

(11) 特許出願公開番号

**特開2018-171663**

(P2018-171663A)

(43) 公開日 平成30年11月8日(2018.11.8)

(51) Int. Cl.	F I	テーマコード (参考)
<b>B 2 5 J 13/00 (2006.01)</b>	B 2 5 J 13/00 Z	3 C 2 6 9
<b>G 0 5 B 19/4155 (2006.01)</b>	G 0 5 B 19/4155 V	3 C 7 0 7

審査請求 有 請求項の数 8 O L (全 18 頁)

(21) 出願番号 特願2017-69866 (P2017-69866)  
(22) 出願日 平成29年3月31日 (2017. 3. 31)

(71) 出願人 390008235  
ファナック株式会社  
山梨県南都留郡忍野村忍草字古馬場358  
〇番地

(74) 代理人 100106002  
弁理士 正林 真之

(74) 代理人 100165157  
弁理士 芝 哲央

(74) 代理人 100160794  
弁理士 星野 寛明

(72) 発明者 山本 知之  
山梨県南都留郡忍野村忍草字古馬場358  
〇番地 ファナック株式会社内

[最終頁に続く](#)

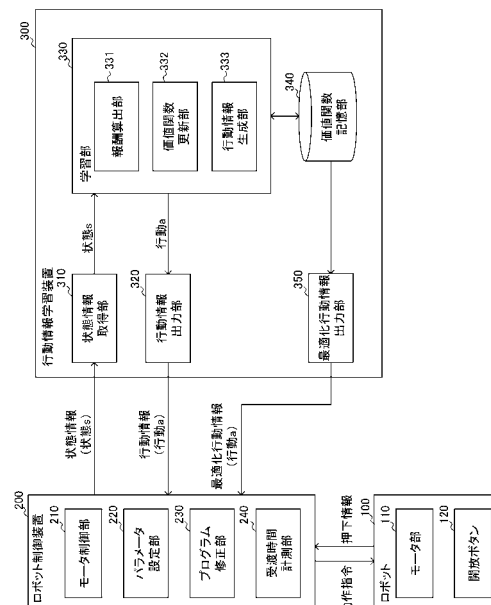
(54) 【発明の名称】 行動情報学習装置、ロボット制御システム及び行動情報学習方法

(57) 【要約】

【課題】作業者がロボットとの協調作業を行いやすくするための行動情報学習装置、ロボット制御システム及び行動情報学習方法を提供する。

【解決手段】行動情報学習装置３００は、ロボット１００がワーク置き場から取得したワーク７を、作業者Ｐに対してワーク７を受け渡す領域である受渡領域８内で受け渡す場合に、ロボット１００の状態ｓを取得する状態情報取得部３１０と、状態ｓの調整情報である行動ａを出力する行動情報出力部３２０と、ワーク７の受け渡しにかかる受渡時間Ｔについての情報である判定情報を取得し、取得した判定情報に基づいて強化学習における報酬の値を算出する報酬算出部３３１と、報酬算出部３３１により算出された報酬の値と、状態ｓと、行動ａとに基づいて強化学習を行うことにより価値関数Ｑを更新する価値関数更新部３３２と、を備える。

【選択図】図2



**【特許請求の範囲】****【請求項 1】**

ロボットがワーク置き場から取得したワークを、作業者に対してワークを受け渡す領域であるワーク受渡領域内で受け渡す場合に、前記ロボットの状態情報を取得する状態情報取得手段と、

前記状態情報の調整情報である行動情報を出力する行動情報出力手段と、

ワークの受け渡しにかかる受渡時間についての情報である判定情報を取得し、取得した前記判定情報に基づいて強化学習における報酬の値を算出する報酬算出手段と、

前記報酬算出手段により算出された前記報酬の値と、前記状態情報と、前記行動情報とに基づいて前記強化学習を行うことにより価値関数を更新する価値関数更新手段と、

を備える行動情報学習装置。

10

**【請求項 2】**

請求項 1 に記載の行動情報学習装置において、

前記状態情報は、前記ロボットの姿勢及び前記ワーク受渡領域内の受渡位置に関する情報を含み、

前記調整情報は、前記状態情報についての調整を行うための情報を含むこと、

を特徴とする行動情報学習装置。

**【請求項 3】**

請求項 2 に記載の行動情報学習装置において、

前記状態情報は、更にワークを取得した位置から前記ワーク受渡領域内への前記ロボットの移動経路を含むこと、

を特徴とする行動情報学習装置。

20

**【請求項 4】**

請求項 1 から請求項 3 までのいずれかに記載の行動情報学習装置において、

前記報酬算出手段は、

前記受渡時間が前回の受渡時間よりも短い場合に、前記報酬の値を正の値とし、

前記受渡時間が前回の受渡時間よりも長い場合に、前記報酬の値を負の値とすること

、

を特徴とする行動情報学習装置。

**【請求項 5】**

30

請求項 1 から請求項 4 までのいずれかに記載の行動情報学習装置において、

前記受渡時間は、ワークを取得してから前記ワーク受渡領域内の位置に移動するまでの移動時間と、ワークを前記ワーク受渡領域内の位置に移動後、前記ワークを作業者が受け取るまでの開放時間とからなり、

前記受渡時間が同じ場合には、前記開放時間が短い場合に、前記移動時間が短い場合より前記報酬の値をより大きな値にすること、

を特徴とする行動情報学習装置。

**【請求項 6】**

請求項 1 から請求項 5 までのいずれかに記載の行動情報学習装置において、

他の行動情報学習装置との間で前記価値関数を共有し、

前記価値関数更新手段が、前記共有した価値関数を更新すること、

を特徴とする行動情報学習装置。

40

**【請求項 7】**

請求項 1 から請求項 6 までのいずれかに記載の行動情報学習装置と、

前記行動情報学習装置に対して通信ネットワークを介して接続され、前記ロボットを制御するロボット制御装置と、

を備えたロボット制御システムであって、

前記行動情報学習装置は、

前記価値関数更新手段により更新された前記価値関数に基づいて、前記ロボットによる前記受渡時間を最短にするための行動情報である最適化行動情報を生成する行動情報生

50

成手段と、

前記行動情報生成手段により生成された前記最適化行動情報を、前記ロボット制御装置に対して出力する行動情報出力手段と、

を備えるロボット制御システム。

【請求項 8】

状態情報取得手段が、ロボットがワーク置き場から取得したワークを、作業者に対してワークを受け渡す領域であるワーク受渡領域内で受け渡す場合に、前記ロボットの状態情報を取得するステップと、

行動情報出力手段が、前記状態情報の調整情報である行動情報を出力するステップと、

報酬算出手段が、ワークの受け渡しにかかる受渡時間についての情報である判定情報を取得し、取得した前記判定情報に基づいて強化学習における報酬の値を算出するステップと、

価値関数更新手段が、算出された前記報酬の値と、前記状態情報と、前記行動情報とに基づいて前記強化学習を行うことにより価値関数を更新するステップと、

を含む行動情報学習方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、行動情報学習装置、ロボット制御システム及び行動情報学習方法に関する。

【背景技術】

【0002】

従来、ロボットが人と共存する空間で作業する場合がある。例えば、ロボットは、予めプログラミングされた位置までワークを運び、作業者にワークを受け渡すと、次のワークを取りに行くという一連の動作を繰り返す場合である。

このように、ロボットと作業者との間で協調作業を行う場合、ロボットは、プログラミングによって定められた位置やタイミング等で作業を行っていた。しかし、このような場合に、受け渡し位置やタイミングによっては、作業者にとって受け取り難い場合や、次の作業を始め難い場合がある。

ここで、ロボットにさせる作業を最適化するための装置が開示されている（例えば、特許文献 1 参照）。

【先行技術文献】

【特許文献】

【0003】

【特許文献 1】特開 2009 - 125920 号公報

【発明の概要】

【発明が解決しようとする課題】

【0004】

ロボットに最適な処理をさせるためには、再度プログラミングを行って受け渡し位置等を教示する必要があるが、試行錯誤して最適な位置等を探し出すのには限度があった。

【0005】

本発明は、係る課題に鑑みなされたものであり、その目的は、作業者がロボットとの協調作業を行いやすくするための行動情報学習装置、ロボット制御システム及び行動情報学習方法を提供することである。

【課題を解決するための手段】

【0006】

（1）本発明の行動情報学習装置（例えば、行動情報学習装置 300）は、ロボット（例えば、ロボット 100）がワーク置き場から取得したワーク（例えば、ワーク 7）を、作業者（例えば、作業者 P）に対してワークを受け渡す領域であるワーク受渡領域（例えば、受渡領域 8）内で受け渡す場合に、前記ロボットの状態情報（例えば、状態 s）を取得する状態情報取得手段（例えば、状態情報取得部 310）と、前記状態情報の調整情報

10

20

30

40

50

である行動情報（例えば、行動 a）を出力する行動情報出力手段（例えば、行動情報出力部 320）と、ワークの受け渡しにかかる受渡時間（例えば、受渡時間 T）についての情報である判定情報を取得し、取得した前記判定情報に基づいて強化学習における報酬の値を算出する報酬算出手段（例えば、報酬算出部 331）と、前記報酬算出手段により算出された前記報酬の値と、前記状態情報と、前記行動情報とに基づいて前記強化学習を行うことにより価値関数（例えば、価値関数 Q）を更新する価値関数更新手段（例えば、価値関数更新部 332）と、を備える。

【0007】

（2）（1）に記載の行動情報学習装置において、前記状態情報は、前記ロボットの姿勢及び前記ワーク受渡領域内の受渡位置に関する情報を含み、前記調整情報は、前記状態情報についての調整を行うための情報を含んでもよい。

10

【0008】

（3）（2）に記載の行動情報学習装置において、前記状態情報は、更にワークを取得した位置から前記ワーク受渡領域内への前記ロボットの移動経路を含んでもよい。

【0009】

（4）（1）から（3）までのいずれかに記載の行動情報学習装置において、前記報酬算出手段は、前記受渡時間が前回の受渡時間よりも短い場合に、前記報酬の値を正の値とし、前記受渡時間が前回の受渡時間よりも長い場合に、前記報酬の値を負の値としてもよい。

【0010】

20

（5）（1）から（4）までのいずれかに記載の行動情報学習装置において、前記受渡時間は、ワークを取得してから前記ワーク受渡領域内の位置に移動するまでの移動時間（例えば、移動時間 T1）と、ワークを前記ワーク受渡領域内の位置に移動後、前記ワークを作業者が受け取るまでの開放時間（例えば、開放時間 T2）とからなり、前記受渡時間が同じ場合には、前記開放時間が短い場合に、前記移動時間が短い場合より前記報酬の値をより大きな値にしてもよい。

【0011】

（6）（1）から（5）までのいずれかに記載の行動情報学習装置において、他の行動情報学習装置との間で前記価値関数を共有し、前記価値関数更新手段が、前記共有した価値関数を更新してもよい。

30

【0012】

（7）本発明によるロボット制御システム（例えば、ロボット制御システム 1000）は、（1）から（5）までのいずれかに記載の行動情報学習装置（例えば、行動情報学習装置 300）と、前記行動情報学習装置に対して通信ネットワーク（例えば、ネットワーク 400）を介して接続され、前記ロボット（例えば、ロボット 100）を制御するロボット制御装置（例えば、ロボット制御装置 200）と、を備え、前記行動情報学習装置が、前記価値関数更新手段（例えば、価値関数更新部 332）により更新された前記価値関数に基づいて、前記ロボットによる前記受渡時間を最短にするための行動情報である最適化行動情報を生成する行動情報生成手段（例えば、最適化行動情報出力部 350）と、前記行動情報生成手段により生成された前記最適化行動情報を、前記ロボット制御装置に対して出力する行動情報出力手段（例えば、最適化行動情報出力部 350）と、を備える。

40

【0013】

（8）本発明による行動情報学習方法は、状態情報取得手段が、ロボットがワーク置き場から取得したワークを、作業者に対してワークを受け渡す領域であるワーク受渡領域内で受け渡す場合に、前記ロボットの状態情報を取得するステップと、行動情報出力手段が、前記状態情報の調整情報である行動情報を出力するステップと、報酬算出手段が、ワークの受け渡しにかかる受渡時間についての情報である判定情報を取得し、取得した前記判定情報に基づいて強化学習における報酬の値を算出するステップと、価値関数更新手段が、算出された前記報酬の値と、前記状態情報と、前記行動情報とに基づいて前記強化学

50

習を行うことにより価値関数を更新するステップと、を含む。

【発明の効果】

【0014】

本発明によれば、作業者がロボットとの協調作業を行いやすくするための行動情報学習装置、ロボット制御システム及び行動情報学習方法を提供できる。

【図面の簡単な説明】

【0015】

【図1】本発明の実施形態全体の基本的構成を示すブロック図である。

【図2】本発明の実施形態に含まれる各装置が備える機能ブロックを表すブロック図である。

【図3】本発明の実施形態におけるロボットの動作を説明するための図である。

【図4】本発明の実施形態における強化学習時の基本的動作を示すフローチャートである。

【図5】本発明の実施形態における最適化行動情報の選択時の基本的動作を示すフローチャートである。

【図6】本発明の変形形態における行動情報学習装置間の連携を示すブロック図である。

【発明を実施するための形態】

【0016】

(実施形態)

まず、本発明の実施形態の概略を説明する。本実施形態において、図1に示すように、ロボット100と、作業者Pとは、作業空間内にて共同で作業をする。そして、ロボット100は、ワーク7(図3参照)を作業者Pに運搬し、作業者Pは、ロボット100からワーク7を受け取って作業をする。そして、ロボット制御システム1000では、ロボット100がワーク7を運搬してから、作業者Pがワーク7を受け取るまでの時間が最短になるように、ロボット100の行動情報を学習する。

【0017】

次に、本実施形態に係るロボット制御システム1000の構成について説明する。ロボット制御システム1000は、ロボット100、ロボット制御装置200、行動情報学習装置300及びネットワーク400を備えている。

ここで、ロボット制御装置200とロボット100とは、1対1の組とされて、通信可能に接続されている。なお、ロボット制御装置200とロボット100とは、接続インターフェースを介して直接接続されても、また、LAN(Local Area Network)等のネットワークを介して接続されてもよい。

【0018】

また、ロボット制御装置200と、行動情報学習装置300とは、それぞれ接続インターフェースを介して直接に接続、又は、それぞれネットワーク400を介して接続されており、相互に通信を行うことが可能である。なお、ネットワーク400は、例えば、工場内に構築されたLANや、インターネット、公衆電話網、或いは、これらの組み合わせである。ネットワーク400における具体的な通信方式や、有線接続及び無線接続のいずれであるか等については、特に限定されない。

【0019】

次に、ロボット制御システム1000に含まれるこれら装置の機能について、図2を参照して説明する。ここで、図2は、各装置に含まれる機能ブロックを表すブロック図である。なお、各装置間に存在するネットワーク400については、その図示を省略する。

【0020】

ロボット100は、ロボット制御装置200に設定されたロボット制御プログラム及びロボット制御装置200に設定されたパラメータの設定値に基づいて生成される動作指令にしたがって、例えば、部品等のワーク7を運搬する。ロボット100は、モータ部110と、開放ボタン120とを備える。

モータ部110は、ロボット100のハンド部13(後述する)等の駆動軸を駆動させ

10

20

30

40

50

るサーボモータである。

開放ボタン１２０は、ハンド部１３に把持したワーク７を取り外す処理を行うためのボタンである。開放ボタン１２０は、作業者Ｐにより操作される。開放ボタン１２０を操作したことによる押下情報は、ロボット制御装置２００に送られる。

#### 【００２１】

ここで、ロボット１００による動作について、図３に基づき説明する。

図３は、本発明の実施形態におけるロボット１００の動作を説明するための図である。

ロボット１００は、例えば、６軸多関節型のロボットである。ロボット１００の各関節部の駆動軸及びハンド部１３の駆動軸は、モータ部１１０によって駆動するが、ロボット制御装置２００によって制御される。

10

ロボット１００は、例えば、ワーク置き場に載置されたワーク７を取得し、作業台上の受渡領域８の所定位置にワーク７を運搬する。このロボット１００がワーク７を取得してから受渡領域８の所定位置までワーク７を運搬するまでの時間を、移動時間Ｔ１とする。

#### 【００２２】

作業者Ｐによるロボット１００への操作、例えば、作業者Ｐがロボット１００のハンド部１３の近傍を掴んで動かす動作をすることによって、ロボット１００は、位置及び姿勢を変える。また、作業者Ｐによる開放ボタン１２０の押下操作によって、ロボット制御装置２００のモータ制御部２１０は、ワーク７をハンド部１３から取り外す制御を行い、作業者Ｐは、ロボット１００からワーク７を受け取る。このロボット１００がワーク７を受渡領域８の所定位置まで運搬してから作業者Ｐがワーク７を受け取るまでの時間を、開放時間Ｔ２とする。そして、移動時間Ｔ１と、開放時間Ｔ２とを合わせた時間を、受渡時間Ｔとする。

20

#### 【００２３】

以上、ロボット１００の機能ブロックについて説明したが、上述した機能ブロックは、本実施形態の動作に特に関連する部分である。ロボット１００は、上述した機能ブロック以外にも、例えば、動作指令を増幅するモータ駆動アンプや、ユーザの操作を受け付けるための操作盤等、一般的な機能ブロックを備えている。しかしながら、これらの一般的な機能ブロックについては、当業者によく知られているので、詳細な説明及び図示を省略する。

#### 【００２４】

30

図２に戻り、ロボット制御装置２００は、ロボット１００を制御することにより、ロボット１００に所定の動作を行わせる装置である。また、ロボット制御装置２００は、行動情報学習装置３００に対して状態情報（「ステータス」ともいう。）を送信する。更に、ロボット制御装置２００は、行動情報学習装置３００から行動情報（「アクション」ともいう。）を受信する。これら各情報の詳細については、行動情報学習装置３００の機能ブロックの説明と併せて説明をする。

#### 【００２５】

ロボット制御装置２００は、モータ制御部２１０と、パラメータ設定部２２０と、プログラム修正部２３０と、受渡時間計測部２４０とを備える。

#### 【００２６】

40

モータ制御部２１０は、ロボット制御プログラム及びパラメータ（例えば、ハンド部１３を含むロボット１００の受渡領域８内での姿勢、受渡領域８内における位置、ワーク７を取得してから受渡領域８内の位置に達するまでのロボット１００の移動経路に関する値）の設定値に基づいて動作指令を生成し、生成した動作指令をロボット１００に送出する。そして、モータ制御部２１０は、ロボット１００に動作指令を送出することにより、ロボット１００のモータ部１１０等の駆動を制御する。この処理により、ロボット１００によるワーク７の運搬動作が実現される。ここで、ロボット制御プログラムには、運搬のための諸条件（例えば、障害物を避けるためのマップ情報、移動速度等）が定義されている。

#### 【００２７】

50

パラメータ設定部 220 は、ロボット 100 の当該ロボット制御プログラムによるワーク 7 の運搬時における移動処理に関するパラメータを設定する部分である。ロボット 100 のパラメータとは、例えば、ハンド部 13 を含むロボット 100 の受渡領域 8 内での姿勢や、受渡領域 8 内の位置や、ロボット 100 の移動経路に関するパラメータである。ハンド部 13 を含むロボット 100 の姿勢を示すパラメータは、例えば、ハンド部 13 の角度に関するデータである。また、受渡領域 8 内の位置を示すパラメータは、例えば、ハンド部 13 の位置を XYZ 座標で表したデータである。ロボット 100 の移動経路に関するパラメータは、例えば、ロボット 100 がワーク 7 を取得したワーク置き場の位置から受渡領域 8 までのロボット 100 の教示点のデータである。ここで、ロボット 100 の教示点とは、ロボット 100 のハンド部 13 の先端部の位置をいう。

10

かかるパラメータの設定値は、行動情報学習装置 300 から出力される行動情報や、最適化行動情報に基づいて調整される。

#### 【0028】

プログラム修正部 230 は、ロボット制御プログラムを直接修正する。具体的には、プログラム修正部 230 は、当該ロボット制御プログラムで記述されたハンド部 13 を含むロボット 100 の受渡領域 8 内での姿勢や、受渡領域 8 内の位置等を、行動情報学習装置 300 から出力される行動情報や、最適化行動情報に基づいて、プログラムコードを直接修正する。

#### 【0029】

受渡時間計測部 240 は、時間を計測する制御部である。受渡時間計測部 240 は、ロボット 100 がワーク置き場に載置されたワーク 7 を取得し、受渡領域 8 の所定位置にワーク 7 を運搬するまでの時間である移動時間 T1 を計測する。また、受渡時間計測部 240 は、受渡領域 8 の所定位置にワーク 7 を運搬してから作業員 P がワーク 7 を受け取るまでの時間である開放時間 T2 を計測する。

20

#### 【0030】

行動情報学習装置 300 は、強化学習を行う装置である。行動情報学習装置 300 に含まれる各機能ブロックの説明に先立って、まずは、強化学習の基本的な仕組みについて説明する。エージェント（本実施形態における行動情報学習装置 300 に相当）は、環境の状態を観測し、ある行動を選択し、当該行動に基づいて環境が変化する。環境の変化に伴って、何らかの報酬が与えられ、エージェントはより良い行動の選択（意思決定）を学習する。

30

教師あり学習が、完全な正解を示すのに対して、強化学習における報酬は、環境の一部の変化に基づく断片的な値であることが多い。

このため、エージェントは、将来にわたっての報酬の合計を最大にするように行動を選択するように学習する。

#### 【0031】

このように、強化学習では、行動を学習することにより、環境に行動が与える相互作用を踏まえて適切な行動を学習、すなわち将来的に得られる報酬を最大にするための学習する方法を学ぶ。これは、本実施形態において、例えば、受渡時間 T を短縮し、更に、開放時間 T2 を短縮するための行動情報を選択するという、未来に影響をおよぼすような行動を獲得できることを表している。

40

#### 【0032】

ここで、強化学習としては、任意の学習方法を用いることができるが、以下の説明では、或る環境の状態  $s$  の下で、行動  $a$  を選択する価値  $Q(s, a)$  を学習する方法である Q 学習 (Q-learning) を用いる場合を例にとって説明をする。

Q 学習では、或る状態  $s$  のとき、取り得る行動  $a$  のなかから、価値  $Q(s, a)$  の最も高い行動  $a$  を最適な行動として選択することを目的とする。

#### 【0033】

しかしながら、Q 学習を最初に開始する時点では、状態  $s$  と行動  $a$  との組み合わせについて、価値  $Q(s, a)$  の正しい値は全く分かっていない。そこで、エージェントは、或

50

る状態  $s$  の下で様々な行動  $a$  を選択し、その時の行動  $a$  に対して、与えられる報酬に基づいて、より良い行動の選択をすることにより、正しい価値  $Q(s, a)$  を学習していく。

【0034】

また、将来にわたって得られる報酬の合計を最大化したいので、最終的に価値  $Q(s, a) = E[(\sum_{t=0}^{\infty} \gamma^t r_t)]$  となるようにすることを目指す。ここで  $E[\ ]$  は期待値を表し、 $t$  は時刻、 $\gamma$  は後述する割引率と呼ばれるパラメータ、 $r_t$  は時刻  $t$  における報酬、 $\sum_{t=0}^{\infty} \gamma^t r_t$  は時刻  $t$  による合計である。この式における期待値は、最適な行動にしたがって状態変化した場合の期待値である。しかし  $Q$  学習の過程において、最適な行動が何であるのかは不明であるので、様々な行動を行うことにより、探索しながら強化学習をする。このような価値  $Q(s, a)$  の更新式は、例えば、次の式 (1) により表すことができる。

10

【0035】

【数 1】

$$Q(s_{t+1}, a_{t+1}) \leftarrow Q(s_t, a_t) + \alpha(r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)) \quad \cdots (1)$$

【0036】

上記の式 (1) において、 $s_t$  は、時刻  $t$  における環境の状態を表し、 $a_t$  は、時刻  $t$  における行動を表す。行動  $a_t$  により、状態は  $s_{t+1}$  に変化する。 $r_{t+1}$  は、その状態の変化により得られる報酬を表している。また、 $\max$  の付いた項は、状態  $s_{t+1}$  の下で、その時に分かっている最も  $Q$  値の高い行動  $a$  を選択した場合の  $Q$  値に  $\alpha$  を乗じたものになる。ここで、 $\alpha$  は、 $0 < \alpha < 1$  のパラメータで、割引率と呼ばれる。また、 $\alpha$  は、学習係数で、 $0 < \alpha < 1$  の範囲とする。

20

【0037】

上述した式 (1) は、試行  $a_t$  の結果、返ってきた報酬  $r_{t+1}$  を元に、状態  $s_t$  における行動  $a_t$  の価値  $Q(s_t, a_t)$  を更新する方法を表している。この更新式は、状態  $s_t$  における行動  $a_t$  の価値  $Q(s_t, a_t)$  よりも、行動  $a_t$  による次の状態  $s_{t+1}$  における最良の行動の価値  $\max_a Q(s_{t+1}, a)$  の方が大きければ、価値  $Q(s_t, a_t)$  を大きくし、逆に小さければ、価値  $Q(s_t, a_t)$  を小さくすることを示している。つまり、或る状態における或る行動の価値を、それによる次の状態における最良の行動の価値に近づける。ただし、その差は、割引率  $\gamma$  と報酬  $r_{t+1}$  のあり方により変わってくるが、基本的には、ある状態における最良の行動の価値が、それに至る一つ前の状態における行動の価値に伝播していく仕組みになっている。

30

【0038】

ここで、 $Q$  学習では、すべての状態行動ペア  $(s, a)$  についての価値  $Q(s, a)$  のテーブルを作成して、学習を行う方法がある。しかし、すべての状態行動ペアの価値  $Q(s, a)$  の値を求めるには状態数が多すぎて、 $Q$  学習が収束するのに多くの時間を要してしまう場合がある。

【0039】

そこで、公知の  $DQN$  (Deep Q - Network) と呼ばれる技術を利用するようにしてもよい。具体的には、価値関数  $Q$  を適当なニューラルネットワークを用いて構成し、ニューラルネットワークのパラメータを調整することにより、価値関数  $Q$  を適当なニューラルネットワークで近似することにより価値  $Q(s, a)$  の値を算出するようにしてもよい。 $DQN$  を利用することにより、 $Q$  学習が収束するのに要する時間を短くすることが可能となる。なお、 $DQN$  については、例えば、以下の非特許文献に詳細な記載がある。

40

【0040】

< 非特許文献 >

「Human - level control through deep reinforcement learning」、Volodymyr Mnih 著 [online]、[平成29年3月17日検索]、インターネット URL : <http://arxiv.org/abs/1312.5602>

50



les.davidquiu.com/research/nature14236.pdf

【0041】

行動情報学習装置300は、上記において説明をしたQ学習を行う。具体的には、行動情報学習装置300は、ロボット100において設定されたロボット制御プログラムの内容及び当該ロボット制御プログラム実行時におけるパラメータの組み合わせを状態sとし、当該状態sに係る当該ロボット制御プログラムの修正及びパラメータの調整を行動aとして、選択する価値関数Qを学習する。

【0042】

行動情報学習装置300は、ロボット100において設定されたロボット制御プログラム及びパラメータ等の状態sを観測して、行動aを決定する。行動情報学習装置300は、行動aをするたびに報酬が返ってくる。行動情報学習装置300は、将来にわたっての報酬の合計が最大になるように、最適な行動aを試行錯誤的に探索する。そうすることで、行動情報学習装置300は、ロボット100において設定されたロボット制御プログラムの内容及び当該ロボット制御プログラム実行時におけるパラメータの組み合わせ等である状態sに対して、最適な行動aを選択することが可能となる。

【0043】

すなわち、行動情報学習装置300により学習された価値関数Qに基づいて、或る状態sに係るロボット制御プログラムの内容及び当該ロボット制御プログラム実行時におけるパラメータの組み合わせに対して、適用される行動aのうち、価値関数Qの値が最大となるような行動aを選択することで、ワーク7の受け渡しに係る時間である受渡時間T及び開放時間T2が最短になるような行動aを選択することが可能となる。

【0044】

以上の強化学習を行うために、行動情報学習装置300は、状態情報取得部310、行動情報出力部320、学習部330及び価値関数記憶部340を備える。

【0045】

状態情報取得部310は、ロボット制御プログラムの内容及び当該ロボット制御プログラム実行時におけるパラメータの組み合わせである状態情報(状態s)を、ロボット制御装置200(及び/又はロボット100)から取得する部分である。この状態sは、Q学習における、環境の状態sに相当する。

【0046】

具体的には、本実施形態における状態sには、ロボット100を制御するためのロボット制御プログラムの内容及び当該ロボット制御プログラム実行時におけるパラメータの組み合わせが含まれ、パラメータには、当該ロボット制御プログラム等によるハンド部13を含むロボット100の受渡領域8内での姿勢、受渡領域8内における位置、ワーク7を取得してから受渡領域8内の位置に達するまでの移動経路に関する情報が含まれる。状態情報取得部310は、取得した状態sを学習部330に対して出力する。

【0047】

また、状態情報取得部310は、Q学習を行うための報酬を算出するための判定情報も取得する。具体的には、状態sに係るロボット制御プログラム及び当該ロボット制御プログラム実行時におけるパラメータの組み合わせを実行後の受渡時間Tを、Q学習を行うための報酬を算出するための判定情報とする。受渡時間Tは、上述したように、ロボット100がワーク7を取得してから受渡領域8内の位置まで移動する時間である移動時間T1と、受渡領域8内の位置に移動してから作業員Pにワーク7が受け渡されるまでの開放時間T2とからなる。

【0048】

行動情報出力部320は、学習部330が生成した行動情報(行動a)を、ロボット制御装置200に対して送信する部分である。ロボット制御装置200は、上述したように、この行動aに基づいて、現在の状態s、すなわち現在設定されているロボット制御プログラム及びパラメータを修正することで、次の状態s'(すなわち修正されたロボット制

10

20

30

40

50

御プログラム、修正されたパラメータ及び当該修正されたロボット制御プログラムと修正されたパラメータによる移動処理を実行した場合の状態)に遷移する。

【0049】

学習部330は、或る環境の状態 $s$ の下で、ある行動 $a$ を選択する場合の価値 $Q(s, a)$ を学習する部分である。具体的には、学習部330は、報酬算出部331、価値関数更新部332及び行動情報生成部333を備える。

【0050】

報酬算出部331は、或る状態 $s$ の下で、行動 $a$ を選択した場合の報酬を、判定情報に基づいて算出する部分である。ここで、本実施形態では、行動 $a$ により修正された状態 $s'$ に係る修正後のロボット制御プログラム及び当該修正後のロボット制御プログラム実行時における修正後のパラメータに基づいて動作したロボット100の受渡時間 $T$ が、行動 $a$ により修正される前の状態 $s$ に係る修正前のロボット制御プログラム及び当該修正前のロボット制御プログラム実行時における修正前のパラメータに基づいて動作したロボット100の受渡時間 $T$ よりも長くなった場合に、報酬の値を負の値とする。

【0051】

一方で、行動 $a$ により修正された状態 $s'$ に係る修正後のロボット制御プログラム及び当該修正後のロボット制御プログラム実行時における修正後のパラメータに基づいて動作したロボット100の受渡時間 $T$ が、行動 $a$ により修正される前の状態 $s$ に係る修正前のロボット制御プログラム及び当該修正前のロボット制御プログラム実行時における修正前のパラメータに基づいて動作したロボット100の受渡時間 $T$ よりも短くなった場合に、報酬の値を正の値とする。

【0052】

また、報酬の値については重みづけを与えるようにすることができる。例えば、移動時間 $T_1$ と、開放時間 $T_2$ とでは、開放時間 $T_2$ が短くなった場合の方が、移動時間 $T_1$ が短くなった場合に比べて、報酬の正の値を大きくすることが好ましい。つまり、開放時間 $T_2$ が短くなった度合いに応じて、正の値が大きくなるようにするとよい。

なお、上記の報酬の値の算出方法は、一例であって、これに限定されない。例えば、状態 $s'$ と状態 $s$ における受渡時間 $T$ の偏差、移動時間 $T_1$ の偏差、及び開放時間 $T_2$ の偏差と、報酬の値を対応付ける報酬対応テーブル(仮称)を予め任意に作成しておき、報酬対応テーブルに基づいて、報酬の値を算出するようにしてもよい。また、受渡時間 $T$ の偏差、移動時間 $T_1$ の偏差、及び開放時間 $T_2$ の偏差を入力とする報酬関数(仮称)を予め任意に作成しておき、報酬関数に基づいて、報酬の値を算出するようにしてもよい。

【0053】

価値関数更新部332は、状態 $s$ と、行動 $a$ と、行動 $a$ を状態 $s$ に適用した場合の状態 $s'$ と、上記のようにして算出された報酬の値とに基づいて $Q$ 学習を行うことにより、価値関数記憶部340が記憶する価値関数 $Q$ を更新する。

【0054】

価値関数 $Q$ の更新は、オンライン学習で行ってもよく、バッチ学習で行ってもよく、ミニバッチ学習で行ってもよい。

オンライン学習とは、或る行動 $a$ を現在の状態 $s$ に適用することにより、状態 $s$ が新たな状態 $s'$ に遷移する都度、即座に価値関数 $Q$ の更新を行うという学習方法である。また、バッチ学習とは、或る行動 $a$ を現在の状態 $s$ に適用することにより、状態 $s$ が新たな状態 $s'$ に遷移することを繰り返すことにより、学習用のデータを収集し、収集したすべての学習用データを用いて、価値関数 $Q$ の更新を行うという学習方法である。更に、ミニバッチ学習とは、オンライン学習と、バッチ学習の中間的な、ある程度学習用データが溜まるたびに価値関数 $Q$ の更新を行うという学習方法である。

【0055】

行動情報生成部333は、 $Q$ 学習の過程において、ロボット100に様々な動作( $Q$ 学習における行動 $a$ に相当する。)を行わせるために、行動 $a$ を生成して、生成した行動 $a$ を行動情報出力部320に対して出力する。

## 【0056】

具体的には、行動情報生成部333は、現在の状態 $s$ に対して、Q学習の過程における行動 $a$ を選択する。本実施形態における行動 $a$ には、現在の状態 $s$ に係るロボット制御プログラムで記述された内容に関する修正情報、及び現在の状態 $s$ に係るパラメータ（例えば、ハンド部13を含むロボット100の姿勢、受渡領域8内における位置、ワーク7を取得してから受渡領域8内の位置に達するまでのロボット100の移動経路に関する値）の設定値が含まれる。

## 【0057】

行動情報生成部333は、例えば、状態 $s$ に含まれるロボット制御プログラム及びパラメータに対して行動 $a$ に含まれるパラメータ（例えば、ハンド部13を含むロボット100の姿勢、受渡領域8内における位置、ワーク7を取得してから受渡領域8内の位置に達するまでのロボット100の移動経路に関する値）の設定値を適用して、状態 $s'$ に遷移して、プラスの報酬（正の値の報酬）が返った場合、次の行動 $a'$ としては、例えば、受渡領域8内の位置を、ワーク7を取り外した位置である開放位置側に少し移動させたり、ロボット100の姿勢を、ワーク7を取り外した姿勢である開放姿勢の方向に少し変化させたりして、受渡時間 $T$ がより短くなるような行動 $a'$ を選択する方策を取るようにしてもよい。

## 【0058】

また、逆に、マイナスの報酬（負の値の報酬）が返った場合、行動情報生成部333は、例えば、状態 $s'$ よりも状態 $s$ 寄りになるように、行動 $a'$ を選択するようにしてもよい。又は、状態 $s'$ 寄りになるような行動 $a'$ を選択することで、マイナスの報酬になると思われる行動を集めるようにしてもよい。

更に、行動情報生成部333は、現在の推定される行動 $a$ の価値の中で、最も価値 $Q(s, a)$ の高い行動 $a'$ を選択するグリーディ法や、ある小さな確率でランダムに行動 $a'$ を選択し、それ以外では最も価値 $Q(s, a)$ の高い行動 $a'$ を選択するグリーディ法といった公知の方法により、行動 $a'$ を選択する方策を取るようにしてもよい。

## 【0059】

価値関数記憶部340は、価値関数 $Q$ を記憶する記憶装置である。価値関数記憶部340に記憶された価値関数 $Q$ は、価値関数更新部332により更新される。

## 【0060】

また、行動情報学習装置300は、価値関数更新部332がQ学習を行うことにより更新した価値関数 $Q$ に基づいて、価値 $Q(s, a)$ が最大となる動作をロボット100に行わせるための行動 $a$ （以下、「最適化行動情報」と呼ぶ。）を生成する。

## 【0061】

行動情報学習装置300は、最適化行動情報出力部350を備えている。

最適化行動情報出力部350は、価値関数記憶部340が記憶している価値関数 $Q$ を取得する。この価値関数 $Q$ は、上述したように、価値関数更新部332がQ学習を行うことにより更新したものである。そして、最適化行動情報出力部350は、価値関数 $Q$ に基づいて、最適化行動情報を生成し、生成した最適化行動情報をロボット制御装置200に対して出力する。この最適化行動情報には、行動情報出力部320がQ学習の過程において出力する行動情報と同様に、修正後のロボット制御プログラム及び当該修正後のロボット制御プログラム実行時における修正後のパラメータが含まれる。

## 【0062】

ロボット制御装置200が、この最適化行動情報に基づいて現在設定されているロボット制御プログラム及びパラメータを修正して、動作指令を生成することにより、ロボット100は、受渡時間 $T$ 及び開放時間 $T_2$ が最短になるように動作することができる。

## 【0063】

以上、ロボット制御装置200や行動情報学習装置300に含まれる機能ブロックについて説明した。

これらの機能ブロックを実現するために、ロボット制御装置200及び行動情報学習装

10

20

30

40

50

置 3 0 0 は、CPU ( C e n t r a l   P r o c e s s i n g   U n i t ) 等の演算処理装置を備える。また、ロボット制御装置 2 0 0 及び行動情報学習装置 3 0 0 は、アプリケーションソフトウェアや OS ( O p e r a t i n g   S y s t e m ) 等の各種の制御用プログラムを格納した HDD ( H a r d   D i s k   D r i v e ) 等の補助記憶装置や、演算処理装置がプログラムを実行する上で一時的に必要とされるデータを格納するための RAM ( R a n d o m   A c c e s s   M e m o r y ) といった主記憶装置も備える。

【 0 0 6 4 】

そして、ロボット制御装置 2 0 0 及び行動情報学習装置 3 0 0 は、演算処理装置が補助記憶装置からアプリケーションソフトウェアや OS を読み込み、読み込んだアプリケーションソフトウェアや OS を主記憶装置に展開させながら、これらのアプリケーションソフトウェアや OS に基づいた演算処理を行う。また、ロボット制御装置 2 0 0 及び行動情報学習装置 3 0 0 は、この演算結果に基づいて、各装置が備える各種のハードウェアを制御する。これにより、本実施形態の機能ブロックは実現される。つまり、本実施形態は、ハードウェアとソフトウェアが協働することにより実現することができる。

10

【 0 0 6 5 】

具体例として、ロボット制御装置 2 0 0 は、一般的なロボット 1 0 0 の制御装置に本実施形態を実現するためのアプリケーションソフトウェアを組み込むことにより実現できる。また、行動情報学習装置 3 0 0 は、一般的なパーソナルコンピュータに、本実施形態を実現するためのアプリケーションソフトウェアを組み込むことにより実現できる。

20

【 0 0 6 6 】

ただし、行動情報学習装置 3 0 0 については、機械学習に伴う演算量が多いため、例えば、パーソナルコンピュータに GPU ( G r a p h i c s   P r o c e s s i n g   U n i t s ) を搭載し、GPGPU ( G e n e r a l - P u r p o s e   c o m p u t i n g   o n   G r a p h i c s   P r o c e s s i n g   U n i t s ) と呼ばれる技術により、GPU を機械学習に伴う演算処理に利用するようにすると、高速処理できるようになるのでよい。更には、より高速な処理を行うために、行動情報学習装置 3 0 0 は、このような GPU を搭載したコンピュータを複数台用いてコンピュータ・クラスターを構築し、このコンピュータ・クラスターに含まれる複数のコンピュータにて並列処理を行うようにしてもよい。

30

【 0 0 6 7 】

次に、図 4 のフローチャートを参照して本実施形態における行動情報学習処理として、Q 学習を行う行動情報学習装置 3 0 0 の動作について説明をする。

【 0 0 6 8 】

まず、ステップ S 1 1 ( 以下、単に「S」という。 ) において、状態情報取得部 3 1 0 がロボット制御装置 2 0 0 から状態情報を取得する。取得した状態情報は、価値関数更新部 3 3 2 や行動情報生成部 3 3 3 に対して出力される。上述したように、この状態情報は、Q 学習における環境の状態 s に相当する情報であり、S 1 1 時点での、ロボット制御プログラムの内容とパラメータの設定値である、ハンド部 1 3 を含むロボット 1 0 0 の受渡領域 8 内での姿勢、受渡領域 8 内の位置、移動経路に関する情報が含まれる。なお、最初に Q 学習を開始する時点でのロボット制御プログラム及びパラメータの設定値は、予めユーザが生成するようにする。つまり、本実施形態では、ユーザが作成したロボット制御プログラム及びパラメータの初期設定値を、強化学習により最適なものに調整する。

40

【 0 0 6 9 】

S 1 2 において、行動情報生成部 3 3 3 が新たな行動情報を生成し、生成した新たな行動情報 ( 行動 a ) を、行動情報出力部 3 2 0 を介してロボット制御装置 2 0 0 に対して出力する。行動情報を受信したロボット制御装置 2 0 0 は、受信した行動情報に基づいて現在の状態 s に係るロボット制御プログラム及びパラメータを修正した状態 s ' により、ロボット 1 0 0 を駆動させて、ワーク 7 の受け渡し処理を行う。上述したように、この行動情報は、Q 学習における行動 a に相当するものである。ここで、行動情報には、例えば、ロボット制御プログラムの修正値と、パラメータの設定値が含まれる点については、上述

50

した通りである。

【0070】

S 1 3において、状態情報取得部310は、新たな状態 $s'$ についての判定情報を取得する。ここで、新たな状態 $s'$ には、状態 $s'$ に係るロボット制御プログラム及びパラメータを含む。また、判定情報は、状態 $s'$ に係る移動処理を行うために要した移動時間 $T_1$ 及び開放時間 $T_2$ からなる受渡時間 $T$ を含む。取得した判定情報は、報酬算出部331に対して出力される。

【0071】

報酬算出部331は、入力された判定情報に基づいて報酬を算出する。そのために、S 1 4において、報酬算出部331は、判定情報に含まれる受渡時間 $T$ が短くなったか否かを判定する。かかる判定は、状態 $s'$ の判定情報に含まれる、状態 $s'$ に係る移動処理を行うために要した受渡時間 $T$ と、状態 $s'$ の前の状態である状態 $s$ の判定情報に含まれる、状態 $s$ に係る移動処理を行うために要した受渡時間 $T$ とを比較することにより行うことができる。受渡時間 $T$ が短くなったのであれば(S 1 4: Yes)、報酬算出部331は、処理をS 1 5に移す。他方、受渡時間 $T$ が長くなったのであれば(S 1 4: NO)、報酬算出部331は、処理をS 1 8に移す。

10

【0072】

S 1 5において、報酬算出部331は、判定情報に含まれる開放時間 $T_2$ が、状態 $s'$ の前の状態である状態 $s$ の判定情報に含まれる、状態 $s$ に係る移動処理を行うために要した開放時間 $T_2$ より短くなったか否かを判定する。開放時間 $T_2$ が短くなったのであれば(S 1 5: Yes)、報酬算出部331は、処理をS 1 6に移す。他方、開放時間 $T_2$ が長くなったのであれば(S 1 5: NO)、報酬算出部331は、処理をS 1 7に移す。

20

【0073】

S 1 6において、報酬算出部331は、報酬を第1の値とする。ここで、第1の値は正の値とする。その後、学習部330は、処理をS 1 9に移す。

S 1 7において、報酬算出部331は、報酬を第2の値とする。ここで、第2の値は正の値とする。また、第2の値は、第1の値より小さいものとする。その後、学習部330は、処理をS 1 9に移す。

S 1 8において、報酬算出部331は、報酬を第3の値とする。ここで、第3の値は負の値とする。

30

なお、第1の値、第2の値及び第3の値については、前回と比較した時間の差の大きさによって、更に重みづけを行うようにしてもよい。

【0074】

S 1 9において、価値関数更新部332は、上述にて算出された報酬の値に基づいて、価値関数記憶部340が記憶している価値関数 $Q$ を更新する。そして、学習部330は、再度S 1 1に戻り、上述した処理を繰り返すことにより、価値関数 $Q$ は、適切な値に収束していく。なお、学習部330は、上述した処理を、所定回数繰り返したことや、所定時間繰り返したことを条件として終了するようにしてもよい。

以上、行動情報学習装置300の動作について説明したが、例えば、S 1 4からS 1 8にかけての報酬の値を算出する処理は、一例であって、これに限定されない。例えば、上述したように、状態 $s'$ と状態 $s$ における受渡時間 $T$ の偏差、移動時間 $T_1$ の偏差、及び開放時間 $T_2$ の偏差と、を予め設定された報酬対応テーブル(仮称)又は報酬関数(仮称)に入力して、報酬の値を算出するようにしてもよい。

40

【0075】

以上、図4を参照して説明した動作により、本実施形態では、受渡時間 $T$ 及び開放時間 $T_2$ を短縮するための行動情報を生成するための価値関数 $Q$ を生成することができる、という効果を奏する。

【0076】

次に、図5のフローチャートを参照して、行動情報学習装置300による最適化行動情報の生成時の動作について説明をする。

50

まず、S 2 1において、行動情報学習装置 3 0 0 の最適化行動情報出力部 3 5 0 は、価値関数記憶部 3 4 0 が記憶している価値関数 Q を取得する。この価値関数 Q は、上述したように価値関数更新部 3 3 2 が Q 学習を行うことにより更新したものである。

【 0 0 7 7 】

S 2 2 において、最適化行動情報出力部 3 5 0 は、この価値関数 Q に基づいて、例えば現在設定されている状態 s において、取り得る行動 a のなかから価値  $Q(s, a)$  の最も高い行動 a を最適な行動として選択することで最適化行動情報を生成し、生成した最適化行動情報をロボット制御装置 2 0 0 に対して出力する。

【 0 0 7 8 】

以上により、ロボット制御装置 2 0 0 が、この最適化行動情報に基づいて現在設定されている状態 s (すなわち、現在設定されているロボット制御プログラム及びパラメータ) を修正して、動作指令を生成する。そして、ロボット制御装置 2 0 0 は、生成した動作指令をロボット 1 0 0 に送ることにより、ロボット 1 0 0 は、受渡時間 T が最短になるように動作することができる、という効果を奏する。

10

【 0 0 7 9 】

また、図 5 を参照して説明した動作により、本実施形態では、行動情報学習装置 3 0 0 は、価値関数 Q に基づいて最適化行動情報を生成し、ロボット制御装置 2 0 0 は、この最適化行動情報に基づいて現在設定されているロボット制御プログラム及びパラメータを修正して、動作指令を生成する。そして、ロボット制御装置 2 0 0 は、生成した動作指令をロボット 1 0 0 に送ることにより、受渡時間 T を短縮して、ロボット 1 0 0 を制御することが可能となる、という効果も奏する。

20

【 0 0 8 0 】

本実施形態では、上述したように、ロボット制御プログラムやパラメータの設定値を調整しながら強化学習を行うことにより、受渡時間 T を短縮することができる。すなわち、本実施形態は、一般的な技術に比べて、有利な効果を奏する。

【 0 0 8 1 】

なお、上記のロボット制御システム 1 0 0 0 に含まれる各装置のそれぞれは、ハードウェア、ソフトウェア又はこれらの組み合わせにより実現することができる。また、上記のロボット制御システム 1 0 0 0 に含まれる各装置のそれぞれの協働により行われる行動情報学習方法も、ハードウェア、ソフトウェア又はこれらの組み合わせにより実現することができる。ここで、ソフトウェアによって実現されとは、コンピュータがプログラムを読み込んで実行することにより実現されることを意味する。

30

【 0 0 8 2 】

プログラムは、様々なタイプの非一時的なコンピュータ可読媒体 (non-transitory computer readable medium) を用いて格納され、コンピュータに供給することができる。非一時的なコンピュータ可読媒体は、様々なタイプの実体のある記録媒体 (tangible storage medium) を含む。非一時的なコンピュータ可読媒体の例は、磁気記録媒体 (例えば、フレキシブルディスク、磁気テープ、ハードディスクドライブ)、光磁気記録媒体 (例えば、光磁気ディスク)、CD-ROM (Read Only Memory)、CD-R、CD-R/W、半導体メモリ (例えば、マスク ROM、PROM (Programmable ROM)、EPROM (Erasable PROM)、フラッシュ ROM、RAM (random access memory)) を含む。また、プログラムは、様々なタイプの一時的なコンピュータ可読媒体 (transitory computer readable medium) によってコンピュータに供給されてもよい。一時的なコンピュータ可読媒体の例は、電気信号、光信号、及び電磁波を含む。一時的なコンピュータ可読媒体は、電線及び光ファイバ等の有線通信路、又は無線通信路を介して、プログラムをコンピュータに供給できる。

40

【 0 0 8 3 】

また、上述した実施形態は、本発明の好適な実施形態ではあるが、上記実施形態のみに

50

本発明の範囲を限定するものではなく、本発明の要旨を逸脱しない範囲において種々の変更を施した形態での実施が可能である。

【0084】

上述した実施形態では、行動情報学習装置300を、ロボット100やロボット制御装置200とは別体の装置により実現することを想定していたが、行動情報学習装置300の機能の一部又は全部を、例えば、ロボット制御装置200により実現するようにしてもよい。

【0085】

上述した実施形態では、行動情報学習装置300を、学習を行う機能と、行動情報を生成する機能とを有するものとしたが、学習を行う機能と、行動情報を生成する機能とを別の装置で行うようにしてもよい。

10

【0086】

上述した実施形態では、行動情報学習装置300が強化学習を行うものを説明した。この点、図6に示すように、ロボット制御システム1000-2が、m個の行動情報学習装置300に対してネットワーク500を介して接続された管理装置600を備えるものとしてもよい。例えば、ロボット制御装置200ごとに行動情報学習装置300を備えた場合には、mは、ロボット制御装置200の数である。

そして、ロボット100と作業員Pの相対的な作業環境が同じ条件（例えば、ロボット100の位置、受渡領域8、ロボット100のハンド部13の移動可能領域等が相対的に同じであること）を満たす場合、複数の行動情報学習装置300-1～300-mに対してネットワーク500を介して管理装置600を接続することで、管理装置600は、各行動情報学習装置300の価値関数Qを集約することができる。そうすることで、価値関数Qは、すべての行動情報学習装置300との間で共有される。価値関数Qを複数の行動情報学習装置300で共有するようにすれば、各行動情報学習装置300にて分散して強化学習を行うことが可能となるので、強化学習の効率を向上させることが可能となる。

20

【0087】

そして、管理装置600が、集約した価値関数Qを、各行動情報学習装置300に対して送信するようにしてもよい。

なお、管理装置600は、各行動情報学習装置300から学習用のデータを収集し、価値関数Qを更新するようにしてもよい。

30

また、管理装置600が、最適化行動情報を、各ロボット制御装置200に対して出力するようにしてもよい。

【符号の説明】

【0088】

7          ワーク  
8          受渡領域  
13        ハンド部  
100       ロボット  
110       モータ部  
120       開放ボタン  
200       ロボット制御装置  
210       モータ制御部  
220       パラメータ設定部  
300       行動情報学習装置  
310       状態情報取得部  
320       行動情報出力部  
330       学習部  
331       報酬算出部  
332       価値関数更新部  
333       行動情報生成部

40

50

3 4 0 価値関数記憶部

3 5 0 最適化行動情報出力部

4 0 0 , 5 0 0 ネットワーク

1 0 0 0 ロボット制御システム

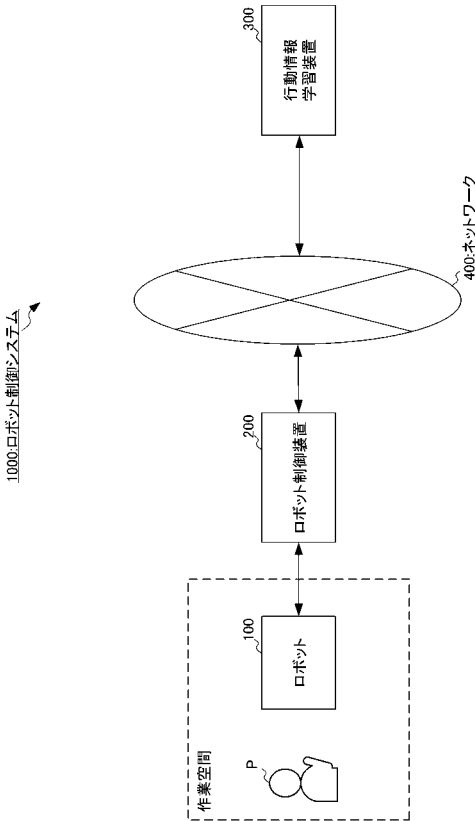
P 作業者

T 受渡時間

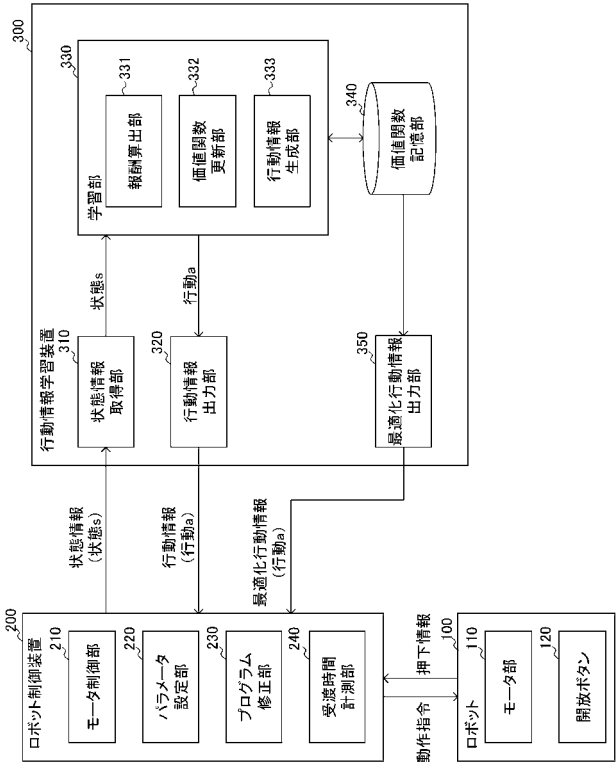
T 1 移動時間

T 2 開放時間

【 図 1 】

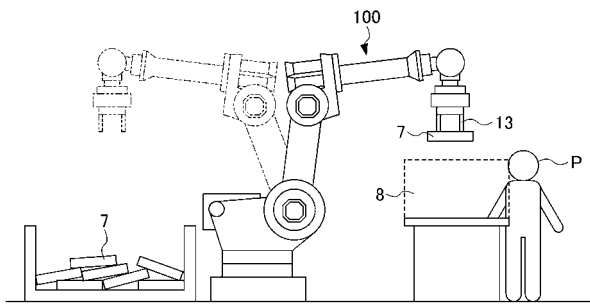


【 図 2 】

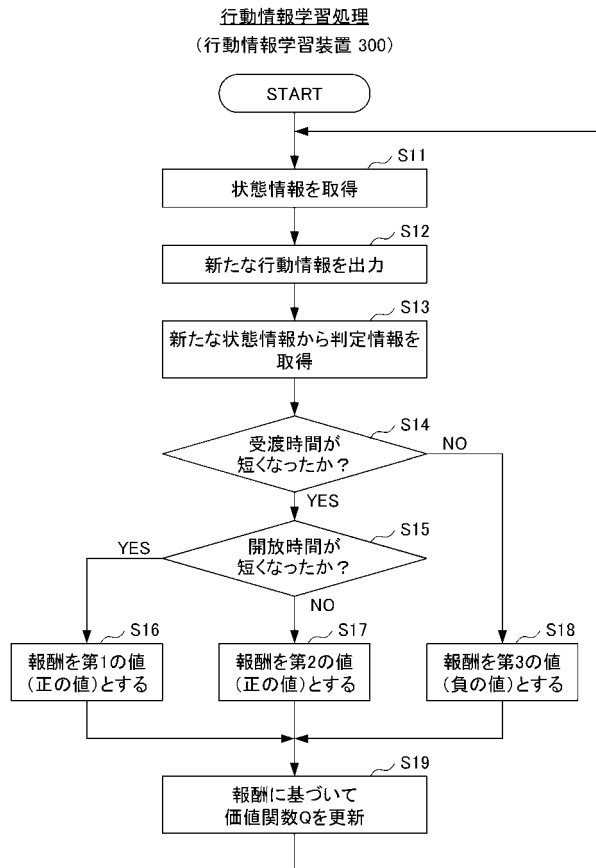




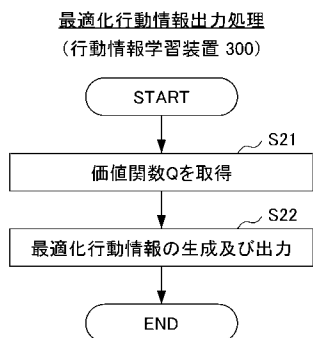
【図3】



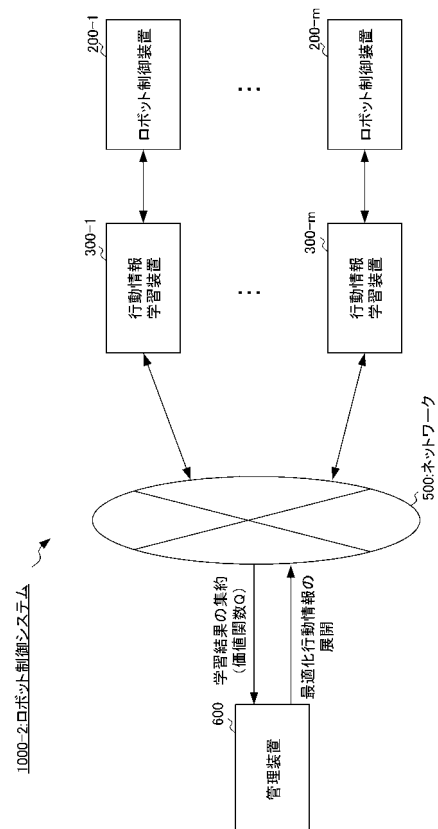
【図4】



【図5】



【図6】



---

フロントページの続き

(72)発明者 栗原 佑典

山梨県南都留郡忍野村忍草字古馬場 3 5 8 0 番地 ファナック株式会社内

F ターム(参考) 3C269 AB33 BB08 EF60 MN27 MN44

3C707 AS01 BS12 KS38 LS14 LV01 LW08 LW12