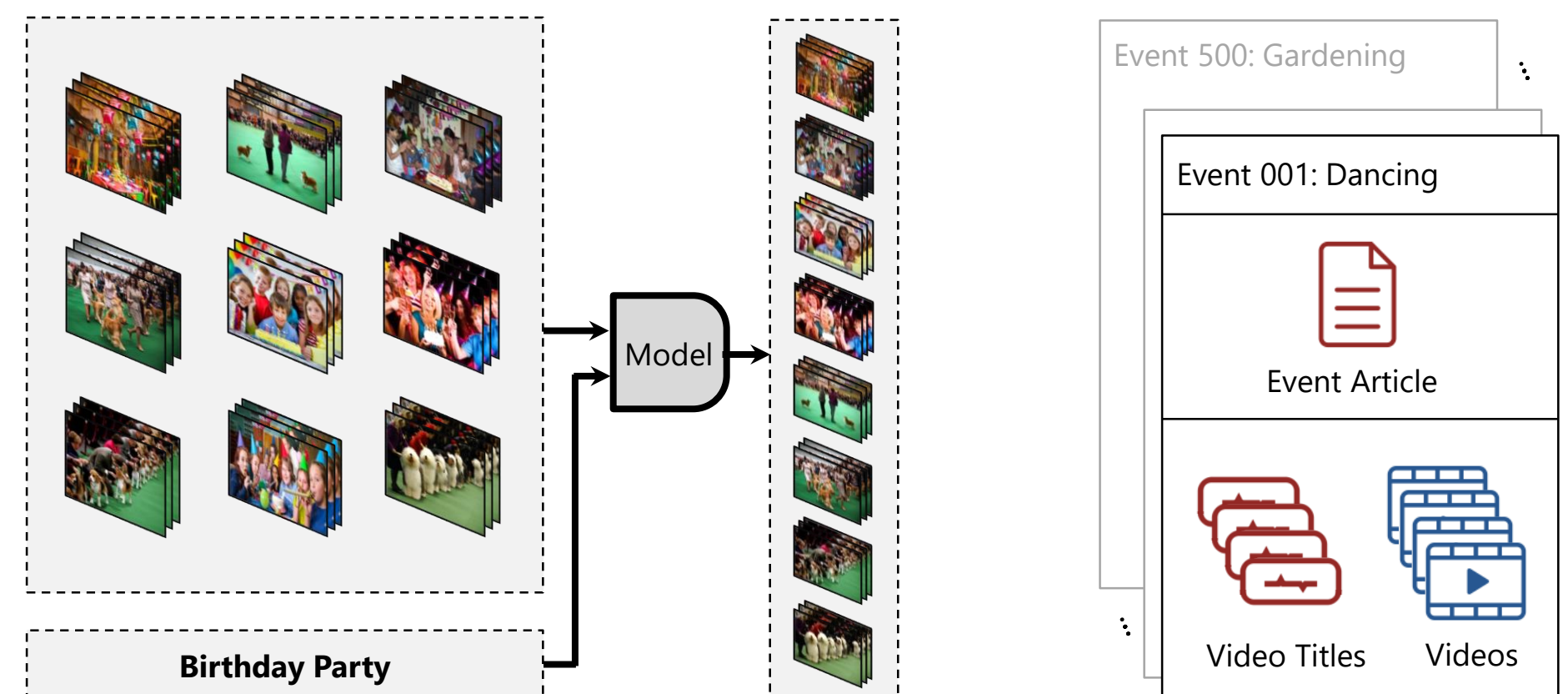




Unified Embedding and Metric Learning for Zero-Exemplar Event Detection

Noureldien Hussein, Efstratios Gavves, Arnold W. M. Smeulders | Quva Lab, University of Amsterdam

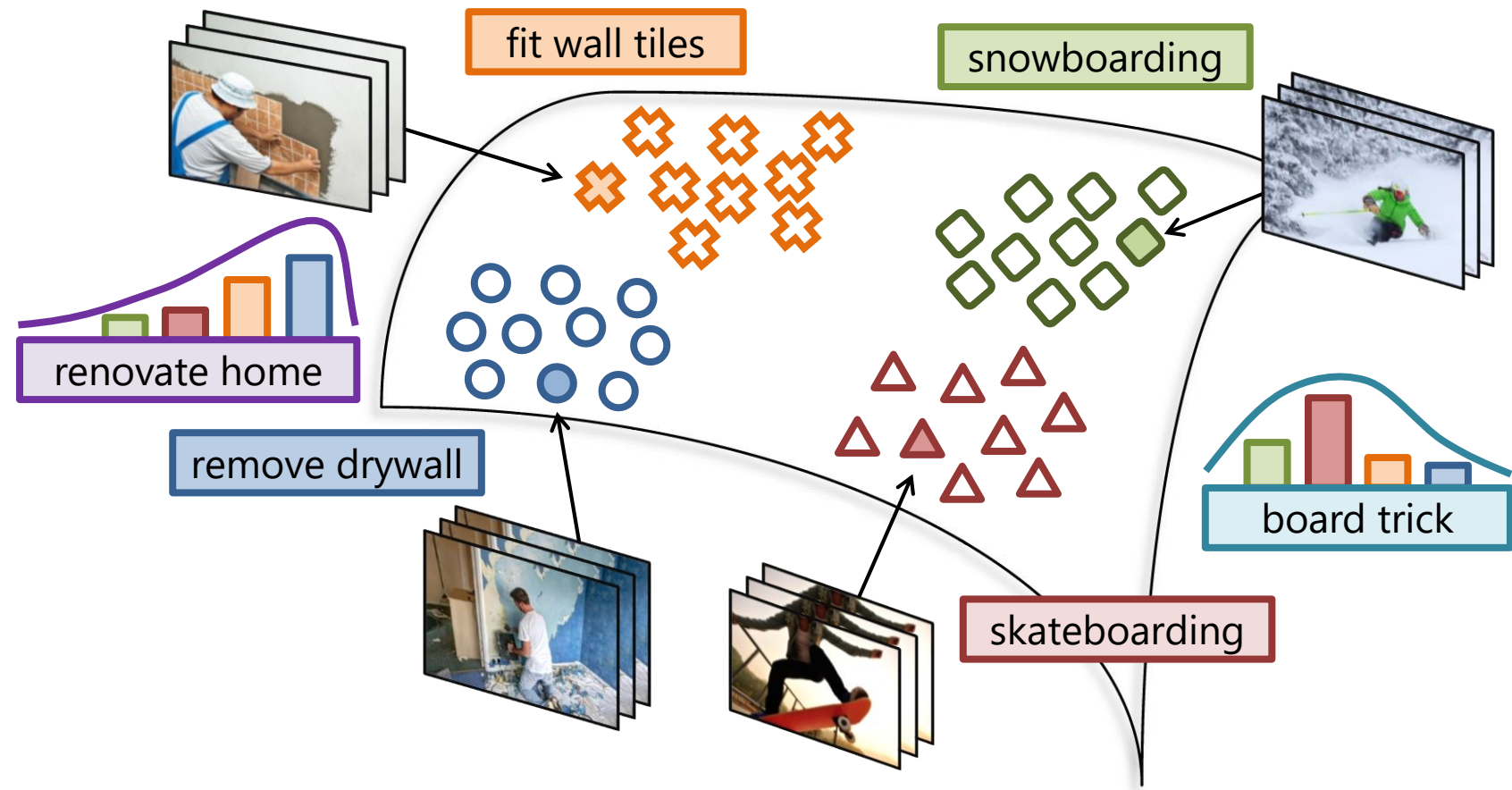
Problem



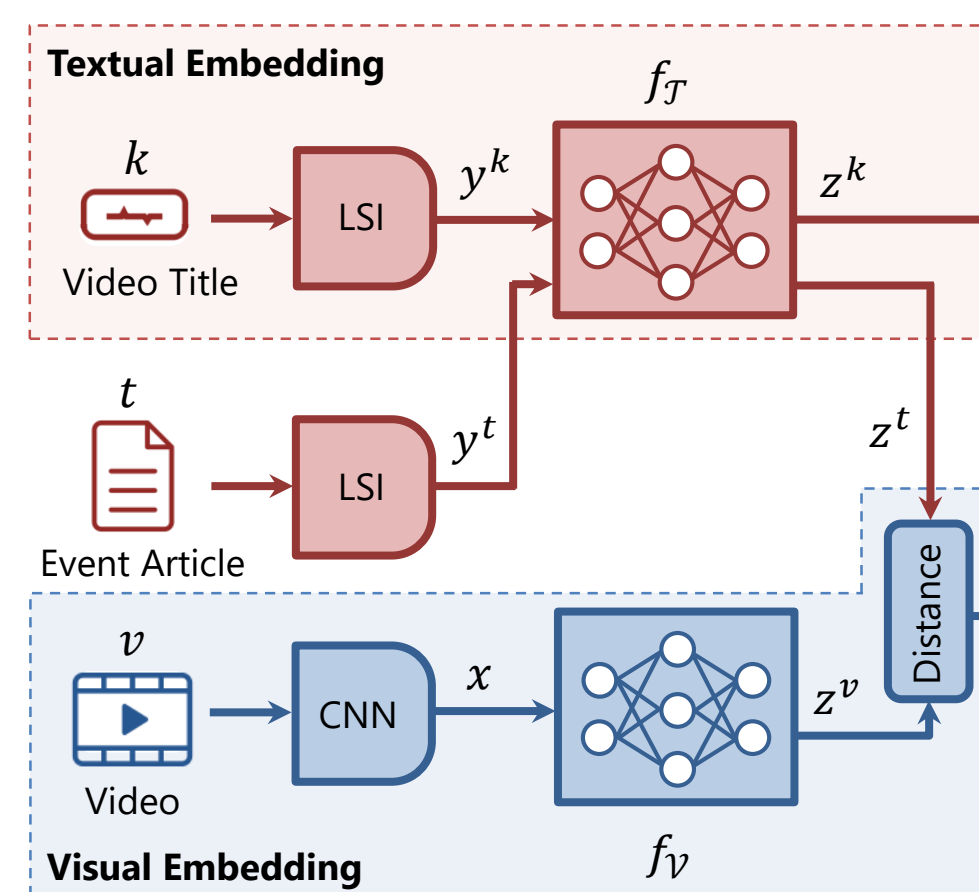
Zero-exemplar Event Detection (ZED) is posed as a video retrieval task. Given test videos and a novel query, the model is required to rank the videos accordingly.

We use video samples from EventNet and event articles from WikiHow.

Approach



We pose ZED as learning from a set of predefined events. Given video exemplars of events "removing drywall" or "fit wall tiles", one may detect a novel event "renovate home" as a probability distribution over the predefined events.



At the top, network f_T learns to classify the title feature y^k into one of M event categories. In the middle, we borrow the network f_T to embed the event articles feature y^t as $z^t \in \mathcal{Z}$. Then, at the bottom, the network f_V learns to embed the video feature x as $z^v \in \mathcal{Z}$ such that the distance between (z^t, z^v) is minimized, in the learned metric space \mathcal{Z} .

$$\mathcal{L}^u = \mathcal{L}_{con} + \mathcal{L}_{log} \quad (1)$$

$$\mathcal{L}_{log} = \sum_{i=1}^N \sum_{j=1}^M -y_i^j \log f_T^j(k_i; W_T)$$

$$\mathcal{L}_{con} = \frac{1}{2N} \sum_{i=1}^N h_i \cdot d_i^2 + (1 - h_i) \max(1 - d_i, 0)^2$$

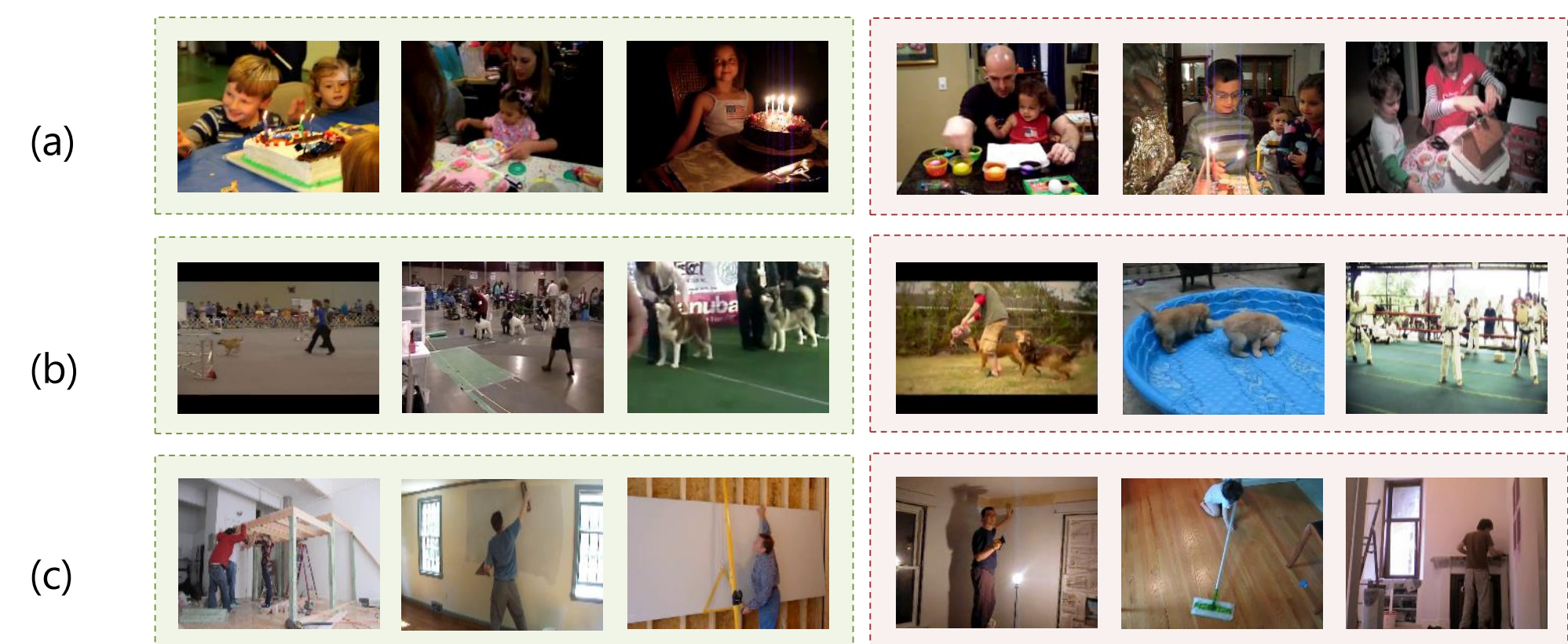
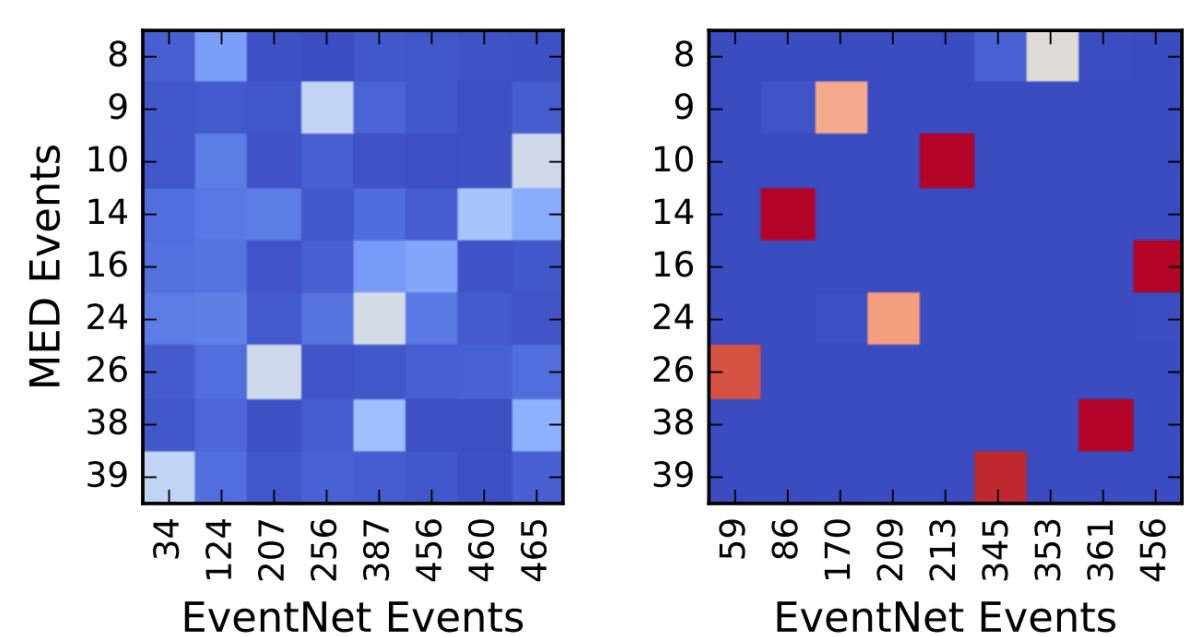
$$d_i = \|f_T(t_i; W_T) - f_V(v_i; W_V)\|_2$$

Novelties

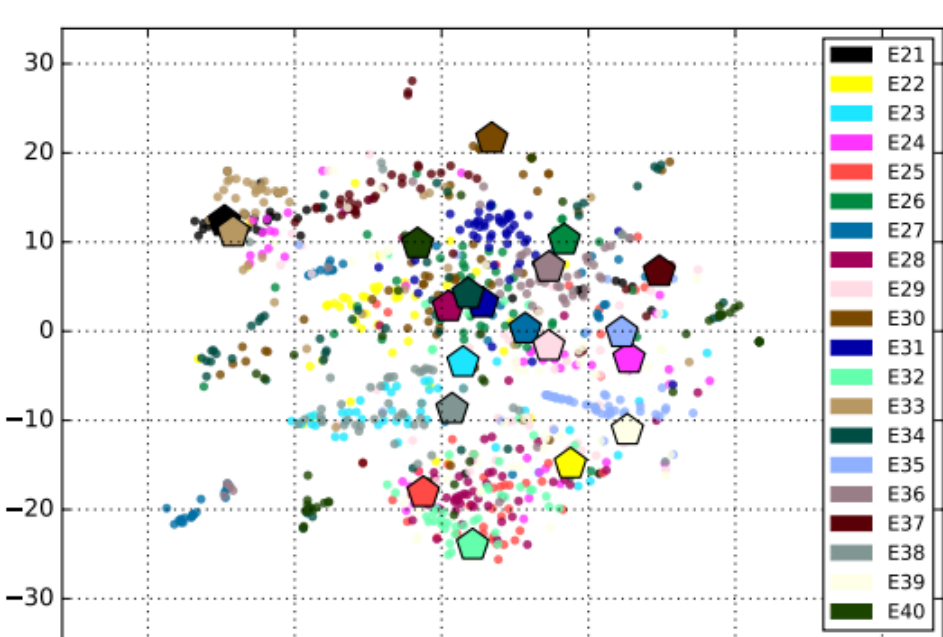
- Unified embedding** for cross-modalities with **metric loss** for maximum discrimination between events.
- Textual embedding** poses a novel query a probability of predefined events.
- External data source**, of event articles and related videos, with end-to-end learning from cross-modal pairs.

Results

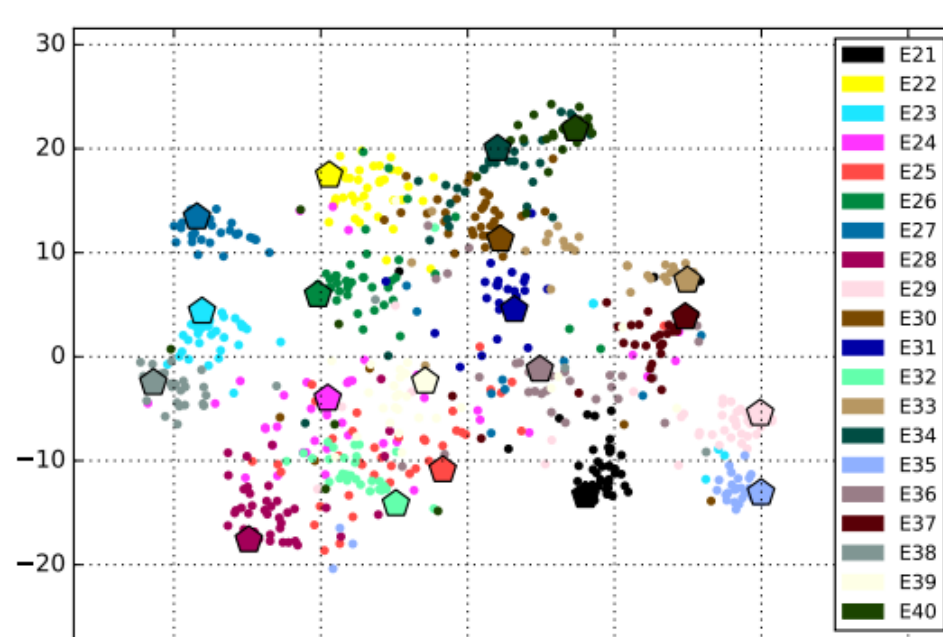
Our textual embedding f_T maps the text description of MED events to EventNet events better than off-the-shelf LDA, LSI or Doc2vec.



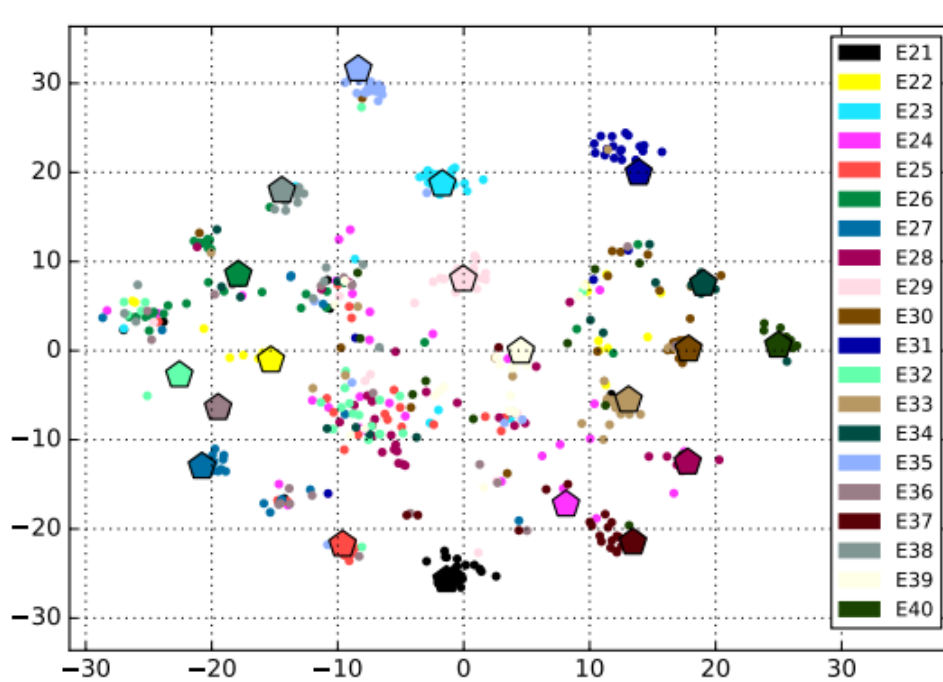
Success (left) and failure (right) video examples of three different events: (a) birthday party, (b) dog show, (c) renovating home.



(a) Visual Embedding ($model^V$)



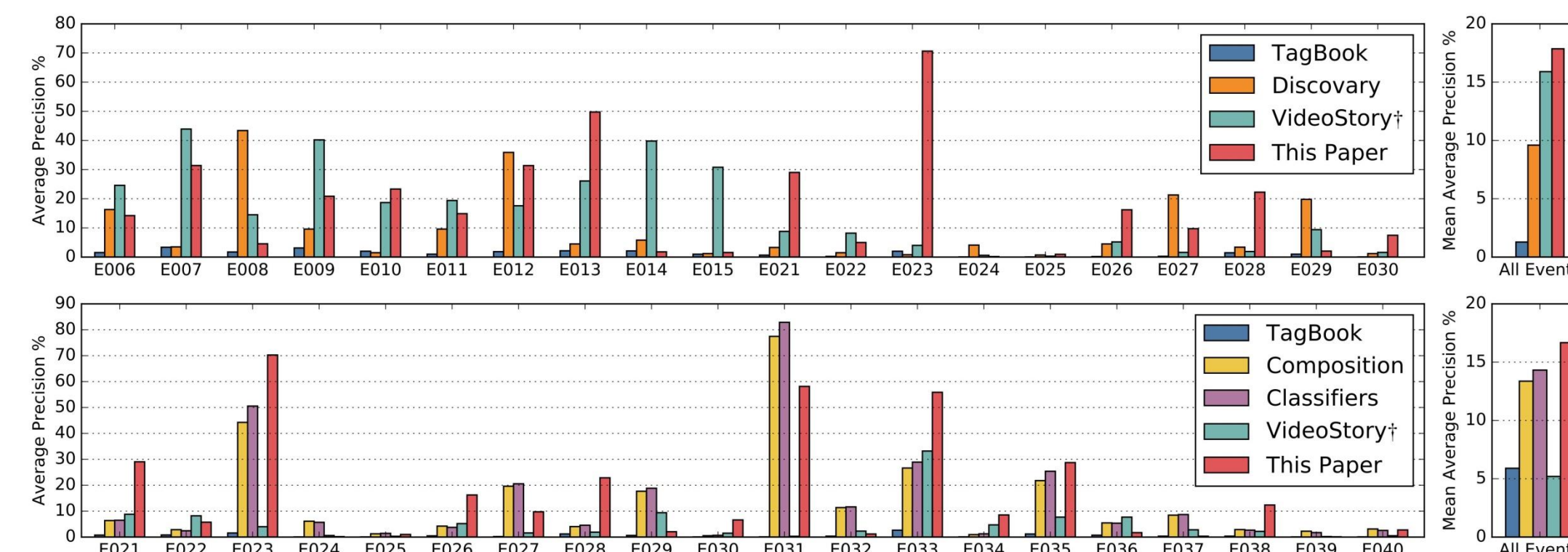
(b) Separate Embedding ($model^S$)



(c) Unified Embedding ($model^U$)

• The unified embedding (c) is doing a better job in discriminating the text articles of the events.

• In (c), projecting the videos on their related events is much better than the baselines (a) and (b).



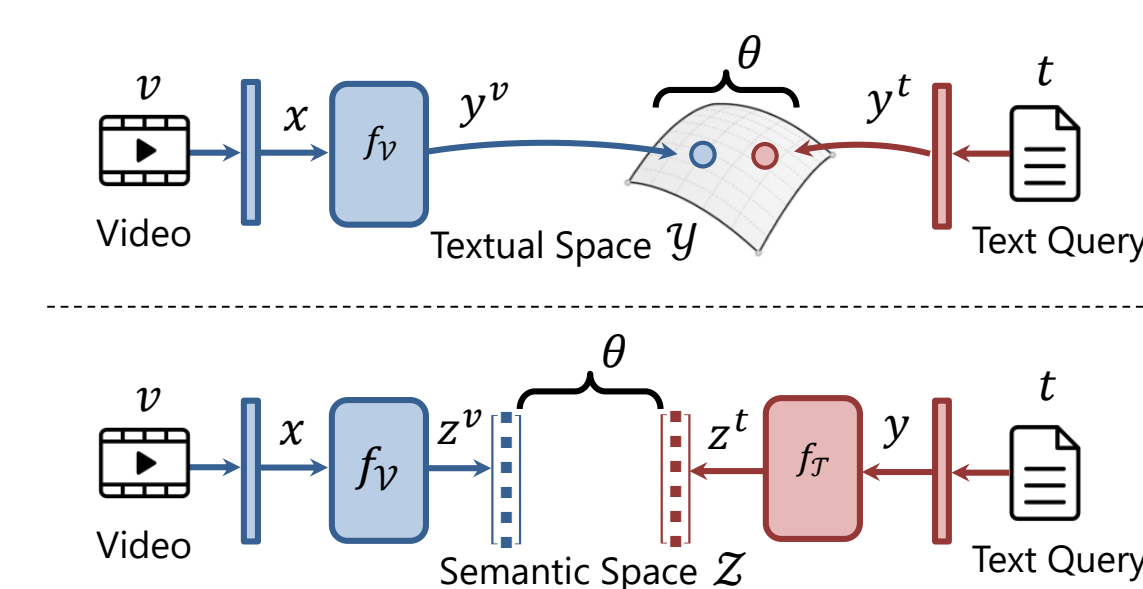
Event detection accuracies: per-event average precision (AP) and per-dataset mean average precision (mAP) for MED-13 and MED-14 datasets.

Method		MED13	MED14	Baseline	Loss	Metric	$f_V(\cdot)$	$f_T(\cdot)$	MED13	MED14
TagBook [18]	ToM '15	12.90	05.90	$model^V$	\mathcal{L}_{mse}^V (2)	✗	✓	✗	11.90	10.76
Discovery [7]	ICAI '15	09.60	—	$model^C$	\mathcal{L}_{con}^C (3)	✓	✓	✓	13.29	12.31
Composition [8]	AAAI '16	12.64	13.37	$model^S$	\mathcal{L}_{log} (4)	✗	✓	✗	15.60	13.49
Classifiers [9]	CVPR '16	13.46	14.32	$model^N$	\mathcal{L}_{mse}^N (5)	✗	✓	✓	15.92	14.36
VideoStory [†] [17]	PAMI '16	15.90	05.20	$model^U$	\mathcal{L}^U (1)	✓	✓	✓	17.86	16.67
VideoStory* [17]	PAMI '16	20.00	08.00							
This Work ($model^U$)		17.86	16.67							

Left: retrieval accuracy (mAP) of our model vs. related works for MED-13 and MED-14 datasets. Right: retrieval accuracy (mAP) of our model (unified embedding) vs. other baselines.

Experiments

Model overview of baseline methods. Top: visual embedding ($model^V$). Bottom: separate embedding ($model^S$).



$$\mathcal{L}_{mse}^V = \frac{1}{N} \sum_{i=1}^N \|y_i - f_V(v_i; W_V)\|_2^2 \quad (2)$$

$$\mathcal{L}_{con}^C = \frac{1}{2N} \sum_{i=1}^N h_i \cdot d_i^2 + (1 - h_i) \max(1 - d_i, 0)^2 \quad (3)$$

$$\mathcal{L}_{log} = \sum_{i=1}^N \sum_{j=1}^M -y_i^j \log f_T^j(k_i; W_T) \quad (4)$$

$$\mathcal{L}_{mse}^N = \frac{1}{N} \sum_{i=1}^N \|f_T(t_i; W_T) - f_V(v_i; W_V)\|_2^2 \quad (5)$$

Loss functions used to train the baseline models: visual embedding $model^V$ (2), contrastive visual $model^C$ (3), separate embedding $model^S$ (4) and non-metric embedding $model^N$ (5).

Take Home

The Good: external knowledge (EventNet, WikiHow) is leveraged for better zero-exemplar event detection.

The Bad: no fine-grained event detection, e.g. "fixing musical instrument" vs. "tuning musical instrument".

The Ugly: is average pooling enough for video representation or temporal modeling is required?

