# TVSum: Summarizing Web Videos using Titles

Yale Song*, Jordi Vallmitjana, Amanda Stent, Alejandro Jaimes
Computer Vision Group, Yahoo Labs NYC

YAHOO! LABS
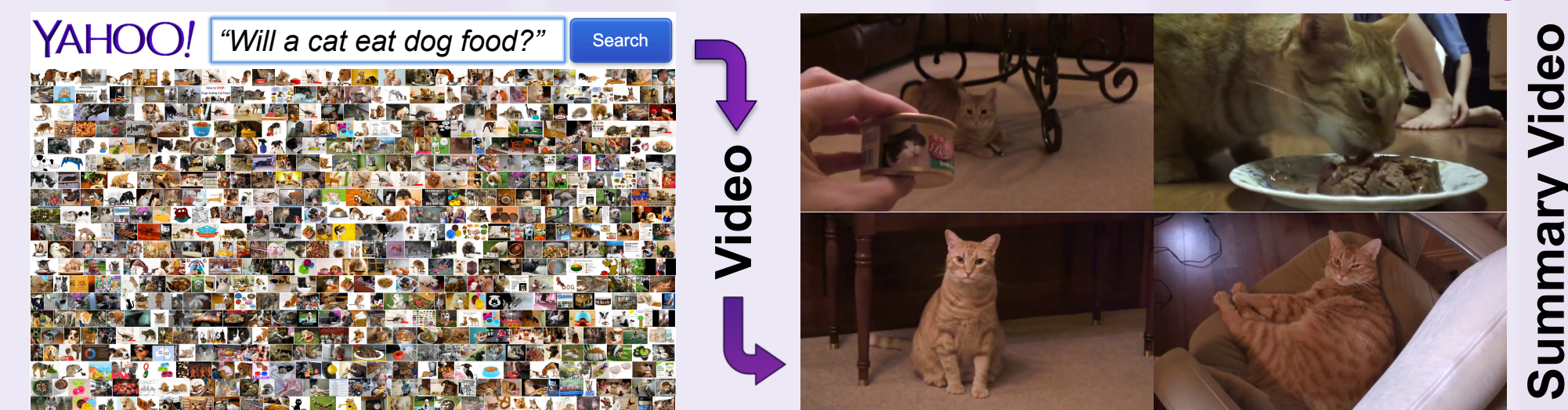
CVPR 2015 BOSTON JUNE 7-12

## Contributions

- ❑ **TVSum**: **T**itle-based **V**ideo **sum**marization system guided by natural language description of the video
- ❑ **Co-archetypal Analysis**: a novel algorithm for **cross-dataset canonical visual concept learning**
- ❑ **TVSum50**: a novel, diverse dataset for video sum., available via Yahoo! Webscope Program
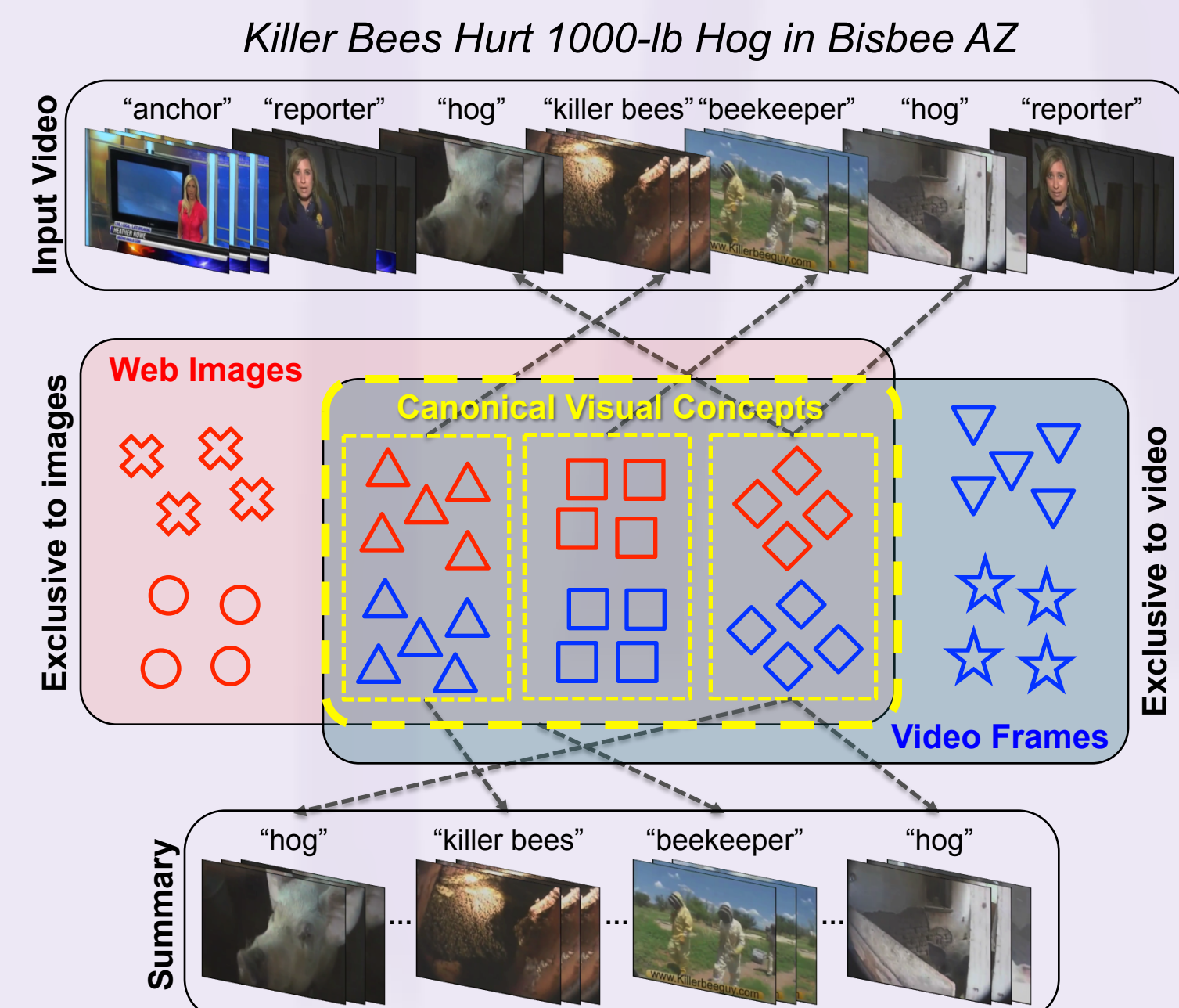
## Overview

- ❑ **Idea:** **a video title hints towards the main story**
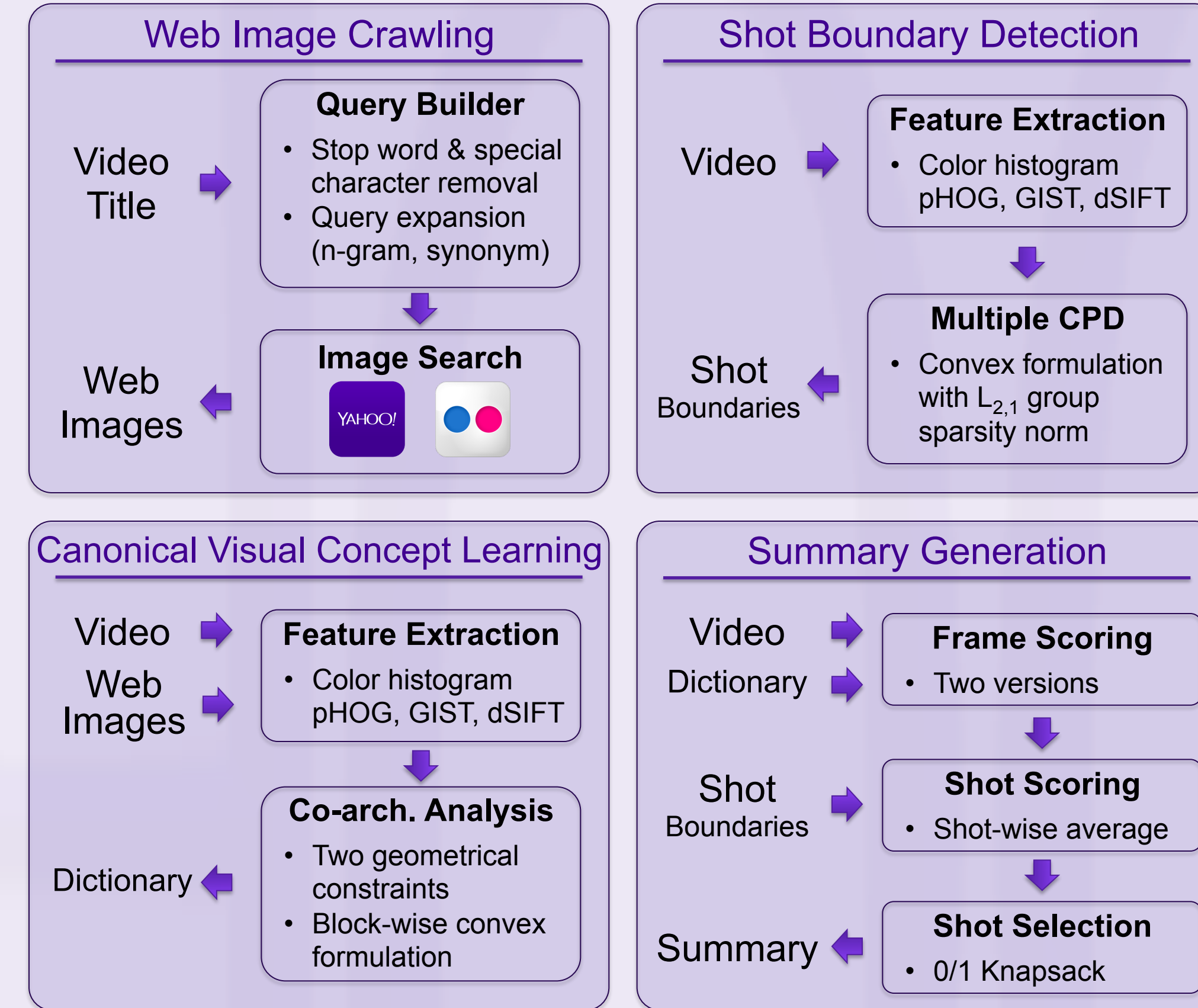
But real-world titles are free-formed, unconstrained, and ambiguous. This increases noise & variance in search results. Our work shows how to leverage title information effectively.

- ❑ **Solution: Co-archetypal Analysis**
- • Learn canonical visual concepts from a combination of video frames and title-based image search results

*Killer Bees Hurt 1000-lb Hog in Bisbee AZ*

## TVSum Pipeline

### Web Image Crawling

Video Title

**Query Builder**
- Stop word & special character removal
- Query expansion (n-gram, synonym)

**Image Search**
YAHOO!  flickr

Web Images

### Shot Boundary Detection

Video

**Feature Extraction**
- Color histogram pHOG, GIST, dSIFT

**Multiple CPD**
- Convex formulation with $L_{2,1}$ group sparsity norm

Shot Boundaries

### Canonical Visual Concept Learning

Video
Web Images

**Feature Extraction**
- Color histogram pHOG, GIST, dSIFT

**Co-arch. Analysis**
- Two geometrical constraints
- Block-wise convex formulation

Dictionary

### Summary Generation

Video
Dictionary

**Frame Scoring**
- Two versions

Shot Boundaries

**Shot Scoring**
- Shot-wise average

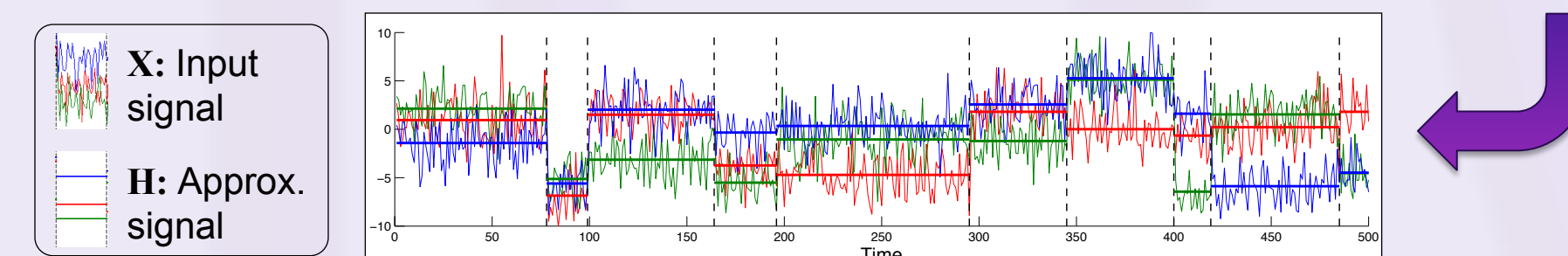Summary

**Shot Selection**
- 0/1 Knapsack

## Shot Boundary Detection

- ❑ **Multiple change point detection**

Convex formulation with an $L_{2,1}$ norm

[Bleakley & Vert, 2011]

$$\min_H \frac{1}{2}\|X - H\|_F^2 + \lambda \sum_{t=1}^{n-1}\|H_{\bullet, t+1} - H_{\bullet, t}\|_2$$

**Reconstruction Error** **Total Variation**

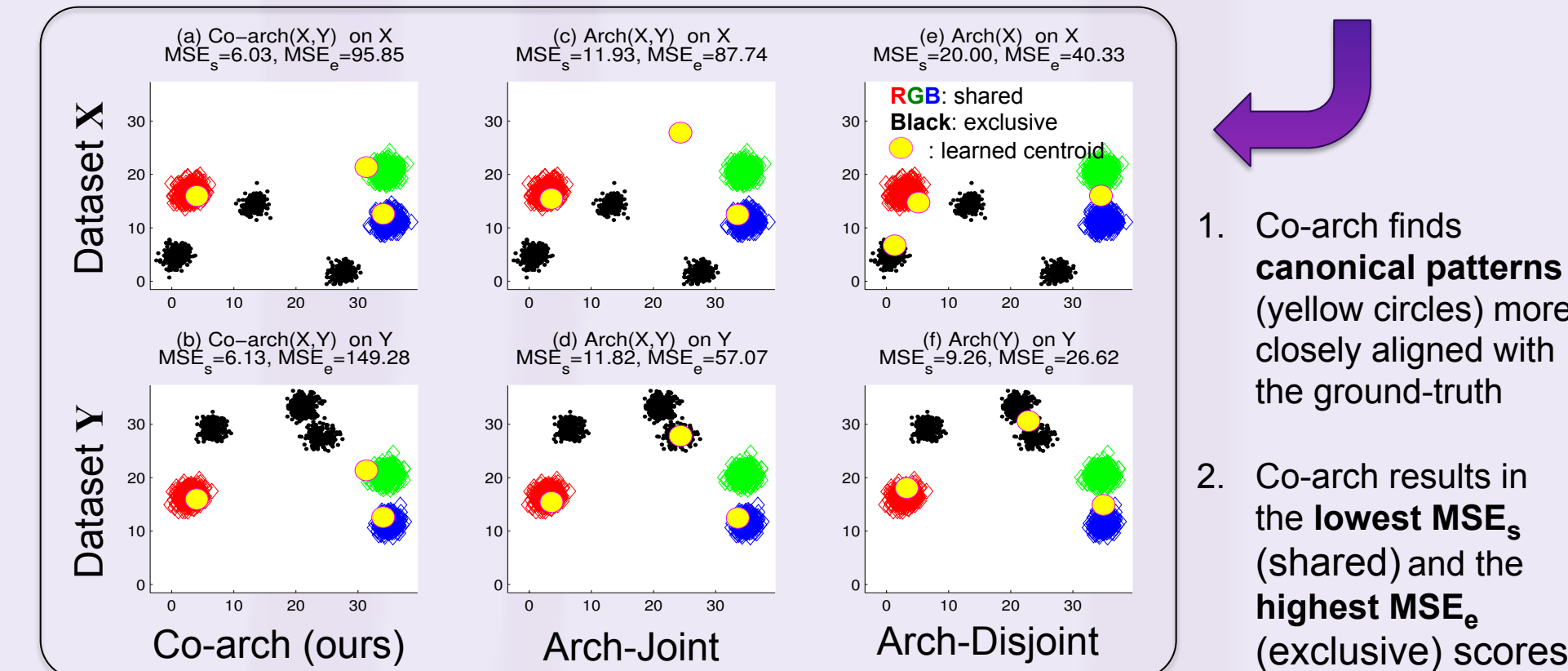X: Input signal
H: Approx. signal

- • **H** is **column-wise sparse** because of the $L_{2,1}$ norm
- • Comparison to DP solution on NASA videos [TREC-01]
  - ☑ **More accurate:** mean F1 score of 0.60 vs. 0.35
  - ☑ **Runs faster:** 7.5s vs. 39.4s on a 7 minute video
  - ☑ **More practical:** no need to specify #CP a priori

## Canonical Visual Concept Learning

- ❑ **Co-archetypal Analysis**
  - • Learn a dictionary of visual patterns $\mathbf{Z}$ shared *jointly* between video frames $\mathbf{X}$ and web images $\mathbf{Y}$
  - • Two geometrical constraints:
    1. Each $\mathbf{x}_i$ and $\mathbf{y}_i$ should be well approx.-ed by a convex combination of $\mathbf{Z}$, i.e., $\mathbf{X} \approx \mathbf{ZA}^X$ and $\mathbf{Y} \approx \mathbf{ZA}^Y$
    2. Each $\mathbf{z}_j$ should be well approx.-ed *jointly* by a C.C. of $\mathbf{X}$ and by a C.C. of $\mathbf{Y}$, i.e., $\mathbf{Z} \approx \mathbf{XB}^X \approx \mathbf{YB}^Y$

$$\min_\Omega \|\mathbf{X} - \mathbf{ZA}^X\|_F^2 + \|\mathbf{Y} - \mathbf{ZA}^Y\|_F^2 + \gamma\|\mathbf{XB}^X - \mathbf{YB}^Y\|_F^2$$

Geometrical constraint #1 Constraint #2
Video Factorization  Web Image Factorization  Video-Image Joint Regularization

Co-arch (ours)   Arch-Joint   Arch-Disjoint

1. Co-arch finds **canonical patterns** (yellow circles) more closely aligned with the ground-truth
2. Co-arch results in the **lowest MSE$_s$** (shared) and the **highest MSE$_e$** (exclusive) scores

- ☑ Compared to Archetypal Analysis [Cutler & Brieman 1994], our Co-archetypal Analysis learns a dictionary Z that captures canonical patterns that appear both in X (video) and Y (images), but not in either alone
- ☑ Effectively deals with noise (images irrelevant to video) and variance (multiple, unknown numbers of visual concepts)

## Summary Generation

- ❑ Compute shot scores by averaging frame scores

$$score_{ver1}(\mathbf{x}_i) = \|\mathbf{x}_i - \mathbf{Z}\alpha_i\|_2, \quad score_{ver2}(\mathbf{x}_i) = \sum_{j=1}^n \mathbf{B}_i\alpha_j$$

- ❑ Given a summary length budget, solve 0/1 knapsack

$$\max \sum_{i=1}^{\#shots} u_i \times score(\mathbf{s}_i) \text{ s.t. } u_i \in \{0,1\}, \sum_i u_i \times length(\mathbf{s}_i) \leq budget$$

## TVSum50 Benchmark Dataset

- ❑ 50 videos, 10 categories, 3.5 hours
- ❑ 1000 crowd ratings, 20 per video
- ❑ Avoiding chronological bias

SumMe [Gygli 2014]
Ours

- • Rating done by visual importance & relevance, without being affected by temporal precedence
- • Highly consistent ratings (Cronbach's α = 0.81)
- ❑ Available via Yahoo! Webscope Program

## Experiments

- ❑ Task: summarize video with 15% length budget
- ❑ Metric: mean pairwise F1 scores w.r.t. human summaries
- ❑ Features: color histograms, pHoG, GIST, dSIFT

| SumMe (25 videos) [Gygli et al '14] | | Our TVSum50 | |
|---|---|---|---|
| Methods | mpF1 | Methods | mpF1 |
| Uniform Sampling* | 0.14 | Uniform Sampling | 0.36 |
| K-means Clustering* | 0.16 | Random Sampling | 0.32 |
| Attention* [Ejaz '13] | 0.17 | K-means Clustering | 0.35 |
| Interestingness* [Gygli '14] | 0.23 | Spectral Clustering | 0.39 |
| LiveLight [Zhao & Xing '14] | 0.24 | LiveLight [Zhao & Xing '14] | 0.46 |
| Web Image Prior [Khosla '13] | 0.24 | Web Image Prior [Khosla '13] | 0.36 |
| **Co-archetypal Analysis** | **0.27** | Arch (Video Only) | 0.33 |
| | | Arch (Video + Web Image) | 0.35 |
| | | **Co-archetypal Analysis** | **0.50** |

* Results are from [Gygli et al. 2014]

- ❑ Qualitative Results

(a) Ruben Sandwich with Corned Beef & Sauerkraut    F$_1$ Score = 0.6247

(b) CASACL – Flashmob in Copenhagen underground – Peer Gynt    F$_1$ Score = 0.6061

(c) David Belle | Fondateur du parkour | Reportage de TF1    F$_1$ Score = 0.3612