

# Supplementary Information for

## AI enabled sign language recognition and VR space bidirectional communication using triboelectric smart glove

Feng Wen, Zixuan Zhang, Tianyiyi He, and Chengkuo Lee\*

\* Corresponding author: elelc@nus.edu.sg (C.L.)

**Supplementary Note 1.** The considerations behind the sensor distribution on gloves

**Supplementary Note 2.** The more detailed discussion about pros and cons of non-segmentation and segmentation methods

**Supplementary Note 3.** The accuracy performance of image recognition by comparing with sensor-based recognition system

**Supplementary Figure 1.** The detailed area information and channel label of sensors on gloves

**Supplementary Figure 2.** The photography of remaining 31 gestures and their corresponding triboelectric signals

**Supplementary Figure 3.** The train and validation accuracy increase with epochs

**Supplementary Figure 4.** The word frequency presented in the investigated 20 sentences

**Supplementary Figure 5.** The schematic diagram of segmentation

**Supplementary Figure 6.** The recognition result of three new sentences

**Supplementary Figure 7.** The radar comparison map of two methods

**Supplementary Figure 8.** The accuracy performance of image recognition when the brightness fades

**Supplementary Table 1.** CNN parameters

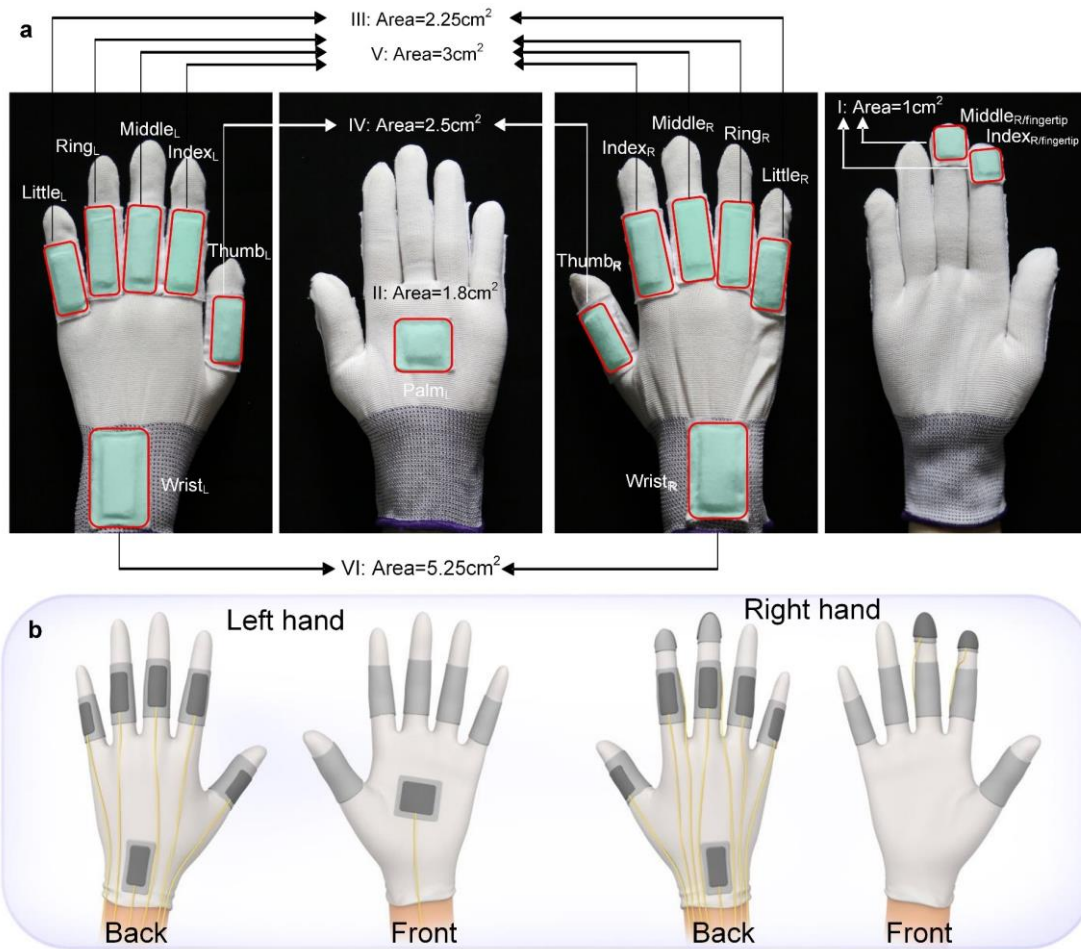
**Supplementary Table 2.** Detailed prediction of three new sentences

**Supplementary Table 3.** Benchmarking with other works

27 **Supplementary Note 1. The considerations behind the sensor distribution on**  
28 **gloves**

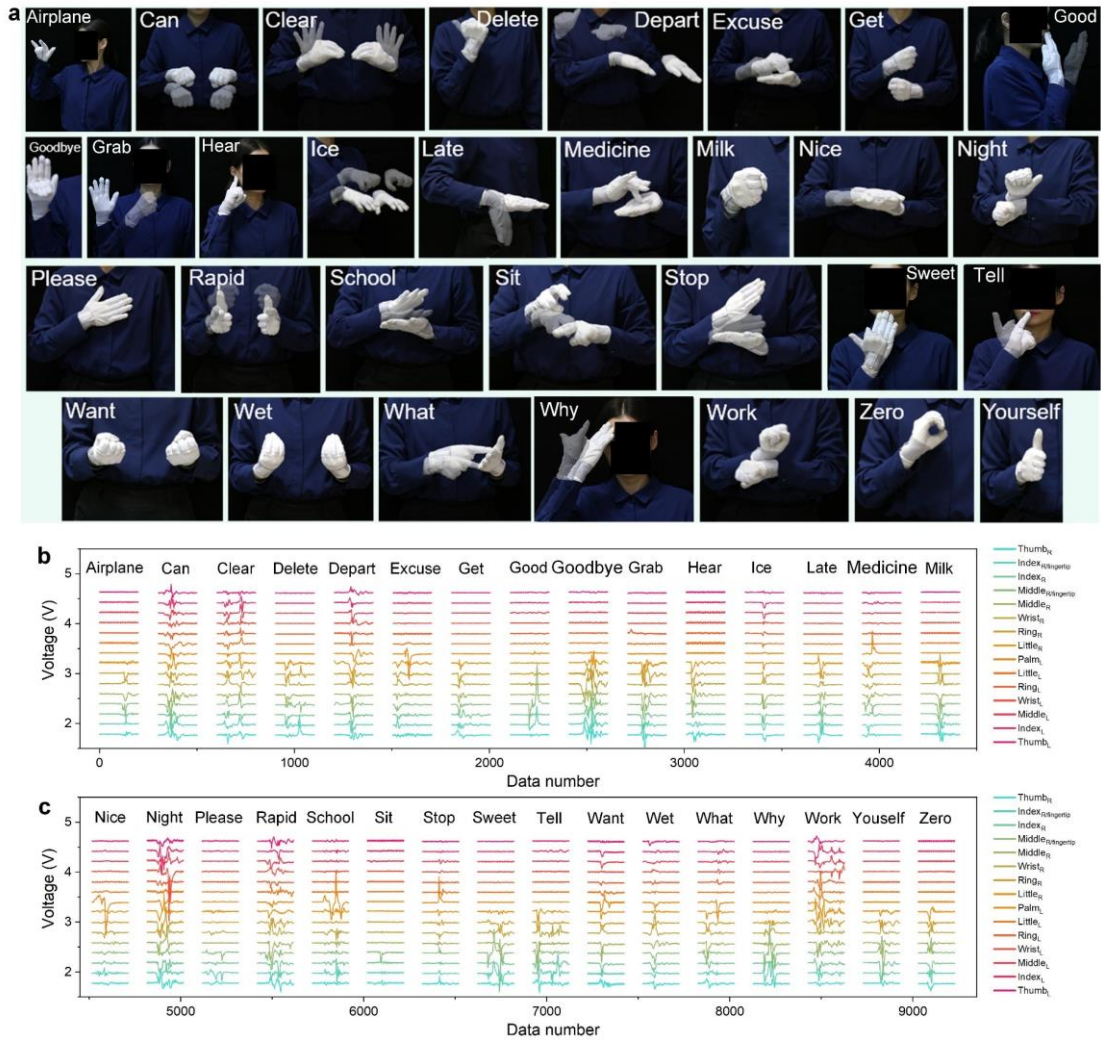
29 As depicted in Fig. 1b, the statistical analysis finds that the daily sign language involves  
30 three major motions, including elbow/shoulder motions, face muscle activities, and  
31 hand movements. The dominant hand motion accounts for 43%. Thus, the hand motion  
32 sensing is inevitable for sign language recognition. As shown in the enlarged pie chart  
33 in the right of Fig. 1b, the hand motions can be subdivided into four categories including  
34 finger bending (56%), wrist motion (18%), touch with fingertips (16%), and interaction  
35 with palm (10%). These detailed hand motions need sensors in different positions of  
36 hands to generate the essential correspondence.

37  
38 Fig. 1c and Supplementary Fig. 1 show the triboelectric sensor is mounted on each  
39 finger for finger bending detection, while two sensors are put on wrists for wrist motion  
40 perception. In addition, the fingertips of index and middle of right hand are also in  
41 frequent use in daily used sign language and hence two sensors are located at fingertips.  
42 Meanwhile, signers often use their palms to interact with other parts of their body to  
43 convey richer information. But we nominally allocate only one sensor on the palm of  
44 left hand rather than two sensors one located on the left hand and one located on right  
45 hand. There are two major considerations behind such arrangement: (1) based on the  
46 minimalist design for reducing system complexity, we expect as few sensors as possible  
47 with the limit of capable of detect necessary hand motions. Thus, only one sensor is  
48 located on the left hand instead of one for left hand and one for right hand. (2) This  
49 sensor is attached on the left hand not right hand. Because the final status of most of  
50 gestures that involves palm end up on the palm of left hand, such as ‘Excuse’,  
51 ‘Medicine’, ‘Nice’, ‘School’, ‘Stop’ and ‘What’ shown in Supplementary Fig. 2. Hence  
52 one palm sensor on the left hand is reasonable to sense the interaction motions.



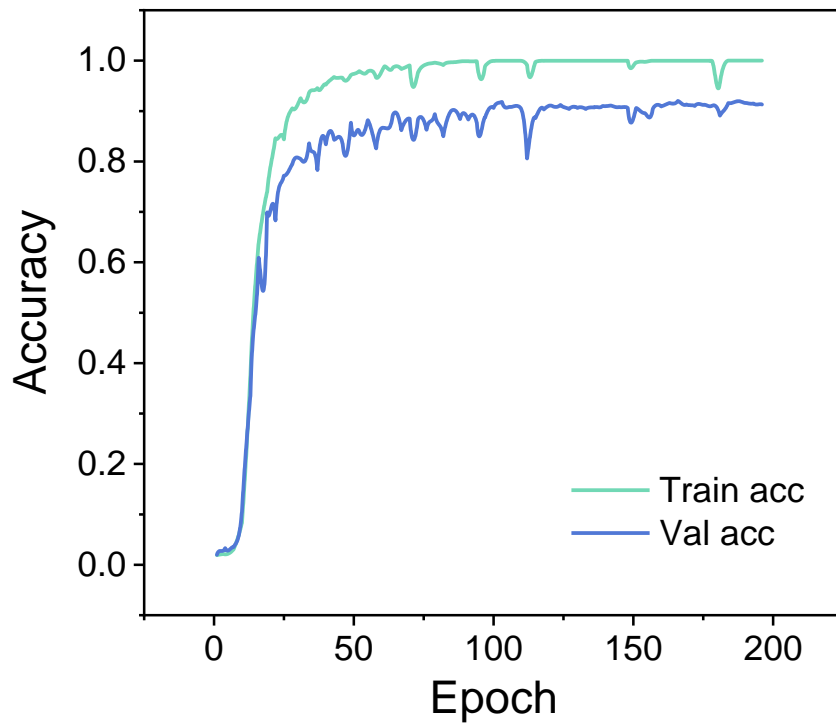
53

54 **Supplementary Figure 1. The detailed area information and channel label of**  
55 **sensors on gloves. a** Fabricated glove photos show detailed sensor area information  
56 and sensor channel label (corresponding with sensor output signal graphs) of sensors in  
57 a different position. **b** Schematic diagram of sensors on hand, corresponding with the  
58 photos of proposed gloves. The hand images are created by the authors via Blender.  
59 Photo credit: Feng Wen, National University of Singapore.



60

61 **Supplementary Figure 2. The photography of remaining 31 gestures and their**  
 62 **corresponding triboelectric signals. a** Photography of the remaining 31 gestures. The  
 63 opaque and translucent gesture images show the starting and final state of the gesture,  
 64 respectively. These photos are of one of the authors. **b-c** Corresponding signals of these  
 65 31 gestures. Photo credit: Feng Wen, National University of Singapore.



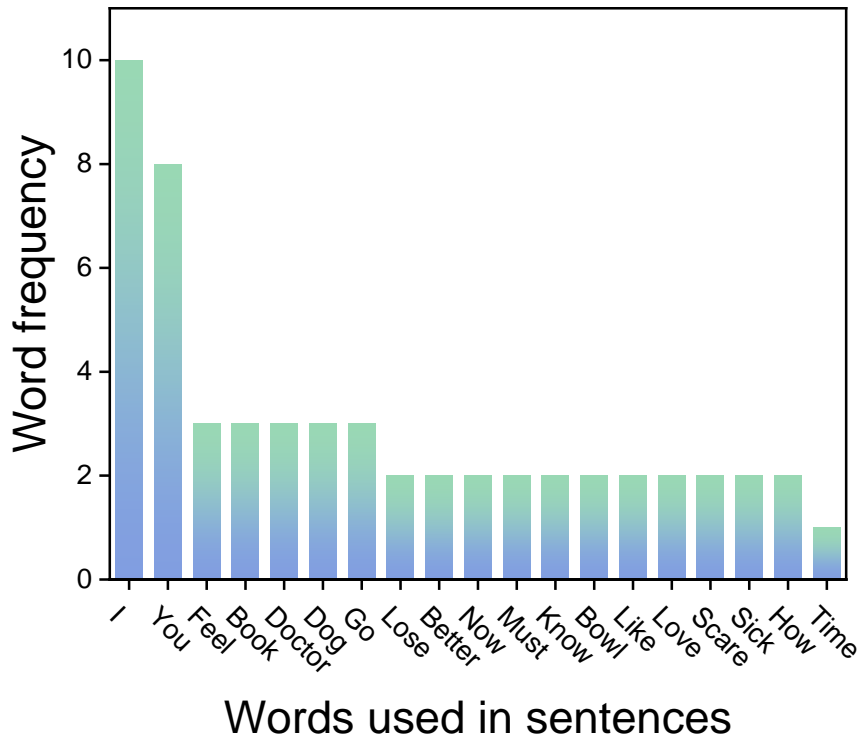
66

67 **Supplementary Figure 3. The train and validation accuracy increase with epochs.**

68 After a small number of training epochs 50, the accuracy achieves an acceptable level,

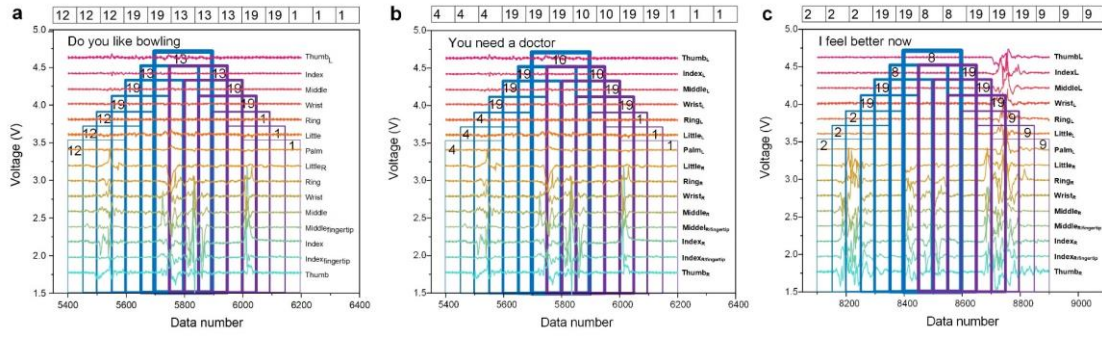
69 proving the good performance of the proposed CNN model for sign language

70 recognition.



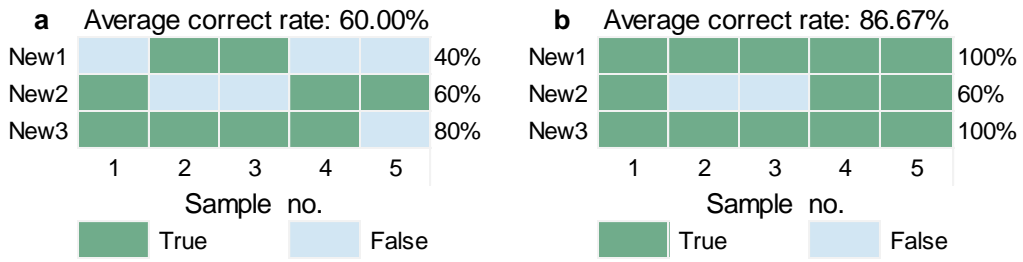
71

72 **Supplementary Figure 4. The word frequency presented in the investigated 20**  
 73 **sentences.** 19 Words are numbered from 0-18 according to the decreased usage  
 74 frequency.



75  
 76  
 77  
 78  
 79

**Supplementary Figure 5. The schematic diagram of segmentation. a** ‘Do you like bowling’, **b** ‘You need a doctor’, and **c** ‘I feel better now’ as examples to show more detailed signal splitting process.



80

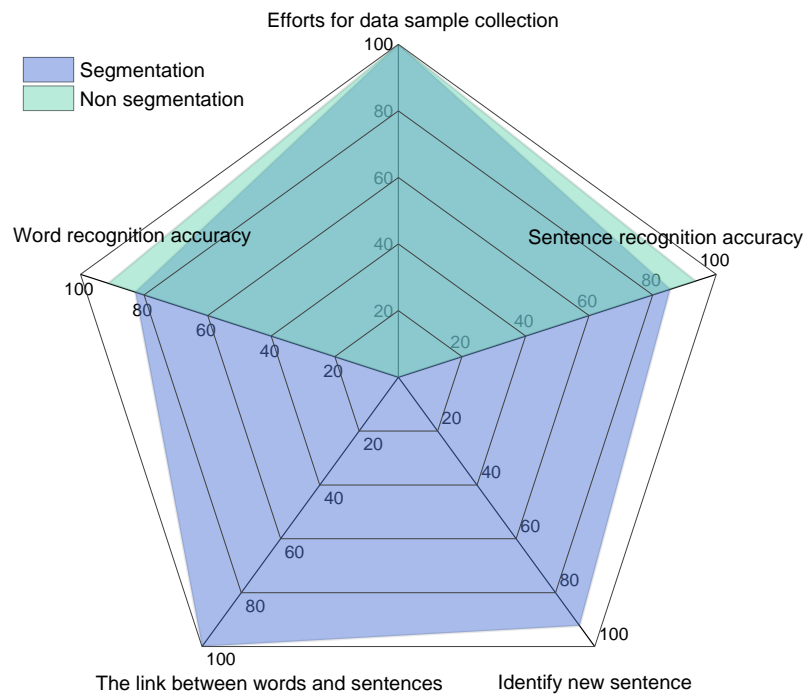
81 **Supplementary Figure 6. The recognition result of three new sentences. a** Using  
 82 single classifier. **b** Using hierarchy classifier. The false prediction area is greatly  
 83 reduced by using hierarchy classifier. Each sentence has been tested for five times with  
 84 five samples.  
 85



86 **Supplementary Note 2. The more detailed discussion about pros and cons of non-**  
87 **segmentation and segmentation methods**

88 To illustrate the pros and cons of non-segmentation and segmentation methods, the  
89 detailed implementation of these two approaches should be discussed first. For non-  
90 segmentation method, each word or sentence is labeled as the independent individual.  
91 Then all the words and sentences will be separately trained in the neural network. Upon  
92 the completion of training, the CNN will recognize words and sentences independently.  
93 With such regime, either words or sentences essentially are different classes with  
94 respect to CNN's cognition, in which there is no built-up relationship between word  
95 units and sentences. For the strategy of segmentation, the data sliding window divides  
96 the entire sentence signal (800 data points) into fragments including intact word signals,  
97 incomplete word signals, and background signals. The label of fragment split from  
98 sentence signal is determined by the principal component. In other words, either word  
99 signal or background noise accounts for more than 50% of the sliding window size, and  
100 the label will be the number of corresponding words or 'empty' 19 as shown in Fig. 4a.  
101 Due to the specific size (200 data points) and sliding step (50 data points) of sliding  
102 window, the entire sentence signal is split into 13 fragments where each one is labeled  
103 with a number. Hence, the label of sentence will be a series of 13 numbers as illustrated  
104 in the top of Fig. 4b. Next, the sentence signals with the label of number sequences are  
105 included the dataset for training. The CNN classifier will go through all the fragments  
106 as well as the fragment sequence in sentences. Ultimately, both fragments and reversely  
107 reconstructed sentences by virtue of fragments can be correctly recognized. In particular,  
108 the CNN classifier is even endowed with the capability of recognizing never-seen  
109 sentences that comprise new-order word fragments, in which the never-seen sentences  
110 are not included in the dataset for training process and hence never learned by the neural  
111 network before.

112  
113 Overall, as the radar map of comparison in Supplementary Fig. 7 shows, the non-  
114 segmentation approach possesses better performance in the aspect of recognition  
115 accuracy either for words or sentences. However, two following limitations of this  
116 means may compromise the universality and practicality of the whole system. Above  
117 all, owing to the independence of words and sentences, the CNN classifier cannot  
118 identify the new sentence although words in the sentence are seen before and only  
119 combined in a new order. In addition, when expanding sentence database, the labor-  
120 intensive data collection of new sentences and successive training are unavoidable. This  
121 kind of independence also leads to increased effort on data collection of words since  
122 the CNN model cannot extract and recognize the word signals in the sentence.  
123 Regarding the segmentation method, in addition to identifying existing sentences in the  
124 dataset, the CNN classifier enables the recognition of new sentences. These new  
125 sentences comprise new-order word series that are different from the order of existed  
126 sentences in the dataset. Nevertheless, the segmentation introduces a large amount of  
127 random and irregular 'empty' signals. It sacrifices the recognition accuracy for both  
128 words and sentences. Further research efforts could be committed to optimizing the  
129 algorithm framework and improve the recognition accuracy.



130

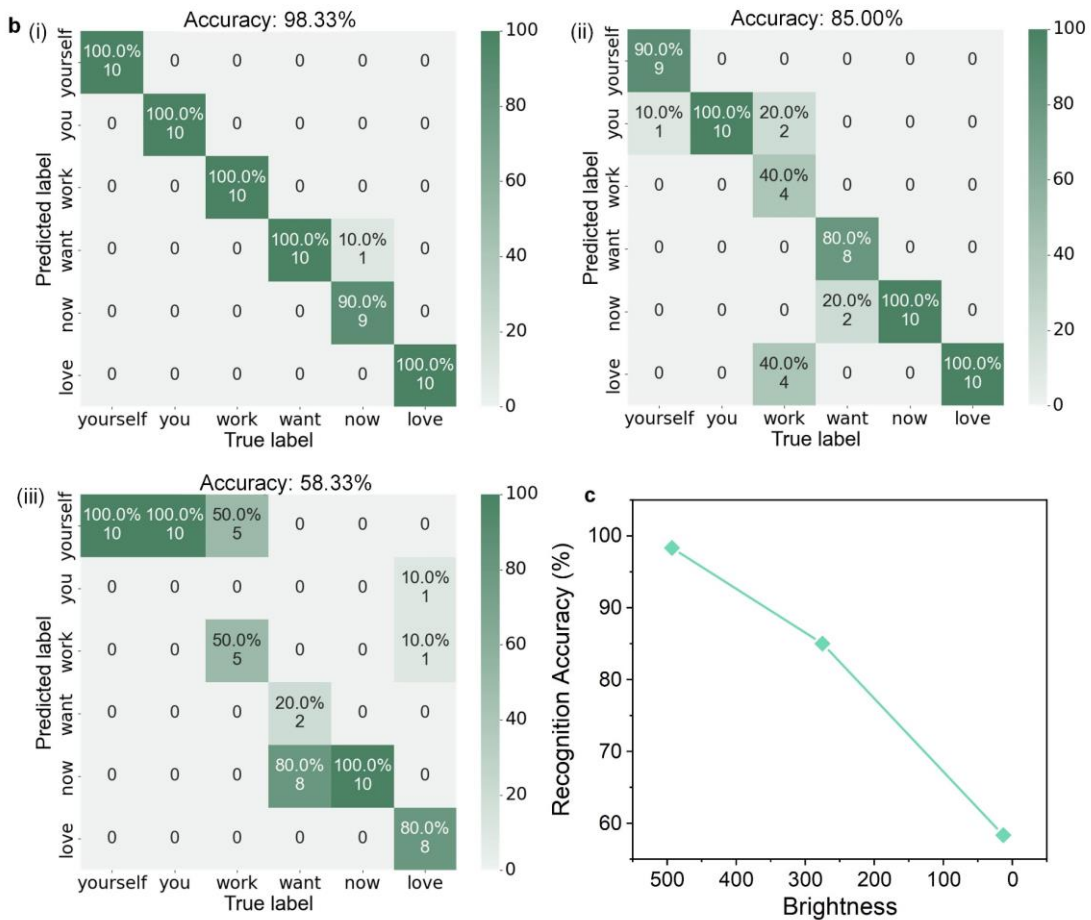
131 **Supplementary Figure 7. The radar comparison map of two methods.** Comparing  
 132 non-segmentation and segmentation recognition methods based on their pros and cons.  
 133

134 **Supplementary Note 3. The accuracy performance of image recognition by**  
135 **comparing with sensor-based recognition system**

136 To clarify the advantages of sensor-based system, the additional test about the accuracy  
137 performance of visual images for gesture recognition is carried out. The recognition  
138 results of six representative gestures are shown in Supplementary Fig. 8 with varying  
139 light conditions (493, 275 and 13 lux). For each light condition, 50 trials of each gesture  
140 (300 trials in total) are carried out for image-based recognition. Supplementary Fig.  
141 8b(i-iii) indicate a dramatically decayed recognition accuracy from 98.33% to 58.33%  
142 when the room light fades. The efficiency of visual images/videos recognition is well-  
143 known limited by the environmental interferences such as occlusions and especially  
144 light conditions. In addition, sign language involves the upper limbs as well as human  
145 faces. When the image-based system captures gesture information, the exposure of  
146 facial information to camera may arise the issue of privacy disclosure.

147

148 For sensor-based human gesture recognition system, wearable sensors, are typically less  
149 bulky, flexible and provide an intimate contact with the user for high-quality data  
150 acquisition and high-accurate recognition that is comparable with its image system  
151 counterpart. The sensor-based system is considered as one of approaches to overcome  
152 the drawbacks of image recognition. On the one hand, such sensor-based systems are  
153 not affected by varying luminance and can work well even under entirely dark condition  
154 with higher environmental tolerance. On the other hand, they can mitigate the privacy  
155 issue in cost-effective way owing to no need for individual information collection such  
156 as facial characteristics.



157

158

**Supplementary Figure 8. The accuracy performance of image recognition when**

159

**the brightness fades. a** Gesture image of ‘Love’ under three light conditions. **b** (i-iii)

160

Accuracy of the image recognition under different light conditions (493, 275 and 13

161

lux). These photos are of one of the authors. **c** Accuracy degradation with decreased

162

brightness. Photo credit: Feng Wen, National University of Singapore.

163

164 **Supplementary Table 1. CNN parameters.** The detailed parameters for constructing  
 165 Convolutional Neural Network (CNN).

166

No	Layer Type	No. of Filters	Kernel/ Pool Size	Stride	Input Size	Output Size	Padding
1	Convolution 2	64	5	1	(None, 100, 32)	(None, 100, 64)	same
2	Max Pooling 2		2	2	(None, 100, 64)	(None, 50, 64)	same
3	Convolution 3	128	5	1	(None, 50, 64)	(None, 50, 128)	same
4	Max Pooling 3		2	2	(None, 50, 128)	(None, 25, 128)	same
5	Convolution 4	256	5	1	(None, 25, 128)	(None, 25, 256)	same
6	Max Pooling 4		2	2	(None, 25, 256)	(None, 13, 256)	same
7	Convolution 5	512	5	1	(None, 13, 256)	(None, 13, 512)	same
8	Max Pooling 5		2	2	(None, 13, 512)	(None, 7, 512)	same
9	Flatten				(None, 7, 512)	(None, 3584)	same
10	Dense (500)				(None, 3584)	(None, 500)	
11	Dense (50)				(None, 500)	(None, 50)	

167

168  
169  
170  
171

**Supplementary Table 2. Detailed prediction of three new sentences.** The predicted and true labels of single classifier and hierarchy classifier for three new/never-seen sentence recognition with numbers in red standing for wrong prediction.

Sentence	Sample	True label	Single classifier predicted label	Hierarchy classifier predicted label
New1	1	[ 0 0 19 19 19 7 7 19 19 19 5 5 5 ]	[ 0 0 19 19 8 7 7 19 19 19 5 5 5 ]	[ 0 0 19 19 19 7 7 19 19 19 19 5 5 ]
	2	[ 0 0 19 19 7 7 7 19 19 19 5 5 5 ]	[ 0 0 19 19 11 7 7 7 19 19 19 5 5 5 ]	[ 0 0 19 19 7 7 7 19 19 19 5 5 5 ]
	3	[ 0 0 19 19 19 7 7 19 19 19 19 5 5 ]	[ 0 0 19 19 6 7 7 19 19 19 19 5 5 ]	[ 0 0 19 19 19 7 7 19 19 19 19 5 5 ]
	4	[ 0 0 19 19 07 7 7 19 19 19 5 5 5 ]	[ 0 0 19 19 7 7 19 19 19 19 5 5 5 ]	[ 0 0 19 19 7 7 19 19 19 19 5 5 5 ]
	5	[ 0 0 19 19 19 7 7 19 19 19 19 5 5 ]	[ 0 0 19 19 19 7 7 7 19 19 19 5 5 ]	[ 0 0 19 19 19 7 7 19 19 19 19 5 5 ]
New2	1	[ 0 0 0 19 19 19 4 4 4 19 19 19 19 ]	[ 19 0 0 0 19 19 4 4 4 19 19 19 19 ]	[ 0 0 0 0 19 19 4 4 4 19 19 19 19 ]
	2	[ 0 0 0 19 19 19 19 4 4 19 19 19 19 ]	[ 0 0 0 19 19 19 4 4 4 19 19 19 19 ]	[ 0 0 0 19 19 19 4 4 4 19 19 19 19 ]
	3	[ 0 0 0 19 19 19 19 4 4 19 19 19 19 ]	[ 0 0 0 19 19 19 4 4 4 4 19 19 19 ]	[ 0 0 0 19 19 19 19 4 4 4 19 19 19 ]
	4	[ 19 0 0 0 19 19 19 4 4 4 19 19 19 ]	[ 19 0 0 18 19 19 19 4 4 4 19 19 19 ]	[ 1 0 0 0 19 19 19 4 4 4 19 19 19 ]
	5	[ 19 0 0 0 19 19 19 19 4 4 4 19 ]	[ 19 3 13 13 19 19 19 19 4 4 4 19 ]	[ 19 0 18 19 19 19 19 19 4 4 4 19 ]
New3	1	[ 17 17 17 19 19 19 1 1 1 19 19 2 2 2 ]	[ 17 17 17 19 19 1 1 1 0 19 19 2 2 2 ]	[ 17 17 19 19 19 1 1 1 19 19 2 2 2 ]
	2	[ 17 17 17 19 19 1 1 1 19 19 2 2 2 ]	[ 17 17 17 7 1 1 1 1 19 19 2 2 2 ]	[ 17 17 17 19 19 1 1 1 19 19 2 2 2 ]
	3	[ 17 17 17 19 1 1 1 1 19 2 2 2 19 19 ]	[ 17 19 17 17 1 1 1 1 19 2 2 2 19 19 ]	[ 17 17 17 19 1 1 1 1 19 2 2 2 19 19 ]
	4	[ 17 17 17 19 19 1 1 1 1 19 2 2 2 19 ]	[ 17 17 17 17 19 1 1 1 1 1 19 2 2 2 19 ]	[ 17 17 17 19 1 1 1 1 1 19 2 2 2 2 ]
	5	[ 17 17 17 19 19 1 1 1 1 19 2 2 2 19 ]	[ 17 19 17 7 19 1 1 1 1 19 2 2 2 19 ]	[ 17 19 19 19 1 1 1 1 1 19 2 2 19 19 ]

172

**Supplementary Table 3. Benchmarking with other works.** The benchmarking table for comparing with other similar works.

Mechanism	Device	Method	Sensor/hand	Gesture	Sentence	Recognize new sentence	Self-powered	Mutual interaction	Ref.
Piezoresistive	Armband	SVM	3	8 letters	-	-	×	×	1
Piezoresistive	Radar+wristsband	HSVM	5	4 letters	0	×	×	×	2
Resistive	Finger sensing	SVM	5	10 letters	0	×	×	×	3
Resistive	Glove	BNN	9	10 letters	0	×	×	×	4
Resistive	Elastomer Glove	Amplitude	5	5	-	-	×	×	5
Resistive	Glove	Amplitude	5	9 words	0	×	×	×	6
Resistive	Glove	Amplitude	9	26 letters	0	×	×	×	7
Capacitive	Glove	Amplitude	5	5 numbers	-	-	×	×	8
Capacitive	Glove	Amplitude	5	36 letters	0	×	×	×	9
Capacitive	Wristband	Amplitude	15	15 words	0	×	×	×	10
Capacitive	Glove	XRF	16	10 numbers	0	×	×	×	11
Capacitive	Finger sensing	Amplitude	4	6 numbers	-	-	×	×	12
Ionic	Glove	Amplitude	8	6 words	0	×	√	×	13
Piezoelectric	Glove	Amplitude	5	6 numbers	0	×	√	×	14
Triboelectric	Glove	Amplitude	5	6 numbers	-	-	√	×	15
Triboelectric	Glove	Amplitude	6	7 numbers	-	-	√	×	16
Triboelectric	Glove	Amplitude	5	5 letters	0	×	√	×	17
Triboelectric	Glove	SVM	5	11 letters	0	×	√	×	18
Triboelectric	Glove	CNN	7	50 words	20	√	√	√	*

Note: (H)SVM: (hierarchical) support vector machine; BNN: binary neural network; XRF: extremely randomized trees; CNN: convolutional neural network. (\*This work)

## 175 **References**

- 176 1. Esposito, D. *et al.* A piezoresistive array armband with reduced number of sensors  
177 for hand gesture recognition. *Front. in neurorob.* **13**, 114 (2020).
- 178 2. Liang, X., Li, H., Wang, W., Liu, Y., Ghannam, R., Fioranelli, F., & Heidari, H.  
179 (2019). Fusion of wearable and contactless sensors for intelligent gesture  
180 recognition. *Adv. Intell. Syst.* **1**, 1900088.
- 181 3. Li, L., Jiang, S., Shull, P. B., & Gu, G. SkinGest: artificial skin for gesture recognition  
182 via filmy stretchable strain sensors. *Adv. Robot.* **32**, 1112-1121 (2018).
- 183 4. Fan, T., *et al.* Analog Sensing and Computing Systems with Low Power  
184 Consumption for Gesture Recognition. *Adv. Intell. Syst.* **3**, 2000184 (2021).
- 185 5. Muth, J. T. *et al.* Embedded 3D printing of strain sensors within highly stretchable  
186 elastomers. *Adv. Mater.* **26**, 6307-6312 (2014).
- 187 6. Ambar, R. *et al.* Development of a wearable device for sign language recognition. *J.*  
188 *Phys. Conf. Ser.* **1019**, 012017 (2018).
- 189 7. O'Connor, T. F. *et al.* The language of glove: wireless gesture decoder with low-  
190 power and stretchable hybrid electronics. *PLOS One* **12**, e0179766 (2017).
- 191 8. Shintake, J., Piskarev, E., Jeong, S. H., & Floreano, D. Ultrastretchable strain sensors  
192 using carbon black-filled elastomer composites and comparison of capacitive  
193 versus resistive sensors. *Adv. Mat. Techno.* **3**, 1700284 (2018).
- 194 9. Abhishek, K. S., Qubeley, L. C. F. & Ho, D. Glove-based hand gesture recognition  
195 sign language translator using capacitive touch sensor. *2016 IEEE Int. Conf.*  
196 *Electron Devices and Solid-State Circuits*, pp. 334-337 (2016).
- 197 10. Truong, H. *et al.* Cap-band: Battery-free successive capacitance sensing wristband  
198 for hand gesture recognition. *In Proc. of the 16th ACM Conf. on Embedded*  
199 *Networked Sens. Syst.*, pp. 54-67 (2018).
- 200 11. Pan, J. *et al.* A wireless multi-channel capacitive sensor system for efficient glove-  
201 based gesture recognition with AI at the edge. *IEEE Trans. on Circuits and Syst.*  
202 *II: Express Briefs* **67**, 1624-1628 (2020).
- 203 12. Bartlett, M. D., Markvicka, E. J., & Majidi, C. Rapid fabrication of soft,  
204 multilayered electronics for wearable biomonitors. *Adv. Funct. Mater.* **26**, 8496-  
205 8504 (2016).
- 206 13. Zhao, J. *et al.* Passive and space-discriminative ionic sensors based on durable  
207 nanocomposite electrodes toward sign language recognition. *ACS Nano* **11**, 8590-  
208 8599 (2017).
- 209 14. Fuh, Y. K., & Ho, H. C. Highly flexible self-powered sensors based on printed  
210 circuit board technology for human motion detection and gesture recognition.  
211 *Nanotechnology* **27**, 095401 (2016).
- 212 15. Lai, Y. C. *et al.* Single-thread-based wearable and highly stretchable triboelectric  
213 nanogenerators and their applications in cloth-based self-powered human-  
214 interactive and biomedical sensing. *Adv. Funct. Mater.* **27**, 1604462 (2017).
- 215 16. He, Q. *et al.* An all-textile triboelectric sensor for wearable teleoperated human-  
216 machine interaction. *Journal of Materials Chemistry A* **7**, 26804-26811(2019).
- 217 17. Maharjan, P. *et al.* A human skin-inspired self-powered flex sensor with thermally



- 218 embossed microstructured triboelectric layers for sign language interpretation.  
219 *Nano Energy* **76**, 105071 (2020).
- 220 18. Zhou, Z. *et al.* Sign-to-speech translation using machine-learning-assisted  
221 stretchable sensor arrays. *Nat. Electron.* **3**, 571-578 (2020).