

Joint Audio-Visual Bi-Modal Codewords for Video Event Detection

Guangnan Ye¹, I-Hong Jhuo², Dong Liu¹, Yu-Gang Jiang³, D. T. Lee^{2,4}, Shih-Fu Chang¹

¹Dept. of Electrical Engineering, Columbia University, New York

²Dept. of Computer Science and Information Engineering, National Taiwan University, Taiwan

³School of Computer Science, Fudan University, Shanghai

⁴Dept. of Computer Science and Engineering, National Chung Hsing University, Taiwan



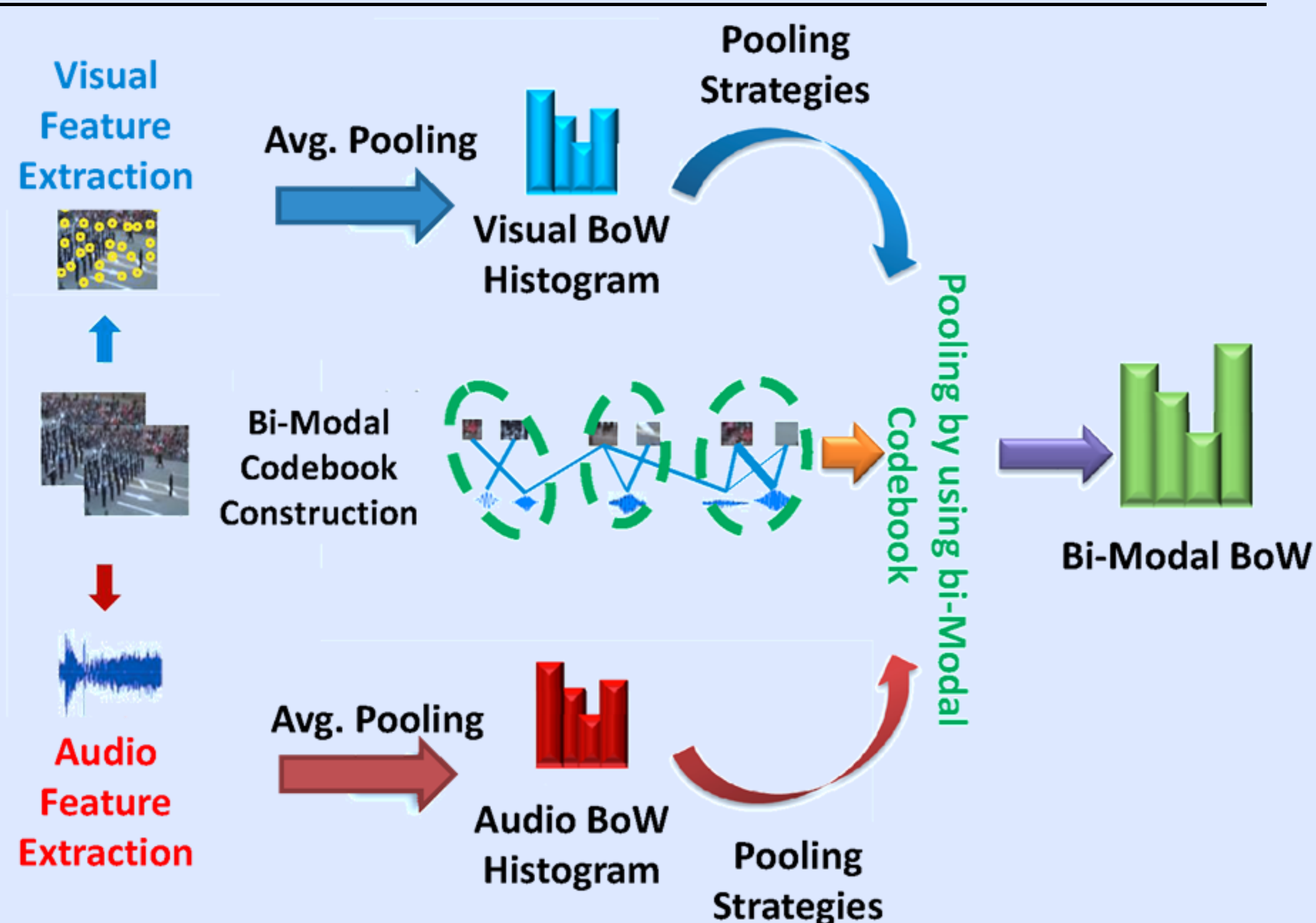
Objective and Overview

Objective: Develop a joint audio-visual bi-modal representation to discover strong audio-visual joint patterns in videos for detecting multimedia events.

- Build a bipartite graph to model relations across the quantized words extracted from the visual and audio modalities;
- Partition the bipartite graph to construct a bi-modal codebook that reveal joint audio-visual patterns;
- Various pooling strategies are employed to re-quantize the visual and audio words into the bi-modal words;
- Bi-modal bag-of-words (BoW) representations are used for event classification.

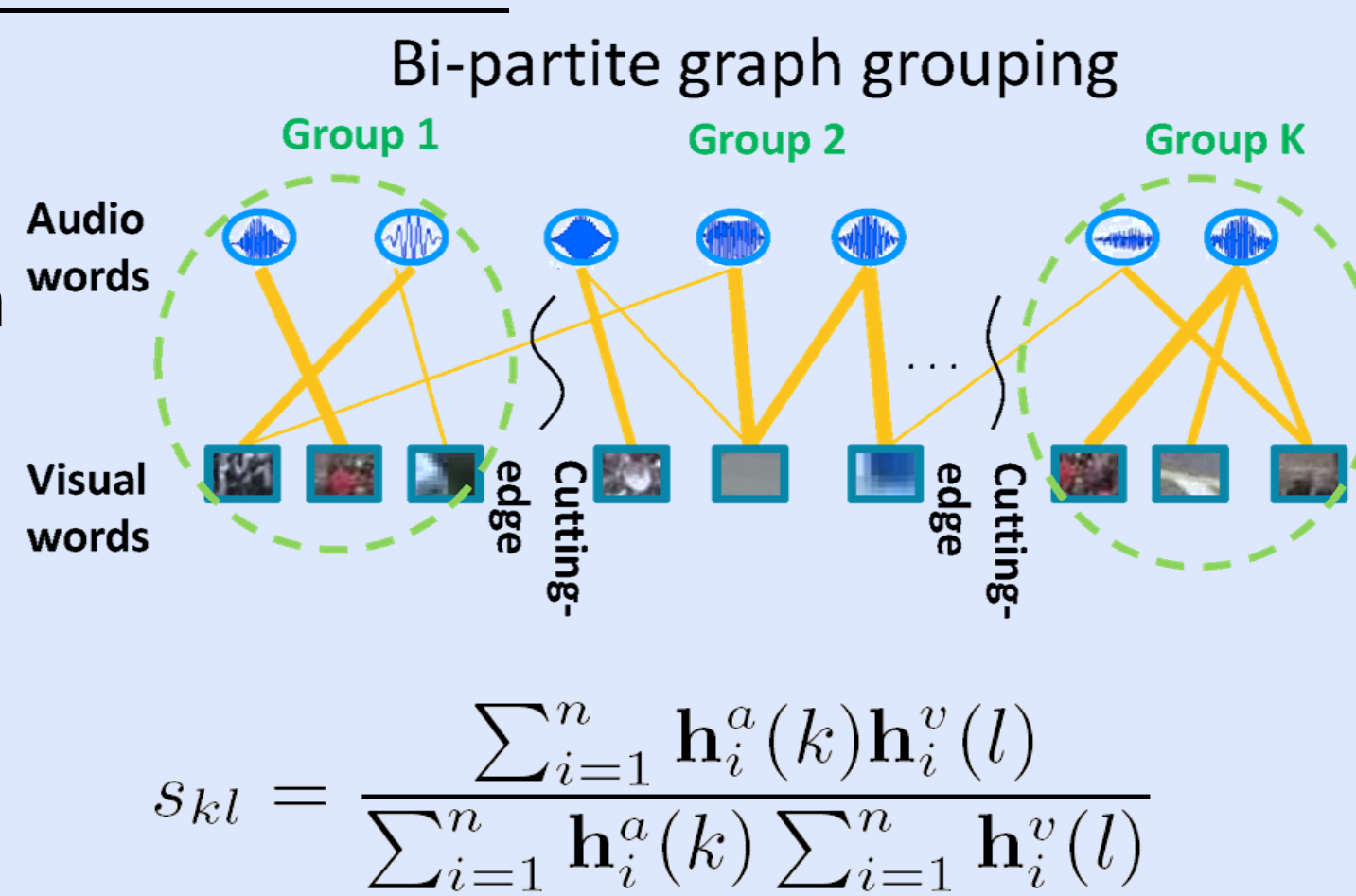
The Proposed Approach and Experiments

• Audio-Visual Bi-Modal BoW Generation



• Bipartite Graph Construction

- Nodes are audio words (h_i^a) and Visual words (h_i^v).
- Edges denote the correlation of audio and visual words.
- Edge weight (line width) is measured by co-occurrence of audio and visual words, defined by s_{kl} :



• The Algorithm

Algorithm 1 Audio-Visual Bi-Modal BoW Representation Generation Procedure

- 1: **Input:** Training video collection $\mathcal{D} = \{d_i\}$ where each d_i is represented as a multi-modality representation $d = \{h_i^a, h_i^v\}$; Size of the audio-visual bi-modal codebook K .
- 2: Produce the correlation matrix \mathbf{S} between the audio and visual words by calculating the co-occurrence probability over \mathcal{D} by Eq. (1).
- 3: Calculate matrix \mathbf{D}_1 , \mathbf{D}_2 and $\hat{\mathbf{S}}$ respectively.
- 4: Apply SVD on $\hat{\mathbf{S}}$ and select $l = \lceil \log_2 K \rceil$ of its left and right singular vectors $\mathbf{U} = [\mathbf{u}_2, \dots, \mathbf{u}_{l+1}]$ and $\mathbf{V} = [\mathbf{v}_2, \dots, \mathbf{v}_{l+1}]$.
- 5: Calculate $\mathbf{Z} = (\mathbf{D}_1^{-1/2} \mathbf{U}, \mathbf{D}_2^{-1/2} \mathbf{V})^\top$.
- 6: Apply k-means clustering algorithm on \mathbf{Z} to obtain K clusters, which form the audio-visual words $\mathcal{B} = \{B_1, \dots, B_K\}$.
- 7: Apply a suitable pooling strategy to re-quantize each video into the audio-visual bi-modal BoW representation.
- 8: **Output:** Audio-visual BoW representation.

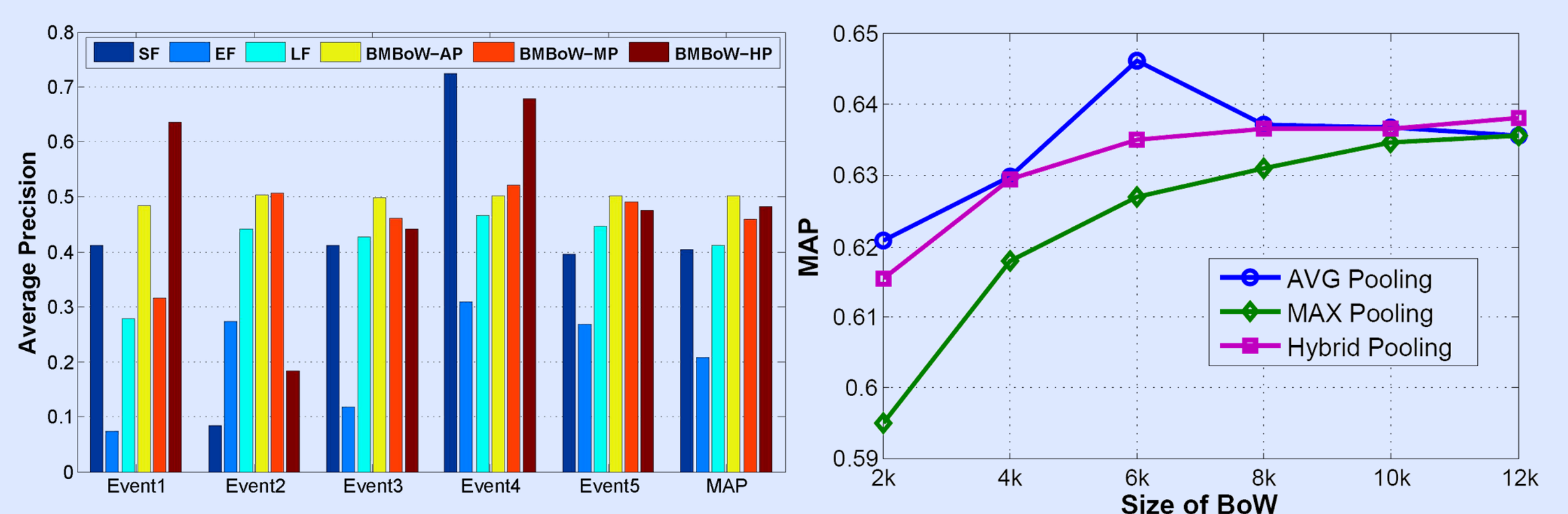
• Pooling Strategies

$$\text{Average Pooling: } h_i^{\text{avg}}(k) = \frac{\sum_{w_p^a \in \mathcal{W}_k^a, w_q^v \in \mathcal{W}_k^v} (h_i^a(p) + h_i^v(q))}{|\mathcal{W}_k^a| + |\mathcal{W}_k^v|}$$

$$\text{Max Pooling: } h_i^{\text{max}}(k) = \max \left(\sum_{w_p^a \in \mathcal{W}_k^a} h_i^a(p), \sum_{w_q^v \in \mathcal{W}_k^v} h_i^v(q) \right)$$

$$\text{Hybrid Pooling: } h_i^{\text{hyb}}(k) = \frac{1}{2} \left(\max_{w_p^a \in \mathcal{W}_k^a} h_i^a(p) + \frac{\sum_{w_q^v \in \mathcal{W}_k^v} h_i^v(q)}{|\mathcal{W}_k^v|} \right)$$

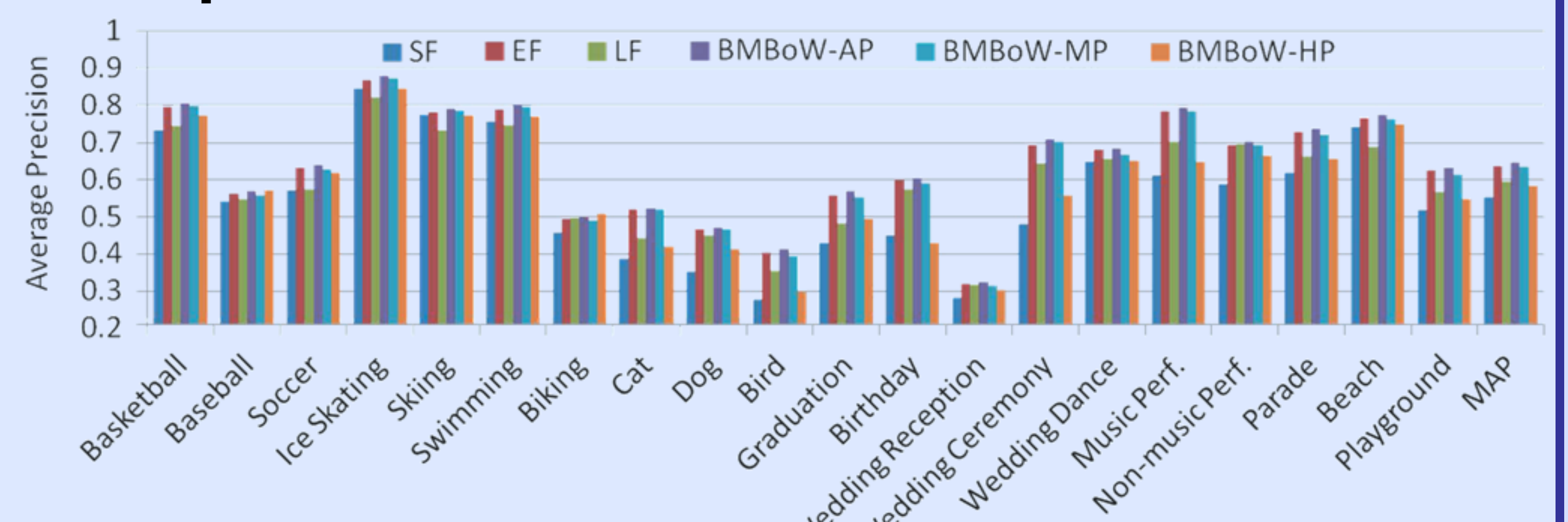
• Experiment on TRECVID MED 2011



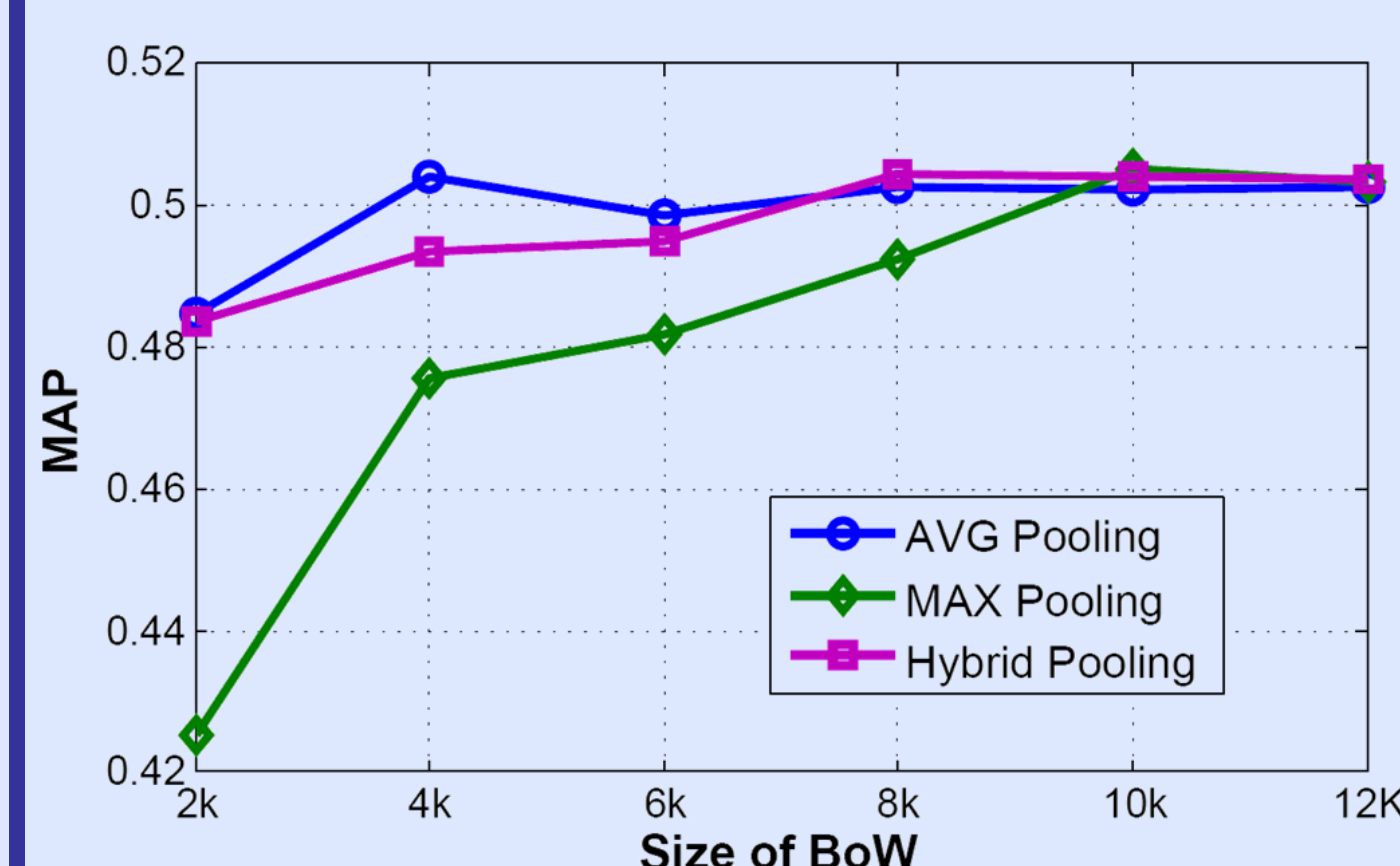
Per-event AP performance, 19.6% gain over LF baseline

Effect of varying bi-modal codebook size; average pooling performs the best

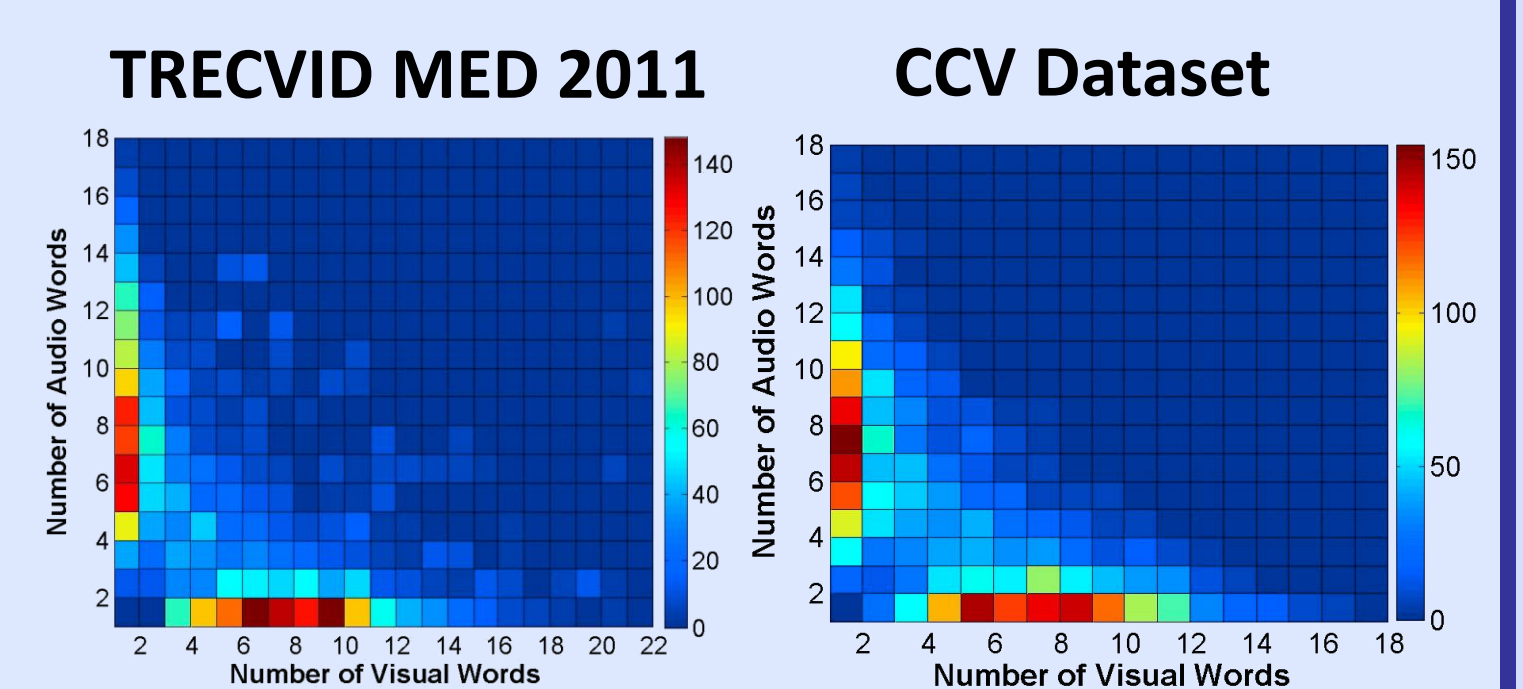
• Experiment on CCV Dataset



Per-event AP performance; 8.6% gain over LF baseline



Effect of varying bi-modal codebook size; average pooling performs the best



Density of audio and visual words within the bi-modal words: 47% and 36% of bi-modal words contain both contributions from audio and visual in TRECVID and CCV dataset respectively.

Summary

- Joint bi-modal codewords achieved 19.6% and 8.6% improvement over LF baseline in TRECVID and CCV, respectively.
- 47% and 36% of bi-modal codewords contain contributions from both modalities in TRECVID and CCV, respectively.
- Among the evaluated pooling strategies, average pooling achieved the best performance.
- Events with multimodal cues, such as “Bird” and “Wedding Ceremony”, show the highest gains.