

# Generating Structured Pseudo Labels for Noise-resistant Zero-shot Video Sentence Localization

Minghang Zheng<sup>1</sup>, Shaogang Gong<sup>2</sup>, Hailin Jin<sup>3</sup>, Yuxin Peng<sup>1, 4</sup>, and Yang Liu<sup>1, 5\*</sup>

<sup>1</sup>Wangxuan Institute of Computer Technology, Peking University

<sup>2</sup>Queen Mary University of London, <sup>3</sup>Adobe Research

<sup>4</sup> National Key Laboratory for Multimedia Information Processing, Peking University

<sup>5</sup> National Key Laboratory of General Artificial Intelligence, BIGAI

{minghang, pengyuxin, yangliu}@pku.edu.cn

s.gong@qmul.ac.uk, hljin@adobe.com



# Task

## Zero-shot Temporal Sentence Localization

- **Inputs:** Video + Sentence query
- **Outputs:** Target video clip (start and end timestamps)
- **Zero-shot Setting:** No manual annotation required

### Video:



### Query:

A man puts on gloves and then clean the snow

5.5s

Output

15.1s

} **Unavailable**



# Motivation

## Existing zero-shot methods:

- Generating **pseudo-events** and **pseudo-queries**
- **Training** with pseudo-event and pseudo-query

## Drawbacks:

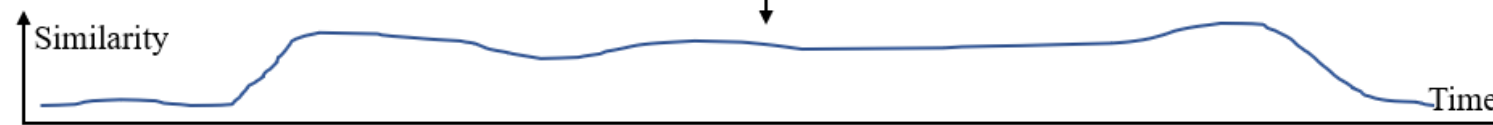
- 1. Pseudo queries are **too simple**
- 2. **Unalignment** between pseudo-events and pseudo-queries
- 3. Ignoring the **noise** in the pseudo labels

### Zero-shot Video:



**Step 1. Pseudo Event:** ← [8s, 11s] →

**Step 2. Pseudo Query:** Person cleans the snow ↑ Unaligned ↓



(c) Existing pipeline



# Method

## Our Structured Pseudo-Label (SPL) generation:

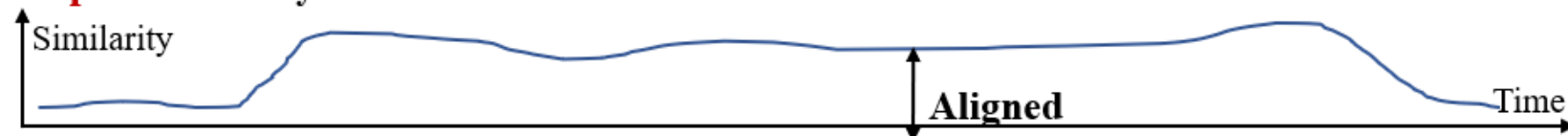
- Generate **free-form** pseudo-queries using image description models
- Generate pseudo-events based on the **event temporal structure**
  - the video **inside** the event has a **high** correlation with the query
  - the video **outside** the event has a **low** correlation with the query
- **Reduce noise** during training
  - Sample **re-weight** and label **refinement**

Video:



**Step 1. Pseudo Query:** A man cleans the snow with a brush

**Step 2. Similarity:**



**Step 3. Pseudo Event:**

[4s, 15s]

(d) Ours pipeline



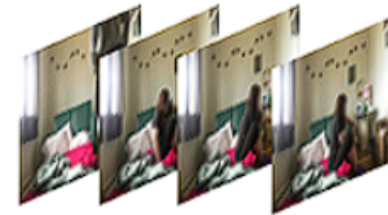
# 1. Pseudo Query Generation

**Aims:** Generate **free-form** pseudo-queries

- Densely sample video frames
- Generate pseudo queries from video frames using pretrained BLIP<sup>1</sup> model

<sup>1</sup>Li, Junnan, et al. "Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation." International Conference on Machine Learning. ICML, 2022.

## 1. Pseudo Query Generation



**Image Caption**

**Query Candidates:**

A woman takes shoes off.  
A woman sitting on a bed.  
A bed in the room.  
...

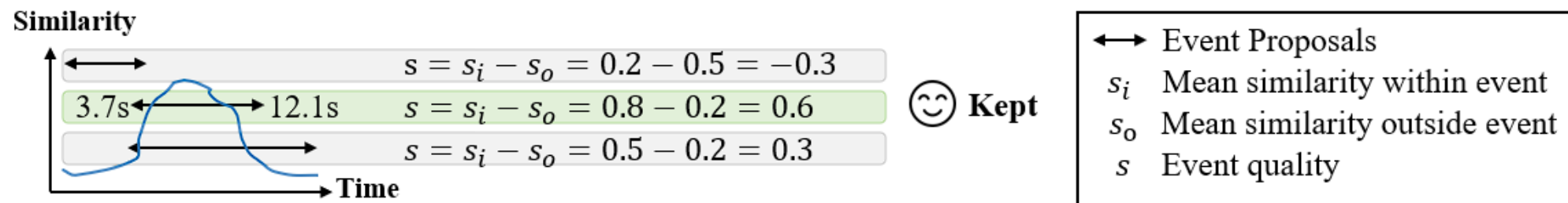


## 2. Pseudo Event Generation

**Aims:** Generate pseudo-events based on **event temporal structure**

- Videos **within** events have **high relevance** to queries
- Videos **outside** events have **low relevance** to queries
- Calculate similarity  $S$  between pseudo-query and video frames
- Event quality:  $Q_{ik} = \frac{1}{N_{p_k}} \sum_{j \in p_k} S_{ij} - \frac{1}{N - N_{p_k}} \sum_{j \notin p_k} S_{ij}$
- Choose the event proposal with highest quality

### Pseudo Event Generation

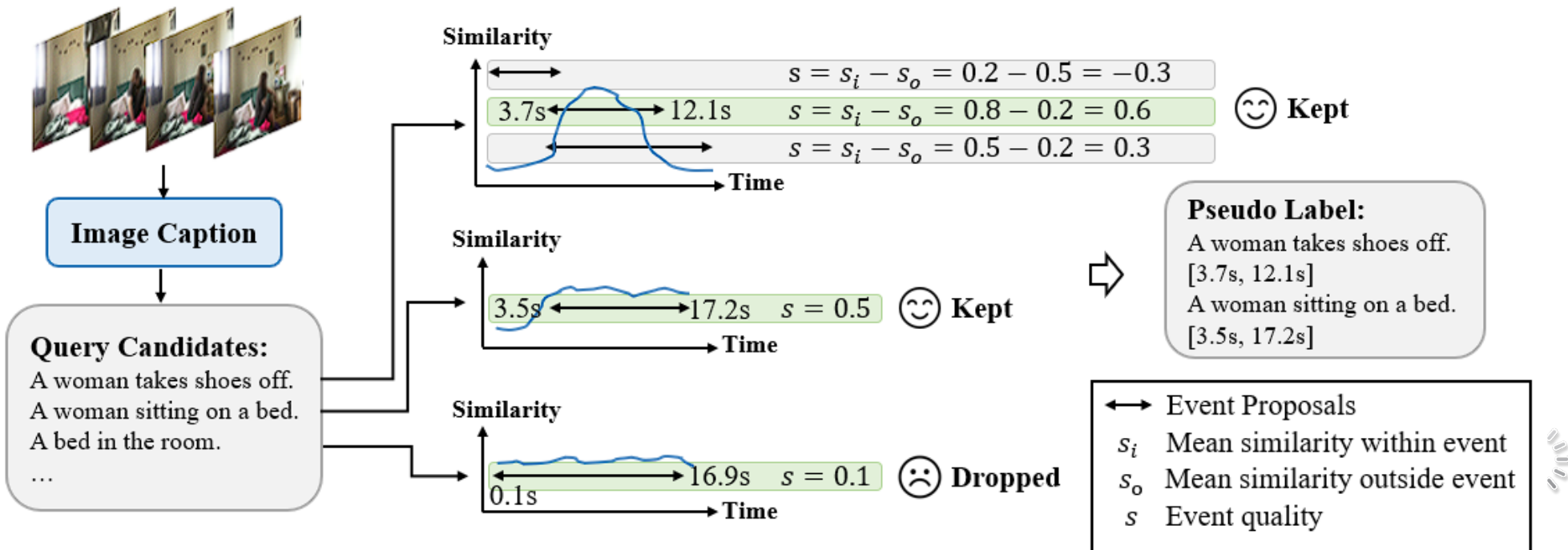




## 2. Pseudo Event Generation

- Filter out **low-quality** pseudo-query event pairs
  - Keep top  $K$  pseudo-query-event pairs with high event quality
  - Use non-maximum suppression to eliminate pseudo-query-event pairs with high event overlap.

### Pseudo Event Generation



### 3. Training with Noisy Pseudo Labels

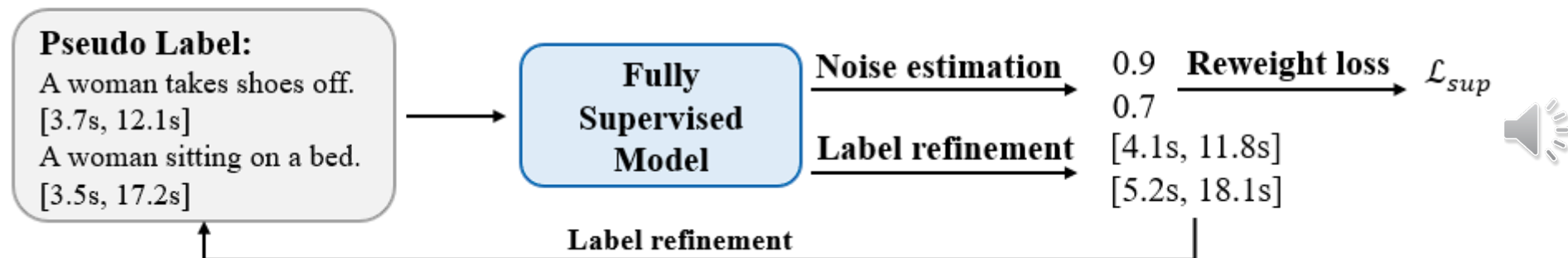
**Aims:** Train model using pseudo-labels and **reduce noise** in the pseudo-labels

- **Sample re-weight:** Estimate noise based on the confidence score  $s^{conf}$  and the IoU  $s^{iou}$  of predictions and pseudo-label and weight the sample loss

$$w = \alpha \frac{1}{1 - s^{iou}} + (1 - \alpha) \frac{1}{1 - s^{conf}}$$

- **Label refinement:** If the model's prediction confidence is high, consider the prediction as a new pseudo-label.

#### Training with Noisy Pseudo Label





# Experiments

## Datasets

- ActivityNet Captions
- Charades-STA

## Metrics

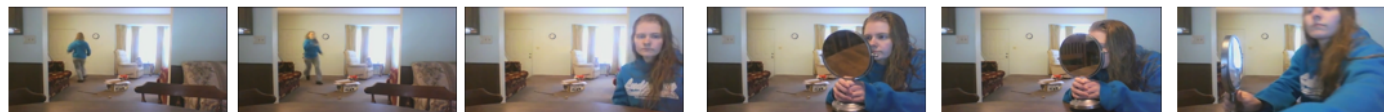
- R@m
- mIoU

**Query:** A man in a red tank top is crossing the monkey bars.



(a) ActivityNet Captions

**Query:** Person runs to a table.



(b) Charades-STA



# Comparing with SOTA

## Results

- **Best zero-shot performance** on most metrics

Method	Sup.	Charades-STA				ActivityNet Captions			
		R@0.3	R@0.5	R@0.7	mIoU	R@0.3	R@0.5	R@0.7	mIoU
2D-TAN (Zhang et al., 2020)	fully	-	39.81	23.25	-	58.75	44.05	27.38	-
EMB (Huang et al., 2022)		<b>72.50</b>	58.33	39.25	<b>53.09</b>	<b>64.13</b>	44.81	26.07	<b>45.59</b>
MGSL-Net (Liu et al., 2022)		-	<b>63.98</b>	<b>41.03</b>	-	-	<b>51.87</b>	<b>31.42</b>	-
CRM (Huang et al., 2021)	weakly	53.66	34.76	16.37	-	55.26	32.19	-	-
CNM* (Zheng et al., 2022a)		60.39	35.43	15.45	-	55.68	<b>33.33</b>	-	-
CPL (Zheng et al., 2022b)		<b>66.40</b>	<b>49.24</b>	<b>22.39</b>	-	<b>55.73</b>	31.37	-	-
Gao et al.* (Gao and Xu, 2021)	no	46.69	20.14	8.27	-	46.15	26.38	11.64	-
PSVL* (Nam et al., 2021)		46.47	31.29	14.17	31.24	44.74	30.08	14.74	29.62
PZVMR* (Wang et al., 2022)		46.83	33.21	18.51	32.62	45.73	31.26	<b>17.84</b>	30.35
Kim et al.* (Kim et al., 2023)		52.95	37.24	19.33	36.05	47.61	<b>32.59</b>	15.42	31.85
SPL* (ours)	no	<b>60.73</b>	<b>40.70</b>	<b>19.62</b>	<b>40.47</b>	<b>50.24</b>	27.24	15.03	<b>35.44</b>

more experiments and ablation studies can be found in paper

# Conclusion

- Propose a zero-shot video sentence localization method based on **structured pseudo-label generation** that is **robust to noise**
- Generate **free-form** pseudo-queries and generate pseudo-events based on **event temporal structure**.
- Reduce the influence of noise in pseudo-labels by **sample reweight** and **label refinement**
- **Best zero-shot performance** on two datasets



# Thank you!

Code

