

IDEAI

Intelligent Data Science and Artificial Intelligence Research Center



### Abstract

Significant progress has been recently made on challenging tasks in automatic sign language (SL) understanding, such as sign language recognition, translation and production. However, these works have focused on datasets with relatively few samples, short recordings and limited vocabulary and signing space.

In this work, we introduce the novel task of sign language topic detection. We base our experiments on How2Sign[1], a large-scale video dataset spanning multiple semantic domains. We provide strong baselines for the task of topic detection, and present a comparison between different visual features commonly used in the domain of sign language.



sign language."CVPR 2021

[2] Duarte, Amanda, Samuel Albanie, Xavier Giró-i-Nieto, and Gül Varol. "Sign Language Video Retrieval with Free-Form Textual Queries." CVPR 2022.

# Topic Detection in Continuous Sign Language Videos











# **Topic Detection**

#### Predict a label that describes the semantic content of a signer's discourse.



- [1] Duarte, Amanda, Shruti Palaskar, Lucas Ventura, Deepti Ghadiyaram, Kenneth DeHaan, Florian Metze, Jordi Torres, and Xavier Giro-i-Nieto. "How2Sign: a large-scale multimodal dataset for continuous American
- [3] Jaegle, Andrew, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda Koppula et al. "Perceiver io: A general architecture for structured inputs & outputs." ICLR 2022.







# Visual Features

We compare four different visual features commonly used in sign language understanding:



#### **Experimental Results**

Test accuracy for each combination of features and architecture.

FAIRSEO

**O** PyTorch

	LSTM	Transformer	PerceiverIO	lstm - spot align - test	lstm - spot_align - test
Cartesian (OP)	$30.35 \pm 3.01$	$34.02 \pm 0.33$	$30.34 \pm 2.58$	1-24 5 0 0 0 0 1 0 1 0	
Cartesian (MP)	$29.43 \pm 1.17$	$33.10\pm2.58$	$33.56 \pm 1.30$	2 - 4 13 0 0 1 1 1 0 0   3 - 0 1 4 0 0 0 0 0 0 0	
Angular (OP)	$31.95 \pm 1.05$	$29.66\pm0.12$	$31.49 \pm 2.34$	4 - 1 1 0 3 1 0 0 0 0 0   a 5 - 1 1 0 0 6 1 3 0 0 1	
Angular (MP)	$32.64 \pm 0.32$	$34.71 \pm 1.42$	$30.80 \pm 1.97$	2 6 0 3 0 0 9 0 0 0 0   7 4 3 0 1 2 1 9 0 0 0	
I3D	$45.75 \pm 1.59$	$46.26 \pm 1.30$	$48.27 \pm 0.33$	8 - 7 2 0 0 1 0 4 2 0   9 - 3 0 0 0 0 0 1 3 0	
Spotted signs	$\textbf{58.03} \pm \textbf{1.40}$	$\textbf{53.33} \pm \textbf{2.18}$	$\textbf{52.88} \pm \textbf{0.32}$	10 - 1 0 0 0 0 0 0 0 1 2 1 2 3 4 5 6 7 8 9 10 Predicted label	
<b>Transcriptions</b> (upper bound)	$70.35 \pm 4.50$	$75.38\pm0.75$	$75.90 \pm 2.48$	Confusion matrix	t-SNE



\_\_\_\_\_

-

\_



https://github.com/imatge-upc/sign-topic



Accessibility, Vision, and Autonomy Workshop

### Neural Architectures

#### We compare three neural architectures:

Visualizations of the spotted signs + LSTM results:

Acknowledgements:



