# Weakly Supervised Dense Video Captioning

## Supplementary Materials

## 1. Region-sequence Generation Algorithm

Algorithm 1 describes the region-sequence generation method, which is based on the CELF (Cost-Effective Lazy Forward selection) algorithm [1]. In this algorithm, $m$ is the number of regions in a sequence, $UC$ and $CB$ are the abbreviation for uniform cost and cost benefit respectively.

## 2. Response Maps

Figure 1 shows some examples of response maps (heatmaps) generated by the Lexical-FCN model. We first associate the response maps to the words in the sentences based on the computed probabilities, and then visualize the best match.
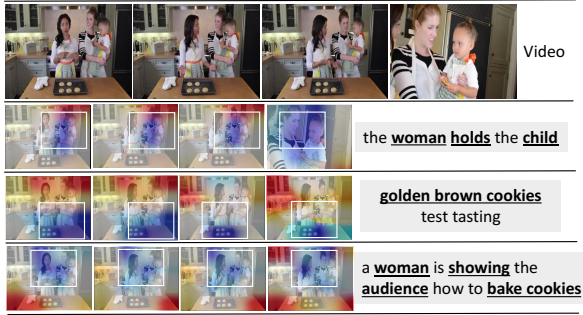


Figure 1. Visualization of learned response maps from the last CNN layer (left), and the corresponding natural sentences (right). The blue areas in the response maps are of high attention, and the region-sequences are highlighted in white bounding-boxes.

## 3. Sentence Re-ranking Module

Figure 2 shows the diagram of our sentence re-ranking module, which re-rank multiple predicted sentences from dense video captioning. This module is similar to [2], which learns the cosine similarity between video features and sentence features with a neural network evaluator.

## 4. More Result Examples

More result examples of our DenseVideoCap system are provided in Figures 3, 4, 5, 6.

---

**Algorithm 1** Region-sequence generation by submodular maximization with the CELF algorithm [1].

```
 1: function LAZYFORWARD(S_v, x_v, R, m, type)
 2:     A ← ∅;                          ▷ Start with the empty sequence
 3:     for each r ∈ S_v do L(w; r) ← ∞   ▷ Init marginal gains
 4:     end for
 5:     while |A| < m do
 6:         for each r ∈ S_v \ A do cur_s ← false;
 7:         end for
 8:         while true do                 ▷ Begin loop
 9:             if type = UC then          ▷ Uniform cost
10:                 r* ← arg max  L(w; r);    ▷ Max gain
                         r∈S_v\A
11:             end if
12:             if type = CB then          ▷ Cost benefit
13:                 r* ← arg max  L(w;r)/R(r);  ▷ Max gain / cost
                         r∈S_v\A
14:             end if
15:             if cue_{r*} then A ← A ∪ {r*}; break;
16:             else                      ▷ Update marginal gain
17:                 L(w; r) ← R(A ∪ {r}) − R(A);
18:                 cur_{r*} ← true;
19:             end if
20:         end while
21:     end while
22:     return A;                        ▷ Return region-sequence
23: end function
24:
25: function MAIN(S_v, x_v, R, m)
26:     A_UC ← LAZYFORWARD(S_v, x_v, R, m, UC)
27:     A_CB ← LAZYFORWARD(S_v, x_v, R, m, CB)
28:     return arg max{R(A_UC), R(A_CB)}
29: end function
```

---

## References

[1] J. Leskovec, A. Krause, and et al. Cost-effective outbreak detection in networks. In *ACM SIGKDD*, 2007. 1

[2] R. Shetty and J. Laaksonen. Frame-and segment-level features and candidate pool evaluation for video caption generation. *arXiv:1608.04959*, 2016. 1
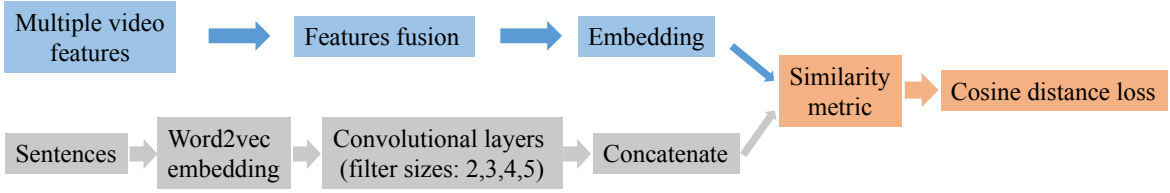
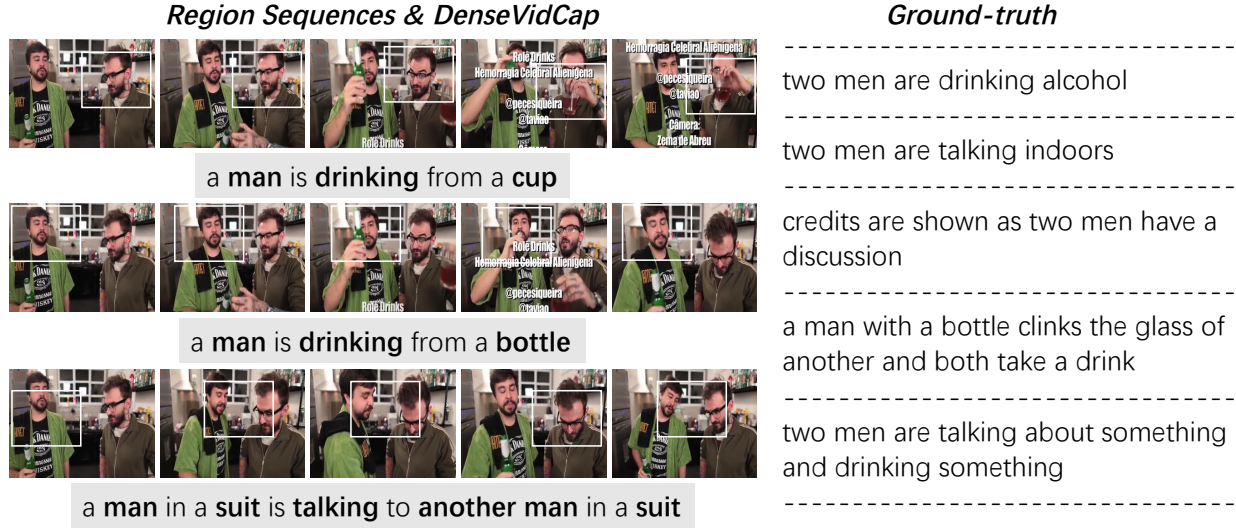Figure 2. Illustration of the sentence re-ranker module.

| Region Sequences & DenseVidCap | Ground-truth |
|---|---|



**Region Sequences & DenseVidCap**

a **man** is **drinking** from a **cup**

a **man** is **drinking** from a **bottle**

a **man** in a **suit** is **talking** to **another man** in a **suit**

**Ground-truth**

----------------------------------

two men are drinking alcohol

----------------------------------

two men are talking indoors

----------------------------------

credits are shown as two men have a discussion

----------------------------------

a man with a bottle clinks the glass of another and both take a drink

----------------------------------

two men are talking about something and drinking something

----------------------------------

Figure 3. Left: Examples of dense sentences produced by our *DenseVidCap* method and corresponding *region sequences*; Right: Ground-truth (video6974).



**Region Sequences & DenseVidCap**

a **man** is **running** on a **track**

a **group** of **men running** in a **race track**

a **track race track** and **field runners run** and **running fast** in a **race**

**Ground-truth**

------------------------------------------

a man is running

------------------------------------------

a group of men are running down a race track

------------------------------------------

athletes are running around a track

------------------------------------------

a group of people are running as fast as they can

------------------------------------------

men run a race around a track while a male voice narrates

------------------------------------------

several young men are racing in a track meet

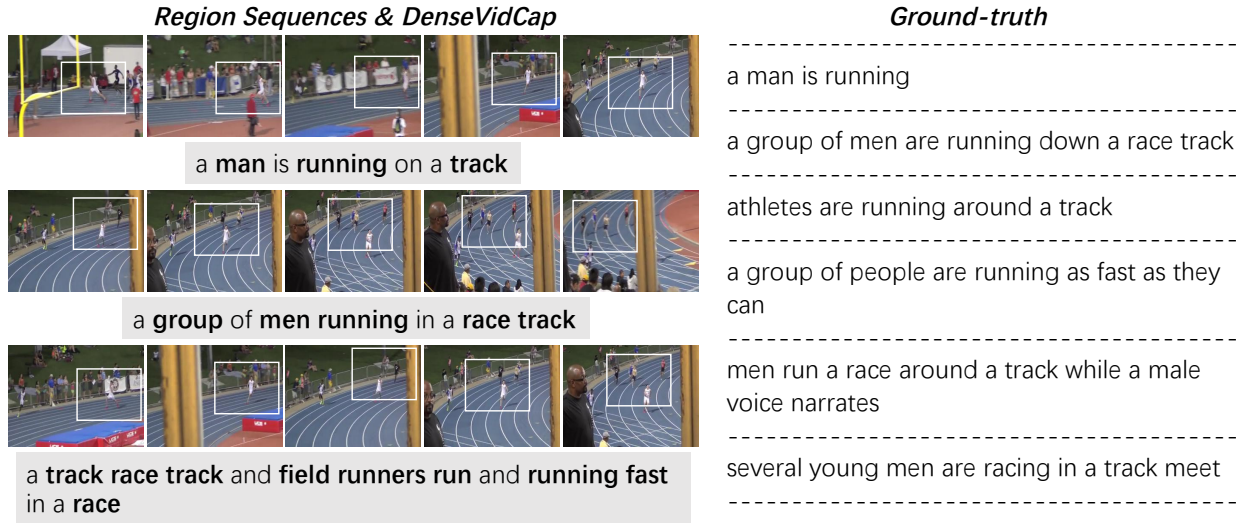------------------------------------------

Figure 4. Left: Examples of dense sentences produced by our *DenseVidCap* method and corresponding *region sequences*; Right: Ground-truth (video6967).

## Region Sequences & DenseVidCap



a **woman** is **cooking food**



a **woman** is **teaching how** to **cook**



a **woman** is **cooking food** in a **pot**



a **woman cooks** a **meal** in a **pot** on the **stove**

## Ground-truth

a person is cooking

a person teaching a recipe

a girl stirring a mixture in a pot

this is a cooking video

a woman is giving a cooking instructional video

creating a thick chocolate pudding

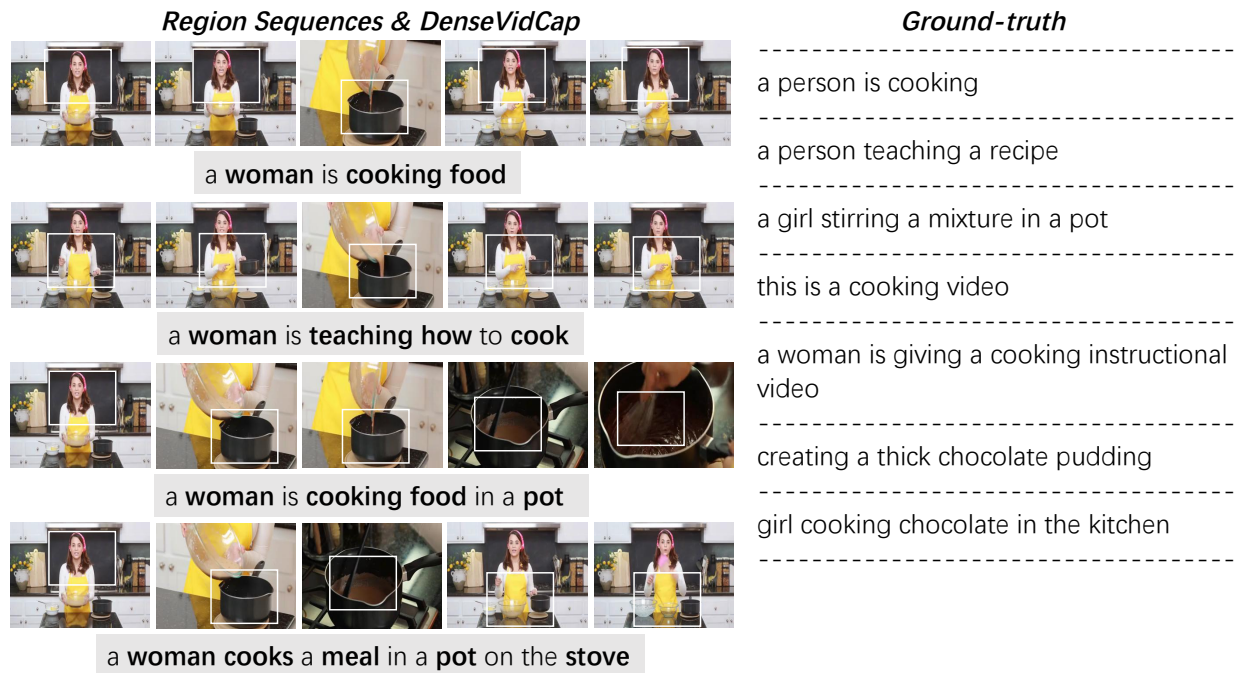girl cooking chocolate in the kitchen

Figure 5. Left: Examples of dense sentences produced by our *DenseVidCap* method and corresponding *region sequences*; Right: Ground-truth (video6911).

## Region Sequences & DenseVidCap



a **man** is **hitting** a **ping pong ball** in a **stadium**



a **man** in a **blue shirt** is **hitting** a **ping pong ball**

## Ground-truth

a couple of men are playing ping pong inside

two men play ping pong in a stadium with people watching

a ping-pong player wins a point runs away from the blue table jumps with knees high runs to gym mats

two asian men playing ping pong in a tournament as the man in the blue shirt won

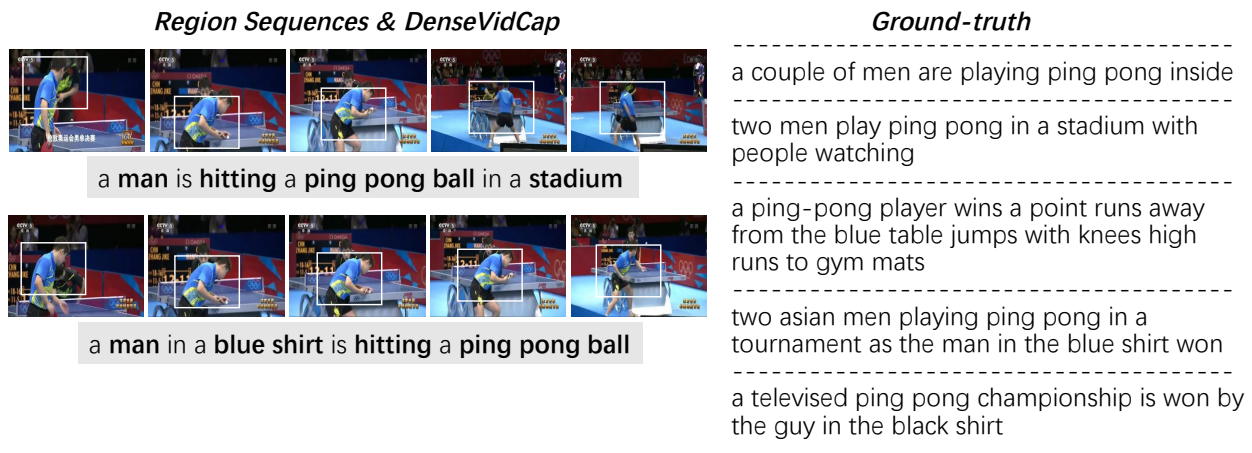a televised ping pong championship is won by the guy in the black shirt

Figure 6. Left: Examples of dense sentences produced by our *DenseVidCap* method and corresponding *region sequences*; Right: Ground-truth (video6973).