BOSTON





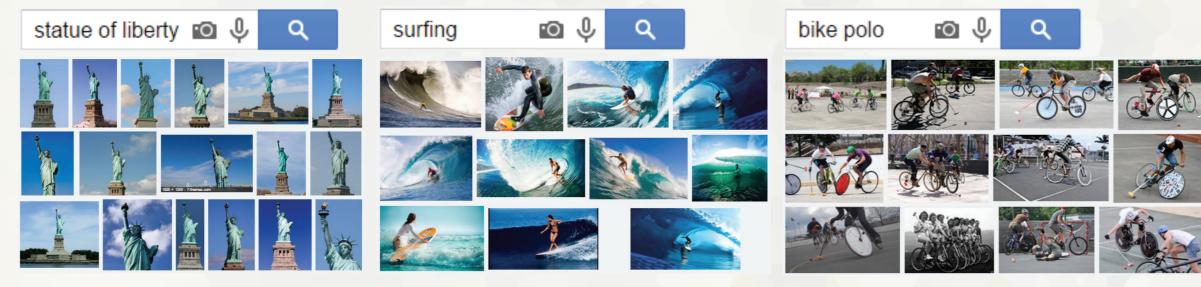
 Δ Require no prior knowledge and annotated data

- Dictionary learning: Cong et al. [TMM12], Zhao & Xing [CVPR14]
- Hierarchical clustering: Mahmoud et al. [ICMLA13]
- Δ Use additional resources
- Human attention during viewing videos: Ngo et al [TCSVT05]
- Web image priors: Khosla et al. [CVPR13], Kim et al. [CVPR14]





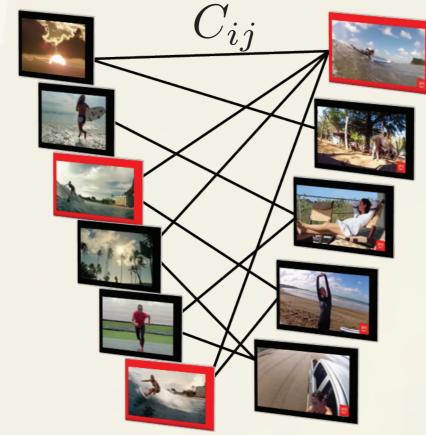
 Observations Δ Videos can share common topics (eg, retrieved by a query string) Δ Important concepts are likely to repeat visually.



Main idea

Input Video





[1] I.S. Dhillon. Co-clustering documents and words using bipartite spectral graph partitioning. In SIGKDD, 2001.

Video Co-summarization: Video Summarization by Visual Co-occurrence



Video Co-summarization

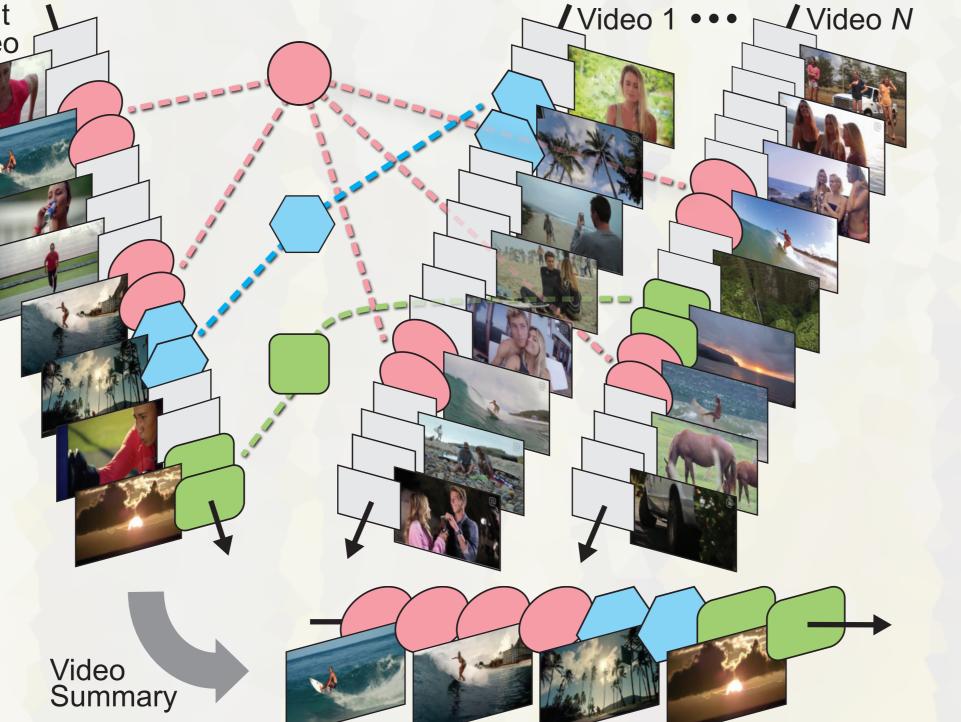
• A desirable video summarization method

Generate adaptive summaries that fits user's interests Scale to large dataset

Require limited/no human supervision



△ Exploits visual co-occurrence across multiple videos Input Query: surfing



Shot segmentation



Visual co-occurrence as bipartite graph and co-clusters [1]

- 1. Compute $\mathbf{D}_1 = \operatorname{diag}(\mathbf{C}\mathbf{e}_n)$
- 2. Compute $\mathbf{D}_2 = \operatorname{diag}(\mathbf{C}^{\top}\mathbf{e}_m)$
- 3. Compute $\widehat{\mathbf{C}} = \mathbf{D}_1^{-1/2} \mathbf{C} \mathbf{D}_2^{-1/2}$
- 4. Perform clustering on $Z = [D_1^{-1/2}U; D_2^{-1/2}V]$ where $\mathbf{C} = \mathbf{U} \Sigma \mathbf{V}^{\dagger}$

Maximal Biclique Finding (MBF)

• Visual co-occurrence as maximal bicliques \triangle Co-clusters could fail when shot co-occurrence is sparse. Δ We formulate co-sum as a maximal bicluque finding problem: $C_{ij}u_iv_j$

$$\max_{\mathbf{u},\mathbf{v}} \sum_{ij} \mathbf{0}$$

subject to $u_i + \mathbf{0}$

$$\mathbf{u} \in \{$$

 Δ Relax to a continous interval and impose sparsity-inducing norm:

$\max_{\mathbf{u},\mathbf{v}}$	\sum_{ij}			
subject to	$u_i +$			
	$\mathbf{u} \in$			

Algorithm 1: Maximal Biclique Finding (MBF)

Input : Bipartite gray
described by
parameters λ
Output: Maximal bic
Initialize $\mathbf{v} \leftarrow \texttt{rand}($
while not converged d
Compute $\widehat{u}_i = \min$
Update $u_i = \min$
Compute $\hat{v}_j = \min$
Update $v_j = \min$

Scalable easily to large dataset

costs $\mathcal{O}(mn^2 + n^3)$ due to an SVD

Closed-form updates

- Parallalizable
- Quality measure of discovered bicliques

$$q(\mathcal{B}) = \cdot$$

- Δ Similarities:
- ACA, TCD, and MBF discover visually similar shots in an unsupervised manner. Δ Differences:
- shots in its objective.
- TCD aims to locate one pair of shots at one time.
- TPAMI, 35(3):582–596, 2013.
- commonality discovery. In ECCV, 2012.

Wen-Sheng Chu¹, Yale Song² and Alejandro Jaimes²

 $v_j \leq 1 + I(C_{ij} \geq \epsilon), \forall i, j,$ $\{0,1\}^m, \mathbf{v} \in \{0,1\}^n$

 $\sum_{i} C_{ij} u_i v_j - \lambda_u \|\mathbf{u}\|_1 - \lambda_v \|\mathbf{v}\|_1$

 $+v_j \leq 1 + I(C_{ij} \geq \epsilon), \forall i, j$, $[0,1]^m, \mathbf{v} \in [0,1]^n$

aph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{W})$, where **W** is the co-occurrence matrix C; $\lambda_u \geq 0, \lambda_v \geq 0, \text{ and } \epsilon.$ clique indicated by **u** and **v** $(n) \in [0,1]^n;$ $\min\{I(\mathbf{C}_{ij} \ge \epsilon) - v_j\}_{j=1}^n;$ $n(I(\mathbf{C}_{i}; \mathbf{v} \ge \lambda_{u}), 1 + (\widehat{u}_{i})_{-});$

 $\min\{I(\mathbf{C}_{ij} \ge \epsilon) - u_i\}_{i=1}^m;$ $\mathbf{n}(I(\mathbf{u}^{\top}\mathbf{C}_{:j} \ge \lambda_v), 1 + (\widehat{v}_j)_{-});$

V Lower complexity: $\mathcal{O}(m+n)$, while co-clustering [1]

$$\frac{1}{|\mathcal{B}||} \sum_{ij} C_{ij} u$$

• Similarities and differences with ACA [2] and TCD [3]:

- ACA is clustering-based method, and by nature consider all

- MBF finds a group of shot pairs at once, and ensures each biclique contains only shots that are similar to each other.

[2] F. Zhou, F. De la Torre, and J. K. Hodgins. Hierarchical aligned cluster analysis for temporal clustering of human motion.

[3] W.-S. Chu, F. Zhou, and F. De la Torre. Unsupervised temporal

Query-specific Video Summarization

Dateset

 Δ Compiled a YouTube dataset using 10 queries from SumMe [4]. Δ SumMe contains one video/category; thus is not suitable for our purpose.

Features

 \triangle CENTRIST, Denst-SIFT, HSV color moments (mean, std, skew) Δ Shot-level representation using bag of temporal words △PCA-reduced dimension: 254+3840+108 --> 400

Subject evaluation

 Δ 3 judges see the query, and select >10% and <50% shots they think relevant. Δ Ground truth are constructed as the shots selected by >1 judges. Δ We used mean average precision (mAP) as an evaluation metric.

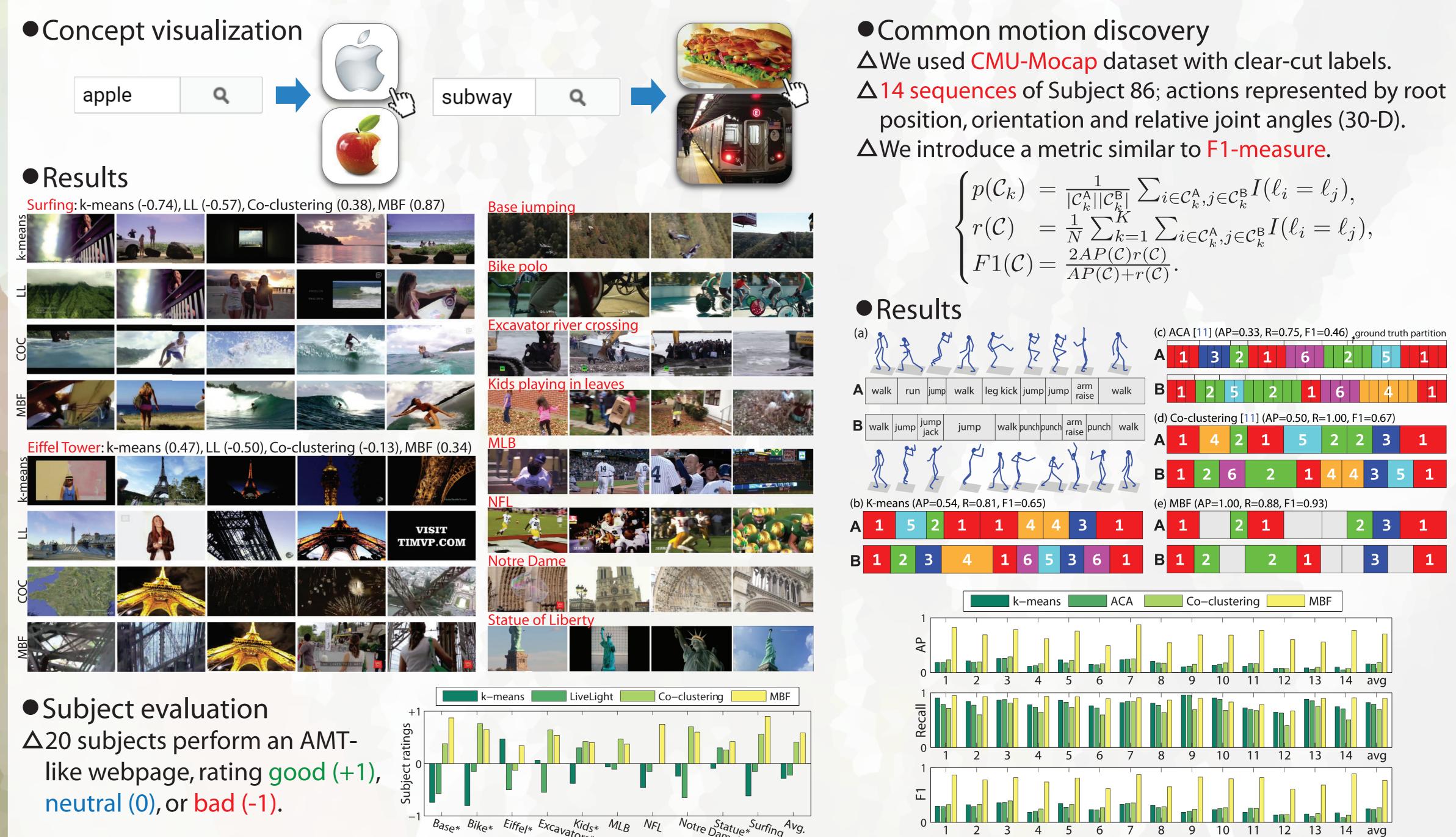
Methods

 Δ k-means (baseline), co-clustering (COC), LiveLight (LL) [4], and MBF.

[4] B. Zhao and E. Xing. Quasi real-time summarization for consumer videos. In CVPR, 2014.

	-	Methods	Base*	Bike*	Eiffel*	Excavators*	Kids*	MLB	NFL	Notre Dame*	Statue*	Surfing	Avg.
		k-means	0.432	0.427	0.422	0.289	0.791	0.556	0.663	0.392	0.543	0.550	0.507
Summaries generated	05	LL	0.226	0.305	0.413	0.667	0.744	0.508	0.710	0.568	0.763	0.334	0.524
tend to match more closely with human generated summaries	Top	COC	0.495	0.802	0.580	0.713	0.859	0.561	0.762	0.803	0.378	0.668	0.662
		MBF	0.680	0.788	0.596	0.690	0.798	0.638	0.680	0.715	0.810	0.684	0.707
		k-means	0.397	0.369	0.422	0.338	0.772	0.485	0.562	0.442	0.597	0.481	0.487
	15	LL	0.318	0.459	0.468	0.671	0.710	0.499	0.737	0.592	0.653	0.337	0.545
	Top	COC	0.496	0.795	0.561	0.656	0.852	0.503	0.823	0.676	0.458	0.586	0.641
		MBF	0.747	0.663	0.562	0.674	0.859	0.755	0.760	0.680	0.661	0.652	0.701
	_						0						

Concept Visualization & Common Motion Discovery



¹Carnegie Mellon ²YAHOO!

t		-					
Base jumping	Bike polo	Eiffel Tower	Excavato	or river Kids	Kids playing in leaves		
MLB	NFL	Notre Dame Cathedral	Statue of	Liberty S	Surfing		
Video quer	у	Length	#Vid	#Frm	#Shot		
Base jumpir	ng	10m54s	5	17960	241		
Bike polo		14m08s	5	22490	341		
Eiffel Tower	r	25m47s	7	43729	381		
Excavators	river xing	10m41s	3	16019	112		
Kids playing	g in leaves	15m40s	6	27972	238		
MLB		12m11s	6	21271	201		
NFL		13m28s	3	23179	405		
Notre Dame	e Cathedral	11m26s	5	20110	196		
Statue of Li	berty	10m44s	5	18542	164		
Surfing		22m40s	6	34790	483		
Total		147m40s	51	246062	2762		