Ph.D. DISSERTATION

# Synaptic and Neuron Devices for Excitatory and Inhibitory Signals in Neuromorphic Systems

신경 모방 시스템에서 흥분 및 억제 신호를 위한 시냅스와 뉴런 소자

by

SUNGYUN WOO

February 2021

DEPARTMENT OF ELECTRICAL ENGINEERING AND
COMPUTER SCIENCE
COLLEGE OF ENGINEERING
SEOUL NATIONAL UNIVERSITY

# Synaptic and Neuron Devices for Excitatory and Inhibitory Signals in Neuromorphic Systems

신경 모방 시스템에서 흥분 및 억제성 신호를 위한 시냅스와 뉴런 소자

지도교수 이 종 호

이 논문을 공학박사 학위논문으로 제출함

2021 년 2 월

서울대학교 대학원

전기정보공학부

우 성 윤

우성윤의 공학박사 학위논문을 인준함

2021 년 2 월

위 원 장 : 박 병 국 (인)

부위원장 : 이 종 호 (인)

위　　원 : 유 승 주 (인)

위　　원 : 김 재 준 (인)

위　　원 : 김 재 하 (인)

# Synaptic and Neuron Devices for Excitatory and Inhibitory Signals in Neuromorphic Systems

by

Sung Yun Woo

Advisor: Jong-Ho Lee

A dissertation submitted in partial fulfillment of

the requirements for the degree of

Doctor of Philosophy

(Electrical and Computer Engineering)

in Seoul National University

February 2021

Doctoral Committee:

Professor Byung-Gook Park, Chair

Professor Jong-Ho Lee, Vice-Chair

Professor Sungjoo Yoo

Professor Jae-Jun Kim

Professor Jaeha Kim

# ABSTRACT

Recently, hardware-based neural networks (HNNs) have emerged since neuromorphic systems can compute complex data efficiently. Various synaptic devices and neuron circuits suitable for architectures and learning algorithms have been researched for high performance in HNNs. Specifically, processing simultaneously both excitatory ($G^+$) and inhibitory ($G^-$) signals transmitted from synaptic arrays are important to process the computation efficiently and improve the performance of HNNs.

In this dissertation, synaptic and neuron devices are proposed for the neuromorphic system with high density and low power consumption. A positive-feedback (PF) device simultaneously processing excitatory and inhibitory signals is used as the neuron device to replace conventional neuron circuits. Owing to the steep switching characteristics of the PF operation, the PF neuron device can reduce the energy consumption during processing integration function of neurons. The PF neuron device is an efficient structure that merges a gated thyristor and a single

MOSFET. By accumulating electrons in an *n* floating body of the PF neuron device,

the integrate-and-fire operation with steep subthreshold swing (SS < 1 mV/dec) is

experimentally implemented. The electrons accumulated in the *n* floating body are

discharged by applying inhibitory signals to the merged FET. Moreover, the

threshold voltage ($V_{th}$) of the proposed PF neuron with a non-volatile memory

function is controlled by program and ease states in a charge storage layer. The PF

neuron circuit that consumes low energy per a spike (~ 0.62 pJ/spike) consists of

one PF device and only five MOSFETs for the integrate-and-fire function and reset

operation. The dual-gate FET with independent two gates (G1 and G2) is proposed

as the synaptic device. Here, G1 turns on and off the synaptic device, and G2 with

the charge storage layer controls the conductance of the dual-gate FET for synaptic

weights. The range of conductance change of the dual-gate FET is very wide (100

pA ~ 1 μA). In the NOR type array based on the dual-gate FETs, program and erase

operations can be implemented with the Fowler-Nordheim (FN) tunneling

mechanism, resulting in low power consumption during the synaptic weight update.

The sum of current (3.63 μA) of eight individual dual-gate FETs is almost the same

(~ 0.87 %) as the $I_{total}$ (3.6 µA) of eight dual-gate FETs in the NOR type synapse array. The variations ($\sigma/\mu$) of the quantized synaptic currents in eight synaptic devices are obtained as 0.023, 0.011, 0.015, and 0.032 for four different synaptic weight states. The PF neuron circuit and synapse array based on the dual-gate FETs provide viable solutions for high-density and low-energy neuromorphic systems.

Keywords: neuromorphic system, positive-feedback (PF) neuron device, dual-gate FET, NOR type synapse array, excitatory/inhibitory signals, hardware-based neural network.

Student number: 2014-21721

# CONTENTS

**Chapter 3**

# Chapter 4

# Chapter 5

# List of Figures

# Chapter 1

# Introduction

## 1.1  Neuromorphic computing

Recently, deep neural networks (DNNs) based on the back-propagation (BP) algorithm have shown excellent performance in many areas including object recognition, internet of things (IoT), autonomous driving, and translation [1]-[6]. The state-of the-art architectures of DNNs include recurrent neural networks (RNNs), generative adversarial networks (GANs) and convolutional neural networks (CNNs). However, to find the optimal weights in DNN based on BP, the vector-to-vector matrix multiplication (VMM) of forward propagation (FP) and BP occupies a large part in the computational task. To solve these problems, many research groups have investigated hardware-based neural networks (HNNs) [7]-[9]. Hardware-based neural networks (HNNs) can easily implement VMM as the current of the synapse array, which is the product of the input voltage and conductance. HNNs are basically composed of synapse array and neuron circuits as

shown in Fig. 1.1. When input signals are applied to the synapse array, the sum of

the currents reflecting the synaptic weights is transmitted to neuron circuits. Neuron

circuits perform an integrate-and-fire operation, generating an output signal to the

next synaptic array. In hardware-based neural networks, the characteristics of

synaptic devices and the stability of neural circuits are very important. Therefore,

many research groups have developed various synaptic devices and neuron circuits

to implement neuromorphic systems for HNNs [10]-[34]. Next, we describe

emerging devices as synaptic devices and the requirements of synaptic devices. And,

we investigate neuron devices to replace a membrane capacitor of neuron circuit

and reduce low energy consumption.

Fig. 1.1. Typical structure of hardware-based neural networks composed of input layer, hidden layers, and output layer.

### 1.1.1  Synaptic devices

Recently, resistive random access memory (RRAM), phase change random access memory (PCRAM), magnetic random access memory (MRAM), ferroelectric material-based devices, and FET-based devices with memory functions are emerged as synaptic devices [10]-[19]. These synaptic devices require high scalability, high reliability, low current in sum and update of synaptic weights, and CMOS compatible technology for neuromorphic systems [18], [19]. The high density of synapse array is very important because large-scale neural networks are composed of thousands or tens of thousands of synaptic devices. Also, the retention characteristics and reproducibility of memory functions in the synaptic devices are important because the conductance change of the synaptic devices with memory function represents the weight of synapses. And, a low current of the synaptic devices is required to prevent fan-in and fan-out in the neuromorphic systems. The synaptic devices should be CMOS compatible because synapse arrays are manufactured with peripheral circuits including neuron circuits. The synapse array is composed of a pair of synaptic devices to improve the performance of HNN, and

the pair of synapses is $G^+$ (excitatory synapse) and $G^-$ (inhibitory synapse), respectively [9], [17]. The excitatory and inhibitory signals reflecting the synaptic weights are transmitted to the neuron circuit.

## 1.1.2 Neuron devices

Conventional neuron circuits consist of a membrane capacitor ($\geq$ 0.1 pF) and many transistors ($\geq$ 11 MOSFETs) to implement the integrate-and-fire function as shown in Fig. 1.2 [19]-[21]. In the neuron circuits, processing excitatory and inhibitory signals simultaneously can reduce memory usage and simplify the peripheral circuitry. The size of the membrane capacitor is 100 $\mu m^2$ when a capacitance of the membrane capacitor is 0.5 pF in 0.35 μm CMOS technology [21]. In a 400–128–10 sized HNN composed of a synaptic array, current mirror circuits, and neuron circuits, neuron circuits consume 415.3 μW, which is about 33.48% of the total power consumption [22]. Most of the power is consumed by the leakage or generation of output spikes. To solve the problems, memristor-based [23-28] and FET-based neuron devices [29-34] with memory functionalities have been researched to implement the integrate-and-fire function of neurons. Memristor-based neuron devices with memory functionalities are used to replace membrane capacitors in neuron circuits, and the structure of memristors with two-terminals is the advantage of high density compared to membrane capacitors and FET-based

neuronal devices. However, generation of output spikes requires infinite endurance of neuron devices during integrate-and-fire operation. Existing memristor-based neuron devices are being studied to improve endurance to generate output spikes [25], [26]. Also, memristor-based neuron devices require a differential amplifier to compare the resistance of the memristor to a reference resistance, and a reset circuit to reset the memristor after generating the output spike [27], [28]. On the other hand, existing FET-based neuron devices can process only one type of synaptic signals or signals transmitted from excitatory and inhibitory of synapses ($G^+$ and $G^-$) sequentially [29]-[34]. So, FET-based neuron devices processing only one type of signals should be paired for processing excitatory and inhibitory signals, and require additional circuits in the neural networks. In neuron devices that process signals transmitted from excitatory and inhibitory synapses sequentially, the computing and inferencing speeds can also be slower than conventional neuron circuits. It is also difficult for the reported neuron devices to control the threshold voltage to reduce variations of the neuron devices. Thus, synaptic and neuron devices with high density and low power consumption should be developed for

neuromorphic systems.

**Integration**  **Fire and reset**



Fig. 1.2. Conventional neuron circuit composed of a membrane, a refractory

capacitor and a comparator [21].

## 1.2　Purpose of research

We propose a neuron device with steep switching characteristics using positive feedback (PF) [35] and a dual-gate FET consisting of independent two gates as a synaptic device. Here, a neuron circuit based on the PF device and a synaptic array based on the dual-gate FETs were fabricated on the same wafer. We investigate the integrate-and-fire function by considering excitatory and inhibitory signals simultaneously. By accumulating (or discharging) electrons into an $n$ floating body, the PF neuron device can implement an integrate-and-fire function of neurons. The PF neuron device with steep switching characteristics can replace a large membrane capacitor and a comparator in conventional neuron circuits. Moreover, a threshold voltage ($V_{th}$) of the PF neuron device is changed by program and erase state of a charge storage layer, thereby adjusting the threshold of neurons in neural networks. The threshold tuning ability of neuron circuits can alleviate the degradation of recognition rate by device variations in neural networks. In NOR type synaptic array based on the dual-gate FETs, selective program and erase operation with Fowler-Nordheim (FN) tunneling mechanism in the charge storage layer can be

implemented by using the independent two gates.

## 1.3  Dissertation outline

Based on the above description, this work mainly focuses on the PF neuron device processing simultaneously excitatory and inhibitory signals and the synaptic array based on the dual-gate FET. The remainder of this dissertation is organized as follows. In Chapter 2, the device structure of the PF neuron device is described and the integrate-and-fire operation of the PF neuron device with excitatory and inhibitory signals is explained. Also, the threshold voltage controllability of the PF neuron device with non-volatile memory function is explained. In Chapter 3, the dual-gate FET fabricated with the PF neuron device on the same wafer are proposed and investigated as a synaptic device. In Chapter 4, a NOR type synaptic array based on the dual-gate FET is described with the selective program and erase operation, and the VMM operation performed by the 8×4 NOR type synaptic array is explained. In Chapter 5, concludes this dissertation with a summary.

# Chapter 2
# Neuron device

## 2.1  Device structure

The nervous system of the brain consists of excitatory and inhibitory synapses. Biological neurons integrate signals transmitted from excitatory and inhibitory synapses, firing spikes to next synapses. Even in neuromorphic systems, it is very important to simultaneously process excitatory and inhibitory signals in neuron circuits to improve the performance of neural networks.

Fig. 2.1 (a) and (b) show 3-D schematic and top views of the proposed PF neuron device, respectively. A structure of the PF neuron device is an efficient structure that merges one FET and a gated PNPN-junction diode. The PF neuron device is consisting of a drain, a cathode, an anode, and three gates (G1, G2, and G3). A gate insulating material of G1 and G1 is a silicon oxide layer ($SiO_2$). A gate stack of G3 consists of a tunneling oxide layer ($SiO_2$), a charge storage layer ($Si_3N_4$), and a blocking oxide layer ($SiO_2$). G1 receives excitatory signals transmitted from

12

the pre-synapse array, which charges electrons in the floating $n$-body region to perform the integrate-and-fire operation of the neuron. On the other hand, G2 receives inhibitory signals transmitted from the pre-synapse array, which discharges electrons in the floating $n$-body region to inhibit the integrate-and-fire operation of the neuron. G3 can modulate an initial potential of the $n$-body region. As a bias applied to G3 ($V_{G3}$) increases, more electrons charged in the $n$-body region are required for the integrate-and-fire operation. Also, the initial potential of the $n$-body region is modulated by the amount of charges in the charge storage layer of the dual-gate FET. By applying the bias or program/erase pulses to G3 with the charge storage layer, the threshold voltage ($V_{th}$) of the PF neuron device can be controlled. And, Fig. 2.2 shows a TEM image of the fabricated PF neuron device. The thicknesses of the gate oxide of G1 ($T_{ox}$) and the $SiO_2/Si_3N_4/SiO_2$ stack of G3 are 10 nm and 3/6/9 nm, respectively. The thickness of Si body ($T_{Si}$) and the channel width (W), $n$-body length ($L_n$), and $p$-body length ($L_p$) are 100 nm, 1 μm, 1.1 μm, and 0.7 μm, respectively. Doping concentrations of $p$ and $n$-body are $1\times10^{18}$ cm-3 and $2\times10^{17}$ cm$^{-3}$, respectively.

- **3D schematic**



- **Top view**



$L_p$ = 0.7 μm / $L_n$ = 1.1 μm $N_n$ = 2×10$^{17}$ cm$^{-3}$
$W$ = 1.0 μm $\quad\quad\quad\quad N_p$ = 1×10$^{18}$ cm$^{-3}$

Fig. 2.1. 3-D schematic and top views of the PF neuron device.

Fig. 2.2. (a) Top SEM view. (b) Cross-sectional TEM images cut along the solid

line in (a) and (c) its magnified views.

## 2.2 Device fabrication

The PF neuron device is fabricated on a 6-inch SOI wafer with 9 masks and conventional CMOS process technology. The used masks are Si channel define (1$^{st}$), $p$-body implantation (2$^{nd}$), $n$-body implantation (3$^{rd}$), G1 formation (4$^{th}$), G2 formation (5$^{th}$), cathode/drain implantation (6$^{th}$), anode implantation (7$^{th}$) contact hole (8$^{th}$), and metal line formation (9$^{th}$).

The main fabrication process diagrams and detailed steps are shown in Fig. 2.3, respectively. Fig 2.3 (a) shows the schematic cross-sectional views of the key fabrication process steps, and Fig 2.3 (b) shows the process flow of the fabrication of the PF neuron device.

After cleaning process, which include sulfuric peroxide mixture (SPM), ammonium hydroxide-hydrogen peroxide mixture (APM), hydrochloric acid-hydrogen peroxide-water mixture (HPM), and diluted hydrogen fluoride (DHF), a 100-nm-thick Si active layer was patterned (1$^{st}$ mask) by a SS03A9 photoresist (PR). A 10-nm-thick sacrificial $SiO_2$ layer was deposited by a low-pressure chemical vapor deposition (LPCVD). A boron and phosphorus ion implantation was

performed for *p*-body and *n*-body doping using $2^{nd}$ and $3^{rd}$ mask, respectively. After the sacrificial $SiO_2$ layer was removed by wet etching in 100:1 DHF, a $SiO_2$ layer is thermally grown as gate oxide by dry oxidation process at 950 °C. Then, a layer of *in situ* $n^+$-doped poly-Si was deposited and patterned as G1 and G2 ($4^{th}$ mask). A layer of tunneling oxide layer/charge storage layer/blocking oxide layer ($SiO_2/Si_3N_4/SiO_2$) stack was deposited by the LPCVD process at 780 °C, after which the layer of *in situ* $n^+$-doped poly-Si was deposited as a G3. After the G3 is defined by the photolithography ($5^{th}$ mask) and etching the $n^+$-doped poly-Si using the RIE process, ion implantation by $As^+$ ions with a dose of $2 \times 10^{15}$ $cm^{-2}$ and energy of 40 keV is performed to form the cathode/drain ($6^{th}$ mask). And, ion implantation by $BF_2^+$ ions with a dose of $2 \times 10^{15}$ $cm^{-2}$ and energy of 40 keV is performed to form the anode ($7^{th}$ mask). This is followed by rapid thermal annealing (RTA) at a temperature of 1000 °C for 10 sec for activation and diffusion of implanted ions. After tetraethyl orthosilicate (TEOS) was deposited by a plasma-enhanced CVD (PECVD) process, contact holes for the G1, G2, G3, cathode, drain, and anode were formed ($8^{th}$ mask) by RIE process. Subsequently,

Ti/TiN/Aluminum (Al)/TiN electrodes were formed by sputtering process and were then patterned ($9^{th}$ mask). Then, hydrogen ($H_2$) annealing at 350 ℃ for 10 min was performed to improve the contact and interface property. The PF neuron device was fabricated with conventional CMOS on the same SOI wafer.

Most of the processes were carried out using the equipment in Inter-University Semiconductor Research Center (ISRC) located in Seoul National University (SNU), Seoul, Korea, and *in situ n$^+$*-doped poly-Si layer was deposited by using the equipment of National NanoFab Center (NNFC) located in Daejeon, Korea.

**(a)**

$n^+$ **poly-Si**

**(d)** G1 G3 G2
$n^+$ p- n- p- $n^+$
SiO₂

p- n- p-
SiO₂

$n^+$ **poly-Si**

**O/N/O**

**(b)** G1 G2
p- n- p-
SiO₂

**(c)** G1 G2
p- n- p-
SiO₂

- **6 inch SOI wafer**
- **Si patterning (active)**
- ***p-*/*n-* implantation  - (a)**
- **$n^+$ poly-Si deposition**
- **Gate1/2 (G1/G2) patterning - (b)**
- **O/N/O stack formation  - (c)**
- **$n^+$ poly-Si deposition**
- **Gate3 (G3) patterning**
- **Cathode/Anode implantation & activation  - (d)**
- **Back end process**

Fig. 2.3. (a) The cross-sectional views and (b) the process flow of the key

fabrication process steps for the PF neuron device.

## 2.3 Integrate-and-fire operation as a neuron device

The proposed PF neuron device performs the integrate-and-fire operation while simultaneously receiving excitatory and inhibitory signals. The integrate-and-fire operation is implemented by the PF mechanism in the floating body of the PF neuron device. To analyze the PF mechanism of the PF neuron device in detail, the integrate-and-fire operation is investigated by G1, G2, and G3.

First, Fig. 2.4 (a) shows a part of the PF neuron device to investigate the device operation by G1. Fig. 2.4 (b) shows the energy band diagram cut along the anode to the cathode in a PF neuron device to illustrate the PF mechanism. As $V_{G1}$ increases, an electron-injection barrier ($V_{e1}$) below G1 decreases (①) and as a result, electrons from $n^+$ cathode accumulate in the $n$-body region (②). As the injected electrons increase, a hole-injection barrier ($V_h$) in the $n$-body region decreases (③). Then, holes from the $p^+$ anode are easily injected into the $p$-body region below G1, which further decreases the height of $V_{e1}$ by the injected holes (④). The repeated positive feedback operation (①-④) in the floating $p$-/$n$-bodies enables steep switching characteristics of the PF neuron device. Since accumulating charges in

the $n$-body region is accelerated as $V_{G1}$ increases, the excitatory signals are transmitted to G1.

Second, Fig. 2.5 (a) and (b) show a top view and an energy band diagram cut along the drain from the cathode in the PF neuron device to investigate the device operation by G2. As $V_{G2}$ increases, an electron injection barrier ($V_{e2}$) to the drain decreases. Electrons accumulated in the $n$-body region can escape to the drain depending on the $V_{G2}$ as shown in Fig. 2.5 (b), which means that G2 prevents charges from accumulating in the $n$-body region. Note that in the inhibitory operation, the $n$-body region acts as the source of a MOSFET. By accumulating and discharging electrons in the $n$-body region, of the PF neuron device can implement the integrate-and-fire operation processing the excitatory and inhibitory signals simultaneously. When the PF neuron device turns off, the amount of electrons accumulated in the $n$-body region decreases over time due to recombination of electrons and holes. The decrease of electrons means a leaky integration, and is related to a retention time, which can be controlled by G3 as a long term [32], [36].

Third, the initial $V_h$ is modulated by $V_{G3}$ as shown in Fig. 3 (c), which means

that more electrons accumulated in the $n$-body region are required for PF operation. Thus, $V_{G3}$ controls the $V_{th}$ of the PF neuron device. Also, the $V_h$ can be controlled by the amount of charges in the charge storage layer ($Si_3N_4$). The $V_{th}$ of the PF neuron device, adjusted by the amount of charges in the charge storage layer, has a non-volatile memory function.

Fig. 2.4. (a) A top view of the PF neuron device with G1. (b) An energy band diagram from the cathode to the anode in the PF neuron device.

Fig. 2.5. (a) A top view of the PF neuron device with G1 and G2. (b) An energy band diagram from the cathode to the drain in the PF neuron device.

Fig. 2.6. (a) A top view of the PF neuron device with G1, G2, and G3. (b) An energy

band diagram cut along the drain from the cathode in the PF neuron device.

## 2.4 Device operation with excitatory and inhibitory signals

### 2.4.1 DC $I$-$V$ characteristics

Fig. 2.7 (a) and (b) show measured $I_A$-$V_{G1}$ and $I_D$-$V_{G1}$ curves of the PF neuron device as a parameter of $V_{G2}$, respectively, As $V_{G1}$ increases at a $V_{G2}$ of 0 V, $I_A$ rapidly increases due to the PF operation. The subthreshold swing (SS) of the PF neuron device is very steep (< 1 mV/dec) as shown in Fig. 2.7 (a). After the PF operation, $I_A$ of the PF neuron device is constant because the constant diode current flows from the anode to the cathode region at a fixed $V_A$. When $I_D$ is higher than $I_A$ by increasing $V_{G2}$ just before the PF action of electrons and holes in the floating $p$-/$n$-bodies occurs, the $V_{th}$ of the PF neuron device increases by preventing electrons from accumulating in the $n$-body region. Fig. 2.8 (a) and (b) show measured $I_A$-$V_{G2}$ and $I_D$-$V_{G2}$ curves of the PF neuron device as a parameter of $V_{G1}$, respectively. Although $V_{G2}$ increases, $I_A$ and $I_D$ are off-state at a $V_{G1}$ of 0 V. As $V_{G2}$ increases at $V_{G1}$ of 0.4 V, $I_D$ increases and $I_A$ decreases. Electrons accumulated in the $n$-body region are discharged to the drain at high $V_{G2}$ of 1 V. When the PF neuron device turns on at a $V_{G1}$ of 0.5 V, the energy band of the PF neuron device is nearly flat.

26

By increasing $V_{G2}$, electrons accumulated in the $n$-body region go to the drain. Then, the $V_h$ of the $n$-body region increases. As a result, $I_A$ rapidly decreases by suppressing the PF operation. At the same time, $I_D$ instantaneously increases a large amount of current flowing to the drain. Then, the $I_D$ decreases by the increased $V_{e1}$, which shows a negative resistance at a $V_{G1}$ of 0.5 V as shown in Fig. 2.8 (b). Since a reverse bias is applied between the $p^+$ anode and $n^+$ drain at the $V_D$ of 1.5 V and the $V_A$ of 1 V, current cannot flow from the anode to drain.

Fig. 2.7. Measured anode current versus gate1 voltage and (b) drain current versus

gate1 voltage curves of the PF neuron device as a parameter of $V_{G2}$, respectively.



Fig. 2.8. Measured anode current versus gate2 voltage and (b) drain current versus

gate2 voltage curves of the PF neuron device as a parameter of $V_{G1}$, respectively.

## 2.4.2   $V_{th}$ variation and controllability

Fig. 2.9 (a) and (b) show measured $I_A$-$V_{G1}$ curves of the PF neuron device as parameters of $V_{G3}$ and program/erase (PGM/ERS) operation, respectively. As $V_{G3}$ increases, the hole-injection barrier ($V_h$) in the $n$-body region increases, and the $V_{th}$ of the PF neuron device increases as show in Fig. 2.9 (a). When program pulses ($V_{G3}$ of 9 V and $t_{PGM}$ of 100 μs) are applied to the G3 to store electrons in the charge storage layer, the $V_{th}$ of the PF neuron device decreases due to decreasing $V_h$ the $n$-body region. Conversely, the $V_{th}$ of the PF neuron device increases by applying erase pulses ($V_{G3}$ of - 8.5 V and $t_{ERS}$ of 1 ms) to the G3 gate as shown in Fig. 2.9 (b). Fig. 2.9 (c) shows the $V_{th}$ retention of the PF neuron device in program (solid circle symbols) and erase (solid square symbols) states. The PF neuron device maintains the non-volatile of two $V_{th}$s. It is confirmed that the $V_{th}$ of the PF neuron device is well maintained in the high (erase) and low (program) states with the non-volatile function until a time of $10^4$ seconds. The variation of the hardware systems, such as conductance of synaptic devices and the $V_{th}$ of neuron circuits, can affect the target spike rates of neurons, which can degrade the performance of HNNs [37],

[38]. In order to improve the performance of the neural networks, it is important to minimize the $V_{th}$ variation by adjusting the $V_{th}$ of the neuron circuit. The method of selectively supplying different voltages to each neuron circuit to reduce the $V_{th}$ variation is very inefficient and practically impossible in large area neural networks consisting of hundreds of neuron circuits. Therefore, the $V_{th}$ of neuron circuits should have a non-volatile memory function and be selectively adjusted by program and erase pulses. This characteristics are very efficient and essential in terms of the size, power consumption and performance of neural networks.

Fig. 2.9 (d) shows the number of dies showing $V_{th}$ difference of two PF neuron devices on the same die. The $V_{th}$ difference of two PF neuron devices on the same die is less than 0.04 V. By adjusting the $V_{th}$ of the PF neuron using $V_{G3}$ control or program/erase operations in the charge storage layer, the $V_{th}$ variation of the PF neuron devices can be reduced. Fig. 2.10 (a) shows the measured $I_A$-$V_{G1}$ of 10 PF neuron devices with the $V_{th}$ variation on the same wafer. And, the $V_{th}$ variation in 10 PF neuron devices can be tuned by applying program and erase pulses to G3 as shown in Fig. 2.10 (b). Also, the threshold tuning ability enables to mimic a

homeostasis function of biological neurons, which is essential in spiking neural

networks (SNNs) based on spike-timing-dependence-plasticity (STDP) to improve

the recognition rate [37], [38].

Fig. 2.9. (a) Measured anode current versus gate1 voltage curves of the PF neuron device as a parameter of $V_{G3}$. (b) Measured anode current versus gate1 voltage curves of the PF neuron device with non-volatile memory function. (c) The threshold voltage ($V_{th}$) retention of the PF neuron device with non-volatile memory function. (d) The number of dies according to the $V_{th}$ difference ($\Delta V_{th}$) between two PF neuron devices on the same die.

32

Fig. 2.10. (a) Measured anode current versus gate1 voltage and (b) tuned anode current versus gate1 voltage of 10 PF neuron devices on the same wafer.

### 2.4.3 Transient characteristics

Fig. 2.11 (a) and (b) show measured step pulse and pulse transients of the PF neuron device as a parameter of $V_{G1}$, respectively. The amount of electrons that accumulate in the $n$-body region increases over time as $V_{G1}$ increases. As a result, the turn-on time ($t_{on}$) of the PF neuron device is shorter as shown in Fig. 2.11 (a). In pulse transients of the PF neuron device, width ($t_{pulse}$), rise time ($t_{rise}$), and fall time ($t_{fall}$) of pules applied to G1 of the PF neuron device are 1 μs, 500 ns, and 500 ns, respectively. The PF neuron device turns on at fewer pulses as $V_{G1}$ increases. Fig. 2.12 (a) and (b) show measured $I$-$t$ plot of the PF neuron device as parameters of $V_{G2}$ and $V_{G3}$, respectively. As $V_{G2}$ increases, the amount of electrons that discharge from the $n$-body region to the drain increases over time. So the $t_{on}$ the PF neuron device is longer at high $V_{G2}$ as shown in Fig. 2.12 (a). Increasing $V_{G3}$ positively deepens the potential well, $V_h$ of the n-body region in the PF neuron device. Since more electrons should be accumulated in the $n$-body region for fire, the $t_{on}$ of the PF neuron device is longer as shown in Fig. 2.12 (b). Fig. 2.13 (a) and (b) show the $t_{on}$s of the PF neuron device for integrate-and-fire function as

parameters of $V_{G1}$ and $V_{G2}$, respectively. Though high $V_G$ is applied to the PF neuron device, a turn-on time delay ($t_{delay}$) of the PF neuron device happens in PF operation [39]. The $t_{on}$ is considered with the $t_{delay}$ of 500 ns. The $t_{on}$ become exponentially short as $V_{G1}$ increases because the amount of electrons accumulated in the $n$-body region increases exponentially. Conversely, the $t_{on}$ to long exponentially as $V_{G2}$ increases because the amount of electrons accumulated in the $n$-body region decreases exponentially.

Fig. 2.11. Measured anode current versus time plots of the PF neuron device for IF function as a parameter of $V_{G1}$. (a) step and (b) pulse transient. Here, $t_{pulse}$, $t_{rise}$ and $t_{fall}$ are 1 μs , 500 ns and 500 ns, respectively.



Fig. 2.12. Measured anode current versus time plots of the PF neuron device with inhibitory signals and threshold tuning ability as parameters of (a) $V_{G2}$ and (b) $V_{G3}$, respectively.

Fig. 2.13. Measured turn-on time ($t_{on}$) of the PF neuron device for IF function as parameters of (a) $V_{G1}$ and (b) $V_{G2}$, respectively.

## 2.5 Neuron circuit

### 2.5.1 Neuron circuit based on PF neuron device

Fig. 2.14 shows a neuron circuit that consists of the proposed PF neuron device, one invertor, and $p$/$n$MOSFETs. The PF neuron device is represented by two $n$MOSFET ($M_{Exc}$ and $M_{Inh}$) and one diode. Excitatory and inhibitory signals ($I_{Exc}$ and $I_{Inh}$) from synaptic arrays are reflected in G1 and G2, respectively. When $M_R$ and $M_{Exc}$ have the same $V_{th}$, $I_{Exc}$ is linear with a current flowing through the $M_{Exc}$ because operation of the $M_{Exc}$ as a current mirror to $I_{Exc}$. Just after the PF neuron device turns on, $V_{out}$ become high state by the invertor, and $V_A$ becomes 0 V by $M_1$. Electrons accumulated in the $n$-body region are discharged at the $V_A$ of 0 V, and the PF neuron device resets. Then, $V_A$ becomes high state by the $M_2$ at $V_{G1}$ of 0 V. Also, a current flowing through the $M_{Inh}$ to the drain is determined by $I_{Inh}$. When $I_{Inh}$ is reflected in G2, electrons accumulated in the $n$-body region are drained by $M_{Inh}$. Fig. 2.15 (a) and (b) show simulated $V$-$t$ plots of the PF neuron circuit as parameters of $V_{G1}$ and $V_{G2}$. The PF neuron circuit was simulated using a mixed-mode option in a TCAD simulator (Sentaurus) of Synopsys. By using the TCAD simulation, IF and

reset operations in the PF neuron circuit are confirmed at $V_D$ of 1.5 V and $V_{DD}$ of 1.1 V, respectively. By increasing the amplitude of $V_{G1}$ pulses, spike rate of the PF neuron circuit increases at a fixed $V_{G2}$ of 0.5 V as shown in Fig. 2.15 (a). By decreasing $V_{G2}$, fire rate of the PF neuron circuit increases at a fixed amplitude (0.48 V) of $V_{G1}$ pulses.

Fig. 2.16 shows a top view of the fabricated neuron circuit that consists of the proposed PF neuron device, $p/n$MOSFETs, and one invertor. Due to steep switching characteristics of the PF neuron device, $V_{out}$ of the PF neuron circuit abruptly changes to high and low states as $V_{G1}$ increases and decreases as shown in Fig. 2.17. Fig. 2.18 shows the measured pulse transients of the PF neuron circuit as a parameter of $V_{G1}$. As $V_{G1}$ increases, $t_{on}$ of the PF neuron circuit is faster. When $V_{out}$ become the high state, M1 turns on, and $V_A$ of the PF neuron device is 0 V. When $V_{G1}$ is applied to G1 at $V_A$ of 0 V, the accumulated electrons and holes in the $p$-/$n$-body regions are discharged. It is demonstrated that the proposed PF neuron circuit can implement the integrate-and-fire function and the reset operation as shown in Fig. 2. 18 (b).

Fig. 2.14. A neuron circuit composed of the proposed PF neuron device, $p$-/$n$-FETs ($M_1$ and $M_2$), and an output invertor. Input pulses from excitatory and inhibitory synapses are applied to G1 and G2, respectively.

Fig. 2.15. Simulated output voltage versus time plot of the PF neuron circuit for IF function as parameters of (a) $V_{G1}$ and (b) $V_{G2}$. Here, $V_D$ and $V_{DD}$ are 1.5 V and 1.1 V, respectively.

Fig. 2.16. Top view of a fabricated neuron circuit consisting of the proposed PF neuron device, $p/n$MOSFETs ($M_1$ and $M_2$), and one invertor.

Fig. 2.17. Measured output voltage versus gate1 voltage plot of the PF neuron

circuit with steep switching characteristics.

Fig. 2.18. Measured output voltage versus time plot of the PF neuron circuit for IF

function as parameters of (a) $V_{G1}$ and (b) $V_{G2}$. Here, $V_D$ and $V_{DD}$ are 1.1 V. (c) IF

function and reset operation of the PF neuron circuit.

## 2.5.2 Energy consumption

Fig. 2.19 shows circuit diagrams and simulated *I-t* plots for the comparison of

energy consumption (J/spike) between the PF neuron circuit and a conventional

neuron circuit, respectively [22]. The conventional neuron circuit consists of

capacitors (membrane and refractory) and minimum number of FETs to mimic

integrate-and-fire function of neurons. $M_1$ is used to fully discharge the input node

of the invertor, and $M_{\text{RESET}}$ is used to reset the membrane potential. In conventional

neuron circuit, the current transmitted from synaptic devices flows into the

membrane capacitor ($C_{\text{mem}}$), and a membrane potential ($V_{\text{mem}}$) changes. When the

$V_{\text{mem}}$ exceeds the $V_{\text{th}}$ of the neuron circuit, the neuron circuit turns on and generates

an output spike. Until the $V_{\text{mem}}$ exceeds the $V_{\text{th}}$ of the neuron circuit during the

integration function, a leakage current flows in the neuron circuit by the relatively

slower SS (> 60 mV/dec) of conventional CMOS, which affects the total energy

consumption of neuromorphic systems as shown in Fig. 2.19 (a). On the other hand,

low leakage current (< 1 nA) flows in the PF neuron circuit before the PF neuron

device with steep switching characteristics (SS < 1 mv/dec) turns on. Thus, the

current in the PF neuron circuit flows only whenever the state of $V_{\text{out}}$ changes by

the fire and reset operations as shown in Fig 2.19 (b). Energy consumption of the

proposed PF neuron circuit is about 0.62 pJ/spike, which is reduced by about 10

times compared to the conventional neuron circuit (~6.14 pJ/spike).

Fig. 2.19. Circuit diagrams and simulated total current versus time plots for the comparison of energy consumption per spike between (a) a conventional neuron circuit and (b) PF neuron circuit.

# Chapter 3
# Synaptic device

## 3.1   Device structure

Fig. 3.1 and 3.2 show a cross-section view and an SEM image of a synaptic device, respectively. The structure of the synaptic device is a dual-gate FET that consists of a gate1 (G1) with a gate oxide layer (SiO$_2$) and a gate2 (G2) with a gate insulator stack (SiO$_2$/Si$_3$N$_4$/SiO$_2$). The thicknesses of the gate oxide and SiO$_2$/Si$_3$N$_4$/SiO$_2$ stack are 10 nm and 3/6/9 nm, respectively. The length of G1, G2, and the thickness of Si body ($T_{Si}$) are 0.6 μm, 0.6 μm, and 100 nm, respectively. Doping concentration of a channel region are $1 \times 10^{18}$ cm$^{-3}$. In the dual-gate FET, G1 is the role of a switch to turn on the synaptic device, and an input voltage is applied to G1. On the other hand, G2 of the dual-gate FET controls conductance of the synaptic device because channel resistance of the dual-gate FET changes by $V_{G2}$ and amount of charges in the charge storage layer. Thus, in the dual-gate FET, G1 determines on and off states of the synaptic device, and G2 controls a synaptic

weight as conductance changes. The dual-gate FET implements a wide range of conductance changes, and can completely turn off the synaptic device when the input voltage is not applied to G1. Also, the dual-get FET is compatible with conventional CMOS technology, and can be fabricated with conventional CMOS and mentioned PF neuron device on the same SOI wafer in Chapter 2.

**$T_{ONO}$ = 3/6/9 nm    $T_{Ox}$ = 10 nm**

**Source**    G2    G1    **Drain**

$n^+$    $p^-$    $n^+$

$p^-$ : $1 \times 10^{18}$ cm$^{-3}$    $L_{G1}$ : 1 um, $L_{G2}$ : 0.7 um

**Source**    G2    G1    **Drain**

$n^+$    $p^-$    $n^+$

$p^-$ : $1 \times 10^{18}$ cm$^{-3}$

**Conductance is changed by gate2 (bias or PGM/ERS).**

Fig. 3.1. (a) A pulse-width modulation (PWM) circuit designed with the assumption

of *n*-type GSDs. (b) Voltage pulses ($V_{in,s}$) with different pulse widths modulated by

the PWM circuit. (c) Synaptic currents corresponding the modulated voltage pulses.

Fig. 3.2. a SEM image of the fabricated dual-gate FETs

## 3.2    Device fabrication

The dual-gate FET is fabricated on a 6-inch SOI wafer with 7 masks and conventional CMOS process technology. The used masks are Si channel define (1st), channel implantation (2nd), G1 formation (3rd), G2 formation (4th), source/drain formation (5th), contact hole (6th), and metal line formation (7th).

The main fabrication process diagrams and detailed steps are shown in the Fig. 3.2 (a) and (b), respectively. Fig. 3.2 (a) shows the schematic cross-sectional views of the key fabrication process steps, and Fig. 3.2 (b) shows the process flow of the fabrication of the dual-gate FET.

After cleaning process, which include sulfuric peroxide mixture (SPM), ammonium hydroxide-hydrogen peroxide mixture (APM), hydrochloric acid-hydrogen peroxide-water mixture (HPM), and diluted hydrogen fluoride (DHF), a 100-nm-thick Si active layer was patterned (1st mask) by a SS03A9 photoresist (PR). A 10-nm-thick sacrificial $SiO_2$ layer was deposited by a low-pressure chemical vapor deposition (LPCVD). A boron ion implantation was performed for channel doping using 2nd mask. After the sacrificial $SiO_2$ layer was removed by wet etching

in 100:1 DHF, a $SiO_2$ layer is thermally grown as gate oxide by dry oxidation process at 950 °C. Then, a layer of *in situ* $n^+$-doped poly-Si was deposited and patterned as G1 (3rd mask). A layer of tunneling oxide/charge storage layer/blocking oxide ($SiO_2/Si_3N_4/SiO_2$) stack was deposited by a LPCVD process at 780 °C, after which a layer of $n^+$-doped poly-Si was deposited as a G2. After the G2 is defined by the photolithography (4th mask) and etching the $n^+$-doped poly-Si using the RIE process, ion implantation by $As^+$ ions with a dose of $2 \times 10^{15}$ cm$^{-2}$ and energy of 40 keV is performed to form the source/drain. This is followed by rapid thermal annealing (RTA) at a temperature of 1000 °C for 10 sec for activation and diffusion of implanted ions. After tetraethyl orthosilicate (TEOS) was deposited by a plasma-enhanced CVD (PECVD) process, contact holes for the G1, G2, sources, drains were formed (6th mask) by RIE process. Subsequently, Ti/TiN/Aluminum (Al)/TiN electrodes were formed by sputtering process and were then patterned (7th mask). Then, hydrogen ($H_2$) annealing at 350 °C for 10 min was performed to improve the contact and interface property. The dual-gate FET was fabricated with conventional CMOS and mentioned PF neuron device on the same SOI wafer.

Most of the processes were carried out using the equipment in Inter-University Semiconductor Research Center (ISRC) located in Seoul National University (SNU), Seoul, Korea, and *in situ* $n^+$-doped poly-Si layer was deposited by using the equipment of National NanoFab Center (NNFC) located in Daejeon, Korea.

**(a)**

$n^+$ **poly-Si**

**(b)**

$n^+$ **poly-Si**

**(c)**

O/N/O

**(d)**

- **6 inch SOI wafer**
- **Si patterning (active)**
- *p*-channel implantation  - (a)
- $n^+$ poly-Si deposition
- **Gate1 (G1) patterning - (b)**
- **O/N/O stack formation**
- $n^+$ poly-Si deposition
- **Gate2 (G2) patterning - (c)**
- **Source/Drain implantation & activation  - (d)**
- **Back end process**

Fig. 3.3. (a) The schematic cross-sectional views of the key fabrication process steps. (b) The process flow of the fabrication of the dual-gate FET.

## 3.3   Device characteristics

The direct current (DC) *I-V* characteristics of the fabricated dual-gate FETs were measured by using a semiconductor parameter analyzer (B1500A, Keysight) and cascade probe station.

Fig. 3.4 (a) shows a schematic view of the device operation by each G1 and G2 in the dual-gate FET. Fig. 3.4 (b) shows the measured $I_D$-$V_{G1}$ curves of the dual-gate FET as a parameter of $V_{G2}$ at $V_D$ of 1 V. Here, the width of gates ($W_G$) and the length of G1 and G2 of the dual-gate FET are 1 μm, 1 μm and 0.6 μm, respectively. Since the dual-gate FET turns on by G1 at a fixed $V_{G2}$, a threshold voltage ($V_{th}$) of $I_D$-$V_{G1}$ curves is constant as shown in Fig 3.4 (b). When the channel region under G1 is in a strong inversion state at a high $V_{G1}$ above a threshold voltage ($V_{th}$), the channel resistance under G2 is the most dominant in the dual-gate FET. Even if $V_{G1}$ increases at the fixed $V_{G2}$, $I_D$ is constant. As $V_{G2}$ increases at a high $V_{G1}$ of 2 V, $I_D$ increases due to the decrease of the channel resistance under G2 as shown in Fig 3.4 (b). Thus, $V_{th}$ of the transfer curves is determined by G1, and on-current ($I_D$) of the device is controlled by G2 as shown in Fig. 3.4 (a). Fig. 3.5 shows the measured

$I_D$-$V_D$ curves as a parameter of $V_{G2}$ at the fixed $V_{G1}$ of 2 V. When $V_D$ is over than

0.5 V and $V_{G2}$ is less than 2 V, $I_D$ saturates. As $V_{G2}$ increases at the fixed $V_{G1}$ of 2 V,

$I_{D,sat}$ increases. It is very important that the current of the device saturates with the

drain and gate voltages in the synaptic array. The gate and drain voltages applied to

each device in the synaptic array can be changed due to problems such as noise at

the input terminals and a voltage drop between devices by the resistance of the metal

lines in the synaptic array. Even if synaptic weights are stored as accurate values in

each synaptic device through software learning, the problems cause errors in the

sum of synaptic weights, which degrades the target recognition rate of neural

networks.

Fig. 3.6 (a) and (b) show the measured $I_D$-$V_{G1}$ curves of the dual-gate FET with

$L_{G2}$ of 0.6 µm and 1 µm, respectively. $I_D$ of the dual-gate FET with $L_{G2}$ of 0.6 µm is

higher than that of the dual-gate FET with $L_{G2}$ of 1 µm at the same $V_{G2}$. Fig. 3.7

shows the measured $I_D$-$V_D$ curves of the dual-gate FET according to the position of

G2 that controls the conductance change. Fig. 3.7 (a) and (b) are the measurement

results of the dual-gate FET where G2 is adjacent to the drain or source region,

respectively. In both cases, $I_D$ saturates as $V_D$ increases at the fixed $V_{G1}$ and $V_{G2}$, and

$I_{D,sat}$ increase as $V_{G2}$ increases at $V_{G1}$ of 2 V.

Fig. 3.4. (a) Device characteristics of the dual-gate FET using G1 (switch in read operation) and G2 (conductance change). (b) Measured $I_D$-$V_{G1}$ curve of the dual-gate FET as a parameter of $V_{G2}$.

Fig. 3.5. Measured $I_D$-$V_D$ curves of the dual-gate FET as a parameter of $V_{G2}$ at $V_{G1}$ of 2 V.

Fig. 3.6. Measured $I_D$-$V_{G1}$ curve of the dual-gate FET with $L_{G2}$ of (a) 0.6 μm and (b)

1 μm as a parameter of $V_{G2}$.

Fig. 3.7. Measured $I_D$-$V_D$ curves of the dual-gate FET according to the position of

G2 that is adjacent to the drain (a) and source (b) region.

## 3.4 Device operation as a synaptic device

To represent the learned weights of excitatory and inhibitory synapses, the synaptic devices with a memory function are needed. In the hardware-based neuromorphic systems, the synaptic memory function requires non-volatile memory characteristics and high reproducibility in repetitive potentiation and depression behavior. Also, synaptic devices with low power consumption are required during synaptic weight update. In the dual-gate FET, the $SiO_2/Si_3N_4/SiO_2$ gate stack with the charge storage layer is used as a gate dielectric of G2 to change the conductance of the synaptic device. The thickness of tunneling oxide layer, charge storage layer, and blocking oxide layer is 3 nm, 6 nm, and 9 nm, respectively. The program and erase operations in the charge storage layer of the dual-gate FET are performed through direct tunneling and Fowler-Nordheim (FN) tunneling mechanism. After the program and erase operations, the channel resistance is changed by amount of charges stored in the charge storage layer. Fig. 3.8 (a) and (b) show the measured $I_D$-$V_{G1}$ curves of the dual-gate FET as the number of erase and program pules, respectively. Here, the width ($t_{ERS}$) and amplitude ($V_{ERS}$) of

erase pulses are 1 ms and – 8.5 V, respectively. The width ($t_{PGM}$) and amplitude ($V_{PGM}$) of program pulses are 100 μs and 5.5 V, respectively. As number of erase pulses applied to G2 increases, $I_D$ increases at $V_{G2}$ of 0 V as shown in Fig. 3.8 (a). Otherwise, as number of program pulses applied to G2 increases, $I_D$ decreases at $V_{G2}$ of 0 V as shown in Fig. 3.8 (b). The $V_{th}$ of the $I_D$-$V_{G1}$ curves are always constant after program and erase operations in the dual-gate FET. Thus, the charges stored in the charge storage layer by the program and the erase pulses applied to G2 only affect the channel resistance under G2. To evaluate the dual-gate FET as the synaptic device, the long-term potentiation (LTP) and depression (LTD) characteristics of a dual-gate FET are analyzed as shown in Fig. 3.9. According to number of program and erase pulses, $I_D$ of the dual-gate FET is measured at $V_{G1}$ of 3 V, $V_{G2}$ of 0 V, and $V_D$ of 1 V. Here, erase pulses ($V_{ERS}$ of -9.5 V and $t_{ERS}$ of 1 ms) and program pulses ($V_{PGM}$ of 6.5 V and $t_{PGM}$ of 100 μs) are applied to G2. Here, $V_{ERS}$ and $t_{ERS}$ of erase pulses are – 9.5 V and 1 ms, respectively. $V_{PGM}$ and $t_{PGM}$ of program pulses are 6.5 V and 100 μs, respectively. As the number of erase pulses increases, the conductance of LTP increases linearly. On the other hand, as the

number of program pulses increases, the conductance of LTD decreases abruptly.

Considering that the dual-gate FET is in subthreshold region by G2, a synaptic

current ($I_{\text{syanpse}}$) is

$$I_{synapse} = I_{G2,sub}. \tag{1}$$

Subthreshold current ($I_{\text{sub}}$) of a conventional MOSFET is

$$I_{sub} = \mu_{eff} \frac{W}{L} (m-1)(\frac{kT}{q})^2 e^{q(V_{gs}-V_t)/mkT}(1 - e^{-qV_{ds}/kT}). \tag{2}$$

The relationship between $I_{\text{syanpse}}$ and $V_{\text{G2}}$ can be obtained by using equation (1) and

(2),

$$I_{synapse} = I_{G2,sub} \propto \exp(|V_{G2}|). \tag{3}$$

The effective $V_{\text{G2}}$ can be regarded as trap charge ($Q_{\text{trap}}$) with relation of

$$Q_{trap} = C_{Oxide} \times V_{G2,eff}. \tag{4}$$

By FN tunneling mechanism during the erase operation in the dual-gate FET, $Q_{\text{trap}}$

trapped in the charge storage layer is

$$Q_{trap} \propto \ln(t_{potentiation}), \tag{5}$$

where $t_{\text{potentiation}}$ is the time of the potentiation operation which is the product of the

pulse width ($t_{\text{pulse}}$) and the number of pulses.

The relationship between $V_{G2,eff}$ and $t_{potentiation}$ can be obtained by using equation (4) and (5),

$$V_{G2,eff} \propto Q_{trap} \propto \ln(t_{pulse} \times number\ of\ pulses). \quad (6)$$

The relationship between $I_{synapse}$ and the number of pulses can be obtained by using equation (3) and (6),

$$I_{synapse} \propto \exp(|V_{G2,eff}|) \propto \exp(\ln(number\ of\ pulses))$$

$$\propto number\ of\ pulses. \qquad (7)$$

Thus, it is demonstrated that the correlation between the conductance change and the number of erase pulses is linear in the subthreshold region of the dual-gate FET.

Fig. 3.10 shows the measured $I_D$-$V_{G2}$ curve of the dual-gate FET at $V_{G1}$ of 3 V during the LTP operation. Since erase pulses are applied to G2, $V_{th}$ of $I_D$-$V_{G2}$ curve decreases as number of erase pulses increases. Fig. 3.11 (a) and (b) show the LTP and LTD characteristics of the dual-gate FET at $V_{G2}$ of 0.2 V and 1 V, respectively. When the conductance of the dual-gate FET are measured at $V_{G2}$ of 1 V, which is higher than $V_{th}$, the LTP characteristic of the dual-gate FET is not satisfied by the above equations (1)-(7). The relationship between the conductance change and

66

number of erase pulses is non-linear at $V_{G2}$ of 1 V as shown in Fig 3.11 (b).

Fig. 3.8. Measured $I_D$-$V_{G1}$ curves of the dual-gate FET as number of erase (a) and

program pules (b), respectively.

Fig. 3.9. The long-term potentiation (LTP) and depression (LTD) characteristics of

the dual-gate FET. Here, $I_D$ of the dual-gate FET is measured at $V_{G1}$ of 3 V, $V_{G2}$ of

0.2 V, and $V_D$ of 1 V.

Fig. 3.10. Measured $I_D$-$V_{G2}$ curve of the dual-gate FET at $V_{G1}$ of 3 V during the LTP

operation.

Fig. 3.11. The LTP and LTD characteristics of the dual-gate FET at $V_{G2}$ of (a) 0.2 V

and (b) 1 V, respectively.

# Chapter 4

# Synaptic array

## 4.1 Synaptic array based on a dual-gate FET

We have successfully demonstrated a dual-gate FET as the synaptic device in Chapter 3. The synaptic current is modulated by $V_{G2}$ and amount of charge in the charge storage layer. Regardless of $V_{G2}$, the dual-gate FET turns off when $V_{G1}$ is 0 V. Fig. 4.1 shows a NOR type synaptic array based on the proposed dual-gate FETs. The G1, G2, source and drain regions of the dual gate FET in Chapter 3 are referred to as word-line (WL), control-gate (CG), source-line (SL), and bit-line (BL) in the synaptic array, respectively. The input signals transmitted from pre-neurons are applied to BLs, and synaptic currents of the dual-gate FETs are summed to SLs. In a conventional flash NOR array, the program operation is performed by a hot carrier injection mechanism and the erase operation is performed by a FN tunneling mechanism. Programming with the hot carrier injection in the charge trap layer accelerates the degradation of the tunneling oxide layer. In addition, high current

flows during the program operation using the hot carrier injection, resulting in higher power consumption compared to the program operation using the FN tunneling. However, program and erase operations can be implemented with the FN tunneling mechanism in the NOR type array based on the dual-gate FETs.

Fig. 4.2 shows bias conditions for selective program operation in 2×2 NOR type synaptic array based on the dual-gate FETs. To program only M1 in the array, a $V_{PGM}$ of 8 V and $V_{inhibition}$ of 4 V are applied to CG1 and SL2, respectively. Since only M1 meets the bias condition for program operation of the dual-gate FETs, that condition enables selective programming in the array. Fig. 4.3 shows the measured $I_{Bit}$-$V_{WL}$ curves of M1, M2, M3, and M4 at $V_{CGS}$ of 2.5 V and $V_{bit}$ of 1 V when selective program bias conditions ($V_{CG1}$ of 8 V and $V_{SL2}$ of 4 V) are applied to the NOR type synaptic array. A width of pulse ($t_{PGM}$) is 100 μs. Under the program bias condition, the $I_{Bit}$ of M1 decreases by 35 nA, and M2 and M4 have little change (~ 1 %). The $I_{Bit}$ of M3 decreases by about 1 nA (~ 3%). This reduction is a reasonable value compared to the overall current change and demonstrates that M2, M3 and M4 are successfully inhibited. Fig. 4.4 shows bias conditions for selective erase

operation in 2×2 NOR type synaptic array based on the dual-gate FETs. To erase

only M1 in the array, a $V_{CG1}$ of $-6$ V and $V_{SL1}$ of 5 V are applied to CG1 and SL1,

respectively. Since only M1 meets the bias condition for erase operation of the dual-

gate FETs, that condition enables selective erasing in the array. Fig. 4.5 shows the

measured $I_{Bit}$-$V_{WL}$ curves of M1, M2, M3, and M4 at $V_{CGS}$ of 2.5 V and $V_{bit}$ of 1 V

when selective erase bias conditions ($V_{CG1}$ of $-6$ V and $V_{SL1}$ of 5 V) are applied to

the NOR type synaptic array. A width of pulse ($t_{ERS}$) is 10 ms. Under the erase bias

condition, the $I_{Bit}$ of M1 increases by 33 nA. The $I_{Bit}$s M2, M3 and M4 have little

change ($\sim 1$ %). The changes are reasonable values compared to the overall current

change and demonstrates that M2, M3 and M4 are successfully inhibited.

Fig. 4.1. Schematic and SEM top views of NOR type synaptic array based on the

dual-gate FETs.

Fig. 4.2. The NOR type synaptic array based on the dual-gate FETs and bias conditions for selective program operation of M1 and inhibition of M2, M3 and M4.



Fig. 4.3. Measured $I_{Bit}$-$V_{WL}$ curve of M1, M2, M3 and M4 for selective program operation in the NOR type synaptic array based on the dual-gate FET.
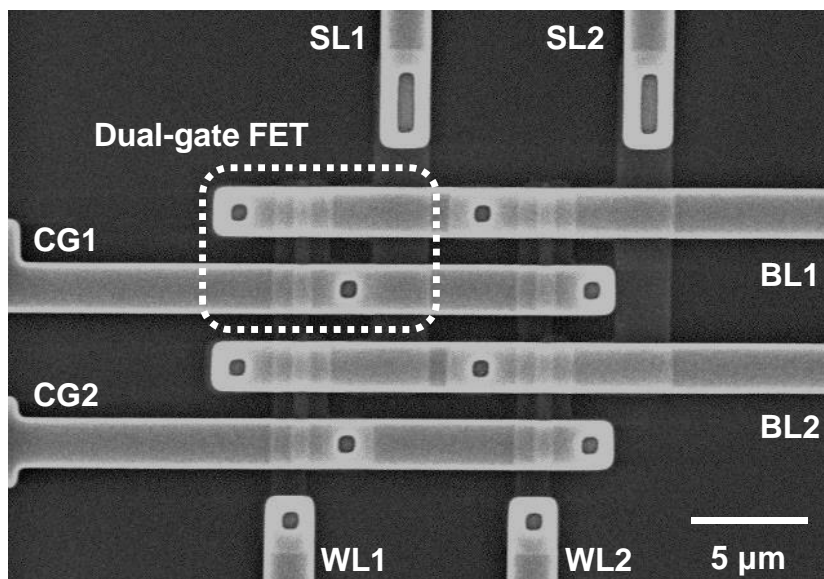
Fig. 4.4. The NOR type synaptic array based on the dual-gate FETs and bias conditions for selective erase operation of M1 and inhibition of M2, M3 and M4.



Fig. 4.5. Measured $I_{Bit}$-$V_{WL}$ curve of M1, M2, M3 and M4 for selective erase operation in the NOR type synaptic array based on the dual-gate FET.

## 4.2 VMM using dual-gate FET array

When synapse devices are configured as a crossbar array, the sneak path and the IR drop problems are main challenges. The unintended current by the sneak path problem can result in inaccurate VMM computation. In addition, the IR drop along metal wires in crossbar array can distort the voltage across the synapse device [40]-[42]. The synapse device at the far-end of the array is most affected by the IR drop, so size of the array can be limited. We assume for simplicity that the resistance of synapse devices are the same, and calculate the voltage across the synapse device at the far-end of the $N{\times}N$ crossbar array ($V_{NN}$) (Fig. 4.6 (a)). Even when the input voltage ($V_{input}$) is applied to the synapse device, the $V_{NN}$ can be changed from $V_{input}$ to a reduced voltage due to the IR drop along metal wires in crossbar array. The resistance of metal wire between adjacent synapse devices is assumed to be 2.5 $\Omega$ [41]. As shown in Fig. 4.6 (b), if the resistance of the synapse device is 5 k$\Omega$ (low resistance state in [41]) and the array size is 64$\times$64, $V_{NN}$ becomes 31% of $V_{input}$ [42]. Therefore, the resistance of the synapse device even in the low resistance state should be sufficiently large considering the array size. However, the synapse array

78

based on the dual-gate FETs is free from these problems. Since the synaptic current is saturated with respect to the input voltage, the distortion of input voltage caused by the IR drop in metal wires does not affect the synaptic current of the dual-gate FETs. These characteristics are very important in synapse array.

Fig. 4.7 shows the SEM image and the layout of the 8×4 NOR type synaptic array. The 8×4 NOR type synaptic array consists of 4 WLs, 4 SLs, 8 BLs, and 8 CGs. The total number of the dual-gate FETs in the array is 24 (8×4). Except for the CGs for program and erase operation, it is similar to a conventional NOR flash array. Fig. 4.8 shows the vector-by-matrix multiplication (VMM) operation performed by the 8×4 NOR type synaptic array. Calculating the VMM in software is a big burden. However, it can be easily implemented by summing the synaptic currents in a synaptic array. The $I_{SL1}$-$V_{WL1}$ curves of eight individual dual-gate FETs are measured at $V_{CGS}$ of 2.5 V and the selected $V_{Bit}$ of 1 V as shown in Fig. 4.8 (a). The $I_{total}$-$V_{WL1}$ curve is measured at all $V_{CGS}$ of 2.5 V and $V_{Bit}$s of 1 V as shown in Fig. 4.8 (b). The sum of current (3.63 μA) of eight individual dual-gate FETs ($I_1$ + $I_2$ + ⋯ $I_8$) is almost the same (~ 0.87 %) as the $I_{total}$ (3.6 μA) of eight dual-gate

FETs. The variation, represented by the standard deviation divided the mean ($\sigma/\mu$) of the synaptic currents in Fig. 4.8 (a), is 0.295. Fig. 4.9 shows the measured $I_{\text{Bit}}$-$V_{\text{WL}}$ curves of 8 dual-gate FETs as the quantized conductance levels (220, 150, 80, and 10 nA) by using program and erase operation in the charge storage layer, respectively. The quantized conductance levels (220, 150, 80, and 10 nA) of 8 dual-gate FETs were measured at $V_{\text{CG}}$ of 2.5 V, $V_{\text{WL}}$ of 2.5 V, and $V_{\text{Bit}}$ of 1 V. The mean values of the synaptic current in 8 devices are 213, 153, 85.6, 10.6 nA, respectively, and the standard deviations are 4.89, 1.77, 1.32, 0.34 nA, respectively. The larger the synaptic current, the larger the standard deviation. The variation, represented by the standard deviation divided the mean ($\sigma/\mu$) of the synaptic currents, at each weight state are obtained as 0.023, 0.011, 0.015 and 0.032, respectively. The dual-gate FET performs well as a synapse device with the help of saturation characteristics. In addition, the power consumption can be low due to controllability of CGs and the device reliability is also good because it is a Si-based device.

Fig. 4.6. (a) A crossbar array showing parasitic resistance along metal wires. (b)

Voltage across the synapse device at the far-end of the array ($V_{NN}$)

Fig. 4.7. The synapse array based on the GSDs.

Fig. 4.8. Vector-by-matrix multiplication (VMM) in the 8×4 NOR type synaptic array. The $I_{total}$-$V_{WL1}$ curve is measured at all $V_{CG}$s of 2.5 V and $V_{Bit}$s of 1 V.

Fig. 4.9. $I_{Bit}$-$V_{WL}$ curves of 8 dual-gate FETs as the quantized conductance levels

(220, 150, 80, and 10 nA), respectively.

Fig. 4.10. Quantized conductance levels (220, 150, 80, and 10 nA) of 8 dual-gate

FETs in the 8×4 NOR type synaptic array with program and erase operations.

# Chapter 5
# Conclusion

In this work, we have investigated a neuron device based on positive feedback (PF) mechanism and a synaptic device based on a dual-gate FET for hardware-based neural networks (HNNs) with high density and low power consumption. The PF neuron device has a structure that efficiently merges a gated PNPN-junction diode and one FET. The PF neuron device implements the integrate-and-fire function by charging electrons in a floating $n$-body region, which can replace a membrane capacitor of neuron circuits. Owing to steep switching characteristics (SS < 1 mV/dec) of the PNPN-junction, the PF neuron device enables the integrate-and-fire operation with low energy consumption. And, the accumulated electrons are discharged to the drain by G2 of the PF neuron device, the PF neuron device can simultaneously process excitatory and inhibitory signals using G1 and G2. A neuron circuit based on the PF neuron device consumes 0.62 pJ of energy per spike, which is about 10 times less than that of a conventional neuron circuit to implement

the integrate-and-fire function. Moreover, the threshold voltage ($V_{th}$) of the PF

neuron device is modulated by adjusting the potential barrier of the $n$-body region

of the PF device ($V_{G3}$ control or modulating the amount of charge stored in the

charge storage layer by the program/erase operation), which controls the firing rate

of neurons. The threshold tuning ability of the PF neuron device can implement the

homeostasis function of biological neurons and compensate for the $V_{th}$ variation of

neurons in HNNs. The dual-gate FET has independent two gates (G1 with a $SiO_2$

layer and G2 with the charge storage layer). Here, G1 turns on and off the synaptic

device in the synapse array, and G2 controls the conductance of the dual-gate FET

for synaptic weights. The range of conductance change of the dual-gate FET is very

wide (100 pA ~ 1 μA). Also, the conductance response of the dual-gate FET shows

the linear potentiation characteristics. In a conventional flash NOR array, the

program operation is performed by a hot carrier injection mechanism, and the erase

operation is performed by an FN tunneling mechanism. However, in the NOR type

array based on the dual-gate FETs, program and erase operations can be

implemented with the FN tunneling mechanism, resulting in low power

consumption during program operation for the synaptic weight update. In an 8×4 NOR type synapse array, the sum of current (3.63 μA) of eight individual dual-gate FETs is almost the same (~ 0.87 %) as the $I_{total}$ (3.6 μA) of eight dual-gate FETs. The variations (σ/μ) of the synaptic currents quantized by four weight states are obtained as 0.023, 0.011, 0.015, and 0.032, respectively. The PF neuron device and the dual-gate synaptic device can be promising solutions for high-density and low-power neuromorphic systems.

# Bibliography

[1] C.S. Poon, K. Zhou, "Neuromorphic silicon neurons and large-scale neural networks: challenges and opportunities," *Frontiers in neuroscience*, 5:108, 2011.

[2] D. Kuzum, S. Yu, and H.-S. P. Wong, "Synaptic electronics: materials, devices and applications", *Nanotechnology*, vol. 24, no. 38, p. 382001, Sep. 2013.

[3] Y. LeCun, "Gradient-based learning applied to document recognition", *Proceeding IEEE*, vol. 86, iss. 11, pp. 2278 – 2324, Nov. 1998.

[4] D. Ciregan, U. Meier, and J. Schmidhuber, "Multi-column deep neural networks for image classification," *in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, pp. 3642–3649, Jun. 2012.

[5] D. Kuzum, Rakesh G. D. Jeyasingh, B. Lee, and H.-S. P. Wong, "Nanoelectronic Programmable Synapses Based on Phase Change Materials for Brain-Inspired Computing," *Nano Letter*, pp. 2179–2186, June 2012.

[6] G. Indiveri and S.-C. Liu, "Memory and information processing in neuromorphic systems," *Proceeding IEEE*, vol. 103, no. 8, pp. 1379–1397, Aug. 2015.

[7] C. Merkel, R. Hasan, N. Soures, D. Kudithipudi, T. Taha, S. Agarwal, M. Marinella, "Neuromemristive Systems: Boosting Efficiency through Brain-Inspired Computing", *Computer*, vol. 49, pp. 56-64, Oct. 2016

[8] EO. Neftci, C. Augustine, S. Paul, G. Detorakis, "Event-driven random back-propagation: enabling neuromorphic deep learning machines", *Frontiers in neuroscience.*, vol. 11, pp. 1-11, June 2017.

[9] S. Lim, J.-H. Bae, J.-H. Eum, S. Lee, C.-H. Kim, D. Kwon, B.-G. Park, J.-H.

Lee, "Adaptive learning rule for hardware-based deep neural networks using electronic synapse devices", *Neural Computing and Applications*, vol. 30, pp 1-16, 2018.

[10] S. Yu, B. Gao, Z. Fang, H. Yu, J. Kang, and H.-S. P. Wong, "A neuromorphic visual system using RRAM synaptic devices with Sub-pJ energy and tolerance to variability: Experimental characterization and large-scale modeling," *IEEE International Electron Devices Meeting (IEDM)*, Dec. 2012.

[11] M. Prezioso, F. Merrikh-Bayat, B. D. Hoskins, G. C. Adam, K. K. Likharev, and D. B. Strukov, "Training and operation of an integrated neuromorphic network based on metal-oxide memristors", *Nature*, vol. 521, pp. 61-64, May 2015.

[12] O. Bichler, M. Suri, D. Querlioz, D. Vuillaume, B. DeSalvo, and C. Gamrat, "Visual Pattern Extraction Using Energy-Efficient "2-PCM Synapse" Neuromorphic Architecture", *IEEE Transections on Electron Devices*, vol. 59, no. 8, pp. 2206-2214, May 2012.

[13] S. Y. Woo, K.-B. Choi, S. Lim, S.-T. Lee, C.-H. Kim, W.-M. Kang, D. Kwon, J.-H. Bae, B.-G. Park, and J.-H. Lee, "Synapse device Using a Floating Fin-Body MOSFET With Memory Functionality for Neural Network," *Solid-State Electronics*, vol. 156, 2019.

[14] J.-H. Bae, S. Lim, B.-G. Park, and J.-H. Lee, "High-Density and Near-Linear Synaptic Device Based on a Reconfigurable Gated Schottky Diode", *IEEE Electron Device Letter*, vol. 38, no. 8, pp. 1153-1156, Aug. 2017.

[15] H. Kim, S. Hwang, J. Park, S. Yun, J.-H. Lee, and B.-G. Park, "Spiking neural network using synaptic transistors and neuron circuits for pattern recognition with noisy images," *IEEE Electron Device Letter*, vol. 39, no. 4, pp. 630-633, Apr. 2018

[16] X. Guo et al., "Fast, energy-efficient, robust, and reproducible mixed-signal neuromorphic classifier based on embedded NOR flash memory technology," *IEEE*

*International Electron Devices Meeting (IEDM)*, Dec. 2017.

[17] G.W. Burr, R.M. Shelby, C. di Nolfo, J.-W. Jang, R.S. Shenoy, P. Narayanan, K. Virwani, E.U. Giacometti, B. Kurdi, H. Hwang, "Experimental demonstration and tolerancing of a large-scale neural network (165,000 synapses), using phasechange memory as the synaptic weight element," *IEEE International Electron Devices Meeting (IEDM)*, Dec. 2014.

[18] C.-H. Kim, S. Lim, S. Y. Woo, W.-M. Kang, Y.-T. Seo, S. T. Lee, S. Lee, D. Kwon, S. Oh, Y. Noh, H. Kim, J. Kim, J.-H. Bae and J.-H. Lee, "Emerging memory technologies for neuromorphic computing," *Nanotechnology*, vol. 30, p. 032001, 2018.

[19] S. Hwang, J. Chang, M.-H. Oh, J.-H. Lee and B.-G. Park, "Impact of the Sub-Resting Membrane Potential on Accurate Inference in Spiking Neural Networks," *Scientific Reports*. vol. 10, no. 3515, pp.1-10, Feb. 2020.

[20] P. Livi, and G. Indiveri, "A current-mode conductance-based silicon neuron for Address-Event neuromorphic systems," *IEEE Int. Symp. Circuits and Systems*, pp. 2898-2901, 2009.

[21] Giacomo Indiveri et al "Neuromorphic silicon neuron circuits", *Frontiers in neuroscience*, vol. 5, p.1, May 2011.

[22] S. Oh, D. Kwon, G. Yeom, W.-M. Kang, S. Lee, S. Y. Woo, J. S. Kim, M. K. Park, and J.-H. Lee, " Hardware Implementation of Spiking Neural Networks Using Time-To-First-Spike Encoding", *arXiv*, June 2020.

[23] K. Moon, E. Cha, D. Lee, J. Jang, J. Park, and H. Hwang, "ReRAM based analog synapse and IMT neuron device for neuromorphic system," *in Proc. Int. Symp. VLSI Technol., Syst. Appl. (VLSI-TSA)*, pp. 1–2, 2016.

[24] J. Lin, Annadi, S. Sonde, C. Chen, L. Stan, K. V. L. V. Achari, S. Ramanathan, and S. Guha, "Low-voltage artificial neuron using feedback engineered insulator-

to-metal-transition devices," *IEEE International Electron Devices Meeting (IEDM)*, Dec. 2016.

[25] J. Lin and J. Yuan, "Capacitor-less RRAM-based stochastic neuron for event-based unsupervised learning," *IEEE Biomedical Circuits and Systems Conference (BioCAS)*, pp. 1-4, 2017.

[26] M. Barci et al., "Bilayer Metal-Oxide Conductive Bridge Memory Technology for Improved Window Margin and Reliability," *IEEE Journal of the Electron Devices Society*, vol. 4, no. 5, pp. 314-320, Sept. 2016.

[27] X. Wu, V. Saxena, and K. Zhu, "Homogeneous spiking neuromorphic system for real-world pattern recognition," *IEEE J. Emerg. Sel. Top. Circuits Syst.*, vol. 5, no. 2, pp. 254–266, Jun. 2015.

[28] X. Shamsi, K. Mohammadi and S. B. Shokouhi, "A Hardware Architecture for Columnar-Organized Memory Based on CMOS Neuron and Memristor Crossbar Arrays," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 26, no. 12, pp. 2795-2805, Dec. 2018.

[29] C. Chen, M. Yang, S. Liu, T. Liu, K. Zhu, Y. Zhao, H. Wang, Q. Huang and R. Huang, "Bio-Inspired Neurons Based on Novel Leaky-FeFET with Ultra-Low Hardware Cost and Advanced Functionality for All-Ferroelectric Neural Network," *Symposium on VLSI Tech.*, pp. T136-T137, June 2019.

[30] J. Luo, L. Yu, T. Liu, M. Yang, Z. Fu, Z. Liang, L. Chen, C. Chen, S. Liu, S. Wu, Q. Huang and R. Huang, "Capacitor-less Stochastic Leaky-FeFET Neuron of Both Excitatory and Inhibitory Connections for SNN with Reduced Hardware Cost", *IEEE International Electron Devices Meeting (IEDM)*, Dec. 2019.

[31] S. Dutta,V. Kumar, A. Shukla, N. R. Mohapatra and U. Ganguly, "Leaky Integrate and Fire Neuron by Charge-Discharge Dynamics in Floating-Body MOSFET", *Scientific Reports.*, vol. 7, no. 8257, pp.1-7, Aug. 2017.

[32] M. Kwon, K. Park, M. Baek, J. Lee and B. Park, "A Low-Energy High-Density Capacitor-Less I&F Neuron Circuit Using Feedback FET Co-Integrated With CMOS," *IEEE Journal of the Electron Devices Society*, vol. 7, pp. 1080-1084, Sep. 2019.

[33] J.-W. Han, M. Meyyappan, "Leaky Integrate-and-Fire Biristor Neuron," *IEEE Electron Device Letter*, vol. 39, no. 9, pp. 1457–1460, Jul. 2018.

[34] J.-K. Han, M. Seo, W.-K. Kim, M.-S. Kim, S.-Y. Kim, M.-S. Kim, G.-J. Yun, G.-B. Lee, J.-M. Yu and Y.-K. Choi, "Mimicry of Excitatory and Inhibitory Artificial Neuron With Leaky Integrate-and-Fire Function by a Single MOSFET", *IEEE Electron Device Letter*, vol. 41, no. 2, pp. 208-211, Feb. 2012.

[35] S. Y. Woo, D. Kwon, N. Choi, W.-M. Kang, Y.-T. Seo, M.-K. Park, J.-H. Bae, B.-G. Park, J.-H. Lee, "Low-Power and High-Density Neuron Device for Simultaneous Processing of Excitatory and Inhibitory Signals in Neuromorphic Systems," *IEEE Access,* vol. 8, 2020.

[36] N. Navlakha, J.-T. Lin, and A. Kranti, "Improved Retention Time in Twin Gate 1T DRAM With Tunneling Based Read Mechanism", *IEEE Electron Device Letter*, vol. 37, no. 9, pp. 1127-1130, Sep. 2016.

[37] D. Querlioz, O. Bichler, P. Dollfus and C. Gamrat, "Immunity to device variations in a spiking neural network with memristive nanodevices," *IEEE Transections on Nanotechnology*, vol.12, no. 3, pp. 288-295, May 2013.

[38] S. Y. Woo, K.-B. Choi, J. Kim, W.-M. Kang, C.-H. Kim, Y.-T. Seo, J.-H. Bae, B.-G. Park, J.-H. Lee, "Implementation of homeostasis functionality in neuron circuit using double-gate device for spiking neural network," *Solid-State Electronics*, vol. 165, article 107741, pp.1-6, Mar. 2020.

[39] K.-B. Choi, S. Y. Woo, W.-M. Kang, S. Lee, C.-H. Kim, J.-H. Bae, S. Lim and J.-H. Lee, "A Split-Gate Positive Feedback Device with an Integrate-and-Fire

Capability for a High-Density Low-Power Neuron Circuit," *Frontiers in Neuroscience*, vol. 12, pp. 1-13, Oct. 2018.

[40] B. Liu, H. Li, Y. Chen, X. Li, T. Huang, Q. Wu, and M. Barnell, "Reduction and IR-drop compensations techniques for reliable neuromorphic computing systems," *in Proc. IEEE/ACM International Confernece on Computer-Aided Design (ICCAD)*, 2014.

[41] J. Liang, H.S. Philip Wong, "Cross-Point Memory Array Without Cell Selectors-Device Characteristics and Data Storage Pattern Dependencies," *IEEE Transactions on Electron Devices*, vol. 57, no. 10, pp. 2531-2537, 2010.

[42] J-H. Lee, S. Y. Woo, S.-T. Lee, S. Lim, W.-M. Kang, Y.-T. Seo, S. Lee, D. Kwon, S. Oh, Y. Noh, H. Kim, J. Kim, J.-H. Bae, "Review of candidate devices for neuromorphic applications," *49th European Solid-State Device Research Conference (ESSDERC)*, Sep. 2019.

# 초    록

하드웨어기반 신경망 기술은 복잡한 데이터를 효율적으로 처리하는 신경 모방 시스템의 활용으로 대두되고 있다. 이러한 하드웨어 기반 신경망 기술의 성능 향상을 위해 아키텍쳐와 학습 알고리즘에 적합한 시냅스 어레이와 뉴런 회로들이 개발되고 있다. 특히, 시냅스 어레이에서 전달되는 흥분 및 억제 신호들을 동시에 처리하는 기술을 하드웨어 기반 신경망의 인지 능력을 향상시키는데 중요하다.

본 논문에서 고집적 및 저전력 신경 모방 시스템을 개발하기 위해 시냅스 및 뉴런 소자를 제안한다. 문턱 전압 조절이 가능하고 흥분 및 억제 신호를 동시에 처리할 수 있는 양의 피드백 뉴런 소자를 기존의 뉴런 회로를 대체하기 위해 제안한다. 양의 피드백 동작으로 인해 가파른 스위칭 특성을 가진 양의 피드백 뉴런 소자는 뉴런의 통합 및 발화 기능을 저전력으로 구현할 수 있다. 양의 피드백 뉴런 소자의 구조는 게이트 PNPN 다이오드와 하나의 MOSFET이 효율적으로 접합된 구조이다. 양의 피드백 뉴런 소자의 부유-바디에 전자를 쌓음으로 뉴런의 통합 및 발화가 실험적으로 구현된다. 부유-바디에 쌓인 전자는 접합된 MOSFET 에 억제 신호가 전달되면 드레인으로 빠져나간다. 더욱이, 제안된 양의 피드백 뉴런 소자는 부유-바디 상단에 존재하는 전자 저장 층에 의해 문턱 전압이 조절된다. 양의 피드백 뉴런 회로는 하나의 양의 피드백 소자와 다섯 개의 MOSFET으로 뉴런의 발화 및 통합 동작과 초기화 동작을 구현 할 수 있으며, 그 에너지 소모는 0.62pJ/spike이다. 두개의 독

립된 게이트를 가지는 듀얼 게이트 FET를 시냅스 소자로 제안한다. G1은 시냅스 소자를 켜고 끄는 역할을 하고, 전하 저장 층을 가지고 있는 G2는 시냅스 가중치를 위해 듀얼 게이트 FET의 전류를 조절하는 역할을 한다. 듀얼 게이트 FET의 컨덕턴스 변화는 100 pA 에서 1 μA까지 그 범위가 아주 넓다. 듀얼 게이트 기반 NOR 시냅스 어레이는 파울러-노르하임 (FN) 터널링으로 프로그램 및 이레이즈 동작을 구현할 수 있으며, 저전력으로 시냅스 가중치를 변화할 수 있다. 시냅스 어레이에서 8개 시냅스 전체 전류의 합은 3.6 μA으로 8개의 각 시냅스 전류를 한 전류(3.63 μA)와 약 0.87 % 차이로 거의 일치한다. 또한, 제작된 시냅스 소자 어레이는 네 가지 시냅스 가중치 상태에 대해 0.023, 0.011, 0.015 및 0.032의 변화를 보여준다. 우리가 제안한 양의 피드백 뉴런 회로와 듀얼 게이트 FET 기반 시냅스 어레이는 고집적 및 저전력 신경 모방 시스템을 구현하는데 해결책을 제공해 줄 것이다.

주요어 : 신경 모방 시스템, 양의 피드백 뉴런 소자, 듀얼 게이트 FET, NOR 타입 시냅스 어레이, 흥분 및 억제 신호, 하드웨어 기반 신경망.

학번 : 2014-21721