

Recognition of Activity

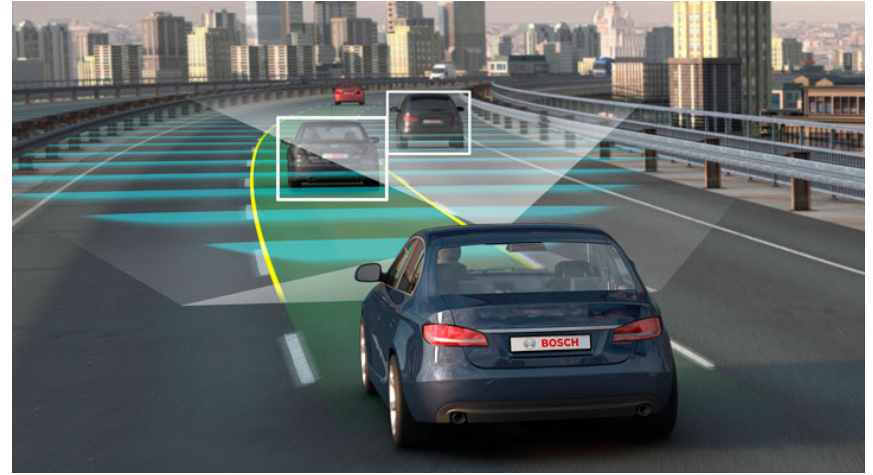
Cees G.M Snoek & Arnold W.M Smeulders
University of Amsterdam



MOTIVATION

The many faces of video.

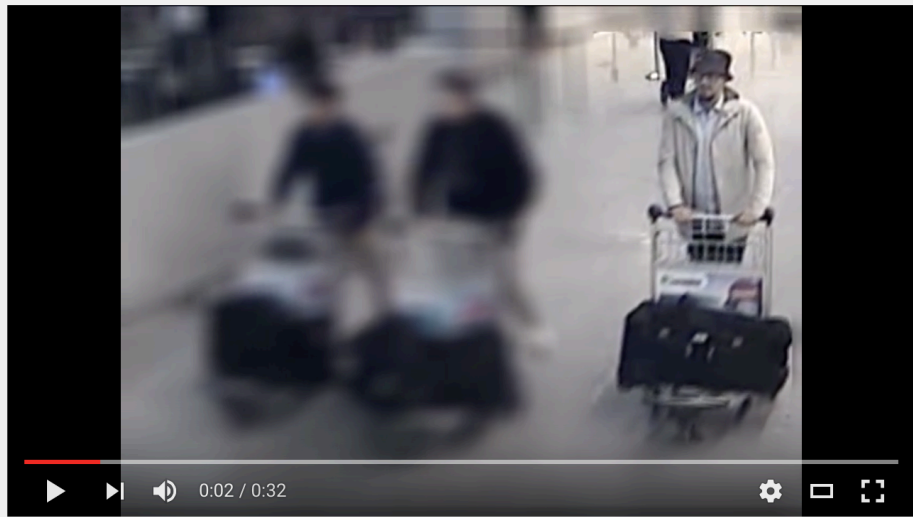
New purposes



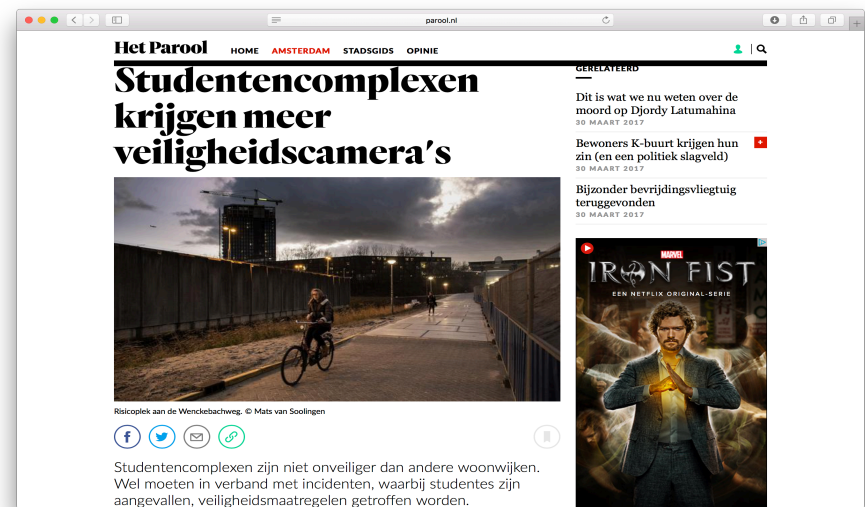
Video events will be used in many more new ways.

New surveillance

reconstruction



prevention



... and reused in old applications ...

New interaction



... and used as part of a loop.

FOCUS

Video is filled with what?

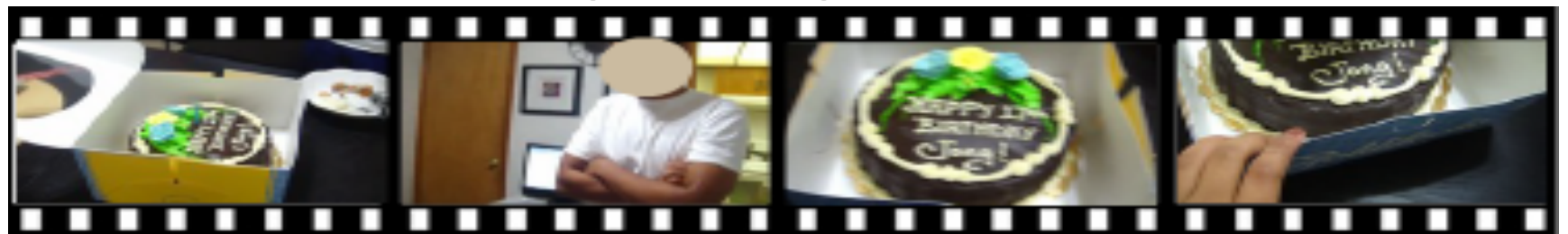
Acts > Actions > Events - ...



Acts: driving a screw



Actions: grooming an animal



Events: birthday party

Act > Action > Event > ...



more degrees of freedom,
rigid schemes can't work

act atomic motion pattern

sleeping

running

driving a screw

action functional pattern

shaking hands

removing a lit

serving the ball

event purposeful pattern

serving an ace

of actors, objects

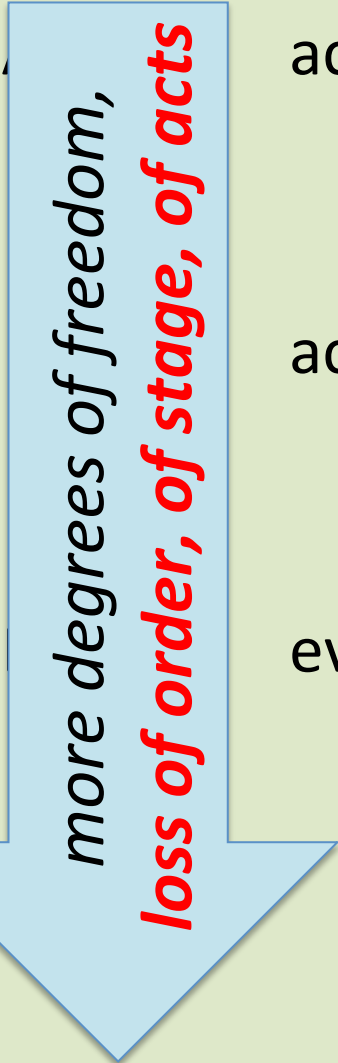
welcoming a friend

and motions

repairing an appliance

Act > Action > Event > ...

time frame *pattern variations*

 more degrees of freedom, loss of order, of stage, of acts	act	<i>sleeping</i>	$\pm 1s$	pose
		<i>running</i>	$\pm 2s$	dress, gait
		<i>driving screw</i>	$\pm 2s$	repetition pace
	action	<i>shaking hands</i>	2-5s	routine, active, solid
		<i>removing a lit</i>	2-5s	size, temperature
		<i>serving ball</i>	2-5s	camera
	event	<i>serving an ace</i>	5-10s	camera, in
		<i>welcoming</i>	1-5m	choice acts, actions
		<i>repairing</i>	1-60m	choice acts, actions

DRIVING FACTORS

What makes success?

Goalgole

Goalgole Demonstrator - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorites History Mail Size Print Edit

Address <http://localhost/VoetbalDemo/> Go

9 hours of video.

Data are the starting point.

Bay. Leverkusen-Man. United
Semi final of UEFA Champions League season 2001/2002, Second half at 25:01

Bay. Leverkusen-Man. United
Semi final of UEFA Champions League season 2001/2002, Second half at 14:49

Bay. Leverkusen-Man. United
Semi final of UEFA Champions League season 2001/2002, First half at 43:41

Bay. Leverkusen-Man. United
Semi final of UEFA Champions League season 2001/2002, Second half at 23:03

Bay. Leverkusen-Man. United
Semi final of UEFA Champions League season 2001/2002, Second half at 22:50

Bay. Leverkusen-Man. United
Semi final of UEFA Champions League season 2001/2002, Second half at 22:50

Paused 37683 / 103245 CC

Play Pause Stop

Done Local intranet

CCV Columbia

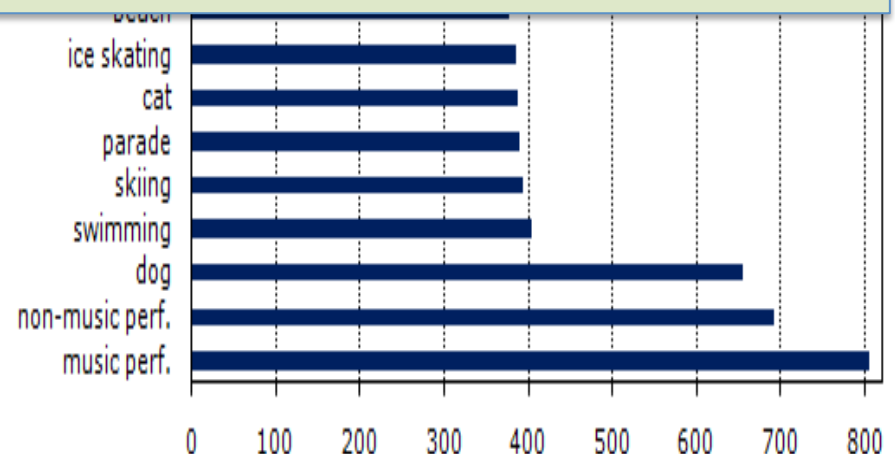
- # videos: 9,317
 - (210 hrs in total)



Everyone their own results.

Progress needs a community who agree.

- average length
 - 80 seconds
- # defined categories
 - 20
- annotation method
 - Amazon Mechanical Turk



TRECVID Internet video collections

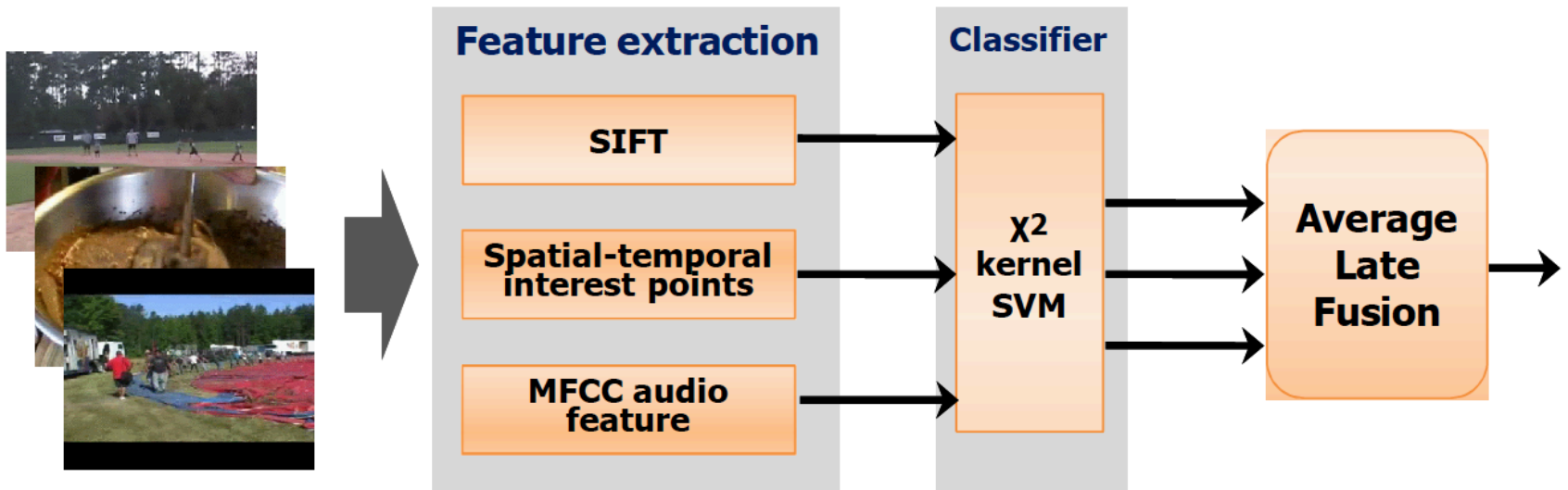
Collection Name	Designated Uses	Target sizes	Annotation
Pilot	<u>2010</u> Development collection	1,723 clips	Clip content annotation for both
<p>The important moment</p> <p><i>The driving factor is a shared & open competition.</i></p>			
	<u>2012-2015</u> (1) and (2) merged to a single training collection		Clip content annotation for the opaque subset
Progress	<u>2012-2015</u> : test collection	120K clips, 4000 hrs	No clip content annotation
Novel 1	<u>2014</u> : test collection	120K clips, 4000 hrs.	No clip content annotation
Novel 2	<u>2015</u> : test collection	120K clips, 4000 hrs.	No clip content annotation

CLASSIFICATION

Giving events a name, step by step TRECvid.

2010 Media diversity

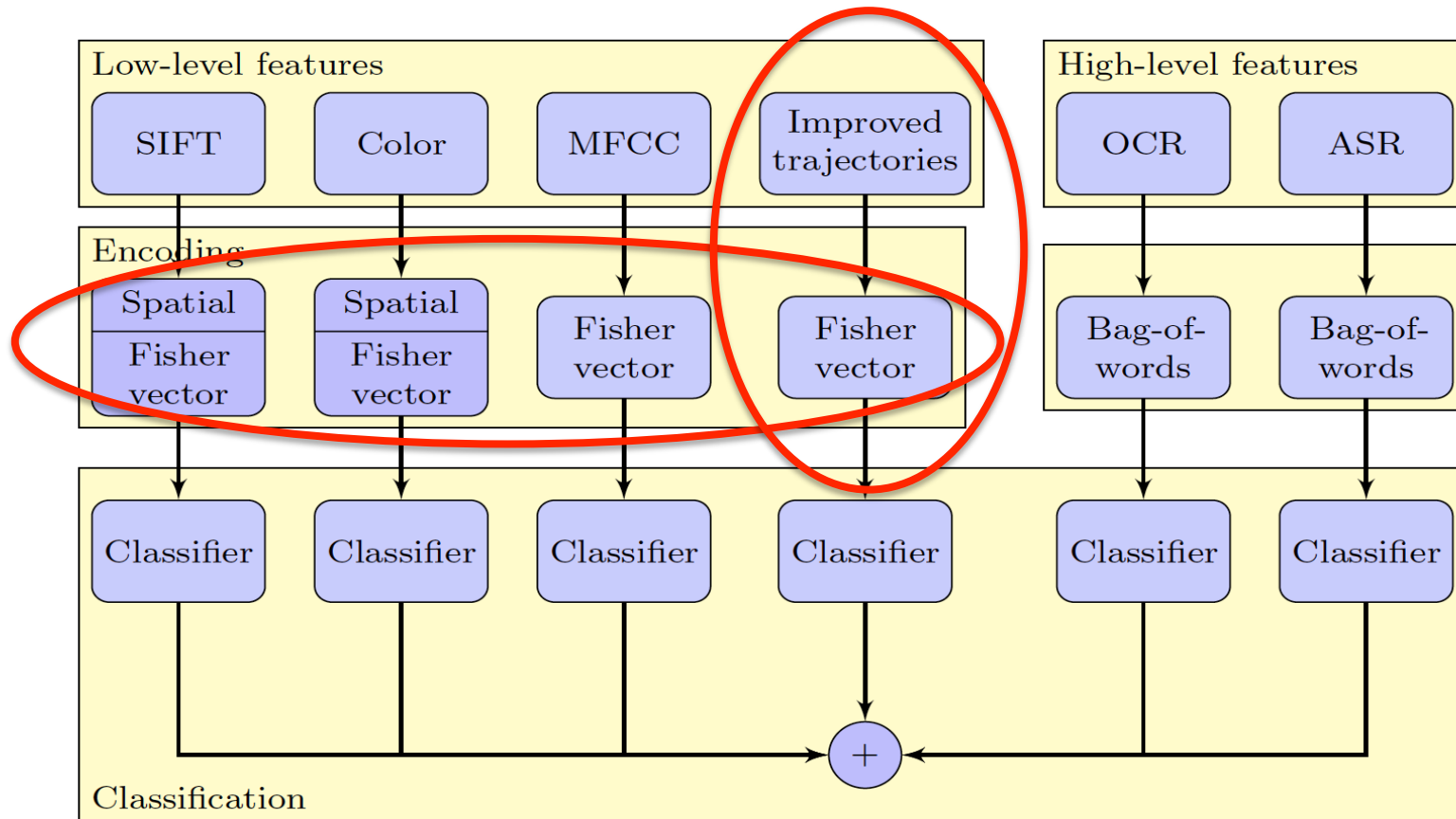
Diverse is better, more is better, fusion is better.



Y.G.Jiang TRECVID10 P.Natarjan CVPR12 Wang ICCV13 *others*

Yu-Gang Jiang, Xiaohong Zeng, Guangnan Ye, Subh Bhattacharya, Dan Ellis, Mubarak Shah, Shih-Fu Chang,
Columbia-UCF TRECVID2010 Multimedia Event Detection: Combining Multiple Modalities, Contextual Concepts, and Temporal Matching, NIST TRECVID Workshop, 2010.

2012/13 Trajectories and aggregation



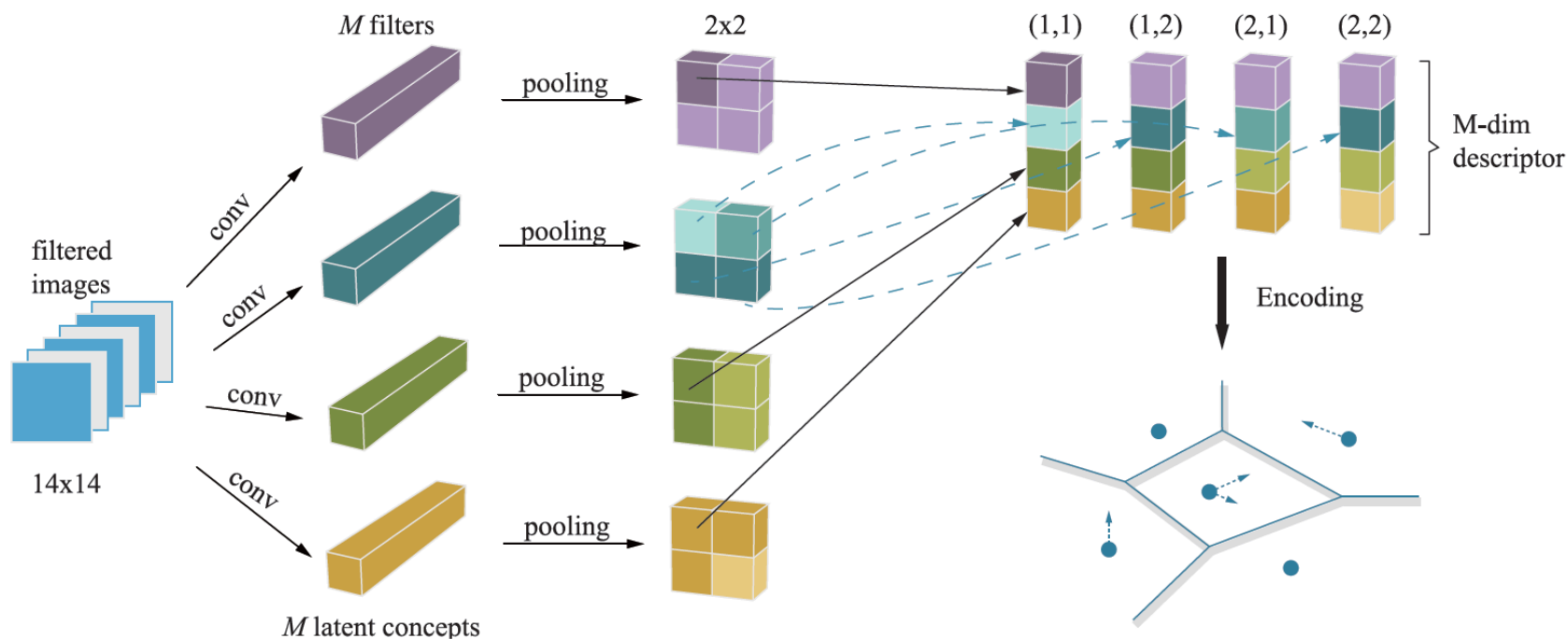
Dense trajectories are more and Fisher aggregation are fusion.

This is the end of hand engineered features.

2014 Deep learning & VLAD

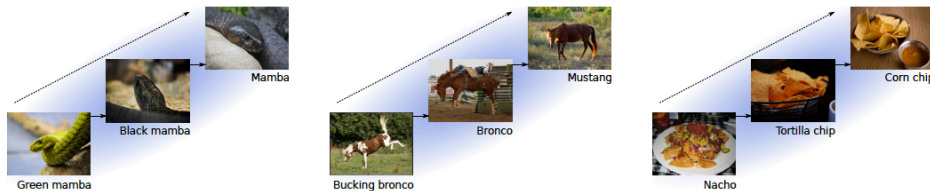
Networks integrate features and classifiers.

Deep learning builds in fusion of diverse, more and late.

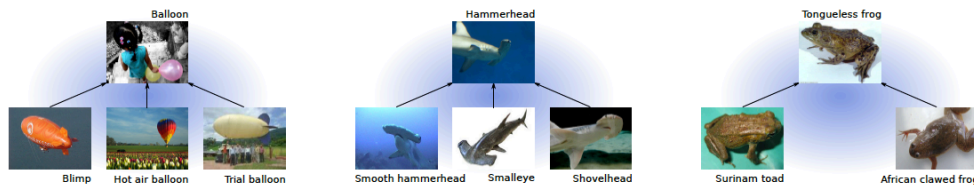


2015 Prior knowledge

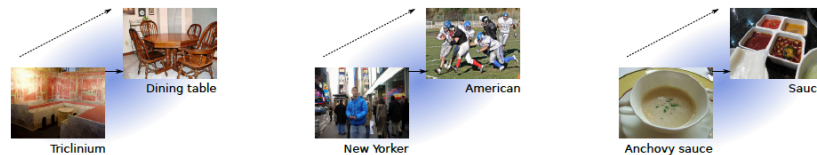
Insert 15000 ImageNet detectors pruned, but first reorganize prior knowledge removing fine semantics and merging small sets.



(a) Roll.



(b) Bind.



(c) Promote.



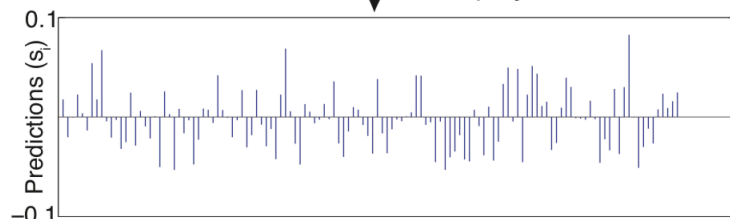
(d) Subsample.

2016 Joint embedding

Fuse media diverse in one embedding to compose stories

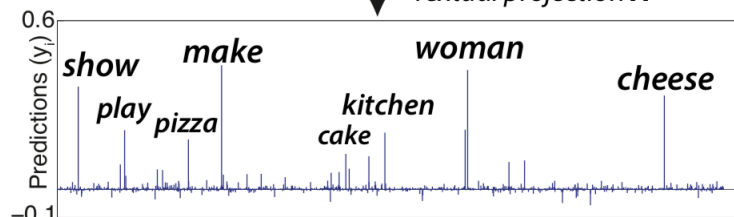


↓ Visual projection W



VideoStory

↓ Textual projection A



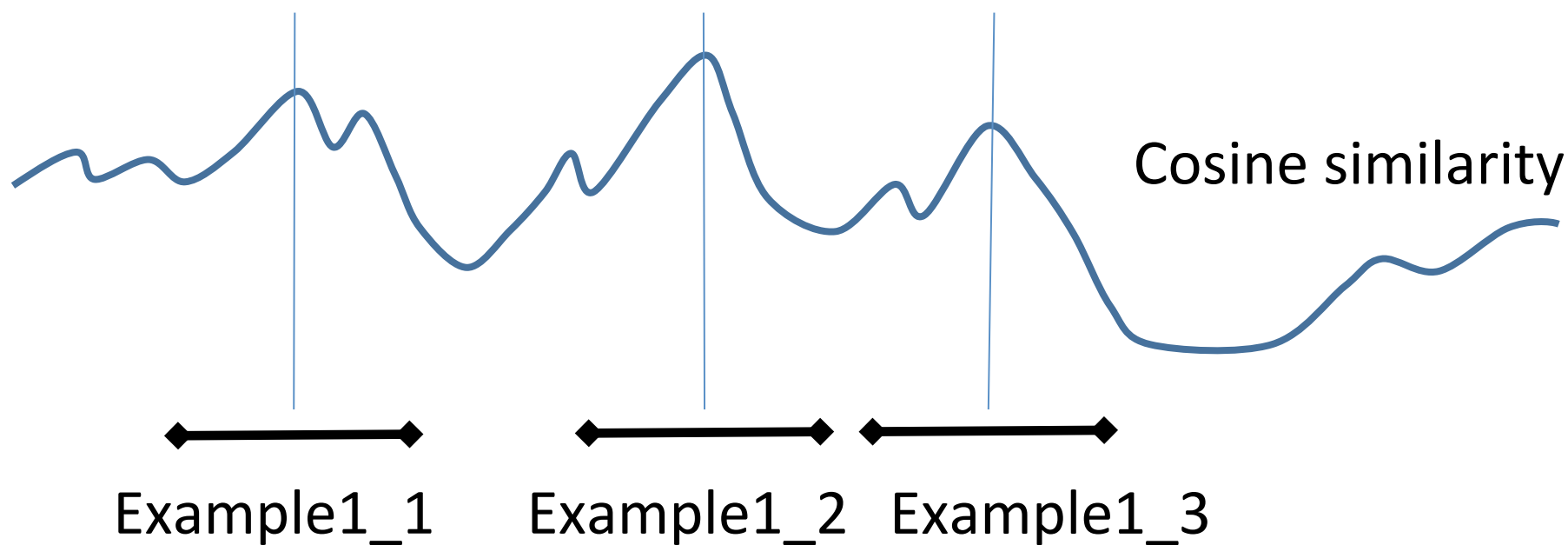
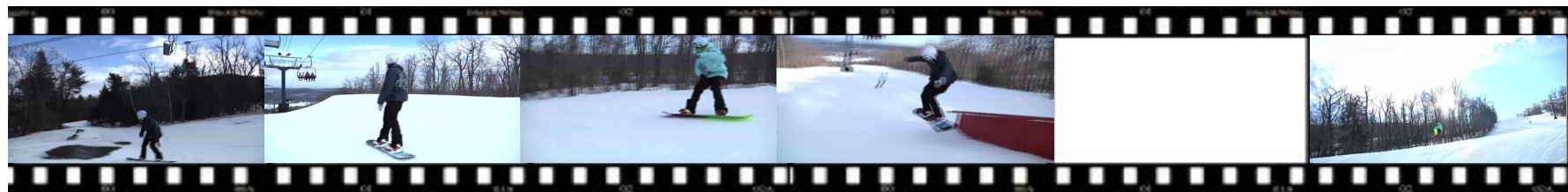
Terms

Pre-train representation
on webly-supervised videos

.. detectors selected for
generality and specificity.

... to achieve stories, even
when that class has few data.

2017 Expand training material



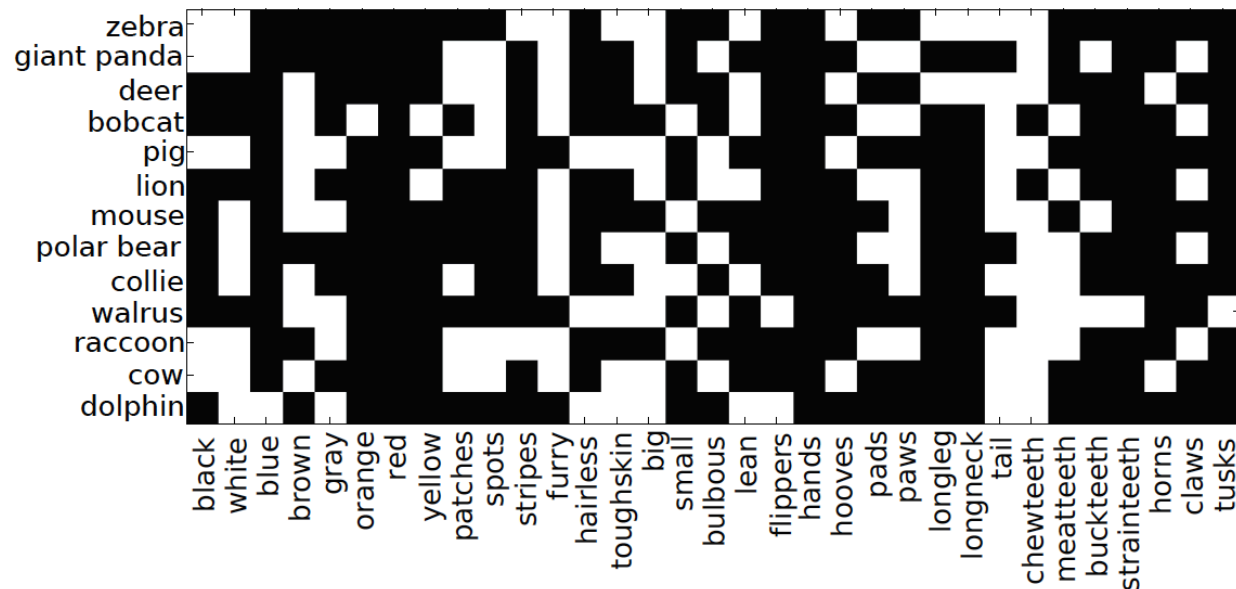
More by less data

RETRIEVAL

Every question is new, so classification not for every day.

Event retrieval: is it zero-shot?

Zero shot aims to classify test videos by predefined mutual relationship using class-to-attribute mappings



We aim for a new event by a text only.

Event recognition without examples

The goal is event recognition without examples, ruling out using videos as histograms or other features.

Event Name: Winning
Definition: An individual or team achieving a specific goal, often in a competitive context. The event is generally defined by the purpose of the activity and the outcome. The only type of race for the purpose of this event is the type of race. The event is completed by a person or animal. Different types of races involve different types of human ...

event
description

term
extraction

term vector



test
videos

**video to
lingual**

term vector

objective, subjective,
cultural, mood, framing

objective, subjective,
cultural, mood, framing

Nouns are easy, propositions are not



Can we have a vote?



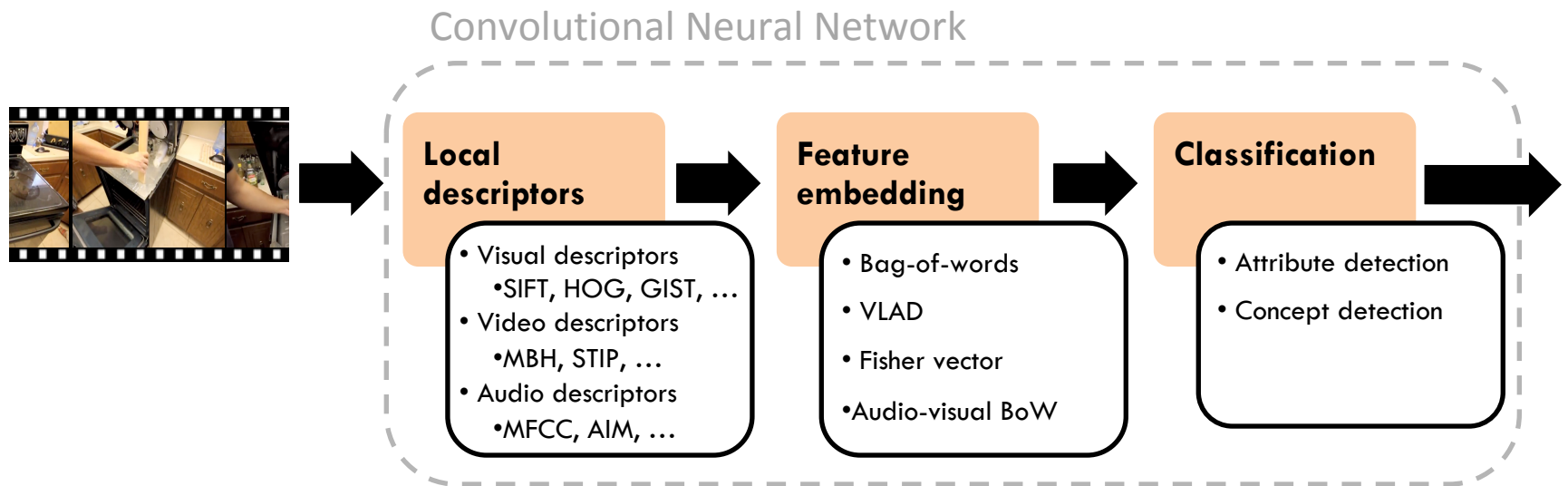
Nouns are stable, adjectives personal



Old is visually different for every notion.

Concept embedding for retrieval without examples

Representing videos as histograms of concept scores



Problem: one classifier against the complexity of the world.

CONCEPT EMBEDDING

Concept embedding label expansion

Expanding the labels by logical combinatorics,

	Ride	Motorcycle	Bike	Bike-AND-Ride	Bike-OR-Motorcycle	Concept Annotations
	0	0	1	0	1	
	1	0	1	1	1	
	1	1	0	0	1	

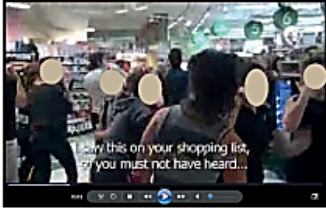


Label expansion expands the vocabulary for free:

bike .and. road for *bicycle trick*, not bike .or. road.

Concept embedding qualitative results

Top ranked videos for *flash mob gathering*.

Most important concepts in their video representation

Detected Videos	Composite Concepts
	Group-AND-Dance-AND-Shopping Celebrating-OR-Marching Performance-OR-Music People-OR-Girl Surprise-OR-Party
	Group-AND-Dance-AND-Shopping Band-OR-Singining Inside-OR-School Performance-OR-Music Surprise-OR-Party
	Group-AND-Dance -AND-Shopping Practice-OR-Gym Living-AND-Room Street-OR-Inside Performance-OR-Music

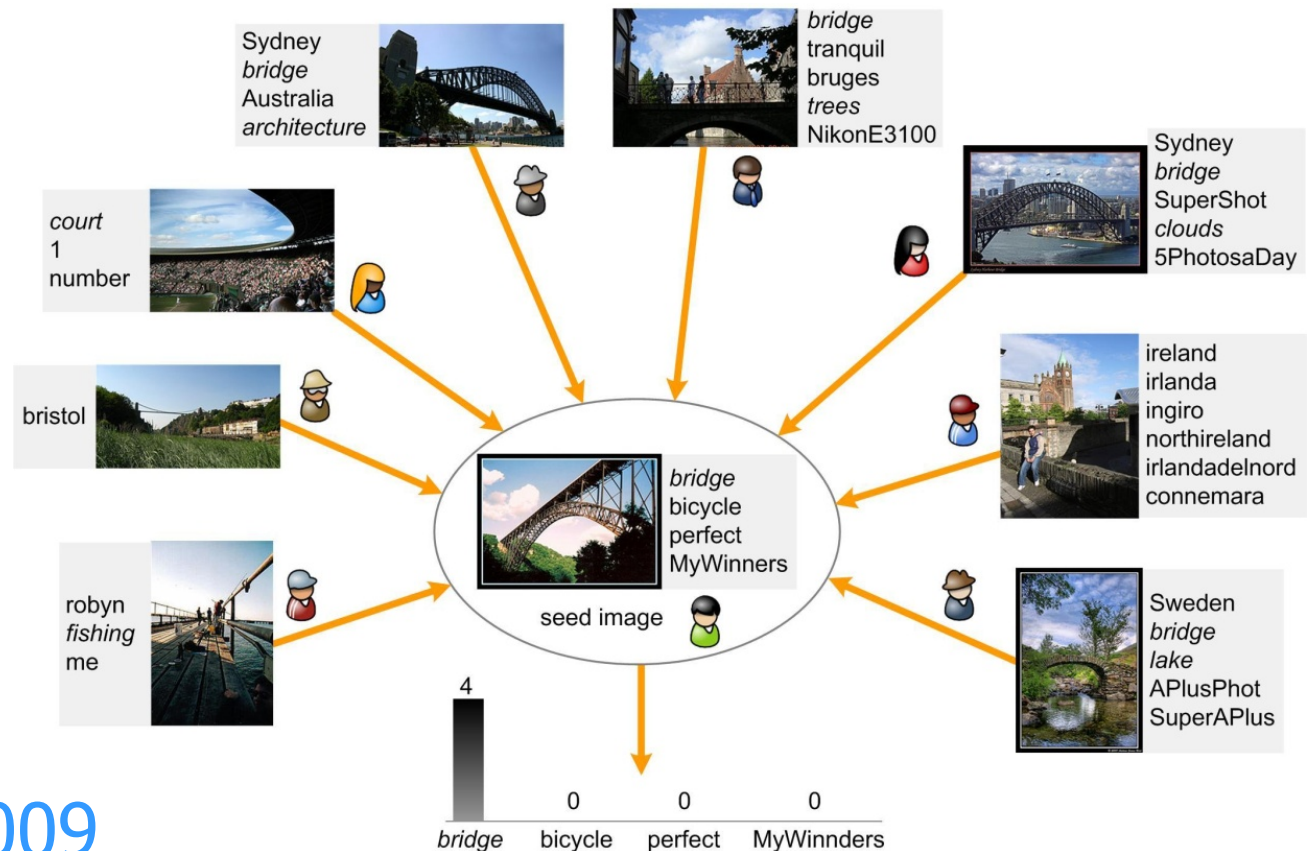
***Still need a labeled basis
for each concept classifier.***

VIDEO TO TAG-TERM EMBEDDING

Embedding inspiration from tags

Embedding based on free social tagged videos only,
without the need for training any intermediate detectors.

Inspired by:



Video2vec embedding

Can we learn the embedding from videos and their stories?

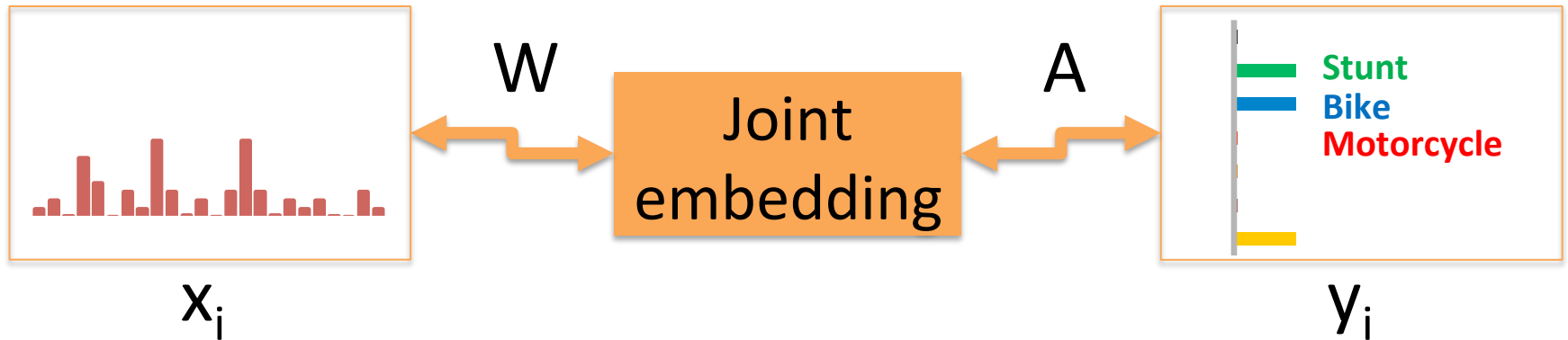
Video



Story

Story usually highlights the key concepts in video jointly.
Videos and stories are freely available on YouTube.

Video2vec embedding



Joint space where $x_i W \approx y_i A$.

Explicitly relate training W and A from multimedia.

W = Visual projection matrix individual term classifiers

A = Textual projection matrix select/group terms

Video2vec embed the video story



Learn W and A such that *descriptiveness* preserves video descriptions and *predictability* recognizes terms from video content

Video2vec key observation



By grouping terms, the number of classes is reduced.
For training classifiers, more positives needed per group.
We can train from freely available web data.

Video2vec joint optimization

S is (the size of) the embedding

L_d Loss function for descriptiveness.

L_p Loss function for predictability.

$$L_{VS}(\mathbf{A}, \mathbf{W}) = \min_{\mathbf{S}} L_d(\mathbf{A}, \mathbf{S}) + L_p(\mathbf{S}, \mathbf{W})$$

Jointly optimize descriptiveness and predictability.

Video2vec objective descriptiveness

The Video2vec embedding should be descriptive.

$$L_d(\mathbf{A}, \mathbf{S}) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{y}_i - \mathbf{A}\mathbf{s}_i\|_2^2 + \lambda_a \Omega(\mathbf{A}) + \lambda_s \Psi(\mathbf{S})$$

Original transcriptions

Reconstructed terms

Regularizers

Latent semantic indexing with L2 norm.

Video2vec objective predictability

The Video2vec embedding should be predictable.

$$L_p(\mathbf{S}, \mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{s}_i - \mathbf{W}^\top \mathbf{x}_i\|_2^2 + \lambda_w \Theta(\mathbf{W})$$

Video2vec embedding

Video feature embedding

Regularizer

Video2vec 46K dataset

Videos and title descriptions from higher quality YouTube,
46K videos, 19K terms in description.

Features x_i any combination.

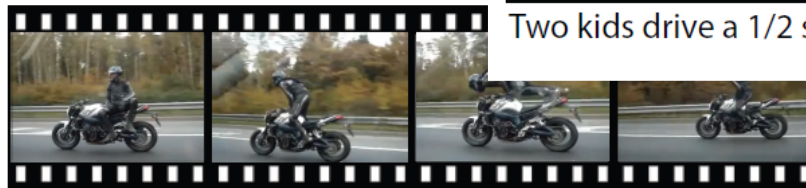
Seeded from video event descriptions y_i in bags.



Cute tabby cat gives her dog a bath



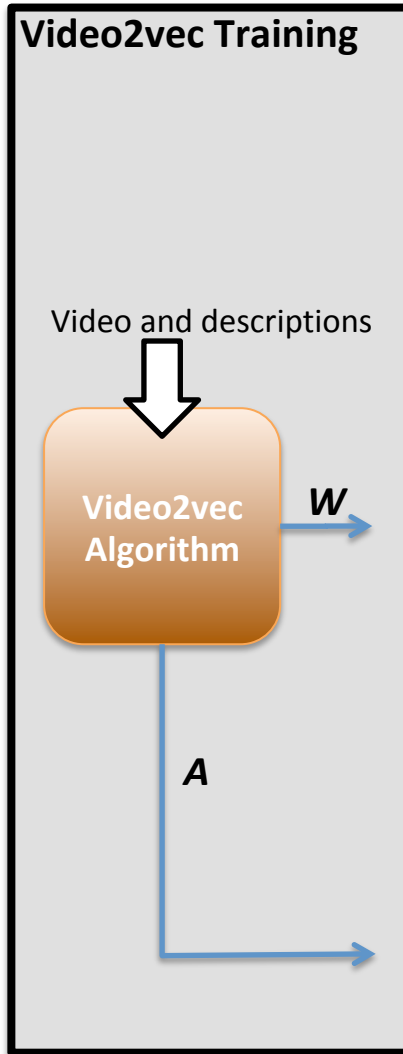
Two kids drive a 1/2 size Jeep through mud



Crazy guy doing insane stunts on bike.

Available for download: www.mediamill.nl

Video2vec training method



Stochastic Gradient Descent starting from a random sample.

The sample gradient wrt objective is:

$$\nabla_{\mathbf{A}} L_{VS} = -2 (\mathbf{y}_t - \mathbf{A} \mathbf{s}_t) \mathbf{s}_t^\top + \lambda_a \mathbf{A},$$

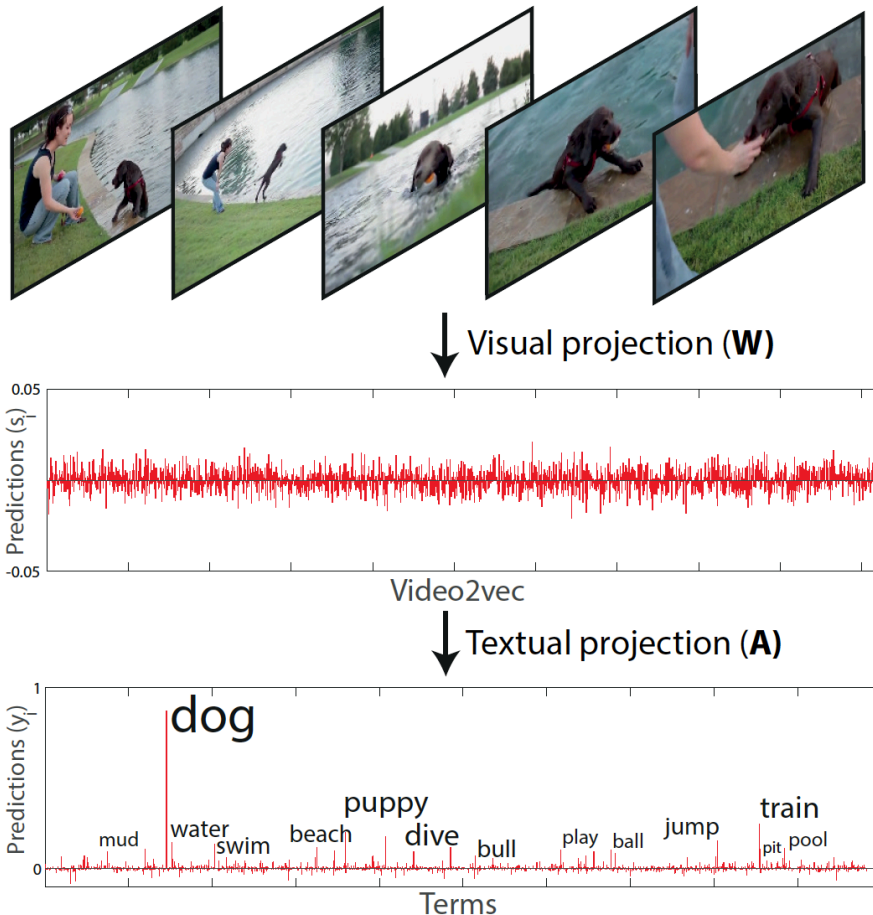
$$\nabla_{\mathbf{W}} L_{VS} = -2 \mathbf{x}_t \left(\mathbf{s}_t - \mathbf{W}^\top \mathbf{x}_t \right)^\top + \lambda_w \mathbf{W}, \text{ and}$$

$$\nabla_{\mathbf{s}_t} L_{VS} = 2 \left[\mathbf{s}_t - \mathbf{W}^\top \mathbf{x}_t - \mathbf{A}^\top (\mathbf{y}_t - \mathbf{A} \mathbf{s}_t) \right] + \lambda_s \mathbf{s}_t.$$

Update parameters with step-size η .

Start \mathbf{A} and \mathbf{S} from SVD of term vectors \mathbf{Y} .

Video2vec at work



1. Project visual features

$$s_i = W^T x_i,$$

2. Translate to text

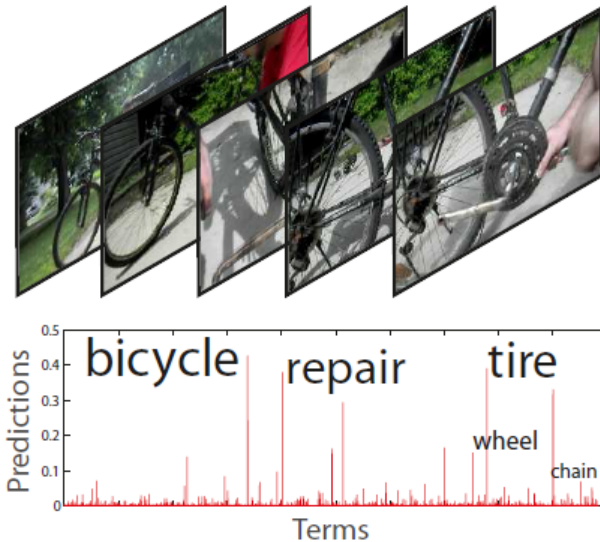
$$\hat{y}_i = A s_i,$$

3. Cosine distance match

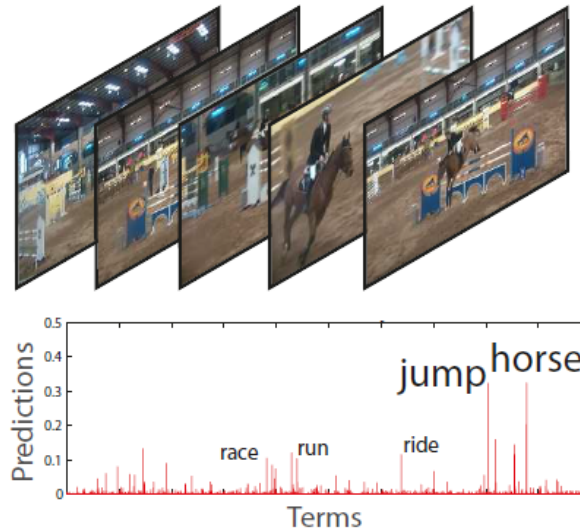
$$s_e(x_i) = \frac{y^e{}^T \hat{y}_i^e}{\|y^e\| \|\hat{y}_i^e\|}$$

Video2vec predicted terms

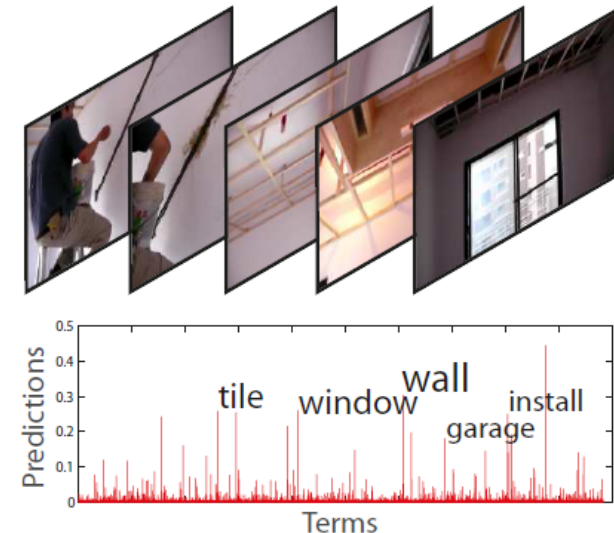
non-motorized vehicle repair



horse riding competition

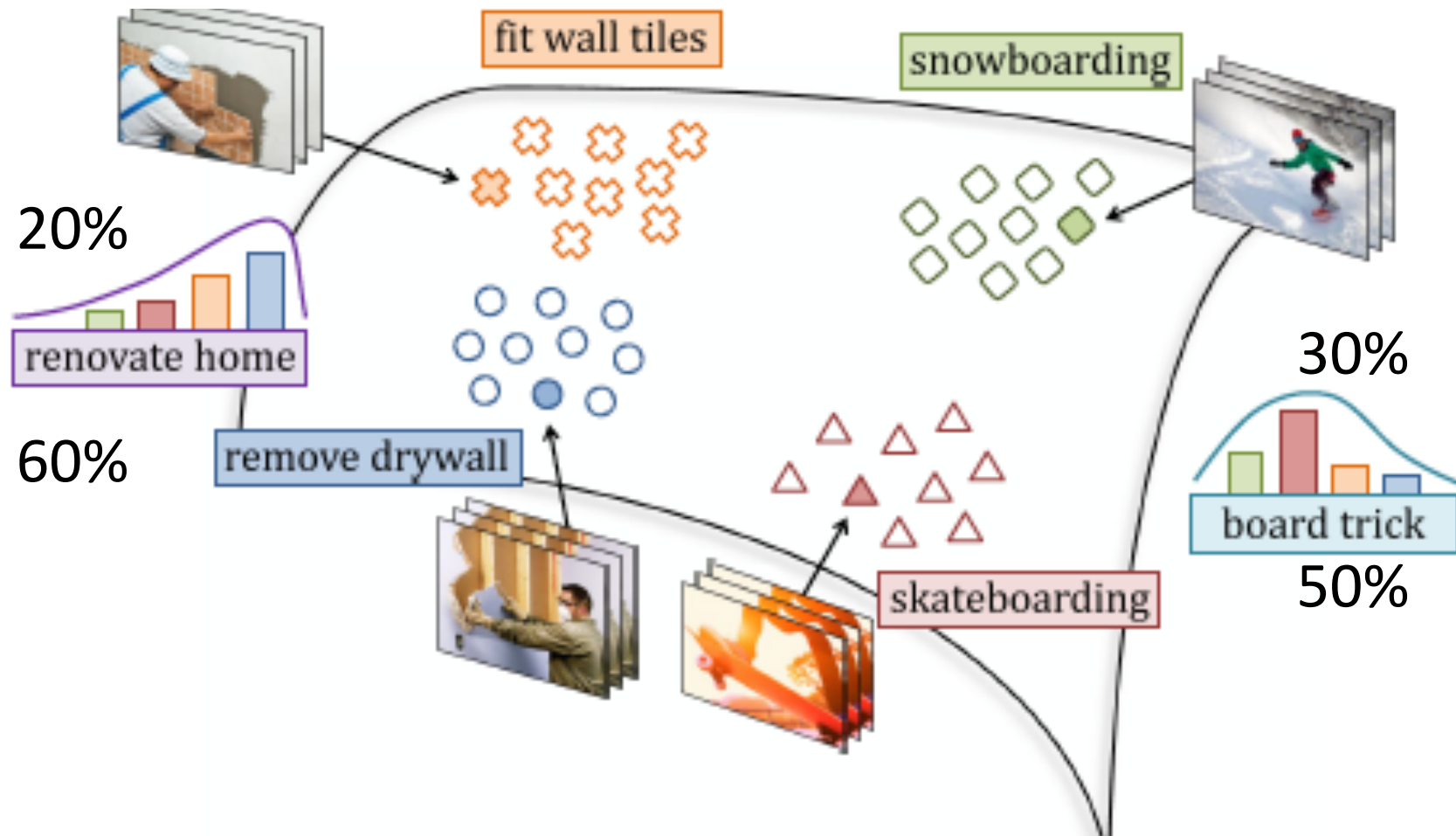


renovating a home



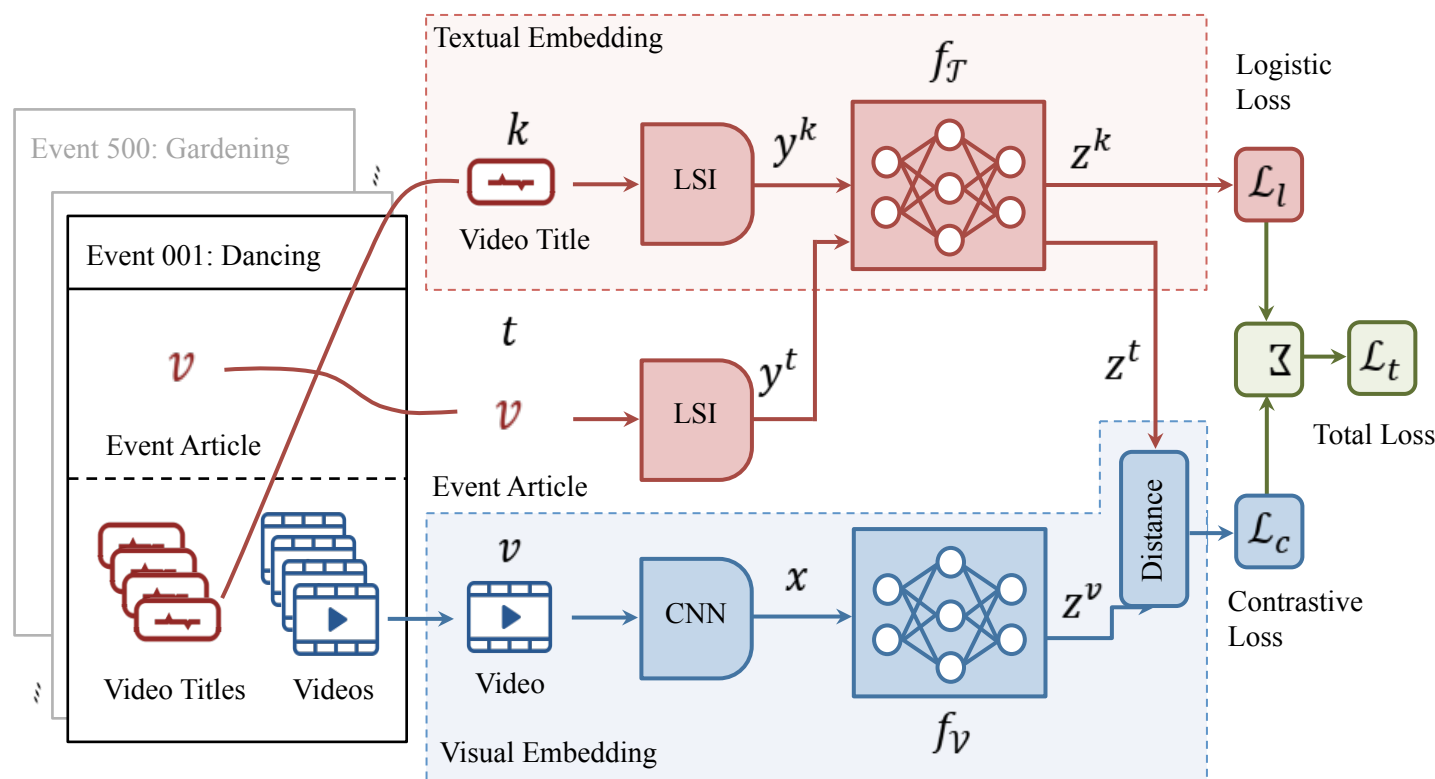
UNIFIED METRIC EMBEDDING

Unified metric embedding

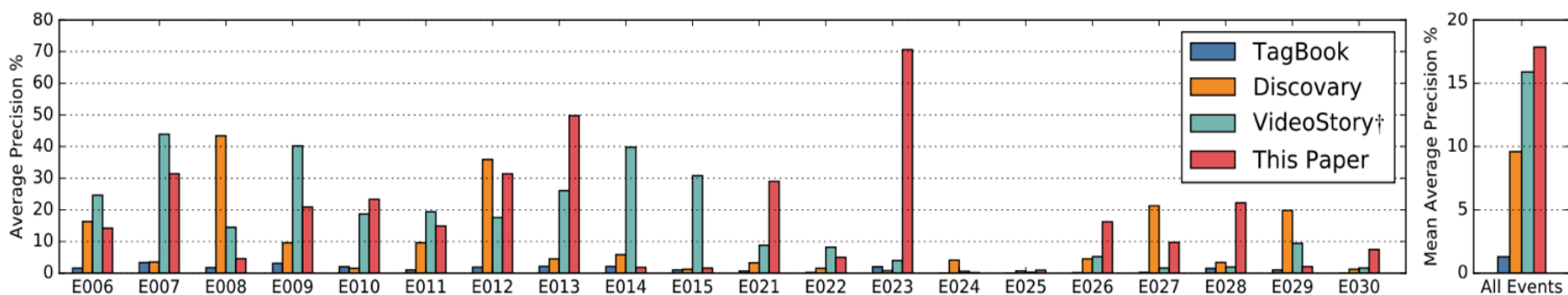


Unified metric embedding

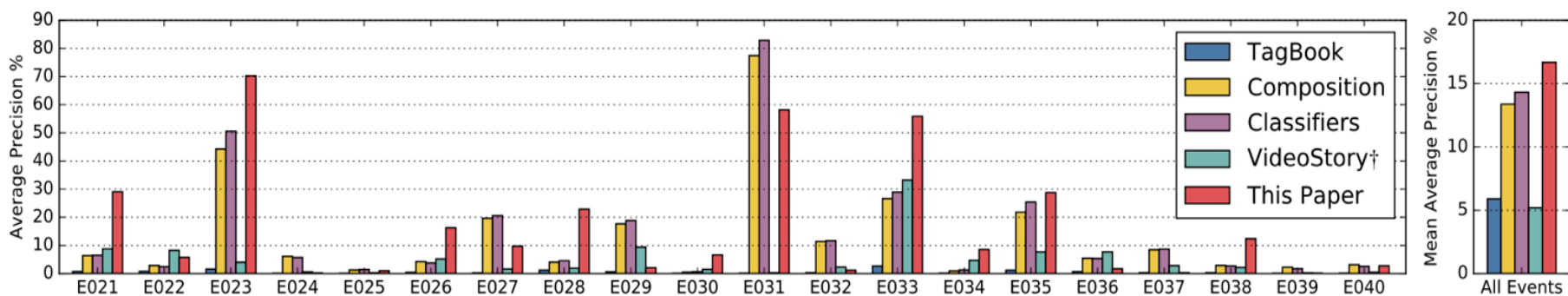
Zero-exemplar is learning from pre-defined events plus novel ones as a probability over the existing events.



Unified metric embedding quantitative



(a) MED-13 Dataset



(b) MED-14 Dataset

Unified metric embedding qualitative

Renovating home improve a home by rebuilding parts of the structures.

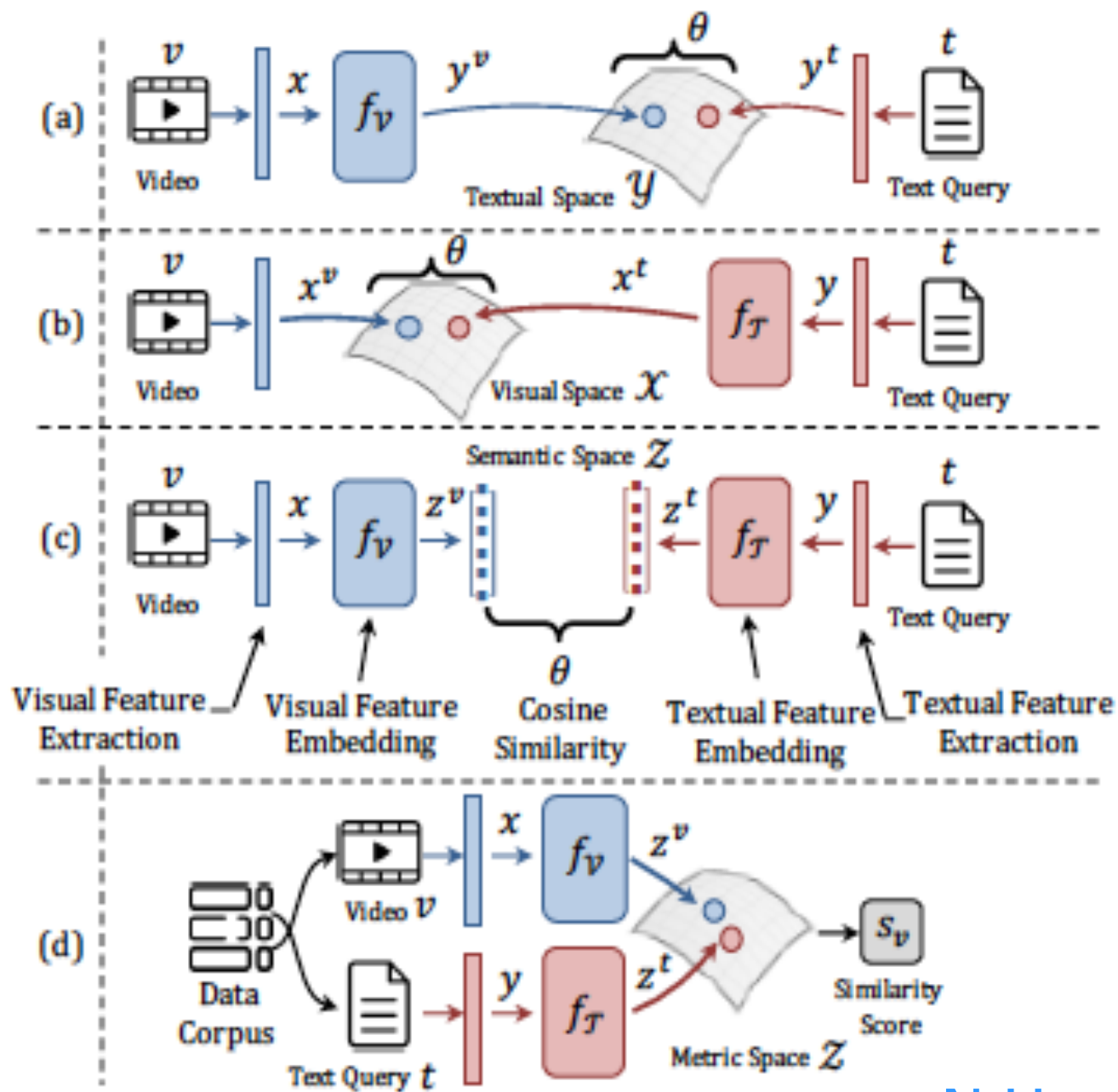
Success



Failure



EVENT RECOGNITION BY EMBEDDING



Retrieval by embedding results

Authors		Published	mAP
Habibian et al.	concept embedding	ICMR 2014	6.4
Ye et al.		MM 2015	9.0
Mazloom et al.		ICMR 2015	11.9
Wu et al.		CVPR 2014	12.7
Jiang et al.		AAAI 2015	12.9
Mazloom et al.	tag embedding	TMM 2016	12.9
Liang et al.	big data & reranking	MM 2015	18.3
Habibian et al.	joint embedding	TPAMI 2017	20.0
Hussein et al.	unified metric embed	CVPR 2017	17.9

N.Hussein CVPR2017 A.Habibian PAMI 2017

OTHER CHALLENGES

In the kitchen of the future.

TRECVID SURVEILLANCE EVENTS > TRECVID ACTIVITIES EXTENDED VIDEO

slides by Jon Fiscus (NIST)

COMPLEX ACTIVITIES IN VIDEO

at the UvA by N.Hussein, S. Gavves, C. Snoek *others*

Multi-cam surveillance from text

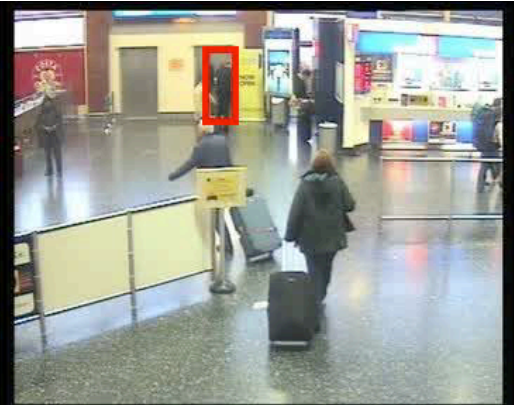
Controlled
Access Door ①



② Waiting Area



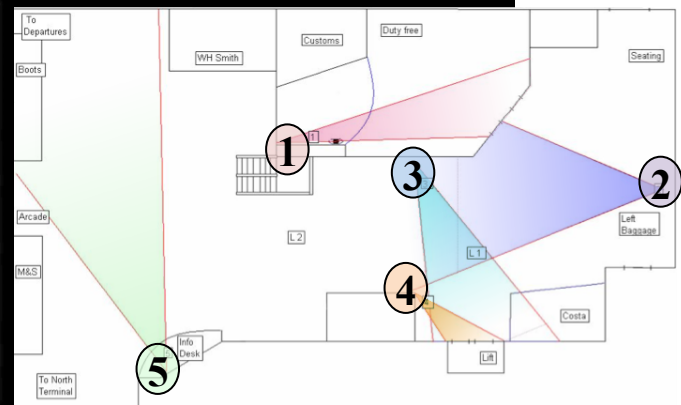
③ Debarcation Area



④ Elevator Close-Up



⑤ Transit Area



Events of Interest



Single Person events

PersonRuns

Someone runs

Pointing

Someone points

Single Person + Object events

CellToEar

Someone puts a cell phone to his/her head or ear

ObjectPut

Someone drops or puts down an object

Multiple People events

Embrace

Someone puts one or both arms at least part way around another person

PeopleMeet

One or more people walk up to one or more other people, stop, and some communication occurs

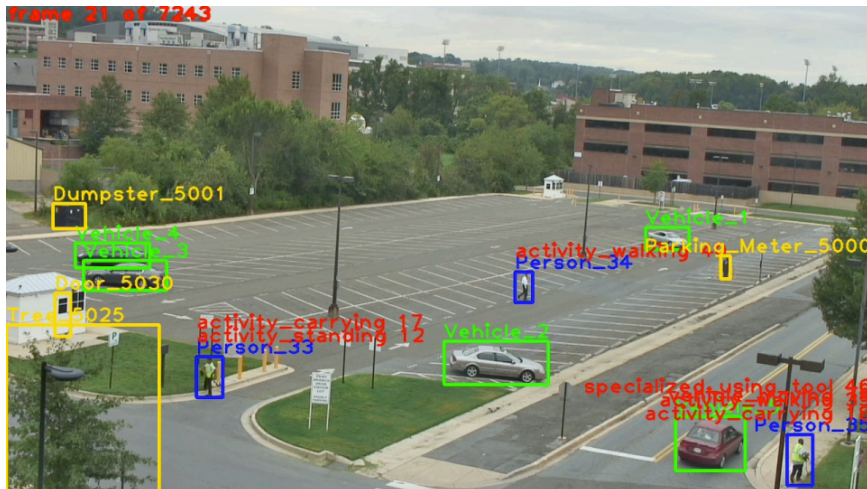
PeopleSplitUp

From two or more people, standing, sitting, or moving together, communicating, one or more people separate themselves and leave the

ActEV new task per 2018

Successor of Surveillance Event Detection by adding a large collection of multi-camera video data, both of simple and complex activities.

ActEV will address activity detection for both forensic applications and for real-time alerting.



Recognizing complex tasks

Strong temporal models are no longer valid.



This depicts *cooking food* regardless frame order.

