

# Fingerspelling recognition with two-steps cascade process of spotting and classification

Masanori Muroi<sup>1</sup>, Naoya Sogi<sup>1</sup>, Nobuko Kato<sup>2</sup>, Kazuhiro Fukui<sup>1</sup>

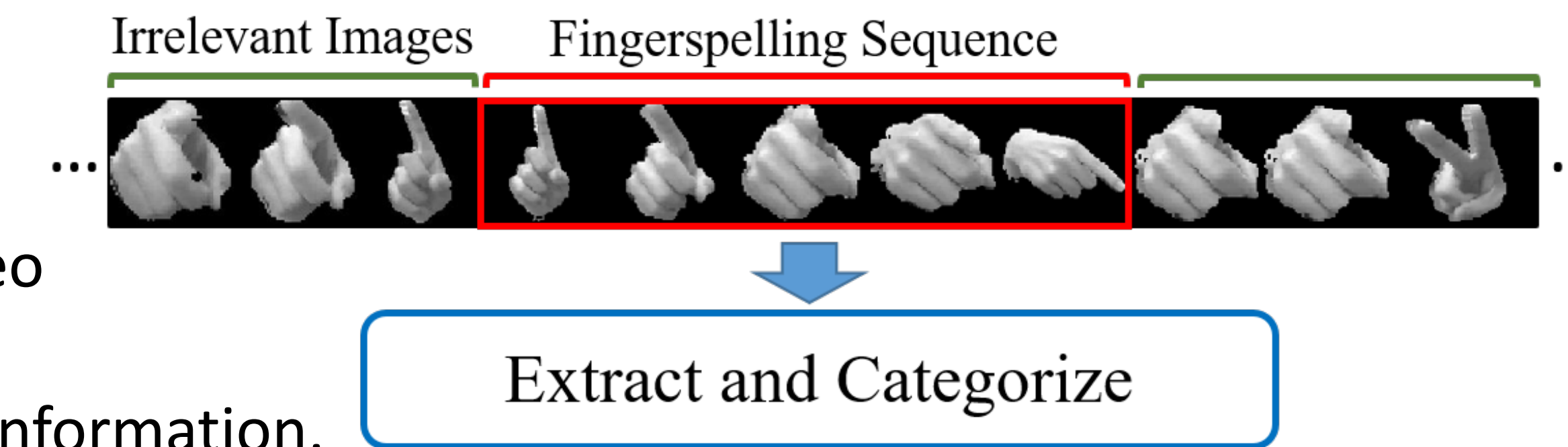
<sup>1</sup> Graduate School of Systems and Information Engineering, University of Tsukuba, Japan

<sup>2</sup> Faculty of Industrial Technology, Tsukuba University of Technology, Japan



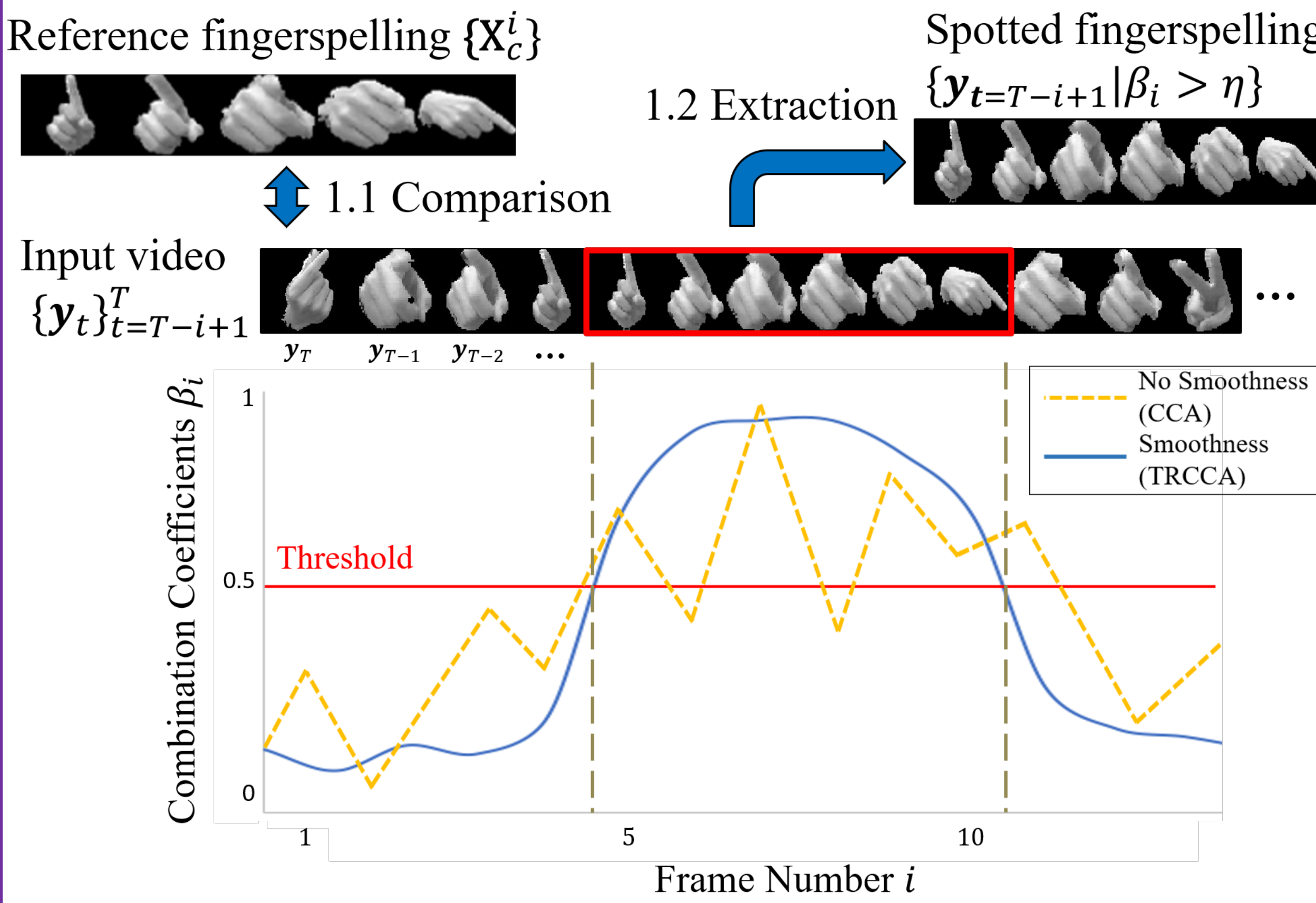
## (1) Introduction

- Fingerspelling is a tool to express a letter by a hand shape. ➤ Used in conjunction with sign language
- Main goal : Detect and categorize fingerspelling in a continuous video as a mixture of fingerspelling sequences and irrelevant images.**
- Basic idea :** Divide a whole process into two-steps cascade process:
  - Spotting :** Segment and extract a fingerspelling sequence in an input video by utilizing temporal dynamic information.
  - Classification :** Classify the spotted sequence by utilizing 3D hand shape information.
- We propose a fingerspelling classification framework based on two types of methods:
  - Temporal Regularized CCA (TRCCA)**[1] for spotting
  - Orthogonal Mutual Subspace Method (OMSM)**[2] with **CNN feature**[3] for classification



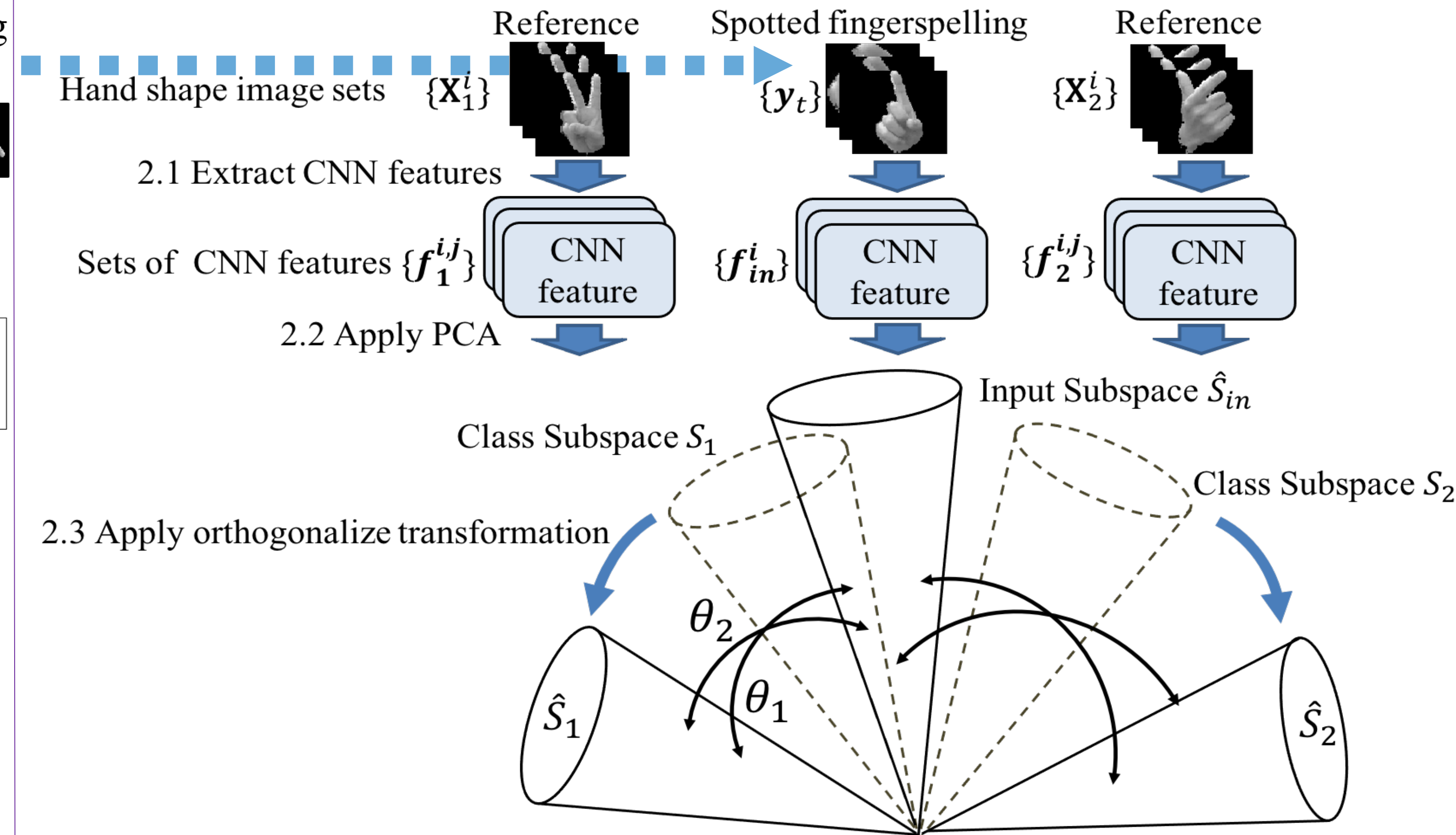
## (2) Proposed framework

### Step1: Spotting process



- Fingerspelling images in an input video are segmented and extracted using TRCCA.
- TRCCA: Extension of the *canonical correlation analysis* with the smoothness on the temporal domain
  - The smoothness efficiently incorporates the temporal information
- The detailed procedure
  - The input image sequence  $\{y_t\}_{t=T-i+1}^T$  is compared with the reference fingerspelling  $\{X_c^i\}$  by TRCCA. If the input sequence has high similarity with the reference fingerspelling, the input sequence is identified to fingerspelling.
  - Fingerspelling images  $\{y_{t=T-i+1} | \beta_i > \eta\}$  are extracted from the input sequence.

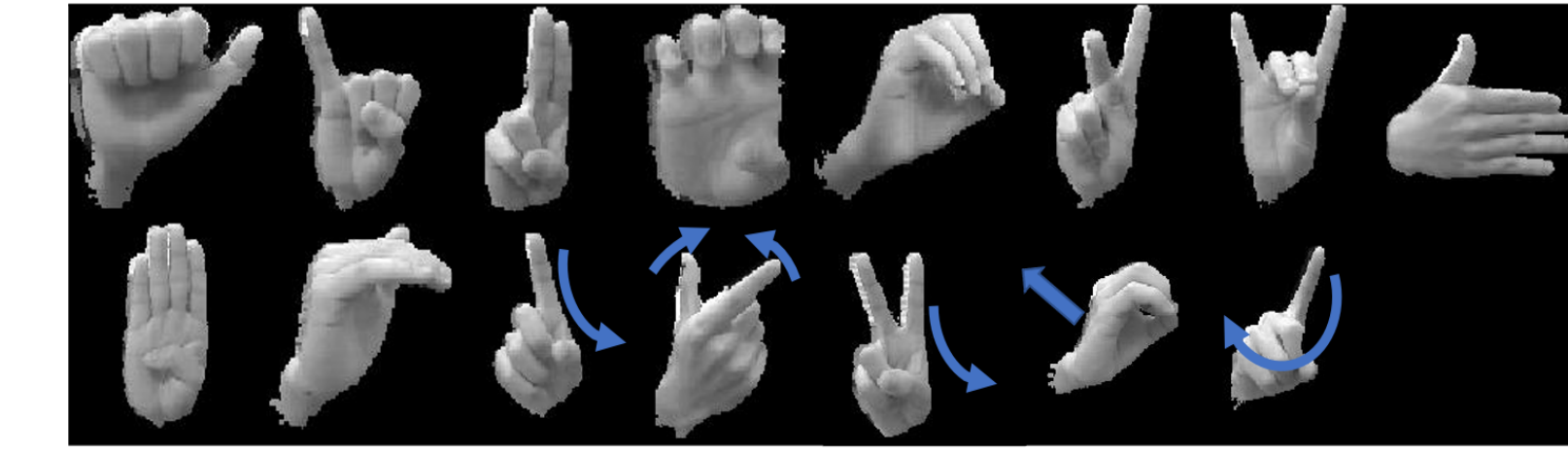
### Step2: Classification process



- Spotted fingerspelling is classified by OMSM with CNN features using hand shape image sets.
- OMSM: Represent 3D shape of a hand by subspace.
  - The subspace representation is effective for 3D object recognition
- The detailed procedure
  - CNN features  $\{f_{in}^i\}$  and  $\{f_c^i\}$  are extracted from  $\{y_{t=T-i+1} | \beta_i > \eta\}$  and  $\{X_c^i\}$ .
  - Each class subspace  $\{S_c\}$  and an input subspace  $S_{in}$  are generated by applying PCA to the sets of CNN features.
  - Orthogonal subspaces  $\{\hat{S}_c\}$  and  $\hat{S}_{in}$  are generated by applying orthogonalize transformation to  $\{S_c\}$  and  $S_{in}$ .
  - The spotted fingerspelling is classified based on similarities between the input subspace  $S_{in}$  and reference subspaces  $\{\hat{S}_c\}$ .

## (3) Experiments

- Dataset:**
  - We recorded **15 fingerspelling classes** by a depth camera.
  - ❑ The hand region is extracted from the whole input image based on the depth map.
  - We synthesized an input video, which continuously inputs fingerspelling and not fingerspelling sequences alternately.
- Evaluation index:**
  - Spotting performance, Classification accuracy, and Recognition time.



Sample images of Japanese fingerspelling dataset.

True Class	Predicted Class	
	Target	Not Target
Target	88 50.0%	0 0.0%
Not Target	15 8.5%	73 41.5%
	85.4% 14.6%	100% 0.0%
	91.5% 8.5%	

Confusion Matrix: Results of the spotting process.

Accuracies and recognition times of different methods.

Framework	Accuracy	Recognition Time
TRCCA[1]	64.1%	39.7 ms
CNN feat - OMSM	68.9%	52.7 ms
KOTRCCA[1]	79.0%	169.0 ms
TRCCA - CNN(softmax)	80.7%	56.9 ms
TRCCA - KOMSM	86.9%	187.3 ms
<b>TRCCA - CNN feat - OMSM(Proposed)</b>	<b>88.2%</b>	<b>91.2 ms</b>

True Class	Predicted Class								Not Target	
	Class 1	Class 2	Class 3	Class 4	Class 5	Class 6	Class 7	Class 8		
Class 1	11 10.8%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
Class 2	0 0.0%	11 10.8%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
Class 3	0 0.0%	0 0.0%	11 10.8%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
Class 4	0 0.0%	0 0.0%	0 0.0%	9 8.8%	1 1.0%	0 0.0%	1 1.0%	0 0.0%	0 0.0%	81.8% 18.2%
Class 5	0 0.0%	0 0.0%	0 0.0%	0 0.0%	11 10.8%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
Class 6	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	11 10.8%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
Class 7	1 1.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	10 9.8%	0 0.0%	0 0.0%	90.9% 9.1%
Class 8	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	11 10.8%	0 0.0%	100% 0.0%
Not Target	1 1.0%	0 0.0%	1 1.0%	2 2.0%	1 1.0%	3 2.9%	1 1.0%	0 0.0%	5 4.9%	35.7% 64.3%
	84.6% 15.4%	100% 0.0%	91.7% 8.3%	81.8% 18.2%	84.6% 15.4%	78.6% 21.4%	83.3% 16.7%	100% 0.0%	100% 0.0%	88.2% 11.8%

Confusion Matrix: Results of the classification process

## (4) Conclusion

- We proposed fingerspelling recognition framework based on a complementary combination of TRCCA and OMSM with CNN features.
- We confirmed that our two-steps process significantly outperforms conventional one-step methods in terms of classification accuracy and recognition time.

## (5) References

[1] S. Tanaka, A. Okazaki, N. Kato, H. Hino and K. Fukui, Spotting fingerspelled words from sign language video by temporally regularized canonical component analysis, *2016 IEEE International Conference on Identity, Security and Behavior Analysis*, 2016, pp. 1-7.  
 [2] K. Fukui and O. Yamaguchi, The kernel orthogonal mutual Subspace method and its application to 3D object recognition, in *Asian Conference on Computer Vision*, 2007, pp. 467-476.  
 [3] N. Sogi, T. Nakayama, and K. Fukui, A method based on convex cone model for image-set classification with cnn features, in *2018 International Joint Conference on Neural Networks*, 2018, pp. 1-8.