# Extraction of Read Text Using a Wearable Eye Tracker for Automatic Video Annotation

Mizuki Matsubara*  Joachim Folz**
Takumi Toyama** Marcus Liwicki**
Andreas Dengel**  Koichi Kise*

*Graduate school of Engineering Osaka Prefecture University
** German Research Center for Artificial Intelligence

1

# Outline

- Motivation
- Approach
- Pilot Study
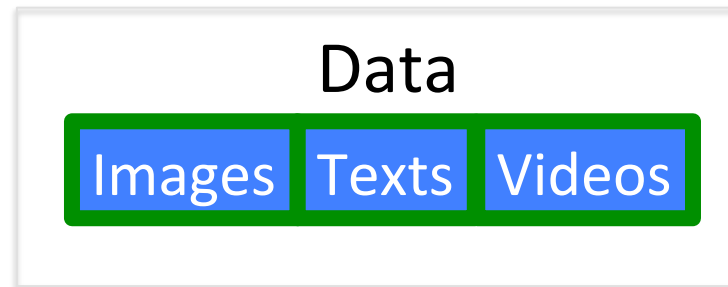- Conclusion & Future Work

# Outline

- <span style="color:red">Motivation</span>
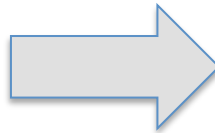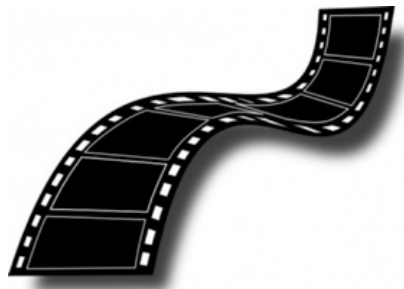- Approach
- Pilot Study
- Conclusion & Future Work

# Life-logging

- Recording a person's life





Data
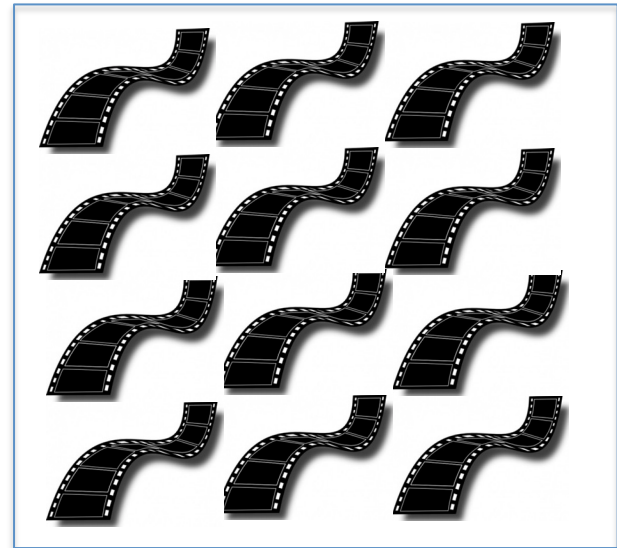
Images  Texts  Videos

# Life log video

- Life log video tends to be very long

Record everyday

Video annotation is needed to index the videos

# Automatic video annotation

Video annotations describe the video contents

Object recognition

Chair

Activity recognition

Playing soccer

Image recognition is a difficult computational problem

J. Thomason, S. Venugopalan, S. Guadarrama, K. Saenko, R. Mooney, "Integrating language and vision to generate natural language descriptions of videos in the wild." Proceedings of the 25th International Conference on Computational Linguistics (2014).

# Proposed method

Our method annotates videos in <u>a particular situation</u>

A user is following a textual manual



Worker following a manual



Cook following a recipe

<u>Texts</u> in the manual can describe the user's actions.
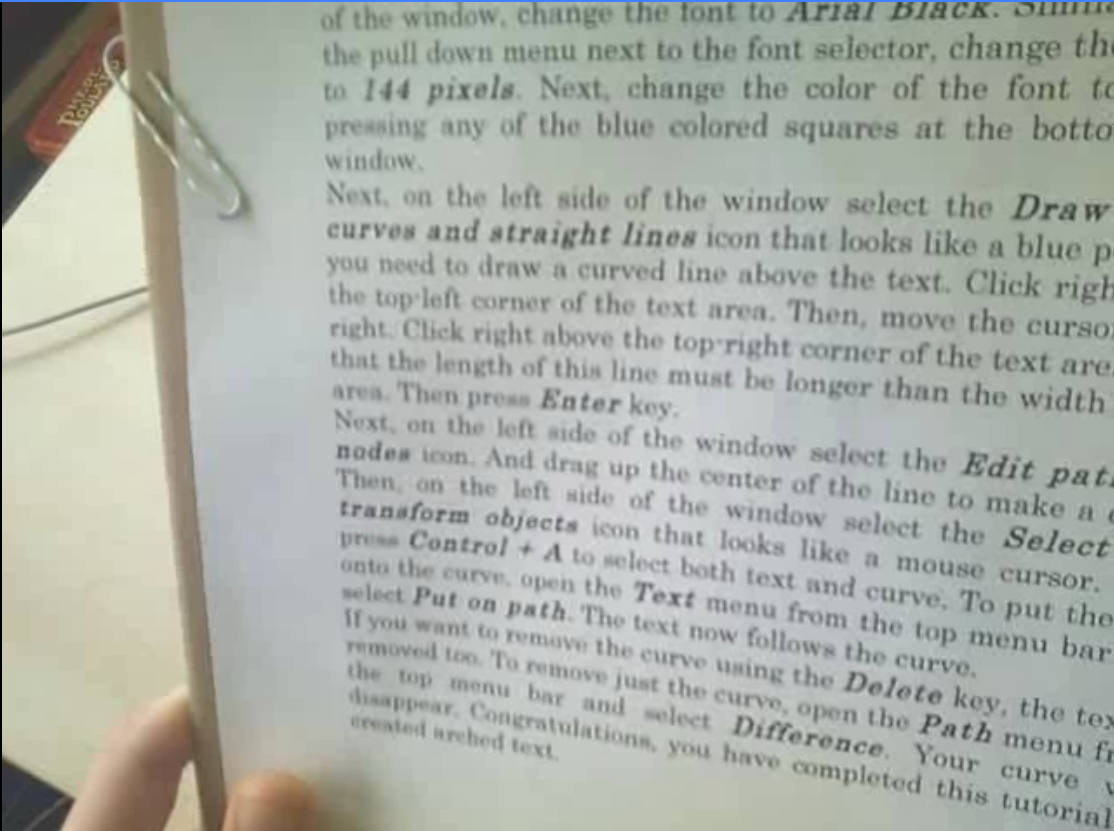
Annotation

# SMI's mobile eye tracker

- A device estimating the user's gaze



Camera

Record life log videos

# Ideal system



Draw a line above the text.

# Example 1(following a manual)

When a worker doesn't understand the manual

Video recording while working



Retrieve

Education of
a new employee

Watching the video
help him

# Example 2 (following a recipe)

When a cook forgets how to cook the meal he/she made before

Video recording while cooking



Retrieve

**Aids of memory**

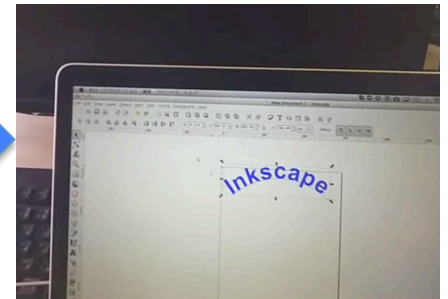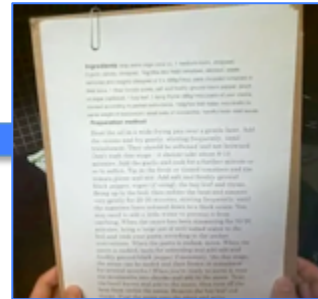The video reminds him/her the way how to cook the meal

# Outline

- Motivation
- Approach
- Pilot Study
- Conclusion & Future Work

# Assumption

Users actions can be annotated with texts read just before



Type
a sentence.

Select the illustration.
Press Enter key.

# Assumption

| Reading | Non-reading | Reading | Non-reading |
|---------|-------------|---------|-------------|

Document

Video

# Retrieval of document image

- LLAH (Locally Likely Arrangement Hashing)



Video → Retrieval → Document image

Kai Kunze, Hitoshi Kawaichi, Kazuyo Yoshimura, and Koichi Kise, Towards inferring language expertise using eye tracking, CHI'13 Extended Abstracts on Human Factors in Computing Systems, 217–222, 2013.

# Retrieval of document image

- LLAH (Locally Likely Arrangement Hashing)

Project

Retrieval

Video                              Document image

Kai Kunze, Hitoshi Kawaichi, Kazuyo Yoshimura, and Koichi Kise, Towards inferring language expertise using eye tracking, CHI'13 Extended Abstracts on Human Factors in Computing Systems, 217–222, 2013.

# Reading detection

We assume that the user is reading the document while his/her gaze point is on the document

Gaze

Reading

Noise
- influence of the user's blink
- error of the eye tracker
- failure of retrieving the document image

# Assumption



Document

| Reading | Non-reading | Reading | Non-reading |

Video

# Extraction of read text

# Outline

- Motivation
- Approach
- <span style="color:red">Pilot Study</span>
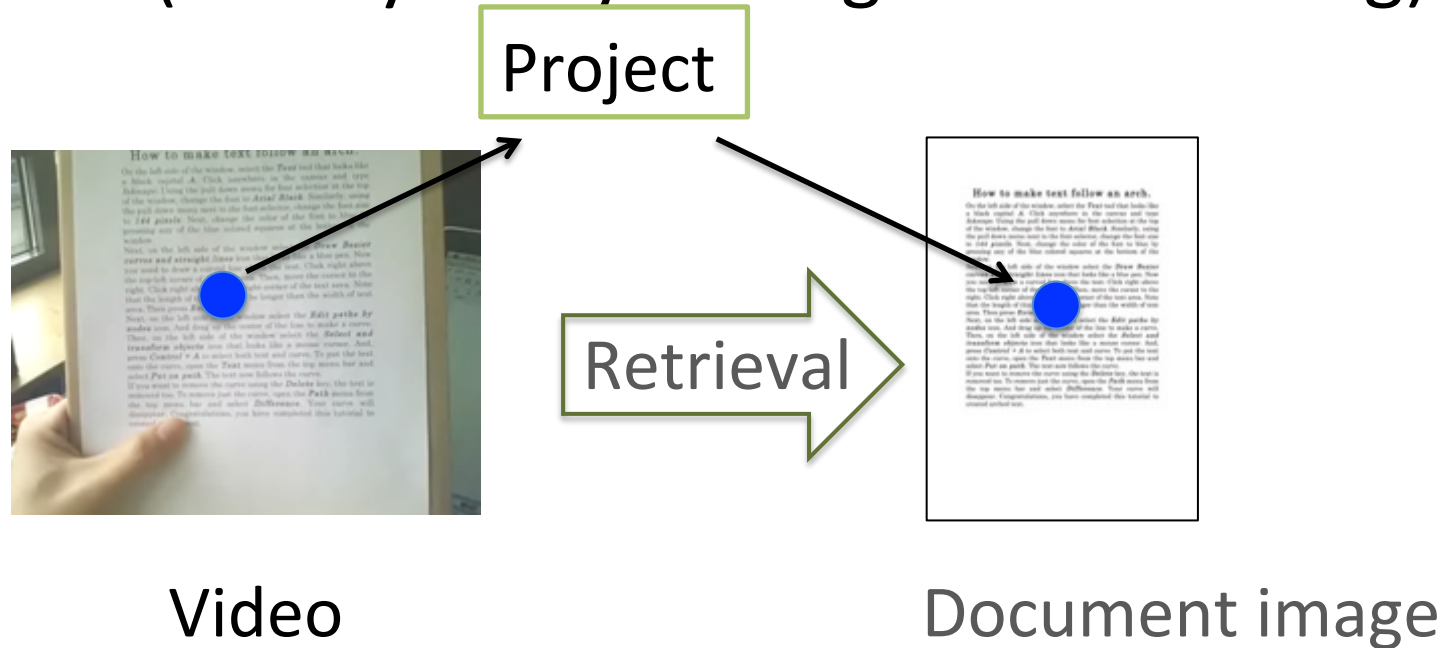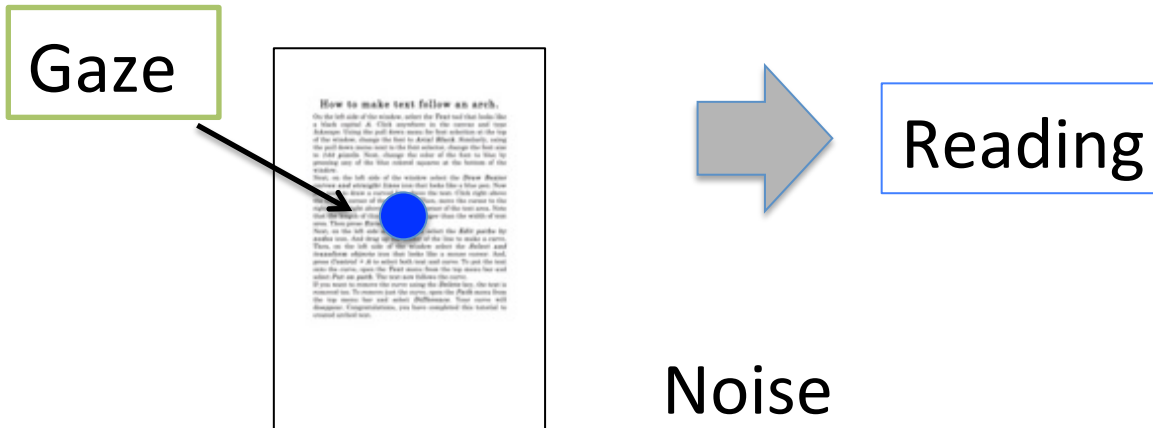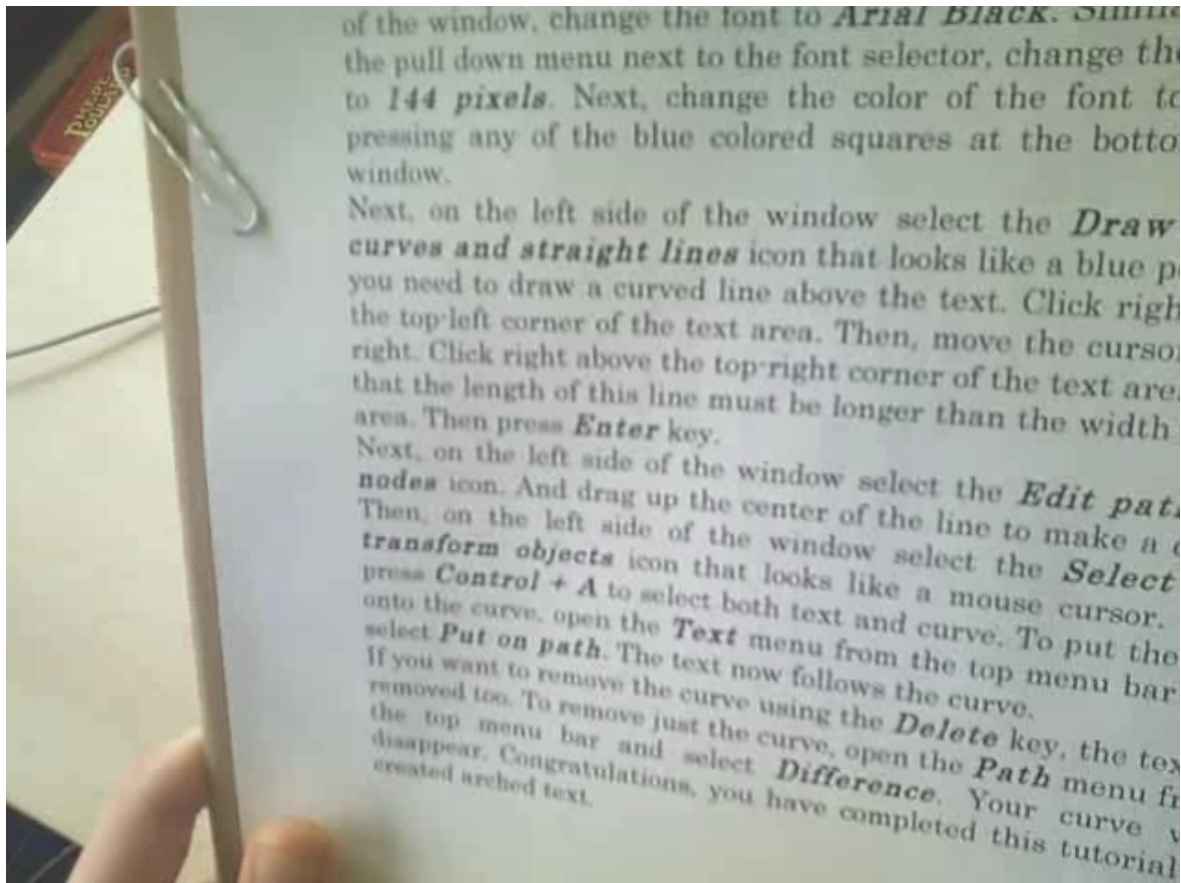- Conclusion & Future Work

# Task

- Follow a tutorial
  - Draw an illustration

# Experimental condition

- Participant : 5 (20〜30's males and females)
- Recording time : 345, 458, 393, 370, 354 [seconds]
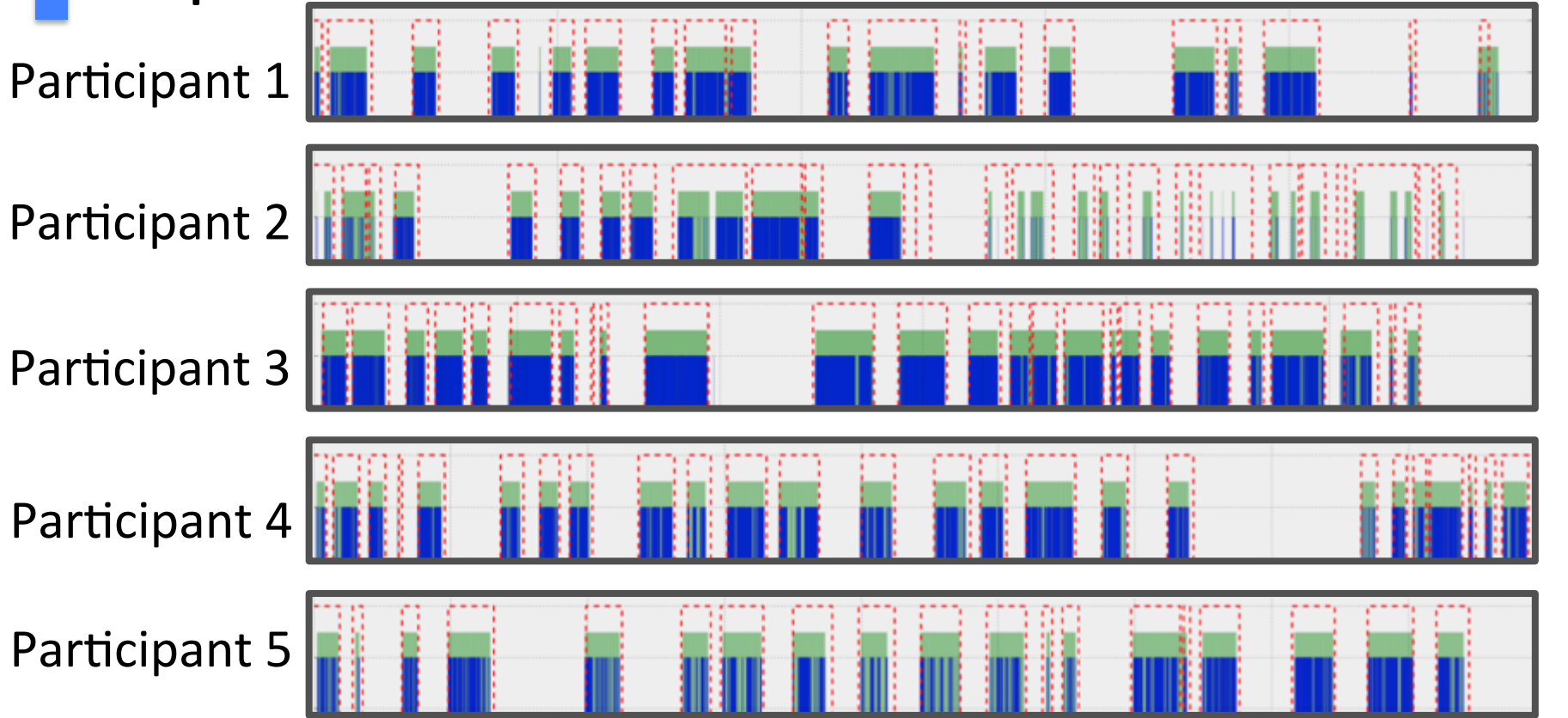- Language : English
- One page

# Evaluation of the reading detection

- We applied the reading detection to the recording videos.
  - examine the accuracy of the reading detection
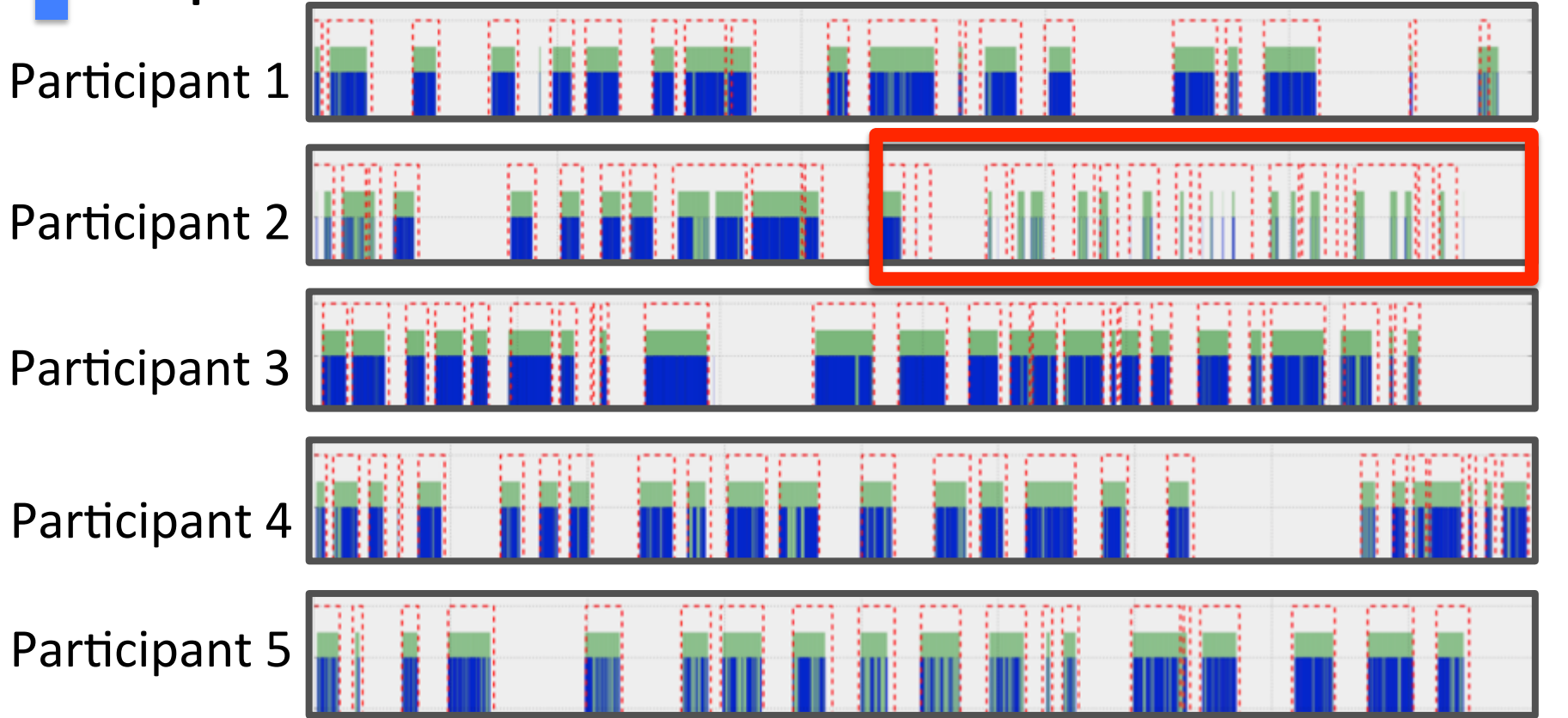
# Experimental result

Participant 1

Participant 2

Participant 3

Participant 4

Participant 5

Time [frame]

—— LLAH          After Smoothing          ----- Ground Truth

# Experimental result



Participant 1

Participant 2

Participant 3

Participant 4

Participant 5
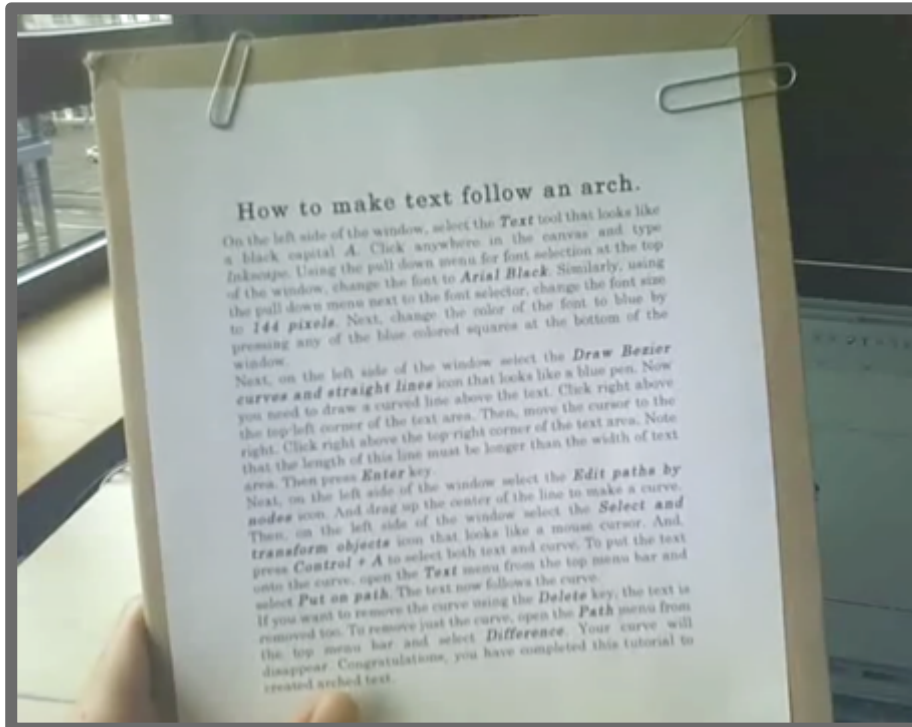
Time [frame]
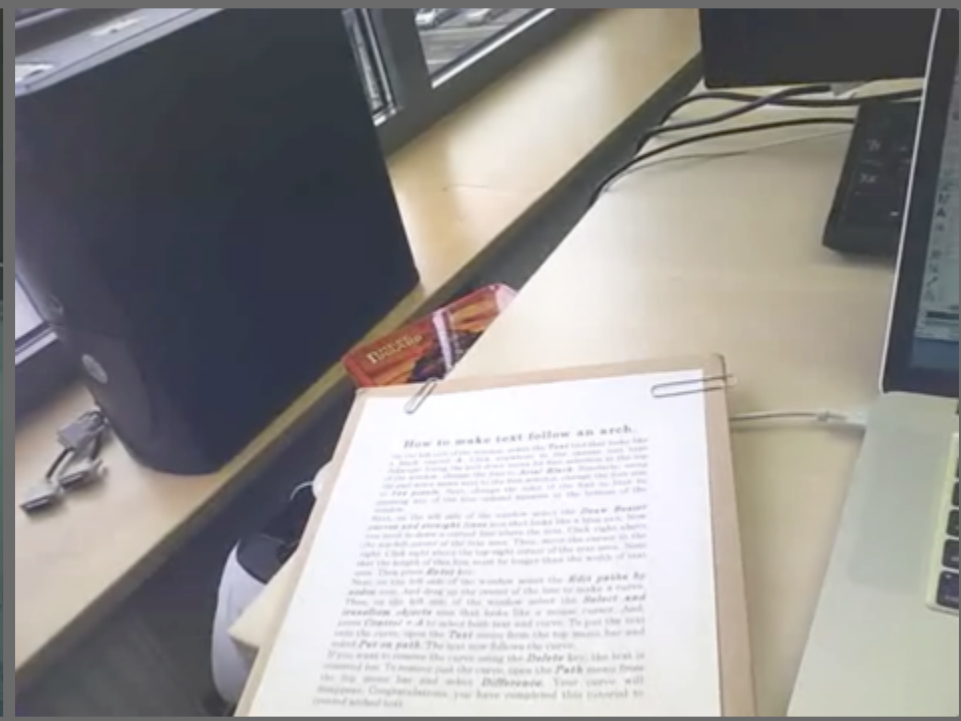
—— LLAH    After Smoothing    ----- Ground Truth

# Examples of recording frames

Success                                              Failure



The distance and angle influence the retrieval of LLAH

# Accuracy of the reading detection

Recall

$$\frac{\text{successfully detected reading frames}}{\text{Ground truth frames}}$$

$$\doteq 80.0\%$$

Precision

$$\frac{\text{successfully detected reading frames}}{\text{Retrieved document frames}}$$

$$\doteq 100\%$$

# Pilot study

Assumption

Users actions can be annotated with texts read just before

Purpose

Examine how our assumption works while participants are following a manual naturally

# Experimental condition

- Participant : 5 (20～30's males and females)
- Recording time : 345, 458, 393, 370, 354 [seconds]
- Language : English
- One page
- Rule out video segmentation error

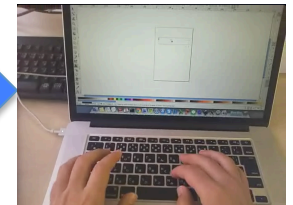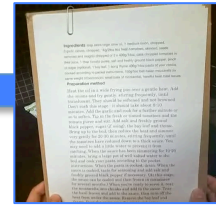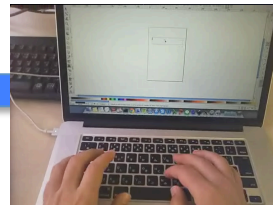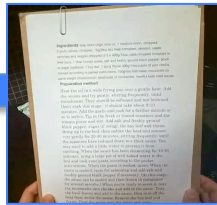# Correct the video segmentation

Correct to "Reading"

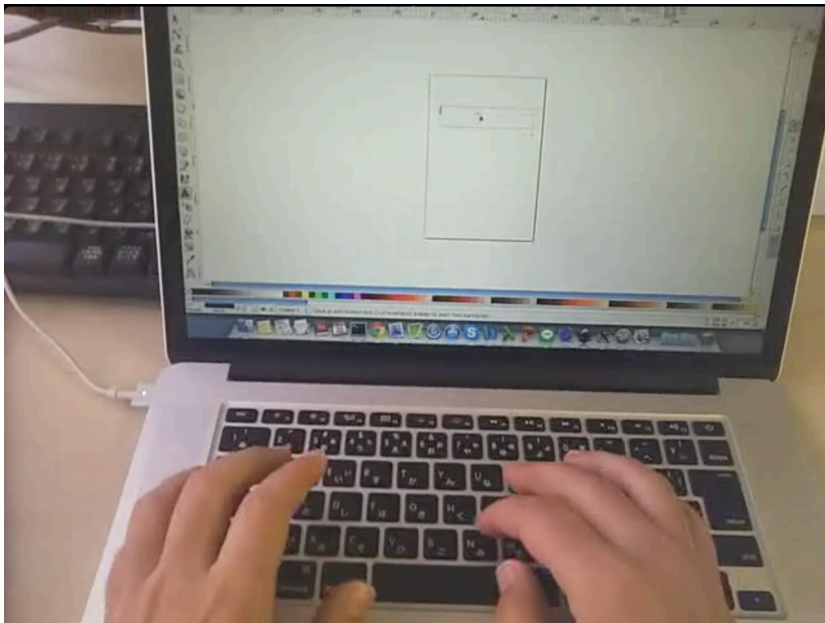| Reading | Not reading | Not reading | Not reading | Reading | Not reading |
|---------|-------------|-------------|-------------|---------|-------------|

# Evaluation of the extraction of the annotations

## Compared the annotations to the correct annotations

Correct annotation: manually annotate



Type a sentence.

# Experimental result

| Recall[%] | Precision[%] |
|---|---|
| 64.5 | 30.8 |

Recall:

Actions which can be annotated by the read texts are correctly annotated at the rate of 64.5%.

Precision:

Annotations are correct at the rate of 30.8%.

# Discussion

Causes for error:

   1. Extraction of read sentences fails

      If we could extract read texts correctly,

         Recall : 86.8 %

         Precision : 61.0%

   2. Assumption sometimes does not work

# Outline

- Motivation
- Approach
- Pilot Study
- Conclusion & Future Work

# Conclusion & Future work

## Conclusion

- Proposed a method for automatic video annotation in the scenario where users are following a textual manual

- Accuracy :

  If the video segmentation is succeeded,
  - Recall : 64.5%   Precision : 30.8%

## Future Work

- Improve the recall and precision of the annotation
- Improve the video segmentation method

35