

## MIT Open Access Articles

*Hardware-Centric AutoML for Mixed-Precision Quantization*

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

**As Published:** <https://doi.org/10.1007/s11263-020-01339-6>

**Publisher:** Springer US

**Persistent URL:** <https://hdl.handle.net/1721.1/131513>

**Version:** Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

**Terms of use:** Creative Commons Attribution-Noncommercial-Share Alike



## Hardware-Centric AutoML for Mixed-Precision Quantization

**Cite this article as:** Kuan Wang, Zhijian Liu, Yujun Lin, Ji Lin and Song Han, Hardware-Centric AutoML for Mixed-Precision Quantization, International Journal of Computer Vision <https://doi.org/10.1007/s11263-020-01339-6>

This Author Accepted Manuscript is a PDF file of an unedited peer-reviewed manuscript that has been accepted for publication but has not been copyedited or corrected. The official version of record that is published in the journal is kept up to date and so may therefore differ from this version.

Terms of use and reuse: academic research for non-commercial purposes, see here for full terms. <https://www.springer.com/aam-terms-v1>

Author accepted manuscript

# Hardware-Centric AutoML for Mixed-Precision Quantization

Kuan Wang\* · Zhijian Liu\* · Yujun Lin\* · Ji Lin · Song Han

Received: date / Accepted: date

**Abstract** Model quantization is a widely used technique to compress and accelerate deep neural network (DNN) inference. Emergent DNN hardware accelerators begin to support *flexible bitwidth* (1-8 bits) to further improve the computation efficiency, which raises a great challenge to find the optimal bitwidth for each layer: it requires domain experts to explore the vast design space trading off accuracy, latency, energy, and model size, which is both time-consuming and usually sub-optimal. There are plenty of specialized hardware accelerators for neural networks, but little research has been done to design specialized neural networks optimized for a particular hardware accelerator. The latter is demanding given the much longer design cycle of silicon than neural nets. Conventional quantization algorithm ignores the different hardware architectures and quantizes all the layers in a uniform way. In this paper, we introduce the **Hardware-Aware Automated Quantization (HAQ)** framework which automatically determine the quantization policy, and we take the hardware accelerator's feedback in the design loop. Rather than relying on proxy signals such as FLOPs and model size,

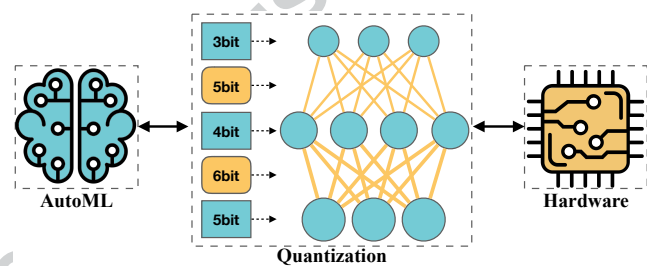


Fig. 1: **Hardware-centric** (right) **automated ML** (left) for **flexible-bitwidth quantization** (middle).

we employ a hardware simulator to generate the direct feedback signals to the RL agent. Compared with conventional methods, our framework is fully automated and can specialize the quantization policy for different neural network and hardware architectures. The learned policy can transfer very well between different neural net architectures. Our framework effectively reduced the latency by **1.4-1.95 $\times$**  and the energy consumption by **1.9 $\times$**  with negligible loss of accuracy compared with the fixed bitwidth (8 bits) quantization. Our framework reveals that the optimal policies on different hardware architectures (*i.e.*, edge and cloud architectures) under different resource constraints (*i.e.*, latency, energy and model size) are drastically different. We interpreted the implication of different quantization policies, which offer insights for both neural network architecture design and hardware architecture design.

Kuan Wang\*  
 Massachusetts Institute of Technology  
 E-mail: kuanwang@mit.edu

Zhijian Liu\*  
 Massachusetts Institute of Technology  
 E-mail: zhijian@mit.edu

Yujun Lin\*  
 Massachusetts Institute of Technology  
 E-mail: yujunlin@mit.edu

Ji Lin  
 Massachusetts Institute of Technology  
 E-mail: jilin@mit.edu

Song Han  
 Massachusetts Institute of Technology  
 E-mail: songhan@mit.edu

\* indicates equal contributions.

**Keywords** Model Quantization · Mixed-Precision · Automated ML · Hardware

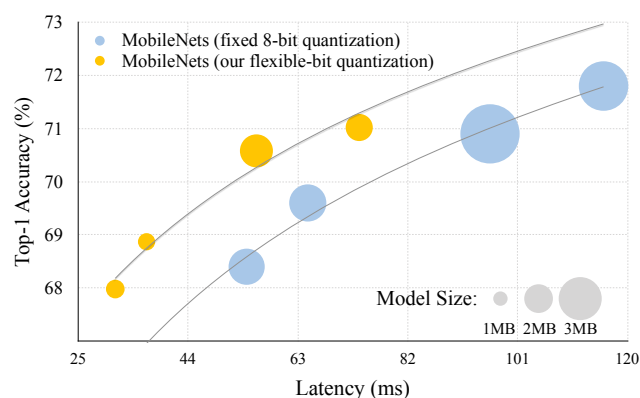


Fig. 2: We need **flexible** number of bits for different layers. We quantize MobileNets (Howard et al 2017) to different number of bits (both weights and activations), and it lies on a better pareto curve (yellow) than fixed-bit quantization (blue). This is because different layers have different redundancy and have different operation intensities (operations per byte) on the hardware, which advocates for using flexible bitwidths for different layers.

## 1 Introduction

In many real-time machine learning applications (such as robotics, autonomous driving, and mobile VR/AR), deep neural networks is strictly constrained by the latency, energy, and model size. In order to improve the hardware efficiency, many researchers have proposed to directly design efficient models (Sandler et al 2018; Howard et al 2017; Cai et al 2019) or to quantize the weights and activations to low precision (Han et al 2016; Zhu et al 2017).

Conventional quantization methods use the same number of bits for all layers (Choi et al 2018; Jacob et al 2018), but as different layers have different redundancy and behave differently on the hardware (computation bounded or memory bounded), it is necessary to use *flexible* bitwidths for different layers (as shown in Figure 2). This flexibility was originally not supported by chip vendors until recently the hardware manufacturers started to implement this feature: Apple released the A12 Bionic chip that supports flexible bits for the neural network inference (Apple 2018); NVIDIA recently introduced the Turing GPU architecture that supports 1-bit, 4-bit, 8-bit and 16-bit arithmetic operations (Nvidia 2018); Imagination launched a flexible neural network IP that supports per-layer bitwidth adjustment for both weights and activations (Imagination 2018). Besides industry, recently academia also works on the bit-level flexible hardware design: BISMO (Umuroglu et al 2018) proposed the bit-serial multiplier to support multiplications of 1 to 8 bits; BitFusion (Sharma et al 2018) supports multiplications of 2, 4, 8 and 16 bits in a spatial manner.

	Inference latency on		
	HW1	HW2	HW3
Best Q. policy for HW1	16.29 ms	85.24 ms	117.44 ms
Best Q. policy for HW2	19.95 ms	64.29 ms	108.64 ms
Best Q. policy for HW3	19.94 ms	66.15 ms	99.68 ms

Table 1: Inference latency of MobileNet-V1 (Howard et al 2017) on three hardware architectures under different quantization policies. The quantization policy that is optimized for one hardware is not optimal for the other. This suggests we need a **specialized** quantization solution for different hardware architectures. (HW1: BitFusion (Sharma et al 2018), HW2: BISMO (Umuroglu et al 2018) edge accelerator, HW3: BISMO cloud accelerator, batch = 16).

A very important missing part is, however, how to **determine the bitwidth of both weights and activations for each layer on different hardware accelerators**. This is a vast design space: with  $M$  different neural network models, each with  $N$  layers, on  $H$  different hardware platforms, there are in total  $O(H \times M \times 8^{2N})$  possible solutions (Here, we assume that the bitwidth is between 1 to 8 for both weights and activations). For a widely used ResNet-50 (He et al 2016) model, the size of the search space is about  $8^{100}$ , which is even larger than the number of particles in the universe. Conventional methods require domain experts (with knowledge of both machine learning and hardware architecture) to explore the huge design space smartly with rule-based heuristics. For instance, we should retain more bits in the first layer which extracts low level features and in the last layer which computes the final outputs; also, we should use more bits in the convolution layers than in the fully-connected layers because empirically, the convolution layers are more sensitive. As the neural network becomes deeper, the exploration space increases exponentially, which makes it infeasible to rely on hand-crafted strategies. Therefore, these *rule-based* quantization policies are usually sub-optimal, and they cannot generalize well from one model to another. In this paper, we would like to *automate* this exploration process by a *learning-based* framework.

Another challenge is how to measure the latency and the energy consumption of a given model on the hardware. A widely adopted approach (Howard et al 2017; Sandler et al 2018) is to rely on some proxy signals (e.g., FLOPs, number of memory references). However, as different hardware behaves very differently, the performance of a model on the hardware cannot always be accurately reflected by these proxy signals. Therefore, it is of great importance to directly *involve the hardware architecture into the loop*. Also, as demonstrated in Table 1, the quantization solution optimized on one hardware might not be optimal on the other,

which raises the demand for *specialized* policies for different hardware architectures.

To this end, we propose the **Hardware-Aware Automated Quantization (HAQ)** framework that leverages reinforcement learning to automatically predict the quantization policy given the hardware's feedback. The RL agent decides the bitwidth of a given neural network in a layer-wise manner. For each layer, the agent receives the layer configuration and statistics as observation, and it then outputs the action which is the bitwidth of weights and activations. We then leverage the hardware accelerator as the environment to obtain the *direct feedback from hardware* to guide the RL agent to satisfy the resource constraints. After all layers are quantized, we fine-tune the quantized model for one more epoch, and feed the validation accuracy after short-term retraining as the reward signal to our RL agent. During the exploration, we leverage the deep deterministic policy gradient (DDPG) (Lillicrap et al 2016) to supervise our RL agent. We studied the quantization policy on multiple hardware architectures: both cloud and edge neural network accelerators, with spatial or temporal multi-precision design.

The contribution of this paper has four aspects:

1. **Automation:** We propose an automated framework for quantization, which does not require domain experts and rule-based heuristics. It frees the human labor from exploring the vast search space of choosing bitwidths.
2. **Hardware-Aware:** Our framework integrates the hardware architecture into the loop so that it can directly reduce the latency, energy and storage on the target hardware instead of relying on some proxy signals.
3. **Specialization:** For different hardware architectures, our framework can offer a specialized quantization policy that's exactly tailored for the hardware architecture.
4. **Design Insights:** We interpreted the different quantization policies learned for different hardware architectures. Taking both computation and memory access into account, the interpretation offers insights on both neural network architecture and hardware architecture design.

## 2 Related Work

**Quantization.** There have been extensive explorations on compressing and accelerating deep neural networks using quantization. Han et al (2016) quantized the network weights to reduce the model size by rule-based strategies: *e.g.*, they used human heuristics to determine the bitwidths for convolution and fully-connected layers. Courbariaux et al (2016) binarized the network weights into  $\{-1, +1\}$ ; Rastegari et al (2016) and Zhou et al (2018) binarized each convolution filter into  $\{-w, +w\}$ ; Zhu et al (2017) mapped the network weights into  $\{-w_N, 0, +w_P\}$  using two bits; Zhou et al (2016)

used one bit for network weights and two bits for activations; Jacob et al (2018) made use of 8-bit integers for both weights and activations. We refer the reader to the survey paper by Krishnamoorthi (2018) for a more detailed overview. These conventional quantization methods either simply assign the same number of bits to all layers or require domain experts to determine the bitwidths for different layers, while our framework automates this design process, and our *learning-based* policy outperforms *rule-based* strategies.

**Automated ML.** Many researchers aimed to improve the performance of deep neural networks by searching the network architectures: Zoph and Le (2017) proposed the Neural Architecture Search (NAS) to explore and design the transformable network building blocks, and their network architecture outperforms several human designed networks; Liu et al (2018) introduced the Progressive NAS to accelerate the architecture search by  $5\times$  using sequential model-based optimization; Pham et al (2018) introduced the Efficient NAS to speed up the exploration by  $1000\times$  using parameter sharing; Cai et al (2018) introduced the path-level network transformation to search the tree-structured architecture space effectively. Motivated by these AutoML frameworks, He et al (2018) leveraged the reinforcement learning to automatically prune the convolution channels. Our framework further explores the automated quantization for network weights and activations, and it takes the hardware architectures into consideration.

**Efficient Models.** To facilitate the efficient deployment, researchers designed hardware-friendly approaches to slim neural network models. For instance, the coarse-grained channel pruning methods (He et al 2017; Liu et al 2017) prune away the entire channel of convolution kernels to achieve speedup. Recently, researchers have explicitly optimized for various aspects of hardware properties, including the inference latency and energy: Yang et al (2016) proposed the energy-aware pruning to directly optimize the energy consumption of neural networks; Yang et al (2018) reduced the inference time of neural networks on the mobile devices through a lookup table. Nevertheless, these methods are still rule-based and mostly focus on pruning. Our framework automates the quantization process by taking hardware-specific metric as direct rewards using a learning based method.

## 3 Approach

The overview of our proposed framework is in Figure 3. We model the quantization task as a reinforcement learning problem. We used the actor critic model with DDPG agent to give action: bits for each layer. We collect hardware counters, together with accuracy as direct rewards to search the optimal quantization policy for each layer. We have three

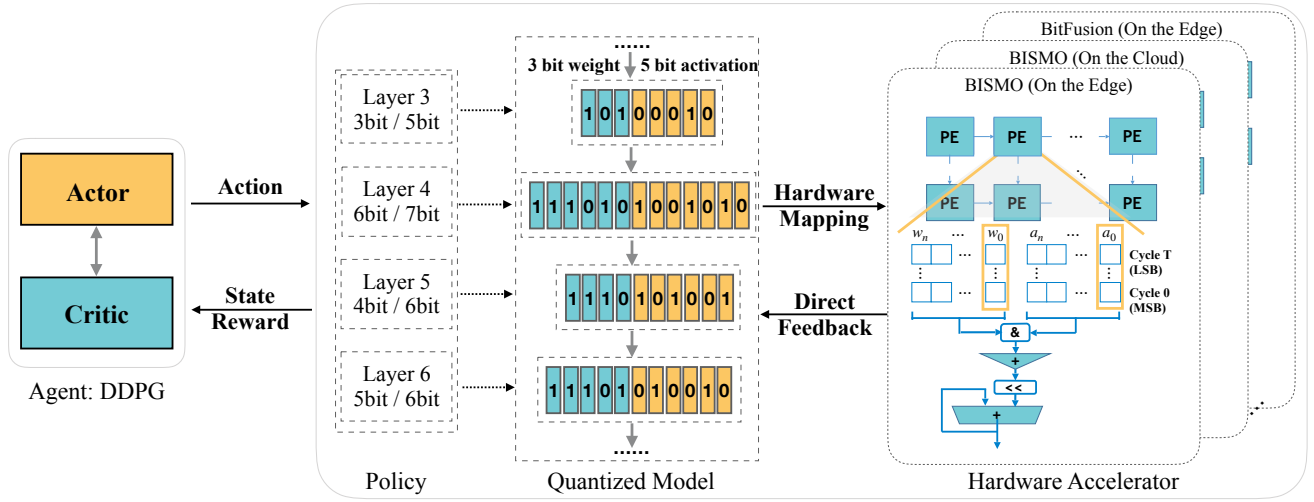


Fig. 3: An overview of our **Hardware-Aware Automated Quantization (HAQ)** framework. We leverage the reinforcement learning to automatically search over the huge quantization design space with hardware in the loop. The agent propose an optimal bitwidth allocation policy given the amount of computation resources (*i.e.*, latency, power, and model size). Our RL agent integrates the hardware accelerator into the exploration loop so that it can obtain the direct feedback from the hardware, instead of relying on indirect proxy signals.

hardware environments that covers edge and cloud, spatial and temporal architectures for multi-precision accelerator. Below describes the details of the RL formulation.

### 3.1 Observation (State Space)

Our agent processes the neural network in a layer-wise manner. For each layer, our agent takes two steps: one for weights, and one for activations. In this paper, we introduce a ten-dimensional feature vector  $O_k$  as our observation:

If the  $k^{\text{th}}$  layer is a convolution layer, the state  $O_k$  is

$$O_k = (k, c_{\text{in}}, c_{\text{out}}, s_{\text{kernel}}, s_{\text{stride}}, s_{\text{feat}}, n_{\text{params}}, i_{\text{dw}}, i_{\text{w/a}}, a_{k-1}), \quad (1)$$

where  $k$  is the layer index,  $c_{\text{in}}$  is #input channels,  $c_{\text{out}}$  is #output channels,  $s_{\text{kernel}}$  is kernel size,  $s_{\text{stride}}$  is the stride,  $s_{\text{feat}}$  is the input feature map size,  $n_{\text{params}}$  is #parameters,  $i_{\text{dw}}$  is a binary indicator for depthwise convolution,  $i_{\text{w/a}}$  is a binary indicator for weight/activation, and  $a_{k-1}$  is the action from the last time step.

If the  $k^{\text{th}}$  layer is a fully-connected layer, the state  $O_k$  is

$$O_k = (k, h_{\text{in}}, h_{\text{out}}, 1, 0, s_{\text{feat}}, n_{\text{params}}, 0, i_{\text{w/a}}, a_{k-1}), \quad (2)$$

where  $k$  is the layer index,  $h_{\text{in}}$  is #input hidden units,  $h_{\text{out}}$  is #output hidden units,  $s_{\text{feat}}$  is the size of input feature vector,  $n_{\text{params}}$  is #parameters,  $i_{\text{w/a}}$  is a binary indicator for weight/activation, and  $a_{k-1}$  is the action from the last step.

For each dimension in the observation vector  $O_k$ , we normalize it into  $[0, 1]$  to make them in the same scale.

### 3.2 Action Space

We use a *continuous* action space to determine the bitwidth. The reason that we do not use a *discrete* action space is because it loses the relative order: *e.g.*, 2-bit quantization is more aggressive than 4-bit and even more than 8-bit. At the  $k^{\text{th}}$  time step, we take the continuous action  $a_k$  (which is in the range of  $[0, 1]$ ), and round it into the discrete bitwidth value  $b_k$ :

$$b_k = \text{round}(b_{\min} - 0.5 + a_k \times (b_{\max} - b_{\min} + 1)), \quad (3)$$

where  $b_{\min}$  and  $b_{\max}$  denote the min and max bitwidth (in our experiments, we set  $b_{\min}$  to 2 and  $b_{\max}$  to 8).

**Resource Constraints.** In real-world applications, we have limited computation budgets (*i.e.*, latency, energy, and model size). We would like to find the quantization policy with the best performance given the constraint.

We encourage our agent to meet the computation budget by limiting the action space. After our RL agent gives actions  $\{a_k\}$  to all layers, we measure the amount of resources that will be used by the quantized model. The feedback is directly obtained from the hardware accelerator, which we will discuss in Section 3.3. If the current policy exceeds our resource budget (on latency, energy or model size), we will sequentially decrease the bitwidth of each layer until the constraint is finally satisfied.



### 3.3 Direct Feedback from Hardware Accelerators

An intuitive feedback to our RL agent can be FLOPs or the model size. However, as these proxy signals are indirect, they are not equal to the performance (*i.e.*, latency, energy consumption) on the hardware. Cache locality, number of kernel calls, memory bandwidth all matters. Proxy feedback can not model these hardware functionality to find the specialized strategies (see Table 1).

Instead, we use direct latency and energy feedback from hardware accelerators to optimize the performance. In simulators, the latency is approximated as the sum of the computation time, the stall caused by the memory access and some other overheads:

$$T = T_{\text{computation}} + T_{\text{stall}} + T_{\text{overhead}}, \quad (4)$$

and the energy consumption is modeled as the sum of the logic circuits and memory:

$$E = E_{\text{memory access per bit}} \times S_{\text{total memory access size}} + P_{\text{dynamic}} \times T_{\text{execution}}. \quad (5)$$

The direct feedback from the hardware simulator is very crucial as it enables our RL agent to determine the bitwidth allocation policy from the subtle differences between different layers: *e.g.*, the vanilla convolution has more data reuse and better locality, while the depthwise convolution (Chollet 2017) has less reuse and worse locality, which makes it memory bounded.

### 3.4 Quantization

We linearly quantize the weights and activations of each layer using the action  $a_k$  given by our RL agent, as linearly quantized model only need fixed point arithmetic unit which is more efficient to implement on the hardware than the  $k$ -means quantization.

Specifically, for each weight value  $w$  in the  $k^{\text{th}}$  layer, we first truncate it into the range of  $[-c, c]$ , and we then quantize it linearly into  $a_k$  bits:

$$\text{quantize}(w, a_k, c) = \text{round}(\text{clamp}(w, c)/s) \times s, \quad (6)$$

where  $\text{clamp}(\cdot, x)$  is to truncate the values into  $[-x, x]$ , and the scaling factor  $s$  is defined as

$$s = c/(2^{a_k-1} - 1). \quad (7)$$

In this paper, we choose the value of  $c$  by finding the optimal value  $x$  that minimizes the KL-divergence between the original weight distribution  $W_k$  and the quantized weight distribution  $\text{quantize}(W_k, a_k, x)$ :

$$c = \arg \min_x D_{\text{KL}}(W_k \parallel \text{quantize}(W_k, a_k, x)), \quad (8)$$

where  $D_{\text{KL}}(\cdot \parallel \cdot)$  is the KL-divergence that characterizes the distance between two distributions. As for activations, we quantize the values similarly except that we truncate them into the range of  $[0, c]$ , not  $[-c, c]$  since the activation values (which are the outputs of the ReLU layers) are non-negative. This calibration based on KL-divergence enables us to make use of the pretrained models rather than training the models from scratch so that it can reduce the training time significantly. As for the overhead, we only use 64 images to calibrate once at the beginning, which is negligible compared to the whole training.

### 3.5 Reward Signal

After quantization, we retrain the quantized model for one more epoch to recover the performance. As we impose the resource constraints by limiting the action space, we define our reward function  $R$  to be only related to the accuracy:

$$R = \lambda \times (\text{accuracy}_{\text{quant}} - \text{accuracy}_{\text{origin}}), \quad (9)$$

where  $\text{accuracy}_{\text{origin}}$  is the top-1 classification accuracy of the full-precision model on the training set,  $\text{accuracy}_{\text{quant}}$  is the top-1 classification accuracy of the quantized model after finetuning, and  $\lambda$  is a scaling factor which is set to 0.1 in our experiments.

### 3.6 Agent

In our environment, one step means that our agent makes an action to decide the number of bits assigned to the weights or activations of a specific layer, while one episode is composed of multiple steps, where our RL agent makes actions to all layers. As for our RL agent, we leverage the deep deterministic policy gradient (DDPG) (Lillicrap et al 2016), which is an off-policy actor-critic algorithm for continuous control problem. We apply a variant form of the Bellman's Equation, where each transition in an episode is defined as

$$T_k = (O_k, a_k, R, O_{k+1}). \quad (10)$$

During exploration, the  $Q$ -function is computed as

$$\hat{Q}_k = R_k - B + \gamma \times Q(O_{k+1}, w(O_{k+1}) \mid \theta^Q), \quad (11)$$

and the gradient signal can be approximated using

$$L = \frac{1}{N_s} \sum_{k=1}^{N_s} (\hat{Q}_k - Q(O_k, a_k \mid \theta^Q))^2, \quad (12)$$

where  $N_s$  denotes the number of steps in this episode, and the baseline  $B$  is defined as an exponential moving average of all previous rewards in order to reduce the variance of the gradient estimation. The discount factor  $\gamma$  is set to 1 since we assume that the action made for each layer should contribute

equally to the final result. Moreover, as the number of steps is always finite (bounded by the number of layers), the sum of the rewards will not explode.

### 3.7 Implementation Details

In this section, we present some implementation details about the RL agent, exploration and finetuning quantized models.

**Agent.** The DDPG agent consists of an actor network and a critic network. Both follow the same network architecture: each network has 3 fully-connected layers with the hidden size of [400, 300]. For the actor network, the input is the state vector, and the output action is normalized to [0, 1] by the sigmoid function; while for the critic network, the input is a vector concatenated by state and its corresponding action produced by actor.

**Exploration.** Optimization of the DDPG agent is carried out using ADAM (Kingma and Ba 2015) with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . We use a fixed learning rate of  $10^{-4}$  for the actor network and  $10^{-3}$  for the critic network. During exploration, we employ the following stochastic process of the noise:

$$w'(O_k) \sim N_{\text{trunc}}(w(O_k | \theta_k^w), \sigma^2, 0, 1), \quad (13)$$

where  $N_{\text{trunc}}(\mu, \sigma, a, b)$  is the truncated normal distribution, and  $w$  is the model weights. The noise  $\sigma$  is initialized as 0.5, and after each episode, the noise is decayed exponentially with a decay rate of 0.99.

**Finetuning.** During exploration, we finetune the quantized model for one epoch to help recover the performance (using SGD with a fixed learning rate of  $10^{-3}$  and momentum of 0.9). We randomly select 100 categories from ImageNet (Deng et al 2009) to accelerate the model finetuning during exploration. After exploration, we quantize the model with our best policy and finetune it on the full dataset.

## 4 Experiments

We conduct extensive experiments to demonstrate the consistent effectiveness of our framework for multiple objectives: *latency, energy, model size, and accuracy*.

**Datasets and Models.** Our experiments are performed on the ImageNet (Deng et al 2009) dataset. As our focus is on more efficient models, we extensively study the quantization of MobileNet-V1 (Howard et al 2017) and MobileNet-V2 (Sandler et al 2018). Both MobileNets are inspired from the depthwise separable convolutions (Chollet 2017) and replace the regular convolutions with the *pointwise* and *depthwise* convolutions: MobileNet-V1 stacks multiple “*depthwise* – *pointwise*” blocks repeatedly; while MobileNet-V2

	Hardware	Batch	PE Array	AXI port	Block RAM
Edge	Zynq-7020	1	8×8	4×64b	140×36Kb
Cloud	VU9P	16	16×16	4×256b	2160×36Kb

Table 2: The configurations of edge and cloud accelerators.

uses the “*pointwise* – *depthwise* – *pointwise*” blocks as its basic building primitives.

### 4.1 Latency-Constrained Quantization

We first evaluate our framework under latency constraints on two representative hardware architectures: spatial and temporal architectures for multi-precision CNN:

**Temporal Architecture.** Bit-Serial Matrix Multiplication Overlay (BISMO) proposed by Umuroglu et al (2018) is a classic temporal design of neural network accelerator on FPGA. It introduces bit-serial multipliers which are fed with one-bit digits from 256 weights and corresponding activations in parallel at one time and accumulates their partial products by shifting over time.

**Spatial Architecture.** BitFusion architecture proposed by Sharma et al (2018) is a state-of-the-art spatial ASIC design for neural network accelerator. It employs a 2D systolic array of Fusion Units which spatially sum the shifted partial products of two-bit elements from weights and activations.

#### 4.1.1 Quantization policy for BISMO Architecture

Inferencing the neural networks on edge devices and cloud servers can be quite different, since the tasks on the cloud servers are more intensive and the edge devices are usually limited to low computation resources and memory bandwidth. We use Xilinx Zynq-7020 FPGA (Xilinx 2018b) as our edge device and Xilinx VU9P (Xilinx 2018a) as our cloud device. Table 2 shows our experiment configurations on these two platforms along with their available resources.

As for comparison, we adopt the PACT (Choi et al 2018) as our baseline, which uses the same number of bits for all layers except for the first layer which extracts the low level features, they use 8 bits for both weights and activations as it has fewer parameters and is very sensitive to errors. We follow a similar setup as PACT: we quantize the weights and activations of the first and last layer to 8 bits and explore the bitwidth allocation policy for all the other layers.

Under the same latency, HAQ consistently achieved better accuracy than the baseline on both the cloud and the edge (Table 3). With similar accuracy, HAQ can reduce the latency by  $1.4\times$  to  $1.95\times$  compared with the baseline.



	Bitwidths	Edge Accelerator						Cloud Accelerator					
		MobileNet-V1			MobileNet-V2			MobileNet-V1			MobileNet-V2		
		Acc.-1	Acc.-5	Latency	Acc.-1	Acc.-5	Latency	Acc.-1	Acc.-5	Latency	Acc.-1	Acc.-5	Latency
PACT	4 bits	62.44	84.19	45.45 ms	61.39	83.72	52.15 ms	62.44	84.19	57.49 ms	61.39	83.72	74.46 ms
Ours	<i>flexible</i>	<b>67.40</b>	<b>87.90</b>	45.51 ms	<b>66.99</b>	<b>87.33</b>	52.12 ms	<b>65.33</b>	<b>86.60</b>	57.40 ms	<b>67.01</b>	<b>87.46</b>	73.97 ms
PACT	5 bits	67.00	87.65	57.75 ms	68.84	88.58	66.94 ms	67.00	87.65	77.52 ms	68.84	88.58	99.43 ms
Ours	<i>flexible</i>	<b>70.58</b>	<b>89.77</b>	57.70 ms	<b>70.90</b>	<b>89.91</b>	66.92 ms	<b>69.97</b>	<b>89.37</b>	77.49 ms	<b>69.45</b>	<b>88.94</b>	99.07 ms
PACT	6 bits	70.46	89.59	70.67 ms	71.25	90.00	82.49 ms	70.46	89.59	99.86 ms	71.25	90.00	127.07 ms
Ours	<i>flexible</i>	<b>71.20</b>	<b>90.19</b>	70.35 ms	<b>71.89</b>	<b>90.36</b>	82.34 ms	<b>71.20</b>	<b>90.08</b>	99.66 ms	<b>71.85</b>	<b>90.24</b>	127.03 ms
Original	8 bits	70.82	89.85	96.20 ms	71.81	90.25	115.84 ms	70.82	89.85	151.09 ms	71.81	90.25	189.82 ms

Table 3: Latency-constrained quantization on BISMO (edge accelerator and cloud accelerator) on ImageNet. Our framework can reduce the latency by  $1.4\times$  to  $1.95\times$  with negligible loss of accuracy compared with the fixed bitwidth (8 bits) quantization.

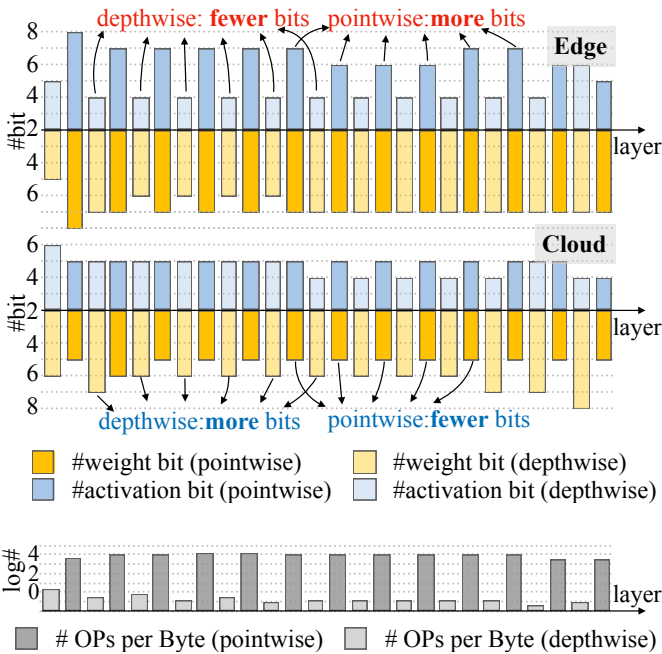


Fig. 4: Quantization policy under latency constraints for MobileNet-V1 on BISMO (57.7 ms for the edge accelerator and 77.5 ms for the cloud accelerator). On edge accelerator, our agent allocates *fewer* activation bits to the depthwise convolutions, which echos that the depthwise convolutions are memory bounded and the activations dominates the memory access. On cloud accelerator, our agent allocates *more* bits to the depthwise convolutions and allocates *fewer* bits to the pointwise convolutions, as cloud device has more memory bandwidth and higher parallelism, the network appears to be computation bounded.

**Interpreting the quantization policy.** Our agent gave quite different quantization policy for edge and cloud accelera-

	Weights	Activations	Acc.-1	Acc.-5	Latency
PACT	4 bits	4 bits	62.44	84.19	7.86 ms
Ours	<i>flexible</i>	<i>flexible</i>	<b>67.45</b>	<b>87.85</b>	7.86 ms
PACT	6 bits	4 bits	67.51	87.84	11.10 ms
Ours	<i>flexible</i>	<i>flexible</i>	<b>70.40</b>	<b>89.69</b>	11.09 ms
PACT	6 bits	6 bits	70.46	89.59	19.99 ms
Ours	<i>flexible</i>	<i>flexible</i>	<b>70.90</b>	<b>89.95</b>	19.98 ms
Original	8 bits	8 bits	70.82	89.85	20.08 ms

Table 4: Latency-constrained quantization on BitFusion (MobileNet-V1 on ImageNet). Our framework can reduce the latency by  $2\times$  with almost no loss of accuracy compared with the fixed bitwidth (8 bits) quantization.

tors. For the activations, the depthwise convolution layers are assigned much less bitwidth than the pointwise layers on the edge; while on the cloud device, the bitwidth of these two types of layers are similar to each other. For weights, the bitwidth of these types of layers are nearly the same on the edge; while on the cloud, the depthwise convolution layers are assigned much more bitwidth than the pointwise convolution layers.

We explain the difference of quantization policy between edge and cloud by the roofline model (Williams et al 2009). Many previous works use FLOPs or BitOPs as metrics to measure computation complexity. However, they are not able to directly reflect the latency, since there are many other factors influencing the hardware performance, such as memory access cost and degree of parallelism (Sandler et al 2018; Liu et al 2017). Taking computation and memory access into account, the roofline model assumes that applications are either computation-bound or memory bandwidth-bound, if not fitting in on-chip caches, depending on their operation inten-

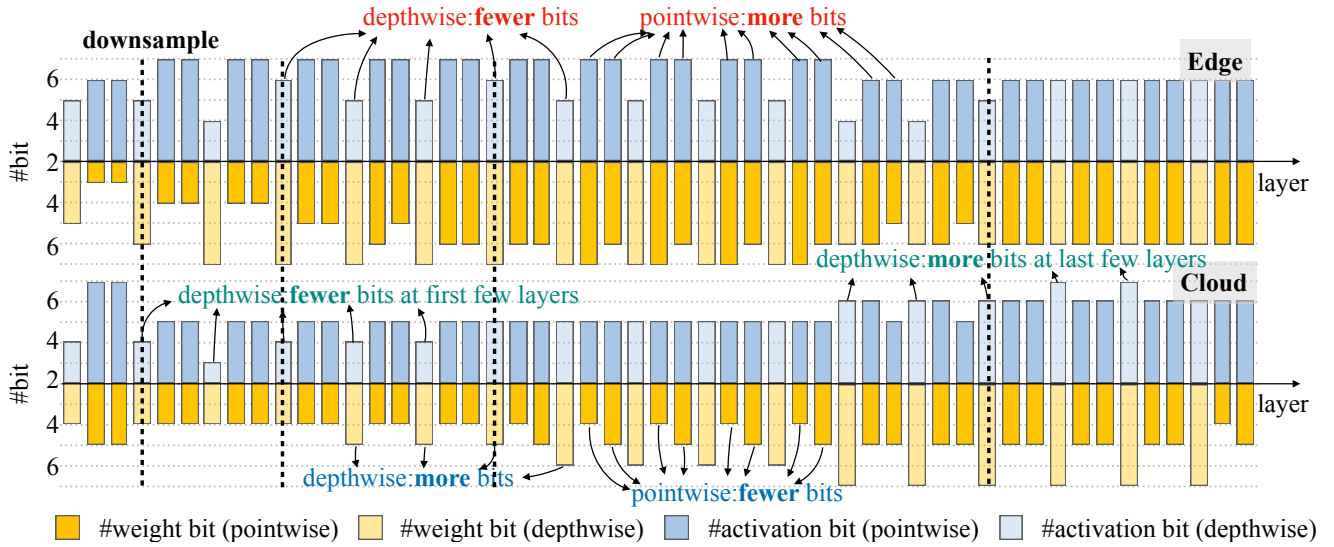


Fig. 5: Quantization policy under latency constraints for MobileNet-V2 on BISMO (66.9 ms for the edge accelerator and 99.1 ms for the cloud accelerator). Similar to Figure 4, depthwise layer is assigned with fewer bits on the edge accelerator, and pointwise layer is assigned with fewer bits on the cloud accelerator.

	Weights	Activations	Acc.-1	Acc.-5	Energy
PACT	4 bits	4 bits	62.44	84.19	13.47 mJ
Ours	<i>flexible</i>	<i>flexible</i>	<b>64.78</b>	<b>85.85</b>	13.69 mJ
PACT	6 bits	4 bits	67.51	87.84	16.57 mJ
Ours	<i>flexible</i>	<i>flexible</i>	<b>70.37</b>	<b>89.40</b>	16.30 mJ
PACT	6 bits	6 bits	70.46	89.59	26.80 mJ
Ours	<i>flexible</i>	<i>flexible</i>	<b>70.90</b>	<b>89.73</b>	26.67 mJ
Original	8 bits	8 bits	70.82	89.95	31.03 mJ

Table 5: Energy-constrained quantization on BitFusion (MobileNet-V1 on ImageNet). Our framework reduces the power consumption by  $2\times$  with nearly no loss of accuracy compared with the fixed bitwidth quantization.

sity. Operation intensity is measured as operations (MACs in neural networks) per DRAM byte accessed. A lower operation intensity indicates that the model suffers more from the memory access.

The bottom of Figure 4 shows the operation intensities (operations per byte) of convolution layers in the MobileNet-V1. Depthwise convolution is a memory bounded operation, and the pointwise convolution is a computation bounded operation. Our experiments show that when running MobileNet-V1 on the edge devices with small batch size, its latency is dominated by the depthwise convolution layers. Since the feature maps take a major proportion in the memory of depthwise convolution layers, our agent gives the activations fewer bits. In contrast, when running MobileNet-V1 on the cloud with large batch size, both two types of layers have nearly

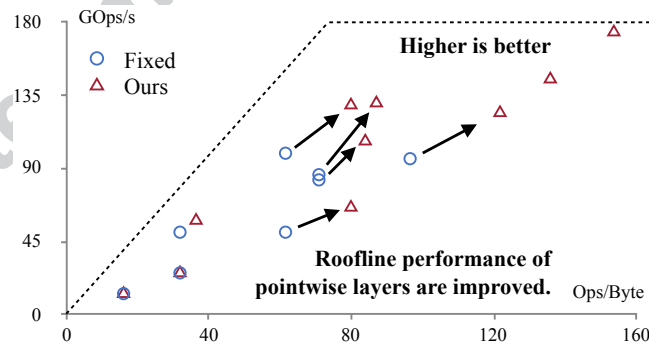


Fig. 6: Roofline model of pointwise layers in MobileNet-V1 (fixed-bitwidth in blue and mixed-precision in red). Our mixed-precision framework improves the roofline performance by a large margin.

the equal influence on the speed. Therefore, our agent tries to reduce the bitwidth of both activation and weights. However, since the weights of the depthwise convolution layers takes a small proportion of the memory, our agent increases their bitwidth to preserve the network accuracy at low memory overhead. Figure 6 shows the roofline model before and after HAQ. HAQ gives more reasonable policy to allocate the bits for each layer and pushes all the points to the upper right corner that is more efficient.

On edge accelerator, which has much less memory bandwidth, our RL agent allocates *fewer* activation bits to the depthwise convolutions since the activations dominates the memory access. On cloud accelerator, which has more memory bandwidth, our agent allocates *more* bits to the depth-

wise convolutions and allocates *fewer* bits to the pointwise convolutions to prevent it from being computation bounded.

A similar phenomenon can be observed in Figure 5 for quantizing MobileNet-V2. Moreover, since the activation size in the deeper layers gets smaller, they get assigned more bits. Another interesting phenomenon we discover in Figure 5 is that the downsample layer gets assigned more activation bits than the adjacent layer. This is because down-sampled layers are more prone to lose information, so our agent learns to assign more bits to the activations to compensate.

#### 4.1.2 Quantization policy for BitFusion Architecture

In order to demonstrate the effectiveness of our framework on different hardware architectures, we further compare our framework with PACT (Choi et al 2018) under the latency constraints on the BitFusion (Sharma et al 2018) architecture. As demonstrated in Table 4, our framework performs much better than the hand-craft policy with the same latency. Also, it can achieve almost no degradation of accuracy with only half of the latency used by the original MobileNet-V1 model (from 20.08 to 11.09 ms). Therefore, our framework is indeed very flexible and can be applied to different hardware platforms.

#### 4.2 Energy-Constrained Quantization

We then evaluate our framework under the energy constraints on the BitFusion (Sharma et al 2018) architecture. Similar to the latency-constrained experiments, we compare our framework with PACT (Choi et al 2018) which uses fixed number of bits for both weights and activations. From Table 5, we can clearly see that our framework outperforms the rule-based baseline: it achieves much better performance while consuming similar amount of energy. In particular, our framework is able to achieve almost no loss of accuracy with nearly half of the energy consumption of the original MobileNet-V1 model (from 31.03 to 16.57 mJ), which suggests that flexible bitwidths can indeed help reduce the energy consumption.

#### 4.3 Model Size-Constrained Quantization

We further evaluate our framework under the model size constraints. Following Han et al (2016), we employ the  $k$ -means algorithm to quantize the values into  $k$  centroids instead of using the linear quantization for compression.

We compare our framework with Deep Compression (Han et al 2016) on MobileNets and ResNet-50. From Table 6, we can see that our framework performs much better than Deep Compression: it achieves higher accuracy with the same model

size. For MobileNets which are already very compactly designed, our framework can preserve the performance to some extent; while Deep Compression significantly degrades the performance especially when the model size is very small. For instance, when Deep Compression quantizes the weights of MobileNet-V1 to 2 bits, the accuracy drops significantly from 70.90 to 37.62; while our framework can still achieve 57.14 of accuracy with the same model size, which is because our framework makes full use of the flexible bitwidths.

**Discussions.** In Figure 7, we visualize the bitwidth allocation strategy for MobileNet-V2. From this figure, we can observe that our framework assigns *more* bitwidths to the weights in depthwise convolution layers than pointwise convolution layers. Intuitively, this is because the number of parameters in the former is much smaller than the latter. Comparing Figure 5 and Figure 7, the policies are drastically different under different optimization objectives (**fewer** bitwidths for depthwise convolutions under *latency* optimization, **more** bitwidths for depthwise convolutions under *model size* optimization). Our framework succeeds in learning to adjust its bitwidth policy under different constraints.

#### 4.4 Accuracy-Guaranteed Quantization

Apart from the resource-constrained experiments, we also evaluate our framework under the accuracy-guaranteed scenario, that is to say, we aim to minimize the resource (*i.e.*, latency and energy) we use while preserving the accuracy.

Instead of using the resource-constrained action space in Section 3.2, we define a new reward function  $R$  that takes both the resource and the accuracy into consideration:

$$R = R_{\text{latency}} + R_{\text{energy}} + R_{\text{accuracy}}. \quad (14)$$

Here, the reward functions  $R_*$  are defined to encourage each term to be as good as possible:

$$\begin{aligned} R_{\text{latency}} &= \lambda_{\text{latency}} \times (\text{latency}_{\text{quant}} - \text{latency}_{\text{origin}}), \\ R_{\text{energy}} &= \lambda_{\text{energy}} \times (\text{energy}_{\text{quant}} - \text{energy}_{\text{origin}}), \\ R_{\text{accuracy}} &= \lambda_{\text{accuracy}} \times (\text{accuracy}_{\text{quant}} - \text{accuracy}_{\text{origin}}), \end{aligned} \quad (15)$$

where  $\lambda_*$  are scaling factors that encourage the RL agent to trade off between the computation resource and the accuracy. We set  $\lambda_{\text{latency}}$  and  $\lambda_{\text{energy}}$  to 1, and  $\lambda_{\text{accuracy}}$  to 20 in our experiments to ensure that our RL agent will prioritize the accuracy over the computation resource.

We choose to perform our experiments on a ten-category subset of ImageNet as it is very challenging to preserve the accuracy while reducing the computation resource. In Figure 8, we illustrate the exploration curves of our RL agents, and we can observe that the exploration process can be divided into three phases. In the first phase, our RL agent puts

	Weights	MobileNet-V1			MobileNet-V2			ResNet-50		
		Acc.-1	Acc.-5	Model Size	Acc.-1	Acc.-5	Model Size	Acc.-1	Acc.-5	Model Size
Han et al (2016)	2 bits	37.62	64.31	1.09 MB	58.07	81.24	0.96 MB	68.95	88.68	6.32 MB
Ours	<i>flexible</i>	<b>57.14</b>	<b>81.87</b>	1.09 MB	<b>66.75</b>	<b>87.32</b>	0.95 MB	<b>70.63</b>	<b>89.93</b>	6.30 MB
Han et al (2016)	3 bits	65.93	86.85	1.60 MB	68.00	87.96	1.38 MB	75.10	92.33	9.36 MB
Ours	<i>flexible</i>	<b>67.66</b>	<b>88.21</b>	1.58 MB	<b>70.90</b>	<b>89.76</b>	1.38 MB	<b>75.30</b>	<b>92.45</b>	9.22 MB
Han et al (2016)	4 bits	71.14	89.84	2.10 MB	71.24	89.93	1.79 MB	<b>76.15</b>	92.88	12.40 MB
Ours	<i>flexible</i>	<b>71.74</b>	<b>90.36</b>	2.07 MB	<b>71.47</b>	<b>90.23</b>	1.79 MB	76.14	<b>92.89</b>	12.14 MB
Original	32 bits	70.90	89.90	16.14 MB	71.87	90.32	13.37 MB	76.15	92.86	97.49 MB

Table 6: Model size-constrained quantization on ImageNet. Compared with Deep Compression (Han 2017), our framework achieves higher accuracy under similar model size (especially under high compression ratio).

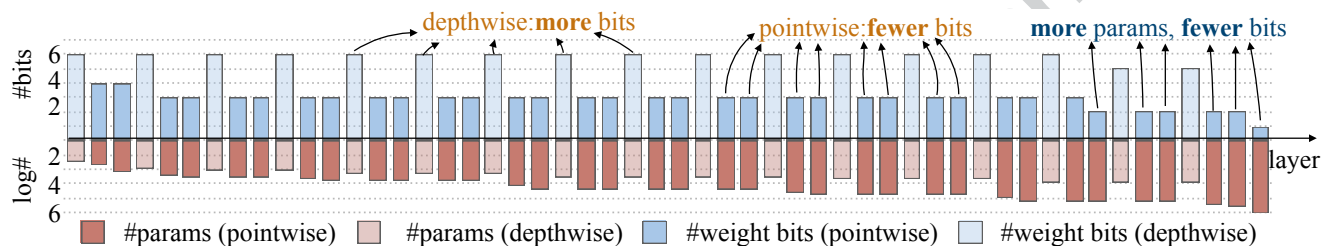


Fig. 7: Quantization policy under model size constraints for MobileNet-V2. Our RL agent allocates *more* bits to the depthwise convolutions, since depthwise convolutions have *fewer* number of parameters.

its focus on the accuracy: it tries to preserve the accuracy while completely ignoring the latency and the energy consumption. In the second phase, the accuracy begins to be more stable, and our RL agent starts to aggressively reduce the latency and the energy. In the third phase, our RL agent converges to the best policy it has found. We conjecture that this interesting behavior is because that the scaling factor  $\lambda_{\text{accuracy}}$  is much larger than the other two, which encourages our agent to first optimize the value of accuracy, and after its value has been stabilized, our agent then attempts to reduce the value of latency and energy to further optimize the reward value (see the reward curve in Figure 8).

#### 4.5 Integration with Architecture Search and Pruning

We integrate the neural architecture search (Cai et al 2019) and automated channel pruning (He et al 2018) with HAQ to demonstrate that our method is orthogonal to other AutoML methods. In Figure 9, we observe significant improvements over the baselines including ProxylessNAS (with 8-bit quantization), ProxylessNAS + AMC (with 8-bit quantization), MobileNetV2 (with 4-bit / 6-bit quantization), and MobileNetV2 + HAQ (with mixed-precision quantization).

	Search Time	Acc.-1	Acc.-5	Latency
ES	<b>17 hours</b>	65.73	86.81	45.45 ms
BO	74 hours	66.28	87.22	45.47 ms
RL (Ours)	<b>17 hours</b>	<b>67.40</b>	<b>87.90</b>	45.51 ms
ES	<b>17 hours</b>	69.11	88.80	57.73 ms
BO	74 hours	70.40	89.56	57.68 ms
RL (Ours)	<b>17 hours</b>	<b>70.58</b>	<b>89.77</b>	57.70 ms

Table 7: Comparison between different optimization methods. RL outperforms EA and BO and is  $4\times$  faster than BO in terms of the total search time.

## 5 Analysis

In this section, we first compare with sample efficiency of different optimization methods; then, we show the generalization and transfer learning ability of our framework; we finally interpret the quantization policy given by HAQ.

### 5.1 Optimization Methods

We leverage the reinforcement learning (RL) as our optimization method. In addition, we also compared with other optimizers including Bayesian optimization (BO) and evolutionary algorithm (EA). Similar to the configurations of

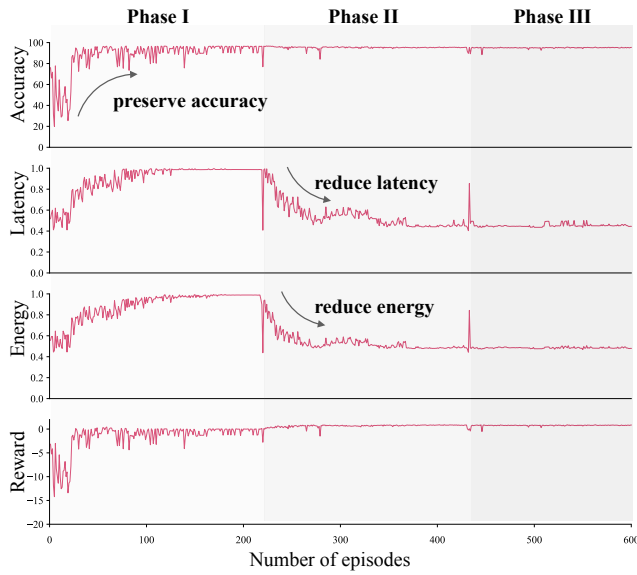


Fig. 8: Exploration curves of accuracy-guaranteed quantization for MobileNet-V1. Our RL agent first tries to preserve the accuracy while completely ignoring the latency and the energy consumption; after the accuracy begins to be more stable, it starts to aggressively reduce the latency and the energy; it finally converges to the best policy it has found.

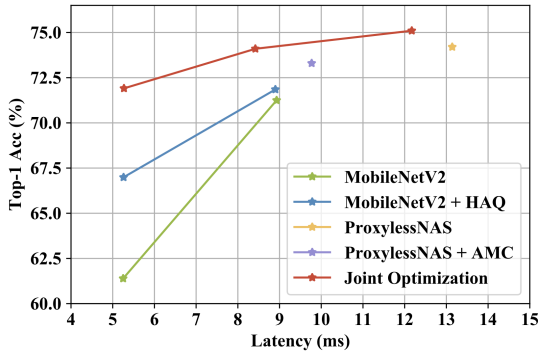


Fig. 9: Integrating NAS and AMC with HAQ together further improves the accuracy-latency trade-off by a significant margin.

RL, we model the outputs of BO and EA as the number of bits of different layers and the objectives of BO and EA as maximizing the validation accuracy of the quantized model.

As a fair comparison, we executed in total 600 runs (samples) for each optimization method. The performance comparison is in Table 7. All experiments are conducted on the BISMO hardware with MobileNet-V1. We observe that RL performs much better than EA and BO and is 4× faster than BO in terms of the total search time. This, we believe, is be-

	Bitwidth	Acc.-1	Acc.-5	Latency
PACT	4 bits	61.39	83.72	52.15 ms
Ours (search for V2)	<i>flexible</i>	66.99	87.33	52.12 ms
Ours (transfer from V1)	<i>flexible</i>	65.80	86.60	52.06 ms

	Bitwidth	Acc.-1	Acc.-5	Latency
PACT	5 bits	68.84	88.58	66.94 ms
Ours (search for V2)	<i>flexible</i>	70.90	89.91	66.92 ms
Ours (transfer from V1)	<i>flexible</i>	69.90	89.24	66.93 ms

Table 8: Comparisons between our agent’s transfer results (from MobileNet-V1 to MobileNet-V2), its direct search results on MobileNet-V2, and the fixed-bitwidth baseline (PACT). Our RL agent is able to generalize well to different network architectures: our quantization policy transferred from V1 to V2 performs better than the fixed-bitwidth baseline and is only slightly worse than the quantization policy directly searched for V2.

Constraint	Bitwidth	BitOPs	latency	Acc.-1	Acc.-5
BitOPs	<i>flexible</i>	8.17 G	85.06 ms	70.29	89.52
Latency	<i>flexible</i>	<b>8.01 G</b>	<b>66.92 ms</b>	<b>70.90</b>	<b>89.91</b>

BitOPs	<i>flexible</i>	11.36 G	97.99 ms	71.41	90.12
Latency	<i>flexible</i>	<b>11.17 G</b>	<b>82.34 ms</b>	<b>71.89</b>	<b>90.36</b>

Table 9: BitOPs-constrained quantization on Bismo (MobileNet-V2 on ImageNet).

cause BO and EA do not make use of the state encoding, which might lead to a worse sample efficiency.

## 5.2 Generalization and Transfer Learning

Another merit of the reinforcement learning is that its agent is able to generalize to different environments (*i.e.*, network architectures). In order to evaluate the transfer ability of our framework, we first train our agent for MobileNet-V1 under the latency constraint, and we then directly evaluate our agent on MobileNet-V2 by feeding its network architecture information in. In Table 8, we compare our agent’s transfer results (from V1 to V2) with its direct search results (for V2) and the fixed-bitwidth baseline (PACT). Our quantization policy transferred from V1 to V2 still performs better than the fixed-bitwidth baseline and is only slightly worse than the quantization policy directly searched for V2. This experiment validates that our RL agent generalize well to different network architectures.

## 5.3 Hardware-Awareness is Important

To evaluate the necessity of involving the hardware in the loop, we replace the Bismo accelerator with theoretical BitOPs



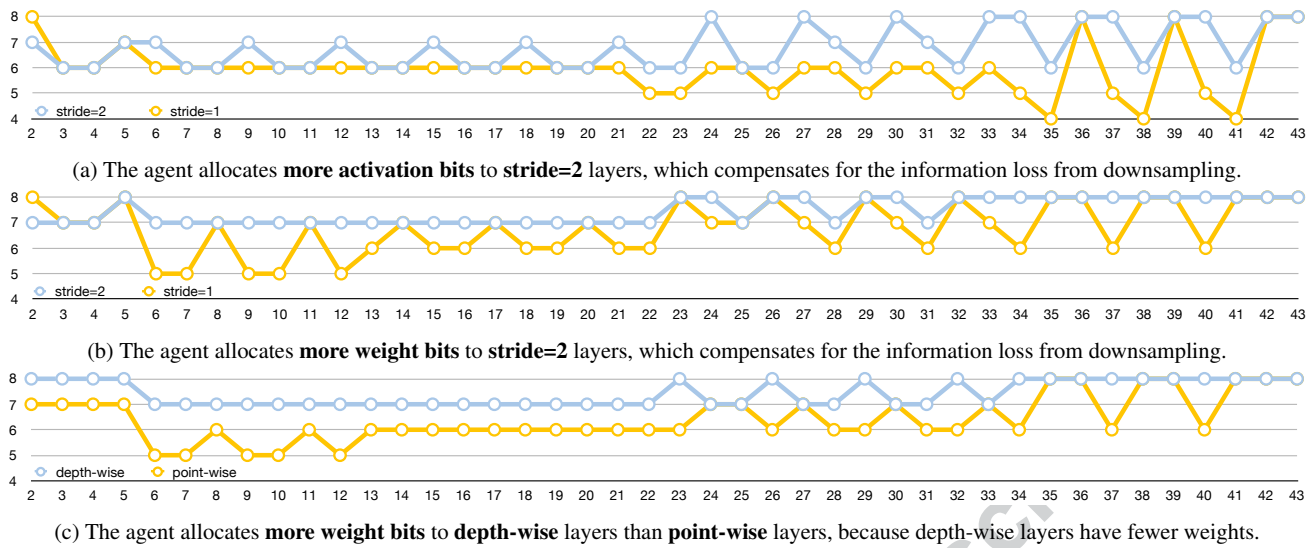


Fig. 10: We change only one dimension of the state vector, and run the actor network again to observe how the action changes across different layers.

analysis, which calculates the latency by  $FLOPs/s \times Bit_{weight} \times Bit_{activation}$  for each layer. The results are listed in Table 9, which shows that under similar BitOPs constraints, BitOPs constrained experiments get worse latency than the experiments with hardware-in-the-loop, and BitOPs constrained experiments could achieve better accuracy when the constrain is tight. The reason is that BitOPs and hardware latency are not linearly correlated, so the similar BitOPs may correspond to totally different latency if the layer is memory bottlenecked. The agent chooses to give more bits to depth-wise layers which have less FLOPs, but more memory access (which means more latency in edge hardware).

	Weights	Activations	Acc.-1	Acc.-5	Latency
PACT	2 bits	4 bits	74.06	91.78	80.03 ms
Ours	<i>flexible</i>	<i>flexible</i>	<b>74.42</b>	<b>91.92</b>	80.38 ms
PACT	3 bits	3 bits	74.73	92.11	85.49 ms
Ours	<i>flexible</i>	<i>flexible</i>	<b>74.98</b>	<b>92.37</b>	84.97 ms
PACT	4 bits	4 bits	76.17	93.03	128.55 ms
Ours	<i>flexible</i>	<i>flexible</i>	<b>76.22</b>	<b>93.15</b>	129.55 ms
Original	8 bits	8 bits	76.64	93.26	446.96 ms

Table 10: Latency-constrained quantization on Bismo (Resnet-50 on ImageNet).

## 5.4 Performance on Large Model

We evaluate our framework on a larger ResNet-50 model with the same search scheme and finetune policy as in Table 3, and the hardware platform is Bismo edge hardware simulator. Table 10 shows that the improvement of HAQ is

not remarkable on the large model like ResNet50. The reason is that ResNet-50 is highly redundant, even PACT can already quantize it to 4 bits without much accuracy loss. Therefore, there is not much room for HAQ to further improve it.

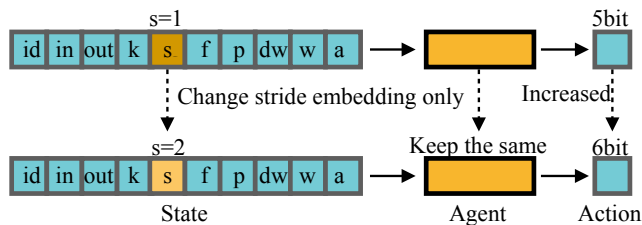


Fig. 11: Policy interpretation. From the model, we first select several layers; then, we only change (flip) one dimension in the state vector; finally, we run our RL agent's actor network again to see how that particular factor affects its decision.

## 5.5 Policy Interpretation

In Section 4, we provided intuitive explanations of our agent's policies. In this section, we quantitatively interpret our agent's quantization policy. As illustrated in Figure 11, we first select several layers from MobileNet-V2; then for each layer, we change (flip) only one dimension in the state vector (in the example, changing the convolution stride from 1 to 2); finally, we run feedforward on our RL agent's actor network again to see how that particular factor (*i.e.*, depthwise, downsample) affects its decision.

From Figure 10, we can clearly observe that some factors will affect the actions. For instance, if we only change the stride embedding of each layer in MobileNet-V2, we could observe that our agent will allocate more activation bits to down-sample layers (stride=2), and this phenomenon is more obviously at deep layers. As for weight bits, our agent also allocate more bits for down-sample layers. Moreover, if we only change the depth-wise/point-wise embedding, we could observe that point-wise layers will be allocated fewer weight bits, the reason may lay on point-wise layer are more computation bounded, fewer weight bits will obviously reduce the computation complexity.

## 6 Conclusion

In this paper, we propose an automated framework for quantization, **Hardware-Aware Automated Quantization (HAQ)**, which does not require any domain experts and rule-based heuristics. We provide a learning based method that can search the quantization policy with hardware feedback. Compared with indirect proxy signals, our framework can offer a specialized quantization solution for different hardware platforms. Extensive experiments demonstrate that our framework performs better than conventional rule-based approaches for multiple objectives: latency, energy and model size. Our framework reveals that the optimal policies on different hardware architectures are drastically different, and we interpreted the implication of those policies. We believe the insights will inspire the future software and hardware co-design for efficient deployment of deep neural networks.

## Acknowledgement

We thank NSF Career Award #1943349, MIT-IBM Watson AI Lab, Samsung, SONY, Xilinx, TI and AWS for supporting this research.

## References

- Apple (2018) Apple describes 7nm A12 bionic chips. URL <http://www.eenewsanalog.com/news/apple-describes-7nm-a12-bionic-chip/page/0/1>
- Cai H, Yang J, Zhang W, Han S, Yu Y (2018) Path-Level Network Transformation for Efficient Architecture Search. In: ICML 3
- Cai H, Zhu L, Han S (2019) ProxylessNAS: Direct Neural Architecture Search on Target Task and Hardware. In: ICLR 2, 10
- Choi J, Wang Z, Venkataramani S, Chuang PIJ, Srinivasan V, Gopalakrishnan K (2018) PACT: Parameterized Clipping Activation for Quantized Neural Networks. arXiv 2, 6, 9
- Chollet F (2017) Xception - Deep Learning with Depthwise Separable Convolutions. In: CVPR 5, 6
- Courbariaux M, Hubara I, Soudry D, El-Yaniv R, Bengio Y (2016) Binarized Neural Networks: Training Deep Neural Networks with Weights and Activations Constrained to +1 or -1. arXiv 3
- Deng J, Dong W, Socher R, Li LJ, Li K, Li FF (2009) ImageNet - A large-scale hierarchical image database. In: CVPR 6
- Han S (2017) Efficient Methods and Hardware for Deep Learning. PhD thesis 10
- Han S, Mao H, Dally W (2016) Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding. In: ICLR 2, 3, 9, 10
- He K, Zhang X, Ren S, Sun J (2016) Deep Residual Learning for Image Recognition. In: CVPR 2
- He Y, Zhang X, Sun J (2017) Channel pruning for accelerating very deep neural networks. In: ICCV 3
- He Y, Lin J, Liu Z, Wang H, Li LJ, Han S (2018) AMC: AutoML for Model Compression and Acceleration on Mobile Devices. In: ECCV 3, 10
- Howard AG, Zhu M, Chen B, Kalenichenko D, Wang W, Weyand T, Andreetto M, Adam H (2017) MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. arXiv 2, 6
- Imagination (2018) Powervr neural network accelerator. URL <https://www.imgtec.com/vision-ai/powervr-series2nx/powervr-ax2145-nna/> 2
- Jacob B, Kligys S, Chen B, Zhu M, Tang M, Howard AG, Adam H, Kalenichenko D (2018) Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference. In: CVPR 2, 3
- Kingma D, Ba J (2015) Adam - A Method for Stochastic Optimization. In: ICLR 6
- Krishnamoorthi R (2018) Quantizing deep convolutional networks for efficient inference - A whitepaper. arXiv 3
- Lillicrap T, Hunt JJ, Pritzel A, Heess N, Erez T, Tassa Y, Silver D, Wierstra D (2016) Continuous control with deep reinforcement learning. In: ICLR 3, 5
- Liu C, Zoph B, Neumann M, Shlens J, Hua W, Li LJ, Fei-Fei L, Yuille A, Huang J, Murphy K (2018) Progressive Neural Architecture Search. In: ECCV 3
- Liu Z, Li J, Shen Z, Huang G, Yan S, Zhang C (2017) Learning efficient convolutional networks through network slimming. In: ICCV 3, 7
- Nvidia (2018) Nvidia tensor cores. URL <https://www.nvidia.com/en-us/data-center/tensorcore/> 2
- Pham H, Guan MY, Zoph B, Le QV, Dean J (2018) Efficient Neural Architecture Search via Parameter Sharing. In: ICML 3
- Rastegari M, Ordonez V, Redmon J, Farhadi A (2016) XNOR-Net - ImageNet Classification Using Binary Convolutional Neural Networks. In: ECCV 3
- Sandler M, Howard A, Zhu M, Zhmoginov A, Chen LC (2018) MobileNetV2: Inverted Residuals and Linear Bottlenecks. In: CVPR 2, 6, 7
- Sharma H, Park J, Suda N, Lai L, Chau B, Chandra V, Esmaeilzadeh H (2018) Bit fusion: Bit-level dynamically composable architecture for accelerating deep neural network. In: ISCA 2, 6, 9
- Umuroglu Y, Rasnayake L, Sjalander M (2018) Bismo: A scalable bit-serial matrix multiplication overlay for reconfigurable computing. In: FPL 2, 6
- Williams S, Waterman A, Patterson D (2009) Roofline: an insightful visual performance model for multicore architectures. Communications of the ACM 52(4):65-76 7
- Xilinx (2018a) Ultrascale architecture and product data sheet: Overview. URL [https://www.xilinx.com/support/documentation/data\\_sheets/ds890-ultrascale-overview.pdf](https://www.xilinx.com/support/documentation/data_sheets/ds890-ultrascale-overview.pdf) 6
- Xilinx (2018b) Zynq-7000 soc data sheet: Overview. URL [https://www.xilinx.com/support/documentation/data\\_sheets/ds190-Zynq-7000-Overview.pdf](https://www.xilinx.com/support/documentation/data_sheets/ds190-Zynq-7000-Overview.pdf) 6

- Yang TJ, Chen YH, Sze V (2016) Designing energy-efficient convolutional neural networks using energy-aware pruning. arXiv [3](#)
- Yang TJ, Howard A, Chen B, Zhang X, Go A, Sandler M, Sze V, Adam H (2018) Netadapt: Platform-aware neural network adaptation for mobile applications. In: ECCV [3](#)
- Zhou A, Yao A, Wang K, Chen Y (2018) Explicit loss-error-aware quantization for low-bit deep neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 9426–9435 [3](#)
- Zhou S, Ni Z, Zhou X, Wen H, Wu Y, Zou Y (2016) DoReFa-Net - Training Low Bitwidth Convolutional Neural Networks with Low Bitwidth Gradients. arXiv [3](#)
- Zhu C, Han S, Mao H, Dally W (2017) Trained Ternary Quantization. In: ICLR [2](#), [3](#)
- Zoph B, Le QV (2017) Neural Architecture Search with Reinforcement Learning. In: ICLR [3](#)

Author accepted manuscript