

(19)日本国特許庁(JP)

(12)公開特許公報(A)

(11)公開番号

特開2023-67686
(P2023-67686A)

(43)公開日 令和5年5月16日(2023.5.16)

(51)国際特許分類	F I	テーマコード(参考)
H 1 0 B 99/00 (2023.01)	H 0 1 L 27/10 4 9 5	5 F 0 8 3
H 1 0 B 43/27 (2023.01)	H 0 1 L 27/11582	5 F 1 0 1
H 1 0 B 43/50 (2023.01)	H 0 1 L 27/11575	
H 0 1 L 21/336(2006.01)	H 0 1 L 29/78 3 7 1	

審査請求 有 請求項の数 18 O L 外国語出願 (全70頁)

(21)出願番号 特願2022-18575(P2022-18575)	(71)出願人 599129074 旺宏電子股 ぶん 有限公司 台湾新竹科学工業園區力行路16號
(22)出願日 令和4年2月9日(2022.2.9)	
(31)優先権主張番号 63/273,876	(74)代理人 100107423 弁理士 城村 邦彦
(32)優先日 令和3年10月29日(2021.10.29)	(74)代理人 100120949 弁理士 熊野 剛
(33)優先権主張国・地域又は機関 米国(US)	(74)代理人 100093997 弁理士 田中 秀佳
(31)優先権主張番号 17/569,419	(72)発明者 葉 騰豪 台湾新竹 縣 竹北市莊敬五街161號 3樓
(32)優先日 令和4年1月5日(2022.1.5)	(72)発明者 呂 函庭 台湾新竹市清華大學東院19號
(33)優先権主張国・地域又は機関 米国(US)	Fターム(参考) 5F083 EP18 EP22 EP76 ER21 最終頁に続く

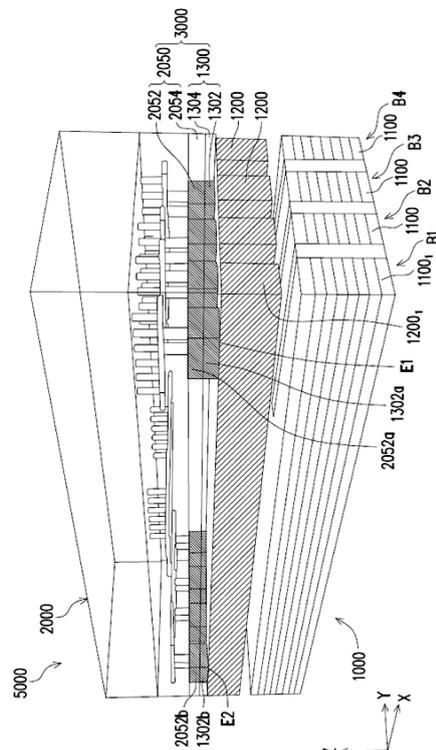
(54)【発明の名称】 3Dフラッシュメモリモジュールチップおよびその製造方法

(57)【要約】 (修正有)

【課題】フラッシュメモリの性能は、複数回の動作後に著しく低下するため、フラッシュメモリに対して修復プロセスを実行する3Dフラッシュメモリモジュール及びその製造方法を提供する。

【解決手段】3Dフラッシュメモリモジュールチップ5000は、メモリチップ1000及び制御チップ2000を含む。メモリチップは、複数のタイル及び複数のヒータ1200を含む。タイルはそれぞれ、複数の3Dフラッシュメモリ構造を含む。ヒータは、タイルのそれぞれの3Dフラッシュメモリ構造の周りに配置されている。制御チップは、ヒータうちの少なくとも1つを駆動するためにメモリチップと接合されている。メモリチップでは、制御チップがヒータを駆動してメモリチップのローカルセクタに対して修復プロセスを行うことができるように制御チップと接合されている。

【選択図】図1A



【特許請求の範囲】**【請求項 1】**

メモリチップであって、
それぞれが複数の 3 Dフラッシュメモリ構造を含む複数のタイル、および
前記タイルのそれぞれの前記 3 Dフラッシュメモリ構造の周りに配置された複数のヒータ、
を含む、メモリチップと、
前記ヒータのうちの少なくとも 1 つを駆動するために前記メモリチップに接合された制御チップと、
を含む、3 Dフラッシュメモリモジュールチップ。

10

【請求項 2】

前記ヒータが前記 3 Dフラッシュメモリ構造の上方に配置され、前記制御チップに隣接している、請求項 1 に記載の 3 Dフラッシュメモリモジュールチップ。

【請求項 3】

前記ヒータが前記 3 Dフラッシュメモリ構造間の複数のスリットトレンチ内に配置されている、請求項 1 に記載の 3 Dフラッシュメモリモジュールチップ。

【請求項 4】

前記メモリチップが、
第 1 の基板上に配置された複数の第 1 のトランジスタと、
前記第 1 のトランジスタの上方に位置する前記 3 Dフラッシュメモリ構造と、
第 1 の相互接続構造であって、前記 3 Dフラッシュメモリ構造が前記第 1 の相互接続構造内に埋め込まれている、第 1 の相互接続構造と、
をさらに含む、請求項 1 に記載の 3 Dフラッシュメモリモジュールチップ。

20

【請求項 5】

前記第 1 の相互接続構造が、
前記 3 Dフラッシュメモリ構造と前記第 1 のトランジスタとの間に位置し、前記 3 Dフラッシュメモリ構造と前記第 1 のトランジスタとを電気的に接続する下部相互接続構造と、
前記 3 Dフラッシュメモリ構造上に位置し、前記 3 Dフラッシュメモリ構造に電気的に接続する上部相互接続構造と、
を含む、請求項 4 に記載の 3 Dフラッシュメモリモジュールチップ。

30

【請求項 6】

前記制御チップが、
それぞれが、
第 2 の基板上に位置する第 2 のトランジスタであって、前記第 2 のトランジスタのソース領域がグローバル電源に電気的に接続されている、第 2 のトランジスタと、
前記第 2 のトランジスタのドレイン領域に電気的に接続され、前記ヒータのうちの 1 つのヒータの第 1 の端部に電気的に接続された第 1 のパッドと、
接地され、前記ヒータのうちの前記 1 つのヒータの第 2 の端部に電気的に接続された第 2 のパッドと、
を含む、複数の駆動行
を含む、請求項 4 に記載の 3 Dフラッシュメモリモジュールチップ。

40

【請求項 7】

前記制御チップが、
前記駆動行の前記第 2 のトランジスタの複数のゲート層に電気的に結合された行デコーダと、
前記第 2 のトランジスタの複数のソース領域および前記グローバル電源に電気的に接続された列デコーダと、
をさらに含む、請求項 6 に記載の 3 Dフラッシュメモリモジュールチップ。

【請求項 8】

50

前記制御チップがアレイ状に配列された複数のタイルを含み、同じ列内の前記タイルの前記第 2 のトランジスタの前記ソース領域が互いに電氣的に接続されている、請求項 6 に記載の 3 D フラッシュメモリモジュールチップ。

【請求項 9】

前記制御チップと前記メモリチップとが接合構造によって接合されている、請求項 6 に記載の 3 D フラッシュメモリモジュールチップ。

【請求項 10】

前記複数の 3 D フラッシュメモリ構造が複数の 3 D AND フラッシュメモリ構造、複数の 3 D NAND フラッシュメモリ構造、または複数の 3 D NOR フラッシュメモリ構造を含む、請求項 6 に記載の 3 D フラッシュメモリモジュールチップ。

10

【請求項 11】

メモリチップを形成するステップであって、

第 1 の基板上に複数のタイルを形成することであり、前記タイルのそれぞれが複数の 3 D フラッシュメモリ構造を含む、第 1 の基板上に複数のタイルを形成すること、および前記タイルのそれぞれの前記 3 D フラッシュメモリ構造の周りに複数のヒータを形成すること、

を含む、メモリチップを形成するステップと、

制御チップを形成するステップと、

前記制御チップと前記メモリチップを接合するステップであって、前記制御チップが前記ヒータを駆動するように構成されている、前記制御チップと前記メモリチップを接合するステップと、

20

を含む、3 D フラッシュメモリモジュールチップの製造方法。

【請求項 12】

前記ヒータが前記 3 D フラッシュメモリ構造の上方に形成されている、請求項 11 に記載の 3 D フラッシュメモリモジュールチップの製造方法。

【請求項 13】

前記ヒータが前記 3 D フラッシュメモリ構造の周りの複数のスリットトレンチ内に形成されている、請求項 11 に記載の 3 D フラッシュメモリモジュールチップの製造方法。

【請求項 14】

前記メモリチップを形成する前記ステップが、

30

前記第 1 の基板上に複数の第 1 のトランジスタを形成するステップと、

前記第 1 のトランジスタの上方に前記 3 D フラッシュメモリ構造を形成するステップと、

を含む、請求項 11 に記載の 3 D フラッシュメモリモジュールチップの製造方法。

【請求項 15】

前記制御チップを形成する前記ステップが、

複数の駆動行を形成するステップであって、前記駆動行のそれぞれの形成が、

第 2 の基板上に第 2 のトランジスタを形成することと、

前記第 2 のトランジスタ上に第 2 の相互接続構造を形成することであり、前記第 2 のトランジスタのソース領域が前記第 2 の相互接続構造を介してグローバル電源に電氣的に結合される、第 2 の相互接続構造を形成することと、

40

前記第 2 の相互接続構造上に第 1 のパッドを形成することであり、前記第 1 のパッドが前記第 2 の相互接続構造を介して前記第 2 のトランジスタのドレイン領域に電氣的に接続される、第 1 のパッドを形成することと、

前記第 2 の相互接続構造上に第 2 のパッドを形成することであり、前記第 2 のパッドが前記第 2 の相互接続構造を介して接地に電氣的に接続される、第 2 のパッドを形成することと、

を含む、複数の駆動行を形成するステップ、

を含む、請求項 14 に記載の 3 D フラッシュメモリモジュールチップの製造方法。

【請求項 16】

50

前記第 1 のパッドを前記ヒータのうちの 1 つのヒータの第 1 の端部に電氣的に接続するステップと、

前記第 2 のパッドを前記ヒータのうちの前記 1 つのヒータの第 2 の端部に電氣的に接続するステップと、

をさらに含む、請求項 15 に記載の 3 D フラッシュメモリモジュールチップの製造方法。

【請求項 17】

前記制御チップと前記メモリチップとが接合構造によってハイブリッド接合されている、請求項 11 に記載の 3 D フラッシュメモリモジュールチップの製造方法。

【請求項 18】

前記複数の 3 D フラッシュメモリ構造が複数の 3 D AND フラッシュメモリ構造、複数の 3 D NAND フラッシュメモリ構造、または複数の 3 D NOR フラッシュメモリ構造を含む、請求項 11 に記載の 3 D フラッシュメモリモジュールチップの製造方法。

【発明の詳細な説明】

【技術分野】

【0001】

本開示の実施形態は、半導体モジュールおよびその製造方法に関し、詳細には 3 D フラッシュメモリモジュールおよびその製造方法に関する。

【背景技術】

【0002】

不揮発性メモリは、電源オフ時に記憶データが消失しないという利点を有するため、パーソナルコンピュータまたは他の電子機器に広く使用されるメモリとなっている。現在、業界で一般的に使用されている 3 次元 (3 D) メモリは、NOR フラッシュメモリおよび NAND フラッシュメモリを含む。加えて、別のタイプの 3 D メモリは、AND フラッシュメモリであり、これは、高集積および高面積利用率を有する多次元メモリアレイに適用することができ、動作速度が速いという利点を有する。したがって、3 D メモリデバイスの開発が徐々に現在のトレンドになってきている。

【発明の概要】

【発明が解決しようとする課題】

【0003】

本開示は、フラッシュメモリに対して局所修復プロセスを行うことができる 3 D フラッシュメモリモジュールチップおよびその製造方法を提供する。

【課題を解決するための手段】

【0004】

本開示の一実施形態において、3 D フラッシュメモリモジュールチップは、メモリチップおよび制御チップを含む。メモリチップは、複数のタイルおよび複数のヒータを含む。タイルはそれぞれ、複数の 3 D フラッシュメモリ構造を含む。ヒータは、タイルのそれぞれの 3 D フラッシュメモリ構造の周りに配置されている。制御チップは、ヒータのうちの少なくとも 1 つを駆動するためにメモリチップと接合されている。

【0005】

本開示の一実施形態において、3 D フラッシュメモリモジュールチップを製造する方法は、以下のステップを含む。メモリチップが形成され、本ステップは、第 1 の基板上に、それぞれが複数の 3 D フラッシュメモリ構造を含む複数のタイルを形成すること、およびタイルのそれぞれの 3 D フラッシュメモリ構造の周りに複数のヒータを形成すること、を含む。制御チップが形成される。制御チップとメモリチップとが接合され、制御チップは、ヒータを駆動するように構成される。

【発明の効果】

【0006】

上記に基づいて、本開示による 3 D フラッシュメモリモジュールチップおよびその製造方法では、追加の制御チップを使用してヒータを駆動し、フラッシュメモリの各セクタに

10

20

30

40

50

対して局所修復プロセスを実行する。制御チップは、ヒータコントローラがメモリチップの面積を占有することを防止するために別個に製造されてもよく、制御チップは、プロセスのコストを削減するためにそれほど高度でないプロセスによって製造されてもよい。

【図面の簡単な説明】

【0007】

【図1A】本開示の一実施形態による3Dフラッシュメモリモジュールチップの概略斜視図である。

【図1B】本開示の一実施形態による3Dフラッシュメモリモジュールチップの概略斜視図である。

【図2A】本開示の一実施形態によるメモリチップの3Dフラッシュメモリ構造の部分上面図である。 10

【図2B】図2AのI-I'線に沿った断面図である。

【図3A】本開示の別の実施形態によるヒータを有するメモリチップの部分上面図である。

【図3B】図3AのI-I'線に沿った断面図である。

【図4A】本開示の別の実施形態によるヒータを有するメモリチップの部分上面図である。

【図4B】本開示の別の実施形態によるメモリチップのヒータおよびパッドの部分上面図である。

【図4C】図4BのII-II'線に沿った断面図である。 20

【図5A】本開示の一実施形態による制御チップの概略斜視図である。

【図5B】本開示の一実施形態による制御チップの概略斜視図である。

【図5C】本開示の一実施形態による制御チップの概略斜視図である。

【図5D】本開示の一実施形態による制御チップの概略斜視図である。

【図5E】本開示の一実施形態による制御チップの概略斜視図である。

【図6A】本開示の一実施形態によるメモリチップおよび制御チップの概略斜視図である。

【図6B】図6Aの概略回路図である。

【図7A】本開示の3Dフラッシュメモリモジュールチップを製造するプロセスの概略断面図である。 30

【図7B】本開示の3Dフラッシュメモリモジュールチップを製造するプロセスの概略断面図である。

【図7C】本開示の3Dフラッシュメモリモジュールチップを製造するプロセスの概略断面図である。

【図8】本開示の別の実施形態による制御チップの概略斜視図である。

【発明を実施するための形態】

【0008】

フラッシュメモリの性能は、複数回の動作後に著しく低下するため、フラッシュメモリに対して修復プロセスを実行する必要がある。修復プロセスでは、フラッシュメモリの電荷蓄積構造（例えば、窒化物層）を修復するために、ヒータを使用してフラッシュメモリを加熱することがある。現在の技術では、ワード線がヒータとして最も一般的に使用されている。しかしながら、ワード線の数が多く、他の構成要素（例えば、ワード線デコーダ）との構成関係が複雑であるため、フラッシュメモリ構造のレイアウト設計がより困難になる場合がある。 40

【0009】

本開示の実施形態は、いくつかの3Dフラッシュメモリモジュールチップを提供しており、メモリチップの3Dフラッシュメモリ構造の上方または側壁の周りにヒータが配置され、メモリチップは、制御チップがヒータを駆動してメモリチップのローカルセクタに対して修復プロセスを行うことができるように制御チップと接合されている。

【0010】

図 1 A および図 1 B は、本開示の一実施形態による 3 D フラッシュメモリモジュールチップの概略斜視図である。図 2 A は、本開示の一実施形態によるメモリチップの 3 D フラッシュメモリ構造の部分上面図である。図 2 B は、図 2 A の I - I ' 線に沿った断面図である。図 3 A は、本開示の別の実施形態によるヒータを有するメモリチップの部分上面図である。図 3 B は、図 3 A の I - I ' 線に沿った断面図である。

【 0 0 1 1 】

図 1 A および図 1 B を参照すると、本開示の一実施形態による 3 D フラッシュメモリモジュールチップ (3 D 集積回路 (3 D I C) と呼ばれる) 5 0 0 0 は、メモリチップ 1 0 0 0 および制御チップ 2 0 0 0 を含む。メモリチップ 1 0 0 0 は、複数の 3 D フラッシュメモリ構造 1 1 0 0 および複数のヒータ 1 2 0 0 を含む。ヒータ 1 2 0 0 は、3 D フラッシュメモリ構造 1 1 0 0 の周りに配置されている。一部の実施形態では、ヒータ 1 2 0 0 は、図 1 A に示すように、3 D フラッシュメモリ構造 1 1 0 0 の上方に配置されている。他の実施形態では、ヒータ 1 2 0 0 は、図 1 B に示すように、3 D フラッシュメモリ構造 1 1 0 0 間のスリットレンチ 1 1 1 0 内に配置されている。制御チップ 2 0 0 0 は、メモリチップ 1 0 0 0 の上方に配置され、メモリチップ 1 0 0 0 内のヒータ 1 2 0 0 を駆動する。制御チップ 2 0 0 0 とメモリチップ 1 0 0 0 は、接合構造 3 0 0 0 によって互いに接合されてもよい。

10

【 0 0 1 2 】

図 1 A および図 1 B を参照すると、メモリチップ 1 0 0 0 の 3 D フラッシュメモリ構造 1 1 0 0 は、(図 2 A および図 2 B に示すような) 3 D A N D フラッシュメモリ構造、3 D N A N D フラッシュメモリ構造 (図示せず) 、または 3 D N O R フラッシュメモリ構造 (図示せず) であってもよい。3 D A N D フラッシュメモリ構造を例として、本開示の 3 D フラッシュメモリ構造 1 1 0 0 を説明するが、本開示の実施形態はこれに限定されない。

20

【 0 0 1 3 】

図 2 A および図 2 B を参照すると、メモリチップ 1 0 0 0 は、複数のタイル T を含むことができる。タイル T は、複数の列および複数の行を含むアレイ状に配置されてもよい。本実施形態では、例示のために 4 つのタイル T (例えば、T 1 ~ T 4) が示されている。4 つのタイル T のうち、タイル T 1 とタイル T 2 とが 1 つの行に配置され、タイル T 3 とタイル T 4 とが別の行に配置されている。タイル T 1 とタイル T 3 とが 1 つの列に配置され、タイル T 2 とタイル T 4 とが別の列に配置されている。各タイル T は、複数のセクタ B (例えば、B 1 ~ B 4) を含むことができる。各セクタ B は、3 D フラッシュメモリ構造 1 1 0 0 を含む。3 D フラッシュメモリ構造 1 1 0 0 は、X 方向に延在し、Y 方向に配列されている。2 つの隣接する 3 D フラッシュメモリ構造 1 1 0 0 は、スリットレンチ 1 1 1 0 によって互いに分離されている。

30

【 0 0 1 4 】

図 2 B を参照すると、3 D フラッシュメモリ構造 1 1 0 0 のそれぞれは、複数のメモリセルによって形成されたメモリアレイを少なくとも含むことができる。具体的には、3 D フラッシュメモリ構造 1 1 0 0 は、第 1 の基板 (例えば、半導体基板) 1 0 1 0 上の 1 つまたは複数の能動デバイス (例えば、第 1 のトランジスタ 1 0 2 0) の上方に配置されてもよい。第 1 のトランジスタ 1 0 2 0 は、例えば、相補型金属酸化膜半導体 (C M O S) 電界効果トランジスタである。したがって、本アーキテクチャは、アレイ下相補型金属酸化膜半導体電界効果トランジスタ (C M O S u n d e r A r r a y (C U A)) アーキテクチャと呼ばれることもある。

40

【 0 0 1 5 】

図 2 B を参照すると、3 D フラッシュメモリ構造 1 1 0 0 は、半導体ダイのバックエンドオブライン (B E O L) において配置されてもよい。例えば、3 D フラッシュメモリ構造 1 1 0 0 は、第 1 の相互接続構造 1 0 3 0 に埋め込まれてもよい。第 1 の相互接続構造 1 0 3 0 は、例えば、下部相互接続構造 1 0 3 2 および上部相互接続構造 1 0 3 4 を含む。下部相互接続構造 1 0 3 2 は、第 1 の基板 (例えば、半導体基板) 1 0 1 0 上の 1 つま

50

たは複数の能動デバイス（例えば、第1のトランジスタ1020）の上方で、3Dフラッシュメモリ構造1100のメモリアレイの下方に配置されている。上部相互接続構造1034は、3Dフラッシュメモリ構造1100のメモリアレイの上方に配置されている。下部相互接続構造1032は、例えば、下部第1金属層BM1、下部第2金属層BM2、および下部第3金属層BM3、ならびにそれらの間のビアBV1およびBV2を含む。上部相互接続構造1034は、例えば、上部第1金属層TM1および上部第2金属層TM2、ならびにそれらの間のビアTV1を含む。下部相互接続構造1032および上部相互接続構造1034の金属層とビアの数は、上記に限定されない。

【0016】

図2Bを参照すると、3Dフラッシュメモリ構造1100は、複数のゲートスタック構造52を含む。ゲートスタック構造52のそれぞれは、下部相互接続構造1032上に形成されている。ゲートスタック構造52のそれぞれは、第1の基板1010のアレイ領域ARから階段領域SRまでX方向に延在する。ゲートスタック構造52は、第1の基板1010の表面上に垂直に積み重ねられた複数のゲート層（ワード線とも呼ばれる）38および複数の絶縁層54を含む。Z方向において、ゲート層38は、間に配置された絶縁層54によって互いに電氣的に絶縁されている。ゲート層38は、タングステンなどの金属層を含む。一部の実施形態では、ゲート層38は、チタン（Ti）、窒化チタン（TiN）、タンタル（Ta）、窒化タンタル（Ta₂N₅）、またはそれらの組合せなどのバリア層37をさらに含む。絶縁層54は、例えば酸化ケイ素である。

【0017】

ゲート層38は、（図2Bに示す）第1の基板1010の表面に平行な方向に延在する。階段領域SR内のゲート層38は、（図2Bに示す）階段構造SCを有することができ、それにより、下側ゲート層38は、上側ゲート層38よりも長く、下側ゲート層38の端部は、上側ゲート層38の端部を越えて横方向に延在する。ゲート層38を接続するためのコンタクトC1は、階段領域SR内のゲート層38の端部上に着地して、ゲート層38のそれぞれを、コンタクトC1および上部相互接続構造1034を介して下部相互接続構造1032の導電線（例えば、下部第3金属層BM3の導電線）に接続することができる。

【0018】

図2Bを参照すると、3Dフラッシュメモリ構造1100は、複数のチャンネルピラー16をさらに含む。チャンネルピラー16は、アレイ領域ARのゲートスタック構造52を貫いて連続的に延びている。一部の実施形態では、チャンネルピラー16は、上面図においてリング形状のプロファイルを有してもよい。チャンネルピラー16の材料は、ドーブされていないポリシリコンなどの半導体であってもよい。

【0019】

図2Bを参照すると、3Dフラッシュメモリ構造1100は、絶縁性充填層24、絶縁性ピラー28、複数の導電性ピラー（例えば、ソースピラーとして機能する）32a、および複数の導電性ピラー（例えば、ドレインピラーとして機能する）32bをさらに含む。導電性ピラー32a、32bおよび絶縁性ピラー28は、チャンネルピラー16内に配置されており、それぞれ、ゲート層38に対して垂直な方向（Z方向）に延在している。導電性ピラー32aおよび32bは、絶縁性充填層24および絶縁性ピラー28によって互いに分離され、チャンネルピラー16に電氣的に接続されている。導電性ピラー32a、32bは、例えばドーブされたポリシリコンである。絶縁性充填層24は、例えば酸化ケイ素であり、絶縁性ピラー28は、例えば窒化ケイ素である。

【0020】

図2Bを参照すると、電荷蓄積構造40は、チャンネルピラー16とゲート層38との間に配置されている。電荷蓄積構造40は、トンネル層（またはバンドギャップ工学設計されたトンネル酸化物層と呼ばれる）14、電荷蓄積層12、およびブロッキング層36を含むことができる。電荷蓄積層12は、トンネル層14とブロッキング層36との間に配置されている。一部の実施形態では、トンネル層14、電荷蓄積層12、およびブロッキ

10

20

30

40

50

ング層 36 は、例えば、酸化ケイ素、窒化ケイ素、および酸化ケイ素である。一部の実施形態では、図 2 B に示すように、電荷蓄積構造 40 の一部（例えば、トンネル層 14）は、ゲート層 38 に垂直な方向（すなわち、Z 方向）に連続的に延在し、電荷蓄積構造 40 の他の部分（例えば、電荷蓄積層 12 およびブロッキング層 36）は、ゲート層 38 を取り囲む。他の実施形態では、電荷蓄積構造 40（例えば、トンネル層 14、電荷蓄積層 12、およびブロッキング層 36）は、ゲート層 38 を取り囲む（図示せず）。ゲート層 38 のそれぞれと、ゲート層 38 によって取り囲まれた電荷蓄積構造 40、チャンネルピラー 16、ソースピラー 32 a、およびドレインピラー 32 b とがメモリセル 20 を画定する。したがって、3D フラッシュメモリ構造 1100 のそれぞれは、複数のメモリセル 20 から構成されるメモリアレイを少なくとも含む。

10

【0021】

3D フラッシュメモリ構造 1100 は、ローカルビット線 LBL_n 、ローカルソース線 LSL_n 、グローバルビット線 GBL_n 、およびグローバルソース線 GSL_n をさらに含む。ローカルビット線 LBL_n およびローカルソース線 LSL_n は、上部相互接続構造 1034 の上部第 1 金属層 TM_1 内に位置し、コンタクト C_2 を介してソースピラー 32 a およびドレインピラー 32 b にそれぞれ電氣的に接続されている。グローバルビット線 GBL_n およびグローバルソース線 GSL_n は、上部相互接続構造 1034 内の上部ビア（図示せず）を介して、ローカルビット線 LBL_n およびローカルソース線 LSL_n にそれぞれ電氣的に接続されている。

20

【0022】

異なる動作方法によると、メモリセル 20 で 1 ビット動作または 2 ビット動作を行うことができる。例えば、ソースピラー 32 a およびドレインピラー 32 b に電圧が印加されると、ソースピラー 32 a およびドレインピラー 32 b がチャンネルピラー 16 に接続されているため、電子は、チャンネルピラー 16 に沿って移動し、電荷蓄積構造 40 全体に蓄積され得る。したがって、メモリセル 20 で 1 ビット動作を行うことができる。加えて、ファウラー・ノルトハイムトンネリングを伴う動作の場合、電子または正孔は、ソースピラー 32 a とドレインピラー 32 b との間の電荷蓄積構造 40 内にトラップされ得る。ソース側注入、チャンネルホットエレクトロン注入、またはバンド間トンネリングホットキャリア注入を伴う動作の場合、電子または正孔は、ソースピラー 32 a およびドレインピラー 32 b の一方に隣接する電荷蓄積構造 40 内に局所的にトラップされ得る。したがって、

30

【0023】

動作中、選択されたワード線（ゲート層）38 に電圧が印加され、例えば、対応するメモリセル 20 の対応する閾値電圧（ V_{th} ）よりも高い電圧が印加されると、選択されたワード線 38 と交差するチャンネルピラー 16 のチャンネル領域がオンになり、電流がローカルビット線 LBL_n からドレインピラー 32 b に入り、オンしたチャンネル領域を介してソースピラー 32 a に流れ、最終的にローカルソース線 LSL_n に流れるようになる。

【0024】

図 3 A および図 3 B を参照すると、メモリチップ 1000 は、複数のヒータ 1200 をさらに含む。ヒータ 1200 は、3D フラッシュメモリ構造 1100 の上方の誘電体層 1040 内に配置されてもよい。誘電体層 1040 の材料は、例えば、酸化ケイ素である。ヒータ 1200 は、銅またはタンゲステンなどの金属層 1202 を含む。一部の実施形態では、ヒータ 1200 は、チタン、タンタル、窒化チタン、窒化タンタル、またはそれらの組合せなどのバリア層 1204 をさらに含む。

40

【0025】

図 3 A を参照すると、一部の実施形態では、1 つのヒータ 1200 が各セクタ B 上に配置され、任意の 2 つの隣接するセクタ B の 2 つのヒータ 1200 は、互いに分離されている。ヒータ 1200 は、X 方向に延在することができる。一実施形態において、ヒータ 1200 は、（図 3 A および図 3 B に示すように）アレイ領域 AR 内に配置され、階段領域

50

S Rまで延在する。一実施形態において、ヒータ1200は、アレイ領域AR内に配置されてもよいが、階段領域SRには配置されない(図示せず)。言い換えれば、ヒータ1200の長さは、3Dフラッシュメモリ構造1100のX方向の長さよりも大きくても、等しくても、または小さくてもよい。

【0026】

加えて、各セクタBに複数のヒータ1200を配置してもよく、例えば、アレイ領域ARと階段領域SRにそれぞれ1つのヒータ1200を配置してもよく、別々に加熱を行ってもよい(図示せず)。しかしながら、本開示の実施形態は、これに限定されない。別の実施形態では、隣接する2つ、3つ、またはそれ以上のセクタBの複数のヒータ1200を1つのヒータ(図示せず)にまとめて、複数のセクタBの3Dフラッシュメモリ構造1100を同時に加熱することもできる。

10

【0027】

図3Aを参照すると、上面図におけるヒータ1200の形状は、例えば、矩形または別の形状である。複数のセクタB上のヒータ1200は、同じ幅または異なる幅を有することができる。アレイ領域ARにおけるヒータ1200の幅W1は、階段領域SRにおけるヒータ1200の幅W2と同じである。しかしながら、本開示は、これに限定されない。ヒータ1200の形状は、実際の要件または設計に従って変更されてもよい。アレイ領域ARにおけるヒータ1200の幅W1は、階段領域SRにおけるヒータ1200の幅W2よりも大きくても、等しくても、または小さくてもよい。

【0028】

図1A、図1B、および図3Bを参照すると、メモリチップ1000は、接合層1300をさらに含む。接合層1300は、パッド1302および絶縁層1304を有する。絶縁層1304は、ヒータ1200上に配置されている。絶縁層1304の材料は、例えば、酸化ケイ素である。パッド1302は、各ヒータ1200の表面の絶縁層1304に配置されている。パッド1302の材料は、例えば、銅である。パッド1302は、パッド1302a、1302bを有する。パッド1302a、1302bは、ヒータ1200の第1の端部E1および第2の端部E2にそれぞれ接続されている。

20

【0029】

上記の実施形態では、3Dフラッシュメモリ構造1100は、3D ANDフラッシュメモリ構造であり、ヒータ1200は、(図3A、図3B、および図6Aに示すように)3D ANDフラッシュメモリ構造の上方に配置されている。他の実施形態では、3Dフラッシュメモリ構造1100は、3D ANDフラッシュメモリ構造であり、ヒータ1200は、(図4A~図4Cに示すように)3D ANDフラッシュメモリ構造間のスリットトレンチ1110内に配置されている。

30

【0030】

図4Aは、本開示の別の実施形態によるヒータを有するメモリチップの部分上面図である。図4Bは、本開示の別の実施形態によるメモリチップのヒータおよびパッドの部分上面図である。図4Cは、図4BのII-II'線に沿った断面図である。

【0031】

図4Aおよび図4Cを参照すると、複数のヒータ1200が、3Dフラッシュメモリ構造1100間のスリットトレンチ1110内に配置されている。ヒータ1200は、ゲートスタック構造52の複数のゲート層38および複数の絶縁層54の周りに配置されている。ヒータ1200は、(図4Cに示すように)絶縁ライナ層1112によってゲート層38および絶縁層54から分離されている。絶縁ライナ層1112は、酸化ケイ素または窒化ケイ素などの絶縁材料を含む。ヒータ1200は、銅またはタングステンなどの金属層1202を(図4Cに示すように)含む。一部の実施形態では、ヒータ1200は、(図4Cに示すように)バリア層1204をさらに含む。バリア層1204は、絶縁ライナ層1112と金属層1202との間に配置されている。バリア層1204は、例えば、チタン、タンタル、窒化チタン、窒化タンタル、またはそれらの組合せである。

40

【0032】

50

一部の実施形態では、1つのヒータ1200がスリットトレンチ1110のそれぞれに配置されている。例えば、ヒータ1200は、X方向に延在することができる。一実施形態において、ヒータ1200は、(図4Aおよび図4Bに示すように)アレイ領域AR内に配置され、階段領域SRまで延在する。一実施形態において、ヒータ1200は、アレイ領域AR内に配置されてもよいが、階段領域SRには配置されない(図示せず)。言い換えれば、ヒータ1200の長さは、3Dフラッシュメモリ構造1100のX方向の長さよりも大きくても、等しくても、または小さくてもよい。

【0033】

あるいは、複数のヒータ1200がスリットトレンチ1110のそれぞれに配置されてもよい。例えば、アレイ領域ARおよび階段領域SRにそれぞれ1つのヒータ1200を設けてもよく、別々に加熱が行われてもよい(図示せず)。しかしながら、本開示の実施形態は、これに限定されない。

10

【0034】

加えて、図4Aを参照すると、上面図におけるヒータ1200の形状は、例えば、矩形または別の形状である。複数のスリットトレンチ1110内のヒータ1200は、同じ幅または異なる幅を有することができる。しかしながら、本開示は、これに限定されない。ヒータ1200の形状は、実際の要件または設計に従って変更されてもよい。

【0035】

図4Bおよび図4Cを参照すると、コンタクトC3が各ヒータ1200の2つの端部(すなわち、E1およびE2)の表面上にそれぞれ配置されている。コンタクトC3は、上部相互接続構造1034を介して上方のパッド1302aおよび1302bに接続されてもよく、それにより、メモリチップ1000のヒータ1200は、上部相互接続構造1034ならびにパッド1302aおよび1302bを介して制御チップ2000に電氣的に接続され得る。パッド1302a、1302bの材料は、例えば、銅である。

20

【0036】

図5A~図5Eは、本開示の一実施形態による制御チップの概略斜視図である。図6Aは、本開示の一実施形態によるメモリチップおよび制御チップの概略斜視図である。図6Bは、図6Aの概略回路図である。

【0037】

図5Aを参照すると、制御チップ2000は、複数のタイルT'を含むことができる。タイルT'は、アレイ状に配置されてもよい。本実施形態では、4つのタイルT'(例えば、T1'~T4')が説明のための例として取り上げられる。4つのタイルT'のうち、タイルT1'とタイルT2'は、1つの行に配置され、タイルT3'とタイルT4'は、別の行に配置されている。タイルT1'とタイルT3'は、1つの列に配置され、タイルT2'とタイルT4'は、別の列に配置されている。

30

【0038】

図5Aおよび図5Eを参照すると、各タイルT'は、複数の駆動行2000Rおよび駆動列2000Cを含む。図5Eに示すように、駆動行2000Rのそれぞれは、第2のトランジスタ2020と、第2の相互接続構造2030と、パッド2052と、を含む。第2のトランジスタ2020は、第2の基板2010の活性領域2012上に配置されている。第2の基板2010は、シリコン基板などの半導体基板であってもよい。第2のトランジスタ2020は、相補型金属酸化膜半導体(CMOS)トランジスタであってもよい。第2のトランジスタ2020は、(図5A~図5Eに示すような)プレーナトランジスタ、または(図8に示すような)フィントランジスタであってもよい。

40

【0039】

図5Eおよび図8を参照すると、第2のトランジスタ2020は、ゲート誘電体層2024、ゲート層2028、ソース領域2022a、およびドレイン領域2022bを含む。ゲート誘電体層2024は、例えば、酸化ケイ素または高誘電率材料である。ゲート層2028は、例えば、ドーパされたポリシリコンまたはタングステンである。ゲート層2028は、ゲート誘電体層2024上に位置する。ゲート層2028は、ストリップ形状

50

を有し、その延在方向は、例えば、図 6 A に示すように、（例えば、X 方向に延在する）ヒータ 1 2 0 0 の延在方向と同じである。一部の実施形態では、図 5 A に示すように、2 つの隣接する行（例えば、タイル T 1 ' と T 2 '、またはタイル T 3 ' と T 4 '）の第 2 のトランジスタ 2 0 2 0 のゲート層 2 0 2 8 は、電氣的に接続されていてもよい。

【0040】

図 5 C および図 5 E を参照すると、第 2 のトランジスタ 2 0 2 0 のソース領域 2 0 2 2 a およびドレイン領域 2 0 2 2 b は、ゲート層 2 0 2 8 の両側の活性領域 2 0 1 2 内に配置されている。ソース領域 2 0 2 2 a およびドレイン領域 2 0 2 2 b は、N 型または P 型ドーパントなどのドーパントを含む。一部の実施形態では、2 つの隣接する第 2 のトランジスタ 2 0 2 0 は、ソース領域 2 0 2 2 a を共有する。

10

【0041】

図 5 B および図 5 C を参照すると、第 2 の相互接続構造 2 0 3 0 は、第 2 のトランジスタ 2 0 2 0 上に配置されている。第 2 の相互接続構造 2 0 3 0 は、（図 5 C に示すような）誘電体層 2 0 3 1 と、誘電体層 2 0 3 1 内に位置する複数のコンタクト 2 0 3 2 および 2 0 3 4 と、複数の導電線 2 0 3 6 および 2 0 4 0 と、複数のビア 2 0 3 8 および 2 0 4 2 と、を含む。コンタクト 2 0 3 2 は、ソース領域 2 0 2 2 a およびドレイン領域 2 0 2 2 b 上にそれぞれ着地し、ソース領域 2 0 2 2 a およびドレイン領域 2 0 2 2 b と電氣的に接続されている。コンタクト 2 0 3 4 は、ゲート層 2 0 2 8 上に着地し、ゲート層 2 0 2 8 に電氣的に接続されている。コンタクト 2 0 3 2 は、図 5 B および図 5 D に示すように、X 方向に沿って延在するストリップ形状を有し、ゲート層 2 0 2 8 と実質的に平行である。コンタクト 2 0 3 4 の形状は、コンタクト 2 0 3 2 の形状とは異なり、例えば、図 5 B に示すように、柱状形状であってもよい。（図 5 C に示すような）導電線 2 0 3 6 および 2 0 4 0 は、それぞれコンタクト 2 0 3 2 および 2 0 3 4 上に配置されている。導電線 2 0 3 6 および導電線 2 0 4 0 は、ビア 2 0 3 8 によって互いに電氣的に絶縁されている。ビア 2 0 4 2 は、導電線 2 0 4 0 上に配置され、導電線 2 0 4 0 を上方の接合層 2 0 5 0 に電氣的に接続する。誘電体層 2 0 3 1 は、例えば、酸化ケイ素である。コンタクト 2 0 3 2 および 2 0 3 4、導電線 2 0 3 6 および 2 0 4 0、ならびにビア 2 0 3 8 および 2 0 4 2 は、タンゲステンまたは銅などの金属層を含む。コンタクト 2 0 3 2 および 2 0 3 4、導電線 2 0 3 6 および 2 0 4 0、ならびにビア 2 0 3 8 および 2 0 4 2 は、チタン、タンタル、窒化チタン、窒化タンタル、またはそれらの組合せなどのバリア層（図示せず）をさらに含むことができる。

20

30

【0042】

図 5 C を参照すると、駆動行 2 0 0 0 R のそれぞれのパッド 2 0 5 2 は、制御チップ 2 0 0 0 の接合層 2 0 5 0 の一部である。接合層 2 0 5 0 は、パッド 2 0 5 2 および絶縁層 2 0 5 4 を有する。絶縁層 2 0 5 4 は、第 2 の相互接続構造 2 0 3 0 上に位置する。パッド 2 0 5 2 は、絶縁層 2 0 5 4 内に位置し、第 2 の相互接続構造 2 0 3 0 のビア 2 0 4 2 に電氣的に接続されている。パッド 2 0 5 2 の材料は、例えば、銅である。絶縁層 2 0 5 4 の材料は、例えば、酸化ケイ素である。

【0043】

図 5 A および図 5 E を参照すると、パッド 2 0 5 2 は、パッド 2 0 5 2 a およびパッド 2 0 5 2 b を含む。具体的には、駆動行 2 0 0 0 R のそれぞれは、X 方向に沿って配置された一対のパッド 2 0 5 2 a、2 0 5 2 b を含む。図 1 A、図 1 B、および図 6 A に示すように、パッド 2 0 5 2 a は、ヒータ 1 2 0 0 の第 1 の端部 E 1 に電氣的に接続され、パッド 2 0 5 2 b は、ヒータ 1 2 0 0 の第 2 の端部 E 2 に電氣的に接続され、接地されている。図 5 C および図 5 D を参照すると、各パッド 2 0 5 2 a は、ビア 2 0 4 2 a を介して下方の導電線 2 0 4 0 a に電氣的に接続されている。同じタイル T ' 内の導電線 2 0 4 0 a は、図 5 A および図 5 C に示すように、第 2 のトランジスタ 2 0 2 0 のドレイン領域 2 0 2 2 b にそれぞれ電氣的に接続されるように、互いに分離され、電氣的に絶縁されている。各パッド 2 0 5 2 b は、図 5 D に示すように、ビア 2 0 4 2 b を介して下方の導電線 2 0 4 0 b に電氣的に接続されている。同じ列内のタイル T '（例えば、タイル T 1 ' と T

40

50

3'、またはタイルT2'とT4')のパッド2052bは、図5Aおよび図5Dに示すように、同じ導電線2040bを介して接地に電氣的に接続されている。

【0044】

図5C、図1A、および図1Bを参照すると、制御チップ2000の接合層2050とメモリチップ1000の接合層1300とが互いに接合されて、接合構造3000が形成されている。具体的には、制御チップ2000の絶縁層2054の位置と、メモリチップ1000の絶縁層1304の位置とが互いに対応し、互いに接合されている。制御チップ2000のパッド2052a、2052bの位置と、メモリチップ1000のパッド1302a、1302bの位置とが互いに対応し、互いに接合されている。

【0045】

10

図5A、図5C、および図5Dを参照すると、制御チップ2000の列2000Cは、同じ列内のタイルT'(例えば、タイルT1'とT3'、またはタイルT2'とT4')の複数の第2のトランジスタ2020の複数の共有ソース領域2022aを、導電線2040cを介してグローバル電源2100に電氣的に結合している。

【0046】

図5A、図5C、および図5Dを参照すると、制御チップ2000の第2のトランジスタ2020のドレイン領域2022bは、図5Cに示すように、第2の相互接続構造2030および接合層2050のパッド2052aに接続されている。パッド2052aは、図6Aに示すように、メモリチップ1000のヒータ1200の第1の端部E1に接続されたパッド1302aと電氣的に接続されている。一実施形態において、制御チップ2000の駆動行2000Rのそれぞれは、図6Aおよび図6Bに示すように、メモリチップ1000の1つの対応するセクタBの1つのヒータ1200を制御することができる。

20

【0047】

図5Eを参照すると、一部の実施形態において、制御チップ2000は、列デコーダ2300および行デコーダ2200をさらに含む。列デコーダ2300は、グローバル電源2100に電氣的に接続されている。列アドレス信号A3およびA4を受信した後、列デコーダ2300は、1つの列(例えば、図5Aの列2000C₁)の複数のタイル(本例では2つのタイル、例えば、図5AのタイルT1'およびT3')を選択する。したがって、グローバル電源2100は、第2の相互接続構造2030の導電線2040c(図5Aに示す)を介して、選択された列(例えば、図5Aの列2000C₁)のタイル(例えば、図5AのタイルT1'およびT3')のそれぞれの第2のトランジスタ2020の共有ソース領域2022aに供給される。行デコーダ2200は、駆動行2000Rの第2のトランジスタ2020のゲート層2028に電氣的に接続されている。行デコーダ2200は、行アドレス信号A0~A2(または制御信号と呼ばれる)を受信した後、入力された行アドレス信号をデコードして、第2のトランジスタ2020のうちの1つ(例えば、図5Aの第2のトランジスタ2020₁)または複数を選択してオンにする。

30

【0048】

一般に、メモリチップ1000は、メモリアレイを制御するための制御ロジックユニットを含み、制御ロジックユニット内のレジスタは、各セクタBのメモリアレイの消去回数のステータス信号を記憶する。消去回数が所定の回数に達すると、ステータス信号が制御チップ2000に送られる。

40

【0049】

図6Aおよび図6Bを参照すると、修復プロセス中に、制御チップ2000は、受信したステータス信号に基づいて、修復を必要とするタイルTおよびセクタB(例えば、図6AのタイルT1のセクタB1)に対応する行アドレス信号および列アドレス信号を生成し、行アドレス信号および列アドレス信号をそれぞれ行デコーダ2200および列デコーダ2300に送信することができる。列デコーダ2300は、受信した列アドレス信号に従って1つの列(例えば、図6Aの列2000C₁)を選択して、この列(例えば、図6Aの列2000C₁)に位置するタイル(例えば、図5AのタイルT1'およびT3')の導電線2040cにグローバル電源2100を供給する。行デコーダ2200は、受信した

50

行アドレス信号に従って、1つの駆動行2000R₁の第2のトランジスタ2020₁を選択してオンにする。したがって、電流は、グローバル電源2100から導電線2040cを介して第2のトランジスタ2020₁のソース領域2022aに流れ、第2のトランジスタ2020₁のチャネルおよびドレイン領域2022bを通過し、次いで、第2の相互接続構造2030およびパッド2052aを介してメモリチップ1000のパッド1302aに流れ込み、次いで、ヒータ1200（例えば、1200₁）の第1の端部E₁に入ることができる。その後、電流は、ヒータ1200₁内を流れ、メモリチップ1000のパッド1302bを介してヒータ1200₁の第2の端部E₂から流れ出て、次いで、制御チップ2000のパッド2052bに入り、次いで、導電線2040bを介して接地に電氣的に接続される。本開示の実施形態では、制御チップ2000の第2のトランジスタ（ドライバ）2020（例えば、2020₁）を用いて、特定のヒータ1200（例えば、1200₁）に高い駆動電流を供給することが可能であり、その結果、ヒータ1200（例えば、1200₁）として機能する導体が加熱されて、特定のタイルT（例えば、T₁）内の特定のセクタB（例えば、B₁）の3Dフラッシュメモリ構造1100内の電荷蓄積層を修復する。

【0050】

図1Aおよび図1Bを参照すると、一部の実施形態では、修復プロセスでは、制御チップ2000は、1つのヒータ1200（例えば、1200₁）を駆動して、1つのセクタB（例えば、B₁）の3Dフラッシュメモリ構造1100（例えば、1100₁）内の電荷蓄積層を修復することができる。図1Bを参照すると、他の実施形態では、修復を行う際に、制御チップ2000は、2つのヒータ1200（例えば、1200₂および1200₃）を同時に駆動して、1つのセクタB（例えば、B₂）の3Dフラッシュメモリ構造1100（例えば、1100₂）内の電荷蓄積層を修復することもできる。

【0051】

図7A～図7Cは、本開示の3Dフラッシュメモリモジュールチップを製造するプロセスの概略断面図を示す。

【0052】

図7Aを参照すると、ウエハ1010Wが用意され、複数のメモリチップ1000がウエハ1010W上に形成される。メモリチップ1000間には、スクライプラインSLが設けられている。メモリチップ1000の形成方法は、以下の通りである。図3Bを参照すると、1つまたは複数の能動デバイス（例えば、第1のトランジスタ）1020が、最初にウエハ1010W上に形成される。次に、下部相互接続構造1032が能動デバイス1020上に形成される。下部相互接続構造1032は、ダマシン、デュアルダマシンなどの任意の既知の方法によって形成することができる。その後、1つの絶縁層（例えば、酸化ケイ素）54と別の絶縁層（図示せず、例えば、窒化ケイ素）を交互に積み重ねることによって形成された絶縁スタック構造（図示せず）が、下部相互接続構造1032上に形成される。次に、任意の既知の方法に従って、電荷蓄積構造40のトンネル層14、チャネルピラー16、ならびに導電性ピラー32aおよび32bが絶縁スタック構造内に形成される。トンネル層14の材料は、酸化ケイ素などの誘電体材料であってもよい。チャネルピラー16の材料は、ドーピングされていないポリシリコンなどの半導体であってもよい。導電性ピラー32a、32bは、例えばドーピングされたポリシリコンである。

【0053】

次いで、リソグラフィおよびエッチングプロセスを実行して、絶縁スタック構造内にスリットレンチ1110を形成し、絶縁スタック構造を複数のセクタBに分割する。

【0054】

その後、ゲート置換プロセスを実行して、ゲートスタック構造52を形成する。最初に、エッチングプロセスを実行して、エッチング液をスリットレンチ1110に注入して絶縁スタック構造内の別の絶縁層を除去し、複数の水平開口部34を形成し、次いで水平開口部34内にゲート層38を形成する。一部の実施形態では、ゲート層38が形成される前に、電荷蓄積層12およびブロッキング層36も水平開口部34内に形成される。電

荷蓄積層 12 は、例えば、窒化ケイ素である。ブロッキング層 36 の材料は、例えば、酸化アルミニウム (Al_2O_3)、酸化ハフニウム (HfO_2)、酸化ランタン (La_2O_5)、遷移金属酸化物、ランタニド酸化物、またはそれらの組合せなどの、7 以上の誘電率を有する高誘電率材料である。ゲート層 38 は、例えば、タングステンである。一部の実施形態では、ゲート層 38 が形成される前に、バリア層 37 が形成される。バリア層 37 の材料は、例えば、チタン (Ti)、窒化チタン (TiN)、タンタル (Ta)、窒化タンタル (Ta_nN)、またはそれらの組合せである。

【0055】

次に、スリットレンチ 1110 にスリット SLT を形成する。スリット SLT を形成する方法は、ゲートスタック構造 52 上およびスリットレンチ 1110 内に絶縁性充填材料を充填し、次いでエッチバックプロセスまたは平坦化プロセスによってゲートスタック構造 52 上の過剰な絶縁性充填材料を除去することを含む。絶縁性充填材料は、例えば、酸化ケイ素である。

10

【0056】

その後、上部相互接続構造 1034 (ローカルビット線 LBL_n 、ローカルソース線 LSL_n 、グローバルビット線 GBL_n 、およびグローバルソース線 GSL_n を含む) が、ゲートスタック構造 52 上に形成される。上部相互接続構造 1034 は、ダマシン、デュアルダマシンなどの任意の既知の方法によって形成することができるが、本明細書では詳細に説明しない。

【0057】

図 3A および図 3B を参照すると、本実施形態では、上部相互接続構造 1034 (ローカルビット線 LBL_n 、ローカルソース線 LSL_n 、グローバルビット線 GBL_n 、およびグローバルソース線 GSL_n を含む) が形成された後、上部相互接続構造 1034 の上方にヒータ 1200 がさらに形成される。ヒータ 1200 を形成する方法は、例えば、最初に上部相互接続構造 1034 の上方に誘電体層 1040 を形成することを含む。誘電体層 1040 の材料は、例えば、酸化ケイ素である。一部の実施形態では、誘電体層 1040 が平坦な表面を有するように、化学機械平坦化プロセスなどの平坦化プロセスがさらに実行される。その後、リソグラフィおよびエッチングプロセスを実行して、誘電体層 1040 内に複数の溝 OP1 を形成する。次いで、誘電体層 1040 上および溝内にバリア材料層および金属材料層を順次形成する。次に、化学機械平坦化プロセスなどの平坦化プロセスを実行して、誘電体層 1040 の表面上のバリア材料層および金属材料層を除去し、溝内にバリア層 1204 および金属層 1202 を形成する。金属材料層は、例えば、銅またはタングステンである。バリア材料層は、例えば、チタン、タンタル、窒化チタン、窒化タンタル、またはそれらの組合せである。

20

30

【0058】

図 3B を参照すると、ヒータ 1200 を形成した後、接合層 1300 が形成される。接合層 1300 の形成方法は、以下の通りである。最初に、ヒータ 1200 および誘電体層 1040 上に絶縁層 1304 を形成した後、リソグラフィおよびエッチングプロセスを実行して、絶縁層 1304 に複数のパッド開口部 OP2 を形成する。パッド開口部 OP2 の底部は、ヒータ 1200 を露出させる。その後、絶縁層 1304 上およびパッド開口部 OP2 内に導電層を形成する。次いで、化学機械平坦化プロセスなどの平坦化プロセスを実行して、絶縁層 1304 上の導電層を除去し、パッド開口部 OP2 内にパッド 1302 を形成する。

40

【0059】

上記の実施形態では、メモリチップ 1000 のヒータ 1200 は、上部相互接続構造 1034 が形成された後に形成される。他の実施形態では、メモリチップ 1000 のヒータ 1200 は、上部相互接続構造 1034 が形成される前に形成されてもよい。

【0060】

図 4C を参照すると、メモリチップ 1000 のヒータ 1200 は、3D フラッシュメモリ構造 1100 のゲートスタック構造 52 が形成された後、上部相互接続構造 1034 (

50

ローカルビット線 LBL_n 、ローカルソース線 LSL_n 、グローバルビット線 GBL_n 、およびグローバルソース線 GSL_n を含む)が形成される前に、ゲートスタック構造 52 間のスリットトレンチ 1110 内に形成される。

【0061】

図 4 A および図 4 C を参照すると、ヒータ 1200 を形成する方法は、例えば、最初にスリットトレンチ 1110 内にライナ材料層を形成することを含む。ライナ材料層は、例えば、酸化ケイ素または窒化ケイ素である。次に、ゲートスタック構造 52 上およびスリットトレンチ 1110 内にバリア材料層および金属材料層を順次形成する。次いで、化学機械平坦化プロセスなどの平坦化プロセスを実行して、ゲートスタック構造 52 の表面上のバリア材料層および金属材料層を除去し、スリットトレンチ 1110 内に絶縁ライナ層 1112、バリア層 1204、および金属層 1202 を形成する。金属材料層は、例えば、銅またはタンゲステンである。バリア材料層は、例えば、チタン、タンタル、窒化チタン、窒化タンタル、またはそれらの組合せである。

10

【0062】

図 4 B および図 4 C を参照すると、ヒータ 1200 が形成された後、上部相互接続構造 1034 (ローカルビット線 LBL_n 、ローカルソース線 LSL_n 、グローバルビット線 GBL_n 、およびグローバルソース線 GSL_n を含む)が形成される。その後、接合層 1300 が、上述の方法に従って上部相互接続構造 1034 上に形成される。

【0063】

図 7 A を参照すると、複数の制御チップ 2000 が設けられている。制御チップ 2000 の形成方法は、以下の通りである。図 5 C を参照すると、第 2 のトランジスタ 2020 が第 2 の基板 (ウエハ) 2010 上に形成される。次いで、第 2 の相互接続構造 2030 が第 2 のトランジスタ 2020 上に形成される。第 2 の相互接続構造 2030 は、ダマシン、デュアルダマシンなどの任意の既知の方法によって形成することができる。その後、接合層 2050 が、上述の方法に従って第 2 の相互接続構造 2030 上に形成される。次に、ダイシングを行って複数の制御チップ 2000 を形成する。

20

【0064】

図 7 B を参照すると、制御チップ 2000 の接合層 2050 とメモリチップ 1000 の接合層 1300 とが接合されて、接合構造 3000 が形成されている。接合方法は、例えば、ハイブリッド接合プロセスである。一部の実施形態では、制御チップ 2000 がウエハ 1010 W 上のメモリチップ 1000 と接合された後、制御チップ 2000 の側壁の周りに封止層 (図示せず) がさらに形成される。

30

【0065】

図 7 C を参照すると、ダイシングプロセスを行って互いに独立した複数の 3D フラッシュメモリモジュールチップ 5000 を形成する。

【0066】

以上をまとめると、本開示では、メモリチップと制御チップを接合して 3D フラッシュメモリモジュールチップを形成する。制御チップのドライバが高い駆動電流を供給してメモリチップ内のヒータを加熱することで、フラッシュメモリの電荷蓄積構造を修復して、より高い消去速度を達成し、フラッシュメモリの耐久性を向上させることが可能である。さらに、制御チップは、メモリチップの制御ロジックユニットのステータス信号に従って、対応するセクタを局所的に加熱することができる。加えて、接合によって形成された 3D フラッシュメモリモジュールチップでは、制御チップを別個に製造することができ、メモリチップ内に大面積のヒータコントローラを形成する必要がない。したがって、ヒータコントローラがメモリチップの面積を占有することを防止することができ、制御チップをそれほど高度でないプロセスで製造することができ、プロセスのコストを削減することができる。

40

【産業上の利用可能性】

【0067】

本発明の 3D フラッシュメモリモジュールチップおよびその製造方法は、3D メモリデ

50

バイスおよびその製造方法に適用することができる。

【符号の説明】

【0068】

14	: トンネル層	
16	: チャネルピラー	
20	: メモリセル	
24	: 絶縁性充填層	
28	: 絶縁性ピラー	
32a	: 導電性ピラー / ソースピラー	
32b	: 導電性ピラー / ドレインピラー	10
34	: 水平開口部	
36	: ブロッキング層	
37, 1204	: バリア層	
38, 2028	: ゲート層	
40	: 電荷蓄積構造	
52	: ゲートスタック構造	
54, 1304, 2054	: 絶縁層	
1000	: メモリチップ	
1010	: 第1の基板	
1010W	: ウエハ	20
1020	: 能動デバイス (第1のトランジスタ)	
1030	: 第1の相互接続構造	
1032	: 下部相互接続構造	
1034	: 上部相互接続構造	
1040, 2031	: 誘電体層	
1100, 1100 ₁ , 1100 ₂	: 3Dフラッシュメモリ構造	
1110	: スリットレンチ	
1112	: 絶縁ライナ層	
1200, 1200 ₁ , 1200 ₂ , 1200 ₃	: ヒータ	
1202	: 金属層	30
1300, 2050	: 接合層	
3000	: 接合構造	
1302, 1302a, 1302b, 2052, 2052a, 2052b	: パッド	
2000	: 制御チップ	
2000C, 2000C ₁	: 列	
2000R, 2000R ₁	: 駆動行	
2010	: 第2の基板	
2012	: 活性領域	
2020, 2020 ₁	: 第2のトランジスタ	
2022a	: ソース領域	40
2022b	: ドレイン領域	
2024	: ゲート誘電体層	
2030	: 第2の相互接続構造	
2032, 2034, C1, C2, C3	: コンタクト	
2036, 2040, 2040a, 2040b, 2040c	: 導電線	
2038, 2042, 2042a, 2042b	: ビア	
2100	: グローバル電源	
2200	: 行デコーダ	
2300	: 列デコーダ	
5000	: 3Dフラッシュメモリモジュールチップ	50

- A 0 , A 1 , A 2 : 行アドレス信号
- A 3 , A 4 : 列アドレス信号
- A R : アレイ領域
- B , B 1 , B 2 , B 3 , B 4 : セクタ
- B M 1 : 下部第 1 金属層
- B M 2 : 下部第 2 金属層
- B M 3 : 下部第 3 金属層
- B V 1 , B V 2 , T V 1 : ピア
- T M 1 : 上部第 1 金属層
- T M 2 : 上部第 2 金属層
- E 1 : 第 1 の端部
- E 2 : 第 2 の端部
- O P 1 : 溝
- O P 2 : パッド開口部
- S C : 階段構造
- S L : スクライプライン
- S L T : スリット
- S R : 階段領域
- T , T ' , T 1 , T 1 ' , T 2 , T 2 ' , T 3 , T 3 ' , T 4 , T 4 ' : タイル
- W 1 , W 2 : 幅
- I - I ' , I I - I I ' : 線
- X , Y , Z : 方向

10

20

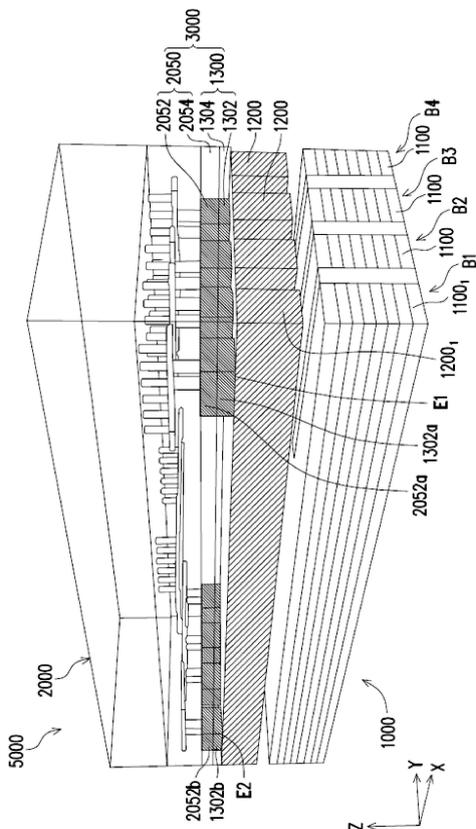
30

40

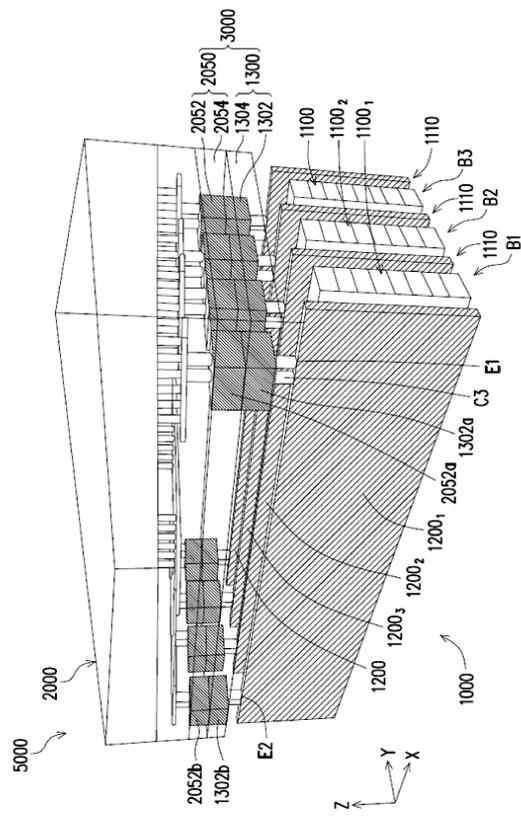
50

【 図 面 】

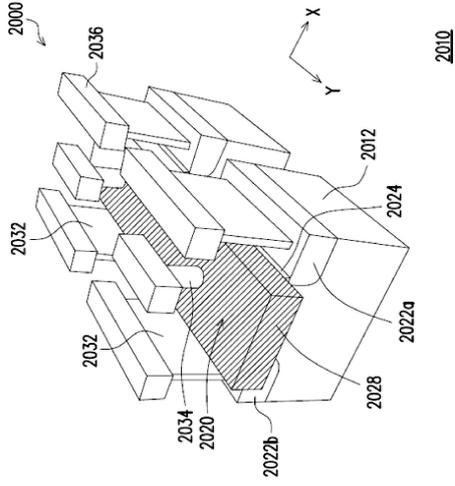
【 図 1 A 】



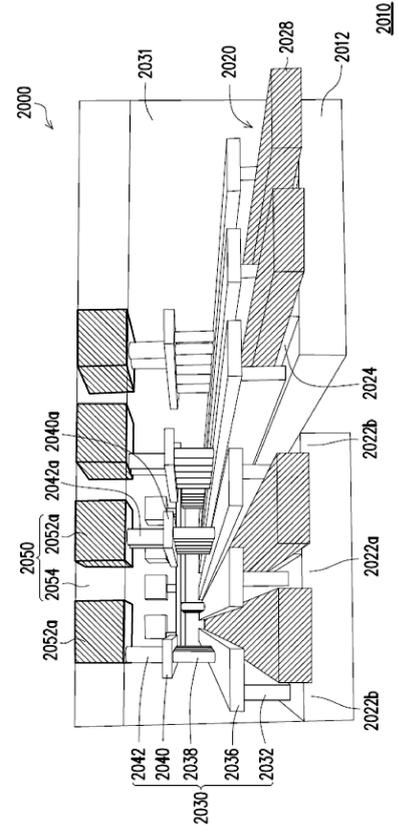
【 図 1 B 】



【 図 5 B 】



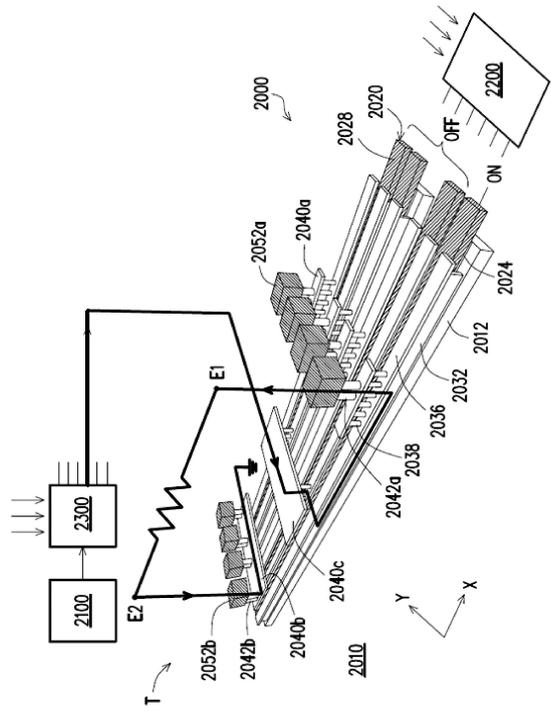
【 図 5 C 】



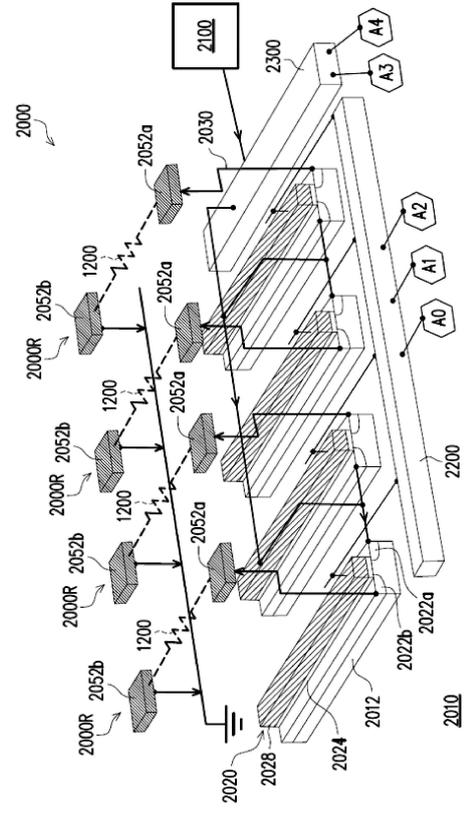
10

20

【 図 5 D 】



【 図 5 E 】

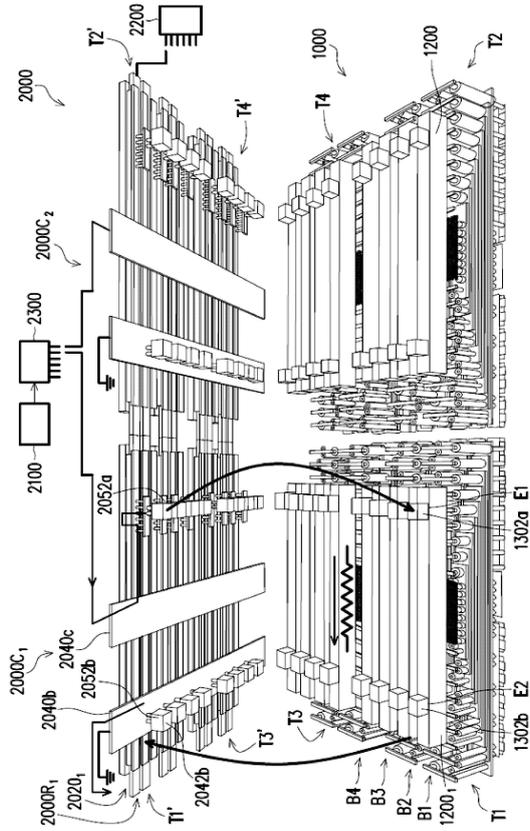


30

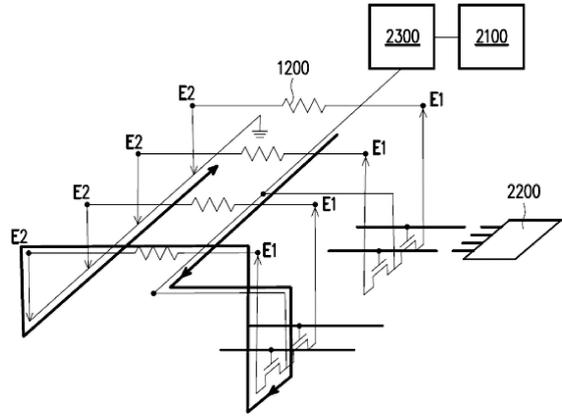
40

50

【 図 6 A 】



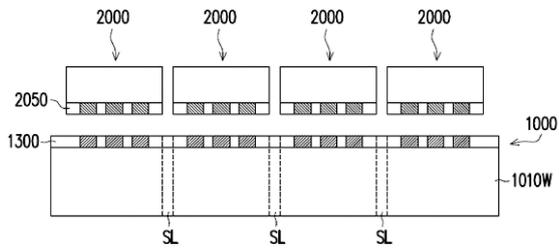
【 図 6 B 】



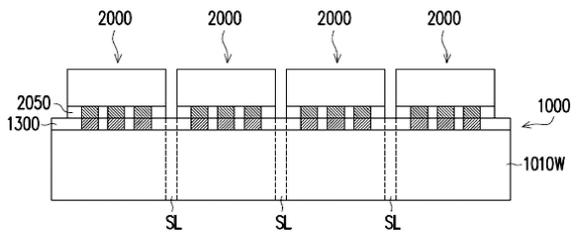
10

20

【 図 7 A 】



【 図 7 B 】

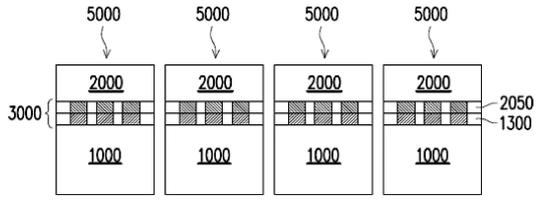


30

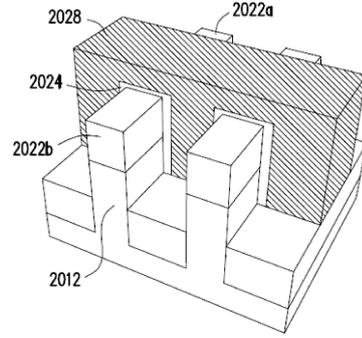
40

50

【 図 7 C 】



【 図 8 】



10

【 外国語明細書 】

20

30

40

50

File: 112718jpf

3D FLASH MEMORY MODULE CHIP AND METHOD OF FABRICATING THE SAME

BACKGROUND

5 [Technical Field] 10

[0001] The embodiment of the disclosure relates to a semiconductor module and a method of fabricating the same, and particularly, to a 3D flash memory module and a method of fabricating the same.

[Description of Related Art]

10 [0002] Since a non-volatile memory has the advantage that stored data does not disappear at power-off, it becomes a widely used memory for a personal computer or other electronics equipment. Currently, the three-dimensional (3D) memory commonly used in the industry includes a NOR flash memory and a NAND flash memory. In addition, another type of 3D memory is an AND flash memory, which can be applied to a multi-dimensional memory array 20

15 with high integration and high area utilization, and has an advantage of a fast operation speed. Therefore, the development of a 3D memory device has gradually become the current trend.

SUMMARY

TECHNICAL PROBLEM

20 [0003] The disclosure provides a 3D flash memory module chip and a method of fabricating the same, which can perform a local healing process on a flash memory. 30

SOLUTION TO PROBLEM

[0004] In an embodiment of the disclosure, a 3D flash memory module chip includes a memory chip and a control chip. The memory chip includes a plurality of tiles and a plurality of heaters.

25 The tiles each include a plurality of 3D flash memory structures. The heaters are disposed around

File: 112718jpf

the 3D flash memory structures of each of the tiles. The control chip is bonded with the memory chip to drive at least one of the heaters.

[0005] In an embodiment of the disclosure, a method of fabricating a 3D flash memory module chip includes the following steps. A memory chip is formed, the step including forming a plurality of tiles on a first substrate, each of the tiles including a plurality of 3D flash memory structures; and forming a plurality of heaters around the 3D flash memory structures of each of the tiles. A control chip is formed. The control chip and the memory chip are bonded, and the control chip is configured to drive the heaters.

10

EFFECTS OF INVENTION

[0006] Based on the above, in the 3D flash memory module chip and the method of fabricating the same according to the disclosure, an additional control chip is used to drive the heater to perform a local healing process on each sector of the flash memory. The control chip may be manufactured separately to prevent the heater controller from occupying the area of the memory chip, and the control chip may be manufactured by a less advanced process to reduce the cost of the process.

20

BRIEF DESCRIPTION OF THE DRAWINGS

[0007] FIG. 1A and FIG. 1B are respectively schematic perspective views of a 3D flash memory module chip according to an embodiment of the disclosure.

[0008] FIG. 2A is a partial top view of a 3D flash memory structure of a memory chip according to an embodiment of the disclosure.

30

[0009] FIG. 2B is a cross-sectional view taken along line I-I' of FIG. 2A.

[0010] FIG. 3A is a partial top view of a memory chip having a heater according to another embodiment of the disclosure.

[0011] FIG. 3B is a cross-sectional view taken along line I-I' of FIG. 3A.

40

50

File: 112718jpf

[0012] FIG. 4A is a partial top view of a memory chip having a heater according to another embodiment of the disclosure.

[0013] FIG. 4B is a partial top view of a heater and a pad of a memory chip according to another embodiment of the disclosure.

5 **[0014]** FIG. 4C is a cross-sectional view taken along line II-II' of FIG. 4B. 10

[0015] FIG. 5A to FIG. 5E are schematic perspective views of a control chip according to an embodiment of the disclosure.

[0016] FIG. 6A is a schematic perspective view of a memory chip and a control chip according to an embodiment of the disclosure.

10 **[0017]** FIG. 6B is a schematic circuit view of FIG. 6A.

[0018] FIG. 7A to FIG. 7C show schematic cross-sectional views of a process of fabricating a 3D flash memory module chip of the disclosure.

[0019] FIG. 8 is a schematic perspective view of a control chip according to another embodiment of the disclosure. 20

15

DESCRIPTION OF THE EMBODIMENTS

[0020] The performance of a flash memory is significantly reduced after multiple operations, so it is necessary to perform a healing process on the flash memory. In the healing process, a heater may be used to heat the flash memory to heal a charge storage structure (e.g., a nitride layer) of the flash memory. In the current art, word lines are most commonly used as the heater. However, due to the large number of word lines and the complicated configuration relationship with other components (e.g., a word line decoder), the layout design of the flash memory structure may be more difficult. 30

25 **[0021]** The embodiments of the disclosure provide several 3D flash memory module chips, in

10

20

30

40

50

File: 112718jpf

which a heater is disposed above, or around sidewalls of, a 3D flash memory structure of a memory chip, and the memory chip is bonded with a control chip so that the control chip can drive a heater to perform a heating process on a local sector of the memory chip.

[0022] FIG. 1A and FIG. 1B are respectively schematic perspective views of a 3D flash memory module chip according to an embodiment of the disclosure. FIG. 2A is a partial top view of a 3D flash memory structure of a memory chip according to an embodiment of the disclosure. FIG. 2B is a cross-sectional view taken along line I-I' of FIG. 2A. FIG. 3A is a partial top view of a memory chip having a heater according to another embodiment of the disclosure. FIG. 3B is a cross-sectional view taken along line I-I' of FIG. 3A.

[0023] Referring to FIG. 1A and FIG. 1B, a 3D flash memory module chip (also referred to as a 3D integrated circuit (3D IC)) 5000 according to an embodiment of the disclosure includes a memory chip 1000 and a control chip 2000. The memory chip 1000 includes a plurality of 3D flash memory structures 1100 and a plurality of heaters 1200. The heaters 1200 are disposed around the 3D flash memory structures 1100. In some embodiments, the heaters 1200 are disposed above the 3D flash memory structures 1100, as shown in FIG. 1A. In other embodiments, the heaters 1200 are disposed in slit trenches 1110 between the 3D flash memory structures 1100, as shown in FIG. 1B. The control chip 2000 is disposed above the memory chip 1000 to drive the heaters 1200 in the memory chip 1000. The control chip 2000 and the memory chip 1000 may be bonded to each other by a bonding structure 3000.

[0024] Referring to FIG. 1A and FIG. 1B, the 3D flash memory structure 1100 of the memory chip 1000 may be a 3D AND flash memory structure (as shown in FIG. 2A and FIG. 2B), a 3D NAND flash memory structure (not shown), or a 3D NOR flash memory structure (not shown). The 3D AND flash memory structure will be taken as an example to illustrate the 3D flash memory structure 1100 of the disclosure, but the embodiment of the disclosure is not limited thereto.

[0025] Referring to FIG. 2A and FIG. 2B, the memory chip 1000 may include a plurality of tiles

File: 112718jpf

T. The tiles T may be arranged in an array including a plurality of columns and a plurality of rows. In this embodiment, four tiles T (e.g., T1 to T4) are shown for illustration. Among the four tiles T, the tile T1 and the tile T2 are arranged in a row, and the tile T3 and the tile T4 are arranged in another row. The tile T1 and the tile T3 are arranged in a column, and the tile T2 and the tile T4 are arranged in another column. Each of the tiles T may include a plurality of sectors B (e.g., B1 to B4). Each of the sectors B includes a 3D flash memory structure 1100. The 3D flash memory structures 1100 extend in the X direction and are arranged in the Y direction. Two adjacent 3D flash memory structures 1100 are separated from each other by a slit trench 1110.

[0026] Referring to FIG. 2B, each of the 3D flash memory structures 1100 may include at least a memory array formed by a plurality of memory cells. Specifically, the 3D flash memory structure 1100 may be disposed above one or more active devices (e.g., first transistors 1020) on a first substrate (e.g., a semiconductor substrate) 1010. The first transistor 1020 is, for example, a complementary metal-oxide-semiconductor (CMOS) field-effect transistor. Therefore, this architecture may also be referred to as a complementary metal-oxide-semiconductor field-effect transistor under array (CMOS under Array (CUA)) architecture.

[0027] Referring to FIG. 2B, the 3D flash memory structure 1100 may be disposed in a back end of line (BEOL) of a semiconductor die. For example, the 3D flash memory structure 1100 may be embedded in a first interconnect structure 1030. The first interconnect structure 1030 includes, for example, a lower interconnect structure 1032 and an upper interconnect structure 1034. The lower interconnect structure 1032 is disposed above one or more active devices (e.g., the first transistors 1020) on the first substrate (e.g., the semiconductor substrate) 1010 and below the memory array of the 3D flash memory structure 1100. The upper interconnect structure 1034 is disposed above the memory array of the 3D flash memory structure 1100. The lower interconnect structure 1032 includes, for example, a lower first metal layer BM1, a lower second metal layer BM2, and a lower third metal layer BM3, and vias BV1 and BV2 therebetween. The

File: 112718jpf

upper interconnect structure 1034 includes, for example, an upper first metal layer TM1 and an upper second metal layer TM2, and vias TV1 therebetween. The numbers of the metal layers and the vias of the lower interconnect structure 1032 and the upper interconnect structure 1034 are not limited to the above.

5 **[0028]** Referring to FIG. 2B, the 3D flash memory structure 1100 includes a plurality of gate stack structures 52. Each of the gate stack structures 52 is formed on the lower interconnect structure 1032. Each of the gate stack structures 52 extends in the X direction, from an array region AR to staircase regions SR of the first substrate 1010. The gate stack structure 52 includes a plurality of gate layers (also referred to as word lines) 38 and a plurality of insulating layers 54
10 vertically stacked on the surface of the first substrate 1010. In the Z direction, the gate layers 38 are electrically isolated from each other by the insulating layers 54 disposed therebetween. The gate layer 38 includes a metal layer such as tungsten. In some embodiments, the gate layer 38 further includes a barrier layer 37, such as titanium (Ti), titanium nitride (TiN), tantalum (Ta),
20 tantalum nitride (TaN), or a combination thereof. The insulating layer 54 is, for example, silicon oxide.
15

[0029] The gate layer 38 extends in a direction parallel to the surface of the first substrate 1010 (shown in FIG. 2D). The gate layers 38 in the staircase region SR may have a staircase structure SC (shown in FIG. 2B), so that a lower gate layer 38 is longer than an upper gate layer 38, and the end of a lower gate layer 38 extends laterally beyond the end of an upper gate layer 38. A contact
20 C1 for connecting the gate layer 38 may land on the end of the gate layer 38 in the staircase region SR to connect each of the gate layers 38 to conductive lines of the lower interconnect structure 1032 (e.g., the conductive line of the lower third metal layer BM3) via the contact C1 and the upper interconnect structure 1034.
30

[0030] Referring to FIG. 2B, the 3D flash memory structure 1100 further includes a plurality of channel pillars 16. The channel pillar 16 continuously extends through the gate stack structure
25

File: 112718jpf

52 in the array region AR. In some embodiments, the channel pillar 16 may have a ring-shaped profile in a top view. The material of the channel pillar 16 may be a semiconductor, such as undoped polysilicon.

[0031] Referring to FIG. 2B, the 3D flash memory structure 1100 further includes an insulating filling layer 24, an insulating pillar 28, a plurality of conductive pillars (e.g., serving as source pillars) 32a, and a plurality of conductive pillars (e.g., serving as drain pillars) 32b. The conductive pillars 32a and 32b and the insulating pillar 28 are disposed in the channel pillar 16 and each extend in a direction (i.e., the Z direction) perpendicular to the gate layer 38. The conductive pillars 32a and 32b are separated from each other by the insulating filling layer 24 and the insulating pillar 28 and are electrically coupled to the channel pillar 16. The conductive pillars 32a and 32b are, for example, doped polysilicon. The insulating filling layer 24 is, for example, silicon oxide, the insulating pillar 28 is, for example, silicon nitride.

[0032] Referring to FIG. 2B, a charge storage structure 40 is disposed between the channel pillar 16 and the gate layers 38. The charge storage structure 40 may include a tunneling layer (or referred to as a bandgap engineered tunneling oxide layer) 14, a charge storage layer 12, and a blocking layer 36. The charge storage layer 12 is located between the tunneling layer 14 and the blocking layer 36. In some embodiments, the tunneling layer 14, the charge storage layer 12, and the blocking layer 36 are, for example, silicon oxide, silicon nitride, and silicon oxide. In some embodiments, a part (e.g., the tunneling layer 14) of the charge storage structure 40 continuously extends in a direction (i.e., the Z direction) perpendicular to the gate layer 38, and the other part (e.g., the charge storage layer 12 and the blocking layer 36) of the charge storage structure 40 surrounds the gate layer 38, as shown in FIG. 2B. In other embodiments, the charge storage structure 40 (e.g., the tunneling layer 14, the charge storage layer 12, and the blocking layer 36) surrounds the gate layer 38 (not shown). Each of the gate layers 38, and the charge storage structure 40, the channel pillar 16, the source pillar 32a, and the drain pillar 32b that are

File: 112718jpf

surrounded by the gate layer 38 define a memory cell 20. Therefore, each of the 3D flash memory structures 1100 includes at least a memory array composed of a plurality of memory cells 20.

[0033] The 3D flash memory structure 1100 further includes a local bit line LBL_n , a local source line LSL_n , a global bit line GBL_n , and a global source line GSL_n . The local bit line LBL_n and the local source line LSL_n are located in the upper first metal layer TM1 of the upper interconnect structure 1034, and are respectively electrically connected to the source pillar 32a and the drain pillar 32b via contacts C2. The global bit line GBL_n and the global source line GSL_n are respectively electrically connected to the local bit line LBL_n and the local source line LSL_n via upper vias (not shown) in the upper interconnect structure 1034.

[0034] According to different operation methods, a 1-bit operation or a 2-bit operation may be performed on the memory cell 20. For example, when a voltage is applied to the source pillar 32a and the drain pillar 32b, since the source pillar 32a and the drain pillar 32b are connected to the channel pillar 16, electrons may be transferred along the channel pillar 16 and stored in the entire charge storage structure 40. Accordingly, a 1-bit operation may be performed on the memory cell 20. In addition, for an operation involving Fowler-Nordheim tunneling, electrons or holes may be trapped in the charge storage structure 40 between the source pillar 32a and the drain pillar 32b. For an operation involving source side injection, channel-hot-electron injection, or band-to-band tunneling hot carrier injection, electrons or holes may be locally trapped in the charge storage structure 40 adjacent to one of the source pillar 32a and the drain pillar 32b. Accordingly, a single level cell (SLC, 1 bit) or multi-level cell (MLC, greater than or equal to 2 bits) operation may be performed on the memory cell 20.

[0035] During operation, a voltage is applied to a selected word line (gate layer) 38; for example, when a voltage higher than a corresponding threshold voltage (V_{th}) of the corresponding memory cell 20 is applied, a channel region of the channel pillar 16 intersecting the selected word line 38

File: 112718jpf

is turned on to allow a current to enter the drain pillar 32b from a bit line BL_n , flow to the source pillar 32a via the turned-on channel region, and finally flow to a source line SL_n .

[0036] Referring to FIG. 3A and FIG. 3B, the memory chip 1000 further includes a plurality of heaters 1200. The heater 1200 may be disposed in a dielectric layer 1040 above the 3D flash memory structure 1100. The material of the dielectric layer 1040 is, for example, silicon oxide. The heater 1200 includes a metal layer 1202, such as copper or tungsten. In some embodiments, the heater 1200 further includes a barrier layer 1204, such as titanium, tantalum, titanium nitride, tantalum nitride, or a combination thereof.

[0037] Referring to FIG. 3A, in some embodiments, one heater 1200 is disposed on each sector B, and two heaters 1200 of any two adjacent sectors B are separated from each other. The heater 1200 may extend in the X direction. In an embodiment, the heater 1200 is disposed in the array region AR and extends to the staircase regions SR (as shown in FIG. 3A and FIG. 3B). In an embodiment, the heater 1200 may be disposed in the array region AR but is not disposed in the staircase regions SR (not shown). In other words, the length of the heater 1200 may be greater than, equal to, or less than the length of the 3D flash memory structure 1100 in the X direction.

[0038] In addition, multiple heaters 1200 may be disposed on each sector B; for example, one heater 1200 may be disposed in the array region AR and the staircase regions SR respectively and may perform heating separately (not shown). However, the embodiment of the disclosure is not limited thereto. In another embodiment, multiple heaters 1200 of adjacent two, three, or more sectors B may also be combined into one heater (not shown) to simultaneously heat the 3D flash memory structures 1100 of multiple sectors B.

[0039] Referring to FIG. 3A, the shape of the heater 1200 in a top view is, for example, a rectangle or another shape. The heaters 1200 on multiple sectors B may have the same width or different widths. A width $W1$ of the heater 1200 in the array region AR is the same as a width $W2$ of the heater 1200 in the staircase region SR. However, the disclosure is not limited thereto.

File: 112718jpf

The shape of the heater 1200 may be changed according to the actual requirements or design. The width W1 of the heater 1200 in the array region AR may be greater than, equal to, or less than the width W2 of the heater 1200 in the staircase region SR.

[0040] Referring to FIG. 1A, FIG. 1B, and FIG. 3B, the memory chip 1000 further includes a bonding layer 1300. The bonding layer 1300 includes a pad 1302 and an insulating layer 1304. The insulating layer 1304 is disposed on the heater 1200. The material of the insulating layer 1304 is, for example, silicon oxide. The pad 1302 is disposed in the insulating layer 1304 on the surface of each of the heaters 1200. The material of the pad 1302 is, for example, copper. The pad 1302 includes pads 1302a and 1302b. The pads 1302a and 1302b are respectively connected to a first end E1 and a second end E2 of the heater 1200.

[0041] In the above embodiment, the 3D flash memory structures 1100 are 3D AND flash memory structures, and the heaters 1200 are disposed above the 3D AND flash memory structures (as shown in FIG. 3A, FIG. 3B, and FIG. 6A). In other embodiments, the 3D flash memory structures 1100 are 3D AND flash memory structures, and the heaters 1200 are disposed in slit trenches 1110 between the 3D AND flash memory structures (as shown in FIG. 4A to FIG. 4C).

[0042] FIG. 4A is a partial top view of a memory chip having a heater according to another embodiment of the disclosure. FIG. 4B is a partial top view of a heater and a pad of a memory chip according to another embodiment of the disclosure. FIG. 4C is a cross-sectional view taken along line II-II' of FIG. 4B.

[0043] Referring to FIG. 4A and FIG. 4C, a plurality of heaters 1200 are disposed in slit trenches 1110 between 3D flash memory structures 1100. The heater 1200 is disposed around a plurality of gate layers 38 and a plurality of insulating layers 54 of a gate stack structure 52. The heater 1200 is separated from the gate layers 38 and the insulating layers 54 by an insulating liner layer 1112 (as shown in FIG. 4C). The insulating liner layer 1112 includes an insulating material such as silicon oxide or silicon nitride. The heater 1200 includes a metal layer 1202 (as shown in FIG.

File: 112718jpf

4C), such as copper or tungsten. In some embodiments, the heater 1200 further includes a barrier layer 1204 (as shown in FIG. 4C). The barrier layer 1204 is located between the insulating liner layer 1112 and the metal layer 1202. The barrier layer 1204 is, for example, titanium, tantalum, titanium nitride, tantalum nitride, or a combination thereof.

5 **[0044]** In some embodiments, one heater 1200 is disposed in each of the slit trenches 1110. For example, the heater 1200 may extend in the X direction. In an embodiment, the heater 1200 is disposed in the array region AR and extends to the staircase regions SR (as shown in FIG. 4A and FIG. 4B). In an embodiment, the heater 1200 may be disposed in the array region AR but is not disposed in the staircase regions SR (not shown). In other words, the length of the heater
10 1200 may be greater than, equal to, or less than the length of the 3D flash memory structure 1100 in the X direction.

[0045] Alternatively, multiple heaters 1200 may be disposed in each of the slit trenches 1110. For example, one heater 1200 may be provided respectively in the array region AR and the staircase region SR, and heating may be performed separately (not shown). However, the
15 embodiment of the disclosure is not limited thereto.

[0046] Referring to FIG. 4A, in addition, the shape of the heater 1200 in a top view is, for example, a rectangle or another shape. The heaters 1200 in multiple slit trenches 1110 may have the same width or different widths. However, the disclosure is not limited thereto. The shape of the heater 1200 may be changed according to the actual requirements or design.

20 **[0047]** Referring to FIG. 4B and FIG. 4C, a contact C3 is disposed respectively on surfaces of two ends (i.e., E1 and E2) of each of the heaters 1200. The contact C3 may be connected to pads 1302a and 1302b above via an upper interconnect structure 1034, so that the heater 1200 of the memory chip 1000 can be electrically connected to the control chip 2000 via the upper interconnect structure 1034 and the pads 1302a and 1302b. The material of the pads 1302a and
25 1302b is, for example, copper.

File: 112718jpf

[0048] FIG. 5A to FIG. 5E are schematic perspective views of a control chip according to an embodiment of the disclosure. FIG. 6A is a schematic perspective view of a memory chip and a control chip according to an embodiment of the disclosure. FIG. 6B is a schematic circuit view of FIG. 6A.

5 **[0049]** Referring to FIG. 5A, the control chip 2000 may include a plurality of tiles T'. The tiles T' may be arranged in an array. In this embodiment, four tiles T' (e.g., T1' to T4') will be taken as an example for illustration. Among the four tiles T', the tile T1' and the tile T2' are arranged in a row, and the tile T3' and the tile T4' are arranged in another row. The tile T1' and the tile T3' are arranged in a column, and the tile T2' and the tile T4' are arranged in another
10 column.

[0050] Referring to FIG. 5A and FIG. 5E, each of the tiles T' includes a plurality of driving rows 2000R and columns 2000C. Each of the driving rows 2000R includes a second transistor 2020, a second interconnect structure 2030, and a pad 2052, as shown in FIG. 5E. The second transistor 2020 is disposed on an active region 2012 of a second substrate 2010. The second
15 substrate 2010 may be a semiconductor substrate, such as a silicon substrate. The second transistor 2020 may be a complementary metal-oxide-semiconductor (CMOS) transistor. The second transistor 2020 may be a planar transistor (as shown in FIG. 5A to FIG. 5E) or a fin transistor (as shown in FIG. 8).

[0051] Referring to FIG. 5E, and FIG. 8, the second transistor 2020 includes a gate dielectric layer 2024, a gate layer 2028, a source region 2022a, and a drain region 2022b. The gate dielectric layer 2024 is, for example, silicon oxide or a high dielectric constant material. The gate dielectric layer 2024 is, for example, silicon oxide or a high dielectric constant material. The gate layer 2028 is, for example, doped polysilicon or tungsten. The gate layer 2028 is located on the gate dielectric layer 2024. The gate layer 2028 has a strip shape, and its extending direction is, for example, the same as the extending direction of the heater 1200 (e.g., extending in the X
20 direction), as shown in FIG. 6A. In some embodiments, the gate layers 2028 of the second

10

20

30

40

50

File: 112718jpf

transistors 2020 in two adjacent rows (e.g., the tiles T1' and T2', or the tiles T3' and T4') may be electrically connected, as shown in FIG. 5A.

[0052] Referring to FIG. 5C and FIG. 5E, the source region 2022a and the drain region 2022b of the second transistor 2020 are disposed in the active region 2012 on two sides of the gate layer 1138. The source region 2022a and the drain region 2022b contain a dopant, such as an N-type or P-type dopant. In some embodiments, two adjacent second transistors 2020 share a source region 2022a.

[0053] Referring to FIG. 5B and FIG. 5C, the second interconnect structure 2030 is located on the second transistors 2020. The second interconnect structure 2030 includes a dielectric layer 2031 (as shown in FIG. 5C), and a plurality of contacts 2032 and 2034, a plurality of conductive lines 2036 and 2040, and a plurality of vias 2038 and 2042 which are located in the dielectric layer 2031. The contacts 2032 respectively land on the source region 2022a and the drain region 2022b, and are electrically connected to the source region 2022a and the drain region 2022b. The contact 2034 lands on the gate layer 2028 and is electrically connected to the gate layer 2028. The contact 2032 has a strip shape which extends along the X direction and is substantially parallel to the gate layer 2028, as shown in FIG. 5B and FIG. 5D. The shape of the contact 2034 is different from the shape of the contact 2032 and may be, for example, a columnar shape, as shown in FIG. 5B. The conductive lines 2036 and 2040 (as shown in FIG. 5C) are respectively disposed on the contacts 2032 and 2034. The conductive line 2036 and the conductive line 2040 are electrically insulated from each other by the via 2038. The via 2042 is disposed on the conductive line 2040 and is electrically connected the conductive line 2040 to a bonding layer 2050 above. The dielectric layer 2031 is, for example, silicon oxide. The contacts 2032 and 2034, the conductive lines 2036 and 2040, and the vias 2038 and 2042 include a metal layer such as tungsten or copper. The contacts 2032 and 2034, the conductive lines 2036 and 2040, and the vias 2038 and 2042 may further include a barrier layer (not shown), such as titanium, tantalum, titanium nitride, tantalum

File: 112718jpf

nitride, or a combination thereof.

[0054] Referring to FIG. 5C, the pad 2052 of each of the driving rows 2000R is a part of the bonding layer 2050 of the control chip 2000. The bonding layer 2050 includes the pad 2052 and an insulating layer 2054. The insulating layer 2054 is located on the second interconnect structure 2030. The pad 2052 is located in the insulating layer 2054 and is electrically connected to the via 2042 of the second interconnect structure 2030. The material of the pad 2052 is, for example, copper. The material of the insulating layer 2054 is, for example, silicon oxide.

[0055] Referring to FIG. 5A and FIG. 5E, the pad 2052 includes a pad 2052a and a pad 2052b. Specifically, each of the driving rows 2000R includes a pair of pads 2052a and 2052b disposed along the X direction. The pad 2052a is electrically connected to the first end E1 of the heater 1200; the pad 2052b is electrically connected to the second end E2 of the heater 1200 and is grounded, as shown in FIG. 1A, FIG. 1B and FIG. 6A. Referring to FIG. 5C and FIG. 5D, each pad 2052a is electrically connected to a conductive line 2040a below via a via 2042a. The conductive lines 2040a in the same tile T' are separated and electrically isolated from each other, so as to be respectively electrically connected to the drain region 2022b of the second transistor 2020, as shown in FIG. 5A and FIG. 5C. Each pad 2052b is electrically connected to a conductive line 2040b below via a via 2042b, as shown in FIG. 5D. The pads 2052b of the tiles T' in the same column (e.g., the tiles T1' and T3', or the tiles T2' and T4') are electrically connected to the ground via the same conductive line 2040b, as shown in FIG. 5A and FIG. 5D.

[0056] Referring to FIG. 5C, FIG. 1A, and FIG. 1B, the bonding layer 2050 of the control chip 2000 and the bonding layer 1300 of the memory chip 1000 are bonded to each other to form a bonding structure 3000. Specifically, the position of the insulating layer 2054 of the control chip 2000 and the position of the insulating layer 1304 of the memory chip 1000 correspond to each other and are bonded to each other. The positions of the pads 2052a and 2052b of the control chip 2000 and the positions of the pads 1302a and 1302b of the memory chip 1000 correspond to

File: 112718jpf

each other and are bonded to each other.

[0057] Referring to FIG. 5A, FIG. 5C and FIG. 5D, the column 2000C of the control chip 2000 electrically couples a plurality of shared source regions 2022a of multiple second transistors 2020 of the tiles T' in the same column (e.g., the tiles T1' and T3', or the tiles T2' and T4') to a global power supply 2100 via a conductive line 2040c.

[0058] Referring to FIG. 5A, FIG. 5C and FIG. 5D, the drain region 2022b of the second transistor 2020 of the control chip 2000 is connected to the second interconnect structure 2030 and the pad 2052a of the bonding layer 2050, as shown in FIG. 5C. The pad 2052a is electrically connected to the pad 1302a connected to the first end E1 of the heater 1200 of the memory chip 1000 as shown in FIG. 6A. In an embodiment, each of the driving rows 2000R of the control chip 2000 may control one heater 1200 of one corresponding sector B of the memory chip 1000, as shown in FIG. 6A and FIG. 6B.

[0059] Referring to FIG. 5E, in some embodiments, the control chip 2000 further includes a column decoder 2300 and a row decoder 2200. The column decoder 2300 is electrically connected to the global power supply 2100. After receiving column address signals A3 and A4, the column decoder 2300 selects multiple tiles (two tiles in this example, e.g., the tiles T1' and T3' in FIG. 5A) of one column (e.g., a column 2000C₁ in FIG. 5A). Accordingly, the global power supply 2100 is provided to the shared source regions 2022a of the second transistors 2020 of each of the tiles (e.g., the tiles T1' and T3' in FIG. 5A) of the selected column (e.g., the column 2000C₁ in FIG. 5A) via the conductive line 2040c (shown in FIG. 5A) of the second interconnect structure 2030. The row decoder 2200 is electrically connected to the gate layers 2028 of the second transistors 2020 of the driving rows 2000R. After receiving row address signals A0 to A2 (or referred to as control signals), the row decoder 2200 decodes the inputted row address signals to select and turn on one (e.g., a second transistor 2020₁ in FIG. 5A) or more of the second transistors 2020.

File: 112718jpf

[0060] Generally, the memory chip 1000 includes a control logic unit for controlling the memory array, and the register in the control logic unit stores a status signal of an erase count of the memory array of each sector B. When the erase count reaches a predetermined count, the status signal is sent to the control chip 2000.

5 **[0061]** Referring to FIG. 6A and FIG. 6B, during the healing process, the control chip 2000 may generate a row address signal and a column address signal corresponding to the tile T and the sector B (e.g., the sector B1 of the tile T1 in FIG. 6A) that need healing based on the received status signal, and transmit the row address signal and the column address signal respectively to the row decoder 2200 and the column decoder 2300. The column decoder 2300 selects one
10 column (e.g., the column 2000C₁ in FIG. 6A) according to the received column address signal to provide the global power supply 2100 to the conductive line 2040c of the tiles (e.g., the tiles T1' and T3' in FIG. 5A) located in this column (e.g., the column 2000C₁ in FIG. 6A). The row decoder 2200 selects and turns on the second transistor 2020_i of one driving row 2000R₁ according
20 to the received row address signal. Therefore, current can flow from the global power supply 2100 into the source region 2022a of the second transistor 2020_i via the conductive line 2040c, pass through the channel and the drain region 2022b of the second transistor 2020_i, then flow into the pad 1302a of the memory chip 1000 via the second interconnect structure 2030 and the pad 2052a, and then enter a first end E1 of a heater 1200 (e.g., 1200_i). Afterwards, the current flow
15 in the heater 1200_i and flows out a second end E2 of the heater 1200_i via the pad 1302b of the memory chip 1000, then enters the pad 2052b of the control chip 2000, and then is electrically
20 connected to the ground via the conductive line 2040b. In the embodiment of the disclosure, with the second transistor (driver) 2020 (e.g., 2020_i) of the control chip 2000, it is possible to provide a high driving current to a specific heater 1200 (e.g., 1200_i), so that the conductor serving as the heater 1200 (e.g., 1200_i) is heated to heal the charge storage layer in the 3D flash memory
25 structure 1100 of a specific sector B (e.g., B1) in a specific tile T (e.g., T1).

File: 112718jpf

[0062] Referring to FIG. 1A and FIG. 1B, in some embodiments, in the healing process, the control chip 2000 may drive one heater 1200 (e.g., 1200₁) to heal the charge storage layer in the 3D flash memory structure 1100 (e.g., 1100₁) of one sector B (e.g., B1). Referring to FIG. 1B, in other embodiments, when healing is performed, the control chip 2000 may also simultaneously
5 drive two heaters 1200 (e.g., 1200₂ and 1200₃) to heal the charge storage layer in the 3D flash
memory structure 1100 (e.g., 1100₂) of one sector B (e.g., B2).

10

[0063] FIG. 7A to FIG. 7C show schematic cross-sectional views of a process of fabricating a 3D flash memory module chip of the disclosure.

[0064] Referring to FIG. 7A, a wafer 1010W is provided, and a plurality of memory chips 1000
10 are formed on the wafer 1010W. Scribe lines SL are provided between the memory chips 1000. The method of forming the memory chip 1000 is as follows. Referring to FIG. 3B, one or more active devices (e.g., first transistors) 1020 are first formed on the wafer 1010W. Next, a lower interconnect structure 1032 is formed on the active devices 1020. The lower interconnect
20 structure 1032 may be formed by any known method, such as damascene, dual-damascene, and the like. Afterwards, an insulating stack structure (not shown) formed by alternately stacking one insulating layer (e.g., silicon oxide) 54 and another insulating layer (not shown, e.g., silicon nitride) is formed on the lower interconnect structure 1032. Next, according to any known method, a tunneling layer 14 of a charge storage structure 40, a channel pillar 16, and conductive pillars 32a and 32b are formed in the insulating stack structure. The material of the tunneling
25 layer 14 may be a dielectric material, such as silicon oxide. The material of the channel pillar 16 may be a semiconductor, such as undoped polysilicon. The conductive pillars 32a and 32b are, for example, doped polysilicon.

20

30

[0065] Then, lithography and etching processes are performed to form slit trenches 1110 in the insulating stack structure to divide the insulating stack structure into a plurality of sectors B.

25 **[0066]** Afterwards, a gate replacement process is performed to form a gate stack structure 52.

40

50

File: 112718jpf

First, an etching process is performed to inject an etchant into the slit trenches 1110 to remove the another insulating layer in the insulating stack structure to form a plurality of horizontal openings 34 and then form gate layers 38 in the horizontal openings 34. In some embodiments, before the gate layer 38 is formed, a charge storage layer 12 and a blocking layer 36 are also formed in the horizontal opening 34. The charge storage layer 12 is, for example, silicon nitride. The material of the blocking layer 36 is, for example, a high dielectric constant material having a dielectric constant greater than or equal to 7, such as aluminum oxide (Al_2O_3), hafnium oxide (HfO_2), lanthanum oxide (La_2O_3), transition metal oxide, lanthanide oxide, or combinations thereof. The gate layer 38 is, for example, tungsten. In some embodiments, before the gate layers 38 are formed, a barrier layer 37 is formed. The material of the barrier layer 37 is, for example, titanium (Ti), titanium nitride (TiN), tantalum (Ta), tantalum nitride (TaN), or a combination thereof.

[0067] Next, slits SLT are formed in the slit trenches 1110. The method of forming the slits SLT includes filling an insulating filling material on the gate stack structure 52 and in the slit trenches 1110, and then removing the excessive insulating filling material on the gate stack structure 52 through an etch-back process or a planarization process. The insulating filling material is, for example, silicon oxide.

[0068] Afterwards, an upper interconnect structure 1034 (including a local bit line LBL_n , a local source line LSL_n , a global bit line GBL_n and a global source line GSL_n) is formed on the gate stack structure 52. The upper interconnect structure 1034 may be formed by any known method, such as damascene, dual-damascene, and the like, which shall not be described in detail herein.

[0069] Referring to FIG. 3A and FIG. 3B, in this embodiment, after the upper interconnect structure 1034 (including the local bit line LBL_n , the local source line LSL_n , the global bit line GBL_n , and the global source line GSL_n) is formed, a heater 1200 is further formed above the upper interconnect structure 1034. The method of forming the heater 1200 includes, for example,

10

20

30

40

50

File: 112718jpf

forming a dielectric layer 1040 above the upper interconnect structure 1034 first. The material of the dielectric layer 1040 is, for example, silicon oxide. In some embodiments, a planarization process such as a chemical mechanical planarization process is further performed, so that the dielectric layer 1040 has a flat surface. Afterwards, lithography and etching processes are performed to form a plurality of grooves OP1 in the dielectric layer 1040. Then, a barrier material layer and a metal material layer are sequentially formed on the dielectric layer 1040 and in the grooves. Next, a planarization process such as a chemical mechanical planarization process is performed to remove the barrier material layer and the metal material layer on the surface of the dielectric layer 1040 and form a barrier layer 1204 and a metal layer 1202 in the groove. The metal material layer is, for example, copper or tungsten. The barrier material layer is, for example, titanium, tantalum, titanium nitride, tantalum nitride, or a combination thereof.

[0070] Referring to FIG. 3B, after the heater 1200 is formed, a bonding layer 1300 is formed. The method of forming the bonding layer 1300 is as follows. First, an insulating layer 1304 is first formed on the heater 1200 and the dielectric layer 1040, and then lithography and etching processes are performed to form a plurality of pad openings OP2 in the insulating layer 1304. The bottom of the pad opening OP2 exposes the heater 1200. Afterwards, a conductive layer is formed on the insulating layer 1304 and in the pad openings OP2. Then, a planarization process such as a chemical mechanical planarization process is performed to remove the conductive layer on the insulating layer 1304 and form pads 1302 in the pad openings OP2.

[0071] In the above embodiment, the heater 1200 of the memory chip 1000 is formed after the upper interconnect structure 1034 is formed. In other embodiments, the heater 1200 of the memory chip 1000 may be formed before the upper interconnect structure 1034 is formed.

[0072] Referring to FIG. 4C, the heater 1200 of the memory chip 1000 is formed in the slit trench 1110 between the gate stack structures 52 after the gate stack structure 52 of the 3D flash memory structure 1100 is formed and before the upper interconnect structure 1034 (including the

10

20

30

40

50

File: 112718jpf

local bit line LBL_n , the local source line LSL_n , the global bit line GBL_n , and the global source line GSL_n) is formed.

[0073] Referring to FIG. 4A and FIG. 4C, the method of forming the heater 1200 includes, for example, forming a liner material layer in the slit trench 1110 first. The liner material layer is, for example, silicon oxide or silicon nitride. Next, a barrier material layer and a metal material layer are sequentially formed on the gate stack structure 52 and in the slit trench 1110. Then, a planarization process such as a chemical mechanical planarization process is performed to remove the barrier material layer and the metal material layer on the surface of the gate stack structure 52 and form an insulating liner layer 1112, a barrier layer 1204, and a metal layer 1202 in the slit trench 1110. The metal material layer is, for example, copper or tungsten. The barrier material layer is, for example, titanium, tantalum, titanium nitride, tantalum nitride, or a combination thereof.

[0074] Referring to FIG. 4B and FIG. 4C, after the heater 1200 is formed, an upper interconnect structure 1034 (including a local bit line LBL_n , a local source line LSL_n , a global bit line GBL_n , and a global source line GSL_n) is formed. Afterwards, a bonding layer 1300 is formed on the upper interconnect structure 1034 according to the above-mentioned method.

[0075] Referring to FIG. 7A, a plurality of control chips 2000 are provided. The method of forming the control chips 2000 is as follows. Referring to FIG. 5C, second transistors 2020 are formed on a second substrate (wafer) 2010. Then, a second interconnect structure 2030 is formed on the second transistors 2020. The second interconnect structure 2030 may be formed by any known method, such as damascene, dual-damascene, and the like. Afterwards, a bonding layer 2050 is formed on the second interconnect structure 2030 according to the above-mentioned method. Next, dicing is performed to form a plurality of control chips 2000.

[0076] Referring to FIG. 7B, the bonding layer 2050 of the control chips 2000 and the bonding layer 1300 of the memory chips 1000 are bonded to form a bonding structure 3000. The bonding

File: 112718jpf

method is, for example, a hybrid bonding process. In some embodiments, after the control chips 2000 are bonded with the memory chips 1000 on the wafer 1010W, an encapsulation layer (not shown) is further formed around the sidewalls of the control chips 2000.

5 [0077] Referring to FIG. 7C, a dicing process is performed to form a plurality of mutually independent 3D flash memory module chips 5000. 10

[0078] In summary of the above, in the disclosure, the memory chip and the control chip are bonded to form the 3D flash memory module chip. With the driver of the control chip providing a high driving current to heat the heater in the memory chip, it is possible to heal the charge storage structure of the flash memory to achieve a higher erase speed and improve the endurance of the 10 flash memory. Furthermore, the control chip can locally heat the corresponding sector according to the status signal of the control logic unit of the memory chip. In addition, in the 3D flash memory module chip formed by bonding, the control chip may be manufactured separately, and it is not required to form a large-area heater controller in the memory chip. Therefore, it is possible to prevent the heater controller from occupying the area of the memory chip, and the 20 control chip may be manufactured by a less advanced process so as to reduce the cost of the process. 15

INDUSTRIAL APPLICABILITY

[0079] The 3D flash memory module chip and the method of fabricating the same of the present invention may be applied to 3D memory devices and methods of fabricating the same.

REFERENCE SIGNS LIST

- 20 [0080] 14: tunneling layer 30
- 16: channel pillar
- 20: memory cell
- 24: insulating filling layer
- 28: insulating pillar
- 25 32a: conductive pillar/source pillar

File: 112718jpf

	32b: conductive pillar/drain pillar	
	34: horizontal opening	
	36: blocking layer	
	37, 1204: barrier layer	
5	38, 2028: gate layer	10
	40: charge storage structure	
	52: gate stack structure	
	54, 1304, 2054: insulating layer	
	1000: memory chip	
10	1010: first substrate	
	1010W: wafer	
	1020: active device (first transistor)	
	1030: first interconnect structure	20
	1032: lower interconnect structure	
15	1034: upper interconnect structure	
	1040, 2031: dielectric layer	
	1100, 1100 ₁ , 1100 ₂ : 3D flash memory structure	
	1110: slit trench	
	1112: insulating liner layer	
20	1200, 1200 ₁ , 1200 ₂ , 1200 ₃ : heater	30
	1202: metal layer	
	1300, 2050: bonding layer	
	3000: bonding structure	
	1302, 1302a, 1302b, 2052, 2052a, 2052b: pad	
25	2000: control chip	

	File: 112718jpf	
	2000C, 2000C ₁ : column	
	2000R, 2000R ₁ : driving row	
	2010: second substrate	
	2012: active region	
5	2020, 2020 ₁ : second transistor	10
	2022a: source region	
	2022b: drain region	
	2024: gate dielectric layer	
	2030: second interconnect structure	
10	2032, 2034, C1, C2, C3: contact	
	2036, 2040, 2040a, 2040b, 2040c: conductive line	
	2038, 2042, 2042a, 2042b: via	
	2100: global power supply	20
	2200: row decoder	
15	2300: column decoder	
	5000: 3D flash memory module chip	
	A0, A1, A2: row address signal	
	A3, A4: column address signal	
	AR: array region	
20	B, B1, B2, B3, B4: sector	30
	BM1: lower first metal layer	
	BM2: lower second metal layer	
	BM3: lower third metal layer	
	BV1, BV2, TV1: via	
25	TM1: upper first metal layer	

File: 112718jpf

	TM2: upper second metal layer	
	E1: first end	
	E2: second end	
	OP1: groove	
5	OP2: pad opening	10
	SC: staircase structure	
	SL: scribe line	
	SLT: slit	
	SR: staircase region	
10	T, T', T1, T1', T2, T2', T3, T3', T4, T4': tile	
	W1, W2: width	
	I-I', II-II': line	
	X, Y, Z: direction	20

30

40

50

File: 112718jpf

WHAT IS CLAIMED IS:

- 1. A 3D flash memory module chip comprising:
 - a memory chip comprising:
 - a plurality of tiles each comprising a plurality of 3D flash memory structures; and
 - 5 a plurality of heaters disposed around the 3D flash memory structures of each of the tiles; and
 - a control chip bonded with the memory chip to drive at least one of the heaters.
 - 2. The 3D flash memory module chip according to claim 1, wherein the heaters are disposed above the 3D flash memory structures and are adjacent to the control chip.
 - 10 3. The 3D flash memory module chip according to claim 1, wherein the heaters are disposed in a plurality of slit trenches between the 3D flash memory structures.
 - 4. The 3D flash memory module chip according to claim 1, wherein the memory chip further comprises:
 - a plurality of first transistors located on a first substrate;
 - 15 the 3D flash memory structures located above the first transistors; and
 - a first interconnect structure, wherein the 3D flash memory structures are embedded in the first interconnect structure.
 - 5. The 3D flash memory module chip according to claim 4, wherein the first interconnect structure comprises:
 - 20 a lower interconnect structure located between the 3D flash memory structures and the first transistors and electrically connecting the 3D flash memory structures and the first transistors; and
 - 30 an upper interconnect structure located on the 3D flash memory structures and electrically connecting the 3D flash memory structures.
 - 25 6. The 3D flash memory module chip according to claim 4, wherein the control chip

File: 112718jpf

comprises:

a plurality of driving rows each comprising:

a second transistor located on a second substrate, wherein a source region of the second transistor is electrically connected to a global power supply;

5 a first pad electrically connected to a drain region of the second transistor and electrically connected to a first end of one of the heaters; and

10

a second pad that is grounded and is electrically connected to a second end of the one of the heaters.

7. The 3D flash memory module chip according to claim 6, wherein the control chip
10 further comprises:

a row decoder electrically coupled to a plurality of gate layers of the second transistors of the driving rows; and

20

a column decoder electrically coupled to a plurality of source regions of the second transistors and the global power supply.

8. The 3D flash memory module chip according to claim 6, wherein the control chip
15 comprises a plurality of tiles arranged in an array, wherein the source regions of the second transistors of the tiles in a same column are electrically connected to each other.

9. The 3D flash memory module chip according to claim 6, wherein the control chip and the memory chip are bonded by a bonding structure.

10. The 3D flash memory module chip according to claim 6, wherein the plurality of 3D
20 flash memory structures comprising a plurality of 3D AND flash memory structure, a plurality of 3D NAND flash memory structure, or a plurality of 3D NOR flash memory structure.

30

11. A method of fabricating a 3D flash memory module chip, comprising:

forming a memory chip, comprising:

25 forming a plurality of tiles on a first substrate, wherein each of the tiles comprises

40

50

File: 112718jpf

a plurality of 3D flash memory structures; and

forming a plurality of heaters around the 3D flash memory structures of each of the tiles;

forming a control chip; and

5 bonding the control chip and the memory chip, wherein the control chip is configured to drive the heaters.

10

12. The method of fabricating a 3D flash memory module chip according to claim 11, wherein the heaters are formed above the 3D flash memory structures.

13. The method of fabricating a 3D flash memory module chip according to claim 11, 10 wherein the heaters are formed in a plurality of slit trenches around the 3D flash memory structures.

14. The method of fabricating a 3D flash memory module chip according to claim 11, wherein the step of forming the memory chip further comprises:

forming a plurality of first transistors on the first substrate; and

forming the 3D flash memory structures above the first transistors.

20

15 15. The method of fabricating a 3D flash memory module chip according to claim 14, wherein the step of forming the control chip comprises:

forming a plurality of driving rows, wherein formation of each of the driving rows comprises:

forming a second transistor on a second substrate;

20 forming a second interconnect structure on the second transistor, wherein a source region of the second transistor is electrically coupled to a global power supply via the second interconnect structure;

30

forming a first pad on the second interconnect structure, wherein the first pad is electrically connected to a drain region of the second transistor via the second interconnect structure; and

25

40

50

File: 112718jpf

forming a second pad on the second interconnect structure, wherein the second pad is electrically connected to ground via the second interconnect structure.

16. The method of fabricating a 3D flash memory module chip according to claim 15, further comprising:

5 electrically connecting the first pad to a first end of one of the heaters; and 10
electrically connecting the second pad that to a second end of the one of the heaters.

17. The method of fabricating a 3D flash memory module chip according to claim 11, wherein the control chip and the memory chip are hybrid bonded by a bonding structure.

18. The method of fabricating a 3D flash memory module chip according to claim 11, 10
wherein the plurality of 3D flash memory structures comprising a plurality of 3D AND flash memory structure, a plurality of 3D NAND flash memory structure, or a plurality of 3D NOR flash memory structure. 20

30

40

50

File: 112718jpf

ABSTRACT OF THE DISCLOSURE

A 3D flash memory module chip includes a memory chip and a control chip. The memory chip includes a plurality of tiles and a plurality of heaters. The tiles each include a plurality of 3D flash memory structures. The heaters are disposed around the 3D flash memory structures of each of the tiles. The control chip is bonded with the memory chip to drive at least one of the heaters.

10

REPRESENTATIVE DRAWING: FIG. 1A

20

30

40

50

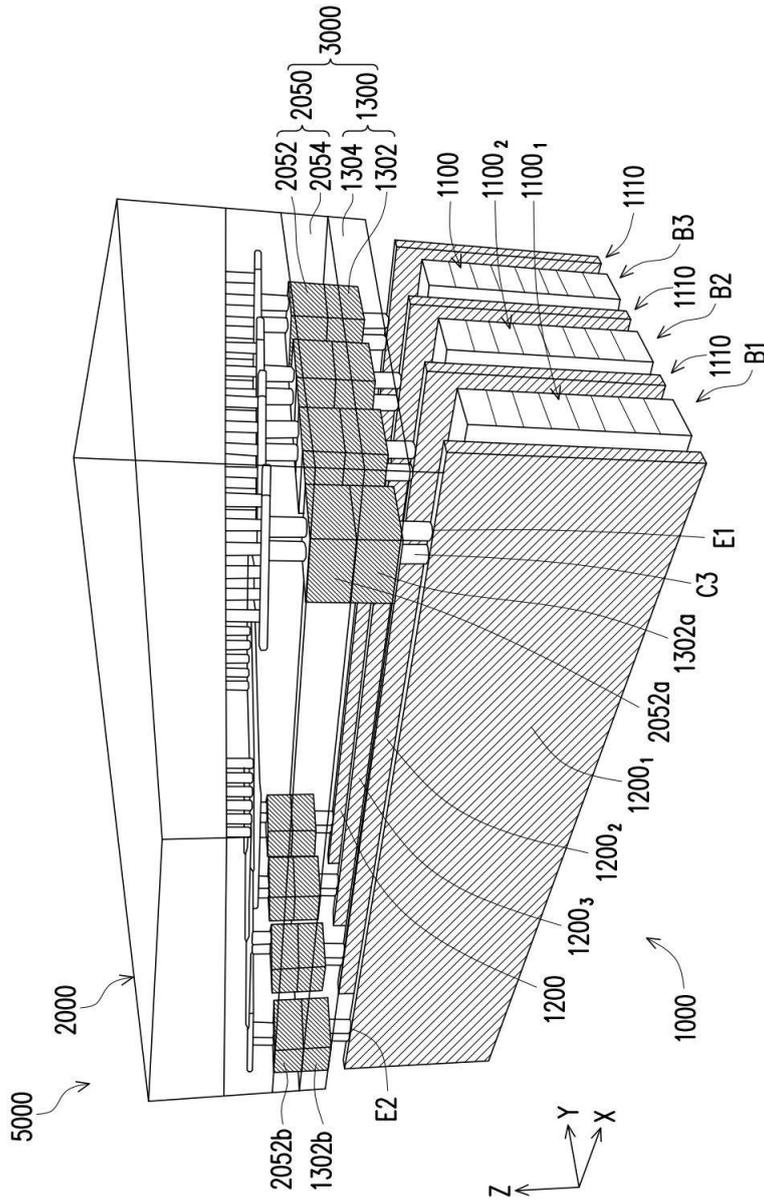


FIG. 1B

10

20

30

40

50

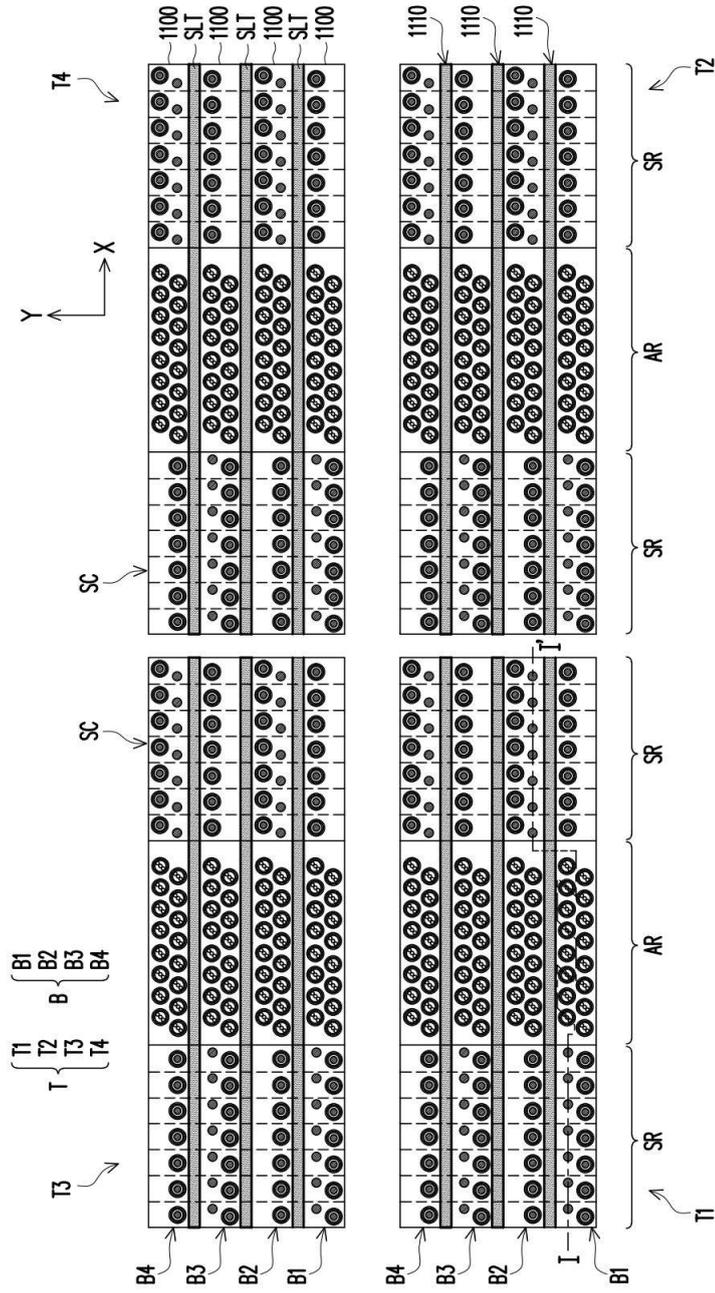


FIG. 2A

10

20

30

40

50

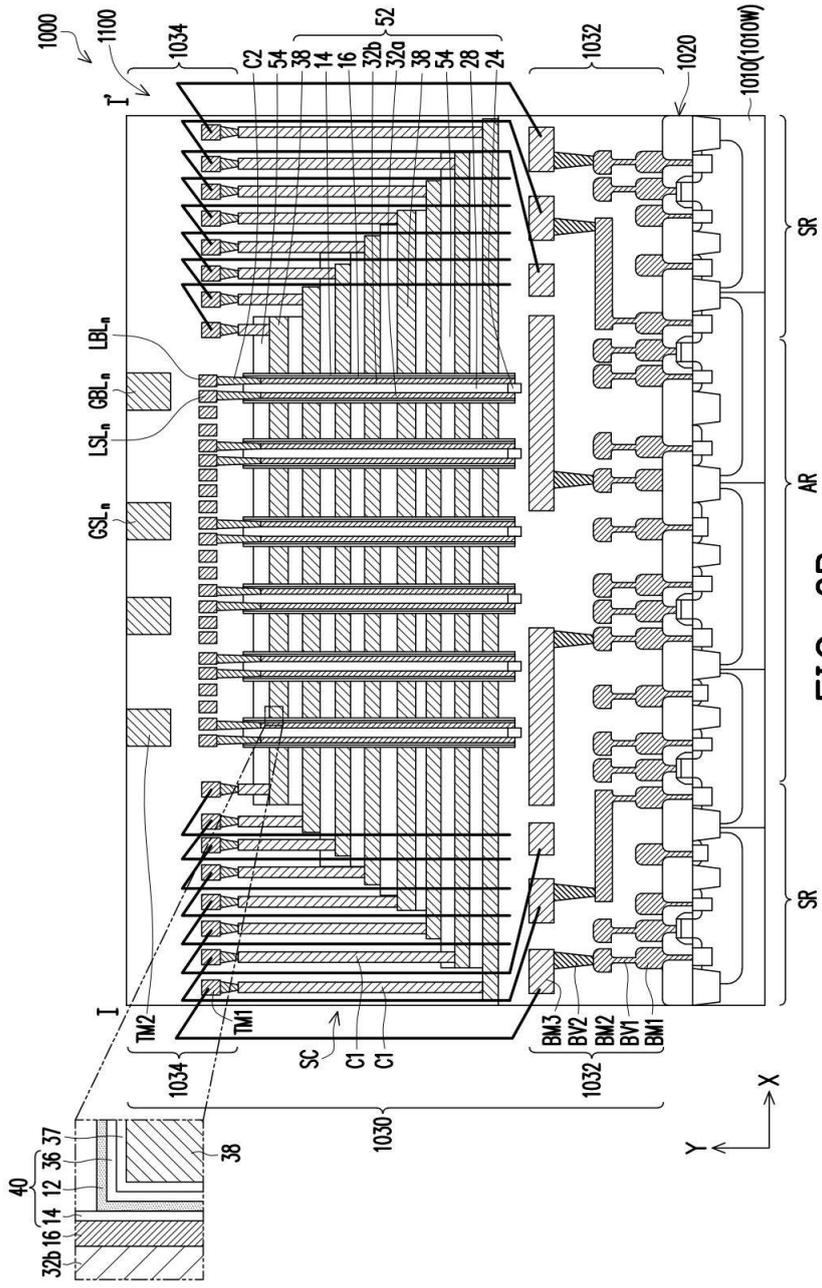


FIG. 2B

10

20

30

40

50

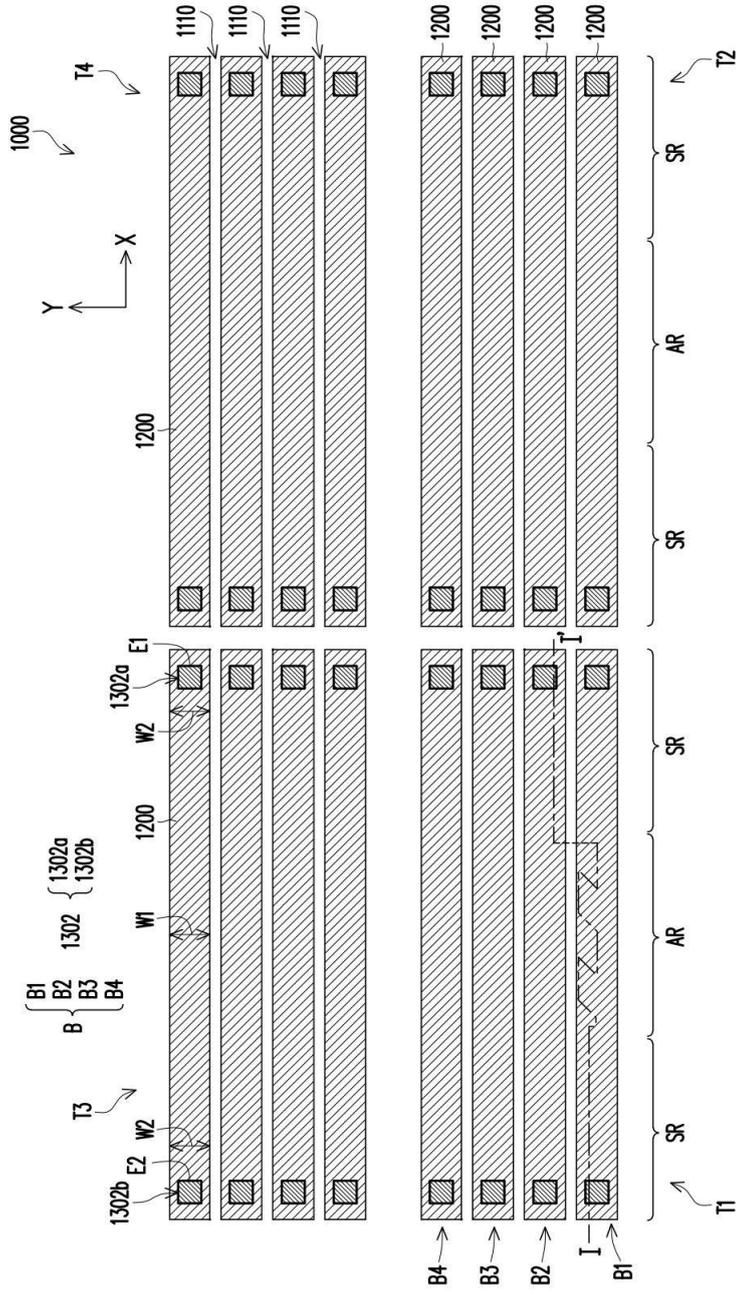


FIG. 3A

10

20

30

40

50

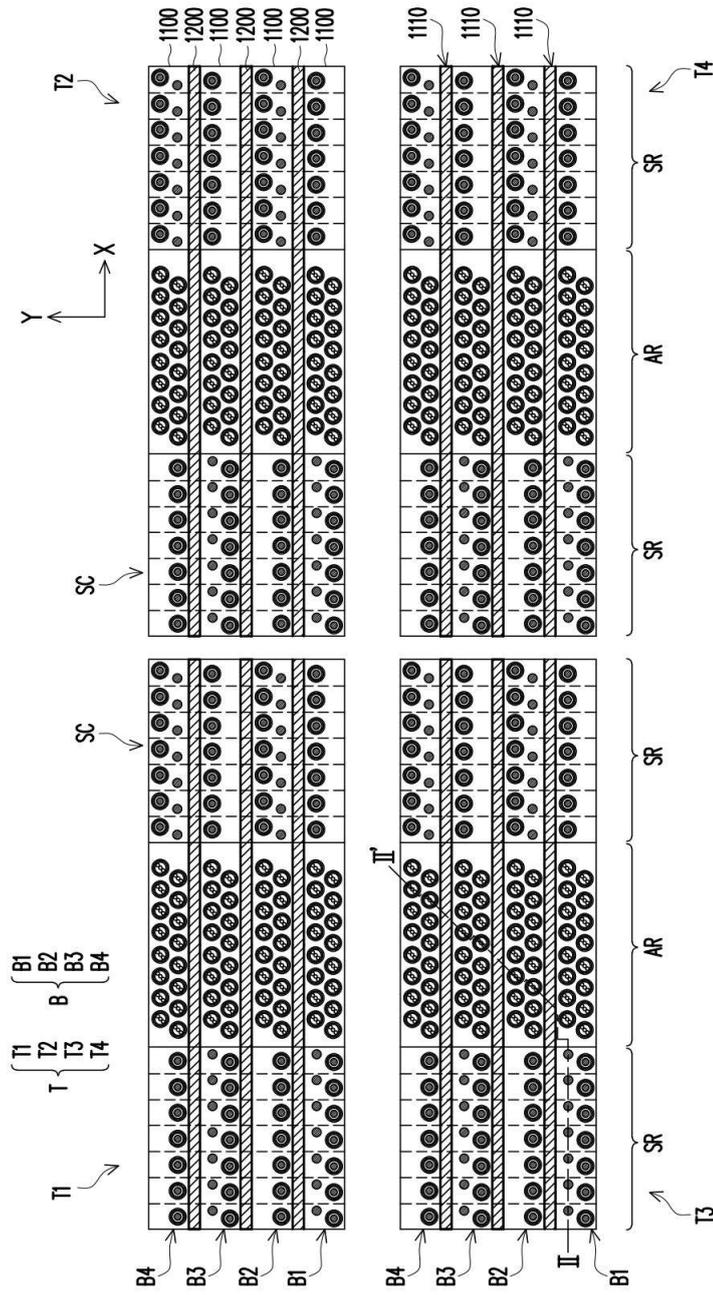


FIG. 4A

10

20

30

40

50

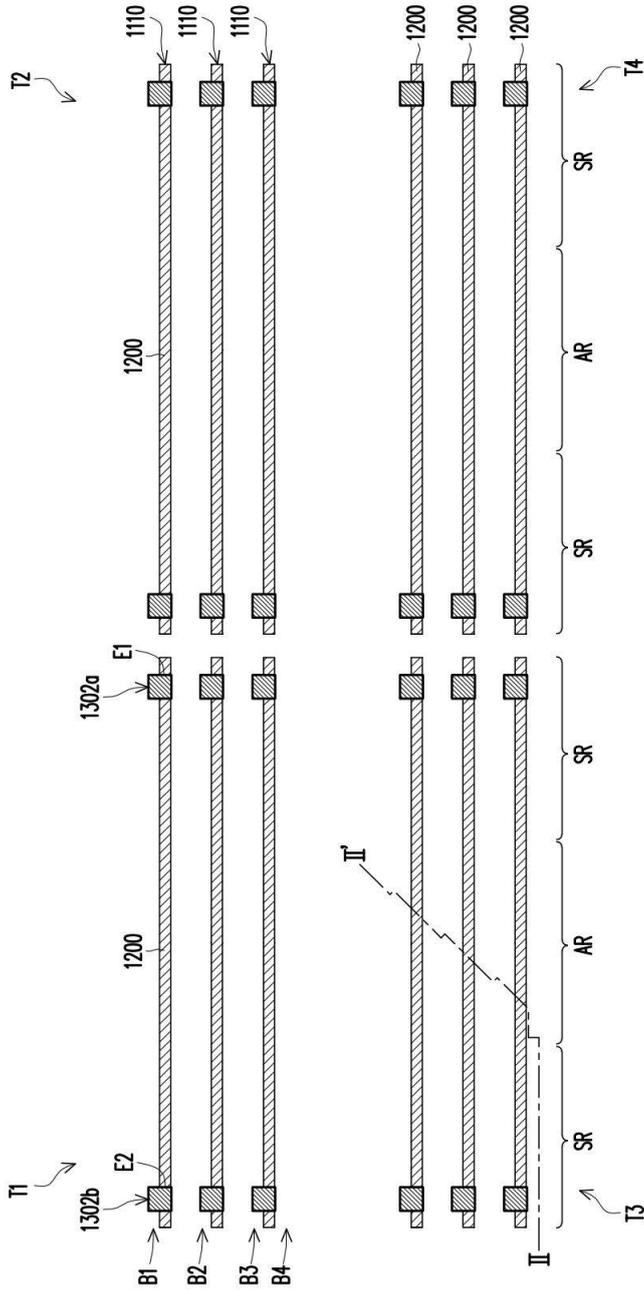


FIG. 4B

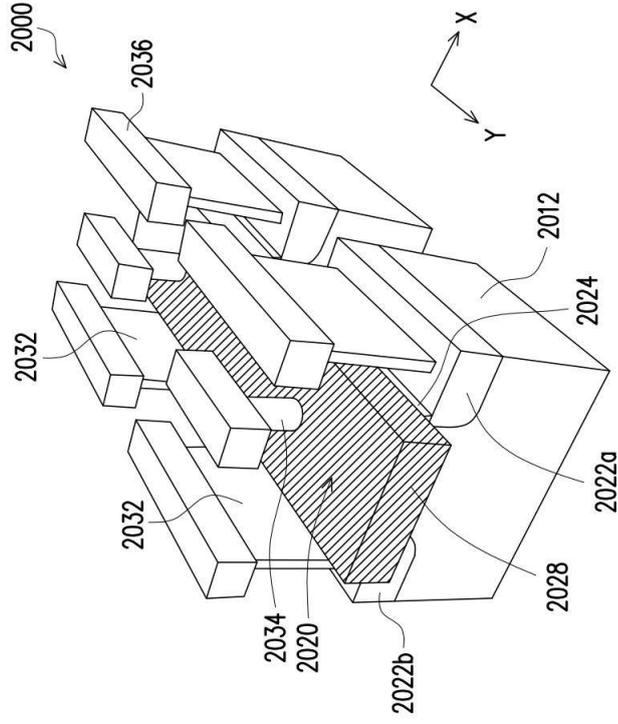
10

20

30

40

50



2000

FIG. 5B

10

20

30

40

50

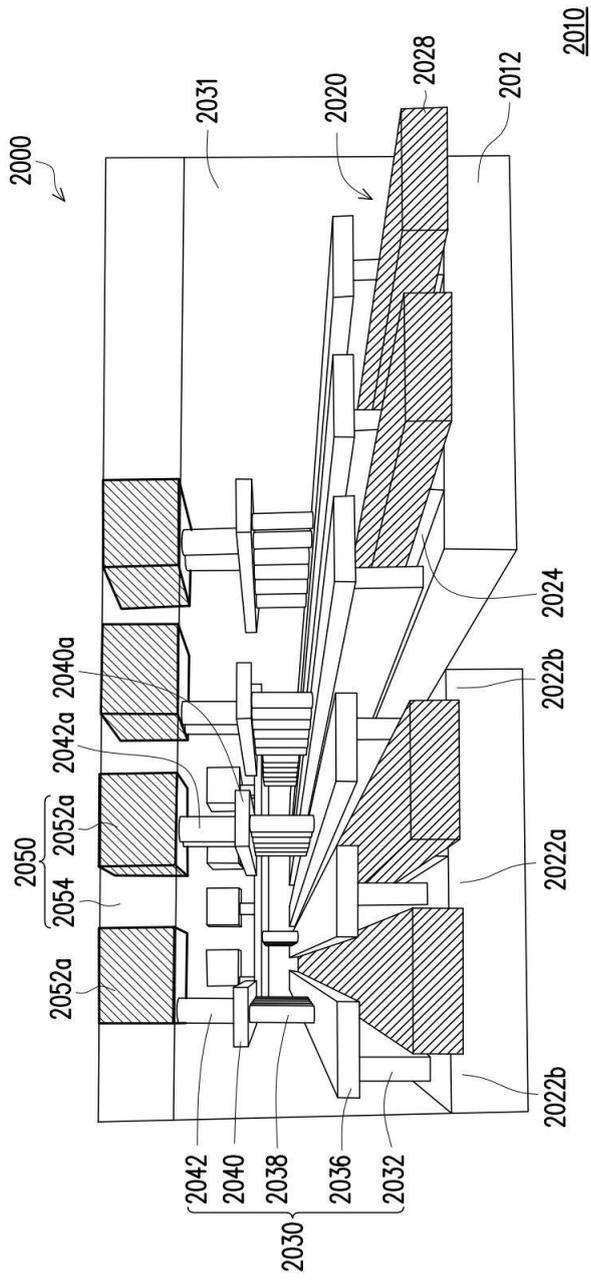


FIG. 5C

10

20

30

40

50

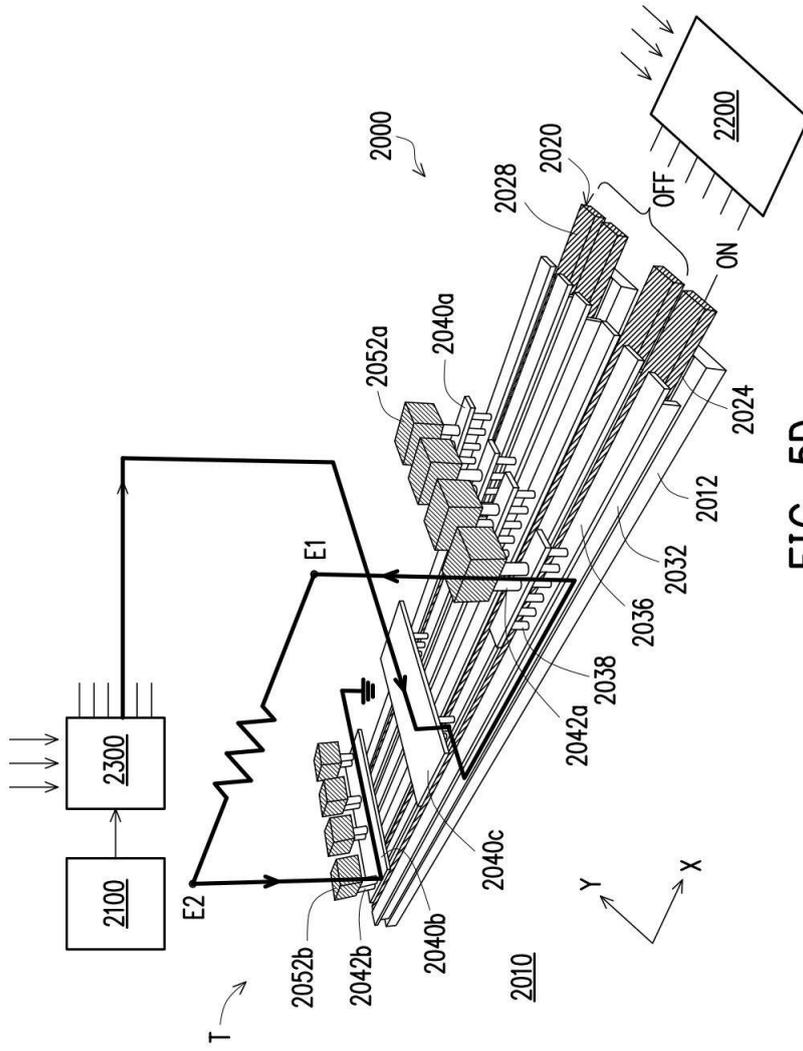


FIG. 5D

10

20

30

40

50

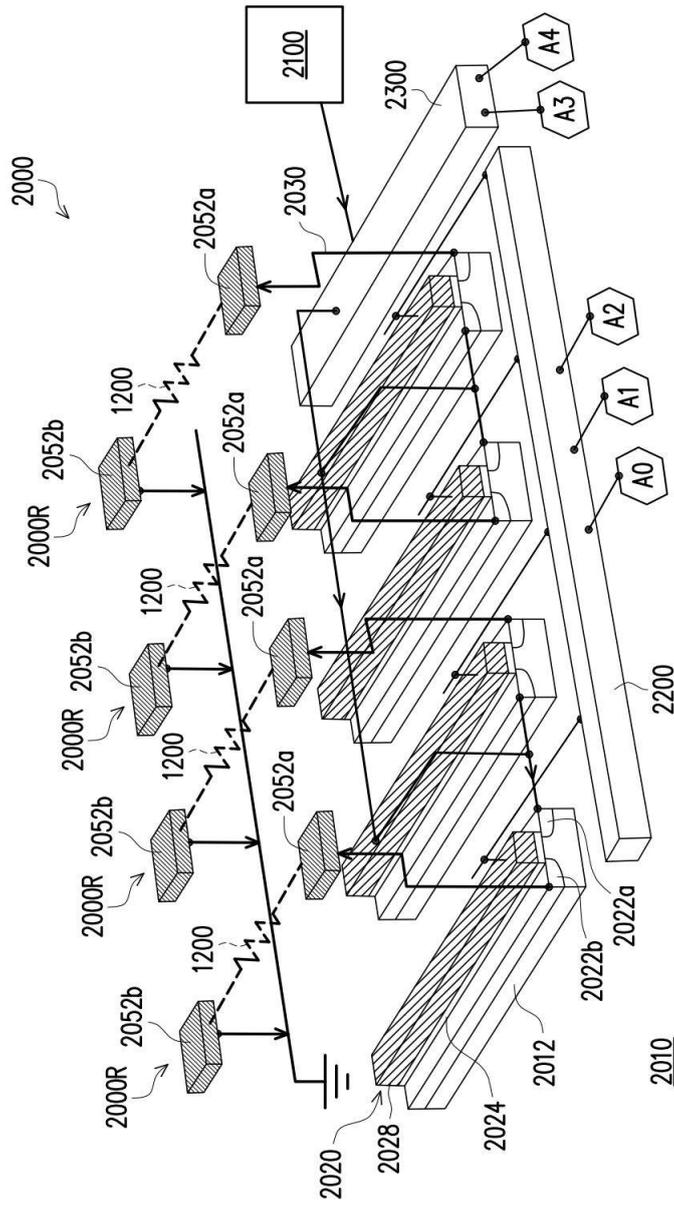


FIG. 5E

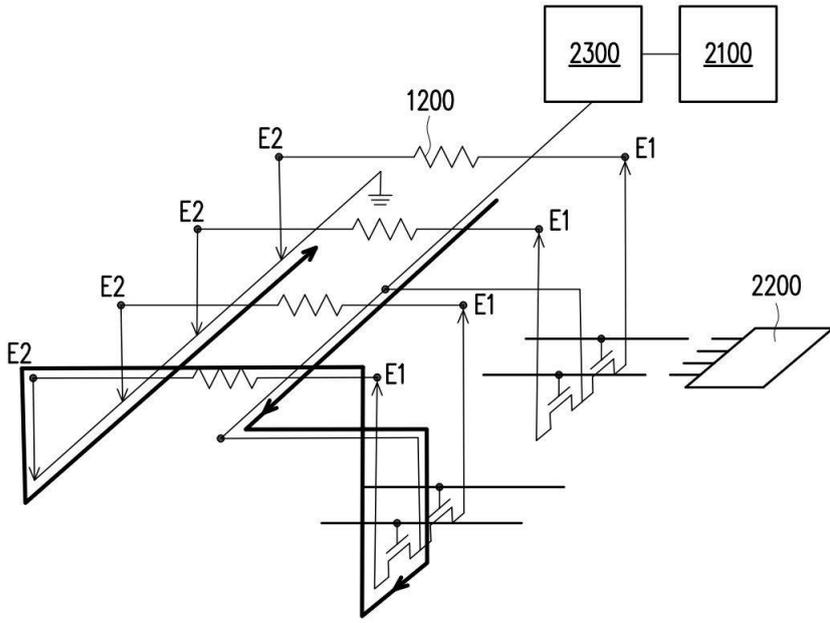
10

20

30

40

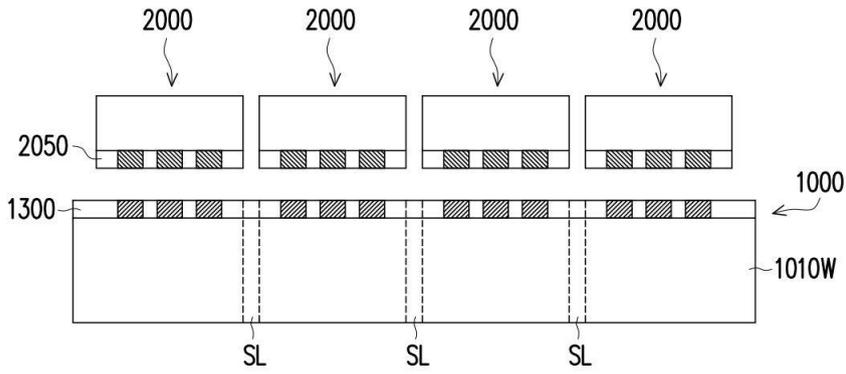
50



10

20

FIG. 6B

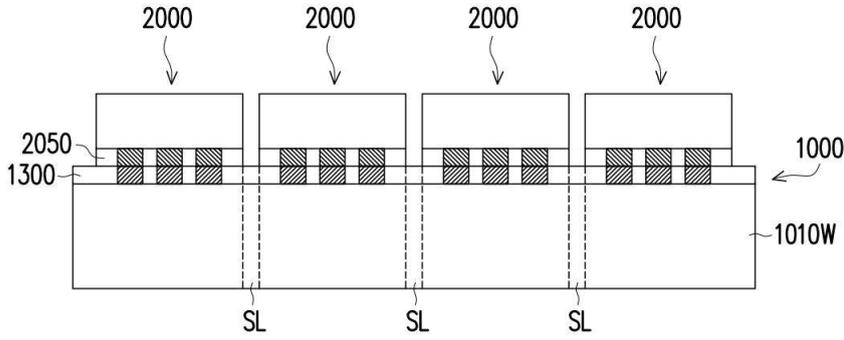


30

FIG. 7A

40

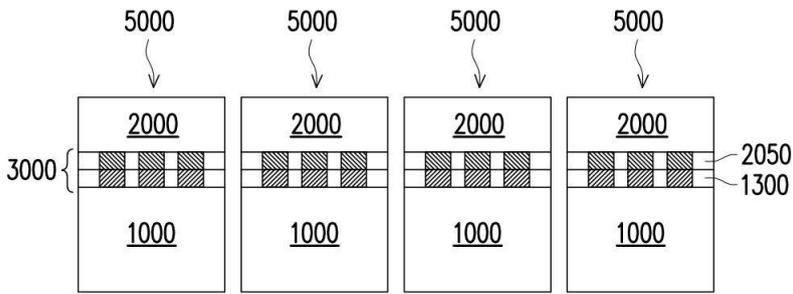
50



10

FIG. 7B

20



30

FIG. 7C

40

50

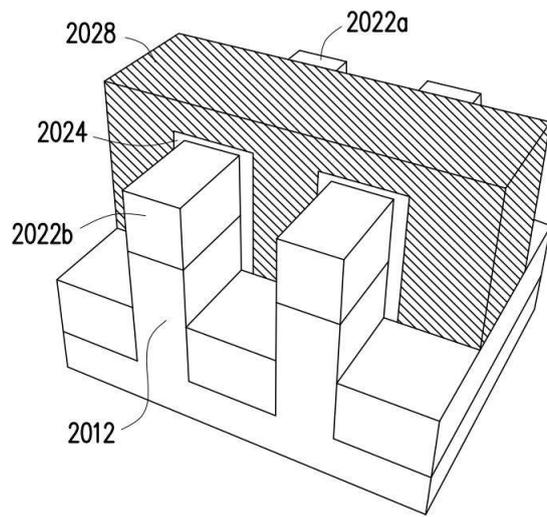


FIG. 8

10

20

30

40

50

フロントページの続き

Fターム(参考) GA10 JA04 JA19 JA37 JA39 JA40 MA06 MA16 MA19 ZA01
5F101 BA46 BA47 BB02 BD16 BD30 BD34 BE07