

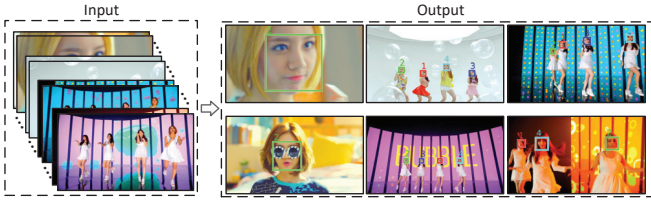
# Tracking Persons-of-Interest via Adaptive Discriminative Features

Shun Zhang<sup>1</sup> Yihong Gong<sup>1</sup> Jia-Bin Huang<sup>2</sup> Jongwoo Lim<sup>3</sup> Jinjun Wang<sup>1</sup> Narendra Ahuja<sup>2</sup> Ming-Hsuan Yang<sup>4</sup>

<sup>1</sup>Xi'an Jiaotong University <sup>2</sup>University of Illinois, Urbana-Champaign <sup>3</sup>Hanyang University <sup>4</sup>University of California, Merced

Code and data available at <http://bit.ly/multi-face-tracking-eccv2016>

## Problem

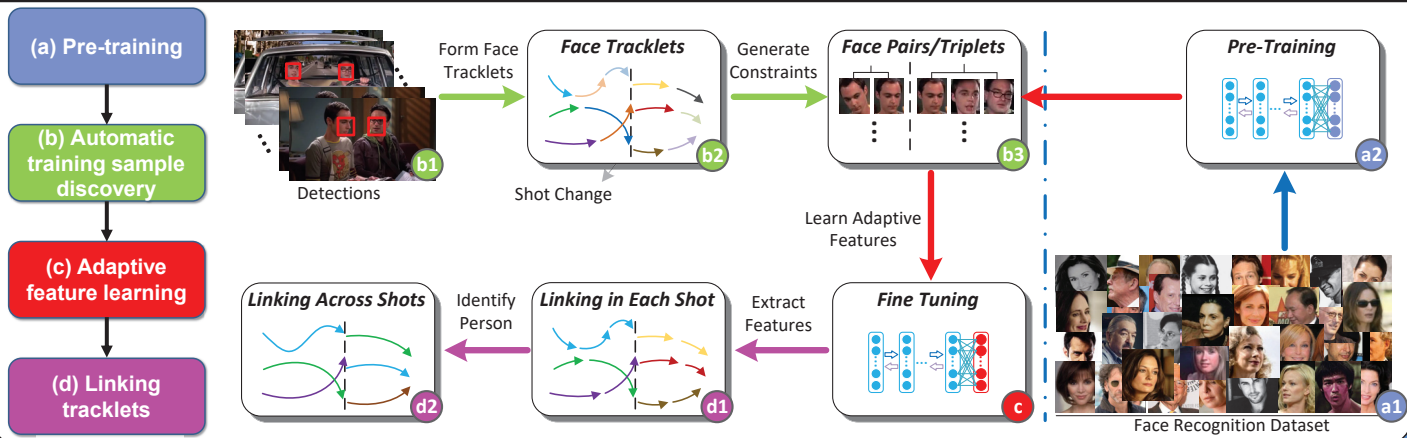


- Frequent shot changes
- Large face appearance variations
- Low resolution, motion blurring, occlusion and so on

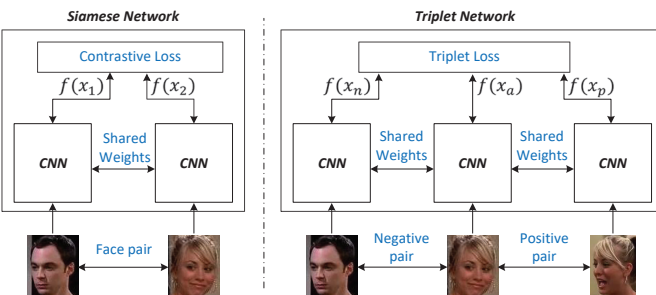
## Contributions

- Learning video-specific features by adapting deep CNN based on contrastive and triplet loss
- An improved symmetric triplet loss function (SymTriplet)
- A fully automatic multi-face tracking algorithm (detection, tracking, clustering, and feature adaptation)
- A new dataset with 8 music videos from YouTube

## Overview



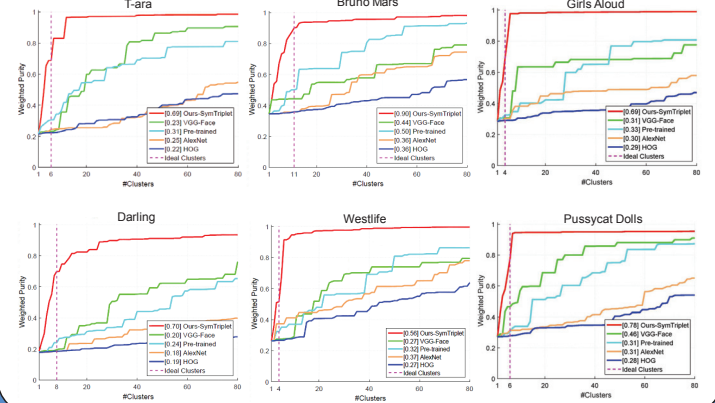
## Adaptive Feature Learning



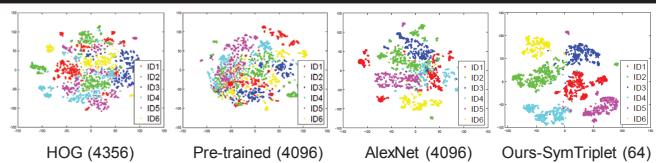
- SymTriplet loss:

$$L_s = \max\left\{0, d(x_a, x_p) - \frac{1}{2}(d(x_a, x_n) + d(x_p, x_n)) + \alpha\right\}$$

## Weighted Purity vs. #Clusters



## t-SNE visualization



## Quantitative Evaluation

Music video dataset							
Method	Recall↑	Precision↑	F1↑	FAF↓	IDS↓	Frag↓	MOTA↑
mTLD [Kalal et al. 2012]	9.7%	36.1%	15.3%	0.39	280	621	68.4%
ADMM [Ayazoglu et al. 2012]	75.5%	61.8%	68.0%	0.50	2382	2959	63.7%
IHTLS [Dicle et al. 2013]	75.5%	68.0%	71.6%	0.41	2013	2880	63.7%
Pre-trained	60.1%	88.8%	71.7%	0.17	931	2140	79.5%
Ours-Siamese	71.5%	89.4%	79.5%	0.12	986	2512	64.0%
Ours-Triplet	71.8%	88.8%	79.4%	0.20	902	2546	64.2%
Ours-SymTriplet	71.8%	89.7%	79.8%	0.19	699	2563	64.3%

## Qualitative Results

