

# Fast Quantum Convolutional Neural Networks for Low-Complexity Object Detection in Autonomous Driving Applications

Hankyul Baek, Donghyeon Kim, and Joongheon Kim, *Senior Member, IEEE*

**Abstract**—Spurred by consistent advances and innovation in deep learning, object detection applications have become prevalent, particularly in autonomous driving that leverages various visual data. As convolutional neural networks (CNNs) are being optimized, the performances and computation speeds of object detection in autonomous driving have been significantly improved. However, due to the exponentially rapid growth in the complexity and scale of data used in object detection, there are limitations in terms of computation speeds while conducting object detection solely with classical computing. Motivated by this, quantum convolution-based object detection (QCOD) is proposed to adopt quantum computing to perform object detection at high speed. The QCOD utilizes our proposed fast quantum convolution that uploads input channel information and re-constructs output channels for achieving reduced computational complexity and thus improving performances. Lastly, the extensive experiments with KITTI autonomous driving object detection dataset verify that the proposed fast quantum convolution and QCOD are successfully operated in real object detection applications.

**Index Terms**—Quantum Machine Learning, Quantum Convolutional Neural Network, Object Detection, Autonomous Driving

## I. INTRODUCTION

With the consistent advancement of deep learning, many deep learning-based applications have improved performance and become practical. These deep learning-based applications require significant computational power due to the expected increase in dataset sizes and algorithm complexity [1]. In particular, the computational complexity becomes more significant in object detection due to the growing complexity of data, which expands from 2D images to 3D point clouds and multi-modal data [2]. To cope with the growing complexity, several research aims to enhance the model architectures and the algorithmic advantages for improving the computation speed and performance of applications [3], [4]. In the era of classical computing, these algorithmic improvements yield highly positive results [5], [6]. However, they encounter fundamental challenges in efficiently conducting highly complex and complicated applications due to the inherent limitations of classical computing resources [7]. A classical convolutional neural network (CNN) is one of the

This research was funded by Institute of Advanced Technology Development (IATD) in Hyundai Motor Company and Kia Corporation. (*Corresponding authors: Donghyeon Kim, Joongheon Kim*)

H. Baek and J. Kim are with the Department of Electrical and Computer Engineering, Korea University, Seoul 02841, Republic of Korea (e-mails: {67back, joongheon}@korea.ac.kr).

D. Kim is with the Institute of Advanced Technology Development (IATD), Hyundai Motor Company (e-mail: donghyeon.kim@hyundai.com)

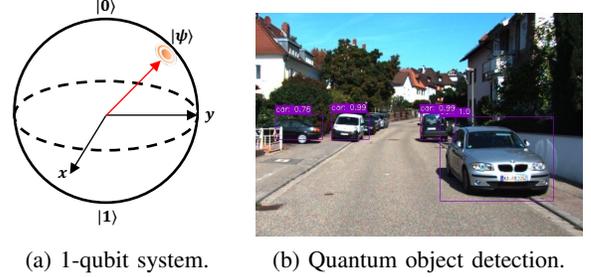


Fig. 1. A brief illustration of quantum computing and its application.

representatives that shows these limitations [8]. While CNN-based algorithms demonstrate rapid execution and viable performance, the computational complexity of CNN is significantly contingent upon the input size, with a computational cost of  $\mathcal{O}(X \cdot C_{in} \cdot C_{out})$ , where  $X$ ,  $C_{in}$ , and  $C_{out}$  denote the product of input data size and kernel size, input channel, and output channel, respectively [9]. This can pose a substantial challenge when applying the convolution process to extensive and intricate datasets, impeding its scalability in dealing with the rapidly expanding dataset. These challenges are exacerbated in complex applications based on CNNs, such as object detection. As the model's architecture grows in complexity and the dataset employed becomes more intricate, it is evident that relying solely on classical computing for such applications for real-time execution is not practical due to computational limitations [10].

Quantum computing is regarded as a promising solution to resolve these computational limitations. The emergence of the noisy intermediate-scale quantum (NISQ) era suggests that the number of available quantum bits, *i.e.*, qubits, will exceed thousands by 2025, potentially achieving quantum advantage [11]. This achievement is based on the intrinsic nature of quantum computing, superposition, which enables quantum computing to excel in classical computing in various complex tasks. Fig. 1 (a) represents an example of a 1-qubit superposition. In contrast to classical bit, the state of a qubit can be depicted as  $|\Phi\rangle = \alpha_0 |0\rangle + \alpha_1 |1\rangle$ , where  $\alpha_0$  and  $\alpha_1$  are the probabilistic complex amplitudes of qubit, satisfying  $\alpha_0^2 + \alpha_1^2 = 1$ . This representational capacity becomes increasingly more extensive as the number of qubits grows, because the number of bases, *e.g.*,  $|00\rangle$  and  $|01\rangle$ , expands to  $2^q$ , where  $q$  denotes the number of qubits. On the other hand, there are challenges in implementing quantum computing-based applications in the NISQ era. As quantum computing and its algorithms are still in their early stages, there is

a lack of optimization methods for tasks using quantum computing [12]. Additionally, quantum computing fails to replace classical computing entirely because it cannot perform structured tasks like convolution in CNNs [13]. Moreover, the limited availability of deep learning techniques, datasets, and optimization tools for classical computing, along with the complexities of quantum computing, make the advancement of deep learning applications through quantum advantage in the realm of quantum machine learning challenging.

Inspired by this, this paper focuses on the object detection, one of the most complicated applications using CNN. To cope with the growing complexity of the object detection and achieve faster operation time, this paper proposes quantum convolution-based object detection (QCOD). Fig. 1 (b) shows a brief example result of QCOD. With our proposed quantum convolution, named fast quantum convolution, which boosts the encoding process. The fast quantum convolution optimizes the advantages of quantum computing. As our fast quantum convolution encodes multi-channel data into an identical quantum system, QCOD achieves fast quantum speed-ups. In addition, to leverage classical optimization schemes and deep learning techniques for QCOD, this paper proposes heterogeneous knowledge distillation, a modified version of knowledge distillation, to train the region proposal layer of QCOD. Knowledge distillation is a well-known training method that transfers the knowledge of the pre-trained teacher model to the un-trained student models. In this paper, heterogeneous knowledge distillation selects the pre-trained classical region proposal network and quantum convolution region proposal network as teacher and student model, respectively. Via heterogeneous knowledge distillation, QCOD addresses the lack of quantum optimization schemes in object detection.

Furthermore, this paper verifies the superiority of fast quantum convolution in QCOD and substantiates the probability of achieving quantum object detection in the near future through extensive experiments and ablation studies. It is difficult to definitively state that quantum object detection is superior to classical object detection in the view of performance. However, this paper observes that quantum object detection can be realized, and the proposed QCOD shows significant speed-ups in object detection. Serving as a foundational step toward the implementation of quantum object detection, this paper provides an outlook for future research in quantum object detection.

**Contributions.** The major contributions of QCOD are as follows. First of all, This paper designs a novel quantum convolution named fast quantum convolution, considering the qubits' representation ability. The fast quantum convolution encodes multiple channels into quantum states and achieves quantum speed-ups. Second, this paper proposes heterogeneous knowledge distillation to leverage classical optimization schemes and the knowledge from classical pre-trained models, addressing the lack of knowledge in the quantum domain. Third, this paper verifies the superiority of the fast quantum convolution when utilizing with quantum random access memory (QRAM). Finally, this paper conducts numerous experiments and, to the best of our knowledge, implements the first quantum object detection.

## II. RELATED WORK

This section introduces previous research that is closely aligned with our fast quantum convolution and QCOD development. The pivotal topics include i) quantum machine learning implementation, ii) data re-uploading, and iii) knowledge distillation for subset model training.

**Quantum machine learning implementation.** The implementation of quantum machine learning hinges on two fundamental categories: i) optimizing qubits' representation abilities and ii) employing fast QRAM searching algorithms to minimize complexity. These properties enable quantum computing to outperform classical counterparts. Among these categories, the research related to our considering quantum convolutional neural network (QCNN) is as follows. *Baek et al.* [14] address the scalability limitation of available qubits in quantum convolution filters by incorporating these filters into massive 3D data classification applications. In this project, they leverage the concept of fidelity to achieve robust performance. *Shen et al.* [15] focus on the architecture of classical CNN, replacing the fourier transform process of the CNN with a quantum circuit, thereby enhancing the speed of the entire CNN. Furthermore, several studies aim to realize quantum advantage by combining QCNN and QRAM. *Oh et al.* [16] implement QCNN on QRAM to store large-sized data. In addition, *Kerenidis et al.* [17] prove the quantum advantage when utilizing QCNN and QRAM with small errors.

**Data re-uploading.** Data re-uploading is an encoding technique based on the quantum information theory that the states of qubits can represent multiple information [18]. *Pérez-Salinas et al.* [19] firstly propose and prove the feasibility of data re-uploading using the 1-qubit system of quantum machine learning. *Friedrich et al.* [20] combine data re-uploading techniques with QCNN to encode multiple data within a few qubits for avoiding barren plateaus<sup>1</sup>. *Schuld et al.* [23] confirm that data re-uploading allows quantum models to represent progressively richer frequency spectra while using a limited number of qubits. In this paper, we propose channel uploading, a modified version of data re-uploading, to cope with numerous number of channels of practical object detection applications.

**Knowledge distillation for subset model training.** Knowledge distillation is a training method to handle variations in deep learning resources and enhance robust training in real-world applications [24]. Knowledge distillation is typically incorporated as a regularizer in the loss function, aiming to minimize the difference between the logits of the teacher model and those of the target student model. The target student model can conduct robust training by transferring pre-trained knowledge from the teacher to the student model. *Cui et al.* [25] adopts a knowledge distillation regularizer as a loss function for semi-supervised learning, aiming to process real-world images. In this paper, we take a step further by employing knowledge distillation between models in a heterogeneous domain. We

<sup>1</sup>The phenomenon of barren plateaus, a characteristic of quantum machine learning, impedes the trainability of quantum machine learning models [21]. Similar to the local minima in classical machine learning, barren plateaus give rise to problems where parameters are not efficiently optimized. In addition, it is well-known that the increase of the number of qubits induces barren plateaus [22].

set a classical CNN-based model as the teacher model and our fast quantum convolution-based model as the student model.

### III. QUANTUM MACHINE LEARNING

Quantum machine learning is a machine learning framework that leverages the quantum advantages of quantum computing to address challenges previously tackled by classical neural networks. The quantum machine learning comprises encoding, parameterized quantum circuits, and decoding. In this section, we dive into each process within quantum machine learning essentials for constructing our proposed fast quantum convolution and QCOD.

**Basic quantum operations.** In contrast to classical bits, which have deterministic values of either 0 or 1, quantum computing allows for the superposition of two states simultaneously. This unique characteristic is expressed using Dirac notation as  $|\Phi\rangle \triangleq \sum_{k=1}^{2^q} \alpha_k |k\rangle$ , where  $|k\rangle$  represents a basis in the Hilbert space, and  $\forall q \in \mathbb{N}[1, \infty)$  and  $\sum_{k=1}^{2^q} |\alpha_k|^2 = 1$ . The initialized  $q$ -qubit system can be expressed as  $|0\rangle^{\otimes q}$ . Similar to classical computer logic gates, quantum gates are operators capable of manipulating the state of qubits. From a physics perspective, operations of the quantum gates can be interpreted as transitions of qubit states on the Bloch sphere from one point to another [26]. Note that these quantum gates are in the form of unitary matrices. In this paper, we denote the unitary matrices used for encoding as  $U_E$  and the unitary matrices used for training as  $U_T$ , depending on their purposes.

**Encoding.** It is essential to convert classical information into quantum information to facilitate the integration of quantum machine learning with classical computing. This transformation is achieved through the implementation of the quantum gates denoted as  $U_E$ . Mathematically,  $U_E$  can be expressed as a  $2^q \times 2^q$  unitary matrix within a  $q$ -qubit quantum system. Therefore, with the classical data  $\mathbf{x}$ , the encoded  $q$ -qubit quantum states can be represented as  $|\psi_{\mathbf{x}}\rangle = U_E(\mathbf{x})|0\rangle^{\otimes q}$ , where the encoded quantum state  $|\psi_{\mathbf{x}}\rangle$  is on  $2^q$ -dimension Hilbert space. In this work, we design quantum gates  $U_E$  without any trainable parameters to encode classical information into quantum states consistently.

**Parameterized quantum circuit.** After encoding the classical information on the target quantum system, the parameterized quantum circuit (PQC) trains the parameter as in attention [27] in classical machine learning. Each PQC has trainable unitary matrices  $U_T$ <sup>2</sup>, which comprises trainable rotation gates  $R_{\Gamma}(\theta)$  and trainable *Controlled- $\Gamma$*  gates  $CT$ , where  $\forall \Gamma \in \{X, Y, Z\}$ . Here,  $\theta$  denotes the trainable parameters where  $\forall \theta \in [0, 2\pi]^{|\theta|}$ . Revisiting the encoded quantum states  $|\psi_{\mathbf{x}}\rangle$ , the output of the PQC can be expressed as  $|\psi_{\mathbf{x},\theta}\rangle = U_T(\theta)|\psi_{\mathbf{x}}\rangle$ . Note that the trainable unitary matrices  $U_T$  use trainable parameters and encoded quantum states as inputs.

**Decoding.** To use the transformed quantum state  $|\psi_{\mathbf{x},\theta}\rangle$  with classical computing, we consider using the expectation quantum value [28]. These expectation value can be designed as  $\langle O_{\mathbf{x},\theta} \rangle = \prod_{M \in \mathcal{M}} \langle \psi_{\mathbf{x},\theta} | M | \psi_{\mathbf{x},\theta} \rangle$ , where  $\langle O_{\mathbf{x},\theta} \rangle$  denotes the

<sup>2</sup>The general expression is a parameterized or variational unitary matrices. In this paper, we emphasize the significance of a PQC, which comprises trainable unitary matrices for a straightforward explanation.

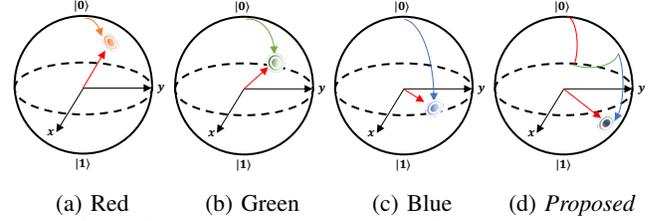


Fig. 2. Comparison between existing channel uploading strategies (a-c) and our proposed channel uploading strategy (d).

expectation of quantum measured values on Hermitian matrices  $M$ . To decode the quantum information of each qubit, this paper designs  $\mathcal{M} = \{M_l\}_{l=1}^q$ , where  $M_l = I^{\otimes l-1} \otimes Z \otimes I^{L-l}$ . Here,  $I$  denotes an identity matrix and  $Z$  denotes a Pauli-Z matrix  $\begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$  [29]. As a result, the output expectation values exist in  $\langle O_{\mathbf{x},\theta} \rangle \in [-1, 1]^q$ .

### IV. FAST QUANTUM CONVOLUTIONAL NEURAL NETWORKS

This section presents the details of our proposed fast quantum convolution, which can mitigate the computational overheads via patch processing, channel uploading, and channel reconstruction.

#### A. Motivation

This paper aims to overcome the structural limitations of CNNs and existing QCNNs to cope with the growing complexity of real-world datasets in the field of object detection. To solve the limitations, we design our fast quantum convolution depicted in Fig. 3. The fast quantum convolution employs i) patch processing, ii) channel uploading, iii) quantum feature extraction, and iv) channel reconstruction layer. Particularly, by employing channel uploading, we successfully mitigate the computational complexity. Fig. 2 briefly illustrates our proposed channel uploading strategy. In existing quantum computing systems, the approach of uploading data from the same channel is employed, as depicted in Fig. 2 (a-c). However, in real-world applications involving computations using numerous channels, a limitation arises where the operations must be repeated, corresponding to the number of channels ( $C$ ). In contrast to the classical encoding scheme, our proposed channel uploading strategy focuses on the qubit's representation ability that is able to contain multiple pieces of information, depicted in Fig. 2 (d). Accordingly, our proposed fast quantum convolution can reduce entire computational complexity.

#### B. Architecture of fast quantum convolutional neural network

This paper employs three steps to design our fast quantum convolution: i) patch pre-processing, ii) quantum state encoding, and iii) quantum state decoding. The first step enables an optimized quantum convolutional process, while the other processes are necessary to attain quantum-induced image features. The overall process of fast quantum convolution is detailed in Algorithm 1.

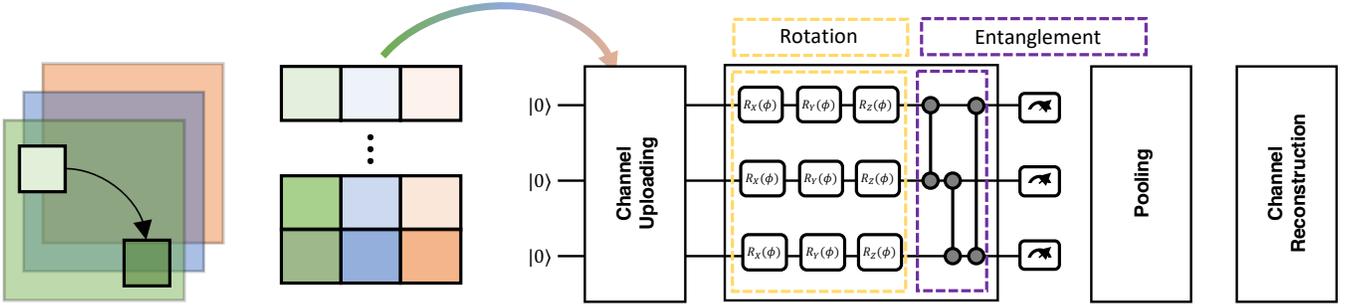
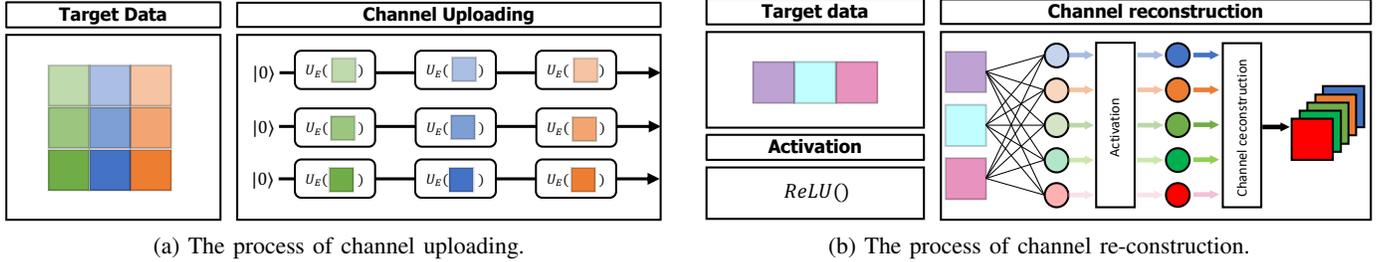


Fig. 3. An illustration of the proposed fast quantum convolutional neural network.



(a) The process of channel uploading.

(b) The process of channel re-construction.

Fig. 4. Detailed illustration of the proposed fast quantum convolution.

**Patch processing.** Inspired by classical image pre-processing optimization [30], this paper employs patch processing operation, named *im2col*, which is widely implemented in classical CNN to fast quantum convolution. Fig. 3 illustrates the patch processing in our fast quantum convolution. In contrast to classical patch processing operations, our methods can reduce the number of operations by uploading the data in different channels in each qubit. A classical 3D tensor image  $X^l \in \mathbb{R}^{H \times W \times C^l}$  is transformed to 2D tensor matrix  $P^l \in \mathbb{R}^{(H^{l+1}W^{l+1}) \times (HWC^l)}$ . Here, each patch in the input image is flattened to the part of each row of the 2D tensor matrix  $P^l$ . Therefore, the number of rows equals the one of operations.

**Encoding via quantum channel uploading.** Based on the ability of qubits to maintain multiple pieces of information simultaneously, in contrast to classical bits, we sequentially upload the channels onto the same qubits. Fig. 4 (a) provides a clear representation of quantum channel uploading. Each row of the 2D tensor matrix  $P^l$  is uploaded on our quantum circuits. We use three different rotation gates  $R_x = \begin{bmatrix} \cos(\frac{\alpha}{2}) & -i \sin(\frac{\alpha}{2}) \\ -i \sin(\frac{\alpha}{2}) & \cos(\frac{\alpha}{2}) \end{bmatrix}$ ,  $R_y = \begin{bmatrix} \cos(\frac{\beta}{2}) & -\sin(\frac{\beta}{2}) \\ \sin(\frac{\beta}{2}) & \cos(\frac{\beta}{2}) \end{bmatrix}$  and  $R_z = \begin{bmatrix} e^{-i\frac{\delta}{2}} & 0 \\ 0 & e^{i\frac{\delta}{2}} \end{bmatrix}$ .  $\alpha$ ,  $\beta$  and  $\delta$  are constant values that can be modified according to the number of uploaded channels. Our proposed quantum channel uploading encoding strategy can be described as

$$|\psi\rangle_i = \prod_{j=0}^{HWC^l} U_E(p_{i,j}^l) |0\rangle^{\otimes q}, \quad (1)$$

where  $p_{i,j}^l$  denotes the components at the  $i$ -th row and  $j$ -th column of the 2D tensor matrix  $P^l$ . Note that  $\forall i \in \mathbb{N}[0, H^{l+1}W^{l+1})$ . Because the number of input components  $HWC^l$  in each row is larger than the number of available qubits

$q$  in the recent NISQ era, where the number of qubits is small, we design the encoding layer  $U_E$  can employ additional data uploading strategy. Note that the quantum channel uploading process occurs in encoding, and the set of encoding layer  $U_E$  doesn't have trainable parameters.

**Quantum convolution.** The fast quantum convolution utilizes the PQC to perform convolution on the information of the encoded quantum feature  $|\psi\rangle_i$ . This use of PQC can be compared with the convolution filters of classical CNNs. Unlike classical CNNs that employ element-wise products, PQC performs convolution through a trainable layer  $U_T$  consisting of trainable *Controlled- $\Gamma$*  gates and rotation gates  $R_\Gamma(\theta)$ , where  $\forall \Gamma \in \{X, Y, Z\}$ . Particularly, by using trainable *Controlled- $\Gamma$* , PQC is designed to induce mutual information referencing between qubits, creating entanglement [31]. This design allows the PQC to incorporate the spatial information of input data into the convolution more effectively. The quantum convolution can be represented as

$$f(|\psi\rangle_i; \theta) : |\psi_\theta\rangle_i \leftarrow U_T(\theta) |\psi\rangle_i, \quad (2)$$

where  $|\psi_\theta\rangle_i$  is the output convoluted quantum states with trainable parameters  $\theta$ .

**Quantum feature extraction via decoding.** Compared to classical CNNs, where convoluted features have discrete values, quantum features possess a probabilistic nature. This paper considers the quantum expectation value  $\langle O_{x,\theta} \rangle$ , represented in Sec. III as quantum features. To ensure stable learning in quantum computing for machine learning, we evaluate and utilize the probabilistic values associated with each basis' amplitude. In addition, as demonstrated in the example in Fig. 1 (a), a quantum state can be represented as  $|\psi\rangle_k = \alpha |0\rangle + \beta |1\rangle$ , where  $k \in \mathbb{N}[1, q]$ , and it satisfies  $\alpha^2 + \beta^2 = 1$ . In this context, we set the output of quantum convolution as the probabilistic

difference in amplitudes for each basis, i.e.,  $\alpha^2 - \beta^2$ . This approach allows us to represent the output of each qubit as  $\langle O_{\mathbf{x},\theta} \rangle_k \in \mathbb{R}[-1, 1]$ . We implement channel reconstruction layer to the entire output  $\langle O_{\mathbf{x},\theta} \rangle \in \mathbb{R}[-1, 1]^{\otimes q}$ . Fig. 4 (b) illustrates our channel reconstruction layers. By employing a linear function and an activation layer on the output, our fast quantum convolution succeeds in achieving scalability.

### C. Strategy of fast quantum convolution

**Quantum backpropagation.** Based on the quantum machine learning theory in [32], the gradients of quantum gates cannot be calculated directly. This is due to the intrinsic nature of quantum computing, where the output corresponds to the expectation of the corresponding computation outcome. Similarly, our fast quantum convolution outputs are also expressed as expectations, making the derivative of such expected outputs an invalid operation within the framework of quantum expectations. To solve this, this paper considers parameter-shift rule [33].

$$\frac{\partial \langle O_{\mathbf{x},\theta} \rangle}{\partial \theta} = \frac{1}{2} [\langle O_{\mathbf{x},\theta+\frac{\pi}{2}} \rangle - \langle O_{\mathbf{x},\theta-\frac{\pi}{2}} \rangle], \quad (3)$$

where  $\langle O_{\mathbf{x},\theta+\frac{\pi}{2}} \rangle$  and  $\langle O_{\mathbf{x},\theta-\frac{\pi}{2}} \rangle$  denote the extracted outputs with modified parameters  $\theta + \frac{\pi}{2}$  and  $\theta - \frac{\pi}{2}$ , respectively.

**QRAM for quantum speed-ups.** As a counterpart to the random-access memory (RAM) in classical computing, quantum random-access memory (QRAM) stores the address of information in the state of the qubit. QRAM technology is a significant part to achieve quantum advantages and has drawn attention. One of the general QRAM structures is the bucket-brigade architecture, which inputs the address of the information in the quantum state and retrieves the data in the quantum state as the output [34]. The QRAM process of our considering fast quantum convolution can be described as

$$\sum_i |i\rangle_{\text{address}} |0\rangle_{\text{data}} \xrightarrow{\text{QRAM}} \sum_i |i\rangle_{\text{address}} |\psi_i\rangle_{\text{data}}, \quad (4)$$

where  $|\cdot\rangle_{\text{address}}$  and  $|\cdot\rangle_{\text{data}}$  denote the storage address of QRAM and corresponding data, respectively. Based on the following Lemma 1 [15], and Lemma 2 [17], we observe quantum speed-ups as depicted in Theorem 1.

**Lemma 1. (Advantages of QRAM)** *Let target input  $P \in \mathbb{R}^{n \times d}$ , there exists a QRAM structure that conducts inserting, deleting, and updating each datum  $p_{i,j}$  in time  $\mathcal{O}(\log(n^2))$ . In addition, there exists a quantum algorithm  $|i\rangle_{\text{address}} |0\rangle_{\text{data}} \rightarrow |i\rangle_{\text{address}} |\psi_i\rangle_{\text{data}}$  in time  $\mathcal{O}(\log^2 n)$ .*

**Lemma 2. (Running time of quantum gates)** *Let symmetric matrix  $M \in \mathbb{R}^{d \times d}$ , datum  $x \in \mathbb{R}^d$  and error  $\delta > 0$ . When the matrix is stored in appropriate QRAM, there exists an algorithm that satisfies  $\| |z\rangle - |Mx\rangle \|_2 \leq \delta$  in time  $\mathcal{O}((\sqrt{d}\kappa(M) + T_x \kappa(M)) \log(1/\delta))$  with probability at least  $1 - \frac{1}{\text{poly}(d)}$ , where  $\kappa(M)$  and  $T_x$  denote the condition number of  $M$  and setting time for  $|x\rangle$ .*

**Theorem 1. (Running time of the fast quantum convolution)** *Let the input classical 3D tensor image  $X^l \in \mathbb{R}^{H^{l+1} \times W^{l+1} \times C^l}$  in*

---

### Algorithm 1: Fast quantum convolution procedure.

---

- 1 **Notation.** Input number of qubits:  $q$ , 3D tensor:  $X^l$ , transformed 2D tensor matrix:  $P^l$ , components at the  $i$ -th row,  $j$ -th column of the  $P^l$ :  $p_{i,j}^l$ , Channel Reconstruction function  $T$  ;
  - 2 **Input:** Input classical image  $X^l$ ;
  - 3 **Patch processing.**  $P^l \leftarrow X^l$ ;
  - 4 **for**  $i \in \{1, 2, \dots, H^{l+1}W^{l+1}\}$  **do**
  - 5     Initialize quantum state  $|0\rangle^{\otimes q}$ ;
  - 6     **for**  $j \in \{1, 2, \dots, HWC^l\}$  **do**
  - 7          $|\psi\rangle_i \leftarrow U_E(p_{i,j}^l)|0\rangle^{\otimes q}$ ;
  - 8          $|\psi_\theta\rangle_i \leftarrow U_T(\theta)|\psi\rangle_i$ ;
  - 9         **for**  $k \in \{1, 2, \dots, q\}$  **do**
  - 10             Achieve  $\langle O \rangle_k \in \mathbb{R}[-1, 1]$ ;
  - 11     Reshaping & Pooling;
  - 12     Channel Reconstruction  $T : \mathbb{R}^q \rightarrow \mathbb{R}^{C^{l+1}}$ ;
  - 13 Achieve extracted features  $X^{l+1}$ ;
  - 14 **Output:** Extracted features
- 

$l$ -th fast quantum convolution layer and the number of qubits  $q = H^{l+1} \times W^{l+1}$ , where  $\forall l \in L$  and  $H^l \times W^l \geq 2$  and let the running time of encoding gates  $U_E$  and  $U_T$  as  $t_E$  and  $t_T$ , respectively. With condition numbers of  $l$ -th encoding  $\kappa(M_E)$  and PQC layer  $\kappa(M_T)$  that satisfy  $\kappa(M_T) = \rho^l \cdot \kappa(M_E)$ , the running time of entire fast quantum convolution process  $T_{\text{total}}$  is conducted in the time of

$$\mathcal{O}(\log^2(H^{l+1}W^{l+1}) + \log(1/\delta)\kappa(M_E)(C^l t_E + \rho^l t_T)). \quad (5)$$

*Proof.* Based on the patch processing, the 2D patch matrix can be generated as  $P^l \in \mathbb{R}^{(H^{l+1}W^{l+1}) \times (HWC^l)}$ . With (4) and Lemma 1, the total setting is conducted in time  $T_S = \mathcal{O}(\log(H^{l+1}W^{l+1})^2 + \log^2(H^{l+1}W^{l+1}))$ . Here, with the assumption  $H^{l+1} \times W^{l+1} \geq 2$ , then the time complexity of  $T_S$  can be expressed as,

$$T_S \approx \mathcal{O}(\log^2(H^{l+1}W^{l+1})). \quad (6)$$

To conduct quantum convolution on the data in QRAM, it is necessary to design  $M_E \in \mathbb{R}^{2^q \times 2^q}$  for encoding and  $M_T \in \mathbb{R}^{2^q \times 2^q}$  for PQC. Based on Lemma 2, the encoding process complexity and convolution process complexity can be described as

$$T_E = \mathcal{O}((\sqrt{2^q}\kappa(M_E) + C^l \cdot t_E \cdot \kappa(M_E)) \log(1/\delta_E)). \quad (7)$$

Similarly, the convolution process of fast quantum convolution can be described as

$$T_T = \mathcal{O}((\sqrt{2^q}\kappa(M_T) + t_T \cdot \kappa(M_T)) \log(1/\delta_T)). \quad (8)$$

Here, we observe that the trainable gate  $t_T$  is called only once due to the advantages of channel uploading. For simplicity, by setting  $\delta_E \approx \delta_T$  and  $C^l t_E + \rho^l t_T \gg \sqrt{2^q}(1 + \rho^l)$ . The total time  $T_{\text{total}} = T_S + T_E + T_T$  is conducted in time

$$\mathcal{O}(\log^2(H^{l+1}W^{l+1}) + \log(1/\delta)\kappa(M_E)(C^l t_E + \rho^l t_T)). \quad (9)$$

On the other hand, with existing quantum convolution method, (7)-(8) is modified as  $T'_E = \mathcal{O}(C^l(\sqrt{2^q}\kappa(M_E) + t_E \cdot \kappa(M_E)) \log(1/\delta_E))$  and  $T'_T = \mathcal{O}(C^l(\sqrt{2^q}\kappa(M_T) + t_T \cdot$

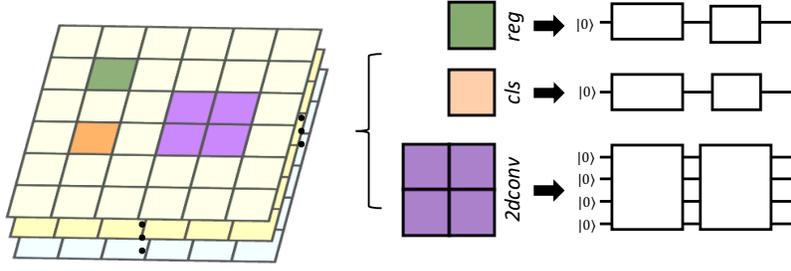


Fig. 5. Proposed QRPN with our fast quantum convolution (The green and orange quantum convolution filter is designed for classification and box regression. The purple quantum convolution filter is designed for 2-dimensional convolution).

$\kappa(M_T) \log(1/\delta_T)$ , respectively. With a large number of channels  $C^l$ , we observe the advantage of proposed fast quantum convolution.  $\square$

## V. QUANTUM OBJECT DETECTION

This section provides a detailed description of our method for implementing our fast quantum convolution in object detection. Note that the quantum version of the region proposal network (RPN) proposed below is a significant component of our QCOD.

### A. Motivation of quantum region proposal network

The RPN is a major network utilized in object detection applications [35]. It is due to the role of RPN that localizes and proposes the target object. The RPN employs convolutional layers composed of spatial filters with dimensions  $n \times n$ , where  $n \geq 1$ . In addition, convolutional layers containing  $1 \times 1$  spatial filters for box regression and classification are employed, respectively. Using these convolutional filters, the RPN takes an extracted feature as input, which is obtained from an extractor, and generates a set of rectangular object proposals, each of them accompanied by an objectness score. The outputs, *i.e.*, object proposals and objectness scores, serve as the foundation for enabling the classifier to categorize objects within the target proposal. However, Despite excellent performance, RPN is still a computationally expensive network owing to the numerous number of proposed regions and channels involved in [36].

### B. Architecture of quantum region proposal network

To solve the limitations, this paper modifies the RPN [37] to a quantum version of the RPN (QRPN) using our fast quantum convolution. QRPN aims to calculate object proposals and objectness scores using our fast quantum convolution. In contrast to RPN, which slides spatial filters on target tensor and utilizes element-wise multiplication, QRPN encodes multiple channel inputs jointly convolute the features via unitary gates (*e.g.*,  $R_x$ ,  $R_y$  and  $CNOT$  gates). By utilizing our fast quantum convolution, QRPN mitigates the time complexity as proved in Theorem 1. In classical RPN structure,  $1 \times 1$  convolution filters are utilized for computing two different loss functions (*i.e.*, classification loss and box regression loss). Here, due to the  $1 \times 1$  scale filter operation strategy, they can be considered as scalar-multiplied fully connected layers. To implement these

characteristics of  $1 \times 1$  filter-based convolution using fast quantum convolution, we design another structure of our fast quantum convolution. With an initialized single qubit  $|0\rangle$  and each  $1 \times 1$  quantum convolution can be expressed as

$$|\psi_\theta\rangle_i = \prod_{j=0}^{HW C^l} U_{1 \times 1}(\theta) \cdot U_E(p_{i,j}^l) |0\rangle^{\otimes 1}, \quad (10)$$

where  $\forall i \in \mathbb{N}[0, H^l W^l]$ , and  $U_{1 \times 1} \subset U_T$  denotes unitary matrices which is activated on each qubit. Note that the difference between  $U_T$  and  $U_{1 \times 1}$  comes from the usage of 2-qubit gates, which can induce the entanglement in the designed quantum circuit.

### C. Heterogeneous knowledge distillation training

Training the QRPN is challenging due to the limited availability of quantum computing resources and optimization tools, especially when considering object detection applications that rely on the pre-trained and well-optimized convolution-based RPN. Thus, to cope with these challenges and optimize QRPN-based object detection, we train QRPN using heterogeneous knowledge distillation. We set the pre-trained classical RPN and QRPN as a teacher and student model, respectively. Accordingly, well-optimized convolution knowledge of pre-trained model can be transferred to QRPN. Here, as the logits of fast quantum convolution are in range  $[-1, 1]$  and the logits of the classical convolution are in the range  $(-\infty, \infty)$ , we normalize both logits of the classical convolution and quantum convolution using  $ReLU$ . Therefore, we make both logits are in the same space  $[0, \infty)$ . With the normalization, we design the classical to quantum (C2Q) loss function as

$$\mathcal{L}_{C2Q}(\theta^Q) = \|\Omega(\mathbf{x}; \theta^Q) - \Omega(\mathbf{x}; \theta^C)\|, \quad (11)$$

where  $\theta^Q$  and  $\theta^C$  denote the trainable parameters of QRPN and classical RPN, respectively.  $\Omega(\mathbf{x}; \theta^Q)$  and  $\Omega(\mathbf{x}; \theta^C)$  denote the outputs of QRPN and RPN when the input tensor  $\mathbf{x}$ , respectively. Note that the dimensions and sizes of the inputs and outputs of QRPN are designed to be identical to those of RPN.

### D. Loss of QRPN

As our QRPN is designed to be activated similarly to the classical RPN, which first regresses the box and then classifies the image within the box, we incorporate the heterogeneous

TABLE I  
NOTATIONS AND IMPLEMENTATION DETAILS.

Notations for quantum computing	
$ \psi_x\rangle$	The quantum state encoded with data $x$ .
$q$	The number of available qubits.
$\langle O_x \rangle$	The observable derived by quantum state $ \psi_x\rangle$ .
$\Gamma$	Pauli- $\Gamma$ gate, <i>e.g.</i> , $\Gamma \in \{X, Y, Z\}$ .
$R_\Gamma$	Rotation- $\Gamma$ gate, <i>e.g.</i> , $\Gamma \in \{X, Y, Z\}$ .
$CT$	Controlled- $\Gamma$ gate, <i>e.g.</i> , $\Gamma \in \{X, Y, Z\}$ .
$\mathcal{M}$	Measurement operator.
$I$	Identity matrix.
Notations for fast quantum convolution	
$X^l$	The classical 3D tensor image of $l$ -th layer.
$P^l$	The transformed 2D tensor matrix of $l$ -th layer.
$ \psi\rangle_i$	The transformed $i$ -th row quantum states.
$U_E$	The un-trainable encoding gates.
$U_T$	The trainable PQC gates.
$U_{1 \times 1}$	The unitary matrices with 1-qubit.
$(H^l, W^l, C^l)$	(Height, width, channels) of $l$ -th layer.
Notations for QCOD.	
$C$	The number of activated channels $\{16, 32, 64\}$ .
$\gamma$	The distillation coefficients $\{0.1, 0.3, \dots, 0.9\}$ .
$\mathcal{L}_{(C2Q, cls, reg)}$	The utilized losses (KD, cls, reg).
$\Omega(\mathbf{x}; \theta^Q)$	Output logit from quantum convolution layers.
$\Omega(\mathbf{x}; \theta^C)$	Output logit from classical convolution layers.

knowledge distillation regularizer into each loss function of the QRPN ( $L_{cls}$  and  $L_{reg}$ ). The total loss functions are designed as

$$L_{total} = \frac{1-\gamma}{N_c} \sum_{c=1}^{N_c} L_{cls} + \frac{\lambda(1-\gamma)}{N_r} \sum_{r=1}^{N_r} L_{reg} + \gamma \mathcal{L}_{C2Q} \quad (12)$$

where  $L_{reg}$  and  $L_{cls}$  are formulated as presented in [37]. In addition,  $\gamma$  denotes the heterogeneous knowledge distillation parameter  $0 \leq \gamma \leq 1$ . Lastly,  $\lambda$  denotes the normalization parameters between  $L_{reg}$  and  $L_{cls}$ .

### E. Implementation details

To realize and simulate our fast quantum convolution and its application, QCOD, we implement QCOD with the following details. Table. I shows our implementation details. Particularly, in all experiments, we set the number of channels  $C = 64$  among 256 total channels utilized in [37]. The *2d-conv* illustrated in Fig. 5 is designed as a 4-qubit PQC. Both *cls* and *reg* in Fig. 5 are designed with 1-qubit PQC. To design PQC, this paper employs trainable *U3CU3* layers which is composed of *CNOT* gates and other rotation gates [38]. The trainable gate is employed for the 1-qubit PQC. As an activation layer, we utilize *ReLU* function.

## VI. PERFORMANCE EVALUATION

First of all, this section introduces following three hypotheses those are the main items which should be verified and discussed.

- **Hypothesis 1.** Our fast quantum convolution has the potential to incorporate multiple pieces of information.

- **Hypothesis 2.** Channel uploading has advantages in reducing computational complexity rather than classical quantum information encoding strategy.
- **Hypothesis 3.** QCOD has trainability, and heterogeneous knowledge distillation has advantages in training QCOD.

To corroborate Hypothesis 1, we visualize quantum states with corresponding encoded channels. As multi-qubit states are challenging to be described [39], we utilize  $1 \times 1$  PQC. In addition, to compare the advantages of channel uploading by adjusting the number of channels uploaded to corroborate Hypothesis 2. Finally, to verify Hypothesis 3, we visualize the results of QCOD and measure the performance of QCOD with various training strategies.

### A. Setup

This paper conducts all the experiments using classical computing with a Linux-based machine which has Intel i9-10990k, NVIDIA Titan X (2ea), and RAM 128GB. Python v3.8.10 and quantum computing simulation libraries (*e.g.*, torchquantum v0.1.5 [40], pytorch v1.8.2 LTS) are employed for simulating the fast quantum convolution. For QCOD simulation, we use KITTI dataset [41] which is a well-known object detection dataset for autonomous driving. All of the comparative models are trained with 15 epochs. This paper utilizes a pre-trained feature extraction module and classification module using VGG-16 [42]. Each module is pre-trained using IMAGENET [42]. The initial learning rate is set to  $10^{-2}$ .

**Comparison techniques.** To corroborate the advantages of the fast quantum convolution and its application, QCOD, we designed various simulations involving classification and object detection. The brief explanation for each model used in simulations is as follows:

- 1) *FQC (proposed)*: *FQC (proposed)* is a simple convolution layer designed using the fast quantum convolution. The input features are encoded via channel uploading and the output features are decoded via channel reconstruction.
- 2) *QCOD (proposed)*: The *QCOD (proposed)* is designed using the fast quantum convolution for object detection. This model is trained via heterogeneous knowledge distillation.
- 3) *QCOD (w/o kd)*: *QCOD (w/o kd)* is a ablation version of *QCOD (proposed)*. The heterogeneous knowledge distillation parameter  $\gamma$  is set to 0.
- 4) *F-RCNN (baseline)*: *F-RCNN (baseline)* is a baseline model to implement QCOD. F-RCNN (baseline) is also used as a teacher network to train QCOD (*proposed*) via heterogeneous knowledge distillation.
- 5) *Quantum convolution*: Quantum convolution is designed as an existing quantum convolution with patch processing. This model utilizes PQC for each channel.

### B. Evaluation Results

This paper corroborates the performance and feasibility of our fast quantum convolution through experiments represented in Fig.6 and Fig.7(a). In addition, this paper conducts object

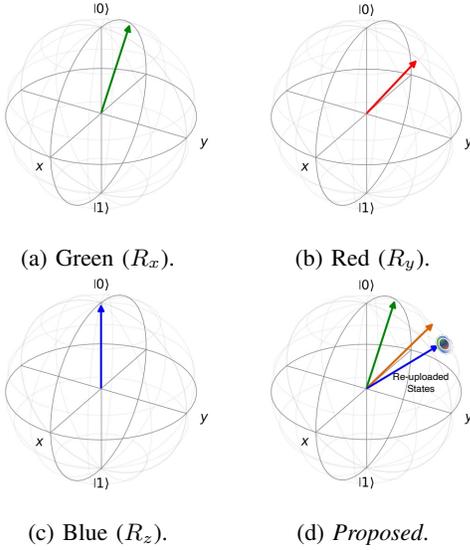


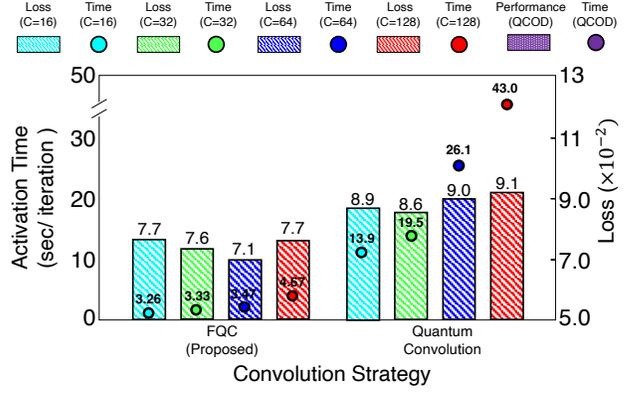
Fig. 6. Visualization of quantum states using various encoding strategies. (a-c) represents quantum states encoded with the classical quantum convolution strategy. (d) represents the states of fast quantum convolution. To visualize the quantum state, a 1-qubit (PQC) is utilized. The (green, orange, blue) vectors in (d) represent the (first, second, final) quantum states with uploaded channels.

TABLE II  
MAP @0.5 (%) OF QCOD AND COMPARATIVE MODELS  
WITH 64 CHANNELS ON KITTI DADASET.

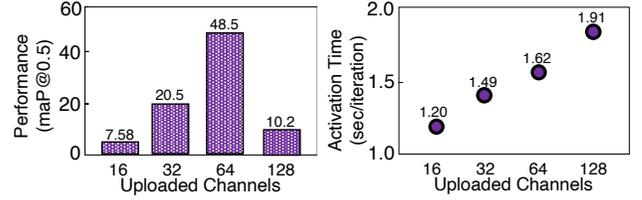
Dataset	Model	Heterogeneous KD ( $\gamma$ )			
		0	0.1	0.3	0.5
KITTI	QCOD (proposed)	-	48.5	<b>51.2</b>	49.1
	QCOD (w/o KD)	21.1	-	-	-
	F-RCNN (baseline)	<b>68.4</b>	-	-	-

detection computation using various training methods and numbers of channels in experiments, as represented in Table II and Fig.7(b-c). Finally, Fig.8 shows the results of object detection computation using our proposed QCOD with our fast quantum convolution on KITTI dataset.

**Performance of fast quantum convolution.** Fig. 6 and Fig. 7 (a) corroborate the feasibility and trainability of fast quantum convolution. We observe that features extracted from our fast quantum convolution in Fig. 7 (a) can be trained using classical optimization techniques. In Fig. 6 (a-c), the encoded states with classical quantum convolution achieve each channel information. Therefore, the number of required qubits equals the number of the input channels. However, As represented in Fig. 6 (d), our fast quantum convolution generates superpositioned quantum states that contain information about the channels. Furthermore, compared with Fig. 6 (c), our proposed model shows more clear representation ability even using same encoding gate. Based on the results presented in Fig. 7 (a), we observe that our fast quantum convolution executes more quickly than other quantum convolution methods. However, a significant decrease in performance is also observed as the number of uploaded channels exceeds the threshold value. This phenomenon is attributed to information loss resulting from the overlap of encoded information on qubits rather than a one-to-one encoding. It is a remaining challenge of our fast quantum convolution, and with the rapid advancement of



(a) Average iteration time and average loss of various convolution strategies.



(b) Performance of QCOD. (c) Iteration time of QCOD.

Fig. 7. Performance and loss comparison of various convolution strategies and their object detection applications (a) is conducted on the  $32 \times 32$  size CIFAR10 dataset. (b) and (c) are measured on the resized  $1382 \times 512$  KITTI dataset with 1 batch-size. In experiments (b) and (c), we set knowledge distillation parameter  $\gamma = 0.1$ . We utilize the set of  $R_y$  gates for uploading.

quantum computing, this challenge can be mitigated as more qubits become available.

**Feasibility of QCOD.** Through extensive experiments, we verify the feasibility of QCOD as an actual object detection application. Fig 7 (b) and (c) represent the performance and activation time of our QCOD according to the number of utilized channels. With the 64 channels, our QCOD shows high performance, even in complex object detection application. As illustrated in Fig.8, QCOD effectively draws bounding boxes tailored to objectives and transfers information, enabling object classification. As the first quantum version of object detection, QCOD shows the feasibility of quantum applications using our fast quantum convolution. To improve our QCOD, finding the optimal number of uploaded channels remains challenging. As QCOD increases the number of channels 32 to 64, the QCOD achieves 38% performance gain. On the other hand, when the number of channels becomes 128, a severe performance degradation is observed. In addition, Fig. 7 (c) represents that the performance and activation time are not proportional. Therefore, considering the performance and activation time, finding the optimal number of uploaded channels is crucial.

**Advantages and disadvantages of heterogeneous knowledge distillation.** We investigate the impact of heterogeneous knowledge distillation on the performance of QCOD, which affects the practical utilization of quantum-based applications. Table II shows the results of heterogeneous knowledge distillation. When we apply heterogeneous knowledge distillation with a parameter value of  $\lambda = 0.3$ , classical knowledge is effectively transferred, leading to an enhancement in QCOD performance.



Fig. 8. QCOD on KITTI dataset with heterogeneous knowledge distillation parameter  $\gamma = 0.3$  and with 64 number of channels. The yellow, green and pink boxes include pedestrians, cyclists and cars, respectively.

However, when we do not utilize heterogeneous knowledge distillation, our QCOD achieves a score of 21.1. These results highlight the substantial performance improvement can be achieved when utilizing heterogeneous distillation. On the other hand, when training with  $\gamma > 0.3$ , we observe a performance degradation. This result indicates a existence of the threshold value for  $\gamma$ , underscoring the significance of finding the appropriate  $\gamma$  value tailored to the objectives of

each QCOD applications.

## VII. CONCLUDING REMARKS

This paper proposes the fast quantum convolution and its practical application in object detection, named QCOD. With our fast quantum convolution, which uploads input channel information and reconstructs output channel information, we observe the feasibility of quantum-based applications. To

implement our fast quantum convolution in QCOD, we design a training method using heterogeneous knowledge distillation. By adopting knowledge distillation to transfer knowledge from the classical object detection domain to the quantum object detection domain, QCOD achieves robustness in object detection and adequately training the PQC. We analyze the complexity of our fast quantum convolution and QCOD to verify the advantages of our fast quantum convolution. Through extensive simulations, this paper corroborates i) the trainability of our fast quantum convolution, ii) the advantages of heterogeneous knowledge distillation, and iii) the feasibility of our QCOD.

## REFERENCES

- [1] X. Chen and X. Lin, "Big data deep learning: Challenges and perspectives," *IEEE Access*, vol. 2, pp. 514–525, 2014.
- [2] S. Whang and J. Lee, "Data collection and quality challenges for deep learning," *Proceedings of the VLDB Endowment*, vol. 13, no. 12, pp. 3429–3432, 2020.
- [3] S. Karimi, A. Karimi, and A. Vahidi, "Level-K reasoning, deep reinforcement learning, and Monte Carlo decision process for fast and safe automated lane change and speed management," *IEEE Transactions on Intelligent Vehicles*, vol. 8, no. 6, pp. 3556–3571, 2023.
- [4] S. Ansari, F. Naghdy, and H. Du, "Human-machine shared driving: Challenges and future directions," *IEEE Transactions on Intelligent Vehicles*, vol. 7, no. 3, pp. 499–519, 2022.
- [5] B. Zhang, T. Chen, B. Wang, and R. Li, "Joint distribution alignment via adversarial learning for domain adaptive object detection," *IEEE Transactions on Multimedia*, vol. 24, pp. 4102–4112, 2022.
- [6] C. Chen, S. Dong, Y. Tian, K. Cao, L. Liu, and Y. Guo, "Temporal self-ensembling teacher for semi-supervised object detection," *IEEE Transactions on Multimedia*, vol. 24, pp. 3679–3692, 2022.
- [7] Z. Wu, S. Li, C. Chen, A. Hao, and H. Qin, "Deeper look at image salient object detection: Bi-stream network with a small training dataset," *IEEE Transactions on Multimedia*, vol. 24, pp. 73–86, 2022.
- [8] Z. Ren, Q. Kong, J. Han, M. D. Plumbley, and B. W. Schuller, "CAA-Net: Conditional atrous CNNs with attention for explainable device-robust acoustic scene classification," *IEEE Transactions on Multimedia*, vol. 23, pp. 4131–4142, 2021.
- [9] K. O'Shea and R. Nash, "An introduction to convolutional neural networks," *CoRR*, vol. abs/1511.08458, 2015.
- [10] Y. Yan, Y. Mao, and B. Li, "Second: Sparsely embedded convolutional detection," *Sensors*, vol. 18, no. 10, p. 3337, 2018.
- [11] F. Arute, K. Arya, R. Babbush, D. Bacon, J. C. Bardin, R. Barends, R. Biswas, S. Boixo, F. G. Brandao, D. A. Buell *et al.*, "Quantum supremacy using a programmable superconducting processor," *Nature*, vol. 574, no. 7779, pp. 505–510, 2019.
- [12] M. Cerezo, G. Verdon, H. Huang, L. Cincio, and P. J. Coles, "Challenges and opportunities in quantum machine learning," *Nature Computational Science*, vol. 2, no. 9, pp. 567–576, 2022.
- [13] N. Aburaed, F. S. Khan, and H. Bhaskar, "Advances in the quantum theoretical approach to image processing applications," *ACM Computing Surveys*, vol. 49, no. 4, pp. 75:1–75:49, 2017.
- [14] H. Baek, W. J. Yun, S. Park, and J. Kim, "Stereoscopic scalable quantum convolutional neural networks," *Neural Networks*, vol. 165, pp. 860–867, August 2023.
- [15] F. Shen and J. Liu, "QFCNN: Quantum Fourier convolutional neural network," *CoRR*, vol. abs/2106.10421, 2021.
- [16] S. Oh, J. Choi, J. Kim, and J. Kim, "Quantum convolutional neural network for resource-efficient image classification: A quantum random access memory (QRAM) approach," in *Proc. of the IEEE International Conference on Information Networking (ICOIN)*, Jeju Island, South Korea, January 2021, pp. 50–52.
- [17] I. Kerentidis, J. Landman, and A. Prakash, "Quantum algorithms for deep convolutional neural networks," in *Proc. of the International Conference on Learning Representations (ICLR)*, Addis Ababa, Ethiopia, April 2020.
- [18] R. LaRose and B. Coyle, "Robust data encodings for quantum classifiers," *Physical Review A*, vol. 102, no. 3, p. 32420, 2020.
- [19] A. Pérez-Salinas, A. Cervera-Lierta, E. Gil-Fuster, and J. I. Latorre, "Data re-uploading for a universal quantum classifier," *Quantum*, vol. 4, p. 226, 2020.
- [20] L. Friedrich and J. Maziero, "Avoiding barren plateaus with classical deep neural networks," *Physical Review A*, vol. 106, no. 4, p. 042433, 2022.
- [21] C. O. Marrero, M. Kieferová, and N. Wiebe, "Entanglement induced barren plateaus," *PRX Quantum*, vol. 2, no. 4, p. 040316, 2021.
- [22] A. Uvarov and J. D. Biamonte, "On barren plateaus and cost function locality in variational quantum algorithms," *Journal of Physics A: Mathematical and Theoretical*, vol. 54, no. 24, p. 245301, 2021.
- [23] M. Schuld, R. Sweke, and J. J. Meyer, "Effect of data encoding on the expressive power of variational quantum-machine-learning models," *Physical Review A*, vol. 103, no. 3, p. 032430, 2021.
- [24] F. Sarfraz, E. Arani, and B. Zonooz, "Knowledge distillation beyond model compression," in *Proc. of the IEEE International Conference on Pattern Recognition (ICPR)*, Milan, Italy, January 2020, pp. 6136–6143.
- [25] X. Cui, C. Wang, D. Ren, Y. Chen, and P. Zhu, "Semi-supervised image deraining using knowledge distillation," *IEEE Transactions on Circuits and Systems and Video Technology*, vol. 32, no. 12, pp. 8327–8341, 2022.
- [26] J. Landman, S. Thabet, C. Dalyac, H. Mhiri, and E. Kashefi, "Classically approximating variational quantum machine learning with random Fourier features," in *Proc. of the International Conference on Learning Representations (ICLR)*, Kigali, Rwanda, May 2023.
- [27] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. of the Advances in Neural Information Processing Systems (Neurips)*, Long Beach, CA, USA, December 2017, pp. 5998–6008.
- [28] A. Czerwinski, "Quantum tomography of entangled qubits by time-resolved single-photon counting with time-continuous measurements," *Quantum Information Processing*, vol. 21, no. 9, p. 332, 2022.
- [29] H. Huang, R. Kueng, G. Torlai, V. V. Albert, and J. Preskill, "Provably efficient machine learning for quantum many-body problems," *Science*, vol. 377, no. 6613, p. eabk3333, 2022.
- [30] C. S. Rohwedder, J. P. L. de Carvalho, J. N. Amaral, G. Araújo, G. Colmenares, and K. A. Wang, "Pooling acceleration in the DaVinci architecture using Im2col and Col2im instructions," in *Proc. of the IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, Portland, OR, USA, June 2021, pp. 46–55.
- [31] W. Xie and W. Li, "Entanglement properties of random invariant quantum states," *Quantum Information and Computing*, vol. 22, no. 11&12, pp. 901–923, 2022.
- [32] L. Banchi and G. E. Crooks, "Measuring analytic gradients of general quantum evolution with the stochastic parameter shift rule," *Quantum*, vol. 5, p. 386, 2021.
- [33] D. Wierichs, J. Izaac, C. Wang, and C. Y. Lin, "General parameter-shift rules for quantum gradients," *Quantum*, vol. 6, p. 677, 2022.
- [34] A. Paler, O. Oumarou, and R. Basmadjian, "Parallelizing the queries in a bucket-brigade quantum random access memory," *Physical Review A*, vol. 102, no. 3, p. 032608, 2020.
- [35] F. Sultana, A. Sufian, and P. Dutta, "A review of object detection models based on convolutional neural network," *Intelligent computing: image processing based applications*, pp. 1–16, 2020.
- [36] J. H. Bappy and A. K. Roy-Chowdhury, "CNN based region proposals for efficient object detection," in *Proc. of the IEEE International Conference on Image Processing (ICIP)*, Phoenix, AZ, USA, September 2016, pp. 3658–3662.
- [37] S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1137–1149, 2017.
- [38] H. Wang, Y. Ding, J. Gu, Y. Lin, D. Z. Pan, F. T. Chong, and S. Han, "Quantumnas: Noise-adaptive search for robust quantum circuits," in *Proc. of the IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, Seoul, South Korea, April 2022, pp. 692–708.
- [39] S. Ruan, Y. Wang, W. Jiang, Y. Mao, and Q. Guan, "VACSEN: A visualization approach for noise awareness in quantum computing," *IEEE Transactions on Visualization and Computer Graphics*, vol. 29, no. 1, pp. 462–472, 2022.
- [40] H. Wang, Z. Li, J. Gu, Y. Ding, D. Z. Pan, and S. Han, "QOC: Quantum on-chip training with parameter shift and gradient pruning," in *Proc. of the IEEE/ACM Design Automation Conference (DAC)*, San Francisco, CA, USA, June 2022, pp. 665–660.
- [41] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The KITTI dataset," *International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.
- [42] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. of the International Conference on Learning Representations (ICLR)*, San Diego, CA, USA, May 2015.



**Hankyul Baek** is currently a Ph.D. student in electrical and computer engineering at Korea University, Seoul, Republic of Korea, since March 2021, where he received his B.S. in electrical engineering in 2020. He was with LG Electronics, Seoul, Republic of Korea, from 2020 to 2021. He was also a visiting scholar at the Department of Electrical and Computer Engineering, The University of Utah, Salt Lake City, UT, USA, in 2023.

His current research interests include quantum machine learning and its applications.

**Dr. Donghyeon Kim** is currently a Principal Researcher at the Institute of Advanced Technology Development (IATD), Hyundai Motor and Kia Corporation, Republic of Korea.



**Prof. Joongheon Kim** (Senior Member, IEEE) has been with Korea University, Seoul, Korea, since 2019, where he is currently an associate professor at the Department of Electrical and Computer Engineering and also an adjunct professor at the Department of Communications Engineering (co-operated by Samsung Electronics) and the Department of Semiconductor Engineering (co-operated by SK Hynix). He received the B.S. and M.S. degrees in computer science and engineering from Korea University, Seoul, Korea, in 2004 and 2006; and the Ph.D. degree

in computer science from the University of Southern California (USC), Los Angeles, CA, USA, in 2014. Before joining Korea University, he was a research engineer with LG Electronics (Seoul, Korea, 2006–2009), a systems engineer with Intel Corporation Headquarter (Santa Clara in Silicon Valley, CA, USA, 2013–2016), and an assistant professor of computer science and engineering with Chung-Ang University (Seoul, Korea, 2016–2019).

He serves as an editor and guest editor for *IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY*, *IEEE INTERNET OF THINGS JOURNAL*, *IEEE TRANSACTIONS ON MACHINE LEARNING IN COMMUNICATIONS AND NETWORKING*, *IEEE COMMUNICATIONS STANDARDS MAGAZINE*, *Computer Networks*, and *ICT Express*. He is also a distinguished lecturer for *IEEE Communications Society (ComSoc)* and *IEEE Systems Council*. He is an executive director of the Korea Institute of Communication and Information Sciences (KICS). He was a recipient of Annenberg Graduate Fellowship with his Ph.D. admission from USC (2009), Intel Corporation Next Generation and Standards (NGS) Division Recognition Award (2015), *IEEE SYSTEMS JOURNAL* Best Paper Award (2020), *IEEE ComSoc Multimedia Communications Technical Committee (MMTC)* Outstanding Young Researcher Award (2020), *IEEE ComSoc MMTC Best Journal Paper Award* (2021), Best Special Issue Guest Editor Award by *ICT Express* (2022), and Best Editor Award by *ICT Express* (2023). He also received several awards from IEEE conferences including *IEEE ICOIN Best Paper Award* (2021), *IEEE Vehicular Technology Society (VTS) Seoul Chapter Awards* (2019, 2021, and 2022), and *IEEE ICTC Best Paper Award* (2022).