

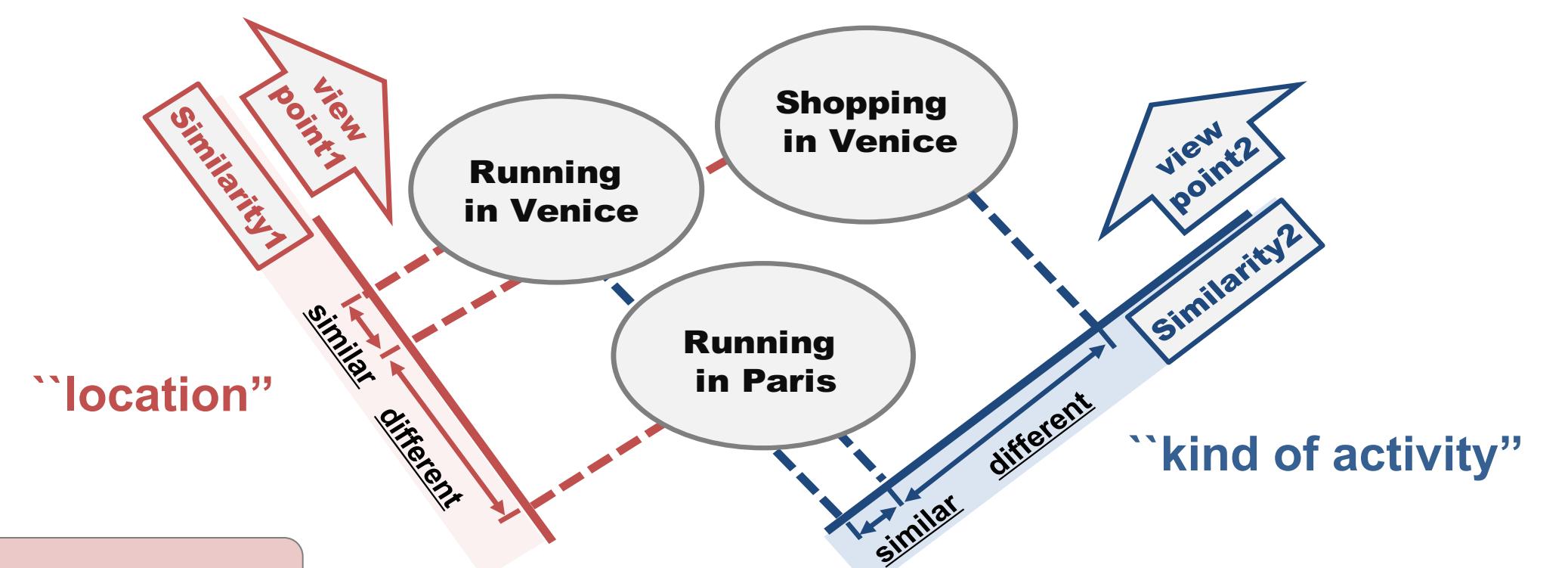
Viewpoint-aware Video Summarization

Atsushi Kanehira¹, Luc Van Gool^{3,4}, Yoshitaka Ushiku¹, Tatsuya Harada^{1,2}

¹ The University Tokyo, ² RIKEN, ³ ETH Zurich, ⁴ KU Leuven

Introduction:

- Viewpoint changes how the video looks.
- This work treats *viewpoint* as ...
 - ✓ A specific aspect of a video
 - e.g., kind of activity, location, ...
 - ✓ Different viewpoint has different semantic similarity.



Basic question

Can we estimate *viewpoint* from given *similarity*?

What is used as *similarity*?

- Pre-determined video groups (easy to obtain, such as by user preference)

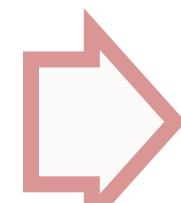
How to show output ?

- Summarizing video (easy to see, contain much information)

Goal

Generating video summary that explains ``why videos are divided in the way they are?''

1. Why the video are in that group?
 - summary of target and other videos in the same group
2. Why the video are not in other groups?
 - summary other videos in the different groups



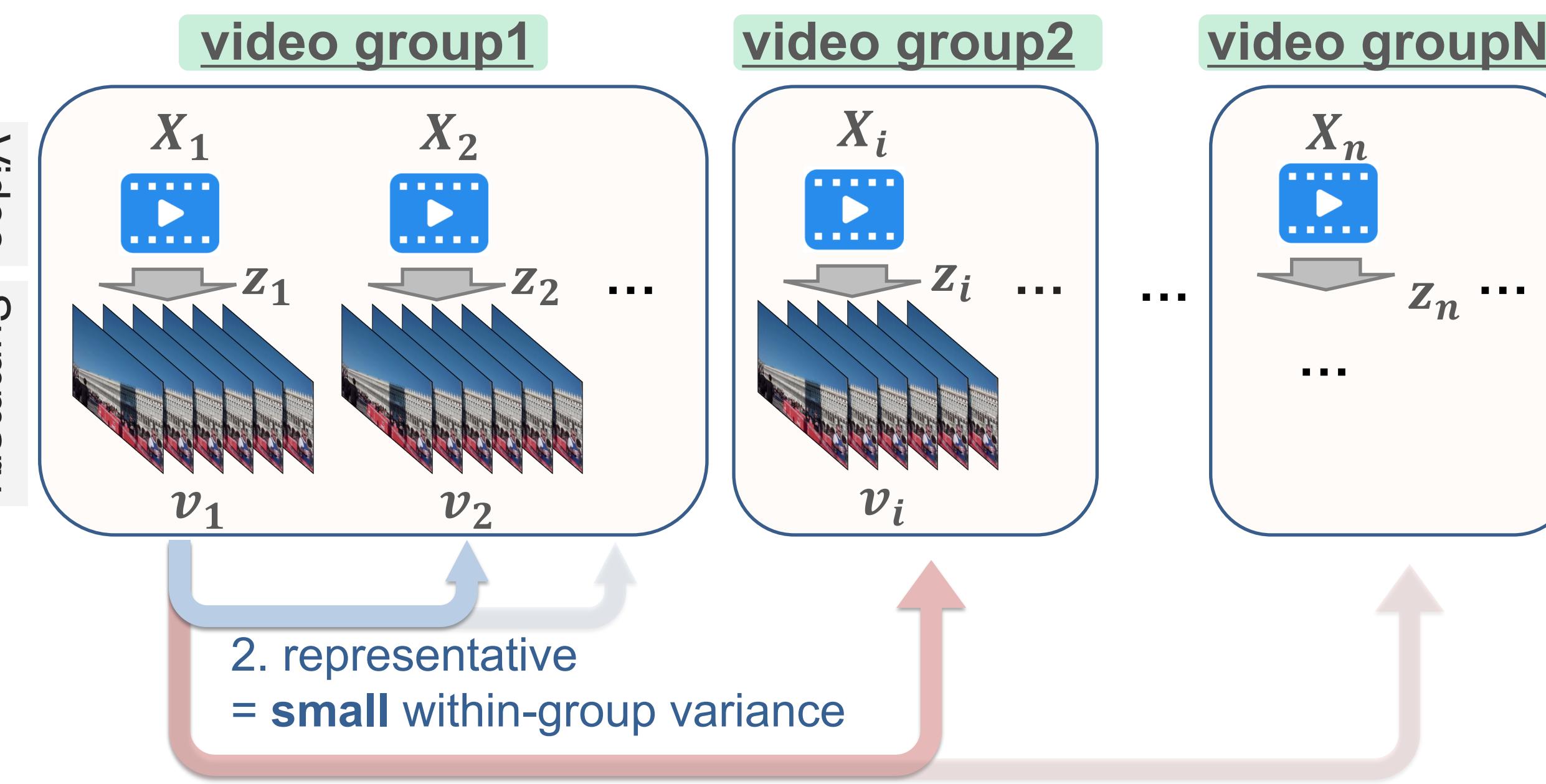
Video Co-summarization

Method:

Video co-summarization reflecting *similarity*

Requirements:

1. **diverse**
2. **representative** to videos in the same group
3. **discriminative** against videos in the different groups



- Visual feature of summaries with shot-indicator $\mathbf{z} = \{0, 1\}^N, \|\mathbf{z}\|_0 = s$

$$\underline{\mathbf{v}} = \mathbf{X}\mathbf{z} \quad \text{summary-feature} \quad \mathbf{X} = [\underline{\mathbf{x}_1}, \underline{\mathbf{x}_2}, \dots, \underline{\mathbf{x}_N}]^T \quad \text{shot-features}$$

- Minimizing the combination of three variances inspired by Fisher Criteria

$$\lambda_1 \text{Tr}(S^V) - \lambda_2 \text{Tr}(S^W) + \lambda_3 \text{Tr}(S^B)$$

inner-summary within-group between-group
 ↓ ↓ ↓
 diverse representative discriminative

- Continuously relaxed problem is DC (difference of convex) programming, which can be optimized with CCCP (concave-convex procedure)

Dataset:

- Collect a set of videos using 3 keywords.

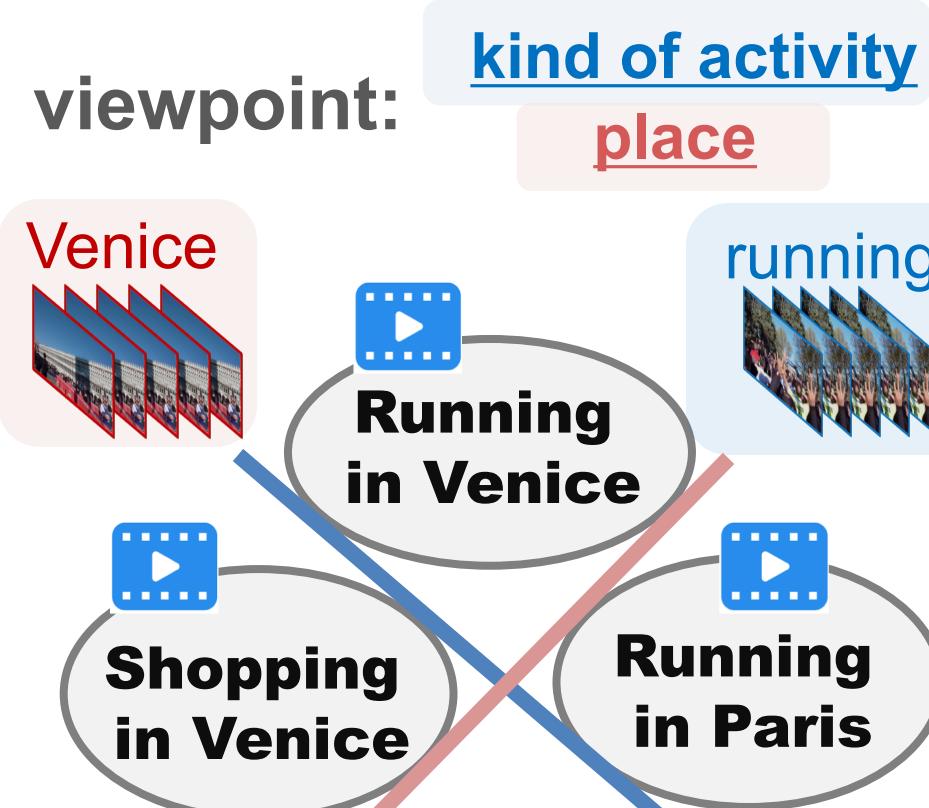
► Videos can be divided with 2 explicit viewpoints.

- Target videos have 2 different annotations.

► Each corresponds to one viewpoint.

- Compare predicted summary with ground-truth.

- 10 types of groups, 150 videos, 4 hours in total.



Examples of dataset

``riding helicopter in NewYork''
with NewYork(↑), helicopter(↓)



``riding horse in the safari''
with riding horse(↑), safari(↓)



Experiment:

Automatic evaluation (top5 mean Average Precision)

method	only representative			representative & discriminative		
	SMRS	CK	CS	MBF	CVS	WSVS
mAP	0.322	0.294	0.312	0.333	0.353	0.322
WSVS(large)						0.319
ours						0.379
ours (feature learning)						0.385

User study

summary quality

method	MBF	CVS	ours
accuracy	0.47	0.60	0.76

explanation quality

method	MBF	CVS	ours
score	1.07	1.22	1.32

example results

``racing in the desert''
with racing(↑), desert(↓)



``riding helicopter in NewYork''
with NewYork(↑), helicopter(↓)

