

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号  
特許第7259978号  
(P7259978)

(45)発行日 令和5年4月18日(2023.4.18)

(24)登録日 令和5年4月10日(2023.4.10)

(51)国際特許分類	F I
H 0 4 L 47/10 (2022.01)	H 0 4 L 47/10
H 0 4 L 41/16 (2022.01)	H 0 4 L 41/16
G 0 6 N 20/00 (2019.01)	G 0 6 N 20/00

請求項の数 10 (全23頁)

(21)出願番号	特願2021-550731(P2021-550731)	(73)特許権者	000004237 日本電気株式会社 東京都港区芝五丁目7番1号
(86)(22)出願日	令和1年9月30日(2019.9.30)	(74)代理人	100141519 弁理士 梶田 邦之
(86)国際出願番号	PCT/JP2019/038454	(72)発明者	沢辺 亜南 東京都港区芝五丁目7番1号 日本電気株式会社内
(87)国際公開番号	WO2021/064766	(72)発明者	岩井 孝法 東京都港区芝五丁目7番1号 日本電気株式会社内
(87)国際公開日	令和3年4月8日(2021.4.8)	審査官	鈴木 肇
審査請求日	令和4年3月28日(2022.3.28)		

最終頁に続く

(54)【発明の名称】 制御装置、方法及びシステム

(57)【特許請求の範囲】

【請求項1】

ネットワークを制御するための行動を学習する、学習手段と、  
前記学習手段が生成した学習情報を記憶する、記憶手段と、を備え、  
前記学習手段は、  
前記ネットワークに対して行われた行動の報酬を、前記行動が行われた後のネットワークを介して提供されたアプリケーションの定常性に基づき定める、制御装置。

【請求項2】

ネットワークを制御するための行動を学習するステップと、  
前記学習により生成された学習情報を記憶するステップと、  
を含み、  
前記学習するステップは、  
前記ネットワークに対して行われた行動の報酬を、前記行動が行われた後のネットワークを介して提供されたアプリケーションの定常性に基づき定める、方法。

【請求項3】

前記学習するステップは、  
前記行動が行われた後のネットワークを介して提供された前記アプリケーションが定常状態であれば、前記ネットワークに対して行われた行動に正の報酬を与え、  
前記行動が行われた後のネットワークを介して提供された前記アプリケーションが非定常状態であれば、前記ネットワークに対して行われた行動に負の報酬を与える、請求項2

に記載の方法。

【請求項 4】

前記学習するステップは、

前記ネットワークに対して行動を起こしたことにより変動するネットワークを介して提供された前記アプリケーションの状態に関する時系列データに基づいて前記ネットワークを介して提供された前記アプリケーションの定常性を判定する、請求項 2 又は 3 に記載の方法。

【請求項 5】

前記学習するステップは、前記ネットワークを介して提供された前記アプリケーションの状態を、前記ネットワークに流れるトラフィックを特徴付ける特徴量、ユーザ体感品質及び制御品質のうち少なくとも 1 つから推定する、請求項 4 に記載の方法。

10

【請求項 6】

前記学習するステップにより生成された学習モデルから得られる行動に基づき、前記ネットワークを制御するステップをさらに含む、請求項 2 乃至 5 のいずれか一項に記載の方法。

【請求項 7】

ネットワークを制御するための行動を学習する、学習手段と、

前記学習手段が生成した学習情報を記憶する、記憶手段と、を含み、

前記学習手段は、

前記ネットワークに対して行われた行動の報酬を、前記行動が行われた後のネットワークを介して提供されたアプリケーションの定常性に基づき定める、システム。

20

【請求項 8】

前記学習手段は、

前記行動が行われた後のネットワークを介して提供された前記アプリケーションが定常状態であれば、前記ネットワークに対して行われた行動に正の報酬を与え、

前記行動が行われた後のネットワークを介して提供された前記アプリケーションが非定常状態であれば、前記ネットワークに対して行われた行動に負の報酬を与える、請求項 7 に記載のシステム。

【請求項 9】

前記学習手段は、

前記ネットワークに対して行動を起こしたことにより変動するネットワークを介して提供された前記アプリケーションの状態に関する時系列データに基づいて前記ネットワークを介して提供された前記アプリケーションの定常性を判定する、請求項 7 又は 8 に記載のシステム。

30

【請求項 10】

前記学習手段は、前記ネットワークを介して提供された前記アプリケーションの状態を、前記ネットワークに流れるトラフィックを特徴付ける特徴量、ユーザ体感品質及び制御品質のうち少なくとも 1 つから推定する、請求項 9 に記載のシステム。

【発明の詳細な説明】

【技術分野】

40

【0001】

本発明は、制御装置、方法及びシステムに関する。

【背景技術】

【0002】

通信技術、情報処理技術の進展に伴い様々なサービスがネットワーク上にて提供される状況にある。例えば、ネットワーク上のサーバから動画データが配信され、端末にて当該動画データを再生することや、サーバから工場等に設置されたロボット等を遠隔制御することが行われている。

【0003】

上記のようなネットワーク上で提供されるサービス、アプリケーションにおいて、エン

50

ドユーザが感じ取る品質（Q o E ; Quality of Experience）や制御品質（Q o C ; Quality of Control）を高める取り組みがなされている。

【0004】

例えば、特許文献1には、個別のwebページの影響が除去された表示待ち時間の品質の推定を可能とする、と記載されている。特許文献1に記載された技術では、任意のエリア及び時間帯におけるトラフィック計測データに基づき当該エリア及び時間帯におけるwebページの表示待ち時間の品質を推定している。

【先行技術文献】

【特許文献】

【0005】

【文献】特開2019-075030号公報

【発明の概要】

【発明が解決しようとする課題】

【0006】

上記特許文献1に開示された技術では、SVM（Support Vector Machine）と称される機械学習が用いられている。ここで、近年、深層学習（ディープラーニング）に代表される機械学習に関する技術が進展し、種々の分野への機械学習の適用が検討されている。

【0007】

例えば、チェス等のゲームやロボット等の制御に機械学習を適用することが検討されている。ゲームの運用に機械学習を適用する場合には、ゲーム内のスコアの最大化が報酬に設定され、機械学習の性能が評価される。また、ロボットの制御では、目標動作の実現が報酬に設定され、機械学習の性能が評価される。通常、機械学習（強化学習）では、即時報酬及びエピソード単位の報酬の総和により学習の性能が議論される。

【0008】

しかし、ネットワークの制御に機械学習を適用する場合には何を報酬に設定するのが問題となる。例えば、ネットワークの制御では、ゲームに機械学習を適用する場合のように最大化するスコアが存在を観念することができない。例えば、ネットワークに含まれる通信機器におけるスループットを最大化することを報酬に設定したとしてもサービス、アプリケーションによっては適切な設定とはいえない。

【0009】

本発明は、機械学習を用いた効率的なネットワークの制御を実現することに寄与する、制御装置、方法及びシステムを提供することを主たる目的とする。

【課題を解決するための手段】

【0010】

本発明の第1の視点によれば、ネットワークを制御するための行動を学習する、学習部と、前記学習部が生成した学習情報を記憶する、記憶部と、を備え、前記学習部は、前記ネットワークに対して行われた行動の報酬を、前記行動が行われた後のネットワークの定常性に基づき定める、制御装置が提供される。

【0011】

本発明の第2の視点によれば、ネットワークを制御するための行動を学習するステップと、前記学習により生成された学習情報を記憶するステップと、を含み、前記学習するステップは、前記ネットワークに対して行われた行動の報酬を、前記行動が行われた後のネットワークの定常性に基づき定める、方法が提供される。

【0012】

本発明の第3の視点によれば、ネットワークを制御するための行動を学習する、学習手段と、前記学習手段が生成した学習情報を記憶する、記憶手段と、を含み、前記学習手段は、前記ネットワークに対して行われた行動の報酬を、前記行動が行われた後のネットワークの定常性に基づき定める、システムが提供される。

【発明の効果】

【0013】

10

20

30

40

50

本発明の各視点によれば、機械学習を用いた効率的なネットワークの制御を実現することに寄与する、制御装置、方法及びシステムが提供される。なお、本発明により、当該効果の代わりに、又は当該効果と共に、他の効果が奏されてもよい。

【図面の簡単な説明】

【0014】

【図1】一実施形態の概要を説明するための図である。

【図2】一実施形態に係る制御装置の動作の一例を示すフローチャートである。

【図3】第1の実施形態に係る通信ネットワークシステムの概略構成の一例を示す図である。

【図4】Qテーブルの一例を示す図である。

10

【図5】ニューラルネットワークの構成の一例を示す図である。

【図6】強化学習により得られる重みの一例を示す図である。

【図7】第1の実施形態に係る制御装置の処理構成の一例を示す図である。

【図8】特徴量とネットワークの状態を対応付ける情報の一例を示す図である。

【図9】行動と制御内容を対応付けたテーブル情報の一例を示す図である。

【図10】特徴量の時系列データの一例を示す図である。

【図11】第1の実施形態に係る制御装置の制御モード時の動作の一例を示すフローチャートである。

【図12】第1の実施形態に係る制御装置の学習モード時の動作の一例を示すフローチャートである。

20

【図13】強化学習実行部の動作を説明するための図である。

【図14】スループットの時系列データの一例を示す図である。

【図15】報酬の与え方を説明するための図である。

【図16】制御装置のハードウェア構成の一例を示す図である。

【発明を実施するための形態】

【0015】

はじめに、一実施形態の概要について説明する。なお、この概要に付記した図面参照符号は、理解を助けるための一例として各要素に便宜上付記したものであり、この概要の記載はなんらの限定を意図するものではない。なお、本明細書及び図面において、同様に説明されることが可能な要素については、同一の符号を付することにより重複説明が省略され得る。

30

【0016】

一実施形態に係る制御装置100は、学習部101と記憶部102を含む(図1参照)。学習部101は、ネットワークを制御するための行動を学習する。記憶部102は、学習部101が生成した学習情報を記憶する。学習部101は、ネットワークに対して行動をする(図2のステップS01)。学習部101は、ネットワークに対して行われた行動の報酬を、行動が行われた後のネットワークの定常性に基づき定め、ネットワークを制御するための行動を学習する(図2のステップS02)。

【0017】

ネットワークにより提供されるサービスやアプリケーションでは、「ネットワークの安定性」が重要視される。制御装置100は、ネットワークに対して行った行動(制御パラメータの変更)により得られる状態の定常性に基づき報酬を定める。即ち、制御装置100は、機械学習(強化学習)の際にネットワークの状態が安定している収束状態に価値が高いものと捉え、そのような状況の場合に高い報酬を与えネットワークを制御するための学習を行う。その結果、機械学習を用いた効率的なネットワークの制御が実現される。

40

【0018】

以下に具体的な実施形態について、図面を参照してさらに詳しく説明する。

【0019】

[第1の実施形態]

第1の実施形態について、図面を用いてより詳細に説明する。

50

## 【 0 0 2 0 】

図 3 は、第 1 の実施形態に係る通信ネットワークシステムの概略構成の一例を示す図である。図 3 を参照すると、通信ネットワークシステムは、端末 1 0 と、制御装置 2 0 と、サーバ 3 0 と、を含んで構成される。

## 【 0 0 2 1 】

端末 1 0 は、通信機能を有する装置である。端末 1 0 には、WEB（ウェブ）カメラ、監視カメラ、ドローン、スマートフォン、ロボット等が例示される。但し、端末 1 0 を上記 WEB カメラ等に限定する趣旨ではない。端末 1 0 は、通信機能を備える任意の装置とすることができる。

## 【 0 0 2 2 】

端末 1 0 は、制御装置 2 0 を介してサーバ 3 0 と通信する。端末 1 0 とサーバ 3 0 により様々なアプリケーション、サービスが提供される。

## 【 0 0 2 3 】

例えば、端末 1 0 が WEB カメラの場合には、サーバ 3 0 が当該 WEB カメラからの画像データを解析し、工場等の資材管理が行われる。例えば、端末 1 0 がドローンの場合には、サーバ 3 0 からドローンに制御コマンドが送信され、ドローンが荷物等を搬送する。例えば、端末 1 0 がスマートフォンの場合には、サーバ 3 0 からスマートフォンに向けて動画が配信され、ユーザはスマートフォンを用いて動画を視聴する。

## 【 0 0 2 4 】

制御装置 2 0 は、例えば、プロキシサーバやゲートウェイ等の通信機器であり、端末 1 0 とサーバ 3 0 からなるネットワークを制御する装置である。制御装置 2 0 は、TCP（Transmission Control Protocol）のパラメータ群やバッファ制御に関するパラメータ群の値を変更し、ネットワークを制御する。

## 【 0 0 2 5 】

例えば、TCP パラメータの制御としては、フローウィンドウサイズの変更が例示される。バッファ制御としては、複数バッファのキュー管理において、最低保証帯域、RED（Random Early Detection）のロス率、ロス開始キュー長、バッファ長に関するパラメータの変更が例示される。

## 【 0 0 2 6 】

なお、以降の説明において、上記 TCP パラメータやバッファ制御に関するパラメータ等、端末 1 0 とサーバ 3 0 の間の通信（トラフィック）に影響を与えるパラメータを「制御パラメータ」と表記する。

## 【 0 0 2 7 】

制御装置 2 0 は、制御パラメータを変更することで、ネットワークを制御する。制御装置 2 0 によるネットワークの制御は、自装置（制御装置 2 0）の packets 転送時に行われてもよいし、端末 1 0 やサーバ 3 0 に制御パラメータの変更を指示することにより行われてもよい。

## 【 0 0 2 8 】

TCP セッションが制御装置 2 0 により終端される場合には、例えば、制御装置 2 0 は、端末 1 0 との間で形成される TCP セッションのフローウィンドウサイズを変更することで、ネットワークを制御する。制御装置 2 0 は、サーバ 3 0 から受信した packets を格納するバッファのサイズを変更したり、当該バッファから packets を読み出す周期を変更したりしてネットワークを制御してもよい。

## 【 0 0 2 9 】

制御装置 2 0 は、ネットワークの制御に「機械学習」を用いる。より具体的には、制御装置 2 0 は、強化学習により得られる学習モデルに基づきネットワークを制御する。

## 【 0 0 3 0 】

強化学習には、種々のバリエーションが存在するが、例えば、制御装置 2 0 は、Q 学習と称される強化学習の結果得られる学習情報（Q テーブル）に基づきネットワークを制御してもよい。

10

20

30

40

50

【 0 0 3 1 】

[ Q 学習 ]

以下、Q 学習について概説する。

【 0 0 3 2 】

Q 学習では、与えられた「環境」における「価値」を最大化するように、「エージェント」を学習させる。当該Q 学習をネットワークシステムに適用すると、端末10やサーバ30を含むネットワークが「環境」であり、ネットワークの状態を最良にするように、制御装置20を学習させる。

【 0 0 3 3 】

Q 学習では、状態（ステート） $s$ 、行動（アクション） $a$ 、報酬（リワード） $r$ の3要素が定義される。

10

【 0 0 3 4 】

状態  $s$  は、環境（ネットワーク）がどのような状態にあるかを示す。例えば、通信ネットワークシステムの場合には、トラヒック（例えば、スループット、平均パケット到着間隔等）が状態  $s$  に該当する。

【 0 0 3 5 】

行動  $a$  は、エージェント（制御装置20）が環境（ネットワーク）に対して取り得る行動を示す。例えば、通信ネットワークシステムの場合には、TCPパラメータ群の設定の変更や機能のオン/オフ等が行動  $a$  として例示される。

【 0 0 3 6 】

報酬  $r$  は、ある状態  $s$  においてエージェント（制御装置20）が行動  $a$  を実行した結果、どの程度の評価が得られるかを示す。例えば、通信ネットワークシステムの場合には、制御装置20が、TCPパラメータ群の一部を変更した結果、スループットが上昇すれば正の報酬、スループットが下降すれば負の報酬の様に定められる。

20

【 0 0 3 7 】

Q 学習では、現時点で得られる報酬（即時報酬）を最大化するのではなく、将来に亘る価値を最大化するように学習が進められる（Qテーブルが構築される）。Q 学習におけるエージェントの学習は、ある状態  $s$  における行動  $a$  を採用した時の価値（Q値、状態行動価値）を最大化するように行われる。

【 0 0 3 8 】

Q値（状態行動価値）は、 $Q(s, a)$ と表記される。Q 学習では、エージェントが行動することによって価値の高い状態に遷移させる行動は、遷移先と同程度の価値を持つことを前提としている。このような前提により、現時点  $t$  におけるQ値は、次の時点  $t + 1$  のQ値により表現することができる（式（1）参照）。

30

【 0 0 3 9 】

【数1】

$$Q(s_t, a_t) = E_{s_{t+1}} \left( r_{t+1} + \gamma E_{a_{t+1}} (Q(s_{t+1}, a_{t+1})) \right) \quad \dots (1)$$

40

【 0 0 4 0 】

なお、式（1）において  $r_{t+1}$  は即時報酬、 $E_{s_{t+1}}$  は状態  $s_{t+1}$  に関する期待値、 $E_{a_{t+1}}$  は行動  $a_{t+1}$  に関する期待値を示す。  $\gamma$  は割引率である。

【 0 0 4 1 】

Q 学習では、ある状態  $s$  において行動  $a$  を採用した結果によりQ値を更新する。具体的には、下記の式（2）に従いQ値を更新する。

【 0 0 4 2 】

50

【数 2】

$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha \left( r_{t+1} + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}) \right) \dots (2)$$

10

【0043】

式(2)において、 $\alpha$  は学習率と称されるパラメータであり、Q値の更新を制御する。また、式(2)における「max」は状態  $S_{t+1}$  の取り得る行動  $a$  のうち最大値を出力する関数である。なお、エージェント(制御装置20)が行動  $a$  を選択する方式には、greedy と称される方式を採用することができる。

【0044】

greedy方式では、確率  $\alpha$  でランダムに行動を選択し、確率  $1 - \alpha$  で最も価値の高い行動を選択する。Q学習の実行により、図4に示すようなQテーブルが生成される。

【0045】

[DQNによる学習]

制御装置20は、DQN(Deep Q Network)と称される深層学習(ディープラーニング)を使った強化学習の結果得られる学習モデルに基づきネットワークを制御してもよい。Q学習では、Qテーブルにより行動価値関数を表現しているが、DQNでは、ディープラーニングにより行動価値関数を表現する。DQNでは、最適行動価値関数を、ニューラルネットワークを使った近似関数により算出する。

20

【0046】

なお、最適行動価値関数とは、ある状態  $s$  時にある行動  $a$  を行うことの価値を出力する関数である。

【0047】

ニューラルネットワークは、入力層、中間層(隠れ層)、出力層を備える。入力層は、状態  $s$  を入力する。中間層の各ノードのリンクには、対応する重みが存在する。出力層は、行動  $a$  の価値を出力する。

30

【0048】

例えば、図5に示すようなニューラルネットワークの構成を考える。図5に示すニューラルネットワークを通信ネットワークシステムに適用すると、入力層のノードは、ネットワークの状態  $S_1 \sim S_3$  に相当する。入力層に入力されたネットワークの状態は、中間層にて重み付けされ、出力層に出力される。

【0049】

出力層のノードは、制御装置20が取り得る行動  $A_1 \sim A_3$  に相当する。出力層のノードは、行動  $A_1 \sim A_3$  のそれぞれに対応する行動価値関数  $Q(s_t, a_t)$  の値を出力する。

40

【0050】

DQNでは、上記行動価値関数を出力するノード間の結合パラメータ(重み)を学習する。具体的には、下記の式(3)に示す誤差関数  $E(s_t, a_t)$  を設定しバックプロパゲーションにより学習を行う。

【0051】

【数 3】

50

$$E(s_t, a_t) = \left( r_{t+1} + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t) \right)^2 \dots (3)$$

## 【 0 0 5 2 】

DQNによる強化学習の実行により、用意されたニューラルネットワークの中間層の構成に対応した学習情報（重み）が生成される（図6参照）。

10

## 【 0 0 5 3 】

ここで、制御装置20の動作モードには、2つの動作モードが含まれる。

## 【 0 0 5 4 】

第1の動作モードは、学習モデルを算出する学習モードである。制御装置20が「Q学習」を実行することで、図4に示すようなQテーブルが算出される。あるいは、制御装置20が「DQN」による強化学習を実行することで、図6に示すような重みが算出される。

## 【 0 0 5 5 】

第2の動作モードは、学習モードにて算出された学習モデルを用いてネットワークを制御する制御モードである。具体的には、制御モードの制御装置20は、現在のネットワークの状態sを算出し、当該状態sの場合に取り得る行動aのうち最も価値の高い行動aを選択する。制御装置20は、当該選択された行動aに対応する動作（ネットワークの制御）を実行する。

20

## 【 0 0 5 6 】

図7は、第1の実施形態に係る制御装置20の処理構成（処理モジュール）の一例を示す図である。図7を参照すると、制御装置20は、パケット転送部201と、特徴量算出部202と、ネットワーク制御部203と、強化学習実行部204と、記憶部205と、を含んで構成される。

## 【 0 0 5 7 】

パケット転送部201は、端末10やサーバ30から送信されたパケットを受信し、当該受信したパケットを対向する装置に転送する手段である。パケット転送部201は、ネットワーク制御部203からの通知された制御パラメータに従い、パケット転送を行う。

30

## 【 0 0 5 8 】

例えば、ネットワーク制御部203からフローウィンドウサイズの設定値が通知されると、パケット転送部201は当該通知されたフローウィンドウサイズにてパケット転送を行う。

## 【 0 0 5 9 】

パケット転送部201は、受信したパケットの複製を特徴量算出部202に引き渡す。

## 【 0 0 6 0 】

特徴量算出部202は、端末10とサーバ30の間の通信トラフィックを特徴付ける特徴量を算出する手段である。特徴量算出部202は、取得したパケットからネットワーク制御の対象となるトラフィックフローを抽出する。なお、ネットワーク制御の対象となるトラフィックフローは、送信元IP（Internet Protocol）アドレス、宛先IPアドレス、ポート番号等が同一のパケットからなるグループである。

40

## 【 0 0 6 1 】

特徴量算出部202は、抽出したトラフィックフローから上記特徴量を算出する。例えば、特徴量算出部202は、スループット、平均パケット到着間隔、パケットロス率、ジッター等を特徴量として算出する。特徴量算出部202は、算出した特徴量を算出時刻と共に記憶部205に格納する。なお、スループット等の算出については既存の技術を用いることができ、且つ、当業者にとって明らかであるのでその詳細な説明を省略する。

## 【 0 0 6 2 】

50

ネットワーク制御部 203 は、強化学習実行部 204 が生成した学習モデルから得られる行動に基づき、ネットワークを制御する手段である。ネットワーク制御部 203 は、強化学習の結果得られる学習モデルに基づきパケット転送部 201 に通知する制御パラメータを決定する。ネットワーク制御部 203 は、主に制御モード時に動作するモジュールである。

【0063】

ネットワーク制御部 203 は、記憶部 205 から最新の（現在時刻の）特徴量を読み出す。ネットワーク制御部 203 は、当該読み出した特徴量から制御対象となっているネットワークの状態を推定（算出）する。

【0064】

例えば、ネットワーク制御部 203 は、特徴量 F とネットワークの状態を対応付けたテーブル（図 8 参照）を参照し、現在の特徴量 F に対応するネットワークの状態を算出する。なお、トラヒックは端末 10 とサーバ 30 の間の通信により生じるものであるから、ネットワークの状態は「トラヒックの状態」と捉えることもできる。即ち、本願開示において、「トラヒックの状態」と「ネットワークの状態」は相互に読み替えが可能である。

【0065】

Q 学習により学習モデルが構築された場合には、ネットワーク制御部 203 は、記憶部 205 に格納された Q テーブルを参照し、現在のネットワーク状態に対応する各行動（アクション）のうち価値 Q が最も高い行動を取得する。例えば、図 4 の例では、算出されたトラヒックの状態が「状態 S1」であり、価値  $Q(S1, A1)$ 、 $Q(S1, A2)$ 、 $Q(S1, A3)$  のうち価値  $Q(S1, A1)$  が最大であれば、行動 A1 が読み出される。

【0066】

あるいは、DNQ により学習モデルが構築された場合には、ネットワーク制御部 203 は、図 5 に示すようなニューラルネットワークに現在のネットワーク状態を入力し、取り得る行動のうち最も価値の高い行動を取得する。

【0067】

ネットワーク制御部 203 は、取得した行動に応じて制御パラメータを決定し、パケット転送部 201 に設定（通知）する。なお、記憶部 205 には、行動と制御内容を対応付けたテーブル（図 9 参照）が格納され、ネットワーク制御部 203 は、当該テーブルを参照してパケット転送部 201 に設定する制御パラメータを決定する。

【0068】

例えば、図 9 に示すように、制御パラメータの変更内容（更新内容）が制御内容として記載されている場合には、ネットワーク制御部 203 は、当該変更内容に応じた制御パラメータをパケット転送部 201 に通知する。

【0069】

強化学習実行部 204 は、ネットワークを制御するための行動（制御パラメータ）を学習する手段である。強化学習実行部 204 は、上記説明した Q 学習や DQN による強化学習を実行し、学習モデルを生成する。強化学習実行部 204 は、主に学習モード時に動作するモジュールである。

【0070】

強化学習実行部 204 は、記憶部 205 に格納された特徴量から現在時刻 t のネットワークの状態 s を算出する。強化学習実行部 204 は、算出した状態 s の取り得る行動 a のなかから上記 *greedy* 方式のような方法で行動 a を選択する。強化学習実行部 204 は、当該選択した行動に対応する制御内容（制御パラメータの更新値）をパケット転送部 201 に通知する。強化学習実行部 204 は、上記行動に応じたネットワークの変化に応じて報酬を定める。その際、強化学習実行部 204 は、ネットワークに対して行われた行動の報酬を、行動が行われた後のネットワークの定常性に基づき定める。

【0071】

具体的には、強化学習実行部 204 は、行動 a を起こした結果、ネットワークが定常状態にあるか否かに基づき報酬を決定する。強化学習実行部 204 は、式（2）や式（3）

10

20

30

40

50

に記載された報酬  $r_{t+1}$  を定める際、ネットワークが定常状態であれば（ネットワークが安定していれば）、正の報酬を与える。対して、ネットワークの状態が非定常状態であれば（ネットワークが不安定であれば）、強化学習実行部 204 は、負の報酬を与える。

【0072】

強化学習実行部 204 は、ネットワークに対して行動を起こしたことにより変動するネットワークの状態に関する時系列データに対して統計処理を実施することで、ネットワークの定常性を判定する。

【0073】

具体的には、強化学習実行部 204 は、上記 *-greedy* 方式のような方法で選択された行動  $a$  に対応するネットワークの制御を実行後の次の時刻  $t+1$  から所定期間前までの特徴量（特徴量の時系列データ）を読み出す。強化学習実行部 204 は、当該読み出された特徴量の時系列データに対して統計処理を施すことで、ネットワークの状態が定常状態か否かを示す評価指標を算出する。

【0074】

具体的には、強化学習実行部 204 は、上記時系列データを自己回帰（Autoregressive model; AR）モデルによりモデル化する。ARモデルは、時系列データ  $x_1$ 、 $x_2$ 、 $\dots$ 、 $x_N$  を下記の式（4）に示すように、現在時刻の値を、重みが付けられた過去の値の加算（線形和）により表現するものである。

【0075】

【数4】

$$x(t) = c + \sum_{i=1}^p w_i x(t-i) + \epsilon(t)$$

・・・(4)

【0076】

式（4）において、 $x(t)$  は特徴量、 $\epsilon(t)$  はノイズ（ホワイトノイズ）、 $c$  は時刻により変化しない定数、 $w_i$  は重みを示す。 $i$  は過去の時刻を指定するためのサフィックスであり、 $p$  は上記所定期間前を指定する整数である。

【0077】

強化学習実行部 204 は、上記式（4）に示される重み  $w_i$  を記憶部 205 から読み出した時系列データを用いて推定する。具体的には、強化学習実行部 204 は、最尤法、ユールウォーカー等のパラメータ推定手法により重み  $w_i$  を推定する。なお、最尤法、ユールウォーカー等のパラメータ推定手法は公知の技術を用いることができるのでその詳細な説明を省略する。

【0078】

次に、強化学習実行部 204 は、時系列データから得られた ARモデルに対して単位根検定を実施する。単位根検定を実施することで、強化学習実行部 204 は、時系列データの定常度（定常度合い）を得る。強化学習実行部 204 は、単位根検定の実行により、「非定常」に対する「定常」の割合を算出することができる。単位根検定は既存のアルゴリズムにより実現でき、且つ、当業者にとって明らかであるのでその詳細な説明を省略する。

【0079】

強化学習実行部 204 は、単位根検定により得られた定常度に対して閾値処理（例えば、取得した値が閾値以上または未満であるかを判定する処理）を実行し、ネットワークの状態が定常状態にあるか否かを判定する。つまり、強化学習実行部 204 は、ネットワークの状態が、定常状態に向かう過渡的な「非定常状態」にあるのか、又は、特定の値を中心に収束している「定常状態」にあるのか判定する。

10

20

30

40

50

## 【 0 0 8 0 】

具体的には、強化学習実行部 2 0 4 は、定常度が閾値以上であればネットワークの状態は「定常」と判定する。強化学習実行部 2 0 4 は、定常度が閾値よりも小さければネットワークの状態は「非定常」と判定する。

## 【 0 0 8 1 】

図 1 0 は、特徴量の時系列データの一部を示す図である。図 1 0 A に示す時系列データに対して、強化学習実行部 2 0 4 が単位根検定を実施すると、ネットワークの状態は「非定常」と判定される。

## 【 0 0 8 2 】

この場合、強化学習実行部 2 0 4 は、式 ( 2 ) や式 ( 3 ) の報酬  $r_{t+1}$  に負の報酬 ( 例  
例えば、 - 1 ) を与え、Q テーブルや重みを更新する。対して、図 1 0 B に示す時系列データ  
に対して、強化学習実行部 2 0 4 が単位根検定を実施すると、ネットワークの状態は「  
定常」と判定される。この場合、強化学習実行部 2 0 4 は、式 ( 2 ) や式 ( 3 ) の報酬  $r_{t+1}$  に正の報酬 ( 例  
例えば、 + 1 ) を与え、Q テーブルや重みを更新する。

10

## 【 0 0 8 3 】

第 1 の実施形態に係る制御装置 2 0 の制御モード時の動作をまとめると図 1 1 に示すフロー  
チャートのとおりとなる。

## 【 0 0 8 4 】

制御装置 2 0 は、パケットを取得し、特徴量を算出する ( ステップ S 1 0 1 ) 。制御装  
置 2 0 は、当該算出された特徴量に基づきネットワークの状態を特定する ( ステップ S 1  
0 2 ) 。制御装置 2 0 は、学習モデルを用いて、ネットワークの状態に応じた最も価値の  
高い行動によりネットワークを制御する ( ステップ S 1 0 3 ) 。

20

## 【 0 0 8 5 】

第 1 の実施形態に係る制御装置 2 0 の学習モード時の動作をまとめると図 1 2 に示すフロ  
ーチャートのとおりとなる。

## 【 0 0 8 6 】

制御装置 2 0 は、パケットを取得し、特徴量を算出する ( ステップ S 2 0 1 ) 。制御装  
置 2 0 は、当該算出された特徴量に基づきネットワークの状態を特定する ( ステップ S 2  
0 2 ) 。制御装置 2 0 は、 - g r e e d y 方式等により現在のネットワーク状態にて取り  
得る行動を選択する ( ステップ S 2 0 3 ) 。制御装置 2 0 は、当該選択された行動により  
ネットワークを制御する ( ステップ S 2 0 4 ) 。制御装置 2 0 は、特徴量の時系列データ  
を用いてネットワークの定常性を判定する ( ステップ S 2 0 5 ) 。制御装置 2 0 は、判定  
結果により報酬を定め ( ステップ S 2 0 6 ) 、学習情報 ( Q テーブル、重み ) を更新する  
( ステップ S 2 0 7 ) 。

30

## 【 0 0 8 7 】

続いて、端末 1 0 の種類ごとに制御装置 2 0 の動作について具体的に説明する。

## 【 0 0 8 8 】

## [ 端末がドローンの場合 ]

端末 1 0 がドローンの場合、ネットワークの状態を示す指標 ( 特徴量 ) として、例えば  
、ドローンからサーバ 3 0 へ向けて送信されるパケットの平均パケット到着間隔が選択さ  
れる。サーバ 3 0 は、ドローンに対して制御パケット ( 制御コマンドを含むパケット ) を  
送信する。当該制御パケットに対するドローンからの応答パケット ( 肯定応答、否定応答  
) の平均パケット到着間隔が特徴量として選択される。

40

## 【 0 0 8 9 】

制御装置 2 0 は、サーバ 3 0 とドローンの間のパケット送受信の間隔が安定するように  
、制御パラメータを決定しネットワークの制御を行う。端末 1 0 がドローンの場合の取り  
得る行動 ( 変更可能な制御パラメータ ) としては、サーバ 3 0 から取得した制御パケット  
を格納するバッファからのパケット読み出し間隔 ( パケット送信間隔 ) が考えられる。

## 【 0 0 9 0 】

強化学習実行部 2 0 4 は、ドローンからサーバ 3 0 に送信される応答パケットの平均パ

50

ケット到着間隔が安定するように、バッファから制御パケットを読み出すパラメータを学習する。サーバ30がドローン（制御対象）を遠隔制御するアプリケーションでは、ドローンとサーバ30間で送受信されるパケット（制御パケット、応答パケット）が安定して相手側に届くことが重視される。

【0091】

ここで、制御パケットや応答パケットのパケットサイズはあまり大きくない。そのため、サーバ30からのスループットが高いが、パケットの送受信が安定しない状況（一度に多くの情報を送れるがパケットの到着にばらつきがある状況）よりも、スループットは低いパケットの送受信が安定する状況の方が、ドローンの制御では価値が高い。

【0092】

第1の実施形態に係る制御装置20は、ネットワークの状態（トラヒックの状態）を特徴付ける特徴量を適切に選択（例えば、平均パケット到着間隔を選択）することで、ドローンの遠隔制御というアプリケーションに適したネットワーク制御を実現できる。

【0093】

[ 端末がWEBカメラの場合 ]

上記説明では、報酬  $r_{t+1}$  を決定する条件（基準）としてネットワークの定常性を用いる場合について説明したが、上記定常性に他の基準を加えて報酬  $r_{t+1}$  を決定してもよい。ここでは、端末10がWEBカメラである場合を例に取り、報酬  $r_{t+1}$  の決定に「ネットワークの定常性」以外の項目を考慮する場合について説明する。

【0094】

端末10がWEBカメラの場合、ネットワークの状態を示す指標（特徴量）として、例えば、WEBカメラからサーバ30に流れるトラヒックのスループットが選択される。強化学習実行部204は、WEBカメラからサーバ30へのスループットが目標値の近傍で安定するように、学習モデルを算出する。

【0095】

例えば、端末10、サーバ30との間で形成されるTCPセッションのフローウィンドウサイズが制御パラメータに設定され、上記目標（スループットが目標値で安定）を実現するような行動が学習される。強化学習実行部204は、特徴量算出部202が算出した特徴量（スループット）の時系列データを用いてネットワークの定常性を判定する。

【0096】

続いて、強化学習実行部204は、特徴量（スループット）の範囲に応じて報酬  $r_{t+1}$  を決定する。例えば、目標値が閾値  $TH_{21}$  以上、且つ、閾値  $TH_{22}$  以下とすれば、強化学習実行部204は、図13に示すような方針（ポリシー）にて報酬  $r_{t+1}$  を決定する。このような報酬の与え方により得られた学習モデルを用いることで、WEBカメラからのスループットが目標とする値近傍で安定するようにネットワークは制御される。

【0097】

具体的には、制御装置20によるネットワーク制御により、図14Aに示すようなネットワークの状態（スループットが目標値近辺で安定）を実現できる。換言すれば、スループットの範囲を考慮して報酬  $r_{t+1}$  を決定することで、図14Bに示すようなネットワークの状態に陥ることが回避される。図14Bでは、最終的にネットワークの状態が安定しているが、定常時のスループットは目標値から大きく乖離している。

【0098】

なお、図13には、スループットが所定の範囲内であれば正の報酬を与える場合を記載したが、スループットが所定の値以上の場合に正の報酬を与えてもよい（図15参照）。図14Bの状況とは逆に、目標値から遠く離れた高い値でスループットが安定することが許容できる場合には、図15に示すように報酬  $r_{t+1}$  が決定されてもよい。

【0099】

スループットに設ける制限に関しては、制御装置20のリソース（通信リソース）を考慮して決定すればよい。例えば、制御パラメータにフローウィンドウサイズを選択した場合、当該ウィンドウサイズを大きくすればスループットは高い値で安定すると考えられる。

10

20

30

40

50

しかしながら、大きなフローウィンドウサイズを用意するためにはメモリ（リソース）の消費が大きくなり、他の端末 10 に割り当て可能なリソースが減少してしまう。制御装置 20 は、上記のようなメリット、デメリットを考慮してテーブル更新ポリシーを決定すればよい。

#### 【0100】

[ 端末がスマートフォンの場合 ]

上記では、1つの特徴量によりネットワークの定常性を判定したりする場合について説明したが、複数の特徴量によりネットワークの定常性の判定等が行われてもよい。以下、端末 10 がスマートフォンである場合を例に取り、ネットワークの定常性が複数の特徴量により判定される場合について説明する。

#### 【0101】

ここでは、サーバ 30 から動画が配信され、スマートフォン（端末 10）にて当該動画が再生される場合を想定する。特徴量算出部 202 は、サーバ 30 からスマートフォンに流れるトラフィックのスループットと平均パケット到着間隔を算出する。

#### 【0102】

強化学習実行部 204 は、当該 2 つの特徴量からネットワークの定常性を判定する。具体的には、強化学習実行部 204 は、スループットの時系列データに基づきスループットが安定しているか否かを判定する。同様に、強化学習実行部 204 は、平均パケット到着間隔の時系列データに基づき平均パケット到着間隔が安定しているか否かを判定する。

#### 【0103】

強化学習実行部 204 は、スループット及び平均パケット到着間隔が共に定常状態にある場合に、ネットワークが定常状態にあると判定し、報酬  $r_{t+1}$  に正の報酬を与え、他の場合には負の報酬を与える。

#### 【0104】

以上のように、第 1 の実施形態に係る制御装置 20 は、ネットワークの状態を、ネットワークに流れるトラフィックを特徴付ける特徴量を用いて推定する。制御装置 20 は、ネットワークに対して行った行動（制御パラメータの変更）により得られる状態の時系列変化に応じて、当該行動に対する報酬を定める。そのため、ネットワークにて提供されるサービスやアプリケーションレベルで求められる、「ネットワークの安定性」に高い報酬が与えられ、アプリケーション等に適したネットワーク品質の向上が実現できる。即ち、本願開示では、強化学習の際にネットワークの状態が安定している収束状態に価値が高いものと捉え、そのような状況の場合に学習器が環境（ネットワーク）に適応できていると考える、報酬を決定している。

#### 【0105】

[ 第 2 の実施形態 ]

続いて、第 2 の実施形態について図面を参照して詳細に説明する。

#### 【0106】

第 1 の実施形態では、ネットワークに流されるトラフィックを特徴付ける特徴量（例えば、スループット）によりネットワークの状態を推定している。第 2 の実施形態では、端末 10 における QoE（ユーザ体感品質）や QoC（制御品質）に基づきのネットワークの状態を決定する場合について説明する。

#### 【0107】

例えば、端末 10 がスマートフォンであって、動画再生アプリケーションが動作している場合を考える。この場合、端末 10 は、再生動画の画質、ビットレート、途絶回数（バッファが空となった回数）、フレームレート等を制御装置 20 に通知する。あるいは、端末 10 は、ITU（International Telecommunication Union）-T 勧告 P. 1203 に規定された MOS（Mean Opinion Score）値を制御装置 20 に送信してもよい。

#### 【0108】

あるいは、スマートフォンにて WEB ページの閲覧（ブラウザが動作）が行われている場合には、端末 10 は、ページ表示までの初期待機時間を制御装置 20 に通知してもよい。

10

20

30

40

50

## 【 0 1 0 9 】

例えば、端末 1 0 がロボットである場合には、ロボットは、制御コマンドの受信間隔、作業完了時間、作業成功回数等を制御装置 2 0 に通知してもよい。

## 【 0 1 1 0 】

あるいは、端末 1 0 が監視カメラである場合には、監視カメラは、監視対象（例えば、人の顔、物体等）の認証率、認証回数等を制御装置 2 0 に通知してもよい。

## 【 0 1 1 1 】

制御装置 2 0 は、端末 1 0 から当該端末 1 0 における Q o E を示す値（例えば、上記初期待機時間等）を取得し、当該値に基づきネットワークの定常性を判定し、報酬  $r_{t+1}$  を決定してもよい。その際、制御装置 2 0 は、第 1 の実施形態にて説明した方法と同様にして、端末 1 0 から取得した Q o E の時系列データに対して単位根検定を実施し、ネットワークの定常性を評価すればよい。

10

## 【 0 1 1 2 】

あるいは、制御装置 2 0 は、端末 1 0 とサーバ 3 0 の間に流れるトラフィックから上記 Q o E を示す値を推定してもよい。例えば、制御装置 2 0 は、スループットからビットレートを推定し、当該推定値に基づきネットワークの定常性を判定してもよい。なお、スループットからビットレートを推定する際には、以下の参考文献 1 に記載された方法を用いればよい。

[ 参考文献 1 ] : 国際公開第 2 0 1 9 / 0 4 4 0 6 5 号

## 【 0 1 1 3 】

以上のように、第 2 の実施形態に係る制御装置 2 0 は、ネットワークの状態を、ユーザ体感品質 ( Q o E ) や制御品質 ( Q o C ) から推定し、ユーザ体感品質等が安定している場合に高い報酬を与えても良い。例えば、ユーザが端末を使用して動画を視聴する場合を考える。この場合、本願開示では、フレームレートが頻繁に変わるネットワーク環境 ( フレームレートが安定しない環境 ) よりも、低いフレームレートであっても一定しているネットワーク環境の方が、ネットワーク品質が高いと判断している。換言すれば、制御装置 2 0 は、このような高いネットワーク品質を実現する制御パラメータを強化学習により学習する。

20

## 【 0 1 1 4 】

続いて、通信ネットワークシステムを構成する各装置のハードウェアについて説明する。図 1 6 は、制御装置 2 0 のハードウェア構成の一例を示す図である。

30

## 【 0 1 1 5 】

制御装置 2 0 は、情報処理装置 ( 所謂、コンピュータ ) により構成可能であり、図 1 6 に例示する構成を備える。例えば、制御装置 2 0 は、プロセッサ 3 1 1、メモリ 3 1 2、入出力インターフェイス 3 1 3 及び通信インターフェイス 3 1 4 等を備える。上記プロセッサ 3 1 1 等の構成要素は内部バス等により接続され、相互に通信可能に構成されている。

## 【 0 1 1 6 】

但し、図 1 6 に示す構成は、制御装置 2 0 のハードウェア構成を限定する趣旨ではない。制御装置 2 0 は、図示しないハードウェアを含んでもよいし、必要に応じて入出力インターフェイス 3 1 3 を備えていなくともよい。また、制御装置 2 0 に含まれるプロセッサ 3 1 1 等の数も図 1 6 の例示に限定する趣旨ではなく、例えば、複数のプロセッサ 3 1 1 が制御装置 2 0 に含まれていてもよい。

40

## 【 0 1 1 7 】

プロセッサ 3 1 1 は、例えば、C P U ( Central Processing Unit )、M P U ( Micro Processing Unit )、D S P ( Digital Signal Processor ) 等のプログラマブルなデバイスである。あるいは、プロセッサ 3 1 1 は、F P G A ( Field Programmable Gate Array )、A S I C ( Application Specific Integrated Circuit ) 等のデバイスであってもよい。プロセッサ 3 1 1 は、オペレーティングシステム ( O S ; Operating System ) を含む各種プログラムを実行する。

## 【 0 1 1 8 】

50

メモリ 312 は、RAM (Random Access Memory)、ROM (Read Only Memory)、HDD (Hard Disk Drive)、SSD (Solid State Drive) 等である。メモリ 312 は、OS プログラム、アプリケーションプログラム、各種データを格納する。

【0119】

入出力インターフェイス 313 は、図示しない表示装置や入力装置のインターフェイスである。表示装置は、例えば、液晶ディスプレイ等である。入力装置は、例えば、キーボードやマウス等のユーザ操作を受け付ける装置である。

【0120】

通信インターフェイス 314 は、他の装置と通信を行う回路、モジュール等である。例えば、通信インターフェイス 314 は、NIC (Network Interface Card) 等を備える。

【0121】

制御装置 20 の機能は、各種処理モジュールにより実現される。当該処理モジュールは、例えば、メモリ 312 に格納されたプログラムをプロセッサ 311 が実行することで実現される。また、当該プログラムは、コンピュータが読み取り可能な記憶媒体に記録することができる。記憶媒体は、半導体メモリ、ハードディスク、磁気記録媒体、光記録媒体等の非トランジエント (non-transitory) なものとすることができる。即ち、本発明は、コンピュータプログラム製品として具現することも可能である。また、上記プログラムは、ネットワークを介してダウンロードするか、あるいは、プログラムを記憶した記憶媒体を用いて、更新することができる。さらに、上記処理モジュールは、半導体チップにより実現されてもよい。

【0122】

なお、端末 10、サーバ 30 も制御装置 20 と同様に情報処理装置により構成可能であり、その基本的なハードウェア構成は制御装置 20 と相違する点はないので説明を省略する。

【0123】

[変形例]

なお、上記実施形態にて説明した通信ネットワークシステムの構成、動作等は例示であって、システムの構成等を限定する趣旨ではない。例えば、制御装置 20 は、ネットワークを制御する装置と学習モデルを生成する装置に分離されていてもよい。あるいは、学習情報 (学習モデル) を記憶する記憶部 205 は、外部のデータベースサーバ等により実現されてもよい。即ち、本願開示は、学習手段、制御手段、記憶手段等を含むシステムとして実施されてもよい。

【0124】

上記実施形態では、特徴量の時系列データに対して単位根検定を実施することとで、ネットワークの定常度を算出している。しかし、ネットワークの定常度は他の指標により算出されてもよい。例えば、強化学習実行部 204 は、データのばらつき度合いを示す標準偏差を計算し、「平均 - 標準偏差」が閾値以上の場合にネットワークは定常状態であると判定してもよい。

【0125】

上記実施形態では、1つの閾値を用いてネットワークの定常性 (安定性) を判定しているが、複数の閾値を用いてより細かくネットワークの定常度合いが算出されてもよい。例えば、「極めて安定」、「安定」、「不安定」、「極めて不安定」のように4段階でネットワークの定常性が判定されてもよい。この場合、ネットワークの定常度合いに応じて報酬が決められていてもよい。

【0126】

なお、端末 10 はセンサ装置である場合がある。センサ装置は、オン/オフモデルに従う通信パターン (通信トラヒック) を発生する。つまり、端末 10 がセンサ装置等であれば、データ (パケット) がネットワークに流れる場合と流れない場合 (無通信状態) が生じ得る。そのため、制御装置 20 が、トラヒック (特徴量) の時系列データそのものを使って定常性判定 (単位根検定) を実施するのではなく、変動パターンにより定常性が判定

10

20

30

40

50

されてもよい。制御装置 20 は、特徴量が上下する時間間隔に関する時系列データを用いてネットワークの定常性を判定してもよい。あるいは、制御装置 20 は、事前にオン/オフモデルに従うアプリケーションを把握している場合には、無通信状態は報酬に反映しない等の対応を行ってもよい。即ち、制御装置 20 は、ネットワークの状態が「通信状態」にある場合に強化学習の報酬を与えるようにしてもよい。

【0127】

上記実施形態では、制御装置 20 は、トラヒックフローを制御の対象（制御単位）とする場合について説明した。しかし、制御装置 20 は、端末 10 単位、又は、複数の端末 10 をまとめたグループを制御の対象としてもよい。つまり、同じ端末 10 であってもアプリケーションが異なればポート番号等が異なり、異なるフローとして扱われる。制御装置 20 は、同じ端末 10 から送信されるパケットには同じ制御（制御パラメータの変更）を適用してもよい。あるいは、制御装置 20 は、例えば、同じ種類の端末 10 を一つのグループとして扱い、同じグループに属する端末 10 から送信されるパケットに対して同じ制御を適用してもよい。

10

【0128】

上述の説明で用いた複数のフローチャートでは、複数の工程（処理）が順番に記載されているが、各実施形態で実行される工程の実行順序は、その記載の順番に制限されない。各実施形態では、例えば各処理を並行して実行する等、図示される工程の順番を内容的に支障のない範囲で変更することができる。また、上述の各実施形態は、内容が相反しない範囲で組み合わせることができる。

20

【0129】

上記の実施形態の一部又は全部は、以下の付記のようにも記載され得るが、以下には限られない。

[付記 1]

ネットワークを制御するための行動を学習する、学習部（101、204）と、前記学習部（101、204）が生成した学習情報を記憶する、記憶部（102、205）と、を備え、前記学習部（101、204）は、前記ネットワークに対して行われた行動の報酬を、前記行動が行われた後のネットワークの定常性に基づき定める、制御装置（20、100）。

30

[付記 2]

前記学習部（101、204）は、前記行動が行われた後のネットワークが定常状態であれば、前記ネットワークに対して行われた行動に正の報酬を与え、前記行動が行われた後のネットワークが非定常状態であれば、前記ネットワークに対して行われた行動に負の報酬を与える、付記 1 に記載の制御装置（20、100）。

[付記 3]

前記学習部（101、204）は、前記ネットワークに対して行動を起こしたことにより変動するネットワークの状態に関する時系列データに基づいて前記ネットワークの定常性を判定する、付記 1 又は 2 に記載の制御装置（20、100）。

40

[付記 4]

前記学習部（101、204）は、前記ネットワークの状態を、前記ネットワークに流れるトラヒックを特徴付ける特徴量、ユーザ体感品質及び制御品質のうち少なくとも一つから推定する、付記 3 に記載の制御装置（20、100）。

[付記 5]

前記学習部（101、204）が生成した学習モデルから得られる行動に基づき、前記ネットワークを制御する、制御部（203）をさらに備える、付記 1 乃至 4 のいずれか一つに記載の制御装置（20、100）。

[付記 6]

50

ネットワークを制御するための行動を学習するステップと、  
前記学習により生成された学習情報を記憶するステップと、  
を含み、  
前記学習するステップは、  
前記ネットワークに対して行われた行動の報酬を、前記行動が行われた後のネットワークの定常性に基づき定める、方法。

[ 付記 7 ]

前記学習するステップは、  
前記行動が行われた後のネットワークが定常状態であれば、前記ネットワークに対して行われた行動に正の報酬を与え、  
前記行動が行われた後のネットワークが非定常状態であれば、前記ネットワークに対して行われた行動に負の報酬を与える、付記 6 に記載の方法。

10

[ 付記 8 ]

前記学習するステップは、  
前記ネットワークに対して行動を起こしたことにより変動するネットワークの状態に関する時系列データに基づいて前記ネットワークの定常性を判定する、付記 6 又は 7 に記載の方法。

[ 付記 9 ]

前記学習するステップは、前記ネットワークの状態を、前記ネットワークに流れるトラフィックを特徴付ける特徴量、ユーザ体感品質及び制御品質のうち少なくとも 1 つから推定する、付記 8 に記載の方法。

20

[ 付記 10 ]

前記学習するステップにより生成された学習モデルから得られる行動に基づき、前記ネットワークを制御するステップをさらに含む、付記 6 乃至 9 のいずれか一つに記載の方法。

[ 付記 11 ]

ネットワークを制御するための行動を学習する、学習手段 ( 101、204 ) と、  
前記学習手段が生成した学習情報を記憶する、記憶手段 ( 102、205 ) と、を含み、  
前記学習手段 ( 101、204 ) は、  
前記ネットワークに対して行われた行動の報酬を、前記行動が行われた後のネットワークの定常性に基づき定める、システム。

30

[ 付記 12 ]

前記学習手段 ( 101、204 ) は、  
前記行動が行われた後のネットワークが定常状態であれば、前記ネットワークに対して行われた行動に正の報酬を与え、  
前記行動が行われた後のネットワークが非定常状態であれば、前記ネットワークに対して行われた行動に負の報酬を与える、付記 11 に記載のシステム。

[ 付記 13 ]

前記学習手段 ( 101、204 ) は、  
前記ネットワークに対して行動を起こしたことにより変動するネットワークの状態に関する時系列データに基づいて前記ネットワークの定常性を判定する、付記 11 又は 12 に記載のシステム。

40

[ 付記 14 ]

前記学習手段 ( 101、204 ) は、前記ネットワークの状態を、前記ネットワークに流れるトラフィックを特徴付ける特徴量、ユーザ体感品質及び制御品質のうち少なくとも 1 つから推定する、付記 13 に記載のシステム。

[ 付記 15 ]

前記学習手段 ( 101、204 ) が生成した学習モデルから得られる行動に基づき、前記ネットワークを制御する、制御手段 ( 203 ) をさらに備える、付記 11 乃至 14 のいずれか一つに記載のシステム。

[ 付記 16 ]

50

コンピュータ（３１１）に、  
 ネットワークを制御するための行動を学習する処理と、  
 前記学習により生成された学習情報を記憶する処理と、  
 を実行させ、  
 前記学習する処理は、  
 前記ネットワークに対して行われた行動の報酬を、前記行動が行われた後のネットワー  
 クの定常性に基づき定める、プログラム。

【０１３０】

なお、引用した上記の先行技術文献の各開示は、本書に引用をもって繰り込むものとする。以上、本発明の実施形態を説明したが、本発明はこれらの実施形態に限定されるもの  
 ではない。これらの実施形態は例示にすぎないということ、及び、本発明のスコップ及び  
 精神から逸脱することなく様々な変形が可能であるということは、当業者に理解されるで  
 あろう。

10

【符号の説明】

【０１３１】

- １０ 端末
- ２０、１００ 制御装置
- ３０ サーバ
- １０１ 学習部
- １０２、２０５ 記憶部
- ２０１ パケット転送装置
- ２０２ 特徴量算出部
- ２０３ ネットワーク制御部
- ２０４ 強化学習実行部
- ３１１ プロセッサ
- ３１２ メモリ
- ３１３ 入出力インターフェイス
- ３１４ 通信インターフェイス

20

30

40

50

【図面】

【図 1】

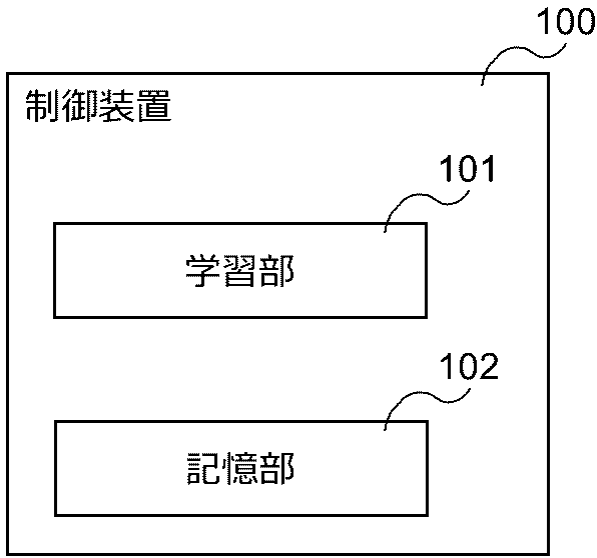


Fig.1

【図 2】

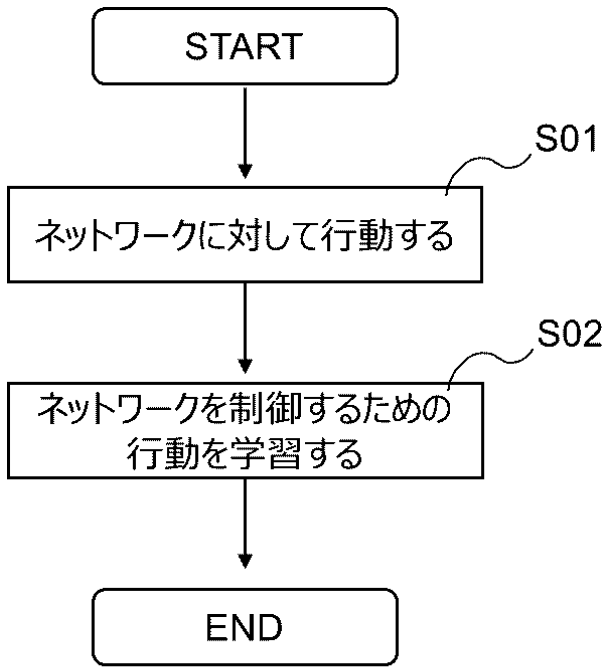


Fig.2

【図 3】

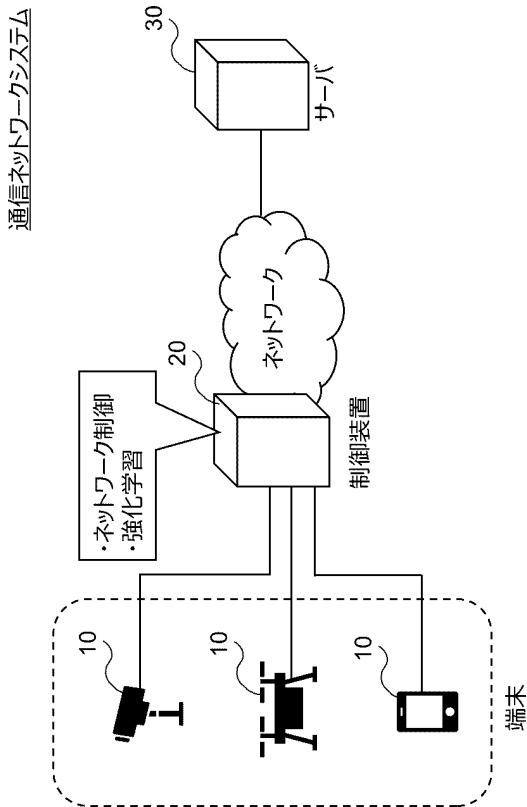


Fig.3

【図 4】

Qテーブル

	行動A1	行動A2	行動A3	...
状態S1	Q(S1,A1)	Q(S1,A2)	Q(S1,A3)	...
状態S2	Q(S2,A1)	Q(S2,A2)	Q(S2,A3)	...
状態S3	Q(S3,A1)	Q(S3,A2)	Q(S3,A3)	...
...	...	...	...	...

Fig.4

10

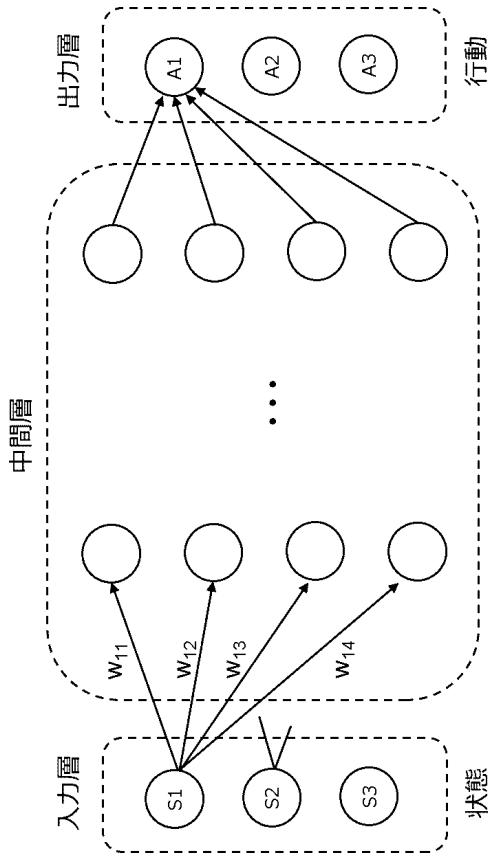
20

30

40

50

【図5】



【図6】

重み

Fig.5

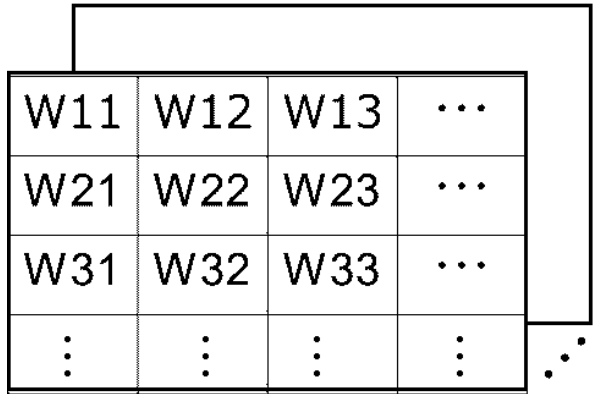
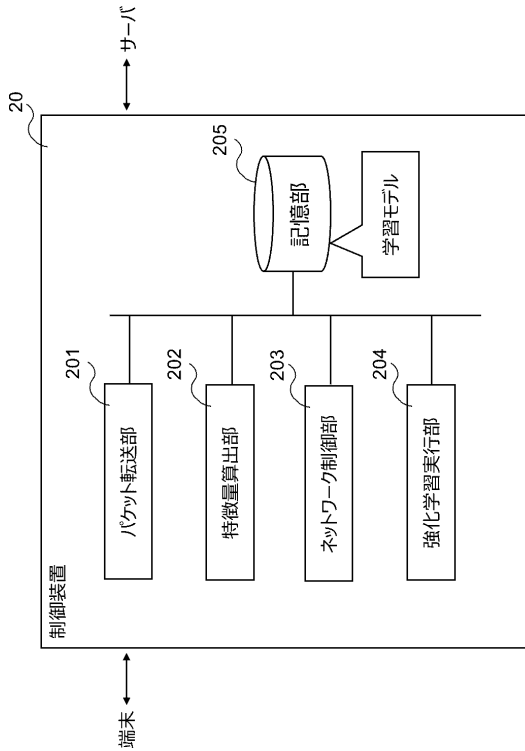


Fig.6

【図7】



【図8】

特徴量	ネットワークの状態
$F < TH11$	状態S1
$TH11 \leq F < TH12$	状態S2
$TH12 \leq F < TH13$	状態S2
⋮	⋮

Fig.7

Fig.8

10

20

30

40

50

【図9】

行動	制御内容
行動A1	ウィンドウサイズをAバイト増加
行動A2	ウィンドウサイズをBバイト増加
行動A3	ウィンドウサイズをCバイト増加
⋮	⋮

Fig.9

【図10】

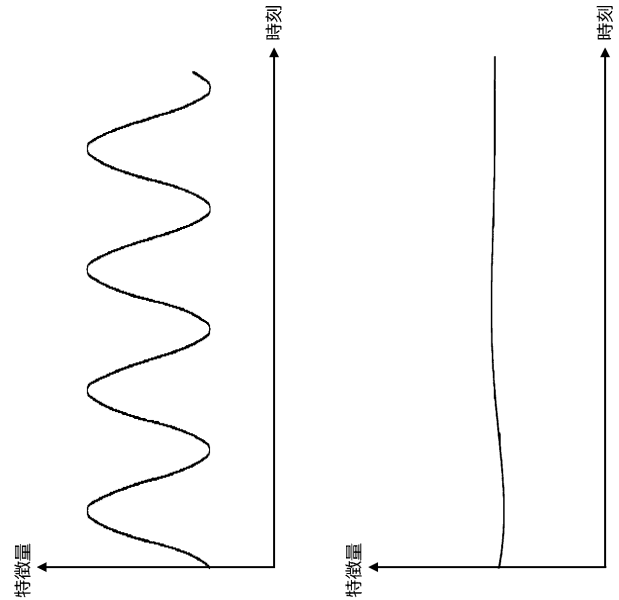


Fig. 10A

Fig. 10B

【図11】

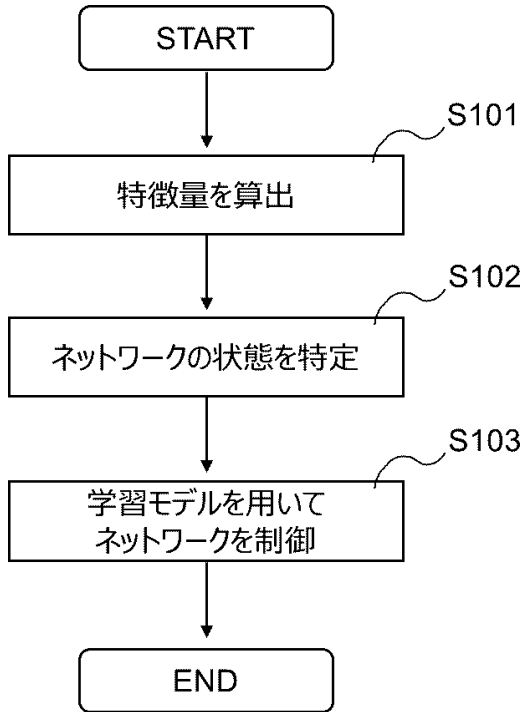


Fig. 11

【図12】

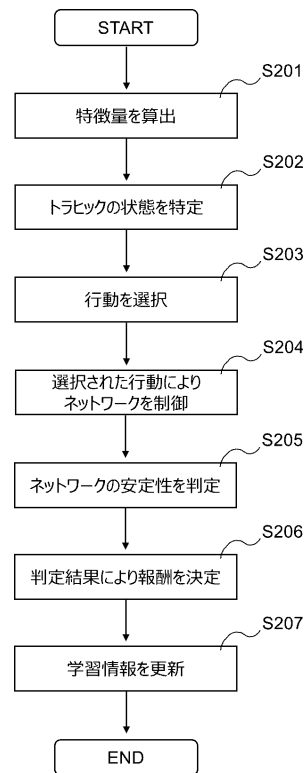


Fig. 12

10

20

30

40

50

【 図 1 3 】

	スループット ( $T < TH21$ )	スループット ( $TH21 \leq T \leq TH22$ )	スループット ( $TH22 < T$ )
ネットワークは安定	負の報酬	正の報酬	負の報酬
ネットワークは不安定	負の報酬	負の報酬	負の報酬

Fig.13

【 図 1 4 】

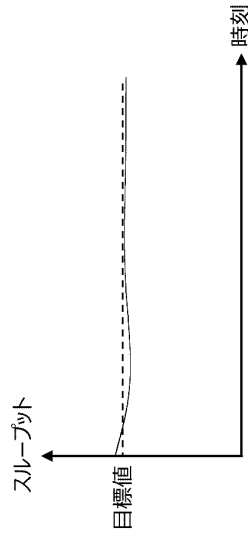


Fig.14A

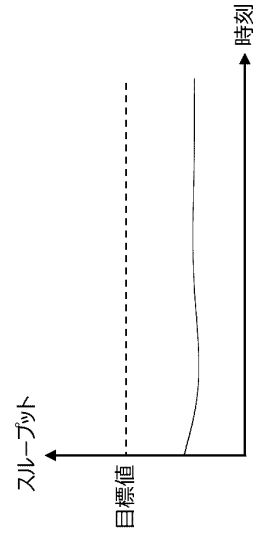


Fig.14B

【 図 1 5 】

	スループット ( $T < TH31$ )	スループット ( $TH31 \leq T$ )
ネットワークは安定	負の報酬	正の報酬
ネットワークは不安定	負の報酬	負の報酬

Fig.15

【 図 1 6 】

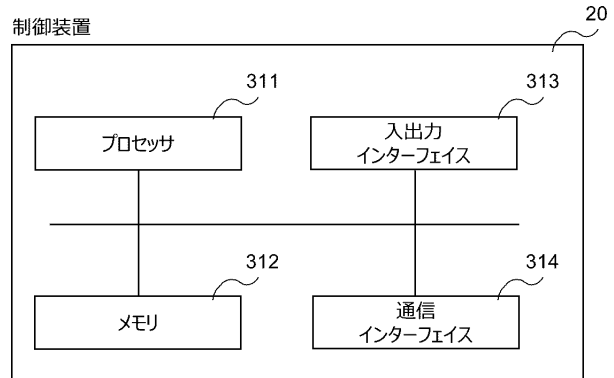


Fig.16

10

20

30

40

50

## フロントページの続き

- (56)参考文献 特開2013-106202(JP,A)  
国際公開第2019/176997(WO,A1)  
特開2019-041338(JP,A)  
米国特許出願公開第2019/0141113(US,A1)  
特開2009-027303(JP,A)
- (58)調査した分野 (Int.Cl., DB名)
- H04L 12/00 - 13/18  
H04L 41/00 - 49/9057  
H04L 61/00 - 65/80  
H04L 69/00 - 69/40  
G06N 3/00 - 3/12  
G06N 7/08 - 99/00