

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7581209号
(P7581209)

(45)発行日 令和6年11月12日(2024.11.12)

(24)登録日 令和6年11月1日(2024.11.1)

(51)国際特許分類		F I	
H 0 1 L	27/088 (2006.01)	H 0 1 L	27/088 3 3 1 E
G 0 6 F	9/38 (2018.01)	G 0 6 F	9/38 3 7 0 C
G 0 6 F	12/00 (2006.01)	G 0 6 F	12/00 5 6 0 F
G 0 6 F	15/78 (2006.01)	G 0 6 F	15/78 5 5 0
G 1 1 C	5/02 (2006.01)	G 1 1 C	5/02 1 0 0
請求項の数 3 (全64頁) 最終頁に続く			
(21)出願番号	特願2021-538510(P2021-538510)	(73)特許権者	000153878
(86)(22)出願日	令和2年7月27日(2020.7.27)		株式会社半導体エネルギー研究所
(86)国際出願番号	PCT/IB2020/057051		神奈川県厚木市長谷398番地
(87)国際公開番号	WO2021/024083	(72)発明者	上妻 宗広
(87)国際公開日	令和3年2月11日(2021.2.11)		神奈川県厚木市長谷398番地 株式会
審査請求日	令和5年7月25日(2023.7.25)		社半導体エネルギー研究所内
(31)優先権主張番号	特願2019-146209(P2019-146209)	(72)発明者	石津 貴彦
(32)優先日	令和1年8月8日(2019.8.8)		神奈川県厚木市長谷398番地 株式会
(33)優先権主張国・地域又は機関			社半導体エネルギー研究所内
	日本国(JP)	(72)発明者	青木 健
(31)優先権主張番号	特願2019-157623(P2019-157623)		神奈川県厚木市長谷398番地 株式会
(32)優先日	令和1年8月30日(2019.8.30)		社半導体エネルギー研究所内
(33)優先権主張国・地域又は機関		(72)発明者	藤田 雅史
	日本国(JP)		神奈川県厚木市長谷398番地 株式会
(31)優先権主張番号	特願2019-216244(P2019-216244)		社半導体エネルギー研究所内
	最終頁に続く		最終頁に続く

(54)【発明の名称】 半導体装置

(57)【特許請求の範囲】

【請求項1】

C P Uと、
アクセラレータと、を有し、
前記アクセラレータは、第1メモリ回路と、駆動回路と、演算回路と、を有し、
前記第1メモリ回路は、第1トランジスタを有し、
前記第1トランジスタは、チャネル形成領域に金属酸化物を有する半導体層を有し、
前記駆動回路は、書き込み回路と、読み出し回路と、を有し、
前記書き込み回路は、切替信号、書き込み制御信号、およびデータ信号に応じて、前記第1メモリ回路に書き込むデータを2値または3値の電圧値に切り替えて出力する機能を有し、
前記読み出し回路は、第1参照電圧および第2参照電圧に応じて、前記第1メモリ回路に保持された電圧レベルに応じた2値または3値のデータを切り替えて読み出す機能を有し、
前記駆動回路および前記演算回路は、第2トランジスタを有し、
前記第2トランジスタは、チャネル形成領域にシリコンを有する半導体層を有し、
前記第1トランジスタと、前記第2トランジスタと、は積層して設けられる、半導体装置。

【請求項2】

請求項1において、

前記CPUは、バックアップ回路が設けられたフリップフロップを有するCPUコアを有し、

前記バックアップ回路は、前記CPUが非動作時において、前記フリップフロップに保持されたデータを電源電圧の供給が停止した状態で保持する機能を有する、半導体装置。

【請求項3】

請求項1または請求項2において、

前記演算回路は、積和演算を行う回路である、半導体装置。

【発明の詳細な説明】

【技術分野】

【0001】

本明細書は、半導体装置等について説明する。

【0002】

なお、本発明の一態様は、上記の技術分野に限定されない。本明細書等で開示する本発明の一態様の技術分野としては、半導体装置、撮像装置、表示装置、発光装置、蓄電装置、記憶装置、表示システム、電子機器、照明装置、入力装置、入出力装置、それらの駆動方法、又はそれらの製造方法、を一例として挙げることができる。

【背景技術】

【0003】

CPU(Central Processing Unit)等を含む半導体装置を有する電子機器が普及している。このような電子機器では、大量のデータを高速に処理するため、半導体装置の性能向上に関する技術開発が活発である。高性能化を実現する技術としては、例えば、GPU(Graphics Processing Unit)等のアクセラレータとCPUとを密結合させた、所謂SoC(System on Chip)化がある。SoC化によって高性能化した半導体装置では、発熱、及び消費電力の増加が問題となってくる。

【0004】

AI(Artificial Intelligence)技術では、計算量とパラメータ数が膨大になるため、演算量が増大する。演算量の増大は、発熱、および消費電力を増加させる要因となるため、演算量を低減するためのアーキテクチャが盛んに提案されている。代表的なアーキテクチャとして、Binary Neural Network(BNN)、およびTernary Neural Network(TNN)があり、回路規模縮小、および低消費電力化に対して特に有効となる(例えば特許文献1を参照)。例えば、BNNでは、もともと32ビット、もしくは16ビット精度で表現されたデータを、「+1」または「-1」の2値に圧縮することで、計算量とパラメータ数を大幅に削減できる。例えば、TNNでは、もともと32ビット、もしくは16ビット精度で表現されたデータを、「+1」、「0」または「-1」の3値に圧縮することで、計算量とパラメータ数を大幅に削減できる。BNNおよびTNNは、回路規模縮小や低消費電力化に有効なため、組み込みチップのように限られたハードウェア資源において低消費電力が求められるアプリケーションと相性が良いと考えられている。

【先行技術文献】

【特許文献】

【0005】

【文献】国際公開第2019/078924号

【発明の概要】

【発明が解決しようとする課題】

【0006】

TNNの演算には3値のデータを用いる。3値のデータをSRAM(Static RAM)に記憶する場合、メモリセル内のトランジスタ数が増えてしまう。そのため、半導体装置の小型化が難しくなるといった虞がある。また、半導体装置が有するアクセラレータでは、メモリが記憶するデータを2値または3値のデータで切り替える場合がある。こ

10

20

30

40

50

の場合、データに応じたメモリセルを用意する構成では、メモリセル内のトランジスタ数が増えてしまう。そのため、半導体装置の小型化が難しくなるといった虞がある。また半導体装置の消費電力は、CPUにおけるデータ転送回数が支配的である。そのため、半導体装置の低消費電力または発熱を抑制するためには、データ転送回数の増加を抑制することが重要となる。

【0007】

本発明の一態様は、半導体装置を小型化することを課題の一とする。または、本発明の一態様は、半導体装置を低消費電力化することを課題の一とする。または、本発明の一態様は、半導体装置の発熱を抑制することを課題の一とする。または、本発明の一態様は、CPUとメモリとして機能する半導体装置との間のデータ転送回数を削減することを課題の一とする。または、新規な構成の半導体装置を提供することを課題の一とする。

10

【0008】

なお、本発明の一態様は、必ずしも上記の課題の全てを解決する必要はなく、少なくとも一の課題を解決できるものであればよい。また、上記の課題の記載は、他の課題の存在を妨げるものではない。これら以外の課題は、明細書、特許請求の範囲、図面などの記載から、自ずと明らかとなるものであり、明細書、特許請求の範囲、図面などの記載から、これら以外の課題を抽出することが可能である。

【課題を解決するための手段】

【0009】

本発明の一態様は、CPUと、アクセラレータと、を有し、アクセラレータは、第1メモリ回路と、演算回路と、を有し、第1メモリ回路は、第1トランジスタを有し、第1トランジスタは、チャンネル形成領域に金属酸化物を有する半導体層を有し、演算回路は、第2トランジスタを有し、第2トランジスタは、チャンネル形成領域にシリコンを有する半導体層を有し、第1トランジスタと、第2トランジスタと、は積層して設けられる、半導体装置である。

20

【0010】

本発明の一態様は、CPUと、アクセラレータと、を有し、アクセラレータは、第1メモリ回路と、駆動回路と、演算回路と、を有し、第1メモリ回路は、第1トランジスタを有し、第1トランジスタは、チャンネル形成領域に金属酸化物を有する半導体層を有し、駆動回路は、書き込み回路と、読み出し回路と、を有し、書き込み回路は、切替信号、書き込み制御信号、およびデータ信号に応じて、第1メモリ回路に書き込むデータを2値または3値の電圧値に切り替えて出力する機能を有し、読み出し回路は、第1参照電圧および第2参照電圧に応じて、第1メモリ回路に保持された電圧レベルに応じた2値または3値のデータを切り替えて読み出す機能を有し、駆動回路および演算回路は、第2トランジスタを有し、第2トランジスタは、チャンネル形成領域にシリコンを有する半導体層を有し、第1トランジスタと、第2トランジスタと、は積層して設けられる、半導体装置である。

30

【0011】

本発明の一態様は、CPUと、アクセラレータと、を有し、アクセラレータは、第1メモリ回路と、演算回路と、を有し、第1メモリ回路は、第1トランジスタを有し、第1トランジスタは、チャンネル形成領域に金属酸化物を有する半導体層を有し、演算回路は、第2トランジスタを有し、第2トランジスタは、チャンネル形成領域にシリコンを有する半導体層を有し、CPUは、バックアップ回路が設けられたフリップフロップを有するCPUコアを有し、バックアップ回路は、第3トランジスタを有し、第3トランジスタは、チャンネル形成領域に金属酸化物を有する半導体層を有し、第1トランジスタと、第2トランジスタと、は積層して設けられる、半導体装置である。

40

【0012】

本発明の一態様は、CPUと、アクセラレータと、を有し、アクセラレータは、第1メモリ回路と、駆動回路と、演算回路と、を有し、第1メモリ回路は、第1トランジスタを有し、第1トランジスタは、チャンネル形成領域に金属酸化物を有する半導体層を有し、駆動回路は、書き込み回路と、読み出し回路と、を有し、書き込み回路は、切替信号、書き

50

込み制御信号、およびデータ信号に応じて、第1メモリ回路に書き込むデータを2値または3値の電圧値に切り替えて出力する機能を有し、読み出し回路は、第1参照電圧および第2参照電圧に応じて、第1メモリ回路に保持された電圧レベルに応じた2値または3値のデータを切り替えて読み出す機能を有し、演算回路は、第2トランジスタを有し、第2トランジスタは、チャネル形成領域にシリコンを有する半導体層を有し、CPUは、バックアップ回路が設けられたフリップフロップを有するCPUコアを有し、バックアップ回路は、第3トランジスタを有し、第3トランジスタは、チャネル形成領域に金属酸化物を有する半導体層を有し、第1トランジスタと、第2トランジスタと、は積層して設けられる、半導体装置である。

【0013】

本発明の一態様において、バックアップ回路は、CPUが非動作時において、フリップフロップに保持されたデータを電源電圧の供給が停止した状態で保持する機能を有する、半導体装置が好ましい。

【0014】

本発明の一態様において、演算回路は、積和演算を行う回路である、半導体装置が好ましい。

【0015】

本発明の一態様において、金属酸化物は、Inと、Gaと、Znと、を含む、半導体装置が好ましい。

【0016】

本発明の一態様において、第1トランジスタは、読出ビット線に電氣的に接続され、読出ビット線は、第2トランジスタが設けられた基板表面に概略垂直に設けられた配線を介して演算回路に電氣的に接続される、半導体装置が好ましい。

【0017】

なおその他の本発明の一態様については、以下で述べる実施の形態における説明、および図面に記載されている。

【発明の効果】

【0018】

本発明の一態様は、半導体装置を小型化することができる。または、本発明の一態様は、半導体装置を低消費電力化することができる。または、本発明の一態様は、半導体装置の発熱を抑制することができる。または、本発明の一態様は、CPUとメモリとして機能する半導体装置との間のデータ転送回数を削減することができる。または、新規な構成の半導体装置を提供することができる。

【0019】

複数の効果の記載は、他の効果の存在を妨げるものではない。また、本発明の一形態は、必ずしも、例示した効果の全てを有する必要はない。また、本発明の一形態について、上記以外の課題、効果、および新規な特徴については、本明細書の記載および図面から自ずと明らかになるものである。

【図面の簡単な説明】

【0020】

図1Aおよび図1Bは、半導体装置の構成例を説明する図である。

図2Aおよび図2Bは、半導体装置の構成例を説明する図である。

図3Aおよび図3Bは、半導体装置の構成例を説明する図である。

図4は、半導体装置の構成例を説明する図である。

図5Aおよび図5Bは、半導体装置の構成例を説明する図である。

図6Aおよび図6Bは、半導体装置の構成例を説明する図である。

図7Aおよび図7Bは、半導体装置の構成例を説明する図である。

図8Aおよび図8Bは、半導体装置の構成例を説明する図である。

図9は、半導体装置の構成例を説明する図である。

図10A、図10Bおよび図10Cは、半導体装置の処理性能と消費電力との関係を説明

10

20

30

40

50

する図である。

図 1 1 A および図 1 1 B は、半導体装置の構成例を説明する図である。

図 1 2 A および図 1 2 B は、半導体装置の構成例を説明する図である。

図 1 3 は、半導体装置の構成例を説明する図である。

図 1 4 A および図 1 4 B は、半導体装置の構成例を説明する図である。

図 1 5 A および図 1 5 B は、半導体装置の構成例を説明する図である。

図 1 6 は、半導体装置の構成例を説明する図である。

図 1 7 は、半導体装置の構成例を説明する図である。

図 1 8 A および図 1 8 B は、半導体装置の構成例を説明する図である。

図 1 9 A および図 1 9 B は、半導体装置の構成例を説明する図である。

10

図 2 0 A および図 2 0 B は、半導体装置の構成例を説明する図である。

図 2 1 は、半導体装置の構成例を説明する図である。

図 2 2 A および図 2 2 B は、半導体装置の構成例を説明する図である。

図 2 3 は、半導体装置の構成例を説明する図である。

図 2 4 は、半導体装置の構成例を説明する図である。

図 2 5 A および図 2 5 B は、半導体装置の構成例を説明する図である。

図 2 6 A および図 2 6 B は、半導体装置の構成例を説明する図である。

図 2 7 A および図 2 7 B は、半導体装置の構成例を説明する図である。

図 2 8 A および図 2 8 B は、半導体装置の構成例を説明する図である。

図 2 9 は、半導体装置の構成例を説明する図である。

20

図 3 0 は、C P U の構成例を説明する図である。

図 3 1 A および図 3 1 B は、C P U の構成例を説明する図である。

図 3 2 は、C P U の構成例を説明する図である。

図 3 3 は、集積回路の構成例を説明する図である。

図 3 4 A および図 3 4 B は、集積回路の構成例を説明する図である。

図 3 5 A および図 3 5 B は、集積回路の適用例を説明する図である。

図 3 6 A および図 3 6 B は、集積回路の適用例を説明する図である。

図 3 7 A、図 3 7 B および図 3 7 C は、集積回路の適用例を説明する図である。

図 3 8 は、集積回路の適用例を説明する図である。

図 3 9 A は、半導体装置の外観写真である。図 3 9 B は、半導体装置の断面 T E M 写真である。

30

図 4 0 は、半導体装置のシステム構成を説明するブロック図である。

図 4 1 A は、メモリセルの回路図である。図 4 1 B は、メモリセルの動作例を示すタイミングチャートである。図 4 1 C は、演算器の構成を示すブロック図である。

図 4 2 A および図 4 2 B は、半導体装置の構成を説明するブロック図である。

図 4 3 A および図 4 3 B は、半導体装置の動作期間中に生じる消費電力の推移を説明する概念図である。

図 4 4 A および図 4 4 B は、情報保持回路の回路図である。

図 4 5 A は、シミュレーション実行後の動作波形を示す図である。図 4 5 B は、シミュレーションで想定したニューラルネットワークモデルを示す図である。

40

【発明を実施するための形態】

【0021】

以下に、本発明の実施の形態を説明する。ただし、本発明の一形態は、以下の説明に限定されず、本発明の趣旨およびその範囲から逸脱することなくその形態および詳細を様々に変更し得ることは、当業者であれば容易に理解される。したがって、本発明の一形態は、以下に示す実施の形態の記載内容に限定して解釈されるものではない。

【0022】

なお本明細書等において、「第 1」、「第 2」、「第 3」という序数詞は、構成要素の混同を避けるために付したものである。従って、構成要素の数を限定するものではない。また、構成要素の順序を限定するものではない。また例えば、本明細書等の実施の形態の

50

一において「第 1」に言及された構成要素が、他の実施の形態、あるいは特許請求の範囲において「第 2」に言及された構成要素とすることもありうる。また例えば、本明細書等の実施の形態の一において「第 1」に言及された構成要素を、他の実施の形態、あるいは特許請求の範囲において省略することもありうる。

【0023】

図面において、同一の要素または同様な機能を有する要素、同一の材質の要素、あるいは同時に形成される要素等には同一の符号を付す場合があり、その繰り返しの説明は省略する場合がある。

【0024】

本明細書において、例えば、電源電位 VDD を、電位 VDD 、 VDD 等と省略して記載する場合がある。これは、他の構成要素（例えば、信号、電圧、回路、素子、電極、配線等）についても同様である。

10

【0025】

また、複数の要素に同じ符号を用いる場合、特に、それらを区別する必要があるときには、符号に“__1”、“__2”、“[n]”、“[m,n]”等の識別用の符号を付記して記載する場合がある。例えば、2 番目の配線 GL を配線 $GL[2]$ と記載する。

【0026】

(実施の形態 1)

本発明の一態様である半導体装置の構成、および動作等について説明する。

【0027】

20

なお、本明細書等において半導体装置とは、半導体特性を利用することで機能し得る装置全般を指す。トランジスタなどの半導体素子をはじめ、半導体回路、演算装置、記憶装置は、半導体装置の一態様である。表示装置（液晶表示装置、発光表示装置など）、投影装置、照明装置、電気光学装置、蓄電装置、記憶装置、半導体回路、撮像装置、電子機器などは、半導体装置を有すると言える場合がある。

【0028】

図 1 A および図 1 B は、本発明の一態様である半導体装置 100 を説明するための図である。半導体装置 100 は、 $CPU10$ 、アクセラレータ 20 およびバス 30 を有する。アクセラレータ 20 は、演算処理部 21 およびメモリ部 22 を有する。演算処理部 21 は、演算回路 23 を有する。メモリ部 22 は、メモリ回路 24 を有する。メモリ部 22 は、デバイスメモリ、共有メモリという場合がある。メモリ回路 24 は、チャンネル形成領域を有する半導体層 29 を有するトランジスタ 25 を有する。演算回路 23 とメモリ回路 24 とは、配線 31 を介して電氣的に接続される。

30

【0029】

$CPU10$ は、オペレーティングシステムの実行、データの制御、各種演算やプログラムの実行など、汎用の処理を行う機能を有する。 $CPU10$ は、1 つまたは複数の CPU コアを有する。 CPU コアはそれぞれ、電源電圧の供給が停止してもデータを保持できるデータ保持回路を有する。電源電圧の供給は、電源ドメイン（パワードメイン）からのパワースイッチ等による電氣的な切り離しによって制御することができる。なお電源電圧は、駆動電圧という場合がある。データ保持回路として、例えば、酸化物半導体（oxide semiconductor）をチャンネル形成領域に有するトランジスタ（OS トランジスタ）を有するメモリが好適である。なお酸化物半導体は、金属酸化物ともいう。OS トランジスタを有するデータ保持回路を備えた CPU コアの構成については、実施の形態 5 で説明する。

40

【0030】

アクセラレータ 20 は、ホストプログラムから呼び出されたプログラム（カーネル、またはカーネルプログラムとも呼ばれる。）を実行する機能を有する。アクセラレータ 20 は、例えば、グラフィック処理における行列演算の並列処理、ニューラルネットワークの積和演算の並列処理、科学技術計算における浮動小数点演算の並列処理などを行うことができる。

50

【0031】

メモリ部22は、アクセラレータ20が処理するデータを記憶する機能を有する。具体的には、ニューラルネットワークの積和演算の並列処理に用いる重みデータ等、演算処理部21に入力するあるいは出力されるデータを記憶することができる。

【0032】

メモリ回路24は、演算処理部21が有する演算回路23と配線31を介して電氣的に接続され、2値または3値のデジタル値を保持する機能を有する。メモリ回路24において、トランジスタ25が有する半導体層29は、酸化物半導体である。つまり、トランジスタ25は、OSトランジスタである。メモリ回路24は、OSトランジスタを有するメモリ（以下、OSメモリともいう。）が好適である。

10

【0033】

金属酸化物のバンドギャップは2.5 eV以上あるため、OSトランジスタは極小のオフ電流をもつ。一例として、ソースとドレイン間の電圧が3.5 V、室温（25℃）下において、チャンネル幅1 μm当たりのオフ電流を 1×10^{-20} A未満、 1×10^{-22} A未満、あるいは 1×10^{-24} A未満とすることができる。すなわち、ドレイン電流のオン/オフ電流比を20桁以上150桁以下とすることができる。そのため、OSメモリは、OSトランジスタを介して保持ノードからリークする電荷量が極めて少ない。従って、OSメモリは不揮発性メモリ回路として機能できるため、アクセラレータのパワーゲーティングが可能となる。

【0034】

20

高密度で集積化された半導体装置は、回路の駆動による熱が発生する場合がある。この発熱により、トランジスタの温度が上がることで、当該トランジスタの特性が変化して、電界効果移動度の変化や動作周波数の低下などが起こることがある。OSトランジスタは、Siトランジスタよりも熱耐性が高いため、温度変化による電界効果移動度の変化が起こりにくく、また動作周波数の低下も起こりにくい。さらに、OSトランジスタは、温度が高くなっても、ドレイン電流がゲート-ソース間電圧に対して指数関数的に増大する特性を維持しやすい。そのため、OSトランジスタを用いることにより、高い温度環境下での安定した動作を行うことができる。

【0035】

OSトランジスタに適用される金属酸化物は、Zn酸化物、Zn-Sn酸化物、Ga-Sn酸化物、In-Ga酸化物、In-Zn酸化物、In-M-Zn酸化物（Mは、Ti、Ga、Y、Zr、La、Ce、Nd、SnまたはHf）などがある。特にMとしてGaを用いる金属酸化物をOSトランジスタに採用する場合、元素の比率を調整することで電界効果移動度等の電気特性に優れたトランジスタとすることができるため、好ましい。また、インジウムおよび亜鉛を含む酸化物に、アルミニウム、ガリウム、イットリウム、銅、バナジウム、ベリリウム、ホウ素、シリコン、チタン、鉄、ニッケル、ゲルマニウム、ジルコニウム、モリブデン、ランタン、セリウム、ネオジム、ハフニウム、タンタル、タングステン、マグネシウムなどから選ばれた一種、または複数種が含まれていてもよい。

30

【0036】

OSトランジスタの信頼性、電気特性の向上のため、半導体層に適用される金属酸化物は、CAAC-OS、CAC-OS、nc-OSなどの結晶部を有する金属酸化物であることが好ましい。CAAC-OSとは、c-axis-aligned crystal line oxide semiconductorの略称である。CAC-OSとは、Cloud-Aligned Composite oxide semiconductorの略称である。nc-OSとは、nanocrystalline oxide semiconductorの略称である。

40

【0037】

CAAC-OSは、c軸配向性を有し、かつa-b面方向において複数のナノ結晶が連結し、歪みを有した結晶構造となっている。なお、歪みとは、複数のナノ結晶が連結する領域において、格子配列の揃った領域と、別の格子配列の揃った領域との間で格子配列の

50

向きが変化している箇所を指す。

【 0 0 3 8 】

C A C - O S は、キャリアとなる電子（または正孔）を流す機能と、キャリアとなる電子を流さない機能とを有する。電子を流す機能と、電子を流さない機能とを分離させることで、双方の機能を最大限に高めることができる。つまり、C A C - O S を O S トランジスタのチャンネル形成領域に用いることで、高いオン電流と、極めて低いオフ電流との双方を実現できる。

【 0 0 3 9 】

金属酸化物は、バンドギャップが大きく、電子が励起されにくいこと、ホールの有効質量が大きいことなどから、O S トランジスタは、一般的な S i トランジスタと比較して、アバランシェ崩壊等が生じにくい場合がある。従って、例えばアバランシェ崩壊に起因するホットキャリア劣化等を抑制できる。ホットキャリア劣化を抑制できることで、高いドレイン電圧で O S トランジスタを駆動することができる。

10

【 0 0 4 0 】

O S トランジスタは、電子を多数キャリアとする蓄積型トランジスタである。そのため、p n 接合を有する反転型トランジスタ（代表的には、S i トランジスタ）と比較して短チャンネル効果の一つである D I B L (D r a i n - I n d u c e d B a r r i e r L o w e r i n g) の影響が小さい。つまり、O S トランジスタは、S i トランジスタよりも短チャンネル効果に対する高い耐性を有する。

【 0 0 4 1 】

O S トランジスタは、短チャンネル効果に対する耐性が高いために、O S トランジスタの信頼性を劣化させずに、チャンネル長を縮小できるので、O S トランジスタを用いることで回路の集積度を高めることができる。チャンネル長が微細化するのに伴いドレイン電界が強まるが、上掲したように、O S トランジスタは S i トランジスタよりもアバランシェ崩壊が起きにくい。

20

【 0 0 4 2 】

また、O S トランジスタは、短チャンネル効果に対する耐性が高いために、S i トランジスタよりもゲート絶縁膜を厚くすることが可能となる。例えば、チャンネル長及びチャンネル幅が 5 0 n m 以下の微細なトランジスタにおいても、1 0 n m 程度の厚いゲート絶縁膜を設けることが可能な場合がある。ゲート絶縁膜を厚くすることで、寄生容量を低減することができるので、回路の動作速度を向上できる。またゲート絶縁膜を厚くすることで、ゲート絶縁膜を介したリーク電流が低減されるため、静的消費電流の低減につながる。

30

【 0 0 4 3 】

以上より、アクセラレータ 2 0 は、O S メモリであるメモリ回路 2 4 を有することで電源電圧の供給が停止してもデータを保持できる。そのため、アクセラレータ 2 0 のパワーゲーティングが可能となり、消費電力の大幅な低減を図ることができる。

【 0 0 4 4 】

O S トランジスタで構成されるメモリ回路 2 4 は、S i C M O S で構成することができる演算回路 2 3 と積層して設けることができる。そのため、回路面積の増加を招くことなく、配置することができる。メモリ回路 2 4 と演算回路 2 3 とは、演算回路 2 3 が設けられる基板表面に対して概略垂直な方向に延在して設けられる配線 3 1 を介して電氣的に接続される。なお「概略垂直」とは、8 5 度以上 9 5 度以下の角度で配置されている状態をいう。

40

【 0 0 4 5 】

メモリ回路 2 4 は、N O S R A M の回路構成とすることができる。「N O S R A M (登録商標) 」とは、「N o n v o l a t i l e O x i d e S e m i c o n d u c t o r R A M」の略称である。N O S R A M は、メモリセルが 2 トランジスタ型 (2 T) 、又は 3 トランジスタ型 (3 T) ゲインセルであり、アクセストランジスタが O S トランジスタであるメモリのことをいう。O S トランジスタはオフ状態でソースとドレインとの間を流れる電流、つまりリーク電流が極めて小さい。N O S R A M は、リーク電流が極めて小さ

50

い特性を用いてデータに応じた電荷をメモリ回路内に保持することで、不揮発性メモリとして用いることができる。特にNOSRAMは保持しているデータを破壊することなく読み出しすること（非破壊読み出し）が可能なため、データ読み出し動作のみを大量に繰り返す、ニューラルネットワークの積和演算の並列処理に適している。

【0046】

演算処理部21は、デジタル値を用いた演算処理を行う機能を有する。デジタル値はノイズの影響を受けにくい。そのためアクセラレータ20は、高い精度の演算結果が要求される演算処理を行うのに適している。なお演算処理部21は、SiCMOS、すなわちシリコンをチャネル形成領域に有するトランジスタ（Siトランジスタ）で構成されること好ましい。当該構成とすることでOSTランジスタと積層して設けることができる。

10

【0047】

演算回路23は、メモリ部22のメモリ回路24のそれぞれに保持されたデジタル値のデータを用いて、整数演算、単精度浮動小数点演算、倍精度浮動小数点演算などの処理のいずれかを行う機能を有する。演算回路23は、積和演算といった同じ処理を繰り返し実行する機能を有する。

【0048】

なお演算回路23は、メモリ回路24の読出ビット線毎、つまり一列（Column）毎に1つの演算回路23を設ける構成とする（Column-Parallel Calculation）。当該構成とすることで、メモリ回路24の1行分（最大で全ビット線）のデータを並列で演算処理することができる。CPU10を用いた積和演算に比べて、CPUとメモリ間のデータバスサイズ（32ビット、など）に制限されないことから、Column-Parallel Calculationでは、演算の並列度を大幅に上げることができるため、AI技術であるディープニューラルネットワークの学習（深層学習）、浮動小数点演算を行う科学技術計算などの膨大な演算処理に係る演算効率の向上を図ることができる。加えてメモリ回路24から出力されるデータの演算を完了させて読み出すことができるため、メモリアクセス（CPUとメモリ間のデータ転送やCPUでの演算）で生じる電力を削減することができ、発熱および消費電力の増加を抑制することができる。さらに、演算回路23とメモリ回路24の物理的な距離を近づけること、例えば積層によって配線距離が短くできることで、信号線に生じる寄生容量を削減できるため、低消費電力化が可能である。

20

30

【0049】

バス30は、CPU10とアクセラレータ20とを電氣的に接続する。つまりCPU10とアクセラレータ20とは、バス30を介してデータ伝送を行うことができる。

【0050】

本発明の一態様は、計算量とパラメータ数が膨大なAI技術などのアクセラレータとして機能する半導体装置を小型化することができる。または、本発明の一態様は、計算量とパラメータ数が膨大なAI技術などのアクセラレータとして機能する半導体装置を低消費電力化することができる。または、本発明の一態様は、計算量とパラメータ数が膨大なAI技術などのアクセラレータとして機能する半導体装置において、発熱を抑制することができる。または、本発明の一態様は、計算量とパラメータ数が膨大なAI技術などのアクセラレータとして機能する半導体装置において、CPUとメモリとして機能する半導体装置との間のデータ転送回数を削減することができる。換言すれば計算量とパラメータ数が膨大なAI技術などのアクセラレータとして機能する半導体装置は非ノイマン型アーキテクチャを有し、処理速度の増加に伴って消費電力が大きくなるノイマン型アーキテクチャと比較して、極めて少ない消費電力で並列処理を行うことができる。

40

【0051】

図2Aは、本発明の半導体装置100が有するメモリ部22に適用可能な回路構成例について説明する図である。図2Aでは、M行N列（M、Nは2以上の自然数）行列方向に並べて配置された書込用ワード線WWL₁乃至WWL_M、読出用ワード線RWL₁乃至RWL_M、書込用ビット線WBL₁乃至WBL_N、および読出用ビット線RBL

50

__ 1 乃至 R B L __ N を図示している。また各ワード線およびビット線に接続されたメモリ回路 2 4 を図示している。

【 0 0 5 2 】

図 2 B は、メモリ回路 2 4 に適用可能な回路構成例について説明する図である。メモリ回路 2 4 は、トランジスタ 2 5、トランジスタ 2 6、トランジスタ 2 7、容量素子 2 8 (キャパシタともいう) を有する。

【 0 0 5 3 】

トランジスタ 2 5 のソースまたはドレインの一方は、書込用ビット線 W B L に接続される。トランジスタ 2 5 のゲートは、書込用ワード線 W W L に接続される。トランジスタ 2 5 のソースまたはドレインの他方は、容量素子 2 8 の一方の電極およびトランジスタ 2 6 のゲートに接続される。トランジスタ 2 6 のソースまたはドレインの一方および容量素子 2 8 の他方の電極は、固定電位たとえばグラウンド電位を与える配線に接続される。トランジスタ 2 6 のソースまたはドレインの他方は、トランジスタ 2 7 のソースまたはドレインの一方に接続される。トランジスタ 2 7 のゲートは、読出用ワード線 R W L に接続される。トランジスタ 2 7 のソースまたはドレインの他方は、読出用ビット線 R B L に接続される。読出用ビット線 R B L は、上述したように、演算回路 2 3 が設けられる基板表面に対して概略垂直な方向に延在して設けられる配線 3 1 等を介して、演算回路 2 3 に接続される。

【 0 0 5 4 】

図 2 B に示すメモリ回路 2 4 の回路構成は、3 トランジスタ型 (3 T) ゲインセルの N O S R A M に相当する。トランジスタ 2 5 乃至トランジスタ 2 7 は、O S トランジスタである。O S トランジスタはオフ状態でソースとドレインとの間を流れる電流、つまりリーク電流が極めて小さい。N O S R A M は、リーク電流が極めて小さい特性を用いてデータに応じた電荷をメモリ回路内に保持することで、不揮発性メモリとして用いることができる。

【 0 0 5 5 】

図 3 A は、本発明の半導体装置 1 0 0 が有する演算処理部 2 1 に適用可能な回路構成例について説明する図である。演算処理部 2 1 は、N 個の演算回路 2 3 __ 1 乃至演算回路 2 3 __ N を有する。N 個の演算回路 2 3 __ 1 乃至演算回路 2 3 __ N はそれぞれ、N 本の読出用ビット線 R B L __ 1 乃至読出用ビット線 R B L __ N のいずれか一の信号が入力され、出力信号 Q __ 1 乃至 Q __ N を出力する。読出用ビット線 R B L __ 1 乃至読出用ビット線 R B L __ N の信号は、センスアンプ等で増幅して読み出す構成としてもよい。出力信号 Q __ 1 乃至 Q __ N は、メモリ回路 2 4 に保持したデータを用いて積和演算を行うことで得られるデータに相当する。

【 0 0 5 6 】

図 3 B は、演算回路 2 3 __ 1 乃至演算回路 2 3 __ N に適用可能な演算回路 2 3 の回路構成例を説明する図である。図 4 は、B i n a r y N e u r a l N e t w o r k (B N N) のアーキテクチャに基づく演算処理を実行するための回路である。演算回路 2 3 は、読出用ビット線 R B L の信号が与えられる読出回路 4 1 と、ビット積和演算器 4 2 と、アキュムレータ 4 3、ラッチ回路 4 4、および出力信号 Q を出力する符号化回路 4 5 を有する。

【 0 0 5 7 】

図 3 B で図示した演算回路 2 3 の構成について、より詳細を示す構成例を図 4 に図示する。図 4 では、8 ビットの信号 (W [0] 乃至 W [7]、A [0] 乃至 A [7]) の積和演算を行い、1 ビットの出力信号 Q、1 1 ビットの出力信号 (a c c o u t [1 0 : 0]) を出力する構成を一例として図示している。図 3 B では、メモリアクセスは 1 クロックで 1 行を選択するため、M 個 (= 1 ビット × M 行) の積とその和を M クロックで実行する。図 4 の演算回路では、同じ M 個の積とその和を 8 並列 × 1 ビット × M / 8 行で実行できるため、M / 8 クロックを要する。したがって、図 4 の構成は並列に積和演算を実行することで演算時間を短縮できるため、演算効率を向上できる。

【 0 0 5 8 】

10

20

30

40

50

図 4 において、ビット積和演算器 4 2 は、8 ビットの信号 ($W[0]$ 乃至 $W[7]$ 、 $A[0]$ 乃至 $A[7]$) が入力される積算器および当該積算器で得られた値が入力される加算器を有する。図 4 に示すように、8 並列で演算される 1 ビットの信号の積を $WA0$ 乃至 $WA7$ 、さらにその和を $WA10$ 、 $WA32$ 、 $WA54$ 、 $WA76$ 、さらにその和を $WA3210$ 、 $WA7654$ として図示している。

【0059】

図 4 において、加算器として機能するアキュムレータ 4 3 は、ビット積和演算器 4 2 の信号とラッチ回路 4 4 の出力信号との和をラッチ回路 4 4 に出力する。なおアキュムレータ 4 3 は、制御信号 $T \times D_EN$ に応じて加算器に入力する信号が切り替えられる。制御信号 $T \times D_EN$ が 0 ($T \times D_EN = 0$) でビット積和演算器 4 2 の信号とラッチ回路 4 4 の出力信号との和をラッチ回路 4 4 に出力する。制御信号 $T \times D_EN$ が 1 ($T \times D_EN = 1$) でロジック回路 4 7 の信号 (11 bit selector) とラッチ回路 4 4 の出力信号との和をラッチ回路 4 4 に出力する。

【0060】

図 4 において、AND 回路で構成されるロジック回路 4 7 は、信号 $A[0]$ 乃至 $A[7]$ と信号 $W[0]$ 乃至 $W[7]$ の積和演算が完了した後、バッチノーマライゼーションのためのデータを足し合わせる。具体的には切替信号 ($th_select[10:0]$) で切り替えながら、信号 $W[7]$ を足し合わせる。なお、バッチノーマライゼーションのためのデータは、例えば信号 $W[7]$ 以外の信号 $W[0]$ 乃至 $W[6]$ から同時に読み出して選択する構成としてもよい。バッチノーマライゼーションは、ニューラルネットワークにおける各層の出力データの分布が一定に収まるように調整するための動作である。例えば、ニューラルネットワークにおける演算によく利用される画像データは、学習に用いるデータの分布がばらつきやすいため、予測データ (入力データ) の分布と異なることがある。バッチノーマライゼーションは、ニューラルネットワークの中間層への入力データの分布を平均 0、分散 1 のガウス分布に正規化することで、ニューラルネットワークにおける学習の精度を高めることができる。Binary Neural Network (BNN) では活性化によって各層の出力結果が 2 値化されるため、しきい値に対してデータ分布の偏りを抑制することで、適切に活性化、つまり情報を分別できるようになる。

【0061】

ラッチ回路 4 4 は、アキュムレータ 4 3 の出力信号 ($account[10:0]$) を保持する。バッチノーマライゼーションによって次のニューラルネットワークにおける層 (NN 層) に渡す 2 値データはラッチ回路 4 4 が保持する積和演算結果の最上位ビットとなる。出力信号 ($account[10:0]$) において、最上位のビットの信号 ($account10$) は、2 の補数で演算されたラッチデータの符号を表し、そのプラスデータを 1、マイナスデータを 0 として次の NN 層に渡すため、符号化回路として機能するインバータ回路 4 6 で反転され、出力信号 Q として出力される。Q は中間層の出力であるため、アクセラレータ 2 0 内のバッファメモリ (入力バッファとも言う) に一時的に保持された後、次層の演算に使用される。

【0062】

図 5 A には、Binary Neural Network (BNN) のアーキテクチャに基づく、階層型のニューラルネットワークを図示する。図 5 A では、ニューロン 5 0、入力層 1 層 ($I1$)、中間層 3 層 ($M1$ 乃至 $M3$)、出力層 1 層 ($O1$) の全結合型のニューラルネットワークを図示している。入力層 $I1$ におけるニューロン数を 786、中間層 $M1$ 乃至 $M3$ におけるニューロン数を 256、出力層 $O1$ におけるニューロン数を 10 とすると、各層 (層 5 1、層 5 2、層 5 3 および層 5 4) の結合数は (784×256) + (256×256) + (256×256) + (256×10) で計 334336 個となる。つまり、ニューラルネットワーク計算に必要な重みパラメータが合計 330 K ビット程度であるため、小規模システムでも十分実装可能なメモリ容量とすることができる。

【0063】

次に、図 5 A に図示するニューラルネットワークの演算ができる、半導体装置 100 の

詳細なブロック図について図 5 B に示す。

【 0 0 6 4 】

図 5 B では、図 1 A および図 1 B で説明した、演算処理部 2 1、演算回路 2 3、メモリ部 2 2、メモリ回路 2 4、および配線 3 1 の他、図 1 A および図 1 B で図示する各構成を駆動するための周辺回路の構成例について図示している。

【 0 0 6 5 】

図 5 B では、コントローラ 6 1、ロウデコーダ 6 2、ワード線ドライバ 6 3、カラムデコーダ 6 4、書き込みドライバ 6 5、プリチャージ回路 6 6、センスアンプ 6 7、セクタ 6 8、入力バッファ 7 1 および演算制御回路 7 2 を図示している。

【 0 0 6 6 】

図 6 A は、図 5 B に図示する各構成について、メモリ部 2 2 を制御するブロックを抜き出した図である。図 6 A では、コントローラ 6 1、ロウデコーダ 6 2、ワード線ドライバ 6 3、カラムデコーダ 6 4、書き込みドライバ 6 5、プリチャージ回路 6 6、センスアンプ 6 7、セクタ 6 8 を抜き出して図示している。

【 0 0 6 7 】

コントローラ 6 1 は、外部からの入力信号を処理して、ロウデコーダ 6 2 およびカラムデコーダ 6 4 の制御信号を生成する。外部からの入力信号は、書き込みイネーブル信号や読み出しイネーブル信号などのメモリ部 2 2 を制御するための制御信号である。またコントローラ 6 1 は、C P U 1 0 との間でバスを介してメモリ部 2 2 に書き込まれるデータあるいはメモリ部 2 2 から読み出されるデータの入出力が行われる。

【 0 0 6 8 】

ロウデコーダ 6 2 は、ワード線ドライバ 6 3 を駆動するための信号を生成する。ワード線ドライバ 6 3 は、書込用ワード線 W W L および読出用ワード線 R W L に与える信号を生成する。カラムデコーダ 6 4 は、センスアンプ 6 7 および書き込みドライバ 6 5 を駆動するための信号を生成する。センスアンプ 6 7 は、読出用ビット線 R B L の電位を増幅する。書き込みドライバは、読出用ビット線 R B L および書込用ビット線 W B L を制御するための信号を生成する。プリチャージ回路 6 6 は、読出用ビット線 R B L などプリチャージする機能を有する。メモリ部 2 2 のメモリ回路 2 4 から読み出される信号は、演算回路 2 3 に入力される他、セクタ 6 8 を介して出力することができる。セクタ 6 8 は、バス幅に応じた分のデータを順次読出し、コントローラ 6 1 を介して必要なデータを C P U 1 0 等に出力することができる。

【 0 0 6 9 】

図 6 B は、図 5 B に図示する各構成について、演算処理部 2 1 を制御するブロックを抜き出した図である。

【 0 0 7 0 】

コントローラ 6 1 は、外部からの入力信号を処理して、演算制御回路 7 2 の制御信号を生成する。またコントローラ 6 1 は、演算処理部 2 1 が有する演算回路 2 3 を制御するための各種信号を生成する。またコントローラ 6 1 は、入力バッファ 7 1 を介して、演算結果に関するデータを入出力する。入力バッファ 7 1 を利用することで、C P U のデータバス幅以上のビット数の並列計算が可能となる。また膨大な数の重みパラメータを C P U 1 0 との間で転送する回数を削減できるため、低消費電力化を図ることができる。

【 0 0 7 1 】

本発明の一態様は、計算量とパラメータ数が膨大な A I 技術などのアクセラレータとして機能する半導体装置を小型化することができる。または、本発明の一態様は、計算量とパラメータ数が膨大な A I 技術などのアクセラレータとして機能する半導体装置を低消費電力化することができる。または、本発明の一態様は、計算量とパラメータ数が膨大な A I 技術などのアクセラレータとして機能する半導体装置において、発熱を抑制することができる。または、本発明の一態様は、計算量とパラメータ数が膨大な A I 技術などのアクセラレータとして機能する半導体装置において、C P U とメモリとして機能する半導体装置との間のデータ転送回数を削減することができる。換言すれば計算量とパラメータ数が

10

20

30

40

50

膨大なAI技術などのアクセラレータとして機能する半導体装置は非ノイマン型アーキテクチャを有し、処理速度の増加に伴って消費電力が大きくなるノイマン型アーキテクチャと比較して、極めて少ない消費電力で並列処理を行うことができる。

【0072】

(実施の形態2)

本発明の一態様である半導体装置の構成、および動作等について説明する。なお本実施の形態において、上記実施の形態と同じ符号が付される構成についての繰り返しの説明を省略する場合がある。

【0073】

図7Aおよび図7Bは、本発明の一態様である半導体装置100Aを説明するための図である。図7Aおよび図7Bにおける半導体装置100Aでは、一例として、CPU10、アクセラレータ20およびバス30を図示している。CPU10は、CPUコア11およびバックアップ回路12を有する。アクセラレータ20は、演算処理部21およびメモリ部22を有する。演算処理部21は、駆動回路15および演算回路23を有する。駆動回路15は、メモリ部22を駆動するための回路である。メモリ部22は、メモリ回路24を有する。メモリ部22は、デバイスメモリ、共有メモリという場合がある。メモリ回路24は、チャンネル形成領域を有する半導体層29を有するトランジスタ25を有する。駆動回路15とメモリ回路24とは、配線31を介して電氣的に接続される。

【0074】

メモリ回路24は、演算処理部21が有する演算回路23と配線31および駆動回路15を介して電氣的に接続される。メモリ回路24は、2値または3値のデータをアナログの電圧値として保持する機能を有する。当該構成とすることで、Binary Neural Network(BNN)、およびTernary Neural Network(TNN)といったアーキテクチャに基づく演算処理を演算処理部21で効率的に行うことができる。

【0075】

駆動回路15は、メモリ部22にデータを書き込むための書き込み回路およびメモリ部22からデータを読み出すための読み出し回路を有する。書き込み回路は、2値または3値のデータ信号の書き込みを切り替える切替信号、書き込み制御信号、およびデータ信号等の各種信号に応じて、メモリ部22にあるメモリ回路24に書き込むデータを2値または3値の電圧値に切り替えて出力する機能を有する。書き込み回路は、複数の信号が入力されるロジック回路で構成される。読み出し回路は、複数の参照電圧を用いて、メモリ部22にあるメモリ回路24に保持された電圧値を2値または3値のデータ信号に切り替えて読み出す機能を有する。読み出し回路は、センスアンプの機能を有する。

【0076】

OSトランジスタで構成されるメモリ回路24と駆動回路15とは、駆動回路15および演算回路23が設けられる基板表面に対して概略垂直な方向に延在して設けられる配線31を介して電氣的に接続される。なお「概略垂直」とは、85度以上95度以下の角度で配置されている状態をいう。なおメモリ回路24に接続されるビット線が書き込みビット線と読み出しビット線の場合、別々の配線を介して接続されることが好ましい。例えば書き込みビット線は、駆動回路15および演算回路23が設けられる基板表面に概略垂直に設けられた配線(第1配線)を介して書き込み回路に接続される。また例えば読み出しビット線は、駆動回路15および演算回路23が設けられる基板表面に概略垂直に設けられた配線(第2配線)を介して読み出し回路に接続される。

【0077】

次いで図8Aでは、図7Aおよび図7Bで説明した半導体装置100Aの構成に加え、バス30に接続されたOSメモリ300の他、DRAMなどで構成されるメインメモリ400を図示している。また図8Aでは、OSメモリ300とCPU10との間のデータをデータDCPUとして図示している。また図8Aでは、OSメモリ300とアクセラレータ20との間のデータをデータDACCとして図示している。

10

20

30

40

50

【 0 0 7 8 】

上述したように本発明の一態様の構成では、アクセラレータ 2 0 において、2 値または 3 値のアナログの電圧値をデータとして保持しつづけることができるとともに、演算回路で演算して得られる演算結果を C P U 1 0 に出力する構成とすることができる。そのため、演算処理のための O S メモリ 3 0 0 からのデータ D A C C を削減することができる。また C P U 1 0 の演算処理量を削減することができるため、O S メモリ 3 0 0 と C P U 1 0 との間のデータ D C P U も削減することができる。つまり本発明の一態様の構成では、バス 3 0 を介したアクセス数の低減、転送するデータ量の削減を図ることができる。

【 0 0 7 9 】

なお C P U 1 0 におけるバックアップ回路 1 2 およびアクセラレータ 2 0 におけるメモリ部 2 2 は、S i C M O S で構成することができる C P U コア 1 1 および演算処理部 2 1 と積層して設けることができる。そのため、回路面積の増加を招くことなく、配置することができる。

10

【 0 0 8 0 】

また O S メモリ 3 0 0 に適用可能な記憶回路としては、D O S R A M あるいは N O S R A M が好ましい。D O S R A M (登録商標)とは、「D y n a m i c O x i d e S e m i c o n d u c t o r R a n d o m A c c e s s M e m o r y (R A M)」の略称であり、1 T (トランジスタ) 1 C (容量) 型のメモリセルを有する R A M を指す。D O S R A M は、N O S R A M と同様に、O S トランジスタのオフ電流が低いことを利用したメモリである。

20

【 0 0 8 1 】

D O S R A M は、O S トランジスタを用いて形成された D R A M であり、D O S R A M は、外部から送られてくる情報を一時的に格納するメモリである。D O S R A M は、O S トランジスタを含むメモリセルと、S i トランジスタ(チャネル形成領域にシリコンを有するトランジスタ)を含む読み出し回路部を有する。上記メモリセルと読み出し回路部は、積層された異なる層に設けることができるため、D O S R A M は、全体の回路面積を小さくすることができる。また、D O S R A M は、メモリセルアレイを細かく分けて、効率的に配置することができる。

【 0 0 8 2 】

なお O S メモリ 3 0 0 において図 8 B に図示するように O S メモリ 3 0 0 は、O S トランジスタを有する層を積層して形成し、D O S R A M を高集積化した O S メモリ 3 0 0 N とすることで、単位面積あたりの記憶容量を大きくすることができる。この場合、半導体装置 1 0 0 A と別に設けるメインメモリ 4 0 0 を省略することも可能である。

30

【 0 0 8 3 】

O S メモリ 3 0 0 N を含め、半導体装置 1 0 0 A が有する C P U 1 0 およびアクセラレータ 2 0 が有する回路の一部を O S トランジスタで構成することで、各回路を一体化した 1 つの集積回路とすることができる。図 9 には、C P U 1 0 、アクセラレータ 2 0 および O S メモリ 3 0 0 N を密結合させた S o C として機能する半導体装置 1 0 0 A の模式図について図示する。

【 0 0 8 4 】

図 9 に図示するように、C P U 1 0 において、C P U コア 1 1 の上層にある O S トランジスタを有する層にバックアップ回路 1 2 を設ける構成とすることができる。また図 9 に図示するように、アクセラレータ 2 0 において、演算処理部 2 1 の上層にある O S トランジスタを有する層にメモリ部 2 2 を設けることができる。また図 9 に図示するように、メモリ部 2 2 と同様に積層された O S メモリ 3 0 0 N を配置することができる。その他、S i トランジスタを有するコントロール回路 5 0 0 、O S トランジスタを有するロジック回路 6 0 0 等を設ける構成とすることができる。なおロジック回路 6 0 0 は、O S トランジスタで代替可能な、切り替えスイッチなどの簡易なロジック回路等が好ましい。

40

【 0 0 8 5 】

図 9 に図示するように、C P U 1 0 、アクセラレータ 2 0 およびメモリ 3 0 0 N 等の各

50

回路を密結合させたS o Cの場合、発熱の問題があるが、O Sトランジスタは熱による電気特性の変動量がS iトランジスタと比べて小さいため、好適である。また、図9に図示するように三次元方向において回路を集積化することによって、シリコン貫通電極(Through Silicon Via:TSV)などを用いた積層構造などと比較して寄生容量を小さくすることができる。各配線の充放電に要する消費電力を削減することができる。そのため、演算処理効率の向上を図ることができる。

【0086】

図10Aは、処理性能(OPS:Operations Per Second)と、消費電力(W)との関係を説明する図である。なお、図10Aにおいて、縦軸が処理能力を、横軸が消費電力を、それぞれ表している。また、図10A中には、演算効率の指標として、0.1TOPS/W(Tera Operations Per Second/W)、1TOPS/W、10TOPS/W、及び100TOPS/Wを、破線にてそれぞれ明示してある。

10

【0087】

また、図10Aにおいて、領域710が従来の汎用AIアクセラレータ(ノイマン型)が含まれる領域を、領域712が本発明の一態様の半導体装置が含まれる領域を、それぞれ示している。なお、領域710には、例えば、CPU(Central Processing Unit)、GPU(Graphics Processing Unit)、FPGA(Field-Programmable Gate Array)などが含まれる。

20

【0088】

図10Aに示すように、本発明の一態様の半導体装置を適用することで、従来の汎用AIアクセラレータ(ノイマン型)よりも、2桁程度の消費電力を低減することができ、且つ処理性能を大幅(例えば1000倍以上)に向上させることができる。なお、本発明の一態様の半導体装置を適用することで、100TOPS/W以上の演算効率が期待できる。

【0089】

ここで、従来構成と、本発明の一態様の半導体装置を適用する構成との具体例について、図10B、及び図10Cを用いて説明する。図10Bが、画像認識における従来構成の半導体装置の消費電力のイメージ図を表し、図10Cが、画像認識における本発明の一態様の構成を用いる半導体装置の消費電力のイメージ図を表している。

30

【0090】

なお、図10B、及び図10Cにおいて、縦軸が電力を、横軸が時間を、それぞれ表している。また、図10Bにおいて、電力714がリーク電力を、電力716がCPU電力を、電力718がメモリ電力を、それぞれ示している。また、図10Cにおいて、電力714がリーク電力を、電力720がCPU電力を、電力722がアクセラレータ電力を、それぞれ示している。なお、電力722には、演算回路、及びメモリ回路に用いられる電力も含まれる。

【0091】

また、図10B、及び図10Cにおいて、矢印a、矢印b、及び矢印cは、それぞれ画像認識における信号を表している。なお、矢印a、矢印b、及び矢印cの信号が入力された際に、半導体装置にて、画像認識などの演算処理が開始されると仮定する。

40

【0092】

図10Bに示すように、従来構成の半導体装置の場合、時間に対して一定のリーク電力(電力714)が生じている。一方で、図10Cに示すように、本発明の一態様の半導体装置を適用する構成の場合、CPU電力(電力720)、及びアクセラレータ電力(電力722)を使用している間はリーク電力(電力714)が生じているが、CPU電力(電力720)、及びアクセラレータ電力(電力722)を使用していない期間は、リーク電力(電力714)が発生しないノーマリーオフ駆動(図10C中に示す期間t1)とすることができる。これにより、消費電力を大幅に低減することが可能となる。すなわち、極低消費電力な半導体装置を提供することができる。

50

【 0 0 9 3 】

図 1 1 A は、本発明の半導体装置 1 0 0 A が有するメモリ部 2 2 に適用可能な回路構成例について説明する図である。図 1 1 A では、M 行 N 列（M、N は 2 以上の自然数）行列方向に並べて配置された書き込み用ワード線 W W L 1 乃至 W W L M、読み出し用ワード線 R W L 1 乃至 R W L M、書き込み用ビット線 W B L 1 乃至 W B L N、および読み出し用ビット線 R B L 1 乃至 R B L N を図示している。また各ワード線およびビット線に接続されたメモリ回路 2 4 を図示している。

【 0 0 9 4 】

図 1 1 B は、メモリ回路 2 4 に適用可能な回路構成例について説明する図である。メモリ回路 2 4 は、トランジスタ 2 5、トランジスタ 2 6、トランジスタ 2 7、容量素子 2 8（キャパシタともいう）を有する。

10

【 0 0 9 5 】

トランジスタ 2 5 のソースまたはドレインの一方は、書き込み用ビット線 W B L に接続される。トランジスタ 2 5 のゲートは、書き込み用ワード線 W W L に接続される。トランジスタ 2 5 のソースまたはドレインの他方は、容量素子 2 8 の一方の電極およびトランジスタ 2 6 のゲートに接続される。トランジスタ 2 6 のソースまたはドレインの一方および容量素子 2 8 の他方の電極は、固定電位たとえばグラウンド電位を与える配線に接続される。トランジスタ 2 6 のソースまたはドレインの他方は、トランジスタ 2 7 のソースまたはドレインの一方に接続される。トランジスタ 2 7 のゲートは、読み出し用ワード線 R W L に接続される。トランジスタ 2 7 のソースまたはドレインの他方は、読み出し用ビット線 R B L に接続される。書き込み用ビット線 W B L および読み出し用ビット線 R B L は、上述したように、演算回路 2 3 が設けられる基板表面に対して概略垂直な方向に延在して設けられる配線等を介して、駆動回路 1 5 に接続される。駆動回路 1 5 は、2 値または 3 値のアナログの電圧値であるデータ信号 S O U T を出力する。また駆動回路 1 5 は、メモリ回路 2 4 から読み出されるデータに応じた読み出し用ビット線 R B L の電圧が与えられ、当該電圧に応じたデータ信号 D O 0、D O 1 を出力する。

20

【 0 0 9 6 】

図 1 1 B に示すメモリ回路 2 4 の回路構成は、3 トランジスタ型（3 T）ゲインセルの N O S R A M に相当する。トランジスタ 2 5 乃至トランジスタ 2 7 は、O S トランジスタである。O S トランジスタはオフ状態でソースとドレインとの間を流れる電流、つまりリーク電流が極めて小さい。N O S R A M は、リーク電流が極めて小さい特性を用いてデータに応じた電荷をメモリ回路内に保持することで、不揮発性メモリとして用いることができる。なお各トランジスタは、バックゲートを有する構成としてもよい。バックゲートを有することで、トランジスタ特性の向上を図ることができる。

30

【 0 0 9 7 】

図 1 2 A は、本発明の半導体装置 1 0 0 A が有する演算処理部 2 1 に適用可能な回路構成例について説明する図である。演算処理部 2 1 は、駆動回路 1 5 および演算回路 2 3 を有する。駆動回路 1 5 は、N 個の駆動回路 1 5 1 乃至駆動回路 1 5 N を有する。演算回路 2 3 は、N 個の演算回路 2 3 1 乃至演算回路 2 3 N を有する。N 個の駆動回路 1 5 1 乃至駆動回路 1 5 N はそれぞれ、N 本の読み出し用ビット線 R B L 1 乃至読み出し用ビット線 R B L N のいずれか一の信号が入力され、データ信号 D O 0 1 乃至 D O 0 N および / またはデータ信号 D O 1 1 乃至 D O 1 N を出力する。データ信号 D O 0 1 乃至 D O 0 N および / またはデータ信号 D O 1 1 乃至 D O 1 N は、演算回路 2 3 1 乃至演算回路 2 3 N に入力され、出力信号 Y 1 乃至 Y N を得る。出力信号 Y 1 乃至 Y N は、メモリ回路 2 4 に保持したデータを用いて積和演算を行うことで得られるデータに相当する。

40

【 0 0 9 8 】

図 1 2 B は、演算回路 2 3 1 乃至演算回路 2 3 N に適用可能な演算回路 2 3 の回路構成例を説明する図である。図 1 3 は、B i n a r y N e u r a l N e t w o r k（B N N）または T e r n a r y N e u r a l N e t w o r k（T N N）のアーキテクチャ

50

に基づく演算処理を実行するための回路である。演算回路 23 は、データ信号 D00 および / またはデータ信号 D01 が入力される読出回路 41 と、ビット積和演算器 42 と、アキュムレータ 43、ラッチ回路 44、および出力信号 Y を出力する符号化回路 45 を有する。

【0099】

図 12B で図示した演算回路 23 の構成について、より詳細を示す構成例を図 13 に図示する。図 13 では、8 ビットの信号 (W[0] 乃至 W[7]、A[0] 乃至 A[7]) の積和演算を行い、出力信号 Y、11 ビットの出力信号 (account[10:0]) を出力する構成を一例として図示している。図 12B では、メモリアクセスは 1 クロックで 1 行を選択するため、M 個 (= 1 ビット × M 行) の積とその和を M クロックで実行する。図 13 の演算回路では、同じ M 個の積とその和を 8 並列 × 1 ビット × M / 8 行で実行できるため、M / 8 クロックを要する。したがって、図 13 の構成は並列に積和演算を実行することで演算時間を短縮できるため、演算効率を向上できる。

【0100】

図 13 において、ビット積和演算器 42 は、8 ビットの信号 (W[0] 乃至 W[7]、A[0] 乃至 A[7]) が入力される積算器および当該積算器で得られた値が入力される加算器を有する。図 13 に示すように、8 並列で演算される 1 ビットの信号の積を WA0 乃至 WA7、さらにその和を WA10、WA32、WA54、WA76、さらにその和を WA3210、WA7654 として図示している。

【0101】

図 13 において、加算器として機能するアキュムレータ 43 は、ビット積和演算器 42 の信号とラッチ回路 44 の出力信号との和をラッチ回路 44 に出力する。なおアキュムレータ 43 は、制御信号 TxD_EN に応じて加算器に入力する信号が切り替えられる。制御信号 TxD_EN が 0 (TxD_EN = 0) でビット積和演算器 42 の信号とラッチ回路 44 の出力信号との和をラッチ回路 44 に出力する。制御信号 TxD_EN が 1 (TxD_EN = 1) でロジック回路 47 の信号 (11 bit selector) とラッチ回路 44 の出力信号との和をラッチ回路 44 に出力する。

【0102】

図 13 において、AND 回路で構成されるロジック回路 47 は、信号 A[0] 乃至 A[7] と信号 W[0] 乃至 W[7] の積和演算が完了した後、バッチノーマライゼーションのためのデータ、具体的には切替信号 (thselect[10:0]) で切り替えながら、信号 W[7] を足し合わせる。なお、バッチノーマライゼーションのためのデータは、例えば信号 W[7] 以外の信号 W[0] 乃至 W[6] から同時に読み出して選択する構成としてもよい。バッチノーマライゼーションは、ニューラルネットワークにおける各層の出力データの分布が一定に収まるように調整するための動作である。例えば、ニューラルネットワークにおける演算によく利用される画像データは、学習に用いるデータの分布がばらつきやすいため、予測データ (入力データ) の分布と異なることがある。バッチノーマライゼーションは、ニューラルネットワークの中間層への入力データの分布を平均 0、分散 1 のガウス分布に正規化することで、ニューラルネットワークにおける学習の精度を高めることができる。Binary Neural Network (BNN) では活性化によって各層の出力結果が 2 値化されるため、しきい値に対してデータ分布の偏りを抑制することで、適切に活性化、つまり情報を分別できるようになる。

【0103】

ラッチ回路 44 は、アキュムレータ 43 の出力信号 (account[10:0]) を保持する。バッチノーマライゼーションによって次のニューラルネットワークにおける層 (NN 層) に渡す 2 値データはラッチ回路 44 が保持する積和演算結果の最上位ビットとなる。出力信号 (account[10:0]) において、最上位のビットの信号 (account10) は、2 の補数で演算されたラッチデータの符号を表し、そのプラスデータを 1、マイナスデータを 0 として次の NN 層に渡すため、符号化回路として機能するインバータ回路 46 で反転され、出力信号 Y として出力される。Y は中間層の出力であるため、ア

10

20

30

40

50

クセラレータ 20 内のバッファメモリ（入力バッファとも言う）に一時的に保持された後、次層の演算に使用される。

【0104】

図14Aには、Binary Neural Network (BNN) または Ternary Neural Network (TNN) のアーキテクチャに基づく、階層型のニューラルネットワークを図示する。図14Aでは、ニューロン50、入力層1層 (I1)、中間層3層 (M1乃至M3)、出力層1層 (O1) の全結合型のニューラルネットワークを図示している。入力層 I1 におけるニューロン数を 786、中間層 M1 乃至 M3 におけるニューロン数を 256、出力層 O1 におけるニューロン数を 10 とすると、例えば Binary Neural Network (BNN) では、各層 (層 51、層 52、層 53 および層 54) の結合数は $(784 \times 256) + (256 \times 256) + (256 \times 256) + (256 \times 10)$ で計 334336 個となる。つまり、ニューラルネットワーク計算に必要な重みパラメータが合計 330 K ビット程度であるため、小規模システムでも十分実装可能なメモリ容量とすることができる。

10

【0105】

次に、図14Aに図示するニューラルネットワークの演算ができる、半導体装置 100 Aの詳細なブロック図について図14Bに示す。

【0106】

図14Bでは、図7Aおよび図7Bで説明した、演算処理部21、演算回路23、メモリ部22、メモリ回路24、および配線31の他、図7Aおよび図7Bで図示する各構成を駆動するための周辺回路の構成例について図示している。

20

【0107】

図14Bでは、コントローラ61、ロウデコーダ62、ワード線ドライバ63、カラムデコーダ64、書き込みドライバ65、プリチャージ回路66、センスアンプ67、セクタ68、入力バッファ71および演算制御回路72を図示している。

【0108】

図15Aは、図14Bに図示する各構成について、メモリ部22を制御するブロックを抜き出した図である。図15Aでは、コントローラ61、ロウデコーダ62、ワード線ドライバ63、カラムデコーダ64、書き込みドライバ65、プリチャージ回路66、センスアンプ67、セクタ68を抜き出して図示している。図7Aおよび図7Bで図示する駆動回路15は、書き込みドライバ65、プリチャージ回路66、およびセンスアンプ67のブロックに相当する。なお駆動回路15には、ワード線ドライバ63およびカラムデコーダ64を含めてもよい。

30

【0109】

コントローラ61は、外部からの入力信号を処理して、ロウデコーダ62およびカラムデコーダ64の制御信号を生成する。外部からの入力信号は、書き込みイネーブル信号や読み出しイネーブル信号などのメモリ部22を制御するための制御信号である。またコントローラ61は、CPU10との間でバスを介してメモリ部22に書き込まれるデータあるいはメモリ部22から読み出されるデータの入出力が行われる。

【0110】

ロウデコーダ62は、ワード線ドライバ63を駆動するための信号を生成する。ワード線ドライバ63は、書き込み用ワード線WWLおよび読み出し用ワード線RWLに与える信号を生成する。カラムデコーダ64は、センスアンプ67および書き込みドライバ65を駆動するための信号を生成する。プリチャージ回路66は、読み出し用ビット線RBLなどをプリチャージする機能を有する。メモリ部22のメモリ回路24から読み出される信号は、演算回路23に inputs される他、セクタ68を介して出力することができる。セクタ68は、バス幅に応じた分のデータを順次読み出しし、コントローラ61を介して必要なデータをCPU10等に出力することができる。

40

【0111】

図15Bは、図14Bに図示する各構成について、演算処理部21を制御するブロック

50

を抜き出した図である。

【 0 1 1 2 】

コントローラ 6 1 は、外部からの入力信号を処理して、演算制御回路 7 2 の制御信号を生成する。またコントローラ 6 1 は、演算処理部 2 1 が有する演算回路 2 3 を制御するための各種信号を生成する。またコントローラ 6 1 は、入力バッファ 7 1 を介して、演算結果に関するデータを入出力する。このバッファメモリを利用することで、CPU のデータバス幅以上のビット数の並列計算が可能となる。また膨大な数の重みパラメータを CPU 1 0 との間で転送する回数を削減できるため、低消費電力化を図ることができる。

【 0 1 1 3 】

図 1 6 では、2 値または 3 値のアナログの電圧値に変換されたデータ信号をメモリ回路に書き込むための、書き込みドライバ 6 5 の構成例について説明する。書き込みドライバ 6 5 は、インバータ回路 6 0 1、NAND 回路 6 0 2、NAND 回路 6 0 3、インバータ回路 6 0 4、トランジスタ 6 0 5、トランジスタ 6 0 6、およびインバータ回路 6 0 7 を有する。書き込みドライバ 6 5 を構成するトランジスタは、Si トランジスタである。トランジスタ 6 0 5 およびトランジスタ 6 0 6 は、図 1 6 に図示するように p チャネル型トランジスタが好ましい。

【 0 1 1 4 】

トランジスタ 6 0 5 およびトランジスタ 6 0 6 のソースまたはドレインの一方には、図 1 6 に図示するように、電位 V_{DD} ($> GND$) または電位 $V_{DD}/2$ ($> GND$) が与えられる。またインバータ回路 6 0 1 には、入力データであるデータ信号 $DI1$ が与えられる。NAND 回路 6 0 2 には、インバータ回路 6 0 1 の出力信号の他、データ信号 $DI0$ 、データの書き込みを制御するための書き込み制御信号 WE および 2 値または 3 値のデータ信号の書き込みを切り替えるための切替信号 B/T が入力される。NAND 回路 6 0 3 には、データ信号 $DI0$ および書き込み制御信号 WE が入力される。インバータ回路 6 0 7 は、2 値または 3 値のデータに応じた電圧値に相当するデータ信号 S_{OUT} を出力する。

【 0 1 1 5 】

図 1 6 に図示する各信号の真理値表は、表 1 のようになる。

【 0 1 1 6 】

【表 1】

WE	B/T	DI1	DI0	S _{out}
1	0	1	1	X
		0	1	VDD
		0	0	GND
	1	1	1	VDD
		0	1	VDD/2
		0	0	GND

【 0 1 1 7 】

つまり 2 値のデータをメモリ回路に書き込む場合、データ信号 S_{OUT} はデータ信号 $DI0$ に応じて、電圧 V_{DD} または電圧 GND に切り替えられる。3 値のデータをメモリ回路に書き込む場合、データ信号 S_{OUT} はデータ信号 $DI0$ および $DI1$ に応じて、電圧 V_{DD} 、電圧 $V_{DD}/2$ または電圧 GND の 3 値に切り替えられる。切り替えられた電圧は、書き込みビット線 WBL を介して、メモリ回路に書き込むことができる。

【 0 1 1 8 】

図 1 7 では、2 値または 3 値のアナログの電圧値に応じたデータ信号を演算回路 2 3 に出力するセンスアンプ 6 7 を含む構成例について説明する。図 1 7 では、入力信号に相当する読み出しビット線 R B L の電位から出力データであるデータ信号 D O 0、D O 1 を生成する、比較回路 6 1 1 および比較回路 6 1 2 が、センスアンプ 6 7 として機能する。比較回路 6 1 1 には、読み出しビット線 R B L の電位の電位および参照電圧 V r e f 1 が与えられる。比較回路 6 1 2 には、読み出しビット線 R B L の電位の電位および参照電圧 V r e f 2 が与えられる。参照電圧 V r e f 2 は、参照電圧 V r e f 1 より大きく、V D D より小さい。参照電圧 V r e f 1 は、G N D より大きく、V D D / 2 より小さい。

【 0 1 1 9 】

二値のデータの場合、バッファ回路 6 1 3 を介して出力される 2 値の出力データであるデータ信号 D O 0 およびデータ信号 B O が得られる。データ信号 D O 0 は、データ信号 B O と同じ論理値である。データ信号 D O 0 と、データ信号 B O と、の各信号の真理値表は、表 2 のようになる。

【 0 1 2 0 】

【表 2】

DO0	BO
0	0
1	1

【 0 1 2 1 】

3 値の出力データの場合、演算回路 2 3 を介して出力されるデータ信号 Y が得られる。データ信号 D O 0、データ信号 D O 1 と、データ信号 Y と、の各信号の真理値表は、表 3 のようになる。

【 0 1 2 2 】

【表 3】

DO1	DO0	W	A	Y=A*X
1	1	0	-1	0
0	1	+1	-1	-1
0	0	-1	-1	+1
1	1	0	+1	0
0	1	+1	+1	+1
0	0	-1	+1	-1

【 0 1 2 3 】

データ信号 Y は、重みデータ A とデータ信号 D O 0、D O 1 (X) とが、積和演算され

ることで積和信号 $Y (= A * X)$ を生成する。

【 0 1 2 4 】

以上説明したように、本発明の一態様は、アクセラレータとCPUを備えた半導体装置において、小型化された半導体装置を提供することができる。または、本発明の一態様は、アクセラレータとCPUを備えた半導体装置において、低消費電力化された半導体装置を提供することができる。または、本発明の一態様は、アクセラレータとCPUを備えた半導体装置において、発熱が抑制された半導体装置を提供することができる。または、本発明の一態様は、CPUにおけるデータ転送回数が削減された半導体装置を提供することができる。または、新規な構成の半導体装置を提供することができる。換言すれば、本発明の一態様の半導体装置は、非ノイマン型アーキテクチャを有し、処理速度の増加に伴って消費電力が大きくなるノイマン型アーキテクチャと比較して、極めて少ない消費電力で並列処理を行うことができる。

10

【 0 1 2 5 】

(実施の形態3)

本発明の一態様である半導体装置の構成、および動作等について説明する。なお本実施の形態において、上記実施の形態と同じ符号が付される構成についての繰り返しの説明を省略する場合がある。

【 0 1 2 6 】

図18Aおよび図18Bは、本発明の一態様である半導体装置100Bを説明するための図である。半導体装置100Bは、CPU10、アクセラレータ20およびバス30を有する。アクセラレータ20は、演算処理部21およびメモリ部22を有する。演算処理部21は、演算回路23を有する。メモリ部22は、メモリ回路24を有する。メモリ部22は、デバイスメモリ、共有メモリという場合がある。メモリ回路24は、チャネル形成領域を有する半導体層29を有するトランジスタ25を有する。演算回路23とメモリ回路24とは、配線31を介して電氣的に接続される。

20

【 0 1 2 7 】

メモリ部22は、アクセラレータ20が処理するデータを記憶および生成する機能を有する。具体的には、ニューラルネットワークの積和演算の並列処理に用いる重みデータ(第1データ信号ともいう)を記憶する機能を有する。またメモリ部22は、入力データ(第2データ信号ともいう)との乗算の結果に応じた出力データ(第3データ信号)を生成する機能を有する。メモリ部は、生成された出力データを演算処理部21に入力する機能を有する。

30

【 0 1 2 8 】

メモリ回路24は、演算処理部21が有する演算回路23と配線31を介して電氣的に接続され、2値で表される重みデータ、つまり1ビットのデジタル信号を保持する機能を有する。またメモリ回路は、重みデータと、入力データと、の乗算結果に相当する排他的論理和によって得られる信号を生成する機能を有する。なおメモリ回路24において、トランジスタ25が有する半導体層29は、酸化物半導体である。つまり、トランジスタ25は、OSトランジスタである。メモリ回路24は、OSトランジスタを有するメモリ(以下、OSメモリともいう。)が好適である。

40

【 0 1 2 9 】

図19Aは、本発明の半導体装置100Bが有するメモリ部22に適用可能な回路構成例について説明する図である。図19Aでは、M行N列(M、Nは2以上の自然数)行列方向に並べて配置された書込用ワード線WWL₁乃至WWL_M、読出用ワード線RWL₁₁乃至RWL_{MN}、読出用反転ワード線RWLB₁₁乃至RWLB_{MN}、書込用ビット線WBL₁乃至WBL_N、書込用反転ビット線WBLB₁乃至WBLB_N、および読出用ビット線RBL₁乃至RBL_Nを図示している。また各ワード線およびビット線に接続された複数のメモリ回路24を図示している。

【 0 1 3 0 】

図19Bは、メモリ回路24に適用可能な回路構成例について説明する図である。メモ

50

り回路 2 4 は、トランジスタ 3 1 A、3 1 B、トランジスタ 3 2 A、3 2 B、トランジスタ 3 3 A、3 3 B、容量素子 3 4 A、3 4 B（キャパシタともいう）の各素子を有する。各素子は、図 1 9 B に図示するように、書込用ワード線 WWL、読出用ワード線 RWL、読出用反転ワード線 RWLB、書込用ビット線 WBL、書込用反転ビット線 WBLB、および読出用ビット線 RBL の各配線に接続される。

【0131】

容量素子 3 4 A、3 4 B の一方の電極、およびトランジスタ 3 2 A、3 2 B のソースまたはドレインの一方は、固定電位たとえばグラウンド電位を与える配線に接続される。読出用ビット線 RBL は、上述したように、演算回路 2 3 が設けられる基板表面に対して概略垂直な方向に延在して設けられる配線 3 1 等を介して、演算回路 2 3 に接続される。

10

【0132】

図 1 9 B に示すメモリ回路 2 4 の回路構成は、トランジスタ 3 1 A、トランジスタ 3 2 A、およびトランジスタ 3 3 A 並びにトランジスタ 3 1 B、トランジスタ 3 2 B、およびトランジスタ 3 3 B で、3 トランジスタ型（3 T）ゲインセルの NOSRAM を構成する。トランジスタ 3 1 A、3 1 B、トランジスタ 3 2 A、3 2 B、トランジスタ 3 3 A、3 3 B は、OSTランジスタである。OSTランジスタはオフ状態でソースとドレインとの間を流れる電流、つまりリーク電流が極めて小さい。NOSRAM は、リーク電流が極めて小さい特性を用いてデータに応じた電荷をメモリ回路内に保持することで、不揮発性メモリとして用いることができる。図 1 9 B で言えば、トランジスタ 3 1 A、3 1 B をオフにすることで、ノード SN1、SN2 に与えられた電荷を保持することができる。なお各トランジスタは、バックゲート電極を有する構成としてもよい。

20

【0133】

図 1 9 B のメモリ回路 2 4 の真理値表は、表 4 のようになる。表 3 において H レベルおよび L レベルの電圧は、論理「1」、「0」で表している。「RWL」、「RWLB」は、入力データとして与えられる読出用ワード線 RWL、読出用反転ワード線 RWLB の電圧に応じた論理に相当する。「SN1」、「SN2」は、重みデータとして書込用ビット線 WBL、書込用反転ビット線 WBLB からノード SN1、SN2 に与えられる電圧に応じた論理に相当する。「RBL」は、出力データとして生成される読出用ビット線 RBL の電圧に応じた論理に相当する。

【0134】

30

【表 4】

RWL	RWLB	SN1	SN2	RBL
1	0	1	0	1
1	0	0	1	0
0	1	1	0	0
0	1	0	1	1

40

【0135】

図 1 9 B の回路構成において、表 4 で示す真理値表のデータを得ることができる。そのため、例えば、表 5 に示す読出用ワード線 RWL（入力データ A）と、ノード SN2（重みデータ W）と、の排他的論理和に基づく出力信号（出力データ $Y = W \times A$ ）を得ることができる。なお表 5 に図示するように、論理「1」、「0」は、Binary Neural Network（BNN）に用いる「+1」または「-1」の 2 値で表されるデータである。

【0136】

50

【表 5】

RWL (A)	SN1 (W)	RBL (Y=W×A)
-1 (0)	-1 (0)	+1 (1)
-1 (0)	+1 (1)	-1 (0)
+1 (1)	-1 (0)	-1 (0)
+1 (1)	+1 (1)	+1 (1)

10

【0137】

読出用ビット線 RBL にデータを読み出す場合の動作について、図 20A を用いて説明する。まずスタンバイ期間 T01 で読出用ワード線 RWL、読出用反転ワード線 RWLB を H レベル、読出用ビット線 RBL を中間電位とする。次いで、プリチャージ期間 T02 で読出用ワード線 RWL、読出用反転ワード線 RWLB を H レベルとし、読出用ビット線 RBL を H レベルとして電氣的に浮遊状態（フローティング）とする。次いで、読み出し期間 T03 で読出用ワード線 RWL、読出用反転ワード線 RWLB を入力データに応じた論理「1」、「0」とすることで、読出用ビット線 RBL の論理が「1」または「0」に変化することで出力データを生成することができる。

20

【0138】

重みデータの保持、および入力データとの排他的論理和に基づく信号を生成可能なメモリ部 22 は、図 20B に図示するような構成とすることができる。つまり複数のメモリ回路 24 において、重みデータである W_{11} 乃至 W_{MN} を記憶部 35 に保持させ、読出用ワード線 RWL_11 乃至 RWL_MN 、読出用反転ワード線 $RWLB_11$ 乃至 $RWLB_MN$ を介して入力データを排他的論理和部 36 (E×OR) に与えることで、重みデータと入力データとの排他的論理和に基づく出力データを読出用ビット線 RBL_1 乃至 RBL_N に入力することができる。

30

【0139】

なお図 19B のメモリ回路 24 は、図 21 の回路構成に変形することができる。図 21 のメモリ回路 24A は、ノード SN1、SN2 の接続先であるトランジスタ 32A、32B のゲートの接続を変更した構成に相当する。図 21 の回路構成において、表 6 で示す真理値表のデータを得ることができる。

【0140】

【表 6】

RWL	RWLB	SN1	SN2	RBL
1	0	1	0	0
1	0	0	1	1
0	1	1	0	1
0	1	0	1	0

40

【0141】

同様に図 22A のメモリ回路 24B は、ノード SN1 の接続先であるトランジスタを同

50

じ極性のトランジスタから、pチャネル型とnチャネル型を組み合わせたトランジスタ32__P、32__Nに変更した構成に相当する。トランジスタ32__P、32__Nは、Siトランジスタ等を用いることができる。当該構成とすることで、図19BにおけるノードSN2に接続されるトランジスタおよび配線を省略することができる。図22Aの回路構成において、表7で示す真理値表のデータを得ることができる。

【0142】

【表7】

RWL	RWLB	SN1	RBL
1	0	1	0
1	0	0	1
0	1	1	1
0	1	0	0

10

20

【0143】

同様に図22Bのメモリ回路24Cは、図19BのノードSN1、SN2の接続先である同じ極性のトランジスタから、異なる極性のトランジスタ32__P、32__Nに変更し、さらにトランジスタ37、38、および容量素子39を追加した構成に相当する。当該構成とすることで、ノードSN2に接続されるトランジスタおよび配線を省略することができる。図22Bの回路構成の真理値表は、表7と同様である。

【0144】

図23は、本発明の半導体装置100Bにおける、複数のメモリ回路24を有するメモリ部22と、演算回路23と、を説明する模式図である。上述したようにメモリ部22におけるメモリ回路24はそれぞれ、記憶部35と乗算部40とを備える。重みデータ W_1 乃至 W_k (k は2以上の自然数)は記憶部35に保持され、読出用ワード線RWL、読出用反転ワード線RWLBを介して入力される入力データ A_1 乃至 A_k と乗算に応じた1ビットのデジタル信号である出力信号($Y_k = A_k \times W_k$)が演算回路23に与えられる。メモリ部22の各トランジスタは、OSTランジスタとすることで、演算回路23と積層して設けることができるため好ましい。

30

【0145】

また図23に示す演算回路23は、アキュムレータ49と符号化回路45を備える。演算回路23は、乗算された出力信号を足し合わせることで、積和演算された信号Qを生成することができる。

【0146】

図23で図示した演算回路23の構成について、より詳細を示す構成例を図24に図示する。図24では、8ビットの信号($WA[0]$ 乃至 $WA[7]$)の加算を行い、1ビットの出力信号Q、11ビットの出力信号($account[10:0]$)を出力する構成を一例として図示している。図24の構成例では、積和演算と、バッチノーマライゼーションのための和の演算と、を切り替えて行う構成を図示している。図24では、メモリアクセスは1クロックで1行を選択するため、M個(=1ビット×M行)の積とその和をMクロックで実行する。図24の演算回路では、同じM個の積とその和を8並列×1ビット×M/8行で実行できるため、M/8クロックを要する。したがって、図24の構成は並列に積和演算を実行することで演算時間を短縮できるため、演算効率を向上できる。

40

【0147】

50

図 2 4 において、ビット加算器 4 2 A は、8 ビットの信号 (WA [0] 乃至 WA [7]) が入力される加算器を有する。図 2 4 に示すように、1 ビットの信号の和を WA 1 0、WA 3 2、WA 5 4、WA 7 6、さらにその和を WA 3 2 1 0、WA 7 6 5 4 として図示している。

【 0 1 4 8 】

図 2 4 において、加算器として機能するアキュムレータ 4 9 は、ビット加算器 4 2 A の信号とラッチ回路 4 4 の出力信号との和をラッチ回路 4 4 に出力する。なお図 2 4 において、アキュムレータ 4 9 に入力される信号は、制御信号 T x D _ E N に応じて切り替えられるセレクト 4 8 を備える。制御信号 T x D _ E N が 0 (T x D _ E N = 0) でビット加算器 4 2 A の信号とラッチ回路 4 4 の出力信号との和をラッチ回路 4 4 に出力する。制御信号 T x D _ E N が 1 (T x D _ E N = 1) でロジック回路 4 7 の信号 (1 1 b i t s e l e c t o r) とラッチ回路 4 4 の出力信号との和をラッチ回路 4 4 に出力する。セレクト 4 8 によって、積和演算と、バッチノーマライゼーションのための和の演算と、を切り替えて行うことができる。

【 0 1 4 9 】

図 2 4 において、AND 回路で構成されるロジック回路 4 7 は、信号 WA 0 乃至 WA 7 の積和演算が完了した後、バッチノーマライゼーションのためのデータ、具体的には切替信号 (t h s e l e c t [1 0 : 0]) で切り替えながら信号 R B L _ t h [1 0 : 0] を足し合わせる。なお、信号 R B L _ t h [1 0 : 0] は、メモリ回路 2 4 に保持される重みデータに相当する。バッチノーマライゼーションは、ニューラルネットワークにおける各層の出力データの分布が一定に収まるように調整するための動作である。例えば、ニューラルネットワークにおける演算によく利用される画像データは、学習に用いるデータの分布がばらつきやすいため、予測データ (入力データ) の分布と異なることがある。バッチノーマライゼーションは、ニューラルネットワークの中間層への入力データの分布を平均 0、分散 1 のガウス分布に正規化することで、ニューラルネットワークにおける学習の精度を高めることができる。Binary Neural Network (BNN) では活性化によって各層の出力結果が 2 値化されるため、しきい値に対してデータ分布の偏りを抑制することで、適切に活性化、つまり情報を分別できるようになる。

【 0 1 5 0 】

ラッチ回路 4 4 は、アキュムレータ 4 9 の出力信号 (a c c o u t [1 0 : 0]) を保持する。ラッチ回路 4 4 は、信号 C L R n でリセットされる。バッチノーマライゼーションによって次のニューラルネットワークにおける層 (NN 層) に渡す 2 値データはラッチ回路 4 4 が保持する積和演算結果の最上位ビットとなる。出力信号 (a c c o u t [1 0 : 0]) において、最上位のビットの信号 (a c c o u t 1 0) は、2 の補数で演算されたラッチデータの符号を表し、そのプラスデータを 1、マイナスデータを 0 として次の NN 層に渡すため、符号化回路として機能するインバータ回路 4 6 で反転され、出力信号 Q として出力される。Q は中間層の出力であるため、アクセラレータ 2 0 内のバッファメモリ (入力バッファとも言う) に一時的に保持された後、次層の演算に使用される。

【 0 1 5 1 】

図 2 5 A には、Binary Neural Network (BNN) のアーキテクチャに基づく、階層型のニューラルネットワークを図示する。図 2 5 A では、ニューロン 5 0、入力層 1 層 (I 1)、中間層 3 層 (M 1 乃至 M 3)、出力層 1 層 (O 1) の全結合型のニューラルネットワークを図示している。入力層 I 1 におけるニューロン数を 7 8 6、中間層 M 1 乃至 M 3 におけるニューロン数を 2 5 6、出力層 O 1 におけるニューロン数を 1 0 とすると、各層 (層 5 1、層 5 2、層 5 3 および層 5 4) の結合数は (7 8 4 × 2 5 6) + (2 5 6 × 2 5 6) + (2 5 6 × 2 5 6) + (2 5 6 × 1 0) で計 3 3 4 3 3 6 個となる。つまり、ニューラルネットワーク計算に必要な重みパラメータが合計 3 3 0 K ビット程度であるため、小規模システムでも十分実装可能なメモリ容量とすることができる。

【 0 1 5 2 】

次に、図 2 5 A に図示するニューラルネットワークの演算ができる、半導体装置 1 0 0

Bの詳細なブロック図について図25Bに示す。

【0153】

図25Bでは、図18Aおよび図18Bで説明した、演算処理部21、演算回路23、メモリ部22、メモリ回路24、および配線31の他、図18Aおよび図18Bで図示する各構成を駆動するための周辺回路の構成例について図示している。

【0154】

図25Bでは、コントローラ61、ロウデコーダ62、ワード線ドライバ63、カラムデコーダ64、書き込みドライバ65、プリチャージ回路66、センスアンプ67、セクタ68、入力バッファ71および演算制御回路72を図示している。

【0155】

図26Aは、図25Bに図示する各構成について、メモリ部22を制御するブロックを抜き出した図である。図26Aでは、コントローラ61、ロウデコーダ62、ワード線ドライバ63、カラムデコーダ64、書き込みドライバ65、プリチャージ回路66、センスアンプ67、セクタ68を抜き出して図示している。

【0156】

コントローラ61は、外部からの入力信号を処理して、ロウデコーダ62およびカラムデコーダ64の制御信号を生成する。外部からの入力信号は、書き込みイネーブル信号や読み出しイネーブル信号などのメモリ部22を制御するための制御信号である。またコントローラ61は、CPU10との間でバスを介してメモリ部22に書き込まれるデータあるいはメモリ部22から読み出されるデータの入出力が行われる。

【0157】

ロウデコーダ62は、ワード線ドライバ63を駆動するための信号を生成する。ワード線ドライバ63は、書込用ワード線WWLおよび読出用ワード線RWLに与える信号を生成する。カラムデコーダ64は、センスアンプ67および書き込みドライバ65を駆動するための信号を生成する。センスアンプ67は、読出用ビット線RBLの電位を増幅する。書き込みドライバは、読出用ビット線RBLおよび書込用ビット線WBLを制御するための信号を生成する。プリチャージ回路66は、読出用ビット線RBLなどをプリチャージする機能を有する。メモリ部22のメモリ回路24から読み出される信号は、演算回路23に入力される他、セクタ68を介して出力することができる。セクタ68は、バス幅に応じた分のデータを順次読出し、コントローラ61を介して必要なデータをCPU10等に出力することができる。

【0158】

図26Bは、図25Bに図示する各構成について、演算処理部21を制御するブロックを抜き出した図である。

【0159】

コントローラ61は、外部からの入力信号を処理して、演算制御回路72の制御信号を生成する。またコントローラ61は、演算処理部21が有する演算回路23を制御するための各種信号を生成する。またコントローラ61は、入力バッファ71を介して、演算結果に関するデータを入出力する。このバッファメモリを利用することで、CPUのデータバス幅以上のビット数の並列計算が可能となる。また膨大な数の重みパラメータをCPU10との間で転送する回数を削減できるため、低消費電力化を図ることができる。

【0160】

また上述したメモリ回路24は、トランジスタ等の構成を追加した回路構成に変形することができる。例えばメモリ回路24に適用可能な図27Aのメモリ回路24Dは、図19Bで図示した構成に加えて、トランジスタ81および容量素子82を追加した構成に相当する。また図27Aでは、ノードSOを図示している。なお図27Aに図示する回路構成は、図21に対応する変形例として、図27Bのメモリ回路24Eの構成とすることもできる。

【0161】

トランジスタ81は、OSトランジスタであることが好ましい。トランジスタ81をO

10

20

30

40

50

ストランジスタとすることで、リーク電流が極めて小さい特性を用いて容量素子 8 2、すなわちノード S O に出力データに応じた電荷を保持させることができる。ノード S O に保持された出力データは、トランジスタ 8 1 のゲートに接続された制御信号 S W に応じて、読出しビット線 R B L に出力させることができる。

【 0 1 6 2 】

図 2 8 A は、図 2 7 A の構成のメモリ回路 2 4 D をメモリ部 2 2 に適用した際の動作を説明するための模式図である。図 2 8 A に示すメモリ回路 2 4 D では、図 2 0 B で説明した記憶部 3 5 および排他的論理和部 3 6 に加えて、図 2 7 A で示したノード S O、およびスイッチとして機能するトランジスタを制御する制御信号 S W を図示している。1 行目のメモリ回路 2 4 D には、読出用ワード線 R W L _ 1 1 乃至 R W L _ 1 N のいずれか一、読出用反転ワード線 R W L B _ 1 1 乃至 R W L B _ 1 N のいずれか一が接続されている。M 行目のメモリ回路 2 4 D には、読出用ワード線 R W L _ M 1 乃至 R W L _ M N のいずれか一、読出用反転ワード線 R W L B _ M 1 乃至 R W L B _ M N のいずれか一が接続されている。また図 2 8 A では、読出しビット線 R B L _ 1 乃至 R B L _ N をプリチャージするためのプリチャージ電圧が与えられる配線に接続されたスイッチを制御する制御信号 P R E、読出しビット線 R B L _ 1 乃至 R B L _ N のノード P A、読出しビット線 R B L _ 1 乃至 R B L _ N と演算回路 2 3 A との間のスイッチを制御する制御信号 O U T を図示している。

10

【 0 1 6 3 】

各行のノード S O に保持される電荷を制御信号 S W で一斉に読出しビット線 R B L _ 1 乃至 R B L _ N にチャージシェアリングさせることで、読出しビット線 R B L _ 1 乃至 R B L _ N は各行のメモリ回路 2 4 D の出力データの和に応じた電位とすることができる。つまり読出しビット線 R B L _ 1 乃至 R B L _ N は、メモリ回路 2 4 D における乗算に応じた電荷の加算に応じたアナログ電圧とすることができる。そのため演算回路 2 3 A では、図 2 3 で説明した加算器の代わりにアナログデジタル変換回路を用いることができる。

20

【 0 1 6 4 】

読出用ビット線 R B L にデータを読み出す場合の動作について、図 2 8 B を用いて説明する。なお各スイッチは、H レベルでオン、L レベルでオフであるとして説明する。

【 0 1 6 5 】

まずスタンバイ期間 T 1 1 で読出用ワード線 R W L、読出用反転ワード線 R W L B を H レベル、制御信号 S W および制御信号 P R E を L レベル、ノード S O およびノード P A を中間電位とする。次いで、プリチャージ期間 T 1 2 で読出用ワード線 R W L、読出用反転ワード線 R W L B を H レベルとし、制御信号 S W および制御信号 P R E を H レベル、ノード S O およびノード P A を H レベルとして電氣的に浮遊状態（フローティング）とする。次いで、乗算を行う期間 T 1 3 で読出用ワード線 R W L、読出用反転ワード線 R W L B を入力データに応じた論理「1」、「0」とすることで、ノード S O の論理が「1」または「0」に変化する。期間 T 1 3 では、制御信号 S W を L レベル、制御信号 P R E およびノード P A を H レベルとする。次いで、加算を行う期間 T 1 4 で読出用ワード線 R W L、読出用反転ワード線 R W L B を H レベル、制御信号 P R E を L レベルとして、制御信号 S W を H レベルとする。ノード S O とノード P A がチャージシェアリングされ、ノード P A の電位は、乗算して得られた複数のメモリ回路におけるノード S O の電荷が加算されていられるアナログ電位とすることができる。当該アナログ電位は、制御信号 O U T によって、演算回路 2 3 A に読み出すことができる。

30

40

【 0 1 6 6 】

本発明の一態様は、計算量とパラメータ数が膨大な A I 技術などのアクセラレータとして機能する半導体装置を小型化することができる。または、本発明の一態様は、計算量とパラメータ数が膨大な A I 技術などのアクセラレータとして機能する半導体装置を低消費電力化することができる。または、本発明の一態様は、計算量とパラメータ数が膨大な A I 技術などのアクセラレータとして機能する半導体装置において、発熱を抑制することができる。または、本発明の一態様は、計算量とパラメータ数が膨大な A I 技術などのアク

50

セラレータとして機能する半導体装置において、CPUとメモリとして機能する半導体装置との間のデータ転送回数を削減することができる。換言すれば計算量とパラメータ数が膨大なAI技術などのアクセラレータとして機能する半導体装置は非ノイマン型アーキテクチャを有し、処理速度の増加に伴って消費電力が大きくなるノイマン型アーキテクチャと比較して、極めて少ない消費電力で並列処理を行うことができる。

【0167】

(実施の形態4)

本実施の形態では、上記実施の形態で説明したCPU10で実行するプログラムの演算の一部をアクセラレータ20で実行する場合の、動作の一例を説明する。

【0168】

図29は、CPUで実行するプログラムの演算の一部をアクセラレータで実行する場合の、動作の一例を説明する図である。

【0169】

CPUにて、ホストプログラムが実行される(ステップS1)。

【0170】

CPUは、アクセラレータを用いて演算を行う際に必要とされるデータ用領域を、メモリ部に確保するとの命令を確認した場合(ステップS2)、該データ用領域を、メモリ部に確保する(ステップS3)。

【0171】

次に、CPUは、メインメモリから上記メモリ部へ入力データを送信する(ステップS4)。上記メモリ部は該入力データを受信し、該入力データを、ステップS2で確保された領域に格納する(ステップS5)。

【0172】

CPUは、カーネルプログラムを起動するとの命令を確認した場合(ステップS6)、アクセラレータは、カーネルプログラムの実行を開始する(ステップS7)。

【0173】

アクセラレータがカーネルプログラムの実行を開始した直後、CPUを、演算を行う状態からPG状態へと切り替えてもよい(ステップS8)。その場合、アクセラレータがカーネルプログラムの実行を終了する直前に、CPUは、PG状態から演算を行う状態へ切り替えられる(ステップS9)。ステップS8からステップS9までの期間、CPUをPG状態にすることで、半導体装置全体として消費電力および発熱を抑制することができる。

【0174】

アクセラレータがカーネルプログラムの実行を終了すると、出力データが上記メモリ部に格納される(ステップS10)。

【0175】

カーネルプログラムの実行が終了した後、CPUは、メモリ部に格納された出力データをメインメモリへ送信するとの命令を確認した場合(ステップS11)、上記の出力データが上記メインメモリへ送信され、上記メインメモリに格納される(ステップS12)。

【0176】

CPUは、メモリ部上に確保されたデータ用領域を解放するとの指示を確認した場合(ステップS13)、上記メモリ部上に確保された領域が解放される(ステップS14)。

【0177】

以上のステップS1からステップS14までの動作を繰り返すことにより、CPUおよびアクセラレータの消費電力および発熱を抑制しつつ、CPUで実行するプログラムの演算の一部をアクセラレータで実行することができる。

【0178】

本実施の形態は、他の実施の形態の記載と適宜組み合わせることができる。

【0179】

(実施の形態5)

本実施の形態では、パワーゲーティングが可能なCPUコアを有するCPUの一例につ

10

20

30

40

50

いて説明する。

【0180】

図30に、CPU10の構成例を示す。CPU10は、CPUコア(CPU Core)200、L1(レベル1)キャッシュメモリ装置(L1 Cache)202、L2キャッシュメモリ装置(L2 Cache)203、バスインターフェース部(Bus I/F)205、パワースイッチ210~212、レベルシフタ(LS)214を有する。CPUコア200はフリップフロップ220を有する。

【0181】

バスインターフェース部205によって、CPUコア200、L1キャッシュメモリ装置202、L2キャッシュメモリ装置203が相互に接続される。

10

【0182】

外部から入力される割り込み信号(Interrupts)、CPU10が発行する信号SLEEP1等の信号に応じて、PMU193はクロック信号GCLK1、各種のPG(パワーゲーティング)制御信号(PG control signals)の生成を行う。クロック信号GCLK1、PG制御信号はCPU10に入力される。PG制御信号は、パワースイッチ210~212、フリップフロップ220を制御する。

【0183】

パワースイッチ210、211は、仮想電源線V__VDD(以下、V__VDD線と呼ぶ)への電圧VDDD、VDD1の供給をそれぞれ制御する。パワースイッチ212は、レベルシフタ(LS)214への電圧VDDHの供給を制御する。CPU10およびPMU193には、パワースイッチを介さずに電圧VSSSが入力される。PMU193には、パワースイッチを介さずに電圧VDDDが入力される。

20

【0184】

電圧VDDD、VDD1はCMOS回路用の駆動電圧である。電圧VDD1は電圧VDDDよりも低く、スリープ状態での駆動電圧である。電圧VDDHはOSトランジスタ用の駆動電圧であり、電圧VDDDよりも高い。

【0185】

L1キャッシュメモリ装置202、L2キャッシュメモリ装置203、バスインターフェース部205それぞれは、少なくとも1つパワーゲーティング可能なパワードメインを有する。パワーゲーティング可能なパワードメインには、1または複数のパワースイッチが設けられている。これらのパワースイッチは、PG制御信号によって制御される。

30

【0186】

フリップフロップ220は、レジスタに用いられる。フリップフロップ220には、バックアップ回路が設けられている。以下、フリップフロップ220について説明する。

【0187】

図31にフリップフロップ220(Flip-flop)の回路構成例を示す。フリップフロップ220はスキャンフリップフロップ(Scan Flip-flop)221、バックアップ回路(Backup Circuit)222を有する。

【0188】

スキャンフリップフロップ221は、ノードD1、Q1、SD、SE、RT、CK、クロックバッファ回路221Aを有する。

40

【0189】

ノードD1はデータ(data)入力ノードであり、ノードQ1はデータ出力ノードであり、ノードSDはスキャンテスト用データの入力ノードである。ノードSEは信号SCEの入力ノードである。ノードCKはクロック信号GCLK1の入力ノードである。クロック信号GCLK1はクロックバッファ回路221Aに入力される。スキャンフリップフロップ221のアナログスイッチは、クロックバッファ回路221AのノードCK1、CKB1に接続される。ノードRTはリセット信号(reset signal)の入力ノードである。

【0190】

50

信号 S C E は、スキャンイネーブル信号であり、P M U 1 9 3 で生成される。P M U 1 9 3 は信号 B K、R C (図示せず) を生成する。レベルシフタ 2 1 4 は信号 B K、R C をレベルシフトし、信号 B K H、R C H を生成する。信号 B K、R C はバックアップ信号、リカバリ信号である。

【 0 1 9 1 】

スキャンフリップフロップ 2 2 1 の回路構成は、図 3 1 に限定されない。標準的な回路ライブラリに用意されているフリップフロップを適用することができる。

【 0 1 9 2 】

バックアップ回路 2 2 2 は、ノード S D _ I N、S N 1 1、トランジスタ M 1 1 ~ M 1 3、容量素子 C 1 1 を有する。

10

【 0 1 9 3 】

ノード S D _ I N は、スキャンテストデータの入力ノードであり、スキャンフリップフロップ 2 2 1 のノード Q 1 に接続される。ノード S N 1 1 は、バックアップ回路 2 2 2 の保持ノードである。容量素子 C 1 1 はノード S N 1 1 の電圧を保持するための保持容量である。

【 0 1 9 4 】

トランジスタ M 1 1 はノード Q 1 とノード S N 1 1 間の導通状態を制御する。トランジスタ M 1 2 はノード S N 1 1 とノード S D 間の導通状態を制御する。トランジスタ M 1 3 はノード S D _ I N とノード S D 間の導通状態を制御する。トランジスタ M 1 1、M 1 3 のオンオフは信号 B K H で制御され、トランジスタ M 1 2 のオンオフは信号 R C H で制御される。

20

【 0 1 9 5 】

トランジスタ M 1 1 ~ M 1 3 は、上述したメモリ回路 2 4 が有するトランジスタ 2 5 乃至 2 7 と同様に、O S トランジスタである。トランジスタ M 1 1 ~ M 1 3 はバックゲートを有する構成を図示している。トランジスタ M 1 1 ~ M 1 3 のバックゲートは、電圧 V B G 1 を供給する電源線に接続されている。

【 0 1 9 6 】

少なくともトランジスタ M 1 1、M 1 2 が O S トランジスタであることが好ましい。オフ電流が極めて小さいという O S トランジスタの特長によって、ノード S N 1 1 の電圧の低下を抑えることができること、データの保持に電力を殆んど消費しないことから、バックアップ回路 2 2 2 は不揮発性の特性をもつ。容量素子 C 1 1 の充放電によってデータを書き換えるため、バックアップ回路 2 2 2 は原理的には書き換え回数に制約はなく、低エネルギーで、データの書き込みおよび読み出しが可能である。

30

【 0 1 9 7 】

バックアップ回路 2 2 2 の全てのトランジスタは O S トランジスタであることが非常に好ましい。図 3 1 B に示すように、シリコン C M O S 回路で構成されるスキャンフリップフロップ 2 2 1 上にバックアップ回路 2 2 2 を積層することができる。

【 0 1 9 8 】

バックアップ回路 2 2 2 は、スキャンフリップフロップ 2 2 1 と比較して素子数が非常に少ないので、バックアップ回路 2 2 2 を積層するためにスキャンフリップフロップ 2 2 1 の回路構成およびレイアウトの変更が必要ない。つまり、バックアップ回路 2 2 2 は、汎用性が非常に高いバックアップ回路である。また、スキャンフリップフロップ 2 2 1 が形成されている領域内にバックアップ回路 2 2 2 を設けることができるので、バックアップ回路 2 2 2 を組み込んでも、フリップフロップ 2 2 0 の面積オーバーヘッドはゼロにすることが可能である。よって、バックアップ回路 2 2 2 をフリップフロップ 2 2 0 に設けることで、C P U コア 2 0 0 のパワーゲーティングが可能となる。パワーゲーティングに必要なエネルギーが少ないため、C P U コア 2 0 0 を高効率にパワーゲーティングすることが可能である。

40

【 0 1 9 9 】

バックアップ回路 2 2 2 を設けることによって、トランジスタ M 1 1 による寄生容量が

50

ノードQ 1に付加されることになるが、ノードQ 1に接続される論理回路による寄生容量と比較して小さいので、スキャンフリップフロップ2 2 1の動作に影響はない。つまり、バックアップ回路2 2 2を設けても、フリップフロップ2 2 0の性能は実質的に低下しない。

【0 2 0 0】

C P Uコア2 0 0の低消費電力状態として、例えば、クロックゲーティング状態、パワーゲーティング状態、休止状態を設定することができる。P M U 1 9 3は、割り込み信号、信号S L E E P 1等に基づき、C P Uコア2 0 0の低消費電力モードを選択する。例えば、通常動作状態からクロックゲーティング状態に移行する場合、P M U 1 9 3はクロック信号G C L K 1の生成を停止する。

10

【0 2 0 1】

例えば、通常動作状態から休止状態に移行する場合は、P M U 1 9 3は、電圧および/または周波数スケーリングを行う。例えば、電圧スケーリングを行う場合、P M U 1 9 3は、電圧V D D 1をC P Uコア2 0 0に入力するため、パワースイッチ2 1 0をオフにし、パワースイッチ2 1 1をオンにする。電圧V D D 1は、スキャンフリップフロップ2 2 1のデータを消失させない電圧である。周波数スケーリングを行う場合、P M U 1 9 3はクロック信号G C L K 1の周波数を低下させる。

【0 2 0 2】

C P Uコア2 0 0を通常動作状態からパワーゲーティング状態に移行する場合には、スキャンフリップフロップ2 2 1のデータをバックアップ回路2 2 2にバックアップする動作が行われる。C P Uコア2 0 0をパワーゲーティング状態から通常動作状態に復帰する際には、バックアップ回路2 2 2のデータをスキャンフリップフロップ2 2 1に書き戻すリカバリ動作が行われる。

20

【0 2 0 3】

図3 2に、C P Uコア2 0 0のパワーゲーティングシーケンスの一例を示す。なお、図3 2において、t 1 ~ t 7は時刻を表している。信号P S E 0 ~ P S E 2は、パワースイッチ2 1 0 ~ 2 1 2の制御信号であり、P M U 1 9 3で生成される。信号P S E 0が“ H ” / “ L ”のとき、パワースイッチ2 1 0はオン / オフである。信号P S E 1、P S E 2についても同様である。

【0 2 0 4】

30

時刻t 1以前は、通常動作状態(N o r m a l O p e r a t i o n)である。パワースイッチ2 1 0はオンであり、C P Uコア2 0 0には電圧V D D Dが入力される。スキャンフリップフロップ2 2 1は通常動作を行う。このとき、レベルシフタ2 1 4は動作させる必要がないため、パワースイッチ2 1 2はオフであり、信号S C E、B K、R Cは“ L ”である。ノードS Eが“ L ”であるため、スキャンフリップフロップ2 2 1はノードD 1のデータを記憶する。なお、図3 2の例では、時刻t 1において、バックアップ回路2 2 2のノードS N 1 1は“ L ”である。

【0 2 0 5】

バックアップ(B a c k u p)時の動作を説明する。時刻t 1で、P M U 1 9 3はクロック信号G C L K 1を停止し、信号P S E 2、B Kを“ H ”にする。レベルシフタ2 1 4はアクティブになり、“ H ”の信号B K Hをバックアップ回路2 2 2に出力する。

40

【0 2 0 6】

バックアップ回路2 2 2のトランジスタM 1 1がオンになり、スキャンフリップフロップ2 2 1のノードQ 1のデータがバックアップ回路2 2 2のノードS N 1 1に書き込まれる。スキャンフリップフロップ2 2 1のノードQ 1が“ L ”であれば、ノードS N 1 1は“ L ”のままであり、ノードQ 1が“ H ”であれば、ノードS N 1 1は“ H ”になる。

【0 2 0 7】

P M U 1 9 3は、時刻t 2で信号P S E 2、B Kを“ L ”にし、時刻t 3で信号P S E 0を“ L ”にする。時刻t 3で、C P Uコア2 0 0の状態はパワーゲーティング状態に移行する。なお、信号B Kを立ち下げるタイミングで信号P S E 0を立ち下げてもよい。

50

【0208】

パワーゲーティング (Power-gating) 時の動作を説明する。信号 PSE0 が “L” になることで、V_{DD} 線の電圧が低下するため、ノード Q1 のデータは失われる。ノード SN11 は、時刻 t3 でのノード Q1 のデータを保持し続ける。

【0209】

リカバリ (Recovery) 時の動作を説明する。時刻 t4 で、PMU193 が信号 PSE0 を “H” にすることで、パワーゲーティング状態からリカバリ状態に移行する。V_{DD} 線の充電が開始され、V_{DD} 線の電圧が V_{DD} になった状態 (時刻 t5) で、PMU193 は信号 PSE2、RC、SCE を “H” にする。

【0210】

トランジスタ M12 はオンになり、容量素子 C11 の電荷がノード SN11 とノード SD とに分配される。ノード SN11 が “H” であれば、ノード SD の電圧は上昇する。ノード SE は “H” であるので、スキャンフリップフロップ 221 の入力側ラッチ回路にノード SD のデータが書き込まれる。時刻 t6 でノード CK にクロック信号 GCLK1 が入力されると、入力側ラッチ回路のデータがノード Q1 に書き込まれる。つまり、ノード SN11 のデータがノード Q1 に書き込まれたことになる。

【0211】

時刻 t7 で、PMU193 は信号 PSE2、SCE、RC を “L” にし、リカバリ動作が終了する。

【0212】

OSTランジスタを用いたバックアップ回路 222 は、動的および静的低消費電力双方が小さいため、ノーマリオフ・コンピューティングに非常に好適である。フリップフロップ 220 を搭載しても、CPUコア 200 の性能低下、動的電力の増加をほとんど発生させないようにできる。

【0213】

なお、CPUコア 200 は複数のパワーゲーティング可能なパワードメインを有してもよい。複数のパワードメインには、電圧の入力を制御するための 1 または複数のパワースイッチが設けられる。また、CPUコア 200 は、1 または複数のパワーゲーティングが行われないパワードメインを有していてもよい。例えば、パワーゲーティングが行われないパワードメインに、フリップフロップ 220、パワースイッチ 210 ~ 212 の制御を行うためのパワーゲーティング制御回路を設けてもよい。

【0214】

なお、フリップフロップ 220 の適用は CPU10 に限定されない。演算装置において、パワーゲーティング可能なパワードメインに設けられるレジスタに、フリップフロップ 220 を適用できる。

【0215】

本実施の形態は、他の実施の形態の記載と適宜組み合わせることができる。

【0216】

(実施の形態 6)

本実施の形態では、上記実施の形態で説明した半導体装置 100 の構成を含む集積回路の構成について図 33 および図 34 を参照しながら説明する。

【0217】

図 33 は、半導体装置 100 の構成を含む集積回路の構成例を説明するためのブロック図の一例である。

【0218】

図 33 に図示する集積回路 390 は、CPU10、アクセラレータ 20、オンチップメモリ 131、DMAC (Direct Memory Access Controller) 141、電源回路 160、パワーマネジメントユニット (PMU) 142、セキュリティー回路 147、メモリコントローラ 143、DDR SDRAM (Double Data Rate Synchronous Dynamic Random Access

10

20

30

40

50

Memory)コントローラ144、USB(Universal Serial Bus)インターフェース回路145、ディスプレイインターフェース回路146、ブリッジ回路150、割り込み制御回路151、インターフェース回路152、バッテリー制御回路153、およびADC(Analog-to-digital converter)/DAC(Digital-to-analog converter)インターフェース回路154を有する。

【0219】

CPU10は、一例として、CPUコア111、命令キャッシュ112、データキャッシュ113、およびバスインターフェース回路114を有する。アクセラレータ20は、メモリ回路121、演算回路122、および制御回路123を有する。

10

【0220】

CPUコア111は、複数のCPUコアを有する。命令キャッシュ112は、CPUコア111で実行する命令を一時的に記憶する回路構成とすればよい。データキャッシュ113は、CPUコア111で処理するデータまたは処理によって得られたデータを一時的に記憶する回路構成とすればよい。バスインターフェース回路114は、CPU10と、半導体装置内の他の回路とを接続するためのバスとデータやアドレス等の信号を送受信することができる回路構成であればよい。

【0221】

メモリ回路121は、実施の形態1で説明したメモリ回路24に相当する。メモリ回路121は、アクセラレータ20で処理するデータを記憶する回路構成とすればよい。演算回路122は、実施の形態1で説明した演算回路23に相当する。演算回路122は、メモリ回路121に保持したデータの演算処理を行う回路構成とすればよい。制御回路123は、図5Bで図示したように、アクセラレータ20内の各回路を制御するための回路構成とすればよい。

20

【0222】

高速バス140Aは、CPU10、アクセラレータ20、オンチップメモリ131、DMAC141、パワーマネジメントユニット142、セキュリティ回路147、メモリコントローラ143、DDR SDRAMコントローラ144、USBインターフェース回路145、およびディスプレイインターフェース回路146の間の各種信号を高速で送受信するためのバスである。一例としては、AMBA(Advanced Microcontroller Bus Architecture)-AHB(Advanced High-performance Bus)をバスとして用いることができる。

30

【0223】

オンチップメモリ131は、集積回路390が有する回路、例えばCPU10またはアクセラレータ20に入出力するデータまたはプログラムを記憶するための回路構成を有する。

【0224】

DMAC141は、ダイレクトメモリアクセスコントローラである。DMAC141を有することで、CPU10以外の周辺機器は、CPU10を介さずにオンチップメモリ131にアクセスすることができる。

40

【0225】

パワーマネジメントユニット142は、集積回路390が有するCPUコア等の回路のパワーゲーティングを制御するための回路構成を有する。

【0226】

セキュリティ回路147は、集積回路390と外部の回路との間で暗号化して信号を送受信するなど、信号の秘匿性を高めるための回路構成を有する。

【0227】

メモリコントローラ143は、集積回路390の外部にあるプログラムメモリからCPU10またはアクセラレータ20で実行するためのプログラムを書き込みまたは読み出しを行うための回路構成を有する。

50

【 0 2 2 8 】

DDR SDRAMコントローラ 1 4 4 は、集積回路 3 9 0 の外部にあるDRAM等のメインメモリとの間でデータを書き込みまたは読み出しを行うための回路構成を有する。

【 0 2 2 9 】

USBインターフェース回路 1 4 5 は、集積回路 3 9 0 の外部にある回路とUSB端子を介してデータの送受信を行うための回路構成を有する。

【 0 2 3 0 】

ディスプレイインターフェース回路 1 4 6 は、集積回路 3 9 0 の外部にあるディスプレイデバイスとデータの送受信を行うための回路構成を有する。

【 0 2 3 1 】

電源回路 1 6 0 は、集積回路 3 9 0 内で用いる電圧を生成するための回路である。例えば、OSトランジスタのバックゲートに与える、電気的特性を安定化するための負電圧を生成する回路である。

【 0 2 3 2 】

低速バス 1 4 0 B は、割り込み制御回路 1 5 1、インターフェース回路 1 5 2、バッテリー制御回路 1 5 3、およびADC/DACインターフェース回路 1 5 4の間の各種信号を低速で送受信するためのバスである。一例としては、AMBA-APB(Advanced Peripheral Bus)をバスとして用いることができる。高速バス 1 4 0 Aと低速バス 1 4 0 Bとの間の各種信号の送受信は、ブリッジ回路 1 5 0 を介して行う。

【 0 2 3 3 】

割り込み制御回路 1 5 1 は、周辺機器から受け取る要求に対して、割り込み処理を行うための回路構成を有する。

【 0 2 3 4 】

インターフェース回路 1 5 2 は、UART(Universal Asynchronous Receiver/Transmitter)や、I2C(Integrated Circuit)、SPI(Serial Peripheral Interface)などのインターフェースを機能させるための回路構成を有する。

【 0 2 3 5 】

バッテリー制御回路 1 5 3 は、集積回路 3 9 0 の外部にあるバッテリーの充放電に関するデータを送受信するための回路構成を有する。

【 0 2 3 6 】

ADC/DACインターフェース回路 1 5 4 は、集積回路 3 9 0 の外部にあるMEMS(Micro Electro Mechanical Systems)デバイス等のアナログ信号を出力するデバイスとの間でデータを送受信するための回路構成を有する。

【 0 2 3 7 】

図 3 4 A、図 3 4 B は、SoC化した際の回路ブロックの配置の一例を示す図である。図 3 4 A に図示する集積回路 3 9 0 のように図 3 3 のブロック図で図示した各構成は、チップ上で領域を区切って配置することができる。

【 0 2 3 8 】

なお図 3 3 で説明したオンチップメモリ 1 3 1 は、OSトランジスタで構成される記憶回路、例えばNOSRAM等で構成することができる。つまりオンチップメモリ 1 3 1 とメモリ回路 1 2 1 とは、同じ回路構成を有する。そのため、SoC化した際、図 3 4 B に図示する集積回路 3 9 0 E のようにオンチップメモリ 1 3 1 とメモリ回路 1 2 1 とを一体化して同じ領域内に配置することも可能である。

【 0 2 3 9 】

以上説明した本発明の一態様により、新規な半導体装置および電子機器を提供することができる。又は、本発明の一態様により、消費電力の小さい半導体装置および電子機器を提供することができる。又は、本発明の一態様により、発熱の抑制が可能な半導体装置および電子機器を提供することができる。

【 0 2 4 0 】

本実施の形態は、他の実施の形態の記載と適宜組み合わせることができる。

【0241】

(実施の形態7)

本実施の形態では、上記実施の形態で説明した集積回路390を適用することが可能な電子機器、移動体、演算システムについて、図35乃至図38を参照しながら説明する。

【0242】

図35Aは、移動体の一例として自動車の外観図を図示している。図35Bは、自動車内でのデータのやり取りを簡略化した図である。自動車590は、複数のカメラ591等を有する。また、自動車590は、赤外線レーダー、ミリ波レーダー、レーザーレーダーなど各種センサ(図示せず)などを備える。

10

【0243】

自動車590において、カメラ591等に上記集積回路390を用いることができる。自動車590は、カメラ591が複数の撮像方向592で得られた複数の画像を上記実施の形態で説明した集積回路390で処理し、バス593等を介してホストコントローラ594等により複数の画像をまとめて解析することで、ガードレールや歩行者の有無など、周囲の交通状況を判断し、自動運転を行うことができる。また、道路案内、危険予測などを行うシステムに用いることができる。

【0244】

集積回路390では、得られた画像データをニューラルネットワークなどの演算処理を行うことで、例えば、画像の高解像度化、画像ノイズの低減、顔認識(防犯目的など)、物体認識(自動運転の目的など)、画像圧縮、画像補正(広ダイナミックレンジ化)、レンズレスイメージセンサの画像復元、位置決め、文字認識、反射映り込み低減などの処理を行うことができる。

20

【0245】

なお、上述では、移動体の一例として自動車について説明しているが、移動体は自動車に限定されない。例えば、移動体としては、電車、モノレール、船、飛行体(ヘリコプター、無人航空機(ドローン)、飛行機、ロケット)なども挙げることができ、これらの移動体に本発明の一態様の半導体装置を適用して、人工知能を利用したシステムを付与することができる。

【0246】

30

図36Aは、携帯型電子機器の一例を示す外観図である。図36Bは、携帯型電子機器内でのデータのやり取りを簡略化した図である。携帯型電子機器595は、プリント配線基板596、スピーカー597、カメラ598、マイクロフォン599等を有する。

【0247】

携帯型電子機器595において、プリント配線基板596に上記集積回路390を設けることができる。携帯型電子機器595は、スピーカー597、カメラ598、マイクロフォン599等で得られる複数のデータを上記実施の形態で説明した集積回路390を用いて処理・解析することで、ユーザの利便性を向上させることができる。また、音声案内、画像検索などを行うシステムに用いることができる。

【0248】

40

集積回路390では、得られた画像データをニューラルネットワークなどの演算処理を行うことで、例えば、画像の高解像度化、画像ノイズの低減、顔認識(防犯目的など)、物体認識(自動運転の目的など)、画像圧縮、画像補正(広ダイナミックレンジ化)、レンズレスイメージセンサの画像復元、位置決め、文字認識、反射映り込み低減などの処理を行うことができる。

【0249】

図37Aに示す携帯型ゲーム機1100は、筐体1101、筐体1102、筐体1103、表示部1104、接続部1105、操作キー1107等を有する。筐体1101、筐体1102および筐体1103は、取り外すことが可能である。筐体1101に設けられている接続部1105を筐体1108に取り付けることで、表示部1104に出力される

50

映像を、別の映像機器に出力することができる。他方、筐体 1 1 0 2 および筐体 1 1 0 3 を筐体 1 1 0 9 に取り付けすることで、筐体 1 1 0 2 および筐体 1 1 0 3 を一体化し、操作部として機能させる。筐体 1 1 0 2 および筐体 1 1 0 3 の基板に設けられているチップなどに先の実施の形態に示す集積回路 3 9 0 を組み込むことができる。

【 0 2 5 0 】

図 3 7 B は U S B 接続タイプのスティック型の電子機器 1 1 2 0 である。電子機器 1 1 2 0 は、筐体 1 1 2 1、キャップ 1 1 2 2、U S B コネクタ 1 1 2 3 および基板 1 1 2 4 を有する。基板 1 1 2 4 は、筐体 1 1 2 1 に収納されている。例えば、基板 1 1 2 4 には、メモリチップ 1 1 2 5、コントローラチップ 1 1 2 6 が取り付けられている。基板 1 1 2 4 のコントローラチップ 1 1 2 6 などに先の実施の形態に示す集積回路 3 9 0 を組み込むことができる。

10

【 0 2 5 1 】

図 3 7 C は人型のロボット 1 1 3 0 である。ロボット 1 1 3 0 は、センサ 2 1 0 1 乃至 2 1 0 6、および制御回路 2 1 1 0 を有する。例えば、制御回路 2 1 1 0 には、先の実施の形態に示す集積回路 3 9 0 を組み込むことができる。

【 0 2 5 2 】

上記実施の形態で説明した集積回路 3 9 0 は、電子機器に内蔵する代わりに、電子機器と通信を行うサーバーに用いることもできる。この場合、電子機器とサーバーによって演算システムが構成される。図 3 8 に、システム 3 0 0 0 の構成例を示す。

【 0 2 5 3 】

システム 3 0 0 0 は、電子機器 3 0 0 1 と、サーバー 3 0 0 2 によって構成される。電子機器 3 0 0 1 とサーバー 3 0 0 2 間の通信は、インターネット回線 3 0 0 3 を介して行うことができる。

20

【 0 2 5 4 】

サーバー 3 0 0 2 には、複数のラック 3 0 0 4 を有する。複数のラックには、複数の基板 3 0 0 5 が設けられ、当該基板 3 0 0 5 上に上記実施の形態で説明した集積回路 3 9 0 を搭載することができる。これにより、サーバー 3 0 0 2 にニューラルネットワークが構成される。そして、サーバー 3 0 0 2 は、電子機器 3 0 0 1 からインターネット回線 3 0 0 3 を介して入力されたデータを用いて、ニューラルネットワークの演算を行うことができる。サーバー 3 0 0 2 による演算の結果は必要に応じて、インターネット回線 3 0 0 3 を介して電子機器 3 0 0 1 に送信することができる。これにより、電子機器 3 0 0 1 における演算の負担を低減することができる。

30

【 0 2 5 5 】

本実施の形態は、他の実施の形態の記載と適宜組み合わせることができる。

【 0 2 5 6 】

(本明細書等の記載に関する付記)

以上の実施の形態、および実施の形態における各構成の説明について、以下に付記する。

【 0 2 5 7 】

各実施の形態に示す構成は、他の実施の形態あるいは実施例に示す構成と適宜組み合わせ、本発明の一態様とすることができる。また、1つの実施の形態の中に、複数の構成例が示される場合は、構成例を適宜組み合わせることが可能である。

40

【 0 2 5 8 】

なお、ある一つの実施の形態の中で述べる内容（一部の内容でもよい）は、その実施の形態で述べる別の内容（一部の内容でもよい）、および／または、一つ若しくは複数の別の実施の形態で述べる内容（一部の内容でもよい）に対して、適用、組み合わせ、または置き換えなどを行うことが出来る。

【 0 2 5 9 】

なお、実施の形態の中で述べる内容とは、各々の実施の形態において、様々な図を用いて述べる内容、または明細書に記載される文章を用いて述べる内容のことである。

【 0 2 6 0 】

50

なお、ある一つの実施の形態において述べる図（一部でもよい）は、その図の別の部分、その実施の形態において述べる別の図（一部でもよい）、および／または、一つ若しくは複数の別の実施の形態において述べる図（一部でもよい）に対して、組み合わせることにより、さらに多くの図を構成させることができる。

【0261】

また本明細書等において、ブロック図では、構成要素を機能毎に分類し、互いに独立したブロックとして示している。しかしながら実際の回路等においては、構成要素を機能毎に切り分けることが難しく、一つの回路に複数の機能が係わる場合や、複数の回路にわたって一つの機能が関わる場合があり得る。そのため、ブロック図のブロックは、明細書で説明した構成要素に限定されず、状況に応じて適切に言い換えることができる。

10

【0262】

また、図面において、大きさ、層の厚さ、または領域は、説明の便宜上任意の大きさに示したものである。よって、必ずしもそのスケールに限定されない。なお図面は明確性を期すために模式的に示したものであり、図面に示す形状または値などに限定されない。例えば、ノイズによる信号、電圧、若しくは電流のばらつき、または、タイミングのずれによる信号、電圧、若しくは電流のばらつきなどを含むことが可能である。

【0263】

また、図面等において図示する構成要素の位置関係は、相対的である。従って、図面を参照して構成要素を説明する場合、位置関係を示す「上に」、「下に」等の語句は便宜的に用いられる場合がある。構成要素の位置関係は、本明細書の記載内容に限定されず、状況に応じて適切に言い換えることができる。

20

【0264】

本明細書等において、トランジスタの接続関係を説明する際、「ソースまたはドレインの一方」（または第1電極、または第1端子）、「ソースまたはドレインの他方」（または第2電極、または第2端子）という表記を用いる。これは、トランジスタのソースとドレインは、トランジスタの構造または動作条件等によって変わるためである。なおトランジスタのソースとドレインの呼称については、ソース（ドレイン）端子や、ソース（ドレイン）電極等、状況に応じて適切に言い換えることができる。

【0265】

また、本明細書等において「電極」や「配線」の用語は、これらの構成要素を機能的に限定するものではない。例えば、「電極」は「配線」の一部として用いられることがあり、その逆もまた同様である。さらに、「電極」や「配線」の用語は、複数の「電極」や「配線」が一体となって形成されている場合なども含む。

30

【0266】

また、本明細書等において、電圧と電位は、適宜言い換えることができる。電圧は、基準となる電位からの電位差のことであり、例えば基準となる電位をグラウンド電圧（接地電圧）とすると、電圧を電位に言い換えることができる。グラウンド電位は必ずしも0Vを意味するとは限らない。なお電位は相対的なものであり、基準となる電位によっては、配線等に与える電位を変化させる場合がある。

【0267】

40

また本明細書等において、ノードは、回路構成やデバイス構造等に応じて、端子、配線、電極、導電層、導電体、不純物領域等と言い換えることが可能である。また、端子、配線等をノードと言い換えることが可能である。

【0268】

本明細書等において、AとBとが接続されている、とは、AとBとが電氣的に接続されているものをいう。ここで、AとBとが電氣的に接続されているとは、AとBとの間で対象物（スイッチ、トランジスタ素子、またはダイオード等の素子、あるいは当該素子および配線を含む回路等を指す）が存在する場合にAとBとの電気信号の伝達が可能である接続をいう。なおAとBとが電氣的に接続されている場合には、AとBとが直接接続されている場合を含む。ここで、AとBとが直接接続されているとは、上記対象物を介すること

50

なく、AとBとの間で配線（または電極）等を介してAとBとの電気信号の伝達が可能である接続をいう。換言すれば、直接接続とは、等価回路で表した際に同じ回路図として見なせる接続をいう。

【0269】

本明細書等において、スイッチとは、導通状態（オン状態）、または、非導通状態（オフ状態）になり、電流を流すか流さないかを制御する機能を有するものをいう。または、スイッチとは、電流を流す経路を選択して切り替える機能を有するものをいう。

【0270】

本明細書等において、チャンネル長とは、例えば、トランジスタの上面図において、半導体（またはトランジスタがオン状態のときに半導体の中で電流の流れる部分）とゲートとが重なる領域、またはチャンネルが形成される領域における、ソースとドレインとの間の距離をいう。

10

【0271】

本明細書等において、チャンネル幅とは、例えば、半導体（またはトランジスタがオン状態のときに半導体の中で電流の流れる部分）とゲート電極とが重なる領域、またはチャンネルが形成される領域における、ソースとドレインとが向かい合っている部分の長さをいう。

【0272】

なお本明細書等において、「膜」、「層」などの語句は、場合によっては、または、状況に応じて、互いに入れ替えることが可能である。例えば、「導電層」という用語を、「導電膜」という用語に変更することが可能な場合がある。または、例えば、「絶縁膜」という用語を、「絶縁層」という用語に変更することが可能な場合がある。

20

【実施例】

【0273】

本発明の一態様に係る半導体装置の一例として、チャンネルが形成される半導体層にIn-Ga-Zn酸化物を用いたトランジスタ（「IGZO-FET」ともいう。）とSiトランジスタ（「Si-FET」ともいう。）を用いたBinary AI Processorを作製した。本実施例では、作製したBinary AI Processorの構成、および動作のシミュレーション結果について説明する。作製したBinary AI Processorは、後述するNoFFコンピューティング可能な半導体装置である。

【0274】

30

近年、IoT（Internet of Things）およびAIなどの技術が注目されている。IoT分野で使用される機器（IoT機器）では、消費電力の低減が求められる一方で、AI処理時は演算性能の高さが求められる。

【0275】

消費電力の低減を目的として、待機状態の回路への電源供給を遮断するパワーゲーティング（PG）技術が知られている。また、IoT機器などの低消費電力化を実現する技術として、PG技術にメモリを組み合わせたNormally-off（NoFF）コンピューティングが提案されている。

【0276】

NoFFコンピューティングでは、システム全体としては動作しているが、一時的に動作不要となる回路に対して当該回路のデータをメモリに退避させた後に、当該回路への電源供給を遮断する動作が行われる。NoFFコンピューティングに用いるメモリとして、ReRAM（抵抗変化型メモリ）、MRAM（磁気メモリ）、PCM（相変化メモリ）などの不揮発性メモリが検討されている。

40

【0277】

OSメモリは、ReRAM、MRAM、およびPCMよりもデータ書き込み時のエネルギー消費が少ないため、NoFFコンピューティングに用いるメモリとして好適である。なお、OSTランジスタは、ReRAM、MRAM、およびPCMなどに用いることも可能である。

【0278】

50

作製した Binary AI Processor Chip (以下、「BAP900」ともいう。)は、130 nm Si CMOSプロセスで形成された演算器 (PE: Processing Element) と PE 上に 60 nm IGZO プロセスで形成された OS メモリを含む。

【0279】

また、BAP900は、IGZO-FETを用いたOSメモリをAI Acceleratorの重みパラメータを格納するメモリ (W-MEM) として使用し、当該メモリの読み出し線を演算器と直結した構成を有する。

【0280】

図39Aに作製したBAP900の外観写真を示す。図39Bに、BAP900の一部を拡大した断面TEM写真を示す。BAP900は、層M1乃至層M8を有する。なお、層M1乃至層M8は、配線または電極などの導電体を含む層である。図39Bより、Si-FETの上方に、IGZO-FETおよびMIM (Metal-Insulator-Metal) 構造の容量 (MIM-Capacitor) が設けられていることがわかる。表8にBAP900の主な仕様を示す。

【0281】

【表8】

Technology		130 nm Si CMOS, 60 nm IGZO-FET (BEOL)
Supply voltage		1.2 V, 3.3 V
CPU		ARM Cortex-M0
AI Accelerator Subsystem	Weight memory	32 KB
	# of PE	128
Scratchpad memory		16 KB

【0282】

BAP900は、回路部901乃至回路部905を有する。回路部901は、32 bitのARM Cortex-M0 CPUと、その周辺回路 (Peripherals) とを含む。回路部902は、AI Accelerator Control Logicを含む。回路部903は、PEアレイ上に設けられた、IGZOプロセスで形成された32 KBのW-MEMを含む (IGZO-based W-MEM (32 KB) on PE Array)。回路部904は、16 KBのScratchpad memoryを含む。回路部905は、Power Switchesを含む。

【0283】

図40は、BAP900の詳細なシステム構成を説明するブロック図である。BAP900は、Cortex-M0サブシステム (Cortex-M0 Subsystem)、AI Acceleratorサブシステム (AI Accelerator Subsystem)、およびCortex-M0サブシステムよりも動作周波数が低い周辺回路 (Low-BW (Band Width) Peripherals) を含む。

【0284】

Cortex-M0サブシステムは、32 bitのARM Cortex-M0 CPU、電源管理ユニット (PMU: Power Management Unit)、2つのGPIO (General purpose input/output)、SYSCTL、記憶容量16 KByteの組み込みIGZOスクラッチメモリ、UARTs (Universal Asynchronous Receiver/Transmitter)、および外部メモリインターフェイス (Ext-MEM IF) を含む。それぞれは、32ビットのAHBバスライン (32b AHB) を介して接続される。

【0285】

AI Acceleratorサブシステムは、AI Accelerator制御回路(AI Accelerator Control Logic)、PEアレイ(PE Array)、およびPEアレイ上に設けられた記憶容量32KbyteのW-MEMを含む。PEアレイは、128個のPEを含む。

【0286】

Low-BW Peripheralsは、パワースイッチ(Power Switches)、SPI(Serial Peripheral Interface)、タイマー(Timers)、Watch dog、およびUARTsを含む。パワースイッチ、SPI、タイマー、Watch dog、およびUARTsは、32ビットのAPBバスライン(32b APB)を介して接続される。パワースイッチは、Cortex-M0サブシステムへの電力供給を制御する機能を有する。

10

【0287】

また、BAP900は、OSCノード、GPIOノード、VDDsノード、Sensorノード、RTCノード、USBノード、およびExt-MEMノードを有する。これらのノードを介して信号の入出力などが行われる。例えば、OSCノードを介して外部からクロック信号(Clock)が入力される。なお、図40に記す「M」はMasterを示し、「S」はSlaveを示している。

【0288】

電源ドメインは、VDDsノードを介して外部から常時供給される電源VDDsと、PG可能な電源PGVDDsの2つがある。PMUは、動作モードに応じて電力供給を制御する機能を有する。待機モードで動作する場合、PMUはPG可能な回路に対してPGを行うことで消費電力を削減する。AI処理(積和演算処理)を行う際にAI Acceleratorサブシステムを用いることで、CPUによる演算よりも高速かつ高効率にAI処理を行うことができる。

20

【0289】

BAP900はPGが可能なため、AI処理を行わない期間は、システム全体として消費電力を低減できる。一方で、Sensorノードからの信号入力が発生すると、元のシステム状態を瞬間に復元し、すぐにAI処理を実行できる。

【0290】

図41Aに、W-MEMに含まれるメモリセル910の回路図を示す。メモリセル910は、3つのIGZO-FETと1つの容量を含むメモリセルである。当該容量はMIM(Metal-Insulator-Metal)構造の容量である。メモリセル910の電源電圧は3.3Vである。メモリセル910はノードSNに電荷を保持するメモリのため、電源遮断時にもデータが消失しない。

30

【0291】

1つのメモリセル910で、1ビットの重み情報Wを保持する。重み情報Wは、配線WBLを介してノードSNに書き込まれる。ノードSNに書き込まれた重み情報Wは、配線RBLを介して読み出される。8つのメモリセルを用いることで、重み情報W[0]から重み情報W[7]で構成される8ビットの重み情報W(「W[7:0]」ともいう。)を保持することができる。

40

【0292】

図41Bは、図41Aに示したメモリセルの動作例を示すタイミングチャートである。図41Bは、書き込みモード(Write)、待機モード(Sleep)、および読み出しモード(Read)における、配線WBL、配線WWL、配線RBL、および配線RWLの電位変化を示している。配線WBLには1.7Vまたは0Vが供給される。配線WWLには3.3Vまたは0Vが供給される。配線RBLには1.0Vまたは0Vが供給される。配線RWLには1.2Vまたは0Vが供給される。また、図41Bでは、ノードSNにdata0を供給する場合の電位変化と、data1を供給する場合の電位変化を示している。

50

【0293】

図41Cは、PE920の構成を示すブロック図である。PE920は、電源電圧1.2VのSiロジックセルで作製した。PE920は、センスアンプ921(SA)、乗算回路922(Multiplier)と加算回路923(Adder tree)を含むbinary積和演算器924(MAC)、アキュムレータ925(Accumulator)を含む。アキュムレータ925は、バッチ正規化用の1ビット(1b)しきい値加算器と11ビットのレジスタ(11bit register)を含む。

【0294】

1つのPE920には8本の配線RBLが並列に接続され、8ビットの重み情報W[7:0]が入力される。入力された重み情報W[7:0]は、センスアンプ(SA)で増幅された後に積和演算処理に使用されるか、積和演算処理に使用されず直接読み出される。どちらが行われるかは、Processing/Read selector信号で決定される。積和演算処理に使用される場合は、重みW[7:0]は乗算回路で信号A[7:0]と乗算され、積信号M[7:0]に変換される。直接読み出す場合、信号readout[7:0]として出力される。

【0295】

積信号M[7:0]はAdder tree回路で加算され、積和信号MAに変換される。MAC/BN selector信号によって、積和信号MAおよびしきい値信号THのどちらをアキュムレータに入力するかが決定される。アキュムレータは、11ビットの信号macout[10:0]を出力する機能と、インバータ回路を介してサインビット(Sign bit)信号を出力する機能と、を有する。

【0296】

図42Aは、回路部903(IGZO-based W-MEM(32KB) on PE Array)の構成を示すブロック図である。回路部903は、1つのGlobal Logic回路と4つのSubarray回路(Subarray0乃至Subarray3)を含む。

【0297】

図42Bは、1つのSubarray回路の構成を示すブロック図である。1つのSubarray回路は、回路部931乃至回路部938を含む。回路部931乃至回路部934は、それぞれが128×128個のメモリセル910を含む記憶容量16kビットのメモリセルアレイ(16k array(128×128))として機能する。1つのメモリセルアレイは、128本の配線RBL(読み出しビット線)を含む。また、1本の配線RBLには128個のメモリセルが接続されている。

【0298】

回路部935および回路部936は、それぞれが、16個のPE920および列ドライバを有する(#16 of PE and Shared Column driver)。回路部935が有する列ドライバは、回路部931および回路部933を駆動する。回路部936が有する列ドライバは、回路部932および回路部934を駆動する。回路部937は、回路部931および回路部932を駆動する行ドライバ(Upper Row driver)を含む。回路部938は、回路部933および回路部934を駆動する行ドライバ(Lower Row driver)を含む。

【0299】

よって、回路部903全体では、PEアレイに1024本の配線RBLが並列に接続される。1024本の配線RBLから読み出された情報は並列演算される。また、行ドライバをメモリセルアレイと重ねて設けることで、情報の読み出しエネルギーとチップ面積を削減できる。

【0300】

図43Aに、作製したBAP900の動作期間中に生じる消費電力の推移とPG期間の概念図を示す(This work)。比較例として、図43Bに従来動作(PGを行わない)の動作期間中に生じる消費電力推移の概念図を示す(Conventional

10

20

30

40

50

）。図43A、図43Bともに、縦軸は消費電力（Power）を示し、横軸は経過時間（Time）を示している。

【0301】

従来動作では、電源供給が停止すると重み情報やニューラルネットワーク構造などの情報が消えるため、再起起動時にこれらの情報をROMなどから読み出して、RAMに書き込む必要があった（ROM/RAM access）。よって、演算処理時間の短縮が難しかった。また、演算処理時間に比例して、メモリアクセス時の消費電力およびCPUの消費電力のみでなく、静的消費電力（Static）も増加する。

【0302】

今回作製したBAP900は、センサノードからBAP900の起動信号Rx（sensor raw data）が入力されると起動し、CPUからAI Acceleratorサブシステムへ生データが転送される。生データはAI Acceleratorサブシステムで演算処理され、演算結果が信号Tx（meaningful data）として出力される。その後、PGが行われる。AI Acceleratorサブシステムでは並列処理が行われるため、従来例よりも演算処理時間が短く（high ops）、消費電力も小さい。よって効率の良い演算処理が実現できる（high efficiency）。

【0303】

また、今回作製したBAP900では、PG開始直前に、重み情報などの復帰時に必要な情報が保持される。このような情報の保持にはOSメモリを用いることが好ましい。図44Aおよび図44Bに、BAP900に用いた情報保持回路の一例を示す。

【0304】

図44Aに示す情報保持回路941は、Siプロセス（Si-FET）で作製したスキャンDフリップフロップ941a（Scan DFF）に、IGZO-FET含むOSメモリ941bを組み合わせた構成を有する。

【0305】

スキャンDフリップフロップ941aは、端子CK、端子D、端子SE、および端子Qと電氣的に接続される。また、スキャンDフリップフロップ941aは、IGZO-FETを介して端子Qと電氣的に接続される。OSメモリ941bは端子BK、端子RE、および端子Qと電氣的に接続される。

【0306】

図44Bに示す情報保持回路942は、Si-FET、IGZO-FET、および容量で構成された、1T1C型のスクラッチメモリセル（IGZO-based Scratchpad memory cell）である。情報保持回路942は、配線WWL、配線RWL、配線WBL、配線RBL、および配線SLと電氣的に接続される。

【0307】

情報保持回路941および情報保持回路942ともに、復帰時に必要な情報をノードSNに保持することができる。

【0308】

回路シミュレーションソフトウェアを用いて、BAP900の動作を検証した。回路シミュレーションソフトウェアとして、SILVACO社SmartSpiceを使用した。

【0309】

シミュレーションでは、以下の動作を検証した。はじめに、W-MEMに学習済の重みデータを格納し（Write trained W-MEM）、次に、電源供給を停止した（PG）。続いて、電源供給を再開し、SPIを介して解像度が28×28の2値イメージデータを入力して（Input 28×28 binary image data from SPI）、推論動作を行なった（AI operation）。その後、推論結果をSPIに出力し（Output inference result to SPI）、再び電源供給を停止した。

【0310】

図 4 5 A に、シミュレーション実行後の動作波形の一例を示す。図 4 5 B に、当該シミュレーションで想定した全結合型のニューラルネットワークモデルを示す。シミュレーションで想定したニューラルネットワークモデルでは、784 個のニューロンを有する入力層と、10 個のニューロンを有する出力層の間に、128 個のニューロンを有する隠れ層を 3 層設定した。なお、図 4 5 A および図 4 5 B では、入力層と 1 つめの隠れ層の全結合を F C 1、1 つめの隠れ層と 2 つめの隠れ層の全結合を F C 2、2 つめの隠れ層と 3 つめの隠れ層の全結合を F C 3、3 つめの隠れ層と出力層の全結合を F C 4 と示している。シミュレーションによって、前述した動作が問題なく行われることが確認できた。

【 0 3 1 1 】

シミュレーションから見積もられた演算効率、消費エネルギーなどを表 9 に示す。

10

【 0 3 1 2 】

【表 9】

Memory Read energy : W-MEM cell	1.4pJ/bit
Computation energy : 1PE	2.46pJ
Energy efficiency (w/ AI Accelerator control Logic)	0.54TOPS/W
Performance	0.82GOPS
W-MEM retention time	> 1 hour @85°C

20

【 0 3 1 3 】

シミュレーションによって、AI Accelerator サブシステムを使用した BAP900 の演算性能が 0.82GOPS であることがわかった。本実施例では記載しないが、別途シミュレーションを行った AI Accelerator サブシステムを使用しない場合の BAP900 の演算性能と比較したところ、約 215 倍の演算性能が得られていた。また、演算効率は 0.54TOPS/W であることがわかった。

【 0 3 1 4 】

IGZO-FET は、極低消費電力、高速復帰が要求されるイベントドリブンシステムと相性がよく、IoT 機器や末端機器での AI アプリケーションに好適に用いることができる。

30

【符号の説明】

【 0 3 1 5 】

10 : CPU、20 : アクセラレータ、21 : 演算処理部、22 : メモリ部、23 : 演算回路、24 : メモリ回路、29 : 半導体層、30 : バス、31 : 配線、100 : 半導体装置

40

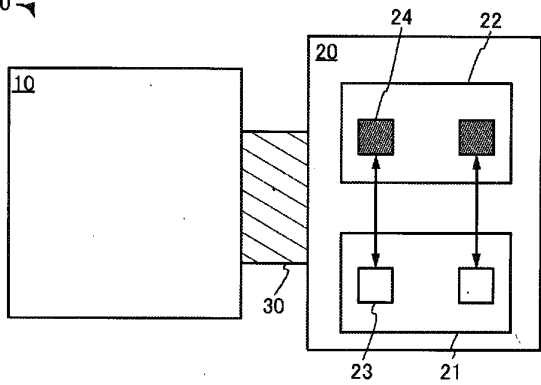
50

【図面】

【図 1 A】

図1A

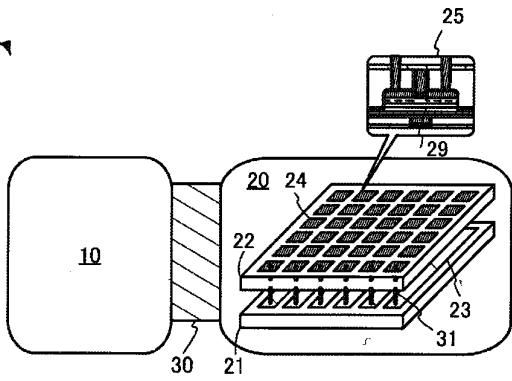
100



【図 1 B】

図1B

100

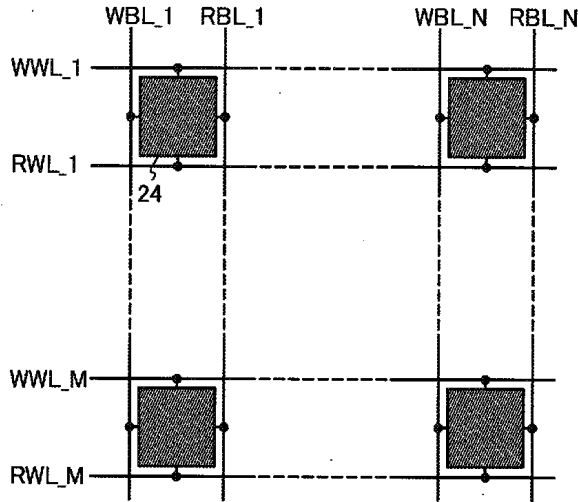


10

【図 2 A】

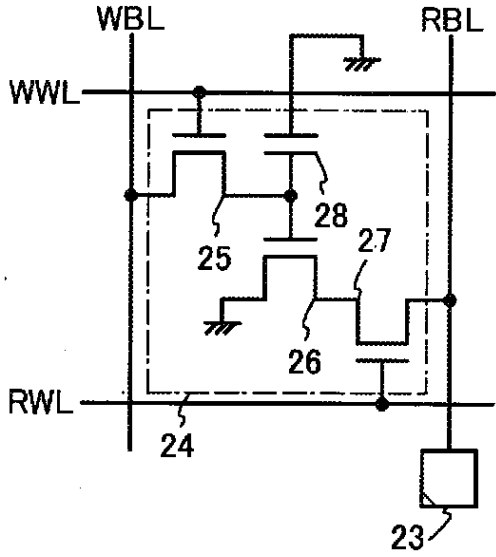
図2A

22



【図 2 B】

図2B



20

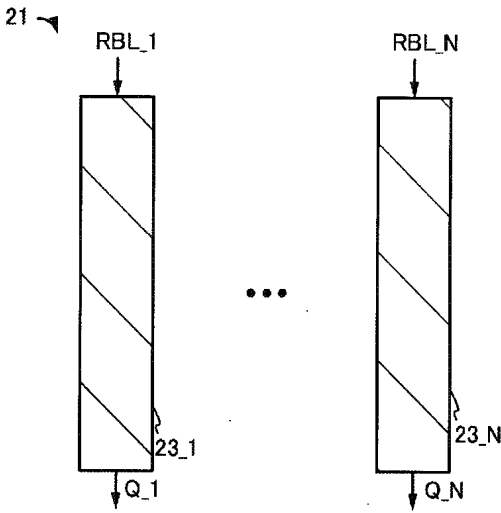
30

40

50

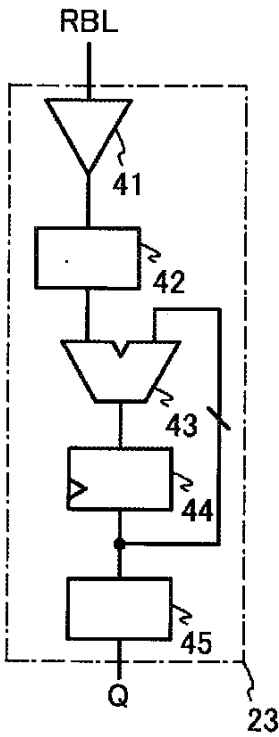
【図 3 A】

図3A

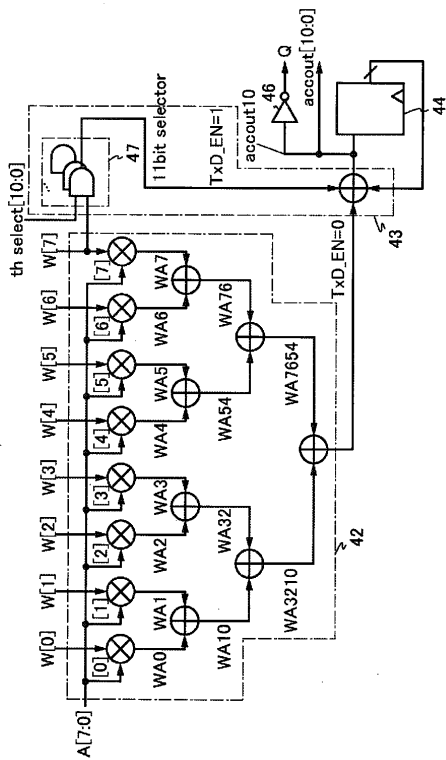


【図 3 B】

図3B



【図 4】



【図 5 A】

図5A

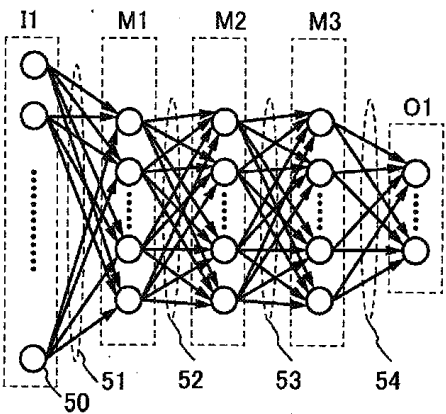


図4

10

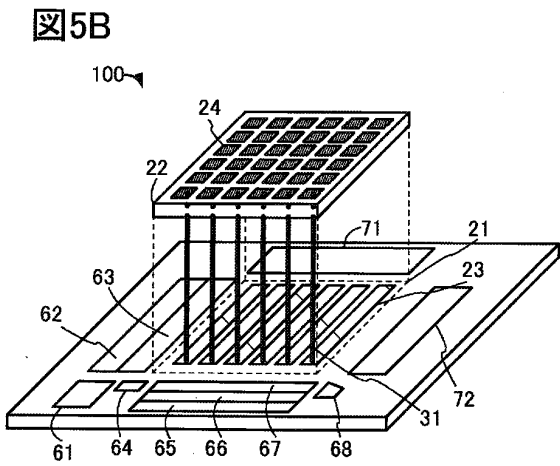
20

30

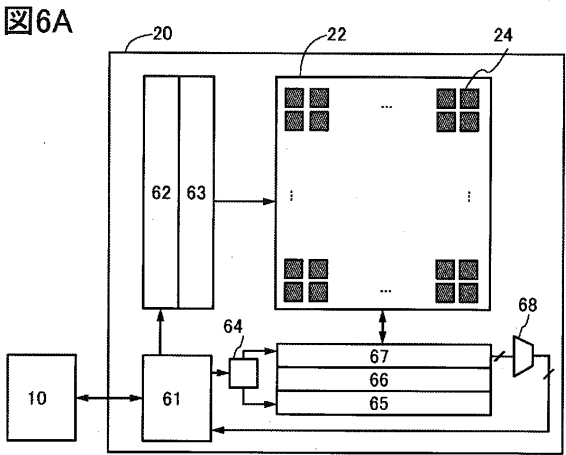
40

50

【図 5 B】

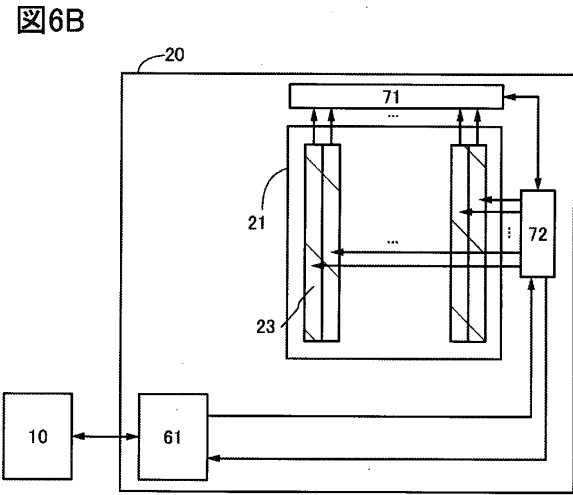


【図 6 A】

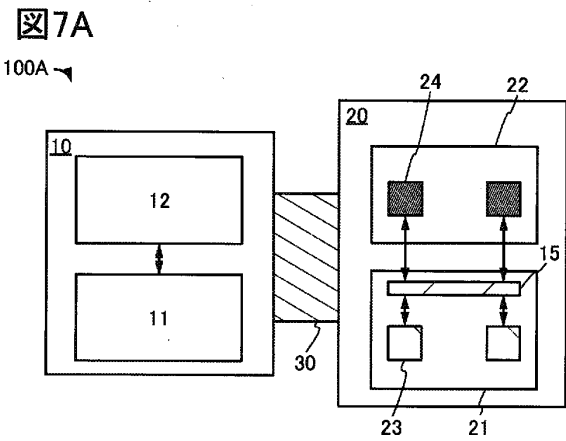


10

【図 6 B】



【図 7 A】



20

30

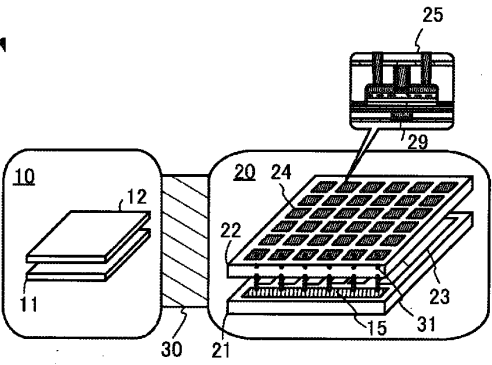
40

50

【図 7 B】

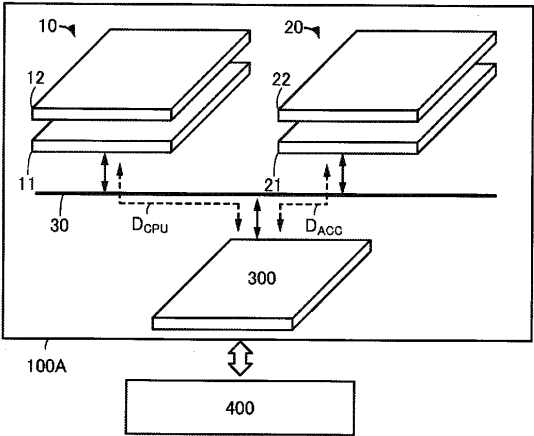
図7B

100A



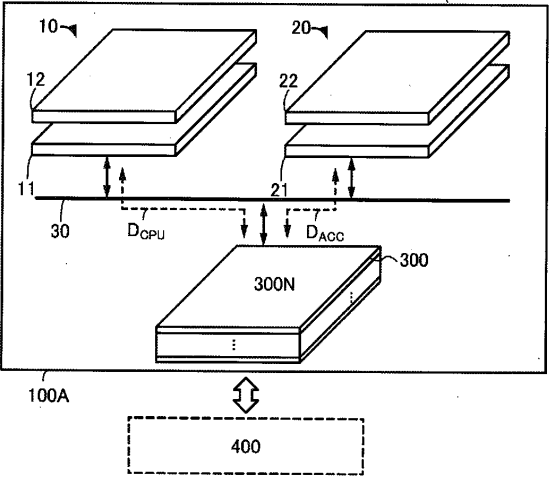
【図 8 A】

図8A



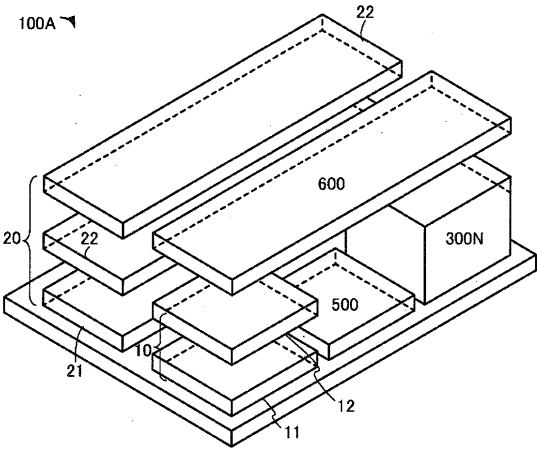
【図 8 B】

図8B



【図 9】

図9



10

20

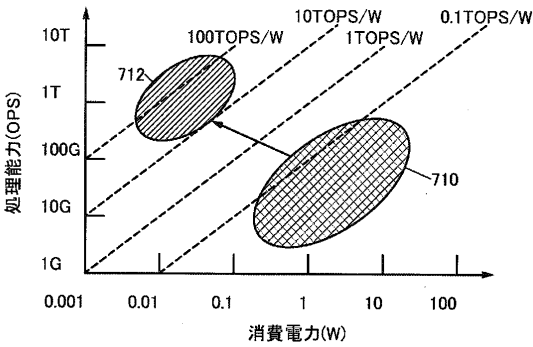
30

40

50

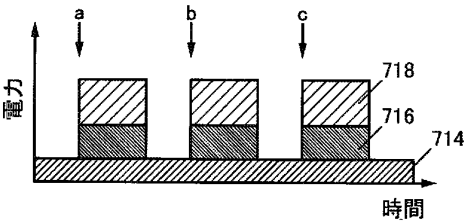
【図10A】

図10A



【図10B】

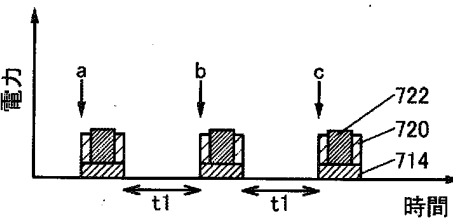
図10B



10

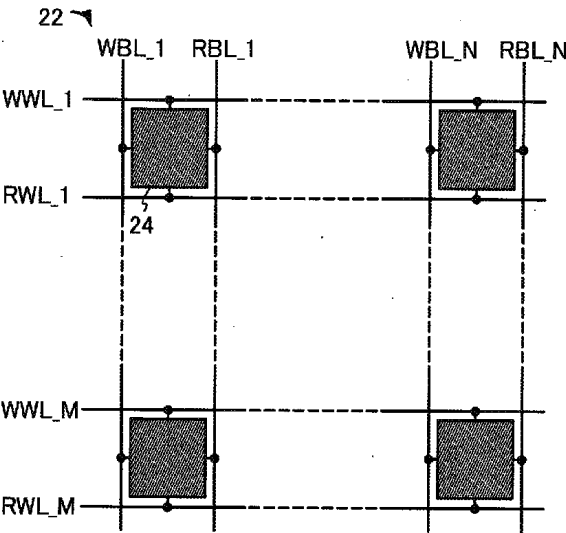
【図10C】

図10C



【図11A】

図11A



20

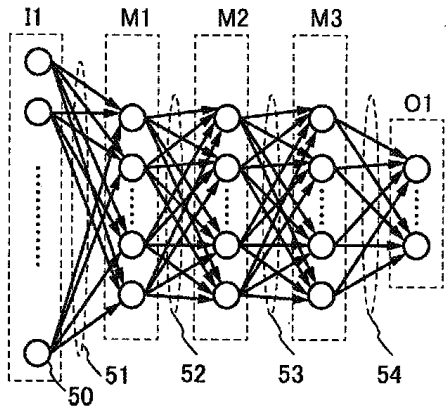
30

40

50

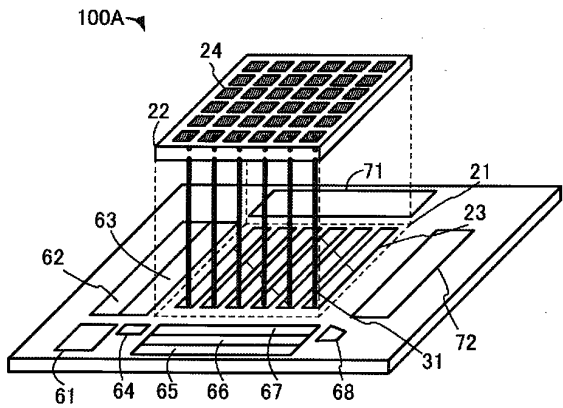
【図14A】

図14A



【図14B】

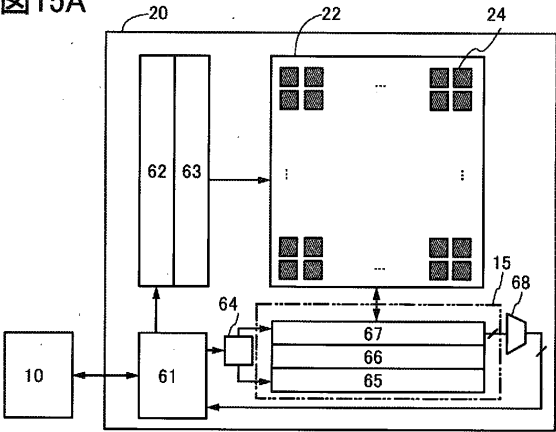
図14B



10

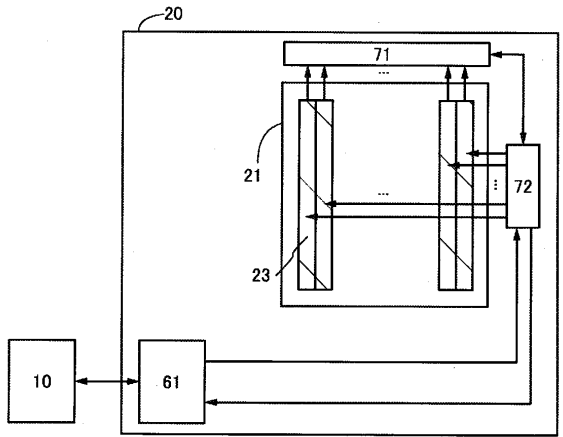
【図15A】

図15A



【図15B】

図15B



20

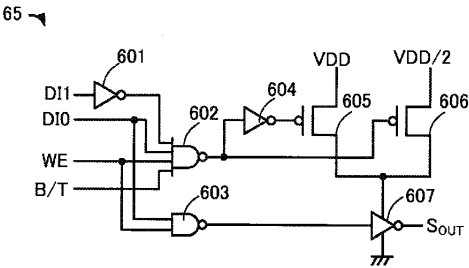
30

40

50

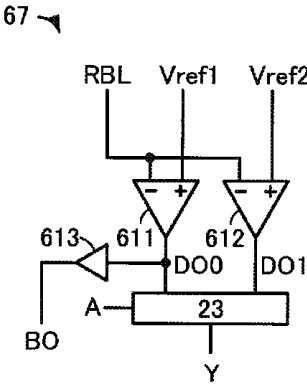
【図16】

図16



【図17】

図17



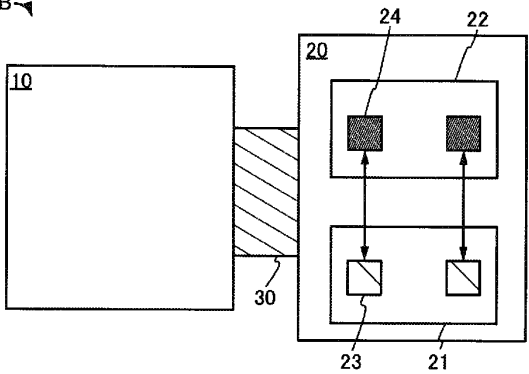
10

20

【図18A】

図18A

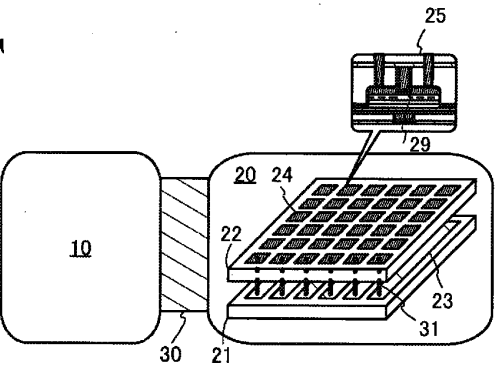
100B ↗



【図18B】

図18B

100B ↗



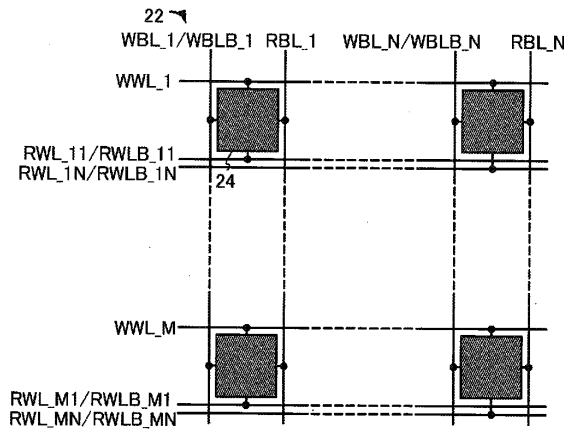
30

40

50

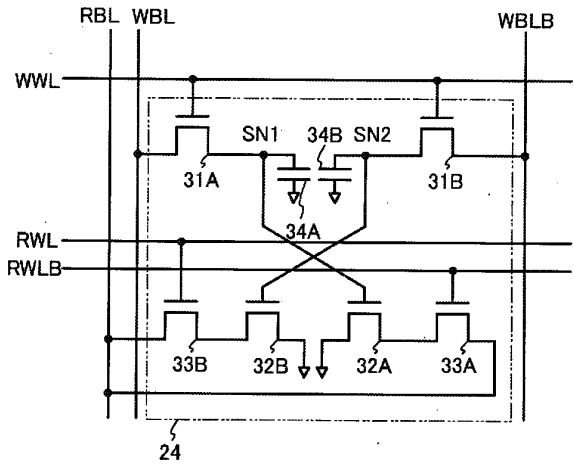
【図 19 A】

図19A



【図 19 B】

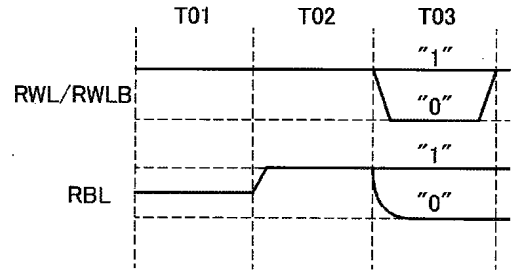
図19B



10

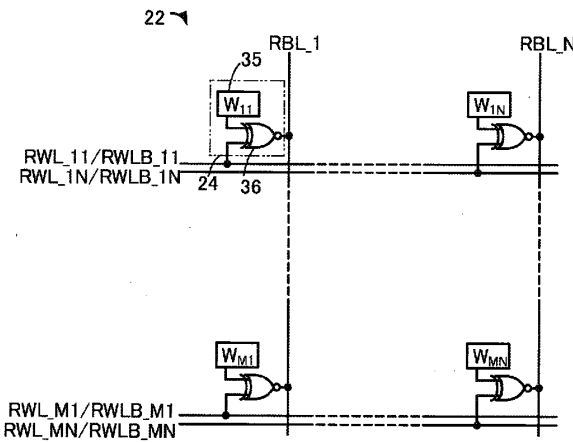
【図 20 A】

図20A



【図 20 B】

図20B



20

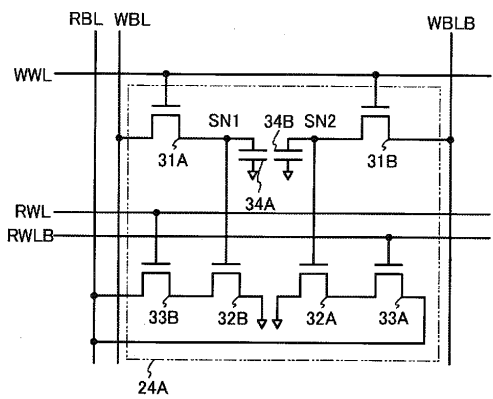
30

40

50

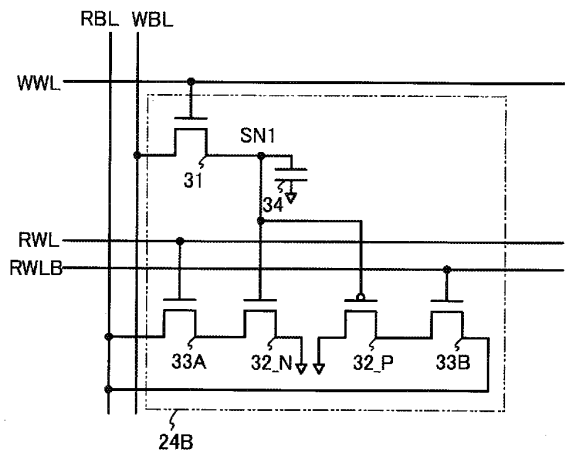
【図 2 1】

図21



【図 2 2 A】

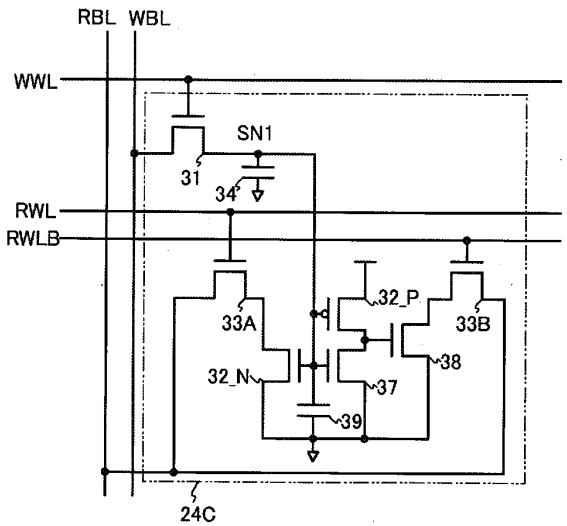
図22A



10

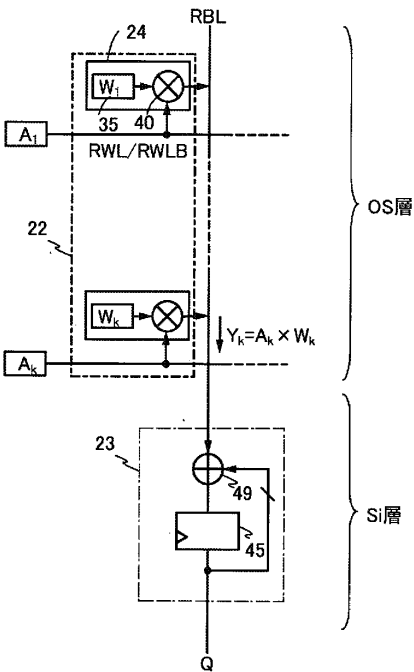
【図 2 2 B】

図22B



【図 2 3】

図23



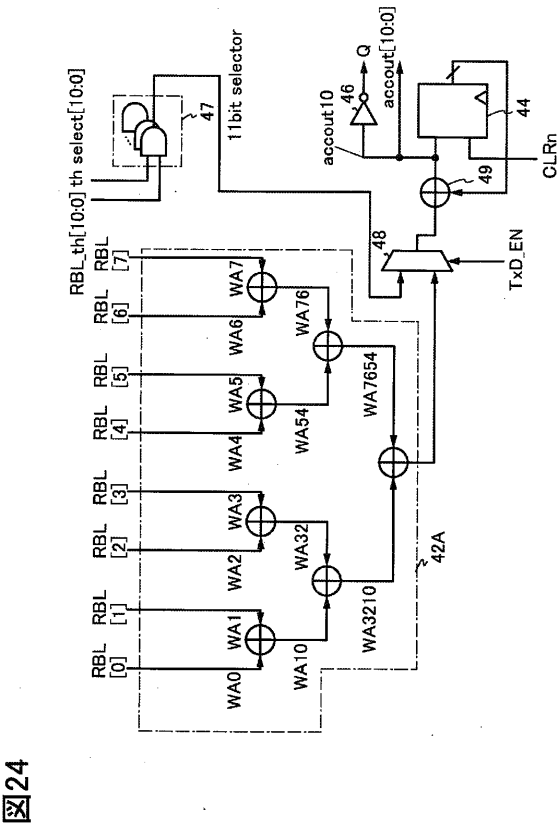
20

30

40

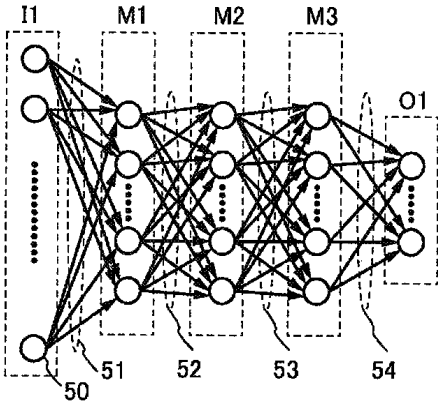
50

【図 2 4】



【図 2 5 A】

図25A

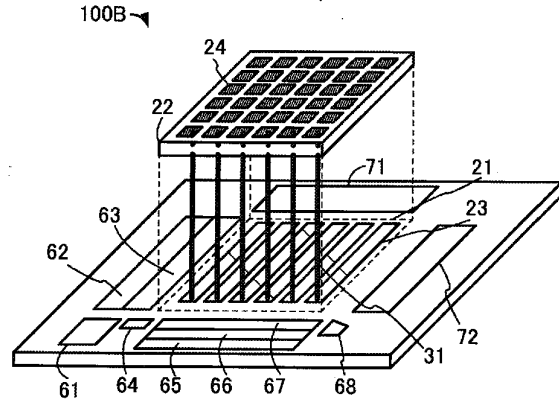


10

20

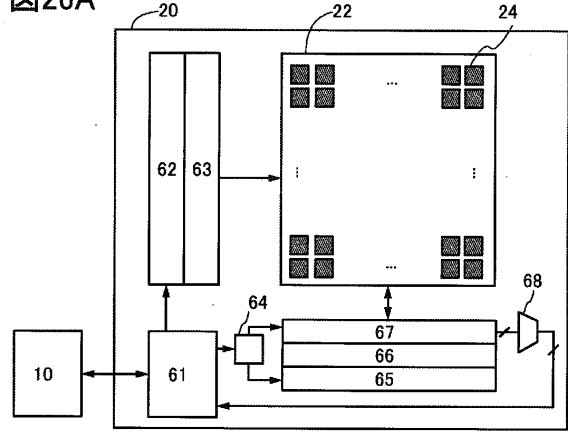
【図 2 5 B】

図25B



【図 2 6 A】

図26A



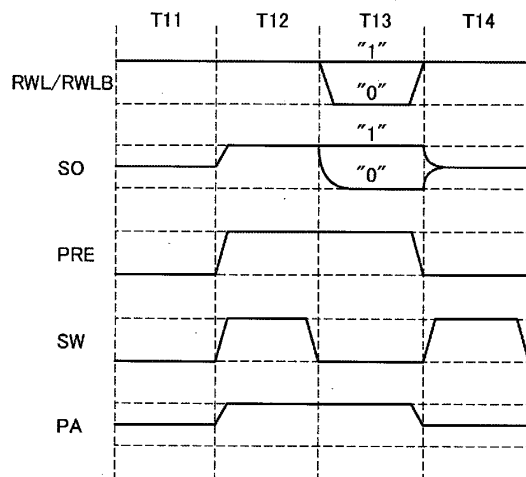
30

40

50

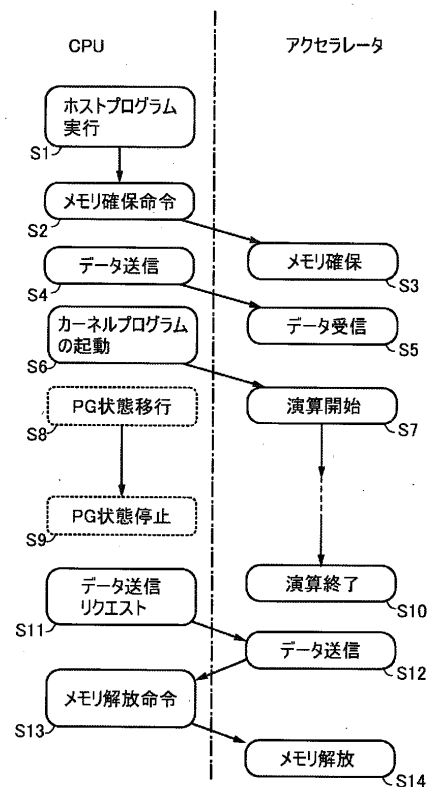
【図 28 B】

図28B



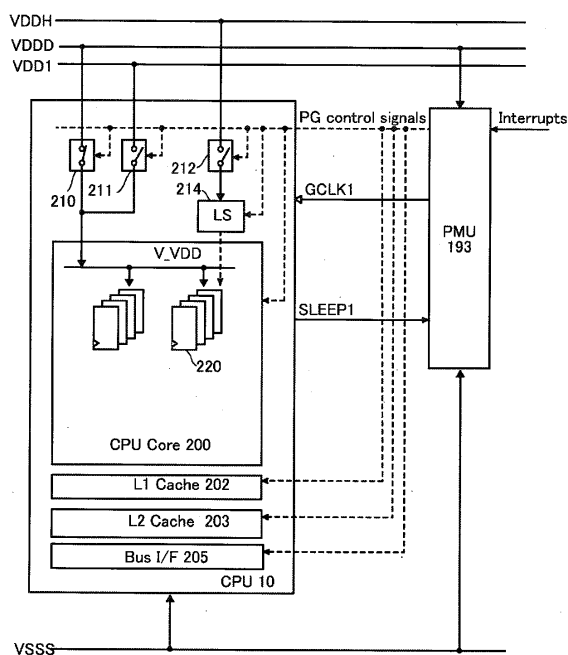
【図 29】

図29



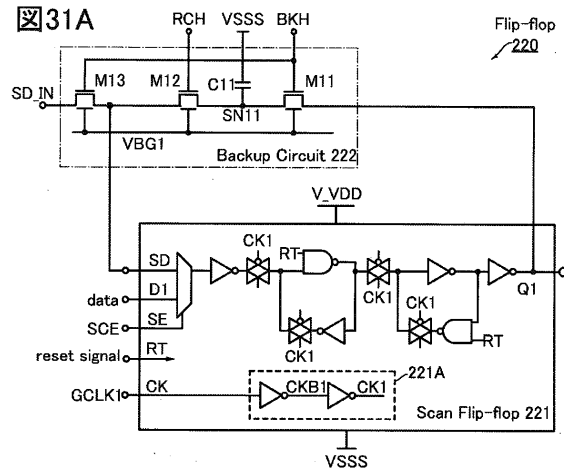
【図 30】

図30



【図 31 A】

図31A



10

20

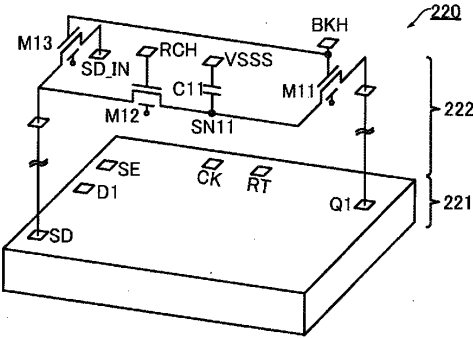
30

40

50

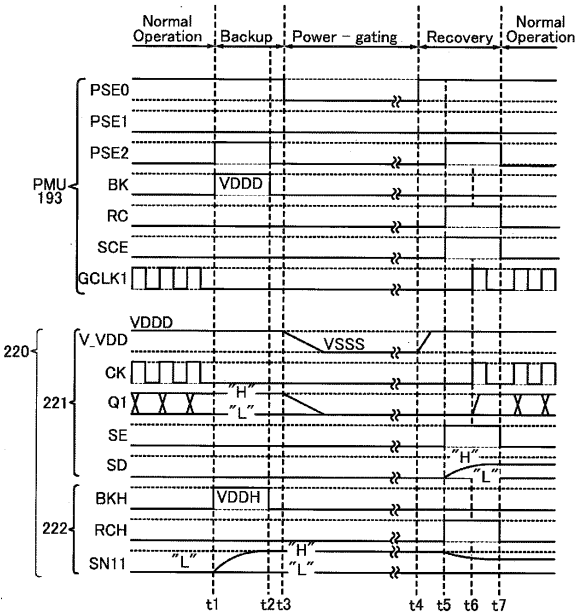
【図 3 1 B】

図31B



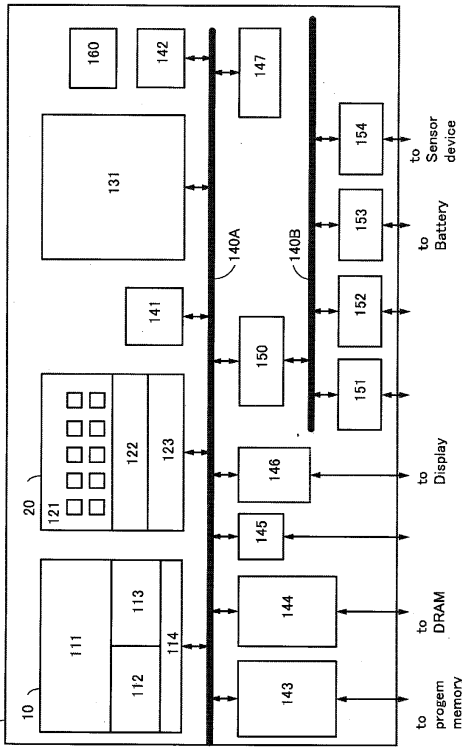
【図 3 2】

図32



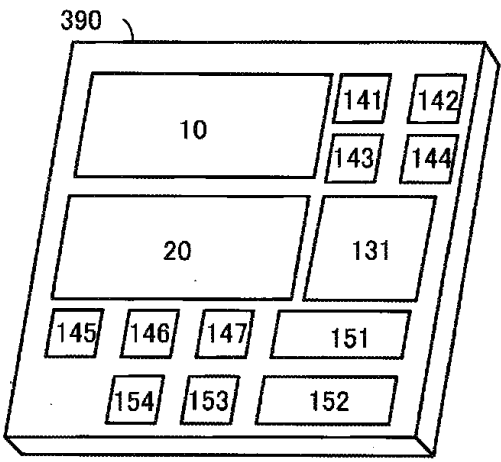
【図 3 3】

図33



【図 3 4 A】

図34A



10

20

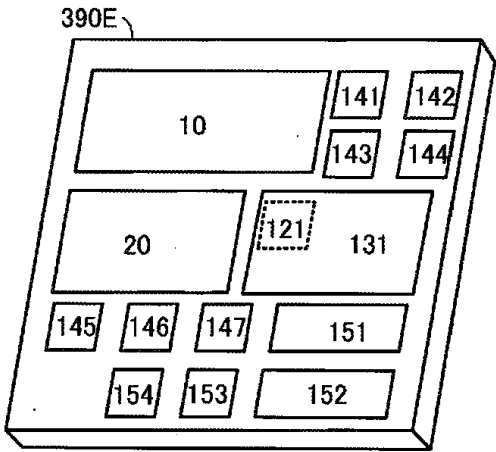
30

40

50

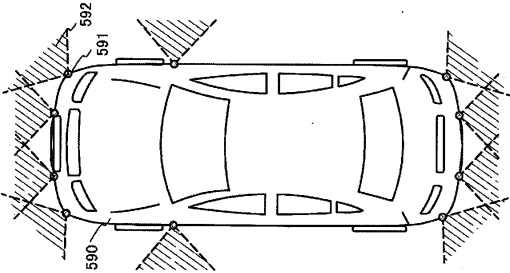
【 3 4 B 】

34B



【 3 5 A 】

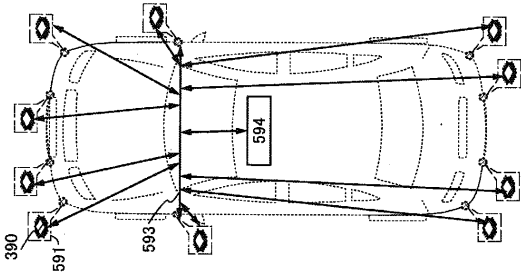
35A



10

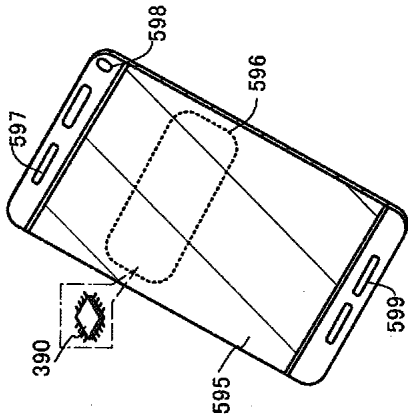
【 3 5 B 】

35B



【 3 6 A 】

36A



20

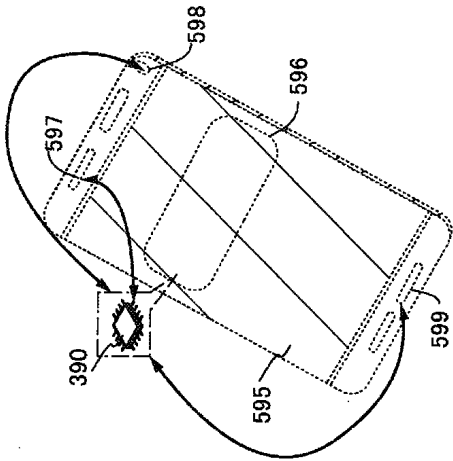
30

40

50

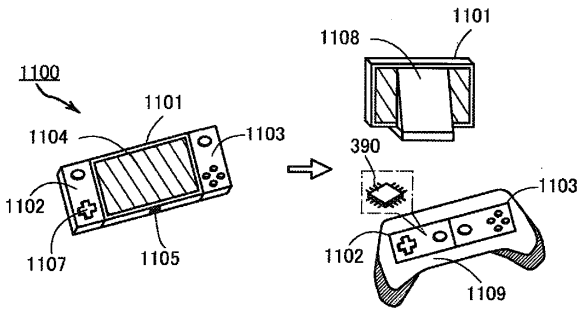
【図 36 B】

図 36B



【図 37 A】

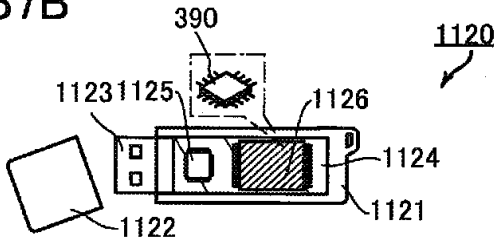
図 37A



10

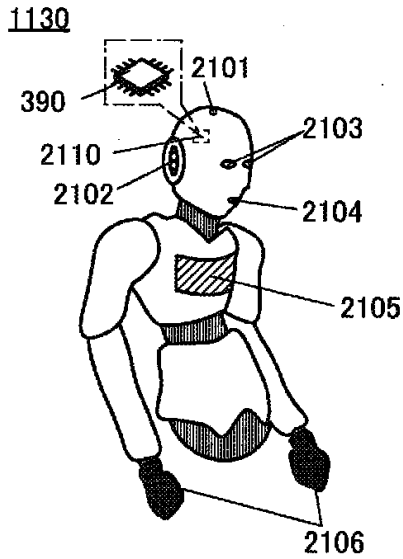
【図 37 B】

図 37B



【図 37 C】

図 37C



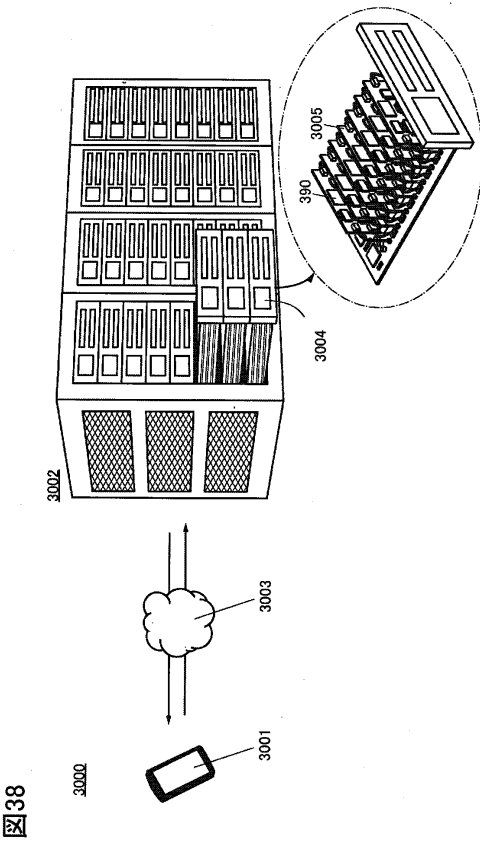
20

30

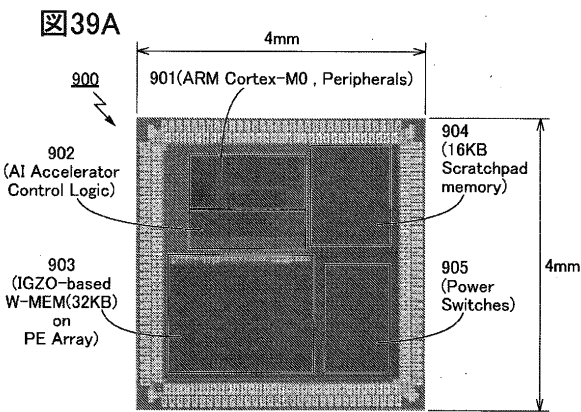
40

50

【図 38】



【図 39 A】

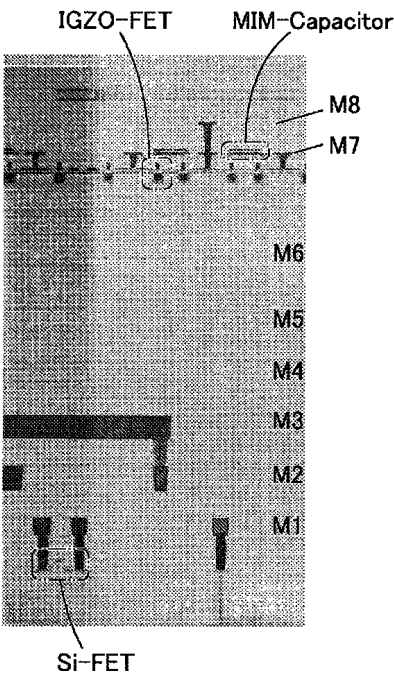


10

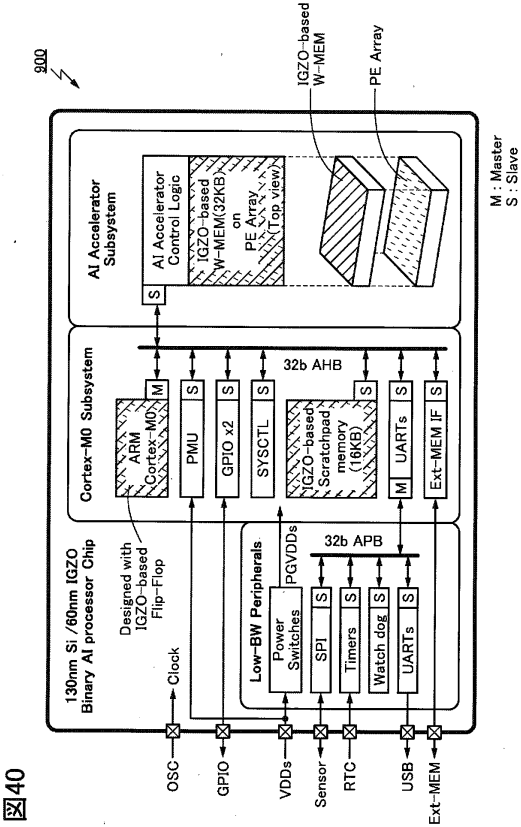
20

【図 39 B】

図 39B



【図 40】



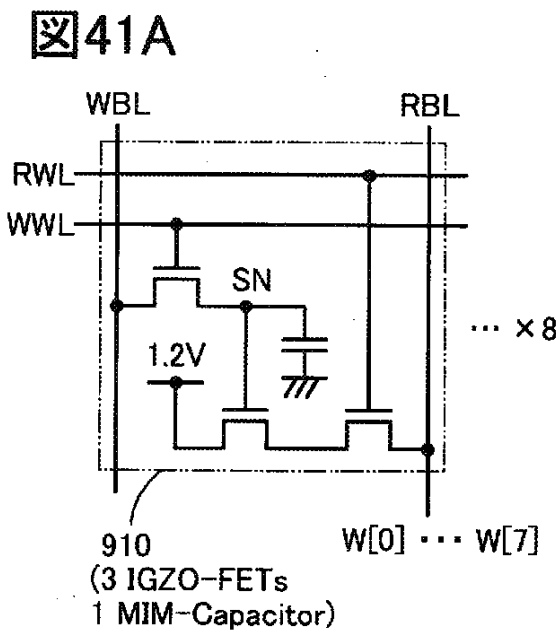
30

40

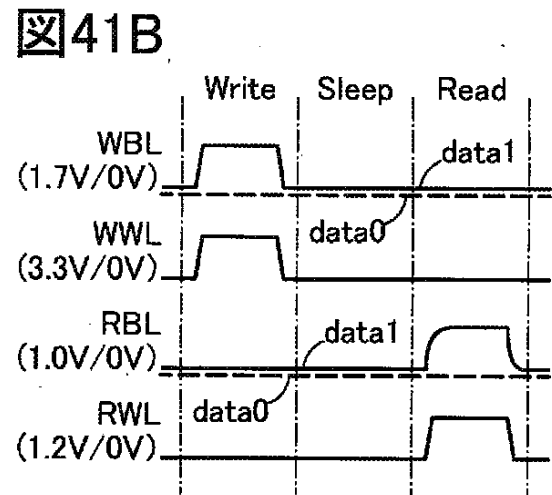
図 40

50

【図 4 1 A】

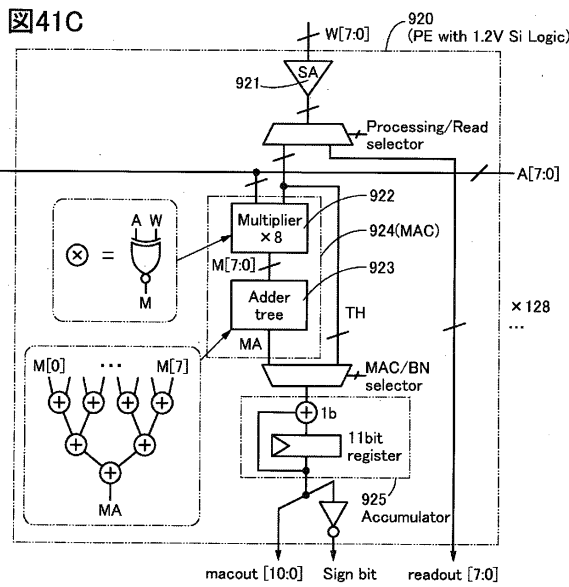


【図 4 1 B】

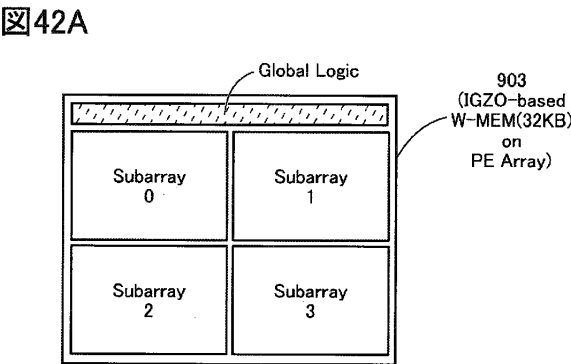


10

【図 4 1 C】



【図 4 2 A】



20

30

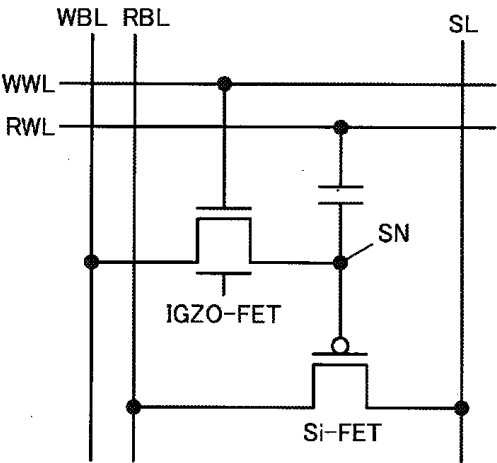
40

50

【 4 4 B 】

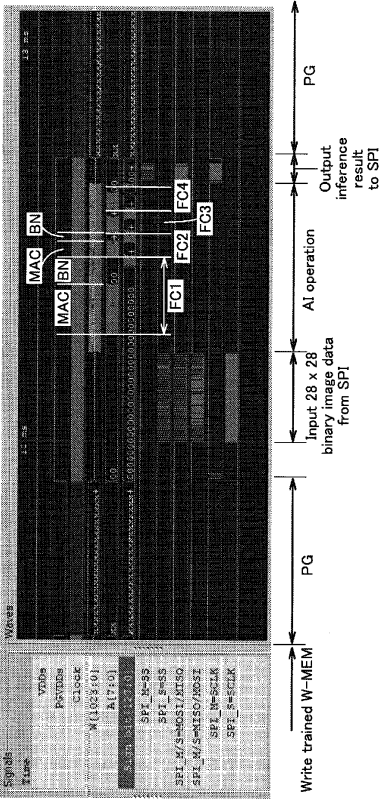
44B

942



【 4 5 A 】

45A

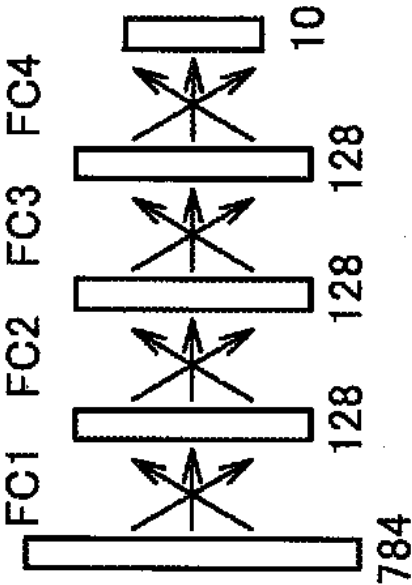


10

20

【 4 5 B 】

45B



30

40

50

フロントページの続き

(51)国際特許分類

F I

G 1 1 C	7/22 (2006.01)	G 1 1 C	7/22	
G 1 1 C	11/404 (2006.01)	G 1 1 C	11/404	
G 1 1 C	11/405 (2006.01)	G 1 1 C	11/405	
G 1 1 C	11/409 (2006.01)	G 1 1 C	11/409	
G 1 1 C	14/00 (2006.01)	G 1 1 C	14/00	
H 0 1 L	21/822 (2006.01)	H 0 1 L	27/04	D
H 0 1 L	27/04 (2006.01)	H 0 1 L	27/04	U
H 0 1 L	21/8234 (2006.01)	H 0 1 L	27/088	E
H 0 1 L	29/786 (2006.01)	H 0 1 L	27/088	H
H 1 0 B	12/00 (2023.01)	H 0 1 L	29/78	6 1 3 B
H 1 0 B	41/70 (2023.01)	H 0 1 L	29/78	6 1 8 B
		H 1 0 B	12/00	8 0 1
		H 1 0 B	41/70	

(32)優先日 令和1年11月29日(2019.11.29)

(33)優先権主張国・地域又は機関

日本国(JP)

(31)優先権主張番号 特願2020-38446(P2020-38446)

(32)優先日 令和2年3月6日(2020.3.6)

(33)優先権主張国・地域又は機関

日本国(JP)

(31)優先権主張番号 特願2020-87645(P2020-87645)

(32)優先日 令和2年5月19日(2020.5.19)

(33)優先権主張国・地域又は機関

日本国(JP)

(72)発明者 古谷 一馬

神奈川県厚木市長谷 3 9 8 番地 株式会社半導体エネルギー研究所内

(72)発明者 佐々木 宏輔

神奈川県厚木市長谷 3 9 8 番地 株式会社半導体エネルギー研究所内

審査官 市川 武宜

(56)参考文献 特開 2 0 1 9 - 3 6 2 8 0 (J P , A)

特開 2 0 1 9 - 4 6 1 9 9 (J P , A)

特開 2 0 1 9 - 4 7 0 0 6 (J P , A)

国際公開第 2 0 1 9 / 0 3 8 6 6 4 (W O , A 1)

(58)調査した分野 (Int.Cl. , D B 名)

H 0 1 L 2 1 / 8 2 2

H 0 1 L 2 1 / 8 2 3 4

H 0 1 L 2 7 / 0 4

H 0 1 L 2 7 / 0 8 8

H 0 1 L 2 9 / 7 8 6

H 1 0 B 1 2 / 0 0

H 1 0 B 4 1 / 7 0

G 0 6 F 9 / 3 8

G 0 6 F 1 2 / 0 0

G 0 6 F 1 5 / 7 8

G 1 1 C 5 / 0 2

G 1 1 C 7 / 2 2

G 1 1 C 1 1 / 4 0 4

G 1 1 C 1 1 / 4 0 5

G 1 1 C 1 1 / 4 0 9

G 1 1 C 1 4 / 0 0