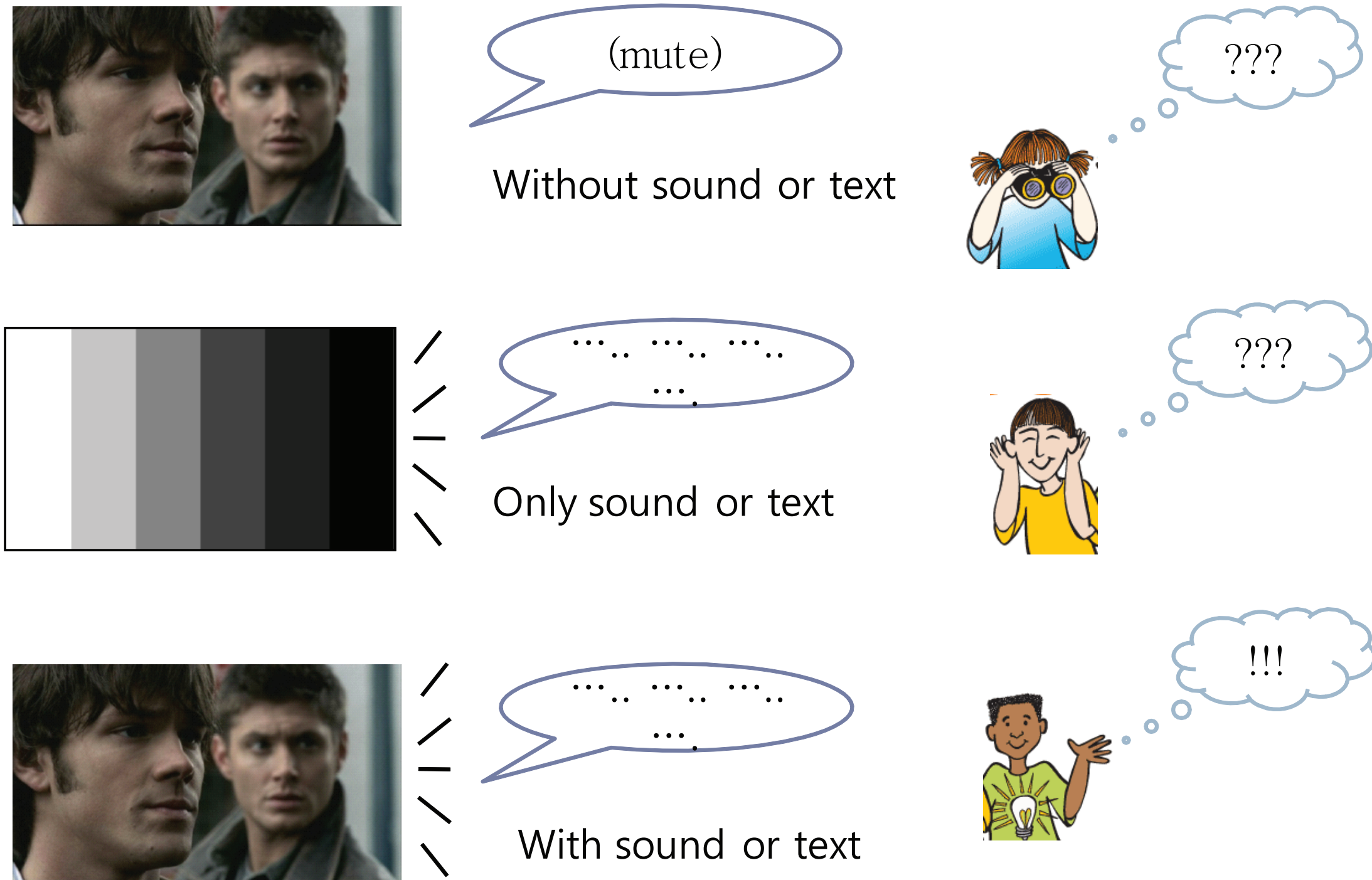


Watch, Listen & Learn: Co-training on Captioned Images and Videos

Sonal Gupta, Joohyun Kim, Kristen Grauman, Raymond Mooney
The University of Texas at Austin
{sonaluta, scimitar, grauman, mooney}@cs.utexas.edu

Introduction



Motivation

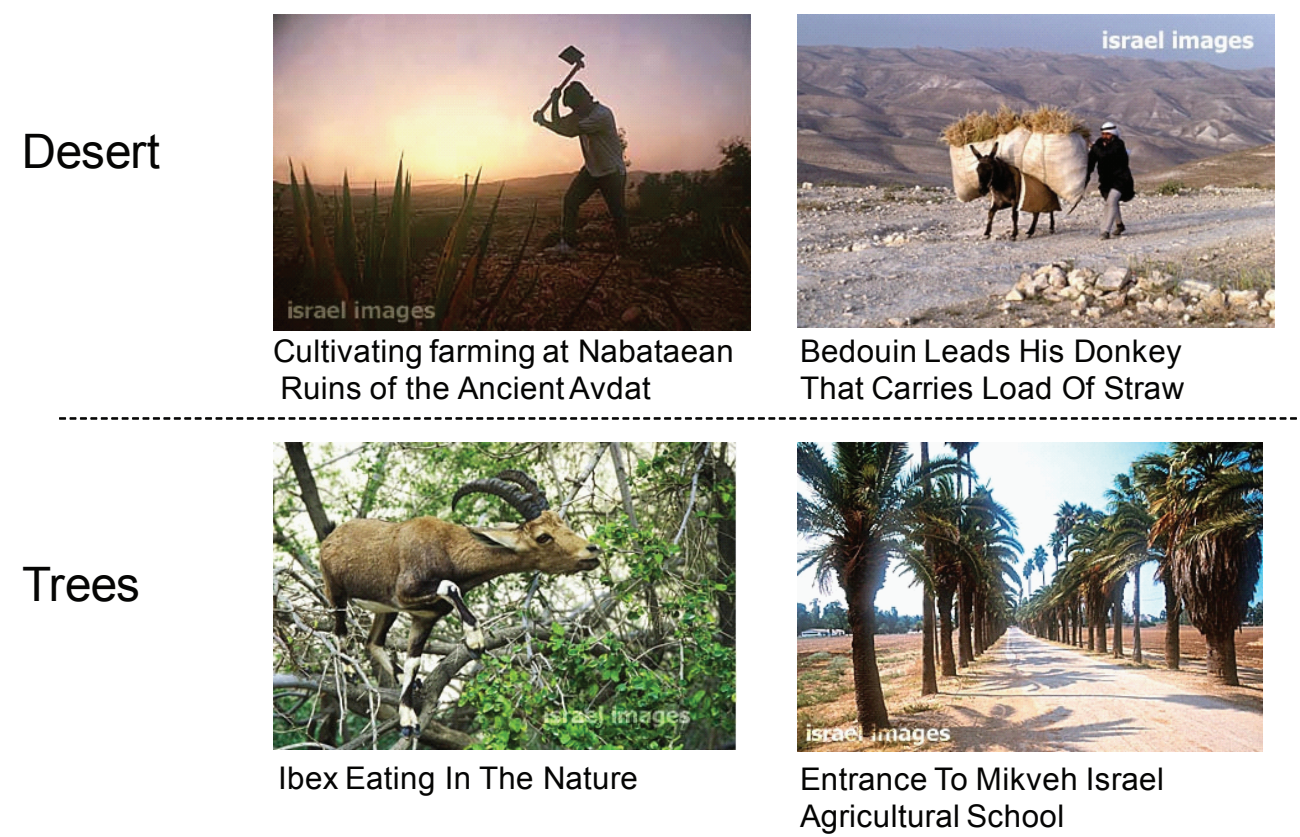
- **Image Recognition & Human Activity Recognition in Videos**
 - Hard to classify, visual cues are ambiguous
 - Expensive to manually label instances
- **Often images and videos have text captions**
 - Leverage multi-modal data
 - Use readily available unlabeled data to improve accuracy

Goals

- Classify images and videos with the help of associated text captions
- Use Co-training to achieve better classification accuracy for image and video classification task

Datasets

- Image : 362 instances with 2 classes



- Video : 221 instances with 4 classes

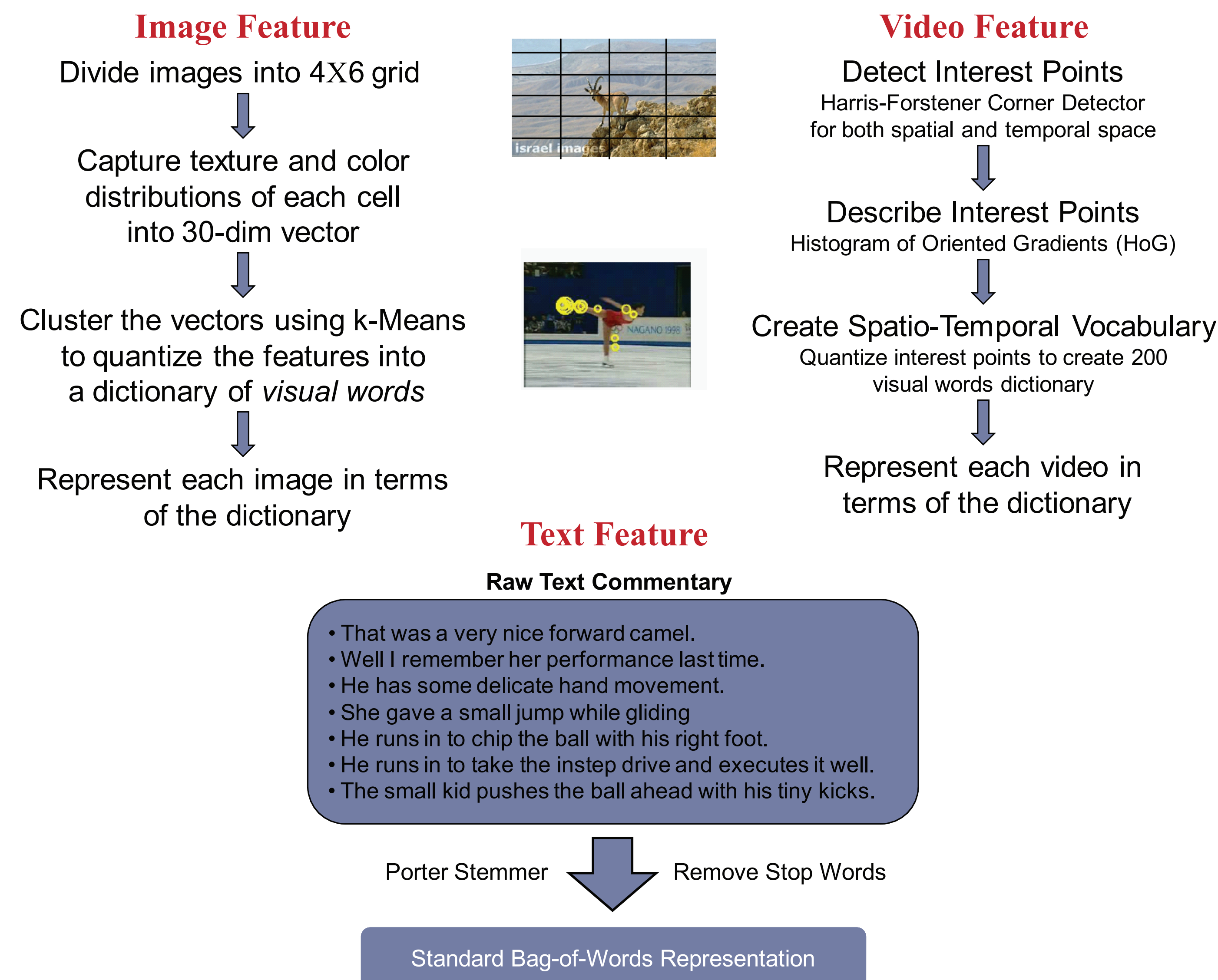


Approach

- Combining two views (Text and Visual) of images and videos using Co-training (Blum and Mitchell '98) learning algorithm

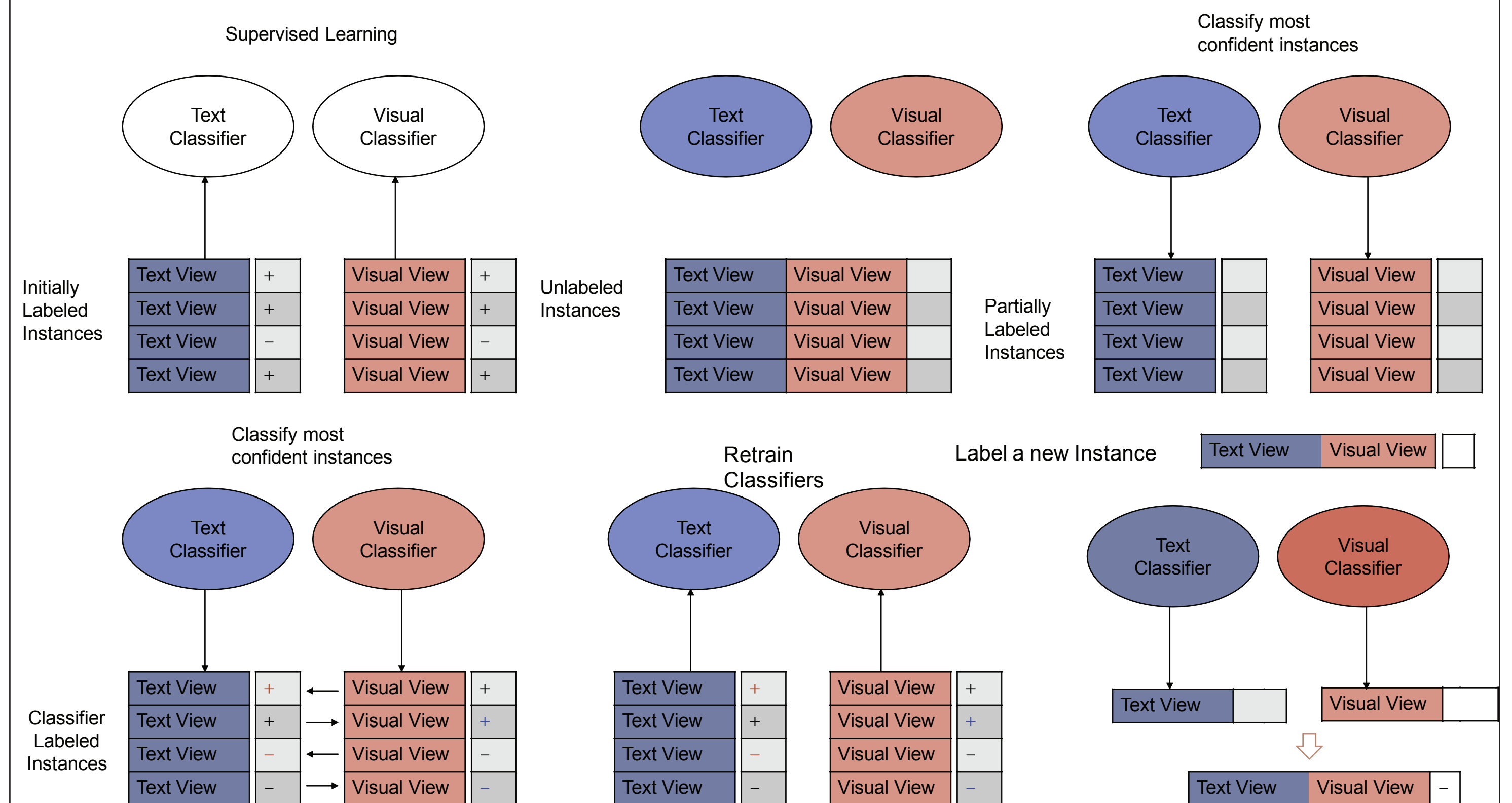
- **Text View**
 - Caption of image or video
 - Readily available
- **Visual View**
 - Color, texture, temporal information in image/video

Feature Extraction



Algorithm

- **Co-training**
 - Semi-supervised learning paradigm that exploits two mutually independent and sufficient views
- **Features of dataset can be divided into two sets:**
 - The instance space: $X = X_1 \times X_2$
 - Each example: $x = (x_1, x_2)$
- **Proven to be effective in several domains**
 - Web page classification (content and hyperlink)
 - E-mail classification (header and body)

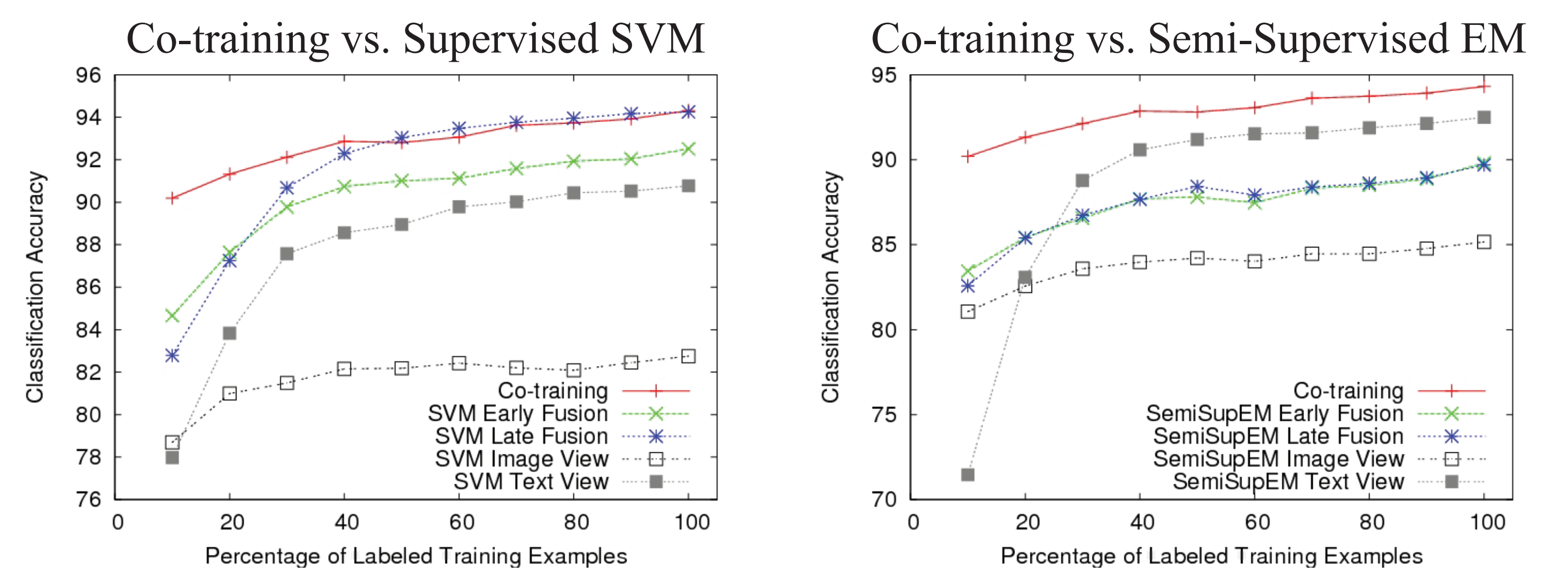


Experimental Results

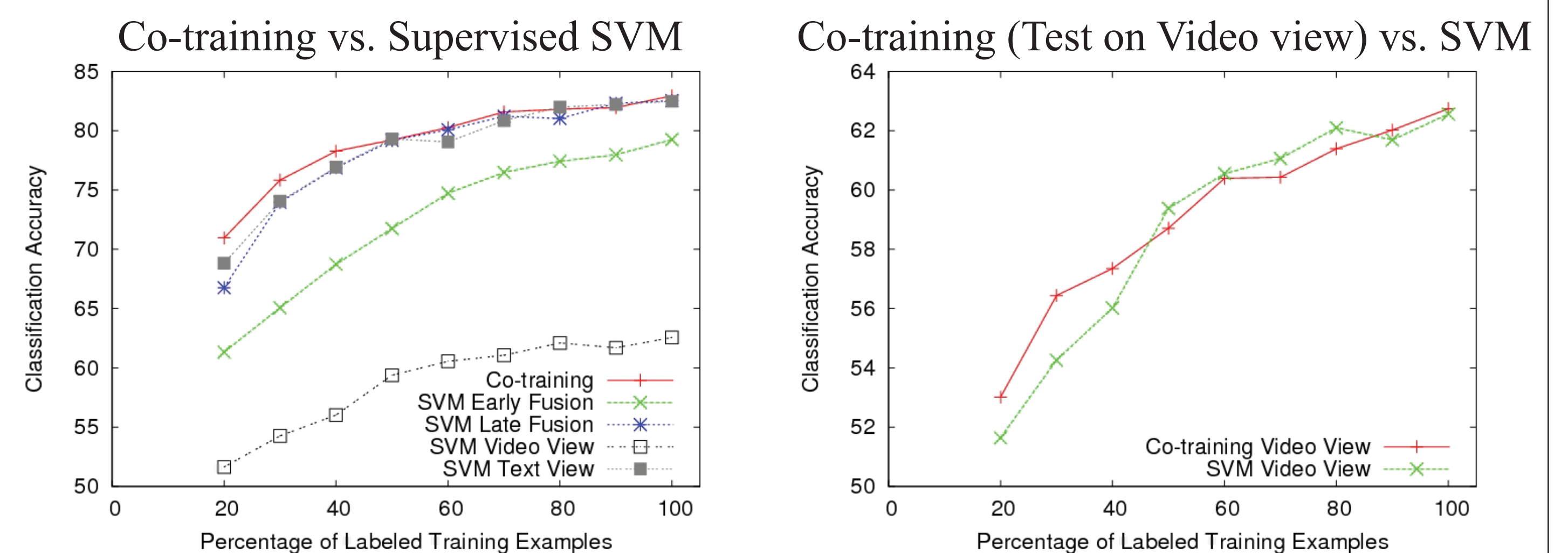
Baselines

- **Uni-modal**
 - Image/Video View : Only image/video features are used
 - Text View : Only textual features are used
- **Multi-modal**
 - Early Fusion : Concatenate visual and textual features and train classifier
 - Late Fusion : Run separate classifiers on each view and concatenate the results

Image Dataset



Video Dataset



Conclusion

- Combining textual and visual features can help improve accuracy
- Co-training can be useful to combine textual and visual features to classify images and videos
- Co-training helps in reducing labeling of images and videos

References

- [1] Bekkerman and Jeon, Multi-modal Clustering for Multimedia Collections. CVPR 2007
- [2] Blum and Mitchell, Combining labeled and unlabeled data with co-training, COLT 1998
- [3] Laptev, On space-time interest points, IJCV 2005
- [4] Weka Data Mining Tool (Witten and Frank)