

MOTIVATION

- **Unprecedented growth in video data**
 - Makes it challenging to store and consume
 - Video summarization attempts to address these challenges, but still an unsolved problem
- **What constitutes a good summary varies from domain to domain**
 - Sports video - 'importance' is more important than 'diversity'
 - Surveillance video - 'diversity' is more important than 'importance'
 - Further, what is 'important' for soccer videos is different from what is 'important' for birthday videos
- **A good video summary of a domain should have both characteristics!**

OUR CONTRIBUTIONS

- **Joint problem** of learning domain specific importance of segments as well as the desired summary characteristic for that domain
- An effective rating mechanism to serve as **indirect ground truth**
- A **novel evaluation measure**, more naturally suited in assessing the quality of video summary for the task at hand than F1 like measures
- A **gold standard dataset** for domain specific video summarization, first known dataset of long videos

METHODOLOGY

- **Weighted mixture** of modular and submodular terms
 - **Modular** terms: capture the domain specific importance of snippets
 - **Submodular** terms like Set Cover, Facility Location etc.: impart certain desired characteristics to the summary
- Learn weighted mixture using **max margin learning framework**
 - Different weights learnt for different domain
- For any given test video of that domain, the **weighted mixture is then maximized** to produce the desired summary video

GOOD SUMMARY?

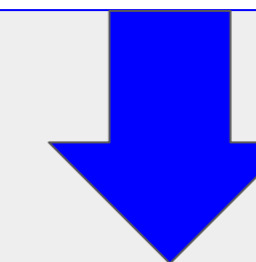


Soccer

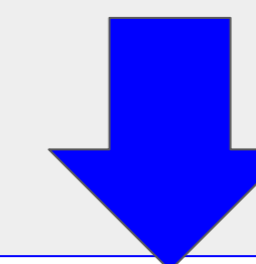
Birthday

Surveillance

Eliminate motionless chunks



Eliminate redundancy



Retain what is important

Naive

Better

Ideal

$$y^* = \operatorname{argmax}_{y \subseteq Y_v, |y| \leq k} o(x_v, y)$$

$$o(x_v, y) = w^T f(x_v, y)$$

$$\min_{w \geq 0} \frac{1}{N} \sum_{n=1}^N L_n(w) + \frac{\lambda_1}{2} \|w_1\|^2 + \frac{\lambda_2}{2} \|w_2\|^2$$

$$L_n(w) = \max_{y \subseteq Y_v^n} (w^T f(x_v^n, y) + l_n(y)) - w^T f(x_v^n, y_{gt}^n)$$

Formulation

Evaluation Measure

$$S_V(y) = \sum_{x_i \in X_P} |y \cap x_i| * \left(1 + \frac{|y \cap x_i|}{|x_i|}\right) * e^{\alpha * \text{rating}(x_i)}$$

$$+ \sum_{x_i \in X_R} \min(|y \cap x_i|, \beta) * \left(1 + \frac{\min(|y \cap x_i|, \beta)}{\min(|x_i|, \beta)}\right) * e^{\alpha * \text{rating}(x_i)}$$

$$- \sum_{x_i \in X_N} |y \cap x_i| * k$$

Positively Rated: Reward

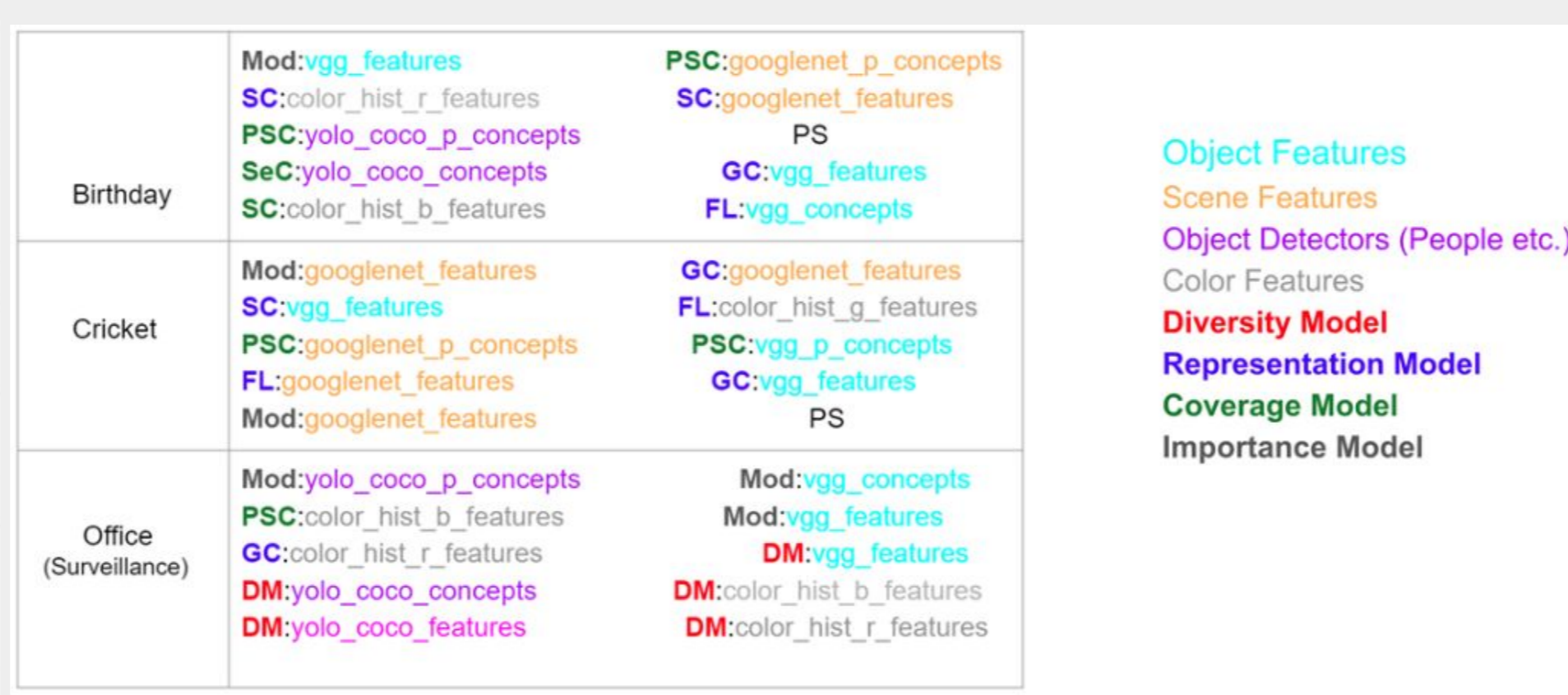
Negatively Rated: Penalize

Repetitive: Saturate

RESULTS & CONCLUSIONS

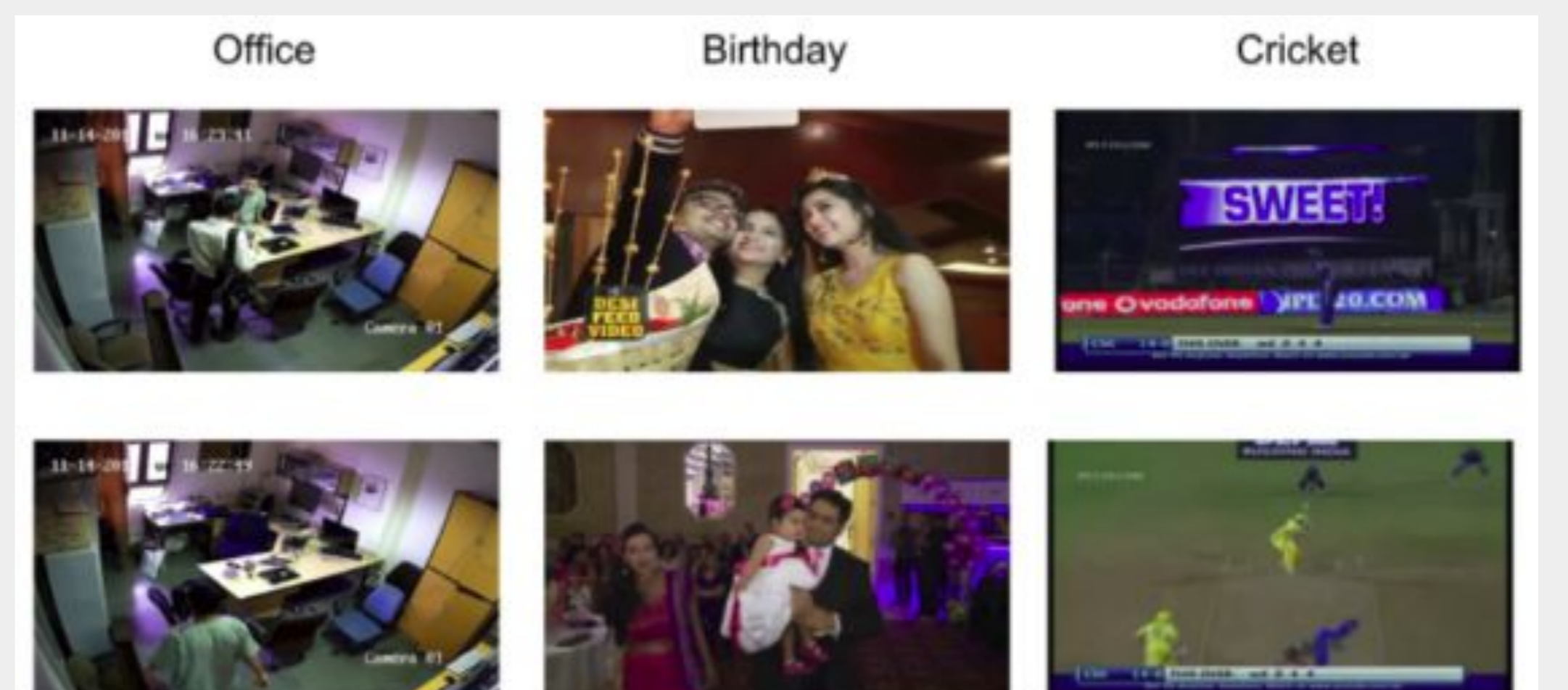
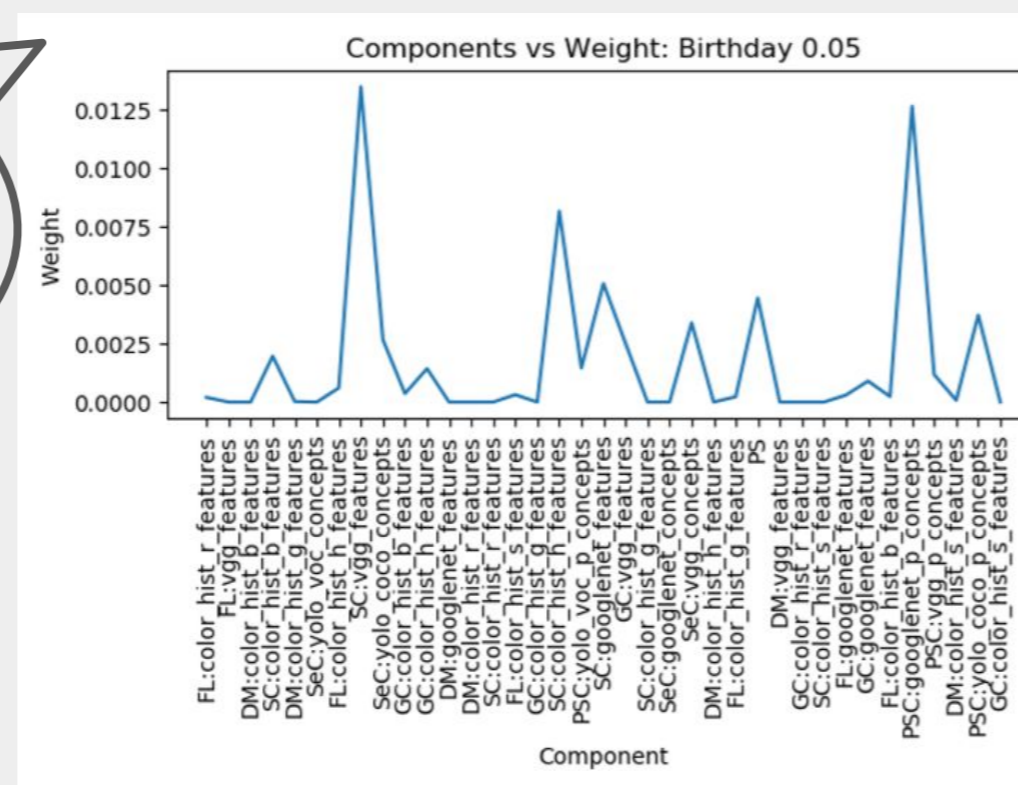
- **Full mixture** works best
- Models trained on one **domain** do not perform well on another
- **Multiple ground truths** help
- **Strong correlation** between components based on learnt weights and components with highest individual score when optimized
- **Intuitive relationship** of individual components with domains

Domain	Method	ScoreLoss
Birthday	All Modular	0.7234
	All Submodular	0.7307
	Full	0.6625
	Random	0.7378
	Uniform	0.7569
Cricket	All Modular	0.7432
	All Submodular	0.5967
	Full	0.5884
	Random	0.7706
	Uniform	0.7785
EntryExit	All Modular	0.6306
	All Submodular	0.6306
	Full	0.5884
	Random	0.7706
	Uniform	0.7785
Office	All Modular	0.8140
	All Submodular	0.8275
	Full	0.7733
	Random	0.8911
	Uniform	0.8979
Soccer	All Modular	0.8849
	All Submodular	0.7645
	Full	0.6533
	Random	0.9217
	Uniform	0.8747
Office (Surveillance)	All Modular	0.9152
	All Submodular	0.9152
	Full	0.6533
	Random	0.9217
	Uniform	0.8747



Different terms relevant for different domains

Full mixture performs the best, as hypothesized



Qualitative verification of video summaries produced