## 스토캐스틱 컴퓨팅 기반의 임베디드용 AI 프로세서

제 23회 대한민국 반도체 설계대전

김정은, 고광현, 정영우, 정원식

Seoul National University of Science & Technology

**Electronic Engineering** 



### **CONTENTS**

- 1 디자인 소개
- 2 성능 이득 분석
- 3 디자인 결과
- 4 데모 영상





## **>>>>** 디자인 소개

## 스토캐스틱 컴퓨팅 기반의 <u>임베디드용</u> AI 프로세서



파워

Fault-tolerance



정확도



창의성

확률적인 연산 회로를 프로세서에 적용

사업성

인공지능 가속기를 포함한 내고장형의 경량 프로세서를 on-chip으로 구현

기술성

높은 정확도의 스토캐스틱 **연산기**와 AI 가속기 및 페리페럴을 구현

완성도

삼성 28nm 공정을 통해 칩 레이아웃 검증 및 제작 완료

기존의 연산 회로는 복잡하고 오류에 취약

→ 확률적인 연산 회로를 이용하여 해결

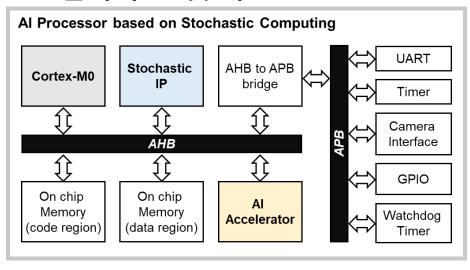




- Stochastic Computing (SC) 확률론적인
  - ✓일반적인 이진수를 [0,1] 범위의 확률 값으로 인코딩
    - 1. 랜덤 수와의 **크기 비교** ex) **100**<sub>(2)</sub> vs 7, 5, 2, 6, 1, 0, 4, 3
    - 2. 숫자 '1'의 비율로 근사 → <u>1의 개수 / 시퀀스의 총 길이</u>



#### 설계 구조 및 목표



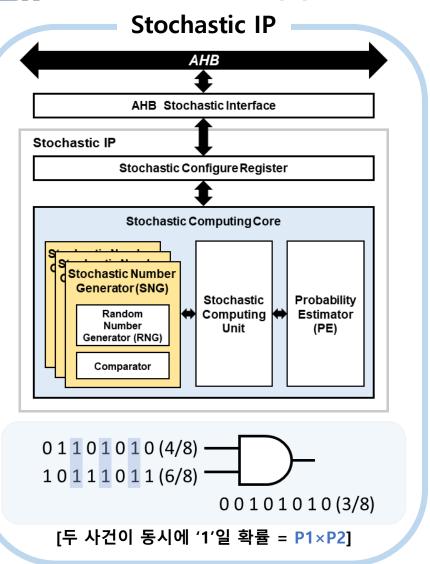
| 항목       | 설계 목표                       |
|----------|-----------------------------|
| 코어 사양    | Arm 사의 Cortex-M 시리즈         |
| 동작 주파수   | 50MHz                       |
| 메모리      | 코드 영역: 16KB, 데이터 영역: 128KB  |
| 인터페이스    | AHB, APB                    |
| 스토캐스틱 IP | 정확도 90% 이상                  |
| 인공지능 가속기 | k-NN <b>알고리즘</b> 적용         |
| 페리페럴     | 카메라 I/F, 시리얼 I/F, GPIO, 타이머 |

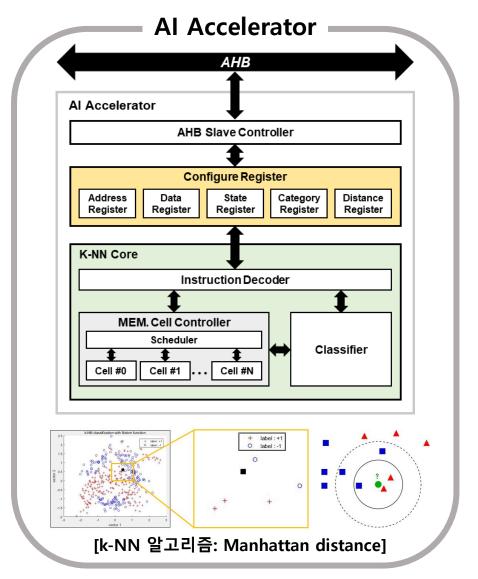
[스토캐스틱 컴퓨팅 기반 AI 프로세서의 전체 구조도]

[초기 설계 목표 및 사양]



### 성능이득 분석(1) : 면적 & 파워

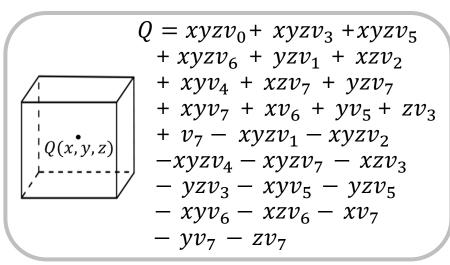




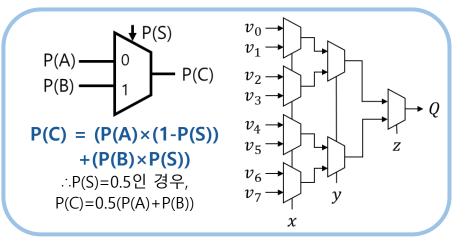




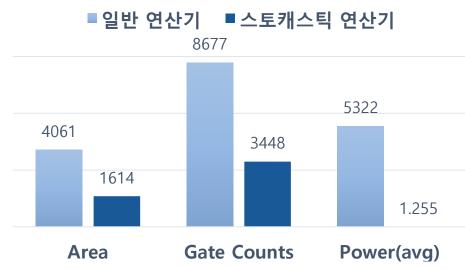
### 성능이득 분석(1) : 면적 & 파워



#### [Tri-linear interpolation 연산식[1]]



[SC 덧셈기 기반의 Tri-linear interpolation 연산기[2]]



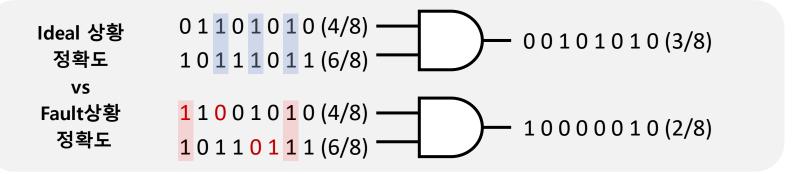
- \* 8bit 연산 기준
- \*Samsung 28nm, Design Compiler

|       |               | 일반 연산기                     | 스토캐스틱 연산기                  |
|-------|---------------|----------------------------|----------------------------|
| La    | tency         | 1cycle                     | 32cycles (8배 가속)           |
| P     | Area          | $4,061 \mu m^2$            | 1,614μm²                   |
|       | Gate<br>ounts | 8,677                      | 3,448                      |
| ŗ     | Max           | 5.3956mW                   | $1.2243 \times 10^{-3} mW$ |
| power | Min           | $1.0172 \times 10^{-3} mW$ | $8.4204 \times 10^{-7} mW$ |
| д     | avg           | $5.322 \times 10^{-3} mW$  | $1.255 \times 10^{-6} mW$  |

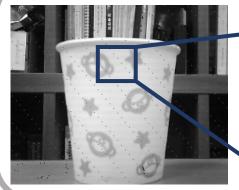


## >>> 성능이득 분석(2) : Fault-tolerance

- 시퀀스의 전체 길이 중 '1'의 비율을 컴퓨팅에 사용 ex) (1,0,0,0) = (0,1,0,0) = (0,1,0,0,0,1,0,0) = 1/4 (시퀀스의 길이 길수록 fault에 강함)
- 시퀀스의 각 비트는 동일한 weight를 가짐 ex)  $3/256 \rightarrow 4/256$  or 2/256 vs  $01000001 \rightarrow 11000001$



Fault-tolerance 분석





| 240 | 242 | 239 |
|-----|-----|-----|
| 255 |     | 255 |
| 255 | 255 | 254 |

(a) C코딩 필터 연산

## [이미지 엣지 검출 연산]

| \            |  |
|--------------|--|
| <i>92}</i> ∵ |  |

| 15  | 17  | 254 |
|-----|-----|-----|
| 15  |     | 248 |
| 247 | 234 | 251 |

\*0(검정색)~255(흰색)

(b) 스토캐스틱 필터 연산



## ὢ 디자인 결과

정확도 분석: 이미지 엣지 검출 연산(Sobel 필터 연산)

| $P_1$ | $P_2$ | $P_3$ |
|-------|-------|-------|
| $P_4$ | $P_5$ | $P_6$ |
| $P_7$ | $P_8$ | $P_9$ |



(a) Sobel 필터 마스크 (b) 원본 이미지

$$|G| = |(P_1 + 2P_2 + P_3) - (P_7 + 2P_8 + 2P_9)| + |(P_3 + 2P_6 + P_9) - (P_1 + 2P_4 + P_7)|$$





(c) C코딩 필터 연산

(d) 스토캐스틱 필터 연산

| 단위 (%) | (c)      | (d)   |
|--------|----------|-------|
| 평균 정확도 | 100%(가정) | 92.4% |

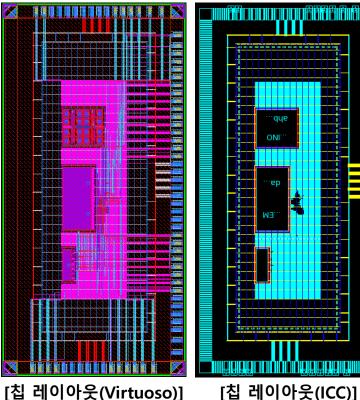
| 항목      | 설계 목표              | 목표 달성 |
|---------|--------------------|-------|
| 코어 사양   | Arm사의 Cortex-M 시리즈 | 0     |
| 동작 주파수  | 50MHz              | 0     |
| 메모리     | 코드 영역: 16KB / 데이터  | 0     |
| -11-11  | 영역: 128KB          |       |
| 인터페이스   | AHB, APB           | 0     |
| 스토캐스틱   | 저희도 02 40/         |       |
| IP      | 정확도 92.4%          | Ο     |
| 인공지능    | k-NN (k-Nearest    | 0     |
| 가속기     | Neighbor) 알고리즘 적용  | O     |
| ᆒᅴᆒᄙ    | 카메라 I/F, 시리얼 I/F,  |       |
| 페리페럴    | GPIO, 타이머          | 0     |
| 치수      | 0                  |       |
| (설계 면적) | 2mm × 4mm          | 0     |

[디자인 결과 및 목표 달성 여부]



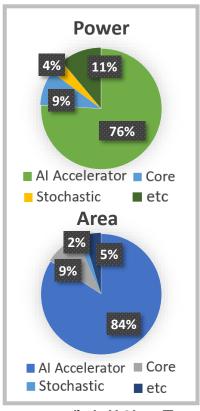
## ὢ 디자인 결과

레이아웃 디자인



| i i i | <b>4</b>    |          |
|-------|-------------|----------|
|       |             |          |
|       | ===         |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       | <u>=::1</u> |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             | <u> </u> |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
|       |             |          |
| (ICC) | 1           |          |
| (ICC) | <i>,</i> 1  |          |

| 항목          | 칩 사양                |
|-------------|---------------------|
| 공정          | Samsung 28nm CMOS   |
| -101        | Core: 1.0V          |
| 파워          | IO: 1.8V            |
| 동작 주파수      | 50MHz               |
| Gate Counts | 1062K @ 50MHz       |
| ==! .!.     | 동작 온도:-40°C~125°C   |
| 동작 사양       | 코어 전압: 0.9V~1.1V    |
|             | code region: 16KB   |
| 메모리         | data region: 128KB  |
|             | Al Accelerator: 1KB |
| 치수          | 2mm × 4mm           |
|             |                     |



[프로세서 하위 모듈 면적 및 소비전력 비율1

내고장형, 경량 회로 → Neural Networks, Image Processing과 같이 복잡하고 많은 연산량이 요구되는 응용 시스템에 적용 예상





### >> 데모 영상

• 데모 순서





카메라 모듈을 통한 **촬영** 및 이미지 데이터 저장

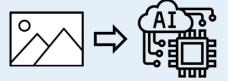
2





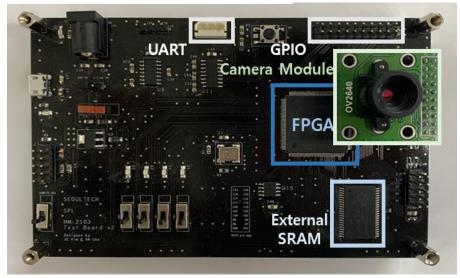
스토캐스틱 연산기를 통한 이미지 **엣지 검출 연산** 

3



출력 이미지를 AI 가속기의 학습/인식 데이터로 활용

#### • 테스트 보드 사양



[테스트 보드 이미지]

| 항목     | 테스트 보드 사양                     |
|--------|-------------------------------|
| FPGA   | Altera MAX10 10M50SCE144C8G   |
| 입력 전압  | 5V / 12V                      |
| 동작 주파수 | 50MHz                         |
| 메모리    | off-chip 512KB SRAM           |
| 카메라 모듈 | OV2640 CMOS                   |
| 페리페럴   | 카메라 I/F, 시리얼 I/F (UART), GPIO |







고광현



정영우

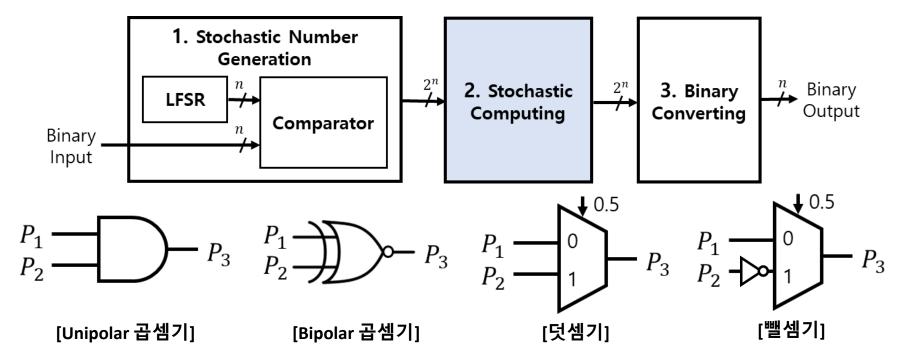


정원식

# 감사합니다



## 성능이득 분석(1) : 면적 & 파워



- 두 개의 스토캐스틱 입력 시퀀스는 서로 **독립적** 1.
- 2. 각 입력 시퀀스에서 **'1'이 나타날 확률**을 P**1, P2** 로 정의
- AND gate의 출력이 '1'일 확률은, 두 사건이 동시에 '1'일 확률 = P1×P2 3.





### Backup slide (reference)

#### Area & Energy efficiency

W. Qian, X. Li, M. D. Riedel, K. Bazargan and D. J. Lilja, "An Architecture for Fault-Tolerant Computation with Stochastic Logic," in IEEE Transactions on Computers, vol. 60, no. 1, pp. 93-105, Jan. 2011, doi: 10.1109/TC.2010.202.

A 16×128 Stochastic-Binary Processing Element Array for Accelerating Stochastic Dot-Product Computation Using 1-16 Bit-Stream Length

> Qian Chen, Yuqi Su, Hyunjoon Kim, Taegeun Yoo, Tony Tae-Hyoung Kim, and Bongjin Kim School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore 50 Nanyang Avenue, Singapore, 639798
> Email: { e170029, yuqi003, kimh0003}@e.ntu.edu.sg, {tgyoo, thkim, bjkim}@ntu.edu.sg

aborace. This work present la-128 actuatio-haus processis demands for energy-mass efficient processing of artificial presents of a constitution of a processing of actual presents of a transport of a constitution of a processing of a maximum processing of an artificial processing of a constitution of a company rather tools for a grant processing of a company rather tools for a constitution of a company rather tools for a company rather than a company rather t stochastic computation using MU dot-product error is 69-to-1.5% from the baseline stochastic meth to-2048. A mean MNIST classifica is 1.19% lower than 8b binary) us The measured energy from a 6tu product, and the energy efficiency

Stochastic computing [1-3] is sethod based on probabilities of samy 1's in a bit-stream). S tochastic computing are perfe logic gates, it is suitable for mo mall footprint. Despite the po has not been employed in pract concerns in the latency and en-



978.3.9819263.4.7/DATE20/6

| Stochasti | c Multiply | Stochas | stic Add      |
|-----------|------------|---------|---------------|
| A-D-c     | A-Ma-c     | A-T)-c  | <u>*</u> —\_c |

|                         | [7] ASSCC'16              | [8] JSSC'19               | Propos ed                                   |
|-------------------------|---------------------------|---------------------------|---|
| Computing Type          | Analog<br>(Deterministic) | Analog<br>(Deterministic) | Digital<br>(Stochastic)                     |
| Technology              | 28nm                      | 65nm                      | 65nm  |
| Precision Control       | Fixed (8b)                | Fixed (6b/1b)             | Reconfigurable                              |
| MAC Circuit Type        | Analog                    | Analog In-Mem.            | Digital                                     |
| ADC/DAC Overhead        | Embedded                  | Required                  | Not Required                                |
| Parallelism             | No                        | 16×                       | 16×   |
| Energy Efficiency       | 9.61TOPS/W                | 51.3TOPS/W                | 25.5TOPS/W @ <sup>D</sup> N=16              |
| Energy FoM <sup>A</sup> | 208fJ                     | 39.0fJ                    | 39.3fJ @ <sup>D</sup> N=16                  |
| Area FoM <sup>B</sup>   | 720.4μm <sup>2</sup>      | 61.5μm²                   | 154.4μm²                                    |
| Accura cy <sup>c</sup>  | N/A                       | 98%<br>(CNN/4-layers)     | 96.1% @ <sup>D</sup> N=16<br>(MLP/3-layers) |

AEnergy FoM=Energy/(# of inputs)×(# of dot-products)

BArea FoM=Area/(# of inputs)×(# of dot-products) CMNIST dataset DBit-stream length

#### [레퍼런스 논문의 이미지]

#### **Error Injection**

Q. Chen, Y. Su, H. Kim, T. Yoo, T. T. -H. Kim and B. Kim, "A 16×128 Stochastic-Binary Processing Element Array for Accelerating Stochastic Dot-Product Computation Using 1-16 Bit-Stream Length," 2020 Design, Automation & Test in Europe Conference & Exhibition (DATE), 2020, pp. 678-681, doi: 10.23919/DATE48585.2020.9116349.

#### An Architecture for Fault-Tolerant Computation with Stochastic Logic

Weikang Qian, Student Member, IEEE, Xin Li, Marc D. Riedel, Member, IEEE, Kia Bazargan, and David J. Lilja, Fellow, IEEE

Abdrace—Abunting concerns over viriability, detects, and note motivate a new approach for digital circularly studiates logge, that is one to yet with more and unchesting. Teachings and proceedings are also supported by mens and unchesting. Teachings also proceedings are unclearly circular and operations are well-exceeded an embodology for circular and systems are well-exceeded an embodology for any embodies, but also loss yet public and personal processing control and the processing control and any embodies and embodies benchmarks for image processing. The stochastic architecture requires less area than conventional hardware implementation. Moreover, I is much more iderant of ord errors (bit flips) than these deterministic implementations. This fault tolerance scale gradefully to very large numbers of errors.

Index Terms—Stochastic logic, reconfigurable hardware, fault-tolerant computation.

#### 1 INTRODUCTION up, the precise Boo uncertainty at the c a major concern, pa environments such tion. A broad dass tion. A broad class graphy and commu complexity if physic [1], [2]. Application physical phenomen-quantum physics, a

Added Error > 5

The authors are will Engineering, Universit MN 55455. E-mail: (a)

10 15

Analysis of Error Distribution of the Gamma Correction Function

Produced Error (%) > 10 > 15 > 20 Conv. Stoc. Conv. Stoc. Conv. Stoc. Conv. Stoc. (%)Conv. Stoc. 0.00.0 0.00.00.0 0.0 0.0 0.0 0.0 4.00.0 3.0 0.0 2.4 0.0 2.0 0.0 0.0 1.6 8.2 0.5 6.1 0.0 5.0 0.0 4.2 0.0 3.2 0.0 5 19.5 27.8 14.71.6 12.10.0 10.3 0.07.9 0.0 22.7 7.9 19.3 24.3 15.1 0.0 38.5 36.6 32.1 21.4 27.3 21.8

TABLE 2

[레퍼런스 논문의 이미지]