

CNNに基づくロボットのリーチング動作に関する End-to-End 学習

CNN-based End-to-End Learning for Reaching Movement of Robotic Arm

○学 久田智己 (宇都宮大) 学 浦山一樹 (宇都宮大)
 宮下隼輔 (宇都宮大) 学 柿木泰成 (宇都宮大)
 正 尾崎功一 (宇都宮大) 正 星野智史 (宇都宮大)

Tomoki HISADA, Utsunomiya University
 Kazuki URAYAMA, Utsunomiya University
 Shunsuke MIYASHITA, Utsunomiya University
 Yasunari KAKIGI, Utsunomiya University
 Koichi OZAKI, Utsunomiya University
 Satoshi HOSHINO, Utsunomiya University

For factory automation, robots are required to substitute for workers. For this issue, researchers have paid attention to humanoid and dual-arm robots. In these researches, autonomous generation of motions for target objects is a challenge. In this paper, we focus on the reaching movement. For this challenge, we propose an end-to-end motion planner based on the convolutional neural network, CNN. In contrast to other related motion planners, we use only right and left images captured from the head camera of robot. The network training is executed by learning from demonstration. This allows the robot to map relationship between the images (input) and joint angles (output). Through the experiments, we show that the robot is enabled to generate the reaching movement toward objects located not only at the demonstrated positions, but also at other unknown position.

Key Words: Convolutional Neural Network, Reaching Movement, End-to-End Learning

1 緒言

労働力人口の減少にともない、工場においてロボットの導入が進んでいる。多品種少量生産に適したセル生産方式では、部品の把持や組み立て等、ロボット1台が行う作業の工程数が増加する。先行研究では、複雑な動作生成のために、人による教示およびロボットによる動作の再生を提案した [1]。しかしながら、動作の教示および再生時において、対象物体の位置が異なる場合は作業に失敗するという問題がある。そこで本研究では、ロボット自身による、対象物体の位置変化に応じた動作生成を目指す。

従来研究では、End-to-End 学習を用いたビジュアルフィードバックによるロボットの動作生成法が提案されている [2]。ビジュアルフィードバックとは、カメラから取得される画像を基にロボットの動作を制御する手法である [3]。Yang らは、複数の深層学習器を組み合わせることにより、タオルの折り畳み動作に成功している [4]。この際、ロボットの関節角度および搭載されたカメラから取得される画像を入力に、関節角度を出力として学習を行っている。しかしながら、異なるセンサを入力に用いて学習を行った場合、サンプリングレートの関係からデータ間での対応付けが困難なことがある。その結果、学習はできても、それに基づいた動作生成に失敗する恐れがある。

そこで本研究では、ロボット頭部のカメラから取得される画像のみを入力に用いて動作の生成を試みる。そして、Convolutional Neural Network (CNN) [5] に基づいた End-to-End 学習によるビジュアルフィードバックを提案する。本稿では、作業に必要な可変なロボットの動作の一つであるリーチングに着目する。リーチングとは、対象物体に向かってロボットが手先を移動させる動作である。そして、ロボットが、提案手法に基づき対象物体の位置変化に応じたリーチング動作を生成可能であることを示す。

2 CNN に基づいたビジュアルフィードバック

ロボットは、搭載されたカメラから取得される画像を基に、自らの動作を生成する。そこで、CNN による End-to-End 動作学習を行い、動作生成器を構築する。CNN は、順伝播型のニューラルネットワークの一種であり、画像処理分野において有効性が示されている。本研究では、動作生成器への入力をロボット頭部

のカメラから得られる RGB 画像、出力を手先位置移動量とする。ただし、動作生成器は移動量を出力し続けるため、対象物体付近で手先を停止し、動作を終了することが困難となる。そこで、CNN に基づき終了判定用の識別器も構築する。そして、識別器によって動作の終了判定を行う。これにより、ロボットは対象物体付近でリーチングを終了することが可能となる。図 1 に、動作生成器および識別器に基づいたビジュアルフィードバックの概要を示す。

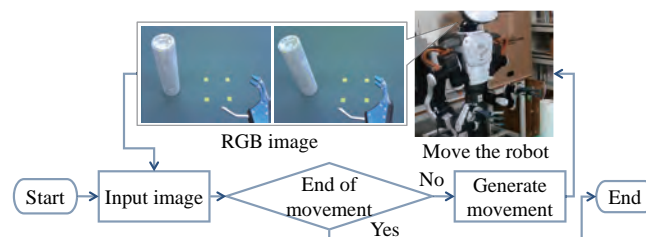


Fig.1 Overview of visual feedback based on CNN for reaching movement of robotic arm

ロボットは、頭部のカメラから画像を取得する。取得された画像を終了判定用の識別器に入力し、動作終了または動作続行の識別を行う。動作続行が識別された際、取得された画像を動作生成器に入力する。動作生成器は、入力された画像に基づき手先位置移動量を出力する。ロボットは、出力された移動量に基づき手先位置を決定する。そして、決定された手先位置へ腕を動かす際の関節角度を求めるため、逆運動学の計算を行う。そして、算出された関節角度に基づき腕を動かす。このフィードバックループは、動作終了の判定が出るまで繰り返される。

CNN への入力に 1 視点のカメラ画像を用いる場合、奥行き方向の情報が含まれない。そのため CNN は、ロボットの正確な手先位置の決定が困難となる。そこで本研究では、CNN への入力に 2 視点のカメラ画像を用いる。これにより、CNN は奥行き方向の情報を考慮して手先位置を決定することができると考えら

れる。

3 CNNに基づく動作生成器および終了判定用識別器

3.1 動作生成器の構築

動作生成器としてのCNNの学習のため、データセットを作成する。対象物体として、円柱を複数か所に置き、ロボットはそれぞれの位置の円柱へリーチングを行う。本研究では、ロボットの右腕のみの動作に焦点を当てる。ロボットの初期手先位置から円柱位置までの動作における逆運動学の計算には、MoveIt!を用いる[6]。そして、算出された関節角度に基づき動作する。このとき、ロボット頭部のステレオカメラから取得されるRGB画像、および次に画像が取得されるまでの右手先位置移動量($\Delta x, \Delta y, \Delta z$)を記録する。

4か所へのリーチングにおいて、右カメラから得られたRGB画像計95枚、および左カメラから得られたRGB画像計95枚、対応した手先位置移動量($\Delta x, \Delta y, \Delta z$)計95個の値をデータセットとする。図2に、作成されたデータセットの例を示す。

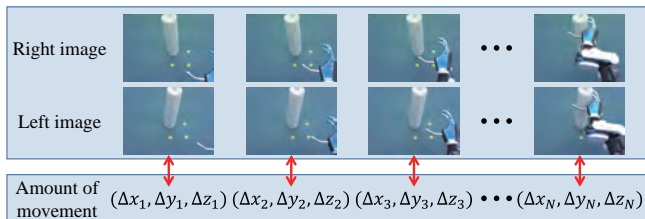


Fig.2 Example of dataset used for End-to-End learning of regression problem

データ数は $N = 95$ である。End-to-End 学習の際、左右画像を入力、手先移動量を出力に対する正解データとして用いる。続いて、入力画像から動作を出力するため、回帰問題に対して用いられるCNNの構造を図3に示す。

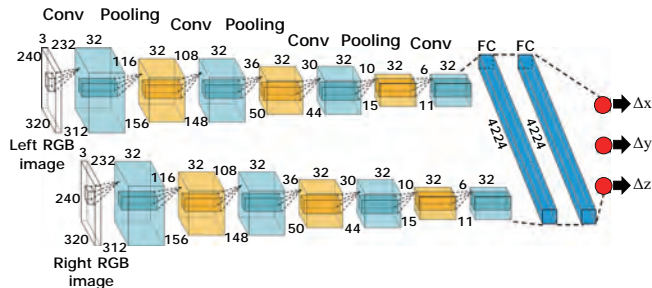


Fig.3 Structure of CNN used for movement generator

入力には、ロボット頭部のステレオカメラから得られる3チャンネルのRGB画像を用いる。ステレオカメラから得られる2視点の画像それぞれについて、入力側から出力側にかけて畳み込み層、プーリング層が繰り返される。

1つ目の畳み込み層のフィルタサイズは $9 \times 9[\text{pix}]$ 、ストライドは $1[\text{pix}]$ である。2つ目の畳み込み層のフィルタサイズは $9 \times 9[\text{pix}]$ 、ストライドは $1[\text{pix}]$ である。3つ目の畳み込み層のフィルタサイズは $7 \times 7[\text{pix}]$ 、ストライドは $1[\text{pix}]$ である。4つ目の畳み込み層のフィルタサイズは $5 \times 5[\text{pix}]$ 、ストライドは $1[\text{pix}]$ である。また、1つ目のプーリング層のフィルタサイズは $2 \times 2[\text{pix}]$ 、ストライドは $2[\text{pix}]$ である。2つ目のプーリング層のフィルタサイズは $3 \times 3[\text{pix}]$ 、ストライドは $3[\text{pix}]$ である。3つ目のプーリング層のフィルタサイズは $3 \times 3[\text{pix}]$ 、ストライドは $3[\text{pix}]$ である。畳み込み層および全結合層における活性化関数には、 ramp 関数を用いる。

そして、2層の全結合層につながり、恒等関数によってロボットの手先位置移動量($\Delta x, \Delta y, \Delta z$)が出力される。学習は、作成したデータセットに基づき、畳み込み層における重みフィルタおよび、全結合層における結合重みを更新することで行う。重みの更新には、Adamを用いる[7]。

3.2 終了判定用識別器の構築

終了判定用識別器としてのCNNの学習のため、3.1節のデータセットを利用する。記録した手先位置移動量($\Delta x, \Delta y, \Delta z$)に基づき、画像に対応した正解ラベルを以下の2クラスとする。

- 動作終了: $\Delta x = 0[\text{m}]$, $\Delta y = 0[\text{m}]$, $\Delta z = 0[\text{m}]$ に対応した画像
- 動作続行: $\Delta x \neq 0[\text{m}]$, $\Delta y \neq 0[\text{m}]$, $\Delta z \neq 0[\text{m}]$ に対応した画像

4か所へのリーチングにおいて、右カメラから得られたRGB画像計95枚、および左カメラから得られたRGB画像計95枚、対応した動作終了および動作続行の2クラスをデータセットとする。図4に、作成されたデータセットの例を示す。

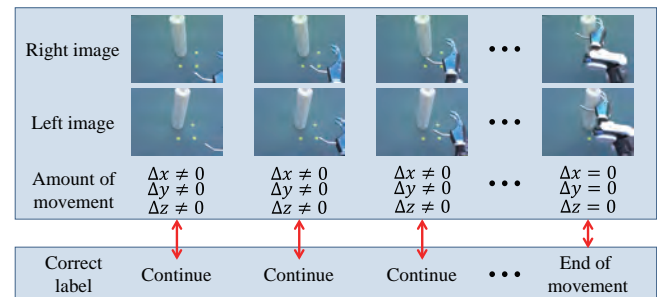


Fig.4 Example of dataset used for End-to-End learning of classification problem

End-to-End 学習の際、一番右の画像に動作終了の正解ラベルを、それ以外の画像に対しては動作続行の正解ラベルと入力に用いる。続いて、入力画像から動作終了または続行の判定結果を出力するため、識別問題に対して用いられるCNNの構造を図5に示す。

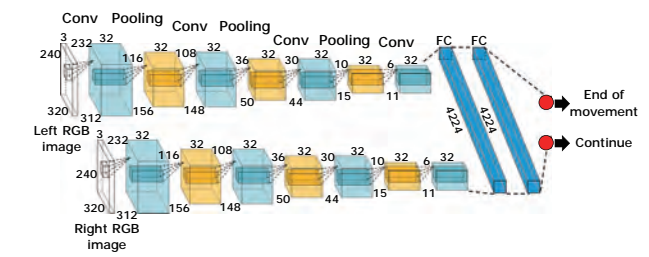


Fig.5 Structure of CNN used for judgement of movement

入力には、ロボット頭部のステレオカメラから得られる3チャンネルのRGB画像を用いる。ステレオカメラから得られる2視点の画像それぞれについて、入力側から出力側にかけて畳み込み層、プーリング層が繰り返される。

1つ目の畳み込み層のフィルタサイズは $9 \times 9[\text{pix}]$ 、ストライドは $1[\text{pix}]$ である。2つ目の畳み込み層のフィルタサイズは $9 \times 9[\text{pix}]$ 、ストライドは $1[\text{pix}]$ である。3つ目の畳み込み層のフィルタサイズは $7 \times 7[\text{pix}]$ 、ストライドは $1[\text{pix}]$ である。4つ目の畳み込み層のフィルタサイズは $5 \times 5[\text{pix}]$ 、ストライドは $1[\text{pix}]$ である。また、1つ目のプーリング層のフィルタサイズは $2 \times 2[\text{pix}]$ 、ストライドは $2[\text{pix}]$ である。2つ目のプーリング層のフィルタサイズは $3 \times 3[\text{pix}]$ 、ストライドは $3[\text{pix}]$ である。3つ目のプーリング層のフィルタサイズは $3 \times 3[\text{pix}]$ 、ストライドは $3[\text{pix}]$ である。畳み込み層および全結合層における活性化関数には、 ramp 関数を用いる。

そして、2層の全結合層につながり、Softmax 関数によって動作終了および動作続行の識別結果が出力される。学習は、作成したデータセットに基づき、畳み込み層における重みフィルタおよび、全結合層における結合重みを更新することで行う。重みの更新には、Adamを用いる。

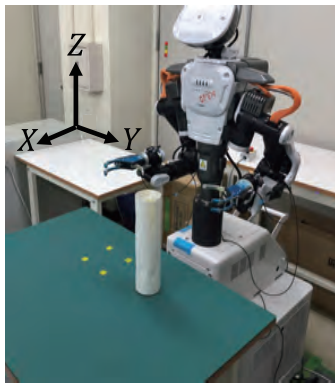
4 ロボットによるリーチング実験

4.1 実験条件

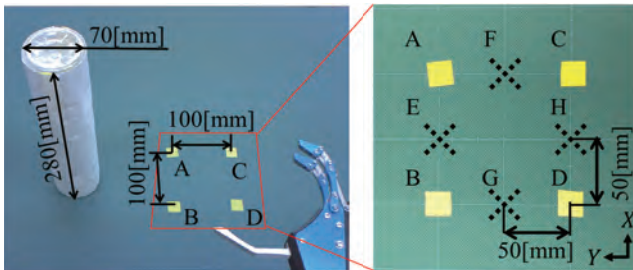
本実験では、以下の2つの手法に基づきロボットにリーチング動作を生成させる。

- 1 視点の画像を入力とする CNN に基づいたビジュアルフィードバック
- 2 2視点の画像を入力とする CNN に基づいたビジュアルフィードバック

比較対象として用いられる1視点画像に対しても、2視点画像と同様に動作生成器および終了判定器を生成し、これをロボットに適用した。そして、それぞれの手法において、動作生成器から出力された手先位置移動量を基にロボットが手先を移動させ、これを繰り返すことにより、リーチング動作を行う。ロボットが手先位置を移動する際、腕の逆運動学の計算には、MoveIt!を用いる。図6に、実験環境を示す。



(a) Dual-arm robot NEXTAGE



(b) Target object on workbench

Fig.6 Experimental environment

図6(a)のロボットは、カワダロボティクス社製のNEXTAGEである[8]。ロボット頭部には、ステレオカメラが搭載されており、ロボット手先には、3本指からなる汎用ハンドである、THK社製のTRK-Sが装着されている。ロボットの正面には作業台が設置され、対象物体として高さ280[mm]、直径70[mm]の円柱が置かれている。図6(b)の黄色の四角で印した位置をそれぞれ位置A~D、破線の×で印した位置をそれぞれ位置E~Hとする。そして、データセット作成に用いた位置A~D、データセット作成時と異なる位置E~Hにそれぞれ円柱を配置した。

4.2 実験結果

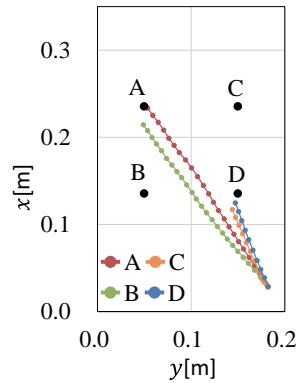
ロボットは、リーチング動作終了後に把持動作を行う。表1に、各入力画像に対するリーチング動作の結果を示す。表中の○は成功、すなわちリーチング動作終了後に円柱を把持、×は失敗、

すなわちリーチング動作終了後に円柱を把持できなかったことを意味している。

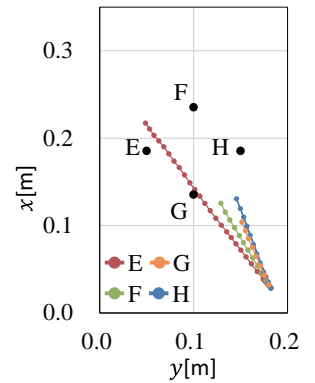
Table 1 Result of reaching movement

Methods	Position of target object							
	A	B	C	D	E	F	G	H
Monocular	○	×	×	○	×	×	×	×
Binocular	○	○	○	○	○	×	×	×

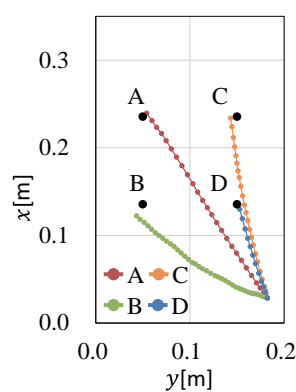
1視点画像によるリーチング動作は、位置Aおよび位置Dの2箇所成功となった。一方、2視点画像によるリーチング動作は位置A~Eの5箇所成功となった。1視点画像および2視点画像によって生成された動作の差異、ならびに2視点画像におけるリーチング動作失敗の要因について考察するため、リーチング動作時の手先位置の軌跡に着目する。図7に、各入力画像に対する位置A~Hの円柱へのリーチング動作におけるロボットの手先の移動軌跡を示す。図中の黒色の点は、円柱を配置した位置A~Hの座標を表している。



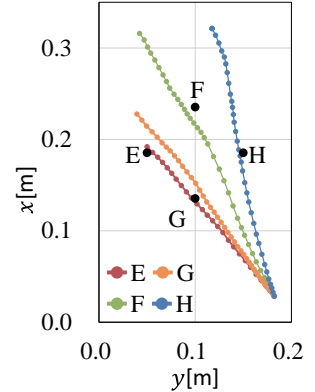
(a) Monocular : A~D



(b) Monocular : E~H



(c) Binocular : A~D



(d) Binocular : E~H

Fig.7 Trajectory of the hand position

図7(a)における位置Bへの軌跡、図7(b)における位置Gへの軌跡は、対応していない位置に接近していることが分かる。これは、1視点画像が対象物体の位置に応じた動作の生成に失敗したことを示している。一方、図7(c)および図7(d)における位置A~Hへの軌跡は、全て対応した位置に接近していることが分

る。これは、2 視点画像が対象物体の位置に応じた動作の生成できたことを示している。以上のことから、ステレオカメラの2 視点画像を入力とすることの有効性が示された。しかしながら、図7(d) の位置 F~H への軌跡は、終了判定用識別器による動作終了の判定がなされなかったため、円柱の位置を通過している。これは、終了判定用識別器の汎化性が不足していることが原因であると考えられる。これに関しては、A~D 以外の位置に置かれた対象物体に対する把持画像を入力として学習を行うことで、解決することができるものと考えられる。

5 結語

本稿では、対象物体の位置に応じたロボットの自律的なリーチング動作の生成を目的に、CNN に基づいた End-to-End 学習によるビジュアルフィードバックを提案した。データセット作成時の位置、およびデータセット作成時と異なる位置に置かれた円柱に対するリーチング実験を行い、ステレオカメラによる2 視点の画像を用いることで、ロボットは画像のみから対象物体の位置に応じた動作を生成できることを示した。しかしながら、終了判定用識別器における汎化性の不足により、データセット作成時と異なる位置の対象物体へのリーチング動作において、適切な位置で動作を終了できず、動作に失敗する結果となった。今後は、終了判定用識別器を構築するためのデータセットを改善することにより、識別器の汎化性を向上させ、より多くの位置へのリーチング動作の成功を目指す。

参考文献

- [1] 浦山一樹 他, “双腕型ロボットに対するティーチングプレイバックシステム,” ロボティクス・メカトロニクス講演会, 2A1-G14, 2018.
- [2] 尾形哲也, “深層学習とマニピュレーション,” 日本ロボット学会誌, Vol. 35, No. 1, pp. 28–31, 2017.
- [3] 橋本浩一, “ビジュアルフィードバック制御と今後,” 日本ロボット学会誌, Vol. 27, No. 4, pp. 400–404, 2009.
- [4] P.-C. Yang *et al.*, “Repeatable Folding Task by Humanoid Robot Worker Using Deep Learning,” *IEEE Robotics and Automation Letters*, Vol. 2, Issue 2, pp. 397–403, 2017.
- [5] Y.Lecun *et al.*, “Gradient-Based Learning Applied to Document Recognition,” *Proceedings of the IEEE*, Vol. 86, Issue 11, pp. 2278–2324, 1998.
- [6] S. Chitta *et al.*, “Moveit![ROS topics],” *IEEE Robotics and Automation Magazine*, Vol. 19, No. 1, pp. 18–19, 2012.
- [7] D. Kingma and J. Ba, “Adam: A Method for Stochastic Optimization,” *Proceedings International Conference for Learning Representations*, pp. 1–15, 2014.
- [8] 五十棲隆勝, 中畑光明, “ヒトとともに未来に向かって進化する次世代産業用ロボット「NEXTAGE」,” 日本機械学会特集, Vol. 30, pp. 12–24, 2011.