# Neural Acceleration for GPU Throughput Processors

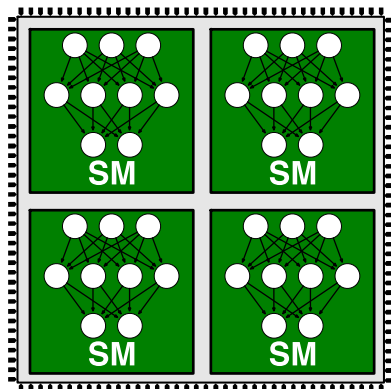**Amir Yazdanbakhsh**        Jongse Park        Hardik Sharma

Pejman Lotfi-Kamran*        Hadi Esmaeilzadeh

**Alternative Computing Technologies (ACT) Lab**
**Georgia Institute of Technology**

***The Institute for Research in Fundamental Sciences**

**NGPU**
**Neurally Accelerated GPU**
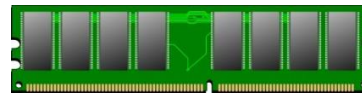
# Approximate computing
## Embracing imprecision

**Relax** the abstraction of "*near perfect*" **accuracy** in



Data Processing          Storage          Communication

Accept **imprecision** to improve

**performance**

**energy dissipation**
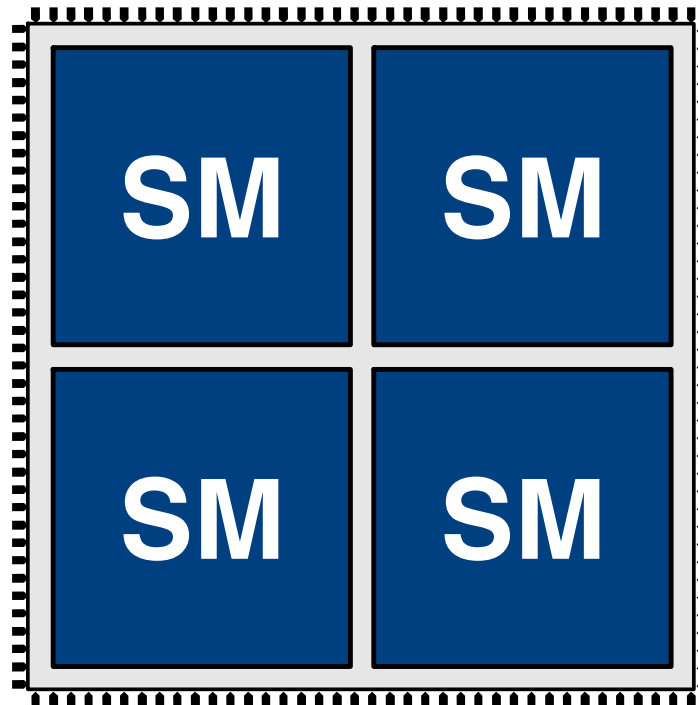
resource utilization **efficiency**

# Opportunity

Many GPU applications are amenable to approximation

Augmented Reality
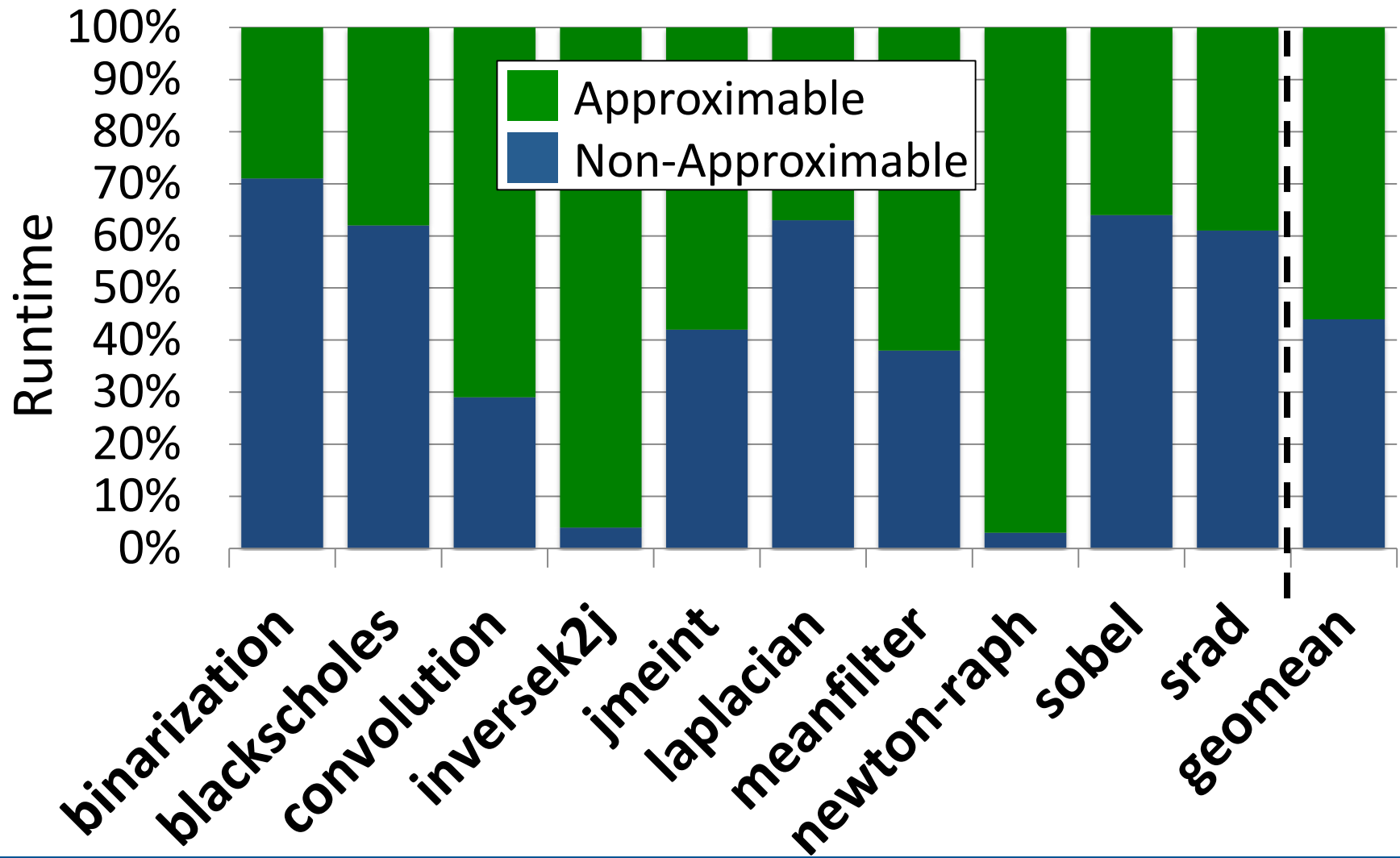
Computer Vision

Robotics

SM  SM

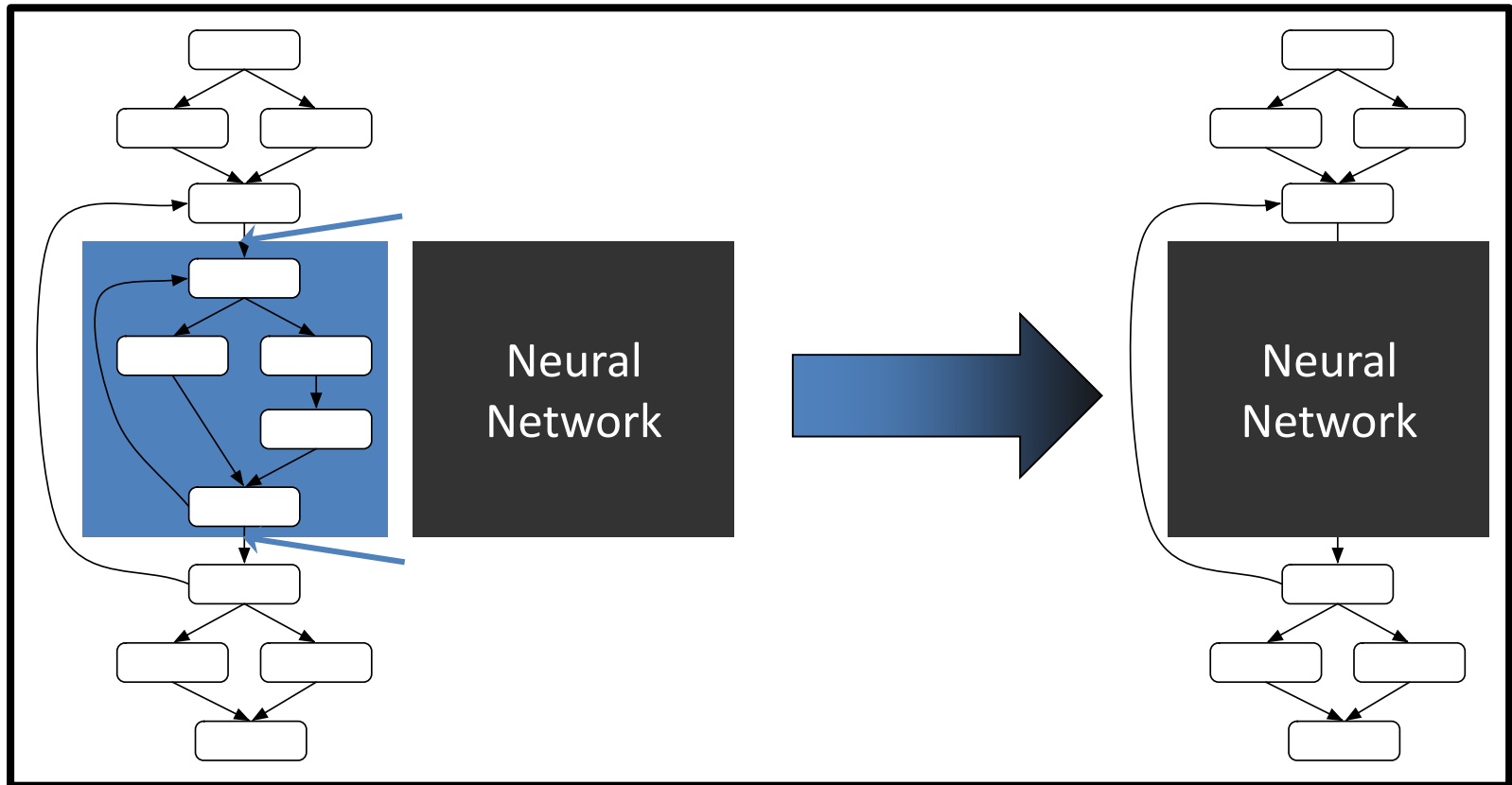SM  SM

Machine Learning

Sensor Processing

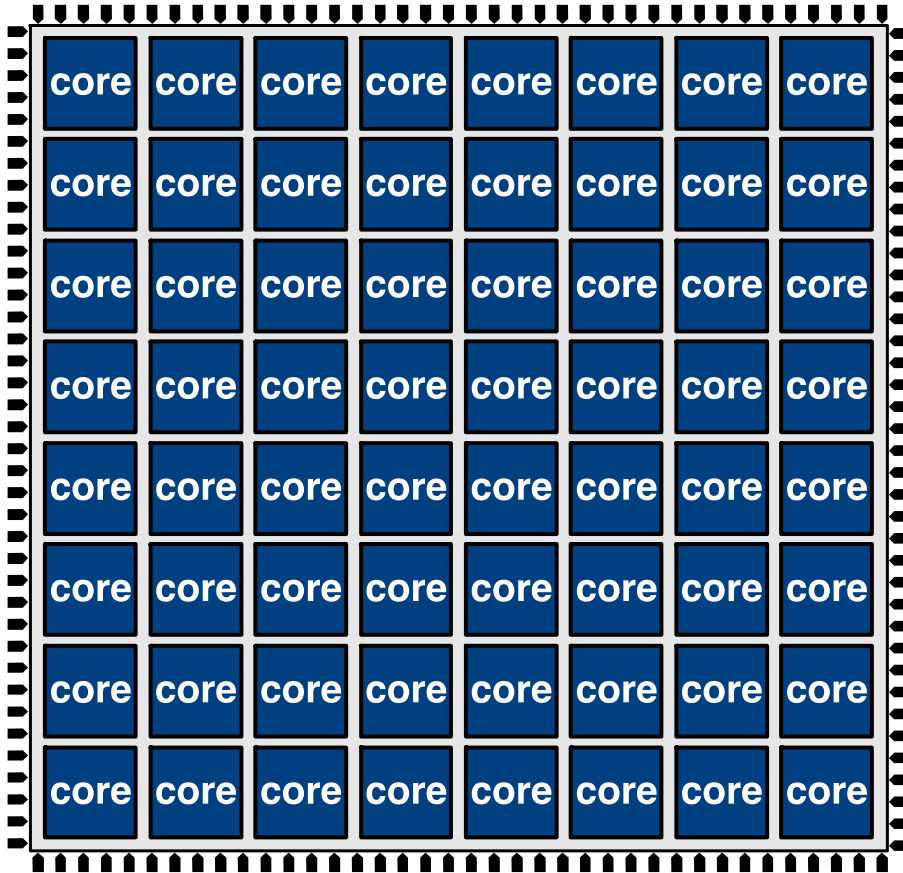Multimedia

# Opportunity



More than **55%** of application **runtime** and **energy** is in **neurally approximable regions**

# Neural Transformation for GPUs

# Challenges

# Challenges

# Challenges

# Challenges



Augmenting the CPU based neural processing units to each SIMD lane imposes **31.2%** area overhead

# NGPU

## Neurally-Accelerated GPU Architecture

# Neuronal Network Operations

$x_{j,0}$ $\cdots$ $x_{j,i}$ $\cdots$ $x_{j,n}$

$w_{j,0}$ $w_{j,i}$ $w_{j,n}$

$y_j$

$y_j =$

$sigmoid($

$w_{j,0} \times x_{j,0} +$

$\cdots$

$w_{j,i} \times x_{j,i} +$

$\cdots$

$w_{j,n} \times x_{j,n} +$

$)$

# NGPU

## Neurally-Accelerated GPU Architecture

# NGPU

## Neurally-Accelerated GPU Architecture



**NGPU** reuses the existing ALU in each SIMD lane

# NGPU

**Neurally-Accelerated GPU Architecture**

**Weight FIFO** is shared among all the SIMD lanes

# **NGPU** Execution Model

ld.global %r0, [addr0];
ld.global %r1, [addr1];
**send.n_data %r0;**
**send.n_data %r1;**
**recv.n_data %r2;**
st.global [addr2], %r2;

**Neurally Accelerated GPU Application**



**Neural Network**

# **NGPU** Execution Model

```
ld.global %r0, [addr0];
ld.global %r1, [addr1];
send.n_data %r0;
send.n_data %r1;
recv.n_data %r2;




st.global [addr2], %r2;
```

$( in_0, in_0, ..., in_0)$
$( in_1, in_1, ..., in_1)$
$w_0 \times ( in_0, in_0, ..., in_0)$
$+ w_2 \times ( in_1, in_1, ..., in_1)$
sigmoid
$w_1 \times ( in_0, in_0, ..., in_0)$
$+ w_3 \times ( in_1, in_1, ..., in_1)$
sigmoid
$w_4 \times ( n_0, n_0, ..., n_0)$
$+ w_5 \times ( n_1, n_1, ..., n_1)$
sigmoid
$( out_0, out_0, ..., out_0)$

$in_0 (\%r_0)$    $in_1 (\%r_1)$

$w_0$  $w_1$    $w_2$  $w_3$

$n_0$    $n_1$

$w_4$    $w_5$

$n_2$

$out_0 (\%r_2)$

# **NGPU** Execution Model

```
ld.global %r0, [addr0];
ld.global %r1, [addr1];
send.n_data %r0;
send.n_data %r1;
recv.n_data %r2;




st.global [addr2], %r2;
```

$( in_0, in_0, …, in_0)$
$( in_1, in_1, …, in_1)$
$w_0 \times ( in_0, in_0, …, in_0)$
$+ w_2 \times ( in_1, in_1, …, in_1)$
sigmoid
$w_1 \times ( in_0, in_0, …, in_0)$
$+ w_3 \times ( in_1, in_1, …, in_1)$
sigmoid
$w_4 \times ( n_0, n_0, …, n_0)$
$+ w_5 \times ( n_1, n_1, …, n_1)$
sigmoid
$( out_0, out_0, …, out_0)$



**SIMD lanes are in normal mode and performs precise computation**

# **NGPU** Execution Model

```
ld.global %r0, [addr0];
ld.global %r1, [addr1];
send.n_data %r0;
send.n_data %r1;
recv.n_data %r2;




st.global [addr2], %r2;
```



$( in_0, in_0, …, in_0)$
$( in_1, in_1, …, in_1)$
$w_0 \times \quad ( in_0, in_0, …, in_0)$
$+ w_2 \times \quad ( in_1, in_1, …, in_1)$
sigmoid
$w_1 \times \quad ( in_0, in_0, …, in_0)$
$+ w_3 \times \quad ( in_1, in_1, …, in_1)$
sigmoid
$w_4 \times \quad ( n_0, n_0, …, n_0)$
$+ w_5 \times \quad ( n_1, n_1, …, n_1)$
sigmoid
$( out_0, out_0, …, out_0)$

$in_0 \ (\%r_0)$ $\qquad$ $in_1 \ (\%r_1)$

$w_0 \quad w_1 \qquad w_2 \quad w_3$

$n_0 \qquad n_1$

$w_4 \qquad w_5$

$n_2$

$out_0 \ (\%r_2)$

**SIMD lanes enter neural mode**
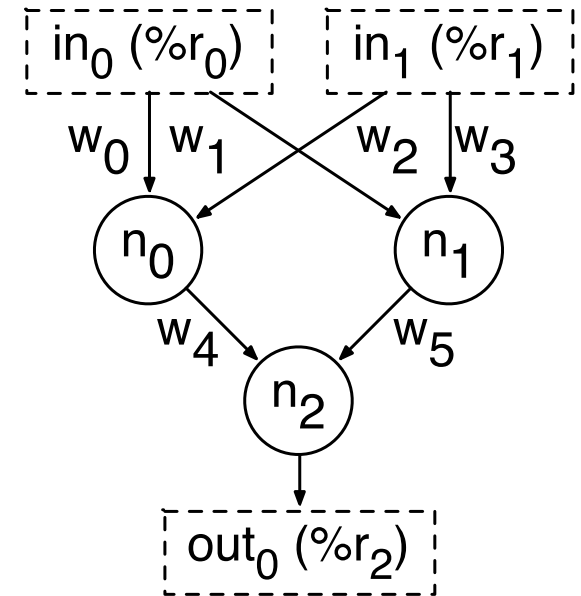
# **NGPU** Execution Model

```
ld.global %r0, [addr0];
ld.global %r1, [addr1];
send.n_data %r0;
send.n_data %r1;
recv.n_data %r2;



st.global [addr2], %r2;
```

$( in_0, in_0, …, in_0)$
$( in_1, in_1, …, in_1)$
$w_0 \times ( in_0, in_0, …, in_0)$
$+ w_2 \times ( in_1, in_1, …, in_1)$
sigmoid
$w_1 \times ( in_0, in_0, …, in_0)$
$+ w_3 \times ( in_1, in_1, …, in_1)$
sigmoid
$w_4 \times ( n_0, n_0, …, n_0)$
$+ w_5 \times ( n_1, n_1, …, n_1)$
sigmoid
$( out_0, out_0, …, out_0)$

$in_0 (\%r_0)$   $in_1 (\%r_1)$

$w_0$  $w_1$   $w_2$  $w_3$

$n_0$   $n_1$

$w_4$   $w_5$

$n_2$

$out_0 (\%r_2)$

**SIMD starts the calculation of the neural network**

# **NGPU** Execution Model

```
ld.global %r0, [addr0];
ld.global %r1, [addr1];
send.n_data %r0;
send.n_data %r1;
recv.n_data %r2;



st.global [addr2], %r2;
```
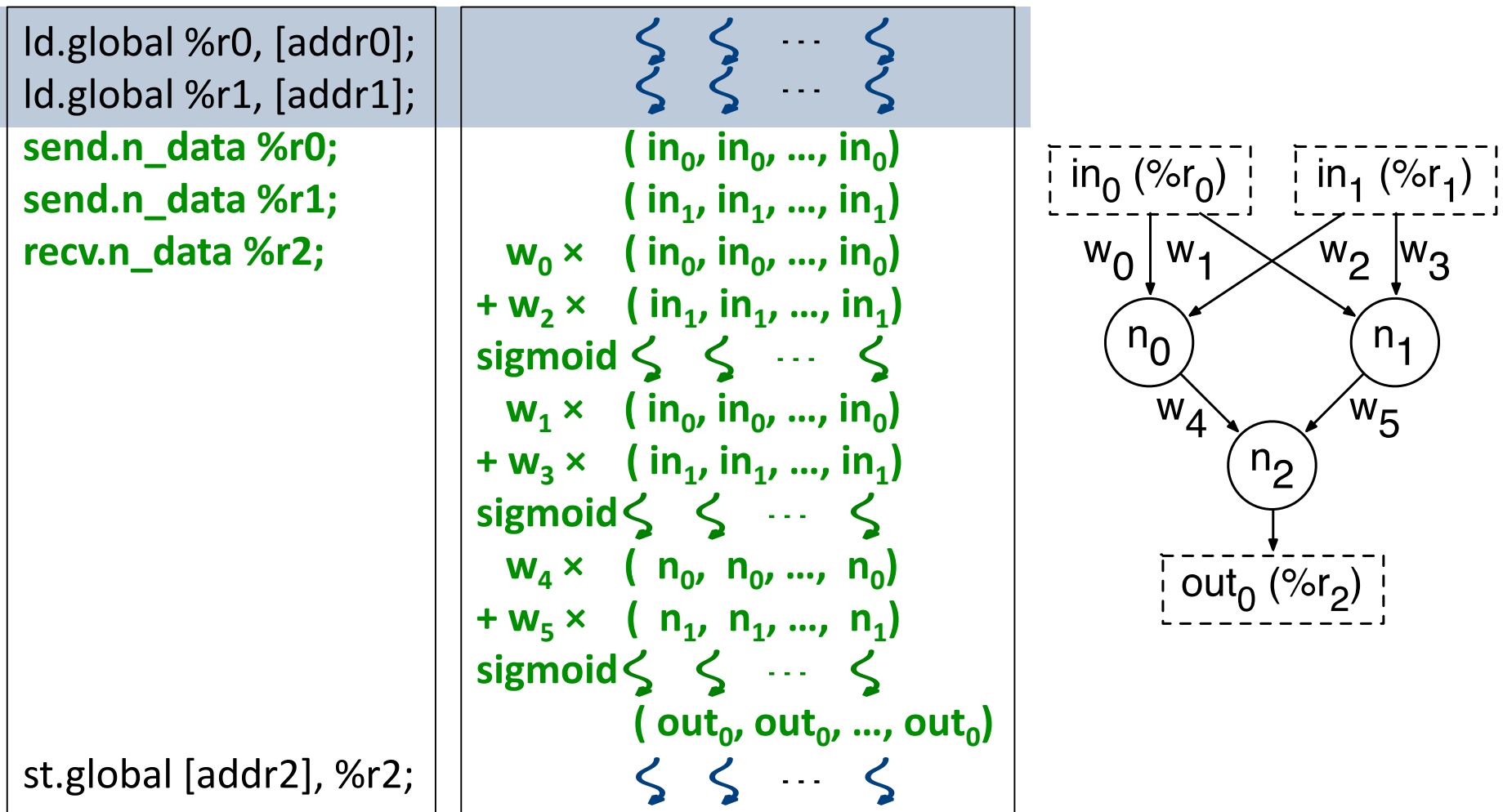
$( in_0, in_0, ..., in_0)$
$( in_1, in_1, ..., in_1)$
$w_0 \times \quad ( in_0, in_0, ..., in_0)$
$+ w_2 \times \quad ( in_1, in_1, ..., in_1)$
sigmoid
$w_1 \times \quad ( in_0, in_0, ..., in_0)$
$+ w_3 \times \quad ( in_1, in_1, ..., in_1)$
sigmoid
$w_4 \times \quad ( n_0, n_0, ..., n_0)$
$+ w_5 \times \quad ( n_1, n_1, ..., n_1)$
sigmoid
$( out_0, out_0, ..., out_0)$

$in_0 (\%r_0)$    $in_1 (\%r_1)$

$w_0 \quad w_1 \qquad w_2 \quad w_3$

$n_0 \qquad n_1$

$w_4 \qquad w_5$

$n_2$

$out_0 (\%r_2)$

**The neurally accelerated SIMD lanes autonomously
calculate the neural outputs in lock-step**

# **NGPU** Execution Model

```
ld.global %r0, [addr0];
ld.global %r1, [addr1];
send.n_data %r0;
send.n_data %r1;
recv.n_data %r2;



st.global [addr2], %r2;
```



$( in_0, in_0, ..., in_0)$
$( in_1, in_1, ..., in_1)$
$w_0 \times \quad ( in_0, in_0, ..., in_0)$
$+ w_2 \times \quad ( in_1, in_1, ..., in_1)$
sigmoid
$w_1 \times \quad ( in_0, in_0, ..., in_0)$
$+ w_3 \times \quad ( in_1, in_1, ..., in_1)$
sigmoid
$w_4 \times \quad ( n_0, n_0, ..., n_0)$
$+ w_5 \times \quad ( n_1, n_1, ..., n_1)$
sigmoid
$( out_0, out_0, ..., out_0)$

$in_0 (\%r_0)$ $\qquad$ $in_1 (\%r_1)$

$w_0$ $w_1$ $\qquad$ $w_2$ $w_3$

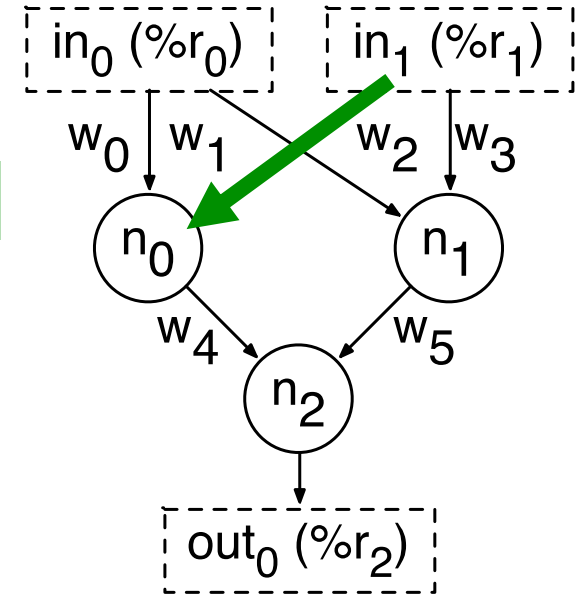$n_0$ $\qquad$ $n_1$

$w_4$ $\qquad$ $w_5$

$n_2$

$out_0 (\%r_2)$

**The accelerated SIMD lanes autonomously calculate the neural outputs in lock-step**

# **NGPU** Execution Model

```
ld.global %r0, [addr0];
ld.global %r1, [addr1];
send.n_data %r0;
send.n_data %r1;
recv.n_data %r2;



st.global [addr2], %r2;
```

$$( in_0, in_0, ..., in_0)$$
$$( in_1, in_1, ..., in_1)$$
$$w_0 \times ( in_0, in_0, ..., in_0)$$
$$+ w_2 \times ( in_1, in_1, ..., in_1)$$
$$sigmoid$$
$$w_1 \times ( in_0, in_0, ..., in_0)$$
$$+ w_3 \times ( in_1, in_1, ..., in_1)$$
$$sigmoid$$
$$w_4 \times ( n_0, n_0, ..., n_0)$$
$$+ w_5 \times ( n_1, n_1, ..., n_1)$$
$$sigmoid$$
$$( out_0, out_0, ..., out_0)$$

$in_0$ (%$r_0$)     $in_1$ (%$r_1$)

$w_0$   $w_1$     $w_2$   $w_3$

$n_0$     $n_1$

$w_4$     $w_5$

$n_2$

$out_0$ (%$r_2$)

**The accelerated SIMD lanes autonomously calculate the neural outputs in lock-step**

# **NGPU** Execution Model

```
ld.global %r0, [addr0];
ld.global %r1, [addr1];
send.n_data %r0;
send.n_data %r1;
recv.n_data %r2;



st.global [addr2], %r2;
```

$( in_0, in_0, ..., in_0)$
$( in_1, in_1, ..., in_1)$
$w_0 \times \quad ( in_0, in_0, ..., in_0)$
$+ w_2 \times \quad ( in_1, in_1, ..., in_1)$
sigmoid
$w_1 \times \quad ( in_0, in_0, ..., in_0)$
$+ w_3 \times \quad ( in_1, in_1, ..., in_1)$
sigmoid
$w_4 \times \quad ( n_0, n_0, ..., n_0)$
$+ w_5 \times \quad ( n_1, n_1, ..., n_1)$
sigmoid
$( out_0, out_0, ..., out_0)$

$in_0 (\%r_0)$   $in_1 (\%r_1)$

$w_0$ | $w_1$    $w_2$ | $w_3$

$n_0$     $n_1$

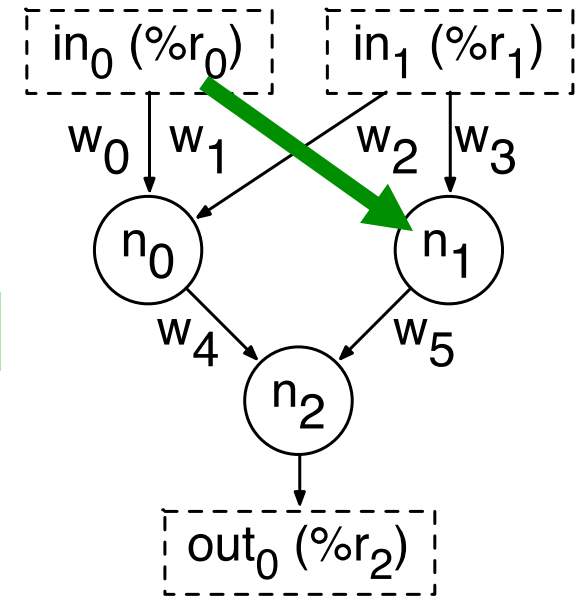$w_4$       $w_5$

$n_2$

$out_0 (\%r_2)$

**The accelerated SIMD lanes autonomously calculate the neural outputs in lock-step**

# **NGPU** Execution Model

```
ld.global %r0, [addr0];
ld.global %r1, [addr1];
send.n_data %r0;
send.n_data %r1;
recv.n_data %r2;



st.global [addr2], %r2;
```

$( in_0, in_0, ..., in_0)$
$( in_1, in_1, ..., in_1)$
$w_0 \times \quad ( in_0, in_0, ..., in_0)$
$+ w_2 \times \quad ( in_1, in_1, ..., in_1)$
sigmoid
$w_1 \times \quad ( in_0, in_0, ..., in_0)$
$+ w_3 \times \quad ( in_1, in_1, ..., in_1)$
sigmoid
$w_4 \times \quad ( n_0, n_0, ..., n_0)$
$+ w_5 \times \quad ( n_1, n_1, ..., n_1)$
sigmoid
$( out_0, out_0, ..., out_0)$

$in_0 (\%r_0)$    $in_1 (\%r_1)$

$w_0 \quad w_1 \qquad w_2 \quad w_3$

$n_0 \qquad\qquad n_1$

$w_4 \qquad\qquad w_5$
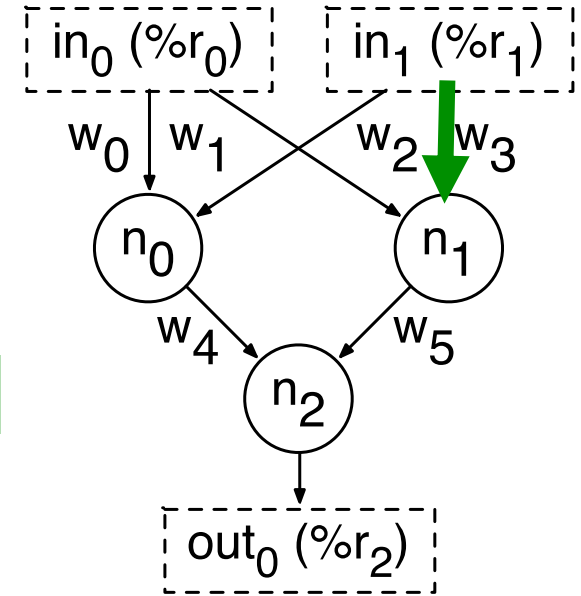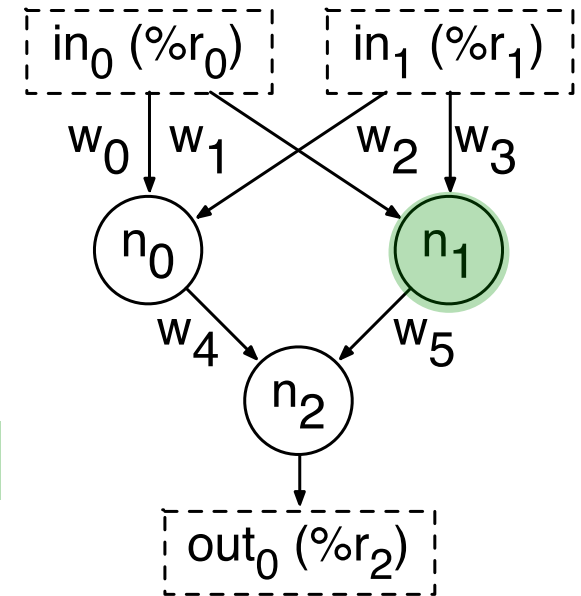
$n_2$

$out_0 (\%r_2)$

**The accelerated SIMD lanes autonomously calculate the neural outputs in lock-step**

# **NGPU** Execution Model

```
ld.global %r0, [addr0];
ld.global %r1, [addr1];
send.n_data %r0;
send.n_data %r1;
recv.n_data %r2;



st.global [addr2], %r2;
```

$( in_0, in_0, ..., in_0)$
$( in_1, in_1, ..., in_1)$
$w_0 \times \quad ( in_0, in_0, ..., in_0)$
$+ w_2 \times \quad ( in_1, in_1, ..., in_1)$
sigmoid
$w_1 \times \quad ( in_0, in_0, ..., in_0)$
$+ w_3 \times \quad ( in_1, in_1, ..., in_1)$
sigmoid
$w_4 \times \quad ( n_0, n_0, ..., n_0)$
$+ w_5 \times \quad ( n_1, n_1, ..., n_1)$
sigmoid
$( out_0, out_0, ..., out_0)$

$in_0 (\%r_0)$   $in_1 (\%r_1)$
$w_0 \quad w_1$   $w_2 \quad w_3$
$n_0$   $n_1$
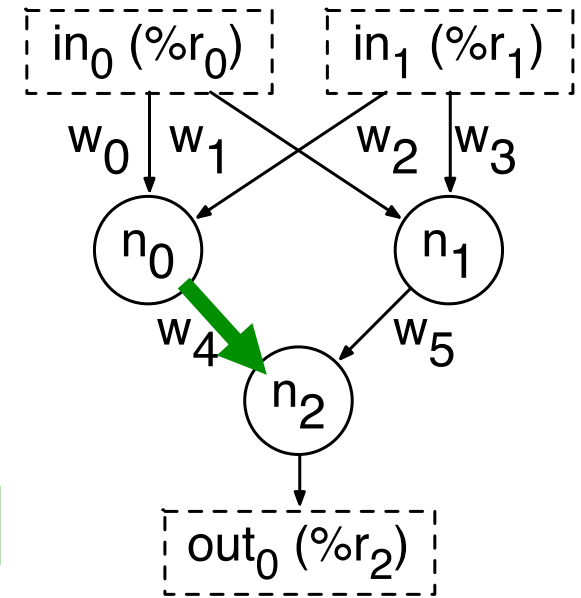$w_4$   $w_5$
$n_2$
$out_0 (\%r_2)$

**The accelerated SIMD lanes autonomously
calculate the neural outputs in lock-step**

# **NGPU** Execution Model

```
ld.global %r0, [addr0];
ld.global %r1, [addr1];
send.n_data %r0;
send.n_data %r1;
recv.n_data %r2;




st.global [addr2], %r2;
```

$$( in_0, in_0, ..., in_0)$$
$$( in_1, in_1, ..., in_1)$$
$$w_0 \times ( in_0, in_0, ..., in_0)$$
$$+ w_2 \times ( in_1, in_1, ..., in_1)$$
sigmoid
$$w_1 \times ( in_0, in_0, ..., in_0)$$
$$+ w_3 \times ( in_1, in_1, ..., in_1)$$
sigmoid
$$w_4 \times ( n_0, n_0, ..., n_0)$$
$$+ w_5 \times ( n_1, n_1, ..., n_1)$$
sigmoid
$$( out_0, out_0, ..., out_0)$$

$in_0 (\%r_0)$    $in_1 (\%r_1)$

$w_0$  $w_1$    $w_2$  $w_3$

$n_0$        $n_1$

$w_4$        $w_5$
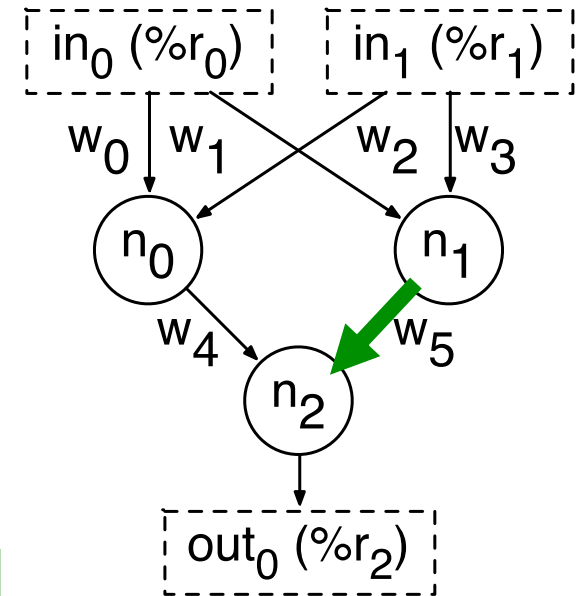
$n_2$

$out_0 (\%r_2)$

**The accelerated SIMD lanes autonomously calculate the neural outputs in lock-step**

# **NGPU** Execution Model

```
ld.global %r0, [addr0];
ld.global %r1, [addr1];
send.n_data %r0;
send.n_data %r1;
recv.n_data %r2;



st.global [addr2], %r2;
```

$( in_0, in_0, ..., in_0)$
$( in_1, in_1, ..., in_1)$
$w_0 \times \quad ( in_0, in_0, ..., in_0)$
$+ w_2 \times \quad ( in_1, in_1, ..., in_1)$
sigmoid
$w_1 \times \quad ( in_0, in_0, ..., in_0)$
$+ w_3 \times \quad ( in_1, in_1, ..., in_1)$
sigmoid
$w_4 \times \quad ( n_0, n_0, ..., n_0)$
$+ w_5 \times \quad ( n_1, n_1, ..., n_1)$
sigmoid
$( out_0, out_0, ..., out_0)$

$in_0 (\%r_0)$  $in_1 (\%r_1)$
$w_0$  $w_1$  $w_2$  $w_3$
$n_0$  $n_1$
$w_4$  $w_5$
$n_2$
$out_0 (\%r_2)$

**The accelerated SIMD lanes autonomously
calculate the neural outputs in lock-step**

# **NGPU** Execution Model



ld.global %r0, [addr0];
ld.global %r1, [addr1];
**send.n_data %r0;**
**send.n_data %r1;**
**recv.n_data %r2;**

st.global [addr2], %r2;

$( in_0, in_0, ..., in_0)$
$( in_1, in_1, ..., in_1)$
$w_0 \times \quad ( in_0, in_0, ..., in_0)$
$+ w_2 \times \quad ( in_1, in_1, ..., in_1)$
sigmoid
$w_1 \times \quad ( in_0, in_0, ..., in_0)$
$+ w_3 \times \quad ( in_1, in_1, ..., in_1)$
sigmoid
$w_4 \times \quad ( n_0, n_0, ..., n_0)$
$+ w_5 \times \quad ( n_1, n_1, ..., n_1)$
**sigmoid**
$( out_0, out_0, ..., out_0)$

$in_0 (\%r_0)$  $in_1 (\%r_1)$
$w_0$ $w_1$ $w_2$ $w_3$
$n_0$  $n_1$
$w_4$ $w_5$
$n_2$
$out_0 (\%r_2)$

**The accelerated SIMD lanes autonomously calculate the neural outputs in lock-step**
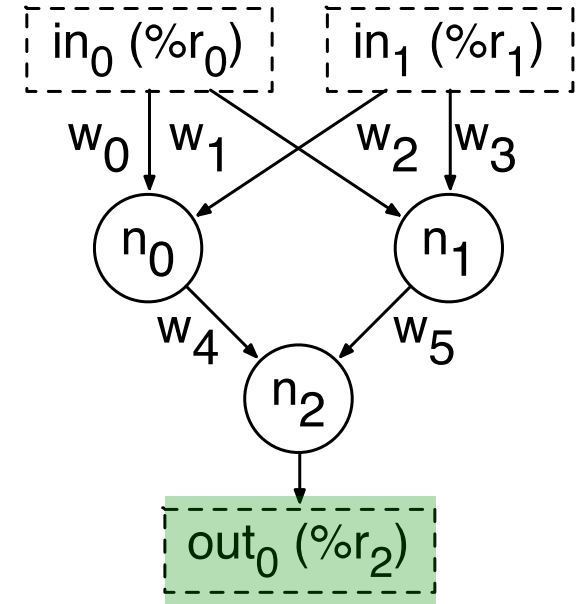
# **NGPU** Execution Model

```
ld.global %r0, [addr0];
ld.global %r1, [addr1];
send.n_data %r0;
send.n_data %r1;
recv.n_data %r2;



st.global [addr2], %r2;
```

$$( in_0, in_0, ..., in_0)$$
$$( in_1, in_1, ..., in_1)$$
$$w_0 \times ( in_0, in_0, ..., in_0)$$
$$+ w_2 \times ( in_1, in_1, ..., in_1)$$
sigmoid
$$w_1 \times ( in_0, in_0, ..., in_0)$$
$$+ w_3 \times ( in_1, in_1, ..., in_1)$$
sigmoid
$$w_4 \times ( n_0, n_0, ..., n_0)$$
$$+ w_5 \times ( n_1, n_1, ..., n_1)$$
sigmoid
$$( out_0, out_0, ..., out_0)$$

$in_0$ (%$r_0$)  $in_1$ (%$r_1$)

$w_0$ | $w_1$  $w_2$ | $w_3$

$n_0$  $n_1$

$w_4$  $w_5$

$n_2$

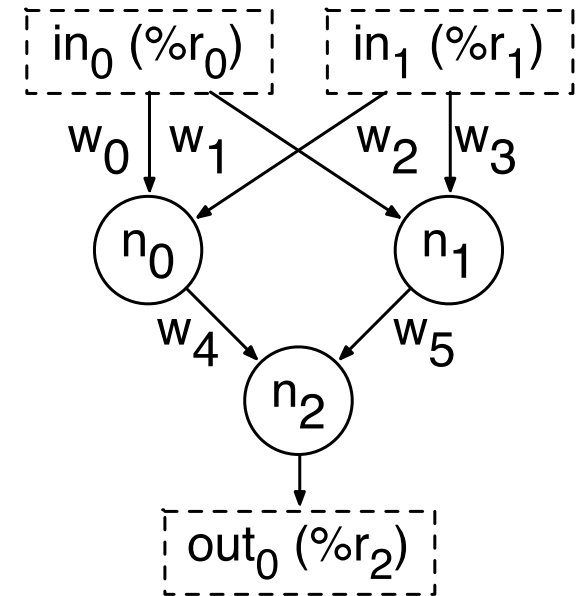$out_0$ (%$r_2$)

**SIMD lanes exit neural mode**

# **NGPU** Execution Model

```
ld.global %r0, [addr0];
ld.global %r1, [addr1];
send.n_data %r0;
send.n_data %r1;
recv.n_data %r2;
```

$( in_0, in_0, ..., in_0)$
$( in_1, in_1, ..., in_1)$
$w_0 \times ( in_0, in_0, ..., in_0)$
$+ w_2 \times ( in_1, in_1, ..., in_1)$
sigmoid
$w_1 \times ( in_0, in_0, ..., in_0)$
$+ w_3 \times ( in_1, in_1, ..., in_1)$
sigmoid
$w_4 \times ( n_0, n_0, ..., n_0)$
$+ w_5 \times ( n_1, n_1, ..., n_1)$
sigmoid
$( out_0, out_0, ..., out_0)$

```
st.global [addr2], %r2;
```

$in_0 (\%r_0)$   $in_1 (\%r_1)$

$w_0$ | $w_1$   $w_2$ | $w_3$

$n_0$   $n_1$

$w_4$   $w_5$

$n_2$

$out_0 (\%r_2)$

**SIMD lanes are in normal mode**

# Experimental Setup

## Machine Learning, Finance, Vision
## 3D Gaming, Medical Imaging
## Numerical Analysis, Image Processing
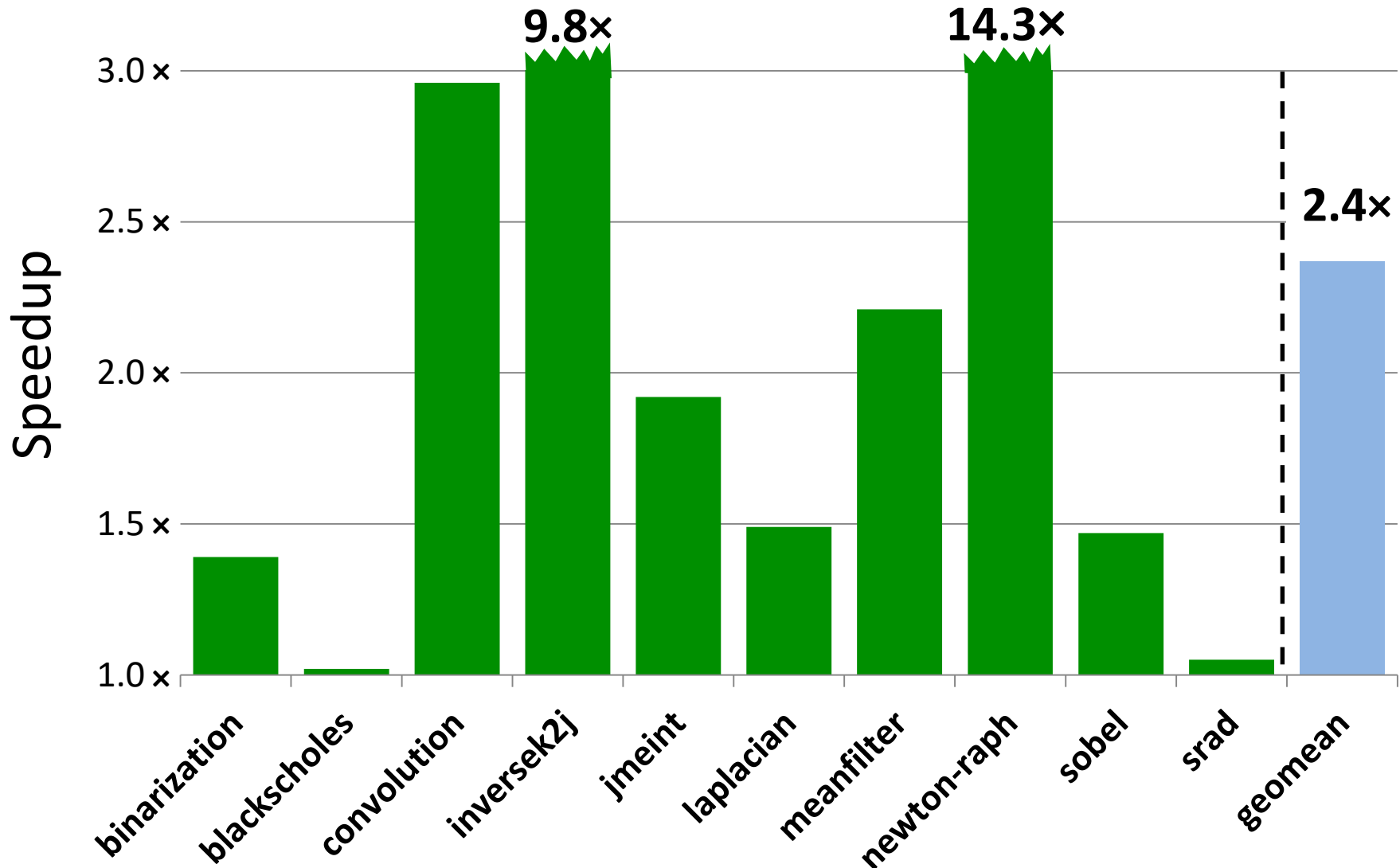
### GPU Simulator

- GPGPUSim Cycle-Level Simulator
- Fermi-based GTX 480, Shader Core Frequency 1.4 GHz
- NVCC Compiler –O3

### Power Model

- Technology Node 40 nm
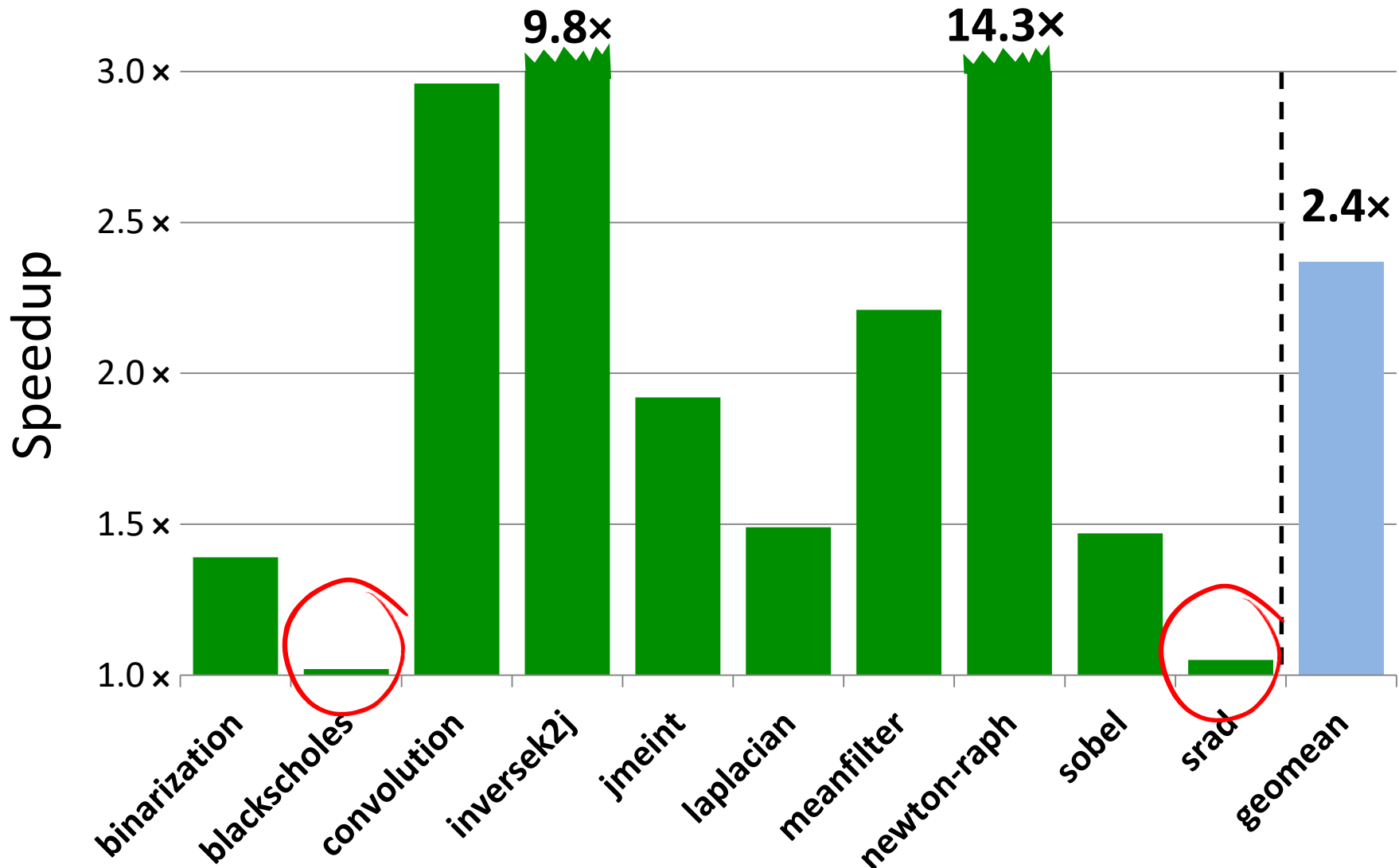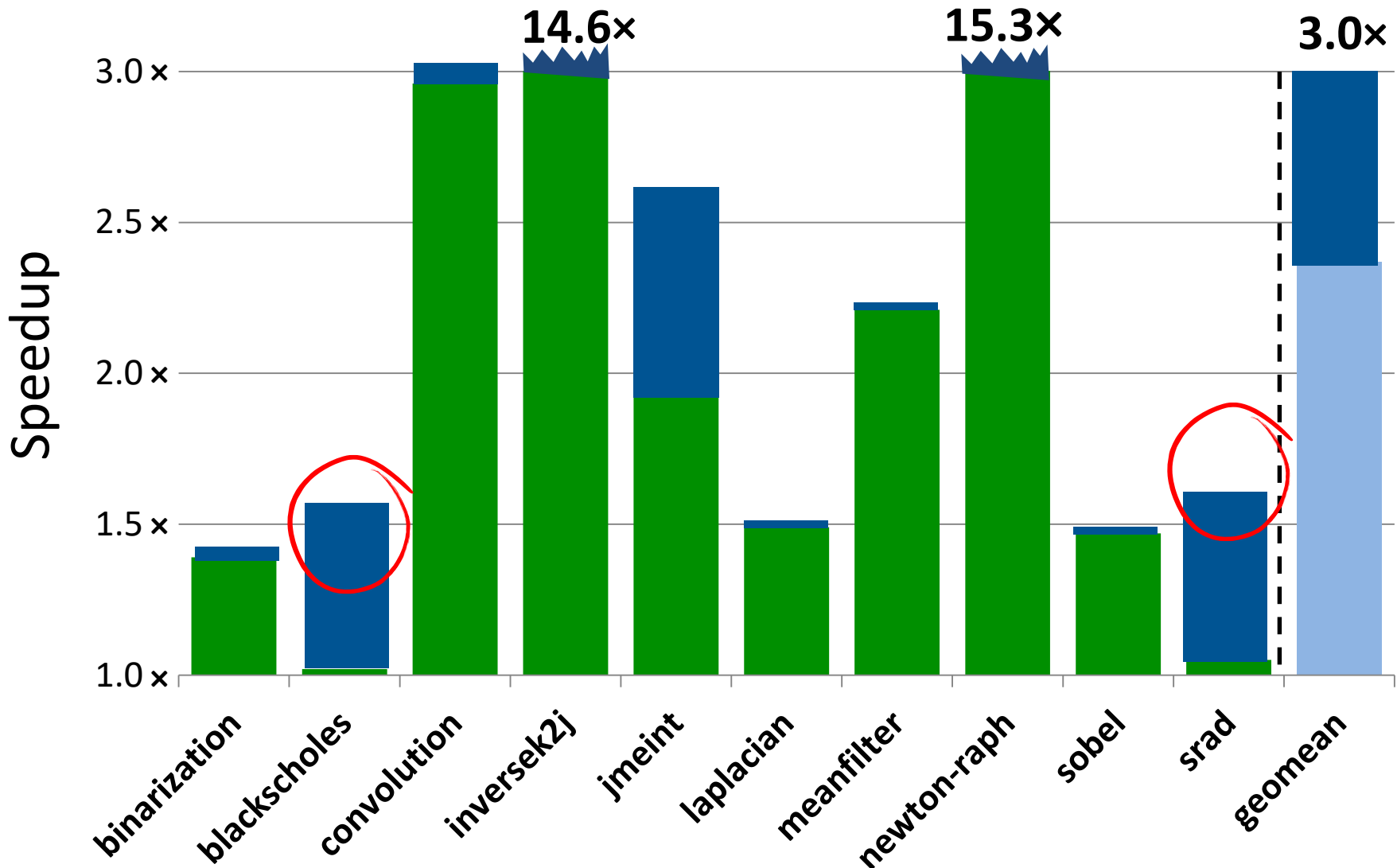- GPUWattch
- McPAT and CACTI, Verilog

# NGPU Speedup



Speedup bar chart showing:
- binarization: ~1.4×
- blackscholes: ~1.0×
- convolution: ~2.95×
- inversek2j: 9.8×
- jmeint: ~1.9×
- laplacian: ~1.5×
- meanfilter: ~2.2×
- newton-raph: 14.3×
- sobel: ~1.5×
- srad: ~1.05×
- geomean: 2.4×
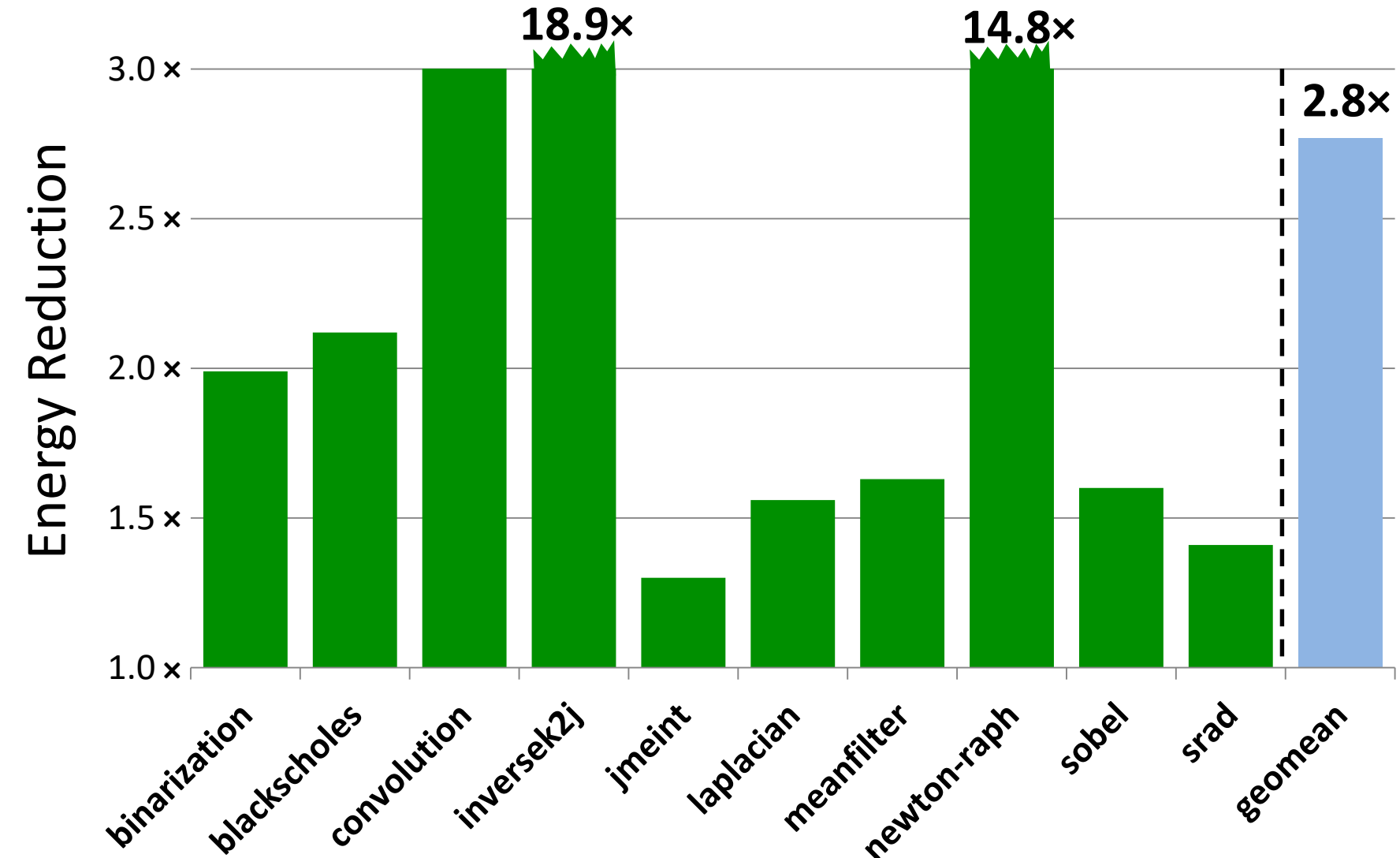
**Most applications** see **speedup** with **NGPU**

# NGPU Speedup

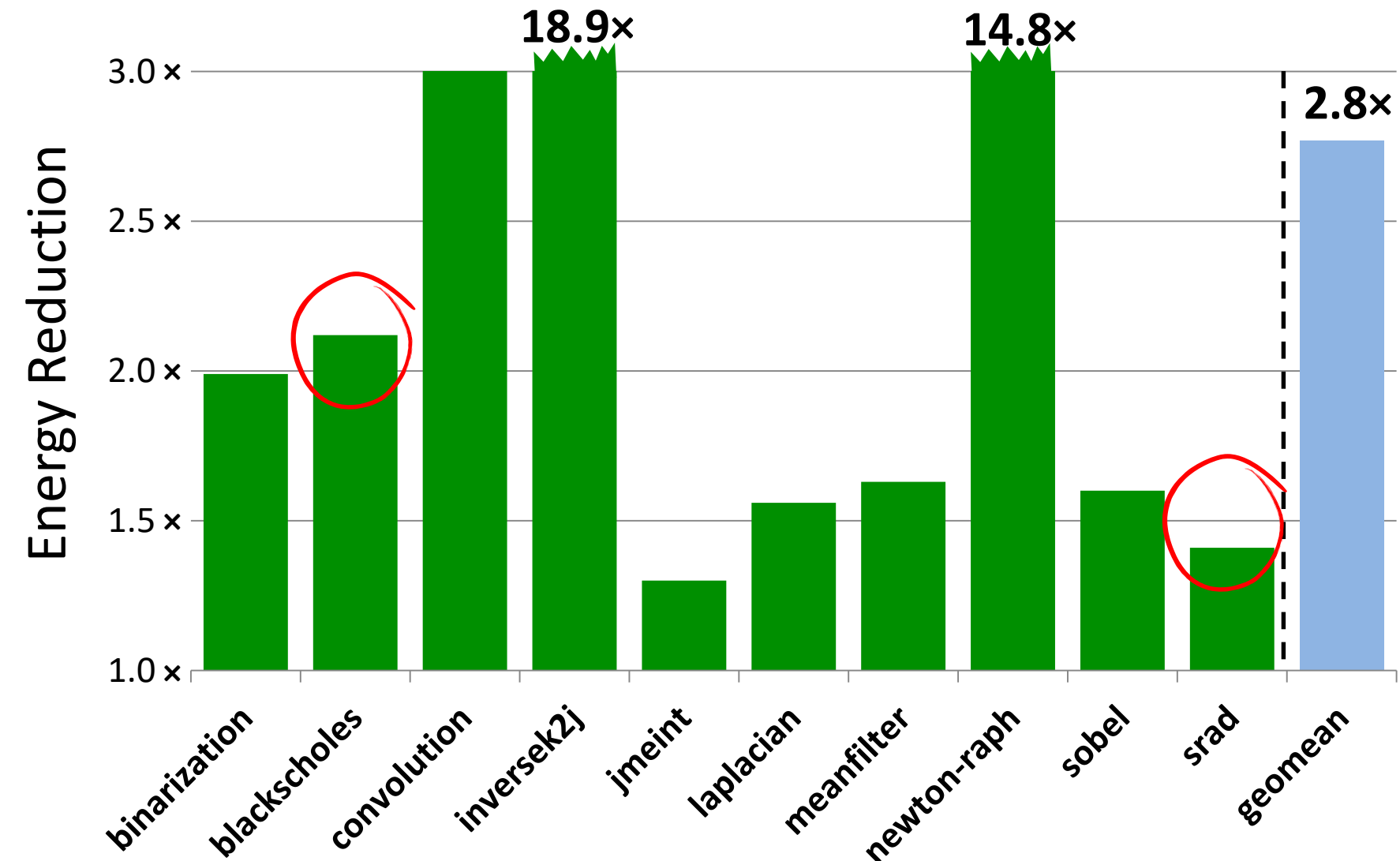The speedup for **bandwidth-sensitive** applications is limited

# NGPU Speedup with 2x Bandwidth

**Bandwidth-sensitive** applications see **speedup** with 2x bandwidth
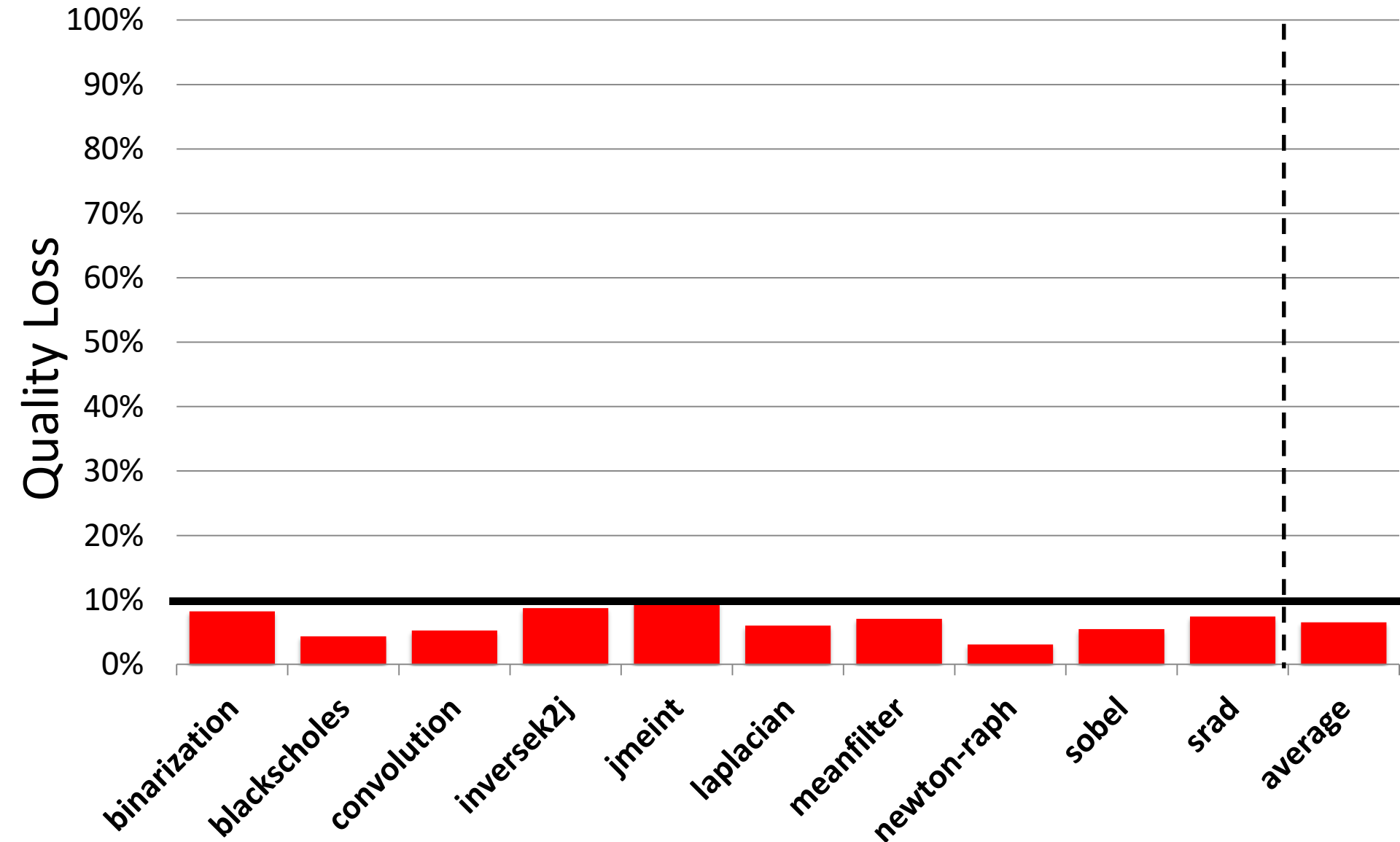
# NGPU Energy Savings with Baseline Bandwidth

**NGPU** eliminates the von Neumann overhead which results in energy reduction

# **NGPU** is a Fair Bargain

**Overhead**

| Area Overhead ≤ 1.0% | |
|---|---|
| Quality ≥ 97.5% | Quality ≥ 90.0% |

**Benefits**

| 1.9× Speedup  2.1× Energy Reduction | 2.4× Speedup  2.8× Energy Reduction |
|---|---|