

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7640552号
(P7640552)

(45)発行日 令和7年3月5日(2025.3.5)

(24)登録日 令和7年2月25日(2025.2.25)

(51)国際特許分類

F I

H 0 4 N 19/126(2014.01)

H 0 4 N 19/13 (2014.01)

H 0 4 N 19/463(2014.01)

H 0 4 N 19/126

H 0 4 N 19/13

H 0 4 N 19/463

請求項の数 46 (全77頁)

(21)出願番号	特願2022-538077(P2022-538077)	(73)特許権者	591037214
(86)(22)出願日	令和2年12月21日(2020.12.21)		フラウンホッファー - ゲゼルシャフト
(65)公表番号	特表2023-507502(P2023-507502 A)		ツァ フェルダールング デア アンゲヴ
(43)公表日	令和5年2月22日(2023.2.22)		アンテン フォアシュンク エー . ファオ
(86)国際出願番号	PCT/EP2020/087489		ドイツ連邦共和国 8 0 6 8 6 ミュンヘ
(87)国際公開番号	WO2021/123438	(74)代理人	100079577
(87)国際公開日	令和3年6月24日(2021.6.24)		弁理士 岡田 全啓
審査請求日	令和4年10月19日(2022.10.19)	(72)発明者	ハーセ ボール
(31)優先権主張番号	19218862.1		ドイツ連邦共和国 1 0 5 8 7 ベルリン
(32)優先日	令和1年12月20日(2019.12.20)		アインシュタインウーファー 3 7 フラ
(33)優先権主張国・地域又は機関	欧州特許庁(EP)		ウンホッファー - インスティチュート
前置審査			フュア ナーハリヒテンデヒニーク ハイ
			ンリッヒ - ヘルツ - インスティチュート
			H H I 内
			最終頁に続く

(54)【発明の名称】 ニューラルネットワークのパラメータを符号化するための概念

(57)【特許請求の範囲】

【請求項 1】

ニューラルネットワーク (1 0) を定義するニューラルネットワークパラメータ (1 3) をデータストリーム (1 4) から復号化するための装置であって、前記装置は、
前記データストリーム (1 4) から復号化された、以前のニューラルネットワークパラメータのための量子化インデックス (5 8) に依存して、現在のニューラルネットワークパラメータ (1 3 ') について、複数 (5 0) の再構成レベルセット (5 2) の中から再構成レベルセット (4 8) を選択 (5 4) することと、
前記現在のニューラルネットワークパラメータについての選択された前記再構成レベルのセット (4 8) のうちの 1 つの再構成レベルを示す、前記現在のニューラルネットワークパラメータ (1 3 ') のための量子化インデックス (5 6) を前記データストリーム (1 4) から復号化することと、

前記現在のニューラルネットワークパラメータのための前記量子化インデックス (5 6) によって示される前記選択された再構成レベルセット (4 8) のうちの前記 1 つの再構成レベルに前記現在のニューラルネットワークパラメータ (1 3 ') を逆量子化 (6 2) すること

により、前記ニューラルネットワークパラメータ (1 3 ') を順次復号化するように構成され、
前記装置は、

前記現在のニューラルネットワークパラメータ (1 3 ') に関連づけられた状態に応じて、

前記現在のニューラルネットワークパラメータ（１３'）について、前記複数（５０）の再構成レベルセット（５２）の中から前記再構成レベルセット（４８）を決定すること、及び、

前記データストリームから復号化された、直前のニューラルネットワークパラメータのための前記量子化インデックス（５８）に応じて、後続のニューラルネットワークパラメータの状態を更新すること

による状態遷移プロセスにより、前記現在のニューラルネットワークパラメータ（１３'）について、前記複数（５０）の再構成レベルセット（５２）の中から前記再構成レベルセット（４８）を選択（５４）し、

前記現在のニューラルネットワークパラメータ（１３'）の状態および以前に復号化されたニューラルネットワークパラメータの前記量子化インデックスに依存する（１２２）確率モデルを用いる算術符号化を用いて、前記データストリーム（１４）から前記現在のニューラルネットワークパラメータ（１３'）のための前記量子化インデックス（５６）を復号化し、

前記現在のニューラルネットワークパラメータ（１３'）に関連付けられた前記状態に依存して、複数の確率モデルの中から確率モデルのサブセットを事前選択し、以前に復号化されたニューラルネットワークパラメータのための前記量子化インデックスに依存して（１２１）、前記確率モデルのサブセットの中から前記現在のニューラルネットワークパラメータのための前記確率モデルを選択し、

前記現在のニューラルネットワークパラメータが関連する部分に隣接する前記ニューラルネットワークの部分に関連する、以前に復号化されたニューラルネットワークパラメータの前記量子化インデックスの値に依存して、前記現在のニューラルネットワークパラメータのための確率モデルを第１、第２および第３の確率モデルを含むサブセットの中から選択する

ように構成され、

ここで、前記以前に復号化されたニューラルネットワークパラメータがゼロよりも小さければ、第１の確率モデルが選択され、

前記以前に復号化されたニューラルネットワークパラメータがゼロよりも大きければ、第２の確率モデルが選択され、

前記以前に復号化されたニューラルネットワークパラメータがゼロに等しければ、前記第３の確率モデルが選択される、

装置。

【請求項２】

前記ニューラルネットワークパラメータ（１３）が、前記ニューラルネットワーク（１０）のニューロン相互接続（１１）の重みに関連する、請求項１に記載の装置。

【請求項３】

前記複数（５０）の再構成レベルセット（５２）のうちの再構成レベルセット（５２）の数は２である、請求項１又は２に記載の装置。

【請求項４】

前記複数（５０）の再構成レベルセット（５２）を既定の量子化ステップサイズ（ＱＰ）によりパラメータ化（６０）し、前記データストリーム（１４）から前記既定の量子化ステップサイズに関する情報を導出するように構成される、請求項１乃至３のいずれかに記載の装置。

【請求項５】

前記ニューラルネットワークは１つ以上のＮＮ層を含み、前記装置は、
それぞれのＮＮ層（ p ； $p-1$ ）について、前記それぞれのＮＮ層のための既定の量子化ステップサイズに関する情報を前記データストリーム（１４）から導出し、
それぞれのＮＮ層について、前記それぞれのＮＮ層について導出された前記既定の量子化ステップサイズを使用して前記複数（５０）の再構成レベルセット（５２）をパラメータ化して、前記それぞれのＮＮ層に属する前記ニューラルネットワークパラメータの逆量

10

20

30

40

50

子化のために使用されるようにする

ように構成される、請求項 1 乃至 4 のいずれかに記載の装置。

【請求項 6】

前記複数 (50) の再構成レベルセット (52) のうちの再構成レベルセット (52) の数は 2 であり、前記複数の再構成レベルセットは、

ゼロと既定の量子化ステップサイズの偶数倍とを含む第 1 の再構成レベルセット (セット 0) と、

ゼロと前記既定の量子化ステップサイズの奇数倍とを含む第 2 の再構成レベルセット (セット 1) と

を含む、請求項 1 乃至 5 のいずれかに記載の装置。

10

【請求項 7】

すべての再構成レベルセットのすべての再構成レベルは、既定の量子化ステップサイズの整数倍を表し、前記装置は、

それぞれのニューラルネットワークパラメータについて、前記それぞれのニューラルネットワークパラメータについての前記選択された再構成レベルセットと、前記それぞれのニューラルネットワークパラメータのための前記量子化インデックスをエントロピー復号化したものに応じて中間整数値を導出すること、及び、

それぞれのニューラルネットワークパラメータについて、前記それぞれのニューラルネットワークパラメータについての前記中間整数値を、前記それぞれのニューラルネットワークパラメータのための前記既定の量子化ステップサイズで乗算すること

20

により、前記ニューラルネットワークパラメータを逆量子化するように構成される、請求項 1 乃至 6 のいずれかに記載の装置。

【請求項 8】

前記複数 (50) の再構成レベルセット (52) のうちの再構成レベルセット (52) の数は 2 であり、前記装置は、

前記それぞれのニューラルネットワークパラメータについての前記選択された再構成レベルセットが第 1 のセットである場合、前記それぞれのニューラルネットワークパラメータのための前記量子化インデックスを 2 倍して、前記それぞれのニューラルネットワークパラメータについての前記中間整数値を得ること、及び、

それぞれのニューラルネットワークパラメータについての前記選択された再構成レベルセットが第 2 のセットであり、かつ前記それぞれのニューラルネットワークパラメータのための前記量子化インデックスがゼロに等しい場合、前記それぞれのニューラルネットワークパラメータについての前記中間整数値をゼロに等しく設定すること、及び、

30

それぞれのニューラルネットワークパラメータについての前記選択された再構成レベルセットが第 2 のセットであり、かつ前記それぞれのニューラルネットワークパラメータのための前記量子化インデックスがゼロより大きい場合、前記それぞれのニューラルネットワークパラメータのための前記量子化インデックスを 2 倍し、該乗算の結果から 1 を引いて前記それぞれのニューラルネットワークパラメータについての前記中間整数値を得ること、及び、

現在のニューラルネットワークパラメータについての前記選択された再構成レベルセットが第 2 のセットであり、前記それぞれのニューラルネットワークパラメータのための前記量子化インデックスがゼロより小さい場合、前記それぞれのニューラルネットワークパラメータのための前記量子化インデックスを 2 倍し、該乗算の結果に 1 を加えて、前記それぞれのニューラルネットワークパラメータについての前記中間整数値を得ることにより、それぞれのニューラルネットワークパラメータについての前記中間整数値を導出するように構成される、請求項 7 に記載の装置。

40

【請求項 9】

前記現在のニューラルネットワークパラメータ (13') について、前記データストリーム (14) から復号化された、以前に復号化されたニューラルネットワークパラメータのための前記量子化インデックス (58) を 2 値化したものの LSB 部分又は以前に復号化

50

されたピンに応じて、前記複数（５０）の再構成レベルセット（５２）の中から前記再構成レベルセット（４８）を選択する（５４）ように構成される、請求項１乃至８のいずれかに記載の装置。

【請求項１０】

前記現在のニューラルネットワークパラメータ（１３'）について、前記データストリーム（１４）から復号化された、以前に復号化されたニューラルネットワークパラメータのための前記量子化インデックス（５８）の２値関数の結果に応じて、前記複数（５０）の再構成レベルセット（５２）の中から前記再構成レベルセット（４８）を選択する（５４）ように構成される、請求項１乃至８のいずれかに記載の装置。

【請求項１１】

前記装置は、前記現在のニューラルネットワークパラメータ（１３'）について、前記データストリーム（１４）から復号化された、以前に復号化されたニューラルネットワークパラメータのための前記量子化インデックス（５８）のパリティに応じて、前記複数（５０）の再構成レベルセット（５２）の中から前記再構成レベルセット（４８）を選択する（５４）ように構成される、請求項１乃至１０のいずれかに記載の装置。

【請求項１２】

前記複数（５０）の再構成レベルセット（５２）のうちの前記再構成レベルセット（５２）の数は２であり、前記装置は、

それぞれのニューラルネットワークパラメータのためのサブセットインデックスを、前記それぞれのニューラルネットワークパラメータについての前記選択された再構成レベルセットと、前記それぞれのニューラルネットワークパラメータのための前記量子化インデックスの２値関数とに基づいて導き出し、前記サブセットインデックスについての４つの可能な値を生じさせ、

前記現在のニューラルネットワークパラメータ（１３'）について、以前に復号化されたニューラルネットワークパラメータのための前記サブセットインデックスに応じて、前記複数（５０）の再構成レベルセット（５２）の中から前記再構成レベルセット（４８）を選択する（５４）

ように構成される、請求項１乃至１１のいずれかに記載の装置。

【請求項１３】

前記装置は、前記現在のニューラルネットワークパラメータ（１３'）について、複数の直前に復号化されたニューラルネットワークパラメータのための前記サブセットインデックスに依存する選択ルールを用いて、前記複数（５０）の再構成レベルセット（５２）の中から前記再構成レベルセット（４８）を選択し（５４）、前記選択ルールを前記ニューラルネットワークパラメータの一部または全部に対して使用するよう構成される、請求項１２に記載の装置。

【請求項１４】

前記選択ルールが依存する、前記直前に復号化されたニューラルネットワークパラメータの数は２である、請求項１３に記載の装置。

【請求項１５】

それぞれのニューラルネットワークパラメータのための前記サブセットインデックスは、前記それぞれのニューラルネットワークパラメータについての前記選択された再構成レベルセットと、前記それぞれのニューラルネットワークパラメータのための前記量子化インデックスのパリティとに基づいて導出される、請求項１２乃至１４のいずれかに記載の装置。

【請求項１６】

前記装置は、

前記現在のニューラルネットワークパラメータ（１３'）について、前記現在のニューラルネットワークパラメータ（１３'）に関連付けられた状態に応じて、前記複数（５０）の再構成レベルセット（５２）の中から前記再構成レベルセット（４８）を決定すること、及び、

10

20

30

40

50

前記データストリームから復号化された、前記直前のニューラルネットワークパラメータのための前記量子化インデックス（５８）に応じて、後続のニューラルネットワークパラメータの状態を更新すること
による状態遷移プロセスにより、前記現在のニューラルネットワークパラメータ（１３'）について、前記複数（５０）の再構成レベルセット（５２）の中から前記再構成レベルセット（４８）を選択する（５４）ように構成される、請求項１乃至１５のいずれかに記載の装置。

【請求項１７】

前記データストリームから復号化された、前記直前のニューラルネットワークパラメータのための前記量子化インデックス（５８）の２値関数を用いて、前記後続の前記ニューラルネットワークパラメータの状態を更新するように構成される、請求項１６に記載の装置。

10

【請求項１８】

前記データストリームから復号化された、前記直前のニューラルネットワークパラメータのための前記量子化インデックス（５８）のパリティを使用して、前記後続のニューラルネットワークパラメータの状態を更新するように構成される、請求項１６に記載の装置。

【請求項１９】

前記状態遷移プロセスは、４つ又は８つの可能な状態の間で遷移するように構成される、請求項１６乃至１８のいずれかに記載の装置。

【請求項２０】

20

前記状態遷移プロセスにおいて偶数の可能な状態の間で遷移し、前記複数（５０）の再構成レベルセット（５２）のうちの再構成レベルセット（５２）の数が２であるように構成され、ここで、前記現在のニューラルネットワークパラメータ（１３'）について、前記現在のニューラルネットワークパラメータ（１３'）に関連づけられた前記状態に応じて、前記再構成レベルセット（５２）のうちの前記再構成レベルセット（４８）を決定することにより、前記状態が前記偶数の可能な状態のうちの前半に属する場合は前記複数（５０）の再構成レベルセット（５２）のうちの第１の再構成レベルセットが決定され、前記状態が前記偶数の可能な状態のうちの後半に属している場合は前記複数（５０）の再構成レベルセット（５２）のうちの第２の再構成レベルセットが決定される、請求項１６乃至１９のいずれかに記載の装置。

30

【請求項２１】

前記データストリームから復号化された、前記直前のニューラルネットワークパラメータの状態及び前記量子化インデックス（５８）のパリティの組み合わせを、前記後続のニューラルネットワークパラメータに関連づけられた別の状態にマッピングする遷移テーブルによって、前記状態の更新を実行するように構成される、請求項１６乃至２０のいずれかに記載の装置。

【請求項２２】

前記現在のニューラルネットワークパラメータ（１３'）について選択された前記再構成レベルセット（４８）に依存する確率モデル（１２３）を用いる算術符号化を用いて、前記現在のニューラルネットワークパラメータ（１３'）のための前記量子化インデックス（５６）を前記データストリーム（１４）から復号化するよう構成される、請求項１乃至２１のいずれかに記載の装置。

40

【請求項２３】

前記量子化インデックス（５６）を２値化したもの（８２）の少なくとも１つのピン（８４）についての前記現在のニューラルネットワークパラメータ（１３'）の状態に依存する（１２２）前記確率モデルを用いることにより、２値算術符号化を用いて、前記現在のニューラルネットワークパラメータのための前記量子化インデックス（５６）を前記データストリーム（１４）から復号化するよう構成される、請求項１乃至２２のいずれかに記載の装置。

【請求項２４】

50

前記少なくとも1つのピンが、前記現在のニューラルネットワークパラメータのための前記量子化インデックス(56)がゼロに等しいか否かを示す有意性ピンを含む、請求項23に記載の装置。

【請求項25】

前記少なくとも1つのピンが、前記現在のニューラルネットワークパラメータのための前記量子化インデックス(56)がゼロより大きいか又はゼロより小さいかを示す符号ピン(86)を含む、請求項23または請求項24に記載の装置。

【請求項26】

前記少なくとも1つのピンが、前記現在のニューラルネットワークパラメータの前記量子化インデックス(56)の絶対値がXより大きいか否かを示す、greater-than-Xピンを含み、ここでXはゼロより大きい整数である、請求項23乃至25のいずれかに記載の装置。

10

【請求項27】

前記確率モデルの依存性には、前記依存性を用いた、前記ニューラルネットワークパラメータについてのコンテキストセットのうちのコンテキスト(87)の選択(103)を含み、それぞれのコンテキストには既定の確率モデルが関連づけられているように構成される、請求項1乃至26のいずれかに記載の装置。

【請求項28】

前記それぞれのコンテキストを用いて算術符号化された前記量子化インデックスに基づいて、前記コンテキストのそれぞれに関連付けられた前記既定の確率モデルを更新するように構成される、請求項27に記載の装置。

20

【請求項29】

前記量子化インデックスを2値化したものの少なくとも1つのピンについて、前記現在のニューラルネットワークパラメータ(13')について選択された前記再構成レベルセット(48)に依存する確率モデルを用いることによって、前記現在のニューラルネットワークパラメータ(13')のための前記量子化インデックス(56)を前記データストリーム(14)から2値算術符号化を用いて復号化するように構成される、請求項1乃至28のいずれかに記載の装置。

【請求項30】

前記少なくとも1つのピンは、前記現在のニューラルネットワークパラメータのための前記量子化インデックス(56)がゼロに等しいか否かを示す有意性ピンを含む、請求項29に記載の装置。

30

【請求項31】

前記少なくとも1つのピンは、前記現在のニューラルネットワークパラメータのための前記量子化インデックス(56)がゼロより大きいか又はゼロより小さいかを示す符号ピンを含む、請求項29又は30に記載の装置。

【請求項32】

前記少なくとも1つのピンが、前記現在のニューラルネットワークパラメータのための前記量子化インデックス(56)の絶対値がXより大きいか否かを示すgreater-than-Xピンを含み、ここでXはゼロより大きい整数である、請求項29乃至31のいずれかに記載の装置。

40

【請求項33】

第1の状態又は再構成レベルセットについて事前選択されたサブセットが、任意の他の状態又は再構成レベルセットについて事前選択されたサブセットと互いに素であるように、前記複数の確率モデルの中から確率モデルの前記サブセットを、前記現在のニューラルネットワークパラメータ(13')について選択された再構成レベルの前記状態又は前記再構成レベルセット(48)に依存して事前選択するように構成される、請求項1乃至32のいずれかに記載の装置。

【請求項34】

前記現在のニューラルネットワークパラメータが関連する部分に隣接する前記ニューラ

50

ルネットワークの部分に関連する、以前に復号化されたニューラルネットワークパラメータの前記量子化インデックスに応じて、前記確率モデルのサブセットの中から前記現在のニューラルネットワークパラメータについての前記確率モデルを選択するように構成される、請求項 1 乃至 3 3 のいずれかに記載の装置。

【請求項 3 5】

前記現在のニューラルネットワークパラメータが関連する部分に隣接する前記ニューラルネットワークの部分に関連する、以前に復号化されたニューラルネットワークパラメータのための前記量子化インデックスの特性に依存して、前記確率モデルのサブセットの中から前記現在のニューラルネットワークパラメータについての前記確率モデルを選択するように構成され、前記特性は、

10

前記現在のニューラルネットワークパラメータが関連する部分に隣接する前記ニューラルネットワークの部分に関連する、以前に復号化されたニューラルネットワークパラメータのゼロでない量子化インデックスの符号と、

前記現在のニューラルネットワークパラメータが関連する部分に隣接する前記ニューラルネットワークの部分に関連する、以前に復号化されたニューラルネットワークパラメータの量子化インデックスの数であって、ゼロでない数と、

前記現在のニューラルネットワークパラメータが関連する部分に隣接する前記ニューラルネットワークの部分に関連する、以前に復号化されたニューラルネットワークパラメータの量子化インデックスの絶対値の合計値と、

前記現在のニューラルネットワークパラメータが関連する部分に隣接する前記ニューラルネットワークの部分に関連する、以前に復号化されたニューラルネットワークパラメータの量子化インデックスの前記絶対値の合計値と、前記現在のニューラルネットワークパラメータが関連する部分に隣接する前記ニューラルネットワークの部分に関連する、以前に復号化されたニューラルネットワークパラメータの量子化インデックスの数であって、ゼロでない数と、の差分と、
のうちの 1 つ以上を含む、請求項 1 乃至 3 3 のいずれかに記載の装置。

20

【請求項 3 6】

前記以前に復号化されたニューラルネットワークパラメータが前記現在のニューラルネットワークパラメータと同じニューラルネットワーク層に関連するように、前記以前に復号化されたニューラルネットワークパラメータを位置づけるように構成される、請求項 3 4 又は 3 5 に記載の装置。

30

【請求項 3 7】

1 つ以上の前記以前に復号化されたニューラルネットワークパラメータが、前記現在のニューラルネットワークパラメータが参照するニューロン相互接続が関連するニューロンまたは該ニューロンに隣接する別のニューロンから出現するニューロン相互接続又はこれらのニューロンに向かうニューロン相互接続に関連するように、前記以前に復号化されたニューラルネットワークパラメータのうちの 1 つ以上のパラメータを位置付けるように構成される、請求項 3 4 又は 3 5 に記載の装置。

【請求項 3 8】

前記ニューラルネットワークパラメータ (1 3) のための前記量子化インデックス (5 6) を復号化し、前記ニューラルネットワークパラメータ (1 3) 間の共通の連続的な順序 (1 4 ') に沿って前記ニューラルネットワークパラメータ (1 3) の前記逆量子化を実行するように構成される、請求項 1 乃至 3 7 のいずれかに記載の装置。

40

【請求項 3 9】

前記量子化インデックスを 2 値化したものの 1 つ以上のリーディングビン (l e a d i n g b i n s) についての以前に復号化されたニューラルネットワークパラメータに依存する確率モデルを用いること、及び、前記 1 つ以上のリーディングビンに後続する、前記量子化インデックスを 2 値化したものの等確率バイパスモードサフィックスビンを用いることによって、前記現在のニューラルネットワークパラメータ (1 3 ') のための前記量子化インデックス (5 6) を前記データストリーム (1 4) から 2 値算術符号化を用いて

50

復号化するように構成される、請求項 1 乃至 3 8 のいずれかに記載の装置。

【請求項 4 0】

前記量子化インデックスを 2 値化したもののサフィックスピンは、絶対値が前記 1 つ以上の前記リーディングピンによって表現可能な最大絶対値を超える前記量子化インデックスの値を 2 値化するためのサフィックス 2 値化の 2 値化コードのピンを表し、前記装置は、以前に復号化されたニューラルネットワークパラメータの前記量子化インデックスに応じて前記サフィックス 2 値化を選択するように構成される、請求項 3 9 に記載の装置。

【請求項 4 1】

前記ニューラルネットワークパラメータは、前記ニューラルネットワーク (1 0) が表現される再構成層のうちの 1 つの再構成層に関連し、前記装置は、

前記ニューラルネットワークパラメータを、ニューラルネットワークパラメータ単位で、1 つ以上の別の再構成層の対応するニューラルネットワークパラメータと組み合わせることによって、前記ニューラルネットワークを再構成するように構成される、請求項 1 乃至 4 0 のいずれかに記載の装置。

【請求項 4 2】

前記現在のニューラルネットワークパラメータに対応する、対応ニューラルネットワークパラメータに依存する確率モデルを用いる算術符号化を用いて、前記現在のニューラルネットワークパラメータ (1 3 ') のための前記量子化インデックス (5 6) を前記データストリーム (1 4) から復号化するように構成される、請求項 4 1 に記載の装置。

【請求項 4 3】

ニューラルネットワークを定義するニューラルネットワークパラメータをデータストリームに符号化するための装置であって、

現在のニューラルネットワークパラメータ (1 3 ') について、前記データストリーム (1 4) に符号化された、以前に符号化されたニューラルネットワークパラメータのための量子化インデックス (5 8) に依存して、複数 (5 0) の再構成レベルセット (5 2) の中から再構成レベルセット (4 8) を選択 (5 4) することと、

選択された前記再構成レベルセット (4 8) の 1 つの再構成レベル上に前記現在のニューラルネットワークパラメータ (1 3 ') を量子化 (6 4) することと、

前記現在のニューラルネットワークパラメータのための量子化インデックス (5 6) が量子化される前記 1 つの再構成レベルを示す、前記現在のニューラルネットワークパラメータ (1 3 ') のための前記量子化インデックス (5 6) を前記データストリーム (1 4) に符号化することと

により、前記ニューラルネットワークパラメータ (1 3 ') を順次符号化するように構成され、

前記装置は、

前記現在のニューラルネットワークパラメータ (1 3 ') に関連づけられた状態に応じて、前記現在のニューラルネットワークパラメータ (1 3 ') について、前記複数 (5 0) の再構成レベルセット (5 2) の中から前記再構成レベルセット (4 8) を決定すること、及び、

前記データストリームに符号化された、直前のニューラルネットワークパラメータのための前記量子化インデックス (5 8) に応じて、後続のニューラルネットワークパラメータの状態を更新すること

による状態遷移プロセスにより、前記現在のニューラルネットワークパラメータ (1 3 ') について、前記複数 (5 0) の再構成レベルセット (5 2) の中から前記再構成レベルセット (4 8) を選択 (5 4) し、

前記現在のニューラルネットワークパラメータ (1 3 ') の状態および以前に符号化されたニューラルネットワークパラメータの前記量子化インデックスに依存する (1 2 2) 確率モデルを用いる算術符号化を用いて、前記データストリーム (1 4) から前記現在のニューラルネットワークパラメータ (1 3 ') のための前記量子化インデックス (5 6) を符号化し、

10

20

30

40

50

前記現在のニューラルネットワークパラメータ（１３'）に関連付けられた前記状態に依存して、複数の確率モデルの中から確率モデルのサブセットを事前選択し、以前に符号化されたニューラルネットワークパラメータのための前記量子化インデックスに依存して（１２１）、前記確率モデルのサブセットの中から前記現在のニューラルネットワークパラメータのための前記確率モデルを選択し、

前記現在のニューラルネットワークパラメータが関連する部分に隣接する前記ニューラルネットワークの部分に関連する、以前に符号化されたニューラルネットワークパラメータの前記量子化インデックスの値に依存して、前記現在のニューラルネットワークパラメータのための確率モデルを第１、第２および第３の確率モデルを含むサブセットの中から選択する

10

ように構成され、

ここで、前記以前に符号化されたニューラルネットワークパラメータがゼロよりも小さければ、前記第１の確率モデルが選択され、

前記以前に符号化されたニューラルネットワークパラメータがゼロよりも大きければ、前記第２の確率モデルが選択され、

前記以前に符号化されたニューラルネットワークパラメータがゼロに等しければ、前記第３の確率モデルが選択される、

装置。

【請求項４４】

ニューラルネットワーク（１０）を定義するニューラルネットワークパラメータ（１３）をデータストリーム（１４）から復号化するための方法（４００）であって、

20

現在のニューラルネットワークパラメータ（１３'）について、前記データストリーム（１４）から復号化された、以前のニューラルネットワークパラメータのための量子化インデックス（５８）に依存して、複数（５０）の再構成レベルセット（５２）の中から再構成レベルセット（４８）を選択する（５４）ことと、

前記現在のニューラルネットワークパラメータについての前記選択された再構成レベルセット（４８）のうちの１つの再構成レベルを示す、前記現在のニューラルネットワークパラメータ（１３'）のための量子化インデックス（５６）を前記データストリーム（１４）から復号化する（４２０）ことと、

前記現在のニューラルネットワークパラメータのための前記量子化インデックス（５６）によって示される前記選択された再構成レベルセット（４８）のうちの前記１つの再構成レベル上に前記現在のニューラルネットワークパラメータ（１３'）を逆量子化する（６２）ことと、

30

によって、前記ニューラルネットワークパラメータ（１３）を順次復号化するステップを含み、

前記方法はさらに、

前記現在のニューラルネットワークパラメータ（１３'）に関連づけられた状態に応じて、前記現在のニューラルネットワークパラメータ（１３'）について、前記複数（５０）の再構成レベルセット（５２）の中から前記再構成レベルセット（４８）を決定すること、及び、

40

前記データストリームから復号化された、直前のニューラルネットワークパラメータのための前記量子化インデックス（５８）に応じて、後続のニューラルネットワークパラメータの状態を更新すること

による状態遷移プロセスにより、前記現在のニューラルネットワークパラメータ（１３'）について、前記複数（５０）の再構成レベルセット（５２）の中から前記再構成レベルセット（４８）を選択（５４）するステップと、

前記現在のニューラルネットワークパラメータ（１３'）の状態および以前に復号化されたニューラルネットワークパラメータの前記量子化インデックスに依存する（１２２）確率モデルを用いる算術符号化を用いて、前記データストリーム（１４）から前記現在のニューラルネットワークパラメータ（１３'）のための前記量子化インデックス（５６）を復号

50

化するステップと、

前記現在のニューラルネットワークパラメータ（１３'）に関連付けられた前記状態に依存して、複数の確率モデルの中から確率モデルのサブセットを事前選択し、以前に復号化されたニューラルネットワークパラメータのための前記量子化インデックスに依存して（１２１）、前記確率モデルのサブセットの中から前記現在のニューラルネットワークパラメータのための前記確率モデルを選択するステップと、

前記現在のニューラルネットワークパラメータが関連する部分に隣接する前記ニューラルネットワークの部分に関連する、以前に復号化されたニューラルネットワークパラメータの前記量子化インデックスの値に依存して、前記現在のニューラルネットワークパラメータのための確率モデルを第１、第２および第３の確率モデルを含むサブセットの中から選択するステップと

10

を含み、

ここで、前記以前に復号化されたニューラルネットワークパラメータがゼロよりも小さければ、前記第１の確率モデルが選択され、

前記以前に復号化されたニューラルネットワークパラメータがゼロよりも大きければ、前記第２の確率モデルが選択され、

前記以前に復号化されたニューラルネットワークパラメータがゼロに等しければ、前記第３の確率モデルが選択される、

方法。

【請求項４５】

20

ニューラルネットワークを定義するニューラルネットワークパラメータをデータストリームに符号化するための方法（５００）であって、

現在のニューラルネットワークパラメータ（１３'）について、前記データストリーム（１４）に符号化された、以前に符号化されたニューラルネットワークパラメータのための量子化インデックス（５８）に依存して、複数（５０）の再構成レベルセット（５２）の中から再構成レベルセット（４８）を選択（５４）することと、

前記現在のニューラルネットワークパラメータ（１３'）を、前記選択された再構成レベルセット（４８）の１つの再構成レベル上に量子化する（６４）ことと、

前記現在のニューラルネットワークパラメータのための量子化インデックス（５６）が量子化される前記１つの再構成レベルを示す、前記現在のニューラルネットワークパラメータのための前記量子化インデックス（５６）を前記データストリーム（１４）に符号化する（５３０）ことと、

30

によって、前記ニューラルネットワークパラメータ（１３）を順次符号化するステップを含み、

前記方法はさらに、

前記現在のニューラルネットワークパラメータ（１３'）に関連づけられた状態に応じて、前記現在のニューラルネットワークパラメータ（１３'）について、前記複数（５０）の再構成レベルセット（５２）の中から前記再構成レベルセット（４８）を決定すること、及び、

前記データストリームに符号化された、直前のニューラルネットワークパラメータのための前記量子化インデックス（５８）に応じて、後続のニューラルネットワークパラメータの状態を更新すること

40

による状態遷移プロセスにより、前記現在のニューラルネットワークパラメータ（１３'）について、前記複数（５０）の再構成レベルセット（５２）の中から前記再構成レベルセット（４８）を選択（５４）するステップと、

前記現在のニューラルネットワークパラメータ（１３'）の状態および以前に符号化されたニューラルネットワークパラメータの前記量子化インデックスに依存する（１２２）確率モデルを用いる算術符号化を用いて、前記現在のニューラルネットワークパラメータ（１３'）のための前記量子化インデックス（５６）を前記データストリーム（１４）に符号化するステップと、

50

前記現在のニューラルネットワークパラメータ（１３'）に関連付けられた前記状態に依存して、複数の確率モデルの中から確率モデルのサブセットを事前選択し、以前に符号化されたニューラルネットワークパラメータのための前記量子化インデックスに依存して（１２１）、前記確率モデルのサブセットの中から前記現在のニューラルネットワークパラメータのための前記確率モデルを選択するステップと、

前記現在のニューラルネットワークパラメータが関連する部分に隣接する前記ニューラルネットワークの部分に関連する、以前に符号化されたニューラルネットワークパラメータの前記量子化インデックスの値に依存して、前記現在のニューラルネットワークパラメータのための確率モデルを第１、第２および第３の確率モデルを含むサブセットの中から選択するステップと

10

を含み、

ここで、前記以前に符号化されたニューラルネットワークパラメータがゼロよりも小さければ、前記第１の確率モデルが選択され、

前記以前に符号化されたニューラルネットワークパラメータがゼロよりも大きければ、前記第２の確率モデルが選択され、

前記以前に符号化されたニューラルネットワークパラメータがゼロに等しければ、前記第３の確率モデルが選択される、

方法。

【請求項４６】

プログラムが１つ以上のコンピュータ上で実行されるときに、請求項４４又は４５に記載の方法を実行するためのコンピュータプログラムを記憶した、非一時的デジタル記憶媒体。

20

【発明の詳細な説明】

【技術分野】

【０００１】

本発明による実施形態は、ニューラルネットワークパラメータの符号化概念に関連する。

【背景技術】

【０００２】

ニューラルネットワークの最も基本的な形態は、アフィン変換の連鎖と、それに続く要素ごとの非線形関数である。図１に示すように、有向非巡回グラフ（directed acyclic graph）として表現することができる。図１は、ニューラルネットワークの一例を示す模式図であり、ここでは例示的に２層フィードフォワードニューラルネットワークである。すなわち、図１は、フィードフォワードニューラルネットワークをグラフで表現したものである。具体的には、この２層構造のニューラルネットワークは、４次元の入力ベクトルを実線にマッピングする非線形関数である。このニューラルネットワークは、ニューラルネットワークの入力となるInput層における４次元の入力ベクトルに応じた４つのニューロン10cと、Hidden層における５つのニューロン10cと、ニューラルネットワークの出力を形成するOutput層における１つのニューロン10cを含む。ニューラルネットワークは、異なる - 又は後続の - 層のニューロンを接続する、ニューロン相互接続11をさらに含む。ニューロン相互接続11は、重みを関連付けてもよく、重みは、互いに接続されたニューロン10cの間の関係に関連付けられる。特に、重みは、後続の層に転送されるとき、ある層のニューロンの活性化を重み付けし、順次、その後続の層の各ニューロンでインバウンドの重み付き活性化の合計が形成され - 線形関数に対応する - 続いて、後続の層の各ニューロン/ノードで形成される重み付け合計に非線形スカラー関数が適用され - 非線形関数に対応する。したがって、各ノード、例えばニューロン10cは、特定の値を伴い、この値は、エッジ、例えばニューロン相互接続11のそれぞれの重み値との乗算によって次のノードに順方向に伝播される。その後、すべての受信値は単純に集計される。

30

40

【０００３】

50

数学的には、図 1 のニューラルネットワークは、次のように出力を計算することになる。

$$\text{output} = \sigma(W_2 \cdot \sigma(W_1 \cdot \text{input}))$$

ここで、 W_2 、 W_1 はニューラルネットワークのパラメータ、例えば、ニューラルネットワークの重みパラメータ（エッジ重み）であり、シグマは何らかの非線形関数である。例えば、[1]にあるように、いわゆる畳み込み層を行列－行列積としてキャストすることで使用することも可能である。以下、与えられた入力から出力を計算する手順を推論と呼ぶことにする。また、例えば上記の第 1 ドット積＋非直線性の計算のような、線形変換＋要素単位の非直線性を構成する中間結果を隠れ層あるいは隠れ活性化値と呼ぶことにする。

10

【発明の概要】

【発明が解決しようとする課題】

【0004】

通常、ニューラルネットワークは数百万個のパラメータを持ち、その表現には数百 MB（例：メガバイト）が必要とされることがある。結果として、それらの推論手順には大きな行列間の多数のドット積演算の計算が含まれるため、実行には高い計算資源が必要となる。そのため、ドット積演算の複雑さを軽減することは非常に重要である。

【0005】

20

また、上記の問題点に加えて、ニューラルネットワークの膨大なパラメータを保存する必要があり、例えばサーバからクライアントへ送信する必要がある場合もある。さらに、例えば連合学習環境（`federated learning environment`）において、または特定の受信者が支払った、または推論にニューラルネットワークを使用するときに対処することができる品質の異なる段階でニューラルネットワークのパラメータ化を提供する場合に、ニューラルネットワークのパラメータ化に関する情報を徐々にエンティティに提供することができることが有利である場合もある。

【0006】

したがって、例えば圧縮の観点からより効率的な、ニューラルネットワークパラメータの効率的な符号化のための概念を提供することが望まれている。さらに、または代替的に、ニューラルネットワークパラメータのためのビットストリーム、ひいては信号化コストを低減することが望まれる。

30

【0007】

この目的は、本願の独立請求項の主題によって達成される。

【0008】

本発明によるさらなる実施形態は、本願の従属請求項の主題によって定義される。

【課題を解決するための手段】

【0009】

本発明の第 1 の側面による実施形態は、ニューラルネットワークを定義するニューラルネットワークパラメータをデータストリームから復号化する装置であって、現在のニューラルネットワークパラメータについて、以前のニューラルネットワークパラメータについてデータストリームから復号化された量子化インデックスに応じて複数の再構成レベルセットから再構成レベルセットを選択することによって、ニューラルネットワークパラメータを順次復号化するように構成されている、装置である。さらに、装置は、データストリームから現在のニューラルネットワークパラメータについての量子化インデックスを復号化することによって、ここで量子化インデックスは、現在のニューラルネットワークパラメータについての選択された再構成レベルセットのうちの 1 つの再構成レベルを示し、及び、現在のニューラルネットワークパラメータについての量子化インデックスによって示される、選択された再構成レベルセットのうちの 1 つの再構成レベル上に現在のニューラルネットワークパラメータを逆量子化することによって、ニューラルネットワークパラメ

40

50

ータを順次復号化するように構成される。

【 0 0 1 0 】

本発明の第 1 の側面によるさらなる実施形態は、ニューラルネットワークを定義するニューラルネットワークパラメータをデータストリームに符号化する装置であって、現在のニューラルネットワークパラメータについて、以前に符号化されたニューラルネットワークパラメータについてデータストリームに符号化された量子化インデックスに応じて複数の再構成レベルセットの中から再構成レベルセットを選択することによって、ニューラルネットワークパラメータを順次符号化するように構成されている、装置である。さらに、装置は、選択された再構成レベルセットのうちの 1 つの再構成レベルに現在のニューラルネットワークパラメータを量子化することによって、及び、現在のニューラルネットワークパラメータについての量子化インデックスが量子化される 1 つの再構成レベルを示す現在のニューラルネットワークパラメータについての量子化インデックスをデータストリームに符号化することにより、ニューラルネットワークパラメータを順次符号化するよう構成される。

10

【 0 0 1 1 】

本発明の第 1 の側面による更なる実施形態は、ニューラルネットワークを定義するニューラルネットワークパラメータをデータストリームから復号化するための方法を備える。この方法は、現在のニューラルネットワークパラメータについて、以前のニューラルネットワークパラメータについてのデータストリームから復号化された量子化インデックスに応じて複数の再構成レベルセットの中から再構成レベルセットを選択することによって、ニューラルネットワークパラメータを順次復号化するステップを備える。さらに、この方法は、データストリームから現在のニューラルネットワークパラメータについての量子化インデックスを復号化することによって、ここで、量子化インデックスが、現在のニューラルネットワークパラメータについての選択された再構成レベルセットのうちの 1 つの再構成レベルを示し、及び、現在のニューラルネットワークパネルについての量子化インデックスが示す選択された再構成レベルセットのうちの 1 つの再構成レベル上に現在のニューラルネットワークパラメータを逆量子化することによって、ニューラルネットワークパラメータを順次符号化するステップを含む。

20

【 0 0 1 2 】

本発明の第 1 の側面による更なる実施形態は、ニューラルネットワークを定義するニューラルネットワークパラメータをデータストリームに符号化するための方法を備える。この方法は、現在のニューラルネットワークパラメータについて、以前に符号化されたニューラルネットワークパラメータについてデータストリームに符号化された量子化インデックスに応じて複数の再構成レベルセットの中から再構成レベルセットを選択することにより、ニューラルネットワークパラメータを順次符号化することを備える。さらに、この方法は、現在のニューラルネットワークパラメータを選択された再構成レベルのうちの 1 つの再構成レベルに量子化することによって、及び、現在のニューラルネットワークパラメータについての量子化インデックスが量子化される 1 つの再構成レベルを示す現在のニューラルネットワークパラメータについての量子化インデックスをデータストリームに符号化することによって、ニューラルネットワークパラメータを順次符号化するステップを含む。

30

40

【 0 0 1 3 】

本発明の第 1 の態様による実施形態は、非定数量子化器を使用するが、ニューラルネットワークパラメータの符号化したものの中に同じものを変化させることによって、すなわち、以前のまたはそれぞれ以前に符号化されたニューラルネットワークパラメータのデータストリームから復号された、またはそれぞれ符号化された量子化インデックスに応じて再構成レベルのセットを選択することによって、ニューラルネットワークパラメータをより効率的に圧縮し得るという考え方に基づくものである。したがって、ニューラルネットワークパラメータの順序付けられたセットを参照することができる再構成ベクトルは、N 次元信号空間においてより密に詰め込むことができ、ここで、N は、処理されるサンプル

50

のセットにおけるニューラルネットワークパラメータの数を示す。このような依存量子化 (dependent quantization) は、復号化するための装置による復号化及び逆量子化、又は符号化するための装置による量子化及び符号化にそれぞれ使用されることができる。

【0014】

本発明の第2の態様による実施形態は、段階、-ニューラル層におけるニューラルネットワークの層構成と区別するために再構成層と呼ばれる-において行われ時に、これらの段階で提供されるパラメータ化を次に、ニューラルネットワークパラメータ単位で組み合わせて、段階のいずれかと比較して改善されたニューラルネットワークのパラメータ化をもたらす場合、より効率のよいニューラルネットワーク符号化が達成できるかもしれないという考え方に基づく。したがって、ニューラルネットワークを定義するニューラルネットワークパラメータを再構成するための装置は、第1のニューラルネットワークパラメータ、例えば、第1の再構成層について第1の再構成層のニューラルネットワークパラメータを導出して、ニューラルネットワークパラメータごとに、第1の再構成層のニューラルネットワークパラメータ値を得ることができる。第1のニューラルネットワークパラメータは、例えば、連合学習プロセス中に以前に送信されている可能性がある。さらに、第1のニューラルネットワークパラメータは、第1の再構成層のニューラルネットワークパラメータ値であってもよい。さらに、装置は、第2のニューラルネットワークパラメータ、例えば第2の再構成層のニューラルネットワークパラメータを、例えば最終ニューラルネットワークパラメータと区別するために、第2の再構成層についてデータストリームから復号化し、ニューラルネットワークパラメータごとに、第2の再構成層ニューラルネットワークパラメータ値をもたらすように構成されている。第2のニューラルネットワークパラメータは、ニューラルネットワーク表現の観点から自己完結した意味を持たず、単に、第1の表現層のパラメータと組み合わせたときに、ニューラルネットワーク表現、すなわち、例えば、最終ニューラルネットワークパラメータをもたらすだけかもしれない。さらに、本装置は、各ニューラルネットワークパラメータについて、第1の再構成層のニューラルネットワークパラメータ値と第2の再構成層のニューラルネットワークパラメータ値とを組み合わせることにより、ニューラルネットワークパラメータを再構成するように構成される。

【0015】

本発明の第2の側面によるさらなる実施形態は、ニューラルネットワークを定義するニューラルネットワークパラメータを、ニューラルネットワークパラメータごとに、第1の再構成層のニューラルネットワークパラメータ値を含む第1の再構成層についての第1のニューラルネットワークパラメータを用いて符号化するための装置を備える。さらに、この装置は、第2の再構成層についての第2のニューラルネットワークパラメータをデータストリームに符号化するように構成され、第2の再構成層は、ニューラルネットワークパラメータごとに第2の再構成層のニューラルネットワークパラメータ値を含み、ニューラルネットワークパラメータは、ニューラルネットワークパラメータごとに、第1の再構成層のニューラルネットワークパラメータ値と第2の再構成層のニューラルネットワークパラメータ値を組み合わせることによって再構成可能である。

【0016】

本発明の第2の態様によるさらなる実施形態は、ニューラルネットワークを定義するニューラルネットワークパラメータを再構成するための方法を備える。この方法は、ニューラルネットワークパラメータごとに、第1の再構成層のニューラルネットワークパラメータ値をもたらすために、第1の再構成層について、例えば連合学習プロセス中に以前に送信された可能性があり、例えば第1の再構成層のニューラルネットワークパラメータと呼ばれる可能性がある第1のニューラルネットワークパラメータを導出するステップを含む。

【0017】

さらに、本方法は、この方法は、例えば最終的な、例えば再構成されたニューラルネットワークパラメータと区別するために第2の再構成層のニューラルネットワークパラメータ

10

20

30

40

50

タと呼ばれ得る第2のニューラルネットワークパラメータを、データストリームから第2再構成層について復号化して、ニューラルネットワークパラメータごとに、第2の再構成層のニューラルネットワークパラメータ値を生成するステップを含む。本方法は、各ニューラルネットワークパラメータについて、第1の再構成層のニューラルネットワークパラメータ値と第2の再構成層のニューラルネットワークパラメータ値とを組み合わせることによって、ニューラルネットワークパラメータを再構成するステップを含む。第2のニューラルネットワークパラメータは、ニューラル表現に関して自己完結した意味を持たず、単に、第1の表現層のパラメータと組み合わせたときに、ニューラル表現、すなわち、例えば最終ニューラルネットワークパラメータを導くかもしれない。

【0018】

10

本発明の第2の側面によるさらなる実施形態は、ニューラルネットワークパラメータごとに、第1の再構成層のニューラルネットワークパラメータ値を含む第1の再構成層についての第1のニューラルネットワークパラメータを用いて、ニューラルネットワークを定義するニューラルネットワークパラメータを符号化する方法を含む。この方法は、第2の再構成層についての第2のニューラルネットワークパラメータをデータストリームに符号化するステップを含み、第2の再構成層は、ニューラルネットワークパラメータごとに、第2の再構成層ニューラルネットワークパラメータ値を含み、ニューラルネットワークパラメータは、ニューラルネットワークパラメータごとに、第1の再構成層のニューラルネットワークパラメータ値と第2の再構成層のニューラルネットワークパラメータ値とを組み合わせることによって再構成可能である。

20

【0019】

本発明の第2の態様による実施形態は、例えばニューラルネットワークパラメータによって定義されるニューラルネットワークが、再構成層、例えばベース層およびエンハンスメント層などのサブ層を用いて、例えばビットストリーム内のデータ量が少ない状態で、効率的に圧縮および/または送信され得るという考えに基づいている。再構成層は、ニューラルネットワークパラメータが、各ニューラルネットワークパラメータについて、第1の再構成層のニューラルネットワークパラメータ値と第2の再構成層のニューラルネットワークパラメータ値とを組み合わせることによって再構成可能であるように、定義されることができる。この分布は、ニューラルネットワークパラメータの効率的な符号化、例えば、符号化および/または復号化、および/または送信を可能にする。したがって、第2の再構成層についての第2のニューラルネットワークパラメータは、データストリームに別々に符号化および/または送信されることができる。

30

【図面の簡単な説明】

【0020】

図面は必ずしも縮尺通りではなく、代わりに一般的に本発明の原理を説明することに重点が置かれている。以下の説明では、本発明の様々な実施形態が、以下の図面を参照して説明され、その中で、本発明の実施形態が説明される。

【図1】図1は、本発明の実施形態と共に使用され得る2層フィードフォワードニューラルネットワークの例示の概略図である。

【図2】図2は、実施形態による、ニューラルネットワークを定義するニューラルネットワークパラメータをデータストリームから復号化するための装置内で実行される逆量子化のための概念を示す概略図である。

40

【図3】図3は、実施形態による、ニューラルネットワークパラメータをデータストリームに符号化するための装置内で実行される量子化のための概念を示す概略図である。

【図4】図4は、実施形態による、ニューラルネットワークを定義するニューラルネットワークパラメータを再構成するための装置内で実行される復号化のための概念を示す概略図である。

【図5】図5は、実施形態による、ニューラルネットワークを定義するニューラルネットワークパラメータを再構成するための装置内で実行される符号化のための概念を示す概略図である。

50

【図 6】図 6 は、本発明による実施形態で使用するための、ニューラルネットワークパラメータについての再構成層を使用する概念の概略図である。

【図 7】図 7 は、本発明の実施形態による均一再構成量子化器のイラストを示す概略図である。

【図 8】図 8 は、本発明の実施形態による 2 つの重みパラメータの単純なケースについて、許容される再構成ベクトルの位置の例を示す図である。

【図 9 a】図 9 a は、本発明の実施形態による単一の量子化ステップサイズ によって完全に決定される再構成レベルの 2 つのセットを有する依存量子化のための例を示す図である。

【図 9 b】図 9 b は、本発明の実施形態による単一の量子化ステップサイズ によって完全に決定される再構成レベルの 2 つのセットを有する依存量子化のための例を示す図である。

10

【図 9 c】図 9 c は、本発明の実施形態による単一の量子化ステップサイズ によって完全に決定される再構成レベルの 2 つのセットを有する依存量子化のための例を示す図である。

【図 1 0】図 1 0 は、本発明の実施形態による、ニューラルネットワークパラメータの再構成プロセスについての好ましい例を示す疑似コードのための例を示す図である。

【図 1 1】図 1 1 は、本発明の実施形態による再構成レベルセットを 2 つのサブセットに分割するための一例を示す図である。

【図 1 2】図 1 2 は、実施形態による層についてのニューラルネットワークパラメータの再構成プロセスについての好ましい例を示す疑似コードの例である。

20

【図 1 3】図 1 3 は、本発明の実施形態による状態遷移表 $s t t a b$ 及び表 $s e t I d$ の好ましい例を示す図であり、これは、状態に関連する量子化セットを指定する。

【図 1 4】図 1 4 は、本発明の実施形態による状態遷移表 $s t t a b$ 及び表 $s e t I d$ の好ましい例を示す図であり、これは、状態に関連する量子化セットを指定する。

【図 1 5】図 1 5 は、本発明の実施形態による、0 に等しい量子化インデックスが状態遷移および依存スカラー量子化から除外される、ニューラルネットワークパラメータレベルの代替再構成プロセスを示す疑似コードである。

【図 1 6】図 1 6 は、本発明の実施形態による、トレリス構造としての依存スカラー量子化における状態遷移の例を示す図である。

30

【図 1 7】図 1 7 は、本発明の実施形態による基本トレリスセルの一例を示す図である。

【図 1 8】図 1 8 は、本発明の実施形態による 8 つのニューラルネットワークパラメータの依存スカラー量子化のためのトレリス例を示す図である。

【図 1 9】図 1 9 は、本発明の実施形態による、コスト量 (ラグランジュコスト量 $D + \cdot R$ など、 $L a g r a n g i a n \ c o s t \ m e a s u r e$) を最小化する量子化インデックスのシーケンス (またはブロック) を決定するために利用することができるトレリス構造例を示す図である。

【図 2 0】図 2 0 は、本発明の実施形態による、ニューラルネットワークを定義するニューラルネットワークパラメータをデータストリームから復号化する方法のブロック図である。

40

【図 2 1】図 2 1 は、本発明の実施形態による、ニューラルネットワークを定義するニューラルネットワークパラメータをデータストリームに符号化するための方法のブロック図である。

【図 2 2】図 2 2 は、本発明の実施形態による、ニューラルネットワークを定義するニューラルネットワークパラメータを再構成する方法のブロック図である。

【図 2 3】図 2 3 は、本発明の実施形態による、ニューラルネットワークを定義する、ニューラルネットワークパラメータを符号化するための方法のブロック図である。

【発明を実施するための形態】

【0 0 2 1】

等しいまたは同等の機能を有する要素または要素は、異なる図で表されていても、以下

50

の説明では、等しいまたは同等の参照数字で示される。

【 0 0 2 2 】

以下の説明では、本発明の実施形態のより詳細な説明を提供するために、複数の詳細が記載されている。しかし、本発明の実施形態は、これらの特定の詳細なしに実施され得ることは、当業者には明らかであろう。他の例では、本発明の実施形態を不明瞭にしないために、周知の構造およびデバイスを詳細ではなくブロック図の形態で示す。加えて、本明細書で後述する異なる実施形態の特徴は、特に断らない限り、互いに組み合わせることができる。

【 0 0 2 3 】

本説明は、本願のいくつかの実施形態の提示から始まる。この説明は、かなり一般的なものであるが、本願の実施形態が基づいている機能性の概要を読者に提供するものである。その後、これらの機能性のより詳細な説明が、実施形態の動機と、それらがどのように上述の効率向上を達成するかと共に提示される。詳細は、現在説明されている実施形態と、個別に、および組み合わせて、組み合わせることが可能である。

【 0 0 2 4 】

図 2 は、実施形態によるデータストリームからニューラルネットワークを定義するニューラルネットワークパラメータを復号化するための装置内で実行される逆量子化についての概念を示す概略図である。ニューラルネットワークは、例えば、相互接続された層のニューロン間のニューロン相互接続を有する、複数の相互接続されたニューラルネットワーク層を含むことができる。図 2 は、データストリーム 1 4 における、例えば符号化されたニューラルネットワークパラメータ 1 3 の量子化インデックス 5 6 を示す。ニューラルネットワークパラメータ 1 3 は、したがって、そのニューロン間の重みの観点など、ニューラルネットワークを定義またはパラメータ化することができる。

【 0 0 2 5 】

本装置は、ニューラルネットワークパラメータ 1 3 を順次復号化するように構成されている。この順次プロセスの間、量子化器（再構成レベルセット）が変化させられる。この変化は、より少ない（またはより良いより密でない）レベルを有する量子化器を使用することを可能にし、したがって、より小さい量子化インデックスを符号化することを可能にし、この量子化から得られるニューラルネットワーク表現の品質が、必要な符号化ビットレートと比較して、一定の量子化器を使用するより向上される。詳細は後述する。特に、装置は、現在のニューラルネットワークパラメータ 1 3 ' について、以前のニューラルネットワークパラメータについてデータストリーム 1 4 から復号化された量子化インデックス 5 8 に応じて複数（5 0）の再構成レベルセット 5 2（セット 0、セット 1）のうち再構成レベルセット 4 8（選択されたセット）を選択 5 4（再構成レベル選択）することによって、ニューラルネットワークパラメータ 1 3 を順次復号する。

【 0 0 2 6 】

さらに、本装置は、データストリーム 1 4 から現在のニューラルネットワークパラメータ 1 3 ' についての量子化インデックス 5 6 を復号化することによって、ここで、量子化インデックス 5 6 が現在のニューラルネットワークパラメータについての再構成レベルの選択されたセット 4 8 のうちの 1 つの再構成レベルを示し、及び、現在のニューラルネットワークパラメータ 1 3 ' を現在のニューラルネットワークパネルについての量子化インデックス 5 6 によって示される再構成レベルの選択されたセット 4 8 の 1 つの再構成レベルに逆量子化 6 2 することによって、ニューラルネットワークパネル 1 3 を連続的に復号化するように構成される。

【 0 0 2 7 】

復号化されたニューラルネットワークパラメータ 1 3 は、一例として、行列 1 5 a で表される。行列は、デシリアライズされた 2 0 b (d e s e r i a l i z a t i o n) ニューラルネットワークパラメータ 1 3 を含んでもよく、これは、ニューラルネットワークのニューロン相互接続の重みに関連してもよい。

【 0 0 2 8 】

10

20

30

40

50

任意選択的に、複数（５０）の再構成レベルセット５２の、本明細書で時々量子化器とも呼ばれる再構成レベルセット５２の数は、図２に示すように、例えばセット０とセット１の２つであってよい。

【００２９】

さらに、装置は、例えば または k で示される所定の量子化ステップサイズ（ $Q P$ ）によって複数（５０）の再構成レベルセット５２（例えば、セット０、セット１）をパラメータ化６０（*parameterization*）し、データストリーム１４から所定の量子化ステップサイズの情報を導出するように構成されることができる。したがって、実施形態による復号化器は、可変ステップサイズ（ $Q P$ ）に適応することができる。

【００３０】

さらに、実施形態によれば、ニューラルネットワークは１つ以上の NN 層を含むことができ、装置は、各 NN 層について、データストリーム１４からそれぞれの NN 層についての所定の量子化ステップサイズ（ $Q P$ ）の情報を導出し、それぞれの NN 層に属するニューラルネットワークパラメータを逆量子化するために用いられるように、各 NN 層について、それぞれの NN 層について導出された所定の量子化ステップサイズを用いて複数５０の再構成レベルセット５２をパラメータ化するように構成され得る。 NN 層に関するステップサイズの適応、したがって再構成レベルセット５２の適応は、符号化効率を向上させ得る。

【００３１】

さらなる実施形態によれば、装置は、現在のニューラルネットワークパラメータ１３'について、複数５０の再構成レベルセット５２のうち再構成レベルセット４８を、以前に復号化されたニューラルネットワークパラメータについてのデータストリーム１４から復号化された量子化インデックス５８を２値化したものの LSB （例えば、最下位ビット）部分または以前に復号化されたピン（例えば、２値決定）に応じて選択５４するよう構成されるとよい。 LSB 比較は、低い計算コストで実行されることができる。特に、状態遷移が用いられることができる。選択５４は、複数５０の再構成レベルセット５２のうちの量子化レベルセット４８のうちの現在のニューラルネットワークパラメータ１３'について、現在のニューラルネットワークパラメータ１３'に関連する状態に応じて複数５０の再構成レベルセット５２のうちの再構成レベルセット４８を現在のニューラルネットワークパラメータ１３'について決定することによって、かつ、直前のニューラルネットワークパラメータについてのデータストリームから復号化された量子化インデックス５８に応じて後続のニューラルネットワークパラメータについての状態を更新することによって、状態推移プロセスにより実行されることができる。また、例えば遷移表を使用することによる状態遷移以外の代替的なアプローチも同様に使用することができ、以下に規定される。

【００３２】

加えて、または代替的に、装置は、例えば、現在のニューラルネットワークパラメータ１３'について、複数５０の再構成レベルセット５２のうち再構成レベルセット４８を、以前に復号化されたニューラルネットワークパラメータについてのデータストリーム１４から復号化された量子化インデックス５８の２値関数の結果に応じて選択５４するように構成されることができる。２値関数は、例えば、量子化インデックス５８が偶数または奇数を表すかどうかを信号化する、ビット単位「*and*」演算を使用するパリティチェックであってよい。これは、量子化インデックス５８を符号化するために使用される再構成レベルセット４８に関する情報を、したがって、例えば、対応する符号化器で使用される再構成レベルセットの所定の順序のために、現在のニューラルネットワークパラメータ１３'を符号化するために使用される再構成レベルのセットに対して提供し得る。パリティは、前述した状態遷移のために使用されることができる。

【００３３】

さらに、実施形態によれば、装置は、例えば、現在のニューラルネットワークパラメータ１３'について、複数５０の再構成レベルセット５２のうち再構成レベルのセット４８を、以前に復号化されたニューラルネットワークパラメータについてデータストリーム１４

10

20

30

40

50

から復号化された量子化インデックス 5 8 のパリティに応じて選択 5 4 するように構成されることができる。パリティチェックは、例えばビット単位「and」演算を用いて、低い計算コストで実行されることができる。

【0034】

任意選択で、装置は、ニューラルネットワークパラメータ 1 3 についての量子化インデックス 5 6 を復号化し、ニューラルネットワークパラメータ 1 3 の間で共通の連続的な順序 1 4 ' に沿ってニューラルネットワークパラメータ 1 3 の逆量子化を実行するように構成されることができる。言い換えれば、両方のタスクに同じ順序が使用されることができる。

【0035】

図 3 は、実施形態による、ニューラルネットワークパラメータをデータストリームに符号化するための装置内で実行される量子化の概念を模式的に示す図である。図 3 は、ニューラルネットワーク層 1 0 a、1 0 b を含むニューラルネットワーク (NN) 1 0 を示し、層はニューロン 1 0 c を含み、相互接続された層のニューロンはニューロン相互接続 1 1 を介して相互接続されている。一例として、NN 層 (p - 1) 1 0 a および NN 層 (p) 1 0 b が示され、p は NN 層のインデックスであり、1 ≤ p ≤ NN の層の数である。ニューラルネットワークは、ニューラルネットワークパラメータ 1 3 によって定義またはパラメータ化され、それは任意に、ニューラルネットワーク 1 0 のニューロン相互接続 1 1 の重みに関連し得る。図 1 の隠れ層のニューロン 1 0 c は、図 3 の層 p (A、B、C、. . .) のニューロンを表してもよく、図 1 の入力層のニューロンは、図 3 に示される層 p - 1 (a、b、c、. . .) のニューロンを表すことができる。ニューラルネットワークパラメータ 1 3 は、図 1 のニューロン相互接続 1 1 の重みに関連づけることができる。

【0036】

異なる層のニューロン 1 0 c の関係は、図 1 において、ニューラルネットワークパラメータ 1 3 の行列 1 5 a によって表される。例えば、ネットワークパラメータ 1 3 がニューロン相互接続 1 1 の重みに関連する場合、行列 1 5 a は、例えば、行列要素が異なる層のニューロン 1 0 c 間の重み (例えば、層 p - 1 の場合は a、b、. . . 、層 p の場合は A、B、. . .) を表すように構成されることができる。

【0037】

この装置は、例えばシリアル 2 0 a (シリアライゼーション、直列化) において、ニューラルネットワークパラメータ 1 3 を順次符号化するように構成されている。この順次処理の間、量子化器 (再構成レベルセット) は変化させられる。この変化により、より少ない (またはより良いより密でない) レベルを有する量子化器を使用することができ、したがって、より小さい量子化インデックスを符号化することができ、この量子化から得られるニューラルネットワーク表現の品質が、必要な符号化ビットレートと比較して、一定の量子化器を使用する場合よりも改善される。詳細は後述する。特に、装置は、現在のニューラルネットワークパラメータ 1 3 ' について、以前に符号化されたニューラルネットワークパラメータについてデータストリーム 1 4 に符号化された量子化インデックス 5 8 に応じて複数 5 0 の再構成レベルセット 5 2 のうち再構成レベルセット 4 8 を選択 5 4 することによって、ニューラルネットワークパラメータ 1 3 を順次に符号化する。

【0038】

さらに、装置は、現在のニューラルネットワークパラメータ 1 3 ' を選択された再構成レベルセット 4 8 の 1 つの再構成レベル上に量子化 6 4 (Q) することによって、かつ、現在のニューラルネットワークパラメータについての量子化インデックス 5 6 がデータストリーム 1 4 中に量子化される 1 つの再構成レベルを示す現在のニューラルネットワークパラメータ 1 3 ' についての量子化インデックス 5 6 を符号化することによって、ニューラルネットワークパラメータ 1 3 を順次符号化するよう構成される。任意選択で、複数 5 0 の再構成レベルセット 5 2 のうち、本明細書で時々量子化器とも呼ばれる再構成レベルセット 5 2 の数は、例えばセット 0 およびセット 1 を用いて示されるように、2 つであることができる。

【0039】

10

20

30

40

50

実施形態によれば、図 3 に示すように、装置は、例えば、所定の量子化ステップサイズ（ $Q P$ ）によって複数 50 の再構成レベルセット 52 をパラメータ化 60 し、データストリーム 14 に所定の量子化ステップサイズに関する情報を挿入するように構成されていることができる。これにより、例えば、量子化効率を向上させるための適応的量子化が可能となり、ニューラルネットワークパラメータ 13 の符号化方法の変化が、所定の量子化ステップサイズに関する情報とともに復号器に伝達されることができる。所定の量子化ステップサイズ（ $Q P$ ）を使用することにより、情報の伝送のためのデータ量を削減することができる。

【0040】

さらに、実施形態によれば、ニューラルネットワーク 10 は、1 つ以上の NN 層 10 a、10 b を含んでいることができる。装置は、各 NN 層（ p ； $p - 1$ ）について、それぞれの NN 層についての所定の量子化ステップサイズ（ $Q P$ ）の情報をデータストリーム 14 に挿入し、それぞれの NN 層に属するニューラルネットワークパラメータの量子化に使用するように、それぞれの NN 層について導かれた所定の量子化ステップサイズを用いて複数 50 の再構成レベルセット 52 のパラメータ化を行うよう構成されることが可能である。先に説明したように、例えば NN 層または NN 層の特性に応じた量子化の適応は、量子化効率を向上させることができる。

【0041】

任意選択的に、装置は、現在のニューラルネットワークパラメータ 13' について、複数 50 の再構成レベルセット 52 のうち再構成レベルのセット 48 を、以前に符号化されたニューラルネットワークパラメータについてデータストリーム 14 に符号化された量子化インデックス 58 を 2 値化したものの LSB 部分または以前に符号化されたピンに応じて選択 54 するよう構成されることができる。LSB 比較は、低い計算コストで実行され得る。

【0042】

図 2 で説明した復号化のための装置と同様に、状態遷移が用いられることができる。現在のニューラルネットワークパラメータ 13' について、現在のニューラルネットワークパラメータ 13' に関連する状態に応じて、複数 50 の再構成レベルセット 52 のうちの再構成レベルのセット 48 を決定することによって、及び、直前のニューラルネットワークパラメータについてのデータストリームに符号化された量子化インデックス 58 に応じて、後続のニューラルネットワークパラメータについての状態を更新することによって、状態遷移プロセスにより、複数 50 の再構成レベルセット 52 のうちの量子化レベルセット 48 のうち現在のニューラルネットワークパラメータ 13' について選択 54 が実行されることができる。例えば遷移表を使用することによる状態遷移以外のアプローチも同様に使用可能であり、以下に規定される。

【0043】

さらに、または代替的に、装置は、現在のニューラルネットワークパラメータ 13' について、以前に符号化されたニューラルネットワークパラメータについてのデータストリーム 14 に符号化された量子化インデックス 58 の 2 値関数の結果に応じて、複数 50 の再構成レベルセット 52 のうち再構成レベルセット 48 を選択 54 するよう構成されることができる。2 値関数は、例えば、量子化インデックス 58 が偶数または奇数を表すかどうかを信号化する、ビット単位「and」演算を使用したパリティチェックであってもよい。これは、量子化インデックス 58 を符号化するために使用される再構成レベルセット 48 に関する情報を提供してもよく、したがって、例えば、所定の順序のために、対応する復号化器が対応する再構成レベルセット 48 を選択できるように、再構成レベルの所定の順序のために、現在のニューラルネットワークのパラメータ 13' についての再構成レベルのセット 48 を決定することができる。パリティは、前述した状態遷移のために使用されることができる。

【0044】

さらに、実施形態によれば、装置は、例えば、現在のニューラルネットワークパラメー

10

20

30

40

50

タ 1 3 ' について、複数 5 0 の再構成レベルセット 5 2 のうち量子化レベルセット 4 8 を、以前に符号化されたニューラルネットワークパラメータについてのデータストリーム 1 4 に符号化された量子化インデックス 5 6 のパリティに応じて、選択 5 4 するように構成されることができる。パリティチェックは、例えばビット単位「a n d」演算を使用して、低い計算コストで実行されることができる。

【 0 0 4 5 】

任意選択的に、装置は、ニューラルネットワークパラメータ (1 3) についての量子化インデックス (5 6) を符号化し、ニューラルネットワークパラメータ (1 3) の量子化を、ニューラルネットワークパラメータ (1 3) 間の共通の連続的な順序 (1 4 ') に沿って実行するように構成されることができる。すなわち、両方のタスクで同じ順序を使用することができる。

10

【 0 0 4 6 】

図 4 は、実施形態による量子化されたニューラルネットワークのパラメータを算術復号化するための概念の概略図である。これは、図 2 の装置内で使用することができる。したがって、図 4 は、図 2 の可能な拡張として見ることができる。それは、現在のニューラルネットワークパラメータ 1 3 ' についての量子化インデックス 5 6 が、算術符号化、例えば 2 値算術符号化の使用によって任意例として示すように、図 4 の装置によって復号化されるデータストリーム 1 4 を示す。例えばあるコンテキストによって定義される確率モデルが使用され、それは、矢印 1 2 3 によって示されるように、現在のニューラルネットワークパラメータ 1 3 ' について選択された再構成レベルセット 4 8 に依存する。詳細は、本明細書において設定される。

20

【 0 0 4 7 】

図 2 に関して説明したように、現在のニューラルネットワークパラメータ 1 3 ' に関連する状態に応じて、現在のニューラルネットワークパラメータ 1 3 ' について複数 5 0 の再構成レベルセット 5 2 のうちの再構成レベルのセット 4 8 を決定することによって、及び、直前のニューラルネットワークパラメータについてのデータストリームから復号化された量子化インデックス 5 8 に応じて、後続のニューラルネットワークパラメータについての状態を更新することによって、状態遷移プロセスにより、複数 5 0 の再構成レベルセット 5 2 のうちの量子化レベルのセット 4 8 を選択する選択 5 4 が現在のニューラルネットワークパラメータ 1 3 ' について実行されることができる。したがって、状態は、現在のニューラルネットワークパラメータ 1 3 ' を符号化 / 復号化するために使用されるべき再構成レベルセット 4 8 へのポイントに準じるものであるが、これは、しかしながら、状態が、過去のニューラルネットワークパラメータまたは過去の量子化インデックスのメモリとして疑似的に機能するように、再構成セットの数に対応する数の状態を区別するだけとして、より細かい要素で更新される。したがって、状態は、ニューラルネットワークパラメータ 1 3 を符号化 / 復号化するために使用される再構成レベルのセットの順序を規定する。図 4 によれば、例えば、現在のニューラルネットワークパラメータ (1 3 ') についての量子化インデックス (5 6) は、現在のニューラルネットワークパラメータ (1 3 ') についての状態に対応する (1 2 2) 確率モデルを用いて、算術符号化を用いてデータストリーム (1 4) から復号化される。状態に応じて確率モデルを適応させることにより、確率モデルの推定が良好となり、符号化効率が向上する可能性がある。さらに、状態に応じて適応させることで、少ない追加データの送信で計算効率のよい適応が可能になる場合がある。

30

40

【 0 0 4 8 】

さらなる実施形態によれば、装置は、例えば、量子化インデックス 5 6 を 2 値化 8 2 したものの少なくとも 1 つのピン 8 4 についての現在のニューラルネットワークパラメータ 1 3 ' についての状態に依存 1 2 2 する確率モデルを用いて、2 値算術コーディングを使ってデータストリーム 1 4 から現在のニューラルネットワークパラメータ 1 3 ' についての量子化インデックス 5 6 を復号化するよう構成されることができる。

【 0 0 4 9 】

さらに、または代替的に、装置は、確率モデルの依存性が、依存性を用いるニューラル

50

ネットワークパラメータについてのコンテキストのセットのうちコンテキスト 87 の選択 103 (導出) を含むように構成されてもよく、各コンテキストは、所定の確率モデルが関連づけられるように構成される。使用される確率推定が優れているほど、圧縮の効率は高くなる。確率モデルは、例えば、コンテキスト適応的 (2 値) 算術符号化を用いて、更新されることができる。

【0050】

任意選択で、装置は、それぞれのコンテキストを用いて算術符号化された量子化インデックスに基づいて、それぞれのコンテキストに関連する所定の確率モデルを更新するように構成されることができる。このように、コンテキストの確率モデルは、実際の統計に適応される。

10

【0051】

さらに、装置は、例えば、量子化インデックスを 2 値化したものの少なくとも 1 つのピンについて、現在のニューラルネットワークパラメータ 13' に対して選択された再構成レベルセット 48 に対応する確率モデルを用いて、2 値算術符号化を用いてデータストリーム 14 から現在のニューラルネットワークパラメータ 13' についての量子化インデックス 56 を復号化するように構成されることができる。

【0052】

任意選択で、少なくとも 1 つのピンは、現在のニューラルネットワークパラメータの量子化インデックス 56 がゼロに等しいか否かを示す有意性ピンを含んでいることができる。さらに、または代替的に、少なくとも 1 つのピンは、現在のニューラルネットワークパラメータの量子化インデックス 56 がゼロより大きい、またはゼロより小さいかを示す符号ピンを含んでもよい。さらに、少なくとも 1 つのピンは、現在のニューラルネットワークパラメータの量子化インデックス 56 の絶対値が X より大きいか否かを示す *greater-than-X* ピンを含んでもよく、ここで X はゼロより大きい整数である。

20

【0053】

以下、図 5 では、図 4 を用いて説明した復号化のための概念の対極を説明することができる。したがって、すべての説明および利点は、以下の符号化のための概念の側面に適宜適用することができる。

【0054】

図 5 は、実施形態によるニューラルネットワークパラメータを算術符号化するための概念を示す概略図である。これは、図 3 の装置内で使用され得る。したがって、図 5 は、図 3 の可能な拡張として見るることができる。それは、現在のニューラルネットワークパラメータ 13' についての量子化インデックス 56 が、算術符号化、例えば任意の例として 2 値算術符号化の使用によって示されるように図 3 の装置によって符号化されるデータストリーム 14 を示す。例えばあるコンテキストによって定義される確率モデルが使用され、それは、矢印 123 によって示されるように、現在のニューラルネットワークパラメータ 13' に対して選択された再構成レベルセット 48 に依存する。詳細は本明細書において設定される。

30

【0055】

図 3 に関して説明したように、現在のニューラルネットワークパラメータ 13' について選択 54 が実行される。選択 54 は、現在のニューラルネットワークパラメータ 13' について、現在のニューラルネットワークパラメータ 13' に関連する状態に応じて、複数 50 の再構成レベルセット 52 のうちの量子化レベルセット 48 を決定することによって、及び、直前のニューラルネットワークパラメータについてのデータストリームに符号化された量子化インデックス 58 に応じて、後続のニューラルネットワークパラメータについての状態を更新することによって、状態遷移プロセスにより複数 50 の再構成レベルセット 52 のうちの量子化レベルセット 48 を選択する。

40

【0056】

したがって、状態は、現在のニューラルネットワークパラメータ 13' を符号化 / 復号化するために使用されるべき再構成レベルセット 48 へのポインタに準ずるものであるが、

50

しかし、状態が疑似的に、過去のニューラルネットワークパラメータまたは過去の量子化インデックスのメモリとして機能するように、再構成セットの数に対応する数の状態を区別するだけとしてより細かい要素で更新されている。したがって、状態は、ニューラルネットワークパラメータ 1 3 を符号化 / 復号化するために使用される再構成レベルセットの順序を定義する。

【 0 0 5 7 】

さらに、現在のニューラルネットワークパラメータ 1 3 ' についての量子化インデックス 5 6 は、現在のニューラルネットワークパラメータ 1 3 ' についての状態 1 2 2 に対応する確率モデルを用いる算術符号化を用いてデータストリーム 1 4 に符号化されることができる。

10

【 0 0 5 8 】

例えば図 3 によれば、量子化インデックス 5 6 は、量子化インデックス 5 6 を 2 値化 8 2 したものの少なくとも 1 つのピン 8 4 に対する現在のニューラルネットワークパラメータ 1 3 ' に対する状態 1 2 2 に対応する確率モデルを用いて、2 値算術コーディングを用いてデータストリーム 1 4 へ符号化される。確率モデルは確率モデルの推定に適している可能性があるため、状態に応じて確率モデルを適応させると、符号化効率が向上する可能性がある。さらに、状態に応じた適応は、送信される追加データの量が少なく、計算効率のよい適応を可能にし得る。

【 0 0 5 9 】

さらに、または代替的に、本装置は、確率モデルの依存性が、依存性を用いたニューラルネットワークパラメータのコンテキストのセットのうちコンテキスト 8 7 の選択 1 0 3 (導出) を含み、各コンテキストは所定の確率モデルが関連づけられているように構成されることができる。

20

【 0 0 6 0 】

任意選択で、装置は、それぞれのコンテキストを使用して算術符号化された量子化インデックスに基づいて、それぞれのコンテキストに関連付けられた所定の確率モデルを更新するように構成されることができる。

【 0 0 6 1 】

さらに、装置は、例えば、量子化インデックスを 2 値化したものの少なくとも 1 つのピンについて、現在のニューラルネットワークパラメータ 1 3 ' について選択された再構成レベルセット 4 8 に対応する確率モデルを用いることによって、2 値算術符号化を用いて、データストリーム 1 4 に現在のニューラルネットワークパラメータ 1 3 ' の量子化インデックス 5 6 を符号化するよう構成されることができる。2 値算術符号化を使用するために、量子化インデックス 5 6 は、2 値化 (b i n a r i z a t i o n) されることができる。

30

【 0 0 6 2 】

任意選択で、少なくとも 1 つのピンは、現在のニューラルネットワークパラメータの量子化インデックス 5 6 がゼロに等しいか否かを示す有意性ピンを含むことができる。さらに、または代替的に、少なくとも 1 つのピンは、現在のニューラルネットワークパラメータの量子化インデックス 5 6 がゼロより大きい、またはゼロより小さいかを示す符号ピンを含むことができる。さらに、少なくとも 1 つのピンは、現在のニューラルネットワークパラメータの量子化インデックス 5 6 の絶対値が X より大きいか否かを示す g r e a t e r - t h a n - X ピンを含むことができ、ここで X はゼロより大きい整数である。

40

【 0 0 6 3 】

次に説明する実施形態は、次のような本願の別の態様に集中している。ニューラルネットワークのパラメータ化が段階または再構成層で符号化され、NN パラメータごとに、各ステージからの 1 つの値は、ニューラルネットワークの改良された / 強化された表現を得るために結合される必要があり、少なくとも 1 つがそれ自体ニューラルネットワークの妥当な表現を表すかもしれないが、低質で貢献する段階のいずれかに強化されているが、後者の可能性は本態様に必須ではないものである。

【 0 0 6 4 】

50

図 6 は、本発明による実施形態で使用するためのニューラルネットワークパラメータのための再構成層を使用する概念の概略図である。図 6 は、例えば第 2 の再構成層である再構成層 i 、例えば第 1 の再構成層である再構成層 $i - 1$ 、および例えば図 3 からの層 10b であるニューラルネットワーク (NN) 層 p が、例えば図 3 からの行列 15a などのアレイまたは行列の形態で表された層であることを示している。

【0065】

図 6 は、ニューラルネットワークを定義するニューラルネットワークパラメータ 13 を再構成するための装置 310 の概念を示している。したがって、この装置は、例えば連合学習プロセスの間に以前に送信された可能性があり、例えば、第 1 の再構成層、例えば再構成層 $i - 1$ について、第 1 の再構成層のニューラルネットワークパラメータと呼ばれてもよい第 1 ニューラルネットワークパラメータ 13a を導出し、ニューラルネットワークパラメータごと、例えば重みごとまたはニューロン間接続ごとに、第 1 の再構成層のニューラルネットワークパラメータ値をもたらすよう構成されている。この導出は、そうでなければ、第 1 のニューラルネットワークパラメータ 13a を復号化すること、または受け取ることを含むことができる。さらに、装置は、ニューラルネットワークパラメータ 13 ごとに第 2 の再構成層のニューラルネットワークパラメータ値を生成するためのデータストリーム 14 から第 2 の再構成層についての例えば最終ニューラルネットワークパラメータ、例えばパラメータ 13 と区別するために第 2 の再構成層のニューラルネットワークパラメータと呼ばれ得る第 2 ニューラルネットワークパラメータ 13b を復号化 312 するよう構成されている。したがって、第 1 および第 2 の再構成層の 2 つの寄与値が、NN パラメータごとに得られてもよく、第 1 および / または第 2 の NN パラメータ値の符号化 / 復号化は、図 2 および図 3 に従った依存量子化 (dependent quantization)、および / または図 4 および図 5 で説明したような量子化インデックスの算術符号化 / 復号化を使用することができる。第 2 のニューラルネットワークパラメータ 13b は、ニューラル表現の観点から自己完結した意味を持たず、単に、第 1 の表現層のパラメータと組み合わせたときに、ニューラルネットワーク表現、すなわち最終的なニューラルネットワークパラメータを導くだけかもしれない。

【0066】

さらに、装置は、各ニューラルネットワークパラメータについて、第 1 の再構成層のニューラルネットワークパラメータ値と第 2 再構成層ニューラルネットワークパラメータ値とを、例えば要素単位の加算および / または乗算を用いて組み合わせる (CB) ことにより、ニューラルネットワークパラメータ 13 を再構成 314 するように構成されている。

【0067】

さらに、図 6 は、第 1 の再構成層、例えば再構成層 $i - 1$ のための第 1 のニューラルネットワークパラメータ 13a を用いて、ニューラルネットワークを定義するニューラルネットワークパラメータ 13 を符号化するための装置 320 の概念を示している。第 1 の再構成層は、ニューラルネットワークパラメータ 13 ごとに、第 1 の再構成層のニューラルネットワークパラメータ値を含む。したがって、装置は、第 2 の再構成層、例えば再構成層 i のための第 2 のニューラルネットワークパラメータ 13b をデータストリームに符号化 322 するように構成される。第 2 の再構成層は、ニューラルネットワークパラメータ 13 ごとに、第 2 の再構成層のニューラルネットワークパラメータ値を含む。ニューラルネットワークパラメータ 13 は、それぞれのニューラルネットワークパラメータについて、例えば要素単位の加算および / または乗算を用いて、第 1 の再構成層のニューラルネットワークパラメータ値と第 2 の再構成層のニューラルネットワークパラメータ値とを組み合わせる (CB) ことによって再構成可能である。

【0068】

任意選択的に、装置 310 は、データストリーム 14 から、または別個のデータストリームから、第 1 の再構成層についての第 1 のニューラルネットワークパラメータを復号化 316 するように構成されることができる。

【0069】

10

20

30

40

50

簡単に言えば、ニューラルネットワークパラメータ13の分解は、パラメータのより効率的な符号化および/または復号化および送信を可能にし得る。

【0070】

以下では、特に、ニューラルネットワーク符号化概念を含む、さらなる実施形態が開示される。以下の説明では、上述した実施形態と個別に、および組み合わせて使用することができる更なる詳細を提供する。

【0071】

まず、本発明の実施形態による依存スカラー量子化(Dependent Scalar Quantization)を伴うニューラルネットワークのパラメータのエントロピー符号化のための方法を提示する。

10

【0072】

依存スカラー量子化を用いたニューラルネットワークパラメータ13(重み、重みパラメータまたはパラメータとも呼ばれる)のセットのパラメータ符号化方法について説明する。本明細書で提示されるパラメータ符号化は、パラメータ13の依存スカラー量子化(例えば、図3の文脈で説明したような)および得られた量子化インデックス56のエントロピー符号化(例えば、図5の文脈で説明したような)から構成される。復号化器側では、量子化インデックス56のエントロピー復号化(例えば、図4の文脈で説明したように)と、ニューラルネットワークパラメータ13の依存再構成(例えば、図2の文脈で説明したように)とによって、再構成されたニューラルネットワークパラメータのセットを得ることができる。独立したスカラー量子化およびエントロピー符号化を伴うパラメータ符号化とは対照的に、ニューラルネットワークパラメータ13についての許容可能な再構成レベルセットは、再構成順序において現在のニューラルネットワークパラメータ13'に先行する送信済み量子化インデックス56に依存する。以下に示す提示は、依存スカラー量子化で使用される再構成レベルを指定する量子化インデックスのエントロピー符号化のための方法を追加的に説明する。

20

【0073】

本説明は、主にニューラルネットワーク圧縮におけるニューラルネットワークパラメータ層の非可逆符号化を対象としているが、他の分野の非可逆符号化にも適用可能である。

【0074】

本装置の方法論は、以下のような異なる主要部分に分けられる。

30

【0075】

1. 量子化
2. ロスレス符号化
3. ロスレス復号化

【0076】

以下に示す実施形態の主な利点を理解するために、まず、ニューラルネットワークの話題と、パラメータ符号化のための関連する方法について、簡単に紹介する。それにもかかわらず、開示されたすべての側面、特徴、および概念は、本明細書に記載される実施形態と別々にまたは組み合わせて使用することができる。

【0077】

40

2 量子化及びエントロピー符号化のための関連する方法

マルチメディアコンテンツの説明及び分析のためのニューラルネットワークの圧縮のためのMPEG-7パート17規格のワーキングドラフト2[2]は、ニューラルネットワークパラメータ符号化に独立スカラー量子化(independent scalar quantization)及びエントロピー符号化を適用している。

【0078】

2.1 スカラー量子化器

ニューラルネットワークのパラメータは、スカラー量子化器を用いて量子化される。量子化の結果、パラメータ13の許容値のセットは減少する。言い換えれば、ニューラルネットワークのパラメータは、いわゆる再構成レベルの可算集合(実際には有限集合)にマッ

50

ピングされる。再構成レベルセットは、可能なニューラルネットワークパラメータ値のセット（集合）の適切なサブセット（部分集合）を表す。以下のエントロピー符号化を単純化するために、許容可能な再構成レベルは、量子化インデックス 5 6 によって表され、これはビットストリーム 1 4 の一部として伝送される。復号化器側では、量子化インデックス 5 6 は、再構成されたニューラルネットワークパラメータ 1 3 にマッピングされる。再構成されたニューラルネットワークパラメータ 1 3 の可能な値は、再構成レベルセット 5 2 に対応する。符号化器側では、スカラー量子化の結果は、1 セットの（整数）量子化インデックス 5 6 である。

【 0 0 7 9 】

このアプリケーションでは、均一再構成量子化器（URQ）が使用される。その基本設計を図 7 に示す。図 7 は、均一再構成量子化器の説明図である。URQ は、再構成レベルが等間隔に配置されるという特性を持つ。隣接する 2 つの再構成レベル間の距離 Δ （QP）を量子化ステップサイズと呼ぶ。再構成レベルの 1 つは 0 に等しい。したがって、利用可能な再構成レベルの完全なセット、例えば $s'_i, i \in \mathbb{N}_0$, は、量子化ステップサイズ Δ （QP）によって一意に特定される。量子化インデックス q 5 6 の再構成された重みパラメータ t' 1 3' への復号化器のマッピングは、原理的に、単純な式で与えられる。

$$t' = q \cdot \Delta.$$

【 0 0 8 0 】

このコンテキストでは、「独立スカラー量子化」という用語は、任意の重みパラメータ 1 3 に対する量子化インデックス q 5 6 が与えられると、関連する再構成された重みパラメータ t' 1 3' が他の重みパラメータに対するすべての量子化インデックスから独立して決定できる、という特性を指す。

【 0 0 8 1 】

2 . 1 . 1 符号化器の動作：量子化

ニューラルネットワークの圧縮に関する標準規格は、ビットストリームのシンタックスと再構成プロセスのみを規定している。与えられたオリジナルのニューラルネットワークパラメータ 1 3 のセットと与えられた量子化ステップサイズ（QP）に対するパラメータ符号化を考える場合、符号化器は多くの自由度を有する。層 1 0 a、1 0 b の量子化インデックス q_k 5 6 が与えられると、エントロピー符号化は、データをビットストリーム 1 4 に書き込む（すなわち、算術符号語（コードワード）を構築する）ための一意に定義されたアルゴリズムに従わなければならない。しかし、重みパラメータのオリジナルセット（例えば層）が与えられた量子化インデックス q_k 5 6 を得るための符号化アルゴリズムは、ニューラルネットワーク圧縮の規格の範囲外である。以下の説明では、各ニューラルネットワークパラメータ 1 3 の量子化ステップサイズ（QP）が既知であると仮定する。それでも、符号化器は、各ニューラルネットワーク（重み）パラメータ t_k 1 3 についての量子化器インデックス q_k 5 6 を選択する自由を有する。量子化インデックスの選択は、歪み（または再構成 / 近似品質）とビットレートの両方を決定するので、使用される量子化アルゴリズムは、生成されるビットストリーム 1 4 のレート歪み性能に実質的な影響を与える。

【 0 0 8 2 】

10

20

30

40

50

最も単純な量子化方法は、ニューラルネットワークパラメータ t_{k13} を最も近い再構成レベルに丸める（最近傍量子化（*nearest neighbor quantization*）とも呼ばれる）。一般的に使用される URQ について、対応する量子化インデックス q_{k56} は、以下に従って決定することができる。

$$q_k = \text{sgn}(t_k) \cdot \left\lfloor \frac{|t_k|}{\Delta_k} + \frac{1}{2} \right\rfloor, \quad 10$$

ここで、 $\text{sgn}()$ は符号関数、演算子 $\lfloor \cdot \rfloor$ はその引数と等しいか小さい最大の整数を返す。この量子化法によって、MSE 歪みは最小になることが保証される。

$$D = \sum_k D_k = \sum_k (t_k - q_k \cdot \Delta_k)^2$$

しかし、結果として得られるパラメータレベル（重みレベル） q_{k56} を送信するために必要とされるビットレートは完全に無視される。なお、この方法は、MSE 歪み測定に限定されるものではなく、他の任意の歪み測定、例えば、以下の方法による MAE 歪みを使用することができる。

$$D^{MAE} = \sum_k D_k^{MAE} = \sum_k |t_k - q_k \cdot \Delta_k|$$

一般的に、丸め方はゼロに偏った方が良い結果が得られる。

$$q_k = \text{sgn}(t_k) \cdot \left\lfloor \frac{|t_k|}{\Delta_k} + a \right\rfloor \quad \text{with} \quad 0 \leq a < \frac{1}{2}. \quad 20$$

【0083】

量子化プロセスはラグランジュ関数 $D + \lambda \cdot R$ を最小化すれば、レート歪み的に良い結果が得られる。ここで、 D はニューラルネットワークパラメータセットの歪み（例えば、MSE 歪み、又は、MAE 歪み）、 R は量子化インデックス 56 を伝送するために必要なビット数、 λ はラグランジュ乗数である。

【0084】

10

20

30

40

50

量子化ステップサイズが与えられると、ラグランジュ乗数 λ と量子化ステップサイズとの間に以下の関係がしばしば用いられる。

$$\lambda = c_1 \cdot \Delta^2,$$

ここで、 c_1 はニューラルネットワークパラメータのセットに対する定数ファクタを表す。歪みとレートとのラグランジュ関数 $D + \lambda \cdot R$ を最小化することを目的とした量子化アルゴリズムは、レート歪み最適化量子化（RDOQ、*rate-distortion optimized quantization*）とも呼ばれる。歪みをMSEまたは重み付きMSE（またはそれぞれMAE）を用いて測定する場合、重みパラメータのセット（例えば層）に対する量子化指標 q_{k56} は、以下のコスト量が最小となるように決定する必要がある。

$$D + \lambda \cdot R = \sum_k \alpha_k \cdot (t_k - \Delta_k \cdot q_k)^2 + \lambda \cdot R(q_k | q_{k-1}, q_{k-2}, \dots).$$

【0085】

このとき、ニューラルネットワークパラメータインデックス k は、ニューラルネットワークパラメータ13の符号化順序（または走査順序）を指定する。 $R(q_k | q_{k-1}, q_{k-2}, \dots)$ の項は、量子化インデックス q_{k56} を送信するために必要なビット数（またはその推定値）を表している。この条件は、（結合確率または条件付き確率の使用により）特定の量子化インデックス q_k についてのビット数が、典型的には、符号化順序、例えば共通の連続的な順番14'において先行する量子化インデックス q_{k-1}, q_{k-2} ,等のための選択値に依存していることを示している。上式における係数 α_k は、個々のニューラルネットワークパラメータ13の寄与度を重み付けするために用いることができる。以下では、一般に、すべての重み付け係数 α_k が1に等しいと仮定する（ただし、異なる重み付け係数を考慮することができるように、アルゴリズムを端的に修正することができる）。

【0086】

実際、近傍量子化は $\lambda = 0$ の些細なケースであり、マルチメディアコンテンツの説明及び分析のためのニューラルネットワークの圧縮に関するMPEG-7 part 17規格のワーキングドラフト2において適用されている。

【0087】

2.2 エントロピー符号化

前のステップで適用された均一量子化の結果として、重みパラメータはいわゆる再構成レベルの有限集合にマッピングされる。これらは、（整数）量子化器インデックス56（パラメータレベルまたは重みレベルとも呼ばれる）と量子化ステップサイズ（QP）によって表すことができ、例えば、全層に対して固定されている場合がある。層のすべての量子化された重みパラメータを復元するために、層のステップサイズ（QP）および次元は、復号化器によって知られてもよい。これらは、例えば、別々に送信されてもよい。

【0088】

2.2.1 コンテキスト適応的2値算術符号化（CABAC、Context-adaptive

ptive binary arithmetic coding)による量子化インデックスの符号化

量子化インデックス56(整数表現)は、次にエントロピー符号化技術を使用して送信される。したがって、重みの層は、スキャンを使用して量子化された重みレベルのシーケンスにマッピングされる。例えば、行列の最上部の行から始めて、含まれる値を左から右へ符号化する、行ファーストスキャン順序(row first scan order)を使用することができる。この方法では、すべての行が上から下へ符号化される。スキャンは、ニューロン相互接続11の重みに関連し得るニューラルネットワークパラメータ13を含む行列15aについて、例えば共通の連続的な順番14'に沿って、図3に示すように実行されてもよい。行列は、重み層、例えば、図3及び図1にそれぞれ示すように、ニューロン相互接続11の層p-1 10aと層p 10bとの間の重み、又は隠れ層及び入力層との間の重みを表してもよい。なお、他の任意のスキャンを適用することができる。例えば、行列(例えば、図2又は図3の行列15a)は、行ファーストスキャンを適用する前に、転置され、又は水平及び/又は垂直に反転され、及び/又は左又は右に90/180/270度だけ回転されることができる。

【0089】

図3および図5に関して説明したように、実施形態による装置は、量子化インデックス56を2値化82したものの少なくとも1つのピン84に対する現在のニューラルネットワークパラメータ13'に対する状態122に対応する確率モデルを用いて、2値算術符号化を用いてデータストリーム14へ現在のニューラルネットワークパラメータ13'についての量子化インデックス56を符号化するように構成されることができる。確率モデルを用いた2値算術符号化は、コンテキスト適応的2値算術符号化(CABAC、Context-adaptive binary arithmetic coding)であることができる。

【0090】

すなわち、実施形態によれば、レベルの符号化のために、CABACが使用される。詳細については、[3]を参照されたい。そこで、量子化された重みレベルq56は、一連の2値記号またはシンタックス要素、例えばピン(2値決定)に分解され、その後、2値算術符号化器(CABAC)に渡されることがある。最初のステップでは、量子化された重みレベルに対して2値シンタックス要素sig_flagが導出され、これは対応するレベルがゼロに等しいかどうかを指定する。言い換えれば、図4に示す量子化インデックス56を2値化82したものの少なくとも1つのピンは、現在のニューラルネットワークパラメータの量子化インデックス56がゼロに等しいか否かを示す有意性ピンを含むことができる。

【0091】

sig_flagが1に等しい場合、さらなる2値シンタックス要素sign_flagが導き出される。このピンは、現在の重みレベルが正(例えば、ピン=0)か負(例えば、ピン=1)かを示す。言い換えれば、図4に示す量子化インデックス56を2値化82したものの少なくとも1つのピンは、現在のニューラルネットワークパラメータの量子化インデックス56がゼロより大きいかまたはゼロより小さいかを示す符号ピン86を含むことができる。

【0092】

次に、ピンの単項シーケンスが符号化され、それに続いて、以下のように固定長シーケンスが次のように符号化される。

【0093】

変数kは負でない整数で初期化され、Xは $1 < k$ で初期化される。

【0094】

量子化された重みレベルの絶対値がXより大きいことを示すabs_level_greater_Xという一つ以上のシンタックス要素が符号化される。abs_level_greater_Xが1に等しい場合、変数kが更新され(例えば、1だけ増加)、次

に $1 < k$ が X に加えられ、さらに $abs_level_greater_X$ が符号化される。この手順は、 $abs_level_greater_X$ が 0 に等しくなるまで続けられる。その後、長さ k の固定長コードで量子化器インデックスの符号化を完了することができる。例えば、変数 $rem = X - |q|$ は、 k ビットを用いて符号化され得る。あるいは、変数 rem' は、 $rem' = (1 < k) - rem - 1$ として定義され、これは k ビットを用いて符号化され得る。また、変数 rem の k ビットの固定長符号への他のマッピングも使用できる。

【0095】

言い換えれば、図 4 に示す量子化インデックス 56 を 2 値化 82 したものの少なくとも 1 つのピンは、現在のニューラルネットワークパラメータの量子化インデックス 56 の絶対値が X より大きいかなを示す $greater_than_X$ ピンを含むことができ、ここで X はゼロより大きい整数である。

10

【0096】

各 $abs_level_greater_X$ の後に k を 1 ずつ増加させる場合、このアプローチは、指数ゴロム符号化 (exponential Golomb coding) を適用することと同一である ($sign_flag$ がみなされていない場合)。

【0097】

また、符号化器側と復号化器側で絶対値の最大値 abs_max が分かっている場合、次に送信する $abs_level_greater_X$ について、 $X \geq abs_max$ が成り立つとき、 $abs_level_greater_X$ シンタックス要素の符号化を終了することができる。

20

【0098】

2. 2. 2 コンテキスト適応的 2 値算術符号化 (CABAC) による量子化インデックスの復号化

量子化された重みレベル 56 (整数表現) の復号化は、符号化と同様に動作する。復号化器は、まず $sign_flag$ を復号化する。もしそれが 1 に等しければ、 $sign_flag$ と $abs_level_greater_X$ の単項シーケンスが続く。ここで、 k の更新 (したがって X の増分) は符号化器と同じルールに従わなければならない。最後に、 k ビットの固定長符号が復号化され、整数値として解釈される (例えば、 rem または rem' として。両者のどちらが符号化されたかに応じている)。そして、復号化された量子化された重みレベルの絶対値 $|q|$ を X から再構成し、固定長部分を形成することができる。例えば、固定長部分として rem が使用された場合、 $|q| = X - rem$ となる。あるいは、 rem' が符号化されていた場合、 $|q| = X + 1 + rem' - (1 < k)$ 。最後のステップとして、復号化された $sign_flag$ に対応して $|q|$ に符号を適用する必要がある、量子化された重みレベル q 56 を得ることができる。最後に、量子化された重みレベル q にステップサイズ Δ (QP) を乗じることにより、量子化された重み w が再構成される。

30

【0099】

言い換えれば、図 2 および図 4 に関して説明したように、実施形態による装置は、量子化インデックス 56 を 2 値化 82 したものの少なくとも 1 つのピン 84 についての現在のニューラルネットワークパラメータ 13' に対する状態 122 に対応する確率モデルを使用することにより、2 値算術符号化を使用してデータストリーム 14 から現在のニューラルネットワークパラメータ 13' に対する量子化インデックス 56 を復号化するよう構成されることができる。

40

【0100】

図 5 に示す量子化インデックス 56 を 2 値化 82 したものの少なくとも 1 つのピンは、現在のニューラルネットワークパラメータの量子化インデックス 56 がゼロに等しいかなを示す有意性ピンを含むことができる。さらに、または代替的に、少なくとも 1 つのピ

50

ンは、現在のニューラルネットワークパラメータの量子化インデックス 56 がゼロより大きい、またはゼロより小さいかを示す符号ビン 86 を含むことができる。さらに、少なくとも 1 つのビンは、現在のニューラルネットワークパラメータの量子化インデックス 56 の絶対値が X より大きい、または小さいかを示す $greater - than - X$ ビンを含むことができ、ここで X はゼロより大きい整数である。

【0101】

好ましい実施形態では、 k は 0 に初期化され、以下のように更新される。各 $abs_level_greater_X$ が 1 に等しくなった後、 k の必要な更新は、以下のルールに従って行われる： $X > X'$ の場合、 k は 1 だけ増分され、 X' はアプリケーションに対応する定数である。例えば、 X' は符号化器が導出し、復号化器に通知する数値（例えば、0 から 100 の間）である。

10

【0102】

2.2.3 コンテキスト・モデリング

CABAC エントロピー符号化では、量子化された重みレベル 56 のほとんどのシンタックス要素は、2 値確率モデリングを用いて符号化される。各 2 値決定 (bin) はコンテキストと関連付けられている。コンテキストは、符号化されたビンのクラスに対する確率モデルを表す。2 つの可能なビン値のうちの 1 つに対する確率は、対応するコンテキストで既に符号化されたビンの値に基づいて、各コンテキストについて推定される。アプリケーションに応じて、異なるコンテキストモデリングアプローチを適用することができる。通常、量子化された重み符号化に関連するいくつかのビンについて、符号化に使用されるコンテキストは、既に送信されたシンタックス要素に基づいて選択される。実際のアプリケーションに応じて、例えば SBMP 0、または HEVC 0 または VTM - 4.0 のものなど、異なる確率推定器が選択され得る。この選択は、例えば、圧縮効率や複雑さに影響を与える。

20

【0103】

言い換えれば、図 5、例えばコンテキスト 87 に関して説明したような確率モデルは、以前に符号化されたニューラルネットワークパラメータの量子化インデックスに追加的に依存する。

【0104】

それぞれ、図 4、例えばコンテキスト 87 に関して説明したような確率モデルは、さらに、以前に復号化されたニューラルネットワークパラメータの量子化インデックスに依存する。

30

【0105】

広範囲のニューラルネットワークに適合するコンテキストモデリング方式は、以下のように説明される。重み行列 (層) の特定の位置 (x, y) で量子化された重みレベル q_{56} を復号化するために、ローカルテンプレートが現在の位置に適用される。このテンプレートは、例えば ($x - 1, y$)、($x, y - 1$)、($x - 1, y - 1$) 等のような多数の他の (順序付けられた) 位置を含んでいる。各位置に対して、ステータス識別子が導出される。

【0106】

40

好ましい実施形態 ($Si1$ と表す) において、位置 (x, y) に対するステータス識別子 $s_{x,y}$ は、以下のように導出される。位置 (x, y) が行列の外側を指す場合、または位置 (x, y) における量子化された重みレベル $q_{x,y}$ がまだ復号化されていないかゼロに等しい場合、ステータス識別子 $s_{x,y} = 0$ とする。それ以外の場合、ステータス識別子 $s_{x,y} = q_{x,y} < 0 ? 1 : 2$ であるものとする。

【0107】

特定のテンプレートに対して、ステータス識別子のシーケンスを導き出し、ステータス

50

識別子の値の可能な各コンステレーション (c o n s t e l l a t i o n) を、使用されるコンテキストを識別するコンテキストインデックスにマッピングする。テンプレートとマッピングは、異なるシンタックス要素に対して異なる場合がある。例えば、(順序) 位置 (x - 1 , y)、(x , y - 1)、(x - 1 , y - 1) を含むテンプレートから、ステータス識別子 $s_{x-1,y}$ 、 $s_{x,y-1}$ 、 $s_{x-1,y-1}$ の順序シーケンスが導出される。例えば、このシーケンスは、コンテキストインデックス $C = s_{x-1,y} + 3 * s_{x,y-1} + 9 * s_{x-1,y-1}$ にマッピングされることがある。例えば、コンテキストインデックス C は、 $s i g_f l a g$ のための多数のコンテキストを識別するために使用されることができる。

【 0 1 0 8 】

好ましい実施形態 (アプローチ 1 とする) において、位置 (x , y) における量子化された重みレベル $q_{x,y}$ の $s i g_f l a g$ のため、または $s i g n_f l a g$ のためのローカルテンプレートは、1つの位置 (x - 1 , y) (すなわち、左隣) からのみ構成される。関連するステータス識別子 $s_{x-1,y}$ は、好ましい実施形態 $S i 1$ に従って導出される。

10

【 0 1 0 9 】

$s i g_f l a g$ については、 $s_{x-1,y}$ の値に応じて3つのコンテキストのうちの1つが選択され、又は $s i g n_f l a g$ については、 $s_{x-1,y}$ の値に応じて他の3つのコンテキストのうちの1つが選択される。

【 0 1 1 0 】

別の好ましい実施形態 (アプローチ 2 とする) では、 $s i g_f l a g$ のためのローカルテンプレートは、3つの順序付けられた位置 (x - 1 , y)、(x - 2 , y)、(x - 3 , y) を含む。ステータス識別子 $s_{x-1,y}$ 、 $s_{x-2,y}$ 、 $s_{x-3,y}$ の関連するシーケンスは、好ましい実施形態 $S i 2$ に従って導出される。

20

【 0 1 1 1 】

$s i g_f l a g$ については、コンテキストインデックス C を以下のように導出する。

【 0 1 1 2 】

$s_{x-1,y} = 0$ ならば $C = 0$ 、それ以外なら $s_{x-2,y} = 0$ ならば $C = 1$ 、それ以外なら $s_{x-3,y} = 0$ ならば $C = 2$ 、そうでなければ、 $s_{x-3,y} = 0$ ならば、 $C = 2$ である。そうでなければ、 $C = 3$ 。

【 0 1 1 3 】

これは、次の式で表すこともできる。

30

$$C = (s_{x-1,y} \neq 0) ? 0 : ((s_{x-2,y} \neq 0) ? 1 : ((s_{x-3,y} \neq 0) ? 2 : 3))$$

【 0 1 1 4 】

同様に、コンテキストインデックス C が左側の次のゼロでない重みまでの距離に等しくなる (テンプレートサイズを超えない) ように、左側に隣接するものの数を増減してもよい。

【 0 1 1 5 】

各 $a b s_l e v e l_g r e a t e r_X_f l a g$ は、例えば、2つのコンテキストの独自のセットを適用することができる。そして、2つのコンテキストのうちの1つが、 $s i g n_f l a g$ の値に応じて選択される。

40

【 0 1 1 6 】

好ましい実施形態では、 X が予め定義された数 X' より小さい $a b s_l e v e l_g r e a t e r_X_f l a g$ について、異なるコンテキストが、 X および / または $s i g n_f l a g$ の値に応じて区別される。

【 0 1 1 7 】

好ましい実施形態では、 X が予め定義された数 X' より大きいか等しい $a b s_l e v e l_g r e a t e r_X_f l a g$ について、異なるコンテキストは、 X に対応してのみ区別される。

50

【 0 1 1 8 】

別の好ましい実施形態では、予め定義された数 X' より大きいか等しい X を有する $abs_level_greater_X_flag$ は、1 の固定コード長を用いて（例えば、算術符号器のバイパスモードを用いて）符号化される。

【 0 1 1 9 】

さらに、シンタックス要素の一部または全部は、コンテキストを使用せずに符号化されることもある。その代わりに、それらは、例えば、C A B A C のいわゆるバイパスピンを使用して、1 ビットの固定長で符号化される。

【 0 1 2 0 】

別の好ましい実施形態では、固定長の余り rem は、バイパスモードを使用して符号化される。

10

【 0 1 2 1 】

別の好ましい実施形態では、符号化器は、予め定義された数 X' を決定し、 $X < X'$ の各シンタックス要素 $abs_level_greater_X$ に対して、符号に応じて2つのコンテキストを区別し、 $X \geq X'$ の各 $abs_level_greater_X$ に対して、1つのコンテキストを使用する。

【 0 1 2 2 】

言い換えれば、図5に関して説明したような確率モデル、例えばコンテキスト87は、現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する以前に符号化されたニューラルネットワークパラメータの量子化インデックスに応じて確率モデルのサブセットの中から、現在のニューラルネットワークパラメータのために選択103されることができる。

20

【 0 1 2 3 】

この部分は、例えば、上記で説明したテンプレートであって、（順序付けられた）位置 $(x-1, y)$ 、 $(x, y-1)$ 、 $(x-1, y-1)$ を含むテンプレートによって定義されることができる。

【 0 1 2 4 】

それぞれ、図5に関して説明したような確率モデルは、現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する以前に復号化されたニューラルネットワークパラメータの量子化インデックスに応じて、確率モデルのサブセットの中から、現在のニューラルネットワークパラメータについて選択されることができる。

30

【 0 1 2 5 】

3 追加の方法

以下では、再構成された層、例えば図6からのニューラルネットワーク層 p が、例えば別々に伝送されてもよい図6からの再構成層 $i-1$ および再構成層 i のような異なるサブ層の構成である、ニューラルネットワーク10の圧縮/送信のための追加の、したがって任意の方法について説明する。

【 0 1 2 6 】

3.1 ベース層とエンハンスメント層の概念

40

この概念では、ベース層とエンハンスメント層と呼ばれる2種類のサブ層を導入している。そして、再構成プロセス（例えば、すべてのサブ層を追加する）は、サブ層からどのように再構成された層を得ることができるかを定義する。ベース層はベース値を含み、例えば、最初のステップで効率的に表現または圧縮/送信できるように選択することができる。エンハンスメント層は、エンハンスメント情報、（例えばオリジナル層に関する）歪み指標を減少させるために例えば（ベース）層の値に追加される差分値を含む。別の例では、ベース層は（小さなトレーニングセットを用いたトレーニングからの）粗い値を含み、エンハンスメント層は（完全なトレーニングセットまたはより一般的には、別のトレーニングセットに基づく）リファインメント値を含む。サブ層は別々に保存/送信されてもよい。

50

【 0 1 2 7 】

好ましい実施形態では、圧縮される層 L_R 、例えばニューラルネットワークパラメータ、例えば図 2 及び図 3 の行列 15 a によって表され得る重みのようなニューラルネットワークの重みは、ベース層 L_B と 1 つ以上のエンハンスメント層 $L_{E,1}$ 、 $L_{E,2}$ 、 \dots 、 $L_{E,N}$ とに分解される。そして、最初のステップでベース層が圧縮 / 送信され、続くステップでエンハンスメント層 $L_{E,1}$ 、 $L_{E,2}$ 、 \dots 、 $L_{E,N}$ が（別個に）圧縮 / 送信される。

【 0 1 2 8 】

別の好ましい実施形態では、再構成された層 L_R は、以下のように、すべてのサブ層 $L_{S,N}$ を（要素単位で）加算することによって得ることができる。

$$L_R = \sum_{i=0}^N L_{S,N}$$

10

【 0 1 2 9 】

さらなる好ましい実施形態において、再構成された層 L_R は、以下に従って、すべてのサブ層 $L_{S,N}$ を（要素単位で）乗算することによって得ることができる。

$$L_R = \prod_{i=0}^N L_{S,N}$$

20

【 0 1 3 0 】

言い換えれば、本発明による実施形態は、再構成された層 L_R の形態で、または例えば再構成された層 L_R を使用して、ニューラルネットワークパラメータごとに、第 1 の再構成層のニューラルネットワークパラメータ値と第 2 の再構成層のニューラルネットワークパラメータ値のパラメータ単位の和またはパラメータ単位の積によって、ニューラルネットワークパラメータ 13 を再構成するように構成された、装置を含む。

【 0 1 3 1 】

それぞれ、実施形態によるニューラルネットワークパラメータ 13 を符号化するための装置の場合、ニューラルネットワークパラメータ 13 は、ニューラルネットワークパラメータごとに、第 1 の再構成層のニューラルネットワークパラメータ値と第 2 の再構成層のニューラルネットワークパラメータ値のパラメータ単位の和またはパラメータ単位の積によって再構成可能である。

30

【 0 1 3 2 】

さらなる好ましい実施形態では、2 . 1 および / または 2 . 2 の方法は、サブセットまたはすべてのサブ層に適用される。

【 0 1 3 3 】

特に好ましい実施形態では、コンテキストモデリングを使用するエントロピー符号化方式（例えば、2 . 2 . 3 に同質または類似）が適用されるが、以下の規則のうちの 1 つ以上に従ってコンテキストモデルの 1 つ以上のセットを追加する。

40

【 0 1 3 4 】

a) 各サブ層は、独自のコンテキストセットを適用する。言い換えれば、本発明による実施形態は、第 1 の再構成層についての第 1 のニューラルネットワークパラメータ 13 a をデータストリームまたは別のデータストリームから / に符号化 / 復号化し、第 2 の再構成層についての第 2 ニューラルネットワークパラメータ 13 b を第 1 および第 2 の再構成層のための別の確率コンテキストを用いたコンテキスト適応的エントロピー符号化によってデータストリームから / に符号化 / 復号化するように構成された、装置を含む。

【 0 1 3 5 】

b) 符号化されるエンハンスメント層のパラメータのために選択されたコンテキスト

50

セットは、符号化順序において先行する層（例えば、ベース層）の同位置のパラメータの値に依存する。コンテキストモデルの第1のセットは、同位置のパラメータがゼロに等しいときは常に選択され、そうでないときは第2のセットが選択される。言い換えれば、本発明による実施形態は、第2の再構成層のニューラルネットワークパラメータ値、例えばエンハンスメント層のパラメータを、第1の再構成層のニューラルネットワークパラメータ値、例えば符号化順序において先行する層（例えばベース層）における同位置のパラメータの値に対応する確率モデルを用いたコンテキスト適応的エントロピー符号化によってデータストリームに符号化するように構成された、装置を備える。さらなる実施形態は、第1の再構成層のニューラルネットワークパラメータ値に応じて確率コンテキストセットの集合から確率コンテキストセットを選択することによって、及び、第1の再構成層のニューラルネットワークパラメータ値に応じて選択された確率コンテキストセットの中から使用する確率コンテキストを選択することによって、コンテキスト適応的エントロピー符号化によって第2の再構成層のニューラルネットワークパラメータ値をデータストリームに符号化するように構成される装置を含む。それぞれに、実施形態によるニューラルネットワークパラメータ13を復号化するための装置について、前記装置は、第1の再構成層のニューラルネットワークパラメータ値に対応する確率モデルを用いたコンテキスト適応的エントロピー復号化によってデータストリームから第2の再構成層のニューラルネットワークパラメータ値を復号化するように構成されることができる。それぞれ、さらなる実施形態は、第1の再構成層のニューラルネットワークパラメータ値に対応する確率コンテキストセットの集合から確率コンテキストセットを選択することによって、及び、第1の再構成層のニューラルネットワークパラメータ値に対応する選択された確率コンテキストセットの中から使用する確率コンテキストを選択することによって、コンテキスト適応的エントロピー復号化によってデータストリームから第2の再構成層のニューラルネットワークパラメータ値を復号化するように構成された、装置を含む。

【0136】

c) 符号化されるエンハンスメント層のパラメータについて選択されたコンテキストセットは、符号化順序において先行する層（例えば、ベース層）の同位置のパラメータの値に対応する。コンテキストモデルの第1のセットは、同位置のパラメータがゼロより小さい（負）場合は常に選択され、第2のセットは、同位置のパラメータがゼロより大きい（正）場合は選択され、それ以外は第3のセットが選択される。言い換えれば、本発明による実施形態は、例えば、符号化するための装置を含む。確率コンテキストセットの集合は3つの確率コンテキストセットを含む。装置は、第1の再構成層のニューラルネットワークパラメータ値が負の場合、確率コンテキストセットの集合から第1の確率コンテキストセットを選択された確率コンテキストセットとして選択するように構成され、第1の再構成層のニューラルネットワークパラメータ値が正である場合、確率コンテキストセットの集合から第2の確率コンテキストセットを選択された確率コンテキストセットとして選択するように構成され、第1の再構成層のニューラルネットワークパラメータ値がゼロである場合、確率コンテキストセットの集合から第3の確率コンテキストセットを選択された確率コンテキストセットとして選択するよう構成される。それぞれ、実施形態によるニューラルネットワークパラメータ13を復号化するための装置については、確率コンテキストセットの集合は、3つの確率コンテキストセットを含むことができ、装置は、第1の再構成層のニューラルネットワークパラメータ値が負の場合、確率コンテキストセットの集合から第1の確率コンテキストセットを選択した確率コンテキストセットとして選択するように構成されることができ、第1の再構成層のニューラルネットワークパラメータ値が正の場合、確率コンテキストセットの集合から第2の確率コンテキストセットを選択した確率コンテキストセットとして選択するように構成されることができ、第1の再構成層のニューラルネットワークパラメータ値がゼロである場合、確率コンテキストセットの集合から第3の確率コンテキストセットを選択した確率コンテキストセットとして選択するよう構成されることができる。

【0137】

d) 符号化されるエンハンスメント層のパラメータのために選択されるコンテキストセットは、符号化順序において先行する層（例えば、ベース層）の同位置のパラメータの値に対応する。コンテキストモデルの第1のセットは、同位置にあるパラメータの（絶対）値が X （ X はパラメータである）より大きいときは常に選択され、そうでないときは第2のセットが選択される。言い換えれば、本発明による実施形態は装置を含み、確率コンテキストセットの集合は、2つの確率コンテキストセットを含み、装置は、第1の再構成層のニューラルネットワークのパラメータ値、例えば、符号化順序において先行する層（例えば、ベース層）における同位置にあるパラメータの値が所定の値、例えば X 、より大きい場合、確率コンテキストセットの集合から第1の確率コンテキストセットを選択された確率コンテキストセットとして選択するように構成され、第1の再構成層のニューラルネットワークパラメータ値が所定値より大きくない場合、確率コンテキストセットの集合から第2の確率コンテキストセットを選択された確率コンテキストセットとして選択するように構成され、あるいは、第1の再構成層のニューラルネットワークパラメータ値の絶対値が所定値より大きい場合、確率コンテキストセットの集合から第1の確率コンテキストセットを選択された確率コンテキストセットとして選択するように構成され、第1の再構成層のニューラルネットワークパラメータ値の絶対値が所定値より大きくない場合、確率コンテキストセットの集合から第2の確率コンテキストセットを選択された確率コンテキストセットとして選択するように構成される。それぞれに、実施形態によるニューラルネットワークパラメータ13を復号化するための装置について、確率コンテキストの集合は2つの確率コンテキストセットを含むことができ、装置は、第1の再構成層のニューラルネットワークパラメータ値が所定の値、例えば X 、より大きい場合、確率コンテキストセットの集合から第1の確率コンテキストセットを選択された確率コンテキストセットとして選択するように構成され、第1の再構成層のニューラルネットワークパラメータ値が所定の値より大きくない場合、確率コンテキストセットの集合から第2の確率コンテキストセットを選択された確率コンテキストセットとして選択するように構成され、又は、第1の再構成層のニューラルネットワークパラメータ値の絶対値が所定値より大きい場合、確率コンテキストセットの集合から第1の確率コンテキストセットを選択された確率コンテキストセットとして選択するように構成され、第1の再構成層のニューラルネットワークパラメータ値の絶対値が所定値より大きくない場合、確率コンテキストセットの集合から第2の確率コンテキストセットを選択された確率コンテキストセットとして選択するように構成されることができる。

【0138】

4 依存スカラー量子化を伴うニューラルネットワークパラメータ符号化

このセクションでは、図2～図4の文脈で説明したような、本発明による概念および実施形態に対するさらなる任意の側面および特徴が開示される。

【0139】

以下では、ニューラルネットワークパラメータ符号化の修正された概念について説明する。先に説明したニューラルネットワークパラメータ符号化に対する主な変更点は、ニューラルネットワークパラメータ13が独立して量子化および再構成されないということである。その代わりに、ニューラルネットワークパラメータ13の許容される再構成レベルは、再構成順序において先行するニューラルネットワークパラメータの選択された量子化インデックス56に対応する。依存スカラー量子化の概念は、ニューラルネットワークパラメータについての確率モデル選択（または、代替的に符号語（コードワード）表選択）が許容される再構成レベルセットに依存する、修正されたエントロピー符号化と組み合わせられる。しかしながら、先に説明した実施形態は、以下に説明する特徴のいずれかを別個にまたは組み合わせて使用および/または組み込みおよび/または拡張することができることに留意されたい。

【0140】

4.1 関連するニューラルネットワークのパラメータ符号化との比較による優位性

ニューラルネットワークパラメータの依存量子化の利点は、許容される再構成ベクトル

10

20

30

40

50

がN次元信号空間（ここで、Nは、処理されるサンプルのセット、例えば層10a、10bにおけるサンプルまたはニューラルネットワークパラメータ13の数を表す）において密に詰め込まれることである。ニューラルネットワークパラメータセットの再構成ベクトルは、ニューラルネットワークパラメータセットの順序付けられた再構成されたニューラルネットワークパラメータ（または、代替的に、順序付けられた再構成されたサンプル）を指す。依存スカラー量子化の効果を、2つのニューラルネットワークパラメータの最も単純なケースについて図8で説明する。図8は、2つの重みパラメータの単純な場合について、許容される再構成ベクトルの位置の一例を示す図である。図8（a）は独立スカラー量子化の例、図8（b）は依存スカラー量子化の例である。図8aは、独立スカラー量子化の場合の許容再構成ベクトル201（2次元平面上の点を表す）を示している。見て分かるように、第2のニューラルネットワークパラメータ t_1' 13に対する許容値のセットは、第1の再構成されたニューラルネットワークパラメータ t_0' 13に対する選択された値には依存しない。図8（b）は、依存スカラー量子化の例を示す。独立スカラー量子化とは対照的に、第2のニューラルネットワークパラメータ t_1' 13に対して選択可能な再構成値は、第1のニューラルネットワークパラメータ t_0' 13に対して選択された再構成レベルに依存することに注意されたい。図8bの例では、第2のニューラルネットワークパラメータ t_1' 13に対する利用可能な再構成レベルの2つの異なるセット52が存在する（異なる色で図示されている）。第1のニューラルネットワークパラメータ t_0' 13に対する量子化インデックス56が偶数（... , -2, 0, 2, ...）であれば、第1のセット（青い点）の任意の再構成レベル201aを第2のニューラルネットワークパラメータ t_1' 13に対して選択することが可能である。そして、第1のニューラルネットワークパラメータ t_0' に対する量子化インデックス56が奇数（... , -3, -1, 1, 3, ...）であれば、第2のセット（赤色点）の任意の再構成レベル201bを第2のニューラルネットワークパラメータ t_1' 13に対して選択することが可能である。この例では、第1セットと第2セットの再構成レベルは、量子化ステップサイズの半分だけシフトされる（第2セットの任意の再構成レベルは、第1セットの2つの再構成レベルの間に位置する）。

【0141】

ニューラルネットワークパラメータ13の依存スカラー量子化は、N次元単位体積あたりの再構成ベクトル201の所定の平均数に対して、ニューラルネットワークパラメータ13の所定の入力ベクトルと最も近い利用可能な再構成ベクトルとの間の距離の期待値が低減されるという効果を有する。その結果、ニューラルネットワークパラメータの入力ベクトルとベクトル再構成されたニューラルネットワークパラメータとの間の平均歪みは、所与の平均ビット数に対して低減され得る。ベクトル量子化では、この効果をスペースフィリングゲイン（space-filling gain）と呼んでいる。ニューラルネットワークパラメータセット13に対して依存スカラー量子化を用いると、高次元ベクトル量子化の潜在的なスペースフィリングゲインの大部分を利用することができる。そして、ベクトル量子化とは対照的に、再構成プロセス（または復号化プロセス）の実装複雑度は、独立スカラー量子化器を用いた関連するニューラルネットワークパラメータ符号化の実装複雑度に匹敵する。

【0142】

4.2 概要

主な変更点は、前述したように依存量子化である。再構成順序インデックス $k > 0$ の再構成されたニューラルネットワークパラメータ t_k' 13は、関連する量子化インデックス q_{k56} だけでなく、再構成順序において先行するニューラルネットワークパラメータの量子化インデックス q_0, q_1, \dots, q_{k-1} にも依存することになる。なお、依存量子化では、ニューラルネットワークパラメータ13の再構成順序を一意に定めなければならない。量子化インデックス q_{k56} に関連する再構成レベルセットに関する知識もエントロピー符号化において利用される場合、ニューラルネットワーク符号化全体の性能は、典型的に改善され得る。つまり、ニューラルネットワークパラメータに適用される再構成レベルセ

ットに基づいて、コンテキスト（確率モデル）またはコードワード表を切り替えることが典型的には好ましい。

【 0 1 4 3 】

エントロピー符号化は、通常、エントロピー復号化プロセスが与えられると一意に規定される。しかし、関連するニューラルネットワークパラメータ符号化と同様に、オリジナルのニューラルネットワークパラメータを与えられた量子化インデックスの選択には多くの自由度が存在する。

【 0 1 4 4 】

本明細書に記載された実施形態は、層単位ニューラルネットワーク符号化に限定されない。任意の有限のニューラルネットワークパラメータ 1 3 の集合のニューラルネットワークパラメータ符号化にも適用可能である。

【 0 1 4 5 】

特に、本方法は、sec. 3. 1 で説明したようなサブ層にも適用することができる。

【 0 1 4 6 】

4. 3 ニューラルネットワークパラメータの依存量子化

ニューラルネットワークパラメータ 1 3 の依存量子化とは、ニューラルネットワークパラメータ 1 3 のために利用可能な再構成レベルセットが、再構成順序において（例えば層またはサブ層のようなニューラルネットワークパラメータの同じセットの内部において）先行するニューラルネットワークパラメータに対する選ばれた量子化インデックスに依存する概念をいう。

【 0 1 4 7 】

好ましい実施形態では、再構成レベルの複数のセットが予め定義され、符号化順序において先行するニューラルネットワークパラメータの量子化インデックスに基づいて、予め定義されたセットのうちの 1 つが、現在のニューラルネットワークパラメータを再構成するために選択される。言い換えれば、実施形態による装置は、現在のニューラルネットワークパラメータ 1 3 に対して、以前の、例えば先行するニューラルネットワークパラメータに対する量子化インデックス（5 8）に応じて、複数 5 0 の再構成レベルセット 5 2 のうち再構成レベルセット 4 8 を選択するように構成されることができる。

【 0 1 4 8 】

再構成レベルのセットを定義するための好ましい実施形態は、4. 3. 1 節で説明されている。選択された再構成レベルの識別とシグナリングは、4. 3. 2 節に記載されている。4. 3. 3 節では、（再構成順序において先行するニューラルネットワークパラメータの選択された量子化インデックスに基づいた）現在のニューラルネットワークパラメータについての再構成レベルの事前定義されたセットの 1 つを選択するための好ましい実施形態について説明する。

【 0 1 4 9 】

4. 3. 1 再構成レベルセット

好ましい実施形態では、現在のニューラルネットワークパラメータのための許容される再構成レベルセットは、再構成レベルの予め定義されたセット 5 2 の集合（2 つ以上のセット、例えば図 2 および図 3 からのセット 0 およびセット 1）の中から（符号化順序において先行するニューラルネットワークパラメータの量子化インデックスに基づいて）選択される。

【 0 1 5 0 】

好ましい実施形態では、パラメータが量子化ステップサイズ（QP）を決定し、すべての（再構成レベルのすべてのセットにおける）再構成レベルが量子化ステップサイズの整数倍を表す。しかし、再構成レベルの各セットが量子化ステップサイズ（QP）の整数倍のサブセットのみを含むことに留意されたい。このような依存量子化のための構成は、再構成レベルのすべてのセットについて可能なすべての再構成レベルが量子化ステップサイズ（QP）の整数倍を表すものであり、均一再構成量子化器（URQs）の拡張と考えることができる。その基本的な利点は、再構成されたニューラルネットワークパラメ

10

20

30

40

50

ータ 1 3 が非常に低い計算複雑度を有するアルゴリズムによって計算され得ることである（より詳細に後述される）。

【 0 1 5 1 】

再構成レベルセットは完全に不連続であることができるが、1 つ以上の再構成レベルが複数のセットに含まれることも可能である（ただし、セットは他の再構成レベルにおいて依然として異なる）。

【 0 1 5 2 】

好ましい実施形態では、ニューラルネットワークパラメータのための依存スカラー量子化は、再構成レベルのちょうど 2 つの異なるセット、例えばセット 0 とセット 1 とを使用する。そして、特に好ましい実施形態では、ニューラルネットワークパラメータ t_{k13} に対する 2 つのセットのすべての再構成レベルは、このニューラルネットワークパラメータ 1 3 に対する量子化ステップサイズ $k(QP)$ の整数倍を表している。なお、量子化ステップサイズ $k(QP)$ は、両セットの許容再構成値に対するスケーリングファクタを表しているに過ぎない。再構成レベルの同じ 2 つのセットは、すべてのニューラルネットワークパラメータ 1 3 に対して使用される。

【 0 1 5 3 】

図 9 では、2 セットの再構成レベル（セット 0 とセット 1）に対する 3 つの好ましい構成（(a) ~ (c)）が図示されている。図 9 は、単一の量子化ステップサイズ (QP) により完全に決定される再構成レベルの 2 つのセットを有する依存量子化の例を示す。再構成レベルの 2 つの利用可能なセットは、異なる色で強調されている（セット 0 は青、セット 1 は赤）。セット内の再構成レベルを示す量子化インデックスの例は、円の下の数字で示されている。中空と塗りつぶしの円は、再構成レベルセット内の 2 つの異なるサブセットを示し、サブセットは再構成順序において次のニューラルネットワークパラメータの再構成レベルセットを決定するために使用することができる。図には、2 つの再構成レベルセットを持つ 3 つの好ましい構成が示されている：(a) 2 つのセットは不一致でゼロに関して対称である、(b) 両方のセットはゼロに等しい再構成レベルを含むがそれ以外は不一致で、セットはゼロの周りで非対称である、(c) 両方のセットはゼロに等しい再構成レベルを含むがそれ以外は不一致で、どちらのセットもゼロの周りで対称的である。なお、すべての再構成レベルは、量子化ステップサイズの整数倍 (IV) で与えられるグリッド上に存在する。さらに、特定の再構成レベルは両方のセットに含まれる可能性があることに注意する必要がある。

【 0 1 5 4 】

図 9 (a) に描かれた 2 つのセットは不連続である。量子化ステップサイズ (QP) の各整数倍は、どちらかのセットにのみ含まれる。第 1 のセット（セット 0）は量子化ステップサイズのすべての偶数整数倍 (IV) を含むが、第 2 のセット（セット 1）は量子化ステップサイズのすべての奇数整数倍を含む。両セットとも、隣接する 2 つの再構成レベル間の距離は、量子化ステップサイズの 2 倍である。これらの 2 つのセットは、通常、高レート量子化、すなわち、ニューラルネットワークパラメータの分散 $(variance)$ 、平方偏差）が量子化ステップサイズ (QP) よりも著しく大きい設定のものに適している。しかし、ニューラルネットワークのパラメータ符号化では、量子化器は通常、低レート領域で動作する。典型的には、多くのオリジナルのニューラルネットワークパラメータ 1 3 の絶対値は、量子化ステップサイズ (QP) の任意のゼロでない倍数よりもゼロに近い値である。その場合、ゼロが両方の量子化セット（再構成レベルセット）に含まれると、典型的には好ましい。

【 0 1 5 5 】

図 9 (b) に示す 2 つの量子化セットは、いずれもゼロを含む。セット 0 では、ゼロに等しい再構成レベルとゼロより大きい第 1 の再構成レベルとの間の距離は、量子化ステップサイズ (QP) に等しく、隣接する 2 つの再構成レベル間の他のすべての距離は、量子化ステップサイズの 2 倍に等しい。同様に、セット 1 では、ゼロに等しい再構成レベルとゼロより小さい第 1 の再構成レベルとの間の距離は、量子化ステップサイズに等しく、一

方、2つの隣接する再構成レベル間の他のすべての距離は、量子化ステップサイズの2倍に等しくなる。両再構成セットは、ゼロの周りで非対称であることに注意が必要である。これは符号の確率を正確に推定することが困難になるため、非効率になる可能性がある。

【0156】

再構成レベルの2つのセットに関する好ましい構成を図9(c)に示す。第1の量子化セット(図ではセット0と表示)に含まれる再構成レベルは、量子化ステップサイズの偶数整数倍を表す(このセットは、実際には図9(a)のセット0と同じであることに注意)。第2の量子化セット(図ではセット1と表示)は、量子化ステップサイズのすべての奇数整数倍を含み、さらに再構成レベルが0に等しい。なお、どちらの再構成セットもゼロについて対称である。ゼロに等しい再構成レベルは両方の再構成セットに含まれ、そうでなければ再構成セットは不連続である。両再構成セットの組み合わせは、量子化ステップサイズのすべての整数倍を含む。

【0157】

言い換えれば実施形態によれば、例えばニューラルネットワークパラメータ13を符号化/復号化するための装置からなり、複数50の再構成レベルセット52の数は2(例えばセット0、セット1)であり、複数の再構成レベルセットは、ゼロおよび所定の量子化ステップサイズの偶数倍を含む第1の再構成レベルセット(セット0)、およびゼロおよび所定の量子化ステップサイズの奇数倍を含む第2の再構成レベルセット(セット1)含む。

【0158】

さらに、すべての再構成レベルセットのすべての再構成レベルは、所定の量子化ステップサイズ(QP)の整数倍(IV)を表してもよく、実施形態による、例えばニューラルネットワークパラメータ13を復号化するための装置は、各ニューラルネットワークパラメータについて、中間整数値、例えば、それぞれのニューラルネットワークパラメータについて選択された再構成レベルセットおよびそれぞれのニューラルネットワークパラメータ13'についてのエントロピー復号化量子化インデックス58に応じた整数倍(IV)を導出することによって、及び、それぞれのニューラルネットワークパラメータ13について、それぞれのニューラルネットワークパラメータの中間値をそれぞれのニューラルネットワークパラメータ13についての所定の量子化ステップサイズに乘じることによって、ニューラルネットワークパネル13を逆量子化するよう構成されることがある。

【0159】

それぞれ、すべての再構成レベルセットのすべての再構成レベルは、所定の量子化ステップサイズ(QP)の整数倍(IV)を表してもよく、装置、例えば、実施形態によれば、ニューラルネットワークパラメータ13を符号化するために、各ニューラルネットワークパラメータについて、それぞれのニューラルネットワークパラメータの選択された再構成レベルセットとそれぞれのニューラルネットワークパラメータのエントロピー符号化量子化インデックスとに応じた中間整数値を導出することによって、及び、それぞれのニューラルネットワークパラメータについて、それぞれのニューラルネットワークパラメータの中間値をそれぞれのニューラルネットワークパラメータの所定の量子化ステップサイズに乘じることによって、同じものが逆量子化可能になるよう、ニューラルネットワークパネルを量子化するよう構成されることができる。

【0160】

本明細書で規定する実施形態は、図9に示す構成に限定されるものではない。他の任意の2つの異なる再構成レベルセットを使用することができる。複数の再構成レベルが両方のセットに含まれてもよい。あるいは、両方の量子化セットの組み合わせが、量子化ステップサイズのすべての可能な整数倍を含んでいない可能性がある。さらに、ニューラルネットワークパラメータの依存スカラー量子化には、2つ以上の再構成レベルのセットを使用することが可能である。

【0161】

4.3.2 選択された再構成レベルのシグナリング

10

20

30

40

50

符号化器が許容される再構成レベルの中から選択する再構成レベルは、ビットストリーム 14 の内部で示されなければならない。従来の独立スカラー量子化と同様に、これは、重みレベルとも呼ばれる、いわゆる量子化インデックス 56 を用いて実現することができる。量子化インデックス 56（または重みレベル）は、量子化セット 52 の内部（すなわち、再構成レベルのセットの内部）で利用可能な再構成レベルを一意に識別する整数値である。量子化インデックス 56 は、（任意のエントロピー符号化技術を用いる）ビットストリーム 14 の一部として復号化器に送られる。復号化器側では、再構成されたニューラルネットワークパラメータ 13 は、再構成レベルの現在のセット 48（これは、符号化 / 再構成順序において先行する量子化インデックスによって決定される）および現在のニューラルネットワークパラメータ 13' に対する送信された量子化インデックス 56 に基づいて一意に計算することが可能である。

10

【0162】

好ましい実施形態では、再構成レベルのセット（または量子化セット）内部の再構成レベルへの量子化インデックス 56 の割り当ては、次のルールに従う。説明のために、図 9 の再構成レベルには、関連する量子化インデックス 56 が付けられている（量子化インデックスは、再構成レベルを表す円の下の数値によって与えられる）。再構成レベルセットが 0 に等しい再構成レベルを含む場合、0 に等しい量子化インデックスは 0 に等しい再構成レベルに割り当てられる。1 に等しい量子化インデックスは、0 より大きい最小の再構成レベルに割り当てられ、2 に等しい量子化インデックスは、0 より大きい次の再構成レベル（すなわち、0 より大きい 2 番目に小さい再構成レベル）に割り当てられる、等々である。あるいは、言い換えれば、0 より大きい再構成レベルには、その値の昇順において 0 より大きい整数値で（すなわち、1、2、3 など）ラベル付けされる。同様に、量子化インデックス - 1 は、0 より小さい最大の再構成レベルに割り当てられ、量子化インデックス - 2 は、0 より小さい次の（すなわち、2 番目に大きい）再構成レベルに割り当てられる、等である。または、言い換えると、0 より小さい再構成レベルには、値の降順において 0 より小さい整数値（つまり、- 1、- 2、- 3 など）がラベル付けされる。図 9 の例では、（0 に等しい再構成レベルを含まない）図 9（a）のセット 1 を除くすべての量子化セットについて、説明した量子化インデックスの割り当てを例示している。

20

【0163】

0 に等しい再構成レベルを含まない量子化セットに対して、量子化インデックス 56 を再構成レベルに割り当てる方法の 1 つは、以下の通りである。0 より大きい全ての再構成レベルには（値の昇順において）0 より大きい量子化インデックスを割り当て、0 より小さい全ての再構成レベルには（値の降順において）0 より小さい量子化インデックスを割り当てる。したがって、量子化インデックス 56 の割り当ては、基本的に 0 に等しい再構成レベルを含む量子化セットと同じ概念に従うが、0 に等しい量子化インデックスが存在しない点が異なる（図 9（a）の量子化セット 1 のラベルを参照）。量子化インデックス 56 のエントロピー符号化では、その点を考慮する必要がある。例えば、量子化インデックス 56 は、その絶対値（0 からサポートされる最大値までの範囲）を符号化し、0 に等しくない絶対値については、量子化インデックス 56 の符号を付加的に符号化することによって送信されることが多い。0 に等しい量子化インデックス 56 がない場合、エントロピー符号化は、絶対レベルから 1 を引いた値が送信され（対応するシンタックス要素の値は 0 から最大サポート値までの範囲）、符号が常に送信されるように修正され得る。代替案として、量子化インデックス 56 を再構成レベルに割り当てるための割り当てルールが修正され得る。例えば、0 に近い再構成レベルの 1 つは、0 に等しい量子化インデックスでラベル付けされ得る。そして、残りの再構成レベルは、以下のルールによってラベル付けされる。0 に等しい量子化インデックスを持つ再構成レベルより大きい再構成レベルには、0 より大きい量子化インデックスが割り当てられる（量子化インデックスは再構成レベルの値と共に増加する）。そして、0 より小さい量子化インデックスは、量子化インデックスが 0 に等しい再構成レベルより小さい再構成レベルに割り当てられる（量子化インデックスは、再構成レベルの値と共に減少する）。このような割り当ての一つの可能性を

30

40

50

、図 9 (a) の括弧内の数字で示す (括弧内の数字が与えられていない場合は、他の数字が適用される)。

【 0 1 6 4 】

上述のように、好ましい実施形態では、再構成レベルの 2 つの異なるセット (これを量子化セットとも呼ぶ) が使用され、両方のセットの内部の再構成レベルは、量子化ステップサイズ (Q P) の整数倍を表す。それは、量子化ステップサイズが、層ベースで (例えば、ビットストリーム 1 4 の内部に層量子化パラメータを送信することによって)、またはニューラルネットワークパラメータ 1 3 の別の有限セット (例えば、ブロック) で (例えば、ビットストリーム 1 4 の内部にブロック量子化パラメータを送信することによって) 変更される、場合を含んでいる。

10

【 0 1 6 5 】

量子化ステップサイズ (Q P) の整数倍を表す再構成レベルの使用は、復号化器側でのニューラルネットワークパラメータ 1 3 の再構成のための計算量の少ない複雑なアルゴリズムを可能にする。これを図 9 (c) の好ましい例に基づいて以下に説明する (他の構成、特に図 9 (a) および図 9 (b) に示す設定についても同様の単純なアルゴリズムが存在する)。図 9 (c) に示す構成では、第 1 の量子化セットは量子化ステップサイズ (Q P) のすべての偶数整数倍を含み、第 2 の量子化セットは量子化ステップサイズのすべての奇数整数倍に 0 に等しい再構成レベルを加えたもの (これは両方の量子化セットに含まれる) が含まれる。ニューラルネットワークパラメータの再構成プロセスは、図 1 0 の疑似コードで規定されるアルゴリズムと同様に実装され得る。図 1 0 は、ニューラルネットワークパラメータ 1 3 の再構成プロセスに関する好ましい例を示す疑似コードである。k は現在のニューラルネットワークパラメータ 1 3 ' の再構成順序を指定するインデックスを表し、現在のニューラルネットワークパラメータの量子化インデックス 5 6 は `level [k] 2 1 0` で表し、現在のニューラルネットワークパラメータ 1 3 ' に適用する量子化ステップサイズ $_k (Q P)$ は `quant__step__size [k]` で表し、`tre c [k] 2 2 0` は再構成したニューラルネットワークのパラメータ t_k ' の値を表している。変数 `set Id [k] 2 4 0` は、現在のニューラルネットワークパラメータ 1 3 ' に適用される再構成レベルのセットを指定する。それは、再構成順序において先行するニューラルネットワークパラメータに基づいて決定される ; `set Id [k]` の可能な値は 0 及び 1 である。変数 n は、量子化ステップサイズ (Q P) の整数係数、例えば中間値 I V を規定する ; それは、選択された再構成レベルのセット (すなわち `set Id [k]` の値) 及び送信された量子化インデックスレベル [k] により与えられる。

20

30

【 0 1 6 6 】

40

50

図10の疑似コードにおいて、 $level[k]$ は、ニューラルネットワークパラメータ t_{k13} について送信される量子化インデックス56を示し、(0または1に等しい) $setId[k]$ は再構成レベルの現在のセットの識別子を指定する(それは以下でより詳細に説明するように再構成順序において先行する量子化インデックス56に基づいて決定される)。変数 n は、量子化インデックス $level[k]$ とセット識別子 $setId[k]$ によって与えられる量子化ステップサイズ(QP)の整数倍を表す。ニューラルネットワークパラメータ13が、量子化ステップサイズ Δ_k (QP)の偶数整数倍を含む再構成レベル($setId[k] == 0$)の第1のセットを用いて符号化される場合、変数 n は、送信された量子化インデックス56の2倍となる。この場合は、図9(c)の第1の量子化セット $Set0$ の再構成レベルで表すことができ、 $Set0$ は量子化ステップサイズ(QP)の偶数整数倍を全て含んでいる。ニューラルネットワークパラメータ13が再構成レベル($setId[k] == 1$)の第2のセットを用いて符号化される場合、(a) $level[k]$ が0に等しければ、 n も0に等しくなり、(b) $level[k]$ が0より大きい場合、 n は量子化インデックス $level[k]$ の2倍から1を引いたものに等しくなり、(c) $level[k]$ が0未満なら n は量子化インデックスレベル $[k]$ の2倍に1を加えたものに等しくなる、の三つの場合が存在する。これは符号関数を用いて指定することができる。

10

20

$$\text{sign}(x) = \begin{cases} 1 & : x > 0 \\ 0 & : x = 0 \\ -1 & : x < 0 \end{cases}$$

【0167】

そして、第2の量子化セットを使用する場合、変数 n は、量子化インデックス $level[k]$ の2倍から量子化インデックスの符号関数 $\text{sign}(level[k])$ を引いた値に等しくなる。この場合は、図9(c)の第2の量子化セット $Set1$ の再構成レベルで表すことができ、 $Set1$ は量子化ステップサイズ(QP)の奇数整数倍を全て含む。

【0168】

(量子化ステップサイズの整数倍を指定する)変数 n が決定されると、 n に量子化ステップサイズ $_k$ を乗じることによって、再構成されたニューラルネットワークパラメータ $t_{k'}$ が求められる。

30

【0169】

言い換えれば、複数50の再構成レベルセット52のうち再構成レベルセット52の数は2であってもよく、本発明の実施形態による、例えばニューラルネットワークパラメータ13を復号化および/または符号化するための装置は、以下のようにしてそれぞれのニューラルネットワークパラメータの中間値を導出するように構成され得る。

それぞれのニューラルネットワークパラメータについての選択された再構成レベルセットが第1のセットである場合、それぞれのニューラルネットワークパラメータに対する量子化インデックスを2倍して、それぞれのニューラルネットワークパラメータに対する中間値を求め；及び、

40

それぞれのニューラルネットワークパラメータについての選択された再構成レベルセットが2番目のセットであり、それぞれのニューラルネットワークパラメータの量子化インデックスがゼロに等しい場合、それぞれのサンプルの中間値をゼロに設定し；及び、

それぞれのニューラルネットワークパラメータについての選択された再構成レベルセットが第2のセットであり、それぞれのニューラルネットワークパラメータに対する量子化インデックスがゼロより大きい場合、それぞれのニューラルネットワークパラメータに対する量子化インデックスを2倍し、その乗算結果から1を引いてそれぞれのニューラルネットワークパラメータに対する中間値を得て、及び、

50

現在のニューラルネットワークパラメータについての選択された再構成レベルセットが第2のセットであり、それぞれのニューラルネットワークパラメータに対する量子化インデックスがゼロより小さい場合、それぞれのニューラルネットワークパラメータの量子化インデックスを2倍し、その乗算結果に1を加えて、それぞれのニューラルネットワークパラメータの中間値を得る。

【0170】

4.3.3 ニューラルネットワークパラメータの依存的再構成

4.3.1節および4.3.2節で説明した再構成レベルのセットの選択に加えて、ニューラルネットワークパラメータ符号化における依存スカラー量子化のもう一つの重要な設計側面は、定義された量子化セット（再構成レベルのセット）間の切り替えに使用するアルゴリズムである。使用されるアルゴリズムによって、ニューラルネットワークパラメータ13のN次元空間（したがって、再構成されたサンプルのN次元空間も同様）で達成できる「パッキング密度」が決まる。パッキング密度が高ければ高いほど、最終的に符号化効率が向上する。

【0171】

次のニューラルネットワークパラメータのための再構成レベルのセットを決定する好ましい方法は、図11に示されるように、量子化セットの分割に基づくものである。図11は、本発明の実施形態による再構成レベルのセットを2つのサブセットに分割するための一例を示す図である。示された2つの量子化セットは、図9(c)の好ましい例の量子化セットである。量子化セット0の2つのサブセットは、「A」及び「B」を用いてラベル付けされ、量子化セット1の2つのサブセットは、「C」及び「D」を用いてラベル付けされている。なお、図11に示す量子化セットは、図9(c)の量子化セットと同じ量子化セットである。2つ（またはそれ以上）の量子化セットのそれぞれは、2つのサブセットに分割される。図11の好ましい例では、第1の量子化セット（セット0と表示）は、2つのサブセット（AおよびBと表示される）に分割され、第2の量子化セット（セット1と表示される）も、2つのサブセット（CおよびDと表示される）に分割される。唯一の可能性ではないにしても、各量子化セットに対する分割は、好ましくは、直接隣接する再構成レベル（したがって、隣接する量子化インデックス）が異なるサブセットと関連付けられるように行われる。好ましい実施形態では、各量子化セットは、2つのサブセットに分割される。図9において、量子化セットのサブセットへの分割は、中空および充填された円によって示されている。

【0172】

図11および図9(c)に示される特に好ましい実施形態では、以下の分割ルールが適用される。

- サブセットAは、量子化セット0のすべての偶数量子化インデックスから構成される。
- サブセットBは、量子化セット0のすべての奇数量子化インデックスから構成される。
- サブセットCは、量子化セット1のすべての偶数量子化インデックスから構成される。
- サブセットDは、量子化セット1のすべての奇数量子化インデックスで構成される。

【0173】

使用されるサブセットは、典型的には、ビットストリーム14の内部で明示的に示されないことに留意されたい。その代わりに、使用される量子化セット（例えば、セット0またはセット1）および実際に送信された量子化インデックス56に基づいて導出することができる。図11に示す好ましい分割の場合、サブセットは、送信された量子化インデックスレベル及び1のビット単位「and」演算によって導出することができる。サブセットAは、 $(level \& 1)$ が0に等しいセット0のすべての量子化インデックスからなり、サブセットBは、 $(level \& 1)$ が1に等しいセット0のすべての量子化インデックスからなり、サブセットCは、 $(level \& 1)$ が0に等しいセット1のすべての量子化インデックスからなり、サブセットDは、 $(level \& 1)$ が1に等しいセット1のすべての量子化インデックスからなる。

【0174】

10

20

30

40

50

好ましい実施形態では、現在のニューラルネットワークパラメータ13'を再構成するために使用される量子化セット（許容再構成レベルのセット）は、最後の2つ以上の量子化インデックス56に関連付けられるサブセットに基づいて決定される。最後の2つのサブセット（これは最後の2つの量子化インデックスによって与えられる）が使用される例を、表1に示す。この表によって指定される量子化セットの決定は、好ましい実施形態を表す。他の実施形態では、現在のニューラルネットワークパラメータ13'に対する量子化セットは、最後の3つ以上の量子化インデックス56と関連するサブセットによって決定される。層の第1のニューラルネットワークパラメータ（またはニューラルネットワークパラメータのサブセット）については、（先行するニューラルネットワークパラメータがないため）先行するニューラルネットワークパラメータのサブセットに関するデータはない。好ましい実施形態では、このような場合に予め定義された値を使用する。特に好ましい実施形態では、利用可能でない全てのニューラルネットワークパラメータについて、サブセットAを推論する。つまり、第1のニューラルネットワークパラメータを再構成する場合、先行する2つのサブセットを「AA」（又は、先行する3つのニューラルネットワークパラメータを考慮する場合は「AAA」）と推測し、したがって、表1に従って、量子化セット0を使用する。また、2番目のニューラルネットワークパラメータについては、直前の量子化インデックスのサブセットをその値によって決定するが（1番目のニューラルネットワークパラメータについてはセット0が使用されるので、サブセットはAまたはBのいずれか）、2番目の最後の量子化インデックス（これは存在しない）についてのサブセットはAに等しいと推測する。もちろん、存在しない量子化インデックスについてのデフォルト値の推測には他の任意のルールを使用することが可能である。また、存在しない量子化インデックスのデフォルトサブセットを導出するために、他のシンタックス要素を使用することも可能である。さらなる代替案として、初期化のために、ニューラルネットワークパラメータ13の先行するセットの最後の量子化インデックス56を使用することも可能である。

【0175】

表1：本発明の実施形態による2つの最後の量子化インデックスに関連するサブセットに基づいて、次のニューラルネットワークパラメータに使用される量子化セット（利用可能な再構成レベルのセット）を決定するための例である。サブセットは、左の表の欄に示されており、それらは、（2つの最後の量子化インデックスについての）使用される量子化セットおよび（量子化インデックスのパリティによって決定されてよい）いわゆるパスによって一義的に決定される。量子化セットと、括弧内にはサブセットのためのパスが左から2番目の列にリストアップされている。3列目には、関連する量子化セットを指定する。最後の列には、いわゆる状態変数の値が示されており、量子化セットを決定するためのプロセスを簡略化するために使用することができる。

【0176】

10

20

30

40

50

【表 1】

最後の 2 つの量子化 インデックスのサブ セット	最後の 2 つの量子化イン デックスについての量子 化セットとパス（括弧 内）	現在のニューラル ネットワークパラ メータに対する量 子化セット	状態変数
A A	0(0), 0(0)	0	0
A B	0(0), 0(1)	0	0
A C	0(0), 1(0)	1	1
A D	0(0), 1(1)	1	1
B A	0(1), 0(0)	1	1
B B	0(1), 0(1)	1	1
B C	0(1), 1(0)	0	0
B D	0(1), 1(1)	0	0
C A	1(0), 0(0)	0	2
C B	1(0), 0(1)	0	2
C C	1(0), 1(0)	1	3
C D	1(0), 1(1)	1	3
D A	1(1), 0(0)	1	3
D B	1(1), 0(1)	1	3
D C	1(1), 1(0)	0	2
D D	1(1), 1(1)	0	2

【 0 1 7 7 】

なお、量子化インデックス 5 6 のサブセット（A、B、C 又は D）は、使用する量子化セット（セット 0 またはセット 1）と量子化セット内部の使用するサブセット（例えば、セット 0 なら A または B、セット 1 なら C または D）により決定されることに注意する必要がある。量子化セットの内部で選択されたサブセットは、（後述するように、依存量子化プロセスをトレリス構造で表現する場合に、パスを指定するため）パスとも呼ばれる。我々の慣例では、パスは 0 か 1 のどちらかに等しい。すると、サブセット A はセット 0 におけるパス 0 に対応し、サブセット B はセット 0 におけるパス 1 に対応し、サブセット C はセット 1 におけるパス 0 に対応し、サブセット D はセット 1 におけるパス 1 に対応する。したがって、次のニューラルネットワークパラメータの量子化セットも、最後の 2 つ（またはそれ以上）の量子化インデックスに関連付けられた量子化セット（セット 0 またはセット 1）とパス（パス 0 またはパス 1）によって一意に決定される。表 1 では、関連する量子化セットとパスが 2 列目に指定されている。

【 0 1 7 8 】

10

20

30

40

50

なお、パスは、例えば2値関数による簡単な演算で決定できることが多い。例えば、図11のような構成の場合、パスは次式で与えられる。

$$\text{path} = (\text{level}[k] \& 1),$$

ここで、 $\text{level}[k]$ は量子化インデックス（重みレベル）56を表し、演算子 $\&$ はビット単位「and」（2補整数演算）を指定する。

10

【0179】

言い換えれば、複数50の再構成レベルセット52のうちの再構成レベルセット52の数は、例えばセット0とセット1との2つであってもよく、本発明の実施形態による、例えばニューラルネットワークパラメータ13を復号化するための装置は、それぞれのニューラルネットワークパラメータについての再構成レベルの選択されたセットとそれぞれのニューラルネットワークパラメータについての量子化インデックスの2値関数とに基づいて、それぞれのニューラルネットワークパラメータについてサブセットインデックスを導き、結果として4つの可能な値、例えばサブセットインデックスについてA、B、C、またはDを得ることができ、現在のニューラルネットワークパラメータ13'について、以前に復号化されたニューラルネットワークパラメータについてのサブセットインデックスに応じて、複数50の再構成レベルセット52のうち再構成レベルのセット48を選択するように構成することができる。

20

【0180】

本発明によるさらなる実施形態は、現在のニューラルネットワークパラメータ13'について、例えば表1の1列目に示すように、直前に復号化された多数のニューラルネットワークパラメータのサブセットインデックスに対応する選択ルールを使用して、複数50の再構成レベルセット52のうちの再構成レベルセット48を選択54し、ニューラルネットワークパラメータのすべて、または一部について選択ルールを使用するように構成される装置である。

30

【0181】

さらなる実施形態によれば、選択ルールが対応する、直前に復号化されたニューラルネットワークパラメータの数は2つであり、例えば表1に示すように、2つの最後の量子化インデックスのサブセットである。

【0182】

追加の実施形態によれば、各ニューラルネットワークパラメータのサブセットインデックスは、それぞれのニューラルネットワークパラメータについての再構成レベルの選択されたセットと、それぞれのニューラルネットワークパラメータについての量子化インデックスの、例えば $\text{path} = (\text{level}[k] \& 1)$ を用いたパリティに基づいて導出される。

40

【0183】

それぞれ、実施形態によるニューラルネットワークパラメータ13を符号化するための装置については、複数50の再構成レベルセット52の再構成レベルセット52の数は、例えばセット0とセット1との2つであってもよく、装置は、それぞれのニューラルネットワークパラメータのための再構成レベルの選択されたセットとそれぞれのニューラルネットワークパラメータについての量子化インデックスの2値関数とに基づいて、それぞれのニューラルネットワークパラメータのためのサブセットインデックスを導出して、サブセットインデックスについて4つの可能な値、例えば、A、B、C、Dを生じさせ、現在のニューラルネットワークパラメータ13'に対して、以前に符号化されたニューラルネットワークパラメータについてのサブセットインデックスに応じて複数50の再構成レベルセ

50

ット 5 2 のうちの再構成レベルセット 4 8 を選択 5 4 するように構成することができる。

【 0 1 8 4 】

本発明によるさらなる実施形態は、現在のニューラルネットワークパラメータ 1 3 ' について、例えば表 1 の第 1 列に示すように、直前に符号化された多数のニューラルネットワークパラメータのサブセットインデックスに対応する選択ルールを用いて、複数 5 0 の再構成レベルセット 5 2 のうち再構成レベルのセット 4 8 を選択 5 4 し、ニューラルネットワークパラメータのすべて、または一部について選択ルールを用いるよう構成される装置を備える。

【 0 1 8 5 】

さらなる実施形態によれば、選択ルールが対応する、直前に符号化されたニューラルネットワークパラメータの数は 2 であり、例えば表 1 に示すように、2 つの最後の量子化インデックスのサブセットである。

【 0 1 8 6 】

追加の実施形態によれば、各ニューラルネットワークパラメータについてのサブセットインデックスは、それぞれのニューラルネットワークパラメータについての再構成レベルの選択されたセットと、それぞれのニューラルネットワークパラメータの量子化インデックスのパリティ、例えば $path = (level[k] \& 1)$ を使用して、導出される。

【 0 1 8 7 】

量子化セット 5 2 (セット 0 とセット 1) の間の遷移は、状態変数によってエレガントに表現することも可能である。そのような状態変数の例を表 1 の最後の列に示す。この例では、状態変数は 4 つの可能な値 (0 、 1 、 2 、 3) を有する。一方、状態変数は、現在のニューラルネットワークパラメータ 1 3 ' に使用される量子化セットを指定する。表 1 の好ましい例では、状態変数が 0 または 2 に等しい場合にのみ、量子化セット 0 が使用され、状態変数が 1 または 3 に等しい場合にのみ、量子化セット 1 が使用される。一方、状態変数は、量子化セット間の可能な遷移も指定する。状態変数を用いることで、表 1 のルールをより小さな状態遷移表で記述することができる。一例として、表 2 は、表 1 で与えられたルールの状態遷移表を規定したものである。これは、好ましい実施形態を表している。現在の状態を指定して、現在のニューラルネットワークのパラメータ (2 列目) に量子化セットを規定した。それはさらに、選択された量子化インデックス 5 6 に関連するパスに基づいて状態遷移を指定する (パスは、量子化セットが与えられた場合、使用されるサブセット A 、 B 、 C 、または D を指定する) 。状態変数の概念を使用することによって、実際に選択されたサブセットを追跡する必要がないことに留意されたい。ある層のニューラルネットワークのパラメータを再構築する際には、状態変数を更新し、使用される量子化インデックスのパスを決定すれば十分である。

【 0 1 8 8 】

表 2 : 本発明の実施形態による、4 つの状態を有する構成の状態遷移表の好ましい例。

【表 2】

現在の状態	電流係数に対する量子化セット	次の状態	
		パス 0	パス 1
0	0	0	1
1	1	2	3
2	0	1	0
3	1	3	2

【 0 1 8 9 】

つまり、実施形態によるニューラルネットワークパラメータを復号化するための装置であって、当該装置は、現在のニューラルネットワークパラメータ 1 3 ' に関連する状態に応

じて、複数50の再構成レベルセット52のうちの量子化レベルのセット48を決定することによって、及び、直前のニューラルネットワークパラメータのデータストリームから復号化された量子化インデックス58に応じて後続のニューラルネットワークパラメータの状態を更新することによって、現在のニューラルネットワークパラメータ13'について、状態遷移プロセスにより複数50の再構成レベルセット52のうちの再構成レベルセット48を選択54のように構成されることができる。

【0190】

それぞれ、実施形態によるニューラルネットワークパラメータ13を符号化するための装置について、前記装置は、現在のニューラルネットワークパラメータ13'に関連する状態に応じて、複数50の再構成レベルセット52のうちの再構成レベルのセット48を決定することによって、及び、直前のニューラルネットワークパラメータのデータストリームに符号化された量子化インデックス58に応じて、後続のニューラルネットワークパラメータについての状態を更新することによって、現在のニューラルネットワークパラメータ13'について、状態遷移プロセスによって複数50の再構成レベルセット52のうちの再構成レベルのセット48を選択54するように構成されてもよい。

【0191】

本発明の好ましい実施形態では、パスは、量子化インデックスのパリティによって与えられる。level[k]を現在の量子化インデックスとすると、次式に従って決定することができる。

$$\text{path} = (\text{level}[k] \& 1),$$

ここで、演算子&は、2補整数演算におけるビット単位「and」を表す。

【0192】

言い換えれば、実施形態による、例えばニューラルネットワークパラメータを復号化するための装置は、直前のニューラルネットワークパラメータについてのデータストリームから復号化された量子化インデックス58の2値関数を用いて、例えば表2に従って、後続のニューラルネットワークパラメータの状態を更新するよう構成されることができる。

【0193】

さらに、実施形態による装置は、直前のニューラルネットワークパラメータについてのデータストリーム14から復号化された量子化インデックス58のパリティ、例えばpath=(level[k]&1)を使用して、後続のニューラルネットワークパラメータのための状態を更新するよう構成されることができる。

【0194】

それぞれ、実施形態によるニューラルネットワークパラメータ13を符号化するための装置について、当該装置は、直前のニューラルネットワークパラメータについてのデータストリームに符号化された量子化インデックス58の2値関数を用いて、後続のニューラルネットワークパラメータについての状態を更新するよう構成されることができる。

【0195】

さらに、実施形態による例えばニューラルネットワークパラメータ13を符号化するための装置は、直前のニューラルネットワークパラメータについてのデータストリームに符号化された量子化インデックス58のパリティを用いて、後続のニューラルネットワークパラメータについての状態を、例えば表2に従って更新するよう構成されることができる。

【0196】

好ましい実施形態では、4つの可能な値を有する状態変数が使用される。他の実施形態では、異なる数の可能な値を有する状態変数が使用される。特に興味深いのは、状態変数

の可能な値の数が2の整数乗、すなわち、4、8、16、32、64などを表す状態変数である。(表1および表2で与えられるように)好ましい構成において、4つの可能な値を有する状態変数は、現在の量子化セットが2つの最後の量子化インデックスのサブセットによって決定されるアプローチと同等であることに注意されたい。8つの可能な値を持つ状態変数は、現在の量子化セットが3つの最後の量子化インデックスのサブセットによって決定される同様のアプローチに対応する。16個の可能な値を持つ状態変数は、現在の量子化セットが最後の4つの量子化インデックスのサブセットによって決定されるアプローチに対応する。一般に、2の整数乗に等しい可能な値の数を有する状態変数を使用することが好ましいにもかかわらず、実施形態は、この設定に限定されない。

【0197】

特に好ましい実施形態では、8つの可能な値(0、1、2、3、4、5、6、7)を有する状態変数が使用される。好ましい例の表3では、状態変数が0、2、4または6に等しい場合にのみ、量子化セット0が使用され、状態変数が1、3、5または7に等しい場合にのみ、量子化セット1が使用される。

【0198】

実施形態による、8つの状態を有する構成の状態遷移表の好ましい例である。

【表3】

現在の状態	電流係数について の量子化セット	次の状態	
		パス0	パス1
0	0	0	2
1	1	7	5
2	0	1	3
3	1	6	4
4	0	2	0
5	1	5	7
6	0	3	1
7	1	4	6

【0199】

すなわち、本発明の実施形態によれば、状態遷移プロセスにおいて、4つ又は8つの可能な状態の間で遷移するように構成される。

【0200】

さらに、実施形態による、ニューラルネットワークパラメータ13を復号化/符号化するための装置は、状態遷移プロセスにおいて、偶数の可能な状態の間で遷移し、複数50の再構成レベルセット52の数が2であるように構成されることができ、ここで、現在のニューラルネットワークパラメータ13'について、現在のニューラルネットワークパラメータ13'に関連する状態に応じて、量子化セット52のうちの量子化レベルのセット48を決定することで、状態が偶数の可能な状態の前半に属する場合、複数50の再構成レベルセット52のうちの第1の再構成レベルセットが決定され、状態が偶数の可能な状態の後半に属する場合、複数50の再構成レベルセット52のうちの第2の再構成レベルセットが決定される。

【0201】

さらなる実施形態による、例えばニューラルネットワークパラメータ13を復号化するための装置は、直前のニューラルネットワークパラメータのデータストリームから復号化された量子化インデックス58の状態およびパリティの組み合わせを、後続のニューラルネットワークパラメータに関連する別の状態上にマッピングする遷移表によって、状態の

更新を実行するように構成されてもよい。

【0202】

それに応じて、実施形態によるニューラルネットワークパラメータ13を符号化するための装置は、直前のニューラルネットワークパラメータについてのデータストリームに符号化された量子化インデックス58の状態とパリティとの組み合わせを後続のニューラルネットワークパラメータに関連する別の状態にマッピングする遷移表によって状態の更新を行うように構成されることができる。

【0203】

状態遷移の概念を用いると、現在の状態、ひいては現在の量子化セットは、（再構成順序において）以前の状態と以前の量子化インデックス56によって一意に決定される。しかし、有限集合（例えば層）の第1のニューラルネットワークパラメータ13については、以前の状態および以前の量子化インデックスが存在しない。したがって、ある層の第1のニューラルネットワークパラメータに対する状態が一意に定義されることが要求される。さまざまな可能性がある。好ましい選択肢は以下の通りである。

- ・層のための第1の状態は、固定された事前定義された値に常に等しく設定される。好ましい実施形態では、第1の状態は0に等しく設定される。

- ・第1の状態の値は、ビットストリーム14の一部として明示的に送信される。これは、可能な状態値のサブセットのみが対応するシンタックス要素によって示され得るアプローチを含む。

- ・第1の状態の値は、その層の他のシンタックス要素に基づいて導出される。つまり、対応するシンタックス要素（またはシンタックス要素）が復号化器への他の態様のシグナリングに使用されても、それらは、依存スカラー量子化のための第1の状態を導出するために追加的に使用されることを意味する。

【0204】

依存スカラー量子化の状態遷移の概念は、復号化器におけるニューラルネットワークパラメータ13の再構成のための複雑さが少ない実装を可能にする。単層のニューラルネットワークパラメータの再構成プロセスに関する好ましい例を、C言語形式の擬似コードを用いて図12に示す。図12は、本発明の実施形態による層のニューラルネットワークパラメータ13の再構成プロセスのための好ましい例を示す擬似コードの一例である。なお、量子化インデックスの導出と、例えば量子化ステップサイズを用いた、あるいは代替的にコードブックを用いた再構成値の導出は、次々と別々のループで行われてもよい。すなわち、言い換えれば、「n」の導出と状態更新を第1のループで行い、「trec」の導出をもう1つの別の第2のループで行うようにすることができる。アレイレベル（array level）210は、その層の送信されたニューラルネットワークパラメータレベル（量子化インデックス56）を表し、アレイtrec（array trec）220は、対応する再構成されたニューラルネットワークパラメータ13を表す。現在のニューラルネットワークパラメータ13'に適用される量子化ステップサイズk（QP）は、quant__step__size[k]で示される。2d表sttab（2d table sttab）230は、例えば表1、表2及び/又は表3のいずれかに従って状態遷移表を指定し、表setId（table setId）240は、状態250に関連付けられる量子化セットを指定する。

【0205】

図12の擬似コードにおいて、インデックスkは、ニューラルネットワークパラメータの再構成順序を指定する。最後のインデックスlayerSizeは、最後に再構成されたニューラルネットワークパラメータの再構成インデックスを指定する。変数layerSizeは、層内のニューラルネットワークパラメータの数と等しく設定されてもよい。各単一のニューラルネットワークパラメータの再構成プロセスは、図10の例と同じである。図10の例と同様に、量子化インデックスはlevel[k]210で表され、関連する再構成されたニューラルネットワークパラメータはtrec[k]220で表される。また、状態変数は状態（state）210で表される。図12の例では、状態は、層

10

20

30

40

50

の始めに 0 に等しく設定されることに留意されたい。しかし、上述したように、他の初期化（例えば、いくつかのシンタックス要素の値に基づいて）も可能である。1 d 表 s e t I d [] (1 d t a b l e s e t I d []) 2 4 0 は、状態変数の異なる値に関連付けられる量子化セットを指定し、2 d 表 s t t a b [] [] (2 d t a b l e s t t a b [] []) 2 3 0 は、現在の状態（第 1 の引数）およびパス（第 2 の引数）を与えられた状態遷移を指定する。この例では、パスは（ビット単位と演算子 & を使用する）量子化インデックスのパリティで与えられるが、他の概念も可能である。C 言語スタイルのシンタックスにおいて、表の例を図 1 3 と図 1 4 に示す（これらの表は表 2、表 3 と同一であり、言い換えれば、表 2、表 3 の表現を提供することができる）。

【 0 2 0 6 】

図 1 3 は、本発明の実施形態による状態遷移表 s t t a b 2 3 0 と、状態 2 5 0 に関連する量子化セットを指定する表 s e t I d 2 4 0 についての好ましい例を示す図である。C スタイルのシンタックスで与えられた表は、表 2 に規定された表を表す。

【 0 2 0 7 】

図 1 4 は、本発明の実施形態による状態遷移表 s t t a b 2 3 0 と、状態 2 5 0 に関連する量子化セットを指定する表 s e t I d 2 4 0 に関する好ましい例を示す図である。C スタイルのシンタックスで与えられた表は、表 3 に規定された表を表している。

【 0 2 0 8 】

別の実施形態では、0 に等しいすべての量子化インデックス 5 6 は、状態遷移および依存再構成プロセスから除外される。量子化インデックス 5 6 が 0 に等しいか等しくないかの情報は、ニューラルネットワークパラメータ 1 3 を 0 と 0 でないニューラルネットワークパラメータに分割するために使用されるだけである。依存スカラー量子化のための再構成プロセスは、ゼロでない量子化インデックス 5 6 の順序付けされたセットにのみ適用される。0 に等しい量子化インデックスに関連するすべてのニューラルネットワークパラメータは、単に 0 に等しく設定される。対応する擬似コードを図 1 5 に示す。図 1 5 は、本発明の実施形態による、0 に等しい量子化インデックスが状態遷移および依存スカラー量子化から除外される、ニューラルネットワークパラメータレベルの代替再構成プロセスを示す擬似コードである。

【 0 2 0 9 】

また、依存量子化における状態遷移は、図 1 6 に示されるように、トレリス構造を用いて表現することも可能である。図 1 6 は、本発明の実施形態によるトレリス構造としての依存スカラー量子化における状態遷移の例を示す図である。横軸は、再構成順序において異なるニューラルネットワークパラメータ 1 3 を表している。縦軸は、依存量子化および再構成プロセスにおける異なる可能な状態 2 5 0 を表す。示された接続は、異なるニューラルネットワークパラメータに対する状態間の利用可能なパスを指定する。この図に示されるトレリスは、表 2 に指定される状態遷移に対応する。各状態 2 5 0 について、現在のニューラルネットワークパラメータ 1 3 ' に対する状態を、再構成順序において次のニューラルネットワークパラメータ 1 3 に対する 2 つの可能な状態と接続する 2 つのパスが存在する。パスは、パス 0 とパス 1 とでラベル付けされ、この番号は、上で紹介したパス変数に対応する（好ましい実施形態では、そのパス変数は、量子化インデックスのパリティに等しい）。各パスは、量子化インデックスに対するサブセット（A、B、C、または D）を一意的に指定することに留意されたい。図 1 6 では、サブセットは括弧で指定されている。初期状態（例えば状態 0）が与えられると、トレリスを通るパスは、送信された量子化インデックス 5 6 によって一意的に指定される。

【 0 2 1 0 】

図 1 6 の例では、状態（0, 1, 2, 3）は次のような性質を持っている。

- ・ 状態 0：前の量子化インデックス l e v e l [k - 1] はセット 0 の再構成レベルを指定し、現在の量子化インデックス l e v e l [k] はセット 0 の再構成レベルを指定する。
- ・ 状態 1：前の量子化インデックス l e v e l [k - 1] はセット 0 の再構成レベルを指定し、現在の量子化インデックス l e v e l [k] はセット 1 の再構成レベルを指定する。

10

20

30

40

50

- ・状態 2：前の量子化インデックス $level[k-1]$ はセット 1 の再構成レベルを指定し、現在の量子化インデックス $level[k]$ はセット 0 の再構成レベルを指定する。
 - ・状態 3：前の量子化インデックス $level[k-1]$ はセット 1 の再構成レベルを指定し、現在の量子化インデックス $level[k]$ はセット 1 の再構成レベルを指定する。
- 【0211】

トレリスは、いわゆる基本トレリスセルを連結したものである。このような基本トレリスセルについての一例を図 17 に示す。図 17 は、本発明の実施形態に係る基本トレリスセルの一例を示す図である。本発明は、4 つの状態 250 を有するトレリスに限定されないことに留意されたい。他の実施形態では、トレリスは、より多くの状態 250 を有することができる。特に、2 の整数乗を表す任意の数の状態が好適である。特に好ましい実施形態では、状態 250 の数は、例えば表 3 と同様に 8 に等しい。トレリスが 2 以上の状態 250 を有する場合でも、現在のニューラルネットワークパラメータ 13' のための各ノードは、典型的には、以前のニューラルネットワークパラメータ 13 のための 2 つの状態および次のニューラルネットワークパラメータ 13 の 2 つの状態と接続される。しかしながら、ノードが、以前のニューラルネットワークパラメータの 2 つ以上の状態、または次のニューラルネットワークパラメータの 2 つ以上の状態と接続されることも可能である。完全に接続されたトレリス（各状態 250 は、以前のニューラルネットワークパラメータ 13 のすべての状態 250 および次のニューラルネットワークパラメータ 13 のすべての状態 250 と接続されている）は、独立したスカラー量子化に対応するであろうことに留意されたい。

【0212】

好ましい実施形態では、（この決定を復号化器に送信するために、何らかのサイド情報レートを必要とするため）初期状態を自由に選択することはできない。その代わりに、初期状態は、予め定義された値に設定されるか、またはその値が他のシンタックス要素に基づいて導出されるかのいずれかである。この場合、第 1 のニューラルネットワークパラメータに対して、すべてのパスと状態 250 が利用できるわけではない。4 状態のトレリスの例として、図 18 は、初期状態が 0 に等しい場合のトレリス構造を示す。図 18 は、本発明の実施形態による 8 個のニューラルネットワークパラメータの依存スカラー量子化のためのトレリス例を示す図である。第 1 の状態（左側）は、初期状態を表し、この例では 0 に等しく設定されている。

【0213】

4.4 エントロピー符号化

依存量子化によって得られた量子化インデックスは、エントロピー符号化方式によって符号化される。これには、任意のエントロピー符号化法が適用可能である。本発明の好ましい実施形態では、コンテキスト適応的 2 値算術符号化 (CABAC) を用いた、第 2.2 節（符号化方法については第 2.2.1 節、復号化方法については第 2.2.2 節参照）によるエントロピー符号化方法が適用される。このために、例えば図 5 に示すように、量子化インデックスを絶対値として送信するために、まず非 2 値が一連の 2 値決定（いわゆるピン）上にマッピングされる（2 値化）。

【0214】

ここで説明したどの概念も、3 節における方法と関連する概念（特にコンテキストモデリング (context modeling) に関するもの）と組み合わせることができることに留意する必要がある。

【0215】

4.4.1 依存スカラー量子化のコンテキストモデリング

依存スカラー量子化の主な態様は、ニューラルネットワークパラメータ 13 のための許容される再構成レベルの異なるセット（量子化セットとも呼ばれる）が存在することである。現在のニューラルネットワークパラメータ 13' に対する量子化セットは、先行するニューラルネットワークパラメータに対する量子化インデックス 56 の値に基づいて決定される。図 11 の好ましい例を考え、2 つの量子化セットを比較すると、ゼロに等しい再構

成レベルと隣接する再構成レベルとの間の距離は、セット 0 においてセット 1 よりも大きいことは明らかである。したがって、量子化インデックス 5 6 が 0 に等しい確率は、セット 0 が使用される場合により大きく、セット 1 が使用される場合により小さくなる。好ましい実施形態では、この効果は、現在の量子化インデックスに使用される量子化セット（または状態）に基づいてコードワード表または確率モデルを切り替えることによって、エントロピー符号化で利用される。

【0216】

コードワード表または確率モデルの好適な切り替えのために、現在の量子化インデックス（または現在の量子化インデックスの対応する 2 値決定）をエントロピー復号化するときに、すべての先行する量子化インデックスのパス（使用する量子化セットのサブセットとの関連）が知られていなければならないことに注意されたい。したがって、ニューラルネットワークパラメータ 1 3 が再構成順序で符号化されることが必要である。したがって、好ましい実施形態では、ニューラルネットワークパラメータ 1 3 の符号化順序は、それらの再構成順序に等しい。その態様のほかに、量子化インデックス 5 6 の任意の符号化 / 再構成順序が可能であり、例えば、2 . 2 . 1 節で規定した順序は、他の任意の一意的に定義された順序である。

【0217】

言い換えれば、本発明による実施形態は、例えば、以前に符号化されたニューラルネットワークパラメータの量子化インデックスに追加的に依存する確率モデルを用いて、ニューラルネットワークパラメータを符号化するための装置を含む。

【0218】

それぞれ、本発明による実施形態は、例えば、ニューラルネットワークパラメータを復号化するための装置であって、以前に復号化されたニューラルネットワークパラメータの量子化インデックスに追加的に依存する確率モデルを使用する装置を含む。

【0219】

絶対レベルのためのピンの少なくとも一部は、典型的には、適応的な確率モデル（コンテキストとも呼ばれる）を使用して符号化される。本発明の好ましい実施形態では、1 つ以上のピンの確率モデルは、対応するニューラルネットワークパラメータの量子化セット（または、より一般的には、対応する状態変数、例えば、表 1 ~ 3 のいずれかによる関係）に基づき選択される。選択された確率モデルは、既に送信された量子化インデックス 5 6 の複数のパラメータまたは特性に依存することができるが、パラメータの 1 つは、符号化される量子化インデックスに適用される量子化セットまたは状態である。

【0220】

言い換えれば、実施形態による、例えばニューラルネットワークパラメータ 1 3 を符号化するための装置は、現在のニューラルネットワークパラメータ 1 3 ' に対して選択された再構成レベルの状態またはセット 4 8 に応じて、複数の確率モデルのうち確率モデルのサブセットを事前に選択し、以前に符号化されたニューラルネットワークパラメータの量子化インデックスに応じて 1 2 1 確率モデルのサブセットのうち現在のニューラルネットワークパラメータの確率モデルを選択するように構成され得る。

【0221】

実施形態による、例えばニューラルネットワークパラメータ 1 3 を復号化するための装置は、現在のニューラルネットワークパラメータ 1 3 ' に対して選択された再構成レベルの状態またはセット 4 8 に応じて、複数の確率モデルのうち確率モデルのサブセットを事前選択し、以前に復号化されたニューラルネットワークパラメータの量子化インデックスに応じて 1 2 1、確率モデルのサブセットのうち現在のニューラルネットワークパラメータの確率モデルを選択するように構成されてもよい。

【0222】

例えば図 9 の文脈で説明したような発明的概念と組み合わせて、例えばニューラルネットワークパラメータ 1 3 の符号化および / または復号化のための、本発明による実施形態は、第 1 の状態または再構成レベルセットに対して事前選択されたサブセットが、任意の

他の状態または再構成レベルセットに対して事前選択されたサブセットと非結合であるように、現在のニューラルネットワークパラメータ 13' に対して選択された再構成レベルの状態またはセット 48 に応じて、複数の確率モデルの中から確率モデルのサブセットを事前選択するように構成される装置を含む。

【0223】

特に好ましい実施形態では、層の量子化インデックスを送信するためのシンタックスは、量子化インデックスが 0 に等しいか、または 0 に等しくないかを指定するピン、例えば前述の `sig_flag` を含む。このピンの符号化に用いる確率モデルは、2 つ以上の確率モデルのセットの中から選択される。使用される確率モデルの選択は、対応する量子化インデックス 56 に適用される量子化セット（即ち、再構成レベルのセット）に依存する。本発明の別の実施形態では、使用される確率モデルは、現在の状態変数に依存する（状態変数は、使用される量子化セットを意味する）。

10

【0224】

さらなる実施形態では、層の量子化インデックスを送信するためのシンタックスは、量子化インデックスがゼロより大きいか、ゼロより小さいかを指定する `bin`、例えば前述の `sign_flag` を含む。すなわち、`bin` は、量子化インデックスの符号を示す。使用される確率モデルの選択は、対応する量子化インデックスに適用される量子化セット（即ち、再構成レベルのセット）に依存する。別の実施形態では、使用される確率モデルは、現在の状態変数に依存する（状態変数は、使用される量子化セットを意味する）。

20

【0225】

さらなる実施形態では、量子化インデックスを送信するためのシンタックスは、量子化インデックスの絶対値（ニューラルネットワークのパラメータレベル）が X より大きいかどうかを指定するピン、例えば前述の `abs_level_greater_X`（詳細はセクション 0 を参照）を含んでいる。このピンの符号化に用いる確率モデルは、2 つ以上の確率モデルのセットの中から選択される。使用される確率モデルの選択は、対応する量子化インデックス 56 に適用される量子化セット（即ち、再構成レベルセット）に依存する。別の実施形態では、使用される確率モデルは、現在の状態変数に依存する（状態変数は、使用される量子化セットを意味する）。

【0226】

本明細書で論じる実施形態の 1 つの有利な態様は、ニューラルネットワークパラメータ 13 の依存量子化がエントロピー符号化と組み合わせられ、量子化インデックスの 2 値表現の 1 つ以上のピン（これは量子化レベルとも呼ばれる）に対する確率モデルの選択が、現在の量子化インデックスに対する量子化セット（許容再構成レベルのセット）または対応する状態変数に依存することである。量子化セット 52（または状態変数）は、符号化および再構成順序における先行するニューラルネットワークパラメータの量子化インデックス 56（または量子化インデックスを表すピンのサブセット）により与えられる。

30

【0227】

好ましい実施形態において、確率モデルの記述された選択は、以下のエントロピー符号化の側面のうちの 1 つ以上と組み合わせられる。

- ・量子化インデックスの絶対値は、適応型確率モデルを用いて符号化される多数のピンと、適応的に符号化されたピンが既に完全に絶対値を指定していない場合、算術符号化エンジンのバイパスモードで符号化されるサフィックス部分（全ピンに対して `pmf`（例えば確率質量関数）（0.5、0.5）を有する非適応型確率モデル）と、からなる 2 値方式を用いて送信される。好ましい実施形態では、サフィックス部分に使用される 2 値化は、既に送信された量子化インデックスの値に依存する。

40

- ・量子化インデックスの絶対値に対する 2 値化は、量子化インデックスが 0 に等しくないかどうかを指定する適応的に符号化されたピンを含む。このピンの符号化に用いられる確率モデル（コンテキストと呼ばれる）は、候補確率モデルのセットの中から選択される。選択された候補確率モデルは、現在の量子化インデックス 56 に対する量子化セット（許容再構成レベルのセット）または状態変数によって決定されるだけでなく、加えて、その

50

層に対する既に送信された量子化インデックスによっても決定される。好ましい実施形態では、量子化セット（または状態変数）は、利用可能な確率モデルのサブセット（コンテキストセットとも呼ばれる）を決定し、既に符号化された量子化インデックスの値は、このサブセット（コンテキストセット）内において使用される確率モデルを決定する。

【0228】

実施形態では、コンテキストセット内の使用される確率モデルは、現在のニューラルネットワークパラメータの局所近傍における既に符号化された量子化インデックスの値、例えば、2.2.3で説明したようなテンプレートに基づいて決定される。以下では、局所近傍の量子化インデックスの値に基づいて導出され、その後、事前に決定されたコンテキストセットの確率モデルを選択するために使用することができるいくつかの例示的な量（*measure*）をリストアップする。

10

- ・局所近傍で0に等しくない量子化インデックスの符号。

- ・局所近傍領域で0に等しくない量子化インデックスの数。この数は最大値にクリップ（短縮化）される可能性がある。

- ・局所近傍における量子化インデックスの絶対値の合計。この数値は最大値にクリップされる可能性がある。

- ・局所近傍における量子化インデックスの絶対値の合計と、局所近傍における0に等しくない量子化インデックスの数との差。この数値は最大値にクリップされる可能性がある。

【0229】

言い換えれば、本発明による実施形態は、例えば、現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する以前に符号化されたニューラルネットワークパラメータの量子化インデックスの特性に応じて確率モデルのサブセットから現在のニューラルネットワークパラメータの確率モデルを選択するように構成されたニューラルネットワークパラメータの符号化のための装置を含み、特性は、以下のもののうち1つまたは複数を含む。

20

現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する、以前に符号化されたニューラルネットワークパラメータのゼロでない量子化インデックスの符号。

現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する、以前に符号化されたニューラルネットワークパラメータの量子化インデックスの数であって、ゼロでない数

30

現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する、以前に符号化されたニューラルネットワークパラメータの量子化インデックスの絶対値の合計値

現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する、以前に符号化されたニューラルネットワークパラメータの量子化インデックスの絶対値の合計と、

現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する、以前に符号化されたニューラルネットワークパラメータの量子化インデックスの数であって、ゼロでない数と、
の間の差。

40

【0230】

それぞれ、本発明による実施形態は、例えば、ニューラルネットワークパラメータの復号化のための装置である。装置は、確率モデルのサブセットのうち、現在のニューラルネットワークパラメータのための確率モデルを、現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する以前に復号化されたニューラルネットワークパラメータの量子化インデックスの特性に応じて選択するように構成されている。特性は、以下のもののうち1つまたは複数を含む。

現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する、以前に復号化されたニューラルネットワークパラメータのゼ

50

ロでない量子化インデックスの符号。

現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する、以前に復号化されたニューラルネットワークパラメータの量子化インデックスの数であって、ゼロでない数。

現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する、以前に復号化されたニューラルネットワークパラメータの量子化インデックスの絶対値の合計値。

現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する、以前に復号化されたニューラルネットワークパラメータの量子化インデックスの絶対値の合計と、

現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する、以前に復号化されたニューラルネットワークパラメータの量子化インデックスの数であって、ゼロでない数と、
の差。

【0231】

・量子化インデックスの絶対値に対する2値化は、量子化インデックスの絶対値がXより大きいかどうかを指定する適応的に符号化されたピン、例えば `abs_level_greater_X` を含む。これらのピンを符号化するために使用される確率モデル（コンテキストと呼ばれる）は、候補確率モデルのセットの中から選択される。選択された確率モデルは、現在の量子化インデックスに対する量子化セット（許容再構成レベルのセット）または状態変数によって決定されるだけでなく、さらに、例えば前述のようなテンプレートを用いて、層に対する既に送信された量子化インデックスによって決定される。好ましい実施形態では、量子化セット（または状態変数）は、利用可能な確率モデルのサブセット（コンテキストセットとも呼ばれる）を決定し、既に符号化された量子化インデックスのデータは、このサブセット（コンテキストセット）内部で使用される確率モデルを決定する、例えば他の言葉で言えば、決定に用いることが可能である。確率モデルの選択には、上述した方法（量子化インデックスが0に等しくないかどうかを指定するピンの場合）のいずれかを使用することができる。

【0232】

さらに、本発明による装置は、以前に符号化されたニューラルネットワークパラメータ13が、現在のニューラルネットワークパラメータ13'と同じニューラルネットワーク層に関連するように、以前に符号化されたニューラルネットワークパラメータ13を位置づけるように構成されることができる。

【0233】

さらに、本発明による例えばニューラルネットワークパラメータを符号化するための装置は、1つ以上の以前に符号化されたニューラルネットワークパラメータが、現在のニューラルネットワークパラメータが参照するニューロン相互接続11が関連するニューロン10cまたは該ニューロンに隣接する別のニューロンから出現するニューロン相互接続、またはこれらのニューロンに向かうニューロン相互接続に関連するように、以前に符号化されたニューラルネットワークパラメータのうちの1つ以上のパラメータを位置付けるように構成されることができる。

【0234】

さらなる実施形態による装置は、量子化インデックスを2値化したものの1つ以上のリーディングピン（`leading bin`）について以前に符号化されたニューラルネットワークパラメータに対応する確率モデルを使用することによって、及び、1つ以上のリーディングピンに続く量子化インデックスを2値化したものの等確率バイパスモードサフィックスピン（`equi-probable bypass mode suffix bins`）を使用することによって、2値算術符号化を使用して、データストリーム14へ現在のニューラルネットワークパラメータ13'に対する量子化インデックス56を符号化するように構成されることができる。

10

20

30

40

50

【 0 2 3 5 】

量子化インデックスを2値化したもののサフィックスピンは、量子化インデックスの値を2値化するためのサフィックス2値化の2値化コードのピンを表し、その絶対値が1つ以上のリーディングピンで表現できる最大絶対値を超えている場合がある。したがって、本発明の実施形態による装置は、以前に符号化されたニューラルネットワークパラメータ13の量子化インデックス56に応じて、サフィックス2値化を選択するように構成されることができる。

【 0 2 3 6 】

それぞれ、本発明による、例えばニューラルネットワークパラメータを復号化するための装置は、以前に復号されたニューラルネットワークパラメータが、現在のニューラルネットワークパラメータ13'と同じニューラルネットワーク層に関連するように、以前に復号化されたニューラルネットワークパラメータ13を位置付けるように構成されてもよい。

【 0 2 3 7 】

さらなる実施形態によれば、例えば本発明によるニューラルネットワークパラメータを復号化するための装置は、1つ以上の以前に復号化されたニューラルネットワークパラメータが、現在のニューラルネットワークパラメータが参照するニューロン相互接続が関連するニューロン10cまたは該ニューロンに隣接する別のニューロンから出現するニューロン相互接続、またはこれらのニューロンに向かうニューロン相互接続に関連するように、以前に復号化されたニューラルネットワークパラメータ13のうちの1つ以上のパラメータを位置付けるように構成されることができる。

【 0 2 3 8 】

さらなる実施形態による装置は、量子化インデックスを2値化したものの1つ以上のリーディングピンについて以前に復号化されたニューラルネットワークパラメータに対応する確率モデルを使用することによって、及び、1つ以上のリーディングピンに続く量子化インデックスを2値化したものの等確率バイパスモードサフィックスピンを使用することによって、2値算術符号化を使用してデータストリーム14から現在のニューラルネットワークパラメータ13'に対する量子化インデックス56を復号化するように構成されることができる。

【 0 2 3 9 】

量子化インデックスの2値化のサフィックスピンは、量子化インデックスの値を2値化するためのサフィックス2値化の2値化符号のピンを表し、その絶対値が1つ以上のリーディングピンで表現可能な最大絶対値を超える。したがって、実施形態に従う装置は、以前に復号化されたニューラルネットワークパラメータの量子化インデックスに応じてサフィックス2値化を選択するように構成されることができる。

【 0 2 4 0 】

4. 5 符号化の方法例

歪み（再構成品質）とビットレートのトレードオフが非常に良いビットストリームを得るためには、ラグランジュコスト量（Lagrangian cost measure）が最小となるように量子化インデックスを選択する必要がある。

$$D + \lambda \cdot R = \sum_k D_k + \lambda \cdot R_k = \sum_k \alpha_k \cdot (t_k - t'_k)^2 + \lambda \cdot R(q_k | q_{k-1}, q_{k-2}, \dots)$$

独立スカラー量子化では、このような量子化アルゴリズム（レート歪み最適化量子化、rate-distortion optimized quantizationまたはRDOQと呼ばれる）が2.1.1節で議論された。しかし、独立スカラー量子化と比較して、さらなる難題がある。再構成されたニューラルネットワークのパラメータ t_k' 、したがって、その歪み $D_k = |t_k - t'_k|$ （または $D_{k,MSE} = (t_k - t'_k)^2$ ）は、関連する量子化インデックス q_k 56に依存するだけでなく、符号化順序において先行する量子化インデックスの値にも依存する。

【 0 2 4 1 】

しかし、4. 3. 3 節で説明したように、ニューラルネットワークのパラメータ 1 3 間の依存関係は、トレリス構造を用いて表現することができる。さらなる説明のために、図 1 1 に与えられた好ましい実施形態を例として用いる。8 個のニューラルネットワークパラメータのセットの例に対するトレリス構造を図 1 9 に示す。図 1 9 は、本発明の実施形態による、コスト量（ラグランジュコスト量 $D + \lambda \cdot R$ など）を最小化する量子化インデックスのシーケンス（またはブロック）を決定するために利用することができるトレリス構造の例を示している。トレリス構造は、4 つの状態を有する依存量子化の好ましい例を表している（図 1 8 参照）。トレリスは、8 つのニューラルネットワークパラメータ（または量子化インデックス）に対して示されている。最初の状態（一番左）は初期状態を表し、0 に等しいと仮定する。トレリスを通るパス（左から右へ）は、量子化インデックス 5 6 に対して可能な状態遷移を表す。2 つのノード間の各接続は、特定のサブセット（A、B、C、D）の量子化インデックスを表していることに注意が必要である。各サブセット（A、B、C、D）から量子化インデックス q_k 5 6 を選び、対応するレート歪コストを 2 つのトレリスノード間の関連する接続に割り当てた場合、以下になる。

$$J_k = D_k(q_k | q_{k-1}, q_{k-2}, \dots) + \lambda \cdot R_k(q_k | q_{k-1}, q_{k-2}, \dots)$$

全体のレート歪みコスト $D + \lambda \cdot R$ を最小化する量子化インデックスのベクトル／ブロックを決定する問題は、トレリスを通る最小コストのパスを見つけることと同等である（図 1 9 において左から右へ）。エントロピー符号化におけるいくつかの依存性を無視すれば、この最小化問題は、よく知られたビタビアルゴリズム（V i t e r b i a l g o r i t h m）を用いて解くことができる。

【 0 2 4 2 】

言い換えれば、本発明による実施形態は、選択および／または量子化を実行するためにビタビアルゴリズムおよびレート歪みコスト量を使用するように構成された装置を含んでいる。

【 0 2 4 3 】

層のための適切な量子化インデックスを選択するための例示的な符号化アルゴリズムは、以下の主要なステップで構成され得る。

1. 初期状態でのレート歪みコストを 0 に設定する。
2. 符号化順序におけるすべてのニューラルネットワークのパラメータ 1 3 に対して、次のようにする。

a. 各サブセット A、B、C、D について、与えられたオリジナルのニューラルネットワークパラメータ 1 3 に対する歪みを最小化する量子化インデックス 5 6 を決定する。

b. 現在のニューラルネットワークパラメータ 1 3 ' に対する全てのトレリスノード（0、1、2、3）に対して、以下を行う。

i 先行するニューラルネットワークパラメータ 1 3 の状態と現在の状態とを結ぶ 2 つのパスのレート歪みコストを計算する。コストは、先行する状態に対するコストと、 $D_k + \lambda \cdot R_k$ の合計として与えられ、ここで、 D_k と R_k は、考慮される接続に関連するサブセット（A、B、C、D）の量子化インデックスを選択するための歪みとレートを表す。

i i 計算されたコストの最小値を現在のノードに割り当て、最小コストのパスを表さない以前のニューラルネットワークパラメータ 1 3 の状態への接続を取り除く。

注：このステップの後、現在のニューラルネットワークパラメータ 1 3 ' に対するすべてのノードは、先行するニューラルネットワークパラメータ 1 3 に対する任意のノードへの単一の接続を有する。

3. (符号化順序における最後のパラメータについての) 4つの最終ノードのコストを比較し、コストが最小のノードを選択する。このノードは、トレリスを通るユニークなパスに関連していることに注意(他のすべての接続は以前のステップで取り除かれた)。

4. (最終ノードで指定される) 選択したパスを逆順にたどり、トレリスノード間の接続に関連する量子化インデックス56を収集する。

【0244】

ビタビアルゴリズムに基づく量子化インデックス56の決定は、独立スカラー量子化のためのレート歪み最適化量子化(RDOQ)よりも実質的に複雑ではないことに留意されたい。それにもかかわらず、依存量子化のためのより単純な符号化アルゴリズムも存在する。例えば、予め定義された初期状態(または量子化セット)から開始して、量子化インデックス56は、現在の量子化インデックスの影響のみを考慮する任意のコスト量を最小化することによって、符号化/再構成の順序で決定されることができる。現在のパラメータに対する決定された量子化インデックス(および先行するすべての量子化インデックス)が与えられると、次のニューラルネットワークパラメータ13に対する量子化セットは既知である。そして、したがって、このアルゴリズムは、符号化順序ですべてのニューラルネットワークパラメータに適用することができる。

【0245】

以下の実施形態による方法は、図20、図21、図22、図23に示されている。

【0246】

図20は、ニューラルネットワークを定義するニューラルネットワークパラメータをデータストリームから復号化するための方法400のブロック図である。方法400は、現在のニューラルネットワークパラメータについて、以前のニューラルネットワークパラメータについてデータストリームから復号化された量子化インデックスに応じて、複数の再構成レベルセットから再構成レベルセットを選択54するステップによって、及び、データストリームから現在のニューラルネットワークパラメータの量子化インデックス420を復号化するステップによって、ここで量子化インデックスが現在のニューラルネットワークパラメータのための再構成レベルの選択されたセットのうちの1つの再構成レベルを示し、及び、現在のニューラルネットワークパラメータのための量子化インデックスによって示される再構成レベルの選択されたセットのうちの1つの再構成レベル上に現在のニューラルネットワークパラメータを逆量子化62するステップによって、ニューラルネットワークパラメータを順次復号化することを含む。

【0247】

図21は、データストリームからニューラルネットワークを定義するニューラルネットワークパラメータを符号化するための方法500のブロック図である。方法500は、現在のニューラルネットワークパラメータについて、以前に符号化されたニューラルネットワークパラメータについてデータストリームに符号化された量子化インデックスに応じて複数の再構成レベルセットの中から再構成レベルセットを選択54するステップによって、及び、現在のニューラルネットワークパラメータを、選択された再構成レベルのセットのうち1つの再構成レベルに量子化64するステップによって、及び、現在のニューラルネットワークパラメータの量子化インデックスが量子化される1つの再構成レベルを示す現在のニューラルネットワークパラメータの量子化インデックスをデータストリームに符号化530するステップによって、ニューラルネットワークパラメータを連続的に符号化することを含む。

【0248】

図22は、本発明の実施形態による、ニューラルネットワークを定義するニューラルネットワークパラメータを再構築する方法のブロック図である。方法600は、ニューラルネットワークパラメータごとに、第1の再構成層のニューラルネットワークパラメータ値をもたらすために、第1の再構成層についての第1のニューラルネットワークパラメータ610を導出するステップを含む。方法600はさらに、データストリームから第2の再構成層の第2のニューラルネットワークパラメータを復号化620(例えば図6の矢印3

10

20

30

40

50

12で示すように)して、ニューラルネットワークパラメータごとに、第2の再構成層のニューラルネットワークパラメータ値をもたらすステップ、及び、ニューラルネットワークパラメータごとに、第1の再構成層のニューラルネットワークパラメータ値および第2の再構成層のニューラルネットワークパラメータ値を組み合わせることによって、ニューラルネットワークパラメータを再構成する630(例えば図6の矢印314で示す)ステップを含んでいる。

【0249】

図23は、本発明の実施形態による、ニューラルネットワークを定義するニューラルネットワークパラメータを符号化するための方法のブロック図である。方法700は、ニューラルネットワークパラメータごとに、第1の再構成層のニューラルネットワークパラメータ値を含む第1の再構成層の第1のニューラルネットワークパラメータを使用し、第2の再構成層の第2のニューラルネットワークパラメータをデータストリームに(例えば、図中に矢印322で示すように)符号化710するステップを含み、第2の再構成層は、ニューラルネットワークパラメータごとに、第2の再構成層ニューラルネットワークパラメータ値を含み、ニューラルネットワークパラメータは、ニューラルネットワークパラメータごとに、第1の再構成層ニューラルネットワークパラメータ値と第2の再構成層ニューラルネットワークパラメータ値とを組み合わせることによって再構成可能である。

【0250】

10

20

30

40

50

以下では、本発明による追加の実施形態が示される。

quant_tensor(dimensions, maxNumNoRem, entryPointOffset) {	
stateId = 0	997
bitPointer = get_bit_pointer()	998
lastOffset = 0	999
for(i = 0; i < Prod(dimensions); i++) {	1000
idx = TensorIndex(dimensions, i, scan_order)	1001
if(entryPointOffset != -1 && GetEntryPointIdx(dimensions, i, scan_order) != - 1) {	1002
lvlCurrRange = 256	1003
j = entryPointOffset + GetEntryPointIdx(dimensions, i, scan_order)	1004
lvlOffset = cabac_offset_list[j]	1005
if(dq_flag)	1006
stateId = dq_state_list[j]	1007
set_bit_pointer(bitPointer + lastOffset + BitOffsetList[j])	1008
lastOffset = BitOffsetList[j]	1009
Invoke initialisation process for probability estimation parameters	1010
}	1011
int_param(idx, maxNumNoRem, stateId)	1012
if(dq_flag) {	1013
nextSt = StateTransTab[stateId][QuantParam[idx] & 1]	1014
if(QuantParam[idx] != 0) {	1015
QuantParam[idx] = QuantParam[idx] << 1	1016
if(QuantParam[idx] < 0)	1017
QuantParam[idx] += stateId & 1	1018
else	1019
QuantParam[idx] += - (stateId & 1)	1020
}	1021
stateId = nextSt	1022
}	
}	

【 0 2 5 1 】

例えば 1 0 1 4 行目に示す 2 次元整数アレイ `StateTransTab [] []` は、
依存スカラー量子化の状態遷移表を指定するもので、次のようになる。

【 0 2 5 2 】

10

20

30

40

StateTransTab[][] = { {0, 2}, {7, 5}, {1, 3}, {6, 4}, {2, 0}, {5, 7}, {3, 1}, {4, 6} }

int_param(i, maxNumNoRem, stateId) {		
QuantParam[i] = 0	5997	
sig_flag	5998	
if(sig_flag) {	5999	
QuantParam[i]++	6000	
sign_flag	6001	10
j = -1	6002	
do {	6003	
j++	6004	
abs_level_greater_x[j]	6005	
QuantParam[i] += abs_level_greater_x[j]	6006	
} while(abs_level_greater_x[j] == 1 && j < maxNumNoRem)	6007	
if(j == maxNumNoRem) {	6008	20
RemBits = 0	6009	
j = -1	6010	
do {	6011	
j++	6012	
abs_level_greater_x2[j]	6013	
if(abs_level_greater_x2[j]) {	6014	
RemBits++	6015	
QuantParam[i] += 1 << RemBits	6016	
}	6017	30
} while(abs_level_greater_x2[j] && j < 30)	6018	
abs_remainder	6019	
QuantParam[i] += abs_remainder	6020	
}	6021	
QuantParam[i] = sign_flag ? -QuantParam[i] : QuantParam[i]	6022	
}		
}		40

【 0 2 5 3 】

このプロセスへの入力は次の通りである。

- ・復号化されるテンソルの次元を指定する変数 `tensorDims`。
- ・復号化のためのエントリポイントが存在するかどうか、およびエントリポイントが存在する場合エントリポイントオフセットを示す変数 `entryPointOffset`。
- ・コードブックの有無と、コードブックが適用されるかどうか、及びコードブックが適用される場合はどのコードブックを使用するかを示す変数 `codebookId`。

このプロセスの出力は、`TENSOR_FLOAT` 型の変数 `recParam` であり、次元は `tensorDims` と等しい。

【 0 2 5 4 】

変数 `stepSize` は以下のように導出される。

```

3001 mul = (1 << QpDensity) + ( (qp_value + QuantizationParameter) & ( ( 1
    << QpDensity ) - 1 ) )
3002 shift = (qp_value + QuantizationParameter) >> QpDensity
3003 stepSize = mul * 2shift - QpDensity
Variable recParam is updated as follows:
4001 recParam = recParam * stepSize

```

注－上記の計算から、`recParam`は常に2進小数で表すことができる。

10

【 0 2 5 5 】

シンタックス要素 `sig_flag` に対して、使用するコンテキストまたは確率推定を示す `ctxInc` の導出プロセスについて。

【 0 2 5 6 】

このプロセスへの入力は、現在の `sig_flag` の前に復号化された `sig_flag`、状態値 `stateId` および関連する `sign_flag`（存在する場合）である。現在の `sig_flag` の前に復号化された `sig_flag` がない場合、それは0と見なされる。また、以前に復号化された `sig_flag` に関連する `sign_flag` が復号化されていない場合、それは0と見なされる。

20

【 0 2 5 7 】

このプロセスの出力は、変数 `ctxInc` である。

変数 `ctxInc` は、以下のように導出される。

- ・ `sig_flag` が0であれば、`ctxInc` には `stateId * 3` が設定される。
- ・ そうでなければ、`sign_flag` が0に等しい場合、`ctxInc` は `stateId * 3 + 1` に設定される。
- ・ そうでなければ、`ctxInc` は `stateId * 3 + 2` に設定される。

【 0 2 5 8 】

上記の例は、ニューラルネットワークパラメータ13をデータストリーム14へ/から符号化/復号化する概念を示し、ニューラルネットワークパラメータ13は、ニューラルネットワーク10のニューロン相互接続11の重み、例えば重みテンソルの重みに関連し得る。ニューラルネットワークパラメータ13の復号化/符号化は、順次行われる。テンソルの次元ごとの重みの数の積と同数の重みでテンソルの重みを循環する `for - next` ループ1000を参照されたい。重みは所定の順序 `TensorIndex (dimensions, i, scan_order)` でスキャンされる。現在のニューラルネットワークパラメータ `idx13` について、2つの再構成レベルセット52のうち再構成レベルセットが、1018および1020において、以前のニューラルネットワークパラメータのデータストリームから復号化された量子化インデックス58に基づいて連続的に更新される量子化状態 `stateId` に応じて選択される。特に、現在のニューラルネットワークパラメータ `idx` に対する量子化インデックスが1012でデータストリームから復号化され、量子化インデックスは、現在のニューラルネットワークパラメータ13' に対する選択された再構成レベルセットのうち1つの再構成レベルを示している。2つの再構成レベルセットは、1016での複製と、それに続く1018および1020での量子化状態インデックスに応じた1またはマイナス1の加算によって定義される。ここで、1018及び1020において、現在のニューラルネットワークパラメータ13' は、現在のニューラルネットワークパラメータ13' に対する量子化インデックス `QuantParam [idx]` によって示される選択された再構成レベルセットのうちの1つの再構成レベル上に実際に逆量子化される。ステップサイズ `stepSize` は、3001～3003で再

30

40

50

構成レベルセットをパラメータ化するために使用される。この所定の量子化ステップサイズ `stepSize` に関する情報は、シンタックス要素 `qp_value` を介してデータストリームから導出される。後者は、それぞれテンソル全体または `NN` 層全体について、あるいは `NN` 全体についてデータストリームに符号化されるかもしれない。すなわち、ニューラルネットワーク 10 は、1 つ以上の `NN` 層 10 a、10 b を含むことができ、各 `NN` 層について、データストリーム 14 からそれぞれの `NN` 層について所定の量子化ステップサイズ (`QP`) の情報を導出し、その後、それぞれの `NN` 層について、それぞれの `NN` 層に属するニューラルネットワークのパラメータ 13 を逆量子化するために用いられるように、それぞれの `NN` 層について導出した所定の量子化ステップサイズを用いて複数の再構成レベルセットのパラメータ化を行うようにしても良い。

10

【0259】

`stateId = 0` に対する第 1 の再構成レベルセットは、ここではゼロおよび所定の量子化ステップサイズの偶数倍を含み、1018 および 1020 で分かるように、`stateId = 1` に対する第 2 の再構成レベルセットは、ゼロおよび所定の量子化ステップサイズ (`QP`) の奇数倍を含む。各ニューラルネットワークパラメータ 13 について、1015 ~ 1021 で、それぞれのニューラルネットワークパラメータ 13 について選択された再構成レベルセットとそれぞれのニューラルネットワークパラメータについてのエントロピー復号化量子化インデックス `QuantParam[idx]` に応じて中間整数値 `QuantParam[idx] (IV)` を導き、次に、各ニューラルネットワークパラメータについて、それぞれのニューラルネットワークパラメータの中間値に、4001 での

20

【0260】

2 つの再構成レベルセット (例えばセット 0、セット 1) のうちの再構成レベルセットの、現在のニューラルネットワークパラメータ 13 ' に対する選択は、1014 に示すように、以前に復号化されたニューラルネットワークパラメータのデータストリームから復号化された量子化インデックスの `LSB` 部分に応じて行われ、遷移表が、`stateId` が既に復号化された量子化インデックス 56 の過去のシーケンスに依存するように `QuantParam[idx]` の `LSB` に応じて `stateId` から次の量子化状態 `nextState` へ遷移される。したがって、状態遷移は、以前に復号化されたニューラルネットワークパラメータのデータストリームから復号化された量子化インデックス 56 の 2 値関数の結果、すなわちそのパリティに依存する。言い換えれば、現在のニューラルネットワークパラメータについて、複数の再構成レベルセットのうちの再構成レベルセットの選択は、現在のニューラルネットワークパラメータについて、1018 及び 1020 において、現在のニューラルネットワークパラメータに関連付けられた状態 `stateId` に応じて、複数の再構成レベルセットのうちの再構成レベルセットを決定することによって、及び、直前のニューラルネットワークパラメータ、すなわちこれまで `stateId` が決定されていたパラメータについてデータストリームから復号化された量子化インデックスに応じて、次に符号化 / 復号化される `NN` パラメータとは限らないが、次に `stateId` が決定されるべき後続のニューラルネットワークパラメータについて 1014 で `stateId` を更新することによって、状態遷移プロセスにより行われる。例えば、ここでは、現在のニューラルネットワークパラメータは、次に符号化 / 復号化されるべき `NN` パラメータの `stateId` をもたらすための更新のために使用される。1014 における更新は、直前の (現在の) ニューラルネットワークパラメータのデータストリームから復号化された量子化インデックスの 2 値関数、すなわちそのパリティを使用して行われる。状態遷移プロセスは、8 つの可能な状態間を遷移するように構成されている。遷移は、表 `StateTransTab[] []` を介して行われる。状態遷移プロセスにおいて、遷移はこれら 8 つの可能な状態の間で行われ、現在のニューラルネットワークパラメータについて、1018 および 1020 における、現在のニューラルネットワークパラメータに関連する状態 `stateId` に応じた量子化セットのうちの再構成レベルのセットの決定することで、状

30

40

50

態が偶数の可能な状態の前半、すなわち奇数の状態に属する場合には2つの再構成レベルセットのうちの第1の再構成レベルセットが決定され、状態が偶数の可能な状態の後半の状態、すなわち偶数の状態に属する場合には、2つの再構成レベルセットのうちの第2の再構成レベルセットが決定される。状態 `stateId` の更新は、遷移表 `StateTransTab[] []` により行われる。遷移表は、直前の（現在の）ニューラルネットワークパラメータのデータストリームから復号化された状態 `stateId` と量子化インデックス（58）のパリティ、`QuantParam[idx] & 1` の組み合わせを、後続のニューラルネットワークパラメータに関連する別の状態にマッピングする。

【0261】

現在のニューラルネットワークパラメータの量子化インデックスは、現在のニューラルネットワークパラメータに対して選択された再構成レベルセット、より正確には量子化状態 `stateId`、すなわち現在のニューラルネットワークパラメータ13'に対する状態に依存する確率モデルを用いて、算術符号化を用いてデータストリームに符号化し、及びデータストリームから復号化する。1012の関数 `int__param` を呼び出すときの第3のパラメータを参照されたい。特に、現在のニューラルネットワークパラメータに対する量子化インデックスは、量子化インデックスを2値化したものの少なくとも1つのピンに対する現在のニューラルネットワークパラメータについての状態に対応する確率モデルを用いて、2値算術符号化/復号化を用いてデータストリームに符号化され、データストリームから復号化されることができる。ここで、2値化の `sig__flag`、`sign__flag`（オプション）、`abs__level__greater__x[j]`、`abs__level__greater__x2[j]`、`abs__remainder` うちピン `sig__flag` は、現在のニューラルネットワークパラメータの量子化インデックス（56）がゼロに等しいか否かを示す有意性ピンである。確率モデルの依存性は、依存性を用いたニューラルネットワークパラメータのコンテキストセットのうちコンテキストを選択することを含み、各コンテキストは、所定の確率モデルが関連づけられるように構成される。ここで、`sig__flag` のコンテキストは、それぞれが2値確率モデルに関連付けられているコンテキストのリストからコンテキストをインデックス化するためのインデックスのインクリメンターとして `ctxInc` を使用して選択される。モデルは、コンテキストに関連付けられたピンを使用して更新されることができる。すなわち、それぞれのコンテキストに関連付けられた所定の確率モデルは、それぞれのコンテキストを用いて算術符号化された量子化インデックスに基づいて更新されることができる。

【0262】

（注）`sig__flag` の確率モデルは、さらに、以前に復号化されたニューラルネットワークパラメータの量子化インデックス、すなわち以前に復号化されたニューラルネットワークパラメータの `sig__flag`、およびその `sign__flag` - その符号を示すもの - に依存することに注意されたい。より正確には、状態 `stateId` に応じて、複数の確率モデルのうち、すなわちコンテキストインクリメンター状態 `0...23` のうち、確率モデルのサブセットが予め選択され、すなわち `{0...23}` のうち連続する3つのコンテキストを含むその8つ、`sig__flag` に対する確率モデルのサブセットのうち、現在のニューラルネットワークパラメータの確率モデルを、以前に復号されたニューラルネットワークパラメータの量子化インデックスに応じて（121）、すなわち以前のNNパラメータの `sig__flag` と `sign__flag` に基づいて選択する。`stateId` の最初の値に対して事前に選択されたサブセットは、`stateId` の他の値に対して事前に選択されたサブセットと不一致である。`sig__flag` と `sign__flag` が使用される以前のNNパラメータは、現在のニューラルネットワークパラメータが関連する部分に隣接するニューラルネットワークの部分に関連する。

【0263】

複数の実施形態が上述されてきた。実施形態の態様および特徴は、個々にまたは組み合わせて使用され得ることに留意されたい。さらに、本発明の第1および第2の側面による実施形態の態様および特徴は、組み合わせて使用されてもよい。

【 0 2 6 4 】

さらなる実施形態は、装置を構成し、ニューラルネットワークパラメータは、ニューラルネットワーク 10 が表される使用する再構成層のうちの 1 つの再構成層、例えばエンハンスメント層に関連する。装置は、ニューラルネットワークパラメータ単位で、対応する、例えば共通のニューロン相互接続に関連するもの、または率直に言って、異なる表現層における NN 層の行列表現に併置されるもの、1 つ以上の別の再構成層のニューラルネットワークパラメータと組み合わせることによって再構成可能であるように、ニューラルネットワークが構成されてもよい。

【 0 2 6 5 】

例えばこの実施形態で説明したように、本発明の第 1 および第 2 の側面の特徴および態様は、組み合わせられてもよい。第 2 の側面による従属請求項の任意の特徴は、さらなる実施形態をもたらすために、ここにも移転可能であるものとする。

10

【 0 2 6 6 】

さらに、本発明の態様による装置は、現在のニューラルネットワークパラメータ 13 ' に対する量子化インデックス 56 を、現在のニューラルネットワークパラメータに対応する対応するニューラルネットワークパラメータに対応する確率モデルを用いた算術符号化によりデータストリーム 14 の中に符号化するように構成されてもよい。

【 0 2 6 7 】

それぞれ、さらなる実施形態は、装置を構成し、ニューラルネットワークパラメータは、ニューラルネットワーク 10 が表現される再構成層のうちの 1 つの再構成層、例えばエンハンスメント層に関連するものである。装置は、ニューラルネットワークパラメータ単位で、ニューラルネットワークパラメータを、対応する例えば共通のニューロン相互接続に関連するもの、率直に言えば、異なる表現層における NN 層の行列表現において併置されるもの、1 つ以上の別の再構成層のニューラルネットワークパラメータと組み合わせることによって、ニューラルネットワークを再構成するように構成されてもよい。

20

【 0 2 6 8 】

例えばこの実施形態で説明したように、本発明の第 1 および第 2 の側面の特徴および態様は、組み合わせられてもよい。第 2 の側面による従属請求項の任意の特徴は、さらなる実施形態をもたらすために、ここにも移転可能であるものとする。

【 0 2 6 9 】

さらに、本発明の態様による装置は、現在のニューラルネットワークパラメータに対応する対応するニューラルネットワークパラメータに依存する確率モデルを用いて、算術符号化を用いてデータストリーム 14 から現在のニューラルネットワークパラメータ 13 ' の量子化インデックス 56 を復号化するように構成されていてもよい。

30

【 0 2 7 0 】

言い換えれば、再構成層のニューラルネットワークパラメータ、例えば説明したような第 2 のニューラルネットワークパラメータは、それぞれ図 3 および図 5 ならびに図 2 および図 4 に関して説明した概念に従って符号化 / 復号化および / または量子化 / 逆量子化することができる。

【 0 2 7 1 】

いくつかの態様は装置の文脈で説明されてきたが、これらの態様は、ブロックまたは装置が方法ステップまたは方法ステップの特徴に対応する、対応する方法の説明も表すことが明らかである。同様に、方法ステップの文脈で説明された側面は、対応するブロックまたは項目または対応する装置の特徴の説明も表している。

40

【 0 2 7 2 】

本発明データストリームは、デジタル記憶媒体に格納することができ、または、無線伝送媒体またはインターネットなどの有線伝送媒体などの伝送媒体で伝送することができる。

【 0 2 7 3 】

特定の実装要件に応じて、本発明の実施形態は、ハードウェアで実装することも、ソフトウェアで実装することも可能である。実装は、デジタル記憶媒体、例えばフロッピー（

50

登録商標) ディスク、DVD、CD、ROM、PROM、EPROM、EEPROMまたはフラッシュメモリであって、その上に格納された電子的に読み取り可能な制御信号を有し、それぞれの方法が実行されるようにプログラム可能なコンピュータシステムと協力する(または協力することができる)ものを使用して実行することができる。

【0274】

本発明によるいくつかの実施形態は、電子的に読み取り可能な制御信号を有するデータキャリアを備え、このデータキャリアは、本明細書に記載の方法の1つが実行されるように、プログラム可能なコンピュータシステムと協働することが可能である。

【0275】

一般に、本発明の実施形態は、プログラムコードを有するコンピュータプログラム製品として実施することができ、プログラムコードは、コンピュータプログラム製品がコンピュータ上で実行されるときに、方法の1つを実行するために動作可能である。プログラムコードは、例えば、機械読み取り可能な担体に格納することができる。

10

【0276】

他の実施形態は、本明細書に記載された方法の1つを実行するためのコンピュータプログラムを、機械可読担体に格納することからなる。

【0277】

言い換えれば、本発明方法の実施形態は、したがって、コンピュータプログラムがコンピュータ上で実行されるときに、本明細書に記載された方法の1つを実行するためのプログラムコードを有するコンピュータプログラムである。

20

【0278】

したがって、本発明方法のさらなる実施形態は、本明細書に記載の方法の1つを実行するためのコンピュータプログラムをその上に記録してなるデータキャリア(またはデジタル記憶媒体、またはコンピュータ読取可能な媒体)である。

【0279】

本発明方法のさらなる実施形態は、したがって、本明細書に記載された方法の1つを実行するためのコンピュータプログラムを表すデータストリームまたは信号のシーケンスである。データストリームまたは信号のシーケンスは、例えば、データ通信接続、例えば、インターネットを介して転送されるように構成されてもよい。

【0280】

30

さらなる実施形態は、本明細書に記載された方法の1つを実行するように構成された、または適合された処理手段、例えばコンピュータ、またはプログラマブルロジックデバイスを備える。

【0281】

さらなる実施形態は、本明細書に記載された方法の1つを実行するためのコンピュータプログラムをその上にインストールしたコンピュータを具備する。

【0282】

いくつかの実施形態では、プログラマブルロジックデバイス(例えば、フィールドプログラマブルゲートアレイ)を使用して、本明細書に記載の方法の機能性の一部または全部を実行してもよい。いくつかの実施形態では、フィールドプログラマブルゲートアレイは、本明細書に記載される方法の1つを実行するためにマイクロプロセッサと協働してよい。一般に、本方法は、好ましくは、任意のハードウェア装置によって実行される。

40

【0283】

上述した実施形態は、本発明の原理を単に例示するものである。本明細書に記載された配置および詳細の修正および変形は、当業者には明らかであろうことが理解される。したがって、差し迫った特許請求の範囲の範囲によってのみ限定され、本明細書における実施形態の説明および解説によって提示される特定の詳細によって限定されないことが意図される。

【0284】

50

References

- [1] C. W. P. V. J. C. J. T. B. C. E. S. Sharan Chetlur, "cuDNN: Efficient Primitives for Deep Learning," arXiv: 1410.0759, 2014
- [2] MPEG, "Working Draft 2 of Compression of neural networks for multimedia content description and analysis", Document of ISO/IEC JTC1/SC29/WG11, w18784, Geneva, Oct. 2019 10
- [3] D. Marpe, H. Schwarz und T. Wiegand, „Context-Based Adaptive Binary Arithmetic Coding in the H.264/AVC Video Compression Standard," *IEEE transactions on circuits and systems for video technology*, Vol. 13, No. 7, pp. 620-636, July 2003.
- [4] H. Kirchhoffer, J. Stegemann, D. Marpe, H. Schwarz und T. Wiegand, „JVET-K0430-v3 - CE5-related: State-based probability estimator," in *JVET*, Ljubljana, 2018. 20
- [5] ITU - International Telecommunication Union, „ITU-T H.265 High efficiency video coding," *Series H: Audiovisual and multimedia systems - Infrastructure of audiovisual services - Coding of moving video*, April 2015.
- [6] B. Bross, J. Chen und S. Liu, „JVET-M1001-v6 - Versatile Video Coding (Draft 4)," in *JVET*, Marrakech, 2019. 30

【 図 5 】

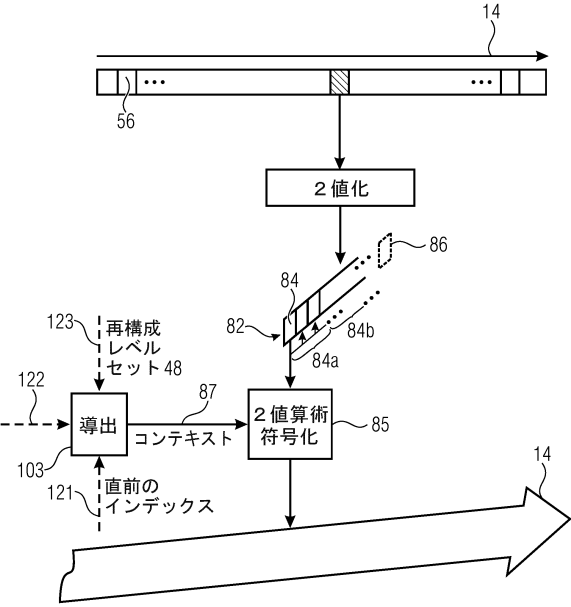


Fig. 5

【 図 6 】

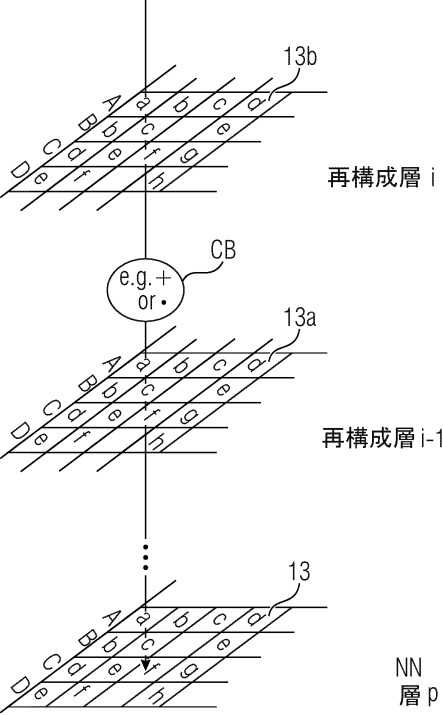


Fig. 6

【 図 7 】

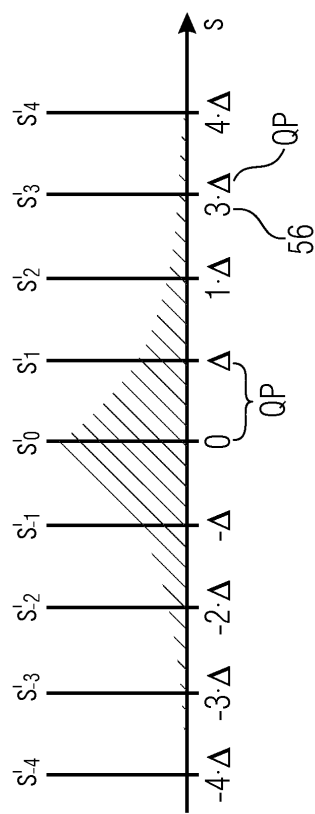


Fig. 7

【 図 8 】

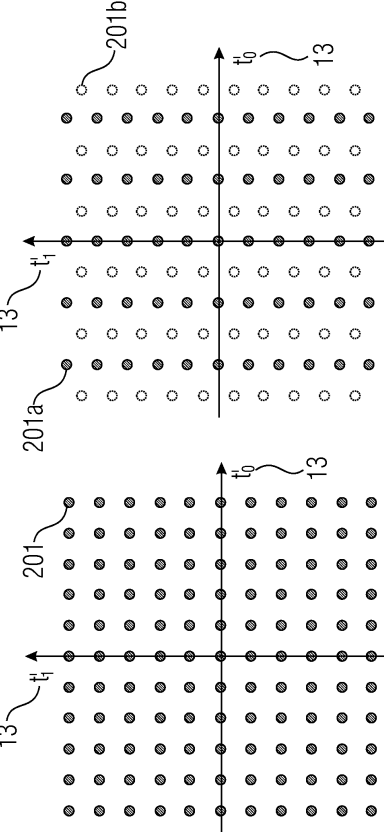


Fig. 8(b)

Fig. 8(a)

10

20

30

40

50

【図 1 1】

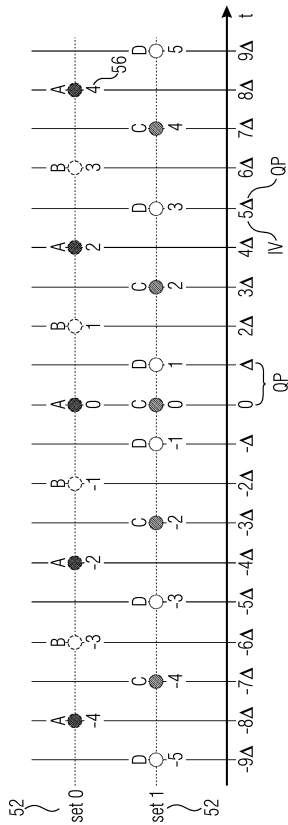


Fig. 11

【図 1 2】

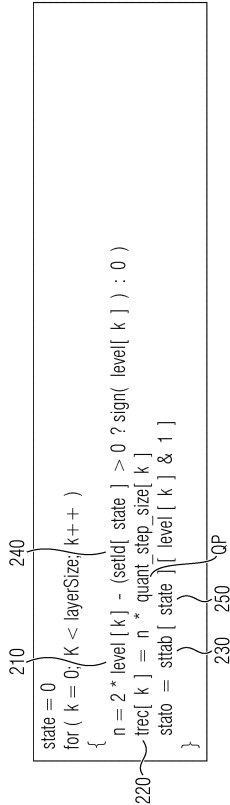


Fig. 12

【図 1 3】

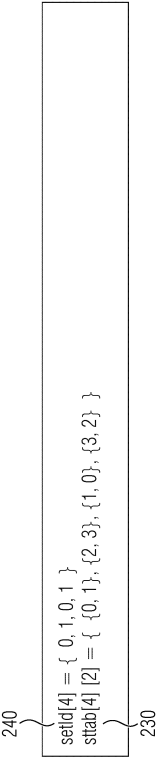


Fig. 13

【図 1 4】

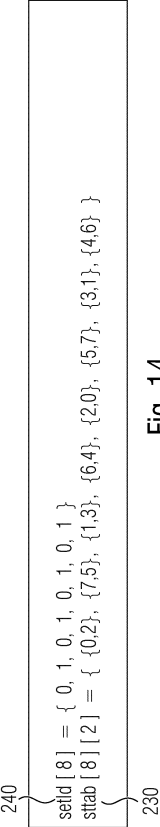


Fig. 14

【図 15】

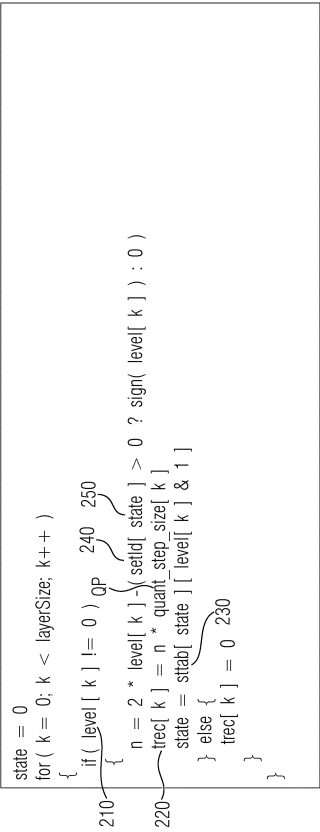


Fig. 15

【図 16】

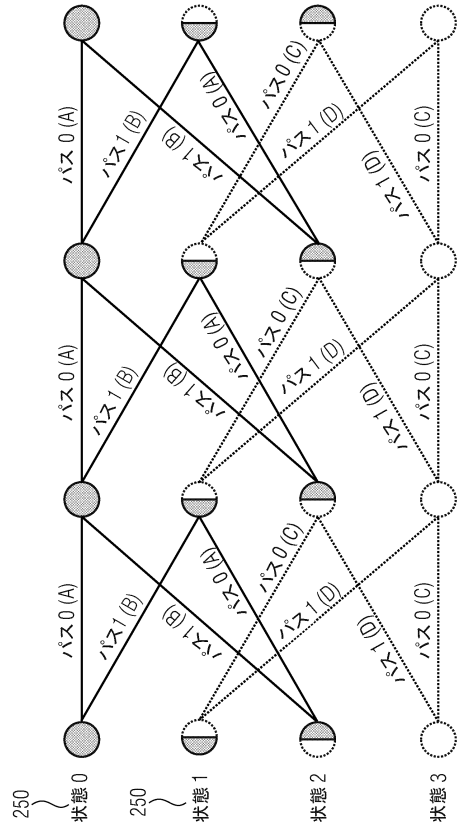


Fig. 16

【図 17】

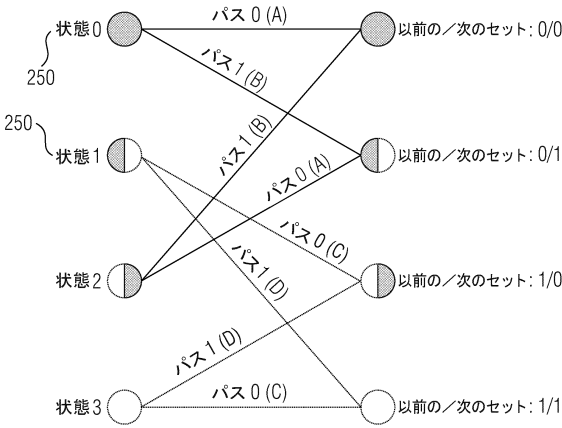


Fig. 17

【図 18】

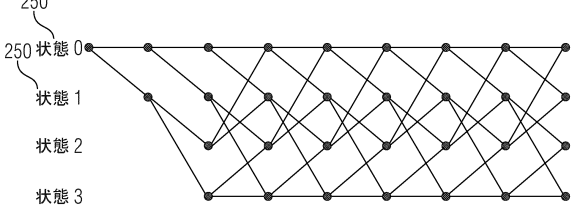


Fig. 18

10

20

30

40

50

【図 19】

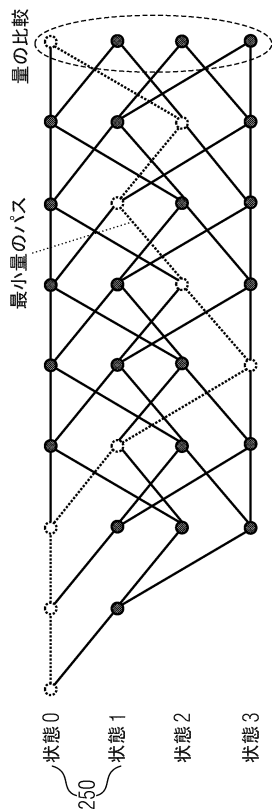


Fig. 19

【図 20】

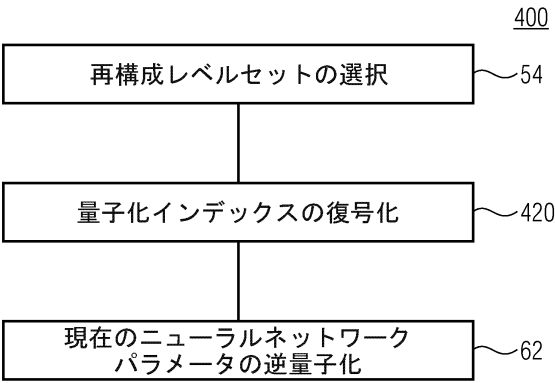


Fig. 20

【図 21】

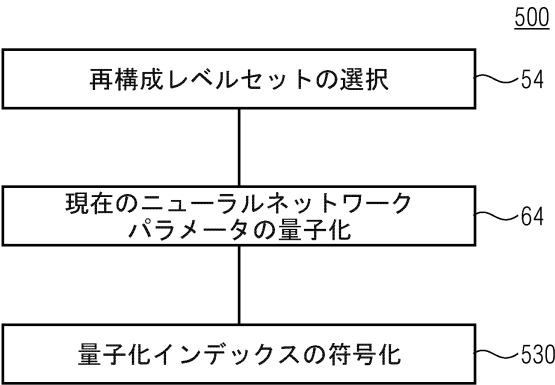


Fig. 21

【図 22】

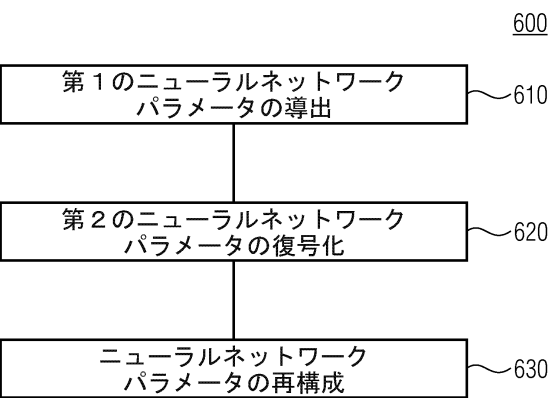


Fig. 22

10

20

30

40

50

【図 23】

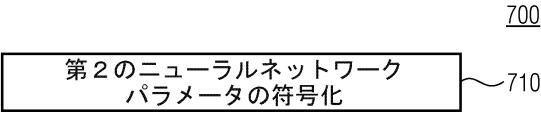


Fig. 23

10

20

30

40

50

フロントページの続き

- (72)発明者 キルヒホフファー ハイナー
ドイツ連邦共和国 1 0 5 8 7 ベルリン アインシュタインウーファー 3 7 フラウンホフファー
- インスティテュート フュア ナーハリヒテンテヒニーク ハインリッヒ - ヘルツ - インスティテ
ュート HHI 内
- (72)発明者 シュヴァルツ ハイコ
ドイツ連邦共和国 1 0 5 8 7 ベルリン アインシュタインウーファー 3 7 フラウンホフファー
- インスティテュート フュア ナーハリヒテンテヒニーク ハインリッヒ - ヘルツ - インスティテ
ュート HHI 内
- (72)発明者 マルベ デトレフ
ドイツ連邦共和国 1 0 5 8 7 ベルリン アインシュタインウーファー 3 7 フラウンホフファー
- インスティテュート フュア ナーハリヒテンテヒニーク ハインリッヒ - ヘルツ - インスティテ
ュート HHI 内
- (72)発明者 ウィーガント トーマス
ドイツ連邦共和国 1 0 5 8 7 ベルリン アインシュタインウーファー 3 7 フラウンホフファー
- インスティテュート フュア ナーハリヒテンテヒニーク ハインリッヒ - ヘルツ - インスティテ
ュート HHI 内
- 審査官 坂東 大五郎
- (56)参考文献 国際公開第 2 0 1 9 / 1 8 5 7 6 9 (WO , A 1)
Simon Wiedemann et al. , DeepCABAC: A universal compression algorithm for deep neural
networks , arXiv , 2019年07月27日 , pp.1-18
Marta Karczewicz et al. , CE8-related: Sign context modelling and level mapping for TS resi
dual coding , Joint Video Experts Team (JVET) , 2019年03月21日 , [JVET-N0455] (version 3)
- (58)調査した分野 (Int.Cl. , D B 名)
H 0 4 N 1 9 / 0 0 - 1 9 / 9 8