

# Integrated Optical Neural Networks Exploiting Light Speed Approximate Parallel Multipliers

---

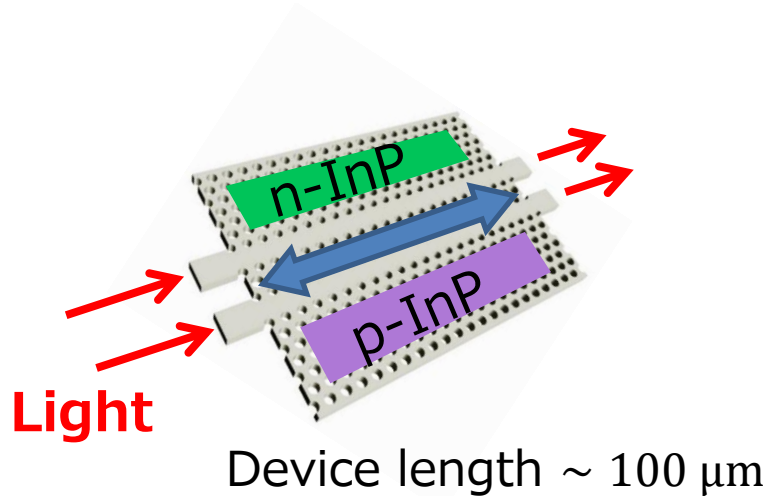
Jun Shiomi<sup>1</sup>, Tohru Ishihara<sup>2</sup>, Hidetoshi Onodera<sup>1</sup>,  
Akihiko Shinya<sup>3</sup>, Masaya Notomi<sup>3</sup>

<sup>1</sup>Graduate School of Informatics, Kyoto University, Japan

<sup>2</sup>Graduate School of Informatics, Nagoya University, Japan

<sup>2</sup>NTT Nanophotonics Center / NTT Basic Research Laboratories, Japan

# Integrated Optical Logic Circuits Using Nanophotonics



E.g. Photonic crystal-based directional coupler

✓ Direction control on the order of the light wavelength

➔ { ☺ On-chip implementation  
☺ **Ultra-low latency**

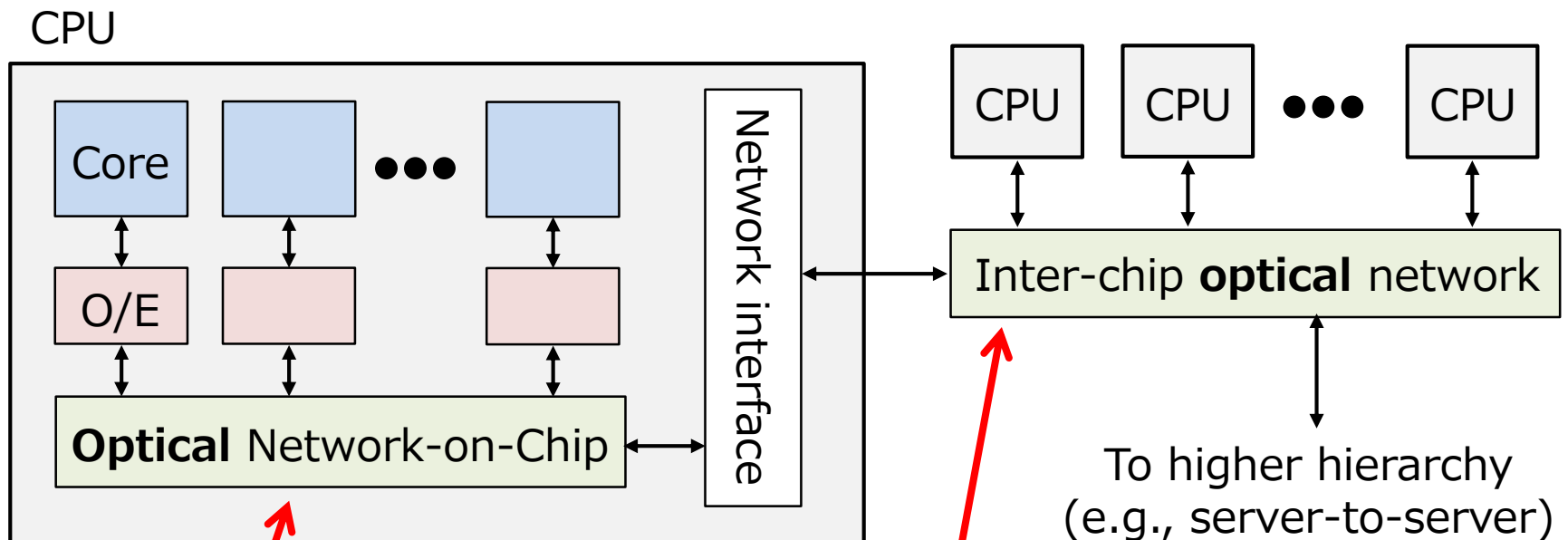
CMOS logic gate  $\sim 10 \text{ ps}$

Nanophotonics-based  $0.1 \text{ ps} \sim 1 \text{ ps}$

Goal: Ultra-fast optical logic circuit design

# Beyond Optical Communication

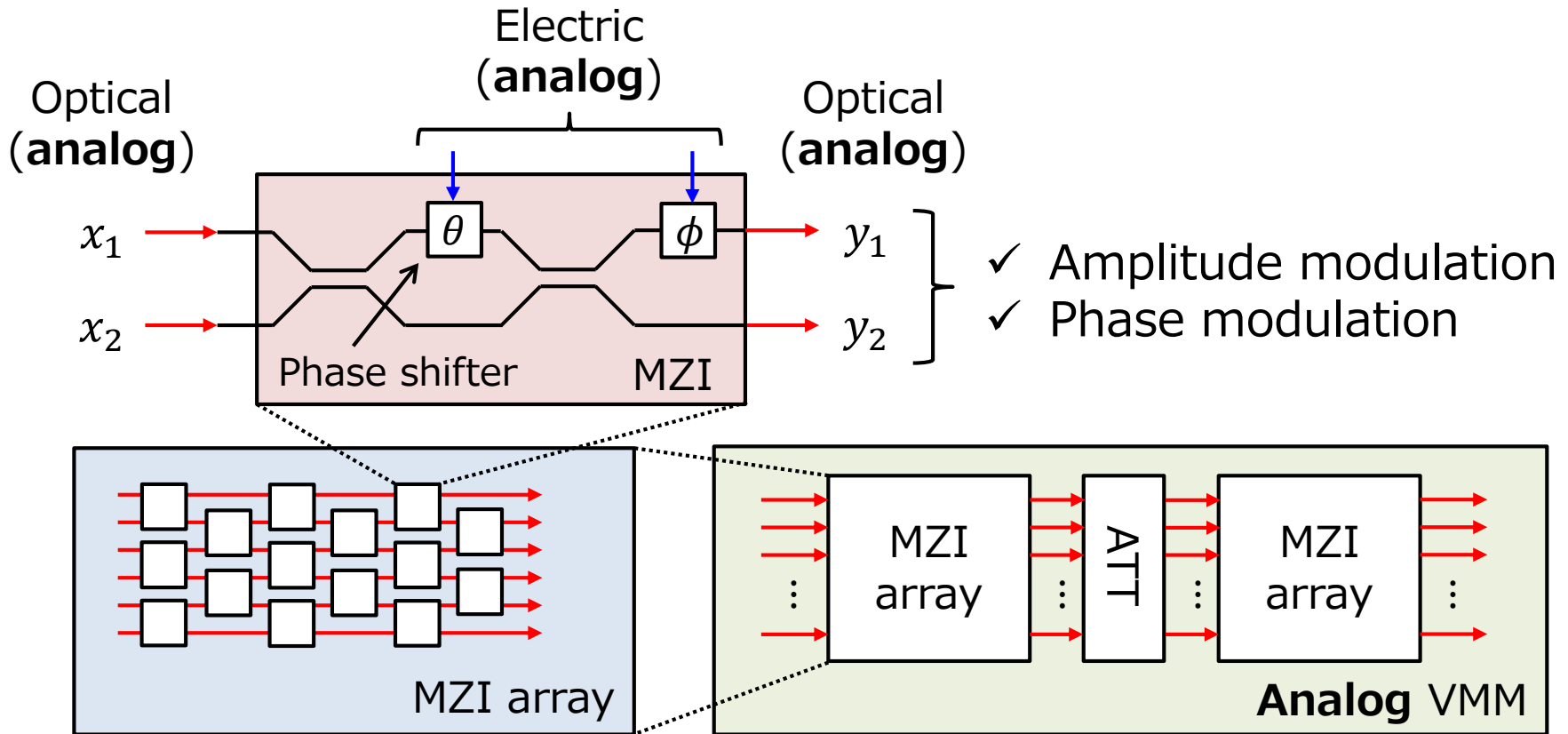
Existing technique's focus: on-chip optical communication



Goal: Add **functional unit** to boost up on-chip communication

➔ E.g. Integrated Neural Network (NN) for packet filtering

# Related Work: Vector Matrix Multiplication (VMM) Exploiting Mach-Zehnder Interferometer (MZI) [1]



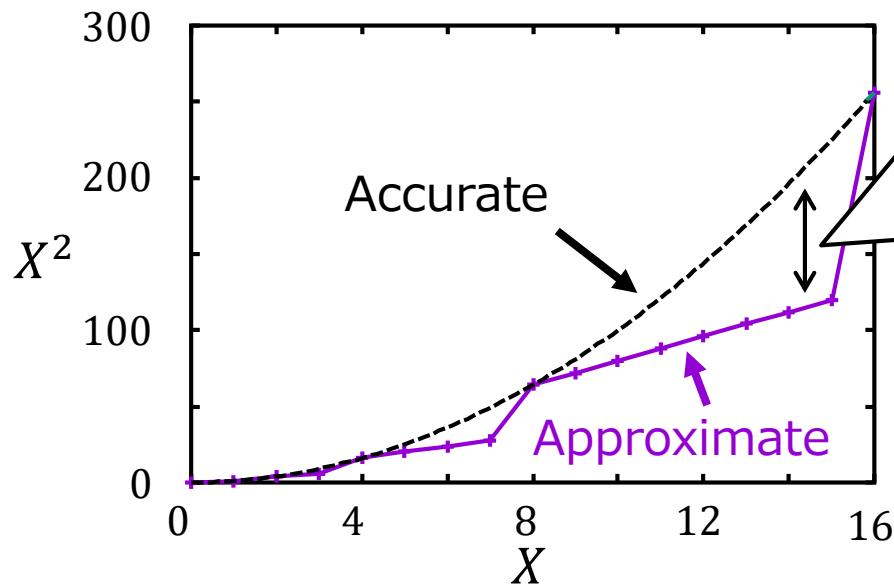
Challenges in analog multiplication:

Environment-dependent error (e.g. noise & temperature)

# Fully Digital Approximate Parallel Multiplier

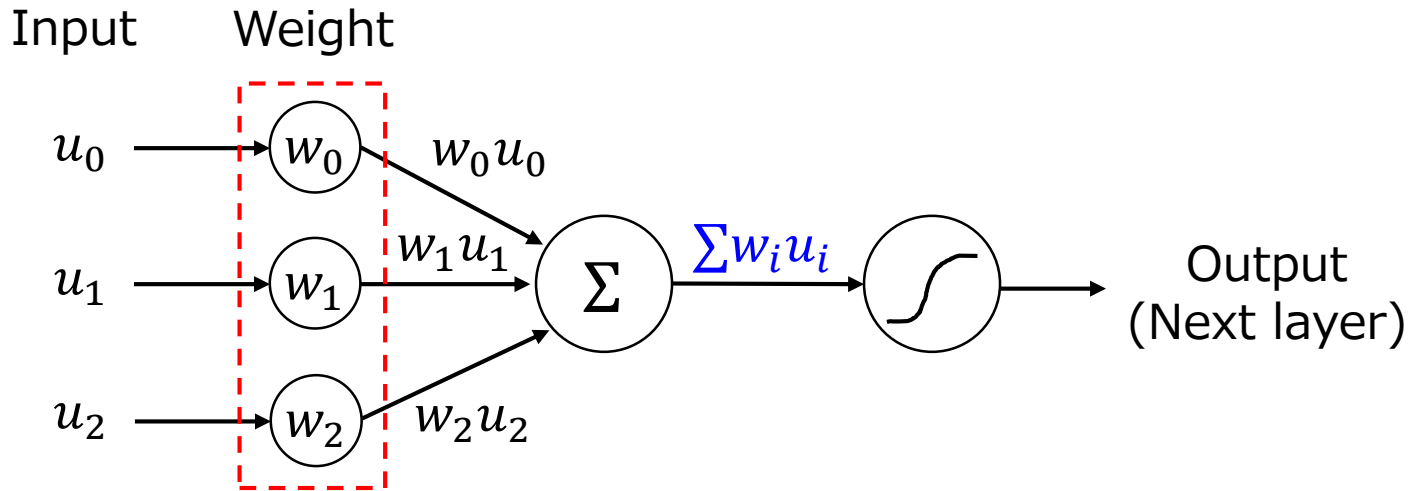
Accurate digital multiplier: larger delay than an analog multiplier

- ➔
- ✓ This work: **Approximate structure** for performance boost
  - ☹ Deterministic multiplication error



- ✓ Dependent on input value only
- ✓ Fixed error @ every calculation
  - Up to  $-50\%$  error

# Neural Network (NN) Exploiting Approximate Multiplier



Performance boost by **approximate multiplier**

- ✓ Pattern recognition: high resiliency to approximation errors
- ✓ Absorb the deterministic error by learning process

$$\frac{\mathbf{w}_0 u_0}{\uparrow} + \frac{\mathbf{w}_1 u_1}{\uparrow} + \frac{\mathbf{w}_2 u_2}{\uparrow}$$

Weight vector optimization to cancel out deterministic error

# Outline

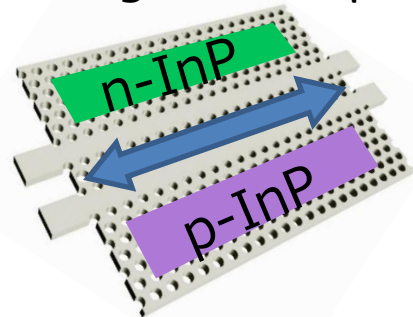
---

- Background
- Parallel Multiplier Using Nanophotonic Devices
- Approximate Parallel Multiplier for NN
- Conclusion

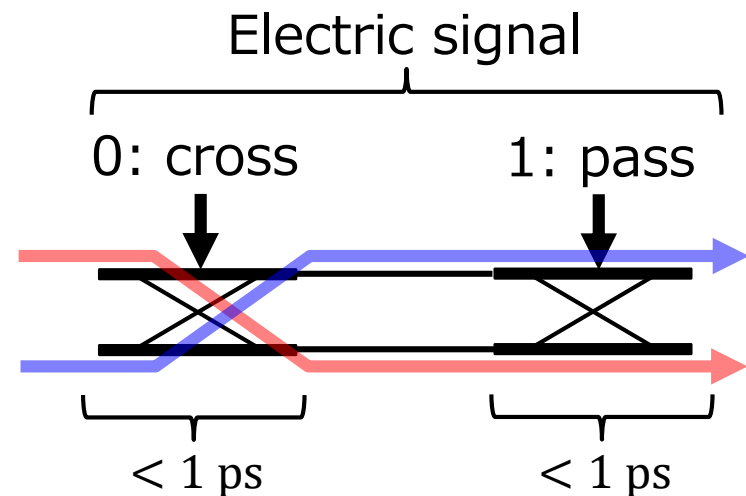
# Photonic Crystal-Based Optical Pass-Gate (OPG)

Directional coupler

Length  $> 100 \mu\text{m}$

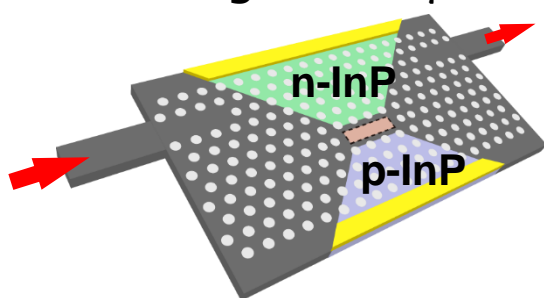


Light  
Light

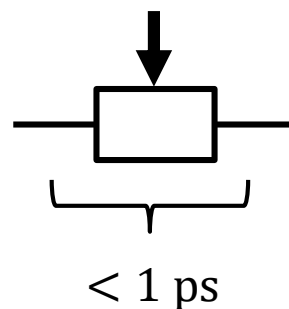


Modulator (switch)

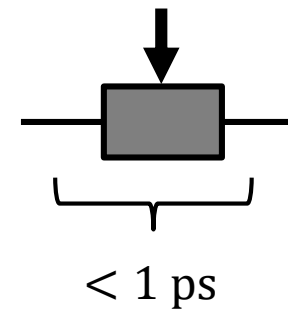
Length  $\sim 1.3 \mu\text{m}$



Electric signal  
0: OFF    1: ON



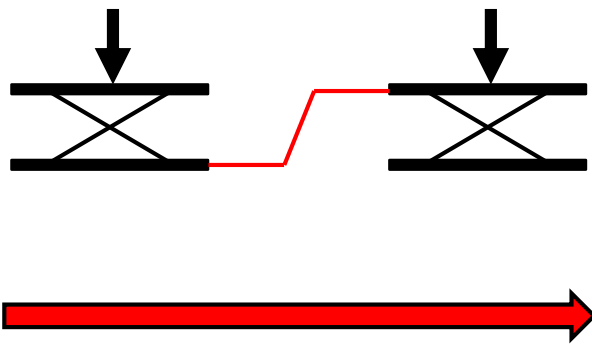
Electric signal  
0: ON    1: OFF





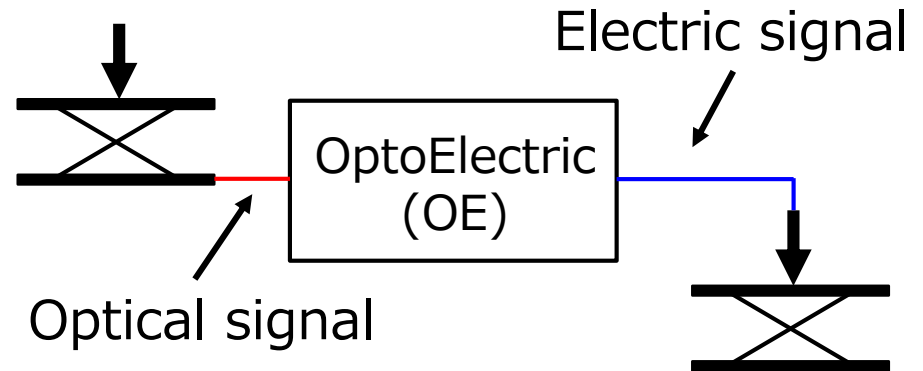
# OptoElectric (OE) Conversion Delay

Serial connection



Light speed ( $\sim 1$  ps/gate)

Cascade connection

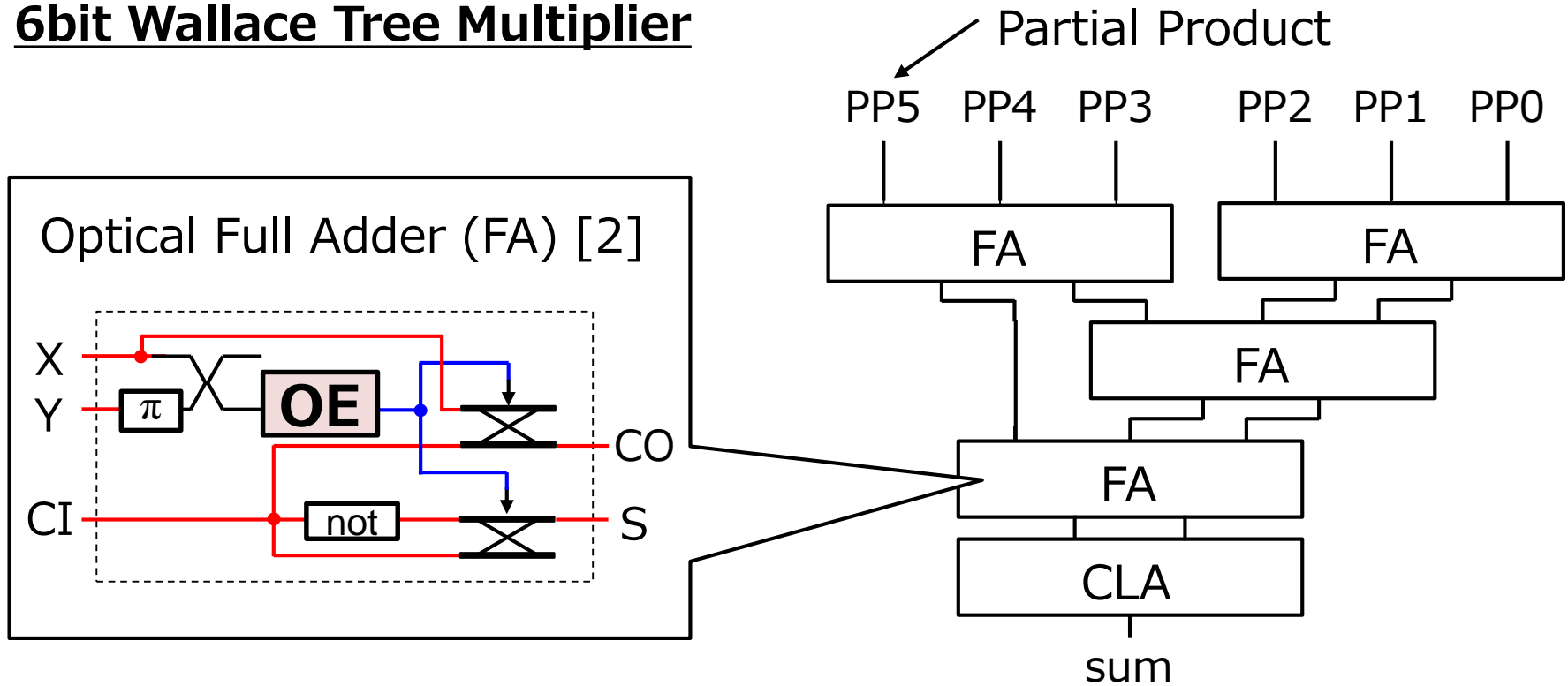


OE conversion ( $\sim 25$  ps/OE)

- ✓ **Reducing OEs on a critical path** is a key to ultra fast operation

# Issues in Conventional Optical Parallel Multiplier

## 6bit Wallace Tree Multiplier



⊗ Large # (OEs on critical path) for large bit width  $n$

➡ Unacceptable latency for large  $n$

# Outline

---

- Background
- Parallel Multiplier Using Nanophotonic Devices
- Approximate Parallel Multiplier for NN
- Conclusion

# Log-quantized Approximate Multiplier

Weight  $\rightarrow$  Data

Key idea:  $W \times U = W \times \text{sign}(U) \times 2^{\log_2|U|}$  ( $W, U \in \mathbb{R}$ )

$$\log_2|U| \simeq \underline{\text{floor}}(\log_2 U) := \tilde{u}$$

Quantization to an integer

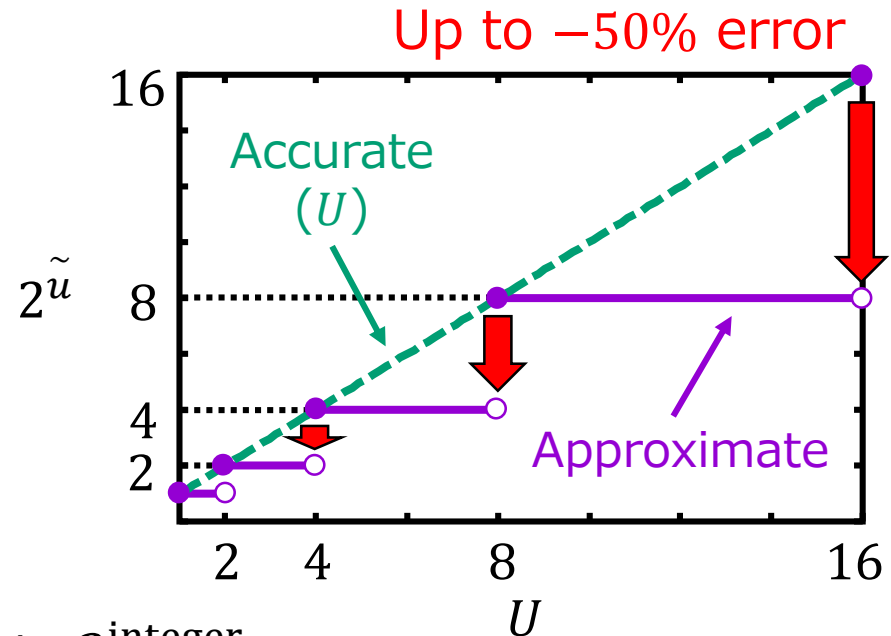
Example ( $U = 42$ )

$U =$  00101010  
MSB LSB

Log-quantization

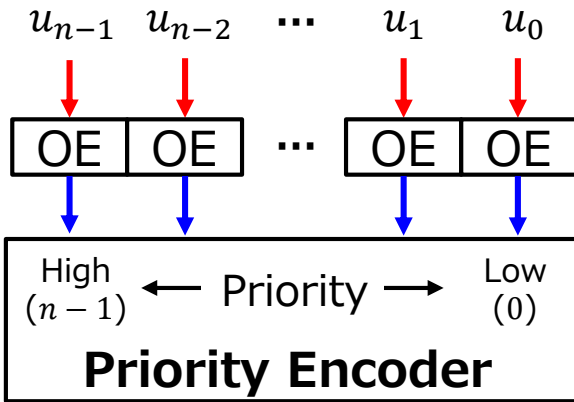
$2^{\tilde{u}} =$  00100000  
MSB LSB

Quantization to  $2^{\text{integer}}$

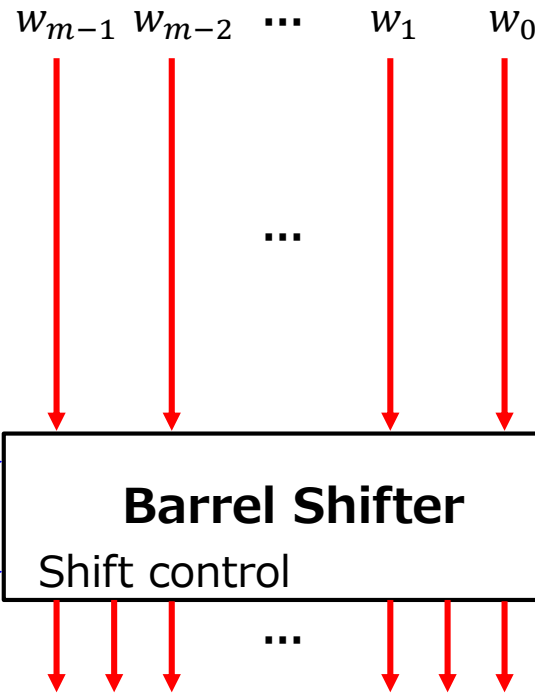


# Overview of the Approximate Multiplier

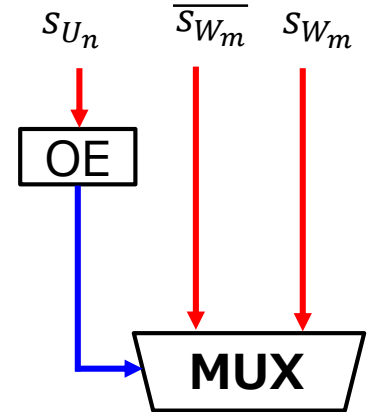
## Input value $U_n$ ( $n$ -bit)



## Weight value $W_m$ ( $m$ -bit)



## Sign bits



## Multiplication Result ( $m + n - 1$ )-bit

$$W_m \times 2^{\tilde{u}}$$

## Sign bit

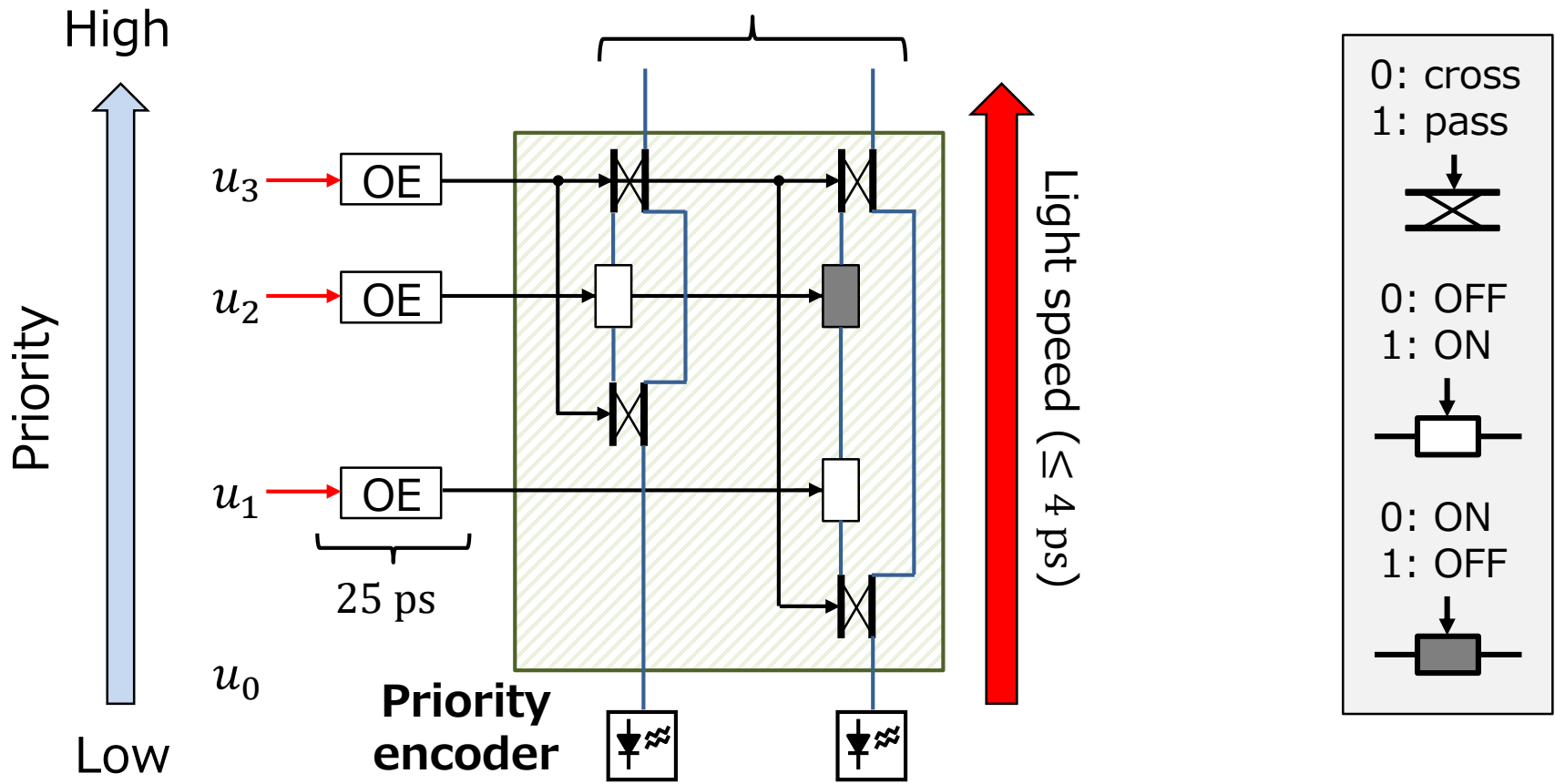
→ Optical    → Electrical

Quantized  $\tilde{u}_n$   
(= floor(log<sub>2</sub> U<sub>n</sub>))

✓ Only two OE converters on critical paths for any bit widths

# Optical Implementation of the Log-Quantizer ( $n = 4$ )

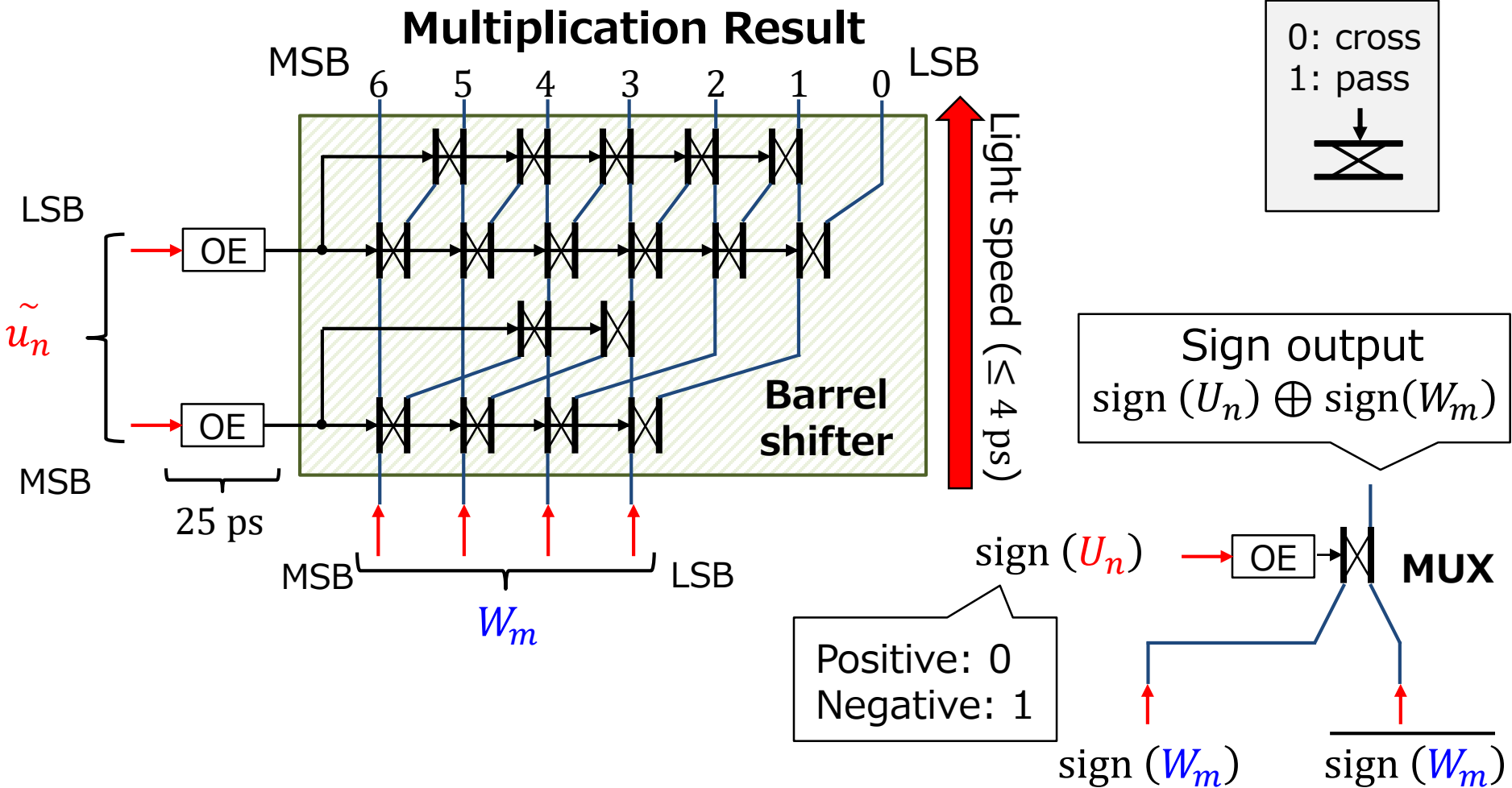
**Quantized log value  $\tilde{u}_n (= \text{floor}(\log_2 U_n))$  (2 bit)**



✓ Only one OE converter on a critical path for any  $n$  14

$$W_m \times 2^{\tilde{u}_n}$$

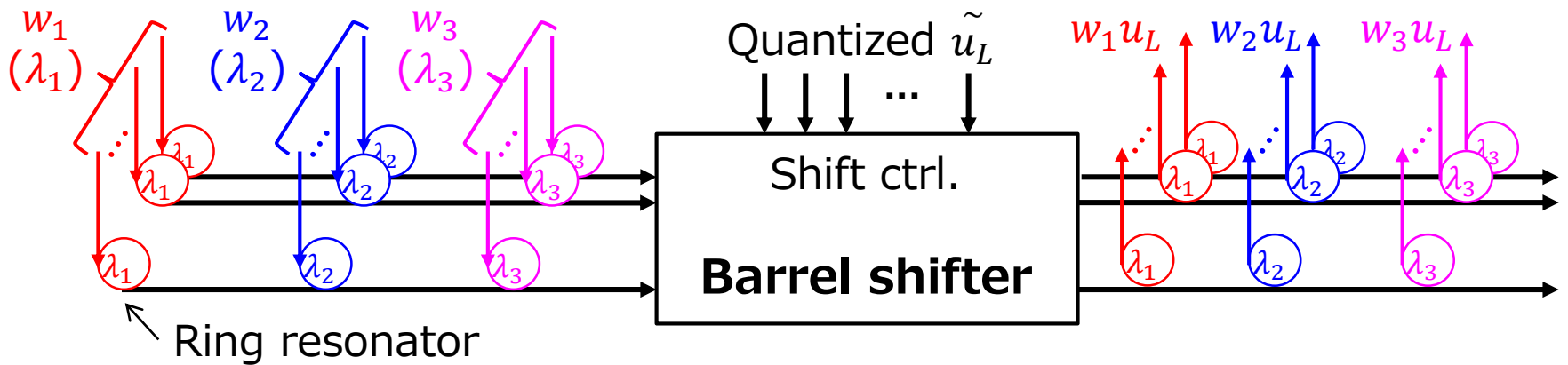
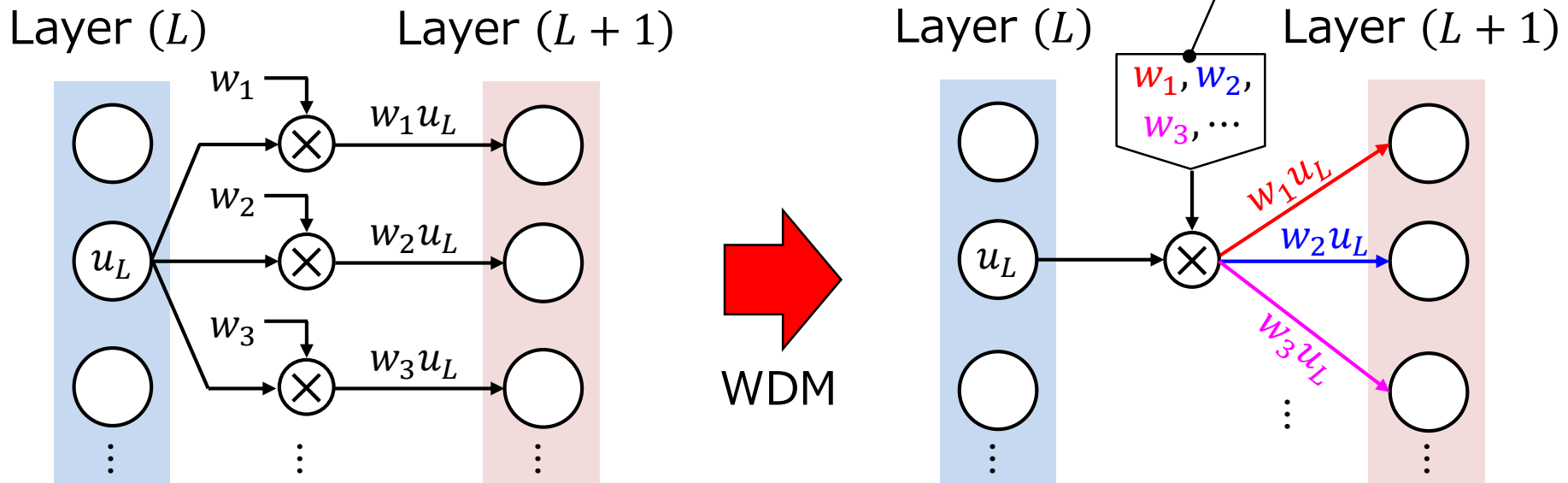
# Barrel Shifter-based Multiplier ( $n = 4$ )



✓ Only one OE converter on a critical path for any  $n$  15

# Area Reduction of Optical NN Circuits Exploiting Wavelength Division Multiplexing (WDM)

✓ Different wavelength for different weight





# Conclusion

---

- Parallel multiplier for neural networks
  - Performance boost by approximate structure
  - Priority encoder & barrel shifter
  - Area reduction exploiting WDM
- Future work
  - Detailed analysis
  - Detailed implementation

# Acknowledgement

---

This work is partly supported by CREST (Core Research for Evolutional science and Technology) of JST (Japan Science and Technology Corporation)