

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2018-126798

(P2018-126798A)

(43) 公開日 平成30年8月16日(2018.8.16)

(51) Int.Cl.
B25J 13/08 (2006.01)F I
B25J 13/08テーマコード (参考)
3C707

審査請求 未請求 請求項の数 13 O L (全 47 頁)

(21) 出願番号 特願2017-19313 (P2017-19313)
(22) 出願日 平成29年2月6日(2017.2.6)(71) 出願人 000002369
セイコーエプソン株式会社
東京都新宿区新宿四丁目1番6号
(74) 代理人 100116665
弁理士 渡辺 和昭
(74) 代理人 100164633
弁理士 西田 圭介
(74) 代理人 100179475
弁理士 仲井 智至
(72) 発明者 長谷川 浩
長野県諏訪市大和3丁目3番5号 セイコーエプソン株式会社内
Fターム(参考) 3C707 KS03 KS04 KS23 KS24 KS33
KT01 KT05 KT11 LU02 LU09
LV24 LW08 LW12 LW15 MT06

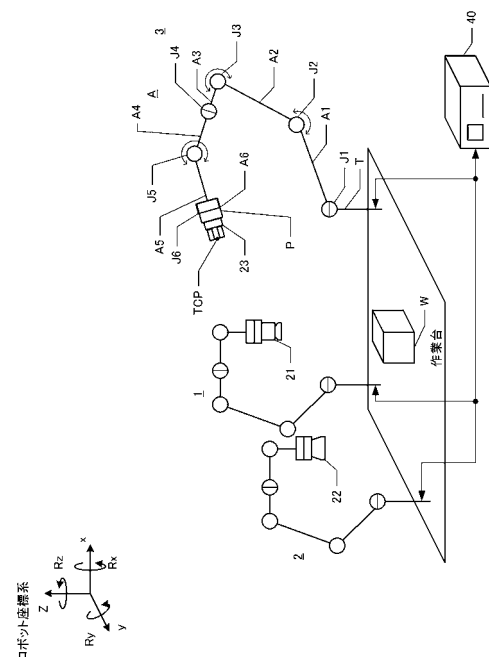
(54) 【発明の名称】 制御装置、ロボットおよびロボットシステム

(57) 【要約】

【課題】ロボットの性能を十分に引き出せる力制御パラメーターを設定することはやはり困難であった。

【解決手段】機械学習を用いて、ロボットの力制御に関する力制御パラメーターを算出する算出部と、算出された前記力制御パラメーターに基づいて前記ロボットを制御する制御部と、を備える制御装置が構成される。

【選択図】図1



【特許請求の範囲】**【請求項 1】**

機械学習を用いて、ロボットの力制御に関する力制御パラメータを算出する算出部と

、

算出された前記力制御パラメータに基づいて前記ロボットを制御する制御部と、
を備える制御装置。

【請求項 2】

前記力制御パラメータは、

前記ロボットがインピーダンス制御で動作する際のインピーダンスパラメータを含む、

請求項 1 に記載の制御装置。

【請求項 3】

前記力制御パラメータは、

前記ロボットの動作の始点と終点との少なくとも一方を含む、

請求項 1 または請求項 2 のいずれかに記載の制御装置。

【請求項 4】

前記力制御パラメータは、

前記ロボットのツールセンターポイントからのオフセット点の位置を含む、

請求項 1 ~ 請求項 3 のいずれかに記載の制御装置。

【請求項 5】

前記算出部は、

状態変数として、少なくとも前記ロボットの位置情報を観測する状態観測部と、

前記状態変数に基づいて前記力制御パラメータを学習する学習部と、を含む、

請求項 1 ~ 請求項 4 のいずれかに記載の制御装置。

【請求項 6】

前記位置情報は、

前記ロボットが備える慣性センサーの出力と、前記ロボットの外部に配置された位置
検出部の出力と、の少なくとも一方に基づいて算出される、

請求項 5 に記載の制御装置。

【請求項 7】

前記学習部は、

前記状態変数に基づいて前記力制御パラメータを変化させる行動を決定し、前記力
制御パラメータを最適化する、

請求項 5 または請求項 6 のいずれかに記載の制御装置。

【請求項 8】

前記学習部は、

前記ロボットが行った作業の良否に基づいて、前記行動による報酬を評価する、

請求項 7 に記載の制御装置。

【請求項 9】

前記学習部は、

前記作業が正常に完了した場合、前記作業の所要時間が基準よりも短い場合、の少な
くとも 1 つにおいて前記報酬を正と評価する、

請求項 8 に記載の制御装置。

【請求項 10】

前記学習部は、

前記ロボットが破損した場合、前記ロボットの作業対象である対象物が破損した場合
、の少なくとも 1 つにおいて前記報酬を負と評価する、

請求項 8 または請求項 9 のいずれかに記載の制御装置。

【請求項 11】

前記算出部は、

10

20

30

40

50

前記状態変数の観測と、当該状態変数に応じた前記行動の決定と、当該行動によって得られる前記報酬の評価とを繰り返すことによって、前記力制御パラメーターを最適化する、

請求項 8 ～ 請求項 10 のいずれかに記載の制御装置。

【請求項 12】

請求項 1 ～ 請求項 11 のいずれかに記載された制御装置によって制御されるロボット。

【請求項 13】

請求項 1 ～ 請求項 11 のいずれかに記載された制御装置と、前記制御装置によって制御される前記ロボットと、
を備えるロボットシステム。

10

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、制御装置、ロボットおよびロボットシステムに関する。

【背景技術】

【0002】

ロボットに作業を行わせるためには、各種の設定が必要であり、従来、各種の設定は人為的に行われている。しかし、当該設定を行うためには高度なノウハウが必要であり、難易度が高い。そこで、従来、予め決められた手順に従ってインピーダンス制御系の慣性パラメーターと粘性パラメーターを調整する技術が開発されている（例えば、特許文献 1）。

20

【0003】

また、従来、工作機械の工具補正の頻度を最適化するために機械学習を利用した技術が知られている（特許文献 2）。

【先行技術文献】

【特許文献】

【0004】

【特許文献 1】特許第 4962551 号公報

【特許文献 2】特許第 5969676 号公報

【発明の概要】

30

【発明が解決しようとする課題】

【0005】

従来の技術を利用しても、パラメーターの設定には高度な知見が必須であった。例えば、特許文献 1 においては、振動回数と閾値とを比較する必要がある。しかし、理想的な閾値を予め決めることは実際には難しい。また、予め決められた手順での調整によってパラメーターが理想的な値になるとは限らない（他の手順であればより理想化できる可能性がある）。従って、ロボットの性能を十分に引き出せる力制御パラメーターを設定することはやはり困難であった。

【課題を解決するための手段】

【0006】

40

上記課題の少なくとも一つを解決するために、制御装置は、機械学習を用いて、ロボットの力制御に関する力制御パラメーターを算出する算出部と、算出された力制御パラメーターに基づいてロボットを制御する制御部と、を備える。この構成によれば、人為的に決められた力制御パラメーターよりも高性能に力制御を行う力制御パラメーターを高い確率で算出することができる。

【0007】

さらに、力制御パラメーターは、ロボットがインピーダンス制御で動作する際のインピーダンスパラメーターを含む構成であっても良い。この構成によれば、人為的な調整によって適切な設定を行うことが困難な、インピーダンスパラメーターを自動的に調整することができる。

50

【 0 0 0 8 】

さらに、力制御パラメータは、ロボットの動作の始点と終点との少なくとも一方を含む構成であっても良い。この構成によれば、人為的に設定された始点や終点を、より高性能に力制御を行うように自動的に調整することができる。

【 0 0 0 9 】

さらに、力制御パラメータは、ロボットのツールセンターポイントからのオフセット点の位置を含む構成であっても良い。この構成によれば、人為的な調整によって適切な設定を行うことが困難な、ツールセンターポイントからのオフセット点の位置を自動的に調整することができる。

【 0 0 1 0 】

さらに、算出部は、状態変数として、少なくともロボットの位置情報を観測する状態観測部と、状態変数に基づいて力制御パラメータを学習する学習部と、を含む構成であっても良い。この構成によれば、高性能な力制御を行う力制御パラメータを容易に算出することができる。

【 0 0 1 1 】

さらに、位置情報は、ロボットが備える慣性センサーの出力と、ロボットの外部に配置された位置検出部の出力と、の少なくとも一方に基づいて算出される構成であっても良い。慣性センサーによれば、ロボットで汎用的に使用されるセンサーに基づいて位置情報を算出することができる。ロボットの外部に配置された位置検出部は、ロボットの動作に影響されることなく位置情報を算出することができる。

【 0 0 1 2 】

さらに、学習部は、状態変数に基づいて力制御パラメータを変化させる行動を決定し、力制御パラメータを最適化する構成であっても良い。この構成によれば、ロボットの使用環境に応じた力制御パラメータとなるように最適化することができる。

【 0 0 1 3 】

さらに、学習部は、ロボットが行った作業の良否に基づいて、行動による報酬を評価する構成であっても良い。この構成によれば、ロボットの作業の質を高めるように力制御パラメータを最適化することができる。

【 0 0 1 4 】

さらに、学習部は、作業が正常に完了した場合、作業の所要時間が基準よりも短い場合、の少なくとも1つにおいて報酬を正と評価する構成であっても良い。作業が正常に完了した場合に報酬を正と評価する構成によれば、ロボットの作業を成功させる力制御パラメータを容易に算出することができる。作業の所要時間が基準よりも短い場合に報酬を正と評価する構成によれば、ロボットを短い時間で作業させる力制御パラメータを容易に算出することができる。

【 0 0 1 5 】

さらに、学習部は、ロボットが破損した場合、ロボットの作業対象である対象物が破損した場合、の少なくとも1つにおいて報酬を負と評価する構成であっても良い。ロボットが破損した場合に報酬を負と評価する構成によれば、ロボットを破損させる可能性が低い力制御パラメータを容易に算出することができる。ロボットの作業対象である対象物が破損した場合に報酬を負と評価する構成によれば、対象物を破損させる可能性が低い力制御パラメータを容易に算出することができる。

【 0 0 1 6 】

さらに、算出部は、状態変数の観測と、当該状態変数に応じた行動の決定と、当該行動によって得られる報酬の評価とを繰り返すことによって、力制御パラメータを最適化する構成であっても良い。この構成によれば、力制御パラメータを自動的に最適化することができる。

【 図面の簡単な説明 】

【 0 0 1 7 】

【 図 1 】 ロボットシステムの斜視図である。

10

20

30

40

50

- 【図 2】制御装置の機能ブロック図である。
- 【図 3】パラメーターを示す図である。
- 【図 4】加減速特性を示す図である。
- 【図 5】ピックアップ処理のフローチャートである。
- 【図 6】算出部に関連する構成のブロック図である。
- 【図 7】光学パラメーターを学習する際の例を示す図である。
- 【図 8】多層ニューラルネットワークの例を示す図である。
- 【図 9】学習処理のフローチャートである。
- 【図 10】動作パラメーターを学習する際の例を示す図である。
- 【図 11】力制御パラメーターを学習する際の例を示す図である。
- 【発明を実施するための形態】
- 【0018】

10

以下、本発明の実施形態について添付図面を参照しながら以下の順に説明する。なお、各図において対応する構成要素には同一の符号が付され、重複する説明は省略される。

- (1) ロボットシステムの構成：
- (2) ロボットの制御：
- (3) ピックアップ処理：
- (4) 学習処理：
- (4-1) 光学パラメーターの学習：
- (4-2) 光学パラメーターの学習例：
- (4-3) 動作パラメーターの学習：
- (4-4) 動作パラメーターの学習例：
- (4-5) 力制御パラメーターの学習：
- (4-6) 力制御パラメーターの学習例：
- (5) 他の実施形態：

20

【0019】

- (1) ロボットシステムの構成：

図 1 は本発明の一実施形態にかかる制御装置で制御されるロボットを示す斜視図である。本発明の一実施例としてのロボットシステムは、図 1 に示すように、ロボット 1 ～ 3 を備えている。ロボット 1 ～ 3 はエンドエフェクターを備える 6 軸ロボットであり、ロボット 1 ～ 3 には異なるエンドエフェクターが取り付けられている。すなわち、ロボット 1 には、撮像部 2 1 が取り付けられ、ロボット 2 には照明部 2 2 が取り付けられ、ロボット 3 にはグリッパー 2 3 が取り付けられている。なお、ここでは、撮像部 2 1 および照明部 2 2 を光学系と呼ぶ。

30

【0020】

ロボット 1 ～ 3 は、制御装置 4 0 によって制御される。制御装置 4 0 はケーブルによりロボット 1 ～ 3 と通信可能に接続される。なお、制御装置 4 0 の構成要素がロボット 1 に備えられていても良い。また、制御装置 4 0 は複数の装置によって構成されても良い（例えば、後述する学習部と制御部とが異なる装置に備えられる等）。また、制御装置 4 0 は、図示しない教示装置をケーブル、または無線通信によって接続可能である。教示装置は、専用のコンピューターであってもよいし、ロボット 1 を教示するためのプログラムがインストールされた汎用のコンピューターであってもよい。さらに、制御装置 4 0 と教示装置とは、一体に構成されていてもよい。

40

【0021】

ロボット 1 ～ 3 は、アームに各種のエンドエフェクターを装着して使用される単腕ロボットであり、本実施形態において、ロボット 1 ～ 3 においてアームや軸の構成は同等である。図 1 においてはロボット 3 においてアームや軸の構成を説明する符号が付されている。ロボット 3 において示されたように、ロボット 1 ～ 3 は、基台 T と、6 個のアーム部材 A 1 ～ A 6 と、6 個の関節 J 1 ～ J 6 を備える。基台 T は作業台に固定されている。基台 T と 6 個のアーム部材 A 1 ～ A 6 は関節 J 1 ～ J 6 によって連結される。アーム部材 A 1

50

～ A 6 とエンドエフェクターは可動部であり、これらの可動部が動作することによってロボット 1 ～ 3 は各種の作業を行うことができる。

【 0 0 2 2 】

本実施形態において、関節 J 2、J 3、J 5 は曲げ関節であり、関節 J 1、J 4、J 6 はねじり関節である。アーム A のうち最も先端側のアーム部材 A 6 には、力覚センサー P とエンドエフェクターとが装着される。ロボット 1 ～ 3 は、6 軸のアームを駆動させることによって、可動範囲内においてエンドエフェクターを任意の位置に配置し、任意の姿勢（角度）とすることができる。

【 0 0 2 3 】

ロボット 3 が備えるエンドエフェクターはグリッパ 2 3 であり、対象物 W を把持することができる。ロボット 2 が備えるエンドエフェクターは照明部 2 2 であり、照射範囲に光を照射することができる。ロボット 1 が備えるエンドエフェクターは撮像部 2 1 であり、視野内の画像を撮像することができる。本実施形態においては、ロボット 1 ～ 3 が備えるエンドエフェクターに対して相対的に固定された位置がツールセンターポイント（TCP）として定義される。TCP の位置はエンドエフェクターの基準の位置となり、TCP が原点となり、エンドエフェクターに対して相対的に固定された 3 次元直交座標系である TCP 座標系が定義される。

【 0 0 2 4 】

力覚センサー P は、6 軸の力検出器である。力覚センサー P は、力覚センサー上の点を原点とした 3 次元直交座標系であるセンサー座標系において互いに直交する 3 個の検出軸と平行な力の大きさと、当該 3 個の検出軸まわりのトルクの大きさとを検出する。なお、本実施例では 6 軸ロボットを例にしているが、ロボットの態様は種々の態様であっても良いし、ロボット 1 ～ 3 の態様が異なってもよい。また、関節 J 6 以外の関節 J 1 ～ J 5 のいずれか 1 つ以上に力検出器としての力覚センサーを備えても良い。

【 0 0 2 5 】

ロボット 1 ～ 3 が設置された空間を規定する座標系をロボット座標系というとき、ロボット座標系は、水平面上において互いに直交する x 軸と y 軸と、鉛直上向きを正方向とする z 軸とによって規定される 3 次元の直交座標系である（図 1 参照）。z 軸における負の方向は概ね重力方向と一致する。また x 軸周りの回転角を R_x で表し、y 軸周りの回転角を R_y で表し、z 軸周りの回転角を R_z で表す。x, y, z 方向の位置により 3 次元空間における任意の位置を表現でき、 R_x , R_y , R_z 方向の回転角により 3 次元空間における任意の姿勢を表現できる。以下、位置と表記した場合、姿勢も意味し得ることとする。また、力と表記した場合、トルクも意味し得ることとする。

【 0 0 2 6 】

なお、本実施形態においてはロボットに作用する力を制御する力制御が実行可能であり、力制御においては、任意の点に作用する当該作用力が目標力になるように制御される。各種の部位に作用する力は、3 次元直交座標系である力制御座標系において定義される。目標力（トルクを含む）は、力制御座標系で表現された力の作用点を起点としたベクトルで表現可能であり、後述する学習が行われる以前において、目標力ベクトルの起点は力制御座標系の原点であり、作用力の方向は力制御座標系の 1 軸方向と一致している。ただし、後述する学習が行われた場合、目標力ベクトルの起点は力制御座標系の原点と異なり得るし、目標力ベクトルの方向は力制御座標系の軸方向と異なり得る。

【 0 0 2 7 】

本実施形態において各種の座標系の関係は予め定義されており、各種の座標系での座標値は互いに変換可能である。すなわち、TCP 座標系、センサー座標系、ロボット座標系、力制御座標系における位置やベクトルは互いに変換可能である。ここでは簡単のため、制御装置 4 0 が TCP の位置および TCP に作用する作用力をロボット座標系で制御する説明をするが、ロボット 1 ～ 3 の位置やロボット 1 ～ 3 に作用する力は、各種の座標系で定義でき、互いに変換可能であるため、位置や力がどの座標系で定義され、制御されても良い。むろん、ここで述べた座標系以外にも他の座標系（例えば対象物に固定されたオブ

10

20

30

40

50

ジェクト座標系等)が定義され、変換可能であっても良い。

【0028】

(2) ロボットの制御：

ロボット1は、教示を行うことにより各種作業が可能となる汎用ロボットであり、図2に示すようにアクチュエーターとしてのモーターM1～M6と、センサーとしてのエンコーダーE1～E6とを備える。アームを制御することはモーターM1～M6を制御することを意味する。モーターM1～M6とエンコーダーE1～E6とは、関節J1～J6のそれぞれに対応して備えられており、エンコーダーE1～E6はモーターM1～M6の回転角度を検出する。また、各モーターM1～M6には電力を供給する電源線が接続されており、各電源線には電流計が備えられている。従って、制御装置40は、各モーターM1～M6に供給された電流を計測することができる。

10

【0029】

制御装置40は、コンピューター等のハードウェア資源と記憶部44に記憶された各種のソフトウェア資源を備え、プログラムを実行可能である。本実施形態において制御装置40は、算出部41、検出部42、制御部43として機能する。なお、ハードウェア資源は、CPU, RAM, ROM等からなる構成であっても良いし、ASIC等によって構成されても良く、種々の構成を採用可能である。

【0030】

本実施形態において検出部42は対象物を検出する処理を実行することが可能であり、制御部43はロボット1～3のアームを駆動することが可能である。検出部42は、光学系20を構成する撮像部21と照明部22とに接続されている。検出部42は、撮像部21を制御し、撮像部21が備える撮像センサーによって撮像された画像を取得することができる。また、検出部42は、照明部22を制御し、出力光の明るさを変化させることができる。

20

【0031】

撮像部21から画像が出力されると、検出部42は、撮像画像に基づいてテンプレートマッチング処理を行い、対象物の位置(位置姿勢)を検出する処理を行う。すなわち、検出部42は、記憶部44に記憶されたテンプレートデータ44cに基づいてテンプレートマッチング処理を実行する。テンプレートデータ44cは複数の位置姿勢毎のテンプレートである。従って、テンプレートデータ44cに対して位置姿勢をID等に対応づけておけば、適合したテンプレートデータ44cの種類によって検出部42から見た対象物の位置姿勢を特定することができる。

30

【0032】

具体的には、検出部42は、複数の位置姿勢毎のテンプレートデータ44cを順次処理対象とし、テンプレートデータ44cの大きさを変化させながら、撮像された画像と比較する。そして、検出部42は、テンプレートデータ44cと画像との差分が閾値以下の像を対象物の像として検出する。

【0033】

対象物の像が検出されると、検出部42は、予め決められた座標系の関係と適合したテンプレートデータ44cの大きさに基づいて対象物の位置姿勢を特定する。すなわち、テンプレートデータ44cの大きさから撮像部21と対象物との光軸方向の距離が判明し、画像内で検出された対象物の位置から光軸に垂直な方向の位置が判明する。

40

【0034】

そこで、例えば、撮像部21の撮像センサーの光軸と撮像平面上の2軸とがTCP座標系の各軸に平行に定義されている場合であれば、検出部42は、テンプレートデータ44cの大きさと、テンプレートデータ44cが画像と適合した位置とに基づいて、TCP座標系において対象物の位置を特定することができる。また、検出部42は、適合したテンプレートデータ44cのIDに基づいて、TCP座標系における対象物の姿勢を特定することができる。このため、検出部42は、上述の座標系の対応関係を利用し、任意の座標系、例えば、ロボット座標系における対象物の位置姿勢を特定することができる。

50

【 0 0 3 5 】

なお、テンプレートマッチング処理は、対象物の位置姿勢を特定するための処理であれば良く、種々の処理を採用可能である。例えば、テンプレートデータ 4 4 c と画像との差分は、階調値の差分によって評価されても良いし、画像の特徴（例えば、画像の勾配等）の差分によって評価されても良い。

【 0 0 3 6 】

検出部 4 2 は、パラメータを参照して当該テンプレートマッチング処理を行う。すなわち、記憶部 4 4 には、各種のパラメータ 4 4 a が記憶されており、当該パラメータ 4 4 a には、検出部 4 2 の検出に関するパラメータが含まれている。図 3 は、パラメータ 4 4 a の例を示す図である。図 3 に示す例において、パラメータ 4 4 a は、光学パラメータと動作パラメータと力制御パラメータとを含んでいる。

10

【 0 0 3 7 】

光学パラメータは、検出部 4 2 の検出に関するパラメータである。動作パラメータと力制御パラメータとはロボット 1 ~ 3 を制御する際のパラメータであり、詳細は後述する。光学パラメータは、撮像部 2 1 に関する撮像部パラメータと、照明部 2 2 に関する照明部パラメータと、撮像部 2 1 によって撮像された対象物の画像に対する画像処理に関する画像処理パラメータとが含まれる。

【 0 0 3 8 】

図 3 においては、これらのパラメータの例が示されている。すなわち、対象物を撮像する際に撮像部 2 1 が配置される位置が撮像部の位置として定義され、撮像部パラメータに含まれている。また、撮像部 2 1 は、露光時間と絞りを調整可能な機構を備えており、対象物を撮像する際の露光時間および絞りの値が撮像部パラメータに含まれている。なお、撮像部の位置は、種々の手法で記述されて良く、例えば、撮像部 2 1 の T C P の位置がロボット座標系で記述される構成等を採用可能である。

20

【 0 0 3 9 】

検出部 4 2 は、撮像部パラメータを参照し、撮像部 2 1 の位置を後述する位置制御部 4 3 a に受け渡す。この結果、位置制御部 4 3 a は、目標位置 L t を生成し、当該目標位置 L t に基づいてロボット 1 を制御する。また、検出部 4 2 は、撮像部パラメータを参照し、撮像部 2 1 の露光時間と絞りを設定する。この結果、撮像部 2 1 においては当該露光時間と絞りによって撮像が行われる状態となる。

30

【 0 0 4 0 】

また、対象物を撮像する際に照明部 2 2 が配置される位置が照明部の位置として定義され、照明部パラメータに含まれている。また、照明部 2 2 は、明るさを調整可能な機構を備えており、対象物を撮像する際の明るさの値が照明部パラメータに含まれている。照明部の位置も、種々の手法で記述されて良く、例えば、照明部 2 2 の T C P の位置がロボット座標系で記述される構成等を採用可能である。

【 0 0 4 1 】

検出部 4 2 は、照明部パラメータを参照し、照明部 2 2 の位置を後述する位置制御部 4 3 a に受け渡す。この結果、位置制御部 4 3 a は、目標位置 L t を生成し、当該目標位置 L t に基づいてロボット 2 を制御する。また、検出部 4 2 は、照明部パラメータを参照し、照明部 2 2 の明るさを設定する。この結果、照明部 2 2 においては当該明るさの光が出力される状態となる。

40

【 0 0 4 2 】

検出部 4 2 は、撮像部 2 1 によって撮像された画像に対してテンプレートマッチング処理を適用する際に、画像処理パラメータを参照する。すなわち、画像処理パラメータには、テンプレートマッチング処理を実行する際の処理順序を示す画像処理シーケンスが含まれている。また、本実施形態において、テンプレートマッチング処理における閾値が可変であり、現在のテンプレートマッチングの閾値が画像処理パラメータに含まれている。さらに、検出部 4 2 は、テンプレートデータ 4 4 c と画像とを比較する前に各種の処理を実行可能である。図 3 においては、各種の処理として平滑化処理と鮮鋭化処理が例示

50

されており、それぞれの強度が画像処理パラメータに含まれている。

【0043】

撮像部 2 1 から画像が出力されると、検出部 4 2 は、画像処理シーケンスに基づいて、画像処理の順序（実行するか否かを含む）を決定し、当該順序で平滑化処理や鮮鋭化処理等の画像処理を実行する。このとき、検出部 4 2 は、画像処理パラメータに記述された強度で平滑化処理や鮮鋭化処理等の画像処理を実行する。また、画像処理シーケンスに含まれる比較（テンプレートデータ 4 4 c と画像との比較）を実行する際には、画像処理パラメータが示す閾値に基づいて比較を行う。

【0044】

なお、以上のように検出部 4 2 は、光学パラメータに基づいて撮像部 2 1 や照明部 2 2 の位置を特定し、ロボット 1、ロボット 2 を動作させることが可能であるが、ロボット 1 およびロボット 2 を駆動する際の位置は、後述する動作パラメータや力制御パラメータによって与えられてもよい。

【0045】

本実施形態において、制御部 4 3 は、位置制御部 4 3 a、力制御部 4 3 b、接触判定部 4 3 c、サーボ 4 3 d を備えている。また、制御部 4 3 においては、モーター M 1 ~ M 6 の回転角度の組み合わせと、ロボット座標系における T C P の位置との対応関係 U 1 が図示しない記憶媒体に記憶され、座標系の対応関係 U 2 が定義され、図示しない記憶媒体に記憶されている。従って、制御部 4 3 や後述する算出部 4 1 は、対応関係 U 2 に基づいて、任意の座標系におけるベクトルを他の座標系におけるベクトルに変換することができる。例えば、制御部 4 3、算出部 4 1 は、力覚センサー P の出力に基づいてセンサー座標系でのロボット 1 ~ 3 への作用力を取得し、ロボット座標系における T C P の位置に作用する力に変換することができる。また、制御部 4 3、算出部 4 1 は、力制御座標系で表現された目標力をロボット座標系における T C P の位置における目標力に変換することができる。むろん、対応関係 U 1、U 2 は記憶部 4 4 に記憶されていても良い。

【0046】

制御部 4 3 は、アームを駆動することによって、ロボット 1 ~ 3 とともに移動する各種の部位の位置や各種の部位に作用する力を制御することができ、位置の制御は主に位置制御部 4 3 a、力の制御は主に力制御部 4 3 b によって実行される。サーボ 4 3 d は、サーボ制御を実行することが可能であり、エンコーダー E 1 ~ E 6 の出力が示すモーター M 1 ~ M 6 の回転角度 D_a と、制御目標である目標角度 D_t とを一致させるフィードバック制御を実行する。すなわち、サーボ 4 3 d は、回転角度 D_a と目標角度 D_t との偏差、当該偏差の積分、当該偏差の微分にサーボゲイン K_{pp} 、 K_{pi} 、 K_{pd} を作用させた P I D 制御を実行することができる。

【0047】

さらに、サーボ 4 3 d は、当該サーボゲイン K_{pp} 、 K_{pi} 、 K_{pd} が作用した出力と、回転角度 D_a の微分との偏差、当該偏差の積分、当該偏差の微分にサーボゲイン K_{vp} 、 K_{vi} 、 K_{vd} を作用させた P I D 制御を実行することができる。当該サーボ 4 3 d による制御は、モーター M 1 ~ M 6 のそれぞれに対して実行可能である。従って、各サーボゲインはロボット 1 ~ 3 が備える 6 軸のそれぞれについて実行可能である。なお、本実施形態において、制御部 4 3 は、サーボ 4 3 d に制御信号を出力し、サーボゲイン K_{pp} 、 K_{pi} 、 K_{pd} 、 K_{vp} 、 K_{vi} 、 K_{vd} を変化させることができる。

【0048】

記憶部 4 4 には、上述のパラメータ 4 4 a に加え、ロボット 1 ~ 3 を制御するためのロボットプログラム 4 4 b が記憶される。本実施形態において、パラメータ 4 4 a およびロボットプログラム 4 4 b は、教示によって生成され、記憶部 4 4 に記憶されるが、後述する算出部 4 1 によって修正され得る。なお、ロボットプログラム 4 4 b は、主に、ロボット 1 ~ 3 が実施する作業のシーケンス（工程の順序）を示し、予め定義されたコマンドの組み合わせによって記述される。また、パラメータ 4 4 a は、主に、各工程を実現するために必要とされる具体的な値であり、各コマンドの引数として記述される。

【 0 0 4 9 】

ロボット 1 ~ 3 を制御するためのパラメータ 4 4 a には、上述の光学パラメータの他に、動作パラメータと力制御パラメータが含まれる。動作パラメータは、ロボット 1 ~ 3 の動作に関するパラメータであり、本実施形態においては、位置制御の際に参照されるパラメータである。すなわち、本実施形態において、一連の作業は複数の工程に分けられ、各工程を実施する際のパラメータ 4 4 a が教示によって生成される。動作パラメータには、当該複数の工程における始点と終点を示すパラメータが含まれている。当該始点と終点は、種々の座標系で定義されて良く、本実施形態においては制御対象のロボットの T C P の始点および終点がロボット座標系で定義される。すなわち、ロボット座標系の各軸についての並進位置と回転位置とが定義される。

10

【 0 0 5 0 】

また、動作パラメータには、複数の工程における T C P の加減速特性が含まれている。加減速特性は、ロボット 1 ~ 3 の T C P が各工程の始点から終点まで移動する際の期間と当該期間内の各時刻における T C P の速度を示している。図 4 は、当該加減速特性の例を示す図であり、始点における T C P の移動開始時刻 t_1 から T C P が終点に到達する時刻 t_4 までの期間内の各時刻において T C P の速度 V が定義されている。また、本実施形態において加減速特性には定速期間が含まれる。

【 0 0 5 1 】

定速期間は時刻 $t_2 \sim t_3$ の期間であり、この期間内に置いて速度は一定である。また、この期間の前後において T C P は加速し、また、減速する。すなわち、時刻 $t_1 \sim t_2$ までの期間において T C P は加速し、時刻 $t_3 \sim t_4$ までの期間において T C P は減速する。当該加減速特性も種々の座標系で定義されて良く、本実施形態においては制御対象のロボットの T C P についての速度であり、ロボット座標系で定義される。すなわち、ロボット座標系の各軸についての並進速度と回転速度（角速度）とが定義される。

20

【 0 0 5 2 】

さらに、動作パラメータには、サーボゲイン K_{pp} , K_{pi} , K_{pd} , K_{vp} , K_{vi} , K_{vd} が含まれている。すなわち、制御部 4 3 は、動作パラメータとして記述された値になるようにサーボ 4 3 d に制御信号を出力し、サーボゲイン K_{pp} , K_{pi} , K_{pd} , K_{vp} , K_{vi} , K_{vd} を調整することができる。本実施形態において当該サーボゲインは、上述の工程毎の値であるが、後述の学習等によってより短い期間毎の値とされても良い。

30

【 0 0 5 3 】

力制御パラメータは、ロボット 1 ~ 3 の力制御に関するパラメータであり、本実施形態においては、力制御の際に参照されるパラメータである。始点、終点、加減速特性、サーボゲインは、動作パラメータと同様のパラメータであり、始点、終点、加減速特性はロボット座標系の 3 軸の並進と回転について定義される。また、サーボゲインはモーター M 1 ~ M 6 のそれぞれについて定義される。ただし、力制御の場合、始点および終点の中の少なくとも一部は定義されない場合（任意とされる場合）もある。例えば、ある方向に作用する力が 0 になるように衝突回避や倣い制御が行われる場合、当該方向における始点および終点は定義されず、当該方向の力を 0 にするように位置が任意に変化し得る状態が定義される場合もある。

40

【 0 0 5 4 】

また、力制御パラメータには、力制御座標系を示す情報が含まれている。力制御座標系は、力制御の目標力を定義するための座標系であり、後述の学習が行われる前においては目標力ベクトルの起点が原点であり、目標力ベクトルの方向に 1 軸が向いている。すなわち、教示において力制御における各種の目標力が定義される際に、各作業の各工程における目標力の作用点が教示される。例えば、対象物の一点を他の物体に当て、両者の接触点で対象物から他の物体に一定の目標力を作用させた状態で対象物の向きを変化させる場合において、対象物が他の物体と接触する点が目標力の作用点となり、当該作用点を原点とした力制御座標系が定義される。そこで、力制御パラメータにおいては、力制御の目標力が作用する点を原点とし、目標力の方向に 1 軸が向いている座標系、すなわち、力制

50

御座標系を特定するための情報を、パラメーターに含んでいる。なお、当該パラメーターは種々の定義が可能であるが、例えば、力制御座標系と他の座標系（ロボット座標系等）との関係を示すデータによって定義可能である。

【 0 0 5 5 】

さらに、力制御パラメーターには、目標力が含まれている。目標力は、各種の作業において、任意の点に作用すべき力として教示される力であり、力制御座標系において定義される。すなわち、目標力を示す目標力ベクトルが、目標力ベクトルの起点と、起点からの 6 軸成分（3 軸の並進力、3 軸のトルク）として定義され、力制御座標系で表現されている。なお、力制御座標系と他の座標系との関係を利用すれば、当該目標力を任意の座標系、例えば、ロボット座標系におけるベクトルに変換することが可能である。

10

【 0 0 5 6 】

さらに、力制御パラメーターには、インピーダンスパラメーターが含まれている。すなわち、本実施形態において力制御部 4 3 b が実施する力制御は、インピーダンス制御である。インピーダンス制御は、仮想の機械的インピーダンスをモーター M 1 ~ M 6 によって実現する制御である。この際、TCP が仮想的に有する質量が仮想慣性係数 m として定義され、TCP が仮想的に受ける粘性抵抗が仮想粘性係数 d として定義され、TCP が仮想的に受ける弾性力のバネ定数が仮想弾性係数 k として定義される。インピーダンスパラメーターはこれらの m , d , k であり、ロボット座標系の各軸に対する並進と回転について定義される。本実施形態において当該力制御座標系、目標力、インピーダンスパラメーターは、上述の工程毎の値であるが、後述の学習等によってより短い期間毎の値とされても良い。

20

【 0 0 5 7 】

本実施形態において、一連の作業は複数の工程に分けられ、各工程を実施するロボットプログラム 4 4 b が教示によって生成されるが、位置制御部 4 3 a は、ロボットプログラム 4 4 b が示す各工程をさらに微小時間 T 毎の微小工程に細分化する。そして、位置制御部 4 3 a は、パラメーター 4 4 a に基づいて微小工程毎の目標位置 L_t を生成する。力制御部 4 3 b は、パラメーター 4 4 a に基づいて一連の作業の各工程における目標力 f_{Lt} を取得する。

【 0 0 5 8 】

すなわち、位置制御部 4 3 a は、動作パラメーターまたは力制御パラメーターが示す始点、終点、加減速特性を参照し、始点から終点まで当該加減速特性で移動する場合（姿勢の場合は姿勢が変化する場合）の微小工程毎の TCP の位置を目標位置 L_t として生成する。力制御部 4 3 b は、各工程についての力制御パラメーターが示す目標力を参照し、力制御座標系とロボット座標系との対応関係 U_2 に基づいて当該目標力をロボット座標系における目標力 f_{Lt} に変換する。当該目標力 f_{Lt} は、任意の点に作用する力として変換され得るが、ここでは、後述の作用力が TCP に作用している力として表現されるため、当該作用力と目標力 f_{Lt} とを運動方程式で解析するため、目標力 f_{Lt} が TCP の位置における力に変換されるとして説明を行う。むろん、工程によっては、目標力 f_{Lt} が定義されない場合もあり、この場合、力制御を伴わない位置制御が行われる。

30

【 0 0 5 9 】

なお、ここで L の文字は、ロボット座標系を規定する軸の方向（ x , y , z , R_x , R_y , R_z ）のなかのいずれか 1 個の方向を表すこととする。また、 L は、 L 方向の位置も表すこととする。例えば、 $L = x$ の場合、ロボット座標系にて設定された目標位置の x 方向成分が $L_t = x_t$ と表記され、目標力の x 方向成分が $f_{Lt} = f_{xt}$ と表記される。

40

【 0 0 6 0 】

位置制御や力制御を実行するため、制御部 4 3 は、ロボット 1 ~ 3 の状態を取得することができる。すなわち、制御部 4 3 は、モーター M 1 ~ M 6 の回転角度 Da を取得し、対応関係 U_1 に基づいて、当該回転角度 Da をロボット座標系における TCP の位置 L （ x , y , z , R_x , R_y , R_z ）に変換することができる。また制御部 4 3 は、対応関係 U_2 を参照し、TCP の位置 L と、力覚センサー P の検出値および位置とに基づいて、力覚

50

センサー P に現実には作用している力を T C P に作用している作用力 f_L に変換してロボット座標系において特定することができる。

【 0 0 6 1 】

すなわち、力覚センサー P に作用している力は、センサー座標系で定義される。そこで、制御部 4 3 は、ロボット座標系における T C P の位置 L と対応関係 U 2 と力覚センサー P の検出値に基づいて、ロボット座標系において T C P に作用する作用力 f_L を特定する。また、ロボットに作用するトルクは、作用力 f_L と、ツール接触点（エンドエフェクターとワークの接触点）から力覚センサー P までの距離とから算出することができ、図示されない f_L トルク成分として特定される。なお、制御部 4 3 は、作用力 f_L に対して重力補償を行う。重力補償とは、作用力 f_L から重力成分を除去する処理である。重力補償は、例えば、T C P の姿勢ごとに T C P に作用する作用力 f_L の重力成分を予め調査しておき、作用力 f_L から当該重力成分を減算するなどして実現可能である。

10

【 0 0 6 2 】

T C P に作用する重力以外の作用力 f_L と、T C P に作用すべき目標力 f_{Lt} とが特定されると、力制御部 4 3 b は、対象物等の物体が T C P に存在し、当該 T C P に力が作用し得る状態において、インピーダンス制御による補正量 L （以後、力由来補正量 L と呼ぶ。）を取得する。すなわち、力制御部 4 3 b はパラメーター 4 4 a を参照して目標力 f_{Lt} とインピーダンスパラメーター m, d, k を取得し、運動方程式（1）に代入して力由来補正量 L を取得する。なお、当該力由来補正量 L は、T C P が機械的インピーダンスを受けた場合に、目標力 f_{Lt} と作用力 f_L との力偏差 $f_L(t)$ を解消するために、T C P が移動すべき位置 L の大きさを意味する。

20

【 数 1 】

$$m\Delta\ddot{S}(t) + d\Delta\dot{S}(t) + k\Delta S(t) = \Delta f_S(t) \quad \cdots (1)$$

【 0 0 6 3 】

（1）式の左辺は、T C P の位置 L の 2 階微分値に仮想慣性係数 m を乗算した第 1 項と、T C P の位置 L の微分値に仮想粘性係数 d を乗算した第 2 項と、T C P の位置 L に仮想弾性係数 k を乗算した第 3 項とによって構成される。（1）式の右辺は、目標力 f_{Lt} から現実の作用力 f_L を減算した力偏差 $f_L(t)$ によって構成される。（1）式における微分とは、時間による微分を意味する。

30

【 0 0 6 4 】

力由来補正量 L が得られると、制御部 4 3 は、対応関係 U 1 に基づいて、ロボット座標系を規定する各軸の方向の動作位置を、各モーター M 1 ~ M 6 の目標の回転角度である目標角度 D_t に変換する。サーボ 4 3 d は、目標角度 D_t からモーター M 1 ~ M 6 の現実の回転角度であるエンコーダー E 1 ~ E 6 の出力（回転角度 D_a ）を減算することにより、駆動位置偏差 $D_e (= D_t - D_a)$ を算出する。サーボ 4 3 d は、パラメーター 4 4 a を参照してサーボゲイン $K_{pp}, K_{pi}, K_{pd}, K_{vp}, K_{vi}, K_{vd}$ を取得し、駆動位置偏差 D_e にサーボゲイン K_{pp}, K_{pi}, K_{pd} を乗算した値と、現実の回転角度 D_a の時間微分値である駆動速度との差である駆動速度偏差に、サーボゲイン K_{vp}, K_{vi}, K_{vd} を乗算した値とを加算することにより、制御量 D_c を導出する。制御量 D_c は、モーター M 1 ~ M 6 のそれぞれについて特定され、各モーター M 1 ~ M 6 の制御量 D_c でモーター M 1 ~ M 6 のそれぞれが制御される。制御部 4 3 がモーター M 1 ~ M 6 を制御する信号は、P W M（Pulse Width Modulation）変調された信号である。

40

【 0 0 6 5 】

以上のように、運動方程式に基づいて目標力 f_{Lt} から制御量 D_c を導出してモーター M 1 ~ M 6 を制御するモードを力制御モードというものとする。また制御部 4 3 は、エンドエフェクター等の構成要素が対象物 W から力を受けない非接触状態の工程では、力制御を行わず、目標位置から線形演算で導出する回転角度でモーター M 1 ~ M 6 を制御する。目標位置から線形演算で導出する回転角度でモーター M 1 ~ M 6 を制御するモードを位置制御モードというものとする。さらに、制御部 4 3 は、目標位置から線形演算で導出する回

50

転角度と目標力を運動方程式に代入して導出する回転角度とを例えば線型結合によって統合し、統合した回転角度でモーターM 1 ~ M 6を制御するハイブリッドモードでもロボット1を制御することができる。これらのモードはロボットプログラム4 4 bによって予め決められる。

【0066】

位置制御モードまたはハイブリッドモードで制御を行う場合、位置制御部4 3 aは、微小工程毎の目標位置 L_t を取得する。微小工程毎の目標位置 L_t が得られると、制御部4 3は、対応関係U 1に基づいて、ロボット座標系を規定する各軸の方向の動作位置を、各モーターM 1 ~ M 6の目標の回転角度である目標角度 D_t に変換する。サーボ4 3 dは、パラメーター4 4 aを参照してサーボゲイン K_{pp} , K_{pi} , K_{pd} , K_{vp} , K_{vi} , K_{vd} を取得し、目標角度 D_t に基づいて、制御量 D_c を導出する。制御量 D_c は、モーターM 1 ~ M 6のそれぞれについて特定され、各モーターM 1 ~ M 6の制御量 D_c でモーターM 1 ~ M 6のそれぞれが制御される。この結果、各工程において、TCPは、微小工程毎の目標位置 L_t を経由し、加減速特性に従って始点から終点まで移動する。

【0067】

なお、ハイブリッドモードでは、制御部4 3は、微小工程毎の目標位置 L_t に、力由来補正量 L を加算することにより動作位置($L_t + L$)を特定し、当該動作位置に基づいて目標角度 D_t を取得し、制御量 D_c を取得する。

【0068】

接触判定部4 3 cは、ロボット1 ~ 3が作業において想定されていない物体と接触したか否かを判定する機能を実行する。本実施形態において、接触判定部4 3 cは、ロボット1 ~ 3のそれぞれが備える力覚センサーPの出力を取得し、出力が予め決められた基準値を超えた場合にロボット1 ~ 3が作業において想定されていない物体と接触したと判定する。この場合において、種々の処理が行われて良いが、本実施形態において接触判定部4 3 cは、ロボット1 ~ 3の制御量 D_c を0としてロボット1 ~ 3を停止させる。なお、停止させる際の制御量は、種々の制御量であって良く、直前の制御量 D_c をキャンセルする制御量でロボット1 ~ 3を動作させる構成等であっても良い。

【0069】

(3) ピックアップ処理：

次に、以上の構成におけるロボット1 ~ 3の動作を説明する。ここでは、ロボット2の照明部2 2で照明され、ロボット1の撮像部2 1で撮像された対象物Wをロボット3のグリッパー2 3でピックアップする作業を例にして説明する。むろん、ロボット1 ~ 3による作業は、ピックアップ作業に限定されず、他にも種々の作業(例えば、ネジ締め作業、挿入作業、ドリルによる穴あけ作業、バリ取り作業、研磨作業、組み立て作業、製品チェック作業等)に適用可能である。ピックアップ処理は、上述のコマンドによって記述されたロボット制御プログラムによって検出部4 2および制御部4 3が実行する処理によって実現される。本実施形態においてピックアップ処理は、作業台に対象物Wが配置した状態で実行される。

【0070】

図5は、ピックアップ処理のフローチャートの例を示す図である。ピックアップ処理が開始されると、検出部4 2は、撮像部2 1が撮像した画像を取得する(ステップS 100)。すなわち、検出部4 2は、パラメーター4 4 aを参照して照明部2 2の位置を特定し、当該位置を位置制御部4 3 aに対して受け渡す。この結果、位置制御部4 3 aは、現在の照明部2 2の位置を始点、パラメーター4 4 aが示す照明部2 2の位置を終点とした位置制御を実行し、パラメーター4 4 aが示す照明部の位置に照明部2 2を移動させる。次に、検出部4 2はパラメーター4 4 aを参照して照明部2 2の明るさを特定し、照明部2 2を制御して照明の明るさを当該明るさに設定する。

【0071】

さらに、検出部4 2は、パラメーター4 4 aを参照して撮像部2 1の位置を特定し、当該位置を位置制御部4 3 aに対して受け渡す。この結果、位置制御部4 3 aは、現在の撮

像部 2 1 の位置を始点、パラメータ 4 4 a が示す撮像部 2 1 の位置を終点とした位置制御を実行し、パラメータ 4 4 a が示す照明部の位置に撮像部 2 1 を移動させる。次に、検出部 4 2 はパラメータ 4 4 a を参照して撮像部 2 1 の露光時間および絞りを特定し、撮像部 2 1 を制御して露光時間および絞りを当該露光時間および絞りに設定する。露光時間および絞りの設定が完了すると、撮像部 2 1 は、画像を撮像し、検出部 4 2 に対して出力する。検出部 4 2 は、当該画像を取得する。

【 0 0 7 2 】

次に、検出部 4 2 は、画像に基づいて、対象物の検出が成功したか否かを判定する（ステップ S 1 0 5）。すなわち、検出部 4 2 は、パラメータ 4 4 a を参照して画像処理シーケンスを特定し、当該画像処理シーケンスが示す各処理をパラメータ 4 4 a が示す強度で実行する。また、検出部 4 2 は、テンプレートデータ 4 4 c を参照し、テンプレートデータ 4 4 c と画像との差分を閾値と比較し、差分が閾値以下である場合に、対象物の検出が成功したと判定する。

【 0 0 7 3 】

ステップ S 1 0 5 において、対象物の検出が成功したと判定されない場合、検出部 4 2 は、テンプレートデータ 4 4 c と画像の相対位置、またはテンプレートデータ 4 4 c の大きさ、の少なくとも一方を変化させ、ステップ S 1 0 0 以降の処理を繰り返す。一方、ステップ S 1 0 5 において、対象物の検出が成功したと判定された場合、制御部 4 3 は、制御目標を特定する（ステップ S 1 1 0）。

【 0 0 7 4 】

本例におけるピックアップ処理は、検出部 4 2 が検出した対象物 W の位置姿勢に合わせてロボット 3 のグリッパー 2 3 を移動させ、姿勢を変化させ、グリッパー 2 3 で対象物 W をピックアップし、所定の位置まで対象物 W を運んでグリッパー 2 3 から対象物 W を離す作業である。そこで、位置制御部 4 3 a および力制御部 4 3 b は、ロボットプログラム 4 4 b に基づいて一連の作業を構成する複数の工程を特定する。

【 0 0 7 5 】

制御目標の特定対象となる工程は、各工程の中で未処理かつ時系列で先に存在する工程である。制御目標の特定対象となる工程が力制御モードの工程である場合、力制御部 4 3 b は、パラメータ 4 4 a の力制御パラメータを参照し、力制御座標系、目標力を取得する。力制御部 4 3 b は、力制御座標系に基づいて、当該目標力をロボット座標系の目標力 f_{Lt} に変換する。また、力制御部 4 3 b は、力覚センサー P の出力を TCP に作用している作用力 f_L に変換する。さらに、力制御部 4 3 b は、パラメータ 4 4 a の力制御パラメータを参照し、インピーダンスパラメータ m, d, k に基づいて、力由来補正量 L を制御目標として取得する。

【 0 0 7 6 】

制御目標の特定対象となる工程が位置制御モードである場合、位置制御部 4 3 a は、当該工程を微小工程に細分化する。そして、位置制御部 4 3 a は、パラメータ 4 4 a の動作パラメータを参照し、始点、終点、および加減速特性に基づいて、微小工程毎の目標位置 L_t を制御目標として取得する。制御目標の特定対象となる工程がハイブリッドモードである場合、位置制御部 4 3 a は、当該工程を微小工程に細分化し、パラメータ 4 4 a の力制御パラメータを参照し、始点、終点、および加減速特性に基づいて、微小工程毎の目標位置 L_t を取得し、力制御座標系、目標力 f_{Lt} 、インピーダンスパラメータ、作用力 f_L に基づいて力由来補正量 L を取得する。これらの目標位置 L_t および力由来補正量 L が制御目標である。

【 0 0 7 7 】

制御目標が特定されると、サーボ 4 3 d は、現在の制御目標でロボット 3 を制御する（ステップ S 1 1 5）。すなわち、現在の工程が力制御モードまたはハイブリッドモードの工程である場合、サーボ 4 3 d は、パラメータ 4 4 a の力制御パラメータを参照し、サーボゲインに基づいて、制御目標に対応する制御量 Dc を特定し、モーター M 1 ~ M 6 のそれぞれを制御する。現在の工程が位置制御モードの工程である場合、サーボ 4 3 d は

、パラメータ 4 4 a の動作パラメータを参照し、サーボゲインに基づいて、制御目標に対応する制御量 D_c を特定し、モーター M 1 ~ M 6 のそれぞれを制御する。

【 0 0 7 8 】

次に、制御部 4 3 は、現在の工程が終了したか否かを判定する（ステップ S 1 2 0）。当該判定は、種々の終了判定条件によって実行されてよく、位置制御であれば、例えば、TCP が目標位置に達したことや目標位置において TCP が整定したこと等が挙げられる。力制御であれば、例えば、作用力が目標力に一致した状態から作用力が指定の大きさ以上、または指定の大きさ以下に変化したことや、TCP が指定の範囲外になったこと等が挙げられる。前者は、例えば、ピックアップ作業における対象物の把持動作の完了や、把持解除動作の完了等が挙げられる。後者は、例えば、ドリルによる対象物の貫通作業においてドリルが貫通した場合等が挙げられる。

10

【 0 0 7 9 】

むろん、他にも各工程が失敗したと推定される場合において、工程が終了したと判定されて良い。ただし、この場合には、作業の中止や中断が行われることが好ましい。工程の失敗を判定するための終了判定条件としては、例えば、TCP の移動速度や加速度が上限値を超えた場合やタイムアウトが発生した場合等が挙げられる。終了判定条件を充足したか否かは、各種のセンサー、力覚センサー P や撮像部 2 1、他のセンサー等が利用されて良い。

【 0 0 8 0 】

ステップ S 1 2 0 において、現在の工程が終了したと判定されない場合、制御部 4 3 は、微小時間 T 後に、次の微小工程についてステップ S 1 1 5 以降の処理を実行する。すなわち、現在の工程が位置制御モードまたはハイブリッドモードである場合、位置制御部 4 3 a は、次の微小工程における目標位置 L_t を制御目標としてロボット 3 を制御する。また、現在の工程が力制御モードまたはハイブリッドモードである場合、力制御部 4 3 b は、再度力覚センサー P の出力に基づいて作用力 f_L を取得し、最新の作用力 f_L に基づいて特定される力由来補正量 L を制御目標としてロボット 3 を制御する。

20

【 0 0 8 1 】

ステップ S 1 2 0 において、現在の工程が終了したと判定された場合、制御部 4 3 は、作業が終了したか否かを判定する（ステップ S 1 2 5）。すなわち、ステップ S 1 2 0 で終了したと判定された工程が最終工程であった場合、制御部 4 3 は、作業が終了したと判定する。ステップ S 1 2 5 で作業が終了したと判定されなかった場合、制御部 4 3 は、作業シーケンスの次の工程を現在の工程に変更し（ステップ S 1 3 0）、ステップ S 1 1 0 以降の処理を実行する。ステップ S 1 2 5 で作業が終了したと判定された場合、制御部 4 3 は、作業が終了したと判定し、ピックアップ処理を終了する。

30

【 0 0 8 2 】

（ 4 ）学習処理：

本実施形態にかかる制御装置 4 0 は、以上のように、パラメータ 4 4 a に基づいてロボット 1 ~ 3 を制御することができる。上述の実施形態において、パラメータ 4 4 a は教示によって生成されたが、人為的な教示によってパラメータ 4 4 a を最適化することは困難である。

40

【 0 0 8 3 】

例えば、検出部 4 2 による対象物 W の検出において、同じ対象物 W であっても、光学パラメータが異なれば対象物 W の位置、画像内での対象物の像や位置、対象物 W に生じる影など、様々な要素が変化し得る。従って、光学パラメータを変化させると、検出部 4 2 による対象物 W の検出精度が変化し得る。そして、光学パラメータを変化させた場合に対象物 W の検出精度がどのように変化するのは、必ずしも明らかではない。

【 0 0 8 4 】

また、動作パラメータや力制御パラメータはロボット 1 ~ 3 の制御に利用されるが、ロボット 1 ~ 3 のように複数の自由度（可動軸）を有するロボットは極めて多数のパターンで動作することが可能である。そして、ロボット 1 ~ 3 においては、振動や異音、オ

50

ーバーシュート等の好ましくない動作が発生しないようにパターンが決められている必要がある。さらに、エンドエフェクターとして各種の装置が取り付けられる場合、ロボット 1 ~ 3 の重心が変化し得るため、最適な動作パラメーター、力制御パラメーターも変化し得る。そして、動作パラメーターや力制御パラメーターを変化させた場合に、ロボット 1 ~ 3 の動作がどのように変化するのは、必ずしも明らかではない。

【 0 0 8 5 】

さらに、力制御パラメーターは、ロボット 1 ~ 3 において力制御が行われる場合に利用されるが、ロボット 1 ~ 3 において実施される各作業において、力制御パラメーターを変化させた場合に、ロボット 1 ~ 3 の動作がどのように変化するのは、必ずしも明らかではない。例えば、どのような方向においてどのようなインピーダンスパラメーターが最適であるのか、全ての作業工程において推定することは困難である。このため、検出部 4 2 の検出精度を高めたり、ロボット 1 ~ 3 の潜在的な性能を引き出ししたりするためには極めて多数の試行錯誤を行う必要がある。

【 0 0 8 6 】

しかし、人為的に極めて多数の試行錯誤を行うことは困難であるため、対象物 W の検出精度が十分に高く、当該検出精度がほぼ上限に達していると推定される状態や、ロボット 1 ~ 3 の潜在的な性能が引き出されている状態（所要時間や消費電力等のパフォーマンスのさらなる向上が困難な状態）を人為的に実現することは困難である。また、パラメーター 4 4 a の調整を行うためには、パラメーター 4 4 a の変化による検出精度の変化やロボット 1 ~ 3 の動作の変化を熟知しているオペレーターが必要になり、熟知していないオペレーターがパラメーター 4 4 a の調整を行うことは困難である。また、常に熟練のオペレーターを必要とするシステムはとても不便である。

【 0 0 8 7 】

そこで、本実施形態においては、人為的なパラメーター 4 4 a の決定作業を行うことなく、自動的にパラメーター 4 4 a を決定するための構成を備えている。なお、本実施形態によれば、多少のパラメーター 4 4 a の変更によって検出精度がより向上しないと推定される（検出精度が極大であると推定される）状態や、多少のパラメーター 4 4 a の変更によってロボット 1 ~ 3 の性能が高性能化することはないと推定される（性能が極大であると推定される）状態を実現することができる。本実施形態においては、これらの状態を最適化された状態と呼ぶ。

【 0 0 8 8 】

本実施形態において制御装置 4 0 は、パラメーター 4 4 a の自動的な決定のために算出部 4 1 を備えている。本実施形態において、算出部 4 1 は、機械学習を用いて、光学パラメーターと動作パラメーターと力制御パラメーターとを算出することができる。図 6 は、算出部 4 1 の構成を示す図であり、図 2 に示す構成の一部を省略し、算出部 4 1 の詳細を示した図である。なお、図 6 に示す記憶部 4 4 は、図 2 に示す記憶部 4 4 と同一の記憶媒体であり、各図においては記憶された情報の一部の図示が省略されている。

【 0 0 8 9 】

算出部 4 1 は、状態変数を観測する状態観測部 4 1 a と、観測された状態変数に基づいてパラメーター 4 4 a を学習する学習部 4 1 b とを備えている。本実施形態において、状態観測部 4 1 a は、パラメーター 4 4 a を変化させたことによって生じた結果を状態変数として観測する。このため、状態観測部 4 1 a は、サーボ 4 3 d の制御結果と、エンコーダー E 1 ~ E 6 の値と、力覚センサー P の出力と、検出部 4 2 が取得する画像とを状態変数として取得可能である。

【 0 0 9 0 】

具体的には、状態観測部 4 1 a は、サーボ 4 3 d の制御結果として、モーター M 1 ~ M 6 に供給される電流値を観測する。当該電流値は、モーター M 1 ~ M 6 で出力されるトルクに相当する。エンコーダー E 1 ~ E 6 の出力は、対応関係 U 1 に基づいてロボット座標系における T C P の位置に変換される。従って、状態観測部 4 1 a は、ロボット 1 であれば撮像部 2 1 の位置、ロボット 2 であれば照明部 2 2 の位置、ロボット 3 であればグリッ

10

20

30

40

50

パー 2 3 の位置を観測することができる。

【 0 0 9 1 】

力覚センサー P の出力は、対応関係 U 2 に基づいてロボット座標系における T C P に作用する作用力に変換される。従って、状態観測部 4 1 a は、ロボット 1 ~ 3 への作用力を状態変数として観測することができる。検出部 4 2 が取得する画像は、撮像部 2 1 で撮像された画像であり、状態観測部 4 1 a は、当該画像を状態変数として観測することができる。状態観測部 4 1 a は、学習対象のパラメーター 4 4 a に応じて観測対象の状態変数を適宜選択することができる。

【 0 0 9 2 】

学習部 4 1 b は、学習によってパラメーター 4 4 a を最適化することができればよく、本実施形態においては、強化学習によってパラメーター 4 4 a を最適化する。具体的には、学習部 4 1 b は、状態変数に基づいてパラメーター 4 4 a を変化させる行動を決定し、当該行動を実行する。当該行動後の状態に応じて報酬を評価すれば、当該行動の行動価値が判明する。そこで、算出部 4 1 は、状態変数の観測と、当該状態変数に応じた行動の決定と、当該行動によって得られる報酬の評価とを繰り返すことによって、パラメーター 4 4 a を最適化する。

【 0 0 9 3 】

本実施形態において、算出部 4 1 は、パラメーター 4 4 a の中から学習対象のパラメーターを選択して学習を行うことができる。本実施形態においては、光学パラメーターの学習と、動作パラメーターの学習と、力制御パラメーターの学習とのそれぞれを独立して実行することができる。

【 0 0 9 4 】

(4 - 1) 光学パラメーターの学習 :

図 7 はエージェントと環境とからなる強化学習のモデルに沿って光学パラメーターの学習例を説明する図である。図 7 に示すエージェントは、予め決められた方策に応じて行動 a を選択する機能に相当し、学習部 4 1 b によって実現される。環境は、エージェントが選択した行動 a と現在の状態 s とに基づいて次の状態 s' を決定し、行動 a と状態 s と状態 s' とに基づいて即時報酬 r を決定する機能に相当し、状態観測部 4 1 a および学習部 4 1 b によって実現される。

【 0 0 9 5 】

本実施形態においては、予め決められた方策によって学習部 4 1 b が行動 a を選択し、状態観測部 4 1 a が状態の更新を行う処理を繰り返すことにより、ある状態 s におけるある行動 a の行動価値関数 $Q(s, a)$ を算出する Q 学習が採用される。すなわち、本例においては、下記の式 (2) によって行動価値関数を更新する。そして、行動価値関数 $Q(s, a)$ が適正に収束した場合には、当該行動価値関数 $Q(s, a)$ を最大化する行動 a が最適な行動であると見なされ、当該行動 a を示すパラメーター 4 4 a が最適化されたパラメーターであると見なされる。

【 数 2 】

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t)) \cdots (2)$$

【 0 0 9 6 】

ここで、行動価値関数 $Q(s, a)$ は、状態 s において行動 a を取った場合において将来にわたって得られる収益 (本例では割引報酬総和) の期待値である。報酬は r であり、状態 s、行動 a、報酬 r の添え字 t は、時系列で繰り返す試行過程における 1 回分のステップを示す番号 (試行番号と呼ぶ) であり、行動決定後に状態が変化すると試行番号がインクリメントされる。従って、式 (2) 内の報酬 r_{t+1} は状態 s_t で行動 a_t が選択され、状態が s_{t+1} になった場合に得られる報酬である。 α は学習率、 γ は割引率である。また、 a' は、状態 s_{t+1} で取り得る行動 a_{t+1} の中で行動価値関数 $Q(s_{t+1}, a_{t+1})$ を最大化する行動であり、 $\max_{a'} Q(s_{t+1}, a')$ は、行動 a' が選択されたことによって最大化された行動価値関数である。

【0097】

光学パラメータの学習においては、光学パラメータを変化させることが行動の決定に相当しており、学習対象のパラメータと取り得る行動とを示す行動情報44dが記憶部44に予め記録される。すなわち、当該行動情報44dに学習対象として記述された光学パラメータが学習対象となる。図7においては、光学パラメータの中の撮像部パラメータと、照明部パラメータと、画像処理パラメータとの一部が学習対象となっている例を示している。

【0098】

具体的には、撮像部パラメータの中で撮像部21のx座標、y座標が学習対象となっている。従って、この例においてz座標やxyz軸に対する回転(姿勢)は学習対象となっておらず、撮像部21は、対象物Wが置かれる作業台に向いている状態であるとともに、撮像部21のx-y平面内での移動が学習対象である。むしろ、他の撮像部パラメータ、例えば、撮像部21の姿勢やz座標、露光時間や絞りが学習対象であっても良い。

【0099】

また、図7に示す例においては、照明部パラメータの中で、照明部22のx座標、y座標および照明部の明るさが学習対象となっている。従って、この例においてz座標やxyz軸に対する回転(姿勢)は学習対象となっておらず、照明部22は、対象物Wが置かれる作業台に向いている状態であるとともに、照明部22のx-y平面内での移動が学習対象である。むしろ、他の照明部パラメータ、例えば、照明部22の姿勢やz座標が学習対象であっても良い。

【0100】

さらに、図7に示す例においては、画像処理パラメータの中で、平滑化処理の強度と鮮鋭化処理の強度とテンプレートマッチングの閾値が学習対象となっている。従って、この例において、画像処理シーケンスは学習対象となっておらず、撮像部21で撮像された画像に対する画像処理の順序は変化しない(むしろ、画像処理シーケンスが学習対象である実施形態も採用可能である)。

【0101】

図7に示す例において行動には値を一定値増加させる行動と、値を一定値減少させる行動とが存在する。従って、図7に示す全8個のパラメータにおいて取り得る行動は全16個である(行動a1~行動a16)。行動情報44dは、学習対象のパラメータと取り得る行動とを示しているため、図7に示す例であれば、図示した8個のパラメータが行動情報44dに学習対象として記述される。また、各行動を特定するための情報(行動のID、各行動での増減量等)が行動情報44dに記述される。

【0102】

図7に示す例において、報酬は対象物Wの検出の成否に基づいて特定される。すなわち、学習部41bは、行動aとして光学パラメータを変化させた後、当該光学パラメータによってロボット1,2を動作させ、検出部42によって撮像部21が撮像した画像を取得する。そして、学習部41bは、当該光学パラメータに基づいてテンプレートマッチング処理を実行し、対象物Wの検出が成功したか否かを判定する。さらに、学習部41bは、検出の成否によって行動a、状態s、s'の報酬を決定する。当該報酬は、対象物Wの検出の成否に基づいて決定されれば良く、例えば、検出の成功に正(例えば+1)、検出の失敗に負(例えば-1)の報酬を与える構成等を採用可能である。この構成によれば、対象物の検出精度を高めるように最適化を行うことができる。

【0103】

現在の状態sにおいて行動aが採用された場合における次の状態s'は、行動aとしてのパラメータの変化が行われた後にロボット1,2を動作させ、状態観測部41aが状態を観測することによって特定可能である。なお、本例にかかる光学パラメータの学習においてロボット3は動作しない。図7に示す例において、状態変数には、撮像部21のx座標、y座標と、照明部22のx座標、y座標、照明部22の明るさと、平滑化処理の強度、鮮鋭化処理の強度、テンプレートマッチングの閾値と、撮像部21で撮像された画

10

20

30

40

50

像とが含まれている。

【0104】

従って、この例において、状態観測部41aは、行動aが実行された後に、ロボット1のエンコーダーE1～E6の出力をU1に基づいて変換して撮像部21のx座標およびy座標を観測する。また、状態観測部41aは、行動aが実行された後に、ロボット2のエンコーダーE1～E6の出力をU1に基づいて変換して照明部22のx座標およびy座標を観測する。

【0105】

本実施形態において、照明部22の明るさは、パラメータ44aによって誤差無く調整可能であると見なされており（または誤差が影響ないと見なされており）、状態観測部41aは、パラメータ44aに含まれる照明部の明るさを取得して状態変数が観測されたと見なす。むしろ、照明部22の明るさは、センサー等によって実測されても良いし、撮像部21が撮像した画像に基づいて（例えば、平均階調値等により）観測されても良い。状態観測部41aは、平滑化処理の強度、鮮鋭化処理の強度、テンプレートマッチングの閾値についても、パラメータ44aを参照して現在の値を取得し、状態変数が観測されたと見なす。

【0106】

さらに、状態観測部41aにおいては、撮像部21が撮像し、検出部42が取得した画像を状態変数として取得する（図7に示す太枠）。すなわち、状態観測部41aは、撮像部21が撮像した画像（対象物が存在し得る注目領域等の画像であっても良い）の画素毎の階調値を状態変数として観測する。撮像部のx座標等は、行動であるとともに観測対象としての状態であるが、撮像部21が撮像した画像は行動ではない。従って、この意味で、撮像された画像は、光学パラメータの変化から直接的に推定することが困難な変化をし得る状態変数である。また、検出部42は、当該画像に基づいて対象物を検出するため、当該画像は検出の成否に直接的に影響を与え得る状態変数である。従って、状態変数として、当該画像を観測することにより、人為的に改善することが困難なパラメータの改善を行い、効果的に検出部42の検出精度を高めるように光学パラメータを最適化することが可能になる。

【0107】

（4-2）光学パラメータの学習例：

次に、光学パラメータの学習例を説明する。学習の過程で参照される変数や関数を示す情報は、学習情報44eとして記憶部44に記憶される。すなわち、算出部41は、状態変数の観測と、当該状態変数に応じた行動の決定と、当該行動によって得られる報酬の評価とを繰り返すことによって行動価値関数 $Q(s, a)$ を収束させる構成が採用されている。そこで、本例において、学習の過程で状態変数と行動と報酬との時系列の値が、順次、学習情報44eに記録されていく。

【0108】

行動価値関数 $Q(s, a)$ は、種々の手法で算出されて良く、多数回の試行に基づいて算出されても良いが、本実施形態においては、行動価値関数 $Q(s, a)$ を近似的に算出する一手法であるDQN(Deep Q-Ne t w o r k)が採用されている。DQNにおいては、多層ニューラルネットワークを用いて行動価値関数 $Q(s, a)$ を推定する。本例においては、状態sを入力とし、選択し得る行動の数N個の行動価値関数 $Q(s, a)$ の値を出力とする多層ニューラルネットワークが採用されている。

【0109】

図8は、本例において採用されている多層ニューラルネットワークを模式的に示す図である。図8において、多層ニューラルネットワークは、M個（Mは2以上の整数）の状態変数を入力とし、N個（Nは2以上の整数）個の行動価値関数Qの値を出力としている。例えば、図7に示す例であれば、撮像部のx座標～テンプレートマッチングの閾値までの8個の状態変数と撮像された画像の画素数との和がM個であり、M個の状態変数の値が多層ニューラルネットワークに入力される。図8においては、試行番号tにおけるM個の状

10

20

30

40

50

態を $s_{1t} \sim s_{Mt}$ として示している。

【0110】

N個は選択し得る行動aの数であり、多層ニューラルネットワークの出力は、入力された状態sにおいて特定の行動aが選択された場合の行動価値関数Qの値である。図8においては、試行番号tにおいて選択し得る行動 $a_{1t} \sim a_{Nt}$ のそれぞれにおける行動価値関数Qを $Q(s_t, a_{1t}) \sim Q(s_t, a_{Nt})$ として示している。当該Qに含まれる s_t は入力された状態 $s_{1t} \sim s_{Mt}$ を代表して示す文字である。図7に示す例であれば、16個の行動が選択可能であるため $N = 16$ である。むろん、行動aの内容や数(Nの値)、状態sの内容や数(Mの値)は試行番号tに応じて変化しても良い。

【0111】

図8に示す多層ニューラルネットワークは、各層の各ノードにおいて直前の層の入力(1層目においては状態s)に対する重みwの乗算とバイアスbの加算とを実行し、必要に応じて活性化関数を経た出力を得る(次の層の入力になる)演算を実行するモデルである。本例においては、層DLがP個(Pは1以上の整数)存在し、各層において複数のノードが存在する。

【0112】

図8に示す多層ニューラルネットワークは各層における重み、とバイアスb、活性化関数、層の順序等によって特定される。そこで、本実施形態においては、当該多層ニューラルネットワークを特定するためのパラメータ(入力から出力を得るために必要な情報)が学習情報44eとして記憶部44に記録される。なお、学習の際には、多層ニューラルネットワークを特定するためのパラメータの中で可変の値(例えば、重みwとバイアスb)を更新していくことになる。ここでは、学習の過程で変化し得る多層ニューラルネットワークのパラメータを と表記する。当該 を使用すると、上述の行動価値関数 $Q(s_t, a_{1t}) \sim Q(s_t, a_{Nt})$ は、 $Q(s_t, a_{1t}; \theta_t) \sim Q(s_t, a_{Nt}; \theta_t)$ とも表記できる。

【0113】

次に、図9に示すフローチャートに沿って学習処理の手順を説明する。光学パラメータの学習処理は、ロボット1,2の運用過程において実施されても良いし、実運用の前に事前に学習処理が実行されてもよい。ここでは、実運用の前に事前に学習処理が実行される構成(多層ニューラルネットワークを示す が最適化されると、その情報が保存され、次回以降の運用で利用される構成)に従って学習処理を説明する。

【0114】

学習処理が開始されると、算出部41は、学習情報44eを初期化する(ステップS200)。すなわち、算出部41は、学習を開始する際に参照される の初期値を特定する。初期値は、種々の手法によって決められて良く、過去に学習が行われていない場合においては、任意の値やランダム値等が の初期値となっても良いし、ロボット1,2や撮像部21、照明部22の光学特性を模擬するシミュレーション環境を準備し、当該環境に基づいて学習または推定した を初期値としてもよい。

【0115】

過去に学習が行われた場合は、当該学習済の が初期値として採用される。また、過去に類似の対象についての学習が行われた場合は、当該学習における が初期値とされても良い。過去の学習は、ロボット1,2を用いてユーザーが行ってもよいし、ロボット1,2の製造者がロボット1,2の販売前に行ってもよい。この場合、製造者は、対象物や作業の種類に応じて複数の初期値のセットを用意しておき、ユーザーが学習する際に初期値を選択する構成であっても良い。 の初期値が決定されると、当該初期値が現在の の値として学習情報44eに記憶される。

【0116】

次に、算出部41は、パラメータを初期化する(ステップS205)。ここでは、光学パラメータが学習対象であるため、算出部41は、光学パラメータを初期化する。すなわち、算出部41は、ロボット1のエンコーダーE1~E6の出力を対応関係U1で

10

20

30

40

50

変換し、撮像部 2 1 の位置を初期値として設定する。また、算出部 4 1 は、予め決められた初期の露光時間（過去に学習が行われた場合には最新の露光時間）を撮像部 2 1 の露光時間の初期値として設定する。さらに、算出部 4 1 は、撮像部 2 1 に制御信号を出力し、現在の絞りの値を初期値として設定する。

【 0 1 1 7 】

さらに、算出部 4 1 は、ロボット 2 のエンコーダー E 1 ~ E 6 の出力を対応関係 U 1 で変換し、照明部 2 2 の位置を初期値として設定する。また、算出部 4 1 は、予め決められた初期の明るさ（過去に学習が行われた場合には最新の明るさ）を照明部 2 2 の明るさの初期値として設定する。さらに、算出部 4 1 は、平滑化処理の強度、鮮鋭化処理の強度、テンプレートマッチングの閾値、画像処理シーケンスについて予め決められた初期値（過去に学習が行われた場合には最新の値）を設定する。初期化されたパラメーターは記憶部 4 4 に現在のパラメーター 4 4 a として記憶される。

【 0 1 1 8 】

次に、状態観測部 4 1 a は、状態変数を観測する（ステップ S 2 1 0）。すなわち、制御部 4 3 は、パラメーター 4 4 a およびロボットプログラム 4 4 b を参照してロボット 1, 2 を制御する。検出部 4 2 は、制御後の状態で撮像部 2 1 が撮像した画像に基づいて対象物 W の検出処理（上述のステップ S 1 0 0, S 1 0 5 に相当）を実行する。この後、状態観測部 4 1 a は、ロボット 1 のエンコーダー E 1 ~ E 6 の出力を U 1 に基づいて変換して撮像部 2 1 の x 座標および y 座標を観測する。また、状態観測部 4 1 a は、ロボット 2 のエンコーダー E 1 ~ E 6 の出力を U 1 に基づいて変換して照明部 2 2 の x 座標および y 座標を観測する。さらに、状態観測部 4 1 a は、パラメーター 4 4 a を参照して照明部 2 2 に設定されるべき明るさを取得して状態変数が観測されたと見なす。

【 0 1 1 9 】

さらに、状態観測部 4 1 a は、平滑化処理の強度、鮮鋭化処理の強度、テンプレートマッチングの閾値についても、パラメーター 4 4 a を参照して現在の値を取得し、状態変数が観測されたと見なす。さらに、状態観測部 4 1 a においては、撮像部 2 1 が撮像し、検出部 4 2 が取得した画像を取得し、各画素の階調値を状態変数として取得する。

【 0 1 2 0 】

次に、学習部 4 1 b は、行動価値を算出する（ステップ S 2 1 5）。すなわち、学習部 4 1 b は、学習情報 4 4 e を参照して $Q(s_t, a_{1t}; \gamma_t)$ を取得し、学習情報 4 4 e が示す多層ニューラルネットワークに最新の状態変数を入力し、N 個の行動価値関数 $Q(s_t, a_{1t}; \gamma_t) \sim Q(s_t, a_{Nt}; \gamma_t)$ を算出する。

【 0 1 2 1 】

なお、最新の状態変数は、初回の実行時においてステップ S 2 1 0、2 回目以降の実行時においてステップ S 2 2 5 の観測結果である。また、試行番号 t は初回の実行時において 0、2 回目以降の実行時において 1 以上の値となる。学習処理が過去に実施されていない場合、学習情報 4 4 e が示す γ は最適化されていないため、行動価値関数 Q の値としては不正確な値となり得るが、ステップ S 2 1 5 以後の処理の繰り返しにより、行動価値関数 Q は徐々に最適化していく。また、ステップ S 2 1 5 以後の処理の繰り返しにおいて、状態 s、行動 a、報酬 r は、各試行番号 t に対応づけられて記憶部 4 4 に記憶され、任意のタイミングで参照可能である。

【 0 1 2 2 】

次に、学習部 4 1 b は、行動を選択し、実行する（ステップ S 2 2 0）。本実施形態においては、行動価値関数 $Q(s, a)$ を最大化する行動 a が最適な行動であると見なされる処理が行われる。そこで、学習部 4 1 b は、ステップ S 2 1 5 において算出された N 個の行動価値関数 $Q(s_t, a_{1t}; \gamma_t) \sim Q(s_t, a_{Nt}; \gamma_t)$ の値の中で最大の値を特定する。そして、学習部 4 1 b は、最大の値を与えた行動を選択する。例えば、N 個の行動価値関数 $Q(s_t, a_{1t}; \gamma_t) \sim Q(s_t, a_{Nt}; \gamma_t)$ の中で $Q(s_t, a_{Nt}; \gamma_t)$ が最大値であれば、学習部 4 1 b は、行動 a_{Nt} を選択する。

【 0 1 2 3 】

10

20

30

40

50

行動が選択されると、学習部 4 1 b は、当該行動に対応するパラメータ 4 4 a を変化させる。例えば、図 7 に示す例において、撮像部の x 座標を一定値増加させる行動 a 1 が選択された場合、学習部 4 1 b は、光学パラメータの撮像部パラメータが示す撮像部の位置において x 座標を一定値増加させる。パラメータ 4 4 a の変化が行われると、制御部 4 3 は、当該パラメータ 4 4 a を参照してロボット 1, 2 を制御する。検出部 4 2 は、制御後の状態で撮像部 2 1 が撮像した画像に基づいて対象物 W の検出処理を実行する。

【 0 1 2 4 】

次に、状態観測部 4 1 a は、状態変数を観測する（ステップ S 2 2 5）。すなわち、状態観測部 4 1 a は、ステップ S 2 1 0 における状態変数の観測と同様の処理を行って、状態変数として、撮像部 2 1 の x 座標および y 座標、照明部 2 2 の x 座標および y 座標、照明部 2 2 に設定されるべき明るさ、平滑化処理の強度、鮮鋭化処理の強度、テンプレートマッチングの閾値、撮像部 2 1 が撮像した画像の各画素の階調値を取得する。なお、現在の試行番号が t である場合（選択された行動が a_t である場合）、ステップ S 2 2 5 で取得される状態 s は s_{t+1} である。

【 0 1 2 5 】

次に、学習部 4 1 b は、報酬を評価する（ステップ S 2 3 0）。本例において、報酬は、対象物 W の検出の成否に基づいて決定される。そこで、学習部 4 1 b は、検出部 4 2 から対象物の検出結果の成否（ステップ S 1 0 5 の成否）を取得し、検出成功であれば既定量の正の報酬、検出失敗であれば既定量の負の報酬を取得する。なお、現在の試行番号が t である場合、ステップ S 2 3 0 で取得される報酬 r は r_{t+1} である。

【 0 1 2 6 】

本実施形態においては式（2）に示す行動価値関数 Q の更新を目指しているが、行動価値関数 Q を適切に更新していくためには、行動価値関数 Q を示す多層ニューラルネットワークを最適化（ θ を最適化）していかななくてはならない。図 8 に示す多層ニューラルネットワークによって行動価値関数 Q を適正に出力させるためには、当該出力のターゲットとなる教師データが必要になる。すなわち、多層ニューラルネットワークの出力と、ターゲットとの誤差を最小化するように θ を改善することによって、多層ニューラルネットワークが最適化されることが期待される。

【 0 1 2 7 】

しかし、本実施形態において、学習が完了していない段階では行動価値関数 Q の知見がなく、ターゲットを特定することは困難である。そこで、本実施形態においては、式（2）の第 2 項、いわゆる TD 誤差（Temporal Difference）を最小化する目的関数によって多層ニューラルネットワークを示す θ の改善を実施する。すなわち、 $(r_{t+1} + \max_{a'} Q(s_{t+1}, a'; \theta_t))$ をターゲットとし、ターゲットと $Q(s_t, a_t; \theta_t)$ との誤差が最小化するように θ を学習する。ただし、ターゲット $(r_{t+1} + \max_{a'} Q(s_{t+1}, a'; \theta_t))$ は、学習対象の θ を含んでいるため、本実施形態においては、ある程度の試行回数にわたりターゲットを固定する（例えば、最後に学習した θ （初回学習時は θ の初期値）で固定する）。本実施形態においては、ターゲットを固定する試行回数である既定回数が予め決められている。

【 0 1 2 8 】

このような前提で学習を行うため、ステップ S 2 3 0 で報酬が評価されると、学習部 4 1 b は目的関数を算出する（ステップ S 2 3 5）。すなわち、学習部 4 1 b は、試行のそれぞれにおける TD 誤差を評価するための目的関数（例えば、TD 誤差の 2 乗の期待値に比例する関数や TD 誤差の 2 乗の総和等）を算出する。なお、TD 誤差は、ターゲットが固定された状態で算出されるため、固定されたターゲットを $(r_{t+1} + \max_{a'} Q(s_{t+1}, a'; \theta_t))$ と表記すると、TD 誤差は $(r_{t+1} + \max_{a'} Q(s_{t+1}, a'; \theta_t)) - Q(s_t, a_t; \theta_t)$ である。当該 TD 誤差の式において報酬 r_{t+1} は、行動 a_t によってステップ S 2 3 0 で得られた報酬である。

【 0 1 2 9 】

また、 $\max_a Q(s_{t+1}, a; \cdot)$ は、行動 a_t によってステップ S 2 2 5 で算出される状態 s_{t+1} を、固定された \cdot で特定される多層ニューラルネットワークの入力とした場合に得られる出力の中の最大値である。 $Q(s_t, a_t; \cdot)$ は、行動 a_t が選択される前の状態 s_t を、試行番号 t の段階の \cdot で特定される多層ニューラルネットワークの入力とした場合に得られる出力の中で、行動 a_t に対応した出力の値である。

【0130】

目的関数が算出されると、学習部 4 1 b は、学習が終了したか否かを判定する（ステップ S 2 4 0）。本実施形態においては、TD 誤差が十分に小さいか否かを判定するための閾値が予め決められており、目的関数が閾値以下である場合、学習部 4 1 b は、学習が終了したと判定する。

【0131】

ステップ S 2 4 0 において学習が終了したと判定されない場合、学習部 4 1 b は、行動価値を更新する（ステップ S 2 4 5）。すなわち、学習部 4 1 b は、TD 誤差の δ による偏微分に基づいて目的関数を小さくするための Q の変化を特定し、 Q を変化させる。むしろ、ここでは、各種の手法で Q を変化させることが可能であり、例えば、RMSProp 等の勾配降下法を採用可能である。また、学習率等による調整も適宜実施されて良い。以上の処理によれば、行動価値関数 Q がターゲットに近づくように Q を変化させることができる。

【0132】

ただし、本実施形態においては、上述のようにターゲットが固定されているため、学習部 4 1 b は、さらに、ターゲットを更新するか否かの判定を行う。具体的には学習部 4 1 b は、既定回数の試行が行われたか否かを判定し（ステップ S 2 5 0）、ステップ S 2 5 0 において、既定回数の試行が行われたと判定された場合に、学習部 4 1 b は、ターゲットを更新する（ステップ S 2 5 5）。すなわち、学習部 4 1 b は、ターゲットを算出する際に参照される Q を最新の Q に更新する。この後、学習部 4 1 b は、ステップ S 2 1 5 以降の処理を繰り返す。一方、ステップ S 2 5 0 において、既定回数の試行が行われたと判定されなければ、学習部 4 1 b は、ステップ S 2 5 5 をスキップしてステップ S 2 1 5 以降の処理を繰り返す。

【0133】

ステップ S 2 4 0 において学習が終了したと判定された場合、学習部 4 1 b は、学習情報 4 4 e を更新する（ステップ S 2 6 0）。すなわち、学習部 4 1 b は、学習によって得られた Q を、ロボット 1, 2 による作業や検出部 4 2 による検出の際に参照されるべきとして学習情報 4 4 e に記録する。当該 Q を含む学習情報 4 4 e が記録されている場合、ステップ 1 0 0 ~ S 1 0 5 のようにロボット 1, 2 による作業が行われる際に、検出部 4 2 はパラメータ 4 4 a に基づいて対象物の検出処理を行う。そして、検出部 4 2 による検出が成功するまで、撮像部 2 1 による撮像が繰り返される工程においては、状態観測部 4 1 a による現在の状態の観測と、学習部 4 1 b による行動の選択が繰り返される。むしろ、この際、学習部 4 1 b は、状態を入力として算出された出力 $Q(s, a)$ の中で最大値を与える行動 a を選択する。そして、行動 a が選択された場合、行動 a が行われた状態に相当する値となるようにパラメータ 4 4 a が更新される。

【0134】

以上の構成によれば、検出部 4 2 は、行動価値関数 Q が最大化される行動 a を選択しながら対象物の検出処理を実行することができる。当該行動価値関数 Q は、上述の処理により、多数の試行が繰り返された結果、最適化されている。そして、当該試行は、算出部 4 1 によって自動で行われ、人為的に実施不可能な程度の多数の試行を容易に実行することができる。従って、本実施形態によれば、人為的に決められた光学パラメータよりも高い確率で対象物を高精度に検出することができる。

【0135】

さらに、本実施形態において検出部 4 2 は、対象物の位置姿勢を検出する構成であるため、本実施形態によれば高精度に対象物の位置姿勢を検出することができる。さらに、本

10

20

30

40

50

実施形態によれば、最適化された行動価値関数 Q に基づいて、光学パラメータである撮像部パラメータを算出することができる。従って、対象物の検出精度を高めるように撮像部 2 1 を調整することができる。さらに、本実施形態によれば、最適化された行動価値関数 Q に基づいて、光学パラメータである照明部パラメータを算出することができる。従って、対象物の検出精度を高めるように照明部 2 2 を調整することができる。

【0136】

さらに、本実施形態によれば、最適化された行動価値関数 Q に基づいて、光学パラメータである画像処理パラメータを算出することができる。従って、対象物の検出精度を高める画像処理を実行することが可能になる。さらに、本実施形態によれば、自動で行動価値関数 Q が最適化されるため、高精度に対象物を検出する光学パラメータを容易に算出することができる。また、行動価値関数 Q の最適化は自動的に行われるため、最適な光学パラメータの算出も自動的に行うことができる。

10

【0137】

さらに、本実施形態において学習部 4 1 b は、状態変数としての画像に基づいて光学パラメータを変化させる行動を決定し、光学パラメータを最適化する。従って、照明部 2 2 によって照明が行われている実環境下において撮像部 2 1 で実際に撮像した画像に基づいて光学パラメータを最適化することができる。従って、ロボット 1, 2 の使用環境に応じた光学パラメータとなるように最適化することができる。

【0138】

本実施形態においては、撮像部 2 1 の位置および照明部 2 2 の位置が行動に含まれており、当該行動に基づいて行動価値関数 Q を最適化することで撮像部 2 1 の位置および照明部 2 2 の位置に関するパラメータ 4 4 a を最適化することができる。従って、学習後においては、少なくとも、撮像部 2 1 と照明部 2 2 の相対位置関係が理想化される。また、対象物 W が作業台の固定位置またはほぼ固定された位置に置かれるのならば、学習後において、撮像部 2 1 と照明部 2 2 のロボット座標系における位置が理想化されと考えることもできる。さらに、本実施形態においては、撮像部 2 1 によって撮像された画像が状態として観測される。従って、本実施形態によれば、各種の画像の状態に対応した撮像部 2 1 の位置や照明部 2 2 の位置が理想化される。

20

【0139】

(4-3) 動作パラメータの学習：

30

動作パラメータの学習においても、学習対象のパラメータを選択することが可能であり、ここでは、その一例を説明する。図 10 は、動作パラメータの学習例を図 7 と同様のモデルで説明した図である。本例も式 (2) に基づいて行動価値関数 $Q(s, a)$ を最適化する。従って、最適化後の行動価値関数 $Q(s, a)$ を最大化する行動 a が最適な行動であると見なされ、当該行動 a を示すパラメータ 4 4 a が最適化されたパラメータであると見なされる。

【0140】

動作パラメータの学習においても、動作パラメータを変化させることが行動の決定に相当しており、学習対象のパラメータと取り得る行動とを示す行動情報 4 4 d が記憶部 4 4 に予め記録される。すなわち、当該行動情報 4 4 d に学習対象として記述された動作パラメータが学習対象となる。図 10 においては、ロボット 3 における動作パラメータの中のサーボゲインと加減速特性が学習対象であり、動作の始点および終点は学習対象となっていない。なお、動作の始点および終点は教示位置であるが、本実施形態においては他の位置は教示されない。従って、本実施形態においては、ロボット 3 に対して教示された教示位置を含まない構成である。

40

【0141】

具体的には、動作パラメータの中のサーボゲイン K_{pp} , K_{pi} , K_{pd} , K_{vp} , K_{vi} , K_{vd} は、モーター $M_1 \sim M_6$ のそれぞれについて定義され、6 軸のそれぞれについて増減可能である。従って、本実施形態においては、1 軸あたり 6 個のサーボゲインのそれぞれを増加または減少させることが可能であり、増加について 3 6 個の行動、減少についても 3

50

6 個の行動、計 7 2 個の行動（行動 a 1 ~ a 7 2）を選択し得る。

【0 1 4 2】

一方、動作パラメーターの中の加減速特性は図 4 に示すような特性であり、モーター M 1 ~ M 6 のそれぞれについて（6 軸について）定義される。本実施形態において加減速特性は、加速域における加速度、減速域における加速度、車速が 0 より大きい期間の長さ（図 4 に示す t_4 ）を変化させることができる。なお、本実施形態において加速域や減速域におけるカーブは加速度の増減によって定義され、例えば、増減後の加速度がカーブ中央の傾きを示し、当該中央の周囲のカーブは予め決められた規則に従って変化する。むしろ、加減速特性の調整法は他にも種々の手法が採用可能である。

【0 1 4 3】

いずれにしても、本実施形態においては、1 軸あたり 3 個の要素（加速域、減速域、期間）で加減速特性を調整可能であり、各要素に応じた数値（加速度、期間長）を増加または減少させることが可能である。従って、増加について 1 8 個の行動、減少についても 1 8 個の行動、計 3 6 個の行動（行動 a 7 3 ~ a 1 0 8）を選択し得る。本実施形態においては、以上のようにして予め定義された行動の選択肢に対応するパラメーターが、行動情報 4 4 d に学習対象として記述される。また、各行動を特定するための情報（行動の ID、各行動での増減量等）が行動情報 4 4 d に記述される。

【0 1 4 4】

図 1 0 に示す例において、報酬はロボット 3 が行った作業の良否に基づいて評価される。すなわち、学習部 4 1 b は、行動 a として動作パラメーターを変化させた後、当該動作パラメーターによってロボット 3 を動作させ、検出部 4 2 によって検出された対象物をピックアップする作業を実行する。さらに、学習部 4 1 b は、作業の良否を観測し、作業の良否を評価する。そして、学習部 4 1 b は、作業の良否によって行動 a、状態 s、s' の報酬を決定する。

【0 1 4 5】

なお、作業の良否は作業の成否（ピックアップ成否等）のみならず、作業の質を含む。具体的には、学習部 4 1 b は、図示しない計時回路に基づいて作業の開始から終了まで（ステップ S 1 1 0 の開始からステップ S 1 2 5 で終了と判定されるまで）の所要時間を取得する。そして、学習部 4 1 b は、作業の所要時間が基準よりも短い場合に正（例えば + 1）、作業の所要時間が基準よりも長い場合に負（例えば - 1）の報酬を与える。なお、基準は種々の要素によって特定されて良く、例えば、前回の作業の所要時間であっても良いし、過去の最短所要時間であっても良いし、予め決められた時間であっても良い。

【0 1 4 6】

さらに、学習部 4 1 b は、作業の各工程において、ロボット 3 のエンコーダー E 1 ~ E 6 の出力を U 1 に基づいて変換してグリッパー 2 3 の位置を取得する。そして、学習部 4 1 b は、各工程の目標位置（終点）と、工程終了の際のグリッパー 2 3 の位置とのずれ量を取得し、グリッパー 2 3 の位置と目標位置とのずれ量が基準よりも小さい場合に正、基準よりも大きい場合の負の報酬を与える。なお、基準は種々の要素によって特定されて良く、例えば、前回のずれ量であっても良いし、過去の最短のずれ量であっても良いし、予め決められたずれ量であっても良い。

【0 1 4 7】

さらに、学習部 4 1 b は、作業の各工程において取得したグリッパー 2 3 の位置を、整定する以前の所定期間にわたって取得し、当該期間における振動強度を取得する。そして、学習部 4 1 b は、当該振動強度の程度が基準よりも小さい場合に正、基準よりも大きい場合の負の報酬を与える。なお、基準は種々の要素によって特定されて良く、例えば、前回の振動強度の程度であっても良いし、過去の最小の振動強度の程度であっても良いし、予め決められた振動強度の程度であっても良い。振動強度の程度は、種々の手法で特定されて良く、目標位置からの乖離の積分値や閾値以上の振動が生じている期間など、種々の手法を採用可能である。なお、所定期間は、種々の期間とすることが可能であり、工程の始点から終点にわたる期間であれば、動作中の振動強度による報酬が評価され、工程の終

10

20

30

40

50

期が所定期間とされれば、残留振動の強度による報酬が評価される。

【0148】

さらに、学習部41bは、作業の各工程の終期において取得したグリッパー23の位置を、整定する以前の所定期間にわたって取得し、当該期間における目標位置からの乖離の最大値をオーバーシュート量として取得する。そして、学習部41bは、当該オーバーシュート量が基準よりも小さい場合に正、基準よりも大きい場合の負の報酬を与える。なお、基準は種々の要素によって特定されて良く、例えば、前回のオーバーシュート量の程度であっても良いし、過去の最小のオーバーシュート量であっても良いし、予め決められたオーバーシュート量であっても良い。

【0149】

さらに、本実施形態においては、制御装置40、ロボット1～3、作業台等の少なくとも1カ所に集音装置が取り付けられており、学習部41bは、作業中に集音装置が取得した音を示す情報を取得する。そして、学習部41bは、作業中の発生音の大きさが基準よりも小さい場合に正、基準よりも大きい場合の負の報酬を与える。なお、基準は種々の要素によって特定されて良く、例えば、前回の作業または工程の発生音の大きさの程度であっても良いし、過去の発生音の大きさの最小値であっても良いし、予め決められた大きさであっても良い。また、発生音の大きさは、音圧の最大値で評価されても良いし、所定期間内の音圧の統計値（平均値等）で評価されても良く、種々の構成を採用可能である。

【0150】

現在の状態sにおいて行動aが採用された場合における次の状態s'は、行動aとしてのパラメータの変化が行われた後にロボット3を動作させ、状態観測部41aが状態を観測することによって特定可能である。なお、本例にかかる動作パラメータの学習は、ロボット1, 2による対象物の検出完了後に、ロボット3に関して実行される。

【0151】

図10に示す例において、状態変数には、モーターM1～M6の電流、エンコーダーE1～E6の値、力覚センサーPの出力が含まれている。従って、状態観測部41aは、サーボ43dの制御結果として、モーターM1～M6に供給される電流値を観測することができる。当該電流値は、モーターM1～M6で出力されるトルクに相当する。また、エンコーダーE1～E6の出力は、対応関係U1に基づいてロボット座標系におけるTCPの位置に変換される。従って、状態観測部41aは、ロボット3が備えるグリッパー23の位置情報を観測することになる。

【0152】

力覚センサーPの出力は積分することによってロボットの位置に変換することができる。すなわち、状態観測部41aは、対応関係U2に基づいてロボット座標系においてTCPへの作用力を積分することでTCPの位置を取得する。従って、本実施形態において状態観測部41aは、力覚センサーPの出力も利用してロボット3が備えるグリッパー23の位置情報を観測する。なお、状態は、各種の手法で観測されて良く、上述の変換が行われない値（電流値やエンコーダー、力覚センサーの出力値）が状態として観測されても良い。

【0153】

状態観測部41aは、行動であるサーボゲインや加減速特性の調整結果を直接的に観測しているのではなく、調整の結果、ロボット3で得られた変化をモーターM1～M6の電流、エンコーダーE1～E6の値、力覚センサーPの出力として観測している。従って、行動による影響を間接的に観測することになり、この意味で、本実施形態の状態変数は、動作パラメータの変化から直接的に推定することが困難な変化をし得る状態変数である。

【0154】

また、モーターM1～M6の電流、エンコーダーE1～E6の値、力覚センサーPの出力は、ロボット3の動作を直接的に示しており、当該動作は作業の良否を直接的に示している。従って、状態変数として、モーターM1～M6の電流、エンコーダーE1～E6の

10

20

30

40

50

値、力覚センサー P の出力を観測することにより、人為的に改善することが困難なパラメータの改善を行い、効果的に作業の質を高めるように動作パラメータを最適化することが可能になる。この結果、人為的に決められた動作パラメータよりも高性能な動作を行う動作パラメータを高い確率で算出することができる。

【0155】

(4-4) 動作パラメータの学習例：

次に、動作パラメータの学習例を説明する。学習の過程で参照される変数や関数を示す情報は、学習情報 44e として記憶部 44 に記憶される。すなわち、算出部 41 は、状態変数の観測と、当該状態変数に応じた行動の決定と、当該行動によって得られる報酬の評価とを繰り返すことによって行動価値関数 $Q(s, a)$ を収束させる構成が採用されている。そこで、本例において、学習の過程で状態変数と行動と報酬との時系列の値が、順次、学習情報 44e に記録されていく。

10

【0156】

なお、本実施形態において、動作パラメータの学習は位置制御モードで実行される。位置制御モードでの学習を実行するためには、位置制御モードのみで構成される作業がロボット 3 のロボットプログラム 44b として生成されても良いし、任意のモードを含む作業がロボット 3 のロボットプログラム 44b として生成されている状況において、その中の位置制御モードのみを用いて学習してもよい。

【0157】

行動価値関数 $Q(s, a)$ は、種々の手法で算出されて良く、多数回の試行に基づいて算出されても良いが、ここでは、DQN によって行動価値関数 Q を最適化する例を説明する。行動価値関数 Q の最適化に利用される多層ニューラルネットワークは、上述の図 8 において模式的に示される。図 10 に示すような状態が観測される本例であれば、ロボット 3 におけるモーター M1 ~ M6 の電流、エンコーダー E1 ~ E6 の値、力覚センサー P の出力 (6 軸の出力) が状態であるため、状態 s の数 $M = 18$ である。また、図 10 に示す行動が選択され得る本例であれば、108 個の行動が選択可能であるため $N = 108$ である。むろん、行動 a の内容や数 (N の値)、状態 s の内容や数 (M の値) は試行番号 t に応じて変化しても良い。

20

【0158】

本実施形態においても、当該多層ニューラルネットワークを特定するためのパラメータ (入力から出力を得るために必要な情報) が学習情報 44e として記憶部 44 に記録される。ここでも学習の過程で変化し得る多層ニューラルネットワークのパラメータを表記する。当該を使用すると、上述の行動価値関数 $Q(s_t, a_{1t}) \sim Q(s_t, a_{Nt})$ は、 $Q(s_t, a_{1t}; \theta_t) \sim Q(s_t, a_{Nt}; \theta_t)$ とも表記できる。

30

【0159】

次に、図 9 に示すフローチャートに沿って学習処理の手順を説明する。動作パラメータの学習処理は、ロボット 3 の運用過程において実施されても良いし、実運用の前に事前に学習処理が実行されてもよい。ここでは、実運用の前に事前に学習処理が実行される構成 (多層ニューラルネットワークを示すが最適化されると、その情報が保存され、次回以降の運用で利用される構成) に従って学習処理を説明する。

40

【0160】

学習処理が開始されると、算出部 41 は、学習情報 44e を初期化する (ステップ S200)。すなわち、算出部 41 は、学習を開始する際に参照される の初期値を特定する。初期値は、種々の手法によって決められて良く、過去に学習が行われていない場合においては、任意の値やランダム値等が の初期値となっても良いし、ロボット 3 や対象物を模擬するシミュレーション環境を準備し、当該環境に基づいて学習または推定した を初期値としてもよい。

【0161】

過去に学習が行われた場合は、当該学習済の が初期値として採用される。また、過去に類似の対象についての学習が行われた場合は、当該学習における が初期値とされても

50

良い。過去の学習は、ロボット 3 を用いてユーザーが行ってもよいし、ロボット 3 の製造者がロボット 3 の販売前に行ってもよい。この場合、製造者は、対象物や作業の種類に応じて複数の初期値のセットを用意しておき、ユーザーが学習する際に初期値を選択する構成であっても良い。の初期値が決定されると、当該初期値が現在の の値として学習情報 4 4 e に記憶される。

【0162】

次に、算出部 4 1 は、パラメータを初期化する（ステップ S 2 0 5）。ここでは、動作パラメータが学習対象であるため、算出部 4 1 は、動作パラメータを初期化する。すなわち、学習が行われていない状態であれば、算出部 4 1 は、教示によって生成されたパラメータ 4 4 a に含まれる動作パラメータを初期値として設定する。過去に何らかの学習が行われた状態であれば、算出部 4 1 は、学習の際に最後に利用されていたパラメータ 4 4 a に含まれる動作パラメータを初期値として設定する。

10

【0163】

次に、状態観測部 4 1 a は、状態変数を観測する（ステップ S 2 1 0）。すなわち、制御部 4 3 は、パラメータ 4 4 a およびロボットプログラム 4 4 b を参照してロボット 3 を制御する（上述のステップ S 1 1 0 ~ S 1 3 0 に相当）。この後、状態観測部 4 1 a は、モーター M 1 ~ M 6 に供給される電流値を観測する。また、状態観測部 4 1 a は、エンコーダ E 1 ~ E 6 の出力を取得し、対応関係 U 1 に基づいてロボット座標系における T C P の位置に変換する。さらに、状態観測部 4 1 a は、力覚センサー P の出力を積分し、T C P の位置を取得する。

20

【0164】

次に、学習部 4 1 b は、行動価値を算出する（ステップ S 2 1 5）。すなわち、学習部 4 1 b は、学習情報 4 4 e を参照して を取得し、学習情報 4 4 e が示す多層ニューラルネットワークに最新の状態変数を入力し、N 個の行動価値関数 $Q(s_t, a_{1t}; \tau_t) \sim Q(s_t, a_{Nt}; \tau_t)$ を算出する。

【0165】

なお、最新の状態変数は、初回の実行時においてステップ S 2 1 0、2 回目以降の実行時においてステップ S 2 2 5 の観測結果である。また、試行番号 t は初回の実行時において 0、2 回目以降の実行時において 1 以上の値となる。学習処理が過去に実施されていない場合、学習情報 4 4 e が示す は最適化されていないため、行動価値関数 Q の値としては不正確な値となり得るが、ステップ S 2 1 5 以後の処理の繰り返しにより、行動価値関数 Q は徐々に最適化していく。また、ステップ S 2 1 5 以後の処理の繰り返しにおいて、状態 s、行動 a、報酬 r は、各試行番号 t に対応づけられて記憶部 4 4 に記憶され、任意のタイミングで参照可能である。

30

【0166】

次に、学習部 4 1 b は、行動を選択し、実行する（ステップ S 2 2 0）。本実施形態においては、行動価値関数 $Q(s, a)$ を最大化する行動 a が最適な行動であると見なされる処理が行われる。そこで、学習部 4 1 b は、ステップ S 2 1 5 において算出された N 個の行動価値関数 $Q(s_t, a_{1t}; \tau_t) \sim Q(s_t, a_{Nt}; \tau_t)$ の値の中で最大の値を特定する。そして、学習部 4 1 b は、最大の値を与えた行動を選択する。例えば、N 個の行動価値関数 $Q(s_t, a_{1t}; \tau_t) \sim Q(s_t, a_{Nt}; \tau_t)$ の中で $Q(s_t, a_{Nt}; \tau_t)$ が最大値であれば、学習部 4 1 b は、行動 a_{Nt} を選択する。

40

【0167】

行動が選択されると、学習部 4 1 b は、当該行動に対応するパラメータ 4 4 a を変化させる。例えば、図 1 0 に示す例において、モーター M 1 のサーボゲイン Kpp を一定値増加させる行動 a 1 が選択された場合、学習部 4 1 b は、動作パラメータが示すモーター M 1 のサーボゲイン Kpp の値を一定値増加させる。パラメータ 4 4 a の変化が行われると、制御部 4 3 は、当該パラメータ 4 4 a を参照してロボット 3 を制御し、一連の作業を実行させる。なお、本実施形態においては、行動選択のたびに一連の作業が実行されるが、行動選択のたびに一連の作業の一部が実行される構成（一連の作業を構成する複数の

50

工程の少なくとも 1 工程が実行される構成)であっても良い。

【0168】

次に、状態観測部 41a は、状態変数を観測する(ステップ S225)。すなわち、状態観測部 41a は、ステップ S210 における状態変数の観測と同様の処理を行って、状態変数として、モーター M1 ~ M6 に供給される電流値、エンコーダー E1 ~ E6 の出力に基づいて特定される TCP の位置、力覚センサー P の出力に基づいて特定される TCP の位置を取得する。なお、現在の試行番号が t である場合(選択された行動が a_t である場合)、ステップ S225 で取得される状態 s は s_{t+1} である。

【0169】

次に、学習部 41b は、報酬を評価する(ステップ S230)。すなわち、学習部 41b は、図示しない計時回路に基づいて作業の開始から終了までの所要時間を取得し、作業の所要時間が基準よりも短い場合に正、作業の所要時間が基準よりも長い場合に負の報酬を取得する。さらに、学習部 41b は、作業の各工程の終了段階におけるグリッパー 23 の位置を取得し、各工程の目標位置とのずれ量を取得する。そして、学習部 41b は、グリッパー 23 の位置と目標位置とのずれ量が基準よりも小さい場合に正、基準よりも大きい場合の負の報酬を取得する。一連の作業が複数の工程で構成される場合、各工程の報酬の和が取得されても良いし、統計値(平均値等)が取得されても良い。

10

【0170】

さらに、学習部 41b は、作業の各工程において取得したグリッパー 23 の位置に基づいて振動強度を取得する。そして、学習部 41b は、当該振動強度の程度が基準よりも小さい場合に正、基準よりも大きい場合の負の報酬を取得する。一連の作業が複数の工程で構成される場合、各工程の報酬の和が取得されても良いし、統計値(平均値等)が取得されても良い。

20

【0171】

さらに、学習部 41b は、作業の各工程の終期において取得したグリッパー 23 の位置に基づいてオーバーシュート量を取得する。そして、学習部 41b は、当該オーバーシュート量が基準よりも小さい場合に正、基準よりも大きい場合の負の報酬を取得する。一連の作業が複数の工程で構成される場合、各工程の報酬の和が取得されても良いし、統計値(平均値等)が取得されても良い。

【0172】

さらに、学習部 41b は、作業中に集音装置が取得した音を示す情報を取得する。そして、学習部 41b は、作業中の発生音の大きさが基準よりも小さい場合に正、基準よりも大きい場合の負の報酬を取得する。なお、現在の試行番号が t である場合、ステップ S230 で取得される報酬 r は r_{t+1} である。

30

【0173】

本例においても、式(2)に示す行動価値関数 Q の更新を目指しているが、行動価値関数 Q を適切に更新していくためには、行動価値関数 Q を示す多層ニューラルネットワークを最適化(を最適化)していかななくてはならない。そして、図 8 に示す多層ニューラルネットワークによって行動価値関数 Q を適正に出力させるためには、当該出力のターゲットとなる教師データが必要になる。すなわち、多層ニューラルネットワークの出力と、ターゲットとの誤差を最小化するようにを改善すると、多層ニューラルネットワークが最適化されることが期待される。

40

【0174】

しかし、本実施形態において、学習が完了していない段階では行動価値関数 Q の知見がなく、ターゲットを特定することは困難である。そこで、本実施形態においては、式(2)の第 2 項、いわゆる TD 誤差を最小化する目的関数によって多層ニューラルネットワークを示すの改善を実施する。すなわち、 $(r_{t+1} + \max_{a'} Q(s_{t+1}, a'; \theta_t))$ をターゲットとし、ターゲットと $Q(s_t, a_t; \theta_t)$ との誤差が最小化するようにを学習する。ただし、ターゲット $(r_{t+1} + \max_{a'} Q(s_{t+1}, a'; \theta_t))$ は、学習対象のを含んでいるため、本実施形態においては、ある程度の試行回数にわたりター

50

ゲットを固定する（例えば、最後に学習した（初回学習時はの初期値）で固定する）。本実施形態においては、ターゲットを固定する試行回数である既定回数が予め決められている。

【0175】

このような前提で学習を行うため、ステップS230で報酬が評価されると、学習部41bは目的関数を算出する（ステップS235）。すなわち、学習部41bは、試行のそれぞれにおけるTD誤差を評価するための目的関数（例えば、TD誤差の2乗の期待値に比例する関数やTD誤差の2乗の総和等）を算出する。なお、TD誤差は、ターゲットが固定された状態で算出されるため、固定されたターゲットを $(r_{t+1} + \max_a Q(s_{t+1}, a; \theta))$ と表記すると、TD誤差は $(r_{t+1} + \max_a Q(s_{t+1}, a; \theta) - Q(s_t, a_t; \theta))$ である。当該TD誤差の式において報酬 r_{t+1} は、行動 a_t によってステップS230で得られた報酬である。

10

【0176】

また、 $\max_a Q(s_{t+1}, a; \theta)$ は、行動 a_t によってステップS225で算出される状態 s_{t+1} を、固定された θ で特定される多層ニューラルネットワークの入力とした場合に得られる出力の中の最大値である。 $Q(s_t, a_t; \theta)$ は、行動 a_t が選択される前の状態 s_t を、試行番号 t の段階の θ で特定される多層ニューラルネットワークの入力とした場合に得られる出力の中で、行動 a_t に対応した出力の値である。

【0177】

目的関数が算出されると、学習部41bは、学習が終了したか否かを判定する（ステップS240）。本実施形態においては、TD誤差が十分に小さいか否かを判定するための閾値が予め決められており、目的関数が閾値以下である場合、学習部41bは、学習が終了したと判定する。

20

【0178】

ステップS240において学習が終了したと判定されない場合、学習部41bは、行動価値を更新する（ステップS245）。すなわち、学習部41bは、TD誤差の γ による偏微分に基づいて目的関数を小さくするための θ の変化を特定し、 θ を変化させる。むろん、ここでは、各種の手法で θ を変化させることが可能であり、例えば、RMSProp等の勾配降下法を採用可能である。また、学習率等による調整も適宜実施されて良い。以上の処理によれば、行動価値関数 Q がターゲットに近づくように θ を変化させることができる。

30

【0179】

ただし、本実施形態においては、上述のようにターゲットが固定されているため、学習部41bは、さらに、ターゲットを更新するか否かの判定を行う。具体的には学習部41bは、既定回数の試行が行われたか否かを判定し（ステップS250）、ステップS250において、既定回数の試行が行われたと判定された場合に、学習部41bは、ターゲットを更新する（ステップS255）。すなわち、学習部41bは、ターゲットを算出する際に参照される θ を最新の θ に更新する。この後、学習部41bは、ステップS215以降の処理を繰り返す。一方、ステップS250において、既定回数の試行が行われたと判定されなければ、学習部41bは、ステップS255をスキップしてステップS215以降の処理を繰り返す。

40

【0180】

ステップS240において学習が終了したと判定された場合、学習部41bは、学習情報44eを更新する（ステップS260）。すなわち、学習部41bは、学習によって得られた θ を、ロボット3による作業の際に参照されるべき θ として学習情報44eに記録する。当該 θ を含む学習情報44eが記録されている場合、ステップ110～S130のようにロボット3による作業が行われる際に、制御部43はパラメータ44aに基づいてロボット3を制御する。そして、当該作業の過程においては、状態観測部41aによる現在の状態の観測と、学習部41bによる行動の選択が繰り返される。むろん、この際、学習部41bは、状態を入力として算出された出力 $Q(s, a)$ の中で最大値を与える行

50

動 a を選択する。そして、行動 a が選択された場合、行動 a が行われた状態に相当する値となるようにパラメータ 4 4 a が更新される。

【 0 1 8 1 】

以上の構成によれば、制御部 4 3 は、行動価値関数 Q が最大化される行動 a を選択しながら作業を実行することができる。当該行動価値関数 Q は、上述の処理により、多数の試行が繰り返された結果、最適化されている。そして、当該試行は、算出部 4 1 によって自動で行われ、人為的に実施不可能な程度の多数の試行を容易に実行することができる。従って、本実施形態によれば、人為的に決められた動作パラメータよりも高い確率でロボット 3 の作業の質を高めることができる。

【 0 1 8 2 】

さらに、本実施形態においては、行動によってパラメータ 4 4 a としてのサーボゲインが変化する。従って、人為的な調整によって適切な設定を行うことが困難な、モーターを制御するためのサーボゲインを自動的に調整することができる。さらに、本実施形態においては、行動によってパラメータ 4 4 a としての加減速特性が変化する。従って、人為的な調整によって適切な設定を行うことが困難な加減速特性を自動的に調整することができる。

【 0 1 8 3 】

さらに、本実施形態においては、行動によってロボットの動作の始点および終点が変わらない。従って、本実施形態においては、ロボット 3 が予定された始点および終点から外れ、利用者の意図しない動作が行われることを防止することができる。さらに、本実施形態においては、行動によってロボットに対する教示位置である始点および終点は変化しない。従って、本実施形態においては、ロボット 3 が教示された位置から外れ、利用者の意図しない動作が行われることを防止することができる。なお、本実施形態において、教示位置は始点および終点であるが、他の位置が教示位置となってもよい。例えば、始点と終点との間で通過すべき位置や取るべき姿勢がある場合、これらが教示位置（教示姿勢）であっても良い。

【 0 1 8 4 】

さらに、本実施形態においては、ロボット 3 が行った作業の良否に基づいて行動による報酬を評価するため、ロボット 3 の作業を成功させるようにパラメータを最適化することができる。さらに、本実施形態においては、作業の所要時間が基準よりも短い場合に報酬を正と評価するため、ロボット 3 を短い時間で作業させる動作パラメータを容易に算出することができる。さらに、本実施形態においては、ロボット 3 の位置と目標位置とのずれ量が基準よりも小さい場合に報酬を正と評価するため、ロボット 3 を目標位置に正確に移動させる動作パラメータを容易に算出することができる。

【 0 1 8 5 】

さらに、本実施形態においては、振動強度が基準よりも小さい場合に報酬を正と評価するため、ロボット 3 の動作による振動を発生させる可能性が低い動作パラメータを容易に算出することができる。さらに、本実施形態においては、ロボット 3 の位置のオーバーシュートが基準よりも小さい場合に報酬を正と評価するため、ロボット 3 がオーバーシュートする可能性が低い動作パラメータを容易に算出することができる。さらに、本実施形態においては、発生音が基準よりも小さい場合に報酬を正と評価するため、ロボット 3 に異常を発生させる可能性が低い動作パラメータを容易に算出することができる。

【 0 1 8 6 】

さらに、本実施形態によれば、自動で行動価値関数 Q が最適化されるため、高性能な動作を行う動作パラメータを容易に算出することができる。また、行動価値関数 Q の最適化は自動的に行われるため、最適な動作パラメータの算出も自動的に行うことができる。

【 0 1 8 7 】

さらに、本実施形態においては、ロボット 3 において汎用的に使用される力覚センサー P によってロボット 3 の位置情報を取得するため、ロボット 3 で汎用的に使用されるセン

10

20

30

40

50

サーに基づいて位置情報を算出することができる。

【0188】

さらに、本実施形態において学習部41bは、状態変数としてのロボット3の動作結果を実測し、動作パラメータを最適化する。従って、ロボット3によって作業が行われている実環境下において合わせて動作パラメータを最適化することができる。従って、ロボット3の使用環境に応じた動作パラメータとなるように最適化することができる。

【0189】

さらに、本実施形態において状態観測部41aは、ロボット3にエンドエフェクターとしてのグリッパ23が設けられた状態で状態変数を観測する。また、学習部41bは、ロボット3にエンドエフェクターとしてのグリッパ23が設けられた状態で行動としてのパラメータ44aの変更が実行される。この構成によれば、エンドエフェクターとしてのグリッパ23を用いた動作を行うロボット3に適した動作パラメータを容易に算出することができる。

【0190】

さらに、本実施形態において状態観測部41aは、エンドエフェクターとしてのグリッパ23が対象物を把持した状態で状態変数を観測する。また、学習部41bは、エンドエフェクターとしてのグリッパ23が対象物を把持した状態で行動としてのパラメータ44aの変更が実行される。この構成によれば、エンドエフェクターとしてのグリッパ23で対象物を把持して動作を行うロボット3に適した動作パラメータを容易に算出することができる。

【0191】

(4-5) 力制御パラメータの学習：

力制御パラメータの学習においても、学習対象のパラメータを選択することが可能であり、ここでは、その一例を説明する。図11は、力制御パラメータの学習例を図7と同様のモデルで説明した図である。本例も式(2)に基づいて行動価値関数 $Q(s, a)$ を最適化する。従って、最適化後の行動価値関数 $Q(s, a)$ を最大化する行動 a が最適な行動であると見なされ、当該行動 a を示すパラメータ44aが最適化されたパラメータであると見なされる。

【0192】

力制御パラメータの学習においても、力制御パラメータを変化させることが行動の決定に相当しており、学習対象のパラメータと取り得る行動とを示す行動情報44dが記憶部44に予め記録される。すなわち、当該行動情報44dに学習対象として記述された力制御パラメータが学習対象となる。図11においては、ロボット3における力制御パラメータの中のインピーダンスパラメータと、力制御座標系と、目標力と、ロボット3の動作の始点および終点が学習対象である。なお、力制御における動作の始点および終点は教示位置であるが、力制御パラメータの学習によって変動し得る。また、力制御座標系の原点は、ロボット3のTCP(ツールセンターポイント)からのオフセット点であり、学習前においては目標力が作用する作用点である。従って、力制御座標系(原点座標と軸回転角)と目標力が変化すると、TCPからのオフセット点の位置が変化することになり、目標力の作用点が力制御座標系の原点ではない場合も生じ得る。

【0193】

力制御パラメータの中のインピーダンスパラメータ m, k, d は、ロボット座標系の各軸に対する並進と回転について定義される。従って、本実施形態においては、1軸あたり3個のインピーダンスパラメータ m, d, k のそれぞれを増加または減少させることが可能であり、増加について18個の行動、減少についても18個の行動、計36個の行動(行動 $a_1 \sim a_{36}$)を選択し得る。

【0194】

一方、力制御座標系は当該座標系の原点座標と、力制御座標系の軸の回転角度と、がロボット座標系を基準として表現されることによって定義される。従って、本実施形態においては、原点座標の3軸方向への増減と、3軸の軸回転角の増減とが可能であり、原点座

10

20

30

40

50

標の増加について3個、減少について3個、軸回転角の増加について3個、減少について3個の行動が可能であり、計12個の行動(行動a37~a48)を選択し得る。目標力は、目標力ベクトルで表現され、目標力の作用点と、力制御座標系の6軸それぞれの成分の大きさ(3軸の並進力、3軸のトルク)によって定義される。従って、本実施形態においては、目標力の作用点の3軸方向への増減について6個、6軸それぞれの成分の増加について6個、減少について6個の行動が可能であり、計18個の行動(行動a49~a66)を選択し得る。

【0195】

ロボット3の動作の始点および終点は、ロボット座標系の各軸方向に沿って座標の増減が可能であり、始点の増減について6個、終点の増減について6個の計12個の行動(行動a67~a78)を選択し得る。本実施形態においては、以上のようにして予め定義された行動の選択肢に対応するパラメーターが、行動情報44dに学習対象として記述される。また、各行動を特定するための情報(行動のID、各行動での増減量等)が行動情報44dに記述される。

10

【0196】

図11に示す例において、報酬はロボット3が行った作業の良否に基づいて評価される。すなわち、学習部41bは、行動aとして力制御パラメーターを変化させた後、当該力制御パラメーターによってロボット3を動作させ、検出部42によって検出された対象物をピックアップする作業を実行する。さらに、学習部41bは、作業の良否を観測し、作業の良否を評価する。そして、学習部41bは、作業の良否によって行動a、状態s、s'の報酬を決定する。

20

【0197】

なお、作業の良否は作業の成否(ピックアップ成否等)のみならず、作業の質を含む。具体的には、学習部41bは、図示しない計時回路に基づいて作業の開始から終了まで(ステップS110の開始からステップS125で終了と判定されるまで)の所要時間を取得する。そして、学習部41bは、作業の所要時間が基準よりも短い場合に正(例えば+1)、作業の所要時間が基準よりも長い場合に負(例えば-1)の報酬を与える。なお、基準は種々の要素によって特定されて良く、例えば、前回の作業の所要時間であっても良いし、過去の最短所要時間であっても良いし、予め決められた時間であっても良い。

30

【0198】

さらに、学習部41bは、作業の各工程において、ロボット3のエンコーダーE1~E6の出力をU1に基づいて変換してグリッパー23の位置を取得する。そして、学習部41bは、作業の各工程において取得したグリッパー23の位置を、整定する以前の所定期間にわたって取得し、当該期間における振動強度を取得する。そして、学習部41bは、当該振動強度の程度が基準よりも小さい場合に正、基準よりも大きい場合の負の報酬を与える。なお、基準は種々の要素によって特定されて良く、例えば、前回の振動強度の程度であっても良いし、過去の最小の振動強度の程度であっても良いし、予め決められた振動強度の程度であっても良い。

【0199】

振動強度の程度は、種々の手法で特定されて良く、目標位置からの乖離の積分値や閾値以上の振動が生じている期間など、種々の手法を採用可能である。なお、所定期間は、種々の期間とすることが可能であり、工程の始点から終点にわたる期間であれば、動作中の振動強度による報酬が評価され、工程の終期が所定期間とされれば、残留振動の強度による報酬が評価される。なお、力制御においては、前者の振動強度による報酬の方が重要である場合が多い。前者の振動強度による報酬の方が重要であれば、後者の残留振動の強度による報酬は評価されない構成とされても良い。

40

【0200】

さらに、学習部41bは、作業の各工程の終期において取得したグリッパー23の位置を、整定する以前の所定期間にわたって取得し、当該期間における目標位置からの乖離の最大値をオーバーシュート量として取得する。そして、学習部41bは、当該オーバーシ

50

ュート量が基準よりも小さい場合に正、基準よりも大きい場合の負の報酬を与える。なお、基準は種々の要素によって特定されて良く、例えば、前回のオーバーシュート量の程度であっても良いし、過去の最小のオーバーシュート量であっても良いし、予め決められたオーバーシュート量であっても良い。

【0201】

さらに、本実施形態においては、制御装置40、ロボット1～3、作業台等の少なくとも1カ所に集音装置が取り付けられており、学習部41bは、作業中に集音装置が取得した音を示す情報を取得する。そして、学習部41bは、作業中の発生音の大きさが基準よりも小さい場合に正、基準よりも大きい場合の負の報酬を与える。なお、基準は種々の要素によって特定されて良く、例えば、前回の作業または工程の発生音の大きさの程度であ

10

【0202】

なお、力制御パラメータの学習においては、動作パラメータの学習において報酬とされていた、目標位置からの乖離は報酬に含まれない。すなわち、力制御パラメータの学習においては、工程の始点や終点が学習に応じて変動し得るため、報酬には含まれていない。

【0203】

現在の状態sにおいて行動aが採用された場合における次の状態s'は、行動aとしてのパラメータの変化が行われた後にロボット3を動作させ、状態観測部41aが状態を観測することによって特定可能である。なお、本例にかかる力制御パラメータの学習は、ロボット1, 2による対象物の検出完了後に、ロボット3に関して実行される。

20

【0204】

図11に示す例において、状態変数には、モーターM1～M6の電流、エンコーダーE1～E6の値、力覚センサーPの出力が含まれている。従って、状態観測部41aは、サーボ43dの制御結果として、モーターM1～M6に供給される電流値を観測する。当該電流値は、モーターM1～M6で出力されるトルクに相当する。また、エンコーダーE1～E6の出力は、対応関係U1に基づいてロボット座標系におけるTCPの位置に変換される。従って、状態観測部41aは、ロボット3が備えるグリッパー23の位置情報を観測することになる。

30

【0205】

本実施形態においては、ロボットの運動中に力覚センサーPによって検出された出力を積分することによってロボットの位置を算出することができる。すなわち、状態観測部41aは、対応関係U2に基づいてロボット座標系において運動中のTCPへの作用力を積分することでTCPの位置を取得する。従って、本実施形態において状態観測部41aは、力覚センサーPの出力も利用してロボット3が備えるグリッパー23の位置情報を観測する。なお、状態は、各種の手法で観測されて良く、上述の変換が行われない値（電流値やエンコーダー、力覚センサーの出力値）が状態として観測されても良い。

【0206】

状態観測部41aは、行動であるインピーダンスパラメータや力制御座標系、工程の始点および終点の調整結果を直接的に観測しているのではなく、調整の結果、ロボット3で得られた変化をモーターM1～M6の電流、エンコーダーE1～E6の値、力覚センサーPの出力として観測している。従って、行動による影響を間接的に観測していることになり、この意味で、本実施形態の状態変数は、力制御パラメータの変化から直接的に推定することが困難な変化をし得る状態変数である。

40

【0207】

また、モーターM1～M6の電流、エンコーダーE1～E6の値、力覚センサーPの出力は、ロボット3の動作を直接的に示しており、当該動作は作業の良否を直接的に示している。従って、状態変数として、モーターM1～M6の電流、エンコーダーE1～E6の

50

値、力覚センサー P の出力を観測することにより、人為的に改善することが困難なパラメータの改善を行い、効果的に作業の質を高めるように力制御パラメータを最適化することが可能になる。この結果、人為的に決められた力制御パラメータよりも高性能な動作を行う力制御パラメータを高い確率で算出することができる。

【0208】

(4-6) 力制御パラメータの学習例：

次に、力制御パラメータの学習例を説明する。学習の過程で参照される変数や関数を示す情報は、学習情報 44e として記憶部 44 に記憶される。すなわち、算出部 41 は、状態変数の観測と、当該状態変数に応じた行動の決定と、当該行動によって得られる報酬の評価とを繰り返すことによって行動価値関数 $Q(s, a)$ を収束させる構成が採用されている。そこで、本例において、学習の過程で状態変数と行動と報酬との時系列の値が、順次、学習情報 44e に記録されていく。

10

【0209】

なお、本実施形態において、力制御パラメータの学習は力制御モードで実行される（位置制御のみが行われる位置制御モードでは力制御パラメータの学習は行われない）。力制御モードでの学習を実行するためには、力制御モードのみで構成される作業がロボット 3 のロボットプログラム 44b として生成されても良いし、任意のモードを含む作業がロボット 3 のロボットプログラム 44b として生成されている状況において、その中の力制御モードのみを用いて学習してもよい。

【0210】

20

行動価値関数 $Q(s, a)$ は、種々の手法で算出されて良く、多数回の試行に基づいて算出されても良いが、ここでは、DQN によって行動価値関数 Q を最適化する例を説明する。行動価値関数 Q の最適化に利用される多層ニューラルネットワークは、上述の図 8 において模式的に示される。図 11 に示すような状態が観測される本例であれば、ロボット 3 におけるモーター M1 ~ M6 の電流、エンコーダー E1 ~ E6 の値、力覚センサー P の出力（6 軸の出力）が状態であるため、状態 s の数 $M = 18$ である。図 11 に示すような行動が選択され得る本例であれば、60 個の行動が選択可能であるため $N = 78$ である。むろん、行動 a の内容や数（ N の値）、状態 s の内容や数（ M の値）は試行番号 t に応じて変化しても良い。

【0211】

30

本実施形態においても、当該多層ニューラルネットワークを特定するためのパラメータ（入力から出力を得るために必要な情報）が学習情報 44e として記憶部 44 に記録される。ここでも学習の過程で変化し得る多層ニューラルネットワークのパラメータをと表記する。当該を使用すると、上述の行動価値関数 $Q(s_t, a_{1t}) \sim Q(s_t, a_{Nt})$ は、 $Q(s_t, a_{1t}; \theta_t) \sim Q(s_t, a_{Nt}; \theta_t)$ とも表記できる。

【0212】

次に、図 9 に示すフローチャートに沿って学習処理の手順を説明する。力制御パラメータの学習処理は、ロボット 3 の運用過程において実施されても良いし、実運用の前に事前に学習処理が実行されてもよい。ここでは、実運用の前に事前に学習処理が実行される構成（多層ニューラルネットワークを示す が最適化されると、その情報が保存され、次回以降の運用で利用される構成）に従って学習処理を説明する。

40

【0213】

学習処理が開始されると、算出部 41 は、学習情報 44e を初期化する（ステップ S200）。すなわち、算出部 41 は、学習を開始する際に参照される の初期値を特定する。初期値は、種々の手法によって決められて良く、過去に学習が行われていない場合においては、任意の値やランダム値等が の初期値となっても良いし、ロボット 3 や対象物を模擬するシミュレーション環境を準備し、当該環境に基づいて学習または推定した を初期値としてもよい。

【0214】

過去に学習が行われた場合は、当該学習済の が初期値として採用される。また、過去

50

に類似の対象についての学習が行われた場合は、当該学習における が初期値とされても良い。過去の学習は、ロボット3を用いてユーザーが行ってもよいし、ロボット3の製造者がロボット3の販売前に行ってもよい。この場合、製造者は、対象物や作業の種類に応じて複数の初期値のセットを用意しておき、ユーザーが学習する際に初期値を選択する構成であっても良い。 の初期値が決定されると、当該初期値が現在の の値として学習情報44eに記憶される。

【0215】

次に、算出部41は、パラメーターを初期化する(ステップS205)。ここでは、力制御パラメーターが学習対象であるため、算出部41は、力制御パラメーターを初期化する。すなわち、学習が行われていない状態であれば、算出部41は、教示によって生成されたパラメーター44aに含まれる力制御パラメーターを初期値として設定する。過去に何らかの学習が行われた状態であれば、算出部41は、学習の際に最後に利用されていたパラメーター44aに含まれる力制御パラメーターを初期値として設定する。

【0216】

次に、状態観測部41aは、状態変数を観測する(ステップS210)。すなわち、制御部43は、パラメーター44aおよびロボットプログラム44bを参照してロボット3を制御する(上述のステップS110~S130に相当)。この後、状態観測部41aは、モーターM1~M6に供給される電流値を観測する。また、状態観測部41aは、エンコーダーE1~E6の出力を取得し、対応関係U1に基づいてロボット座標系におけるTCPの位置に変換する。さらに、状態観測部41aは、力覚センサーPの出力を積分し、TCPの位置を取得する。

【0217】

次に、学習部41bは、行動価値を算出する(ステップS215)。すなわち、学習部41bは、学習情報44eを参照して を取得し、学習情報44eが示す多層ニューラルネットワークに最新の状態変数を入力し、N個の行動価値関数 $Q(s_t, a_{1t}; \gamma) \sim Q(s_t, a_{Nt}; \gamma)$ を算出する。

【0218】

なお、最新の状態変数は、初回の実行時においてステップS210、2回目以降の実行時においてステップS225の観測結果である。また、試行番号tは初回の実行時において0、2回目以降の実行時において1以上の値となる。学習処理が過去に実施されていない場合、学習情報44eが示す は最適化されていないため、行動価値関数Qの値としては不正確な値となり得るが、ステップS215以後の処理の繰り返しにより、行動価値関数Qは徐々に最適化していく。また、ステップS215以後の処理の繰り返しにおいて、状態s、行動a、報酬rは、各試行番号tに対応づけられて記憶部44に記憶され、任意のタイミングで参照可能である。

【0219】

次に、学習部41bは、行動を選択し、実行する(ステップS220)。本実施形態においては、行動価値関数 $Q(s, a)$ を最大化する行動aが最適な行動であると見なされる処理が行われる。そこで、学習部41bは、ステップS215において算出されたN個の行動価値関数 $Q(s_t, a_{1t}; \gamma) \sim Q(s_t, a_{Nt}; \gamma)$ の値の中で最大の値を特定する。そして、学習部41bは、最大の値を与えた行動を選択する。例えば、N個の行動価値関数 $Q(s_t, a_{1t}; \gamma) \sim Q(s_t, a_{Nt}; \gamma)$ の中で $Q(s_t, a_{Nt}; \gamma)$ が最大値であれば、学習部41bは、行動 a_{Nt} を選択する。

【0220】

行動が選択されると、学習部41bは、当該行動に対応するパラメーター44aを変化させる。例えば、図11に示す例において、ロボット座標系のx軸に関するインピーダンスパラメーターmを一定値増加させる行動a1が選択された場合、学習部41bは、力制御パラメーターが示すx軸に関するインピーダンスパラメーターmを一定値増加させる。パラメーター44aの変化が行われると、制御部43は、当該パラメーター44aを参照してロボット3を制御し、一連の作業を実行させる。なお、本実施形態においては、行動

選択のたびに一連の作業が実行されるが、行動選択のたびに一連の作業の一部が実行される構成（一連の作業を構成する複数の工程の少なくとも１工程が実行される構成）であっても良い。

【０２２１】

次に、状態観測部４１ａは、状態変数を観測する（ステップＳ２２５）。すなわち、状態観測部４１ａは、ステップＳ２１０における状態変数の観測と同様の処理を行って、状態変数として、モーターＭ１～Ｍ６に供給される電流値、エンコーダーＥ１～Ｅ６の出力に基づいて特定されるＴＣＰの位置、力覚センサーＰの出力に基づいて特定されるＴＣＰの位置を取得する。なお、現在の試行番号が t である場合（選択された行動が a_t である場合）、ステップＳ２２５で取得される状態 s は s_{t+1} である。

10

【０２２２】

次に、学習部４１ｂは、報酬を評価する（ステップＳ２３０）。すなわち、学習部４１ｂは、図示しない計時回路に基づいて作業の開始から終了までの所要時間を取得し、作業の所要時間が基準よりも短い場合に正、作業の所要時間が基準よりも長い場合に負の報酬を取得する。さらに、学習部４１ｂは、作業の各工程におけるグリッパー２３の位置を取得し、作業の各工程において取得したグリッパー２３の位置に基づいて振動強度を取得する。そして、学習部４１ｂは、当該振動強度の程度が基準よりも小さい場合に正、基準よりも大きい場合の負の報酬を取得する。一連の作業が複数の工程で構成される場合、各工程の報酬の和が取得されても良いし、統計値（平均値等）が取得されても良い。

【０２２３】

20

さらに、学習部４１ｂは、作業の各工程の終期において取得したグリッパー２３の位置に基づいてオーバーシュート量を取得する。そして、学習部４１ｂは、当該オーバーシュート量が基準よりも小さい場合に正、基準よりも大きい場合の負の報酬を取得する一連の作業が複数の工程で構成される場合、各工程の報酬の和が取得されても良いし、統計値（平均値等）が取得されても良い。

【０２２４】

さらに、学習部４１ｂは、作業中に集音装置が取得した音を示す情報を取得する。そして、学習部４１ｂは、作業中の発生音の大きさが基準よりも小さい場合に正、基準よりも大きい場合の負の報酬を取得する。なお、現在の試行番号が t である場合、ステップＳ２３０で取得される報酬 r は r_{t+1} である。

30

【０２２５】

本例においても、式（２）に示す行動価値関数 Q の更新を目指しているが、行動価値関数 Q を適切に更新していくためには、行動価値関数 Q を示す多層ニューラルネットワークを最適化（を最適化）していかななくてはならない。そして、図８に示す多層ニューラルネットワークによって行動価値関数 Q を適正に出力させるためには、当該出力のターゲットとなる教師データが必要になる。すなわち、多層ニューラルネットワークの出力と、ターゲットとの誤差を最小化するようにを改善すると、多層ニューラルネットワークが最適化されることが期待される。

【０２２６】

40

しかし、本実施形態において、学習が完了していない段階では行動価値関数 Q の知見がなく、ターゲットを特定することは困難である。そこで、本実施形態においては、式（２）の第２項、いわゆるＴＤ誤差を最小化する目的関数によって多層ニューラルネットワークを示すの改善を実施する。すなわち、 $(r_{t+1} + \max_{a'} Q(s_{t+1}, a'; \theta_t))$ をターゲットとし、ターゲットと $Q(s_t, a_t; \theta_t)$ との誤差が最小化するようにを学習する。ただし、ターゲット $(r_{t+1} + \max_{a'} Q(s_{t+1}, a'; \theta_t))$ は、学習対象のを含んでいるため、本実施形態においては、ある程度の試行回数にわたりターゲットを固定する（例えば、最後に学習した（初回学習時はの初期値）で固定する）。本実施形態においては、ターゲットを固定する試行回数である既定回数が予め決められている。

【０２２７】

50

このような前提で学習を行うため、ステップ S 2 3 0 で報酬が評価されると、学習部 4 1 b は目的関数を算出する（ステップ S 2 3 5）。すなわち、学習部 4 1 b は、試行のそれぞれにおける T D 誤差を評価するための目的関数（例えば、T D 誤差の 2 乗の期待値に比例する関数や T D 誤差の 2 乗の総和等）を算出する。なお、T D 誤差は、ターゲットが固定された状態で算出されるため、固定されたターゲットを $(r_{t+1} + \max_a Q(s_{t+1}, a; \cdot))$ と表記すると、T D 誤差は $(r_{t+1} + \max_a Q(s_{t+1}, a; \cdot) - Q(s_t, a_t; \cdot))$ である。当該 T D 誤差の式において報酬 r_{t+1} は、行動 a_t によってステップ S 2 3 0 で得られた報酬である。

【0228】

また、 $\max_a Q(s_{t+1}, a; \cdot)$ は、行動 a_t によってステップ S 2 2 5 で算出される状態 s_{t+1} を、固定された \cdot で特定される多層ニューラルネットワークの入力とした場合に得られる出力の中の最大値である。 $Q(s_t, a_t; \cdot)$ は、行動 a_t が選択される前の状態 s_t を、試行番号 t の段階の \cdot で特定される多層ニューラルネットワークの入力とした場合に得られる出力の中で、行動 a_t に対応した出力の値である。

【0229】

目的関数が算出されると、学習部 4 1 b は、学習が終了したか否かを判定する（ステップ S 2 4 0）。本実施形態においては、T D 誤差が十分に小さいか否かを判定するための閾値が予め決められており、目的関数が閾値以下である場合、学習部 4 1 b は、学習が終了したと判定する。

【0230】

ステップ S 2 4 0 において学習が終了したと判定されない場合、学習部 4 1 b は、行動価値を更新する（ステップ S 2 4 5）。すなわち、学習部 4 1 b は、T D 誤差の \cdot による偏微分に基づいて目的関数を小さくするための \cdot の変化を特定し、 \cdot を変化させる。むろん、ここでは、各種の手法で \cdot を変化させることが可能であり、例えば、RMSProp 等の勾配降下法を採用可能である。また、学習率等による調整も適宜実施されて良い。以上の処理によれば、行動価値関数 Q がターゲットに近づくように \cdot を変化させることができる。

【0231】

ただし、本実施形態においては、上述のようにターゲットが固定されているため、学習部 4 1 b は、さらに、ターゲットを更新するか否かの判定を行う。具体的には学習部 4 1 b は、既定回数の試行が行われたか否かを判定し（ステップ S 2 5 0）、ステップ S 2 5 0 において、既定回数の試行が行われたと判定された場合に、学習部 4 1 b は、ターゲットを更新する（ステップ S 2 5 5）。すなわち、学習部 4 1 b は、ターゲットを算出する際に参照される \cdot を最新の \cdot に更新する。この後、学習部 4 1 b は、ステップ S 2 1 5 以降の処理を繰り返す。一方、ステップ S 2 5 0 において、既定回数の試行が行われたと判定されなければ、学習部 4 1 b は、ステップ S 2 5 5 をスキップしてステップ S 2 1 5 以降の処理を繰り返す。

【0232】

ステップ S 2 4 0 において学習が終了したと判定された場合、学習部 4 1 b は、学習情報 4 4 e を更新する（ステップ S 2 6 0）。すなわち、学習部 4 1 b は、学習によって得られた \cdot を、ロボット 3 による作業の際に参照されるべき \cdot として学習情報 4 4 e に記録する。当該 \cdot を含む学習情報 4 4 e が記録されている場合、ステップ 1 1 0 ~ S 1 3 0 のようにロボット 3 による作業が行われる際に、制御部 4 3 はパラメータ 4 4 a に基づいてロボット 3 を制御する。そして、当該作業の過程においては、状態観測部 4 1 a による現在の状態の観測と、学習部 4 1 b による行動の選択が繰り返される。むろん、この際、学習部 4 1 b は、状態を入力として算出された出力 $Q(s, a)$ の中で最大値を与える行動 a を選択する。そして、行動 a が選択された場合、行動 a が行われた状態に相当する値となるようにパラメータ 4 4 a が更新される。

【0233】

以上の構成によれば、制御部 4 3 は、行動価値関数 Q が最大化される行動 a を選択しな

10

20

30

40

50

がら作業を実行することができる。当該行動価値関数 Q は、上述の処理により、多数の試行が繰り返された結果、最適化されている。そして、当該試行は、算出部41によって自動で行われ、人為的に実施不可能な程度の多数の試行を容易に実行することができる。従って、本実施形態によれば、人為的に決められた力制御パラメータよりも高い確率でロボット3の作業の質を高めることができる。

【0234】

さらに、本実施形態においては、行動によってパラメータ44aとしてのインピーダンスパラメータが変化する。従って、人為的な調整によって適切な設定を行うことが困難な、インピーダンスパラメータを自動的に調整することができる。さらに、本実施形態においては、行動によってパラメータ44aとしての始点と終点が変わる。従って、人為的に設定された始点や終点を、より高性能に力制御を行うように自動的に調整することができる。

10

【0235】

さらに、本実施形態においては、行動によってパラメータ44aとしての力制御座標系が変化する。この結果、ロボット3のTCPからのオフセット点の位置が変わる。従って、人為的な調整によって適切な設定を行うことが困難な、TCPからのオフセット点の位置を自動的に調整することができる。さらに、本実施形態においては、行動によってパラメータ44aとしての目標力が変化し得る。従って、人為的な調整によって適切な設定を行うことが困難な、目標力を自動的に調整することができる。特に、力制御座標系と目標力との組み合わせを人為的に理想化することは困難であるため、これらの組が自動的に調整される構成は、有用である。

20

【0236】

さらに、本実施形態においては、ロボット3が行った作業の良否に基づいて行動による報酬を評価するため、ロボット3の作業を成功させるようにパラメータを最適化することができる。さらに、本実施形態においては、作業の所要時間が基準よりも短い場合に報酬を正と評価するため、ロボット3を短い時間で作業させる力制御パラメータを容易に算出することができる。

【0237】

さらに、本実施形態においては、振動強度が基準よりも小さい場合に報酬を正と評価するため、ロボット3の動作による振動を発生させる可能性が低い力制御パラメータを容易に算出することができる。さらに、本実施形態においては、ロボット3の位置のオーバーシュートが基準よりも小さい場合に報酬を正と評価するため、ロボット3がオーバーシュートする可能性が低い力制御パラメータを容易に算出することができる。さらに、本実施形態においては、発生音が基準よりも小さい場合に報酬を正と評価するため、ロボット3に異常を発生させる可能性が低い力制御パラメータを容易に算出することができる。

30

【0238】

さらに、本実施形態によれば、自動で行動価値関数 Q が最適化されるため、高性能な力制御を行う力制御パラメータを容易に算出することができる。また、行動価値関数 Q の最適化は自動的に行われるため、最適な力制御パラメータの算出も自動的に行うことができる。

40

【0239】

さらに、本実施形態においては、ロボット3において汎用的に使用される力覚センサーPによってロボット3の位置情報を取得するため、ロボット3で汎用的に使用されるセンサーに基づいて位置情報を算出することができる。

【0240】

さらに、本実施形態において学習部41bは、状態変数としてのロボット3の動作結果を実測し、力制御パラメータを最適化する。従って、ロボット3によって作業が行われている実環境下において合わせて力制御パラメータを最適化することができる。従って、ロボット3の使用環境に応じた力制御パラメータとなるように最適化することができ

50

る。

【 0 2 4 1 】

さらに、本実施形態において状態観測部 4 1 a は、ロボット 3 にエンドエフェクターとしてのグリッパー 2 3 が設けられた状態で状態変数を観測する。また、学習部 4 1 b は、ロボット 3 にエンドエフェクターとしてのグリッパー 2 3 が設けられた状態で行動としてのパラメーター 4 4 a の変更が実行される。この構成によれば、エンドエフェクターとしてのグリッパー 2 3 を用いた動作を行うロボット 3 に適した力制御パラメーターを容易に算出することができる。

【 0 2 4 2 】

さらに、本実施形態において状態観測部 4 1 a は、エンドエフェクターとしてのグリッパー 2 3 が対象物を把持した状態で状態変数を観測する。また、学習部 4 1 b は、エンドエフェクターとしてのグリッパー 2 3 が対象物を把持した状態で行動としてのパラメーター 4 4 a の変更が実行される。この構成によれば、エンドエフェクターとしてのグリッパー 2 3 で対象物を把持して動作を行うロボット 3 に適した力制御パラメーターを容易に算出することができる。

【 0 2 4 3 】

(5) 他の実施形態：

以上の実施形態は本発明を実施するための一例であり、他にも種々の実施形態を採用可能である。例えば、制御装置は、ロボットに内蔵されていても良いし、ロボットの設置場所と異なる場所、例えば外部のサーバー等に備えられていても良い。また、制御装置は、複数の装置で構成されていても良く、制御部 4 3 と算出部 4 1 とが異なる装置で構成されても良い。また、制御装置は、ロボットコントローラー、ティーチングペンダント、P C、ネットワークにつながるサーバー等であっても良いし、これらが含まれていても良い。さらに、上述の実施形態の一部の構成が省略されてもよいし、処理の順序が変動または省略されてもよい。さらに、上述の実施形態においては、T C P について目標位置や目標力の初期ベクトルが設定されたが、他の位置、例えば力覚センサー P についてのセンサー座標系の原点やネジの先端等について目標位置や目標力の初期ベクトルが設定されても良い。

【 0 2 4 4 】

ロボットは、任意の態様の可動部で任意の作業を実施できれば良い。エンドエフェクターは、対象物に関する作業に利用される部位であり、任意のツールが取り付けられて良い。対象物は、ロボットによる作業対象となる物体であれば良く、エンドエフェクターによって把持された物体であっても良いし、エンドエフェクターが備えるツールで扱われる物体であっても良く、種々の物体が対象物となり得る。

【 0 2 4 5 】

ロボットに作用させる目標力は、当該ロボットを力制御によって駆動する際にロボットに作用させる目標力であれば良く、例えば、力覚センサー等の力検出部によって検出される力（または当該力から算出される力）を特定の力に制御する際に、当該力が目標力となる。また、力覚センサー以外のセンサー、例えば加速度センサーで検出される力（または当該力から算出される力）が目標力になるように制御されても良いし、加速度や角速度が特定の値になるように制御されても良い。

【 0 2 4 6 】

さらに、上述の学習処理においては、試行のたびに の更新によって行動価値を更新し、既定回数の試行が行われるまでターゲットを固定したが、複数回の試行が行われてからの更新が行われてもよい。例えば、第 1 既定回数の試行が行われるまでターゲットが固定され、第 2 既定回数（＜第 1 既定回数）の試行が行われるまで を固定する構成が挙げられる。この場合、第 2 既定回数の試行後に第 2 既定回数分のサンプルに基づいて を更新し、さらに試行回数が第 1 既定回数を超えた場合に最新の でターゲットを更新する構成となる。

【 0 2 4 7 】

10

20

30

40

50

さらに、学習処理においては、公知の種々の手法が採用されてよく、例えば、体験再生や報酬のClipping等が行われてもよい。さらに、図8においては、層DLがP個（Pは1以上の整数）存在し、各層において複数のノードが存在するが、各層の構造は、種々の構造を採用可能である。例えば、層の数やノードの数は種々の数を採用可能であるし、活性化関数としても種々の関数を採用可能であるし、ネットワーク構造が畳み込みニューラルネットワーク構造等になっていても良い。また、入力や出力の態様も図8に示す例に限定されず、例えば、状態sと行動aとが入力される構成や、行動価値関数Qを最大化する行動aがone-hotベクトルとして出力される構成が少なくとも利用される例が採用されても良い。

【0248】

上述の実施形態においては、行動価値関数に基づいてgreedy方策で行動を行って試行しながら、行動価値関数を最適化することにより、最適化された行動価値関数に対するgreedy方策が最適方策であると見なしている。この処理は、いわゆる価値反復法であるが、他の手法、例えば、方策反復法によって学習が行われてもよい。さらに、状態s、行動a、報酬r等の各種変数においては、各種の正規化が行われてよい。

【0249】

機械学習の手法としては、種々の手法を採用であり、行動価値関数Qに基づいたgreedy方策によって試行が行われてもよい。また、強化学習の手法としても上述のようなQ学習に限定されず、SARSA等の手法が用いられてもよい。また、方策のモデルと行動価値関数のモデルを別々にモデル化した手法、例えば、Actor-Criticアルゴリズムが利用されてもよい。Actor-Criticアルゴリズムを利用するのであれば、方策を示すactorである $\mu(s; \cdot)$ と、行動価値関数を示すcriticである $Q(s, a; \cdot)$ とを定義し、 $\mu(s; \cdot)$ にノイズを加えた方策に従って行動を生成して試行し、試行結果に基づいてactorとcriticを更新することで方策と行動価値関数とを学習する構成であっても良い。

【0250】

算出部は、機械学習を用いて、学習対象のパラメータを算出することができればよく、パラメータとしては、光学パラメータ、画像処理パラメータ、動作パラメータ、力制御パラメータの少なくとも1個であれば良い。機械学習は、サンプルデータを用いてよりよいパラメータを学習する処理であれば良く、上述の強化学習以外にも、教師あり学習やクラスタリングなど種々の手法によって各パラメータを学習する構成を採用可能である。

【0251】

光学系は、対象物を撮像することができる。すなわち、対象物が含まれる領域を視野にした画像を取得する構成を備える。光学系の構成要素としては上述のように、撮像部や照明部を含むことが好ましく、他にも種々の構成要素が含まれていて良い。また、上述のように、撮像部や照明部はロボットのアームによって移動可能であっても良いし、2次元的な移動機構によって移動可能であっても良いし、固定的であっても良い。むしろ、撮像部や照明部は交換可能であっても良い。また、光学系で用いる光（撮像部による検出光や照明部の出力光）の帯域は可視光帯域に限定されず、赤外線や紫外線、X線等の任意の電磁波が用いられる構成が採用可能である。

【0252】

光学パラメータは、光学系の状態を変化させ得る値であれば良く、撮像部や照明部等で構成される光学系において状態を直接的または間接的に特定するための数値等が光学パラメータとなる。例えば、撮像部や照明部等の位置や角度等を示す値のみならず、撮像部や照明部の種類を示す数値（IDや型番等）が光学パラメータとなり得る。

【0253】

検出部は、算出された光学パラメータによる光学系での撮像結果に基づいて、対象物を検出することができる。すなわち、検出部は、学習された光学パラメータによって光学系を動作させて対象物を撮像し、撮像結果に基づいて対象物の検出処理を実行する構成

10

20

30

40

50

を備える。

【0254】

検出部は対象物を検出することができればよく、上述の実施形態のように、対象物の位置姿勢が検出される構成の他、対象物の有無が検出される構成であっても良く、種々の構成を採用可能である。なお、対象物の位置姿勢は、例えば、3軸における位置と3軸に対する回転角とによる6個のパラメーターによって定義可能であるが、むしろ、必要に応じて任意の数のパラメーターが考慮されなくても良い。例えば、平面上に設置された対象物であれば、少なくとも1個の位置に関するパラメーターが既知であるとして検出対象から除外されても良い。また、平面に固定的な向きで設置された対象物であれば、姿勢に関するパラメーターが検出対象から除外されても良い。

10

【0255】

対象物は、光学系で撮像され、検出される対象となる物体であればよく、ロボットの作業対象となるワークや、ワークの周辺の物体、ロボットの一部など、種々の物体が想定可能である。また、撮像結果に基づいて対象物を検出する手法としても種々の手法を採用可能であり、画像の特徴量抽出によって対象物が検出されても良いし、対象物の動作（人等の可動物体等の検出）によって対象物が検出されても良く、種々の手法が採用されてよい。

【0256】

制御部は、対象物の検出結果に基づいてロボットを制御することができる。すなわち、制御部は、対象物の検出結果に応じてロボットの制御内容を決定する構成を備える。従って、ロボットの制御は、上述のような対象物をつかむための制御の他にも種々の制御が行われてよい。例えば、対象物に基づいてロボットの位置決めをする制御や、対象物に基づいてロボットの動作を開始または終了させる制御など、種々の制御が想定される。

20

【0257】

ロボットの態様は、種々の態様であって良く、上述の実施形態のような垂直多関節ロボット以外にも直交ロボット、水平多関節ロボット、双腕ロボット等であって良い。また、種々の態様のロボットが組み合わせられても良い。むしろ、軸の数やアームの数、エンドエフェクターの態様等は種々の態様を採用可能である。例えば、撮像部21や照明部22がロボット3の上方に存在する平面に取り付けられ、当該平面上で撮像部21や照明部22が移動可能であっても良い。

30

【0258】

状態観測部は、行動等の試行に応じて変化した結果を観測することができればよく、各種のセンサー等によって状態が観測されても良いし、ある状態から他の状態に変化させる制御が行われ、制御の失敗（エラー等）が観測されなければ当該他の状態が観測されたと見なされる構成であっても良い。前者のセンサーによる観測は、位置等の検出の他にも撮像センサーによる画像の取得も含まれる。

【0259】

さらに、上述の実施形態における行動や状態、報酬は例であり、他の行動や状態、報酬を含む構成や任意の行動や状態が省略された構成であっても良い。例えば、撮像部21や照明部22が交換可能であるロボット1,2において、撮像部21や照明部22の種類の変更を行動として選択可能であり、状態として種類を観測可能であっても良い。接触判定部43cによる判定結果に基づいて報酬が決定されても良い。すなわち、学習部41bにおける学習過程において、接触判定部43cが作業において想定されていない物体とロボットとが接触したと判定した場合、当該直前の行動による報酬を負に設定する構成を採用可能である。この構成によれば、ロボットが想定外の物体に接触しないようにパラメーター44aを最適化することができる。

40

【0260】

また、例えば、光学パラメーターの最適化に際して、ロボット1~3によって対象物の検出結果に基づいた作業（例えば、上述のピックアップ作業等）を行い、学習部41bが、対象物の検出結果に基づいてロボット1~3が行った作業の良否に基づいて、行動によ

50

る報酬を評価する構成であってもよい。この構成は、例えば、図 7 に示す報酬の中で、対象物の検出の替わりに、または、対象物の検出に加えて作業の成否（例えば、ピックアップの成否）を報酬とする構成が挙げられる。

【0261】

作業の成否は、例えば、作業の成否を判定可能な工程（ピックアップの工程等）におけるステップ S 120 の判定結果等で定義可能である。この場合、行動や状態において、ロボット 1～3 の動作に関する行動や状態が含まれても良い。さらに、この構成においては、ロボット 1～3 の作業対象である対象物を撮像部 21 および照明部 22 を備える光学系で撮像した画像を状態とすることが好ましい。この構成によれば、ロボットの作業を成功させるように光学パラメーターを最適化することができる。なお、光学パラメーターや動作パラメーター、力制御パラメーターを学習するために観測される状態としての画像は、撮像部 21 で撮像された画像そのものであっても良いし、撮像部 21 で撮像された画像に対して画像処理（例えば、上述の平滑化処理や鮮鋭化処理等）が行われた後の画像であっても良い。

10

【0262】

さらに、光学パラメーター、動作パラメーター、力制御パラメーターのそれぞれを別個に最適化するのではなく、これらのパラメーターの中の 2 種以上を最適化する構成が採用されてもよい。例えば、図 7 に示す例において、動作パラメーターや力制御パラメーターを変化させる行動が含まれる構成であれば、光学パラメーターとともに、動作パラメーターや力制御パラメーターを最適化することが可能である。この場合、最適化された動作パラメーターや力制御パラメーターに基づいてロボット 1～3 が制御される。この構成によれば、対象物の検出を伴う作業を行うパラメーターを最適化することができ、対象物の検出精度を高める学習を実行することができる。

20

【0263】

画像処理パラメーターは、対象物の撮像結果としての画像を変化させ得る値であれば良く、図 3 に示す例に限定されず、追加または削除されてよい。例えば、画像処理の有無や画像処理の強度、画像処理の順序など、実行される画像処理アルゴリズムを特定するための数値（処理順序等を示すフラグ等を含む）等が画像処理パラメーターとなり得る。より具体的には、画像処理としては、二値化処理、直線検出処理、円検出処理、色検出処理、OCR 処理等があげられる。

30

【0264】

さらに、画像処理は、複数の種類の画像処理を組み合わせた処理であってもよい。例えば、円検出処理と OCR 処理を組み合わせて、「円内の文字を認識する処理」という処理が行われてもよい。いずれにしても、各画像処理の有無や強度を示すパラメーターが画像処理パラメーターとなり得る。また、これらの画像処理パラメーターの変化が行動となり得る。

【0265】

動作パラメーターは、上述の実施形態に挙げられたパラメーターに限定されない。例えば、学習対象となる動作パラメーターに、ロボット 1～3 が備える慣性センサーに基づいて制御を行うためのサーボゲインが含まれていても良い。すなわち、慣性センサーの出力に基づいた制御ループでモーター M 1～M 6 が制御される構成において、当該制御ループにおけるサーボゲインが行動によって変化する構成であっても良い。例えば、ロボット 1～3 に取り付けられたエンコーダー E 1～E 6 に基づいてロボット 1～3 の特定の部位の角速度を算出し、慣性センサーの一種であるジャイロセンサーによって当該特定の部位の角速度を検出し、両者の差分にジャイロサーボゲインを乗じてフィードバック制御を行う構成において、当該ジャイロサーボゲインが行動によって変化する構成が挙げられる。この構成であれば、ロボットの特定の部位に生じる角速度の振動成分を抑制する制御を行うことができる。むろん、慣性センサーはジャイロセンサーに限定されず、加速度センサー等において同様のフィードバック制御が行われる構成において加速度ゲインが行動によって変化する構成であっても良い。以上の構成によれば、人為的な調整によって適切な設定

40

50

を行うことが困難な、慣性センサーに基づいて制御を行うためのサーボゲインを自動的に調整することができる。なお、加速度センサーはロボットの運動によって生じる加速度を検知するセンサーであり、上述の力覚センサーはロボットに作用する力を検知するセンサーである。通常、加速度センサーと力覚センサーとは異なるセンサーであるが、一方が他方の機能を代替できる場合には、一方が他方として機能しても良い。

【0266】

むろん、力制御パラメーターも上述の実施形態に挙げられたパラメーターに限定されないし、学習対象となるパラメーターも適宜選択されてよい。例えば、目標力に関し、6軸中の全成分または一部の成分が行動として選択し得ない(すなわち固定である)構成であっても良い。この構成は、固定された固定対象物(細い筒等)に、ロボットが把持した対象物を挿入する作業において、目標力は固定対象物のある点に対して固定的な成分を有するが、ロボットの挿入作業に応じて力制御座標系が変化するように学習する構成等を想定する事ができる。

10

【0267】

学習部41bは、ロボット3が把持した対象物を作業完了前に落とした場合、ロボット3の作業対象である対象物の一部が作業完了前に分離した場合、ロボット3が破損した場合、ロボット3の作業対象である対象物が破損した場合、の少なくとも1つにおいて報酬を負と評価する構成であっても良い。ロボット3が把持した対象物を作業完了前に落とした場合に報酬を負と評価する構成によれば、対象物を落とさずに作業を完了させる可能性が高い動作パラメーターや力制御パラメーターを容易に算出することができる。

20

【0268】

ロボット3の作業対象である対象物の一部が作業完了前に分離した場合に報酬を負と評価する構成によれば、対象物を分離させることなく作業を完了させる可能性が高い動作パラメーターや力制御パラメーターを容易に算出することができる。ロボット3が破損した場合に報酬を負と評価する構成によれば、ロボット3を破損させる可能性が低い動作パラメーターや力制御パラメーターを容易に算出することができる。

【0269】

ロボット3の作業対象である対象物が破損した場合に報酬を負と評価する構成によれば、対象物を破損させる可能性が低い動作パラメーターや力制御パラメーターを容易に算出することができる。なお、ロボット3が把持した対象物を作業完了前に落としたか否か、ロボット3の作業対象である対象物の一部が作業完了前に分離したか否か、ロボット3が破損したか否か、ロボット3の作業対象である対象物が破損したか否かは、各種のセンサー、例えば撮像部21等によって検出される構成を採用可能である。

30

【0270】

さらに、学習部41bは、ロボット3による作業が正常に完了した場合において報酬を正と評価する構成であっても良い。ロボット3による作業が正常に完了した場合に報酬を正と評価する構成によれば、ロボット3の作業を成功させる動作パラメーターや力制御パラメーターを容易に算出することができる。

【0271】

さらに、ロボット3の位置を検出するための位置検出部は、上述の実施形態のようなエンコーダー、力覚センサーに限定されず、他のセンサー、専用の慣性センサーや撮像部21等の光学センサー、距離センサー等であっても良い。また、センサーはロボットに内蔵されていても良いが、ロボットの外部に配置されても良い。ロボットの外部に配置された位置検出部を利用すれば、ロボットの動作に影響されることなく位置情報を算出することができる。

40

【0272】

さらに、算出部41は、ロボットの異なる複数の動作に基づいて、複数の動作に共通の動作パラメーターや力制御パラメーターを算出する構成であっても良い。複数の動作は、最適化された動作パラメーターを利用して実行される動作を含んでいれば良い。従って、複数の動作は、異なる種類の複数の作業(ピックアップ作業、研磨作業、ネジ締め作業な

50

ど)である構成や、同種の作業(ネジの大きさが異なる複数のネジ締め作業等)である構成等が挙げられる。この構成によれば、各種の動作に適用可能な汎用的な動作パラメータや力制御パラメータを容易に算出することができる。

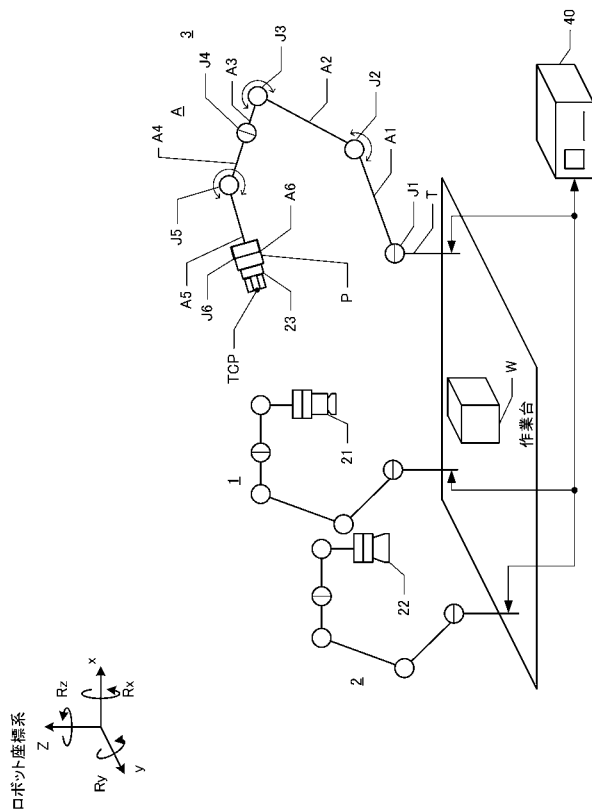
【符号の説明】

【0273】

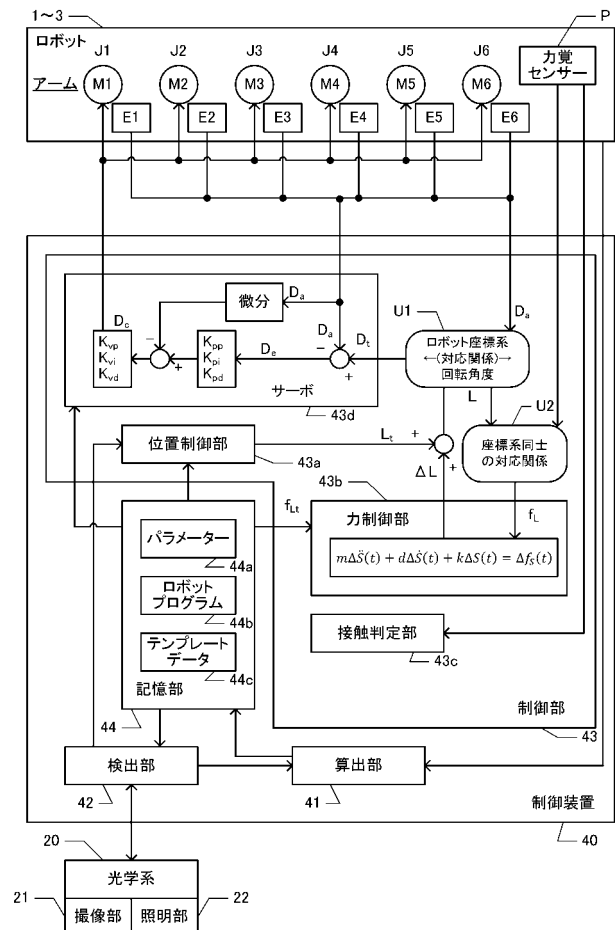
1～3…ロボット、20…光学系、21…撮像部、22…照明部、23…グリッパー、40…制御装置、41…算出部、41a…状態観測部、41b…学習部、42…検出部、43…制御部、43a…位置制御部、43b…力制御部、43c…接触判定部、43d…サーボ、44…記憶部、44a…パラメータ、44b…ロボットプログラム、44c…テンプレートデータ、44d…行動情報、44e…学習情報

10

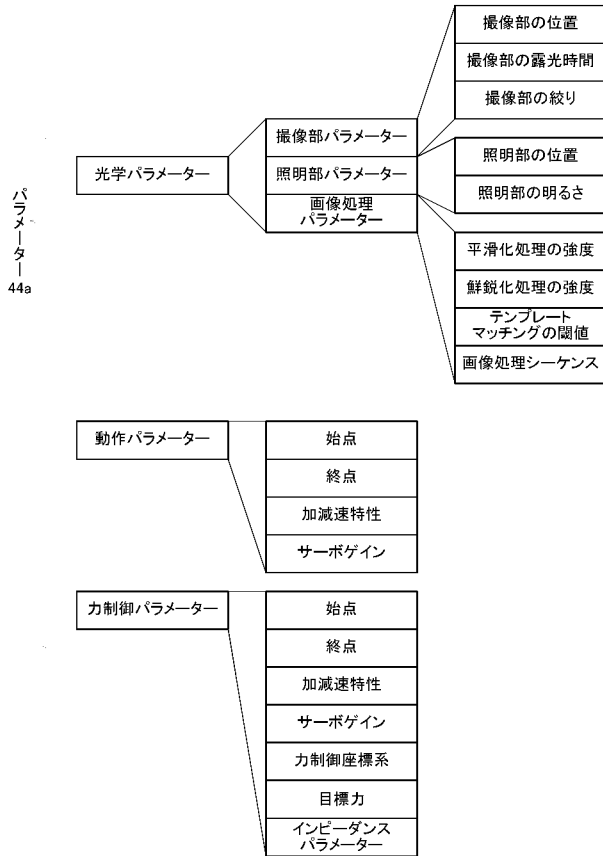
【図1】



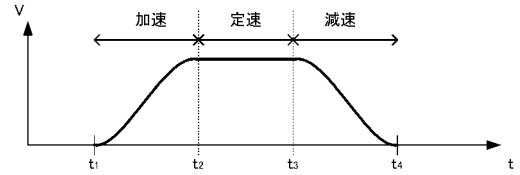
【図2】



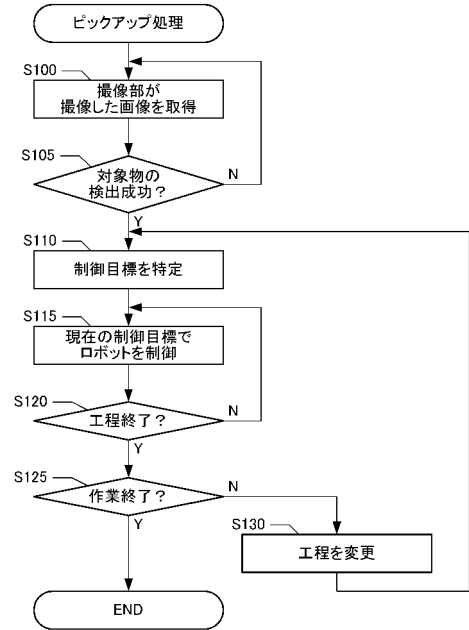
【図 3】



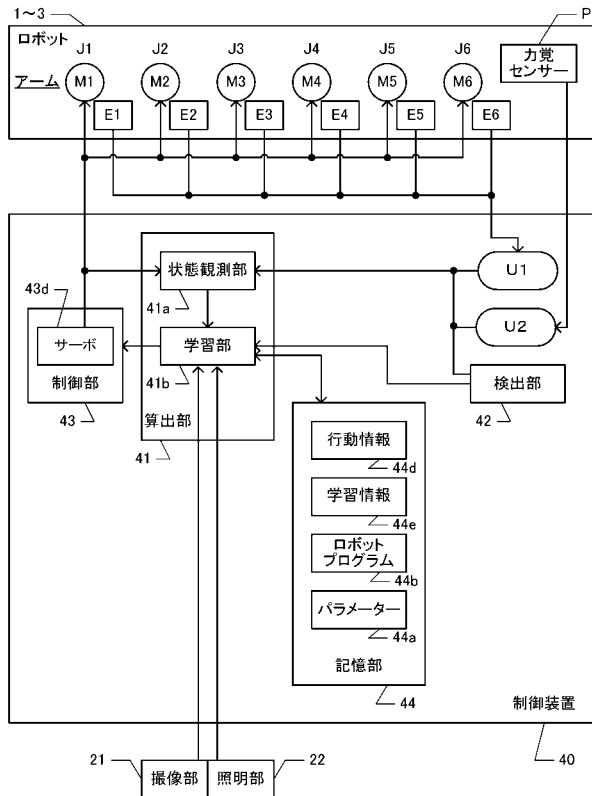
【図 4】



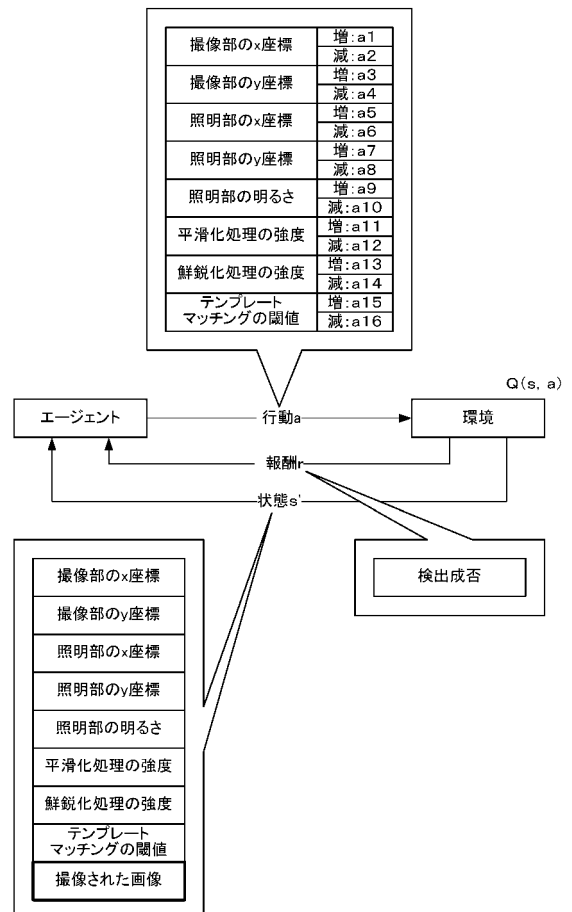
【図 5】



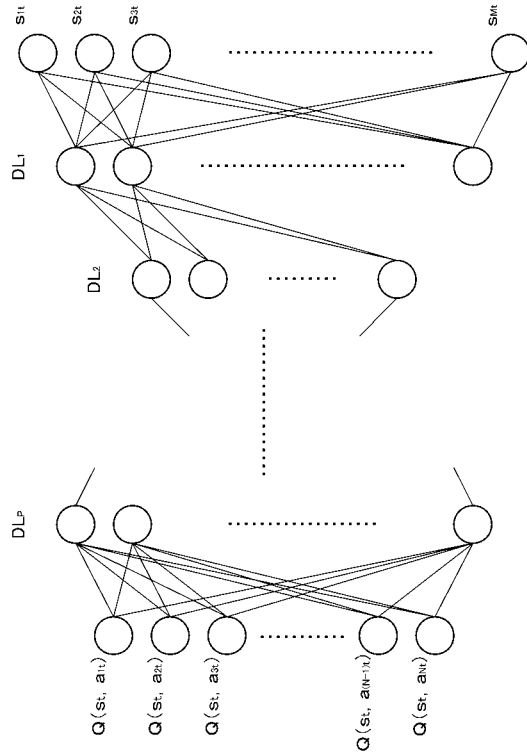
【図 6】



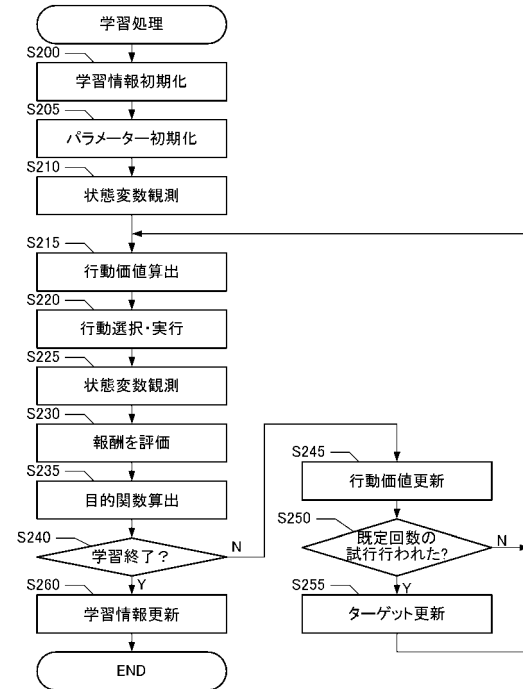
【図 7】



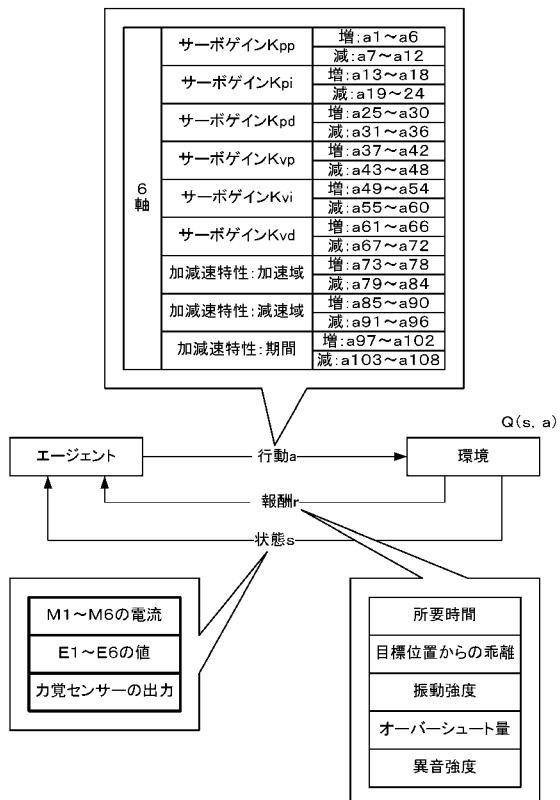
【図 8】



【図 9】



【図 10】



【図 11】

