|\mathcal{B}| \times |\mathcal{G}|$ , where  $|\mathcal{G}|$  is the number of face clusters obtained after performing DBSCAN on the  $F$  face detections.

Captionset 2 is an example with multiple blanks and 4 characters. We observe that some confusion in attention scores leads to errors in the predicted person id labels. In captionset 4, we also see 6 blanks, now with 3 characters. In the last row, while the model wrongly predicts P1, the model does look at cluster 3 (corresponding to P3) correctly. Captionset 1 and 2 are examples of perfect attention scores

and clusters. P1 and C1, and P2 and C2 go together strongly in these examples.

**Impact of number of clusters on FITB.** Fig. 3 shows the results on FITB class-accuracy for varying the DBSCAN epsilon parameter. These results indicate the importance of clustering across videos and choosing an appropriate number of clusters. Qualitatively, we adopt 0.75 and it is unlikely to merge characters incorrectly.

## B. Qualitative Results

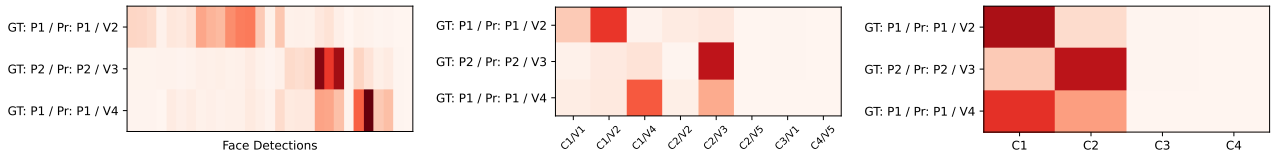
**iSPICE validation examples.** To validate our new metric, we propose an experiment that measures similarity between captions when identity names are added, removed, or replaced (Sec. 4 of the main paper). While the quantitative results favor iSPICE, as seen in Tab. 1 of the main paper, we illustrate with examples the process of metric computation in Fig. 4. We observe that the small difference in identity names is captured correctly by iSPICE, due to the focus on tuples containing identities, while other metrics do not show this sensitivity.

**FITB examples.** While Fig. 2 clearly shows the importance of cross-attention scores of detected faces and computed clusters, the challenging visual scenarios are not evident. We pick two examples (captionset 3 and 4) from Fig. 2 and pair them together with one frame from each video of the videosets. Fig. 5 shows the challenging nature of the videos where characters are often not looking at the camera (example 1 video 1, 3), the scene is dark, or the face may not even be visible (example 1 video 4 or example 2 video 3). MI-Cap leverages the ability to look at faces and clusters across videos to improve results on the FITB task.

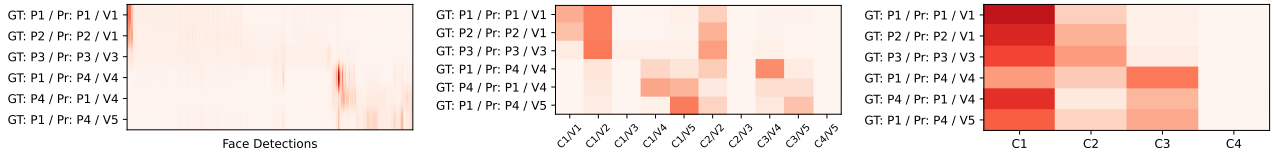
**Id-aware captioning examples.** Fig. 6 shows 2 examples where our model does relatively well, while Fig. 7 shows 2 difficult examples where our model makes mistakes.

In the left column of Fig. 6 we see that the model rightly

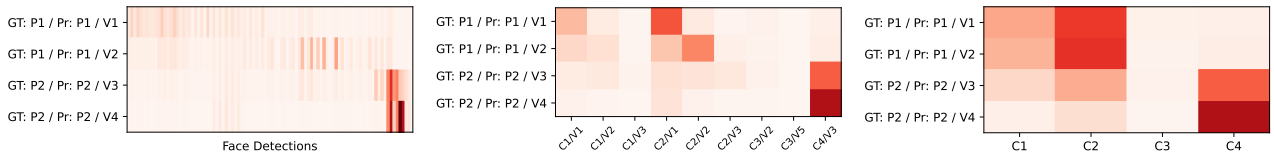
**Captionset 1:** Someone watches the aliens draw closer. \_\_\_\_\_ sits back in the doorway clutching a radio. \_\_\_\_\_ watches from his position several yards away. \_\_\_\_\_ squeezes the detonator the bus blows apart.



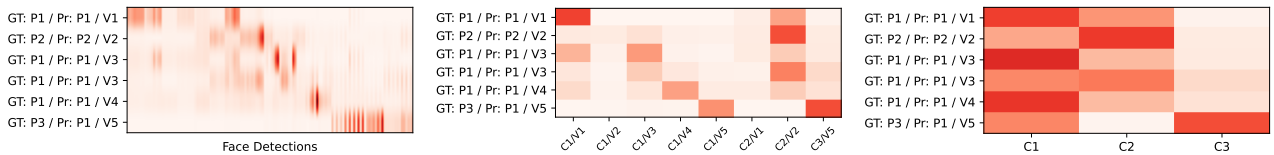
**Captionset 2:** \_\_\_\_\_ and \_\_\_\_\_ killed their first witch. They advance cautiously. Suddenly \_\_\_\_\_ is thrown to the ground with a jolt. \_\_\_\_\_ whips around a weapon poised to find \_\_\_\_\_ holding her wand to neck. \_\_\_\_\_ begins to put the gun on the ground.



**Captionset 3:** \_\_\_\_\_ pulls her phone from her bag and answers. \_\_\_\_\_ frowns uncertainly. \_\_\_\_\_ leans on a wall and slips. \_\_\_\_\_ lowers his phone and folds it shut. The next morning two women stroll across the street in front of apartment building.



**Captionset 4:** \_\_\_\_\_ scrutinizes his earnest face. His eyes gleaming in the dim light. \_\_\_\_\_ abruptly gets to his feet and heads for the door now. \_\_\_\_\_ talks on his cell as \_\_\_\_\_ steps into the daylight silhouetted against the sunny day. \_\_\_\_\_ faces the door frame and leans his head against it now. In a hotel suite a woman applies makeup to \_\_\_\_\_.



**Captionset 5:** \_\_\_\_\_ turns and spots the brown chevy 4x4 parked on a short driveway. \_\_\_\_\_ approaches the vehicle cautiously across a lawn leaning over to get a view of its occupant. The passenger side window is lowered. \_\_\_\_\_ puts both hands on the sill and leans in with an inquisitive frown. \_\_\_\_\_, the asian man who in town sits with one hand clamped to the steering wheel rocking nervously and staring numbly ahead.

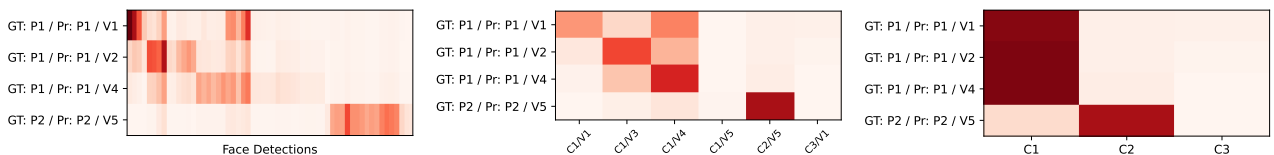


Figure 2. We show 5 examples of our model’s attention scores on the FITB task. For each example (row), we show the captionset (with blanks) and the attention scores grouped in various ways. The **left** column shows the attention score for each blank across all face detections in the video. The **middle** column shows attention scores for face detections grouped by clusters in each video. C1/V1 indicates faces appearing in cluster 1 and video 1, while C1/V2 indicates faces of the same cluster 1 appearing in video 2. The **right** column shows attention scores of each blank for face clusters (across videos). For each row in the attention scores, we indicate the ground-truth (GT) and predicted (Pr) person id label and the video index (V1 .. V5) in which this blank appeared. See Appendix A.2 for a discussion.

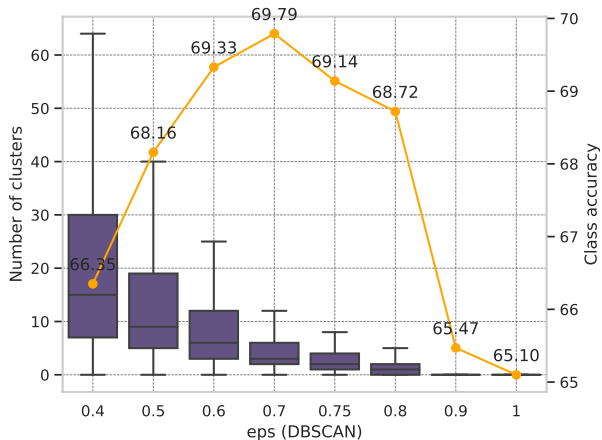


Figure 3. Class-accuracy for the FITB task by varying the DBSCAN eps distance threshold. We also show a box-plot for the number of clusters created at each threshold across samples of the validation set.

identifies P1 as the male character and P2 as the female. The last caption is quite interesting – while the GT points to P1 giving P2 a bowl, our model predicts that P2 gives a sad smile, which is not wrong. This also illustrates some of the challenges of evaluating captioning. In the right column of Fig. 6, the predicted caption uses P2 to refer to the man, and is consistent across videos 3, 4, and 5 in the videoset.

In the complex visual example of Fig. 7 (left), our model assigns P1 to all blanks. Similarly, in the multi-character example of Fig. 7 (right), we observe some confusion between characters. Nevertheless, P2, identified as the man on the left in video 3, is correctly identified for the first 3 videos. The model is also able to predict that they are on a plane (caption for video 2). Nevertheless, these examples illustrate the challenges of id-aware captioning. As future work, they also highlight the need to evaluate visual grounding of the identities beyond captioning performance.

### C. Limitation and Future Work

One limitation of our work, inherited from the task definition in LSMDC, is restricting videosets to local groups of 5 videos. In the future, we would like to extend this to larger videosets, perhaps spanning the entire movie. However, the approach will need to be modified to operate on full movies as: (i) providing features of all movie frames as decoder memory creates a huge number of embeddings; (ii) face clustering across the entire movie could be error-prone; and (iii) auto-regressively generating one caption at a time for hundreds of clips seems challenging, as the model needs to be cognizant of all previously generated captions. We believe that a hierarchical model that builds from shots to scenes to the full movie may be more appropriate here.

Second, the tasks for FITB and full captioning do not

learn at the same pace, and choosing a single best checkpoint for both may be difficult. We posit that the user may choose two checkpoints, one for each task. Furthermore, we observe that by weighing the FITB and full captioning losses appropriately, additional performance improvements can be achieved for one task at the cost of the other task.

We have also not considered using external knowledge or pre-trained large language models (LLMs) or vision-language models (VLMs) built for captioning. We believe that it is interesting to learn what can be achieved by training on LSMDC alone. As seen in multiple examples throughout Appendix B, MICap does perform quite well given the challenging scenarios.

### References

[1] Jae Sung Park, Trevor Darrell, and Anna Rohrbach. Identity-aware multi-sentence video description. In *European Conference on Computer Vision (ECCV)*, 2020. 2

### Add Example

Candidate : A path leads from the side of the circle splitting into two prongs. A third crop circle has two straight lines at either side and a circle of maize remaining in the center with another path leading off from its side. It splits into three larger prongs the central one of which points towards a smaller circle. P1 is on the phone as P2 looks out of his window at the yard. P2 bows his head.

Tuples : [[path, side, lead from], ..., [prong], [p2, window, look out of], [p1, phone, on], [window, yard, at], [phone, window, have], [window], [yard], [p1], [phone], [p2], [p2, head, bow], [p2, head, have], [head], [p2]]

Reference : A path leads from the side of the circle splitting into two prongs. A third crop circle has two straight lines at either side and a circle of maize remaining in the center with another path leading off from its side. It splits into three larger prongs the central one of which points towards a smaller circle. P1 is on the phone as P1 looks out of his window at the yard. P1 bows his head.

Tuples : [[path, side, lead from], ..., [prong], [p1, window, look out of], [p1, phone, on], [window, yard, at], [phone, window, have], [window], [yard], [p1], [phone], [p1], [p1, head, bow], [p1, head, have], [head], [p1]]

CIDER: 92.4 | METEOR: 64.0 | BLEU: 92.0 | SPICE : 91.76 | iSPICE : 16.66

(Term1) = ([[p1, 'on', 'phone'], [p2, 'bow', 'head'], [p2, 'have', 'head'], [p2, 'look out of', 'window']]) = 4  
(Term2) = ([[p2]]) = 2

Common = 1  
P = 1/4 = 0.25  
R = 1/4 = 0.25  
F1 = (2\*0.25\*0.25)/(0.25+0.25) = 0.25

(Term1) = ([[p1, 'on', 'phone'], [p1, 'bow', 'head'], [p1, 'have', 'head'], [p1, 'look out of', 'window']]) = 4  
(Term2) = ([[p1]]) = 1

Common = 1  
P = 1/2 = 0.5  
R = 1/1 = 1  
F1 = (2\*0.5\*1)/(0.5+1) = 0.66

$$F1 * F2 = 0.25 * 0.66 \sim 0.16$$

### Remove Example

Candidate : Opening a small chest filled with personal items P1 takes out a pair of green drawstring pants. In the common sleeping area P1 sets his bags on a lower bunk. A rat runs along a shelf by the headboard. P1 springs up and hits his head on the top bunk. P1 scrambles wildly off the bed grabbing his duffel and peers after the rat with a fearful stare.

Tuples : [[p1, pants, take out], [p1, chest, take out opening], ..., [p1], [chest], [item], [p1, area, set in], [p1, bag, set], [p1, bunk, set on], ..., [p1], ..., [p1, bunk, hit on], [p1, head, hit], [p1, spring], [bunk, top], [p1, head, have], [head], [p1], [bunk], [p1, bed, scramble off], ..., [p1, duffel, have], ..., [p1], [rat], [stare]]

(Term1) = ([[p1, 'take out', 'pants'], [p1, 'take out opening', 'chest'], [p1, 'have', 'duffel'], [p1, 'have', 'head'], [p1, 'hit', 'head'], [p1, 'hit on', 'bunk'], [p1, 'scramble off', 'bed'], [p1, 'set', 'bag'], [p1, 'set in', 'area'], [p1, 'set on', 'bunk'], [p1, 'spring']]) = 11  
(Term2) = ([[p1]]) = 1

CIDER: 91.5 | METEOR: 63.0 | BLEU : 89.0 | SPICE : 79.12 | iSPICE : 12.12

Common = 2  
P = 2/11 = 0.18  
R = 2/11 = 0.18  
F1 = (2\*0.18\*0.18)/(0.18+0.18) = 0.18

Reference : Opening a small chest filled with personal items P1 takes out a pair of green drawstring pants. In the common sleeping area P2 sets his bags on a lower bunk. A rat runs along a shelf by the headboard. P2 springs up and hits his head on the top bunk. P2 scrambles wildly off the bed grabbing his duffel and peers after the rat with a fearful stare.

Tuples : [[p1, pants, take out], [p1, chest, take out opening], [chest, item, fill with], ..., [pants], [pair], [p1], [chest], [item], [p2, area, set in], [p2, bag, set], [p2, bunk, set on], ..., [p2], ..., [p2, bunk, hit on], [p2, head, hit], [p2, spring], [bunk, top], [p2, head, have], [head], [p2], [bunk], [p2, bed, scramble off], ..., [p2, duffel, have], [bed], [duffel], [peer], [p2], [rat], [stare]]

(Term1) = ([[p1, 'take out', 'pants'], [p1, 'take out opening', 'chest'], [p2, 'have', 'duffel'], [p2, 'have', 'head'], [p2, 'hit', 'head'], [p2, 'hit on', 'bunk'], [p2, 'scramble off', 'bed'], [p2, 'set', 'bag'], [p2, 'set in', 'area'], [p2, 'set on', 'bunk'], [p2, 'spring']]) = 11  
(Term2) = ([[p1], [p2]]) = 2

Common = 1  
P = 1/1 = 1  
R = 1/2 = 0.5  
F2 = (2\*1\*0.5)/(1+0.5) = 0.66

$$F1 * F2 = 0.18 * 0.66 \sim 0.12$$

### Replacement Example

Candidate : Meanwhile P1 races to his car in the airport parking lot. P2 stows his bags in the trunk then climbs in. As P2 starts the engine his wipers clear a layer of dirt off the windshield. In an exam room at the clinic the dark haired nurse draws his blood. P2 winces.

Tuples : [[p1, car, race to], [p1, lot, race in], [p1, car, have], [lot, airport], [lot, parking], [car], [p1], [p2, stow], [bag, climb], [bag, trunk, in], [p2, bag, have], [bag], [p2], [trunk], [wiper, layer, clear], [wiper, windshield, clear off], [p2, engine, start], [layer, dirt, of], [engine], [p2], ..., [nurse, haired], ..., [p2, wince], [p2]]

(Term1) = ([[p1, car, race to], [p1, lot, race in], [p1, car, have], [p2, stow], [p2, bag, have], [p2, engine, start], [p2, wince]]) = 7  
(Term2) = ([[p1], [p2]]) = 2

Common = 4  
P = 4/7 = 0.57  
R = 4/7 = 0.57  
F1 = (2\*0.57\*0.57)/(0.57+0.57) = 0.57

Reference : Meanwhile P1 races to his car in the airport parking lot. P1 stows his bags in the trunk then climbs in. As P1 starts the engine his wipers clear a layer of dirt off the windshield. In an exam room at the clinic the dark haired nurse draws his blood. P2 winces.

Tuples : [[p1, car, race to], [p1, lot, race in], [p1, car, have], [lot, airport], [lot, parking], [car], [p1], [p1, stow], [bag, climb], [bag, trunk, in], [p1, bag, have], [bag], [p1], [trunk], [wiper, layer, clear], [wiper, windshield, clear off], [p1, engine, start], [layer, dirt, of], [engine], [p1], [windshield], [layer], [dirt], [wiper], ..., [nurse, haired], ..., [p2, wince], [p2]]

(Term1) = ([[p1, car, race to], [p1, lot, race in], [p1, car, have], [p1, stow], [p1, bag, have], [p1, engine, start], [p2, wince]]) = 7  
(Term2) = ([[p1], [p2]]) = 2

Common = 2  
P = 2/2 = 1  
R = 2/2 = 1  
F2 = (2\*1\*1)/(1+1) = 1

$$F1 * F2 = 0.57 * 1 \sim 0.57$$

Figure 4. We show the effect of identity on captioning metrics using add, remove, and replacement examples. This corresponds to the validation experiment conducted in Tab. 1 of the main paper. For each example, the identity labels are underlined in the candidate and reference captionsets. We also show how iSPICE works by illustrating the tuples, highlighting tuples with identities, and showing the computation of term 1 (left) and term 2 (right) corresponding to tuples with size  $\geq 1$  and  $= 1$  respectively. iSPICE takes into account the identity whereas the other metrics show a high score due to high number of n-gram matches.



[...] pulls her phone from her bag and answers.

[...] frowns uncertainly.

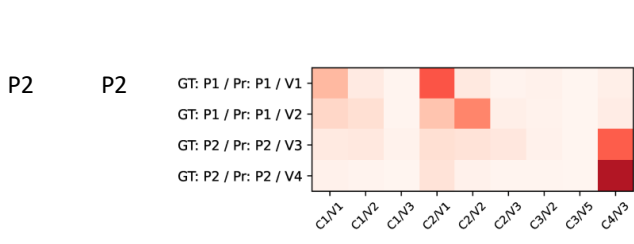
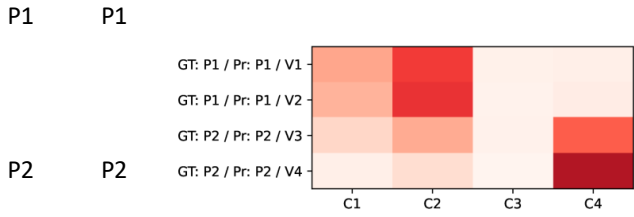
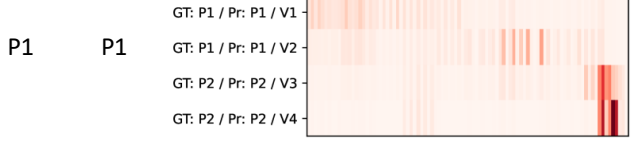
[...] unfolds a map and lifts it up to study it.

[...] lowers his phone and folds it shut.

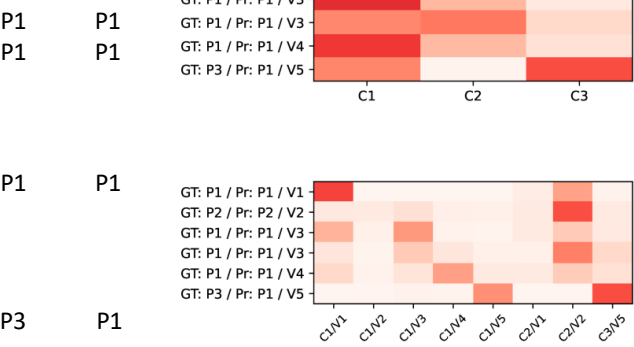
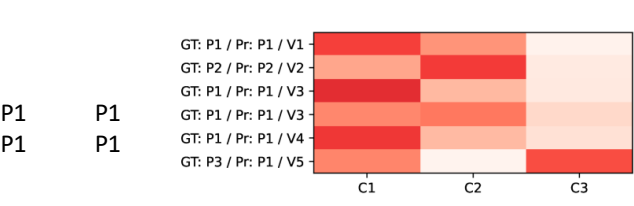
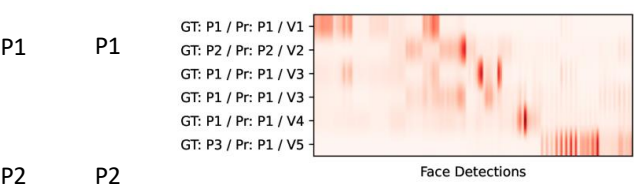
The next morning two women stroll across the street in front of [UNK] apartment building.



**GT**    **Pred**



**GT**    **Pred**



[...] scrutinizes his [UNK] earnest face.

His eyes gleaming in the dim light [...] abruptly gets to his feet and heads for the door

Now [...] talks on his cell as [...] steps into the daylight.

Silhouetted against the sunny day [...] faces the [UNK] door frame and leans his head against it.

Now in a hotel suite a woman applies makeup to [...].

Figure 5. Examples from the Fill-in-the-blanks (FITB) task. On the left, we show one frame from each video of the videoset and the corresponding caption with blanks. In the middle, we show the ground-truth and predicted person id labels. On the right, we show the cross-attention maps (face detections, clusters, and clusters by video ids), presented in Fig. 2. We pick the examples corresponding to captionset 3 and 4 of Fig. 2 for better understanding. In general, we observe that person predictions depend strongly on the cluster features and their attention. In some cases, the identity may be difficult to predict as seen in the last row of the second example, where our model predicts P1 instead of P3, even though the attention masks are correctly focusing on C3/V5.



GT : P1 pours Cheerios.  
Pred : P1 hands her a box.



GT : P1 adds Life cereal to their bowls.  
Pred : P1 takes a bite of food from a box.



GT : P2 gives him a nod, in a t-shirt and sweatpants.  
Pred : P2 nods.



GT : P1 pours blueberries over their cereal.  
Pred : P1 takes a bite.



GT : P1 gives her a bowl.  
Pred : P2 gives a sad smile



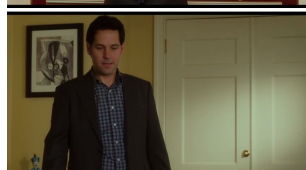
GT : P1 buries his face in his hand and P2 wraps her arm around him.  
Pred : P1 sits on the bench.



GT : Nighttime at the Bowlen Building.  
Pred : Now at the entrance.



GT : Now P1 stands tensely in an elevator.  
Pred : P2 enters the apartment.



GT : Now with his father.  
Pred : P2 returns to his room.



GT : P1 sits on a velvet couch facing his father.  
Pred : P2 sits down.

Figure 6. The above examples showcase MICap’s ability to perform id-aware captioning. We see that the predicted captions are quite good, although terse. While the GT captions tend to be more descriptive in nature, we believe that such behavior may be introduced in the future by incorporating Large Language Models for captioning.

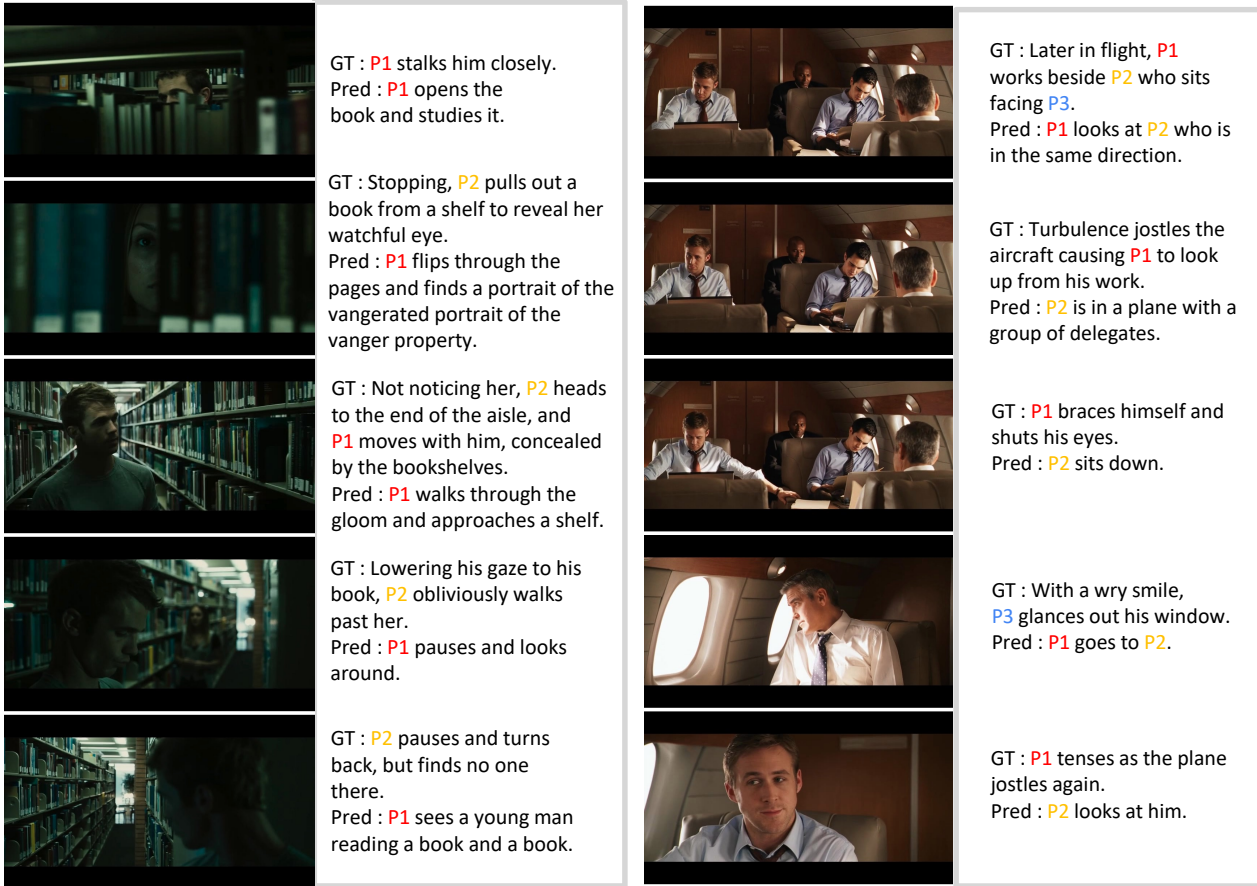


Figure 7. The above examples are relatively difficult cases where there are multiple characters involved with lot of drama or action happening in quick succession. The characters faces are also occluded or partly visible (left example) making it harder to predict identity. We observe that the predicted captions do not capture the tension (e.g. plane turbulence) and the identities.