

# Multimedia Event Detection using GS-SVMs and Audio-HMMs

Nakamasa Inoue, Yusuke Kamishima,  
Koichi Shinoda,  
*Department of Computer Science,  
Tokyo Institute of Technology*

Shunsuke Sato  
*Canon Inc.*

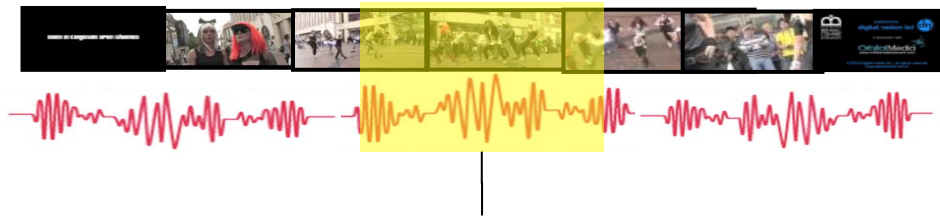
# Outline

- Motivation
- System Overview
- Method
  - Features extraction
  - GS-SVM
  - Audio HMMs
- Results
  - Best result: Minimum NDC = 0.525

# Motivation

- Two event feature categories:
  - Features that appear **in every frame**
  - Features that appear **only in some frames**
- Their combination can improve the detection performance.

ex.) Flash Mob Gathering clips



Some frames:

- Dancing
- Dance music
- Cheering voice

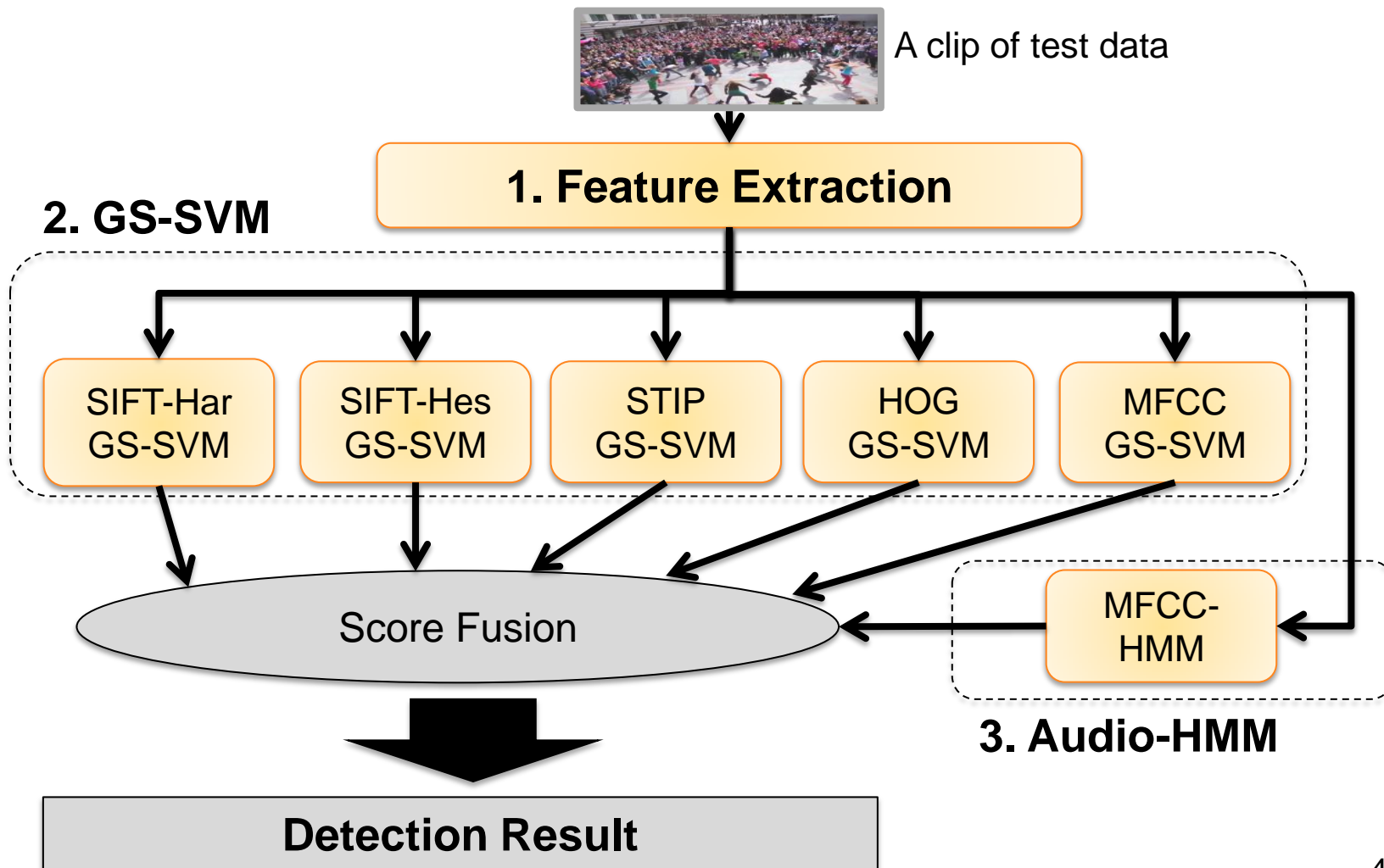
Every frame:

- Outdoor
- Dancers
- Road
- Crowd
- Crowd buzz
- ...

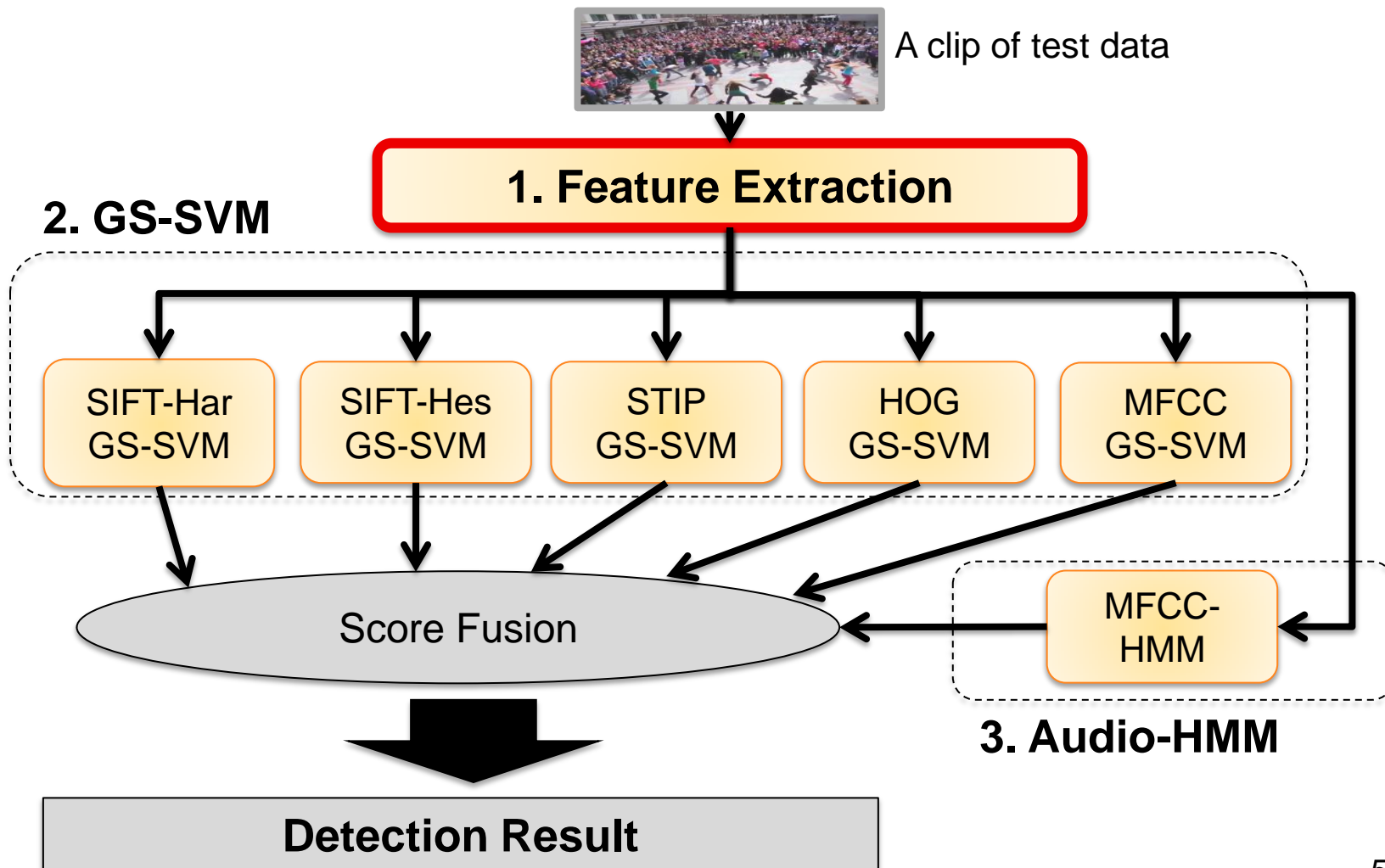
# Method Overview

- For every-frame features: **GS-SVM**  
(**GMM-Supervector Support Vector Machine**)
  - Use several visual and audio features
  - Soft clustering - robust against quantization errors
  - Based on our system of TRECVID 2010 SIN task
  
- For some-frame features: **HMM**  
(**Hidden Markov model**)
  - Model temporal features in sound
  - Apply word-spotting in speech recognition
  - Use only audio, not video

# System Overview

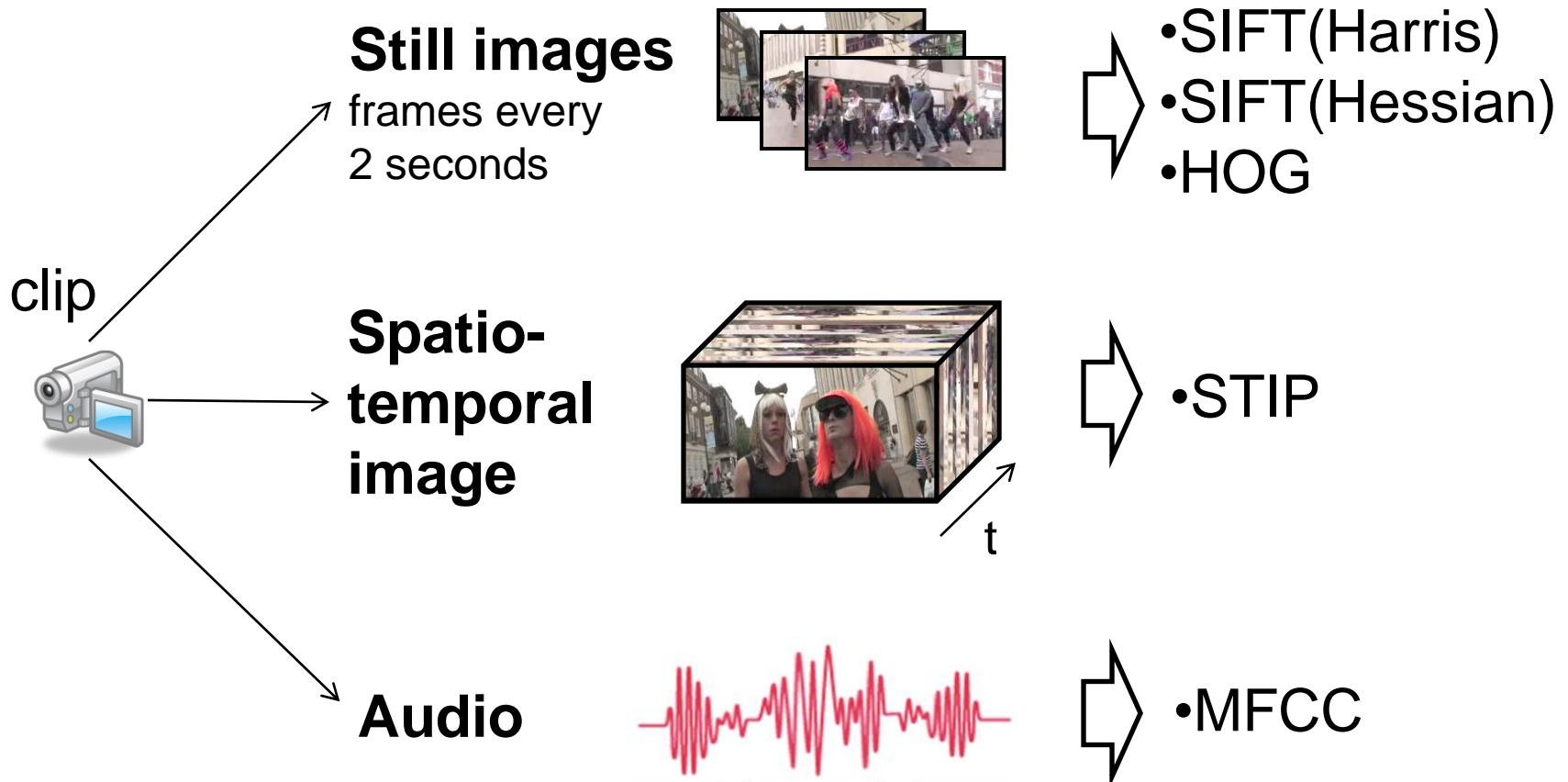


# System Overview



# Feature Extraction

- 5 types of features, from 3 kinds of sources

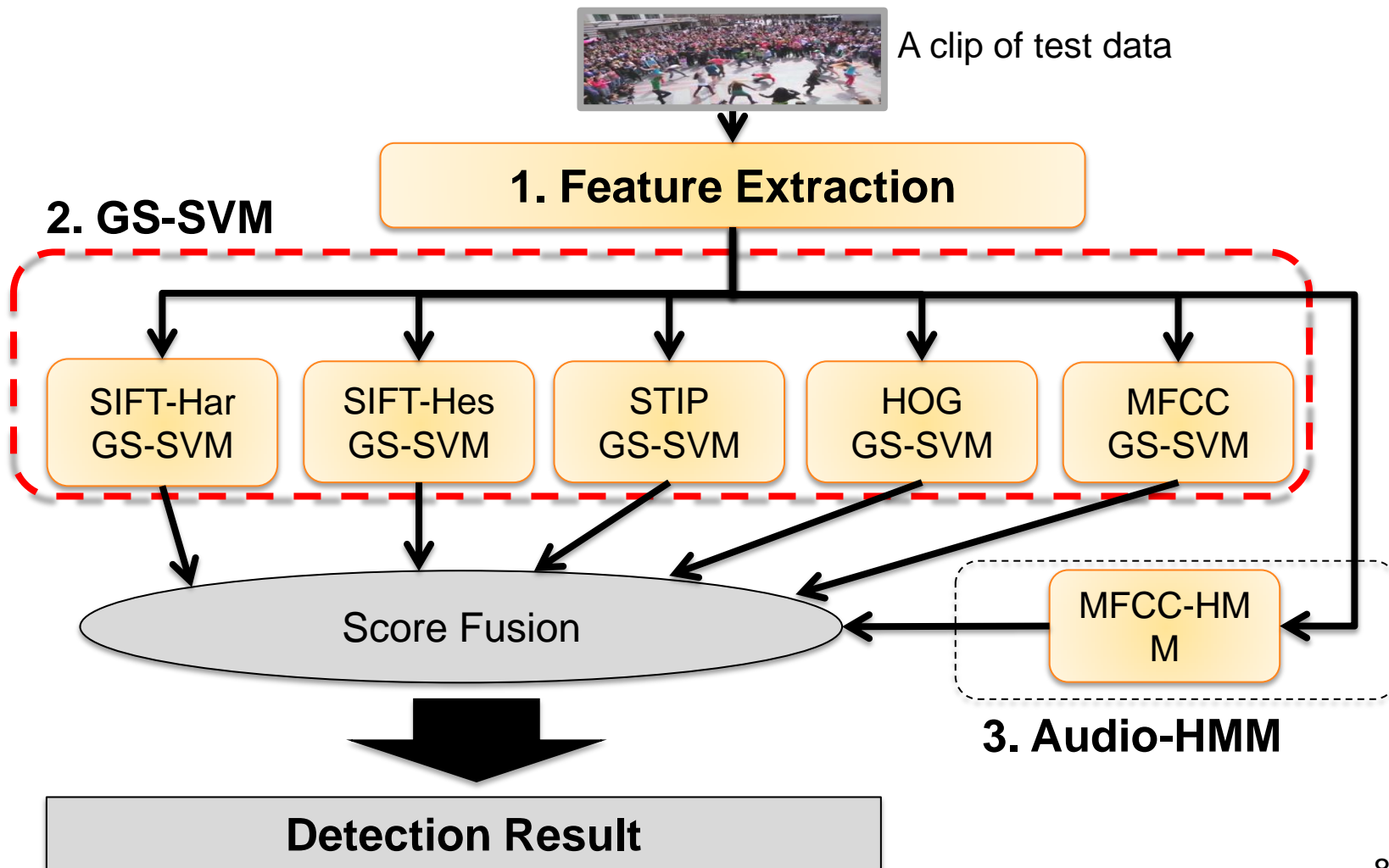


# List of Features

source	feature	description
Still images	<b>SIFT (Harris)</b>	Scale-Invariant Feature Transform with Harris-affine regions and Hessian-affine regions [Mikolajczyk, 2004]
	<b>SIFT (Hessian)</b>	
	<b>HOG</b>	32 dimensional HOG Dense sampling (every 4 pixels)
Spatio-temporal images	<b>STIP</b>	Space-Time Interest Points HOG and HOF features extracted [Laptev, 2005]
Audio	<b>MFCC</b>	Mel-frequency cepstral coefficients Audio features for speech recognition

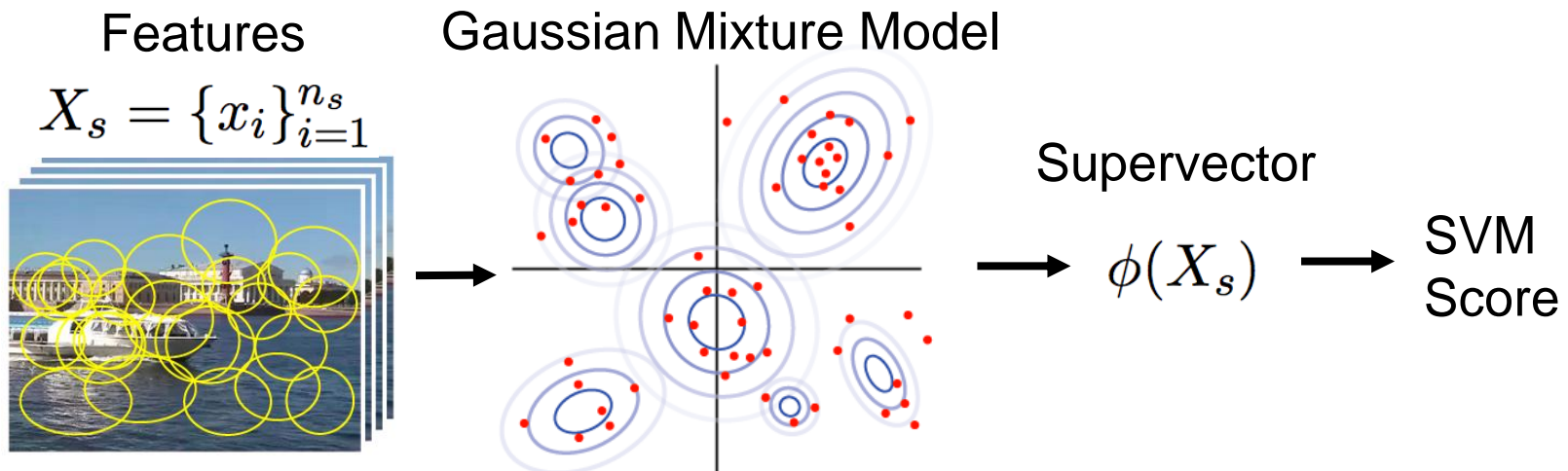


# System Overview



# GMM Supervector SVM (GS-SVM)

- Represent **the distribution of each feature**
  - Each clip is modeled by a **GMM (Gaussian Mixture Model)**
  - Derive **a supervector** from the GMM parameters
  - Train **SVM (Support Vector Machine)** of the supervectors

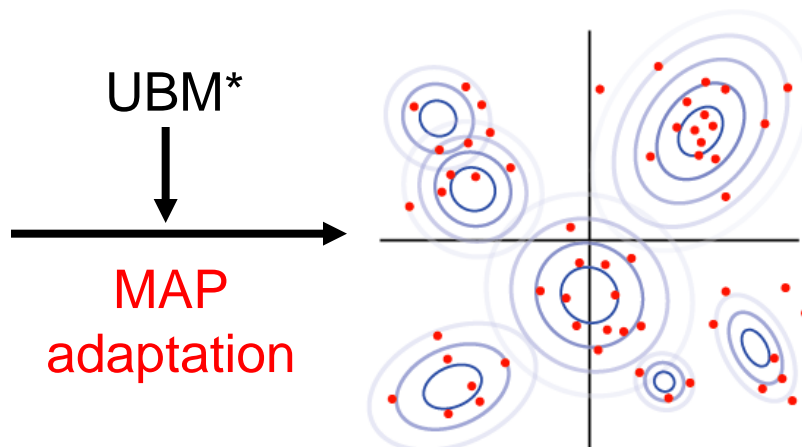
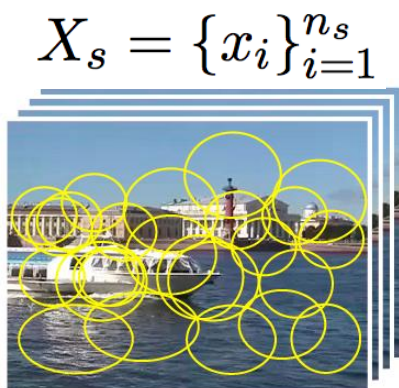


# GMM Estimation

- Estimated by using **maximum a posteriori (MAP) adaptation** for mean vectors:

$$\hat{\mu}_k^{(s)} = \frac{\tau \mu_k^{(U)} + \sum_{i=1}^{n_s} c_{ik} x_i}{\tau + C_k} \quad \left[ \begin{array}{l} \text{where} \\ c_{ik} = \frac{w_k \mathcal{N}(x_i | \mu_k^{(U)}, \Sigma_k^{(U)})}{\sum_{k=1}^K w_k \mathcal{N}(x_i | \mu_k^{(U)}, \Sigma_k^{(U)})}, \quad C_k = \sum_{i=1}^{n_s} c_{ik} \end{array} \right]$$

adapted mean      UBM's mean



\*Universal background model (UBM): a prior GMM which is estimated by using all video data.

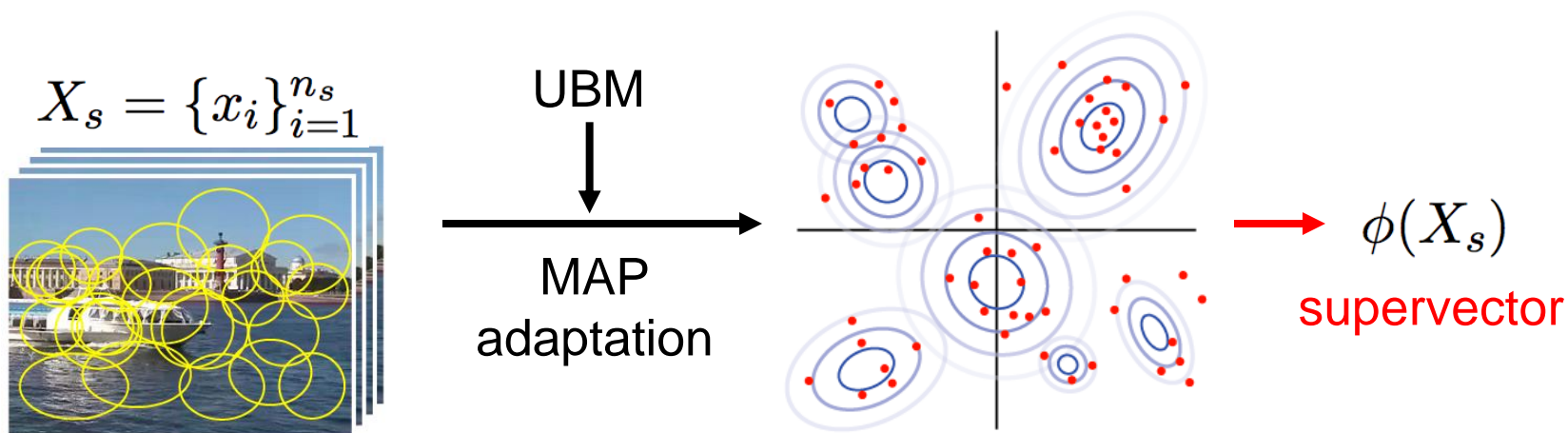
# GMM Supervector

- **GMM Supervector**: combination of the mean vectors.

$$\phi(X_s) = \begin{pmatrix} \tilde{\mu}_1^{(s)} \\ \tilde{\mu}_2^{(s)} \\ \vdots \\ \tilde{\mu}_K^{(s)} \end{pmatrix}$$

where

$$\tilde{\mu}_k^{(s)} = \frac{\sqrt{w_k^{(U)}} (\Sigma_k^{(U)})^{-\frac{1}{2}} \hat{\mu}_k^{(s)}}{\text{normalized mean}}$$



# Score Fusion in GS-SVM

- GS-SVMs use RBF-kernels:

$$k(X_F, X'_F) = \exp(-\gamma \|\phi(X_F) - \phi(X'_F)\|_2^2),$$

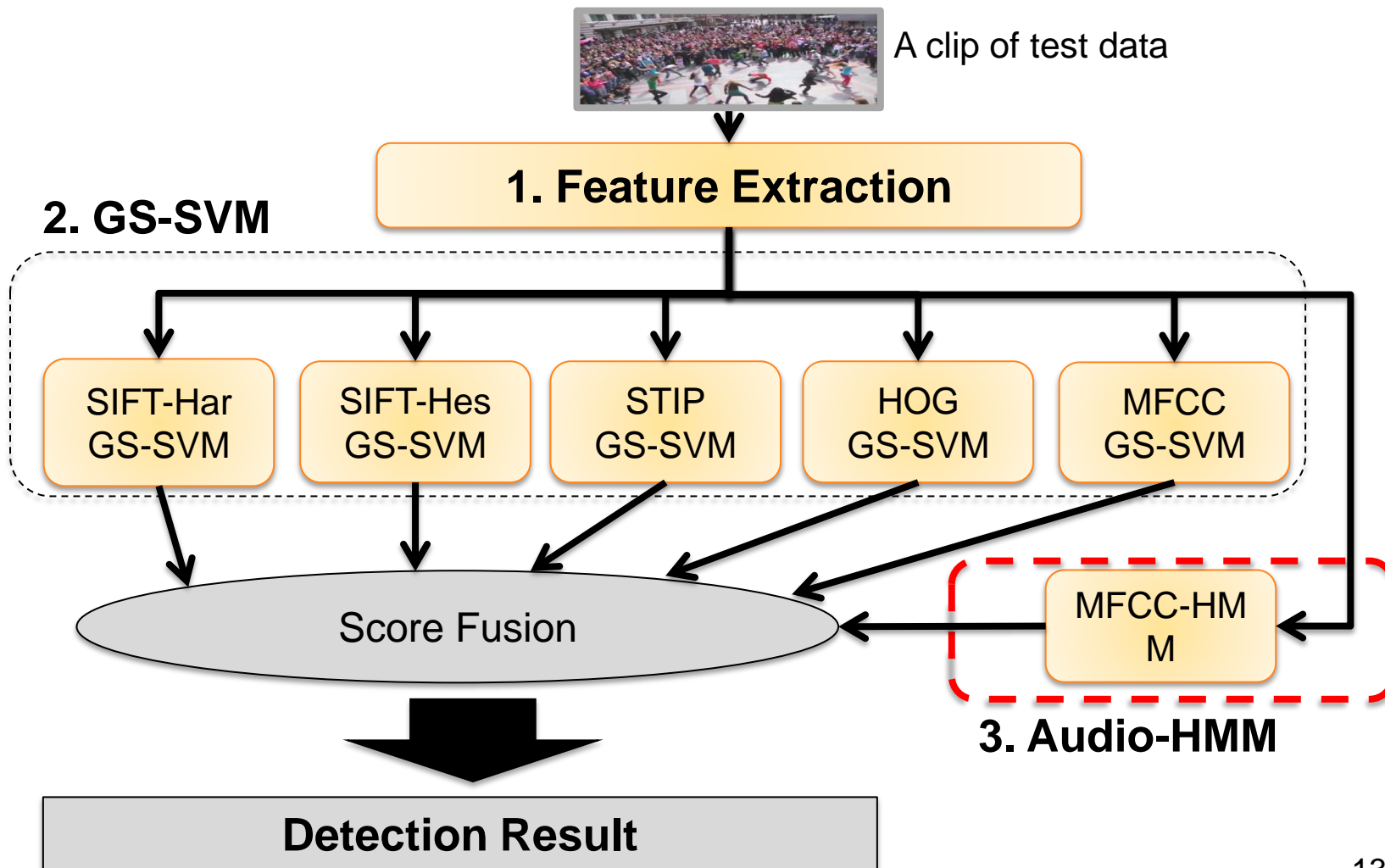
- **Score:** Weighted Average of SVM outputs:

$$f(X) = \sum_{F \in \mathcal{F}} \alpha_F f_F(X_F), \quad 0 \leq \alpha_F \leq 1, \quad \sum_F \alpha_F = 1$$

where  $\mathcal{F} = \{\text{SIFT-Her, SIFT-Hes, HOG, STIP, MFCC}\}$

- $\alpha_F$  are decided by 2-fold cross validation based on
  - Minimum Normalized Detection Cost - Run 1 & Run 2
  - Average Precision - Run 3
  - In Run 4,  $\alpha_F$  is equal for all features

# System Overview



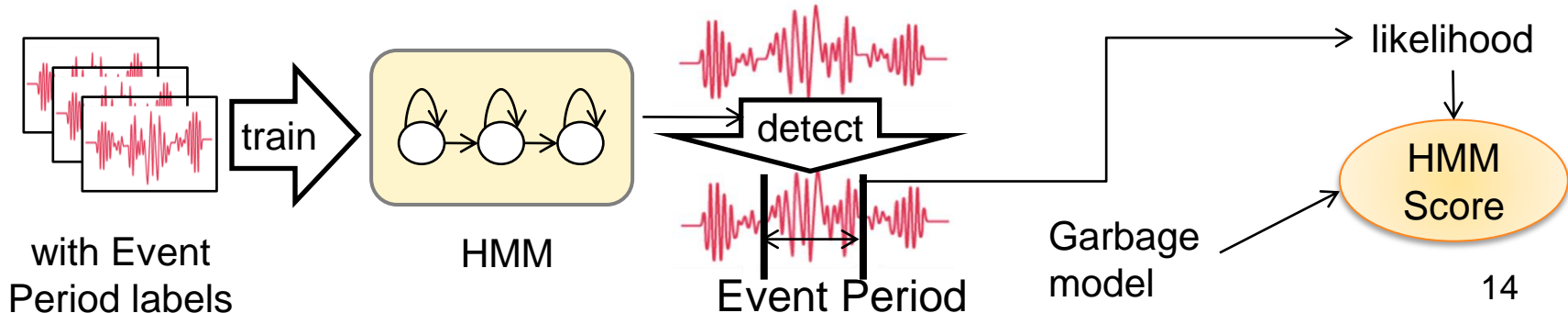
# Audio HMM

Training:

1. Label an **event period** manually for each event clip
2. Train an **event HMM** using MFCC

Test:

1. Find likelihood  $L_E$  of the event period **by word-spotting**
2. Find likelihood  $L_G$  of the event period for a **garbage model** estimated from all video data
3. Calculate likelihood ratio  $L_E/L_G$  as the detection score



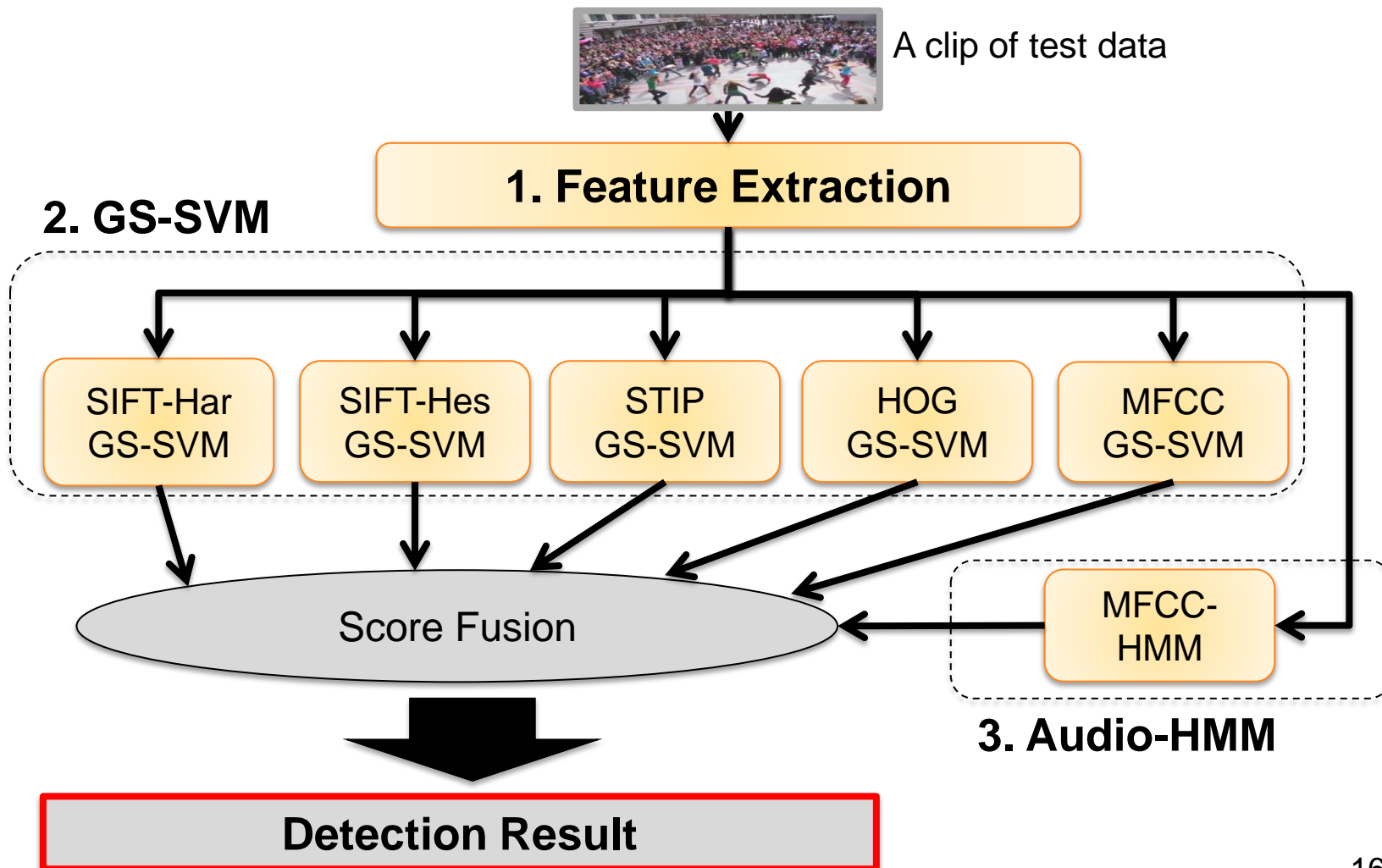
# Preliminary result of Audio HMMs

- Fuse HMM score with GS-SVM by weighted average.
- Audio HMMs are effective in 3 events – Use them in Run1.



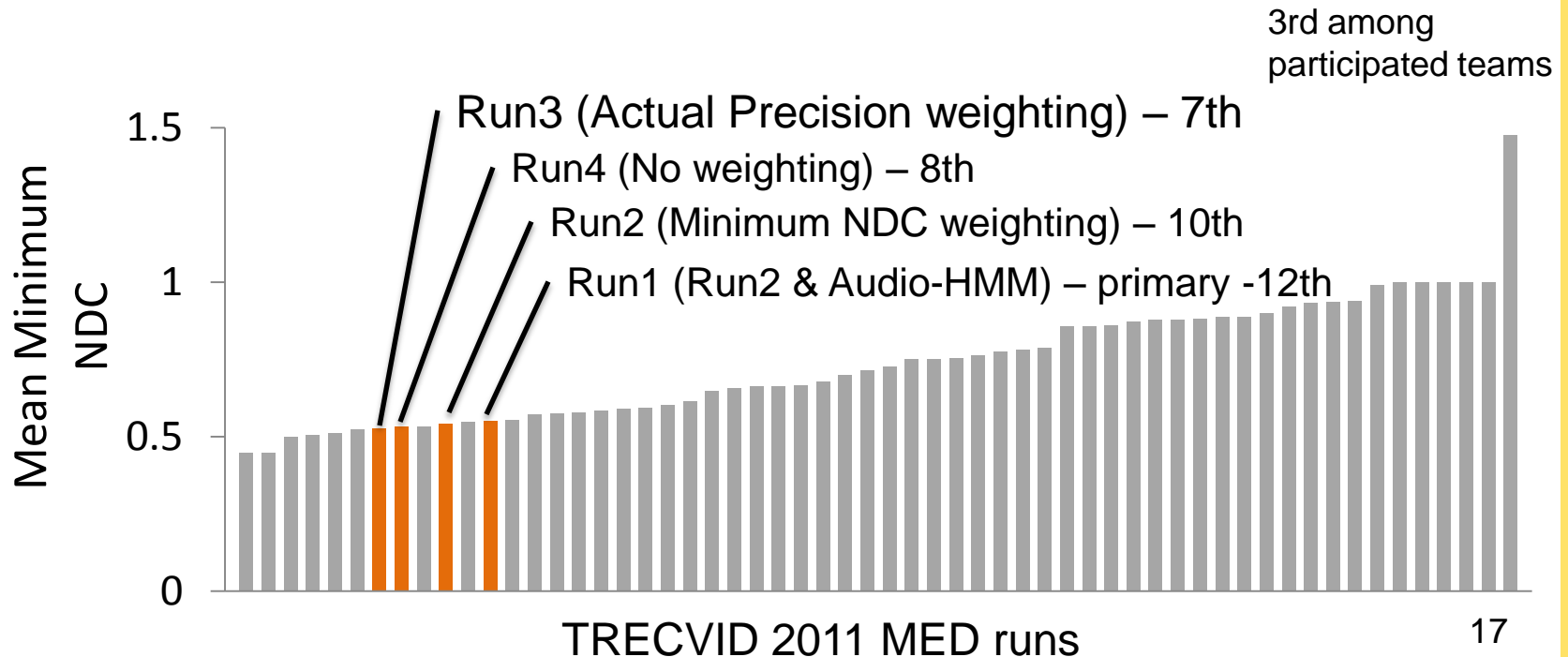


# System Overview



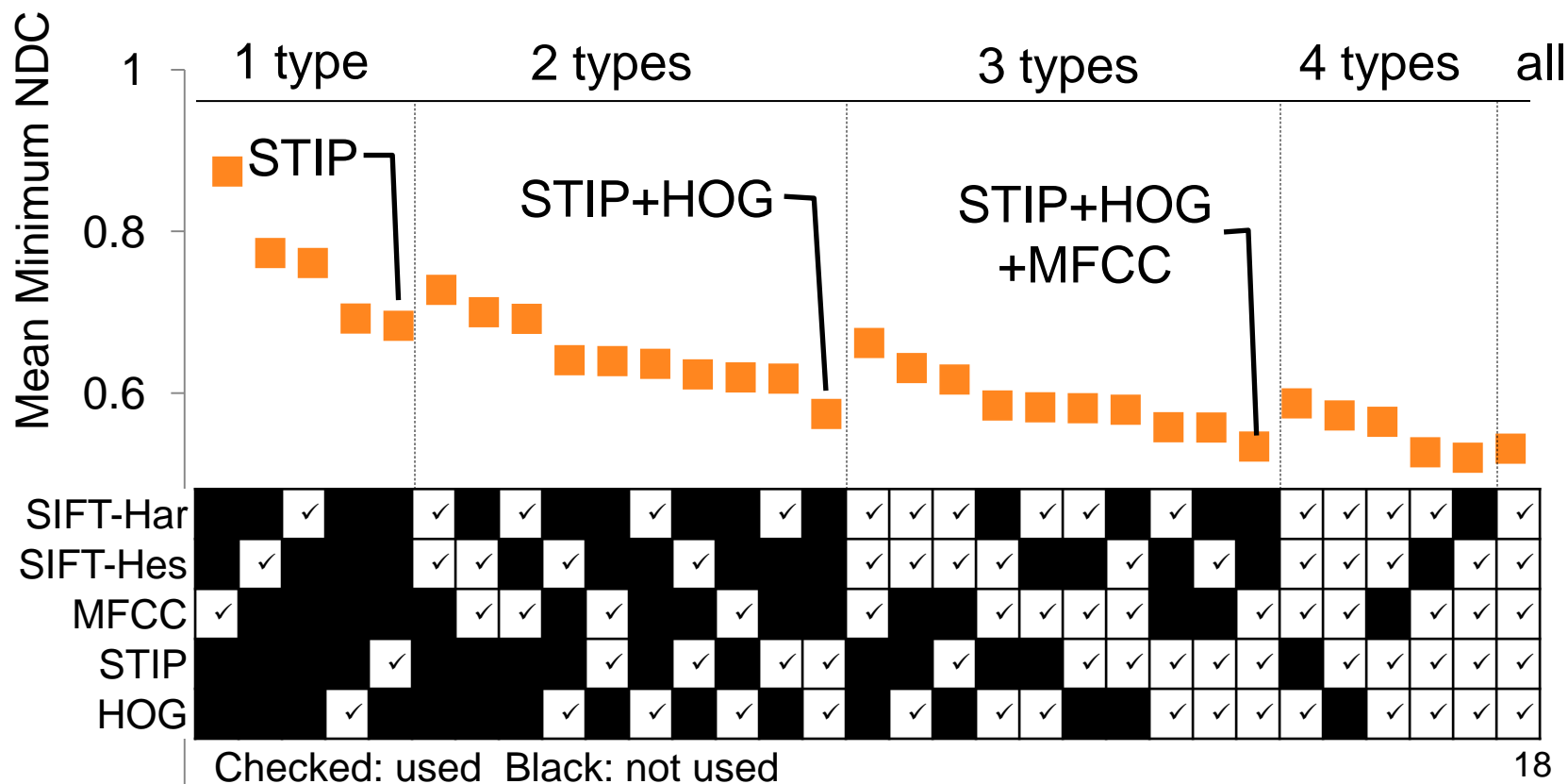
# Experiments

- Run3 was the best. GS-SVM was effective
- Run1 (Audio-HMM) did not show good performance
- Run2, weights decided by Minimum NDC, is not good
  - Simple cross validation may have failed.



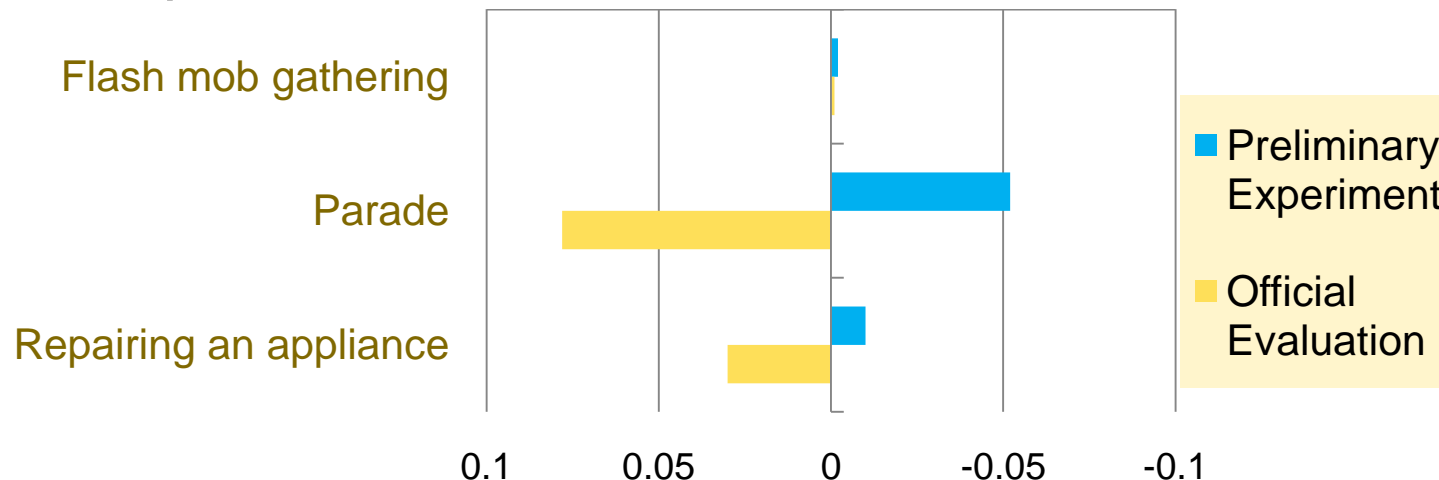
# Effect of each feature in GS-SVM

- STIP and HOG had better performance.
- MFCC was effective when combined with STIP and HOG.



# Why Audio HMM did not work?

- It failed to capture **temporal features**
  - Each state represents a specific sound such as drum, cheering, which may appear in non-event and/or at random.
- Test data include many sounds **not appear** in training and development data



**Difference of Minimum NDC between  
with and without Audio HMMs**

# Conclusion

- We combine GS-SVM and Audio HMM
- GS-SVMs are effective for MED.
  - STIP, HOG, and MFCC are important
- Audio HMMs are not effective
  - It cannot capture temporal features
  - Variety of sounds are larger than expected
- Future works
  - Include other features, such as Dense SIFT
  - Improve the HMM-based sound detection
    - Model event subclasses and their relationship