

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7621675号
(P7621675)

(45)発行日 令和7年1月27日(2025.1.27)

(24)登録日 令和7年1月17日(2025.1.17)

(51)国際特許分類	F I
G 0 6 F 21/60 (2013.01)	G 0 6 F 21/60 3 2 0
G 0 6 F 21/62 (2013.01)	G 0 6 F 21/62 3 1 8
G 0 6 N 3/02 (2006.01)	G 0 6 N 3/02
G 0 6 N 3/0464(2023.01)	G 0 6 N 3/0464

請求項の数 6 (全25頁)

(21)出願番号 特願2023-129058(P2023-129058)	(73)特許権者 514274672 延世大学校 産学協力団 UIF (University Industry Foundation), Yonsei University 大韓民国 03722 ソウル、ソデム -グ、ヨンセイ-ロ 50 50, YONSEI-RO, SEOD AEMUN-GU, SEOUL 037 22, REPUBLIC OF KOR EA
(22)出願日 令和5年8月8日(2023.8.8)	(74)代理人 110000051 弁理士法人共生国際特許事務所
(65)公開番号 特開2024-160912(P2024-160912 A)	(72)発明者 キム, ヨン ソク 大韓民国 ソウル特別市 西大門区 延世路 最終頁に続く
(43)公開日 令和6年11月15日(2024.11.15) 審査請求日 令和5年8月8日(2023.8.8)	
(31)優先権主張番号 10-2023-0058605	
(32)優先日 令和5年5月4日(2023.5.4)	
(33)優先権主張国・地域又は機関 韓国(KR) 特許法第30条第2項適用 ウェブサイトの掲載日 2 022年10月24日 ウェブサイトのアドレス ht t p s : / / d l . a c m . o r g / d o i / a b s / 1 0 . 1 1 4 5 / 3 5 2 8 5 3 5 . 3 5 3 1 5 1 3	

(54)【発明の名称】 トラスト環境ベースの人工知能装置

(57)【特許請求の範囲】

【請求項1】

暗号化入力データを送信し、暗号化出力データを受信する第1タイプのメモリと、
トラスト信頼空間で動作し、前記暗号化入力データ及び前記暗号化出力データの人工知
能演算を行うトラスト人工知能処理部と、を備え、

前記トラスト人工知能処理部は、

前記暗号化入力データの復号化を介して復号化入力データを生成し、前記暗号化出力デ
ータを生成するために非暗号化出力データの暗号化を行う暗号処理フロント-エンドプロ
セッサと、

前記復号化入力データ及び前記非暗号化出力データに対するバッファを提供し、前記第
1タイプのメモリよりも相対的に速い動作速度及び少ない格納容量を有する第2タイプのメモ
リと、

前記復号化入力データに基づいてニューラルネットワーク演算を行って前記非暗号化出
力データを生成するプロセッサと、を含み、

前記プロセッサは、

直接畳み込みベースのニューラルネットワーク演算を行って前記第1タイプのメモリのア
クセス数を減らし、

前記第2タイプのメモリを、畳み込みレイヤの入力アクティベーション、復号化フィルタ
及び非暗号出力アクティベーションのアクセス領域により管理し、

前記第2タイプのメモリに畳み込みレイヤの実行全体で繰り返しアクセスするときに、全

10

20

ての前記復号化入力データを前記入力アクティベーションの領域に固定的に格納し、残りの前記第2種類のメモリを、前記復号化フィルタ及び前記非暗号出力アクティベーションの領域の2つの循環キューで構成された循環キュー方式で格納し、

前記暗号処理フロント-エンドプロセッサの暗号化及び復号化演算の実行にオーバーラップするように前記ニューラルネットワーク演算を実行してイントラ-レイヤパイプライニングを実現し、

前記イントラ-レイヤパイプライニングのために、データ復号化ステップ、演算ステップ、及びデータ暗号化ステップに細分化して、前記ニューラルネットワーク演算を途切れることなく行うことを特徴とするトラスト環境ベースの人工知能装置。

【請求項2】

前記暗号処理フロント-エンドプロセッサは、前記第1種類のメモリから前記暗号化入力データとして暗号化入力アクティベーション及び暗号化フィルタの入力を受けて、前記第2種類のメモリに復号化入力アクティベーション及び復号化フィルタを格納することを特徴とする請求項1に記載のトラスト環境ベースの人工知能装置。

【請求項3】

前記暗号処理フロント-エンドプロセッサは、前記人工知能演算の過程で前記プロセッサによるオン-デマンド要求を受信し、前記第1種類のメモリにアクセスして前記暗号化入力データを取り込むことを特徴とする請求項1に記載のトラスト環境ベースの人工知能装置。

【請求項4】

前記プロセッサは、前記暗号処理フロント-エンドプロセッサの割り込み駆動のオフロード(`interrupt driven offloading`)を介して、前記第1種類のメモリ及び前記第2種類のメモリとデータ送受信を行うことを特徴とする請求項1に記載のトラスト環境ベースの人工知能装置。

【請求項5】

前記プロセッサは、DMA(`Direct Memory Access`)コントローラのDMA駆動のオフロードを介して、前記暗号処理フロント-エンドプロセッサと前記第1種類のメモリ及び前記第2種類のメモリとデータ送受信を行うことを特徴とする請求項1に記載のトラスト環境ベースの人工知能装置。

【請求項6】

前記プロセッサは、前記ニューラルネットワーク演算の実行中に前記暗号処理フロント-エンドプロセッサが前記暗号化入力データの復号化演算を行うようにして、前記ニューラルネットワーク演算を途切れることなく行うことを特徴とする請求項1に記載のトラスト環境ベースの人工知能装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、人工知能ニューラルネットワークフレームワーク技術に関し、より詳細には、トラスト(`Trust`)環境で人工ニューラルネットワークの実行を安全に加速化することができるトラスト環境ベースの人工知能装置に関する。

【背景技術】

【0002】

DNN(`Deep Neural Network`)は、入力層及び出力層とその間のいくつかの隠れ層とからなる人工ニューラルネットワークである。DNNは、モバイル及び組み込みアプリケーションで広く使用されている。特に、DNNは、ユーザの身元を検証するためにユーザの生物学的特性(例えば、指紋、虹彩、顔など)を使用する生体認証実行のアプリケーションに有用である。

【0003】

DNNには敏感なユーザデータが多く含まれているため、モバイル及び組み込み機器はセキュリティ攻撃からユーザ及びDNNデータを安全に保護することができる安全なDN

10

20

30

40

50

Nの実行環境を実現しなければならない。

【0004】

従来は、ARM (Advanced RISC Machine) プロセッサで使用されるハードウェアベースのセキュリティ技術であるトラストゾーン (Trust Zone) を介してトラスト実行環境 (Trusted Execution Environment) でDNNを実行することが提案されていた。トラストゾーンは、プロセッサ (processor) 内に独立したセキュリティ区域を別々に配置して重要な情報を保護するハードウェアベースのセキュリティ技術である。しかし、トラストゾーンでDNNを実行するだけではデータを完全に保護することはできない。トラストゾーンはメモリ保護機能が制限されているためである。トラストゾーンを使用すると、ハードウェア及びソフトウェアリソースを安全な一般区域に分割して、DNNの実行を他のプロセスから分離することができる。一般区域でセキュリティ区域のデータにアクセスすることは、ハードウェアによって防止される。しかし、トラストゾーンはデータを揮発性メモリに暗号化していない状態で維持するため、コールドブート攻撃 (Cold Boot Attack) のような物理的なセキュリティ攻撃は、トラストゾーンでDNNの実行にも拘らず、敏感なユーザ及びDNNデータを取得することができる。

10

【0005】

物理的攻撃からデータを保護するために、データを暗号化して安全なオンチップメモリでのみを復号化するように選択することができる。このような方式は、トラストゾーンを使用してDNNの実行を他のプロセスから隔離することができるだけでなく、暗号化を介して物理的攻撃からユーザ及びDNNデータを保護することができる。しかし、メモリで暗号化されたデータを保護すると、遅いメモリアクセスが大きく増加し、プロセッサに課される高いデータ暗復号化のオーバーヘッドにより、DNNの実行時間が大きく増加するという問題が発生する。

20

【0006】

従って、敏感なユーザ及びDNNデータを物理的攻撃から保護するだけでなく、遅いメモリアクセス及び高いデータ暗復号化のオーバーヘッドを克服し、DNNの実行時間を減らす新しいセキュリティDNNフレームワークが必要になった。

【先行技術文献】

【特許文献】

30

【0007】

【文献】韓国登録特許第10-2474875号公報

【発明の概要】

【発明が解決しようとする課題】

【0008】

本発明は、上記従来の問題点に鑑みてなされたものであって、本発明の目的は、トラスト (Trust) 環境で人工ニューラルネットワークの実行を安全に加速化することができるトラスト環境ベースの人工知能装置を提供することにある。

【課題を解決するための手段】

【0009】

40

上記目的を達成するためになされた本発明の一態様によるトラスト環境ベースの人工知能装置は、暗号化入力データを送信し、暗号化出力データを受信する第1タイプのメモリと、トラスト信頼空間で動作し、前記暗号化入力データ及び前記暗号化出力データの人工知能演算を行うトラスト人工知能処理部と、を備え、前記トラスト人工知能処理部は、前記暗号化入力データの復号化を介して復号化入力データを生成し、前記暗号化出力データを生成するために非暗号化出力データの暗号化を行う暗号処理フロント-エンドプロセッサと、前記復号化入力データ及び前記非暗号化出力データに対するバッファを提供する第2タイプのメモリと、前記復号化入力データに基づいてニューラルネットワーク演算を行って前記非暗号化出力データを生成するプロセッサと、を含む。

【0010】

50

前記暗号処理フロント - エンドプロセッサは、前記第 1 類型のメモリから前記暗号化入力データとして暗号化入力アクティベーション及び暗号化フィルタの入力を受けて、前記第 2 類型のメモリに復号化入力アクティベーション及び復号化フィルタを格納することができる。

前記暗号処理フロント - エンドプロセッサは、前記人工知能演算の過程で前記プロセッサによるオン - デマンド要求を受信し、前記第 1 類型のメモリにアクセスして前記暗号化入力データを取り込むことができる。

前記第 2 類型のメモリは、前記第 1 類型のメモリよりも相対的に速い動作速度及び少ない格納容量を有することができる。

前記プロセッサは、直接畳み込みベースのニューラルネットワーク演算を行って前記第 1 類型のメモリのアクセス数を減らすことができる。

10

前記プロセッサは、前記第 2 類型のメモリに復号化入力アクティベーションを固定的に格納し、復号化フィルタ及び非暗号出力アクティベーションを循環キュー方式で格納することができる。

前記プロセッサは、前記暗号処理フロント - エンドプロセッサの割り込み駆動のオフロード (interrupt driven offloading) を介して、前記第 1 類型のメモリ及び前記第 2 類型のメモリとデータ送受信を行うことができる。

前記プロセッサは、DMA (Direct Memory Access) コントローラの DMA 駆動のオフロードを介して、前記暗号処理フロント - エンドプロセッサと前記第 1 類型のメモリ及び前記第 2 類型のメモリとデータ送受信を行うことができる。

20

前記プロセッサは、前記暗号処理フロント - エンドプロセッサの暗号化及び復号化演算の実行にオーバーラップするように前記ニューラルネットワーク演算を実行してイントラ - レイヤパイプラインングを実現することができる。

前記プロセッサは、前記ニューラルネットワーク演算の実行中に前記暗号処理フロント - エンドプロセッサが前記暗号化入力データの復号化演算を行うようにして、前記ニューラルネットワーク演算を途切れることなく行うことができる。

前記プロセッサは、前記イントラ - レイヤパイプラインングのために、データ復号化ステップ、演算ステップ、及びデータ暗号化ステップに細分化して、前記ニューラルネットワーク演算を途切れることなく行うことができる。

【0011】

30

上記目的を達成するためになされた本発明の他の態様によるトラスト環境ベースの人工知能装置は、暗号化入力データを送信する第 1 類型のメモリと、トラスト信頼空間で動作し、前記暗号化入力データの人工知能演算を行うトラスト人工知能処理部と、を備え、前記トラスト人工知能処理部は、前記暗号化入力データの復号化を介して復号化入力データを生成する暗号処理フロント - エンドプロセッサと、前記復号化入力データ及び非暗号化出力データに対するバッファを提供する第 2 類型のメモリと、前記復号化入力データに基づいてニューラルネットワーク演算を行って前記非暗号化出力データを生成するプロセッサと、を含む。

【0012】

前記プロセッサは、直接畳み込みベースのニューラルネットワーク演算を行って前記第 1 類型のメモリのアクセス数を減らすことができる。

40

前記プロセッサは、前記第 2 類型のメモリに復号化入力アクティベーションを固定的に格納し、復号化フィルタ及び非暗号出力アクティベーションを循環キュー方式で格納することができる。

前記プロセッサは、前記暗号処理フロント - エンドプロセッサの暗号化及び復号化演算の実行にオーバーラップするように前記ニューラルネットワーク演算を実行してイントラ - レイヤパイプラインングを実現することができる。

【発明の効果】

【0013】

本明細書で開示する技術は、次の効果を有する。但し、特定の実施形態が次の効果を全

50

て含まなければならないとか、又は次の効果のみを含まなければならないという意味ではないため、開示する技術の権利範囲はこれによって制限されるものと理解してはならない。

【0014】

本発明のトラスト環境ベースの人工知能装置によれば、トラスト(Trust)環境で人工ニューラルネットワークの実行を安全に加速化することができる。

【0015】

また、トラスト信頼空間でデータの暗号化を行って物理的攻撃からセキュリティを強化し、直接畳み込みベースのニューラルネットワーク演算を行ってメモリのアクセス数を減らすことができる。

【0016】

更に、暗号ハードウェアでオフロードを行ってニューラルネットワークの実行作業が制限されたプロセッサリソースを活用することができ、イントラ-レイヤパイプラインングを介してデータの暗号化及び復号化にオーバーラップして人工ニューラルネットワークの実行時間を短縮することができる。

【図面の簡単な説明】

【0017】

【図1】本発明によるトラスト環境ベースの人工知能装置を説明する図である。

【図2】図1の人工知能装置でトラスト環境ベースの人工知能演算の実行方法を説明する図である。

【図3】トラスト実行環境を説明する図である。

【図4a】従来と比べた本発明によるDNNフレームワークの作業モデルを説明する図である。

【図4b】従来と比べた本発明によるDNNフレームワークの作業モデルを説明する図である。

【図5】従来のDNNフレームワークの実行速度、帯域幅、及び復号化のスループットをそれぞれ示す図である。

【図6】従来のDNNフレームワークの実行速度、帯域幅、及び復号化のスループットをそれぞれ示す図である。

【図7】従来のDNNフレームワークの実行速度、帯域幅、及び復号化のスループットをそれぞれ示す図である。

【図8a】ニューラルネットワーク演算の実行に使用される畳み込みを説明する図である。

【図8b】ニューラルネットワーク演算の実行に使用される畳み込みを説明する図である。

【図9a】DNNに優しいSRAM管理過程を説明する図である。

【図9b】DNNに優しいSRAM管理過程を説明する図である。

【図10】CPUと暗号ハードウェアとの間のスループットを示す図である。

【図11】暗号ハードウェアのオフロードを説明する図である。

【図12】ニューラルネットワーク演算の実行過程でイントラ-レイヤパイプラインングの実現を説明する図である。

【図13】ニューラルネットワーク演算の実行過程でイントラ-レイヤパイプラインングの実現を説明する図である。

【図14】本発明のソフトウェアの構成例を説明する図である。

【図15】本発明による実験結果を説明する図である。

【図16】本発明による実験結果を説明する図である。

【図17】本発明による実験結果を説明する図である。

【発明を実施するための形態】

【0018】

本発明は、下記の研究課題をもって支援を受けて出願された。

[この発明を支援した国家研究開発事業]

[課題固有番号] 1711193986

[課題番号] 2020-0-01361-004

10

20

30

40

50

[省庁名] 科学技術情報通信部

[課題管理 (専門) 機関名] 情報通信企画評価院

[研究事業名] 情報通信放送革新人材養成

[研究課題名] 人工知能大学院支援 (延世大学校)

[寄与率] 1 / 2

[課題実行機関名] 延世大学校産学協力団

[研究期間] 2 0 2 3 . 0 1 . 0 1 ~ 2 0 2 3 . 1 2 . 3 1

[この発明を支援した国家研究開発事業]

[課題固有番号] 1 4 1 5 1 7 7 6 5 9

[課題番号] P 0 0 1 6 1 5 0

10

[省庁名] 産業通商資源部

[課題管理 (専門) 機関名] 韓国産業技術振興院

[研究事業名] 産業技術国際協力 (R & D)

[研究課題名] スマート I o T ネットワークセキュリティのための攻撃探知

及び防御技術開発

[寄与率] 1 / 2

[課題実行機関名] 韓国電子技術研究院

[研究期間] 2 0 2 1 . 1 2 . 0 1 ~ 2 0 2 2 . 1 1 . 3 0

【 0 0 1 9 】

本発明に関する説明は、構造的又は機能的説明のための実施形態に過ぎないため、本発明の権利範囲は、本明細書に説明している実施形態によって制限されるものと解釈してはならない。即ち、実施形態は、様々な変更が可能であり、種々の形態を有することができるため、本発明の権利範囲は、技術的思想を実現することができる均等物を含むものと理解しなければならない。また、本発明で提示している目的又は効果は、特定の実施形態がこれを全て含むべきであるとか、そのような効果だけを含むべきであるという意味ではないため、本発明の権利範囲は、これによって制限されるものと理解してはならない。

20

【 0 0 2 0 】

一方、本明細書で述べる用語の意味は、次のように理解されなければならない。

【 0 0 2 1 】

「第 1」、「第 2」などの用語は、1つの構成要素を他の構成要素から区別するためのものであって、これらの用語によって権利範囲が限定されてはならない。例えば、第 1 構成要素は第 2 構成要素と名付けられてもよく、同様に第 2 構成要素も第 1 構成要素と名付けられてもよい。

30

【 0 0 2 2 】

ある構成要素が他の構成要素に「連結されて」いると言及する場合には、その他の構成要素に直接連結されることもあるが、中間に他の構成要素が存在することもあると理解しなければならない。それに対し、ある構成要素が他の構成要素に「直接連結されて」いると言及する場合には、中間に他の構成要素が存在しないものと理解しなければならない。一方、構成要素間の関係を説明する他の表現、即ち「~間に」と「すぐ~間に」、又は「~に隣り合う」と「~に直接隣り合う」なども同様に解釈されなければならない。

40

【 0 0 2 3 】

単数の表現は、文脈上明白に異なる意味ではない限り、複数の表現を含むものと理解しなければならない。「含む」又は「有する」などの用語は、実施される特徴、数字、ステップ、動作、構成要素、部分品、又はこれらを組み合わせたものが存在することを指定しようとするものであり、1つ又はそれ以上の他の特徴や数字、ステップ、動作、構成要素、部分品、又はこれらを組み合わせたものの存在又は付加可能性を予め排除しないものと理解しなければならない。

【 0 0 2 4 】

各ステップにおいて、識別符号 (例えば、a、b、c など) は、説明の便宜のために使用されるものであって、識別符号は、各ステップの順序を説明するものではなく、各ステ

50

ップは、文脈上明白に特定の順序を記載しない限り、明記された順序と異なって生じ得る。即ち、各ステップは、明記された順序と同一に生じることもあり、実質的に同時に行われることもあり、反対の順序で行われることもある。

【0025】

本発明は、コンピュータ読み取り可能な記録媒体にコンピュータ読み取り可能なコードとして実現され、コンピュータ読み取り可能な記録媒体は、コンピュータシステムによって読み取られるデータが格納される全ての種類の記録装置を含む。コンピュータ読み取り可能な記録媒体の例としては、ROM、RAM、CD-ROM、磁気テープ、フロッピー（登録商標）ディスク、光データ格納装置などがある。また、コンピュータ読み取り可能な記録媒体は、ネットワークで連結されたコンピュータシステムに分散されて、分散方式

10

【0026】

ここで使用する全ての用語は、異なって定義されない限り、本発明が属する分野における通常の知識を有する者により一般的に理解されるものと同じ意味を有する。一般的に使用される辞書に定義されている用語は、関連技術の文脈上有する意味と一致するものと解釈されなければならない。本明細書において明白に定義しない限り、理想的であるか又は過度に形式的な意味を有するものと解釈することはできない。

【0027】

以下、本発明を実施するための形態の具体例を、図面を参照しながら詳細に説明する。

【0028】

図1は、本発明によるトラスト環境ベースの人工知能装置を説明する図である。

20

【0029】

図1を参照すると、トラスト環境ベースの人工知能装置100は、モバイル又は組み込み機器で人工ニューラルネットワークの実行に対してセキュリティ強化及び加速化のためのフレームワークであり、第1タイプのメモリ110及びトラスト人工知能処理部130を含んで実現される。

【0030】

第1タイプのメモリ110は、暗号化入力データを送信し、暗号化出力データを受信する。ここで、第1タイプのメモリ110はDRAM(Dynamic Random Access Memory)で構成されるが、必ずしもこれに限定されるものではない。

30

【0031】

トラスト人工知能処理部130は、トラスト信頼空間で動作し、暗号化入力データ及び暗号化出力データの人工知能演算を行う。このために、トラスト人工知能処理部130は、暗号処理フロント-エンドプロセッサ131、第2タイプのメモリ133、及びプロセッサ135を含む。

【0032】

暗号処理フロント-エンドプロセッサ131は、暗号ハードウェア(Cryptographic Hardware)として暗号化入力データの復号化を介して復号化入力データを生成し、暗号化出力データを生成するために非暗号化出力データの暗号化を行う。暗号処理フロント-エンドプロセッサ131は、第1タイプのメモリ110から、暗号化入力データである暗号化入力アクティベーション及び暗号化フィルタの入力を受けて、第2タイプのメモリ133に復号化入力アクティベーション及び復号化フィルタを格納する。暗号処理フロント-エンドプロセッサ131は、人工知能演算の過程でプロセッサ135によるオン-デマンド(On-Demand)要求を受信し、第1タイプのメモリ110にアクセスして暗号化入力データを取り込む。

40

【0033】

第2タイプのメモリ133は、復号化入力データ及び非暗号化出力データに対するバッファを提供する。ここで、第2タイプのメモリ133は、第1タイプのメモリ110よりも相対的に速い動作速度及び少ない格納容量を有する。例えば、第1タイプのメモリ110がDRAMで構成された場合、第2タイプのメモリ133はSRAM(Static Random

50

Access Memory)で構成される。第1タイプのメモリ110と第2タイプのメモリ133との間でデータを送受信するとき、暗号処理フロント-エンドプロセッサ131でデータを暗号化及び復号化して、正確性及びセキュリティ性を保証する。

【0034】

プロセッサ135は復号化入力データに基づいてニューラルネットワーク演算を行って非暗号化出力データを生成する。プロセッサ135は直接畳み込みベースのニューラルネットワーク演算を行って第1タイプのメモリ110のアクセス数を減らす。即ち、プロセッサ135は、直接畳み込みを使用してニューラルネットワーク演算の作業セットサイズを最小化し、減少した作業セットサイズにより、ニューラルネットワーク演算を行っている間、第2タイプのメモリ133の動作速度よりも相対的に遅い第1タイプのメモリ110に対するアクセス及びデータ暗号化及び復号化の呼び出しを減らす。

10

【0035】

プロセッサ135は、第2タイプのメモリ133に復号化入力アクティベーションを固定的に格納し、復号化フィルタ及び非暗号出力アクティベーションを循環キュー方式で格納する。プロセッサ135は、ニューラルネットワーク演算の実行で畳み込みレイヤの入力アクティベーション、フィルタ、及び出力アクティベーションのアクセスパターンに基づいて、第2タイプのメモリ133を効率的に管理する。プロセッサ135は、入力アクティベーションの場合、データの再使用を増加させる傾向があるため固定的に格納し、フィルタの場合、出力チャンネルを生成するとき同じフィルタを複数回使用するが、サイズが入力アクティベーションよりもはるかに小さく、1つの出力チャンネルのみが生成するのに必要であるため、循環キュー(Circular Queue)方式で格納し、出力アクティベーションの場合、空間的集約性の高い1回の書き込みデータであるため、循環キュー方式で格納する。

20

【0036】

プロセッサ135は暗号処理フロント-エンドプロセッサ131の割り込み駆動のオフロード(interrupt driven offloading)を介して第1タイプのメモリ110及び第2タイプのメモリ133とデータ送受信を行う。また、プロセッサ135は、DMA(Direct Memory Access)コントローラのDMA駆動のオフロードを介して、暗号処理フロント-エンドプロセッサ131並びに第1タイプのメモリ110及び第2タイプのメモリ133とデータ送受信を行う。プロセッサ135は割り込み駆動のオフロード及びDMA駆動のオフロードの中から選択される。割り込み駆動のオフロードは待ち時間が重要な作業を行うのに適しており、DMA駆動のオフロードは一度により多くの量のデータを処理するのに適している。ここで、両オフロードのメカニズムは、ニューラルネットワーク演算を加速化する。

30

【0037】

プロセッサ135は、暗号処理フロント-エンドプロセッサ131の暗号化及び復号化演算の実行にオーバーラップするようにニューラルネットワーク演算を実行して、イントラ-レイヤパイプラインングを実現する。プロセッサ135は、ニューラルネットワーク演算の実行中に暗号処理フロント-エンドプロセッサ131が暗号化入力データの復号化演算を行うようにして、ニューラルネットワーク演算を途切れることなく行う。プロセッサ135は、イントラ-レイヤパイプラインングのために、データ復号化ステップ、演算ステップ、及びデータ暗号化ステップに細分化して、ニューラルネットワーク演算を途切れることなく行う。プロセッサ135は、イントラ-レイヤパイプラインングを介して、データ復号化ステップ、演算ステップ、及びデータ暗号化ステップを並列化する。

40

【0038】

図2は、図1の人工知能装置でトラスト環境ベースの人工知能演算の実行方法を説明するフローチャートである。

【0039】

図2において、人工知能装置100は、第1タイプのメモリ110から暗号化入力データの入力を受けて、復号化を介して復号化入力データを生成する(ステップS210)。人

50

人工知能装置 100 は暗号処理フロント - エンドプロセッサ 131 を介して暗号化入力データに対して復号化を行う。人工知能装置 100 は復号化入力データを第 2 種類のメモリ 133 に格納する。

【0040】

人工知能装置 100 は、復号化入力データに基づいてニューラルネットワーク演算を行って非暗号化出力データを生成する（ステップ S230）。ここで、人工知能装置 100 はプロセッサ 135 を介して直接畳み込みベースのニューラルネットワーク演算を行って暗号処理フロント - エンドプロセッサ 131 が第 1 種類のメモリ 110 にアクセスして暗号化入力データを取り込むアクセス数を減らす。人工知能装置 100 は、非暗号化出力データを第 2 種類のメモリ 133 に格納する。

10

【0041】

人工知能装置 100 は、非暗号化出力データの暗号化を行って暗号化出力データを生成する（ステップ S250）。人工知能装置 100 は、暗号処理フロント - エンドプロセッサ 310 を介して非暗号化出力データに対して暗号化を行う。人工知能装置 100 は暗号化出力データを第 1 種類のメモリ 110 に出力する。

【0042】

人工知能装置 100 は、暗号処理フロント - エンドプロセッサ 131 の暗号化及び復号化演算の実行にオーバーラップするようにプロセッサ 135 のニューラルネットワーク演算を実行して、データ復号化ステップ、演算ステップ、及びデータ暗号化ステップでイントラ - レイヤパイプラインングを実現することによって、ニューラルネットワーク演算を加速化する。

20

【0043】

以下、図 3 ~ 17 を参照して、本発明によるトラスト環境ベースの人工知能装置についてより詳しく説明する。

【0044】

図 3 は、トラスト実行環境を説明する図である。

【0045】

図 3 を参照すると、トラスト実行環境 (Trusted Execution Environment: 以下、TEE という) は、処理、メモリ、及び格納機能を備えたセキュリティ処理環境であり、敏感な作業及びデータが当該環境を外れないように制限し、高いセキュリティを達成する。TEE 内の作業及びデータは、OS (Operating System) 及びアプリケーションが実行されるリッチ実行環境 (Rich Execution Environment: 以下、REE という) で隔離される。

30

【0046】

モバイル及び組み込み装置で使用されるトラストゾーンサポートの TEE を実現するために、トラストゾーンは CPU が指定された時点で TEE 及び REE のうちの 1 つでのみ独占的に作動するようにするセキュリティプロセッサモードを実現する。TEE と REE との間の切替え及び相互作用は、SMC (Secure Monitor Calls) で呼び出せるセキュリティモニターにより管理される。また、トラストゾーンは、DRAM をセキュリティ及び非セキュリティ領域に分割し、TEE の敏感なデータを保護するために REE でセキュリティ領域にアクセスすることを許容しない。

40

【0047】

トラストゾーンサポートの TEE で実行される既存の DNN フレームワークは、1) 周辺装置から入力データ (例えば、指紋センサの指紋イメージ) を受信して REE から開始される。2) REE で DNN のいくつかの初期レイヤを実行した後、3) 出力アクティベーションを暗号化し、SMC を介して TEE に送信する。その次に、4) TEE 内で送信されたアクティベーションを復号化し、5) 復号化されたアクティベーション及び事前送信されたフィルタを使用して残りのレイヤを実行する。残りのレイヤの実行を完了すると、6) DNN が作った予測を、SMC を介して REE に返還する。

【0048】

50

このような方式でDNNの実行を隔離して、いくつかのセキュリティ攻撃からDNNを保護する。

【0049】

図4a及び4bは、従来と比べた発明によるDNNフレームワークの作業モデルを説明する図であって、図4aは従来のDNNフレームワークであり、図4bは、本発明で提案するDNNフレームワーク（以下、GuardianNNという）である。

【0050】

図4aの従来のDNNフレームワークの場合、最初にDNNの入力データ及びフィルタが全て暗号化されてDRAMに格納される。フレームワークがDNNのレイヤの実行を開始すると、1)暗号化された入力アクティベーション及びフィルタをSRAMにロードし、2)CPUを使用してデータを復号化する。その次に、フレームワークは、3)復号化されたデータを使用して各レイヤを実行し、出力アクティベーションをSRAMに格納する。その後、フレームワークは、4)出力アクティベーションを暗号化し、5)これをDRAMに格納する。入力アクティベーション及びフィルタは必要に応じてSRAMに置き換えられる。SRAMがレイヤの作業集合よりも小さい場合、データスワッピング及び暗復号化がよく発生する。暗号化されたDRAMを使用してTEE内でDNNを実行すると、従来のDNNフレームワークがセキュリティ攻撃から敏感なユーザ及びDNNデータを完全に保護できるが、遅いDNNの実行で困難をきたすことができる。

10

【0051】

図5乃至図7は、従来のDNNフレームワークの実行速度、帯域幅、及び復号化のスループットをそれぞれ示す図である。

20

【0052】

図5に示すように、図4aの従来のDNNフレームワーク(OP-TEE with Pager)とDRAMデータを暗号化しない安全ではないDarknetZ作業モデルとのDNNの実行速度を比較した結果であって、従来のDNNフレームワークのDNNの実行速度がDarknetZよりも顕著に遅いことが示された。従来のDNNフレームワークはDarknetZよりもAlexNetを実行するのに3.42倍長くかかった。

【0053】

従来のDNNフレームワークは、組み込みSRAMの制限された容量により遅いDRAMのアクセス数が増加し、CPUに課される高いデータの暗号化及び復号化のオーバーヘッドによる性能のボトルネック現象がDNNの実行速度に影響を与える。具体的に、従来のDNNフレームワークは、TEEが暗号化されたDRAMデータを安全にロードし、ロードされたデータを復号化する安全なオンチップメモリとして組み込みSRAMを活用する。しかし、SRAMの容量(数百KB)は、一般的にモバイル及び組み込み装置のオンチップCPUキャッシュ(数百KB)よりも小さいため、SRAMを安全なオンチップバッファとして使用するが、効果的なオンチップメモリのサイズを一桁に下げる。これにより、DNNを実行するときに、遅いDRAMのアクセス数が増加して、DNNの実行速度が遅くなる。

30

【0054】

図6に示すように、作業セットサイズが多様なセキュリティ内蔵型SRAMとオフチップDRAMとの帯域幅を比較した結果、SRAMが多様な作業セットサイズに対して持続的にDRAMよりも高い帯域幅を提供することが示された。これは、増加したDRAMアクセスがDNNの実行速度に否定的な影響を与えることを意味する。結果は、またより高い帯域幅を達成するためにシーケンシャルアクセスが任意のアクセスよりも好まなければならないということを示唆する。結果として、遅いDRAMアクセスがDNNの実行速度に及ぼす否定的な影響を最小化するためには、速いSRAMデータの再使用を最大化して、DNNの実行のメモリアccessをシーケンシャルアクセスで構成しなければならない。

40

【0055】

組み込みSRAMと暗号化されたDRAMとの間でデータを交換するときに、機能的正確性及び高いセキュリティ性を保証するために、CPUでデータを暗号化及び復号化しな

50

なければならない。しかし、CPUベースの暗号化は速度が遅いだけでなく、制限されたCPU帯域幅のかなりの部分を消費するため、制限されたCPU帯域幅がコンピューティング集約的なDNNの実行と暗号化及び復号化とで全て共有され、DNNの実行速度が遅くなる。

【0056】

図7に示すように、ペイロードサイズが多様なCPUベースのデータ復号化のスループットは、データの暗号化及び復号化が制限されたCPUリソースによりも困難をきたしており、2KBのペイロードで4.51MB/sのスループットを達成することを示す。2KBのDRAMデータにシーケンシャルアクセスするスループット(609.24MB/s、図6の(a))と比較した場合、CPUベースの2KBデータの暗号化及び復号化はシーケンシャルDRAMアクセスよりも135.09倍遅い。低いデータの暗号化及び復号化のスループットは、データの暗号化及び復号化のオーバーヘッドが相当であり、従来のDNNフレームワークが速いDNNの実行を達成することは難しいということを示す。結果として、従来のDNNフレームワークは制限されたCPUリソースに課される高いDRAMデータの暗号化及び復号化のオーバーヘッドによりかなりの困難をきたしており、速く安全なDNNの実行のためには、オーバーヘッドを解決しなければならない。

10

【0057】

従って、本発明は、従来のDNNフレームワークの遅いDNNの実行問題を解決し、速く安全なモバイル及び組み込み装置用DNNフレームワークであるGuardiANNを提案する。本発明で提案したGuardiANNは、次のような特徴を有する。

20

【0058】

- 直接畳み込みを介してDNN実行の作業セットサイズを最小化することができる。作業セットサイズを減少させて、DNNの実行中の、遅いDRAMアクセス及びデータの暗号化及び復号化要求を減らすことができる。

【0059】

- DNNに優しいSRAM管理を使用し、畳み込みレイヤのデータの再使用を最大化することができる。入力アクティベーションをSRAMに固定し、フィルタ及び出力アクティベーションを循環キュー方式で格納して、SRAMを効率的に管理することができる。

【0060】

- 暗号ハードウェアでデータの暗号化及び復号化をオフロードすることができる。データの暗号化及び復号化をオフロードすると、制限されたCPUリソースをDNN専用を使用することができる。

30

【0061】

- DNNレイヤとデータの暗号化及び復号化作業とをオーバーラップして、DNNの実行を更に加速化することができる。これは、CPUと暗号ハードウェアとがそれぞれDNNとデータの暗号化及び復号化作業とを同時に行うことができるため可能である。

【0062】

図4bを見ると、本発明で提案したDNNフレームワークの作業モデルは、DNNレイヤを実行するときに、1) モバイル及び組み込み装置で使用される暗号ハードウェア(図1の暗号処理フロント-エンドプロセッサ)によりDRAMで暗号化された入力アクティベーション及びフィルタをロードしてデータを復号化し、復号化データをSRAMに格納する。その次に、2) CPU(図1のプロセッサ)は、SRAMデータを使用してDNNレイヤの作業を行い、出力アクティベーションをSRAM(図1の第2種類のメモリ)に格納する。その後、3) 暗号ハードウェアは、SRAMに格納された出力アクティベーションを暗号化し、暗号化された出力アクティベーションをDRAM(図1の第1種類のメモリ)に格納する。

40

【0063】

図4aの従来のDNNフレームワークの作業モデルと比較すると、直接畳み込み及びDNNに優しいSRAM管理を使用し、遅いDRAMアクセス及び暗号ハードウェアを使用して、CPUに課される高いデータの暗号化及び復号化のオーバーヘッドを大幅に減らす

50

ことができる。本発明で提案した主要な特徴がDNNの実行に与える影響は下記表1にまとめることができる。

【0064】

【表1】

Key Idea	Beneficial to DNNs Having
Direct Convolutions	Large filter sizes, high input channel counts, small stride sizes
DNN-Friendly SRAM Mgmt.	High output channel counts
Cryptographic Hardware	Large memory footprint sizes
Intra-Layer Pipelining	Balanced compute & encryption/decryption

10

【0065】

本発明は、直接畳み込みを使用して畳み込みレイヤの作業セットサイズを減らし、DNNに優しいSRAM管理を使用してSRAMデータの再使用を最大化し、DNNの実行中の遅いDRAMのアクセスを大幅に減らすことができる。

【0066】

図8a及び8bは、ニューラルネットワーク演算の実行に使用される畳み込みを説明する図であって、図8aはim2col(Image to Column)の畳み込みであり、図8bは直接畳み込みである。

20

【0067】

従来のモバイル及び組み込みDNNフレームワーク(例えば、DarknetZ、TensorFlow Lite)は、im2col(Image to Column)畳み込みを使用して、時間が多くかかる畳み込みレイヤを実行する。im2colは、多次元データを行列に変換して、行列演算を行うようにする関数をいう。多次元データの畳み込みは、im2colを介して行列に変換されたデータの内積と同じである。im2colの畳み込みは、図8aのように各パッチ(即ち、要素がフィルタの要素にマッピングされる入力アクティベーションの集合)を2次元行列に平面化し、平面化されたパッチとフィルタとの間の行列乗算を行って速い畳み込みレイヤの実行を達成する。しかし、パッチを併合すると、併合されたパッチを格納するための追加のバッファを割り当てなければならないため、割り込みレイヤの作業セットサイズが大幅に増加する。SRAMの制限された容量(最大数百KB)により作業セットサイズが増加すると、DNNの実行速度が大幅に遅くなる遅いDRAMのアクセスが多く発生する。従って、速いDNNの実行を達成するために、畳み込みレイヤの作業セットサイズを最小化しなければならない。

30

【0068】

畳み込みレイヤの作業セットサイズを最小化するために、本発明ではim2col畳み込みの代わりに直接畳み込みを使用する。直接畳み込みは、入力に畳み込みレイヤのフィルタを相関演算して、フィルタをスライディングウィンドウ方式で全ての領域で相関値を計算した結果値が出力となる。直接畳み込みは、図8bのように必要な入力アクティベーションを要求する際に取り込むため、作業セットサイズを大きくしない。減少した作業セットサイズ以外にも、直接畳み込みは、入力アクティベーションの固有の時間及び空間的地域性により、SRAMデータの再使用を増加させる傾向がある。直接畳み込みが入力アクティベーションにかけてフィルタを空間的にスライドして出力アクティベーションのチャンネルを生成するため、最近アクセスした入力アクティベーションに隣接する入力アクティベーションが近い将来にアクセスする可能性が高い。フィルタの場合、出力チャンネルを生成するときと同じフィルタを複数回使用するため、フィルタは時間的地域性が高い。需要ページング内蔵のSRAMに結合された直接畳み込みの高い空間的及び時間的集約性は、遅いDRAMアクセスを大幅に減らすのに役立つ。

40

【0069】

図9a及び9bは、DNNに優しいSRAM管理過程を説明する図である。

50

【 0 0 7 0 】

図 9 a に示すように、ARM トラストゾーン技術が適用された既存のオープンソース TEE のフレームワークである Pager は、デマンドページング (Demand Paging) を使用してセキュリティ組み込み S R A M の制限された容量を管理する。デマンドページングは、必要に応じて暗号化された D R A M から S R A M にデータをロードし、最も近くに使用されていない S R A M データを D R A M にプッシュして S R A M 空間を回収する。デマンドページングはデータの時間的地域性を効率的に活用するが、畳み込みレイヤの入力アクティベーション、フィルタ、及び出力アクティベーションの様々なアクセスパターンを認識できない。また、S R A M のサイズは、モバイル及び組み込み装置で数百 K B に過ぎないため、畳み込みレイヤに必要な全てのデータを S R A M に格納することはできない。このような非効率的な S R A M 管理は、遅い D R A M のアクセス数を大幅に増加させ、速い D N N の実行を妨げる。従って、S R A M が多数の遅い D R A M アクセスを誘発しないように、畳み込みレイヤの様々なデータアクセスパターンを統合して、小さい S R A M を効率的に管理しなければならない。

10

【 0 0 7 1 】

図 9 b に示すように、本発明は、S R A M データの再使用を最大化するために、D N N に優しい S R A M 管理を実現する。ここで、D N N に優しい S R A M 管理は、畳み込みレイヤの入力アクティベーション、フィルタ、及び出力アクティベーションの様々なアクセスパターンを活用する。本発明は、畳み込みレイヤの実行全体で繰り返しアクセスするときに、先ず全ての入力アクティベーションを S R A M に固定する。出力チャンネルを生成するためにアクセスされる入力アクティベーションは、他の全ての出力チャンネルに対して再度アクセスされなければならない。その次に、残りの S R A M 空間を 2 つの循環キューで構成する。1 つはフィルタ用であり、もう 1 つは出力アクティベーション用である。フィルタは入力アクティベーションのようによくアクセスするデータである。しかし、サイズが入力アクティベーションよりもはるかに小さく、1 つの出力チャンネルのみを生成するのにフィルタが必要であるため、フィルタを循環キューに格納するだけでも高い時間的地域性を十分活用することができる。出力アクティベーションは、空間的集約性の高い 1 回の書き込みデータであるため、循環キューに適している。D N N に優しい方式で S R A M を管理することによって、S R A M データの再使用を最大化して遅い D R A M アクセスを最小化し、D N N の実行を加速化することができる。

20

30

【 0 0 7 2 】

本発明は、暗号ハードウェアを使用してデータの暗号化及び復号化を行うようにして、制限された C P U リソースを D N N の実行に完全に使用することができる。これにより、C P U リソースを D N N の実行に完全に専用するだけでなく、暗号ハードウェアの高性能を活用して速い D N N の実行を達成することができる。データの暗号化及び復号化は、S R A M が暗号化された D R A M でデータをロード及び格納して、制限された C P U リソースのかなりの量を消費する度に発生する。データの暗号化及び復号化の高いオーバーヘッドを克服するために、暗号ハードウェアを使用して暗号化及び復号化作業を行う。オーバーヘッドの高いデータの暗号化及び復号化を暗号ハードウェアにオフロードし、制限された C P U リソースを完全に専用して、D N N の実行を加速化することができる (図 1 0 参照) 。

40

【 0 0 7 3 】

図 1 1 は、暗号ハードウェアのオフロードを説明する図であって、(a) は割り込みベースのオフロードであり、(b) は D M A (D i r e c t M e m o r y A c c e s s) 駆動のオフロードである。

【 0 0 7 4 】

本発明は、図 1 1 の (a) 割り込みベースのオフロード、及び (b) D M A 駆動のオフロードの 2 つのハードウェアのオフロードメカニズムから選択される。2 つのオフロードメカニズムは、互いに異なる周辺装置の使用シナリオで加速化される。割り込みベースのオフロードは、周辺装置が与えられた作業処理を完了するとすぐに割り込みを生成するた

50

め、待ち時間が重要な作業を周辺装置で行うのに適している。これに対し、DMA駆動のオフロードは、一度により多くの量のデータを処理するのに適している。DMAは、CPUが介入することなくデータを処理するための全てのデータ転送及び周辺の呼び出しを処理し、全てのデータが処理された後に割り込みを発生させる。DNNに優しいSRAM管理は、SRAMとDRAMとの間で大量のデータ転送（一般に数十KB）を要求するため、割り込みベースのオフロードよりも大きなペイロードサイズでより高いスループットを達成するDMAベースのオフロードが有利なことがある。

【0075】

暗号ハードウェアを使用すると、制限されたCPUリソースをDNNの実行に完全に専らし、DNNの実行を加速化することができ、更なる性能の最適化が可能である。CPU及び暗号ハードウェアでそれぞれ実行される重複DNNの実行及びデータの暗号化及び復号化である。畳み込みレイヤの1つの属性は、互いに異なる出力チャンネルを生成する作業が互いに独立していることである。畳み込みレイヤは1つのフィルタを使用して1つの出力チャンネルを生成し、入力アクティベーションは全ての出力チャンネル間で読み取り専用データとして共有される。この属性はプーリングレイヤにも適用される。プーリングレイヤの各出力チャンネルは、当該入力チャンネルのみを要求するため、出力チャンネルの作業をデータ並列にする。これに基づいて、畳み込みレイヤの大量の出力チャンネルの実行をデータの復号化ステップ、演算ステップ、及びデータの暗号化ステップの3つのパイプラインステップに分割する。その次に、出力チャンネルの互いに異なるバルクの3つのステップをパイプラインして、より速いDNNの実行を達成する。これをイントラ-レイヤパイプラインングという。これは、畳み込みレイヤの様々な出力チャンネル作業をパイプラインで接続するためである。

【0076】

図12は、ニューラルネットワーク演算の実行過程でイントラ-レイヤパイプラインングの実現を説明する図であり、256個の出力チャンネル（SRAMの制限された容量により大量の32個の出力チャンネル含む）を生成し、AlexNetの実行待ち時間に最も大きく寄与するAlexNetの6番目のレイヤにイントラ-レイヤパイプラインングを適用する場合の利点を示す。

【0077】

図12において、データの復号化ステップは、先ず大量の出力チャンネルを生成するのに必要な入力アクティベーション及び/又はフィルタを復号化してSRAMにロードする。その次に、演算ステップでは入力アクティベーション及びフィルタを使用して大量の出力チャンネルを演算する。その後、データの暗号化ステップで大量の出力チャンネルを暗号化してDRAMに格納する。イントラ-レイヤパイプラインングがない場合、(a)のように3つのステップが直列化され、レイヤを実行するのに57msがかかる。これに対し、出力チャンネルの互いに異なるバルクの3つのステップをパイプラインするイントラ-レイヤパイプラインングは、(b)のように3つのステップが並列化されてレイヤを実行するのに(a)の場合よりも24.6%更に速い43msしかかからない。この例は、イントラ-レイヤパイプラインングがより速いDNNの実行を達成することができることを示す。図13は、畳み込みレイヤにイントラ-レイヤパイプラインングを適用するための擬似コードを示すアルゴリズムである。

【0078】

しかし、畳み込みレイヤにイントラ-レイヤパイプラインングを適用すると、パイプラインを使用しないレイヤの実行よりもSRAMの容量がより多く消費される。イントラ-レイヤパイプラインングの機能的正確性を保証するためには、より大きな容量が必要である。大量の出力チャンネルの演算ステップと次の大量の出力チャンネルのデータ復号化ステップとをオーバーラップするためには、2つの大量に対するSRAMバッファが同時に割り当てられなければならない。より大きなSRAMバッファが必要な場合、レイヤの出力チャンネルはより多くのバルクにグループ化されるため、DNNの実行速度が遅くなる可能性がある。しかし、より大きなSRAMバッファが必要であるにも拘らず、イントラ

10

20

30

40

50

- レイヤパイプラインの性能の利点は、バルク当たりの出力チャンネル数が少なく、潜在的な性能低下を超えることにある。従って、基本的に畳み込み及びプーリングレイヤに対するイントラ - レイヤパイプラインをアクティベーションする。しかし、非常に小さな畳み込みSRAMが装着されたモバイル及び畳み込み装置の場合、潜在的な性能低下を防止するために、イントラ - レイヤパイプラインを非アクティベーションすることができる。

【0079】

図14は、セキュリティ畳み込みSRAM及び暗号ハードウェアが装着されたモバイル及び畳み込み装置上に本発明を実現するのに必要なソフトウェアの構成例を説明する図である。

【0080】

図14において、最も広く使用されるTEE実現のうちの1つであるオープンソースTrustZoneベースの実現であるOP-TEE上に本発明を実現しているが、必ずしもこれに限定されるものではなく、全てのTrustZoneベースのTEE実現上に実現され得る。本発明の核心には、TEE内の与えられた入力に対する予測を生成するためにDNNのレイヤを実行する信頼できるアプリケーションであるGuardiaNNランタイムがある。DNNの実行はREEから開始される。まず、REEは説明(例えば、レイヤ数、レイヤ当たりの入力及び出力アクティベーションサイズ)、暗号化入力データ、及びDNNのフィルタをTEEメモリに送信し、GuardiaNNランタイムを呼び出す。その次に、GuardiaNNランタイムはそれぞれ異なる種類のDNNのレイヤを実現するヘルパー関数を使用してDNNの各レイヤを実行する。その後、GuardiaNNランタイムはDNNを実行して作られた予測をREEに再度返還する。

【0081】

直接畳み込みは、GuardiaNNランタイム内でのみ実現することができ、GuardiaNNランタイムが信頼できるOS(例えば、OP-TEE)と相互作用してSRAMを割り当て、暗号ハードウェア及びDMAを使用しなければならない。DNNに優しいSRAM管理及びイントラ - レイヤパイプラインを実現するためには、GuardiaNNランタイムがDRAMに置き換えられないバッファをSRAMに割り当てなければならない。このために、TEEメモリの割り当てのためのTEE Internal Core API関数であるTEE_Malloc()を拡張し、isSRAMという追加入力因数を使用する。isSRAMの値がtrueである場合、信頼できるOSはSRAMバッファを割り当て、バッファがDRAMにスワップアウトされることを防止する。また、信頼できるOS(例えば、Pager)のメモリ管理を拡張し、isSRAMをDRAMにtrueに設定した状態でTEE_Malloc()を呼び出し、割り当てられたSRAMバッファを置き換えることを防止する。isSRAMの基本値はfalseに設定され、isSRAMを認識できない既存の信頼できるアプリケーションの機能的正確性を保証する。例えば、TEE_Malloc(1024, hint, true)を呼び出すと、DRAMに除去されない1KB SRAMバッファが割り当てられる。ここで、ヒントはバッファの特性に関するいくつかのヒントを提供する(例えば、0で満たされる)。

【0082】

DMA駆動のデータ暗号化及び復号化のオフロードの場合、GuardiaNNランタイムは信頼できるOSのDMA装置ドライバで定義されたユーザ指定システム呼び出しを呼び出す。GuardiaNNの実現は、信頼できるOSを拡張して、EncryptData()及びDecryptData()という2つのユーザ指定システム呼び出しを提供する。EncryptData()システム呼び出しは、暗号コンテキスト(暗号類型、キーサイズ等含む)、SRAM開始アドレス、DRAM開始アドレス、及びデータサイズを入力として使用する。それから、システム呼び出しは、CPUキャッシュをフラッシュしてダーティキャッシュラインをSRAMから除去し、SRAM開始アドレスで暗号化されていないSRAMデータを読み取り、暗号化コンテキスト及び暗号ハードウェアを

10

20

30

40

50

使用してデータを暗号化し、暗号化されたデータをDRAM開始アドレスでDRAMに格納する。類似の方式で、DecryPtData()システム呼び出しは、暗号コンテキスト、DRAM開始アドレス、SRAM開始アドレス、及びデータサイズの四つの因数を入力として使用する。その次に、EncryptData()システム呼び出しと類似する手順によって暗号化されたDRAMデータを読み取り、データを復号化し、復号化されたデータをSRAMに格納する。

【0083】

拡張されたTEE_Malloc()API機能及び既存のGlobalPlatform APIと共に2つのユーザ指定システム呼び出しを使用すると、本発明をモバイル及び組み込み装置上に充実に実現することができる。

10

【0084】

評価(EVALUATION)

【0085】

<実験セットアップ(Experimental Setup)>

【0086】

速く安全なDNNの実行に対するGuardiaNNの効果を評価するために、STM32MP157C-DK2開発ボード上にGuardiaNNをプロトタイプにして、DNNの実行速度及びエネルギー消費を基本セキュリティDNNフレームワークと比較した。開発ボードは、オープンソーストラストゾーン(TrustZone)ベースのTEE実現であるOP-TEEで公式的にサポートし、最新の組み込み装置の一般的なハードウェアの構成を反映した。デュアルコアARM Cortex-A7 CPU、256KBのセキュリティ組み込みSRAM、暗号ハードウェア、及び512MB DDR3L DRAMで構成される。信頼できるOSとしてGlobalPlatform TEE Client API v1.1及びInternal Core API v1.0をサポートするPagerと共にOP-TEE v3.11.0を使用する。GuardiaNN及び基本フレームワークの実現は、OP-TEEが現在単一のTEEインスタンス内でマルチスレッドをサポートしないため、1つのCPUコアのみ使用する。ARM NEON単一コマンドの多重データコマンドでDarknetZのDNNレイヤの実現を拡張し、各DNNレイヤ作業間のデータ並列性を活用した。拡張DNNレイヤの実現はGuardiaNN及び基本フレームワークに全て適用される。信頼できるOSが全てのSRAM(TEEとREEとの間の共有メモリでOP-TEEが予約した4KBを除外)をGuardiaNN及び基本フレームワークに全て割り当てると仮定する。エネルギー消費の比較のために、Monsoon HVPM(High Voltage Power Monitor)を使用して、装置全体のエネルギー消費を測定する。

20

30

【0087】

ベンチマークとして8ビットの整数量子化を使用して量子化された8個のDNNを選択し、敏感なユーザ及びDNNデータに関連する5個の代表的なモバイル及び組み込みアプリケーションのドメインを取り扱う。5個のドメインは、イメージ分類、顔認識、指紋認識、視線追跡、及び感情認識である。各ドメインに対して多様なDNNが存在する。しかし、ここで、実行するのに合理的に短い実行待ち時間を有する各ドメインで代表的な軽量のDNNを選択する。例えば、STM32MP157C-DK2開発ボードで基本フレームワークがDNNを実行するのに192秒がかかったため、イメージ分類のためのDNNであるResNet-18をベンチマークに含めない。下記の表2には8個のDNN及びその特性が並べられている。

40

【0088】

50

【表 2】

Domain	Name	# Layers	Input Size
Image classification	AlexNet [35]	11	3x32x32
	MobileNetV1 [25]	29	3x32x32
	Tiny Darknet [63]	21	3x32x32
	VGG-7 [37]	11	3x32x32
Face recognition	DeepID [69]	9	3x31x31
Fingerprint recognition	FgptAuth [66]	11	1x128x128
Gaze tracking	GAZEL [58]	10	1x64x64
Emotion recognition	Smart Doll [11]	11	3x50x50

10

【0089】

< 速いDNNの実行 (Fast DNN Execution) >

【0090】

まず、選択された全てのDNNの実行待ち時間を測定し、GuardiaNNのDNNの実行速度を評価する。提案された技術の寄与度を分析するために、基本フレームワークから開始して、各提案された技術を漸進的に適用し、DNNの実行待ち時間を測定する。直接畳み込み及びDNNに優しいSRAM管理、暗号ハードウェア、及びイントラ-レイヤパイプラインに対する作業をすると、図15に5個の棒で表示された計5個の構成が提供される。実験でイントラ-レイヤパイプラインのために8個の出力チャンネルのバルクサイズを使用する。

20

【0091】

図15の(a)は待ち時間の実験結果を示す。殆どのDNNの実行は、基本フレームワークで極めて遅く、最大11.4秒の待ち時間が発生する。GuardiaNNは全てのDNNで待ち時間を1秒未満に減らす。図15の(b)は、直接畳み込み及びDNNに優しいSRAM管理技術(三番目の棒)が適用された構成で正規化されたGuardiaNNが提供する相対的速度の向上を示す。それぞれの提案された技術が効果的であることを観察することができる。直接畳み込みを適用すると、im2colを使用する基本フレームワークに比べて幾何平均速度が2.58倍向上する。DNNに優しいSRAM管理は3.19倍の追加の幾何平均速度の向上をもたらす。暗号化及び復号化を暗号ハードウェアでオフロードすると、1.73倍の追加の幾何平均速度の向上が提供される。イントラ-レイヤパイプラインを適用すると、これを更に改善し、幾何平均速度を1.07倍高めることができる。GuardiaNNはベースラインに比べて15.3倍の幾何平均速度の向上を達成する。ここで評価する8個のDNNのうち、GuardiaNNはSmartDollを最も加速化し、ベースラインよりも31.4倍更に速くする。このような速度向上の殆どは、DNNに優しいSRAM管理を直接畳み込みと共に適用した結果、遅いDRAMのアクセス数が減少したためである。AlexNetはGuardiaNNの暗号ハードウェア及びイントラ-レイヤパイプラインの使用で最も大きい利点を得る。レイヤ演算に比べて、AlexNetの暗号化及び復号化のオーバーヘッドは、暗号ハードウェア及びイントラ-レイヤパイプラインを使用して、大幅に減った他のDNNよりも大きい。下記の表3は、GuardiaNNがDNNの実行速度に与える影響を示す。

30

40

【0092】

50

【表 3】

Name	Pager	+Direct Conv	+SRAM Mgmt	+Crypt HW	+IntraLayer Pipelining
AlexNet [35]	1.00x	1.28x	2.81x	8.26x	9.92x
DeepID [69]	1.00x	1.12x	7.93x	13.25x	13.80x
FgptAuth [66]	1.00x	2.98x	7.89x	15.54x	16.76x
GAZEL [58]	1.00x	8.68x	20.59x	26.21x	26.54
MobileNetV1 [25]	1.00x	1.59x	3.31x	5.90x	6.54x
Smart Doll [11]	1.00x	8.44x	22.20x	30.58x	31.45x
Tiny Darknet [63]	1.00x	1.10x	4.11x	8.32x	9.35x
VGG-7 [37]	1.00x	3.59x	19.28x	25.26x	25.94x

10

【0093】

上記表3に示すように、本発明で提案したGuardiaNNはセキュリティ保証を損傷させることなく、広範囲なDNNの実行を加速化することができる。

【0094】

<高エネルギーの効率(High Energy Efficiency)>

【0095】

ここで、ベンチマークで各DNNに対する全ての構成でDNN実行のエネルギー消費を調査する。まず、遊休状態及びDNNの実行中に装置全体の平均電力を測定し、DNNの実行による平均電力の増加を計算するために両値を引く。その次に、平均電力の増加と待ち時間とを掛けてDNN実行のエネルギー消費を計算する。正規化された結果は、図15の(b)に出ている。提案された各技法を適用することによってエネルギー消費が減少することを観察することができる。基本フレームワークと比較して、GuardiaNNのエネルギー消費は、幾何平均92.3%減少し、エネルギー効率性が15.2倍向上した。これは、GuardiaNNが提供するかなりの待ち時間の減少のためである。提案された技術を適用した後にも、DNNの実行中に装置の全体電力は同じレベルに維持される。減少した待ち時間と結合されたGuardiaNNは、基本フレームワークよりもはるかに高いエネルギー効率性を達成する。

20

【0096】

<敏感度研究(Sensitivity Studies)>

【0097】

イントラ-レイヤパイプラインのバルクサイズがDNNの実行速度に与える影響を研究するために、4個のバルクサイズ(4、8、16、及び32の出力チャンネル)でGuardiaNNのDNNの実行待ち時間を測定する。図16は、バルクサイズ4に正規化された測定値を示す。4個のバルクサイズ以外にも、各DNNを最低の待ち時間に導く最適なレイヤ別バルクサイズの集合を計算して、次のように待ち時間を測定した。イントラ-レイヤパイプラインが最適なバルクサイズで待ち時間を最大14.24%向上させることを観察することができる(図16の5番目の棒)。待ち時間が最も高いバルクサイズと比較する。殆どのDNNの場合、バルクサイズが増加するにつれて、待ち時間が減少する傾向がある。バルクサイズが大きいほど暗号ハードウェアに対する暗号化及び復号化要求(及び割り込み)数が減少するためである。これは、パイプライン演算ステップの待ち時間を効果的に減らして、全体待ち時間を減らす。しかし、バルクサイズはSRAM容量に拘束される。例えば、FgptAuthに8よりも大きいバルクサイズを使用すると、必要なメモリサイズがSRAMサイズを超える。これが、バルクサイズが16及び32であるFgptAuthの実行待ち時間が図16で空いている理由である。

40

【0098】

GuardiaNNで使用するDMA可能暗号ハードウェアは、複数のブロック暗号及び作業モードをサポートする。GuardiaNNはセキュリティ強化のためにAESを使用し、ここではAES-ECB及びAES-CBCの2つの作動モードを比較する。(

50

a) 暗号ハードウェア、及び(b) OP - TEEの基本CPUベースの暗号化ライブラリである LibTomCrypt のある CPU で実行し、多様なキーサイズで AES - ECB 及び AES - CBC の暗号化及び復号化のスループットを測定する。図 17 は測定されたスループットを示す。LibTomCrypt を使用する CPU でよりも、暗号ハードウェアで AES - ECB 及び AES - CBC がはるかに速く実行されることが観察された。また、キーサイズが長いほど AES の暗号化及び復号化でより多くのラウンドが発生するため、LibTomCrypt を使用する CPU のスループットは、キーサイズが増加するにつれて減少する。動作モードを比較すると、AES - CBC はチェーンによって AES - ECB よりもスループットが少し少ない。逆に、暗号ハードウェアでは、キーサイズや作動モードがスループットに顕著な影響を与えない。これは、暗号化及び復号化性能が暗号ハードウェア及び DMA バッファ管理コストによって制限されるためである。それでも、GuardianN で使用される暗号ハードウェアは、全てのケースで CPU よりもはるかに高い暗号化及び復号化の効率性を提供する。

10

【0099】

モバイル及び組み込み機器で DNN の実行に対する関心が高まるにつれて、当該機器で DNN の実行を加速化する様々な技術が登場している。しかし、連合学習のような装置内で処理されるデータは差分プライバシーに基づいているため、依然として多くのプライバシーの問題がある。従って、DNN にトラスト実行環境を活用することが合理的な方向である。例えば、SecureTF は、トラスト実行環境を活用する TensorFlow ベースの分散セキュリティ機械学習フレームワークである。PPFL は、ローカル教育及び集計、多重パート ML にトラスト実行環境を活用して、安全な連合学習を加速化する。Chiron 及び Myelin は、機械学習内でトラスト実行環境をサービスとして活性化する。

20

【0100】

DarknetZ、Infenclave、及び Slalom は、いずれもトラスト実行環境内で DNN の一部を実行するように提案する。しかし、残りのレイヤを攻撃にさらして素直にアクセス方式を適用すると、かなりの量の性能オーバーヘッドが発生する。このような問題を解決するために、HybridTEE はリモートサーバーのトラスト実行環境に DNN の実行を要求することを提案する。しかし、HybridTEE の加速効果は、ローカルトラスト実行環境の DNN の実行を最適化していないため、あまり重要ではない。

30

【0101】

本発明は、ローカルトラスト実行環境で DNN の完全な保護を提案するだけでなく、実際のモバイル及び組み込み環境で実行できるようにする驚くべき速度の向上を提供することができる。

【0102】

< DNN 暗号化 (DNN Encryption) >

【0103】

DNN を保護するまた別の方向は、DNN データを暗号化することである。例えば、SecureML は、拡張可能な個人情報保護の DNN フレームワークを構築するために、安全なマルチパーティ計算を活用した。SoftME は信頼できる環境を提供し、暗号化と復号化及び演算ステップとで構成された信頼できる作業を実行する。SoftME で DNN を実行すると機密性が保証されるが、データの暗号化及び復号化に CPU を使用して、大きな性能オーバーヘッドが発生する。その中で、準同型暗号を適用すると、暗号化されたデータに対する演算が可能であるため有望である。Cryptonets は、そのようなアイデアが実現可能であることを示し、アイデアをセキュリティ教育に拡張する。MiniONN は、事前訓練された DNN を認識できないように変換する技術を提案する。また、準同型暗号化を活用して安全な連合遷移学習プロトコルを構築する。しかし、準同型暗号は演算スループットが低く、モバイル及び組み込み装置には実用的ではないと考えられることが多い。

40

50

【0104】

TEEは、ARM TrustZone及びIntel SGXが広く使用される商用実装であって、高いセキュリティ保証のために注目を集めている。ソフトウェアベースのセキュリティソリューションが適用されるが、様々なアプリケーションを保護するために成功的に悪用された。しかし、キャッシュアーキテクチャ、二重インスタンスアプリ、又は重畳アプリのようなTEEシステムを対象とする多くの脅威がある。これにより、セキュリティ強化のために、最近いくつかの提案がある。また、TEEの活用の困難を緩和するための多くの作業が提案されている。最小カーネルは、TEEの制限されたメモリ問題を解決するために小さなカーネルを構築する。CoSMIXは、アプリケーションレベルのセキュリティページエラー処理器を許容する。TEEMonはTEEのための性能モニタリングフレームワークである。

10

【0105】

本発明によるトラスト環境ベースの人工知能装置は、トラスト実行環境で人工ニューラルネットワークの実行を隔離して、動作速度の遅いDRAMに格納されたデータを暗号化してセキュリティを強化することができる。また、直接畳み込み及びSRAM管理を通じてDRAMのアクセス数を減らすことができ、暗号ハードウェアでデータの暗号化及び復号化をオフロードし、パイプラインングを実現してニューラルネットワーク演算の実行をデータの暗号化及び復号化演算にオーバーラップするように実行して人工ニューラルネットワークの実行を加速化することができる。

20

【0106】

以上、本発明の実施形態について図面を参照しながら詳細に説明したが、本発明は、上述の実施形態に限定されるものではなく、本発明の技術思想から逸脱しない範囲内で多様に変更実施することが可能である。

【符号の説明】

【0107】

- 100 (トラスト環境ベースの)人工知能装置
- 110 第1タイプのメモリ
- 130 トラスト人工知能処理部
- 131 暗号処理フロント-エンドプロセス
- 133 第2タイプのメモリ
- 135 プロセッサ

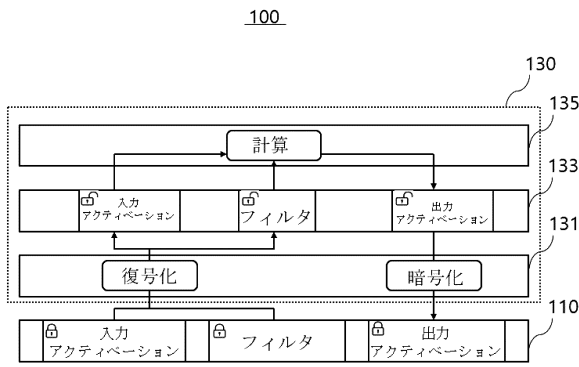
30

40

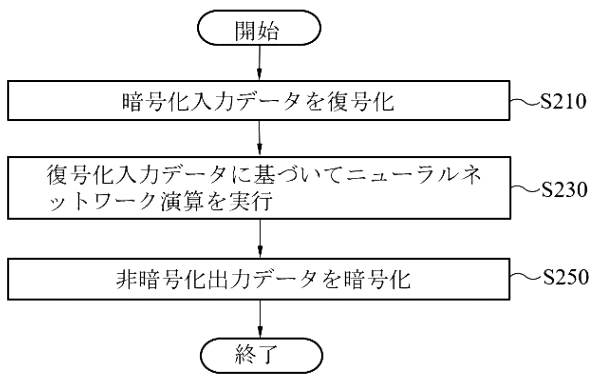
50

【図面】

【図 1】

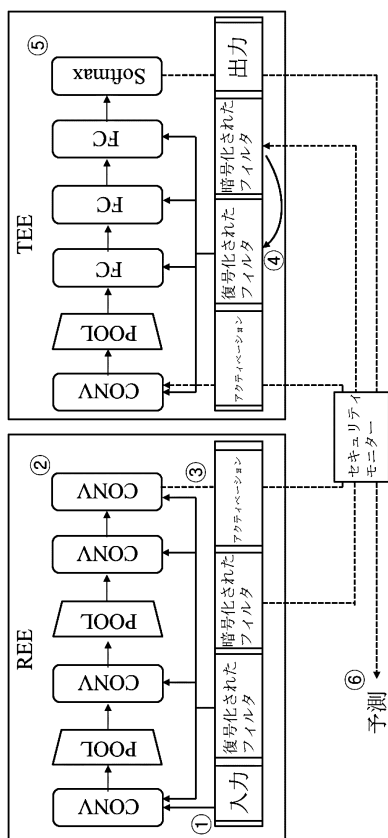


【図 2】

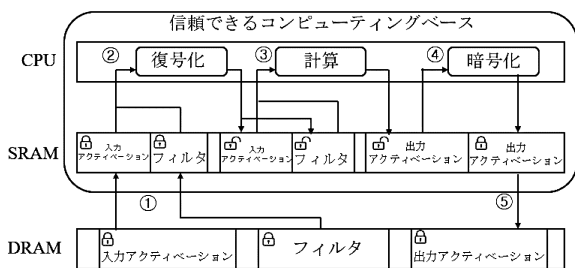


10

【図 3】



【図 4 a】



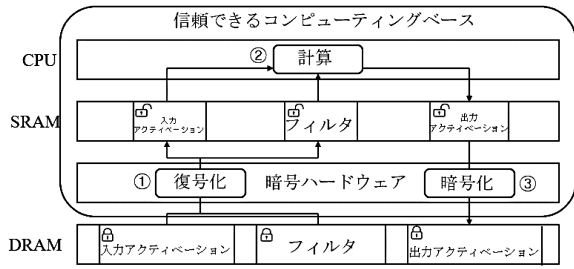
20

30

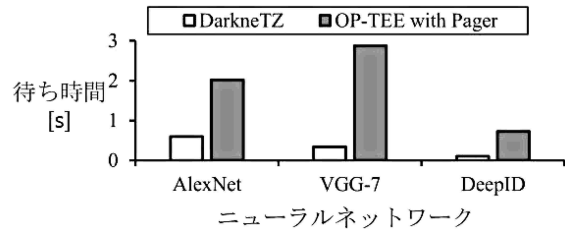
40

50

【図 4 b】

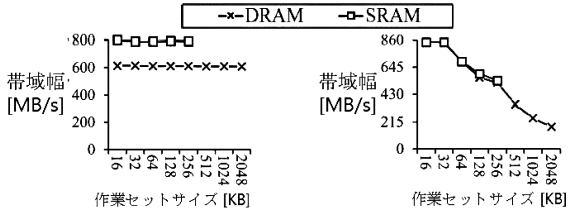


【図 5】



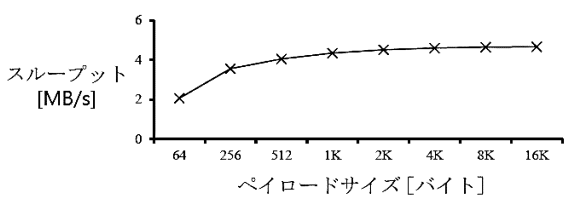
10

【図 6】



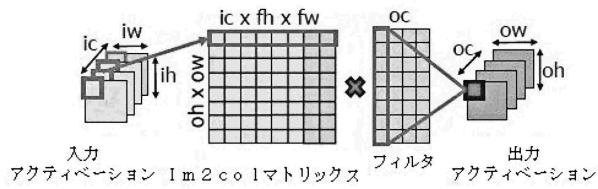
(a) シークンシャル 4-バイトアクセス (b) ランダム 4-バイトアクセス

【図 7】

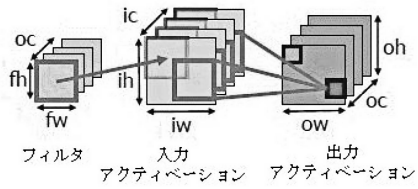


20

【図 8 a】



【図 8 b】

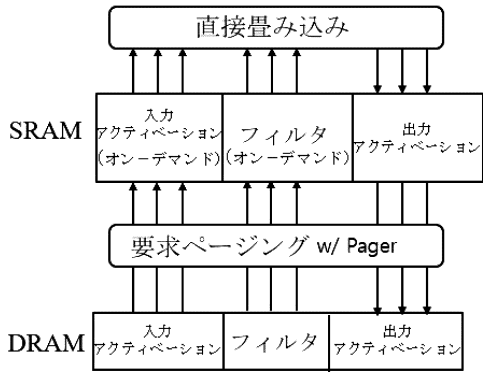


30

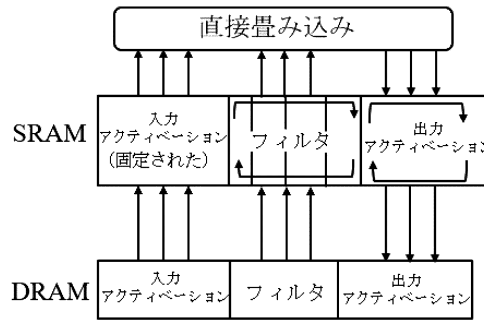
40

50

【図 9 a】

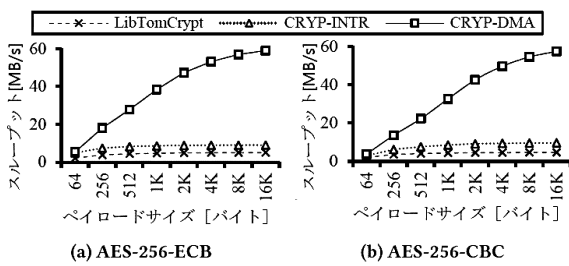


【図 9 b】

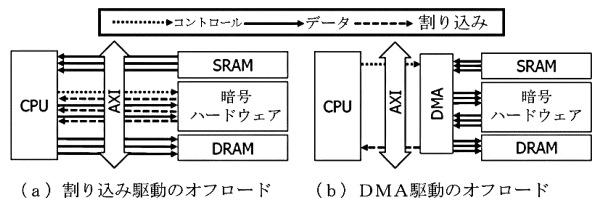


10

【図 10】

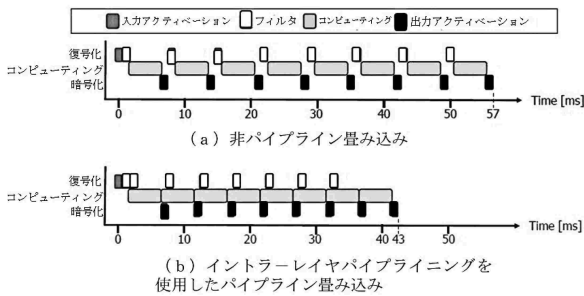


【図 11】



20

【図 12】



【図 13】

アルゴリズム 1: イントラレイヤパイプラインングのための畳み込みレイヤ

```

1: DOC: The first output channel for each intra-layer pipelining stage
2: CrypBlkSize: Cryptographic hardware invocation granularity
3: DfilterSoact: DRAM addresses for filters and output activations; CHW format
4: Siact: SRAM address for input activations; CHW format
5: SfilterSoact: In-SRAM circular queues for filters and output activations
6: f, o: # filters / # output activations of a bulk of output channels
7: F, O: # filters / # output activations to be requested to cryptographic hardware
8: crypCtx: the cryptographic context for the convolutional layer
9:
10: procedure PIPELINEDCONV_LAYER
11:   F = [f/CrypBlkSize] * CrypBlkSize
12:   DecryptData(crypCtx, Dfilter, Sfilter.tail, F)
13:   Sfilter.tail += F
14:   for oc in DOC do
15:     if (Sfilter.tail - Sfilter.head) == 0 then
16:       F = [f/CrypBlkSize] * CrypBlkSize
17:     else
18:       leftover = (Sfilter.tail - Sfilter.head - f)
19:       F = [(f + leftover)/CrypBlkSize] * CrypBlkSize
20:     end if
21:     DecryptData(crypCtx, Dfilter[oc], Sfilter.tail, F)
22:     Sfilter.tail = Sfilter.tail + F
23:     DirectConv(Siact, Sfilter.head, Soact.tail)
24:     Sfilter.head = Sfilter.head + f
25:     Soact.tail = Soact.tail + o
26:     O = [(Soact.tail - Soact.head)/CrypBlkSize] * CrypBlkSize
27:     EncryptData(crypCtx, Soact.head, Doact[oc], O)
28:     Soact.head = Soact.head + O
29:   end for
30: end procedure

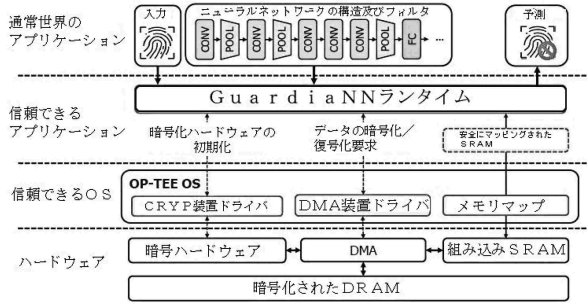
```

30

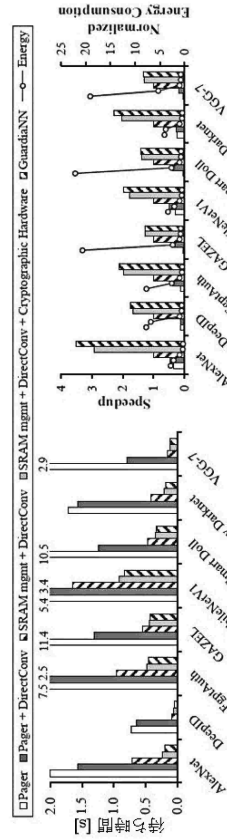
40

50

【図14】



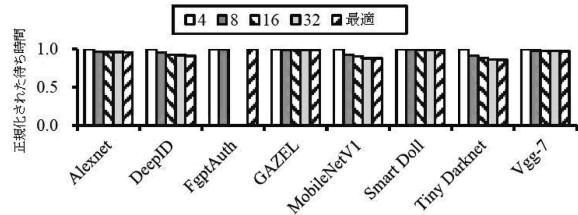
【図15】



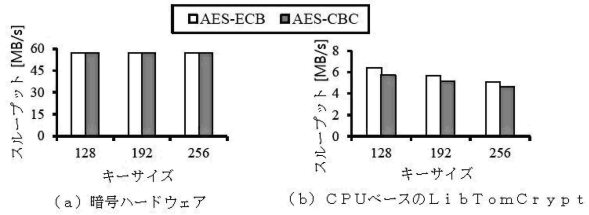
(b) 速度向上及び正規化されたエネルギー消費

(a) DNNの実行待ち時間

【図16】



【図17】



(a) 暗号ハードウェア

(b) CPUベースのLibTomCrypt

10

20

30

40

50

フロントページの続き

- 5 0
- (72)発明者 チェ , ジン ウ
大韓民国 ソウル特別市 西大門区 延世路 5 0
- (72)発明者 イム , チェ ミン
大韓民国 ソウル特別市 西大門区 延世路 5 0
- (72)発明者 イ , ス ヒョン
大韓民国 ソウル特別市 西大門区 延世路 5 0
- (72)発明者 ソン , ド キョン
大韓民国 ソウル特別市 西大門区 延世路 5 0
- (72)発明者 イ , ジン ホ
大韓民国 ソウル特別市 麻浦区 ワールド カップ ブク - ロ 2 3 5
- 審査官 青木 重徳
- (56)参考文献 米国特許第 1 0 9 5 6 5 8 4 (U S , B 1)
特開 2 0 0 5 - 0 1 8 4 3 4 (J P , A)
中井 綱人 , " P - T E E を用いた隔離 A I ハードウェアアクセラレーションの実装評価"
 , 2 0 2 2 年 暗号と情報セキュリティシンポジウム予稿集 , 2022年01月21日 , pp.1-8
Yunjie Deng et al. , StrangBox: A GPU TEE on Arm Endpoints , CCS'22: Proceedings of the
2022 ACM SIGSAC Conference on Computer and Communications Security , 2022年11月0
7日 , pp. 769-783
Meiyu Zhang et al. , SoftME: A Software-Based Memory Protection Approach for TEE Syste
m to Resist Physical Attacks , Security and Communication Networks , Hindawi [オンライ
ン] , 2019年01月 , Volume 2019, Issue 1 , (2024.10.29 検索)、インターネット , URL: htt
ps://onlinelibrary.wiley.com/doi/epdf/10.1155/2019/8690853
- (58)調査した分野 (Int.Cl. , D B 名)
G 0 6 F 2 1 / 6 0
G 0 6 F 2 1 / 6 2
G 0 6 N 3 / 0 2
J S T P l u s / J M E D P l u s / J S T 7 5 8 0 (J D r e a m I I I)
I E E E X p l o r e
T H E A C M D I G I T A L L I B R A R Y