

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号

特許第7699791号
(P7699791)

(45)発行日 令和7年6月30日(2025. 6. 30)

(24)登録日 令和7年6月20日(2025. 6. 20)

(51)Int. Cl. F I
 G 0 6 N 3/063 (2023. 01) G 0 6 N 3/063
 G 0 6 F 9/50 (2006. 01) G 0 6 F 9/50 1 5 0 D

請求項の数 5 (全 18 頁)

(21)出願番号	特願2020-194733(P2020-194733)	(73)特許権者	506301140
(22)出願日	令和2年11月24日(2020. 11. 24)		公立大学法人会津大学
(65)公開番号	特開2022-83341(P2022-83341A)		福島県会津若松市一箕町大字鶴賀字上居合
(43)公開日	令和4年6月3日(2022. 6. 3)		9 0 番地
審査請求日	令和5年10月31日(2023. 10. 31)	(74)代理人	100094525
			弁理士 土井 健二
		(74)代理人	100094514
			弁理士 林 恒徳
		(72)発明者	ベン アブダラ アブデラゼク
			福島県会津若松市一箕町大字鶴賀字上居合
			9 0 番地 公立大学法人会津大学内
		(72)発明者	ファンクン ホアン
			福島県会津若松市一箕町大字鶴賀字上居合
			9 0 番地 公立大学法人会津大学内

最終頁に続く

(54)【発明の名称】 AI プロセッサ

(57)【特許請求の範囲】

【請求項 1】

複数の演算コアを有し、

前記複数の演算コアの少なくともいずれかが、畳み込み層と全結合層とを有する CNN (Convolutional Neural Network) の機械学習モデルに含まれる複数のニューロンのそれぞれに対応付けられた計算プログラムを分割して前記複数の演算コアのそれぞれに割り当てるマッピング処理を実行し、

前記複数の演算コアのそれぞれが、前記マッピング処理によって割り当てられた前記計算プログラムを実行し、

前記マッピング処理では、前記複数の演算コア間における通信コストが所定の閾値以下になるように、遺伝的アルゴリズムによって前記計算プログラムを前記複数の演算コアに割り当て、さらに、

前記機械学習モデルにおけるパラメータと他の機械学習モデルにおけるパラメータとから生成されたグローバルパラメータを用いることによって、前記機械学習モデルにおけるパラメータを更新し、

前記機械学習モデルは、個人情報を含む画像データの入力に伴って前記個人情報に対応する人物が所定の状態にあるか否かについての情報を出力する機械学習モデルであり、

前記マッピング処理では、前記複数の演算コアに含まれる第 1 演算コアが、前記マッピング処理の結果を示すマッピングテーブルを生成し、

前記機械学習モデルにおけるパラメータを更新する処理では、前記第 1 演算コアが、前

記マッピングテーブルを参照し、前記複数の演算コアのうちの複数の他の演算コアに対して前記機械学習モデルにおけるパラメータを送信することによって、前記機械学習モデルにおけるパラメータを更新する、

ことを特徴とするA Iプロセッサ。

【請求項2】

請求項1において、

前記画像データは、患者が映る画像データであり、

前記機械学習モデルは、前記画像データの入力に伴って前記患者が前記所定の状態にあるか否かについての情報を出力する機械学習モデルである、

ことを特徴とするA Iプロセッサ。

10

【請求項3】

請求項1において、

前記CNNは、プーリング層をさらに有する、

ことを特徴とするA Iプロセッサ。

【請求項4】

請求項1において、

前記マッピング処理では、

前記計算プログラムについての前記複数の演算コアに対するN個のマッピング解をランダムに生成し、

前記N個のマッピング解のそれぞれを採用した場合における前記通信コストを算出し、

前記N個のマッピング解から、算出した前記通信コストが小さい順にM個のマッピング解を特定し、

20

前記M個のマッピング解を交差させることによってN - M個の新たなマッピング解を生成し、

前記M個のマッピング解と前記N - M個の新たなマッピング解とを含むN個の新たなマッピング解において突然変異を発生させ、

前記突然変異を発生させた前記N個の新たなマッピング解のそれぞれを採用した場合における前記通信コストを再度算出し、

前記N個の新たなマッピング解のうち、再度算出した前記通信コストが最小である特定のマッピング解を特定し、

30

前記特定のマッピング解の前記通信コストが前記所定の閾値以下であるか否かを判定し

、前記特定のマッピング解の前記通信コストが前記所定の閾値以下であると判定した場合、前記計算プログラムを前記特定のマッピング解に従って前記複数の演算コアに割り当てる、

ことを特徴とするA Iプロセッサ。

【請求項5】

請求項4において、

前記マッピング処理では、

前記特定のマッピング解の前記通信コストが前記所定の閾値以下でないと判定した場合、前記N個の新たなマッピング解について、前記通信コストを算出する処理と、前記M個のマッピング解を特定する処理と、前記N - M個の新たなマッピング解を生成する処理と、前記突然変異を発生させる処理と、前記通信コストを再度算出する処理と、前記特定のマッピング解を特定する処理と、前記通信コストが所定の閾値であるか否かを判定する処理とを再度行う、

40

ことを特徴とするA Iプロセッサ。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、A Iプロセッサに関する。

【背景技術】

【0002】

昨今、SARS-CoV2ウイルス（以下、新型コロナウイルス）によって引き起こされるCOVID-19（以下、新型コロナウイルス感染症とも呼ぶ）が流行している。この新型コロナウイルスに感染しているか否かの検査を行うための標準的な方法は、例えば、患者由来の採取サンプルを使用する逆転写ポリメラーゼ連鎖反応（RT-PCR：Reverse Transcription Polymerase Chain Reaction）であり、60（%）から97（%）程度の感度を有している。また、新型コロナウイルスに感染しているか否かの検査を行うための別の方法としては、例えば、患者の肺を撮影したX線画像の解析があり、80（%）から90（%）程度の精度を有している。

10

【0003】

ここで、上記のようなX線画像の解析が行われる場合、医師は、患者のX線画像を手動で1枚ずつ診断する必要があるが、非効率的な診断処置の原因となっている。

【0004】

そこで、各病院では、例えば、患者が新型コロナウイルスに感染しているか否かの診断や患者が肺炎になっているか否かの診断を機械学習モデル（以下、診断モデルとも呼ぶ）に行わせることによって、多くの患者についての効率的な診断処置を行う場合がある。

【先行技術文献】

【特許文献】

20

【0005】

【特許文献1】米国特許出願公開20200026992A1号明細書

【非特許文献】

【0006】

【非特許文献1】Kun-Chih (Jimmy) Chen, Ting-Yi Wang, “NN-Noxim: High-Level Cycle-Accurate NoC-based Neural Networks Simulator”, 2018 11th International Workshop on Network on Chip Architectures (NoCArc).

【非特許文献2】Xiaoxiao Liu, Wei Wen, Xuehai Qian, Hai Li, Yiran Chen, “Neu-NoC: A High-efficient Interconnection Network for Accelerated Neuromorphic Systems”, 2018 23rd Asia and South Pacific Design Automation Conference (ASP-DAC)

30

【非特許文献3】Kun-Chih (Jimmy) Chen, Masoumeh Ebrahimi, Ting-Yi Wang, Yuch-Chi Yang, “NoC-based DNN Accelerator: A Future Design Paradigm”, NOCS '19, October 17-18, 2019, New York, NY, USA.

【発明の概要】

【発明が解決しようとする課題】

【0007】

しかしながら、上記のような診断モデルを用いた診断は、診断モデルを実行するコンピュータの性能によって効率性が大きく異なる。そのため、病院等の医療現場では、例えば、上記のような診断モデルを用いた診断をより効率的に行うことが可能なAI（Artificial Intelligence）プロセッサの開発が求められている。

40

【0008】

そこで、本発明の目的は、診断モデルを用いた診断をより効率的に行うAIプロセッサを提供することにある。

【課題を解決するための手段】

【0009】

本発明の一態様におけるAIプロセッサは、複数の演算コアを有し、前記複数の演算コアの少なくともいずれかが、畳み込み層と全結合層とを有するCNN（Convolutional Neural Network）の機械学習モデルに含まれる複数のニューロンのそれぞれに対応付けられた計算プログラムを分割して前記複数の演算コアのそれぞれに割り当てるマッピング処理を実行し、前記複数の演算コアのそれぞれが、前記マッピ

50

ング処理によって割り当てられた前記計算プログラムを実行し、前記マッピング処理では、前記複数の演算コア間における通信コストが所定の閾値以下になるように、遺伝的アルゴリズムによって前記計算プログラムを前記複数の演算コアに割り当てる。

【発明の効果】

【0010】

本発明の一態様によれば、診断モデルを用いた診断をより効率的に行うことが可能になる。

【図面の簡単な説明】

【0011】

【図1】図1は、従来の情報処理システム200の構成について説明する図である。 10

【図2】図2は、本実施の形態における情報処理システム100の構成について説明する図である。

【図3】図3は、本実施の形態におけるユーザインターフェースの例について説明する図である。

【図4】図4は、本実施の形態におけるユーザインターフェースの例について説明する図である。

【図5】図5は、本実施の形態におけるユーザインターフェースの例について説明する図である。

【図6】図6は、診断モデルMDの具体例について説明する図である。

【図7】図7は、診断モデルMDにおける更新処理について説明する図である。 20

【図8】図8は、本実施の形態におけるAIプロセッサPRの具体例を説明する図である。

【図9】図9は、本実施の形態におけるマッピング処理について説明する図である。

【図10】図10は、本実施の形態におけるマッピング処理について説明する図である。

【図11】図11は、本実施の形態におけるマッピング処理について説明する図である。

【図12】図12は、本実施の形態におけるマッピング処理について説明する図である。

【図13】図13は、本実施の形態におけるマッピング処理について説明する図である。

【図14】図14は、本実施の形態におけるマッピング処理について説明する図である。

【図15】図15は、本実施の形態におけるマッピング処理について説明する図である。

【図16】図16は、本実施の形態におけるマッピング処理について説明する図である。 30

【図17】図17は、本実施の形態におけるマッピング処理について説明する図である。

【図18】図18は、本実施の形態におけるマッピング処理について説明する図である。

【図19】図19は、本実施の形態におけるマッピング処理について説明する図である。

【図20】図20は、本実施の形態におけるマッピング処理について説明する図である。

【図21】図21は、本実施の形態におけるマッピング処理について説明する図である。

【発明を実施するための形態】

【0012】

以下、図面を参照して本発明の実施の形態について説明する。各実施の形態は、本発明のより良い理解のために準備されている。ただし、かかる実施の形態は、本発明の技術的範囲を限定するものではない。また、本発明の範囲は、特許請求の範囲及びこれと同等のものを網羅している。 40

【0013】

[従来の情報処理システム]

初めに、従来の情報処理システム200について説明を行う。図1は、従来の情報処理システム200の構成について説明する図である。なお、以下、3つの病院(病院11、病院12及び病院1n)が存在する場合について説明を行うが、3つ以外の数の病院が存在するものであってもよい。

【0014】

図1に示す情報処理システム200において、病院11では、医師が患者のX線画像21を分析することによって診断を行う。そして、医師は、例えば、病院11に設定された 50

情報処理装置（図示しない）を用いることにより、各患者が新型コロナウイルス感染症に感染しているか否かの診断結果（例えば、感染者数、各感染者の重篤度及び死者数等）を示す患者情報 3 1 を定期的に生成し、生成した患者情報 3 1 を政府の情報処理装置 4 に送信する。

【 0 0 1 5 】

同様に、図 1 に示す例において、病院 1 2 では、医師が患者の X 線画像 2 2 についての診断結果を示す患者情報 3 2 を生成し、生成した患者情報 3 2 を政府の情報処理装置 4 に送信する。また、病院 1 n では、医師が患者の X 線画像 2 n についての診断結果を示す患者情報 3 n を生成し、生成した患者情報 3 n を政府の情報処理装置 4 に送信する。

【 0 0 1 6 】

ここで、情報処理システム 2 0 0 では、病院 1 1、病院 1 2 及び病院 1 n のそれぞれが他の病院と協調することができない。この場合、各医師は、X 線画像についての診断を効率的に行うことができない。

【 0 0 1 7 】

[本実施の形態における情報処理システム]

次に、本実施の形態における情報処理システム 1 0 0 について説明を行う。図 2 は、本実施の形態における情報処理システム 1 0 0 の構成について説明する図である。また、図 3 から図 5 は、本実施の形態におけるユーザインターフェースの例について説明する図である。

【 0 0 1 8 】

図 2 に示す情報処理システム 1 0 0 において、病院 1 1 では、患者の携帯端末 5 1 から送信された X 線画像を病院 1 1 に設置された診断システム 6 1 に入力する。検知システム 6 1 は、例えば、携帯端末 5 1 から送信された X 線画像が新型コロナウイルスに感染した患者の画像である否かを判定する診断モデルを実行する情報処理装置である。そして、診断システム 6 1 は、携帯端末 5 1 から送信された X 線画像についての診断結果を携帯端末 5 1 に送信する。

【 0 0 1 9 】

同様に、図 2 に示す例において、病院 1 2 に設置された診断システム 6 2 では、携帯端末 5 2 から送信された X 線画像についての診断を行い、その診断結果を携帯端末 5 2 に送信する。また、病院 1 n に設置された診断システム 6 n では、携帯端末 5 n から送信された X 線画像についての診断を行い、その診断結果を携帯端末 5 n に送信する。

【 0 0 2 0 】

これにより、病院 1 1、病院 1 2 及び病院 1 n では、X 線画像についての診断に伴う医師の負担を軽減させることが可能になる。また、病院 1 1、病院 1 2 及び病院 1 n では、X 線画像についての誤った診断の発生を防止することが可能になる。また、各医師は、例えば、図 3 に示すユーザインターフェース U 1 を閲覧することにより、各患者の診断状況をリアルタイムで確認することが可能になる。

【 0 0 2 1 】

また、病院 1 1、病院 1 2 及び病院 1 n では、診断モデルによる診断を行うことで X 線画像についての診断を迅速に行うことが可能になり、携帯端末 5 1、携帯端末 5 2 及び携帯端末 5 n に対して診断結果を迅速に通知することが可能になる。そのため、各患者は、例えば、図 4 に示すユーザインターフェース U 2 を閲覧することにより、X 線画像についての診断結果を迅速に確認することが可能になる。

【 0 0 2 2 】

そして、診断システム 6 1 は、診断モデルによる診断結果（例えば、感染者数、各感染者の重篤度及び死者数等）を示す患者情報をリアルタイムに生成し、生成した患者情報を政府の情報処理装置 4 に送信する。

【 0 0 2 3 】

同様に、図 2 に示す例において、診断システム 6 2 及び診断システム 6 n のそれぞれは、診断モデルによる診断結果を示す患者情報をリアルタイムに生成して政府の情報処理装

10

20

30

40

50

置 4 に送信する。

【 0 0 2 4 】

具体的に、診断システム 6 1、診断システム 6 2 及び診断システム 6 n のそれぞれは、例えば、クラウドサーバ 7 に対して患者情報を送信する。そして、クラウドサーバ 7 は、例えば、診断システム 6 1、診断システム 6 2 及び診断システム 6 n からリアルタイム情報 8 を生成し、生成したリアルタイム情報 8 を政府の情報処理装置 4 に送信する。

【 0 0 2 5 】

これにより、病院 1 1、病院 1 2 及び病院 1 n は、政府の情報処理装置 4 に対して、最新の患者情報を迅速に送信することが可能になる。そのため、政府の担当者は、例えば、図 5 に示すユーザインターフェース U 3 を閲覧することにより、各患者についての診断結果を迅速に確認することが可能になる。

10

【 0 0 2 6 】

[診断モデルの具体例]

次に、診断システム 6 1、診断システム 6 2 及び診断システム 6 n において実行される診断モデル M D の具体例について説明を行う。図 6 は、診断モデル M D の具体例について説明する図である。

【 0 0 2 7 】

図 6 に示す診断モデル M D では、入力層から X 線画像が入力された場合、X 線画像が新型コロナウイルスに感染していないこと（肺炎になっていないこと）を示すカテゴリである「Normal」、または、X 線画像が新型コロナウイルスに感染している疑いがあること（肺炎になっている疑いがあること）を示すカテゴリである「Suspect」を出力する。

20

【 0 0 2 8 】

具体的に、図 6 に示す診断モデル M D は、例えば、ソフトマックス関数を用いることによって X 線画像の分類を行う。また、図 6 に示す診断モデル M D は、例えば、以下の式 (1) で表されるクロスエントロピー損失を損失関数 E として用いる。

【 0 0 2 9 】

【 数 1 】

$$E(\mathbf{w}_o, \mathbf{b}_o) = \frac{1}{N} \sum_{i=1}^N y_i \times \log(y'_i(\mathbf{w}_o, \mathbf{b}_o))$$

【 0 0 3 0 】

式 (1) において w_o 及び b_o は、診断モデル M D におけるパラメータであり、 y 及び y' のそれぞれは、実際のラベル（正解のラベル）及び予測されたラベルをそれぞれ示す。図 6 に示す診断モデル M D では、式 (1) を利用することによって、実際のラベルと予測されたラベルとの間における損失関数 E を算出する。そして、診断モデル M D では、確率的勾配縮小アルゴリズムや逆伝搬アルゴリズムを用いた損失関数 E の最小化が行われ、さらに、パラメータの最適化が行われる。

【 0 0 3 1 】

[診断モデルの更新]

次に、診断モデル M D の更新を行う処理（以下、更新処理とも呼ぶ）について説明を行う。図 7 は、診断モデル M D における更新処理について説明する図である。具体的に、図 7 (A) は、更新処理を説明するフローチャート図である。また、図 7 (B) は、更新処理の実行時における情報の送受信を説明する図である。なお、以下、図 2 で説明した診断システム 6 1 とクラウドサーバ 7 との間において行われる更新処理について説明を行う。

40

【 0 0 3 2 】

診断システム 6 1 は、病院 1 1 の患者から送信された X 線画像を含むデータセットを用いることによって、病院 1 1 における診断モデル M D の生成を予め行う。そして、診断システム 6 1 は、生成した診断モデル M D についてのパラメータ（例えば、以下の式 (2) に示す勾配 g_L ）をクラウドサーバ 7 に送信する (S 1)。

【 0 0 3 3 】

【 数 2 】

$$\nabla g_L = \frac{\delta E(W)}{\delta W}$$

【 0 0 3 4 】

続いて、クラウドサーバ7は、例えば、診断システム61、診断システム62及び診断システム6nのそれぞれからパラメータを受信したことに応じて、受信したパラメータのそれぞれからグローバルパラメータ（例えば、以下の式（3）に示すグローバル勾配 g_G ）を算出する（S2）。

10

【 0 0 3 5 】

さらに、クラウドサーバ7は、算出したグローバルパラメータを診断システム61、診断システム62及び診断システム6nのそれぞれに送信する（S3）。

【 0 0 3 6 】

【 数 3 】

$$\nabla g_G = \frac{1}{n} \sum_{i=1}^n \nabla g_L^i$$

【 0 0 3 7 】

その後、診断システム61は、クラウドサーバ7から送信されたグローバルパラメータを受信したことに応じて、以下の式（4）及び式（5）に示すように、病院11における診断モデルMDのパラメータを更新する（S4）。

20

【 0 0 3 8 】

【 数 4 】

$$W^{r+1} = W^r - \eta \nabla g_G$$

【 0 0 3 9 】

【 数 5 】

$$b^{r+1} = b^r - \eta \nabla g_G$$

【 0 0 4 0 】

式（4）及び式（5）において、 W^r 及び b^r のそれぞれは、r回目に行われたS1からS4の処理（r番目のトレーニングラウンド）における重み及びバイアスをそれぞれ示している。また、式（4）及び（5）において、 η は学習率を示している。

【 0 0 4 1 】

すなわち、情報処理システム100では、フェデレーテッドモデルラーニング（FML：Federated Model Learning）によって、各病院の診断モデルMD（診断モデルMDのパラメータ）の生成を行う。

【 0 0 4 2 】

これにより、情報処理システム100は、各病院の患者についての個人情報を含むデータセットを他の病院等の外部に送信することなく、各病院の診断モデルMDの精度を高めることが可能になる。そのため、各病院の診断モデルMDは、各患者のプライバシーを守りつつ、新型コロナウイルスの診断を精度良く行うことが可能になる。

40

【 0 0 4 3 】

なお、S1からS4の処理は、各病院の診断モデルMDのそれぞれの判定精度が必要な条件を上回るまで繰り返し行われるものであってよい。

【 0 0 4 4 】

また、各病院の診断モデルMDは、例えば、診断システム61、診断システム62及び診断システム6n以外のコンピュータ（例えば、後述する図8に示すホストコンピュータHC）において生成されるものであってもよい。

50

【 0 0 4 5 】

[診断モデルのマッピング]

次に、従来の A I プロセッサ (以下、A I チップとも呼ぶ) P R に対する診断モデル M D のマッピングを行う処理 (以下、マッピング処理とも呼ぶ) について説明を行う。

【 0 0 4 6 】

A I プロセッサ P R は、例えば、診断システム 6 1、診断システム 6 2 及び診断システム 6 n に搭載されたプロセッサである。そして、A I プロセッサ P R に対する診断モデル M D のマッピングが行われる場合、A I プロセッサ P R では、例えば、診断モデル M D を構成するニューロンを複数のグループにクラスタリングした後、A I プロセッサ P R に含まれる複数の演算コアのそれぞれに対する各グループのマッピングを行う。

10

【 0 0 4 7 】

しかしながら、従来のマッピング処理では、複数の演算コア間における通信コストについて考慮されていない場合がある。そのため、従来のマッピング処理では、マッピング処理が必要な時間内に終了しない場合があった。

【 0 0 4 8 】

また、ニューロンのクラスタリングと複数の演算コアへのマッピングは、一般的に、どちらも N P 困難な問題であり、多項式時間において最適に解決することができない場合がある。

【 0 0 4 9 】

そこで、本実施の形態における A I プロセッサ P R では、遺伝的アルゴリズムを用いることにより、複数の演算コア間における通信コストを抑えるように、診断モデル M D の各層を構成するニューロンのマッピングを行う。

20

【 0 0 5 0 】

[本実施の形態における A I プロセッサの具体例]

次に、本実施の形態における A I プロセッサ P R の具体例について説明を行う。図 8 は、本実施の形態における A I プロセッサ P R の具体例を説明する図である。以下、診断モデル M D がホストコンピュータ H C において生成されるものとして説明を行う。

【 0 0 5 1 】

図 8 に示す A I プロセッサ P R は、15 個の演算コアを有している。具体的に、図 8 に示す A I プロセッサ P R は、畳み込み層と対応付けられた 10 個の演算コア C と、全結合層と対応付けられた 3 個の演算コア F を有している。また、各演算コアは、それぞれルータ R と接続している。なお、以下、A I プロセッサ P R が 15 個の演算コアを有している場合について説明を行うが、A I プロセッサ P R は、これ以外の数の演算コアを有するものであってもよい。

30

【 0 0 5 2 】

また、図 8 に示す A I プロセッサ P R は、プーリング層やアクティベーション機能と対応付けられた 1 個の演算コア U と、各演算コアに対して重み係数を送信する 1 個の演算コア I / O を有している。なお、演算コア I / O は、例えば、他の A I プロセッサ P R との間の通信 (チップ間通信) や各 A I プロセッサ P R の制御を行うホストコンピュータ H C との通信を行うものであってもよい。

40

【 0 0 5 3 】

また、図 8 に示す A I プロセッサ P R は、A I プロセッサ P R に対する入力を記憶するオンチップメモリ M (以下、単にメモリ M と呼ぶ) を有する。演算コア C 等の各演算コアは、メモリ M から入力をロードして処理を実行する。そして、各演算コアは、次の層に対応する処理の実行を可能とするために、処理の実行に伴う各演算コアの出力をメモリ M に記憶する。

【 0 0 5 4 】

さらに、図 8 に示す A I プロセッサ P R は、演算コア及びメモリ M のそれぞれに対応するルータ R と、External DRAM (Dynamic Random Access Memory) とを有する。

50

【 0 0 5 5 】

[本実施の形態におけるマッピング処理]

次に、本実施の形態におけるマッピング処理について説明を行う。図 9 は、本実施の形態におけるマッピング処理について説明するフローチャート図である。図 1 0 から図 2 1 は、本実施の形態におけるマッピング処理について説明する図である。

【 0 0 5 6 】

A I プロセッサ P R の演算コア U は、診断モデル M D (ニューラルネットワーク) を構成するニューロンについての N 個のマッピング解をランダムに決定する (S 1 1) 。

【 0 0 5 7 】

具体的に、A I プロセッサ P R がそれぞれ 8 個のニューロンと対応付けることが可能な 4 個の演算コアを有しており、かつ、診断モデル M D を構成するニューロンの数が 3 0 個である場合、マッピング解の組合せは、3 2 個のニューロンにおいて 3 0 個のニューロンを配置する組合せである 4 9 6 通りになる。そのため、演算コア U は、この場合、その 4 9 6 通りのうちの N 通りに対応する N 個のマッピング解をランダムに決定する。

10

【 0 0 5 8 】

以下、A I プロセッサ P R がそれぞれ 8 個のニューロンと対応付けることが可能な 4 個の演算コアを有しており、かつ、診断モデル M D を構成するニューロンの数が 3 0 個であるものとして説明を行う。さらに、診断モデル M D が 1 0 個のニューロンからなる層 L 1、層 L 2 及び層 L 3 をそれぞれ有するニューラルネットワークであるものとして説明を行う。

20

【 0 0 5 9 】

[マッピング結果の具体例]

次に、本実施の形態におけるマッピング結果の具体例について説明を行う。図 1 0 は、各演算コアにマッピングされたニューロンの数を示す図である。また、図 1 1 は、各演算コアにマッピングされたニューロンの識別情報を示す図である。

【 0 0 6 0 】

図 1 0 及び図 1 1 に示す例は、1 行目 1 列目の演算コアと対応付けられているニューロンが、層 L 1 に対応する 3 つのニューロン (ニューロン 2、3 及び 5) と、層 L 2 に対応する 2 つのニューロン (ニューロン 1 1 及び 1 2) と、層 L 3 に対応する 3 つのニューロン (ニューロン 2 2、2 3 及び 2 8) とであることを示している。

30

【 0 0 6 1 】

また、図 1 0 及び図 1 1 に示す例は、1 行目 2 列目の演算コアと対応付けられているニューロンが、層 L 1 に対応する 3 つのニューロン (ニューロン 6、7 及び 8) と、層 L 2 に対応する 1 つのニューロン (ニューロン 1 3) と、層 L 3 に対応する 3 つのニューロン (ニューロン 2 1、2 9 及び 3 0) とであることを示している。なお、図 1 0 及び図 1 1 に示す例は、1 行目 2 列目の演算コアにおいてさらに対応可能なニューロンの数 (F) が 1 であることを示している。

【 0 0 6 2 】

また、図 1 0 及び図 1 1 に示す例は、2 行目 1 列目の演算コアと対応付けられているニューロンが、層 L 1 に対応する 2 つのニューロン (ニューロン 1 及び 4) と、層 L 2 に対応する 3 つのニューロン (ニューロン 1 5、1 7 及び 1 9) と、層 L 3 に対応する 3 つのニューロン (ニューロン 2 4、2 5 及び 2 6) とであることを示している。

40

【 0 0 6 3 】

また、図 1 0 及び図 1 1 に示す例は、2 行目 2 列目の演算コアと対応付けられているニューロンが、層 L 1 に対応する 2 つのニューロン (ニューロン 9 及び 1 0) と、層 L 2 に対応する 4 つのニューロン (ニューロン 1 4、1 6、1 8 及び 2 0) と、層 L 3 に対応する 1 つのニューロン (ニューロン 2 7) とであることを示している。なお、図 1 0 及び図 1 1 に示す例は、2 行目 2 列目の演算コアにおいてさらに対応可能なニューロンの数 (F) が 1 であることを示している。

【 0 0 6 4 】

50

図 9 に戻り、演算コア U は、S 1 1 の処理で決定した N 個のマッピング解から不適切なマッピング解を削除する (S 1 2)。

【 0 0 6 5 】

具体的に、演算コア U は、例えば、図 1 2 に示すように、9 個のニューロンと対応付けられた演算コア (1 行目 1 列目の演算コア) が存在するマッピング解を削除する。また、演算コア U は、例えば、図 1 3 に示すように、同一のニューロン (ニューロン 8) が複数の演算コアと対応付けられたマッピング解を削除する。

【 0 0 6 6 】

そして、演算コア U は、S 1 2 の処理で削除されなかったマッピング解のそれぞれに対応する通信コストを算出する (S 1 3)。

【 0 0 6 7 】

具体的に、各マッピング解の通信コストは、以下の式 (6) によって表現される。

【 0 0 6 8 】

【 数 6 】

$$F_{cost} = \sum_{i=0, j=0}^W d_{i,j} \times c_{i,j}$$

【 0 0 6 9 】

式 (6) において、 $d_{i,j}$ は、ニューロン i とニューロン j との間の距離を示しており、 $c_{i,j}$ は、ニューロン i とニューロン j との間の接続状況を示している。

【 0 0 7 0 】

具体的に、 $d_{i,j}$ は、例えば、ニューロン i と対応付けられている演算コアとニューロン j と対応付けられている演算コアとの間に存在するルータ R の数に 1 を加算した値である。また、 $c_{i,j}$ は、例えば、ニューロン i とニューロン j とが直接接続している場合に 1 になり、ニューロン i とニューロン j とが直接接続していない場合に 0 になる。

【 0 0 7 1 】

さらに具体的に、図 1 1 で説明した例において、ニューロン 1 及びニューロン 2 は、それぞれ層 L 1 と対応付けられている。そのため、この場合、 $d_{1,2}$ は 1 になり、 $c_{1,2}$ は 0 になる。また、図 1 1 で説明した例において、ニューロン 1 は、層 L 1 と対応付けられており、ニューロン 1 4 は、層 L 2 と対応付けられている。そのため、この場合、 $d_{1,14}$ は 4 になり、 $c_{1,2}$ は 1 になる。

【 0 0 7 2 】

続いて、演算コア U は、S 1 2 の処理で削除されなかったマッピング解のうち、S 1 3 の処理で算出した通信コストが条件を満たす M 個のマッピング解を特定する (S 1 4)。

【 0 0 7 3 】

具体的に、演算コア U は、例えば、S 1 2 の処理で削除されなかったマッピング解から、S 1 3 の処理で算出した通信コストが高い順に M 個のマッピング解を特定する。

【 0 0 7 4 】

その後、演算コア U は、S 1 4 の処理で特定した M 個のマッピング解を交差 (クロスオーバー) させることによって N - M 個の新たなマッピング解を決定する (S 1 5)。

【 0 0 7 5 】

具体的に、演算コア U は、例えば、図 1 4 及び図 1 5 に示す親 1 及び親 2 が S 1 4 の処理で特定した M 個のマッピング解に含まれている場合、図 1 6 に示すように、親 1 及び親 2 のそれぞれの割合を 5 0 (%) とした新たな子孫を作成する。

【 0 0 7 6 】

さらに具体的に、図 1 4 に示す例において、親 1 の場合の 1 行目 1 列目の演算コアには、層 L 1 に対応する 3 つのニューロンと、層 L 2 に対応する 2 つのニューロンと、層 L 3 に対応する 3 つのニューロンとが対応付けられている。また、図 1 5 に示す例において、親 2 の場合の 1 行目 1 列目の演算コアには、層 L 1 に対応する 1 つのニューロンと、層 L 2 に対応する 4 つのニューロンと、層 L 3 に対応する 3 つのニューロンとが対応付けられ

10

20

30

40

50

ている。そのため、この場合、新たな子孫の場合の1行目1列目の演算コアには、図16に示すように、層L1に対応する $2(3 * 0.5 + 1 * 0.5)$ つのニューロンと、層L2に対応する $3(4 * 0.5 + 2 * 0.5)$ つのニューロンと、層L3に対応する $3(3 * 0.5 + 3 * 0.5)$ つのニューロンとが対応付けられる。

【0077】

なお、例えば、図17に示すように、各演算コアと対応付けられたニューロンの数に小数が含まれている場合、演算コアUは、図18に示すように、各演算コアと対応付けられたニューロンの数のそれぞれが整数になるように調整を行うものであってもよい。

【0078】

図9に戻り、演算コアUは、N個のマッピング解(S14の処理で特定したM個のマッピング解とS15の処理で決定したN-M個の新たなマッピング解との合計)において突然変異を発生させる(S16)。

【0079】

具体的に、演算コアUは、図19に示すように、例えば、図18に示す複数の演算コアから1行目1列目の演算コアと2行目2列目の演算コアとを特定し、さらに、診断モデルMDを構成する複数の層から層L1と層L2とを特定する。そして、演算コアUは、例えば、特定した1行目1列目の演算コアにおける層L1に対応付けられたニューロンの数である2と、2行目2列目の演算コアにおける層L2に対応付けられたニューロンの数である3とのうちの最小値である2を特定する。その後、演算コアUは、例えば、1行目1列目の演算コアにおける層L1に対応付けられたニューロンの数である2から、最小値として特定した値である2を減算し、さらに、1行目1列目の演算コアにおける層L2に対応付けられたニューロンの数である3に対して、最小値として特定した値である2を加算する。

【0080】

同様に、演算コアUは、例えば、2行目2列目の演算コアにおける層L2に対応付けられたニューロンの数である3から2を減算し、さらに、2行目2列目の演算コアにおける層L1に対応付けられたニューロンの数である4に対して2を加算する。

【0081】

そして、演算コアUは、例えば、S12の処理と同様に、S16の処理が行われた後のN個のマッピング解が制約を満たしているか否かを判定する。その結果、制約を満たしていないマッピング解が存在した場合、演算コアUは、存在したマッピング解を削除する。さらに、演算コアUは、S13の処理と同様に、削除されなかったマッピング解の通信コストを算出する(S17)。

【0082】

その後、S17の処理で算出した通信コストのうちの最適なコスト(最小のコスト)が予め定められた条件を満たしているか否かを判定する(S18)。

【0083】

その結果、S17の処理で算出した通信コストのうちの最適なコストが予め定められた条件を満たしていると判定した場合、演算コアUは、マッピング処理を終了(正常終了)する。

【0084】

一方、S17の処理で算出した通信コストのうちの最適なコストが予め定められた条件を満たしていないと判定した場合、演算コアUは、例えば、S12以降の処理の実行回数(すなわち、世代数)が予め定められた所定回数に到達したか否かを判定する(S19)。

【0085】

その結果、S12以降の処理の実行回数が予め定められた所定回数に到達していないと判定した場合、演算コアUは、例えば、S12以降の処理を再度行う。

【0086】

一方、S12以降の処理の実行回数が予め定められた所定回数に到達していると判定し

10

20

30

40

50

た場合、演算コアUは、例えば、ニューロンのマッピングを終了（異常終了）する。

【0087】

このように、本実施の形態におけるAIプロセッサPRは、複数の演算コアを有し、複数の演算コアの少なくともいずれかが、畳み込み層と全結合層とを有する診断モデルMDに含まれる複数のニューロンのそれぞれに対応付けられた計算プログラムを分割して複数の演算コアのそれぞれに割り当てるマッピング処理を実行する。

【0088】

そして、本実施の形態におけるAIプロセッサは、複数の演算コアのそれぞれが、マッピング処理によって割り当てられた計算プログラムを実行する。

【0089】

具体的に、本実施の形態におけるマッピング処理では、複数のコア間における通信コストが所定の閾値以下になるように、遺伝的アルゴリズムによって計算プログラムを複数の演算コアに割り当てる。

【0090】

さらに具体的に、演算コアUは、マッピング処理が正常終了した場合、マッピング処理の結果を示すマッピングテーブル（図示しない）を生成する。続いて、演算コアUは、ホストコンピュータHCから診断モデルMDのパラメータをダウンロードする。さらに、演算コアUは、マッピングテーブルを参照し、診断モデルMDのパラメータを演算コアC及び演算コアFのそれぞれに送信する。また、演算コアUは、マッピングテーブルを各ルータRに対しても送信する。

【0091】

その後、演算コアUは、例えば、入力データ（例えば、患者のX線画像）が入力された場合、マッピングテーブルを参照し、入力データを最初の層に含まれるニューロンに対応する演算コアのそれぞれに送信する。そして、各ルータRは、最初の層に対応する処理が完了したことに応じて、最初の層からの出力データを次の層に含まれるニューロンに対応する演算コアに送信する。さらに、各ルータRは、最後の層に対応する処理が完了するまでの間、処理対象の層の次の層に含まれるニューロンに対応する演算コアに対する送信を繰り返し行う。そして、各ルータRは、最後の層からの出力データ（診断モデルMDの出力データ）をDRAMに格納する。

【0092】

すなわち、本実施の形態におけるAIプロセッサPRは、従来よりも簡易なアルゴリズムであって、かつ、予測可能な時間内に結果を得ることが可能なアルゴリズムである遺伝的アルゴリズムを用いることによってマッピング処理を行う。

【0093】

これにより、本実施の形態におけるAIプロセッサPRは、各ニューロンの演算コアへのマッピングが行われる際の通信コストを抑えることが可能になる。そのため、本実施の形態におけるAIプロセッサPRは、診断モデルを用いた診断をより効率的に行うことが可能になる。

【符号の説明】

【0094】

- 4：情報処理装置
- 7：クラウドサーバ
- 8：リアルタイム情報
- 11：病院
- 12：病院
- 1n：病院
- 51：携帯端末
- 52：携帯端末
- 5n：携帯端末
- 61：診断システム

10

20

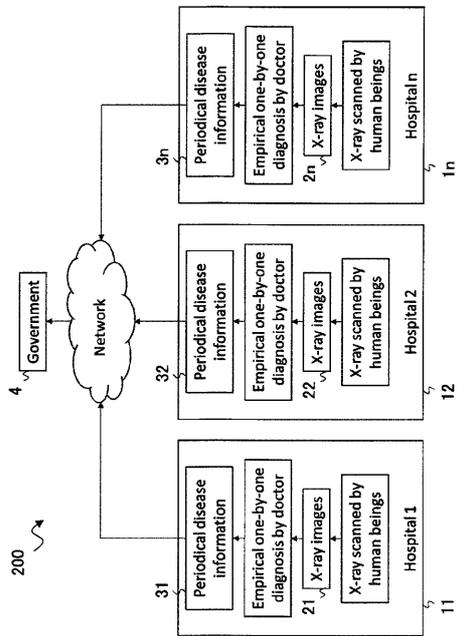
30

40

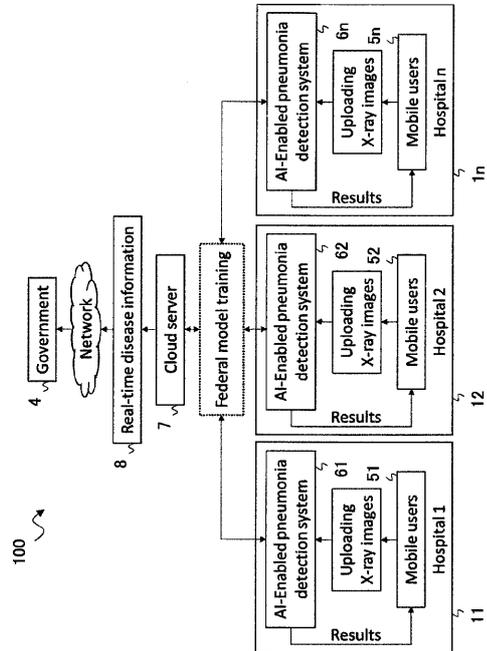
50

6 2 : 診断システム
 6 n : 診断システム
 1 0 0 : 情報処理システム
 M D : 診断モデル
 P R : A I プロセッサ

【 図 1 】



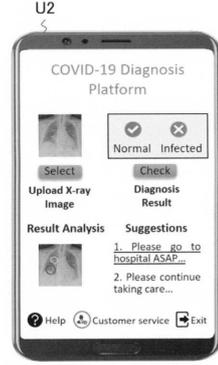
【 図 2 】



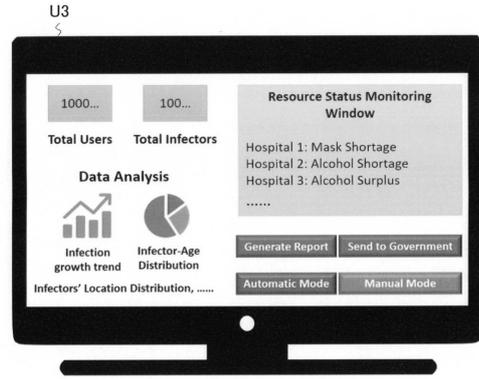
【 図 3 】



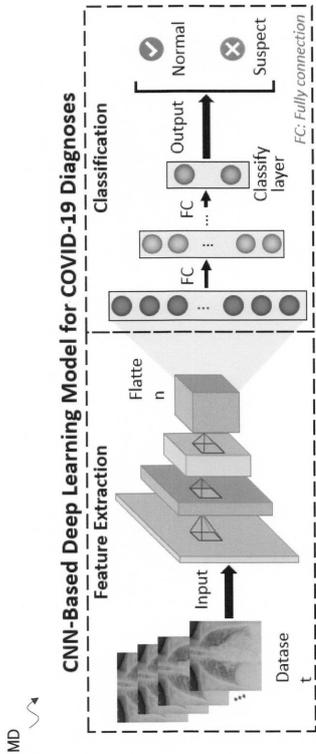
【 図 4 】



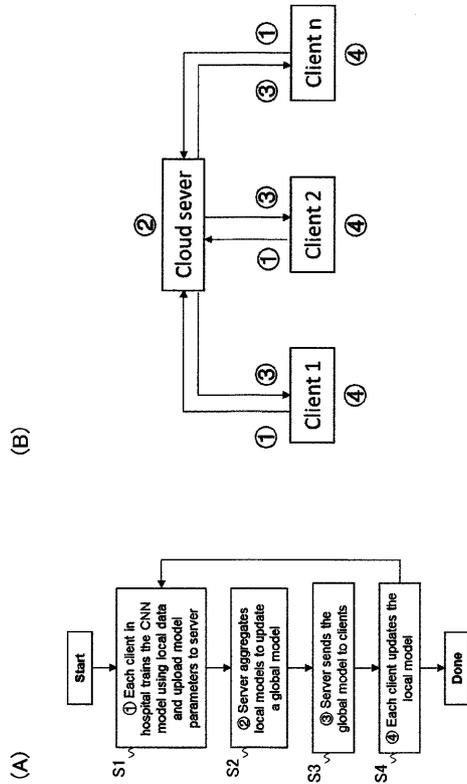
【 図 5 】



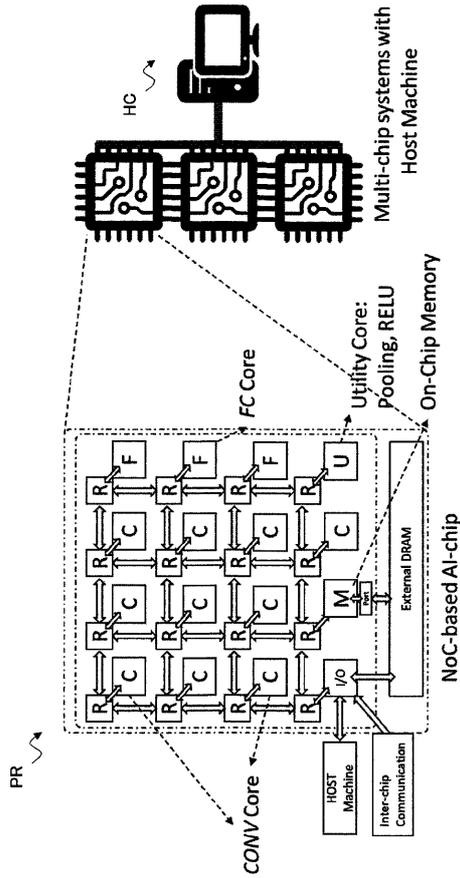
【 図 6 】



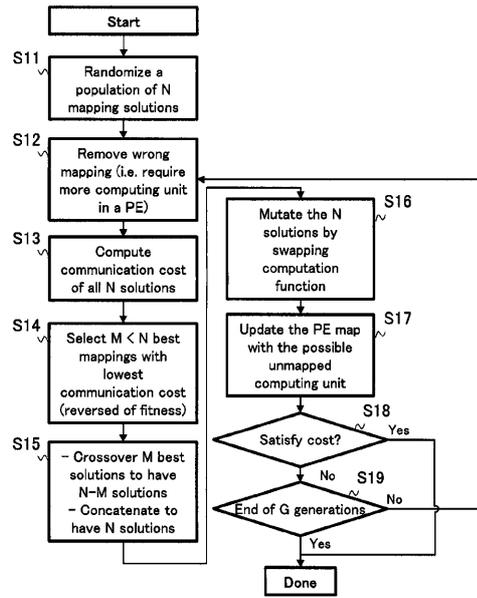
【 図 7 】



【 図 8 】



【 図 9 】



【 図 1 0 】

Mapping-1		Column	
		1	2
Row	1	L1 : 3 neurons L2 : 2 neurons L3 : 3 neurons	L1 : 3 neurons L2 : 1 neurons L3 : 3 neurons F : 1
	2	L1 : 2 neurons L2 : 3 neurons L3 : 3 neurons	L1 : 2 neurons L2 : 4 neurons L3 : 1 neurons F : 1

【 図 1 2 】

Mapping-1		Column	
		1	2
Row	1	L1 : 2, 3, 5, 8 L2 : 11, 12 L3 : 22, 23, 28	L1 : 6, 7 L2 : 13 L3 : 21, 29, 30 F : 1
	2	L1 : 1, 4 L2 : 15, 17, 19 L3 : 24, 25, 26	L1 : 9, 10 L2 : 14, 16, 18, 20 L3 : 27 F : 1

【 図 1 1 】

Mapping-1		Column	
		1	2
Row	1	L1 : 2, 3, 5 L2 : 11, 12 L3 : 22, 23, 28	L1 : 6, 7, 8 L2 : 13 L3 : 21, 29, 30 F : 1
	2	L1 : 1, 4 L2 : 15, 17, 19 L3 : 24, 25, 26	L1 : 9, 10 L2 : 14, 16, 18, 20 L3 : 27 F : 1

【 図 1 3 】

Mapping-1		Column	
		1	2
Row	1	L1 : 2, 3, 8 L2 : 11, 12 L3 : 22, 23, 28	L1 : 6, 7, 8 L2 : 13 L3 : 21, 29, 30 F : 1
	2	L1 : 1, 4 L2 : 15, 17, 19 L3 : 24, 25, 26	L1 : 9, 10 L2 : 14, 16, 18, 20 L3 : 27 F : 1

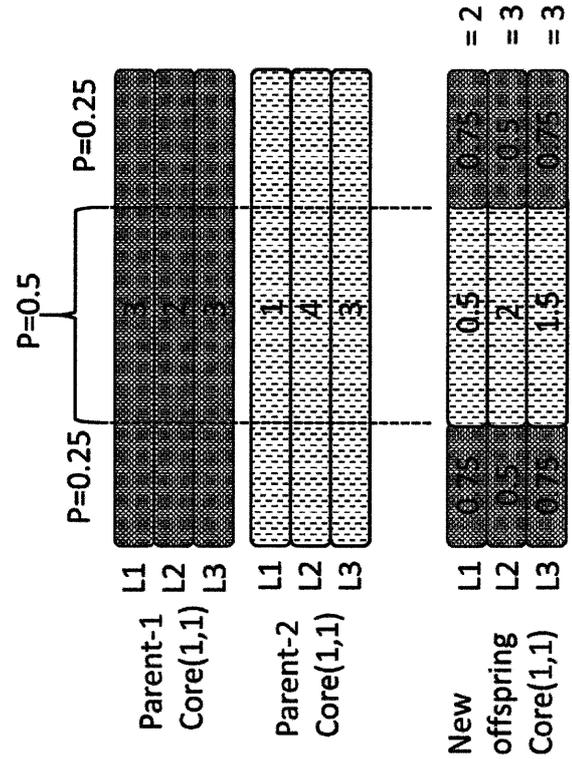
【 図 1 4 】

Parent-1		Column	
		1	2
Row	1	L1 : 3 neurons L2 : 2 neurons L3 : 3 neurons	L1 : 3 neurons L2 : 1 neurons L3 : 3 neurons F : 1
	2	L1 : 2 neurons L2 : 3 neurons L3 : 3 neurons	L1 : 2 neurons L2 : 4 neurons L3 : 1 neurons F : 1

【 図 1 5 】

Parent-2		Column	
		1	2
Row	1	L1 : 1 neurons L2 : 4 neurons L3 : 3 neurons	L1 : 0 neurons L2 : 3 neurons L3 : 3 neurons F : 2
	2	L1 : 4 neurons L2 : 1 neurons L3 : 3 neurons	L1 : 5 neurons L2 : 2 neurons L3 : 1 neurons

【 図 1 6 】



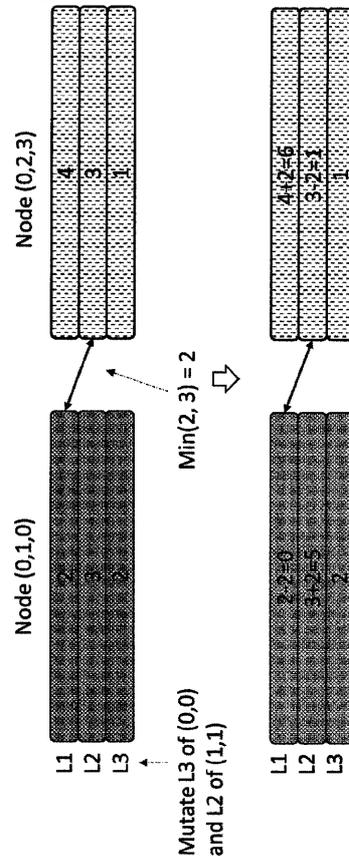
【 図 1 7 】

Parent-2		Column	
		1	2
Row	1	L1 : 2 neurons L2 : 3 neurons L3 : 3 neurons	L1 : 1.5 neurons L2 : 2 neurons L3 : 3 neurons F : 1.5
	2	L1 : 3 neurons L2 : 2 neurons L3 : 3 neurons	L1 : 3.5 neurons L2 : 3 neurons L3 : 1 neurons F : 0.5

【 図 1 8 】

Parent-2		Column	
		1	2
Row	1	L1 : 2 neurons L2 : 3 neurons L3 : 3 neurons	L1 : 1 neurons L2 : 2 neurons L3 : 3 neurons F : 2
	2	L1 : 3 neurons L2 : 2 neurons L3 : 3 neurons	L1 : 4 neurons L2 : 3 neurons L3 : 1 neurons F : 0

【 図 1 9 】



【 2 0 】

Parent-2		Column	
		1	2
Row	1	L1 : 0 neurons L2 : 5 neurons L3 : 3 neurons	L1 : 1 neurons L2 : 2 neurons L3 : 3 neurons F : 2
	2	L1 : 3 neurons L2 : 2 neurons L3 : 3 neurons	L1 : 6 neurons L2 : 1 neurons L3 : 1 neurons F : 0

【 2 1 】

Parent-2		Column	
		1	2
Row	1	L2 : 11, 12, 13, 17, 18 L3 : 28, 29, 30	L1 : 10 L2 : 15, 16 L3 : 25, 26, 27 F : 2
	2	L1 : 4, 5, 6 L2 : 14, 20 L3 : 21, 22, 23	L1 : 1, 2, 3, 7, 8, 9 L2 : 19 L3 : 24

フロントページの続き

(72)発明者 ダン ナム カイン

ベトナム社会主義共和国 ハノイ カウ ザイ ディストリクト スアン トゥイ ストリート
144 ベトナム国家大学ハノイ校内

(72)発明者 ジャンニン ソン

オーストラリア連邦 3800 ビクトリア クレイトン ウェリントン ロード モナシュ大学
内

審査官 大倉 峻吾

(56)参考文献 米国特許出願公開第2020/0089534 (US, A1)

特開2004-185271 (JP, A)

中国特許出願公開第111652863 (CN, A)

国際公開第2020/053887 (WO, A1)

特開2020-046821 (JP, A)

特開2020-160564 (JP, A)

国際公開第2019/220692 (WO, A1)

国際公開第2020/225879 (WO, A1)

特開2003-058520 (JP, A)

樽林亮介 ほか, "データ駆動プロセッサによる実時間処理のためのプログラム割当手法", 電子
情報通信学会論文誌, Vol. J86-D-I, No. 10, 2003年10月, p. 721-732

(58)調査した分野(Int.Cl., DB名)

G06N 3/00 - 99/00

G06F 9/50

G06F 9/54