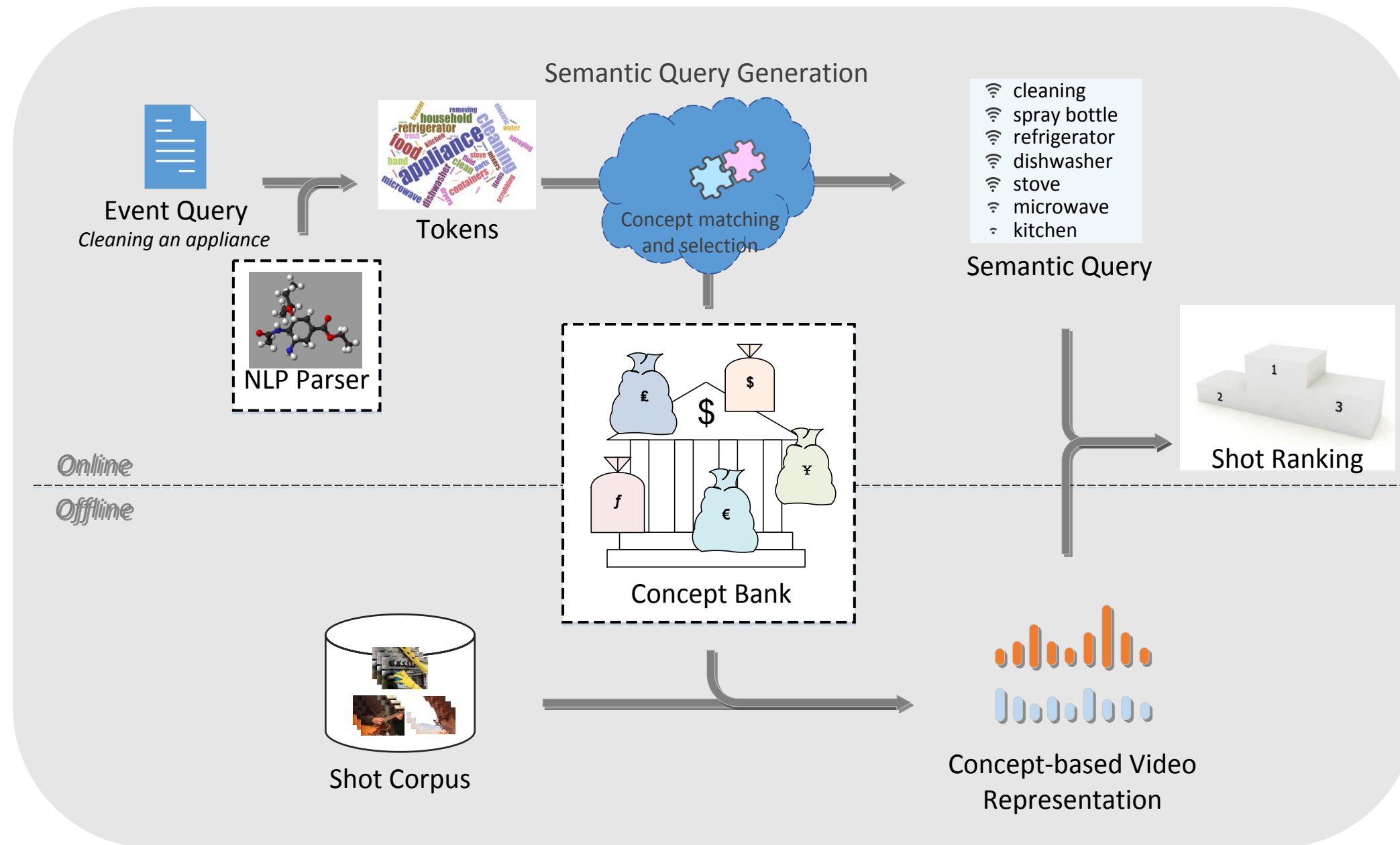


Ad-hoc Video Search (AVS)

Framework



Highlighted Concept Reranking (HCR)

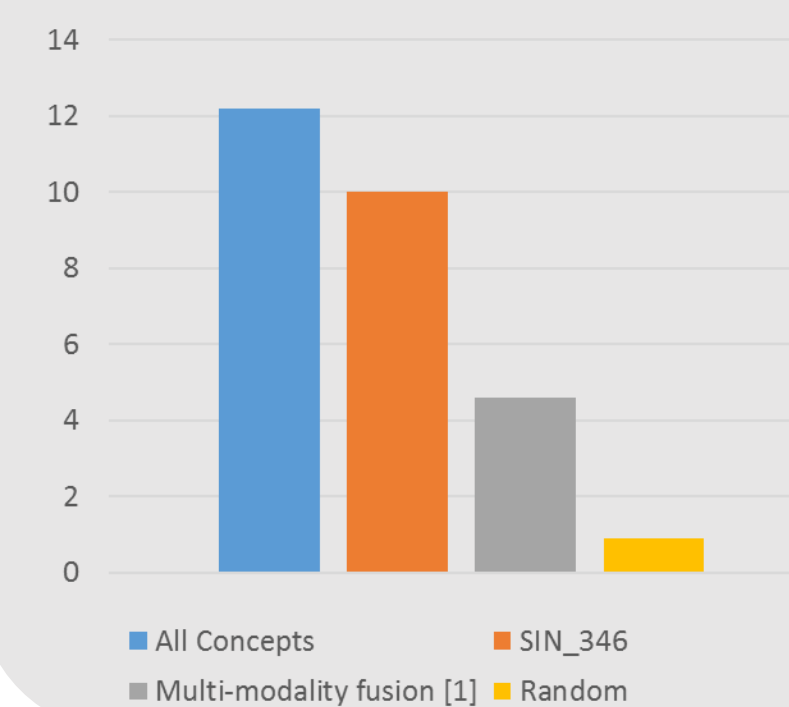
- After **Semantic Query Generation**, human experts can highlight a few important concepts.
- The videos on the top which contain the highlighted concepts are given high weights and their ranking are largely boosted.

Vocabulary

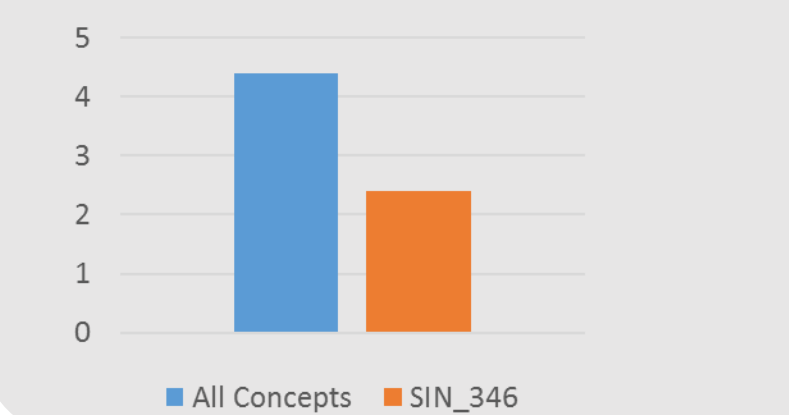
ImageNetShuffle-13K
 ImageNet-1000
 Places-205
 SIN-346
 Research Collection-497

Experiments

Performance for different concept settings in Video-Search 2008 dataset (MinfAP%)



Performance for different concept settings in AVS 2016 evaluation dataset (MinfAP%)

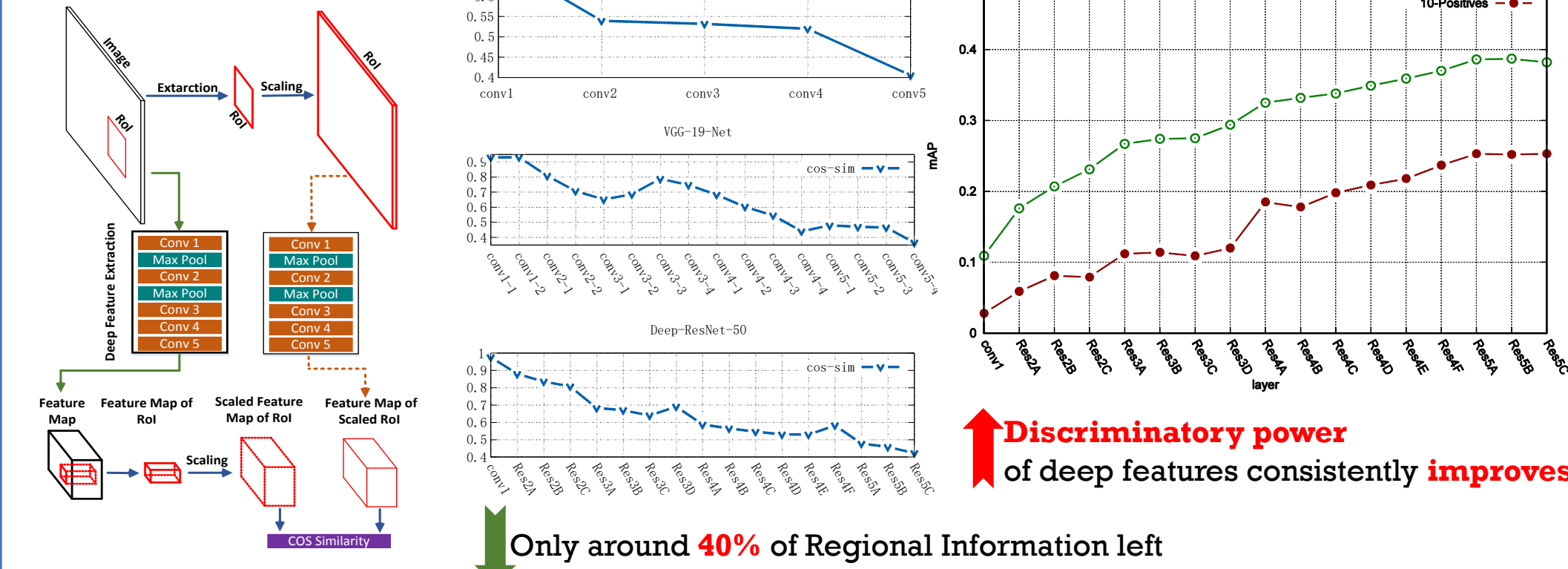


Performance with and without highlighted concept ranking (HCR)

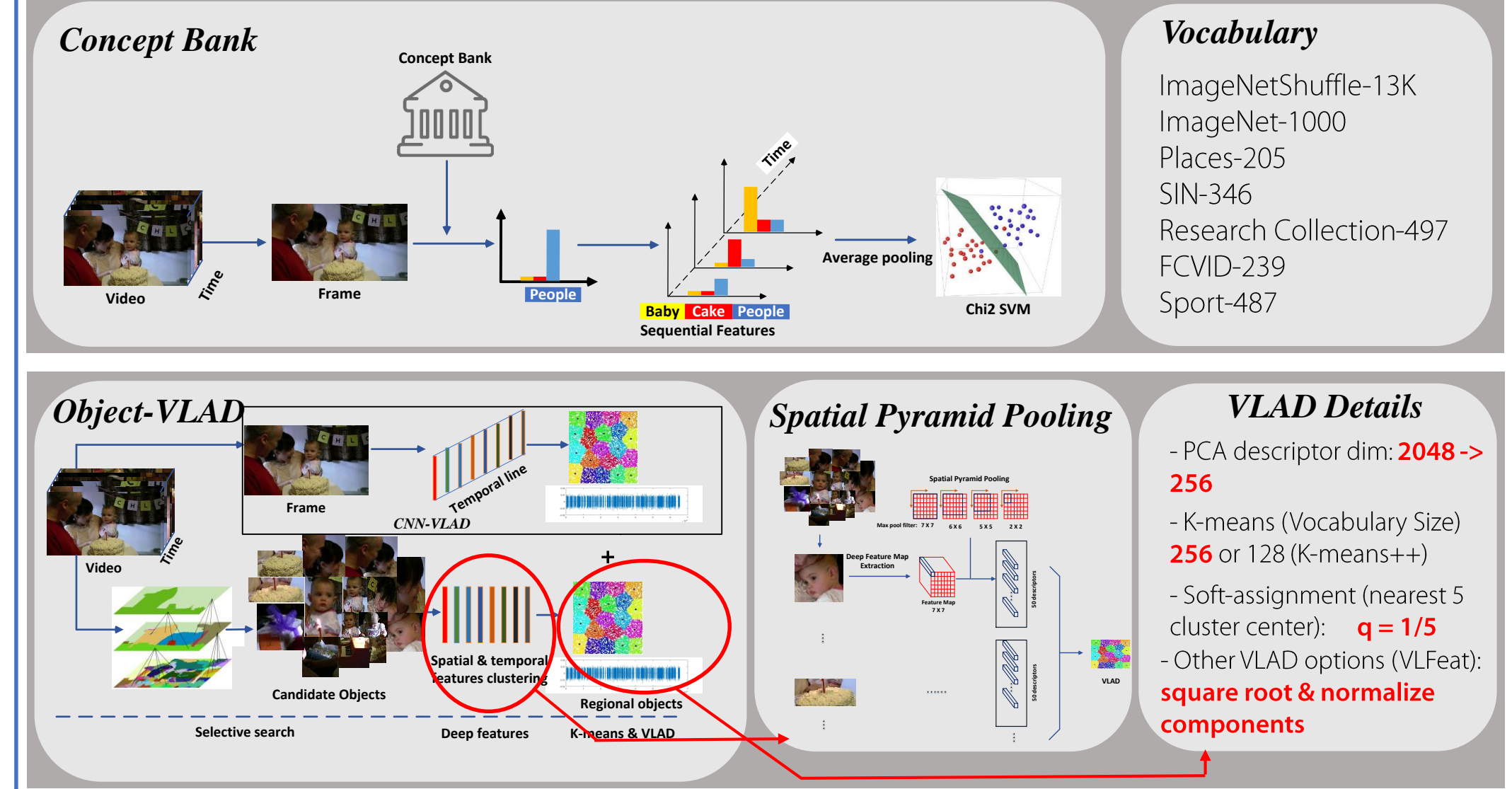


Multimedia Event Detection (MED)

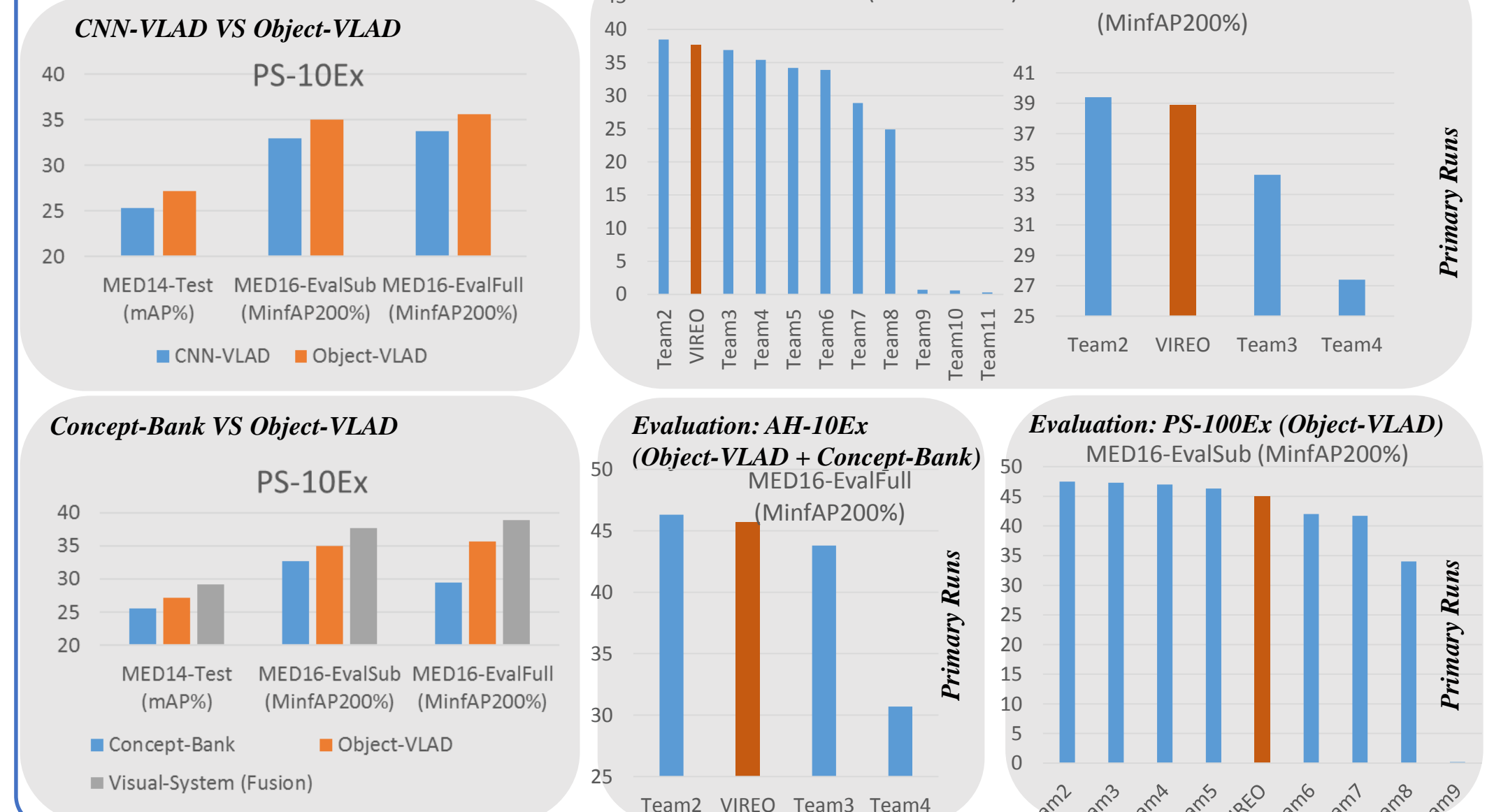
Motivation



Framework



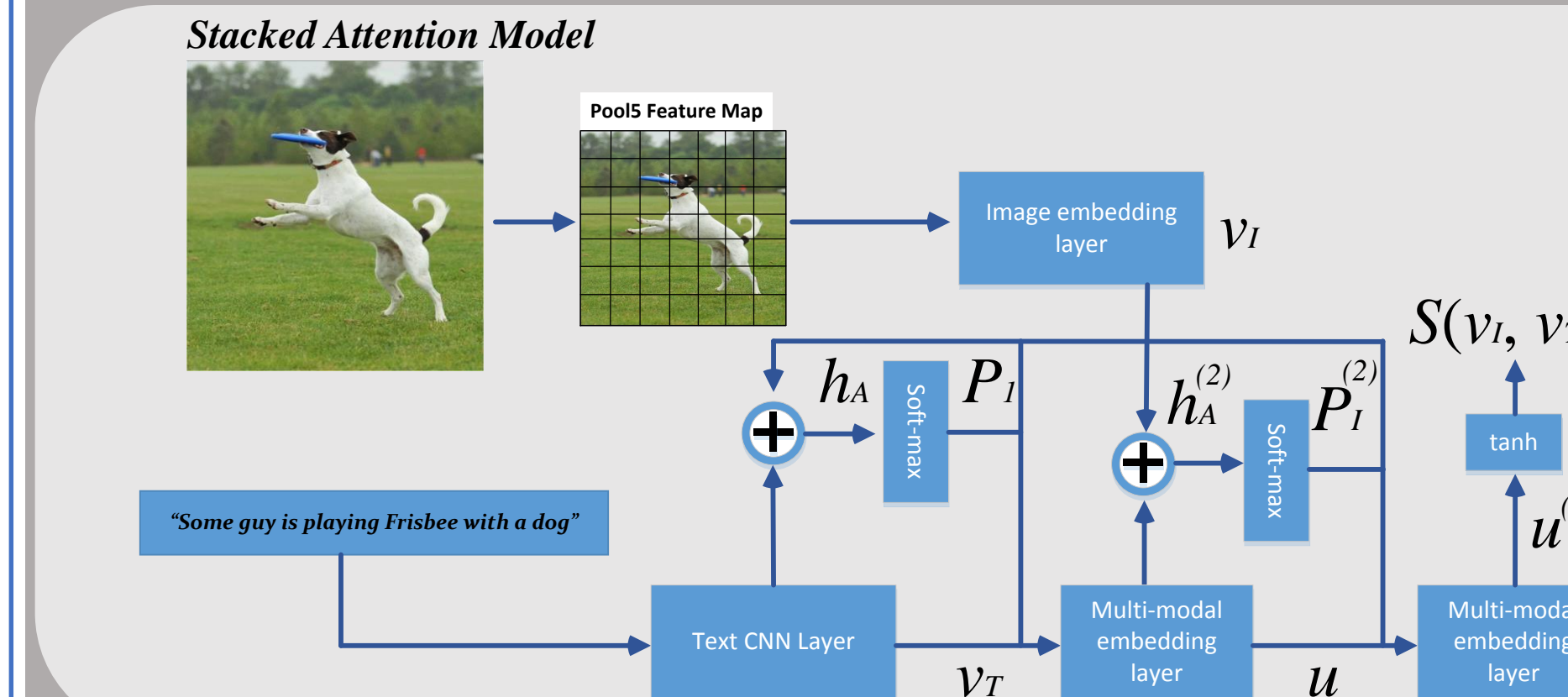
Experiments



Video to Text Description (VTT)

Framework

Concept-based Text Description Matching
 Zero-example MED system in reversed pipeline for concept-based video to text retrieval.
 - Each video is represented as concept vector (max pooling in temporal domain)
 - Query sentences are feed in SQG module for generating the concept-based sentence representation (WordNet to find synonyms).
 - Similarity calculation between video and query representation



f_i : feature for each region

1st Attention Layer

$$v_i = \tanh(W_i f_i + b_i)$$

$$h_A = \tanh(W_{i,A} v_i \oplus (W_{R,A} v_T + b_A))$$

$$P_i = \text{softmax}(W_P h_A + b_P)$$

$$\tilde{v}_i = \sum_1^m p_i v_i$$

$$u = \tilde{v}_i + v_T$$

2nd Attention Layer

$$h_A^{(2)} = \tanh(W_{i,A}^{(2)} v_i \oplus (W_{R,A}^{(2)} u + b_A^{(2)}))$$

$$P_i^{(2)} = \text{softmax}(W_P^{(2)} h_A^{(2)} + b_P^{(2)})$$

$$\tilde{v}_i^{(2)} = \sum_1^m p_i^{(2)} v_i$$

$$u^{(2)} = \tilde{v}_i^{(2)} + u$$

Objective-Function

$$S < v_i, v_T > = \tanh(W_{u,S} u^{(2)} + b_S)$$

$$l(W, D_{trn}) = \sum_{(v_i, v_T^+, v_T^-) \in D_{trn}} \max(0, \delta + S < v_i, v_T^- > - S < v_i, v_T^+ >)$$

Training Dataset
 Flickr30K
 Coco
 TGIF

Deep Features
 - Visual Content: VGG Feature
 Motion: C3D Feature
 - Sentence: Textual CNN

Experiments

