

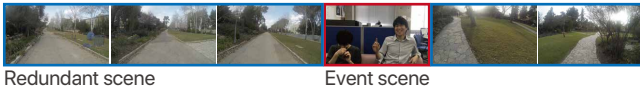
# Dynamic Object Scanning: Object-Based Elastic Timeline for Quickly Browsing First-Person Videos

Seita Kayukawa<sup>1</sup> Keita Higuchi<sup>2</sup> Ryo Yonetani<sup>2</sup> Masanori Nakamura<sup>1</sup> Yoichi Sato<sup>2</sup> Shigeo Morishima<sup>3</sup>  
<sup>1</sup>Waseda University <sup>2</sup>University of Tokyo <sup>3</sup>Waseda Research Institute for Science and Engineering

## Introduction

### First-person videos captured by wearable cameras

- Large and diverse collection of **long and untrimmed** videos
- Important events can be distributed sparsely

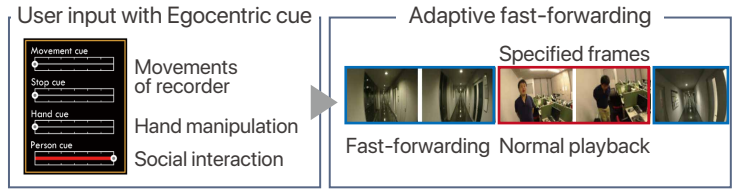


### Our Goal

To present adaptive video fast-forwarding technique that helps users to browse long and untrimmed videos quickly

## Related Work

### Adaptive Video Fast-forwarding based on User Input<sup>[1]</sup>

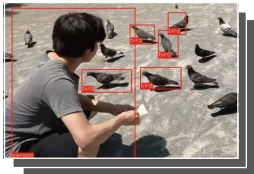


**Limitation:** The type of cues have been **fixed for any given video**  
 ▶ Limits the variety of videos that can get the benefit of Egocentric cue

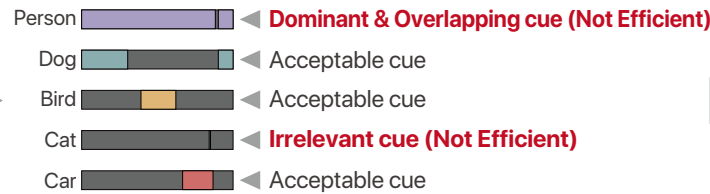
## Our Approach

**Key Idea:** To generate a efficient set of cues arranged adaptively tailored to the contents of a given video

### Object Detection (YOLOv2<sup>[2]</sup>)

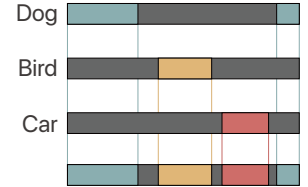


### Detected Objects



### Cue Selection

### Arranged set of Object cues



## Greedy Algorithm

### Objective function

$$F(C) = A(C) - B(C)$$

$C = \{c_1, \dots, c_N\}$ : set of detected objects

$A(C)$ : the number of frames where at least one of the categories in  $C$  is observed  
 To avoid an **irrelevant category** and a **temporally-overlapping category**

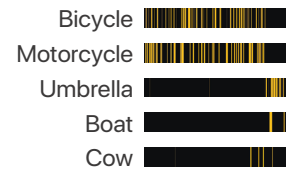
$B(C)$ : the number of frames with the object category observed most frequently  
 To avoid a **temporally-dominant category**

## Cue Selection Results

Applied our cue selection algorithm to diverse scenes



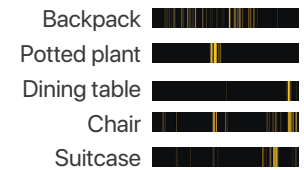
### Selected cues



### Omitted cue



### Selected cues



### Omitted cue



Pedestrians were **detected in nearly every frames** ▶ Person was **not selected** as cue

**Our algorithm selects a set of cues while excluding temporally dominant categories**

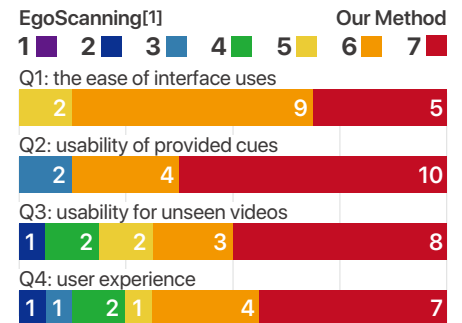
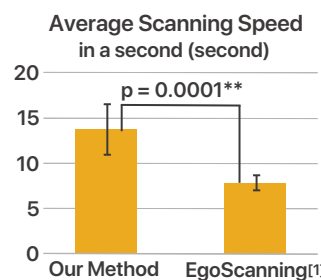
## User Study

Task: Finding predefined event from long first-person videos

Participants: 16 university students

Our Interface: Object cues selected by our method

Baseline (EgoScanning<sup>[1]</sup>): **Fixed** Egocentric cues



## Feedback

- Object cues helped several users to **infer the contents of given videos**.
- Object cues are useful to find important events since they often **emphasized only a limited part of videos**.