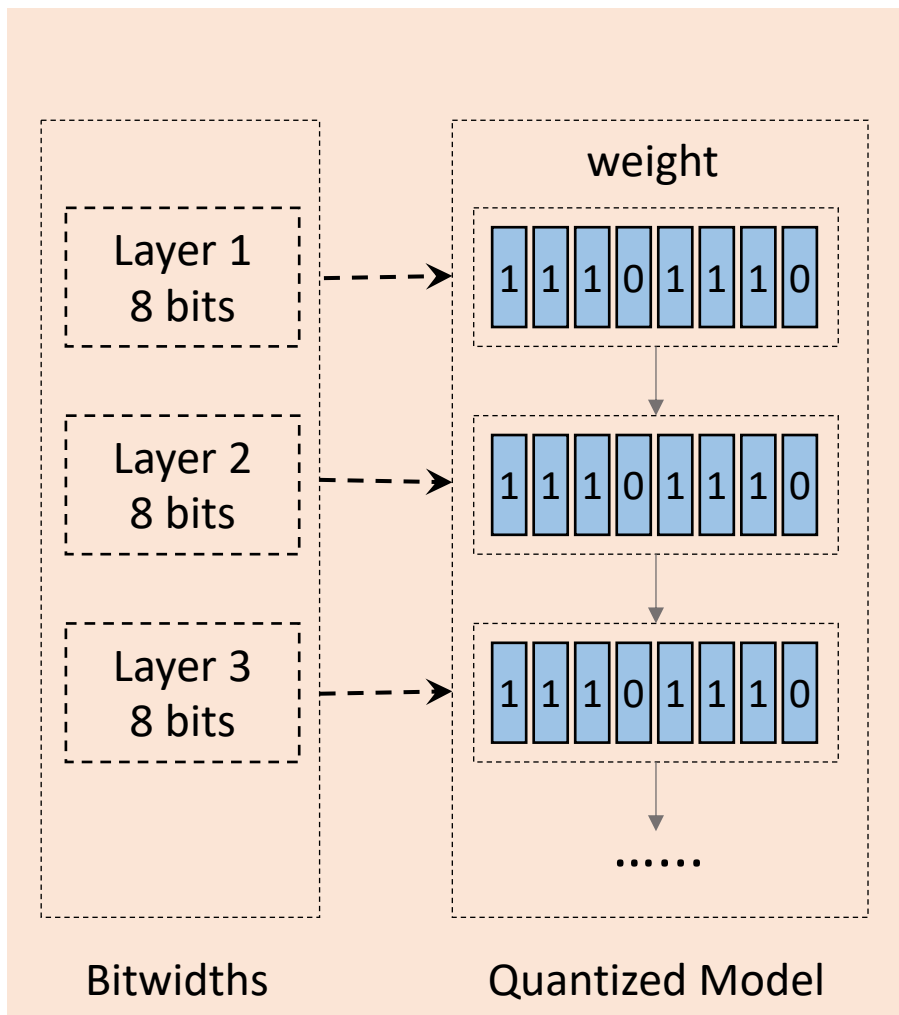# Simple Augmentation Goes A Long Way: ADRL for DNN Quantization

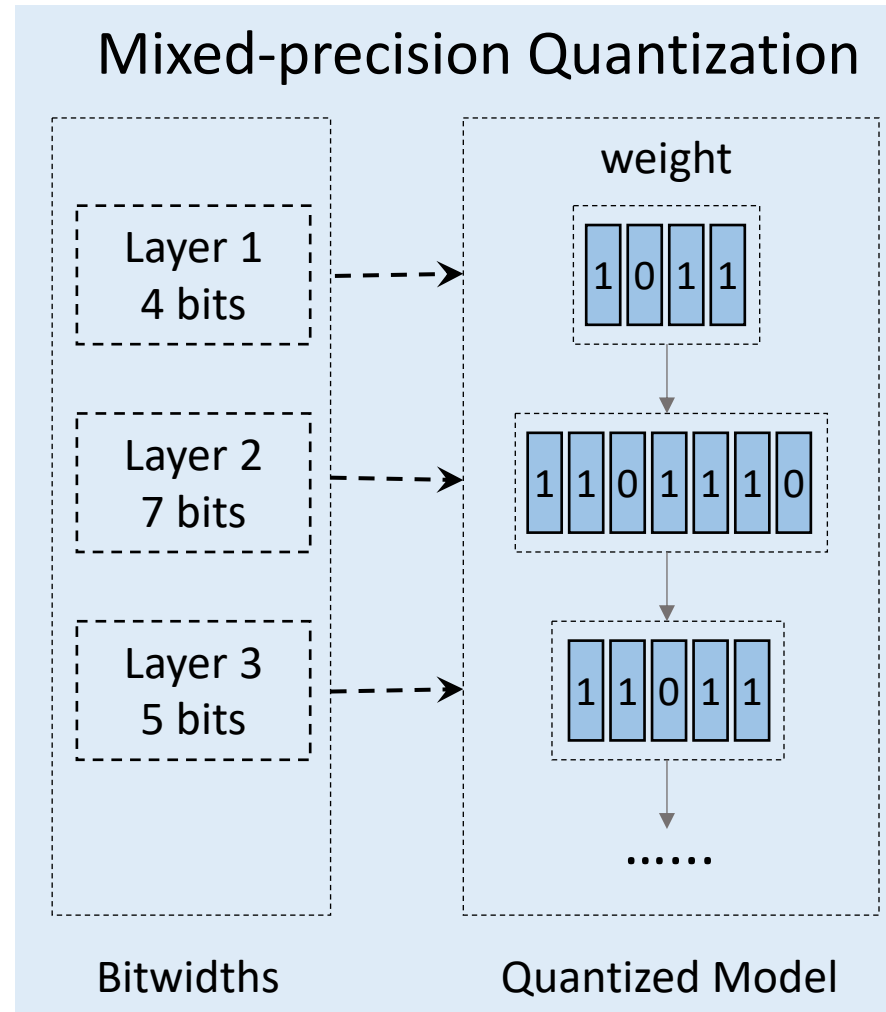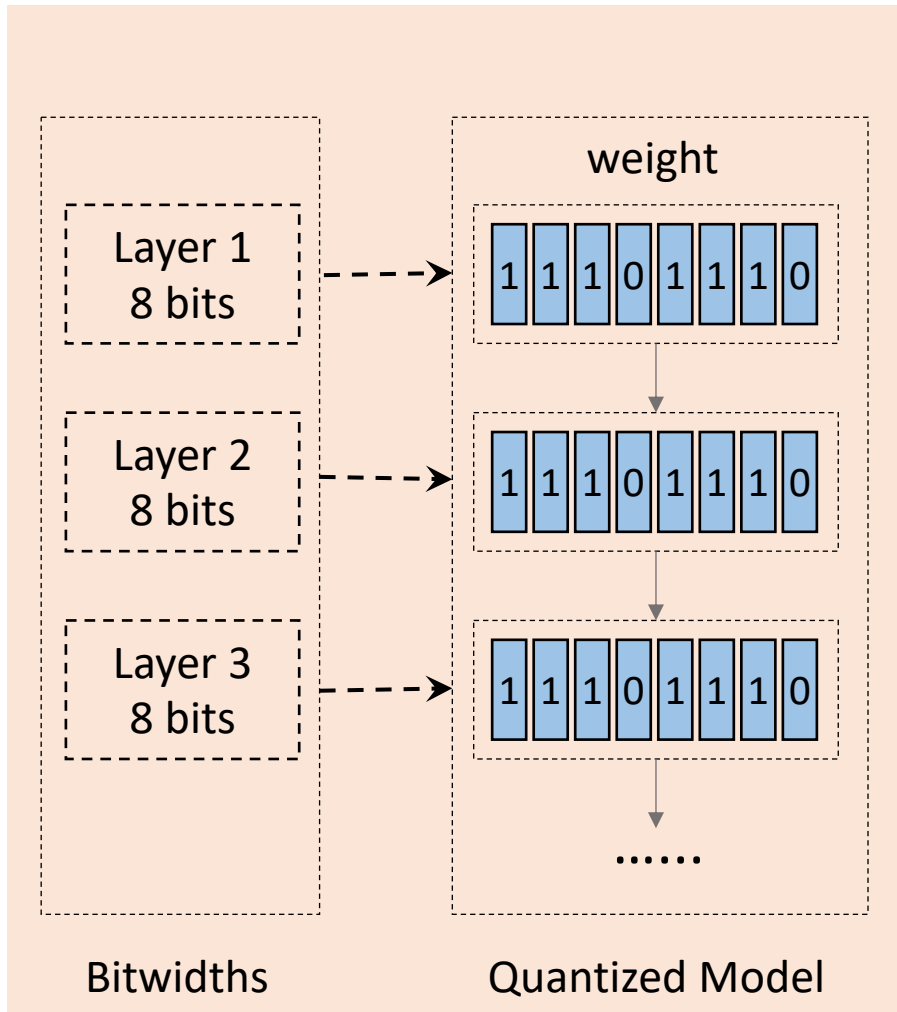Lin Ning[1], Guoyang Chen[2], Weifeng Zhang[2], Xipeng Shen[1]

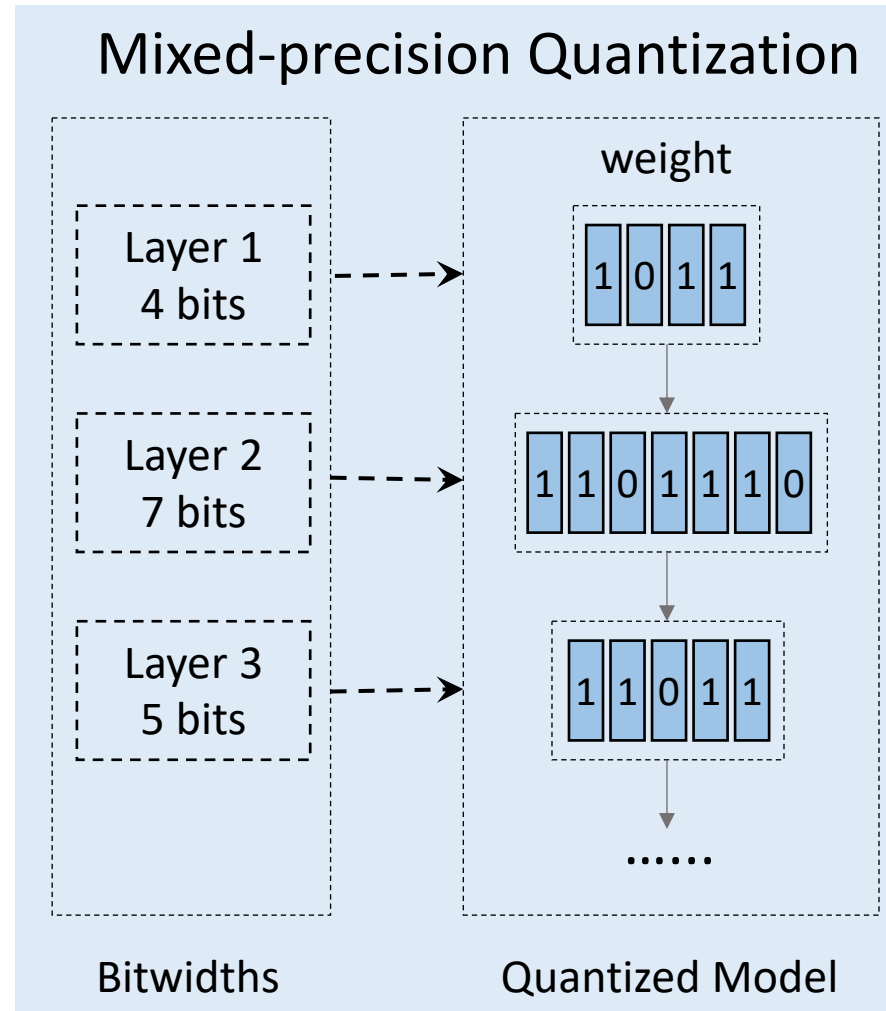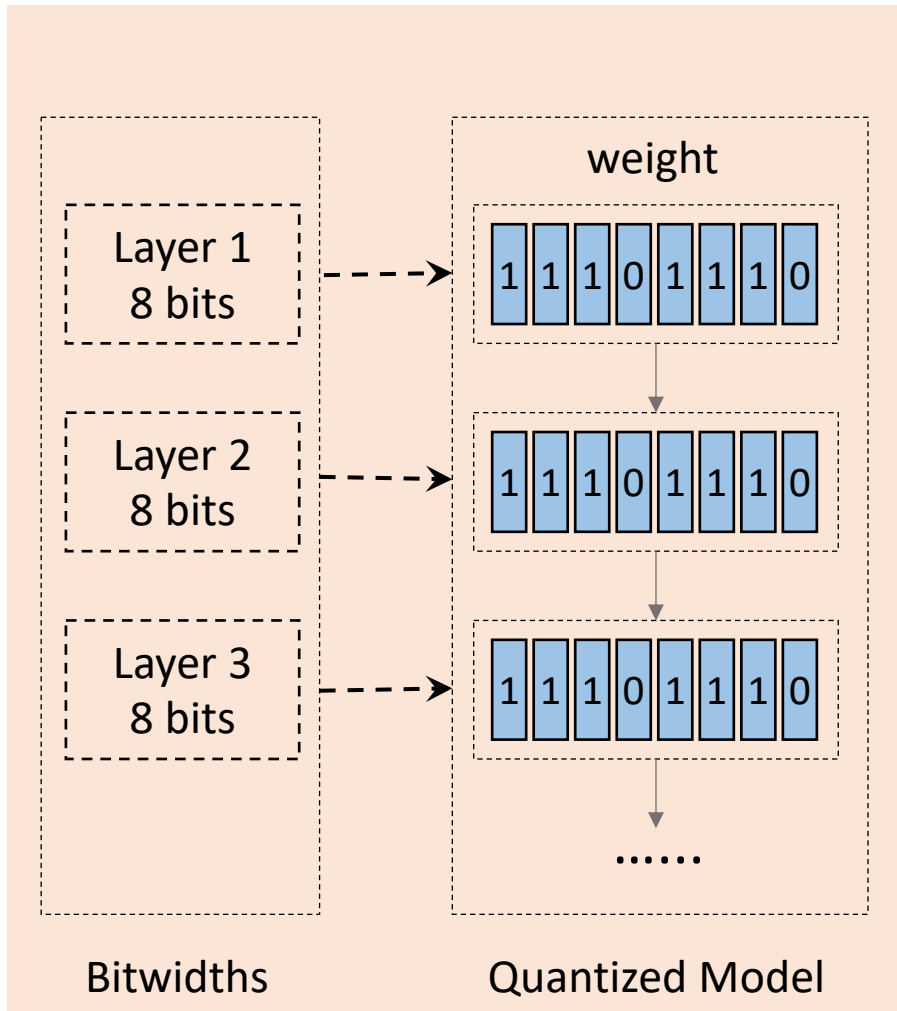[1] North Carolina State University, Raleigh, NC, USA

[2] Alibaba Group, Sunnyvale, CA, USA

# Mixed-Precision Quantization Problem

# Mixed-Precision Quantization Problem
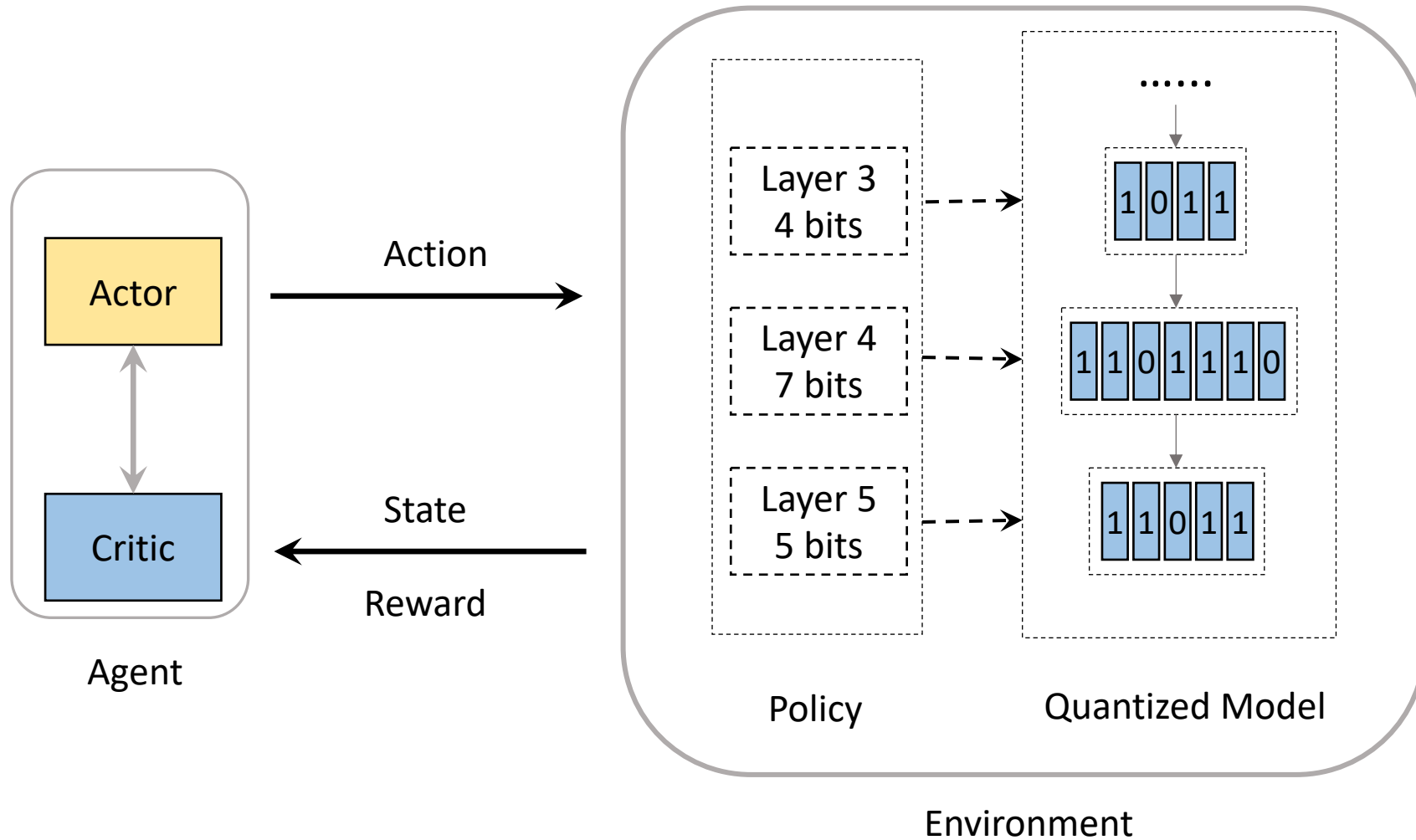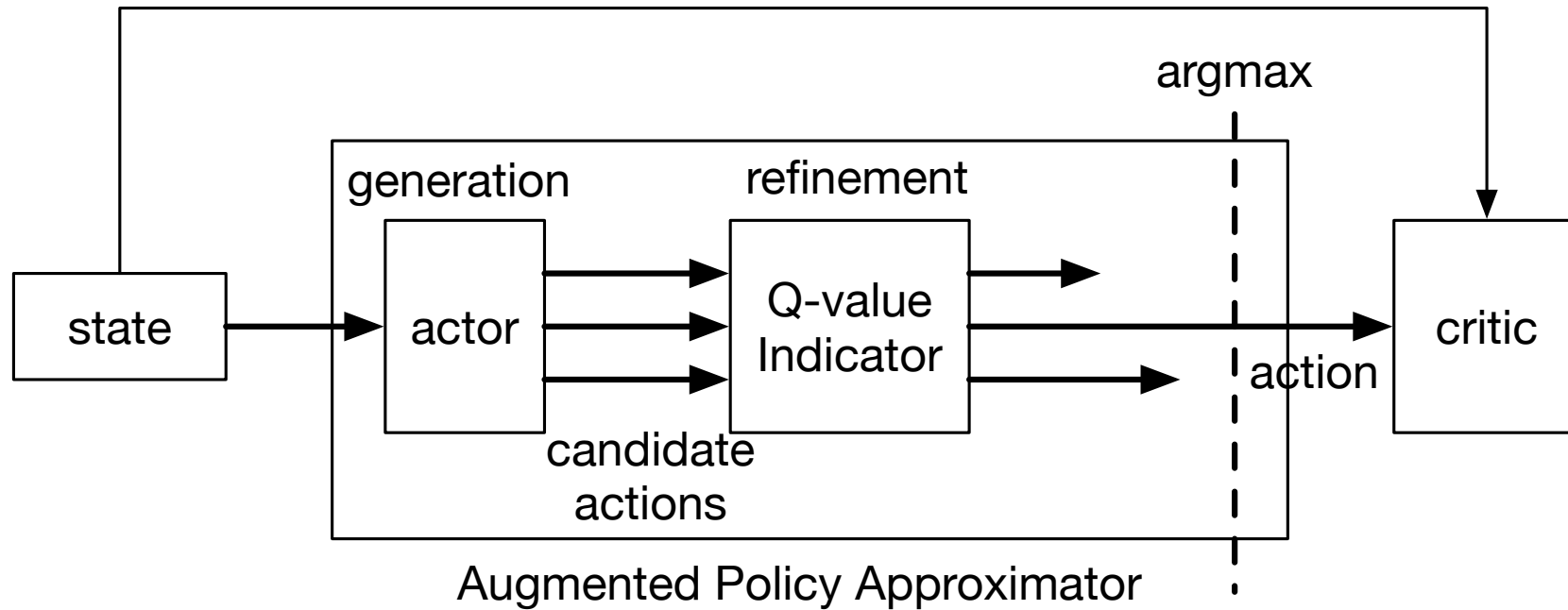
# Mixed-Precision Quantization Problem
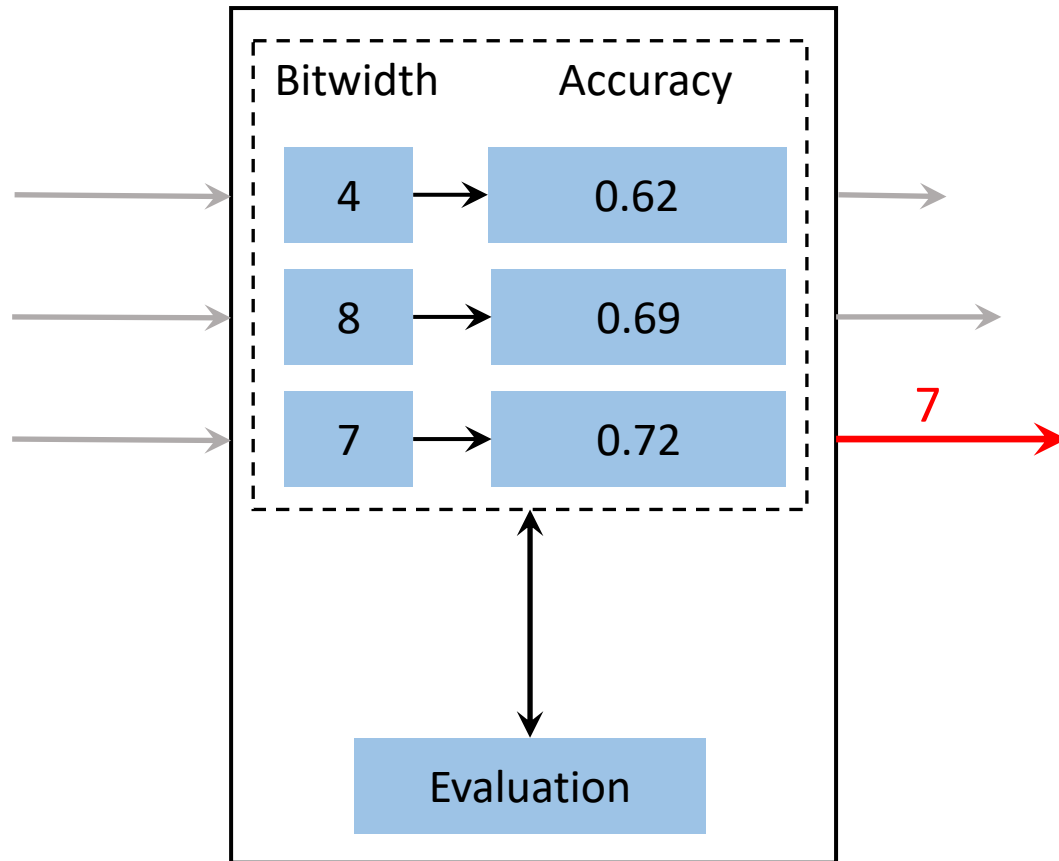


**Design Space: $8^n$**

# DRL for Mixed-Precision Quantization

# ADRL: Augmented Policy Approximator

# ADRL: Q-value Indicator

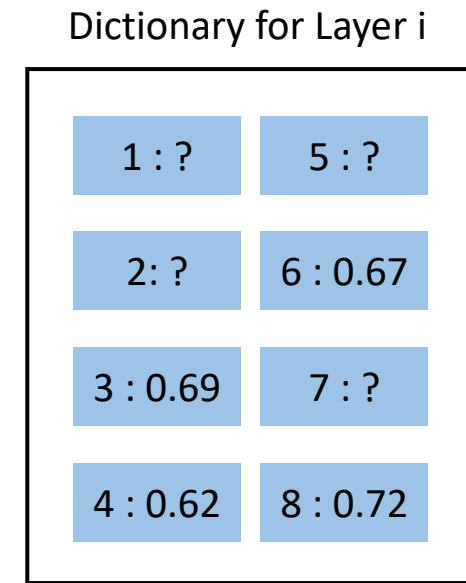Profiling-Based Indicator

# ADRL: Q-value Indicator

# ADRL: Evaluation

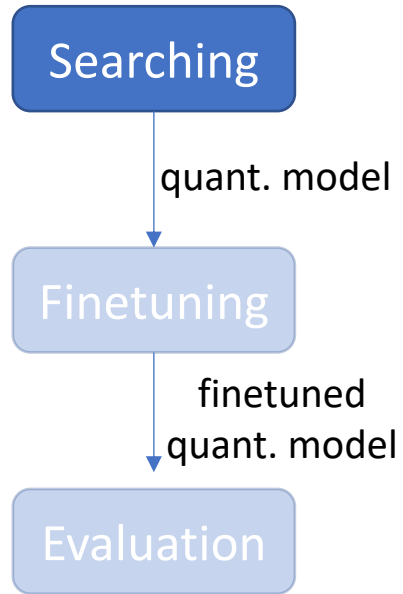| Network | Dataset |
|---------|---------|
| CifarNet | Cifar10 |
| ResNet20 | Cifar10 |
| AlexNet | ImageNet |
| ResNet50 | ImageNet |

## Server

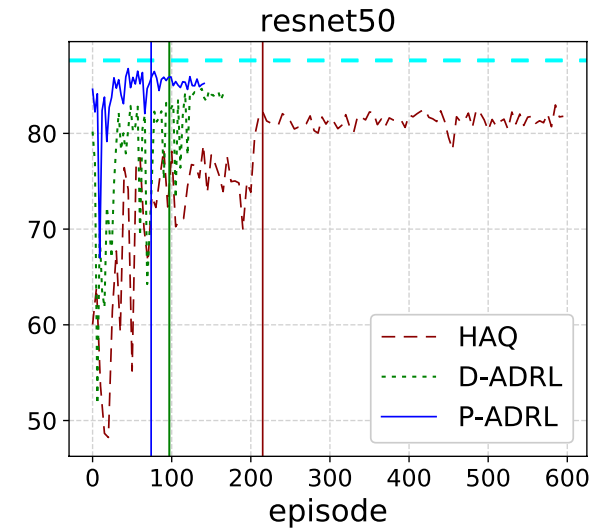- Intel(R) Xeon(R) Platinum 8168 Processor
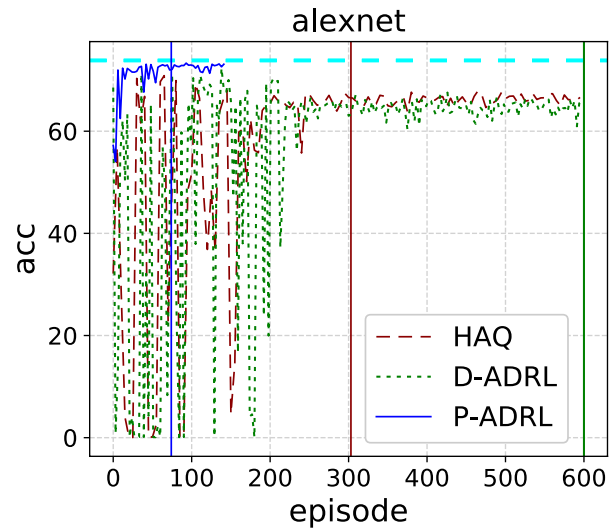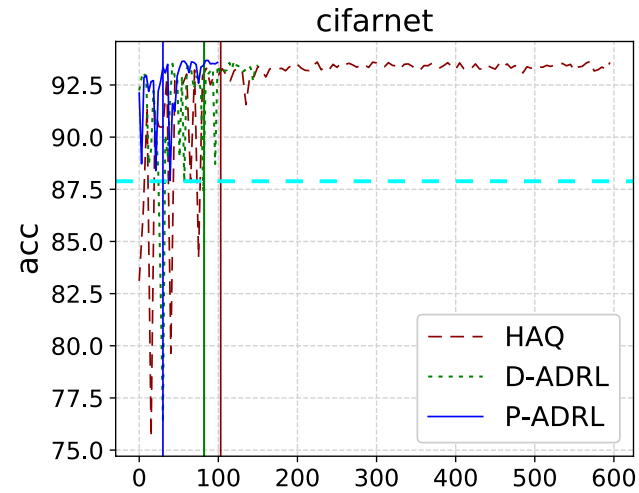- 32GB memory
- 4 NVIDIA Tesla V100 32GB GPUs.

## DRL based mix-precision quantization:

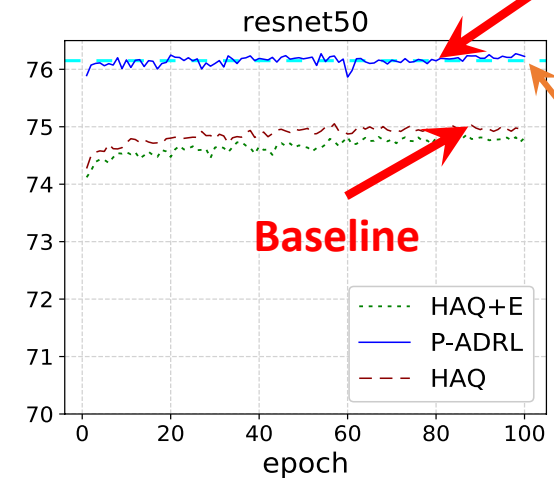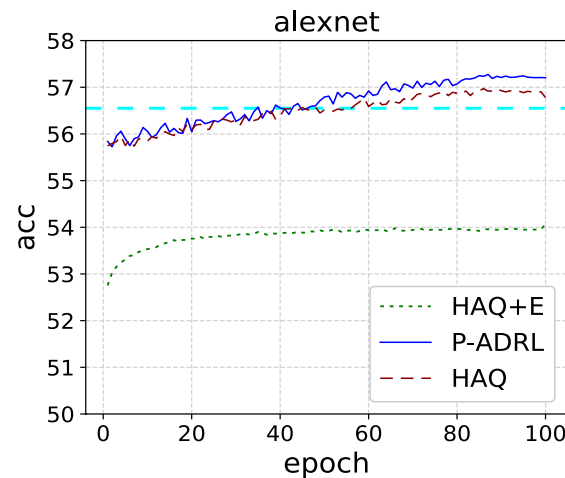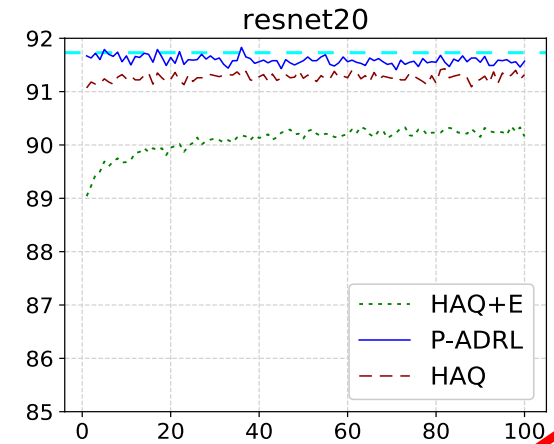Searching → quant. model → Finetuning → finetuned quant. model → Evaluation

# ADRL: Evaluation
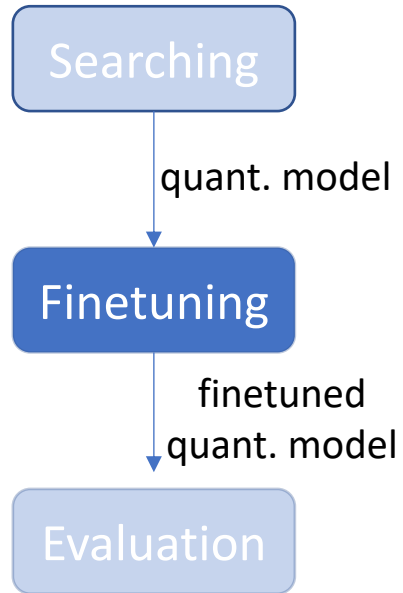
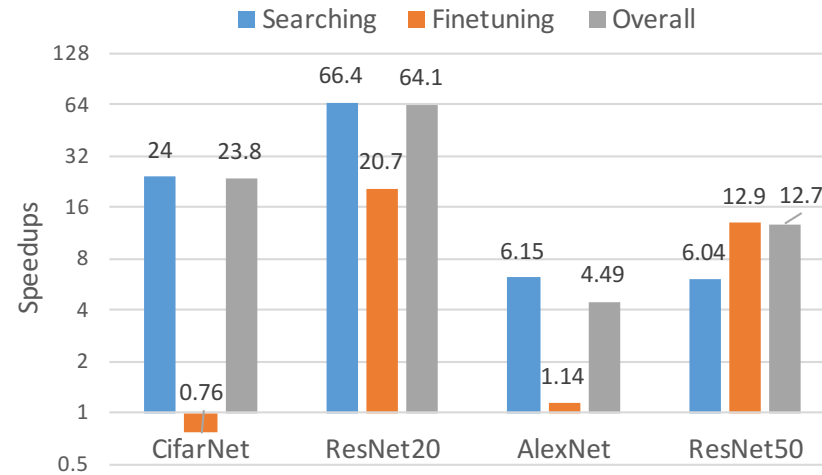Searching: Accuracy

# ADRL: Evaluation

Finetuning: Accuracy

# ADRL: Evaluation

Searching

↓ quant. model

Finetuning

↓ finetuned quant. model

Evaluation

## Speedup



## Summarize

Augmented DRL

➤ Produces more accurate quantized models than the state of the art DRL-based quantization

➤ Improves the learning speed

➤ Significantly magnifies the potential of DRL for DNN quantization.