



(12) 发明专利

(10) 授权公告号 CN 110008440 B

(45) 授权公告日 2021.07.27

(21) 申请号 201910299610.9

(51) Int.Cl.

(22) 申请日 2019.04.15

G06F 17/16 (2006.01)

G06N 3/063 (2006.01)

(65) 同一申请的已公布的文献号

申请公布号 CN 110008440 A

审查员 李兵兵

(43) 申请公布日 2019.07.12

(73) 专利权人 恒烁半导体(合肥)股份有限公司

地址 230000 安徽省合肥市庐阳区天水路  
与太和路交叉口西北庐阳中科大校友企  
业创新园11号楼

(72) 发明人 任军 徐伟民 蒋明峰 李政达

吕向东 徐培

(74) 专利代理机构 北京中政联科专利代理事务

所(普通合伙) 11489

代理人 刘艳

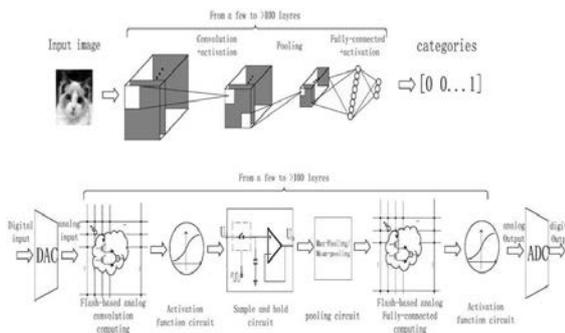
权利要求书3页 说明书13页 附图17页

(54) 发明名称

一种基于模拟矩阵运算单元的卷积运算及其应用

(57) 摘要

本发明涉及电路设计技术领域,公开了一种基于模拟矩阵运算单元的卷积运算及其应用,所述卷积运算将卷积核经过转换拼接成P个长度为Q的横向量,并映射到模拟矩阵运算单元中,输入特征图按照滑动窗口进行切片转换成N个长度为Q的横向量,在脉冲波下,依次将N个长度为Q的横向量映射到向量中,同时输出N个长度为P的运算结果经采样保持器得到完整的卷积结果输出。本发明有效提高了卷积计算的速率,减少了功耗和电路面积,具有高度计算并行性,大大提高了计算的密度和效率,具有较高的实用价值和广泛的应用前景。



1. 一种基于模拟矩阵运算单元的卷积运算方法,其特征在于,所述模拟矩阵运算单元

能够实现向量 $\{a_1 \cdots a_Q\}$ 与矩阵 $\begin{pmatrix} b_{11} & \cdots & b_{1Q} \\ \vdots & \ddots & \vdots \\ b_{P1} & \cdots & b_{PQ} \end{pmatrix}^T$ 的乘法输出;

所述卷积运算包括以下步骤:

将卷积核经过转换拼接成P个长度为Q的横向量,并映射到矩阵 $\begin{pmatrix} b_{11} & \cdots & b_{1Q} \\ \vdots & \ddots & \vdots \\ b_{P1} & \cdots & b_{PQ} \end{pmatrix}$ 中;

将输入特征图按照滑动窗口进行切分转换成N个长度为Q的横向量;

在N个脉冲波下,依次将N个长度为Q的横向量映射到向量 $\{a_1 \cdots a_Q\}$ 中,矩阵运算单元的模拟输出端口按照脉冲波序列,输出N个长度为P的运算结果;

通过采样保持器将所有输出的运算结果采样保持至一个时间点,得到完整的卷积结果;

所述模拟矩阵运算单元由P个模拟乘加单元构成,所述模拟乘加单元能够实现一个行向量 $[V_1, V_2, \dots, V_n]$ 与列向量 $[W_1, W_2, \dots, W_n]^T$ 的乘法输出;

所述模拟乘加单元由Q个模拟乘法电路横构成;

所述模拟乘法电路包括一对浮栅场效应管M1、M2和一个差分电流检测电路;所述M<sub>1</sub>和M<sub>2</sub>共栅极并接入电压源,共漏极或共源极并接入模拟电压输入,所述差分电流检测电路包括两个电流输入端和一个输出端,所述两个电流输入端对应接入M<sub>1</sub>和M<sub>2</sub>各自源极或漏极,所述差分电流检测电路能保持两个电流输入端电压不变且输出端结果为两个电流输入端输入电流差值的函数;所述M<sub>1</sub>和M<sub>2</sub>的栅源电压大于两者阈值电压的最大值,漏源电压小于栅源电压分别与两者阈值电压差值的最小值;

所述模拟乘加单元中的Q个模拟乘法电路共栅极且共用一个差分电流检测电路,所述每个模拟乘法电路中该对浮栅场效应管的共漏极或共源极分别接入对应的输入电压信号;所述差分电流检测电路的输出端结果为每个模拟乘法电路输出的基于该对浮栅场效应管输出电流差值的函数之和;

所述模拟乘法电路执行乘法运算的步骤包括:

对两个浮栅场效应管M<sub>1</sub>、M<sub>2</sub>中的浮置栅极执行擦除和写操作,将乘数以所述一对浮栅场效应管的阈值电压差值的形式进行存储;

对两个浮栅场效应管源极、漏极和选栅极执行电压施加,以所述模拟电压输入与所述两个电流输入端电压差值的形式输入被乘数;

通过差分电流检测电路输出基于所述一对浮栅场效应管输出电流差值的函数作为所述被乘数与乘数的乘积结果。

2. 根据权利要求1所述的一种基于模拟矩阵运算单元的卷积运算方法,其特征在于,所述模拟矩阵运算单元为一个P行Q列的模块电路,所述模块电路中的每列浮栅场效应管的共漏极或共源极接入同一的输入电压信号;

通过控制P个模拟乘加单元的栅极电压,确定参加计算的乘加单元的个数,实现向量

$[a_1 \cdots a_p]$ 与矩阵  $\begin{pmatrix} b_{11} & \cdots & b_{1Q} \\ \vdots & \ddots & \vdots \\ b_{J1} & \cdots & b_{JQ} \end{pmatrix}^T$  的乘法输出,其中 $J \leq P$ 。

3. 根据权利要求2所述的一种基于模拟矩阵运算单元的卷积运算方法,其特征在于,所述模拟矩阵运算单元在确定参加计算的乘加单元的个数后执行以下步骤:

对每对浮栅场效应管中的浮置栅极执行擦除和写操作,按照一对浮栅场效应管的阈值电压差值对应一个乘数进行存储,完成大小为 $J \times Q$ 的乘数矩阵的存储;

对每对浮栅场效应管源极、漏极和控制栅极执行电压施加,以每个乘法电路中模拟电压输入与所述两个电流输入端电压差值对应一个被乘数,完成元素个数为 $Q$ 的横向被乘数向量的输入;

通过每行中设置差分电流检测电路输出得到被乘数向量与乘数矩阵的乘积结果。

4. 根据权利要求1所述的一种基于模拟矩阵运算单元的卷积运算方法,其特征在于,所述差分电流检测电路包括两个电流-电压转换器和一个电压差分放大器;所述电流-电压转换器由一个运算放大器和一个阻性负载构成,运算放大器的负极输入端为所述电流-电压转换器的电流输入端口,运算放大器的输出端为所述电流-电压转换器的电压输出端口,运算放大器的正极输入端外加直流电平 $V_b$ ,所述阻性负载连接在运算放大器的负极输入端与运算放大器的输出端之间。

5. 根据权利要求1所述的一种基于模拟矩阵运算单元的卷积运算方法,其特征在于,所述电压差分放大器由一个运算放大器、两个阻值相同的电阻 $R_1$ 、 $R_3$ 和两个阻值相同的电阻 $R_2$ 、 $R_4$ 构成,所述电压差分放大器的两个电压输入端口分别通过电阻 $R_1$ 、 $R_3$ 连接到运算放大器的正、负极输入端,运算放大器的正极输入端通过电阻 $R_2$ 连接到直流电平 $V_b$ ,运算放大器的负极输入端通过电阻 $R_4$ 连接到运算放大器的输出端,运算放大器的输出端为所述电压差分放大器的电压输出端口。

6. 根据权利要求4所述的一种基于模拟矩阵运算单元的卷积运算方法,其特征在于,所述阻性负载为浮栅场效应管或MOS管,所述浮栅场效应管或MOS管的栅极外加栅极电压,漏极或源极连接到所述运算放大器的输出端,源极或漏极连接到所述运算放大器的负极输入端;

所述浮栅场效应管或MOS管的漏源电压 $V_{DS}$ 、栅源电压 $V_{GS}$ 和阈值电压 $V_{TH}$ 满足:  $V_{DS} \ll 2(V_{GS} - V_{TH})$ 。

7. 一种采用如权利要求1-6任意一项所述的基于模拟矩阵运算单元的卷积运算方法实现卷积神经网络的方法,其特征在于,包括:

DAC电路、卷积层电路、激活函数电路、采样保持电路、池化层电路、全连接层电路和ADC电路,各个电路间的级联构成卷积神经网络;

所述的DAC电路用于将输入层的数字输入转化为模拟输入;

所述卷积层电路与全连接层电路均采用所述基于模拟矩阵运算单元的卷积运算实现功能并构建电路;

所述的激活函数电路将卷积层的输出电流转换为电压,同时完成激活函数作用;

所述的采样保持电路用于将卷积层的序列计算结果采样并保持在同一时间点;

所述的池化层电路用于池化层的实现,包括平均池化电路和最大池化电路;

所述的ADC电路用于将输出层的模拟输出转换为数字输出；

具体步骤包括：

输入特征图经过DAC电路转换为模拟信号输入；

经过卷积层电路得到第一次卷积计算结果，然后经过激活函数电路，将计算结果从电流转变为电压，并且进行非线性化；

经过采样保持电路，得到一个完整的中间特征图，再经过池化电路，将模拟信号降维；

继续经过全连接层电路进行全连接运算，最后经过激活函数电路得到模拟输出特征图，并经ADC电路转换为数字输出。

8. 根据权利要求7所述的一种实现卷积神经网络的方法，其特征在于，所述全连接层电路为一个P行Q列的模拟矩阵运算单元，具体实现方法为：

将输入层 $[a_1 \cdots a_Q]$ 按照每个乘法电路中模拟电压输入与两个电流输入端电压差值对应一个 $a_i$ 其中 $i \in (1 \cdots Q)$ ，从模拟矩阵运算单元模拟电压输入端输入；

将全连接层的 $P \times Q$ 个权重拆分为P个横向量，即 $[w_{11}, w_{21}, \cdots, w_{Q1}], [w_{12}, w_{22}, \cdots, w_{Q2}], \cdots, [w_{1P}, w_{2P}, \cdots, w_{QP}]$ ，按照一对浮栅场效应管的阈值电压差值对应一个乘数进行存储，将其映射到模拟矩阵运算单元的P行中，完成大小为 $P \times Q$ 的权重矩阵的存储；

经过一个计算周期，通过采样保持器将所有输出的运算结果采样保持至一个时间点，得到完整的卷积结果，输出层向量 $[b_1, b_2, \cdots, b_P]$ 。

## 一种基于模拟矩阵运算单元的卷积运算及其应用

### 技术领域

[0001] 本发明涉及电路设计技术领域,具体涉及一种基于模拟矩阵运算单元的卷积运算及其应用。

### 背景技术

[0002] 卷积神经网络在图像识别,目标检测和许多机器学习应用领域显示出巨大的优势。卷积神经网络主要由卷积层,池化层,全连接层级联组成,主要有输入层像素块和卷积核之间的卷积操作、为引入非线性而进行的激活操作、为减少特征值而对特征图进行的下采样操作(即池化)以及卷积之后的全连接操作,其中,绝大部分计算量都在卷积层和全连接层。

[0003] 大型的卷积神经网络,具有庞大的参数集和计算量。为完成庞大的计算量,一般芯片的设计思路是基于冯诺依曼架构上大量增加并行的运算单元,从早期的GPU,再到现在的FPGA,ASIC,NPU,TPU,都是由控制单元,存储单元,计算单元构成的。运算过程中,首先要将权重和输入特征存至片外的存储器,然后将需要运算的数据通过片内的一二级缓存,再到寄存器,最后送入ALU单元进行运算。这种架构存在两大问题:运算过程中,片内片外的数据往返传输消耗了大量的运算时间和功耗,计算单元和存储器之间的数据搬运消耗了大量的资源;为满足算力需要不断增加的并行计算单元的数量与存储单元带宽之间存在矛盾,成为AI芯片算力提升的瓶颈。

### 发明内容

[0004] 针对现有技术的不足,本发明提供一种模拟乘法电路、模拟乘法方法及其应用,用以解决背景技术中提出的问题。

[0005] 本发明解决技术问题采用如下技术方案:

[0006] 一种基于模拟矩阵运算单元的卷积运算,所述模拟矩阵运算单元能够实现向量

$(a_1 \cdots a_q)$  与矩阵  $\begin{pmatrix} b_{11} & \cdots & b_{1Q} \\ \vdots & \ddots & \vdots \\ b_{P1} & \cdots & b_{PQ} \end{pmatrix}^T$  的乘法输出;

[0007] 所述卷积运算包括以下步骤:

[0008] 将卷积核经过转换拼接成P个长度为Q的横向向量,并映射到矩阵  $\begin{pmatrix} b_{11} & \cdots & b_{1Q} \\ \vdots & \ddots & \vdots \\ b_{P1} & \cdots & b_{PQ} \end{pmatrix}$  中;

[0009] 将输入特征图按照滑动窗口进行切分转换成N个长度为Q的横向向量;

[0010] 在N个脉冲波下,依次将N个长度为Q的横向向量映射到向量  $(a_1 \cdots a_q)$  中,矩阵运算单元的模拟输出端口按照脉冲波序列,输出N个长度为P的运算结果;

[0011] 通过采样保持器将所有输出的运算结果采样保持至一个时间点,得到完整的卷积

结果。

[0012] 优选地,所述模拟矩阵运算单元由P个模拟乘加单元构成,所述模拟乘加单元能够实现一个行向量 $[V_1, V_2, \dots, V_n]$ 与列向量 $[W_1, W_2, \dots, W_n]^T$ 的乘法输出。

[0013] 优选地,所述模拟乘加单元由Q个模拟乘法电路横构成;

[0014] 所述模拟乘法电路包括一对浮栅场效应管M<sub>1</sub>、M<sub>2</sub>和一个差分电流检测电路;所述M<sub>1</sub>和M<sub>2</sub>共栅极并接入电压源,共漏极或共源极并接入模拟电压输入,所述差分电流检测电路包括两个电流输入端和一个输出端,所述两个电流输入端对应接入M<sub>1</sub>和M<sub>2</sub>各自源极或漏极,所述差分电流检测电路能保持两个电流输入端电压不变且输出端结果为两个电流输入端输入电流差值的函数;所述M<sub>1</sub>和M<sub>2</sub>的栅源电压大于两者阈值电压的最大值,漏源电压小于栅源电压分别与两者阈值电压差值的最小值;

[0015] 所述模拟乘加单元中的Q个模拟乘法电路共栅极且共同一个差分电流检测电路,所述每个模拟乘法电路中该对浮栅场效应管的共漏极或共源极分别接入对应的输入电压信号;所述差分电流检测电路的输出端结果为每个模拟乘法电路输出的基于该对浮栅场效应管输出电流差值的函数之和。

[0016] 优选地,所述模拟乘法电路执行乘法运算的步骤包括:

[0017] 对两个浮栅场效应管M<sub>1</sub>、M<sub>2</sub>中的浮置栅极执行擦除和写操作,将乘数以所述一对浮栅场效应管的阈值电压差值的形式进行存储;

[0018] 对两个浮栅场效应管源极、漏极和选栅极执行电压施加,以所述模拟电压输入与上述两个电流输入端电压差值的形式输入被乘数;

[0019] 通过差分电流检测电路输出基于所述一对浮栅场效应管输出电流差值的函数作为所述被乘数与乘数的乘积结果。

[0020] 优选地,所述模拟矩阵运算单元为一个P行Q列的模块电路,所述模块电路中的每列浮栅场效应管的共漏极或共源极接入同一的输入电压信号;

[0021] 通过控制P个模拟乘加单元的栅极电压,确定参加计算的乘加单元的个数,实现向

量 $[a_1 \dots a_Q]$ 与矩阵 $\begin{pmatrix} b_{11} & \dots & b_{1Q} \\ \vdots & \ddots & \vdots \\ b_{J1} & \dots & b_{JQ} \end{pmatrix}^T$ 的乘法输出,其中 $J \leq P$ 。

[0022] 优选地,所述模拟矩阵运算单元在确定参加计算的乘加单元的个数后执行以下步骤:

[0023] 对每对浮栅场效应管中的浮置栅极执行擦除和写操作,按照一对浮栅场效应管的阈值电压差值对应一个乘数进行存储,完成大小为 $J \times Q$ 的乘数矩阵的存储;

[0024] 对每对浮栅场效应管源极、漏极和控制栅极执行电压施加,以每个乘法电路中模拟电压输入与上述两个电流输入端电压差值对应一个被乘数,完成元素个数为Q的横向被乘数向量的输入;

[0025] 通过每行中设置差分电流检测电路输出得到被乘数向量与乘数矩阵的乘积结果。

[0026] 优选地,所述差分电流检测电路包括两个电流-电压转换器和一个电压差分放大器;所述电流-电压转换器由一个运算放大器和一个阻性负载构成,运算放大器的负极输入端为所述电流-电压转换器的电流输入端口,运算放大器的输出端为所述电流-电压转

换器的电压输出端口,运算放大器的正极输入端外加直流电平 $V_b$ ,所述阻性负载连接在运算放大器的负极输入端与运算放大器的输出端之间。

[0027] 优选地,所述电压差分放大器由一个运算放大器、两个阻值相同的电阻 $R_1$ 、 $R_3$ 和两个阻值相同的电阻 $R_2$ 、 $R_4$ 构成,所述电压差分放大器的两个电压输入端口分别通过电阻 $R_1$ 、 $R_3$ 连接到运算放大器的正、负极输入端,运算放大器的正极输入端通过电阻 $R_2$ 连接到直流电平 $V_b$ ,运算放大器的负极输入端通过电阻 $R_4$ 连接到运算放大器的输出端,运算放大器的输出端为所述电压差分放大器的电压输出端口。

[0028] 优选地,所述阻性负载为浮栅场效应管或MOS管,所述浮栅场效应管或MOS管的栅极外加栅极电压,漏极或源极连接到所述运算放大器的输出端,源极或漏极连接到所述运算放大器的负极输入端;

[0029] 所述浮栅场效应管或MOS管的漏源电压 $V_{DS}$ 、栅源电压 $V_{GS}$ 和阈值电压 $V_{TH}$ 满足: $V_{DS} < 2(V_{GS} - V_{TH})$ 。

[0030] 本发明还提供一种采用所述的基于模拟矩阵运算单元的卷积运算实现卷积神经网络的方法,包括:

[0031] DAC电路、卷积层电路、激活函数电路、采样保持电路、池化层电路、全连接层电路和ADC电路,各个电路间的级联构成卷积神经网络;

[0032] 所述的DAC电路用于将输入层的数字输入转化为模拟输入;

[0033] 所述卷积层电路与全连接层电路均采用所述基于模拟矩阵运算单元的卷积运算实现功能并构建电路;

[0034] 所述的激活函数电路将卷积层的输出电流转换为电压,同时完成激活函数作用,包括ReLU激活函数实现和sigmoid函数实现;

[0035] 所述的采样保持电路用于将卷积层的序列计算结果采样并保持在同一时间点;

[0036] 所述的池化层电路用于池化层的实现,包括平均池化电路和最大池化电路;

[0037] 所述的ADC电路用于将输出层的模拟输出转换为数字输出;

[0038] 具体步骤包括:

[0039] 输入特征图经过DAC电路转换为模拟信号输入;

[0040] 经过卷积层电路得到第一次卷积计算结果,然后经过激活函数电路,将计算结果从电流转变为电压,并且进行非线性化;

[0041] 经过采样保持电路,得到一个完整的中间特征图,再经过池化电路,将模拟信号降维;

[0042] 继续经过全连接层电路进行全连接运算,最后经过激活函数电路得到模拟输出特征图,并经ADC电路转换为数字输出。

[0043] 优选地,所述全连接层电路为一个P行Q列的的模拟矩阵运算单元,具体实现方法为:

[0044] 将输入层 $[a_1 \cdots a_Q]$ 按照每个乘法电路中模拟电压输入与所述两个电流输入端电压差值对应一个 $a_i$ 其中 $i \in (1 \cdots Q)$ ,从模拟矩阵运算单元模拟电压输入端输入;

[0045] 将全连接层的 $P \times Q$ 个权重拆分为P个横向量,即 $[w_{11}, w_{21}, \cdots, w_{Q1}]$ ,  $[w_{12}, w_{22}, \cdots, w_{Q2}]$ ,  $\cdots$ ,  $[w_{1P}, w_{2P}, \cdots, w_{QP}]$ ,按照一对浮栅场效应管的阈值电压差值对应一个乘数进行存储,将其映射到模拟矩阵运算单元的P行中,完成大小为 $P \times Q$ 的权重矩阵的存储;

[0046] 经过一个计算周期,通过采样保持器将所有输出的运算结果采样保持至一个时间点,得到完整的卷积结果,输出层向量 $[b_1, b_2, \dots, b_p]$ 。

[0047] 与现有技术相比,本发明具有如下的有益效果:

[0048] (1) 本发明采用模拟矩阵运算单元实现存储与计算一体化进行卷积运算,模拟运算单元在存储权重的同时完成该权重的相关运算,省去了数据的往返传输;采用多脉冲下向量依次输入和依次输出,并对输出进行采样保持,相对于冯诺依曼结构的存储运算分隔体系,节省了90%以上的能效,另一方面模拟计算单元能够完成多位计算且模拟运算单元之间具有高度的计算并行性,相对于数字计算单元实现,极大地减少了面积,大大的提高了的计算密度和效率,此外对于组成模拟矩阵运算单元的浮栅场效应管阵列进行栅极控制选择,可以选定模拟乘加单元的个数,从而快速决定计算的规模,使得卷积运算的实用性和适用性更强。

[0049] (2) 本发明构成模拟矩阵运算单元的模拟乘法电路通过阈值电压差值存储乘数,模拟电压输入与上述两个电流输入端电压差值的形式输入被乘数,以电流差值的函数读出乘积的方式实现乘法运算,由于乘法的乘数是预先存储在闪存单元中的,通过预先存储能够多次复用的乘数,能极大的减少运算数据的读取,并且通过不同的电路拓扑结构可拓展为向量的点积以及向量与矩阵的乘法。

[0050] (3) 本发明用一对浮栅场效应管的阈值电压差值存储一个乘数,以电流差值的形式获取乘积的方法,比仅采用一个浮栅场效应管进行存储的线性度更好,可以达到更高的精度。

[0051] (4) 本发明对于差分电流检测电路的改进型设计,一方面能够使得电流输入端的电压保持稳定,避免了现有技术中当电流改变时,电流流过负载成为负载上的电压明显会改变而引起的计算误差,另一方面,特别的采用工作在深三极管区的浮栅场效应管或MOS管代替常用的电阻作为电流-电压转换器中的阻性负载,实现了输出电流流经作为负载的浮栅场效应管或MOS管,可以抵消工艺参数的影响,确保在不同温度和不同工艺下的乘法结果一致性。

[0052] (5) 本发明更优异的地方在于输出电流流经作为负载的浮栅场效应管,可以通过控制调整作为负载的浮栅场效应管的阈值电压,实现对乘法结果的比例进行灵活缩放的效果。

[0053] (6) 本发明成对出现的浮栅场效应管,能够有效抵消阈值电压受到的体效应的影响,从而保证乘积结果的一致性,因此乘数与被乘数都可以为正值、负值或零,扩大了乘积使用范围,能够完成多比特运算。

[0054] (7) 本发明在浮栅阵列的基础上实现存算一体化,用于卷积神经网络的推理过程,浮栅单元存储权重参数同时完成和此权重相关的乘加法运算,在此基础上完成了卷积神经网络中卷积层和全连接层的运算实现,并配合具体的池化层模块和激活函数模块电路,实现了多层卷积神经网络,具有较好的通用性。

[0055] 关于本发明相对于现有技术,其他突出的实质性特点和显著的进步在实施例部分进一步详细介绍。

## 附图说明

[0056] 通过阅读参照以下附图对非限制性实施例所作的详细描述,本发明的其它特征、目的和优点将会变得更明显:

[0057] 图1是本发明的模拟乘法电路的结构图;

[0058] 图2a是本发明采用的一般差分电流检测电路的结构图;

[0059] 图2b是电流-电压转换器结构图;

[0060] 图2c是电压差分放大器的结构图;

[0061] 图3a是本发明采用MOS管作为阻性负载的改进型差分电流检测电路的结构图;

[0062] 图3b是本发明采用浮栅场效应管作为阻性负载的改进型差分电流检测电路的结构图;

[0063] 图3c是本发明采用MOS管作为阻性负载的电流-电压转换器结构图;

[0064] 图3d是本发明采用浮栅场效应管作为阻性负载的电流-电压转换器结构图;

[0065] 图4是本发明采用改进型差分电流检测电路的模拟乘法电路的结构图;

[0066] 图5是本发明乘累加电路的结构图;

[0067] 图6是本发明卷积运算采用的模拟矩阵运算单元电路结构图;

[0068] 图7是本发明的带数据选择器的模拟矩阵运算单元电路结构图;

[0069] 图8是本发明的采用改进型差分电流检测电路的模拟矩阵运算单元电路结构图;

[0070] 图9是本发明的采用改进型差分电流检测电路的带数据选择器的模拟矩阵运算单元结构图;

[0071] 图10是本发明模拟乘法方法的流程图;

[0072] 图11是本发明的一种应用在模拟矩阵运算单元中的乘法方法的流程图;

[0073] 图12是本发明的基于模拟矩阵运算单元实现2D卷积运算的示意图;

[0074] 图13是本发明的基于模拟矩阵运算单元实现2D卷积运算的原理图;

[0075] 图14是本发明的基于模拟矩阵运算单元实现一个具体的3D卷积运算的示意图;

[0076] 图15是本发明的基于模拟矩阵运算单元实现一个具体的3D卷积运算的原理图;

[0077] 图16是本发明的基于模拟矩阵运算单元实现一个具体的3D卷积运算的时序图;

[0078] 图17是本发明的基于模拟矩阵运算单元实现一般3D卷积操作的示意图;

[0079] 图18是本发明的基于模拟矩阵运算单元实现一般3D卷积操作中将3D卷积核的权重参数映射到模拟矩阵计算单元的流程图;

[0080] 图19是本发明的基于模拟矩阵运算单元实现一般3D卷积操作的时序图;

[0081] 图20是本发明的一种采用上述的基于模拟矩阵运算单元的卷积运算实现卷积神经网络的方法流程示意图;

[0082] 图21是本发明提供的一种ReLU激活函数电路图;

[0083] 图22是本发明提供的一种sigmoid激活函数电路图;

[0084] 图23是本发明提供的一种输入特征图分别经过最大值池化和平均值池化的示意图;

[0085] 图24是本发明提供的一种输入特征图最大值池化的电路原理图;

[0086] 图25是本发明提供的一种全连接层的输入输出示意图;

[0087] 图26是本发明提供的基于模拟矩阵运算单元进行全连接运算的示意图。

## 具体实施方式

[0088] 下面将结合本发明实施例中的附图,对本发明实施例中的技术方案进行清楚、完整地描述,显然,所描述的实施例仅仅是本发明一部分实施例,而不是全部的实施例。基于本发明中的实施例,本领域普通技术人员在没有做出创造性劳动前提下所获得的所有其他实施例,都属于本发明保护的范围。

[0089] 参照说明书附图1-26,对本发明一种基于模拟矩阵运算单元的卷积运算及其应用进行具体实施例描述:

[0090] 实施例1

[0091] 如图1所示,为本实施例中构成模拟矩阵运算单元的单个模拟乘法电路,包括一对浮栅场效应管 $M_1$ 、 $M_2$ 和一个差分电流检测电路;

[0092] 其中 $M_1$ 和 $M_2$ 共栅极并接入电压源,共漏极或共源极并接入模拟电压输入,差分电流检测电路包括两个电流输入端和一个输出端,所述两个电流输入端对应接入 $M_1$ 和 $M_2$ 各自源极或漏极;在这里说明的是如果 $M_1$ 和 $M_2$ 共漏极接入模拟电压输入,则两个电流输入端对应接入 $M_1$ 和 $M_2$ 各自源极,如果 $M_1$ 和 $M_2$ 共源极接入模拟电压输入,则两个电流输入端对应接入 $M_1$ 和 $M_2$ 各自漏极,这是由于浮栅场效应管的源漏极可以互换的结果;

[0093] 本实施例中采用的差分电流检测电路能保持两个电流输入端电压不变,所述输出端结果为两个电流输入端输入电流差值的函数;

[0094] 如图2a所示是满足本实施例功能的一般差分电流检测电路:

[0095] 所述差分电流检测电路包括两个电流-电压转换器和一个电压差分放大器;所述电流-电压转换器由一个运算放大器和一个阻性负载构成,运算放大器的负极输入端为所述电流-电压转换器的电流输入端口,运算放大器的输出端为所述电流-电压转换器的电压输出端口,运算放大器的正极输入端外加直流电平 $V_b$ ,所述阻性负载连接在运算放大器的负极输入端与运算放大器的输出端之间;所述电压差分放大器由一个运算放大器、两个阻值相同的电阻 $R_1$ 、 $R_3$ 和两个阻值相同的电阻 $R_2$ 、 $R_4$ 构成,所述电压差分放大器的两个电压输入端口分别通过电阻 $R_1$ 、 $R_3$ 连接到运算放大器的正、负极输入端,运算放大器的正极输入端通过电阻 $R_2$ 连接到直流电平 $V_b$ ,运算放大器的负极输入端通过电阻 $R_4$ 连接到运算放大器的输出端,运算放大器的输出端为所述电压差分放大器的电压输出端口;

[0096] 具体的,图2a的差分电流检测电路是由图2b的电流-电压转换器和图2c的电压差分放大器构成的。图2b的电流-电压转换器的输入电流为 $I_{in}$ ,通过负载 $R$ 转换为输出电压 $V_{out}$ ,输出与输入有如下关系: $V_{out} = V_b + I_{in}R$ ,实现了电流-电压的转换;图2c的电压差分放大器的输出与输入有如下关系: $V_{out} = V_b + \frac{R_2}{R_1}(V_{in1} - V_{in2})$ ,实现了差分电压放大。

[0097] 参考图10,在上述模拟乘法电路的实现乘法方法,包括以下步骤:

[0098] 步骤 $S_1$ 、对两个浮栅场效应管 $M_1$ 、 $M_2$ 中的浮置栅极执行擦除和写操作,将乘数以及所述一对浮栅场效应管的阈值电压差值的形式进行存储;

[0099] 步骤 $S_2$ 、对两个浮栅场效应管源极、漏极和选栅极执行电压施加,以所述模拟电压输入与所述两个电流输入端电压差值的形式输入被乘数;

[0100] 步骤 $S_3$ 、通过差分电流检测电路输出基于所述一对浮栅场效应管输出电流差值的函数作为所述被乘数与乘数的乘积结果。

[0101] 具体说明上述乘法方法：

[0102] 以图1所述电路结构图为例，对于 $M_1$ 和 $M_2$ 以读取操作进行乘法运算时， $M_1$ 和 $M_2$ 须确保工作在三极管区，而栅源极电压 $V_{GS}$ 、漏源电压 $V_{DS}$ 均相同，在此处 $V_{DS} = V_{in} - V_b$ ，读出 $M_1$ 和 $M_2$

的输出电流 $I_{D1,2}$ 为 $I_{D1,2} = \mu C_{ox} \frac{W}{L} \left[ (V_{GS} - V_{TH1,2}) V_{DS} - \frac{1}{2} V_{DS}^2 \right]$ ， $M_1$ 和 $M_2$ 的电流差值

$\Delta I_D = I_{D1} - I_{D2} = \mu C_{ox} \frac{W}{L} (V_{TH2} - V_{TH1}) V_{DS}$ 可以表示为阈值电压差值 $V_{TH2} - V_{TH1}$ 与漏源电压 $V_{DS}$

的乘积，以阈值电压差值 $V_{TH1} - V_{TH2}$ 为乘数，漏源电压 $V_{DS}$ 为被乘数即可实现乘法。

[0103] 在本实施例中用一对浮栅场效应管的阈值电压差值存储一个乘数，以电流差值的形式获取乘积的方法，乘积结果有很好的线性度，可以达到较高的精度。

[0104] 为了使 $M_1$ 和 $M_2$ 都工作在三极管区，需要满足 $M_1$ 和 $M_2$ 的栅源电压大于两者阈值电压的最大值也即 $V_{GS} > \max \{V_{TH1}, V_{TH2}\}$ ，漏源电压小于栅源电压分别与两者阈值电压差值的最小值，也即 $V_{DS} < \min \{V_{GS} - V_{TH1}, V_{GS} - V_{TH2}\}$ 。

[0105] 在 $V_{GS}$ 固定的情况下，需要对阈值电压 $V_{TH}$ 的动态范围与漏源电压 $V_{DS}$ 的动态范围进行折中选择。对于乘数的存储方式，所述 $M_1$ 和 $M_2$ 的栅源电压 $V_{GS}$ 为固定值，漏源电压 $V_{DS}$ 取值范围为 $[0, V_{DS(max)}]$ ，阈值电压 $V_{TH1}$ 、 $V_{TH2}$ 的取值范围为 $[0, V_{TH(max)}]$ ：

[0106] 当 $V_{TH1} - V_{TH2} = \Delta V_{TH} > 0$ 时，选择 $V_{TH1} = V_{TH(max)}$ ， $V_{TH2} = V_{TH(max)} - \Delta V_{TH}$ ，当 $V_{TH1} - V_{TH2} = -\Delta V_{TH} < 0$ 时，选择 $V_{TH2} = V_{TH(max)}$ ， $V_{TH1} = V_{TH(max)} - \Delta V_{TH}$ ；这种选择是为了选择较大的阈值电压可以进一步减小浮栅场效应管的电流以减小功耗。

[0107] 场效应管的体效应：阈值电压会受源极-衬底的电压 $V_{SB}$ 影响，一般施加的衬底电压 $V_b$ 是固定的，p衬底接地，n衬底接最高电位，所以只看 $V_s$ 当被乘数 $V_{in} - V_b > 0$ 时， $V_s = V_b$ 是恒定值；当被乘数 $V_{in} - V_b < 0$ 时， $V_s = V_{in}$ 是随输入改变的，所以阈值电压也随输入改变。具体的在本实施例中，如图2所示输入的被乘数为漏源电压即 $V_{in} - V_b$ 。当 $V_{in} > V_b$ 时，电流由电压输入端流向差分电流检测电路的输入端口，两个闪存单元的上端为漏极，下端为源极；当 $V_{in} < V_b$ 时，电流由差分电流检测电路的输入端口流向电压输入端，两个闪存单元的上端为源极，下端为漏极，此时由于源极的电位由外部输入，闪存单元的阈值电压受到体效应的影响，变化量为 $\Delta V_{TH} = \gamma(\sqrt{2\Phi_F + V_{SB}} - \sqrt{2\Phi_F})$ ，由于两个闪存单元的源极-衬底电压 $V_{SB}$ 始终相同，阈值电压受到体效应影响的变化量 $\Delta V_{TH1,2}$ 相同，因此存储的阈值电压差值 $V_{TH1} - V_{TH2}$ 不变，乘积结果不受影响。因此，被乘数可以为正值、负值或零。

[0108] 在本实施例中，可以继续对差分电流检测电路进行进一步改进，如图3a和3b所示，主要是阻性负载采用浮栅场效应管或MOS管，所述浮栅场效应管或MOS管的栅极外加栅极电压，漏极连接到所述运算放大器的输出端，源极连接到所述运算放大器的负极输入端；

[0109] 为了保证浮栅场效应管或MOS管工作在深三极管区，所述浮栅场效应管或MOS管的漏源电压 $V_{DS}$ 、栅源电压 $V_{GS}$ 和阈值电压 $V_{TH}$ 满足： $V_{DS} \ll 2(V_{GS} - V_{TH})$ 。

[0110] 更进一步具体的说明，参照图3c与图3d中的电流-电压转换器分别采用MOS管或

浮栅场效应管作为负载，工作在深三极管区的阻抗为 $R_{M_r} = \frac{1}{\mu C_{ox} \frac{W}{L} (V_{GS} - V_{TH})}$ 。

[0111] 在此改进的差分电流检测电路中输出电流流经作为负载的浮栅场效应管或MOS管，可以抵消工艺参数的影响，确保在不同温度和不同工艺角下的乘法结果一致性。

[0112] 继续参考图4，乘数以浮栅场效应管的阈值电压差值 $V_{TH1} - V_{TH2}$ 存储，可以为正值、负值或零，被乘数以输入电压值 $V_{in}$ 输入。闪存单元M1、M2的电流为

$$I_{D1,2} = \mu C_{ox} \frac{W}{L} \left[ (V_G - V_{in} - V_{TH1,2})(V_b - V_{in}) - \frac{1}{2}(V_b - V_{in})^2 \right],$$

$$V_{1,2} = V_b + \frac{(V_G - V_{in} - V_{TH1,2})(V_b - V_{in}) - \frac{1}{2}(V_b - V_{in})^2}{V_{GS(MR)} - V_{TH(MR)}},$$

$$V_{out} = V_b + \frac{R_2 (V_{TH2} - V_{TH1})(V_b - V_{in})}{R_1 V_{GS(MR)} - V_{TH(MR)}} = V_b + K(V_{TH2} - V_{TH1})(V_b - V_{in}),$$

$$= K(V_{TH2} - V_{TH1})(V_b - V_{in}),$$

系数 $K = \frac{R_2}{R_1} \frac{1}{V_{GS(MR)} - V_{TH(MR)}}$ 与工艺参数 $\mu C_{ox}$ 无关，且电阻为比值

的形式。

[0113] 在此改进型差分电流检测电路中，输出电流流经作为负载的浮栅场效应管，可以通过控制调整作为负载的闪存单元的阈值电压，实现对乘法结果的比例进行灵活缩放的效果。

[0114] 参考图5，本实施例中直接构成模拟矩阵运算单元一种乘累加电路，包括若干个模拟乘法电路，所述若干个模拟乘法电路共栅极且共同一个差分电流检测电路，所述每个模拟乘法电路中该对浮栅场效应管的共漏极或共源极分别接入对应的输入电压信号；

[0115] 所述差分电流检测电路的输出端结果为每个模拟乘法电路输出的基于该对浮栅场效应管输出电流差值的函数之和。

[0116] 在本实施例中差分电流检测电路可以获取多对所述浮栅场效应管电流差值之和，或者，可以获取多对所述浮栅场效应管中的第一个浮栅场效应管电流之和，以及多对所述浮栅场效应管中的第二个浮栅场效应管电流之和，再得到其差值，以实现乘累加运算的效果。

[0117] 参考图6是本实施例中卷积运算采用的模拟矩阵运算单元，所述模拟矩阵运算单元为一个P行Q列的模块电路，所述每一行均为一个包括Q个模拟乘法电路的乘累加电路；

[0118] 所述每列浮栅场效应管的共漏极或共源极接入同一的输入电压信号。

[0119] 具体而言本实施例中P行Q列的模块电路，参照图11中的乘法流程图，按照一对浮栅场效应管的阈值电压差值对应一个乘数进行存储，完成大小为 $P \times Q$ 的乘数矩阵的存储

$$\begin{pmatrix} b_{11} & \cdots & b_{1Q} \\ \vdots & \ddots & \vdots \\ b_{P1} & \cdots & b_{PQ} \end{pmatrix};$$

每列浮栅场效应管的一端源或漏极相连构成位线(BL)信号，以每个乘法

电路中模拟电压输入与所述两个电流输入端电压差值对应一个被乘数，完成元素个数为Q的横向被乘数向量的输入 $(a_1 \cdots a_Q)$ ；模块电路中Q行中每行的栅极接入同一个字线(WL)信号，另一端源或漏极每行相连构成两条源/漏极线(SDL)，接到一个差分电流检测电路，产生一个乘累加输出，所有的Q个输出可表示为输出向量

$(c_1 \cdots c_Q) = K(a_1 \cdots a_Q) \begin{pmatrix} b_{11} & \cdots & b_{1Q} \\ \vdots & \ddots & \vdots \\ b_{P1} & \cdots & b_{PQ} \end{pmatrix}^T$ , 实现了向量与矩阵的乘法。

[0120] 参照图8,是本实施例中模拟矩阵运算单元也可以采用改进型差分电流检测电路的实现的电路示意图,其差分电流检测电路为采用浮栅场效应管作为阻性负载的改进型

差分电流检测电路,实现向量与矩阵的乘法:  $(c_1 \cdots c_Q) = K(a_1 \cdots a_Q) \begin{pmatrix} b_{11} & \cdots & b_{1Q} \\ \vdots & \ddots & \vdots \\ b_{P1} & \cdots & b_{PQ} \end{pmatrix}^T$ , 被乘数

$a_j = V_{in(j)} - V_b$ , 乘数  $b_{ij} = V_{TH(i,j)} - V'_{TH(i,j)}$ , 系数  $K = \frac{R_2}{R_1} \frac{1}{V_{GS(MR)} - V_{TH(MR)}}$ , 乘法结果  $c_i =$

$V_{out(i)} - V_b$ 。

[0121] 参照图7,本实施例中进一步改进的模拟矩阵运算单元,包括第一数据选择器和第二数据选择器,所述第一数据选择器设置在每列浮栅场效应管的共漏极或共源极与其对应的输入电压信号之间,所述第二数据选择器设置在每行浮栅场效应管的共源极或共漏极与差分电流检测电路的电流输入端之间;

[0122] 所述第一、二数据选择器组合用于从P行Q列的模块电路任意选取X行、Y列单元形成新的工作模块;以减少电路中输入信号、差分电流检测电路与输出信号的数量。

[0123] 参照图9,本实施例中进一步改进的模拟矩阵运算单元电路结构示意图,其差分电流检测电路为采用浮栅场效应管作为阻性负载的改进型差分电流检测电路。

[0124] 本实施例中卷积运算包括以下步骤:

[0125] 将卷积核经过转换拼接成P个长度为Q的横向量,并映射到矩阵  $\begin{pmatrix} b_{11} & \cdots & b_{1Q} \\ \vdots & \ddots & \vdots \\ b_{P1} & \cdots & b_{PQ} \end{pmatrix}$  中;

[0126] 将输入特征图按照滑动窗口进行切分转换成N个长度为Q的横向量;

[0127] 在N个脉冲波下,依次将N个长度为Q的横向量映射到向量  $(a_1 \cdots a_Q)$  中,矩阵运算单元的模拟输出端口按照脉冲波序列,输出N个长度为P的运算结果;

[0128] 通过采样保持器将所有输出的运算结果采样保持至一个时间点,得到完整的卷积结果。

[0129] 下面以具体维度的卷积计算解释实施例1中的如何基于上述模拟矩阵运算单元实现卷积方法:

[0130] 实施例2

[0131] 参考图12,是本实施例的一个基于模拟矩阵运算单元实现2D卷积运算的示意图:其中输入特征图的维度为 $5 \times 5$ ,输出特征图的维度为 $3 \times 3$ ,卷积核的维度为 $3 \times 3$ ,未考虑偏置bias。

[0132] 图13是该2D卷积用模拟矩阵运算单元实现的原理图;将大小为 $3 \times 3$ 卷积核转换为横向量  $[w_{11}, w_{12}, w_{13}, w_{21}, w_{22}, w_{23}, w_{31}, w_{32}, w_{33}]$  将其映射到模拟矩阵运算单元的一行;将输入特征图按照滑动窗口转换为9个横向量,即  $[a_{11}, a_{12}, a_{13}, a_{21}, a_{22}, a_{23}, a_{31}, a_{32}, a_{33}]$ ,  $[a_{12},$

$a_{13}, a_{14}, a_{22}, a_{23}, a_{24}, a_{32}, a_{33}, a_{34}, \dots, [a_{33}, a_{34}, a_{35}, a_{43}, a_{44}, a_{45}, a_{53}, a_{54}, a_{55}]$ , 将其转换为模拟电压, 然后在9个脉冲波下从模拟输入端口并行序列依次输入; 在模拟输出端口依次输出9个运算结果,  $[b_{11}, b_{12}, b_{13}, b_{21}, b_{22}, b_{23}, b_{31}, b_{32}, b_{33}]$ 。通过采样保持器, 将所有运算结果采样保持至一个时间点, 得到一个 $3 \times 3$ 的输出特征图。1个脉冲进行一次滑窗运算, 9个脉冲构成一个完整的卷积周期, 进行一次完整的卷积计算, 输出一个完整2D卷积结果;

[0133] 实施例3

[0134] 参考图14, 是本实施例基于模拟矩阵运算单元实现一个具体的3D卷积运算的示意图:

[0135] 其中输入特征图维度为 $3 \times 3 \times 3$ , 输出特征图维度为 $3 \times 3 \times 2$ , 采用了2个 $1 \times 1 \times 3$ 的卷积核, 2个 $1 \times 1 \times 1$ 的偏置bias。

[0136] 参考图15, 是该3D卷积用模拟矩阵运算单元实现的原理图, 将2个大小为 $1 \times 1 \times 3$ 卷积核和2个 $1 \times 1 \times 3$ 的偏置bias转换为2个横向量 $[w_{11}, w_{12}, w_{13}, b_1], [w_{21}, w_{22}, w_{23}, b_2]$ , 将其映射到模拟矩阵运算单元的两行; 将输入特征图按照滑动窗口转换为9个横向量, 即 $[a_{111}, a_{112}, a_{113}, 1], [a_{121}, a_{122}, a_{123}, 1], \dots, [a_{331}, a_{332}, a_{333}, 1]$ 。其转换为模拟电压, 然后在9个脉冲波下从模拟输入端口并行序列依次输入; 在模拟输出端口依次输出9个运算结果 $[b_{111}, b_{112}], [b_{121}, b_{122}], \dots, [b_{331}, b_{332}]$ , 通过采样保持器, 将所有运算结果采样保持至一个时间点, 得到一个 $3 \times 3 \times 2$ 的输出特征图。一个脉冲进行1次滑窗运算, 9个脉冲构成一个完整的卷积周期, 进行一次完整的卷积计算, 输出一个完整3D卷积结果;

[0137] 参考图16是该3D卷积运算的时序图。

[0138] 实施例4

[0139] 参照图17-19, 本实施例是基于模拟矩阵运算单元的一般3D卷积操作;

[0140] 图17是该操作的示意图, 其中输入特征图的维度是 $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$  (即第 $l-1$ 层的特征图维度(长 $\times$ 宽 $\times$ 通道)), 输出特征图的维度是 $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$  (即第 $l$ 层的特征图维度(长 $\times$ 宽 $\times$ 通道)), 卷积操作使用了 $n_C^{[l]}$ 个维度为 $f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$  (即第 $l$ 层的3D卷积核维度(长 $\times$ 宽 $\times$ 通道))的卷积核;

[0141] 图18是将3D卷积核的权重参数映射到模拟矩阵计算单元的流程; 每个3D卷积核的维度为 $f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$ , 将其切片成 $n_C^{[l-1]}$ 个大小为 $f^{[l]} \times f^{[l]}$ 的2D卷积核, 然后将这些2D卷积核转换为横向量, 拼接成一个横向量, 即 $[w_{111}, w_{121}, \dots, w_{ff1}, w_{112}, \dots, w_{ff2}, \dots, w_{ffn}]$ , 向量长度为 $f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$ ; 一共有 $n_C^{[l]}$ 个大小为 $f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$ 的3D卷积核, 将其全部切片展开为 $n_C^{[l]}$ 个长度为 $f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$ 的横向量, 然后将它们拼接为一个维度为 $n_C^{[l]} \times (f^{[l]} \times f^{[l]} \times n_C^{[l-1]})$ 的权重矩阵, 将其映射到矩阵运算模块, 具体为: 权重矩阵第一行映射模拟矩阵计算单元第一行, 矩阵第二行映射模拟矩阵计算单元第二行, 矩阵第三行映射模拟矩阵计算单元第三行……

[0142] 图19是3D卷积操作进行运算的时序图; 进行一次完整的卷积操作需要 $n_H^{[l]} \times n_W^{[l]}$ 个脉冲周期, 将输入特征图按照滑动窗口进行切分,

[0143] 可以转化为 $n_H^{[l]} \times n_W^{[l]}$ 个横向量, 每个横向量的长度为 $f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$ , 即 $[a_{11},$

$a_{12}, \dots, a_{1f}, a_{21}, \dots, a_{2f}, \dots, a_{ff}] \dots [a_{12}, a_{13}, \dots, a_{ff+1}, a_{22}, \dots, a_{2f+1}, \dots, a_{f+1f+1}]$ 。将其通过DAC转换为模拟电压,然后在 $n_H^{[l]} \times n_W^{[l]}$ 个脉冲波下从模拟输入端口并行序列依次输入;矩阵运算单元的模拟输出端口按照脉冲波序列输出 $n_H^{[l]} \times n_W^{[l]}$ 个长度为 $n_c^{[l]}$ 的运算结果,通过采样保持器,将其输出保持在一个时间点,于是得到一个大小为 $n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$ 的3D输出矩阵;一个脉冲进行一次滑窗运算, $n_H^{[l]} \times n_W^{[l]}$ 个脉冲构成一卷积运算周期,进行一次完整的卷积计算,输出一个完整3D卷积结果。

[0144] 实施例5

[0145] 参照图20,本实施例提供一种采用上述的基于模拟矩阵运算单元的卷积运算实现卷积神经网络的方法,实现了一个输入层,卷积层,池化层,全连接层,输出层共5层得卷积神经网络,具体包括DAC电路、卷积层电路、激活函数电路、采样保持电路、池化层电路、全连接层电路和ADC电路,各个电路间的级联构成卷积神经网络;

[0146] 所述的DAC电路用于将输入层的数字输入转化为模拟输入;

[0147] 所述卷积层电路与全连接层电路均采用所述基于模拟矩阵运算单元的卷积运算实现功能并构建电路;

[0148] 所述的激活函数电路将卷积层的输出电流转换为电压,同时完成激活函数作用;

[0149] 所述的采样保持电路用于将卷积层的序列计算结果采样并保持在同一时间点;

[0150] 所述的池化层电路用于池化层的实现,包括平均池化电路和最大池化电路;

[0151] 所述的ADC电路用于将输出层的模拟输出转换为数字输出;

[0152] 具体步骤包括:

[0153] 输入特征图经过DAC电路转换为模拟信号输入;

[0154] 经过卷积层电路得到第一次卷积计算结果,然后经过激活函数电路,将计算结果从电流转变为电压,并且进行非线性化;

[0155] 经过采样保持电路,得到一个完整的中间特征图,再经过池化电路,将模拟信号降维;

[0156] 继续经过全连接层电路进行全连接运算,最后经过激活函数电路得到模拟输出特征图,并经ADC电路转换为数字输出。

[0157] 在本实施例中激活函数可以采用ReLU激活函数或者sigmoid激活函数,具体的:

[0158] 参考图21,为满足本实施例的一种ReLU激活函数电路图,该ReLU激活函数的电路包括:电流电压转换器,反相器,电压限幅器;电流电压转换器用于将来自矩阵计算单元的模拟输出端口的电流转换为电压;反相器起到电压反向缓冲的作用;电压限幅器采用二极管方式,将大于0的电压输出,小于0的电压保持为0;输入电流与输出电压的关系为

$$V_{\text{out}} = \begin{cases} R_1 \cdot I_{\text{in}} & I_{\text{in}} > 0 \\ 0 & I_{\text{in}} < 0 \end{cases}$$

[0159] 参考图22,为满足本实施例的一种sigmoid激活函数电路图,该sigmoid激活函数的电路包括:带偏置的电流电压转换器,反相器,电压限幅器;电流电压转换器用于将来自矩阵计算单元的模拟输出端口的电流转换为电压,在电流电压转换器的负极性端串联电阻并提供偏置电压;反相器起到电压反向缓冲的作用;电压限幅器将输出电压限制在一

定的电压范围内,输入电流与输出电压的关系为

$$[0160] \quad V_{\text{out}} = \begin{cases} 0 & I_{\text{in}} < -\frac{V_{\text{bias}}}{R_2} \\ R_1 \cdot I_{\text{in}} + \frac{V_{\text{bias}} \cdot R_1}{R_2} & -\frac{V_{\text{bias}}}{R_2} < I_{\text{in}} < \frac{V_{\text{bias}}}{R_2} ; \\ V_{DD} & I_{\text{in}} > \frac{V_{\text{bias}}}{R_2} \end{cases}$$

[0161] 本实施例中还提供了—个卷积神经网络的池化层电路池化操作的示例,包括平均值池化和最大值池化具体的:

[0162] 参考图23是一个维度为 $4 \times 4$ 的输入特征图分别经过最大值池化和平均值池化的示意图;

[0163] 参考图24是最大值池化的电路原理图,包括配置单元和一些平均池化单元和最大池化单元。配置单元用来根据池化降维的大小配置平均池化单元和最大池化单元的个数,一个池化单元(平均或者最大)能够进行一个池化窗口的操作,就是将多个模拟信号取最大值或者平均值输出。最大池化单元由模拟信号比较器和模拟多路复用器构成,模拟信号比较器比较池化窗口的多路模拟输入信号的大小,找出最大模拟输入信号,然后通过模拟多路复用器选通最大模拟输入信号,过滤其余模拟信号;平均池化单元由模拟反向加法器和一个电压反向器串联组成,池化窗口的多路模拟输入信号从 $v_i^1, v_i^2, \dots, v_i^n$ 输入,池化结果

从 $v_0^1$ 输出。当 $R_1 = R_2 = \dots = R_n = nR_i$ 时,  $V_0 = \frac{v_i^1 + v_i^2 + \dots + v_i^n}{n}$ ,起到一个平均模拟输入电压的作用。

[0164] 本实施例中给出一个全连接层电路均采用所述基于模拟矩阵运算单元的卷积运算实现的示例,具体的:

[0165] 参考图25是一个全连接层的示意图,输入层为 $[a_1 \dots a_Q]$ ,输出层为 $[b_1, b_2, \dots, b_P]$ ,共有 $P \times Q$ 个权重参数,未考虑偏置;

[0166] 参考图26,是用模拟矩阵运算单元进行全连接运算的示意图;将输入层 $[a_1 \dots a_Q]$ 按照每个乘法电路中模拟电压输入与所述两个电流输入端电压差值对应一个 $a_i$ 其中 $i \in (1, 2, \dots, Q)$ ,从模拟矩阵运算单元模拟电压输入端输入;

[0167] 将全连接层的 $P \times Q$ 个权重拆分为 $P$ 个横向量,即 $[w_{11}, w_{21}, \dots, w_{Q1}]$ ,  $[w_{12}, w_{22}, \dots, w_{Q2}]$ ,  $\dots, [w_{1P}, w_{2P}, \dots, w_{QP}]$ ,按照一对浮栅场效应管的阈值电压差值对应一个乘数进行存储,将其映射到模拟矩阵运算单元的 $P$ 行中,完成大小为 $P \times Q$ 的权重矩阵的存储;

[0168] 经过一个计算周期,通过采样保持器将所有输出的运算结果采样保持至一个时间点,得到完整的卷积结果,输出层向量 $[b_1, b_2, \dots, b_P]$ 。

[0169] 本发明的一种模拟乘法电路、模拟乘法方法及其应用能通过以乘数预先存储、被乘数在运算时输入的方式实现两个数的乘法、两个向量的点积以及向量与矩阵的乘法,运算速度快、电路功耗低、方法简单易行,具有较高的实用价值和广泛的应用前景。

[0170] 对于本领域技术人员而言,显然本发明不限于上述示范性实施例的细节,而且在不背离本发明的精神或基本特征的情况下,能够以其他的具体形式实现本发明。因此,无

论从哪一点来看,均应将实施例看作是示范性的,而且是非限制性的,本发明的范围由所附权利要求而不是上述说明限定,因此旨在将落在权利要求的等同要件的含义和范围内的所有变化囊括在本发明内。不应将权利要求中的任何附图标记视为限制所涉及的权利要求。

[0171] 此外,应当理解,虽然本说明书按照实施方式加以描述,但并非每个实施方式仅包含一个独立的技术方案,说明书的这种叙述方式仅仅是为清楚起见,本领域技术人员应当将说明书作为一个整体,各实施例中的技术方案也可以经适当组合,形成本领域技术人员可以理解的其他实施方式。

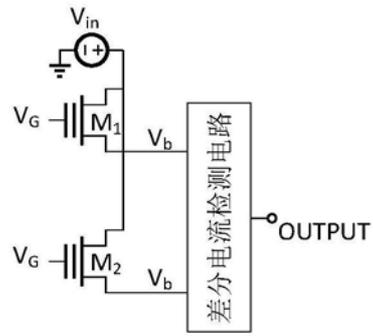


图1

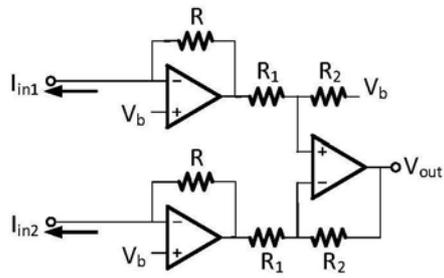


图2a

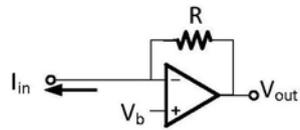


图2b

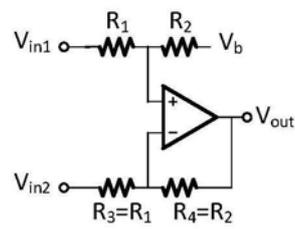


图2c

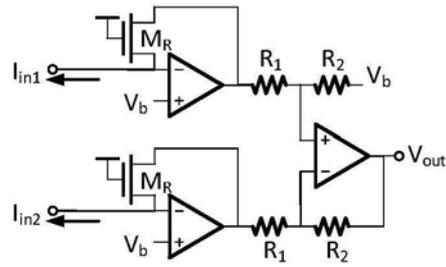


图3a

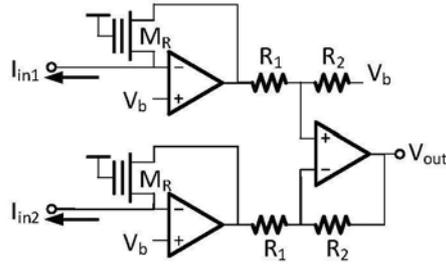


图3b

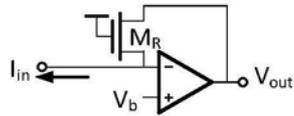


图3c

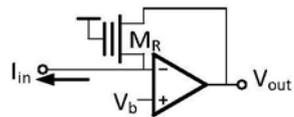


图3d

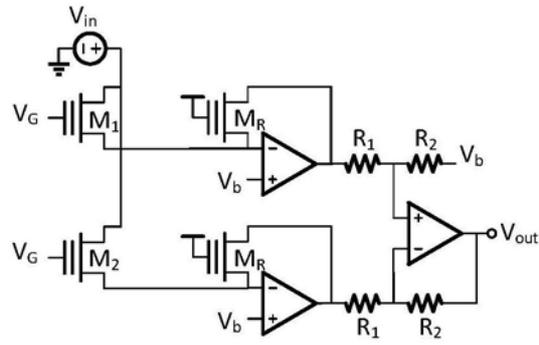


图4

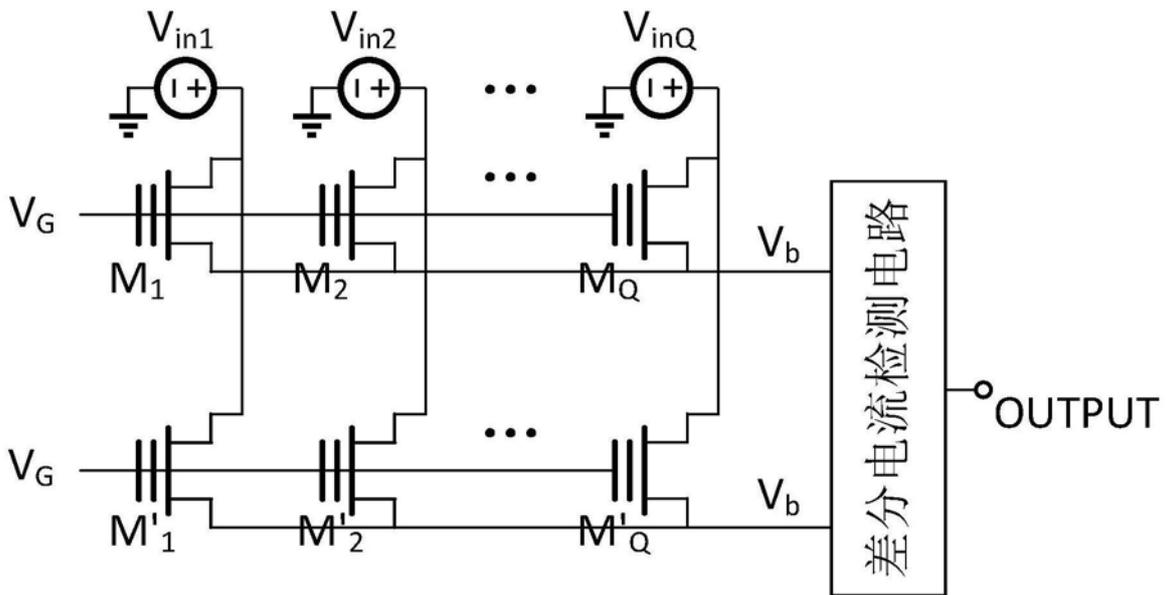


图5

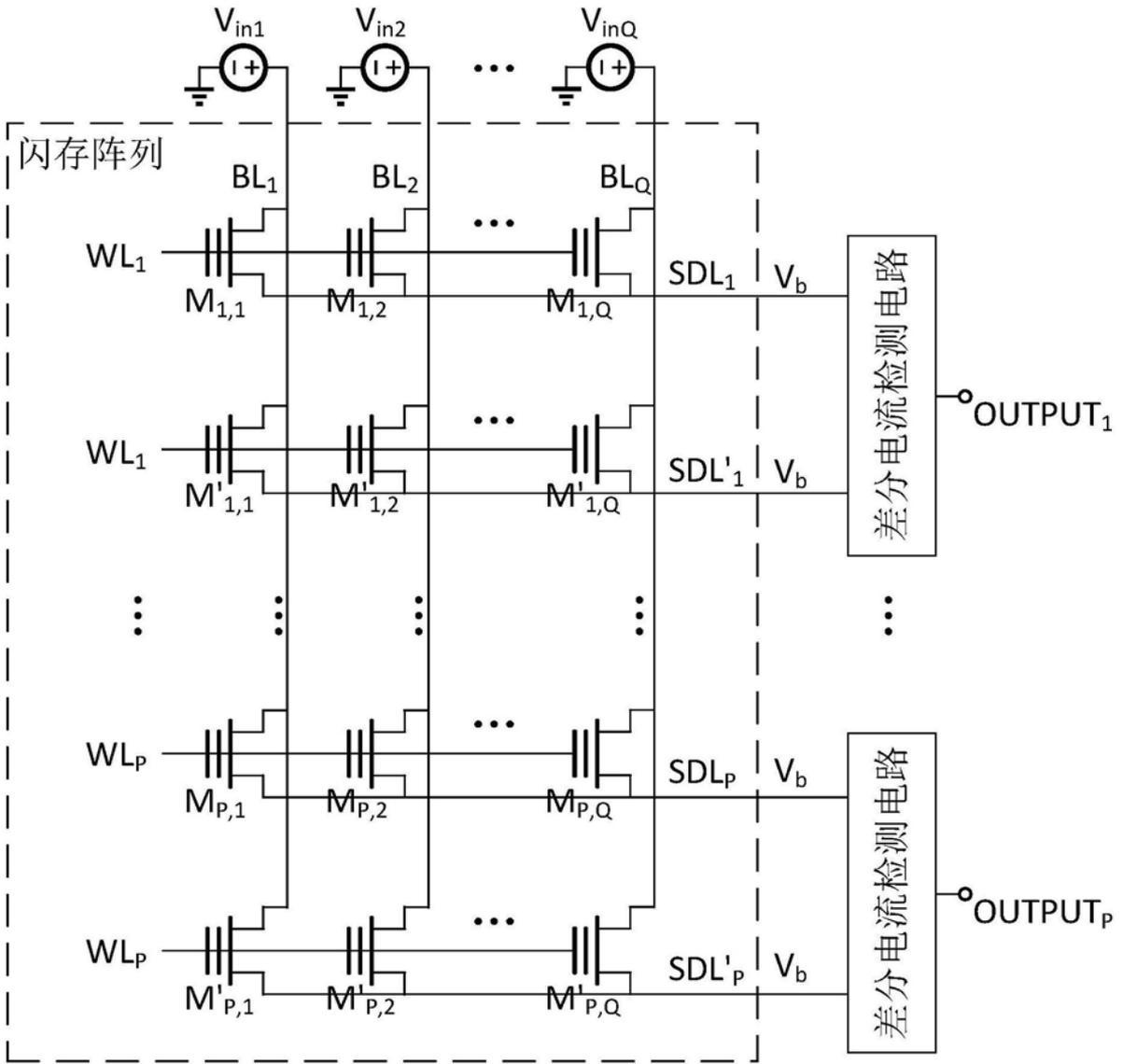


图6

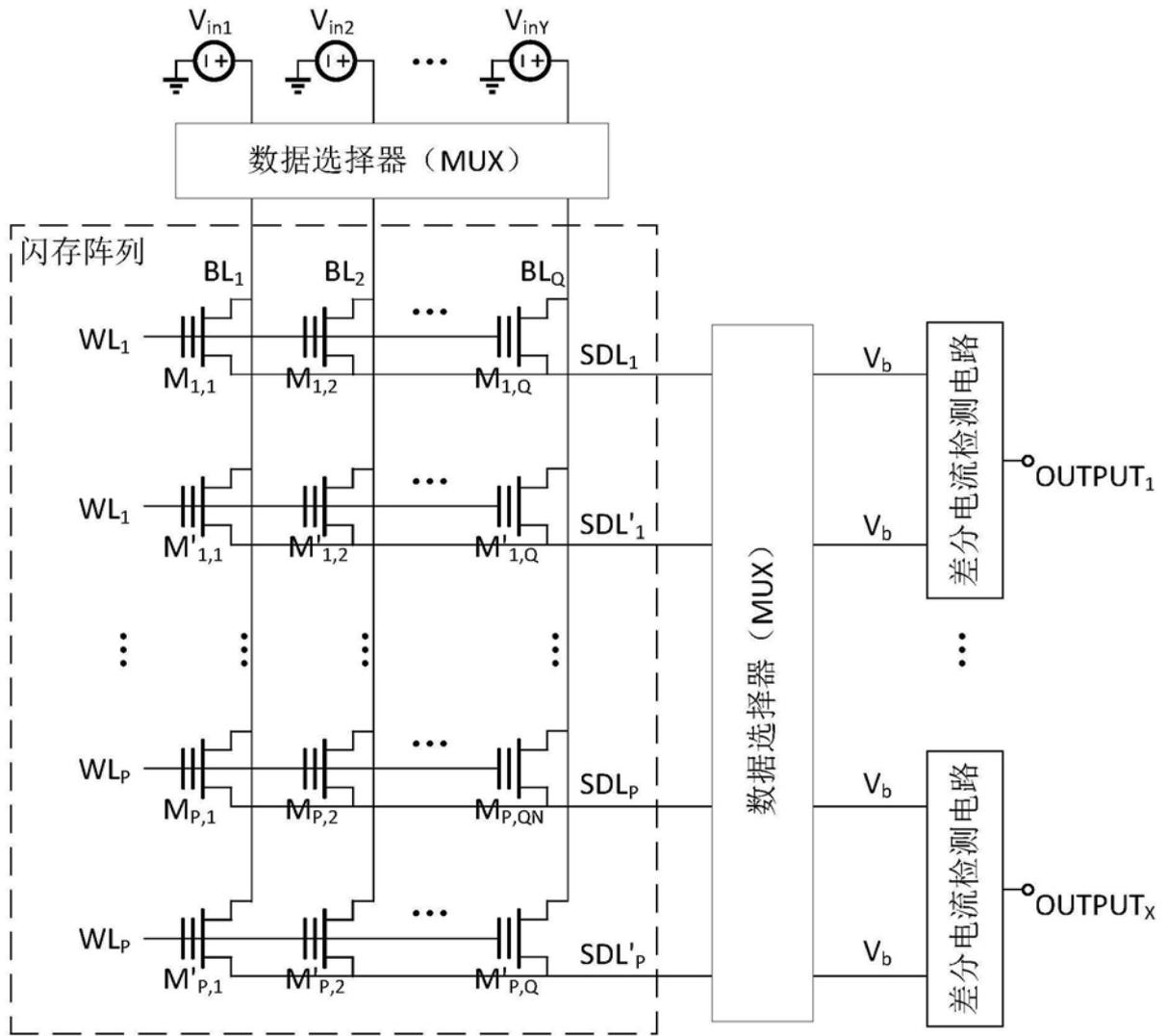


图7

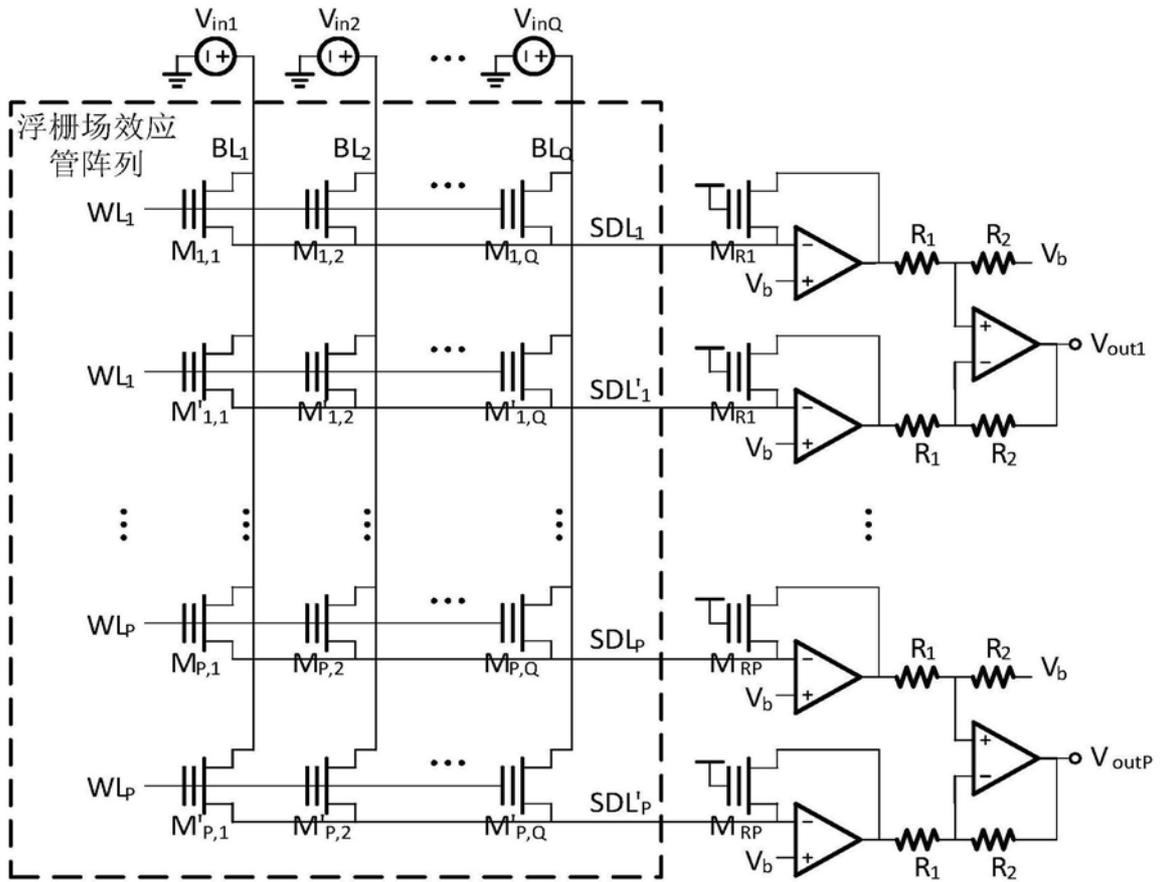


图8

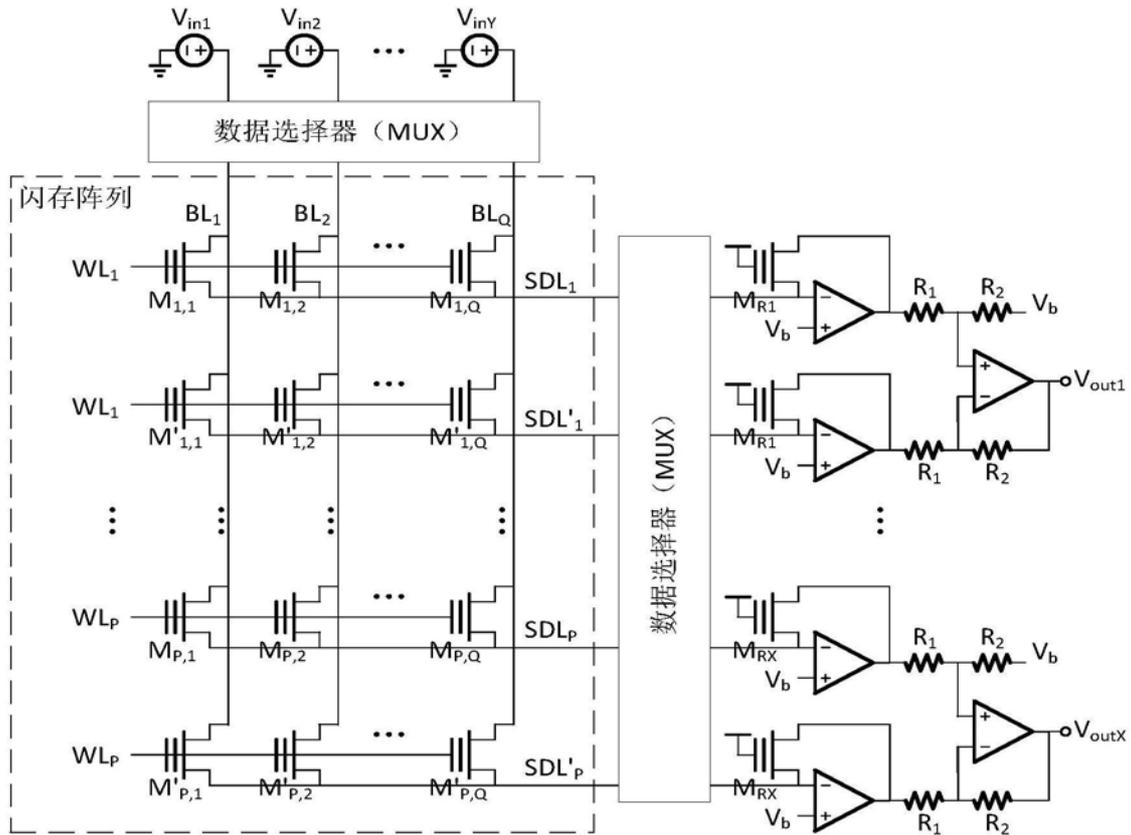


图9

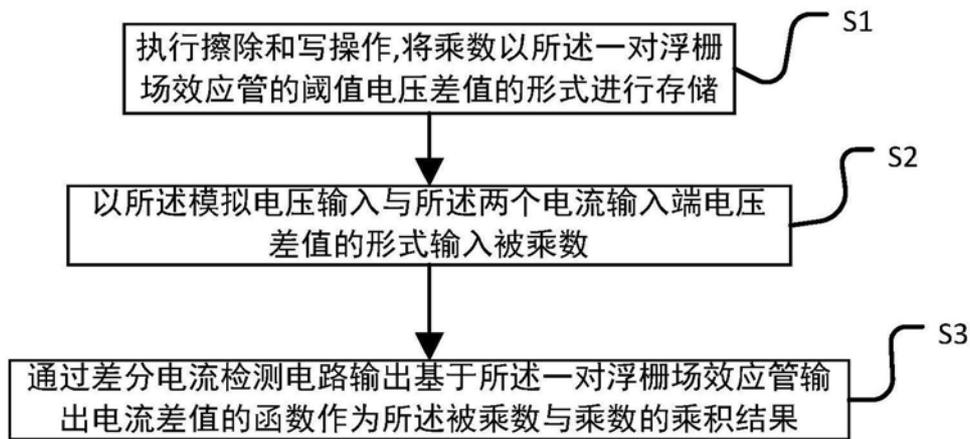


图10

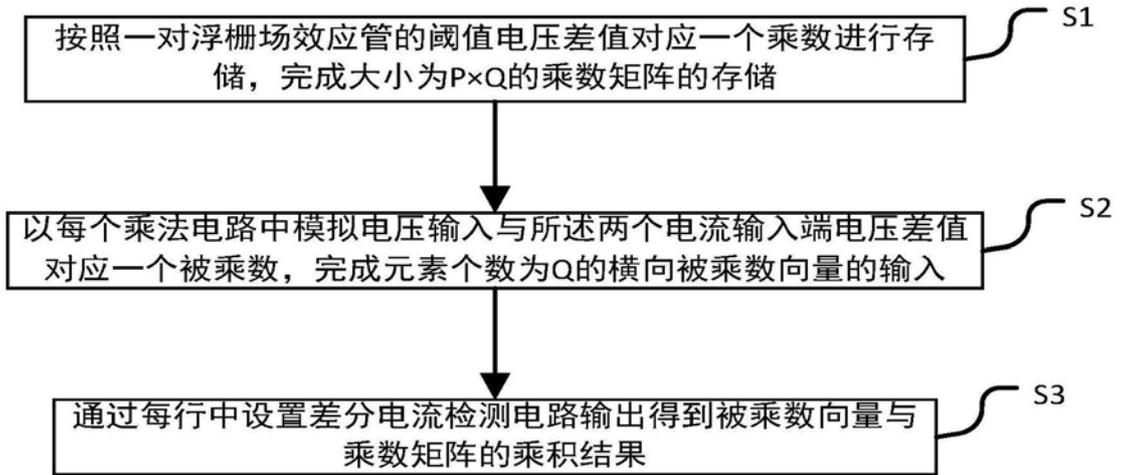


图11

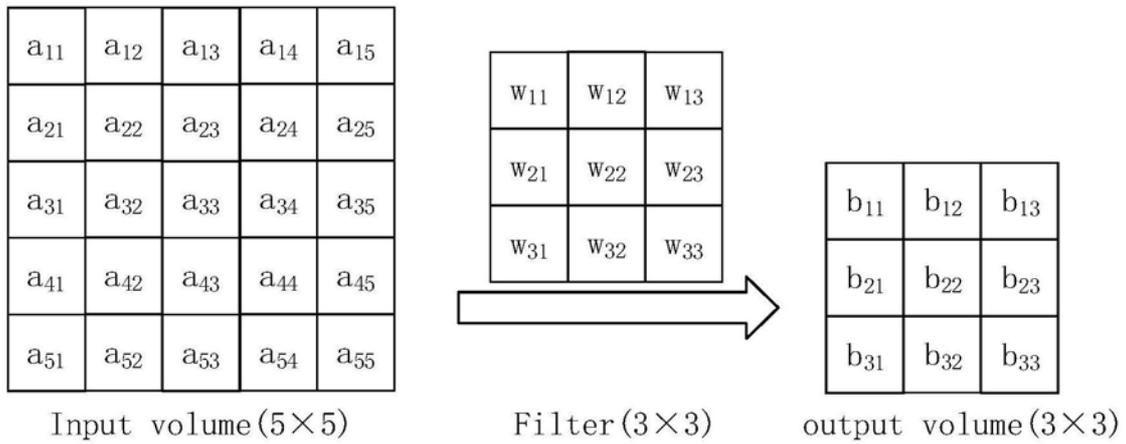


图12

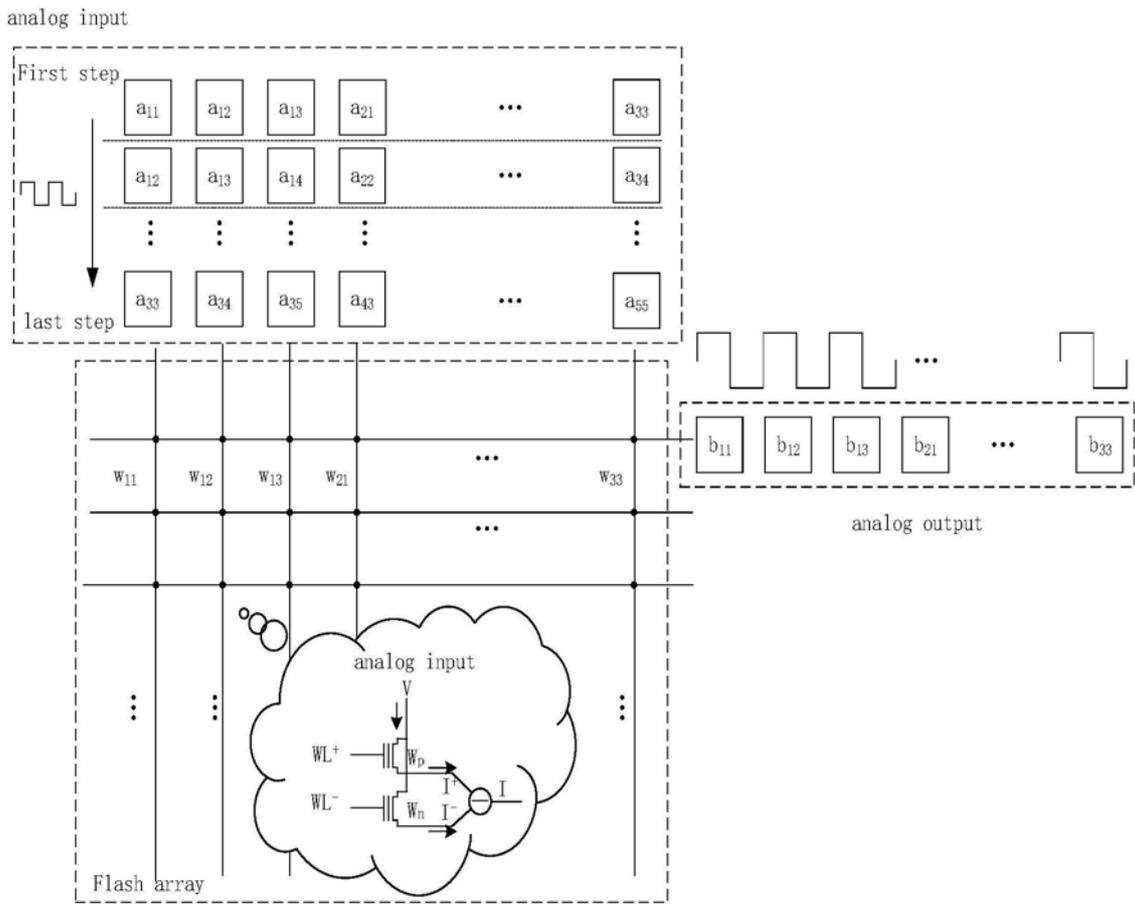


图13

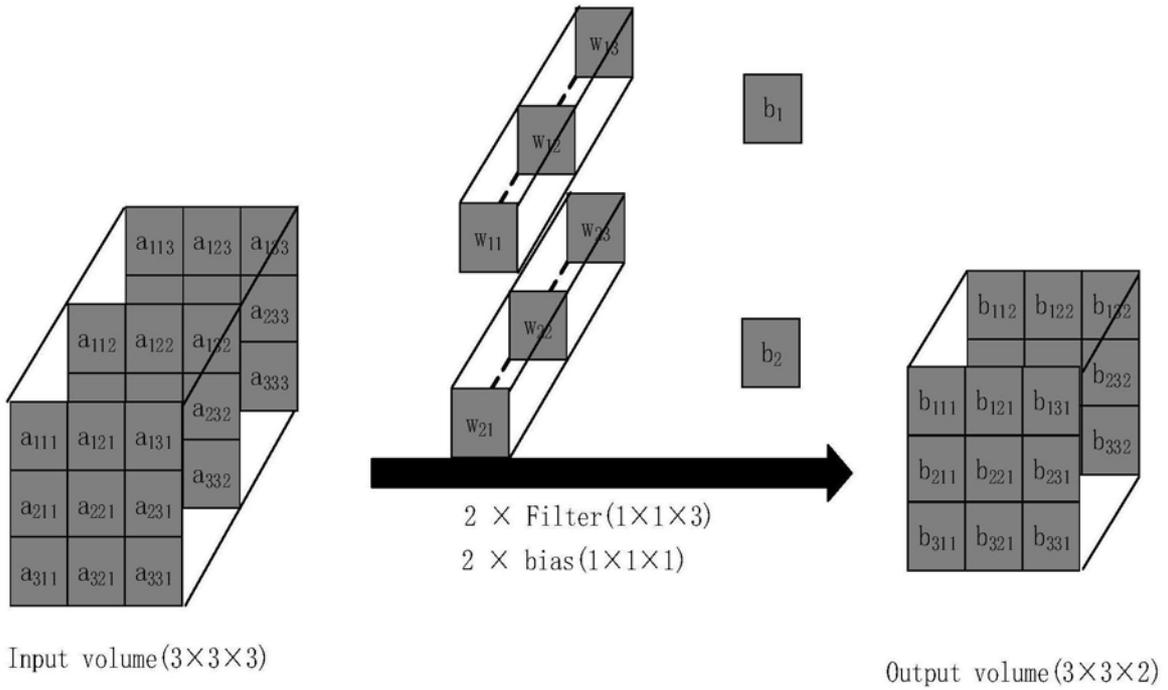


图14

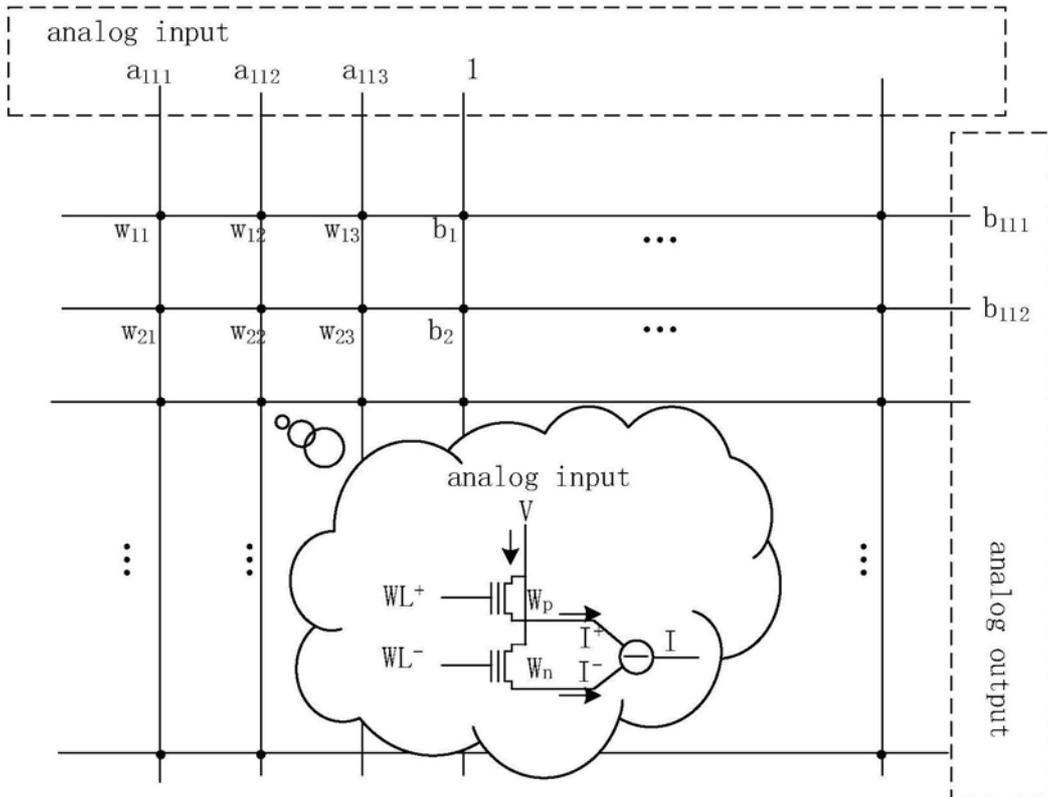


图15

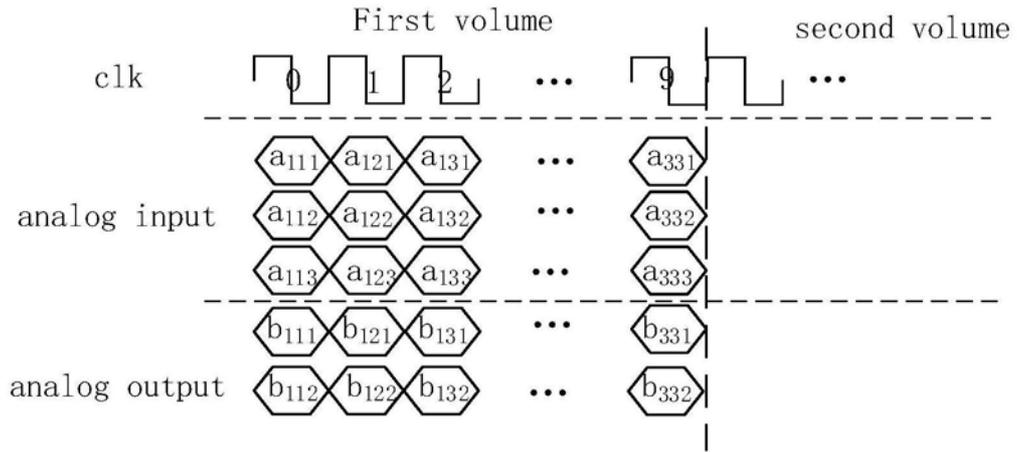


图16

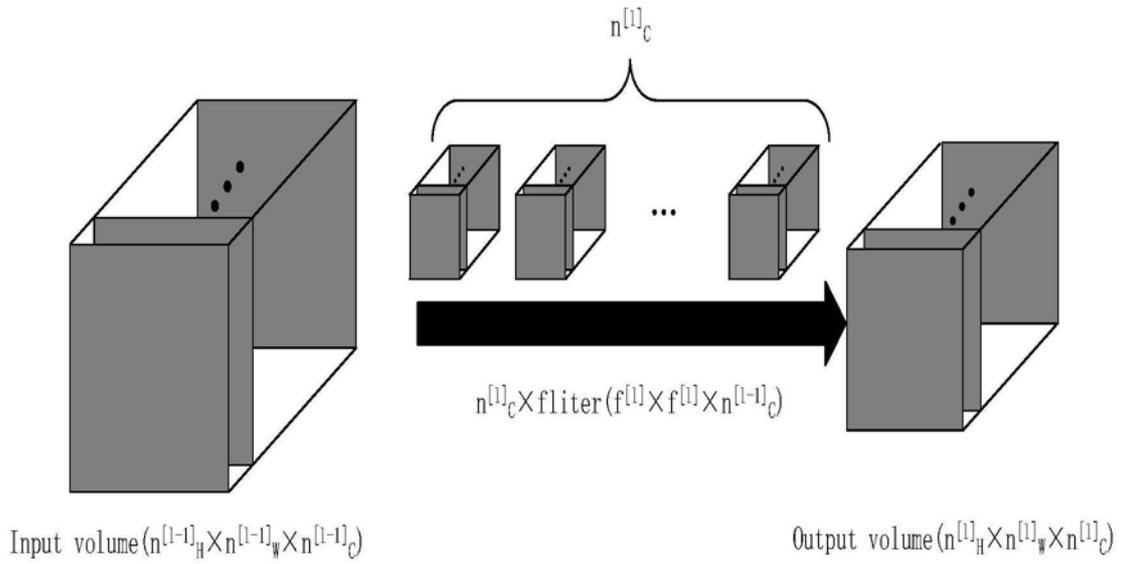


图17

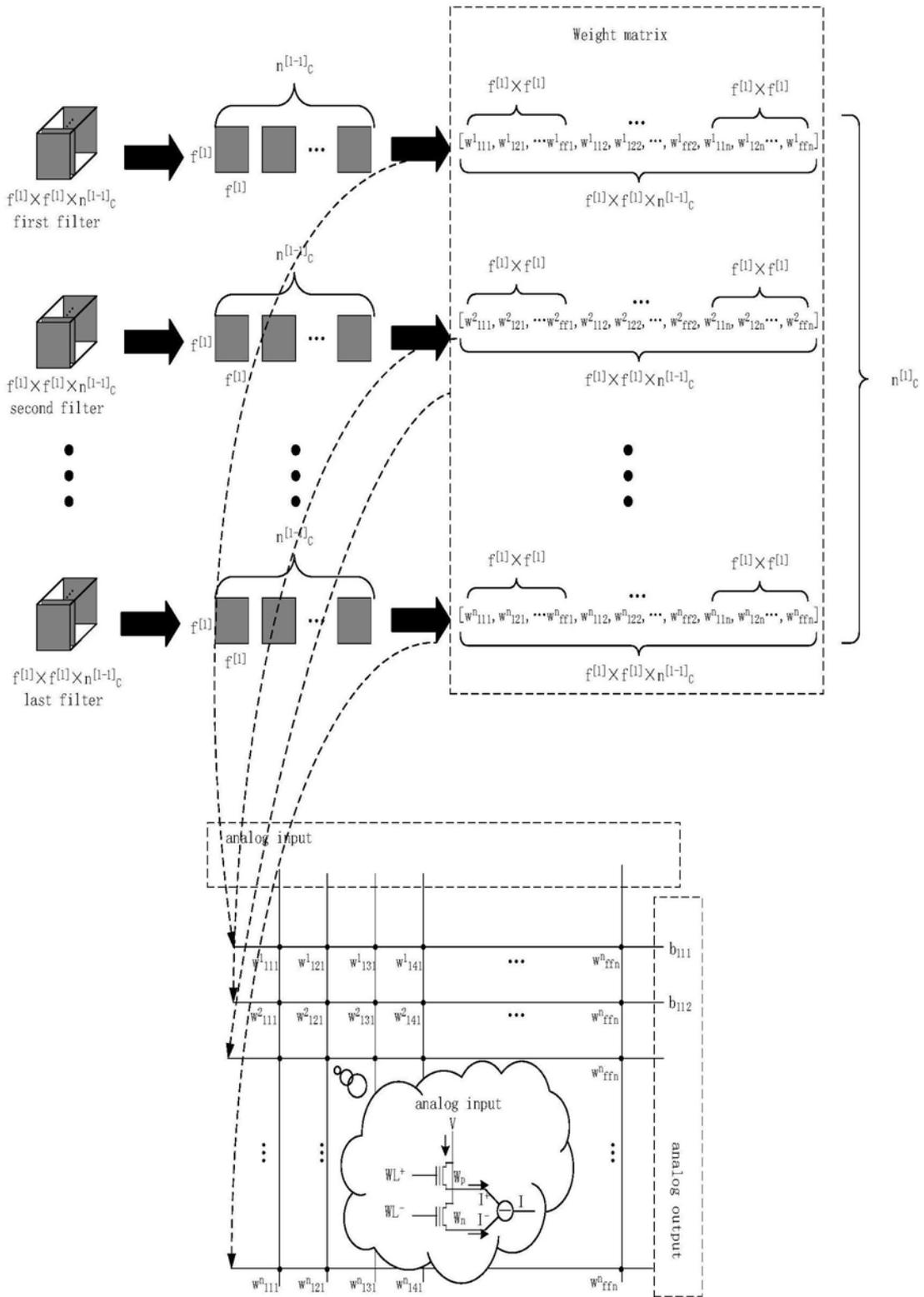


图18

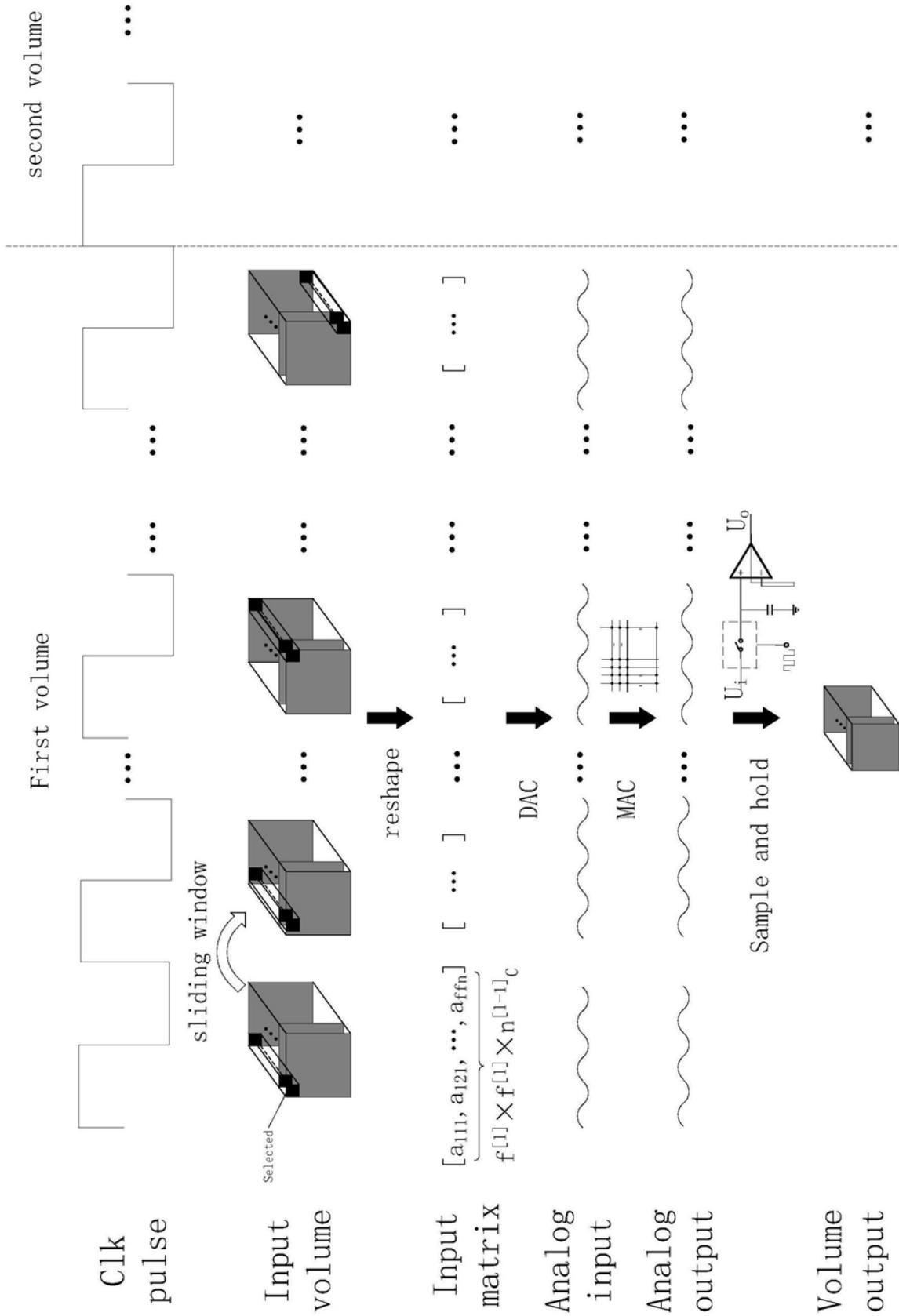


图19

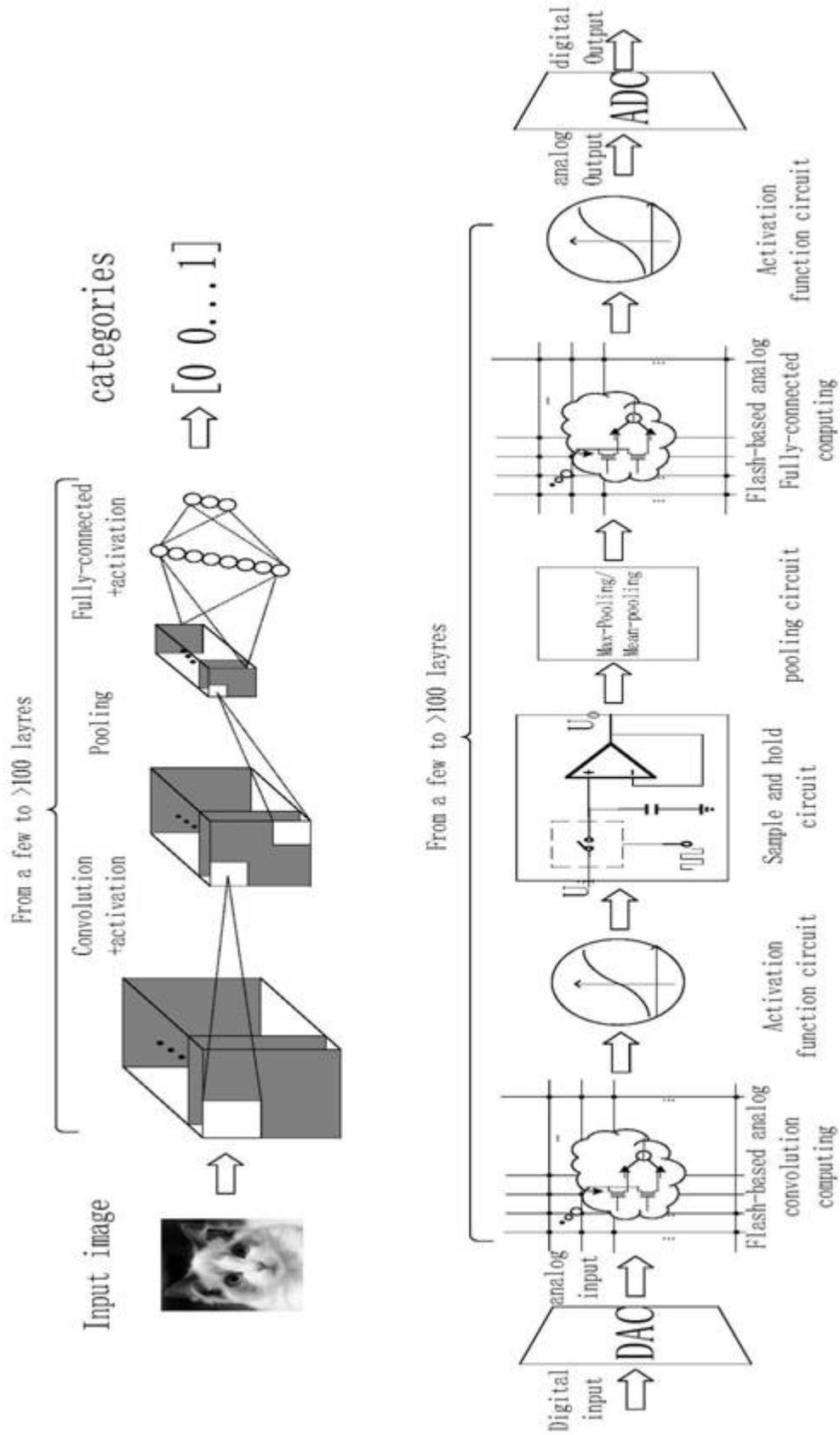


图20

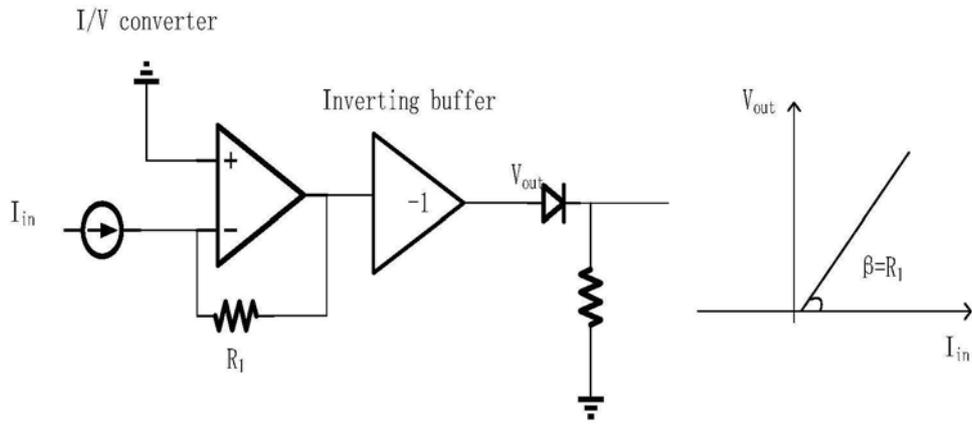


图21

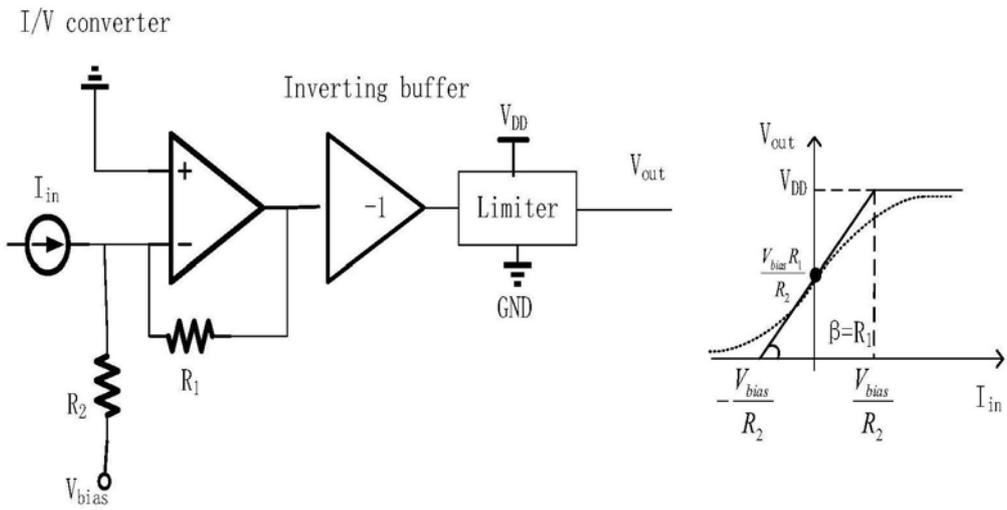


图22

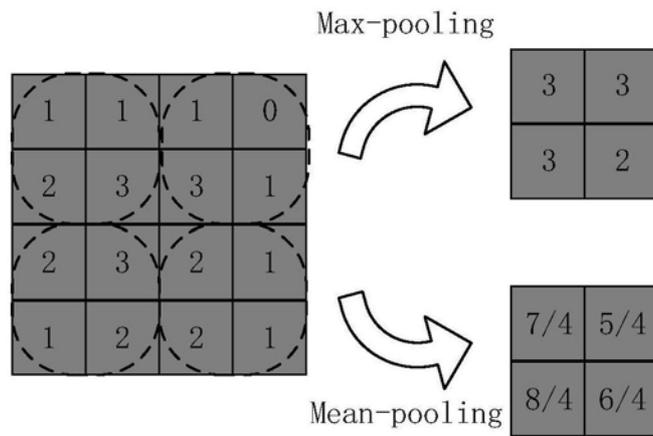


图23

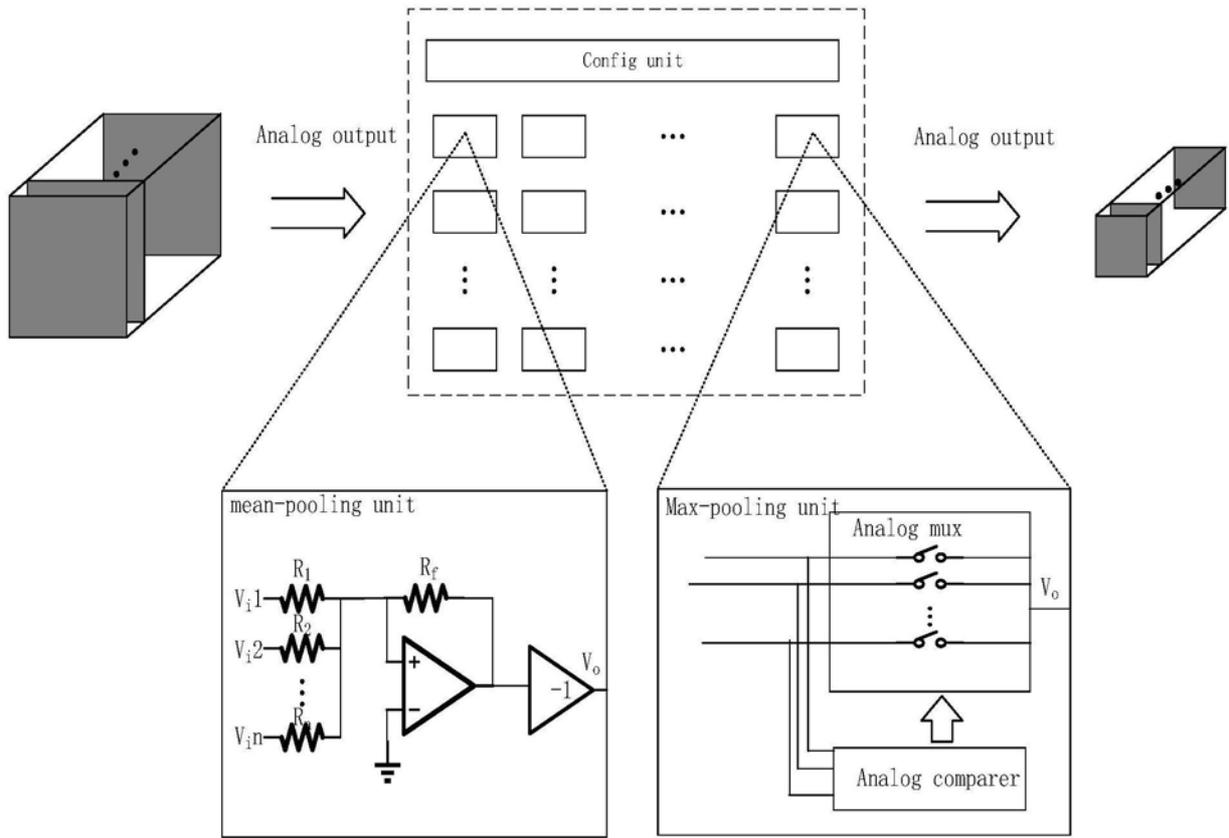


图24

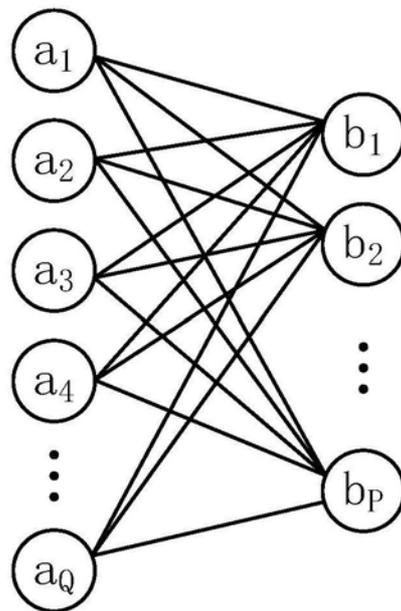


图25

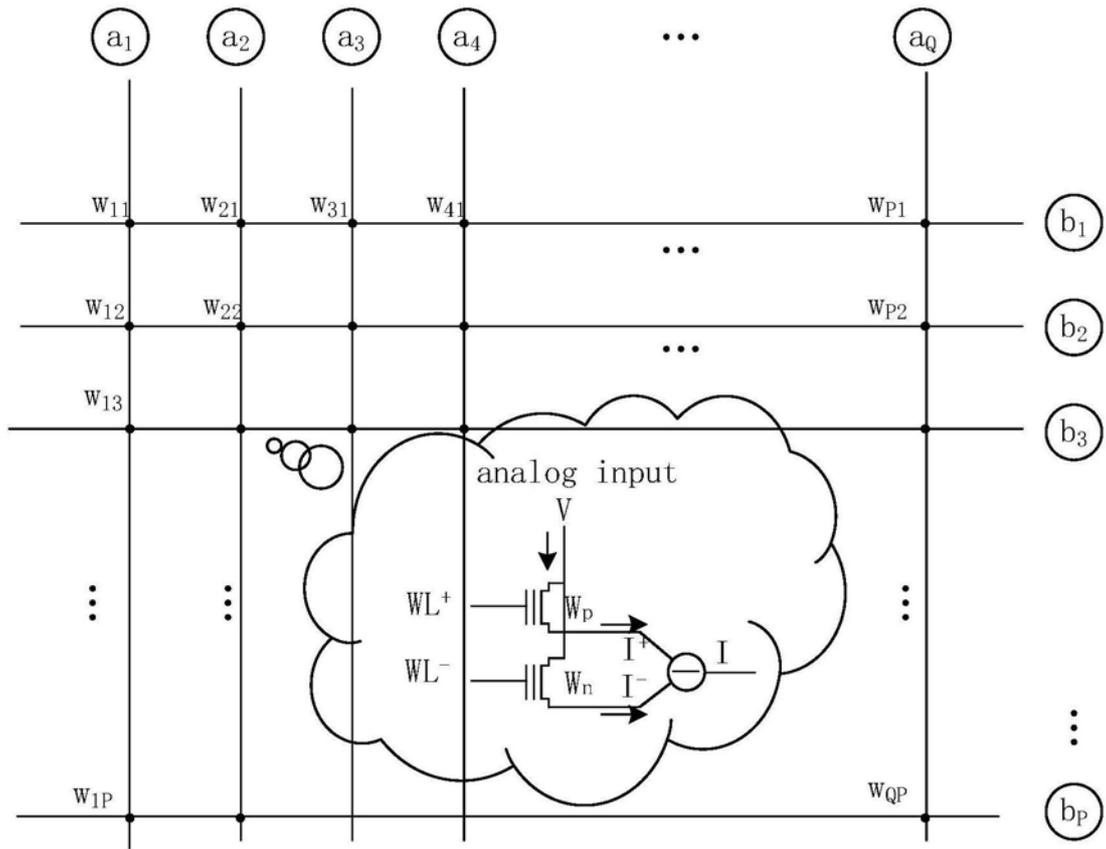


图26