# Modeling of Non-verbal Behaviors of Students in Cooperative Learning by Using OpenPose

Eiji Watanabe[1(✉)] 📙, Takashi Ozeki[2], and Takeshi Kohama[3]

[1] Konan University, Kobe 658-8501, Japan
`e_wata@konan-u.ac.jp`
[2] Fukuyama University, Fukuyama, Hiroshima 729-0292, Japan
[3] Kindai University, Kinokawa, Wakayama 649-6493, Japan

**Abstract.** In this paper, we discuss the modeling for the interactions between non-verbal behaviors of "teaching/learning" students in the cooperative learning. First, we adopt the positions eyes, face and hands detected by OpenPose [7] as a skeleton detection algorithm by a single camera. Next, we propose a modeling method for non-verbal behaviors based on neural networks. Furthermore, we discuss the modeling results for the interactions between non-verbal behaviors of students based on the internal presentations in the above models.

**Keywords:** Cooperative learning · Student · Non-verbal behavior · Modeling · Neural network · OpenPose

## 1 Introduction

The cooperative learning is an effective approach which aims to understand the content as the group with helping each other. In the cooperative learning, students teach other students and vice versa. It is becoming one of the hot topics to be researched [1]. The object of the cooperative learning is to improve the cooperation of the group and the understanding for given contents. Moreover, in [2], the following fundamental factors to be effective among the learning group which is listed as follows; (i) mutually beneficial cooperation, (ii) roles and responsibilities of the individual, and (iii) stimulatory interaction. However, one teacher can not grasp the cooperation and understanding of multiple groups and can not evaluate the above fundamental factors in real time. Therefore, it is important to construct methods for the estimation of the cooperation and the understanding in the group based on the non-verbal behaviors [3,4]. On the other hand, a conceptual model TSCL (Tabletop-Supported Collaborative Learning) has been proposed for understanding of the collaborative learning process [5].

The authors have already discussed the relationship between non-verbal behaviors and understandings of students in the cooperative learning [6]. However, non-verbal behaviors of "teaching/learning" students are represented by

the facial size detected by the camera and the exact direction of the facial movement has not been discussed. In this paper, we discuss the modeling for the interactions between non-verbal behaviors of "teaching/learning" students in the cooperative learning. First, we adopt the positions of eyes, face, and hands detected by OpenPose [7] as a skeleton detection algorithm by a single camera. Next, we propose a modeling method for non-verbal behaviors based on neural networks [10]. Furthermore, we discuss the modeling results for the interactions between non-verbal behaviors of students based on the internal presentations in the above models.

## 2    Detection of Non-verbal Behaviors of Students in Cooperative Learning Environment

### 2.1    Non-verbal Behaviors

In this paper, we treat the learning environment using a whiteboard and a table as shown in Fig. 1. In this learning environment, we have to consider two cases; (Case-1) using a whiteboard and (Case-2) not using a whiteboard. In Case-1, the "teaching" student standing in the front of a whiteboard explains and writes the content in a whiteboard. At the same time, other "learning" students sitting around a table take their notes, look at the whiteboard and listen to the explanation. In this case, the role "teaching/learning" of each student is clear. In Case-2, all students sit around a table and they teach and learn with each other by the explanation, listening and taking notes. In this case, the role "teaching/learning" of each student is unclear.
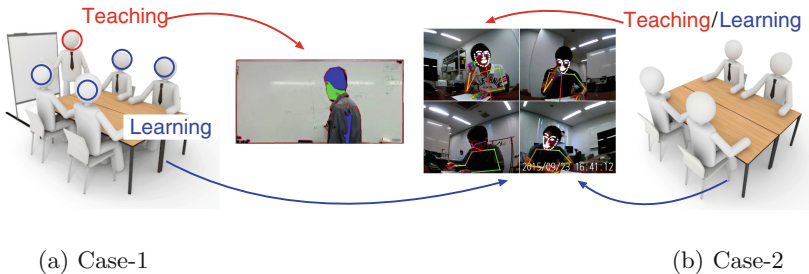


(a) Case-1                                              (b) Case-2

**Fig. 1.** Non-verbal behaviors of "teaching/learning" students in cooperative learning.

### 2.2    Detection of Non-verbal Behaviors by OpenPose [7]

In this learning environment, we can detect the following non-verbal behaviors; (i) writing on the whiteboard, (ii) the explanation to "learning" students, listening to "teaching" student, and (iii) taking notes by using OpenPose [7].

OpenPose can detect a human body, hand, facial, and foot keypoints (in total 135 keypoints) on single images as follows;

– Body: "Neck", "Shoulder", "Elbow", "Wrist", $\cdots$
– Face: "Nose", "Eye", "Mouth", $\cdots$
– Hand: "Finger", "Palm", $\cdots$

### 2.3 Detection of Non-verbal Behaviors in the Case of Using a Whiteboard

In this case, the "teaching" student in the front of a whiteboard has the following behaviors; (i) writing the content on a whiteboard, (ii) explanation to "learning" students as shown in Fig. 2(a). On the other hand, the "learning" students surrounding a table have the following behaviors; taking notes and listening to the explanation as shown in Fig. 2(b). Therefore, in this case, we use some body parts "Neck", "Eye", and "Finger" for the detection of each behavior as follows; (i) writing the content on a whiteboard, (ii) explanation to "learning" students, (iii) taking notes and listening to the explanation. Here, we define the positions $\boldsymbol{p}^{WB}(t)$ of some body parts for non-verbal behaviors of "teaching" students in the front of a whiteboard as follows;

$$
\begin{aligned}
\boldsymbol{p}^{WB}(t) = (&x_{Neck}(t), y_{Neck}(t), x_{Eye}^{L}(t), y_{Eye}^{L}(t), x_{Eye}^{R}(t), y_{Eye}^{R}(t), \\
&x_{Hand}^{L}(t), y_{Hand}^{L}(t), x_{Hand}^{R}(t), y_{Hand}^{R}(t))^{T},
\end{aligned}
\tag{1}
$$

where $(x_{Neck}(t), y_{Neck}(t))$ denotes the coordinates of "Neck". $(x_{Eye}^{L}(t), y_{Eye}^{L}(t))$ and $(x_{Eye}^{R}(t), y_{Eye}^{R}(t))$ denote the coordinates of "Left Eye" and "Right Eye" respectively. $(x_{Hand}^{L}(t), y_{Hand}^{L}(t))$ and $(x_{Hand}^{R}(t), y_{Hand}^{R}(t))$ denote the coordinates of "Left Hand" and "Right Hand" respectively. For example, we can summarize the relationships between the behaviors and the positions $\boldsymbol{p}^{WB}(t)$ of some parts detected by OpenPose as follows;

– Writing the content on a whiteboard: $x_{Neck}(t) \neq 0$, $x_{Eye}^{L}(t) = 0$, $x_{Eye}^{R}(t) = 0$.
– Explanation to "learning" students by speech: $\{x_{Eye}^{L}(t) \neq 0$ and $x_{Eye}^{R}(t) \neq 0\}$, $\{x_{Hand}^{L}(t) \neq 0$ and $x_{Hand}^{R}(t) \neq 0\}$.
– Explanation to "learning" students with a whiteboard: $\{x_{Eye}^{L}(t) \neq 0$ or $x_{Eye}^{R}(t) \neq 0\}$, $\{x_{Hand}^{L}(t) \neq 0$ or $x_{Hand}^{R}(t) \neq 0\}$.

$(x_{Eye}^L(t), y_{Eye}^L(t))$
$(x_{Eye}^R(t), y_{Eye}^R(t))$

$(x_{Neck}(t), y_{Neck}(t))$

$(x_{Hand}^L(t), y_{Hand}^L(t))$
$(x_{Hand}^R(t), y_{Hand}^R(t))$

[Writing contents]

[Explanation]

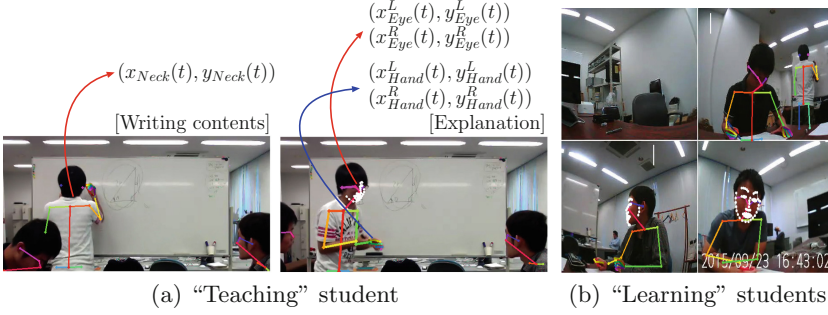(a) "Teaching" student          (b) "Learning" students

**Fig. 2.** Positions of some body parts and non-verbal behaviors of "teaching/learning" students in the case of using a whiteboard.

## 2.4 Detection of Non-verbal Behaviors in the Case of Not Using a Whiteboard

In this case, the "teaching" and "learning" students sitting around a table have the following behaviors; (i) taking notes, (ii) explanation to "learning" students, and (iii) listening to the explanation as shown in Fig. 3. Furthermore, Fig. 3(a) and (b) shows the number of "teaching/learning" students. We can evaluate the number of "teaching/learning" students as the activity for the cooperative learning. Therefore, in this case, we use some body parts "Neck", "Eye", and "Finger" for the detection of each behavior as follows; (i) taking notes, (ii) explanation to "learning" students, and (iii) listening to the explanation. Similarly, we define the positions $\boldsymbol{p}_i^{Table}(t)$ of some body parts for non-verbal behaviors of "teaching" and "learning" students sitting around a table as follows;

$$
\begin{aligned}
\boldsymbol{p}_i^{Table}(t) = (&x_{Neck,i}(t), y_{Neck,i}(t), x_{Eye,i}^L(t), y_{Eye,i}^L(t), x_{Eye,i}^R(t), y_{Eye,i}^R(t), \\
&x_{Hand,i}^L(t), y_{Hand,i}^L(t), x_{Hand,i}^R(t), y_{Hand,i}^R(t))^T,
\end{aligned} \tag{2}
$$

where $i$ denotes the student number. For example, we can summarize the relationships between the behaviors of "learning/teaching" students and the positions $\boldsymbol{p}_i^{Table}(t)$ of some body parts detected by OpenPose as follows;

– Taking notes: $x_{Neck,i}(t) \neq 0$, $x_{Eye,i}^L(t) = x_{Eye,i}^R(t) = 0$, $\{x_{Hand,i}^L(t) \neq 0$ or $x_{Hand,i}^R(t) \neq 0\}$,
– Looking at a whiteboard: $x_{Neck,i}(t) \neq 0$, $\{x_{Eye,i}^L(t) \neq 0$ or $x_{Eye,i}^R(t) \neq 0\}$, $|y_{Eye,i}^R(t) - y_{Eye,i}^L(t)| > 0$,
– Speaking, listening and and/or other students: $x_{Neck,i}(t) \neq 0$, $\{x_{Eye,i}^L(t) \neq 0$ or $x_{Eye,i}^R(t) \neq 0\}$, $|y_{Eye,i}^R(t) - y_{Eye,i}^L(t)| > 0$.

$(x_{Eye}^L(t), y_{Eye}^L(t))$
$(x_{Eye}^R(t), y_{Eye}^R(t))$

$(x_{Neck}(t), y_{Neck}(t))$

$(x_{Hand}^L(t), y_{Hand}^L(t))$
$(x_{Hand}^R(t), y_{Hand}^R(t))$

(a) $n_{Act}$: 0          (b) $n_{Act}$: 3

**Fig. 3.** Positions of some body parts and non-verbal behaviors of "teaching/learning" students in the case of not using a whiteboard ($n_{Act}$ denotes the number of students looking at other students).

## 3 Modeling of Non-verbal Behaviors of Students

### 3.1 Modeling of Non-verbal Behaviors of Students Based on Neural Networks

In the cooperative learning, the non-verbal behaviors of students have strong relations with the understandings and interests for the given contents and the explanation by "teaching" students. Therefore, we have to discuss the interactions between the non-verbal behaviors of all students. In this section, we discuss a modeling method for the interaction among students.

First, we can convert the positions $\boldsymbol{p}^{WB}(t) = (p_m^{WB}(t))$ and $\boldsymbol{p}_i^{Table}(t) = (p_{m,i}^{Table}(t))$ of some body parts represented by the coordinate into the features $\boldsymbol{x}_m(t) = \{x_{m,i}(t)\} = \{x_m^{WB}(t), x_{m,1}^{Table}(t), \cdots, x_{m,P}^{Table}(t)\}$ represented by the binary for the $m$-th event (e.g: Whether there is each behavior?) as follows;

$$x_m^{WB}(t) = \begin{cases} 1 & p_m^{WB}(t) \neq 0, \\ 0 & \text{Otherwise.} \end{cases} \qquad x_{m,i}^{Table}(t) = \begin{cases} 1 & p_{m,i}^{Table}(t) \neq 0, \\ 0 & \text{Otherwise.} \end{cases}$$

where $i(= 1, 2, \cdots, P)$ and $m(= 1, 2, \cdots, M)$ denote the student number and the event number respectively.

Next, non-verbal behaviors of "teaching/learning" students have relationships with each other in the cooperative learning. Therefore, we evaluate the strength of the interactions between behaviors of "teaching/learning" students by using the models with the time-delay. We introduce the following non-linear time-series model for the features $\boldsymbol{x}_m(t) = \{x_{m,i}(t)\}$ represented by the binary for the event (e.g., Whether eyes were detected?) for "teaching/learning" students. Concretely, this model can predict the $m$-th feature of the $i$-th student by the past features $x_{n,k}(t - \ell)$ of all students.

$$x_{m,i}(t) = f\left(\sum_{j=1}^{J} \alpha_{m,i,j} h_{m,j}(t-\ell)\right) + e(t),$$

$$= f\left(\sum_{j=1}^{J} \alpha_{m,i,j} f\left(\sum_{n=1}^{N}\sum_{k=0}^{P}\sum_{\ell=1}^{L} w_{n,j,k,\ell} x_{n,k}(t-\ell)\right)\right) + e(t), \quad (3)$$

where $i$ and $k$ denote the student number ($i = 0$: for the student standing in the front of a whiteboard, $i = 1, \cdots, P$: for the student sitting around a table). $m$ and $n$ denote the event numbers. $e(t)$ denotes a Gaussian noise and $\alpha_{m,i,j}$ denotes the influence of the non-verbal behavior by other students. $J$ denotes the number of hidden units and $h_{m,j}(t-\ell)$ denotes the output of the hidden unit. Moreover, $w_{n,j,k,\ell}$ denotes the time-correlation for the non-verbal behavior of the $k$-th student. The weights $\alpha_{m,i,j}$ and $w_{n,j,k,\ell}$ are initialized by the random number. $f(\cdot)$ denotes the sigmoid function $f(x) = \tanh x$.

Furthermore, the non-linear time-series model defined by Eq. 3 can be represented by the neural network [8] model shown in Fig. 4. The learning object for a neural network model shown in Fig. 4 is to minimize the following error function.

$$E = \sum_{t=1}^{T} E_t = \sum_{t=1}^{T}\sum_{m=1}^{M}\sum_{i=0}^{P} (x_{m,i}(t) - \hat{x}_{m,i}(t))^2, \quad (4)$$

where $T$ denotes the length for the modeling section and $\hat{x}_{m,i}(t)$ denotes the prediction value for the feature $x_{m,i}(t)$ for the $m$-th event of the $i$-th student. Here, we use the forgetting learning algorithm [9] for the purpose of the clarifying of the internal representations of neural networks by the elimination of unnecessary weights.
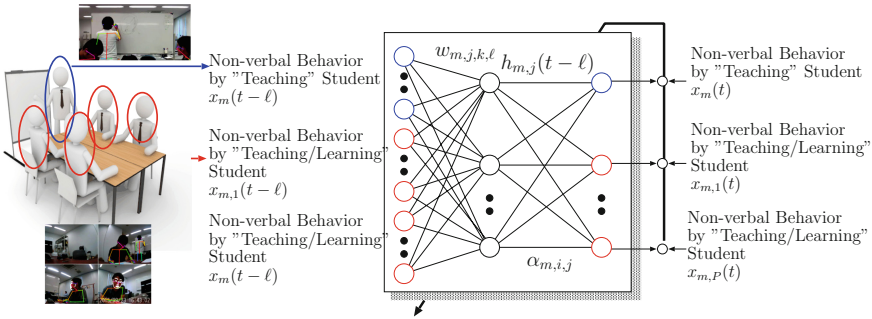


**Fig. 4.** Neural network model for Eq. 3.

## 3.2 Evaluation of the Interaction Based on the Differential Coefficient

We represented the interaction between the behaviors of "teaching/learning" students in Eq. 3. In this equation, the weights $\{\alpha_{m,i,j}\}$ and $\{w_{m,j,k,\ell}\}$ play

important roles on the interaction. Here, we evaluate that the change of the output $x_{m,k}(t)$ (e.g., Whether eyes were detected?) in the output units by the input $x_{m,i}(t - \ell)$ in the input units with the time-delay $\ell$ as follows;

$$\frac{\partial x_{m,k}(t)}{\partial x_{m,i}(t - \ell)} = x'_k(t) \sum_{j=1}^{J} \alpha_{m,i,j} w_{m,j,k,\ell} h'_{m,j}(t - \ell), \tag{5}$$

where $'$ denotes the differential operator. Here, we assume that $x'_k(t) \approx 0$ under the condition which the error function $E$ becomes small.

Here, we define the index $\Delta_{k,i}$ which can evaluate the change of the output $x_k(t)$ (e.g., Whether eyes were detected?) in the output units by the input $x_i(t-\ell)$ in the input units with the time-delay $\ell$.

$$\Delta_{k,i} = \frac{1}{TJLM} \sum_{m=1}^{M} \sum_{t=1}^{T} \sum_{j=1}^{J} \sum_{\ell=1}^{L} \left( \alpha_{m,i,j} w_{m,j,k,\ell} h'_{m,j}(t - \ell) \right)^2. \tag{6}$$

## 4   Experimental Results

### 4.1   Outline of Experiments

We had the two video lectures concerning on the derivation of the formula for two trigonometric functions (law of sines and law of cosines). Four "teaching" and "learning" students are undergraduates. Moreover, we recorded movies for "learning" students sitting around a table by "Meeting Recorder" (Kingjim Co. Ltd., $640 \times 480$ [pixel], 30 [fps]) which is equipped with four cameras and an omni directional microphone. Moreover, we recorded movies for "teaching" students standing in the front of a whiteboard by "MacBook Air" (Apple Co. Ltd., $1280 \times 720$ [pixel], 30 [fps]). The procedure of this experiment is as follows;

1. before-test (about 10 [min]),
2. taking video lectures and taking notes (about 10 [min]),
3. **cooperative learning using a whiteboard (about 10 [min])**,
4. after-test (about 10 [min]).

Table 1 shows scores of before- and after-test and evaluation of notes taken by students for given video lectures. Such scores for tests and notes are evaluated by the three authors. From this table, we can see that Student-A and B have higher scores. The test scores by Student-A, B, and C are improved through the cooperative learning in Lecture-1 and Lecture-2. However, the score of after-test of Student-D is lower ($1.67 \rightarrow 4.00$) than that of before-test in Lecture-2.

**Table 1.** Scores (1:best, 4:worst) of before/after tests and evaluation of notes.

| Student | Lecture-1 | | | Lecture-2 | |
|---|---|---|---|---|---|
| | Test | Note | | Test | Note |
| | (Before/After) | | | (Before/After) | |
| A | 1.33/1.33 | 1.67 | | 2.00/1.67 | 1.33 |
| B | 3.67/1.00 | 2.67 | | 4.00/1.00 | 1.67 |
| C | 3.67/2.00 | 3.00 | | 3.33/2.67 | 1.33 |
| D | 3.33/3.33 | 2.33 | | 1.67/4.00 | 2.33 |
| Ave. | 3.00/1.92 | 2.42 | | 2.75/2.34 | 1.67 |

## 4.2    Features for Behaviors of "Teaching/Learning" Students

Fig. 5(a) shows the behaviors (Writing/Explaining) and the features $\boldsymbol{x}(t)$ for body parts "Neck", "Eyes" and "Hands" of "teaching" students in Lecture-1. In Fig. 5(a), student-A moves to a whiteboard at 350 [sec]. If $x_{Eye}^{L/R}(t) = 0$ and $x_{Hand}^{L/R}(t) \neq 0$, then this "teaching" student is writing the content in a whiteboard. If $x_{Eye}^{L/R}(t) \neq 0$ and $x_{Hand}^{L/R}(t) \neq 0$, then this "teaching" student is explaining to "learning" students sitting around a table.

On the other hand, Fig. 5(b) shows the behaviors (Taking Notes/Looking at other students) and the features $\boldsymbol{x}_B(t)$ for the "learning" student (Student-B) in Lecture-1. In Fig. 5(b), the features change according to the non-verbal behaviors of Student-B ("learning" student). If $x_{Neck,B}(t) \neq 0$, the student is sitting around a table. If $x_{Eye,B}^{L/R}(t) \neq 0$, then the student is looking at other students. If $x_{Eye,B}^{L/R}(t) = 0$ and $x_{Hand,B}^{L/R}(t) \neq 0$, then the student is taking notes.
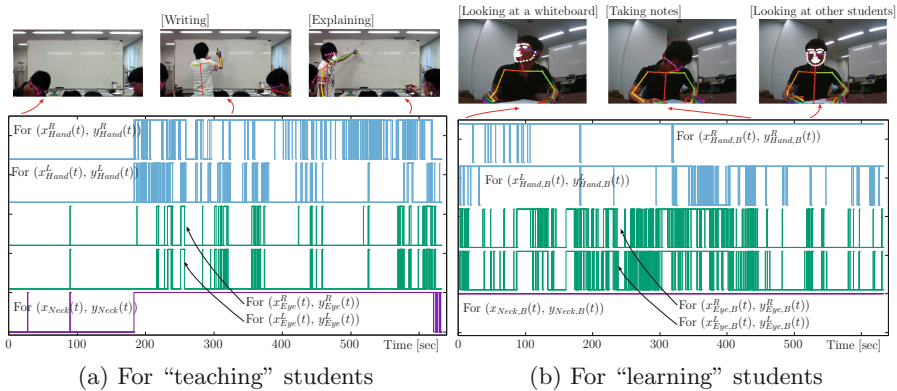


(a) For "teaching" students          (b) For "learning" students

**Fig. 5.** Behaviors and features of "teaching/learning" students (Lecture-1).

### 4.3   Modeling Results of the Non-verbal Behaviors of Students

We used the non-linear time-series model defined by Eq. 3. This model can be represented by the neural network model shown in Fig. 4. Here, we use the following parameters for the size of neural networks;

- the number of students: $P = 4$, the length for the modeling: $L = 10$ [sec],
- the number of body parts: $M = 5$ (Neck, LR-Eyes, LR-Hands),
- the numbers of input, hidden and output units: $L \times M \times (P + 1)$, $J = 10$ and $M \times (P + 1)$ (Here, $P + 1$ includes the standing student).

Table 2 shows the index $\Delta_{k,i}$ (Influence of the $i$-th student on the $k$-th student). Here, "0" denotes the standing student and "A", "B", "C" and "D" denote the sitting student. In Lecture-1, $\Delta_{k,0}$ (0: 2.017, A: 2.349, B: 0.000, C: 0.000, D: 0.000) in the case of $i = 0$ means that Student-0 is influenced by oneself and Student-A. Moreover, Student-A and Student-B are influenced by Student-C ($\Delta_{C,A} = 3.943$ and $\Delta_{C,B} = 8.144$). Furthermore, Student-C and Student-D are influenced by the behavior of oneself ($\Delta_{C,C} = 13.850$ and $\Delta_{D,D} = 84.053$). Similarly, we can discuss the interactions among "teaching/learning" students In Lecture-2.

**Table 2.** $\Delta_{k,i}$: Influence of the $i$-th student on the $k$-th student.

| (a) Lecture-1 | | | | | | (b) Lecture-2 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| $i\backslash k$ | 0 | A | B | C | D | $i\backslash k$ | 0 | A | B | C | D |
| 0 | 2.017 | 2.076 | 3.590 | 1.451 | 2.883 | 0 | 2.629 | 1.128 | 2.240 | 3.908 | 3.763 |
| A | 2.349 | 1.841 | 1.916 | 1.513 | 1.331 | A | 2.888 | 0.012 | 0.010 | 4.445 | 1.412 |
| B | 0.000 | 1.125 | 0.801 | 3.944 | 0.104 | B | 0.000 | 1.181 | 3.209 | 1.023 | 2.830 |
| C | 0.000 | 3.943 | 8.144 | 13.850 | 0.560 | C | 1.400 | 2.493 | 3.528 | 0.675 | 2.260 |
| D | 0.000 | 3.427 | 4.256 | 7.664 | 84.053 | D | 0.000 | 2.381 | 7.937 | 4.009 | 3.332 |

In Fig. 6 shows the sum $X_i(t) = \sum_m x_{m,i}(t)$ of the features $\boldsymbol{x}_m(t) = \{x_{m,i}(t)\}$ for the $m$-th event defined in Sect. 3.1. Here, $i$ denotes the student number ($i = 0$: standing student, $i = A, B, C, D$: sitting students). When the sum $X_i(t)$ becomes large, it means that many body parts of Student-$i$ are recorded by the camera. On the other hand, when the sum $X_i(t)$ becomes small, it means that the movement of the student is small. In Fig. 6(a), Student-A is standing at a whiteboard at 180 [sec] and the sum $X_C(t)$ changes according to the non-verbal behavior. Moreover, $X_C(t)$ and $X_D(t)$ of Student-C and Student-D are comparatively large and they are no influenced by the behavior of standing student (Student-0). Similarly, $\Delta_{C,C}$ and $\Delta_{C,C}$ of Student-C and Student-D in Table 2(a) becomes large. It is shown that the index $\Delta_{k,i}$ has a relation with the sum $X_i(t)$. In Fig. 6(b), Student-A becomes "teaching" student in the sections of [240–405] and [450–470] and Student-C becomes "teaching" student in the sections of [20–238] and [430–450]. Furthermore, the sum of features of Student-B changes at 240 [sec]. It is shown that "learning" Student-B is given the different influences by Student-A and Student-C as the standing "teaching" students.
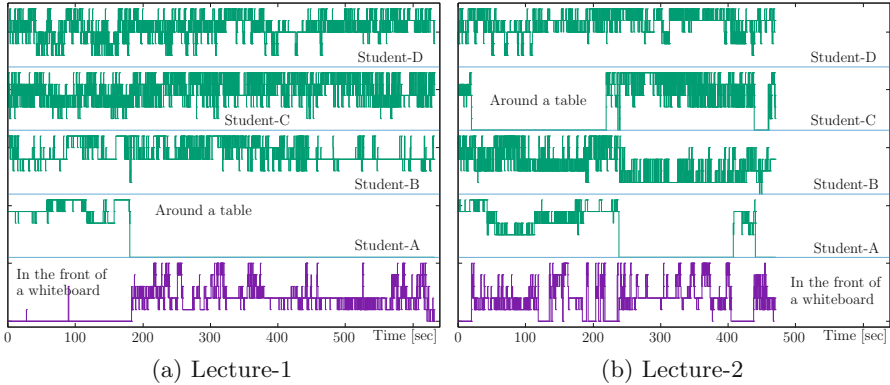
(a) Lecture-1          (b) Lecture-2

**Fig. 6.** $X_i(t) = \sum_m x_{m,i}(t)$ of features of the $i$-th students.

## 5   Conclusions

In this paper, we have discussed the modeling for the interactions between non-verbal behaviors of "teaching/learning" students in a cooperative learning environment. First, we have used the positions of eyes, face and hands detected by OpenPose [7]. Next, we have proposed a modeling method for the interactions of non-verbal behaviors of students based on neural networks. Furthermore, we have proposed the index $\Delta_{k,i}$ for representing "interactions" based on the internal representations of neural networks. From experimental results, we have shown that the index $\Delta_{k,i}$ has give a strong influence on the modeling of the non-verbal behaviors of students. As future work, we would like to discuss carefully the followings; (i) the application to other cases, the usage of various behaviors of students, (ii) the relationships among the modeling results, the progress of the cooperative learning, and the understandings of students.

## References

1. Sugie, S.: An invitation to cooperative learning. Nakanishiya (2011)
2. Johnson, D.W., Johnson, R.T.: Circles of Learning: Cooperation in the Classroom. Interaction Book Co., Edina (1993)
3. Otsuka, K., Araki, S., Ishizuka, K., Fujimoto, M., Heinrich, M., Yamato, J.: A realtime multimodal system for analyzing group meetings by combining face pose tracking and speaker diarization. In: Proceedings of ICMI, pp. 257–264 (2008)
4. Shinnishi, M., Kasuya, Y., Inamoto, H.: Wi-Wi-Meter: a prototype system of evaluating meeting by measuring of activity. IEICE Technical report. HCS2014-63, pp. 19–24 (2014)

5. Martinez-Maldonado, R., Yacef, K., Kay, J.: TSCL: a conceptual model to inform understanding of collaborative learning processes at interactive tabletops. Int. J. Hum.-Comput. Stud. **83**, 62–82 (2015)
6. Watanabe, E., Ozeki, T., Kohama, T.: Analysis of non-verbal behaviors by students in cooperative learning. In: Yoshino, T., Chen, G.-D., Zurita, G., Yuizono, T., Inoue, T., Baloian, N. (eds.) CollabTech 2016. CCIS, vol. 647, pp. 203–211. Springer, Singapore (2016). https://doi.org/10.1007/978-981-10-2618-8_16
7. Cao, Z., Simon, S., Wei, S., Sheikh, Y.: Realtime multi-person 2D pose estimation using part affinity fields. https://arxiv.org/abs/1611.08050. Accessed 12 Dec 2018
8. Rumelhart, D.E., McClelland, J.L.: The PDP Research Group: Parallel Distributed Processing. MIT Press, Cambridge (1986)
9. Ishikawa, M.: Structural learning with forgetting. Neural Netw. **9**(3), 509–521 (1996)
10. Watanabe, E., Ozeki, T., Kohama, T.: Analysis of interactions between lecturer and students. In: Proceedings of the 8th International Conference on Learning Analytics and Knowledge, 5 pages (2018)