SMART FACTORY WITH ARTIFICIAL INTELLIGENCE AT ENDPOINT

200

ATSUSHI HASEGAWA

IOT AND INFRASTRUCTURE BUSINESS UNIT RENESAS ELECTRONICS CORPORATION

MPSOC'19 July 8th, 2019

BIG IDEAS FOR EVERY SPACE RENESAS

/// _lus _lus

© 2019 Renesas Electronics Corporation. All rights reserved.

6.677

SMART FACTORY CHALLENGE





FURTHER HIGH-QUALITY: FDC CHALLENGE FDC : Fault Detection and Classification



EFFECT OF AUTOMATED FDC SYSTEM



ISSUE OF FDC AUTOMATED DETECTION



Overlooking by improper threshold

Trade off between error detection and false detection

New approach



③Integrate parameters





ADDITIONAL RESULT : HIGHER SAMPLING RATE



Analog value and higher sampling rate makes Invisible to visible

Can we find more valuable fact by higher sampling rate as 1ms, sub msec ??

USE CASE CLASS-1: e-AI ANOMALY DETECTION FOR HUNDREDS MILLION MOTORS

Benefits:

- Improve service quality
- Avoid downtimes
- Reduce maintenance cost



e-AI detects and pre-warns anomalies everywhere every time



Renesas is shipping 200M+ motor control MCU per year.

New product series "RX66T" (Nov.) will enable e-AI Anomaly Detection.

Renesas e-AI Solution



e² studio MCU/MPU Development Environment Apply learned AI network to all embedded systems

Significant performance improvement under low power and reduction of memory usage are required for the factory.

Page 11

TARGET`<u>10x Performance</u>' for End-point Intelligence



DRP: Dynamically Reconfigurable Processor



DRP: State Transition Controller (STC) + Coarse grained reconfigurable data-path

Data-path Resources (16bit width)

- Array of ALU:
 - 96x PE
- Linear memory and multiplier units:
 - 16x VMU (512word 2port mem, 2x {MUL16, FP16 ALU})
 - 10x HMU (8Kword 2port mem)
 - 2x {DIV32, FP32 ALU}
- Programmable Wires (not shown)

e-AI Accelerator DRP: Dynamically Reconfigurable Processor



*DRP; Dynamically Reconfigurable Processor

Page 14



Spatial + Time-multiplexed computing



Control vs Data Processing



TARGET <u>`100x Performance</u>' for End-point Intelligence





AI-MAC Data Processing



Multi Layer Processing: Layer 1



(a) CPU prepares a DRP descriptor

(b) DRP starts processing the Layer 1

- By using optimized Layer 1 context

(c) AI-MAC is configured by AI-Mac Layer 1 descriptor

(d) Weight data 1 is transferred to prepare the calculation

Multi Layer Processing: Layer 3



(a) DRP repeats dynamic reconfiguration and layer processing.

(b) Finally, DRP generate an interrupt signal to CPU to notice the end of all layer processing.



AI-MAC Architecture

- AI-MAC has 1024 MACs
 - MACs are Divided into four MAC groups.
 - A MAC group has four of 64-MAC.
- MAC group is connected with DRP by 64bit in/out data port.
- AI-MAC can handle up to four independent matrix.



AI-MAC Architecture



- MAC has three precision modes
 - FP16: FP16 weight and data
 - Binary Weight: 1bit weight, FP16 data
 - Binary Net: 1bit weight and data
- A MAC has local memory data input and shared 16bit input/output data port.



Machine Learning, Deep Learning adaption



Al-Mac

Prototype Chip



Chip Specification Technology: UMC 28HPC, 10-Metal Package: 324pin BGA Core: Dual STP3-AI, ARM CPU, etc. On-Chip Memory: 128Mbit SRAM AI-MAC: max 500MHz DRP: max 333MHz One STP3-AI core peak performance (designed): 1,000 GOPS

Built upon micro-computer chip

Measured 960 GOPS

Evaluation Board

Evaluation (Some VGG Layers)

Layer	Data Size (in, out)	Mode	Batch	[GOPS]	Bottleneck
Convolution, 3x3	56 x 56 x 256ch	FP16	1	524.5	MAC Usago Pato
Convolution, 3x3	28 x 28 x 512ch	FP16	1	754.5	(DRP Program)
Convolution, 3x3	28 x 28 x 512ch	Binary Weight	1	754.5	
Full Connection	4096	FP16	1	19.3	
Full Connection	4096	FP16	16	268.9	Weight Data
Full Connection	4096	FP16	64	687.1	Transfer Time
Full Connection	4096	Binary Weight	1	461.5	

One STP3-AI core used, Exclude initial mem transfer time , Ideal: 1,000GOPS

- Convolution throughput depends on the utilization rate of 1024 MACs.
 - Data processing design and DRP program determine.
- Full Connection throughput depends on weight memory transfer time from ext. mem.
 - Binary Weight mode (1/16 weight data size) or batch data processing are effective.

Conclusion

- There's a certain application to match various range of ML performance. In smart factory and other embedded segment have good opportunity to utilize relatively lower performance Machine Learning solution.
- Dynamic reconfigurable Processor architecture can provide 10x performance acceleration compare to standard MCUs for AI inference.
- Dedicated AI-Mac accelerator and DRP architecture can improve another 10x boost of AI inference performance.
- Part of these performance improvement came from limited flexibility of programing, less generality compare to generic CPU or GPU. But we think such restriction will not harm the adaption for future Machine Learning application.

Thank you for your attention

Atsushi Hasegawa <u>Atsushi.hasegawa.gx@Renesas.com</u> IoT and Infrastructure Business Unit **Renesas Electronics Corporation**



© 2019 Renesas Electronics Corporation. All rights reserved.