

## OpenPose を用いる簡便な手話認識手法

青野 由崇<sup>†</sup> 花泉 弘<sup>†</sup>法政大学情報科学部<sup>†</sup>

## 1. まえがき

手話を行う者が意思疎通を行う上で重要な手段であるが、手話の習得は難しく、健聴者の多くはその意味が理解できない。こうした背景から、手話の自動翻訳が望まれており、様々な取り組みが行なわれてきた[4][5]。文献[1]では、Kinect から得られた骨格ポーズデータの手の座標を取得し、それらをディープラーニングによって各手話単語に対応付けることで認識を行っていた。しかし、認識を行っているのは単体の手話単語のみであり、手話の文章を認識させるには再びその手話の文の映像からディープラーニングで学習させる必要があった。本研究は、特殊なカメラやディープラーニングによる学習を用いないより簡便な手話認識システムの実現を目的としている。ここでは OpenPose[6]を使用して簡単に手首や指の関節点群の動きを得た。パターンマッチングに基づいて認識する手法を提案する。

## 2. 原理と処理手順

使用する手話の映像は CG による手話単語の映像を教師データとし、手話者が行っている手話の映像を判別・翻訳させる。映像は 60FPS の環境で撮影したものである。まず、手話の動作を OpenPose で取得できる関節座標の座標として抽出し、体格の違いを補正して各手話単語のデータベースを構築する。次に、CG と手話者が行っている手話に含まれている単語間の分割するための処理が必要になる。手首の動作速度から次の手話単語の動作に移るときに生じる極小点、極大点を用いて分割を行う[2]。分割したデータをマッチングさせる際に、人によって動作速度の違いがあるので Dynamic Time Warping (DTW) [3]を用いてマッチング処理を施す。分割した各手話単語の手首の移動量からお手本となる CG の手話単語のデータベースよ DTW によるパターンマッチングの類似度を求め手話の判別を行うという手法で手話の判別を行う。

## 2.1 関節点位置

OpenPose で取得できる関節座標は体 25 点、手指 42 点顔 70 点であるが本研究では、上半身の首、右の肩・右ひじ、右手首、左肩、左ひじ、左手首の計 7 点と手指の座標を用いる。上半身のそれぞれ部位を  $P_1, P_2, P_3, P_4, P_5, P_6, P_7$  とし、この中で最も動きの少ない首  $P_1$  を基準にし、各部位までのベクトル  $P_{12}, P_{13}, P_{14}, P_{15}, P_{16}, P_{17}$  を求める。手首の動きが体格の影響をなくすためこれらのベクトルを、首から両肩までの距離の平均値

$$l = \frac{|P_{12}| + |P_{15}|}{2} \quad (1)$$

で割って正規化して体格による大きさの違いを吸収する。各ベクトルに  $k$  点移動平均をかけ、ノイズを低減させる。

次に、単語の分割のため、文献[2]の手法より手話の次単語への手の移動に生じる速度の極小点、極大点を検出のため手話の映像の映像から得られた手首の  $x, y$  座標の移動量を求め、各移動量データに  $\Delta t$  秒で割ってジェスチャーの速度を算出

$$\Delta t = t_{n+1} - t_n \quad (2)$$

$n$  はフレーム番号であり、

$$P(t_n) = (x_n, y_n) \quad (3)$$

の  $n$  フレーム目の手首の位置を用いて

$$v = \frac{|P(t_{n+1}) - P(t_n)|}{\Delta t} \quad (4)$$

と図 1 (a) で示している手首の速度を算出する。手首の速度データ  $v$  から手話単語間に分割する極大値、極小値を求める。手話の単語分割に極大値、極小値より、最大となる極大値の位置を区切りの開始地点とし 0 近い値となる極小値、区切りの終了地点とする。手話に複数の手話単語が含まれている場合、最初の単語の分割が終了時、次の最大となる極大値の位置を区切りの開始地点とし 0 に近い値となる極小値を区切りの終了地点と繰り返し行う。

## 2.4 パターンマッチング

DTW は時間のずれを考慮して、二つの時系列データの距離を算出する方法である。以下に概要を説明する。比較する時系列データ  $X, Y$  をそれぞれ  $N$  個と  $M$  個のデータ量を要素として持っている。

$$X = (x[1], x[2], \dots, x[N]),$$

$$Y = (y[1], y[2], \dots, y[M]) \quad (5)$$

まず、 $X$  の要素  $x[n]$  と  $Y$  の要素  $y[m]$  の類似を表す局コスト値をすべての時刻の組合せで計算する。

$$C(x[n], y[m]) = f(x[n] - y[m])$$

$$(n = 1, \dots, N, m = 1, \dots, M) \quad (6)$$

ここで、 $f$  は要素間の相違を計算する評価関数である。

$$f(x[n] - y[m]) = \|x[n] - y[m]\|_2 \quad (7)$$

次に、再帰的に蓄積コスト値を計算し距離を求める。

$$D(n, m) = \min \begin{cases} D(n-1, m-1) \\ D(n-1, m) \\ D(n, m-1) \end{cases} + c(x[n], y[m]) \quad (8)$$

## 4. 実験および考察

一般のカメラで撮影した CG の手話と人が行っている手話の映像に OpenPose により上半身の関節点座標を取得し、DTW による手話判別を行った。本研究で、使用した手話の映像はテンプレートの CG 手話単語①「私」、②「自転車」、③「壊す」、④「怒られる」、⑤「パン

コン」、⑥「ハワイ」、⑦「行く」、⑧「先生」、⑨「妹」の教師データで構築したデータベースを用いた。判別させる手話は被験者 A,B の2人で行った「私は自転車を壊す」、「私の妹がハワイに行く」、「私が先生に怒られる」を用いた。まずこれらのデータから手話の単語部分を分割するため文献[2]の手法を用いて、手首の速度から極小点、極大点より単語間分割を行った。図1より被験者 A 「私は自転車を壊す」の結果例を示す。

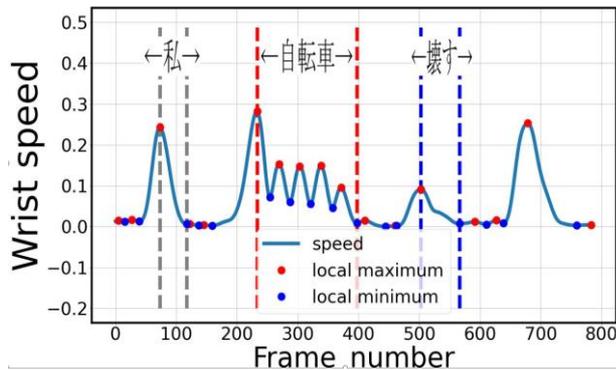


図1: 「私は自転車を壊す」手首の速度および「私」、「自転車」「壊す」部分の分割領域

図1の赤点線の範囲は「私」、黒点線の範囲は「自転車」、青点線の範囲は「壊す」の手話を行っている部分の分割範囲を表している。「自転車」の範囲内に極大点4つ入っているがこの範囲に0に近い値となる極小値がないため手話を続けて行っていると判別して範囲内の極大点4つは区切り位置の場所ではないとしている。

教師データのCG手話単語データと被験者A,Bが行った手話の分割後、これらのデータをパターンマッチングさせる。評価方法として分割した手話単語データと各テンプレートのCG手話単語データをDTWによる類似度で評価した表1から表3より、同じような手首の動作する①「私」や⑨「妹」の場合でも、それぞれ対応した手話単語に対して倍以上の高い類似度が出た。これは、①「私」や⑨「妹」は手の形に違いがあり、手と指の関節点を用いたことでそれぞれ明確に区別できたのではないかと考えられる。その他の手話も表1から表3より太字箇所に対応した手話に対して他の手話よりも倍以上の高い類似度を示している。これら結果から手話に手と指の関節点を用いたことでそれぞれより明確に区別しマッチングできたのではないかと考えられる

表1 「私は自転車を壊す」の類似度

	①	②	③	④	⑤	⑥	⑦	⑧	⑨
A①	<b>4*10<sup>-1</sup></b>	1*10 <sup>-1</sup>	8*10 <sup>-4</sup>	8*10 <sup>-4</sup>	3*10 <sup>-4</sup>	1*10 <sup>-1</sup>	8*10 <sup>-4</sup>	1*10 <sup>-1</sup>	6*10 <sup>-1</sup>
A②	8*10 <sup>-1</sup>	<b>7*10<sup>-1</sup></b>	1*10 <sup>-4</sup>	2*10 <sup>-4</sup>	2*10 <sup>-4</sup>	8*10 <sup>-4</sup>	1*10 <sup>-1</sup>	1*10 <sup>-1</sup>	3*10 <sup>-4</sup>
A③	2*10 <sup>-1</sup>	3*10 <sup>-1</sup>	<b>3*10<sup>-1</sup></b>	1*10 <sup>-4</sup>	1*10 <sup>-4</sup>	2*10 <sup>-4</sup>	8*10 <sup>-4</sup>	9*10 <sup>-4</sup>	3*10 <sup>-4</sup>
B①	<b>6*10<sup>-1</sup></b>	8*10 <sup>-1</sup>	<b>5*10<sup>-8</sup></b>	7*10 <sup>-4</sup>	2*10 <sup>-4</sup>	1*10 <sup>-1</sup>	1*10 <sup>-4</sup>	1*10 <sup>-1</sup>	4*10 <sup>-7</sup>
B②	3*10 <sup>-1</sup>	<b>5*10<sup>-1</sup></b>	1*10 <sup>-4</sup>	1*10 <sup>-4</sup>	1*10 <sup>-4</sup>	8*10 <sup>-4</sup>	1*10 <sup>-1</sup>	1*10 <sup>-1</sup>	2*10 <sup>-4</sup>
B③	1*10 <sup>-1</sup>	2*10 <sup>-1</sup>	<b>4*10<sup>-1</sup></b>	1*10 <sup>-4</sup>	1*10 <sup>-4</sup>	1*10 <sup>-4</sup>	1*10 <sup>-4</sup>	7*10 <sup>-4</sup>	1*10 <sup>-4</sup>

表2 「私の妹がハワイに行く」の類似度

	①	②	③	④	⑤	⑥	⑦	⑧	⑨
A①	<b>8*10<sup>-1</sup></b>	3*10 <sup>-1</sup>	8*10 <sup>-4</sup>	7*10 <sup>-4</sup>	3*10 <sup>-4</sup>	1*10 <sup>-1</sup>	1*10 <sup>-4</sup>	3*10 <sup>-1</sup>	4*10 <sup>-1</sup>
A⑨	5*10 <sup>-7</sup>	3*10 <sup>-4</sup>	2*10 <sup>-4</sup>	4*10 <sup>-4</sup>	6*10 <sup>-4</sup>	5*10 <sup>-1</sup>	7*10 <sup>-1</sup>	1*10 <sup>-1</sup>	<b>7*10<sup>-1</sup></b>
A⑥	1*10 <sup>-1</sup>	4*10 <sup>-4</sup>	4*10 <sup>-4</sup>	3*10 <sup>-4</sup>	4*10 <sup>-1</sup>	<b>5*10<sup>-1</sup></b>	3*10 <sup>-4</sup>	1*10 <sup>-1</sup>	1*10 <sup>-4</sup>
A⑦	1*10 <sup>-4</sup>	4*10 <sup>-4</sup>	2*10 <sup>-4</sup>	4*10 <sup>-4</sup>	2*10 <sup>-1</sup>	4*10 <sup>-1</sup>	<b>8*10<sup>-1</sup></b>	1*10 <sup>-1</sup>	4*10 <sup>-1</sup>
B①	<b>1*10<sup>-1</sup></b>	3*10 <sup>-1</sup>	9*10 <sup>-4</sup>	8*10 <sup>-4</sup>	2*10 <sup>-4</sup>	1*10 <sup>-1</sup>	1*10 <sup>-4</sup>	3*10 <sup>-1</sup>	3*10 <sup>-7</sup>
B⑨	4*10 <sup>-7</sup>	3*10 <sup>-4</sup>	2*10 <sup>-4</sup>	4*10 <sup>-4</sup>	7*10 <sup>-4</sup>	6*10 <sup>-1</sup>	1*10 <sup>-1</sup>	1*10 <sup>-1</sup>	<b>5*10<sup>-1</sup></b>
B⑥	1*10 <sup>-1</sup>	4*10 <sup>-4</sup>	4*10 <sup>-4</sup>	2*10 <sup>-4</sup>	4*10 <sup>-1</sup>	<b>6*10<sup>-1</sup></b>	4*10 <sup>-4</sup>	1*10 <sup>-1</sup>	1*10 <sup>-4</sup>
B⑦	1*10 <sup>-4</sup>	3*10 <sup>-4</sup>	2*10 <sup>-4</sup>	4*10 <sup>-4</sup>	2*10 <sup>-1</sup>	4*10 <sup>-1</sup>	<b>9*10<sup>-1</sup></b>	1*10 <sup>-1</sup>	4*10 <sup>-1</sup>

表3 「私は先生に怒られる」の類似度

	①	②	③	④	⑤	⑥	⑦	⑧	⑨
A①	<b>9*10<sup>-1</sup></b>	3*10 <sup>-1</sup>	9*10 <sup>-4</sup>	8*10 <sup>-4</sup>	3*10 <sup>-4</sup>	1*10 <sup>-1</sup>	8*10 <sup>-4</sup>	2*10 <sup>-1</sup>	4*10 <sup>-1</sup>
A⑧	3*10 <sup>-4</sup>	2*10 <sup>-4</sup>	4*10 <sup>-4</sup>	6*10 <sup>-4</sup>	5*10 <sup>-1</sup>	8*10 <sup>-1</sup>	3*10 <sup>-1</sup>	<b>8*10<sup>-1</sup></b>	8*10 <sup>-4</sup>
A④	3*10 <sup>-4</sup>	1*10 <sup>-1</sup>	1*10 <sup>-4</sup>	<b>8*10<sup>-1</sup></b>	1*10 <sup>-1</sup>	4*10 <sup>-4</sup>	4*10 <sup>-4</sup>	3*10 <sup>-4</sup>	4*10 <sup>-1</sup>
B①	<b>1*10<sup>-1</sup></b>	3*10 <sup>-1</sup>	9*10 <sup>-4</sup>	8*10 <sup>-4</sup>	3*10 <sup>-4</sup>	1*10 <sup>-1</sup>	7*10 <sup>-4</sup>	2*10 <sup>-1</sup>	4*10 <sup>-7</sup>
B⑧	3*10 <sup>-4</sup>	2*10 <sup>-4</sup>	4*10 <sup>-4</sup>	6*10 <sup>-4</sup>	5*10 <sup>-1</sup>	8*10 <sup>-1</sup>	3*10 <sup>-1</sup>	<b>7*10<sup>-1</sup></b>	8*10 <sup>-4</sup>
B④	3*10 <sup>-4</sup>	2*10 <sup>-1</sup>	1*10 <sup>-4</sup>	<b>9*10<sup>-1</sup></b>	1*10 <sup>-1</sup>	4*10 <sup>-4</sup>	4*10 <sup>-4</sup>	3*10 <sup>-4</sup>	4*10 <sup>-1</sup>

#### 4. まとめ

ここでは、画像から人物の関節点と手話の手首の動作速度による単語間の分割、DTWを用いて手話判別手法を提案した。実験では、手首の動作速度の極小点、極大点を用いることで単語部分の分割を行えた。分割した教師データのCG手話データと被験者が行った判別させる手話データとのマッチングはDTWを用いた。今回、使用した判別させる手話は教師データの9つのCG手話単語のデータベース内にある手話単語データに高い類似度を示した。しかし、日常の会話で扱う手話単語数はおよそ500語が目安とされている。そのため手話単語数を増やす必要があるがパターンマッチングの処理時間大きくなるのが問題となる。今後の課題としては、お手本の手話のデータベースを増やしたとき、似ている手話をまとめてグループ化し、マッチングの処理時間の対策することである。

#### 文献

- [1] Al Amin Hosain, Panneer Selvam Santhalingam, Parth Pathak, Huzefa Rangwala, Jana Košecá, "Sign Language Recognition Analysis using Multimodal Data", Proc. DSAA-2019, pp. 1-8, 2019
- [2] 喜安千弥, 藤村貞夫, "動画像による連続手話の認識方式", 計測自動制御学会論文集, Vol.35, No.8, pp.1-6, 1999
- [3] Benjamin Johnen, Bernd Kuhlenkoetter, "A Dynamic Time Warping Algorithm for Industrial Robot Motion Analysis", Proc. CISS-2016, pp. 1-6, 2016
- [4] Shogo Okada; Osamu Hasegawa, "Motion recognition based on Dynamic-Time Warping method with Self-Organizing Incremental Neural Network", Proc. ICPR-2008, pp. 1-4, 2008
- [5] J. Huang, W. Zhou, H. Li, and W. Li, "Sign language recognition using 3d convolutional neural networks," in 2015 IEEE International Conference on Multimedia and Expo (ICME), June 2015, pp. 1-6.
- [6] Z. Cao, T. Simon, S-E Wei, Y.Sheikh, "Realtim e Multi-Person 2D Pose Estimation using Part Affinity Fields", In 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), No.121, pp.1302-1310, July 2017