# LRADNN: High-Throughput and Energy-Efficient Deep Neural Network Accelerator using Low Rank Approximation

Jingyang Zhu[1], Zhiliang Qian[2] , and Chi-Ying Tsui[1]

[1] The Hong Kong University of Science and Technology, Hong Kong
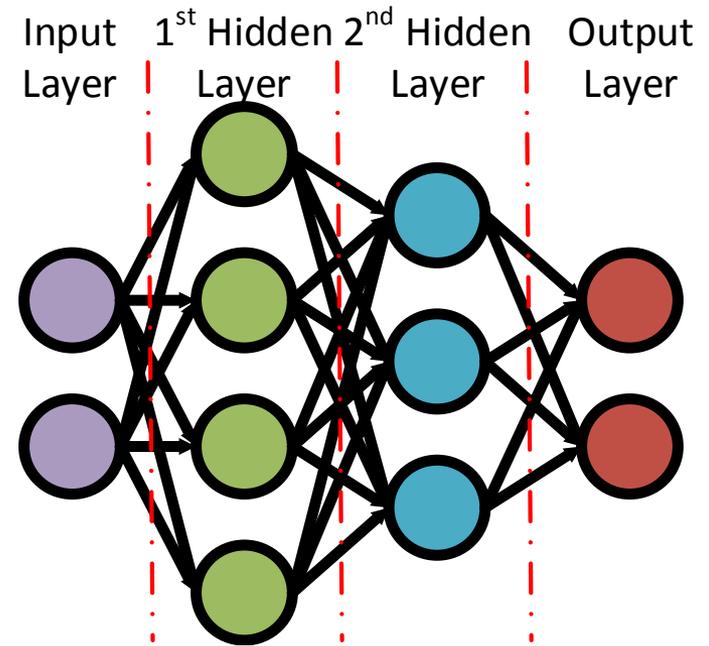[2] Shanghai Jiao Tong University, Shanghai, China

# Outline

- **Introduction**
- Related work
- Motivation of LRADNN
- Low rank approximation (LRA) predictor
- LRADNN: hardware architecture
- Experiment results
- Conclusion

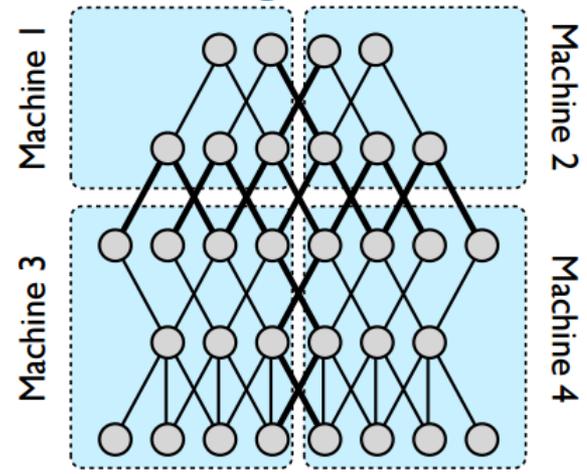# Deep neural network (DNN) and hardware acceleration

- Layer-wise organization with hierarchical feature extractions
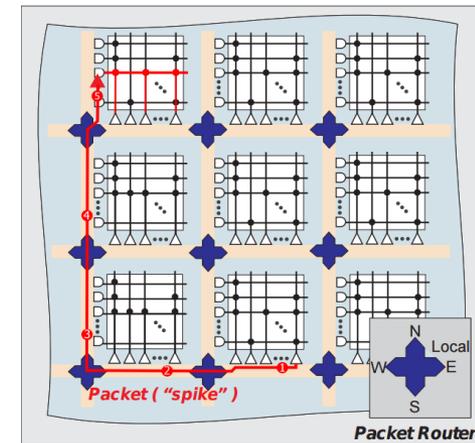
# Deep neural network (DNN) and hardware acceleration

- Layer-wise organization with hierarchical feature extractions

- Hardware acceleration
  - CPU clusters: Google brain
  - GPU clusters: AlexNet
  - ASIC: IBM TrueNorth
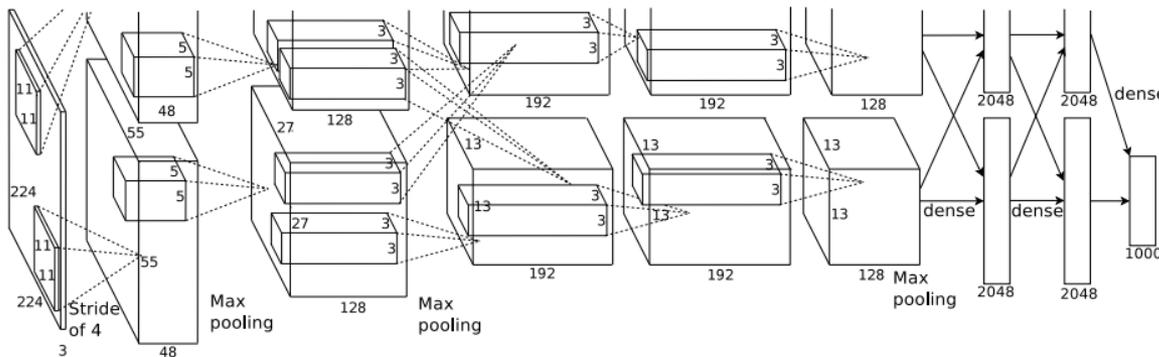
**Google brain**



**AlexNet**



**IBM TrueNorth**

# Outline

- Introduction

- **Related work**

- Motivation of LRADNN

- Low rank approximation (LRA) predictor

- LRADNN: hardware architecture

- Experiment results

- Conclusion

# Related work

- DianNao: a general accelerator for DNN
  - Large energy consumption in memory access
- AxNN: energy-efficient accelerator by approximating resilient neurons
  - Only reduce power in datapath!

**Energy consumption in DianNao**

**AxNN datapath**

# Outline

# Motivation: sparsity in DNN

- Intrinsic sparsity in DNN to avoid overfitting / better feature extractions
  - Sparsity: proportion of inactive neurons (activation = 0)

**Sparsity in DNN targeted for MNIST**

# Motivation: sparsity in DNN

- Intrinsic sparsity in DNN to avoid overfitting / better feature extractions
  - Sparsity: proportion of inactive neurons (activation = 0)
- Conventional computation

$$a_i^{(l+1)} = \sigma\left(\sum_{j=1}^{s_l} W_{ij}^{(l)} a_j^{(l)}\right), i \in [1, s_{l+1}]$$

Regardless of the activeness of neurons

  - Arithmetic ops / neuron: $s_l$ mults, $s_l - 1$ adds, 1 nonlinear
  - Memory accesses / neuron: $2s_l$ (weights & acts)

# Motivation: bypass unnecessary operations

- Dynamically bypass the inactive neurons



- Dedicated predictor for the activeness of neurons
  - Simple: less arithmetic operations involved
  - Accurate: small performance degradation

# Outline

- Introduction
- Related work
- Motivation of LRADNN
- Low rank approximation (LRA) predictor
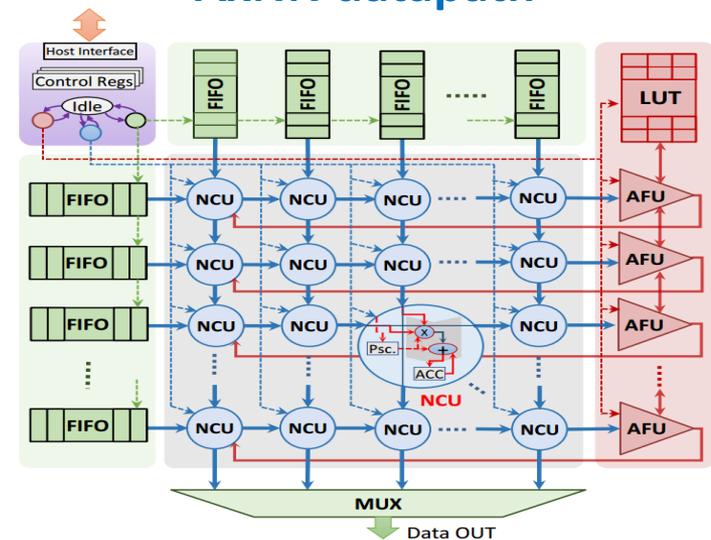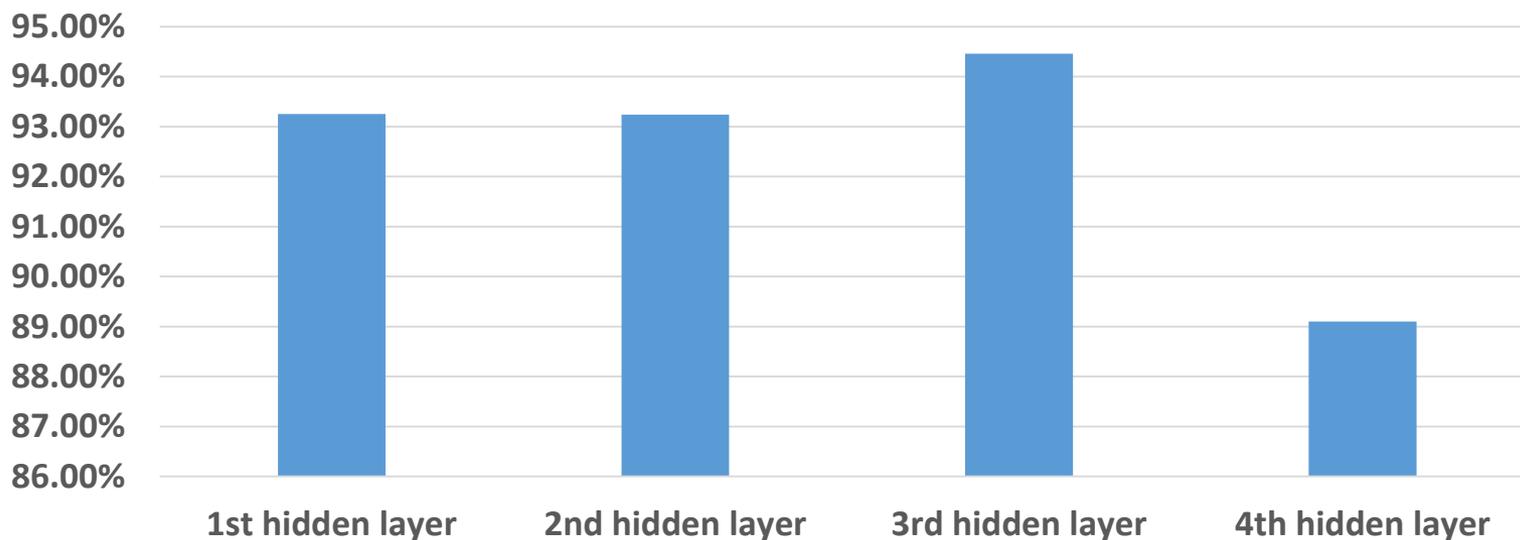- LRADNN: hardware architecture
- Experiment results
- Conclusion

# Low rank approximation (LRA) predictor [1]

- Approximation for synaptic weights W

$$\min_{\widetilde{W}} |\widetilde{W} - W|_F$$

$$s.t. \ \mathrm{rank}(\widetilde{W}) = r$$

  where r is a hyper-parameter controlling the prediction accuracy vs. computation complexity

- A nice closed form solution based on SVD



$U = U_r$

$V = \Sigma_r V_r^T$

[1] Davis A, Arel I. Low-Rank Approximations for Conditional Feedforward Computation in Deep Neural Networks[J]. arXiv preprint arXiv:1312.4461, 2013.

# Feedforward pass with LRA

- **Algorithm 1** DNN feed-forward with LRA

  1: **procedure** FEED-FORWARD($a^{(1)}, \{W^{(l)}\}_{l=1}^{L-1}$)
  2:      **for** $l := 1$ **to** $L - 2$ **do**      ▷ Over all hidden layers
  3:         $y^{(l+1)} = V^{(l)}a^{(l)}$         ▷ V calculation stage
  4:         $p^{(l+1)} = U^{(l)}y^{(l+1)}$         ▷ U calculation stage
  5:         **for** $i := 1$ **to** $s_{l+1}$ **do**         ▷ W calculation stage
  6:            **if** $p_i^{(l+1)} \geq 0$ **then**     ▷ Compute active neurons
  7:               $z_i^{(l+1)} = \sum_j W_{ij}^{(l)}a_j^{(l)}$
  8:            **else**         ▷ Gate inactive neurons
  9:               $z_i^{(l+1)} = 0$
  10:         $a^{(l+1)} = g(z^{(l+1)})$         ▷ ReLU transformation
  11:       $a^{(L)} = g(W^{(L-1)}a^{(L-1)})$    ▷ No LRA for final output

- Take LRA into offline training (backpropagation) to improve the performance

# Outline

- Introduction
- Related work
- Motivation of LRADNN
- Low rank approximation (LRA) predictor
- **LRADNN: hardware architecture**
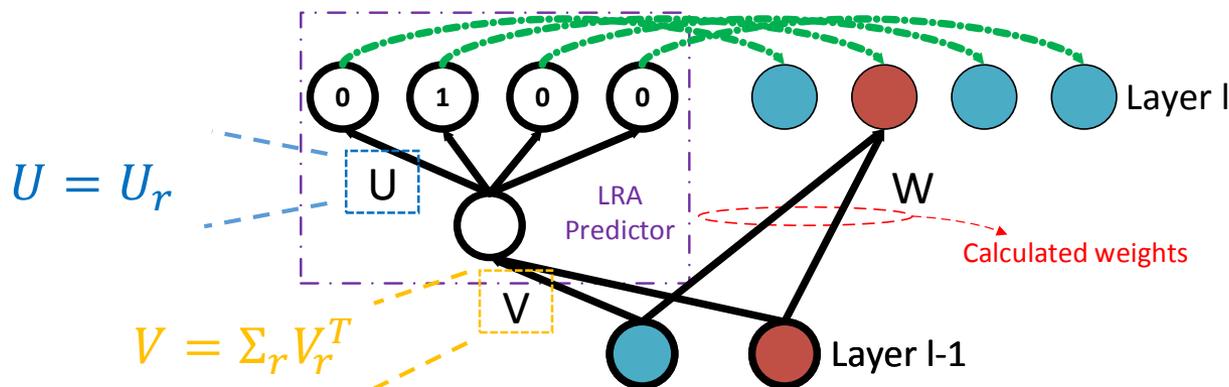- Experiment results
- Conclusion

# LRADNN: hardware architecture for LRA predictor

- Top view of the architecture: 5-stage pipeline
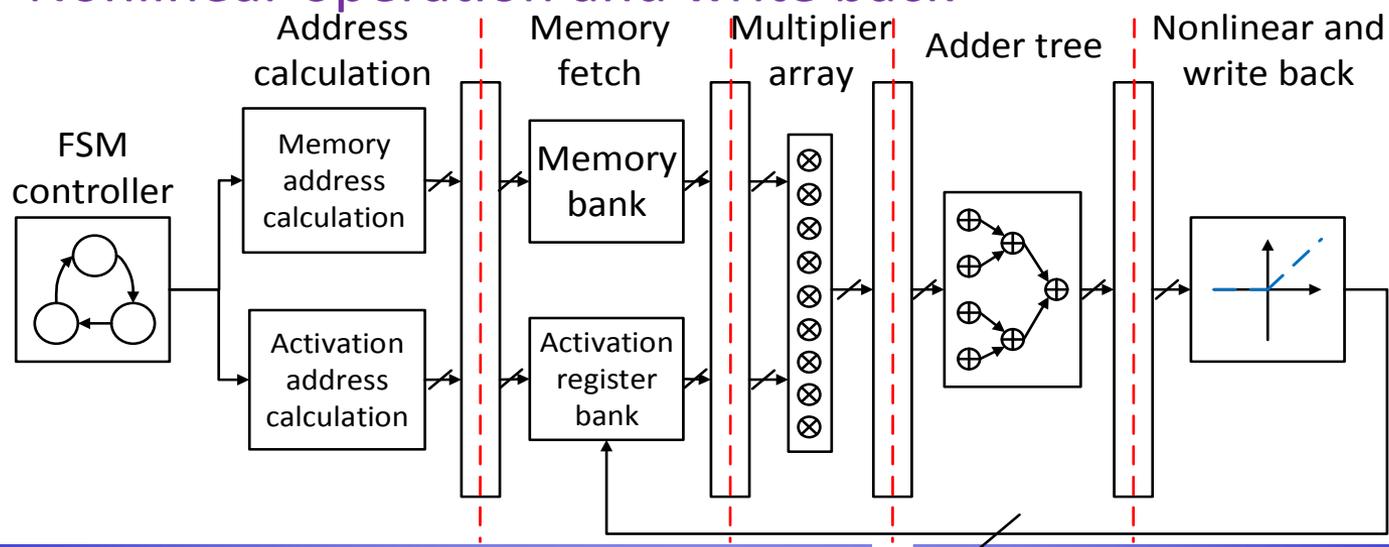  - Address calculation
  - Memory fetch
  - Multiplication
  - Addition
  - Nonlinear operation and write back

$$a_i^{(l+1)} = \sigma \left( \sum_{j=1}^{s_l} W_{ij}^{(l)} \times a_j^{(l)} \right)$$

# Status controller

- 3 calculation statuses
  - V computation
  - U computation
  - W computation
    - Hidden layer
    - Output layer (no nonlinear operation)
  - Data dependency (delay 3 CCs)

$p^{(l)}=U^{(l)}V^{(l)}z^{(l)}$

$a^{(l+1)}=LRA(p^{(l)}, a^{(l)})$



State diagram nodes: Idle → V → Stall (V finishes) → U (U finishes) → Stall → Hid W (Hidden W finishes) → Stall → Out W; Out W → Idle (Only output layers); Out W ↔ Idle (Output W finishes); V → (Still has hidden layers).

| ADDR | MEM | MUL | ADD | NON/W B |
|------|-----|-----|-----|---------|

| ADDR | MEM | MUL | ADD | NON/W B |
|------|-----|-----|-----|---------|

# Address calculation: memory organization

- Memory word width: parallelism of the accelerator

$$\begin{bmatrix} W_{11} & \cdots & W_{1n} \\ W_{21} & \cdots & W_{2n} \\ \vdots & \ddots & \vdots \\ W_{m1} & \cdots & W_{mn} \end{bmatrix}$$

W(1, :)

W(2, :)

One word: parallelism of the accelerator

- For a specified weight $W_{ij}$

$$\text{addr}(W_{ij}) = i\lceil n/p \rceil + \lfloor j/p \rfloor$$

\# word lines for an output neuron

Column block $j_{blk}$

# Memory access: activation register banks

- Two physical register banks are interleaved to two logical registers: input and output activations

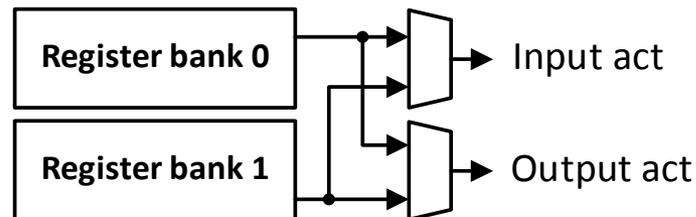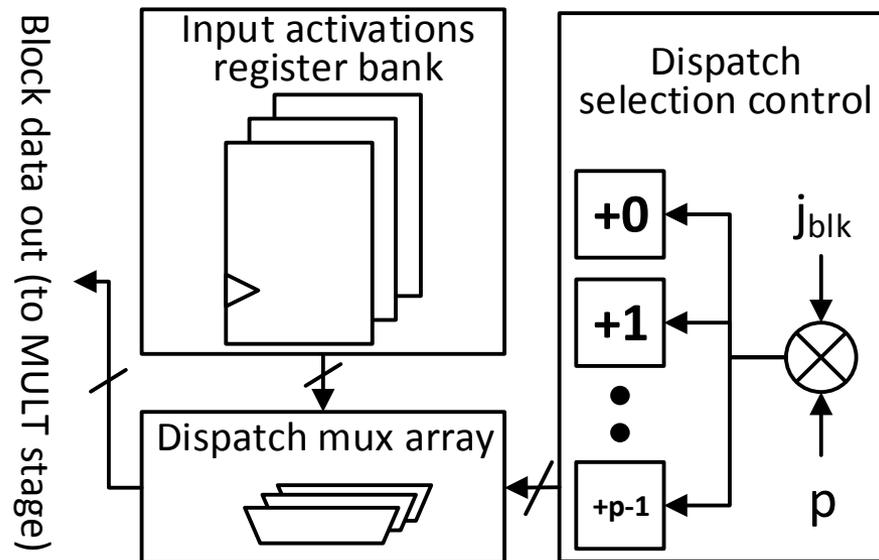

| | |
|---|---|
| **Register bank 0** | → Input act |
| **Register bank 1** | → Output act |

- Input activations

# Memory access: activation register banks (cont.)

- Output activations

- Extra register for LRA scheme ($\leq 5\%$)
  - Tmp1: $\text{tmp1} = V^{(l)}z^{(l)}$
  - Tmp2: $\text{tmp2} = U^{(l)}\text{tmp1}$
  - Predictor: $p^{(l)} = \text{tmp2} > 0$

Output activations register bank

Write demux

i → Read out activations

To WB stage (for accumulation)

| Logical variable | Tmp1 | Tmp2 | Predictor |
|---|---|---|---|
| Physical location | Tmp1 reg | Output act | Predictor reg |
| Size (depth x width) | rank x FP width | # acts x FP width | # acts x 1 |

# Computational stages

- Parallel multiplication
  - p multipliers, determining the parallelism of the accelerator

- Merging operation
  - Adder tree

- Nonlinear operation
  - ReLU
  - Dropout

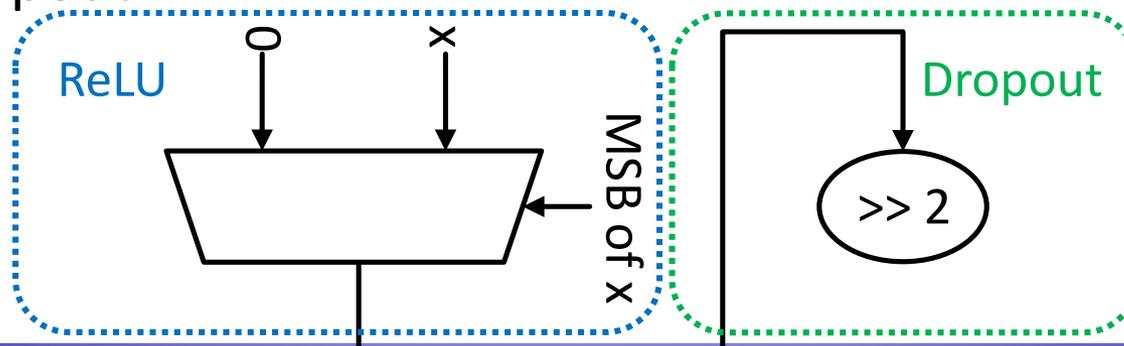# adders @ level 1: p

Depth: $\log_2(p)$

ReLU

0    x

MSB of x

Dropout

>> 2

# Active neurons search: behavior level

- Priority encoder based search

**Algorithm 2** Active neurons search

1: **procedure** ACTIVENEURONSEARCH($i, p^{(l)}$)
2:      **for** searchId := $i + 1$ **to** $i +$ SCAN_WINDOW_SIZE **do**
3:          **if** $p^{(l)}_{\text{searchId}} == 1$ **then**
4:              nextActiveId = searchId
5:              **break**
6:      **if** searchId $>$ $i +$ SCAN_WINDOW_SIZE **then**
7:          nextActiveId = $i +$ SCAN_WINDOW_SIZE $+ 1$
     **return** nextActiveId

*(Priority encoder search: lines 2–5)*

*(Search miss: lines 6–7)*

- Search miss penalty: 1 CC

# Active neurons search: hardware level

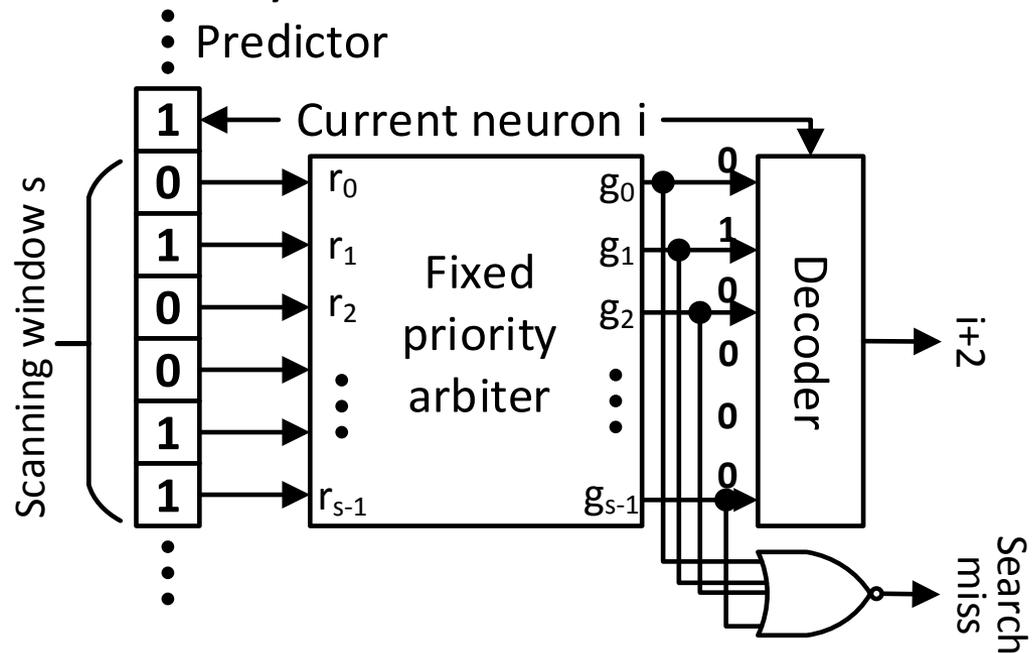- Priority encoder based search
  - Higher priority assigned to LSB

- Decoder
  - One hot to binary decoder

# Outline

- Introduction
- Related work
- Motivation of LRADNN
- Low rank approximation (LRA) predictor
- LRADNN: hardware architecture
- **Experiment results**
- Conclusion

# Simulation setup

- Behavior level simulation on MATLAB
  - Offline training
  - Fixed point simulation
- RTL implementation using Verilog
  - Technology node: TSMC 65nm LP standard cell
  - Memory model: HP CACTI 6.5
- Comparison of behavior, pre-synthesis and post-synthesis simulations
- Power evaluation based on the post-synthesis simulation (extracted switching activities)
- 3 real benchmarks are tested on LRADNN: MNIST, Caltech101, and SVHN

# Summary of the implementation

| Micro-architecture parameters | Value |
|---|---|
| # multipliers (parallelism) | 32 |
| Depth of the activation register bank | 1024 |
| Layer index width (max layer no.) | 3 |
| Fixed point for the internal data path | Q7.12 |
| Fixed point for W memory | Q2.11 |
| Fixed point for U, V memory | Q5.8 |
| W memory size | 3.5MB |
| U, V memory size | 448KB |
| Scanning window size | 16 |
| # V calculation registers (max rank) | 128 |

# Area and timing results

- Memory dominant for area and timing

| | LRADNN (direct) | LRADNN |
|---|---|---|
| Total area (mm2) | 51.94 | 64.18 |
| Critical path (ns) | 8.94 | 9.03 |

- Area overhead caused by U, V memory
$$448KB \times 2/3.5MB = 25\%$$

- Timing overhead caused by extra U, V, and W memory data MUX selection (~ 1% increase)

# Training results on real applications

- Prediction loss: test error rate for LRA feedforward – test error rate for plain feedforward

| | Caltech 101 silhouettes | MNIST | SVHN |
|---|---|---|---|
| Architecture | 784-1023-101 | 784-1000-600-400-10 | 1023-1000-700-400-150-10 |
| # connections | 0.90M | 1.63M | 2.06M |
| Rank | 50 | 50-35-25 | 100-70-50-25 |
| Prediction loss (fixed point) | -0.09% | 0.07% | -0.93% |

# Theoretical lower bounds for accelerators

- Number of cycles

$$\frac{\text{\# synptic connections in DNN}}{\text{parallelism}} = \frac{\sum_{l=1}^{L-1} s_{l+1}(1 + s_l)}{\text{parallelism}}$$

- Power consumption (only consider memory access power)

$$\frac{\text{E}_{\text{FF}}}{t_{FF}} = \frac{\text{E}_{\text{MEM}} \times \text{n}_{\text{cyc}}}{T_{\text{cyc}} \times \text{n}_{\text{cyc}}}$$

$\text{E}_{\text{FF}}$: energy consumption during feedforward
$\text{t}_{\text{FF}}$: elapsed time for feedforward
$\text{n}_{\text{cyc}}$: number of ideal cycles
$\text{E}_{\text{MEM}}$: memory read energy per access

# Timing and power results on real applications

- Power consumption: average post-synthesis simulations on the first 10 testing samples

## Number of cycles

|  | Ideal | LRADNN (direct) | LRADNN |
|---|---|---|---|
| Caltech 101 silhouettes | 28327 | 29840 | 23141 |
| MNIST | 50938 | 52971 | 30105 |
| SVHN | 64586 | 66245 | 49371 |

## Power consumption (mW) / Energy consumption (mJ)

|  | Ideal | LRADNN (direct) | LRADNN |
|---|---|---|---|
| Caltech 101 silhouettes | 517.76 / 0.15 | 551.61 / 0.16 | 487.88 / 0.11 |
| MNIST | 517.76 / 0.26 | 557.98 / 0.30 | 459.73 / 0.14 |
| SVHN | 517.76 / 0.33 | 561.42 / 0.37 | 438.37 / 0.22 |

# Scalability for high parallelism of LRADNN

- Not fully utilize the hardware (multipliers) due to the memory alignment

$$\mathrm{U} = \frac{m \times n}{m \times p \times \lceil n/p \rceil}$$

**Hardware utilization under different parallelsims**

# Outline

- Introduction
- Related work
- Motivation of LRADNN
- Low rank approximation (LRA) predictor
- LRADNN: hardware architecture
- Experiment results
- Conclusion

# Conclusion

- A general hardware accelerator LRADNN for DNN is proposed

- A time and power-saving accelerator based on LRA
  - 31% ~ 53% energy reduction
  - 22% ~ 43% throughput increase

- A better scheme compared to the existing work

- A better scheme compared to the existing work

|  | AxNN | LRADNN |
|---|---|---|
| Prediction loss | 0.5% | < 0.1% |
| Energy improvement (w/o memory) | 1.14x – 1.92x | 1.18x – 1.61x |
| Energy improvement (w/ memory) | N.A. | 1.45x – 2.13x |

- Thank you

- Q & A