

# ooc\_cuDNN : A Deep Learning Library

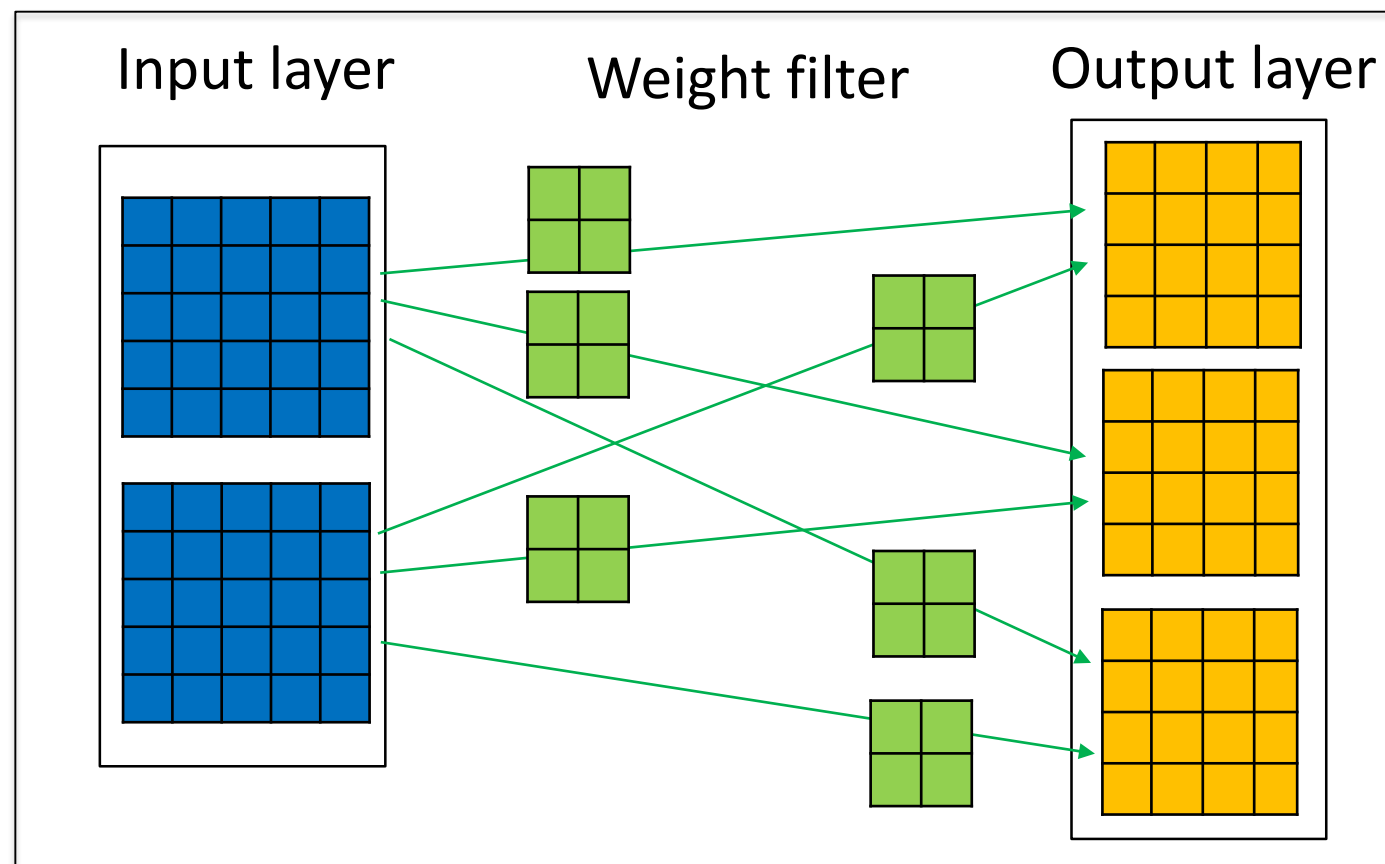
## Supporting CNNs over GPU Memory Capacity

Yuki Ito, Ryo Matsumiya, Toshio Endo (Tokyo Institute of Technology)

### Background

• Convolutional neural networks (CNNs) are used in many fields.

- Image recognition, Image processing, speech recognition, etc...



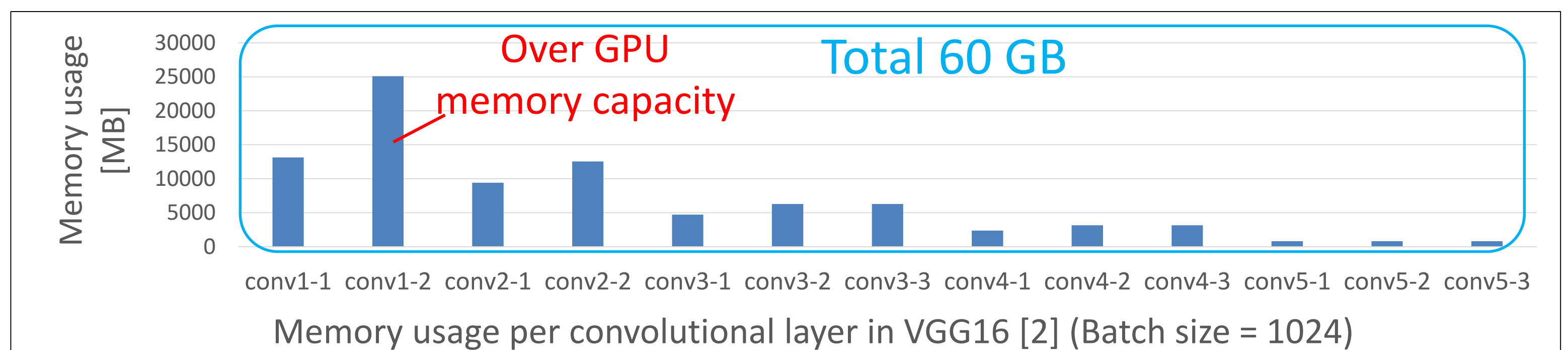
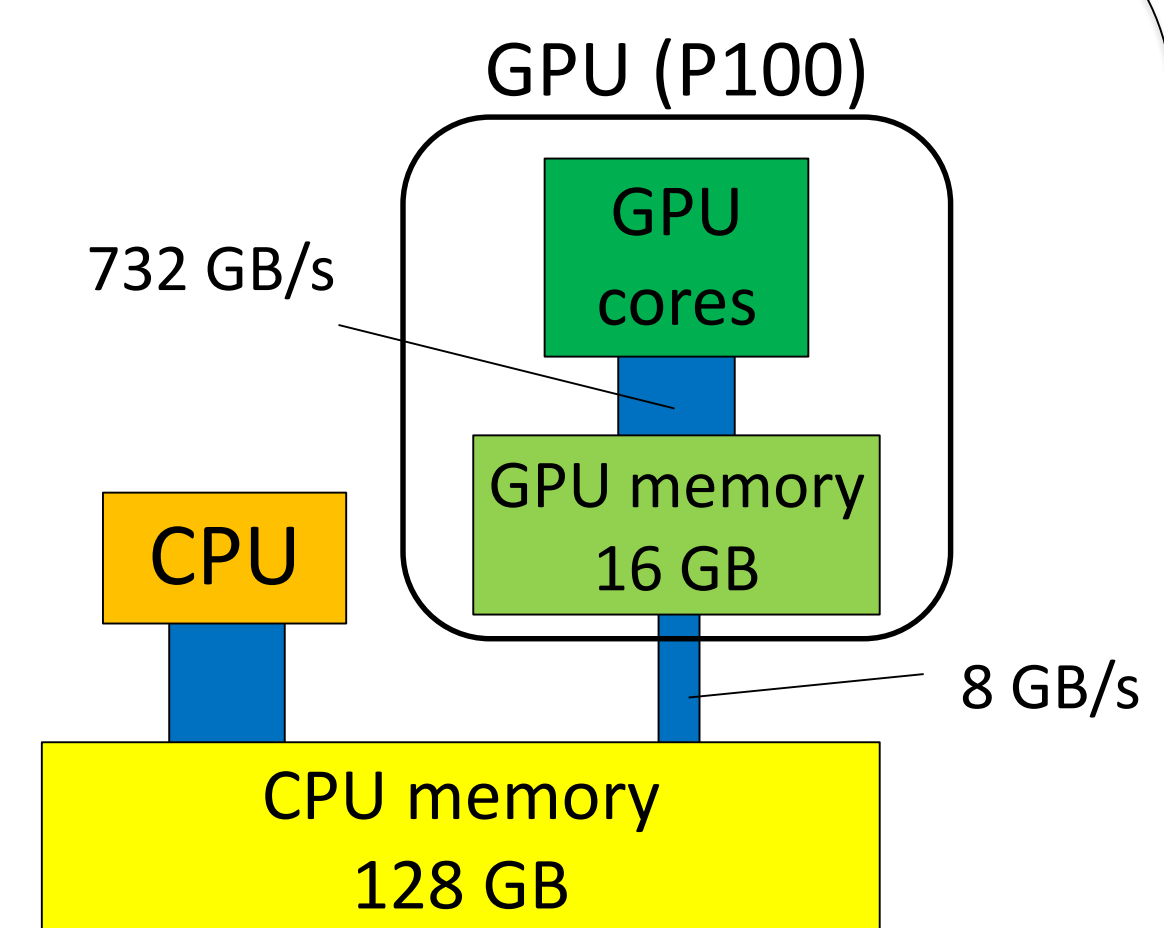
Convolution computation

• cuDNN [1] library can accelerate computation of CNNs

- Developed by NVIDIA
- Used by many deep learning frameworks
- Use **graphic processing units (GPUs)** effectively

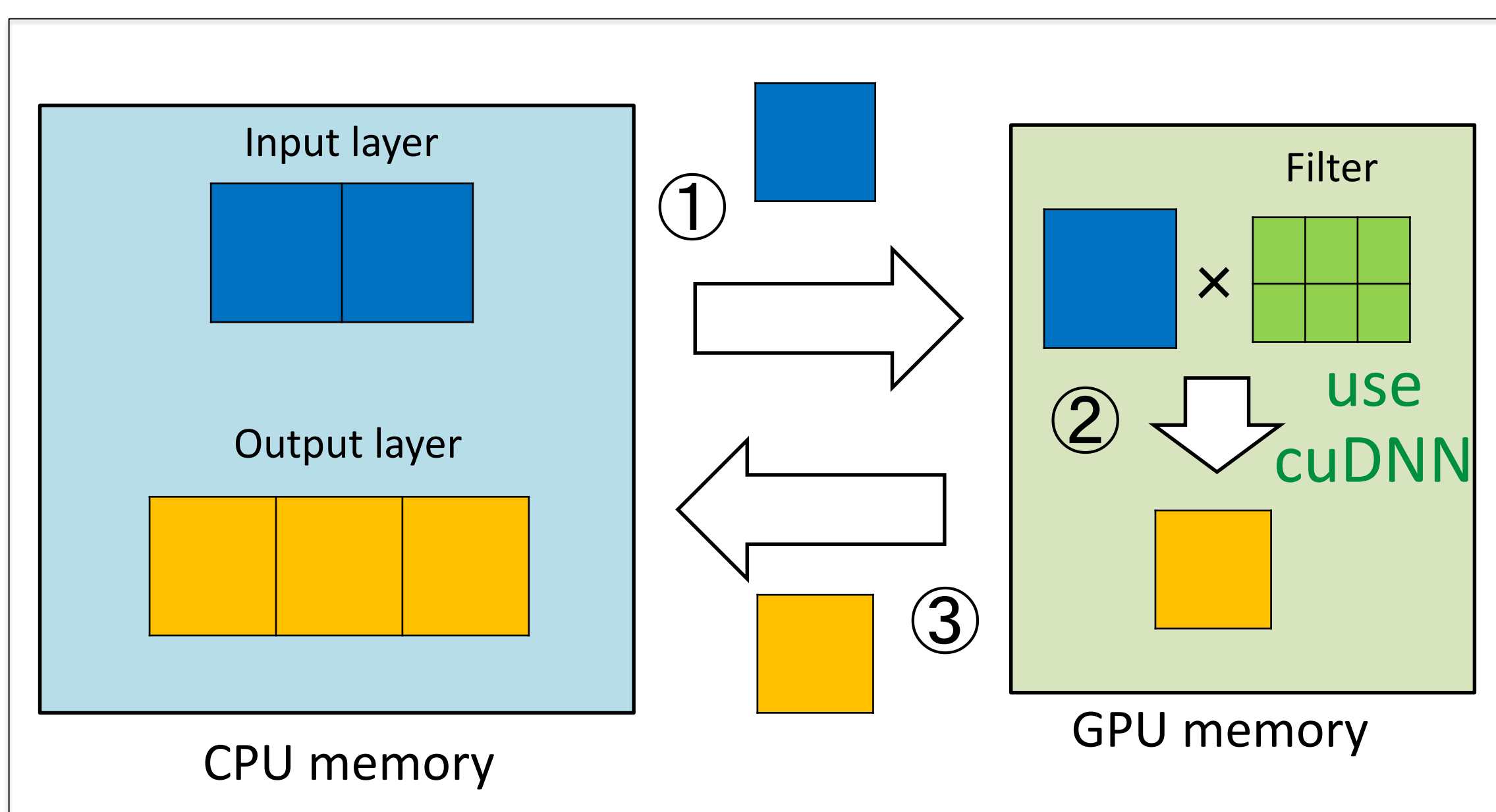
### • Motivation

- It is hard for large scale CNNs to be computed using cuDNN
- ❖ cuDNN can use GPU memory only
- ❖ GPU memory capacity is limited
- ✓ Even computation of one layer may run out of GPU memory



### Our solution

- We designed and implemented **ooc\_cuDNN** [2] library.
- **ooc\_cuDNN (out-of-core cuDNN)** supports large scale CNNs
  - Compatible with cuDNN
  - Enable to compute CNNs that exceed GPU memory capacity
    - ❖ Use both GPU and CPU memory
  - Divide layers and filters
    - ❖ Each layer (or filter) is put on GPU or CPU memory
  - Divided data are used for computation on GPU with cuDNN.
    - ❖ Swap data between CPU and GPU memory
    - ❖ Overlap CPU-GPU communication and computation



Use CPU and GPU memory

### Optimization(1) : Auto-tuning division sizes

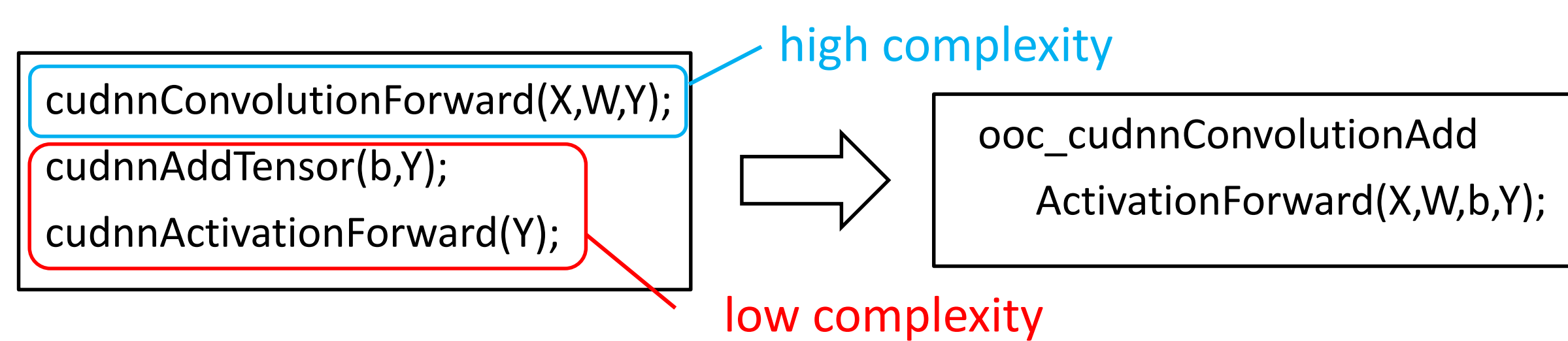
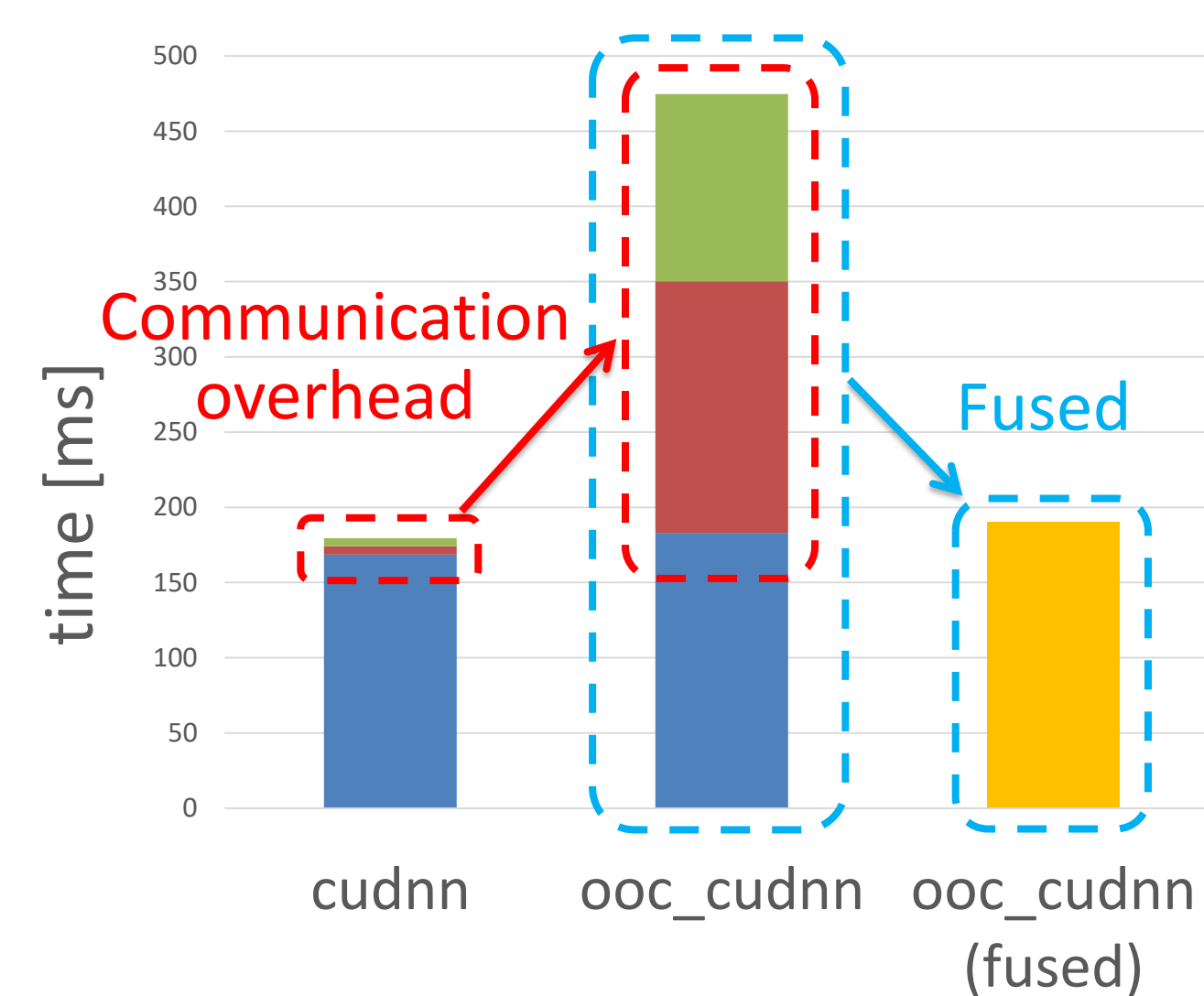
- Performance of **ooc\_cuDNN** is affected by each division size.
  - Make performance model
  - Optimize division size based on the model.

$$T_{conv} = t_{HtoD} + t_{conv} + t_{DtoH} + \left(\left\lceil \frac{c_x}{d_{c_x}} \right\rceil - 1\right) \max(t_{HtoD}, t_{conv}) + \left(\left\lceil \frac{n}{d_n} \right\rceil \left\lceil \frac{c_y}{d_{c_y}} \right\rceil \left\lceil \frac{h_y}{d_{h_y}} \right\rceil - 1\right) \max\left(\left\lceil \frac{c_x}{d_{c_x}} \right\rceil t_{HtoD}, \left\lceil \frac{c_x}{d_{c_x}} \right\rceil t_{conv}, t_{DtoH}\right)$$

Performance model of convolution

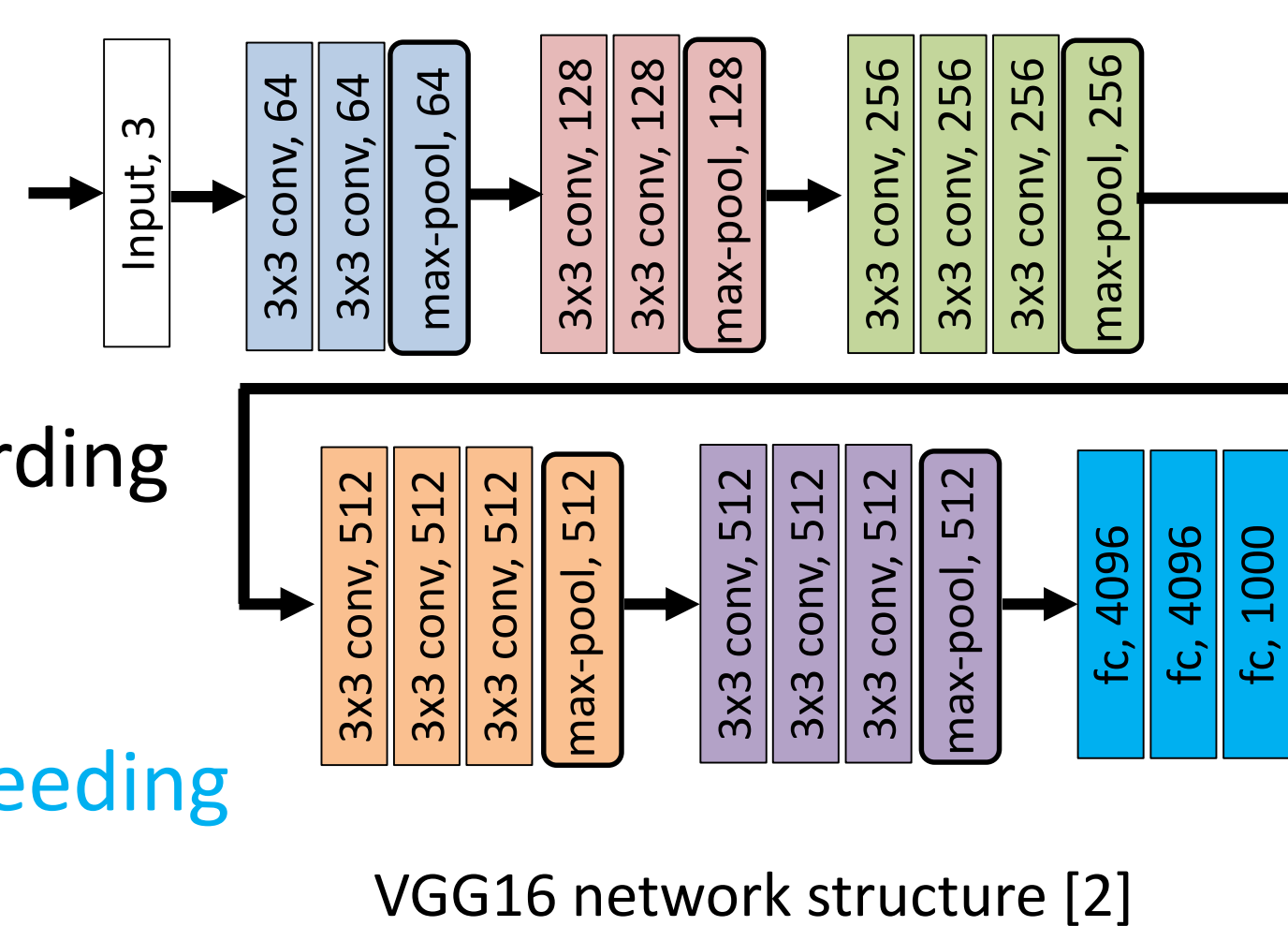
### Optimization(2) : Fusion of computations

- Performance of low complexity computations is too low in **ooc\_cuDNN**.
  - In those computations, communication can not be hidden completely.
  - Provide **fused functions** that perform high complexity computations and low complexity computations at once.

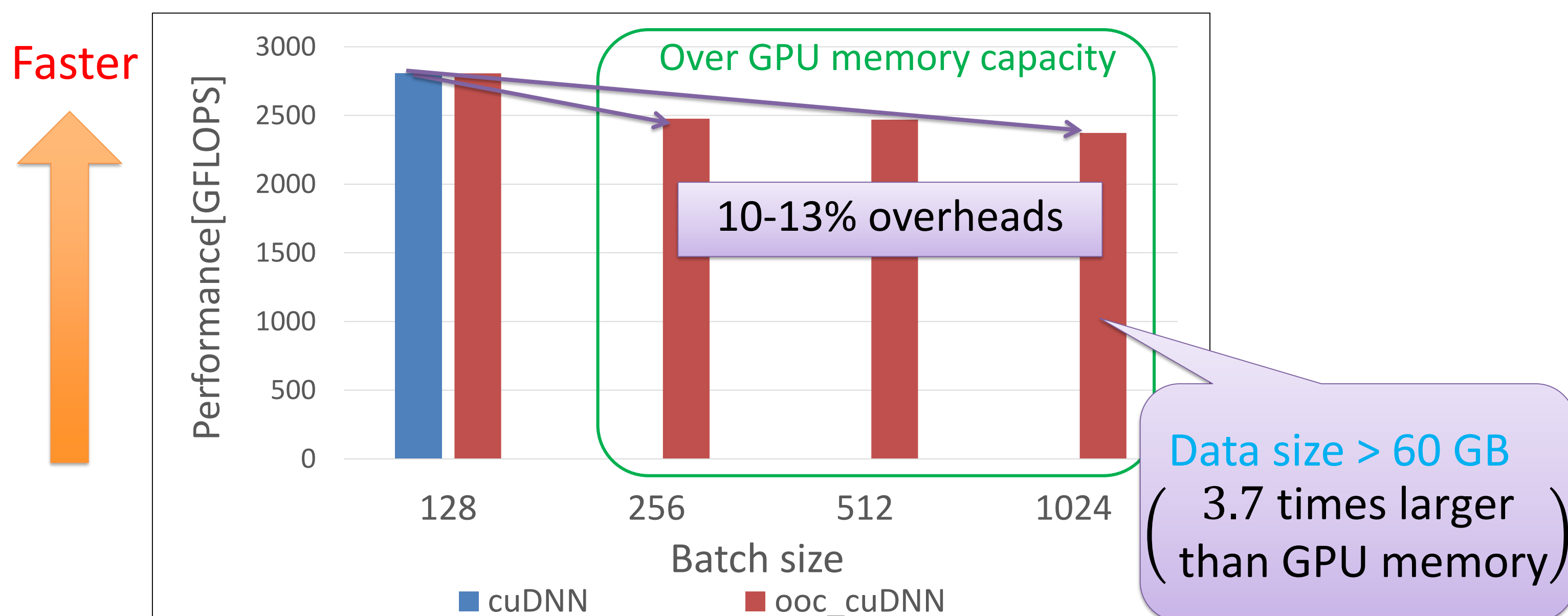


### Evaluation

- Apply **ooc\_cuDNN** to CNN application
  - Forward and Backward of VGG16[3]
  - The required memory size increases according to batch size.
- Experiment with Tesla P100
  - **ooc\_cuDNN** enables to compute CNN exceeding GPU memory capacity.

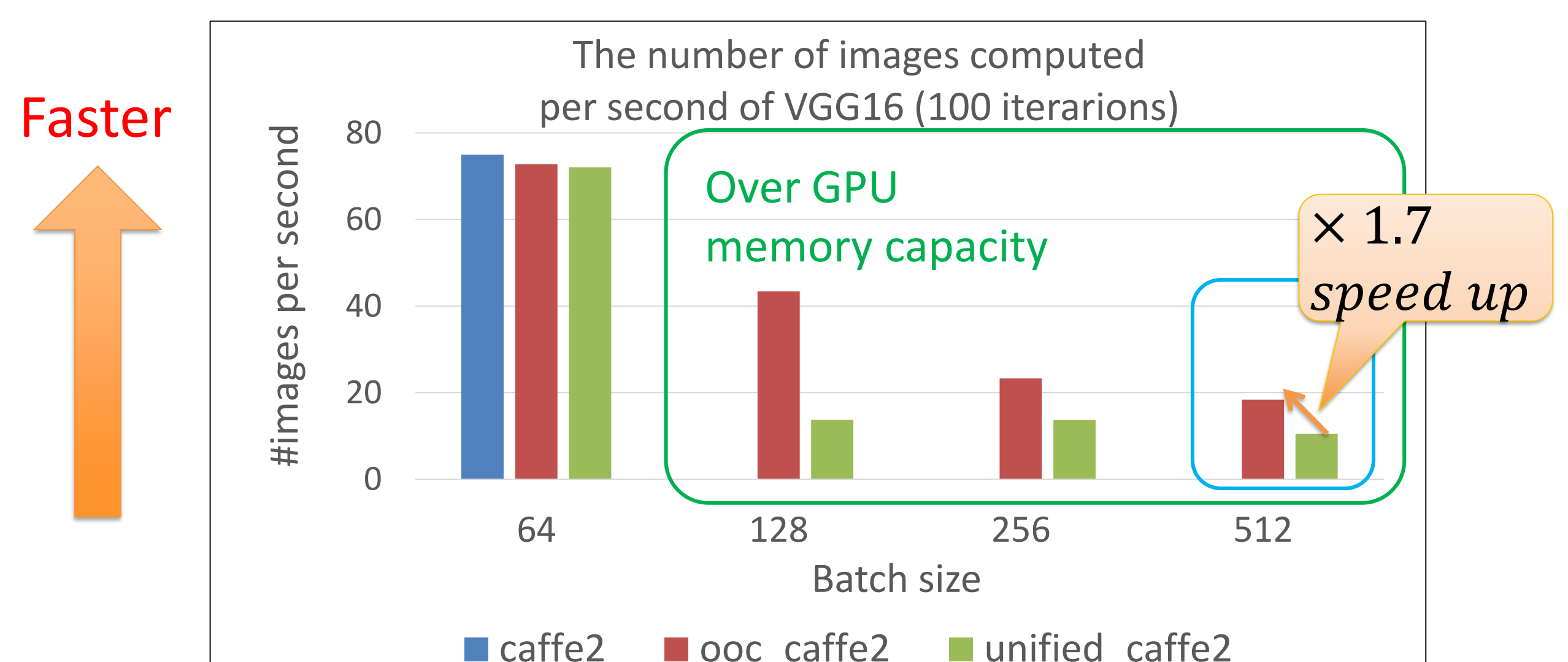


VGG16 network structure [2]



### Integrating with deep learning framework

- We implemented **ooc\_Caffe2** (Caffe2 with ooc\_cuDNN).
  - Caffe2[4] is a deep learning framework developed by Facebook.
  - Not support ooc\_cuDNN's fused functions in current design.
- For comparison, we implemented **unified\_Caffe2**.
  - Use original cuDNN, and allocate data as **Unified Memory**.
    - ❖ Unified memory supports data exceeding GPU memory capacity by swapping mechanism between CPU and GPU.
  - ooc\_Caffe2 is **> x1.7 faster** in out-of-core cases.



### Future work

- Optimization considering the entire CNN
  - Which data should be put on CPU memory?
  - Which computation should be fused?
- Improve ooc\_Caffe2
  - Use fused functions
  - Support distributed computation

[1] NVIDIA Corporation, NVIDIA cuDNN, <https://developer.nvidia.com/cudnn>  
 [2] Yuki Ito et al. ooc\_cuDNN: Accommodating Convolutional Neural Networks over GPU Memory Capacity, IEEE BigData, 2017  
 [3] Karen Simonyan et al. Very Deep Convolutional Networks for Large-Scale Image Recognition, ICLR, 2015  
 [4] Facebook, Caffe2 | A New Lightweight, Modular, and Scalable Deep Learning Framework, <https://caffe2.ai/>