# Learning from Video and Text via Large-Scale Discriminative Clustering

**Antoine Miech¹², Jean-Baptiste Alayrac¹², Piotr Bojanowski¹, Ivan Laptev¹², Josef Sivic¹²³**

¹Inria, Paris

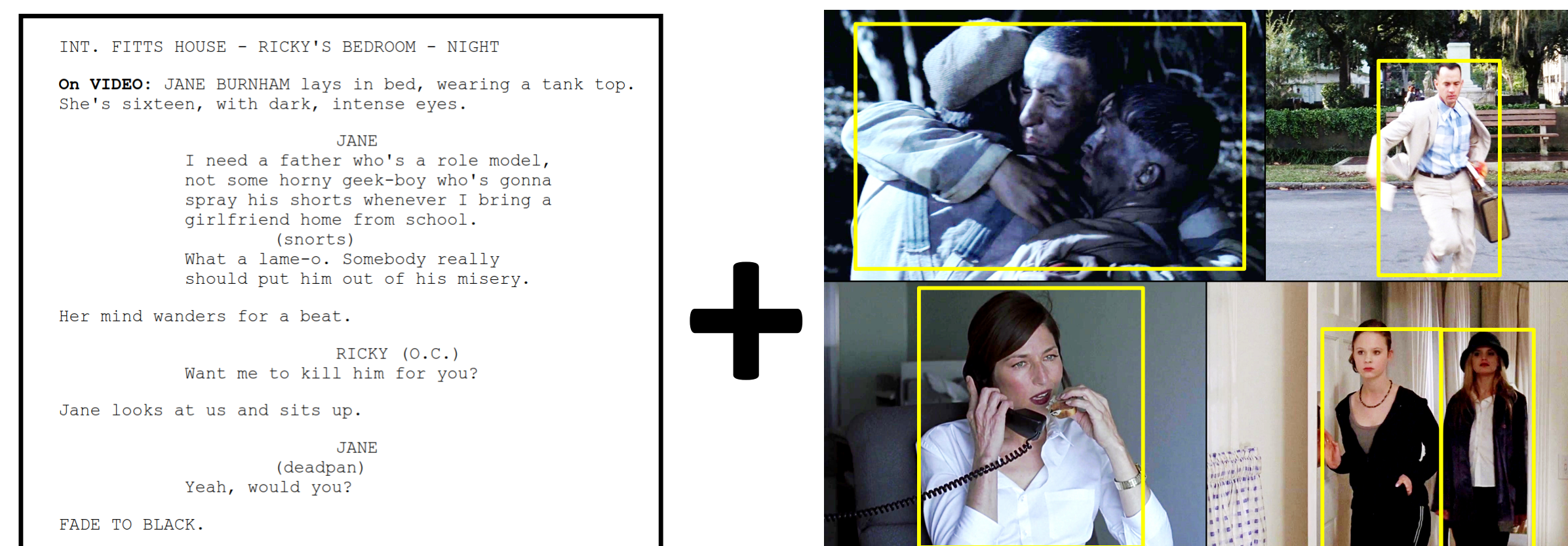²DI ENS, École Normale Supérieure, PSL Research University, Paris

³CIIRC, CTU, Prague

**Keywords:** Text-Video, Weak-Supervision, Discriminative Clustering Person-Action Recognition, Block-Coordinate Frank-Wolfe

## Goal

• **Scale-up** discriminative clustering for weakly supervised learning

• Demonstrate **weakly supervised** learning of actors and actions on large-scale dataset of movies
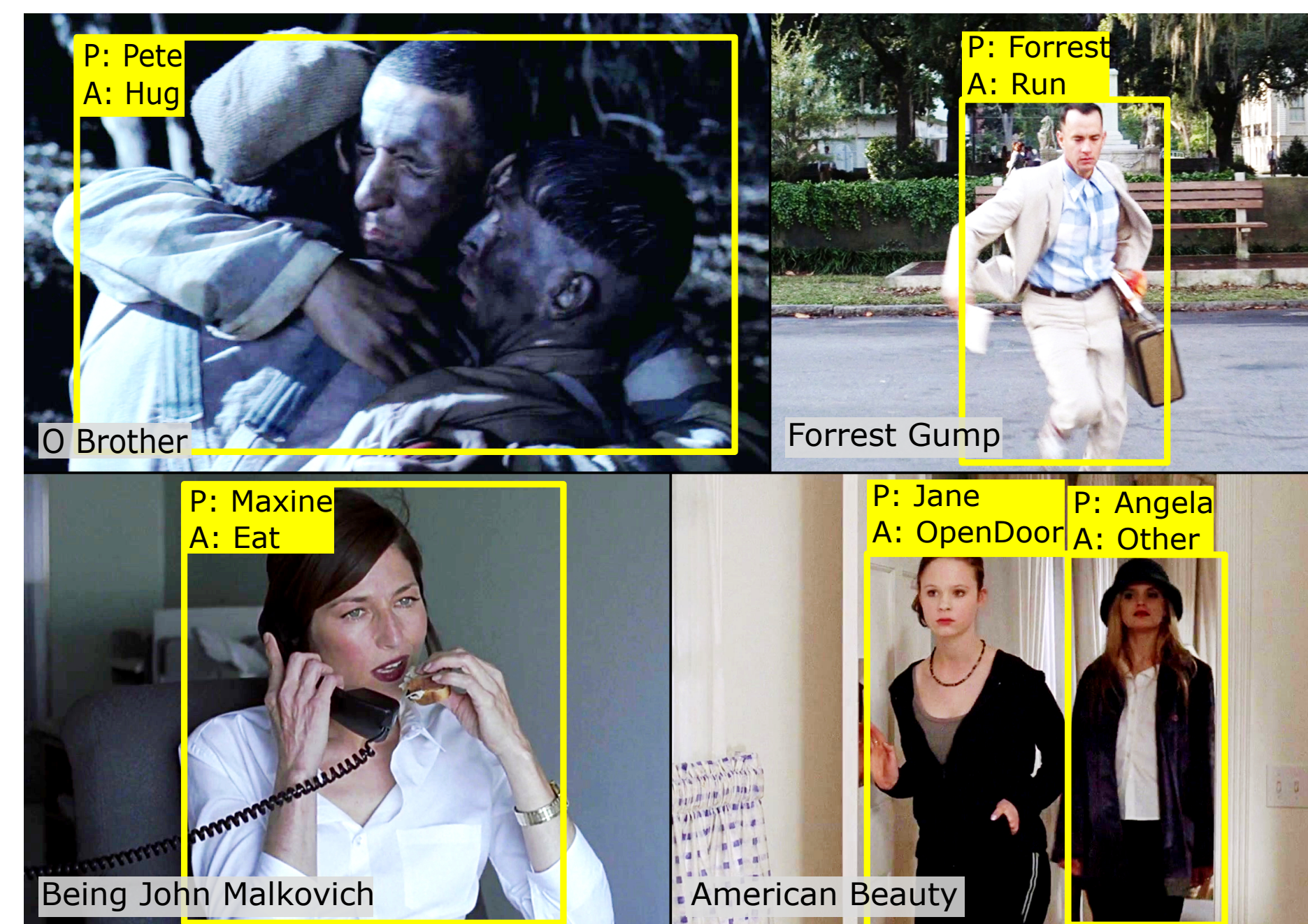
**INPUT**



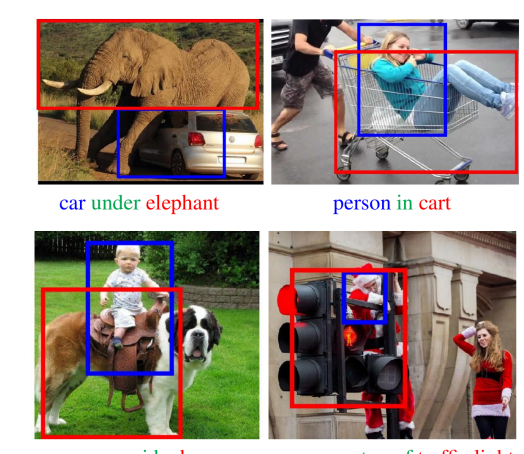Movie script "Free Annotation" + Pre-extracted Person tracks
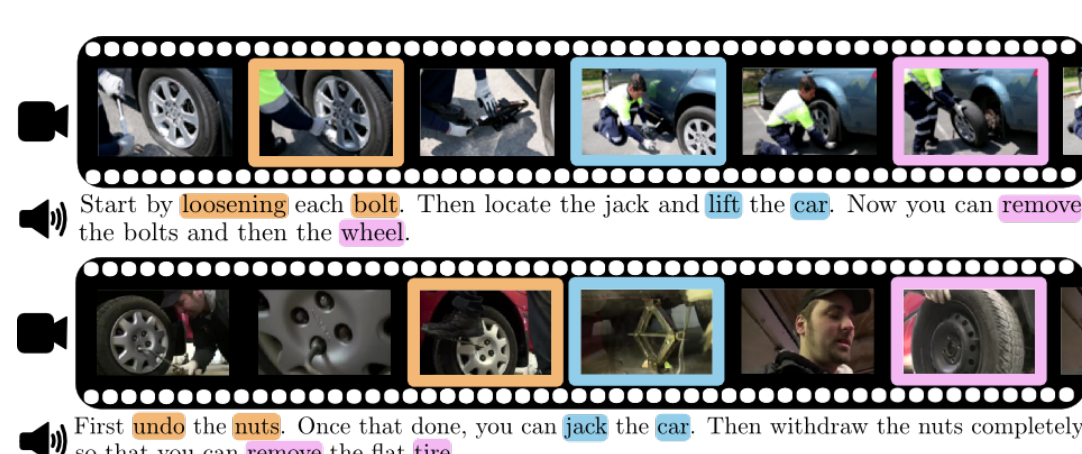
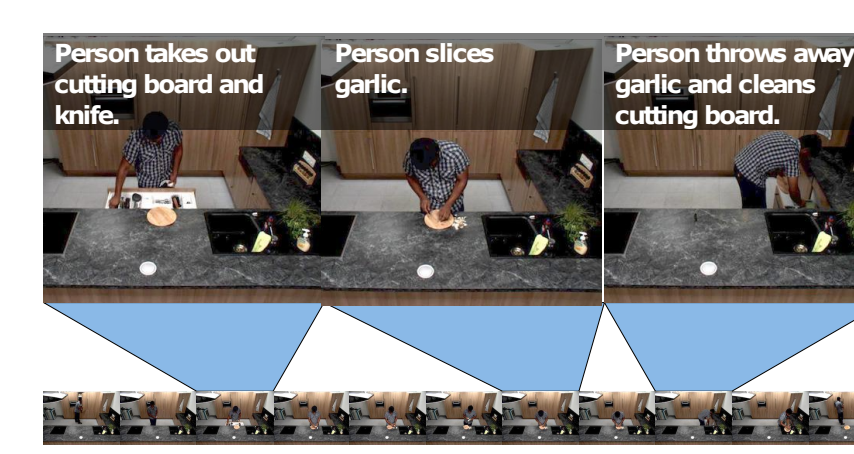**OUTPUT**



Person + Action recognition

## Motivation

Scale-up other weakly-supervised applications:



*Weakly-supervised learning of visual relations* [Peyre et al. ICCV17]

*Unsupervised learning from narrated instructional videos* [Alayrac et al. CVPR16]

*Weakly-Supervised Alignment of Video with Text* [Bojanowski et al. ICCV15]
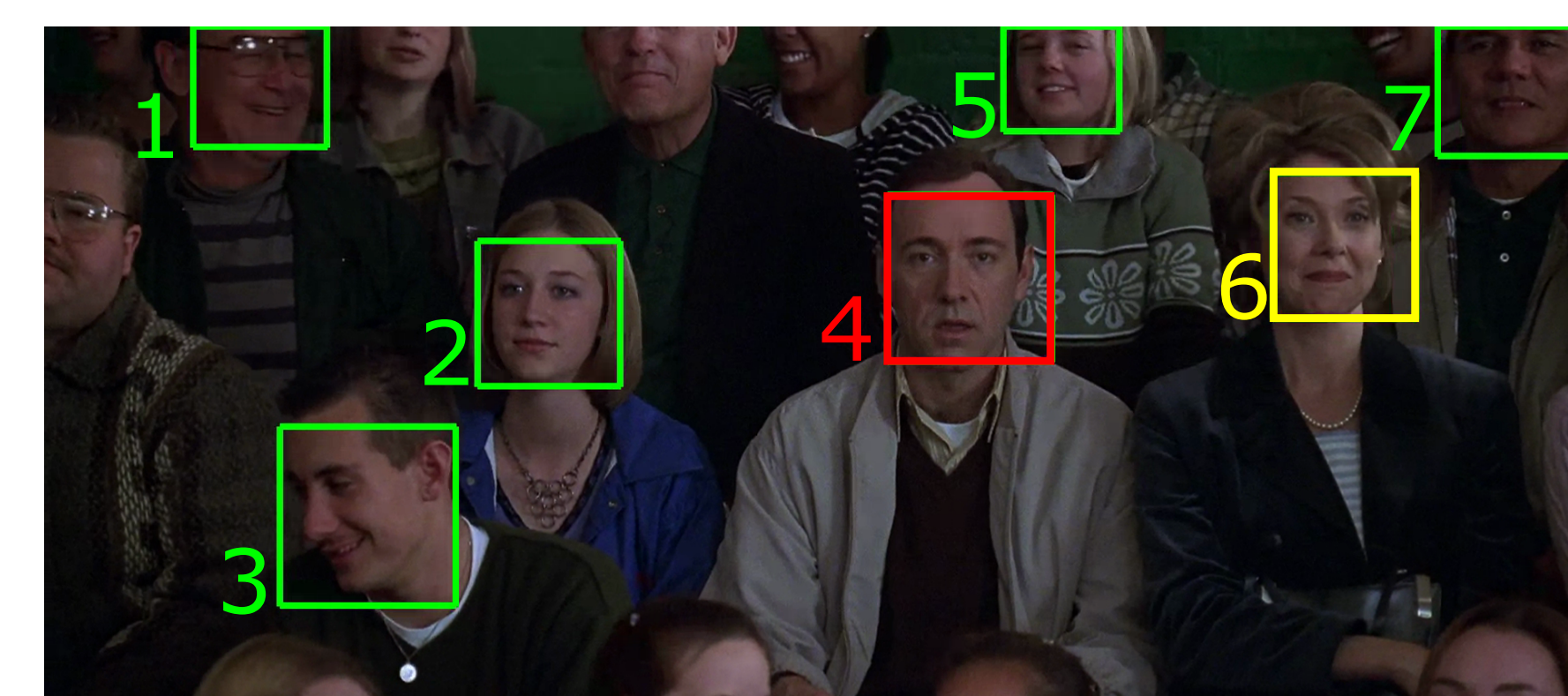
## Contributions

• New online optimization algorithm based on Block-Coordinate Frank-Wolfe (BCFW) for scaling-up discriminative clustering

• Improved model of the background class

## Discriminative Clustering

$$\min_{Z,W} \frac{1}{N}\|XW - Z\|_F^2 + \lambda\|W\|_F^2$$

• $Z \in \mathbb{R}^{N \times K}$: Assignment matrix (e.g Person name or Action class)

• $X \in \mathbb{R}^{N \times d}$: Person tracks features (e.g VGG-face features for face recognition and Improved Dense Trajectories for Action Recognition)

• $W \in \mathbb{R}^{d \times K}$: Linear model to learn



**American Beauty**
Scene i Person name parsed from Script: **LESTER**

**Weak-supervision as Linear Constraints on Z**

• **At Least One Constraint**

• **Background Class Constraint**

• **Mutual Exclusion Constraint**
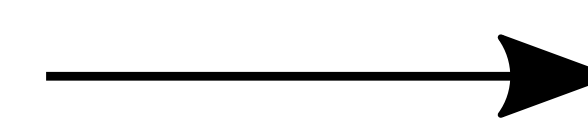
[Bach and Harchaoui, 2007] [Bojanowski et al., 2013]

## Large-Scale Optimization

**An algorithm based on the Block-Coordinate Frank-Wolfe method for efficient online optimization**

$$Z = \begin{matrix} Z^1 \\ Z^2 \\ \vdots \\ Z^{N_{movies}} \end{matrix}$$

• Z is a **block constraint separable** variable

• Exploit the Block-Coordinate Frank-Wolfe algorithm to treat each block in **an online manner**

• Efficient **Time and Space complexity** of block gradient computation via **smart update rules**
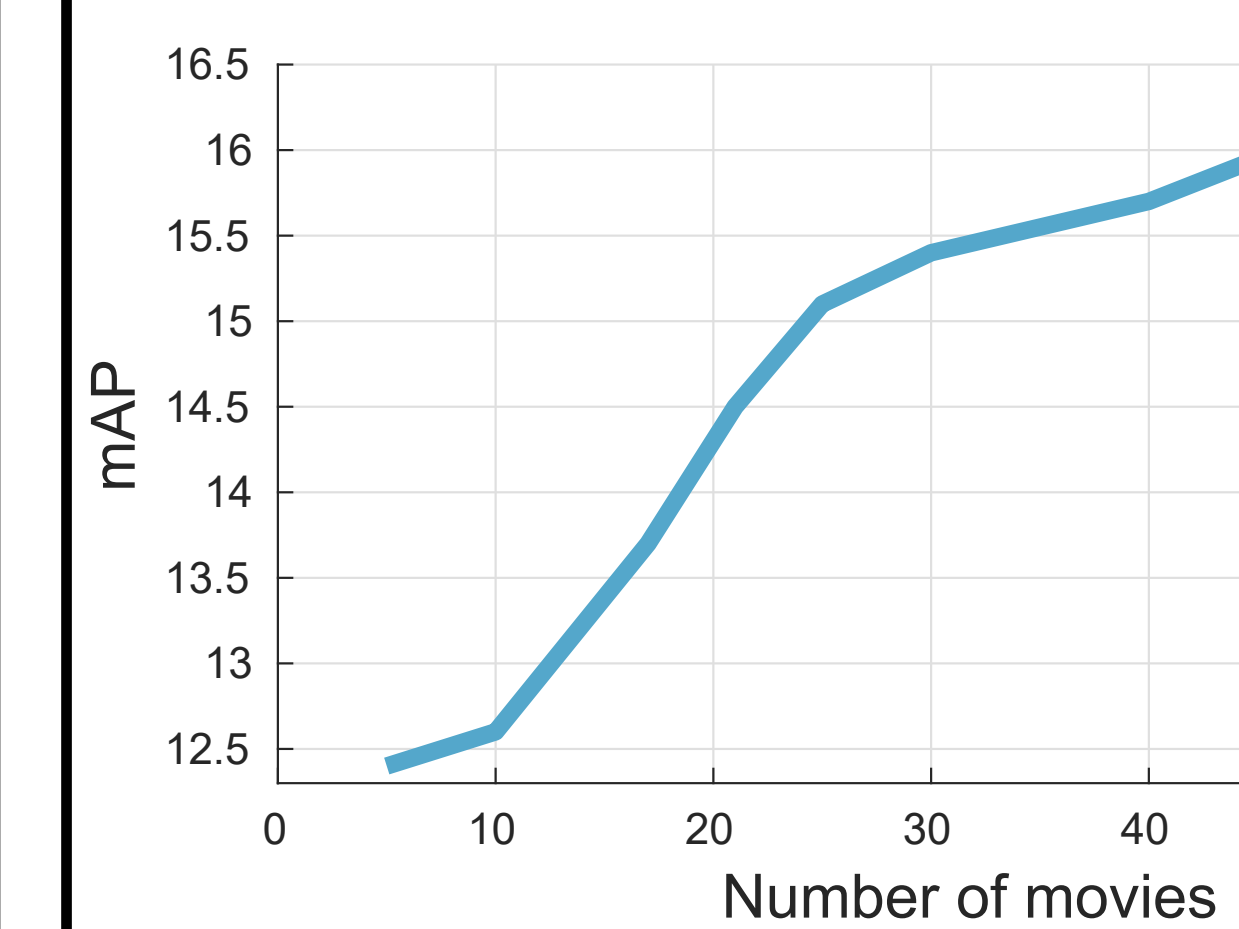
• **Convex relaxation** of the problem

**Standard FW step**
Space complexity: $\mathcal{O}(N^2)$
Time complexity: $\mathcal{O}(N)$

→

**Our optimized BCFW step**
Space complexity: $\mathcal{O}(N_{block})$
Time complexity: $\mathcal{O}(N_{block})$

## Results

**Dataset:** 66 feature-length movies together with scripts

**Actions:** A vocabulary of 14 different actions

Comparison of different method for Person recognition on *Casablanca*

| Method | Acc. | Multi-Class AP | Background AP |
|---|---|---|---|
| Cour et al. [8] | 48% | 63% | – |
| Sivic et al. [35] | 49% | 63% | – |
| Bojanowski et al. [4] | 57% | 75% | 51% |
| Parkhi et al. [27] | 74% | 93% | 75% |
| Our method | **83%** | **94%** | **82%** |

Performance when varying the number of training movies

Comparison of different method for the Action Sit Down on *Casablanca*



- Ours (66 movies) AP=0.25
- Ours (1 movie) AP=0.13
- Bojanowski et al. AP=0.06