# pvCNN: Privacy-Preserving and Verifiable Convolutional Neural Network Testing

Jiasi Weng, Jian Weng*, *Member, IEEE*, Gui Tang, Anjia Yang, Ming Li, Jia-Nan Liu

*Abstract*—We propose a new approach for privacy-preserving and verifiable convolutional neural network (CNN) testing in a distrustful multi-stakeholder environment. The approach is aimed to enable that a CNN model *developer* convinces a *user* of the truthful CNN performance over non-public data from *multiple testers*, while respecting model and data privacy. To balance the security and efficiency issues, we appropriately integrate three tools with the CNN testing, including collaborative inference, homomorphic encryption (HE) and zero-knowledge succinct non-interactive argument of knowledge (zk-SNARK).

We start with strategically partitioning a CNN model into a private part kept locally by the model developer, and a public part outsourced to an outside server. Then, the private part runs over the HE-protected test data sent by a tester, and transmits its outputs to the public part for accomplishing subsequent computations of the CNN testing. Second, the correctness of the above CNN testing is enforced by generating zk-SNARK based proofs, with an emphasis on optimizing proving overhead for two-dimensional (2-D) convolution operations, since the operations dominate the performance bottleneck during generating proofs. We specifically present a new quadratic matrix program (QMP)-based arithmetic circuit with *a single multiplication gate* for expressing 2-D convolution operations between multiple filters and inputs in a batch manner. Third, we aggregate multiple proofs with respect to a same CNN model but different testers' test data (*i.e.*, different statements) into one proof, and ensure that the validity of the aggregated proof implies the validity of the original multiple proofs. Lastly, our experimental results demonstrate that our QMP-based zk-SNARK performs nearly $13.9\times$ faster than the existing quadratic arithmetic program (QAP)-based zk-SNARK in proving time, and $17.6\times$ faster in Setup time, for high-dimension matrix multiplication. Besides, the limitation on handling a bounded number of multiplications of QAP-based zk-SNARK is relieved.

## I. INTRODUCTION

Convolutional neural networks (CNNs) [1], [2] have been widely applied in various application scenarios, such as healthcare analysis, autonomous vehicle and face recognition. But real-life reports demonstrate that neural networks often exhibit erroneous decisions, leading to disastrous consequences, *e.g.*, self-driving crash due to the failure of identifying unexpected driving environments [3], and prejudice due to racial biases embedded in face recognition systems [4]. The reports emphasize that when applying CNN models in security-critical

J.s. Weng, J. Weng, G. Tang, A. Yang and M. Li are with the College of Cyber Security of Jinan University, Guangzhou 510632, and Pazhou Lab, Guangzhou 510335, China. J.-N. Liu is affiliated with the School of Computer Science and Technology, Dongguan University of Technology, Dongguan 523808, and Guangzhou Fongwell Data Limited Company, Guangzhou 511400, and Pazhou Lab, 510335, China. E-mail: wengjiasi@gmail.com, cryptjweng@gmail.com, guitang001@gmail.com, anjiayang@gmail.com, limjnu@gmail.com, j.n.liu@foxmail.com. Jian Weng is the corresponding author.

scenarios, model users should be sufficiently cautious of benchmarking the CNN models to obtain the truthful multi-faceted performance, like robustness and fairness, not limited to natural accuracy. For example, when a CNN model is deployed in a self-driving car for identifying camera images, a user should be assured that the CNN model is always highly accurate, even for perturbed images; for face recognition applications, users want to know that the underlying CNN models can accurately recognize face images without discrimination.

A widely adopted approach to benchmarking a CNN model is black-box testing [5], [6]. Black-box testing enables users to have a black-box access to the CNN model, that is, feeding the model with a batch of test data and merely observing its outputs, *e.g.*, the proportion of correct predictions. Such approach is suitable for the setting where model users and model developers are not the same entities. The main reasons are that the parameters of CNN models are intellectual properties and are vulnerable to privacy attacks, such that developers are unwilling to reveal the parameters. We are also aware that existing excellent white-box testing approaches [7] can test CNN models in a fine-granularity fashion using model parameters, but black-box testing is preferred in the above setting where model parameters are inaccessible.

Starting by black-box testing, two essential issues need consideration to probe a CNN model's truthful and multi-faceted performance. (***i***) **Black-box testing strongly relies on the test dataset [7], which in turn requires multi-source data support in reality.** In order to support testing a model's multi-faceted performance, the available test dataset should be as various as possible. One evidence is that multiple benchmarks are built by embedding dozens of types of noises into ImageNet to measure the robustness of CNN models [8]. Despite the necessity, building such benchmarks consumes many manpower and multi-party efforts, even for large companies, *e.g.*, a recent project named Crowdsourcing Adverse Test Sets for Machine Learning (CATS4ML) launched by Google Research [9]. Thus, the off-the-shelf test datasets for supporting multi-faceted CNN testing are likely from multiple sources. (***ii***) **Black-box testing opens a door for untrusted model developers to cheat in testing.** They might forge untruthful outputs without correctly running the processes of CNN testing on given test datasets. Additionally, if developers can in advance learn a given test dataset, they might craft a CNN which typically adapts to the test dataset, which deviates from our initial goal of testing their truthful multi-faceted performance. From this point, test datasets cannot be public, not merely due to that datasets themselves are privacy-sensitive in many security-critical scenarios. Motivated by the

two issues above, there needs an approach for users to validate the correctness of the black-box CNN testing over multi-source test datasets, in which test datasets are not public and the CNN model is privacy-preserving.
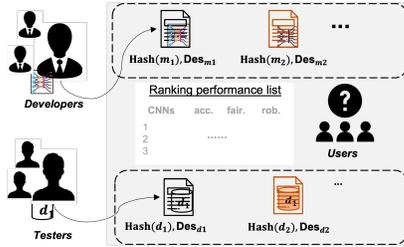


Fig. 1: Scenario example ($m_1, m_2$ refer to CNN models, and $d_1, d_2$ refer to test datasets; Des means non-private descriptions for the CNN models or test datasets).

While many awesome efforts have been done to make CNN prediction/testing verifiable [10]–[17] using cryptographic proof techniques [18]–[23], we still need a new design satisfying our scenario requirements. We now clarify our scenario requirements, which cannot be fully satisfied by the previous work (see explanation in Section II). As shown in Fig. 1, a publicly accessible platform (in gray color), *e.g.*, ModelZoo [24], Kaggle [25] and AWS Marketplace [26], can allow third-party *developers* to advertise their CNN models towards users for earning profits. For *users* who will pay for a CNN model, they may be understandably concerned about its performance, given previous reports of misleading claims. The platform might have a motivation to maintain a good reputation for sustainability by publishing good models. To address the concerns, the platform can announce a crowdsourcing task to test the CNN model, such that *testers* participate in probing the model performance using their test datasets. Note that the model parameters and test datasets are not disclosed on the platform, while the model architectures can be public. Our design empowers untrusted CNN developers by providing them with a way to prove the truthful performance of their models over multiple test datasets to convince potential users. The design needs to satisfy three requirements as follows:

(a) **Public verifiability and no need of interaction.** Since the user is later-coming and not pre-designated, proof generation and verification w.r.t the CNN testing should not share private information between the developer (*i.e.*, prover) and the user (*i.e.*, verifier), and thereby requiring public verifiability. The proving ability should also be non-interactive, due to that the CNN developer, the user and testers are not always simultaneously online.

(b) **Privacy preservation.** During the CNN testing, model parameters should not be exposed to users, since they may be intellectual properties for profits, and the parameters memorize sensitive training data [27]. Also, the test data should also be protected against the CNN developer, considering the CNN developer might strategically craft the CNN model based on the test data. Besides, any verifier should fail to learn private information from the final CNN performance and proofs.

(c) **Batch proving and verification.** A CNN model usually receives test data in batches, so it is a natural requirement to support batch operations in proof generation. In terms of verification in our scenario, a later-coming user has to verify multiple proofs w.r.t a single CNN model which is tested with the multi-tester test data. Hence, enabling the user to verify multiple proofs in a batch manner can be another desirable requirement.

### A. Our Designs

To satisfy the three requirements, we populate our designs by strategically integrating fully homomorphic encryption (FHE), zero-knowledge Succinct Non-interactive ARgument of Knowledge (zk-SNARK) and an idea of collaborative inference [28], [29]. We adopt zk-SNARK-based proof systems, mainly due to that the class of proof technique supports many properties like public delegatability and public verifiability which suit our scenario involving multiple distrustful entities. To summarize, our designs include three components: (1) privacy-preserving CNN testing based on an FHE algorithm and collaborative inference, (2) generating zk-SNARK proofs for the above privacy-preserving CNN testing, with an emphasis on proposing a new quadratic matrix programs (QMP)-based zk-SNARK for proving 2-D convolutional relations, and (3) enabling zk-SNARK proof aggregation for verifying multiple proofs in a batch manner.

Firstly, we start by letting a developer locally run the CNN testing process over the FHE-protected test data from every tester. Our starting point is to enable performing the whole CNN testing processing over ciphertext domain. While FHE can provide stronger data confidentiality, compared to other secure computing techniques, such as secure multi-party computation and differential privacy [30], the straightforward adoption is computation-expensive, nearly 4 to 5 orders slower than computing on unencrypted data. Moreover, proving overhead on encrypted data (in ring field $\mathcal{R}_q$) and encoded CNN parameters (in ring field $\mathcal{R}_t$) becomes unaffordable, compared to generating proofs over plaintext domain. We hence find a middle point for balancing privacy and efficiency, by introducing a collaborative inference strategy, and thereby making partial testing process run on ciphertext domain. Particularly, the CNN developer can split his CNN model into two parts: one is private, called PriorNet, and another one is less private, called LaterNet. Then, the developer can locally evaluate PriorNet on the FHE-protected test data, and delegate LaterNet to a computation-powerful server for accomplishing the subsequent computation by feeding it with the plaintext outputs of PriorNet. Next, due to that the plaintext outputs of PriorNet might be susceptible to model inversion attacks [28] by the malicious server, an off-the-shelf and generic adversarial training strategy [29] is able to protect the PriorNet's outputs against the attacks.

The second step is to prove the correctness of the above splitting CNN testing on FHE-encrypted test data (*i.e.*, polynomial ring elements). We concretely adopt a proving roadmap recently proposed by Fiore *et al.* [22], in which Step 1 is to prove the simultaneous evaluation of multiple polynomial ring elements in the same randomly selected point (for PriorNet), and Step 2 is to prove the satisfiability of quadratic arithmetic program (QAP)-based arithmetic circuits (for PriorNet and LaterNet) whose inputs and outputs are committed.

But we do not directly apply Step 2 into our case, since the direct adoption causes high proving and storage overhead due to a large number of multiplication gates for handling convolutional relations [12]. We thus present a new method to reduce the number of multiplication gates during proof generation, for improving the proving time of the convolutional relations. At first, we represent the 2-D convolution operations between $M$ filters which are $m \times m$ matrices and $Mn^2$ test inputs which are $n \times n$ matrices into *a single matrix multiplication (MM)* computation. The MM computation is conducted between a $Mn^2 \times Mn^2$ filter matrix which is strategically assigned with the $M$ filters and zeros as padding for operation correctness, and another $Mn^2 \times Mn^2$ input matrix which is assigned with the $Mn^2$ test inputs. After that, we start from the QAP-based zk-SNARK and present a new QMP formula [31]. We then express the above single MM computation as a circuit over $Mn^2 \times Mn^2$ matrices *by an only one-degree QMP*. As a result, the number of multiplication gates always is 1 and the proving time linearly increases by the matrix dimension, *i.e.*, $Mn^2 \times Mn^2$. With the effort, we prove the satisfiability of the QMP-based circuits via checking the divisibility in randomly selected point of the set of matrix polynomials of the QMP. Note that there are heterogeneous operations in each neural network layer during CNN testing. We express the convolution and full connection operations as QMP-based circuits, while expressing the activation and max pooling functions as QAP-based circuits, and then separately generate specialized proofs for them. Furthermore, we add *commit-and-prove* (CaP) components based on LegoSNARK [21] for gluing the separately generated specialized proofs, aiming to prove that the current layer indeed takes as inputs the previous layer's outputs.

Thirdly, we enable batch verification via proof aggregation, and only a single proof is generated for a CNN model tested by multi-tester inputs. We aggregate multiple proofs with regard to the same CNN but different test inputs from multiple testers based on Snarkpack [32]. We note that the aggregation computation and proof validation can be delegated to computation-powerful parties in a competition manner, *e.g.*, decentralized nodes who run a secure consensus protocol to maintain the publicly accessible platform previously described in Fig. 1. The validation result will be elected according to the consensus protocol, and then uploaded to the public platform and attached to the corresponding CNN model. To the end, a future user who is interested in the CNN model can enjoy a lightweight verification by merely checking the validation result, via browsing the platform.

In summary, this work makes the following contributions:

- We generate zk-SNARK proofs for the CNN testing based on FHE and collaborative inference, which protects both of the model and data privacy and ensures the integrity.
- We present a new QMP-based arithmetic circuit to express convolutional relations for efficiency improvement.
- We aggregate multiple proofs respective to a same CNN model and different testers' inputs for reducing the entire verification cost.
- We give a proof-of-concept implementation, and the experimental results demonstrate that our QMP-based method for matrix multiplication performs about $17.6\times$

faster than the existing QAP-based method in Setup time, and $13.9\times$ faster in proving time. The code is available at https://github.com/muclover/pvCNN.

### B. Organization

In Section II, we revisit related work. In Section III, we introduce intrinsic preliminaries. Section IV elaborates the scenario problem this work is concerned about. Section V presents our solutions with customized proof designs. In Section VI and Section VII, we show the security analysis and experimental results, respectively.

## II. RELATED WORK

TABLE I: Comparison of verifiable CNN prediction/testing schemes. Consider applying a two-dimensional (2-D) convolution to a filter matrix of $m \times m$ and a test data matrix of $n \times n$. The column of Proving Time is for such 2-D convolutions on $M$ filter matrices and $Mn^2$ test data matrices (under plaintext). For simplicity, we note $M > m$ and the channel number of test data is $1$. Herein, the requirements include (a) public verifiability and no need of interaction, (b) privacy preservation of data and model and (c) batch proving and verification, respectively. See concrete description in Section I.

| Schemes | Requirements | | | Proving Time |
|---|---|---|---|---|
| | (a) | (b) | (c) | |
| SafetyNets [10] | ○ | ○ | ○ | $O(M \cdot Mn^2 \cdot (n^2 \cdot m^2))$ |
| Keuffer's [14] | ● | ◑ | ○ | $O(M \cdot Mn^2 \cdot (n^2 \cdot m^2))$ |
| SafeTPU [11] | ○ | ○ | ◑ | $O(M \cdot Mn^2 \cdot (n^2 \cdot m^2 \cdot \frac{m^2}{n^2}))$ |
| Madi's [16] | ○ | ● | ○ | $O(M \cdot Mn^2 \cdot (n^2 \cdot m^2))$ |
| vCNN [12] | ● | ◑ | ○ | $O(M \cdot Mn^2 \cdot (n^2 + m^2))$ |
| VeriML [15] | ● | ● | ○ | $O(M \cdot Mn^2 \cdot (n^2 \cdot m^2))$ |
| ZEN [13] | ● | ◑ | ◑ | $O(M \cdot Mn^2 \cdot (n^2 \cdot m^2 \cdot \frac{1}{\mathsf{fac}}))$ |
| zkCNN [17] | ○ | ◑ | ○ | $O(M \cdot Mn^2 \cdot (n^2 + m^2))$ |
| Ours | ● | ● | ● | $O(Mn^2 \cdot Mn^2)$ |

○ denotes the requirement is not satisfied; ◑ denotes the scheme partially supports the requirement; ● denotes the scheme fully supports the requirement. Besides, fac depends on the real size of filters. ZEN's optimization method may not be effective for small filters. See Table 4 in [13].

**Verifiable CNN Prediction/Testing.** Many prior schemes [10]–[17] fail to fully satisfy our requirements, as summarized in TABLE I. The previous schemes can be roughly classified into two major groups, according to their underlying cryptographic proof techniques, such as Sum-Check like protocols [18], [33], [34] and Groth16 zk-SNARK-based systems [19]–[23]. Based on the underlying techniques, some schemes clearly elaborate how they support privately verifiable and interactive proofs [10], [11], [16]; some schemes support partial privacy protection [12], [13], [16]. A few schemes [35], [36] do not consider CNN. In addition, other promising directions are designing an efficient and memory-scalable proof protocol to support complicated neural network inference [37], and using clever verification methods to convert any semi-honest secure inference into a malicious secure one [38], when entities are allowed to be simultaneously online.

We next differentiate this paper from previous schemes in terms of our requirements. ***Regarding requirement (a).*** We adopt a class of zk-SNARK proof systems that provide public delegatability and public verifiability, which is suitable to our scenario. We depart from previous schemes like SafetyNets,

SafeTPU and Madi et al.'s work, since they generally need not the properties. Compared to the line of Groth16 zk-SNARK-based schemes [12]–[15], we make new contributions considering the following requirements (b) and (c). ***Regarding requirement (b).*** We protect data and model privacy by enabling CNN inference on encrypted test data while ensuring the correctness of the result. Specifically, we protect the test data from the entity who runs the CNN inference, preserve the model privacy from the entity who provides the test data, and prevent any verifier who checks the result correctness from learning private information either of the test data or of the model. Most existing work protects the privacy either of model or of test data depending on which one is treated as the witnesses on the side of a prover. For example, Keuffer's [14], vCNN [12] and zkCNN [17] treat model parameters as the witnesses to be proven, while in ZEN [13], test data is witness and model parameters should be shared between a prover and a verifier. Although Madi's [16] provides privacy protection on the both sides, by leveraging homomorphic encryption and homomorphic message authenticator, the work requires a designated verifier. As requirement (a) mentioned, our work considers a non-designated verifier, since each user as a verifier is not designated in advance. ***Regarding requirement (c).*** We improve the proving time for handling 2-D convolutional operations on batches of filters and test data, as well as aggregate multiple proofs for reducing verification cost. In terms of proving time, Keuffer's and VeriML directly apply the Groth16 zk-SNARK [19] to the 2-D convolutional operations between each filter and each input, resulting in $O(n^2 \cdot m^2)$ proving time. For such convolutional operations, vCNN reduces the proving time to $O(n^2 + m^2)$ by firstly transforming the original convolutional representations of a sum of products into a product of sums representations, and then employing the quadratic polynomial program (QPP) in the original QAP-based zk-SNARK. ZEN also makes effort to reduce the proving time by reducing the number of constraints in the underlying zk-SNARK system based on a new stranded encoding method. Essentially, ZEN's encoding method can be complementary to other zk-SNARK-based applications, including ours, when encoding the original numbers into finite field elements. *Despite the effort, they do not consider the batch case as we previously described. When handling convolutional operations between $M$ filters and $Mn^2$ input data, we have $O(Mn^2 \cdot Mn^2)$ proving complexity. It derives from our two-step effort: first, we strategically represent the $M$ filters into a matrix of $Mn^2 \times Mn^2$ and meanwhile reshape the $Mn^2$ inputs into another matrix; second, we employ a new quadratic matrix program (QMP) in the QAP-based zk-SNARK, making proving time depend on the dimension of matrices,* i.e., *$Mn^2 \times Mn^2$. Besides, the previous work does not consider proof aggregation as our work.*

**Secure Outsourced ML Inference.** Our work delegates partial computation of CNN inference to a strong but untrusted server, which relates to a long line of work about secure outsourced ML inference. One of most relevant work starts from CryptoNets, enabling CNN inference on FHE-encrypted data. CryptoNets [39] inspires many follow-up schemes [38], [40]–[44] that aim to improve accuracy, communication or

computational efficiency by carefully leveraging cryptographic primitives, other than FHE [39], [45]. For example, Huang *et al.* [44] achieve efficient and accurate CNN feature extraction using a secret sharing-based encryption technique, which avoids approximating the ReLU function with a low-degree polynomial that is needed by CryptoNets, and thereby ensures accuracy. However, the integrity of outsourced computation is out of the consideration of the above schemes. Similarly targeted at the scenario of privacy-preserving outsourced computation, another line of work studies a group of toolkits, which may pave the way for achieving more complex privacy-preserving ML. The toolkits support general operations of integer numbers [46] and rational numbers [47], as well as enable large-scale computation [48]; many ML-based applications usually contain such computation characteristics.

## III. PRELIMINARIES

### A. CNN Prediction

We present here a prediction process of a CNN model on a step-by-step basis. The CNN model basically contains two convolution (conv) layers and three full connection (fc) layers in order; other layers between them include activation (act) and max pooling or average pooling (pool) layers, and the output layer is softmax layer. Notice, complex CNN models generally contain the above layers. Specifically, taking as inputs a single-channel test data $\mathbf{X}$ which is a $n \times n$ matrix, *e.g.*, gray-scale image, the layer-by-layer operations can be conducted sequentially: $y_o = f^o(f^{fc}(f^{fc}(f^{fc}(f^{pool}(f^{act}(f^{conv}(f^{pool}(f^{act}(f^{conv}(\mathbf{X})))))))))),$ where $f^o$ is the output function which selects out the maximal value among the values outputted by the last $f^{fc}$, determining the prediction output.

We proceed to elaborate the layer-by-layer operations in detail. Denote a $m \times m$ weight matrix $\mathbf{W}^{(k)}$ in the $k_{th}$ layer. *(1) **Convolutional Layer** $f^{conv}$.* Two dimensional (2-D) convolution operations will be applied to input matrices and weight matrices (also called *filters*). Here, we show a 2-D convolution operation applying to two matrices $\mathbf{X}$ and $\mathbf{W}^{(1)}$ with a $(n - m + 1) \times (n - m + 1)$ matrix $\mathbf{Y}^{(1)}$ as output:

$$\mathbf{Y}^{(1)}[i][j] = \sum_{l_r^{(1)}=0, l_c^{(1)}=0}^{m-1, m-1} \mathbf{X}[i + l_r^{(1)}][j + l_c^{(1)}] \times \mathbf{W}^{(1)}[l_r^{(1)}][l_c^{(1)}],$$

where $i, j \in [0, n - m + 1)$ and the row (resp. column) index of the weight is $l_r^{(1)}$ (resp. $l_c^{(1)}$), and the stride size is 1. The subsequent $f^{conv}$ are done similarly, applying to the previous layer's output matrices and the weight matrices in the current layer. *(2) **Activation Layer** $f^{act}$.* There are two widely used activation functions, applying to each element of the output matrices of $f^{conv}$ in layer $l$, in order to catch non-linear relationships. Concretely, the two functions are $f^{act}(\mathbf{Y}^{(k)}[i][j]) = max(\mathbf{Y}^{(k)}[i][j], 0)$ named the ReLU function, and $f^{act}(\mathbf{Y}^{(k)}[i][j]) = \frac{1}{1+e^{-\mathbf{Y}^{(k)}[i][j]}}$ named the Sigmoid function. *(3) **Pooling Layer** $f^{pool}$.* Average pooling and max pooling are two common pooling functions for reducing the dimension of the output matrices in certain layer $k$. They are applied to each region covered by a

$m \times m$ filter, within each output matrix of the previous $f^{act}$ layer. Specifically, an average pooling function is done by $(\mathbf{Y}^{(k)}[0][0],...,\mathbf{Y}^{(k)}[m-1][m-1])/m^2$, and a max pooling function is by $\max(\mathbf{Y}^{(k)}[0][0],...,\mathbf{Y}^{(k)}[m-1][m-1])$. *(4) Full Connection Layer $f^{fc}$.* In this layer, each output matrix of the previous layer multiplies by each weight matrix, and then add with a bias in the current layer. We will omit the bias, for simplicity. *(5)* After executing the mentioned sequential operations, $f^o$ outputs a prediction label $l_{test}$ for $\mathbf{X}$.

### B. Non-Interactive Zero-Knowledge Arguments

Non-interactive zero-knowledge (NIZK) arguments allow a prover to convince any verifier of the validity of a statement without revealing other information. Groth [19] proposed the most efficient zk-SNARK scheme, with small constant size and low verification time. Our work leverages a CaP Groth16 variant to prove arithmetic circuit satisfiability with committed inputs, parameters and outputs. We here recall some essential preliminaries.

**Bilinear Groups.** We review Type-3 bilinear group $(p, \mathbb{G}_1, \mathbb{G}_2, \mathbb{G}_T, e)$, where $p$ is a prime, and $\mathbb{G}_1, \mathbb{G}_2$ are cyclic groups of prime order $p$. Note that $g \in \mathbb{G}_1$, $h \in \mathbb{G}_2$ are the generators. Then, $e : \mathbb{G}_1 \times \mathbb{G}_2 \to \mathbb{G}_T$ is a bilinear map, that is, $e(g^a, h^b) = e(g, h)^{ab}$, where $a, b \in \mathbb{Z}_p, e(g, h)^{ab} \in \mathbb{G}_T$.

**Arithmetic Circuits.** Arithmetic circuits are the widely adopted computational models for expressing the computation of polynomials over finite fields $\mathbb{F}_p$. An arithmetic circuit is formed with a directed acyclic graph, where each vertex of fan-in two called *gate* is labelled by an operation $+$ or $\times$, and each edge called *wire* is labelled by an operand in the finite field. Besides, there are usually two measures for the complexity of a circuit, such as its size and depth. Specifically, the number of wires of the circuit determines its size; the longest path from inputs to the outputs determine its depth.

**Quadratic Arithmetic Program (QAP) [49]** For an arithmetic circuit with $n_{in}$ input variables, $n_{ot}$ output variables in $\mathbb{F}_p$ ($p$ is a large prime) and $n_*$ multiplication gates, a QAP is consisting of three sets of polynomials $\{l_i(X), r_i(X), o_i(X)\}_{i=0}^m$ and a $n_*$-degree target polynomial $t(X)$. The QAP holds and $(a_0, ..., a_{n_{in}}, a_{m-n_{ot}+1}, ..., a_m) \in \mathbb{F}_p^{n_{in}+n_{ot}}$ is a valid assignment for the input/output variables of the arithmetic circuit *iff* there is $(a_{n_{in}+1}, ..., a_{m-n_{ot}}) \in \mathbb{F}_p^{m-n_{in}-n_{ot}}$ such that $\sum_{i=0}^m a_i l_i(X)) \cdot \sum_{i=0}^m a_i r_i(X) \equiv \sum_{i=0}^m a_i o_i(X) \mod t(X)$. Herein, the size of the QAP is $m$ and the degree is $n_*$.

**zk-SNARK.** Groth16 [19] is a QAP-based zk-SNARK scheme, enabling proving the satisifiability of QAP via conducting a divisibility check between polynomials, which is equivalent to prove the satisifiability of the corresponding circuits. We figure out some major notations which can pave the way for presenting Groth16 and our design later. Concretely, we can specify the mentioned assignment $(a_0, ..., a_{n_{in}}, a_{m-n_{ot}+1}, ..., a_m) \in \mathbb{F}_p^{n_{in}+n_{ot}}$ as **statement** st to be proven, and specify $(a_{n_{in}+1}, ..., a_{m-n_{ot}}) \in \mathbb{F}_p^{m-n_{in}-n_{ot}}$ as **witness** wt only known by a prover. Then, we define the polynomial time computable binary **relation** R that comprises the pairs of (st, wt) satisfying the denoted QAP $\{\{l_i(X), r_i(X),$ $o_i(X)\}_{i \in [0,m]}, t(X)\}$. The degree of $\{l_i(X), r_i(X), o_i(X)\}$ is lower than that of $t(X)$. If R holds on (st, wt), R = 1; otherwise, R = 0. Formally, the relation is denoted as R = $\{(\mathsf{st}, \mathsf{wt}) \mid \mathsf{st} := (a_0, ..., a_{n_{in}}, a_{m-n_{ot}+1}, ..., a_m), \mathsf{wt} :=$ $(a_{n_{in}+1}, ..., a_{m-n_{ot}})$. Associated to the relation R, we denote that **language** L contains the statements that the corresponding witnesses exist in R, that is, L = $\{\mathsf{st} \mid \exists \mathsf{wt} \text{ s.t. } \mathsf{R}(\mathsf{st}, \mathsf{wt}) = 1\}$.

Based on the above notations, we are proceeding to review the definition of a zk-SNARK scheme for the relation R. Specifically, a zk-SNARK scheme is consisting of a quadruple of probabilistic polynomial time (PPT) algorithms (Setup, Prove, Verify, Sim), which satisfies three properties:

Setup$(1^\lambda, \mathsf{R}) \to (\mathsf{crs}, \mathsf{td})$: the Setup algorithm takes a security parameter $\lambda$ and a relation R as inputs, which returns a common reference string crs and a simulation trapdoor td.

Prove$(\mathsf{crs}, \mathsf{st}, \mathsf{wt}) \to \pi$: the Prove algorithm takes crs, statement st and witness wt as inputs, and outputs an argument $\pi$.

Verify$(\mathsf{R}, \mathsf{crs}, \mathsf{st}, \pi) \to 0/1$: the Verify algorithm takes the relation R, common reference string crs, statement st and argument $\pi$ as inputs, and returns 0 as Reject or 1 as Accept.

Sim$(\mathsf{R}, \mathsf{td}, \mathsf{st}) \to \pi'$: the Sim algorithm takes the relation R, simulation trapdoor td and statement st as inputs, and returns an argument $\pi'$.

· **Completeness.** Given a security parameter $\lambda$, $\forall(\mathsf{st}, \mathsf{wt}) \in \mathsf{R}$, a honest prover can convince a honest verifier the validity of a correctly generated proof with an overwhelming probability, that is, Pr$\{$Setup$(1^\lambda, \mathsf{R}) \to (\mathsf{crs}, \mathsf{td})$, Prove$(\mathsf{crs}, \mathsf{st}, \mathsf{wt}) \to \pi|$ Verify$(\mathsf{R}, \mathsf{crs}, \mathsf{st}, \pi) \to 1\} = 1 - \mathrm{negl}(1^\lambda)$.

· **Computational soundness.** For every computationally bounded adversary $\mathcal{A}$, if an invalid proof is successfully verified by a honest verifier, there exists a PPT extractor $\mathcal{E}$ who can extract wt with a non-negligible probability, that is, Pr$\{\exists(\mathsf{st}, \mathsf{wt}) \notin \mathsf{R}, \mathcal{A}(\mathsf{crs}, \mathsf{st}, \mathsf{wt}) \to \pi|$Verify$(\mathsf{R}, \mathsf{crs}, \mathsf{st}, \pi) \to 1$ $\wedge \mathcal{E}(\mathsf{crs}) \to \mathsf{wt}\} = 1 - \mathrm{negl}(1^\lambda)$.

· **Zero knowledge.** For every computationally bounded distinguisher $\mathcal{D}$, there exists the Sim algorithm such that $\mathcal{D}$ successfully distinguishes a honestly generated proof from a simulation proof with a negligible probability, that is, Pr$\{$Setup$(1^\lambda, \mathsf{R}) \to (\mathsf{crs}, \mathsf{td})$, Prove$(\mathsf{crs}, \mathsf{st}, \mathsf{wt}) \to \pi,$ Sim$(\mathsf{R}, \mathsf{td}, \mathsf{st}) \to \pi'|\mathcal{D}(\pi, \pi') = 1\} \leq \mathrm{negl}(1^\lambda)$.

**Commit-and-prove zk-SNARK.** We particularly emphasize the capability of CaP zk-SNARK that is useful for our scenario. The formal definition can be found in Definition 3.1 of [21]. Essentially, the CaP capability allows a prover to convince a verifier of "$C_{y^1}$ commits to $y^1$ such that $y^1 = F_1(x)$ and meanwhile $y^2 = F_2(y^1)$".

Intuitively, we suppose $C_m$ committed to an intermediate model $m$ of a computation and meanwhile the model $m$ is taken as input for completing subsequent computation, that is $m = F(x) \wedge y = G(m)$. We next explain how such capability matches the features of our scenario: (*a*) **Compression**. Before proving the correctness of a CNN testing, the CNN developer (resp. testers) store commitments to the CNN model (resp. test inputs). Herein, commitment is a lightweight method to compress large-size CNN models and test inputs for saving storage overhead. In terms of security, both the models and test inputs are sealed as well. (*b*) **Flexibility**. Committing to

the CNN model in advance provides the developer with certain flexibility that enables proving later-defined statements on the previously committed CNN model over unexpected test inputs. *(c)* ***Interoperability***. The proofs corresponding to the process of CNN testing are generated separately due to that a CNN model is partitioned into separate parts, so commitments to the CNN model will provide the interoperability between the separately generated proofs (see Fig. 6 and Fig. 7).

## IV. PROBLEM STATEMENT

### A. Collaborative Inference

We take the CNN model described in Section III-A for image classification as an example, and partition it into two parts. We choose the first pooling layer as a split point, such that the model can be partitioned into a part named PriorNet $F_{\mathbf{M}_1} = f^{pool_1}(f^{act_1}(f^{conv_1}(\cdot)))$ and another part named LaterNet $F_{\mathbf{M}_2} = f^o(f^{fc_4}(f^{fc_3}(f^{pool_2}(f^{act_2}(f^{conv_2}(\cdot))))))$. Suppose that the chosen split point is sufficiently optimized [28] or the model is pre-processed [29] for privacy protection, which makes the outputs of the PriorNet not reveal private information of its input data. The choice of the split point depends on the model architecture. A general principle is that the deeper layer a split point locates at, the less privacy will be leaked. Besides, we note that PriorNet $\mathbf{M}_1$ will be locally evaluated and LaterNet $\mathbf{M}_2$ will be delegated.
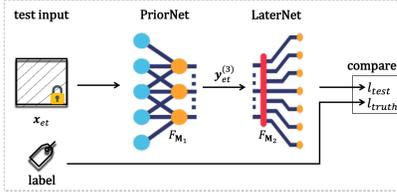


Fig. 2: Pipeline of PriorNet and LaterNet inference.

We now elaborate in Fig. 2 the pipeline of PriorNet and LaterNet inference over a given test input encrypted under a tester's public key of homomorphic encryption. We note that the corresponding label of the test input is not encrypted. Concretely, PriorNet is privately performed by the developer. It takes as input an encrypted test $\mathbf{X}_{et} \in \mathcal{R}_q$, and produces an FHE-encrypted intermediate result $\mathbf{Y}_{et}^{(3)} \in \mathcal{R}_q$ as the output. Since the homomorphic encryption cannot efficiently support the operations beyond additions and multiplications [30], [39], the activation layer within the PriorNet is approximated by a 2-degree polynomial function.

Successively, $\mathbf{Y}_{et}^{(3)}$'s plaintext $\mathbf{Y}^{(3)}$ will be fed into the LaterNet $F_{\mathbf{M}_2}$ whose computations are delegated to a service provider. We can leverage the re-encryption technique to let the service provider obtain the plaintext $\mathbf{Y}^{(3)}$. Since the computations are conducted on plaintext domain, the later activation functions remain unchanged. Finally, LaterNet outputs the prediction result $l_{\text{test}}$. The service provider proceeds to compare the previously given true label $l_{\text{truth}}$ with $l_{\text{test}}$, and returns an equality indicator $0$ or $1$. For a batch of test data, he would count the number of $1$ indicator, and return a proportion of correct prediction. For the case of using a posterior probability vector as a prediction output, we determine the indicator value according to the distance between the prediction output and the corresponding truth.

### B. Basic Workflow

We consider four main entities—*public platform (PP)*, *model tester (MT)*, *model developer (MD)* and *service provider (SP)* as demonstrated in Fig. 3.
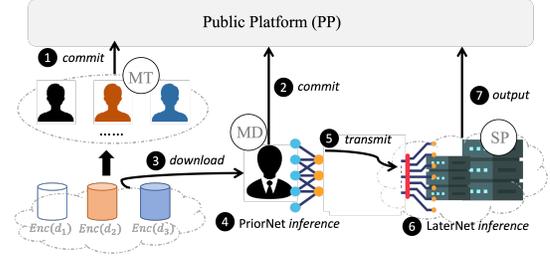


Fig. 3: Scenario overview of pvCNN

Specifically, the *PP* in real world may refer to a publicly accessible and online platform, which can be maintained by a set of computational nodes. On the platform, an *MD* is allowed to advertise his/her pre-trained CNN model by announcing some usage descriptions and model performance, so as to attract users of interest for earning profits, while an *MT* can provide test data for measuring the model performance. In order to establish trust among the *MD*, the *MT* and future users in the scenario, the pvCNN wants to enable verifiable announcements on the public platform, that is, any future user is assured of the truthfulness of model performance without learning any information of either private model parameters or test data. We note the *MD* can delegate computation-consuming tasks to the powerful *SP*.

The basic workflow around the entities is shown below: ❶ *MTs* freely participate in the platform and commit to the test data. They also store encrypted test data by using a leveled FHE (L-FHE) scheme [50] on the publicly accessible server. ❷ The *MD* splits his CNN model into a private PriorNet which is kept at local device, as well as, a public LaterNet which is sent to the *SP*; After that, he generates two commitments to the PriorNet and LaterNet, and sends them to the *PP*. ❸ The *MD* downloads the ciphered test data. ❹ The *MD* evaluates the plaintext PriorNet on the encrypted test data and generates a proof of executing the inference computation as promised. ❺ The *MD* sends the encrypted output of the Step ❹. Note that the encrypted output can be transformed into the ciphertext under the *SP*'s public key via re-encryption. ❻ The *SP* successively runs LaterNet inference in the clear by taking the plaintext output of PriorNet inference as inputs, and proves that LaterNet is executed correctly with the true inputs and meanwhile LaterNet is consistent to the committed one. ❼ The *SP* submits the final test results (*e.g.,* classification correctness rate) to the platform, and distributes proofs to the decentralized nodes for proofs aggregation and validation.

### C. Threat Assumptions

**Model Developers.** We consider CNN developers are untrusted; they can provide untruthful prediction results, *e.g.,* outputting meaningless predictions or using arbitrary test data or CNN model. They can also be curious about test data during a CNN testing. Note that our verifiable computing design does not focus on poisoned or backdoored models [51], but it

greatly relies on multi-tester test data to probe the performance of CNN models in a black-box manner. Recently, SecureDL [52] and VerIDeep [53] resort to sensitive samples to detect model changes, which can be complementary to our work.

**Service Providers.** They can be computationally powerful entities, and accept outsourced LaterNet from developers. Despite service providers' powerful capabilities on computation and storage, we do not trust them. They can try their best to steal private information with regard to PriorNet and test inputs. They even incorrectly run LaterNet and produce untruthful results due to machine disruption. Service providers can also collude with a developer to fake correctness proofs, with respect to the independent inference processes based on a splitting CNN, aiming to evade verification by later-coming users. But the service providers cannot obtain PriorNet by colluding with the developer, since the developer has no motivation. Similar to previous studies [15], [16], we do not account for hyper-parameter extraction attacks through side-channel analysis [54], [55] from service providers in our research. Instead, we assume hyper-parameters to be publicly accessible in our setting. Moreover, there are defensive methods available [56] to mitigate such side-channel analysis attacks .

**Model Testers.** Testers are organized via a crowdsourcing manner. We consider that the testers are willing to contribute their test inputs for measuring the quality of the CNNs. A tester's test inputs concretely contain a batch of data on an input-label basis. We assume the original inputs need protected but their labels can be public (refer to encrypted inputs and plaintext labels in Section IV-A). We also assume each label is consistent with the corresponding input's true label, considering currently popular crowdsourcing-empowered label platforms.

**Public Platform and Future Users.** We assume both the public platform and future users are honest. Users have access to all of data (not including model parameters or test data) uploaded to the platform and the data cannot be tampered. Finally, we assume that all entities communicate via a secure authenticated channel.

## V. CONCRETE DESIGN

This section will introduce how to generate publicly verifiable zk-SNARK proofs for CNN testing based on L-FHE and collaborative inference. We will present high-level ideas in the subsection V-A. We then illustrate the optimization for 2-D convolutional relations in the subsection V-B. In the subsequent subsection V-C, we generate zk-SNARK proofs for a whole CNN testing process in a divide-and-conquer fashion, and aggregate multiple proofs in the subsection V-D.

### A. High-level Ideas

On top of the scenario of testing PriorNet and LaterNet with encrypted test data from testers, we want to generate zk-SNARK proofs for convincing any future user that (1) the model developer correctly runs PriorNet on the ciphered test data and returns the ciphered intermediate result as output; and (2) the service provider exactly takes the intermediate result as

input and correctly accomplishes the successive computations of LaterNet, which produces a correct result. Meanwhile, any future user does not learn any information about the model parameters, the intermediate result and the test data.
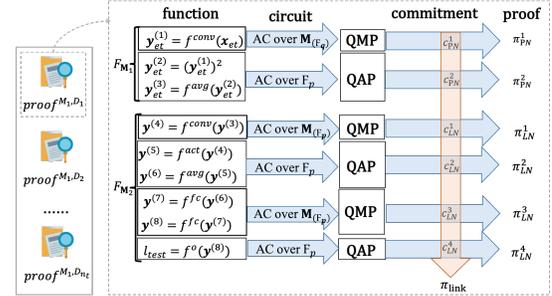


Fig. 4: Overview of proof generation. See the column of function, we use a 2-degree polynomial to approximate the ReLU function in the PriorNet $F_{\mathbf{M}_1}$.

Most importantly, we make two efforts to improve proving and verification costs, by reducing the number of multiplication gates and compressing multiple proofs, respectively. For ease of understanding, we firstly introduce a roadmap of proof generation, and then particularly emphasize our two efforts during proof generation. As shown in Fig. 4, the roadmap of generating QAP-based zkSNARK proofs is: a) compiling the *functions* of the splitting CNN testing into circuits, b) expressing the *circuits* as a set of QAPs, c) constructing the zk-SNARK proofs for the QAPs, and d) generating *commitments* for linking together the zk-SNARK proofs which are constructed separately. Based on the roadmap, our main efforts include: **(A) Optimizing convolutional relations for b) (see Section V-B).** We allow matrices in constructing the QAP, leading to quadratic matrix programs (QMP) [31], which is inspired by vCNN [12]'s efforts of applying polynomials into QAP (QPPs, quadratic polynomial programs [57]). Our QMP can optimize convolutional relations, considering a convolution layer with multiple filters and inputs can be strategically represented by a single matrix multiplication operation. The operation of matrix multiplication then is complied into the arithmetic circuits in QMP with only a single multiplication gate. As a consequence, the number of multiplication gates dominating proving costs is greatly reduced, compared to the previous case of using the original expression of QAP. **(B) Aggregating proofs for c) (see Section V-D).** We enable aggregating multiple proofs for different statements respective to test datasets over the same QAP/QMP for the same CNN model. The single proof after aggregation then is validated by a secure committee and the validity result based on the majority voting will be submitted and recorded on our public platform. In such way, any future user only needs to check the validation result to determine whether or not the statements over the CNN are valid.

### B. Optimizing Convolutional Relations

As demonstrated in Fig. 5, we consider 2-D convolution operations between $Mn^2$ amount of $n \times n$ input matrices and $M$ amount of $m \times m$ convolution filters (herein, $M > n$), which produces $M^2n^2$ filtered matrices of
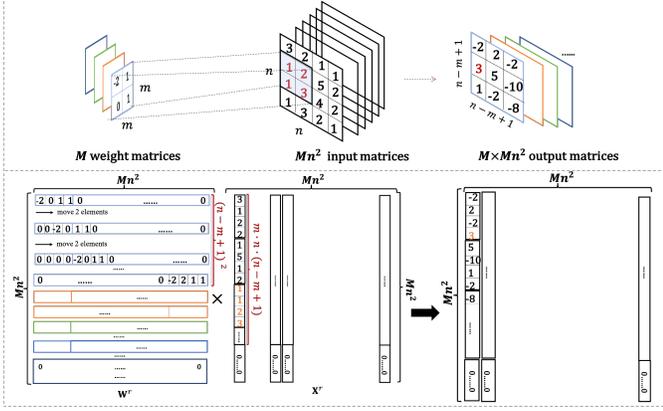
Fig. 5: Matrix multiplication for convolution operations. We transform the convolution operations of $M$ weight matrices and $Mn^2$ input matrices (demonstrated on the top) into the matrix multiplication operation between two square matrices of $Mn^2 \times Mn^2$ (demonstrated on the bottom).

$(n - m + 1) \times (n - m + 1)$, by setting the stride 1. Specifically, by computing the inner products between each input matrix and each filter matrix, $\forall i, j \in [0, n - m]$ and $k \in [0, M - 1]$, it can be expressed as: $\mathbf{Y}^k[i][j] = \sum_{l_r=0, l_c=0}^{m-1, m-1} \mathbf{X}[i + l_r][j + l_c] \times \mathbf{W}^k[l_r][l_c]$. Obviously, the multiplication complexity of such inner product computation is $O(n^2 \cdot m^2)$. Note that there are $M \cdot Mn^2$ amount of such set of $\mathbf{Y}$, with respect to $Mn^2$ inputs and $M$ filters. Thus, $O(M \cdot Mn^2 \cdot (n^2 \cdot m^2))$ multiplication gates in QAP-based circuits are needed for handling $Mn^2$ inputs and $M$ filters. Since the proving time depends on the number of multiplication gates, a proof generation process for 2-D convolution operations becomes impractical.

Now, we seek for improving proving efficiency by reducing the multiplication complexity to $O(Mn^2 \cdot Mn^2)$. The very first step is to represent the inner products between the $Mn^2$ inputs and $M$ filters with *a single matrix multiplication*. Concretely, we show a concrete example of the representation in Fig. 5. First, given $Mn^2$ amount of $n \times n$ input matrices $\{\mathbf{X}\}$, they are reshaped into a square matrix $\mathbf{X}^r$ of $Mn^2 \times Mn^2$. Second, another square matrix $\mathbf{W}^r$ of the same dimension can be constructed by packing the $M$ filter matrices in a moving and zero padding manner. As a result, the multiplication between the two square matrices produces a square matrix $\mathbf{Y}^r$ of $Mn^2 \times Mn^2$, which can also be regarded as a reshaped matrix from the original $M \cdot Mn^2$ amount of output matrices. Based on such representation, the two square matrices $\mathbf{W}^r$ and $\mathbf{X}^r$ of the matrix multiplication can later be the left and right inputs of an arithmetic circuit with only one multiplication gate. We proceed to express the above circuit with our new QMP by using matrices in QAP, which follows the similar principle of using polynomials in QAP [57].

**QMP Definition [31].** For an arithmetic circuit with $n_{in}$ input variables, $n_{ot}$ output variables in $\mathbf{M}^{s \times s}_{(\mathbb{F}_p)}$ and $n_*$ multiplication gates, a QMP is consisting of three sets of polynomials $\{L_i(x), R_i(x), O_i(x)\}_{i=0}^{m}$ with coefficients in $\mathbf{M}^{s \times s}_{(\mathbb{F}_p)}$ and a $n_*$-degree target polynomial $t(x) \in \mathbb{F}_p[x]$. $\mathbb{F}_p$ is a large finite field, *e.g.*, $2^{254}$. The QMP computes the arithmetic

circuit and $(\mathbf{A}_0, ..., \mathbf{A}_{n_{in}}, \mathbf{A}_{m-n_{ot}+1}, ..., \mathbf{A}_m) \in \mathbf{M}^{s \times s}_{(\mathbb{F}_p)}$ is valid assignment for the $n_{in} + n_{ot}$ input/output variables *iff* there exists $m - n_{in} - n_{ot}$ coefficients $(\mathbf{A}_{n_{in}+1}, ..., \mathbf{A}_{m-n_{ot}}) \in \mathbf{M}^{s \times s}_{(\mathbb{F}_p)}$ for arbitrary $\mathbf{X} \in \mathbf{M}^{s \times s}_{(\mathbb{F}_p)}$ such that $t(x)$ divides $p(x, \mathbf{X}) = \text{tr}\{\mathbf{X}^T \sum_{i=0}^m \mathbf{A}_i L_i(x)) \cdot \sum_{i=0}^m \mathbf{A}_i R_i(x)\} - \text{tr}\{\mathbf{X}^T \sum_{i=0}^m \mathbf{A}_i O_i(x)\}$. Herein, $\text{tr}$ means the trace of a square matrix, and the degree of the QMP is $deg(t(x))$. We note that each wire of the arithmetic circuit is labeled by a square matrix $\mathbf{A}_i \in \mathbf{M}^{s \times s}_{(\mathbb{F}_p)}$. Suppose that the arithmetic circuit contains $n_*$ multiplication gates, and then $t(x) = \prod_{r=0}^{n_*-1} t(x - x_r)$. With regard to a multiplication gate $x_r$, its left input wire is labeled by $\sum_{i=0}^m \mathbf{A}_i L_i(x_r))$, right input wire is labeled by $\sum_{i=0}^m \mathbf{A}_i R_i(x_r))$, and the output wire is $\sum_{i=0}^m \mathbf{A}_i O_i(x_r))$. We note that the multiplication gate is constrained by $\text{tr}\{\mathbf{X}^T \sum_{i=0}^m \mathbf{A}_i L_i(x_r)) \cdot \sum_{i=0}^m \mathbf{A}_i R_i(x_r)\} = \text{tr}\{\mathbf{X}^T \sum_{i=0}^m \mathbf{A}_i O_i(x_r)\}$.

**QMP for Matrix Multiplication Relations.** As we introduced before, we can represent convolutional operations as matrix multiplication $\mathbf{Y}^r = \mathbf{W}^r \times \mathbf{X}^r$, where $\mathbf{Y}^r, \mathbf{X}^r$ and $\mathbf{W}^r \in \mathbf{M}^{Mn^2 \times Mn^2}_{(\mathbb{Z}_p)}$. It is noteworthy that we can use advanced quantization techniques to transform floating-point numbersm of 8-bit precision into unsigned integer numbers in $[0, 255]$, which adapts to a large finite field. Now according to the QMP definition, when the matrix multiplication computation is expressed by the arithmetic circuit over matrices, the corresponding QMP is $\{L_{\mathbf{W}}(x), R_{\mathbf{X}}(x), O_{\mathbf{Y}}(x), t(x)\}$ of the arithmetic circuit. Herein, there are three input/output variables, and the degree of the QMP is one.

Similarly, we can apply our above idea into full connection operations, and finally the operations can also be expressed with similarly low-degree QMP.

### C. Proof Generation for a Whole CNN Testing

This section elaborates how to generate proofs based on the CaP zk-SNARK for ensuring the result correctness of executing PriorNet and LaterNet. The proofs convince an arbitrary user that PriorNet is correctly evaluated on the ciphered test inputs, and subsequently, LaterNet is correctly performed by taking the clear outputs of PriorNet as inputs, yielding correct results. We note that the processes of proof generation for the PriorNet and LaterNet are slightly different. For PriorNet, the values, such as weights, test inputs and outputs, are represented as polynomial ring elements, and we follow Fiore *et al.*'s work [22] to generate proofs for the computation over polynomial rings. For LaterNet, we generate zk-SNARK proofs for the computation over scalars. We will begin with relation definition, and then present concrete realization for generating proofs. In light of the nature of layer-by-layer computation of model inference, we will use a divide-and-conquer method to generate zk-SNARK proofs for each separate layer, and combine them as a whole via a commit-and-prove methodology.

Let us define some notations at the beginning. We let MPoly.Com be a polynomial commitment scheme from [22]. We denote $ck_{\text{PN}}$ and $ck_{\text{LN}}$ are the commitment keys used for committing to the PriorNet and LaterNet; $\mathbf{r}_{\text{PN}}, \mathbf{r}_{\text{LN}}$ are non-reused commitment randomnesses. $C_{\mathbf{M}_{1,ed}}$ means the commit-

---

**Part 1: Layer-wise prediction correctness for PriorNet.**

Step 1: to prove that the prover knows the openings of the commitments in $\mathsf{st}_{\mathsf{PN}}^{S1}$

$R_{\mathsf{PN}}^{S1} := \{ (\mathsf{st}_{\mathsf{PN}}^{S1} = (C_{\mathsf{L}_1}, C_{\mathsf{X}}, C_{\mathsf{PN}}^{inter}, C_{\mathsf{L}_1,k}, C_{\mathsf{L}_2,k}, C_{\mathsf{L}_3,k}, k, ck^{S1}, ck_{\mathsf{PN}}), \mathsf{wt}_{\mathsf{PN}}^{S1} = (\mathbf{W}_{ed}^{r(1)}, \mathbf{X}_{et}^{r(1)}, \mathbf{Y}_{et}^{r(1)}, \mathbf{Y}_{et}^{r(2)}, \mathbf{Y}_{et}^{r(3)}, \underline{\mathbf{W}_{ed,k}^{r(1)}, \mathbf{X}_{et,k}^{r(1)}, \mathbf{Y}_{et,k}^{r(1)}, \mathbf{Y}_{et,k}^{r(2)}, \mathbf{Y}_{et,k}^{r(3)}}, \mathbf{r}_{\mathsf{PN}}) \mid$

$\mathbf{W}_{ed,k}^{r(1)} = \mathbf{W}_{ed}^{r(1)}(k) \wedge \mathbf{X}_{et,k}^{r(1)} = \mathbf{X}_{et}^{r(1)}(k) \wedge \mathbf{Y}_{et,k}^{r(1)} = \mathbf{Y}_{et}^{r(1)}(k) \wedge \mathbf{Y}_{et,k}^{r(2)} = \mathbf{Y}_{et}^{r(2)}(k) \wedge \mathbf{Y}_{et,k}^{r(3)} = \mathbf{Y}_{et}^{r(3)}(k)$

$\wedge C_{\mathsf{L}_1} = \mathsf{MPoly.Com}(\mathbf{W}_{ed}^{r(1)}, ck^{S1}, \mathbf{r}_{\mathsf{PN}}) \wedge C_{\mathsf{X}} = \mathsf{MPoly.Com}(\mathbf{X}_{et}^{r(1)}, ck^{S1}, \mathbf{r}_{\mathsf{PN}}) \wedge C_{\mathsf{PN}}^{inter} = \mathsf{MPoly.Com}(\mathbf{Y}_{et}^{r(1)}, \mathbf{Y}_{et}^{r(2)}, \mathbf{Y}_{et}^{r(3)}, ck_{\mathsf{PN}}, \mathbf{r}_{\mathsf{PN}})$

$\wedge C_{\mathsf{L}_1,k} = \mathsf{MPoly.Com}(\mathbf{W}_{ed,k}^{r(1)}, \mathbf{X}_{et,k}^{r(1)}, \mathbf{Y}_{et,k}^{r(1)}, ck_{\mathsf{PN}}, \mathbf{r}_{\mathsf{PN}}) \wedge \underline{C_{\mathsf{L}_2,k} = \mathsf{MPoly.Com}(\mathbf{Y}_{et,k}^{r(2)}, ck_{\mathsf{PN}}, \mathbf{r}_{\mathsf{PN}})} \wedge \underline{C_{\mathsf{L}_3,k} = \mathsf{MPoly.Com}(\mathbf{Y}_{et,k}^{r(3)}, ck_{\mathsf{PN}}, \mathbf{r}_{\mathsf{PN}})} \}.$

*//$ck_{\mathsf{PN}}$ is the commitment key used before the evaluation while $ck^{S1}$ is used during the prediction*
*//$C_{\mathsf{X}}$ commits to $\mathbf{X}_{et}^{r(1)}$; $C_{\mathsf{L}_1}$ commits to $\mathbf{W}_{ed}^{r(1)}$*

Step 2: to prove that the equations with respect to $\mathsf{L}_1, \mathsf{L}_2, \mathsf{L}_3$ hold

$R_{\mathsf{PN}}^{S2}[\mathsf{L}_1] := \{ (\mathsf{st}_{\mathsf{PN}}^{S2} = (C_{\mathsf{L}_1,k}', k, ck_{\mathsf{L}_1}), \mathsf{wt}_{\mathsf{PN}}^{S2} = (\mathbf{W}_{ed,k}^{r(1)}, \mathbf{X}_{et,k}^{r(1)}, \mathbf{Y}_{et,k}^{r(1)}, \mathbf{r}_{\mathsf{PN}})) \mid \mathbf{Y}_{et,k}^{r(1)} = \mathbf{W}_{ed,k}^{r(1)} \odot \mathbf{X}_{et,k}^{r(1)} \wedge C_{\mathsf{L}_1,k}' = \mathsf{MPoly.Com}(\mathbf{W}_{ed,k}^{r(1)}), \mathbf{X}_{et,k}^{r(1)}, \mathbf{Y}_{et,k}^{r(1)}, ck_{\mathsf{L}_1}, \mathbf{r}_{\mathsf{PN}}) \}.$

*//$ck_{\mathsf{L}_1}$ is built depending on the relation; $C_{\mathsf{L}_1,k}'$ commits to the consistent $\mathbf{W}_{ed,k}^{r(1)}, \mathbf{X}_{et,k}^{r(1)}, \mathbf{Y}_{et,k}^{r(1)}$ with the above $C_{\mathsf{L}_1,k}$*

$R_{\mathsf{PN}}^{S2}[\mathsf{L}_2, \mathsf{L}_3] := \{ (\mathsf{st}_{\mathsf{PN}}^{S2} = (C_{\mathsf{L}_2,k}', C_{\mathsf{L}_2}[\mathbf{Y}_{et,k}^{r(1)}], C_{\mathsf{L}_3,k}', C_{\mathsf{L}_3}[\mathbf{Y}_{et,k}^{r(2)}], k, ck_{\mathsf{L}_2}, ck_{\mathsf{L}_3}), \mathsf{wt}_{\mathsf{PN}}^{S2} = (\mathbf{Y}_{et,k}^{r(2)}, \mathbf{Y}_{et,k}^{r(3)}, \mathbf{r}_{\mathsf{PN}}) \mid$

$\{ \mathbf{Y}_{et,k}^{r(2)}[i][j] = (\mathbf{Y}_{et,k}^{r(1)}[i][j])^2 + \mathbf{Y}_{et,k}^{r(1)}[i][j] \}_{i=0,j=0}^{Mn^2-1, Mn^2-1} \wedge \mathbf{Y}_{et,k}^{r(3)} = f^{avg}(\mathbf{Y}_{et,k}^{r(2)}) \wedge C_{\mathsf{L}_2,k}' = \mathsf{MPoly.Com}(\mathbf{Y}_{et,k}^{r(2)}, ck_{\mathsf{L}_2}, \mathbf{r}_{\mathsf{PN}})$

$\wedge C_{\mathsf{L}_2}[\mathbf{Y}_{et,k}^{r(1)}] = \mathsf{MPoly.Com}(\mathbf{Y}_{et,k}^{r(1)}, ck_{\mathsf{L}_2}, \mathbf{r}_{\mathsf{PN}}) \wedge C_{\mathsf{L}_3,k}' = \mathsf{MPoly.Com}(\mathbf{Y}_{et,k}^{r(3)}, ck_{\mathsf{L}_3}, \mathbf{r}_{\mathsf{PN}}) \wedge C_{\mathsf{L}_3}[\mathbf{Y}_{et,k}^{r(2)}] = \mathsf{MPoly.Com}(\mathbf{Y}_{et,k}^{r(2)}, ck_{\mathsf{L}_3}, \mathbf{r}_{\mathsf{PN}}) \}.$

*//$C_{\mathsf{L}_2,k}'$ (and $C_{\mathsf{L}_3,k}'$) commit to the same $\mathbf{Y}_{et,k}^{r(2)}$ (resp. $\mathbf{Y}_{et,k}^{r(3)}$) with $C_{\mathsf{L}_2,k}$ (resp. $C_{\mathsf{L}_3,k}$)*
*//$C_{\mathsf{L}_2}[\mathbf{Y}_{et,k}^{r(1)}]$ should commit to the same $\mathbf{Y}_{et,k}^{r(1)}$ that $C_{\mathsf{L}_1,k}'$ commits to; similarly, $C_{\mathsf{L}_3}[\mathbf{Y}_{et,k}^{r(2)}]$ should commit to the same $\mathbf{Y}_{et,k}^{r(2)}$ that $C_{\mathsf{L}_2,k}'$ commits to*

---

Fig. 6: Layer-by-layer relation definition for PriorNet.

---

**Part 2: Layer-wise prediction correctness for LaterNet.**

$R_{\mathsf{LN}}[\mathsf{L}_4] := \{ (\mathsf{st}_{\mathsf{LN}} = (C_{\mathbf{Y}^{r(3)}}, C_{\mathsf{L}_4}, C_{\mathbf{Y}^{r(4)}}, ck_{\mathsf{L}_4}), \mathsf{wt}_{\mathsf{LN}} = (\mathbf{W}^{r(4)}, \mathbf{Y}^{r(3)}, \mathbf{Y}^{r(4)}, \mathbf{r}_{\mathsf{LN}})) \mid \mathbf{Y}^{r(4)} = \mathbf{W}^{r(4)} \times \mathbf{Y}^{r(3)}$

$\wedge C_{\mathbf{Y}^{r(3)}} = \mathsf{MPoly.Com}(\mathbf{Y}^{r(3)}, ck_{\mathsf{L}_4}, \mathbf{r}_{\mathsf{LN}}) \wedge C_{\mathsf{L}_4} = \mathsf{MPoly.Com}(\mathbf{W}^{r(4)}, ck_{\mathsf{L}_4}, \mathbf{r}_{\mathsf{LN}}) \wedge C_{\mathbf{Y}^{r(4)}} = \mathsf{MPoly.Com}(\mathbf{Y}^{r(4)}, ck_{\mathsf{L}_4}, \mathbf{r}_{\mathsf{LN}}) \}.$

*//$C_{\mathbf{Y}^{r(3)}}$ should commit to $\mathbf{Y}^{r(3)}$ whose encryption is committed in $C_{\mathsf{PN}}^{inter}$*

$R_{\mathsf{LN}}[\mathsf{L}_5, \mathsf{L}_6] := \{ (\mathsf{st}_{\mathsf{LN}} = (C_{\mathsf{L}_5}[\mathbf{Y}^{r(4)}], C_{\mathbf{Y}^{r(5)}}, C_{\mathsf{L}_6}[\mathbf{Y}^{r(5)}], C_{\mathbf{Y}^{r(6)}}, ck_{\mathsf{L}_5}, ck_{\mathsf{L}_6}), \mathsf{wt}_{\mathsf{LN}} = (\mathbf{Y}^{r(4)}, \mathbf{Y}^{r(5)}, \mathbf{Y}^{r(6)}, \mathbf{r}_{\mathsf{LN}})) \mid$

$\{ \mathbf{Y}^{r(5)}[i][j] = max(\mathbf{Y}^{r(4)}[i][j], 0) \}_{i=0,j=0}^{Mn^2-1, Mn^2-1} \wedge \mathbf{Y}^{r(6)} = f^{avg}(\mathbf{Y}^{r(5)})$

$\wedge C_{\mathsf{L}_5}[\mathbf{Y}^{r(4)}] = \mathsf{MPoly.Com}(\mathbf{Y}^{r(4)}, ck_{\mathsf{L}_5}, \mathbf{r}_{\mathsf{LN}}) \wedge C_{\mathbf{Y}^{r(5)}} = \mathsf{MPoly.Com}(\mathbf{Y}^{r(5)}, ck_{\mathsf{L}_6}, \mathbf{r}_{\mathsf{LN}})$ *//$C_{\mathsf{L}_5}[\mathbf{Y}^{r(4)}]$ should commit to $\mathbf{Y}^{r(4)}$ that is committed in $C_{\mathbf{Y}^{r(4)}}$ of $\mathsf{L}_4$*

$\wedge C_{\mathsf{L}_6}[\mathbf{Y}^{r(5)}] = \mathsf{MPoly.Com}(\mathbf{Y}^{r(5)}, ck_{\mathsf{L}_6}, \mathbf{r}_{\mathsf{LN}}) \wedge C_{\mathbf{Y}^{r(6)}} = \mathsf{MPoly.Com}(\mathbf{Y}^{r(6)}, ck_{\mathsf{L}_6}, \mathbf{r}_{\mathsf{LN}}) \}.$ *//$C_{\mathsf{L}_6}[\mathbf{Y}^{r(5)}]$ should commit to $\mathbf{Y}^{r(5)}$ that is committed in $C_{\mathbf{Y}^{r(5)}}$ of $\mathsf{L}_5$*

$R_{\mathsf{LN}}[\mathsf{L}_7, \mathsf{L}_8] := \{ (\mathsf{st}_{\mathsf{LN}} = (C_{\mathsf{L}_{7,8}}, C_{\mathsf{L}_7}[\mathbf{Y}^{r(6)}], C_{\mathsf{L}_8}[\mathbf{Y}^{r(8)}], ck_{\mathsf{L}_{7,8}}), \mathsf{wt}_{\mathsf{LN}} = (\mathbf{W}^{r(7)}, \mathbf{W}^{r(8)}, \mathbf{Y}^{r(6)}, \mathbf{Y}^{r(7)}, \mathbf{Y}^{r(8)}, \mathbf{r}_{\mathsf{LN}})) \mid \mathbf{Y}^{r(7)} = \mathbf{Y}^{r(6)} \times \mathbf{W}^{r(7)}$

$\wedge C_{\mathsf{L}_{7,8}} = \mathsf{MPoly.Com}(\mathbf{W}^{r(7)}, \mathbf{W}^{r(8)}, ck_{\mathsf{L}_{7,8}}, \mathbf{r}_{\mathsf{LN}}) \wedge C_{\mathsf{L}_7}[\mathbf{Y}^{r(6)}] = \mathsf{MPoly.Com}(\mathbf{Y}^{r(6)}, ck_{\mathsf{L}_{7,8}}, \mathbf{r}_{\mathsf{LN}}) \wedge C_{\mathsf{L}_8}[\mathbf{Y}^{r(8)}] = \mathsf{MPoly.Com}(\mathbf{Y}^{r(8)}, ck_{\mathsf{L}_{7,8}}, \mathbf{r}_{\mathsf{LN}}) \}.$

*//$C_{\mathsf{L}_7}[\mathbf{Y}^{r(6)}]$ should commit to $\mathbf{Y}^{r(6)}$ that is committed in $C_{\mathbf{Y}^{r(6)}}$ of $\mathsf{L}_6$*

$R_{\mathsf{LN}}[\mathsf{L}_o] := \{ (\mathsf{st}_{\mathsf{LN}} = (C_{\mathsf{L}_o}[\mathbf{Y}^{r(8)}], \mathbf{l}_{test}, ck_{\mathsf{L}_o}), \mathsf{wt}_{\mathsf{LN}} = (\mathbf{Y}^{r(8)})) \mid \mathbf{l}_{test}[j] = \mathsf{argmax}(\mathbf{Y}^{r(8)}[j]) \wedge C_{\mathsf{L}_o}[\mathbf{Y}^{r(8)}] = \mathsf{MPoly.Com}(\mathbf{Y}^{r(8)}, ck_{\mathsf{L}_o}, \mathbf{r}_{\mathsf{LN}}) \}.$

*//$C_{\mathsf{L}_o}[\mathbf{Y}^{r(8)}]$ should commit to $\mathbf{Y}^{r(8)}$ that is committed in $C_{\mathbf{Y}^{r(8)}}$ of $\mathsf{L}_8$*

---

Fig. 7: Layer-by-layer relation definition for LaterNet.

---

ments to the inputs, weights, intermediates of $\mathbf{M}_1$ while $C_{\mathbf{M}_2}$ represents the commitments to the weights and intermediates of $\mathbf{M}_2$.

**Relation Definition.** We define the relations $R_{\mathsf{PN}}$ and $R_{\mathsf{LN}}$ with regard to PriorNet and LaterNet execution correctness:

$R_{\mathsf{PN}} := \{ (\mathsf{st}_{\mathsf{PN}} = (C_{\mathbf{M}_{1,ed}}, C_{\mathbf{X}_{et}^{r(1)}}, \mathbf{Y}_{et}^{r(3)}, ck_{\mathsf{PN}}),$

$\mathsf{wt}_{\mathsf{PN}} = (\mathbf{M}_{1,ed}, \mathbf{X}_{et}, \mathbf{r}_{\mathsf{PN}})) \mid \mathbf{Y}_{et}^{r(1)} = \mathbf{W}_{ed}^{r(1)} \odot \mathbf{X}_{et}^{r(1)}$

$\wedge \{ \mathbf{Y}_{et}^{r(2)}[i][j] = (\mathbf{Y}_{et}^{r(1)}[i][j])^2 + \mathbf{Y}_{et}^{r(1)}[i][j] \}_{i=0,j=0}^{Mn^2-1, Mn^2-1}$

$\wedge \mathbf{Y}_{et}^{r(3)} = f^{avg}(\mathbf{Y}_{et}^{r(2)})$

$\wedge C_{\mathbf{M}_{1,ed}} = \mathsf{MPoly.Com}(\mathbf{M}_{1,ed}, ck_{\mathsf{PN}}, \mathbf{r}_{\mathsf{PN}})$

$\wedge C_{\mathbf{X}_{et}} = \mathsf{MPoly.Com}(\mathbf{X}_{et}^{r(1)}, ck_{\mathsf{PN}}, \mathbf{r}_{\mathsf{PN}}) \},$

$R_{\mathsf{LN}} := \{ (\mathsf{st}_{\mathsf{LN}} = (C_{\mathbf{Y}^{r(3)}}, C_{\mathbf{M}_2}, \mathbf{l}_{test}, ck_{\mathsf{LN}}),$

$\mathsf{wt}_{\mathsf{LN}} = (\mathbf{Y}^{r(3)}, \mathbf{M}_2, \mathbf{r}_{\mathsf{LN}})) \mid \mathbf{Y}^{r(4)} = \mathbf{W}^{r(4)} \times \mathbf{Y}^{r(3)}$

$\wedge \{ \mathbf{Y}^{r(5)}[i][j] = max(\mathbf{Y}^{r(4)}[i][j], 0) \}_{i=0,j=0}^{Mn^2-1, Mn^2-1}$

$\wedge \mathbf{Y}^{r(6)} = f^{avg}(\mathbf{Y}^{r(5)})$

$\wedge \mathbf{Y}^{r(7)} = \mathbf{Y}^{r(6)} \times \mathbf{W}^{r(7)}$

$\wedge \mathbf{Y}^{r(8)} = \mathbf{Y}^{r(7)} \times \mathbf{W}^{r(8)}$

$\wedge \mathbf{l}_{test}[j] = \mathsf{argmax}(\mathbf{Y}^{r(8)}[j])$

$\wedge C_{\mathbf{Y}^{r(3)}} = \mathsf{MPoly.Com}(\mathbf{Y}^{r(3)}, ck_{\mathsf{LN}}, \mathbf{r}_{\mathsf{LN}})$

$\wedge C_{\mathbf{M}_2} = \mathsf{MPoly.Com}(\mathbf{M}_2, ck_{\mathsf{LN}}, \mathbf{r}_{\mathsf{LN}}) \}.$

Besides the relations, we further define layer-by-layer relations for proving the layer-wise computations of the $R_{\mathsf{PN}}$ and $R_{\mathsf{LN}}$ in Fig. 6 and 7. In Part 1, the layer-wise computations of the PriorNet performs in an encryption manner, so we define the relations to be proven with Setp 1 and Setp 2, as guided by the work [22]. In Part 2, the LaterrNet runs in a plaintext setting, so we define the corresponding relations similar to the above Setp 2 (not needing Setp 1). Lastly, in order to combine the layer-wise computations as a whole, we also denote the relations for ensuring that starting from the second layer, each layer's inputs are consistent with the former layer's outputs.

With the relations, we proceed to generate zk-SNARK proofs accordingly. At a high level, the computations of the PriorNet and LaterNet are proved correct *iff* both of binary relations $R_{\mathsf{PN}}$ and $R_{\mathsf{LN}}$ return 1, which means $R_{\mathsf{PN}}$ and $R_{\mathsf{LN}}$ hold on the corresponding pairs of statement and witness $(\mathsf{st}, \mathsf{wt})$. For $R_{\mathsf{PN}}$, we generate publicly verifiable proofs that $\mathbf{Y}_{et}^{r(3)}$ is computed correctly in an encryption manner, with the trained parameters of PriorNet $\mathbf{M}_{1,ed}$ and ciphered test inputs $\mathbf{X}_{et}^{r(1)}$ through the sequential computations mentioned in the $R_{\mathsf{PN}}$. In the relation, the witnesses $\mathbf{M}_{1,ed}$ and $\mathbf{X}_{et}^{r(1)}$ are committed using the key $ck_{\mathsf{PN}}$, and the generated commitments plus $\mathbf{Y}_{et}^{r(3)}$ are treated as the public statement. For $R_{\mathsf{LN}}$, we generate proofs that LaterNet takes the correct plaintext $\mathbf{Y}^{r(3)}$ and returns the correct inference result $\mathbf{l}_{test}$ by conducting the sequential operations included in the $R_{\mathsf{LN}}$. Similarly, the witnesses $\mathbf{M}_2$ and $\mathbf{Y}^{r(3)}$ are committed using the key $ck_{\mathsf{LN}}$,

and then the commitments and the execution result are public statements. The processes of proof generation for $R_{PN}$ and $R_{LN}$ are conducted separately, since the two parts PriorNet and LaterrNet are executed by different entities. After completing the proofs for the two parts, we link the proofs together by additionally proving that the output of the PriorNet is exactly equal to the input of the LaterNet.

---

**Step 1 CaP zk-SNARK Proofs of Knowledge of Commitments**

1: $\mathsf{MUniEv\text{-}\Pi.Setup}(\lambda, d_c, n_c) \to ck^{S1}$ :
2: $//d_c$ refers to the degree of ciphers; $n_c$ is the number of ciphers
3: $g, h \xleftarrow{\$} \mathsf{G}, g^* \xleftarrow{\$} \mathcal{G}, \alpha, s, t \xleftarrow{\$} \mathbb{Z}_p$
4: $\hat{g} = g^\alpha, \hat{h} = h^\alpha, \hat{g^*} = g^{*\alpha}, g_1^* = g^{*s}, h_1 = h^s,$
5: $\{g_{i,j} = g^{s^i t^j}, \hat{g}_{i,j} = \hat{g}^{s^i t^j}\}_{i=0,j=0}^{d_c,n_c}$, HASH: $\mathsf{G}^* \times \mathbb{Z}_q^* \to \mathbb{Z}_p$
6: $ck^{S1} := (g, h, g^*, \{g_{i,j}, \hat{g}_{i,j}\}_{i=0,j=0}^{d_c,n_c}, \hat{g}, \hat{h}, \hat{g^*}, g_1^*, h_1, \mathsf{HASH})$
7: $\mathsf{MUniEv\text{-}\Pi.Prove}(ck^{S1}, \mathbf{W}_{ed}^{r(1)}, \mathbf{X}_{et}^{r(1)}) \to \pi^{S1}$ :
8: $//\mathbf{W}_{ed}^{r(1)}, \mathbf{X}_{et}^{r(1)}$ are square matrices of $l_{new} = Mn^2$
9: $//l_{new}^2$ is smaller that $n_c$
10: //call the commitment scheme MPoly.Com of [22]
11: call MPoly.Com($ck^{S1}, \mathbf{W}_{ed}^{r(1)}) \to C_{\mathbf{W}_1}$:
  - $r_{\mathbf{W}_1} \xleftarrow{\$} \mathbb{Z}_p$
  - $C_{\mathbf{W}_1}^1 = h^{r_{\mathbf{W}_1}} g_{i,j}^{\sum_{i=0,j=0}^{d_c,l_{new}^2-1} w_{i,j}}, C_{\mathbf{W}_1}^2 = \hat{h}^{r_{\mathbf{W}_1}} \hat{g}_{i,j}^{\sum_{i=0,j=0}^{d_c,l_{new}^2-1} w_{i,j}},$
  - $C_{\mathbf{W}_1} := (C_{\mathbf{W}_1}^1, C_{\mathbf{W}_1}^2)$
12: call MPoly.Com($ck^{S1}, \mathbf{X}_{et}^{r(1)}) \to C_{\mathbf{X}}$:
  - $r_{\mathbf{X}} \xleftarrow{\$} \mathbb{Z}_p, C_{\mathbf{X}}^1 = h^{r_{\mathbf{X}}} g_{i,j}^{\sum_{i=0,j=0}^{d_c,l_{new}^2-1} x_{i,j}}, C_{\mathbf{X}}^2 = \hat{h}^{r_{\mathbf{X}}} \hat{g}_{i,j}^{\sum_{i=0,j=0}^{d_c,l_{new}^2-1} x_{i,j}},$
  - $C_{\mathbf{X}} := (C_{\mathbf{X}}^1, C_{\mathbf{X}}^2)$
13: call MPoly.Com($ck^{S1}, \mathbf{Y}_{et}^{r(1)}) \to C_{\mathbf{Y}}$:
  - $r_{\mathbf{Y}} \xleftarrow{\$} \mathbb{Z}_p, C_{\mathbf{Y}}^1 = h^{r_{\mathbf{Y}}} g_{i,j}^{\sum_{i=0,j=0}^{d_c,l_{new}^2-1} y_{i,j}}, C_{\mathbf{Y}}^2 = \hat{h}^{r_{\mathbf{Y}}} \hat{g}_{i,j}^{\sum_{i=0,j=0}^{d_c,l_{new}^2-1} y_{i,j}},$
  - $C_{\mathbf{Y}} := (C_{\mathbf{Y}}^1, C_{\mathbf{Y}}^2)$
14: $k = \mathsf{HASH}(C_{\mathbf{W}_1}, C_{\mathbf{X}}, \mathbf{W}_{ed}^{r(1)}, \mathbf{X}_{et}^{r(1)}, \mathbf{Y}_{et}^{r(1)})$
15: $//\mathbf{X}_{et}^{r(1)} \in \mathbf{M}_{\mathbb{Z}_q[x](x^{d_c}+1)}^{l_{new}^2}, \mathbf{X}(x,y) = \sum_{i=0,j=0}^{d_c,l_{new}^2-1} x_{i,j} x^i y^j$
16: $//\mathbf{W}_{et}^{r(1)} \in \mathbf{M}_{\mathbb{Z}_q[x](x^{d_c}+1)}^{l_{new}^2}, \mathbf{W}(x,y) = \sum_{i=0,j=0}^{d_c,l_{new}^2-1} w_{i,j} x^i y^j$
17: $//\mathbf{Y}_{et}^{r(1)} \in \mathbf{M}_{\mathbb{Z}_q[x](x^{d_c}+1)}^{l_{new}^2}, \mathbf{Y}(x,y) = \sum_{i=0,j=0}^{d_c,l_{new}^2-1} y_{i,j} x^i y^j$
18: $L_1 := L_1(x,y) = \mathbf{W}(x,y) + \mathbf{X}(x,y) + \mathbf{Y}(x,y) = \sum_{i=0,j=0}^{d_c,l_{new}^2-1} l_{i,j} x^i y^j$
19: call MPoly.Com($ck^{S1}, L_1(s,t)) \to C_{L_1} = (C_{L_1}^1, C_{L_1}^2)$
20: $L_{1,k} := L_1(k,y) = \mathbf{W}(k,y) + \mathbf{X}(k,y) + \mathbf{Y}(k,y) = \sum_{i=0,j=0}^{d_c,l_{new}^2-1} l_{i,j} k^i y^j$
21: call MPoly.Com($ck^{S1}, \mathbf{W}_{ed,k}^{r(1)}) \to C_{\mathbf{w},k} = (C_{\mathbf{w},k}^1, C_{\mathbf{w},k}^2)$
22: call MPoly.Com($ck^{S1}, \mathbf{X}_{ed,k}^{r(1)}) \to C_{\mathbf{x},k} = (C_{\mathbf{x},k}^1, C_{\mathbf{x},k}^2)$
23: call MPoly.Com($ck^{S1}, \mathbf{Y}_{ed,k}^{r(1)}) \to C_{\mathbf{y},k} = (C_{\mathbf{y},k}^1, C_{\mathbf{y},k}^2)$
24: $C_{L_1,k}^1 = (C_{\mathbf{w},k}^1 \cdot C_{\mathbf{x},k}^1 \cdot C_{\mathbf{y},k}^1); C_{L_1,k}^2 = (C_{\mathbf{w},k}^2 \cdot C_{\mathbf{x},k}^2 \cdot C_{\mathbf{y},k}^2)$
25: $\mathbf{T}(x,y) = \frac{L_1(x,y) - L_1(k,y)}{(x-k)}$
26: $\bar{g} = \frac{h_1}{h^k}, a, b \xleftarrow{\$} \mathbb{Z}_p,$
27: call MPoly.Com($ck^{S1}, \mathbf{T}(x,y)) \to C_{\mathbf{T}}$:
  - $r_{\mathbf{T}} \xleftarrow{\$} \mathbb{Z}_p, C_{\mathbf{T}}^1 = h^{r_{\mathbf{T}}} \bar{g}^{\sum_{t=0,i=0,j=0}^{d_c-1,d_c-t,l_{new}^2-1} l_{i+t,j} k^{i-1} s^i t^j},$
    $C_{\mathbf{T}}^1 = h^{r_{\mathbf{T}}} \hat{g}^{\sum_{t=0,i=0,j=0}^{d_c-1,d_c-t,l_{new}^2-1} l_{i+t,j} k^{i-1} s^i t^j},$
  - $C_{\mathbf{T}} := (C_{\mathbf{T}}^1, C_{\mathbf{T}}^2)$
28: $U = e(h^a \bar{g}^b, g^*), \mathsf{e} \leftarrow \mathsf{HASH}(C_{L_1}, C_{L_1,k}, C_{\mathbf{T}}, U, k)$
29: $\sigma = a - (r_{L_1,k} - r_{L_1}) \cdot \mathsf{e} \mod q, \tau = b - r_{\mathbf{T}} \cdot \mathsf{e} \mod q, \pi^{S1} = (C_{\mathbf{T}}, \mathsf{e}, \sigma, \tau)$
30: $\mathsf{MUniEv\text{-}\Pi.Verify}(ck^{S1}, C_{L_1}, C_{\mathbf{X}}, C_{L_1,k}, C_{\mathbf{T}}, k, \pi^{S1}) \to 0/1$ :
31: $b_1 := (e(C_{L_1}^1, \hat{g^*}) == e(C_{L_1}^2, g^*)); b_2 := (e(C_{\mathbf{X}}^1, \hat{g^*}) == e(C_{\mathbf{X}}^2, g^*));$
32: $b_3 := (e(C_{\mathbf{w},k}^1, \hat{g^*}) == e(C_{\mathbf{w},k}^2, g^*) \wedge e(C_{\mathbf{x},k}^1, \hat{g^*}) == e(C_{\mathbf{x},k}^2, g^*)$
    $\wedge e(C_{\mathbf{y},k}^1, \hat{g^*}) == e(C_{\mathbf{y},k}^2, g^*));$
33: $b_4 := (e(C_{\mathbf{T}}^1, \hat{g^*}) == e(C_{\mathbf{T}}^2, g^*));$
34: $H = e(C_{L_1}^1, g_1^*/g^{*k}) \cdot e(C_{L_1}^1/C_{L_1,k}^1, g^*)^{-1}, \bar{g} = h_1/h^k,$
35: $U = e(h^\sigma \bar{g}^\tau, g^*) \cdot H^{\mathsf{e}},$
36: $b_5 := (\mathsf{e} == \mathsf{HASH}(C_{L_1}, C_{L_1,k}, C_{\mathbf{T}}, U, k));$
37: $b = (b_1 \wedge b_2 \wedge b_3 \wedge b_4 \wedge b_5) \to 0/1.$

Fig. 8: Generating zk-SNARK proofs for relation $R_{PN}^{S1}$.

---

**Step 2 CaP zk-SNARK Proofs over QMP**

$//\mathsf{st} = (\mathbf{W}_{ed,k}^{r(1)}, \mathbf{X}_{et,k}^{r(1)}, \mathbf{Y}_{et,k}^{r(1)}); \quad io = m = 3; \mathsf{QMP} = (L(x), R(x), O(x), t(x))$
$//p(x, \mathbf{Z}) = \mathrm{tr}\{\mathbf{Z}^T \cdot L(x) \cdot R(x)\} - \mathrm{tr}\{\mathbf{Z}^T \cdot O(x)\} = h(x, \mathbf{Z}) t(x), \forall \mathbf{Z} \in \mathbf{M}_{(\mathbb{Z}_p)}^{l_{new}}$
1: $\mathsf{AC\text{-}\Pi.Setup}(1^\lambda, \mathsf{QMP}) \to \mathsf{crs}$:
2: $(\alpha, \beta, \gamma, \delta, \eta, z) \leftarrow \mathbb{Z}_p, \mathbf{Z} \leftarrow \mathbf{M}_{(\mathbb{Z}_p)}^{l_{new}^2}, g, h \leftarrow \mathsf{G}$
3: $\mathsf{crs} := (g, h, g^\alpha, g^\beta, h^\beta, h^\gamma, g^\delta, h^\delta, g^{\frac{\eta}{\delta}}, g^{\frac{\eta}{\gamma}}, \{\mathsf{OP}_i = g^{\frac{\beta L_i(z)}{\gamma}} \cdot g^{\frac{\alpha R_i(z)}{\gamma}} \cdot g^{\frac{O_i(z)}{\gamma}}\}_{i=0}^t)$
4: $\mathsf{AC\text{-}\Pi.Prove}(\mathsf{crs}, \mathsf{st}, \mathsf{wt}) \to \pi_{L_1}^{S2}$:
5: $L(x) = \cdot \mathbf{W}_{ed,k}^{r(1)} \cdot L_{\mathbf{W}}(z); R(x) = \mathbf{X}_{et,k}^{r(1)} \cdot R_{\mathbf{X}}(z); O(x) = \cdot \mathbf{Y}_{et,k}^{r(1)} \cdot O_{\mathbf{Y}}(z)$
6: $(t, s, v) \xleftarrow{\$} \mathbb{Z}_p, A := g^\alpha \cdot g^{\mathbf{Z}^T L(z)} \cdot g^{\delta t}, B := h^\beta \cdot h^{R(z)} \cdot h^{\delta s},$
7: $C := g^{h(z, \mathbf{Z}) t(z)/\delta} \cdot A^s \cdot (g^\beta \cdot g^{R(z)} \cdot g^{\delta s})^t \cdot g^{-ts\delta} \cdot g^{-v\eta/\delta},$
8: $D := g^{\frac{\beta \mathbf{Z}^T L(z)}{\gamma}} \cdot g^{\frac{\alpha R(z)}{\gamma}} \cdot g^{\frac{\mathbf{Z}^T O(z)}{\gamma}} \cdot g^{v\eta/\gamma},$
9: $d_1 = g^{\frac{\beta L(z)}{\gamma}} \cdot g^{v\eta/\gamma}, d_2 = g^{\frac{\alpha R(z)}{\gamma}} \cdot g^{v\eta/\gamma}, d_3 = g^{\frac{O(z)}{\gamma}} \cdot g^{v\eta/\gamma},$
10: $C_{L_1,k}' := (d_1, d_2, d_3), \quad \pi_{L_1}^{S2} := (A, B, C, D, C_{L_1,k}')$
11: $\mathsf{AC\text{-}\Pi.Verify}(\mathsf{crs}, \pi_{L_1}^{S2}) \to 0/1$:
12: $b := (\mathrm{tr}\{e(A, B)\} == \mathrm{tr}\{e(g^\alpha, h^\beta) \cdot e(D, h^\gamma) \cdot e(C, h^\delta)\}) \to 0/1.$

Fig. 9: Generating zk-SNARK proofs for relation $R_{PN}^{S2}[L_1]$.

---

**Proof Generation for $R_{PN}$.** As described before, the relation $R_{PN}$ defines that PriorNet is performed in an encryption manner, by taking the encrypted inputs $\mathbf{X}_{et}^{r(1)}$ in ring field $\mathbb{Z}_q[x](x^{d_c} + 1)$. To prove it, we adopt the methodology of efficient verification computation on encrypted data by Fiore *et al.* [22]. We then integrate it with the aforementioned CaP zk-SNARK from Groth scheme [21] for proving the satisfiability of QMP-based arithmetic circuits that are used to express convolutional operations.

Following Fiore *et al.*, there are two steps for completing the proofs. With respective to the relation $R_{PN}^{S1}$, Step 1 is to prove the knowledge of the commitments to test inputs, trained weights of $L_1$, and intermediate outputs, by proving simultaneous evaluation of multiple ring field polynomials on a same point, as demonstrated in Fig. 8. With respective to the relation $R_{PN}^{S2}[L_1]$, Step 2 shown in Fig. 9 is to prove the arithmetic correctness of the matrix multiplication computation with our idea of using QMP. The proofs in Fig. 9 contains the commitments to the QMP, and they are $C_{L_1,k}' = (d_1, d_2, d_3)$. For another relation $R_{PN}^{S2}[L_2, L_3]$, the computation conducted in the two layers involves a 2-degree polynomial evaluated on $\mathbf{Y}_{et}^{r(1)}$ and subsequently an average pooling function. Hence, we remain using the QAP-based circuit to express them, instead of QMP-based one, and generate proofs by directly leveraging the CaP variant of the Groth scheme [21].

**Proof Generation for $R_{LN}$.** We previously describe in the relation $R_{LN}$ that LaterNet is executed over scalars instead of polynomial rings, due to all of the data here are plaintext. Hence, Step 1 is not needed and we directly adopt the above Step 2 to prove the validity of the relations regarding the LaterNet. Notice, before proof generation, we express matrix multiplication operations in $R_{LN}[L_4]$ and $R_{LN}[L_7, L_8]$ using QMP-based arithmetic circuits, while expressing non-linear operations in $R_{LN}[L_5, L_6]$ and $R_{LN}[L_o]$ using QAP-based arithmetic circuits.

**Proofs Composition.** We now need to combine the previous layer-wise proofs together by following the commit-and-prove methodology. Concretely, we need to prove two
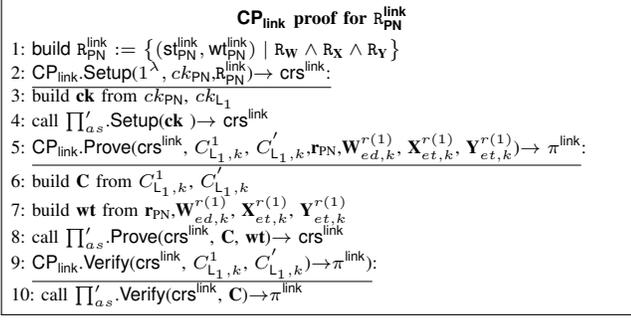
Fig. 10: Generating zk-SNARK proofs for relation $\mathsf{R}_{\mathsf{PN}}^{\mathsf{link}}$.

commitments of two successive layer proofs open to the same values. One of the commitments may commit to the inputs of a current layer while another one may commit to its previous layer's outputs. Taking two separate proofs $\pi^{\mathsf{S1}}$ and $\pi_{\mathsf{L}_1}^{\mathsf{S2}}$ as an example, we ensure the commitments $C_{\mathsf{L}_1,k}^1$ and $C_{\mathsf{L}_1,k}'$ are opened to the same value, *i.e.*, $\mathsf{R}_{\mathsf{PN}}^{\mathsf{link}} := \{(\mathsf{st}_{\mathsf{PN}}^{\mathsf{link}} = (C_{\mathsf{L}_1,k}^1, C_{\mathsf{L}_1,k}', ck_{\mathsf{PN}}, ck_{\mathsf{L}_1}), \mathsf{wt}_{\mathsf{PN}}^{\mathsf{link}} = (\mathbf{W}_{et,k}^{r(1)}, \mathbf{X}_{et,k}^{r(1)}, \mathbf{Y}_{et,k}^{r(1)}, r_{\mathsf{PN}}, r_{\mathsf{PN}}')) \mid C_{\mathsf{L}_1,k}^1 = \mathsf{MPoly.Com}(\mathbf{W}_{et,k}^{r(1)}, \mathbf{X}_{et,k}^{r(1)}, \mathbf{Y}_{et,k}^{r(1)}, r_{\mathsf{PN}}, ck_{\mathsf{PN}}) \wedge C_{\mathsf{L}_1,k}' = \mathsf{MPoly.Com}(\mathbf{W}_{et,k}^{r(1)}, \mathbf{X}_{et,k}^{r(1)}, \mathbf{Y}_{et,k}^{r(1)}, r_{\mathsf{PN}}', ck_{\mathsf{L}_1})\}$. Note that they are Pedersen-like commitments. We make use of the previously proposed CaP zk-SNARK $\mathsf{CP}_{\mathsf{link}}$ [21] to composite two separate proofs by proving the validity of the relation $\mathsf{R}_{\mathsf{PN}}^{\mathsf{link}}$. Also, the similar methodology can be applied into other-layer relations for linking the proofs regarding any two subsequent layers. The key idea of the methodology is to reshape the computation of Pedersen commitment into the linear subspace computation. Concretely, the relation of proving $C_{\mathsf{w},k}^1 = h^{r_{w,k}} g_{0,i}^{\sum_{i=0}^{l_{new}^2-1} w_i}$ (line 21 of Fig. 8) and $d_1 = g^{\frac{v\eta}{\gamma}} \cdot g^{\frac{\beta}{\gamma}\sum_{i=0}^{l_{new}^2-1} w_i}$ (line 9 of Fig. 9) committing to $\mathbf{W}_{ed,k}^{r(1)}$ is transformed into the linear subspace relation

$$\mathsf{R}_{\mathbf{W}} : \mathbf{C} = \mathbf{ck} \cdot \mathbf{wt} = \overbrace{\begin{bmatrix} g^{\frac{\beta}{\gamma}\mathsf{w}_0} \cdot g^{v\eta/\gamma} \\ g^{\frac{\beta}{\gamma}\mathsf{w}_1} \cdot g^{v\eta/\gamma} \\ ... \\ h^{r_{\mathsf{w},k}} g_{0,i}^{\sum_{i=0}^{l_{new}^2-1} w_i} \end{bmatrix}}^{\mathbf{C} \in \mathsf{G}^{l_{new}^2}}$$

$$= \overbrace{\begin{bmatrix} 0 & g^{\frac{\eta}{\gamma}} & g^{\frac{\beta}{\gamma}} & 0 & ... & 0 \\ 0 & g^{\frac{\eta}{\gamma}} & 0 & g^{\frac{\beta}{\gamma}} & ... & 0 \\ ... & ... & ... & ... & ... & ... \\ 0 & g^{\frac{\eta}{\gamma}} & 0 & 0 & ... & g^{\frac{\beta}{\gamma}} \\ h & 0 & g_{0,0} & ... & ... & g_{0,(l_{new}^2-1)} \end{bmatrix}}^{\mathbf{ck} \in \mathsf{G}^{(l_{new}^2) \times (l_{new}^2+2)}} \times \overbrace{\begin{bmatrix} r_{\mathsf{w},k} \\ v \\ \mathsf{w}_0 \\ \mathsf{w}_1 \\ ... \\ \mathsf{w}(l_{new}^2-1) \end{bmatrix}}^{\mathbf{wt} \in \mathsf{Z}_p^{(l_{new}^2+2)}}.$$

Similarly, we can deduce relations $\mathsf{R}_{\mathbf{X}}$ and $\mathsf{R}_{\mathbf{Y}}$ for proving $C_{\mathsf{x},k}^1$ and $d_2$ committing to $\mathbf{X}_{et,k}^{r(1)}$, and $C_{\mathsf{y},k}^1$ and $d_3$ committing to $\mathbf{Y}_{et,k}^{r(1)}$, respectively. As a result, we represent the former $\mathsf{R}_{\mathsf{PN}}^{\mathsf{link}}$ as $\mathsf{R}_{\mathsf{PN}}^{\mathsf{link}} := \{(\mathsf{st}_{\mathsf{PN}}^{\mathsf{link}}, \mathsf{wt}_{\mathsf{PN}}^{\mathsf{link}}) \mid \mathsf{R}_{\mathbf{W}} \wedge \mathsf{R}_{\mathbf{X}} \wedge \mathsf{R}_{\mathbf{Y}}\}$. Then, as demonstrated in Fig. 10, the proof $\pi^{\mathsf{link}}$ for such a relation is generated by leveraging $\mathsf{CP}_{\mathsf{link}}$ that calls a scheme for proving linear subspace relations $\mathsf{R}_{\mathbf{W}}$, $\mathsf{R}_{\mathbf{X}}$ and $\mathsf{R}_{\mathbf{Y}}$.

### D. Aggregating multiple proofs

Due to that a CNN model is tested with the test data from multiple testers, there are multiple proofs with respect to proving the correctness of each inference process. For ease of storage overhead and verification cost, it is desirable to aggregate multiple proofs into a single proof. Besides, ensuring the verification computation as simple and low as possible can make our public platform more easily accessible to a later-coming user. To the end, this section proceeds to introduce the algorithm of aggregating multiple proofs.

**Relation Definition.** Suppose there are $n_t$ CaP proofs $\{\pi^{\mathsf{m},i} = (A_i, B_i, C_i, D_i)\}_{i \in [n_t]}$ w.r.t a CNN $\mathsf{m}$ which is tested by $n_t$ testers. We define the aggregation relation that the multiple generated proofs are valid at the same time: $\mathsf{R}_{\mathsf{agg}} = \{(\mathsf{st}_{\mathsf{agg}}, \mathsf{wt}_{\mathsf{agg}}) \mid \{\mathsf{AC}\text{-}\prod.\mathsf{Verify}(\mathsf{vk}^{\mathsf{m}}, \pi^{\mathsf{m},i}, \mathsf{st}^{\mathsf{m},i}) \to 1\}_{i \in [n_t]}\}$. Here, $\mathsf{vk}^{\mathsf{m}}$ is the verification key, extracted from the common randomness string in the $\mathsf{Setup}$ algorithm, and $\mathsf{vk}^{\mathsf{m}} = (g^{\alpha}, h^{\beta}, \{\mathsf{OP}_i\}_{i=0}^t, h^{\gamma}, h^{\delta})$. We note that $\mathsf{vk}^{\mathsf{m}}$ is the same for the different statements $\{\mathsf{st}^{\mathsf{m},i}\}_{i \in [n_t]}$, each of which contains $\mathsf{st}^{\mathsf{m},i} = (a_{i,0}, ..., a_{i,t})$. $\pi^{\mathsf{m},i}$ is the proof for the statement $\mathsf{st}^{\mathsf{m},i}$ based on the different test data and the different intermediates during the computation of $\mathsf{m}$.

**Proof Aggregation.** With the defined relation, we are ready to provide a proof for it. A crucial tool to generate such a proof is recently proposed SnarkPack [32]. The external effort for us is to extend this work targeted at the original Groth16 scheme [19] without a commit-and-prove component to our case, that is, the proofs based on the CaP Groth16 scheme [21] as well as our QMP-based proofs (see Fig. 9). As shown in line 18 of Fig. 11, we additionally compute a multi-exponentiation inner product for $\{D_i\}_{i \in [n_t]}$ of the multiple proofs, which is similar to computing the original $\{C_i\}_{i \in [n_t]}$. As a result, we enable a single verification on the aggregated proof $\pi_{\mathsf{agg}}$ for a later-coming user who is interested in the model $\mathsf{m}$ that is tested by $n_t$ testers.

## VI. SECURITY ANALYSIS

**Theorem 1.** *If the underlying CaP zk-SNARK schemes are secure commit-and-prove arguments of knowledge $\prod_s$, the used leveled FHE scheme L-FHE is semantically secure and commitments Com are secure, and moreover a model is partitioned into PriorNet and LaterNet with an optimal splitting strategy, then our publicly verifiable model evaluation satisfies the following three properties:*
**Correctness.** *If PriorNet $F_{M_1}$ runs on correctly encrypted inputs $\mathbf{X}_{et}$ and outputs $\mathbf{Y}_{et}^{(3)}$, and meanwhile, LaterNet $F_{M_2}$ runs on the correctly decrypted $\mathbf{Y}^{(3)}$ and returns $l_{test}$, then $l_{test} = F_{M_2}(\mathbf{Y}^{(3)})$ and $\mathbf{Y}_{et}^{(3)} = F_{M_1}(\mathbf{X}_{et})$ are verified;*
**Security.** *If a PPT $\mathcal{A}$ knows all public parameters and public outputs of $\prod_s$, L-FHE and Com, as well as, the computation of $F_{M_1}$ and $F_{M_2}$, but has no access to the model parameters of $F_{M_1}$, it cannot generate $l_{test}^*$ which passes verification but $l_{test}^* \neq F_{M_2}(\mathbf{Y}^{(3)})$;*
**Privacy.** *If a PPT $\mathcal{A}$ knows all public parameters and public outputs of $\prod_s$, L-FHE and Com, as well as, the computation*

Fig. 11: Generating zk-SNARK proofs for relation $R_{agg}$.

*of $F_{M_1}$ and $F_{M_2}$, but has no access to the model parameters of $F_{M_1}$, it cannot obtain any information about the clear $X$.*

We recall that a secure CaP zk-SNARK scheme should satisfy the properties of completeness, knowledge soundness and zero-knowledge; a secure commitment should be correct, hiding and binding, and a secure L-FHE scheme satisfies semantic security and correctness. Now we are ready to analyze how the three properties of correctness, security and privacy can be supported by the underlying used cryptographic components of our publicly verifiable model evaluation method. Notice, as we mentioned in Section III-A, the proof generation for $R_{PN}$ needs an additional step (namely Step 1 in Fig. 8), since $R_{PN}$ relates to encrypted computation. $R_{LN}$ only needs Step 2, relying on the completeness, knowledge soundness and zero-knowledge of the underlying CaP zk-SNARK scheme [21]. Our following analysis is exactly for $R_{PN}$.

The correctness property relies on the correctness of L-FHE and commitments, and the completeness of the underlying CaP zk-SNARK schemes. We employ a widely-adopted L-FHE scheme [50] which has been proved correct. For CaP zk-SNARK schemes, our designs use (*a*) Fiore et al.'s CaP zk-SNARK for simultaneous evaluation of multiple ciphertexts generated by the L-FHE scheme [22] (see Fig. 8), (*b*) the CaP Groth16 scheme [19], [21], (*c*) our QMP-based zk-SNARK derived from the CaP Groth16 (see Fig. 9), and (*d*) the CaP zk-SNARK $CP_{link}$ for compositing separate proofs [21] (see Fig. 10).

Based on *Theorem 10* of scheme (*a*) [22], we deduce the correctness of zk-SNARK proofs for relation $R_{PN}^{S1}$ by direct verification (see Fig. 8). Specifically, if $C_{L_1}, C_X$ are correct commitments, and for $L_1(x,y) = W(x,y) + X(x,y) + Y(x,y) = \sum_{i=0,j=0}^{d_c, l_{new}^2} l_{i,j} x^i y^j$ and random point $k \in Z_p$,

$s, t \in Z_p$, the following formula holds,

$$\mathbf{T}(s,t) = \frac{\mathsf{L}_1(s,t) - \mathsf{L}_1(k,t)}{(s-k)}$$

$$= \frac{[(\mathsf{L}_{1,0}(s) \cdot t^0 + ... + \mathsf{L}_{1,l_{new}^2-1}(s) \cdot t^{l_{new}^2-1}) - }{(s-k)}$$
$$\frac{(\mathsf{L}_{1,0}(k) \cdot t^0 + ... + \mathsf{L}_{1,l_{new}^2-1}(k) \cdot t^{l_{new}^2-1})]}{(s-k)}$$

$$= \frac{[(l_{0,0}(s-k) + ... + l_{d_c,0}(s^{d_c} - k^{d_c})) \cdot t^0 + ... +}{(s-k)}$$
$$\frac{(l_{0,l_{new}^2-1}(s-k) + ... + l_{d_c,l_{new}^2-1}(s^{d_c} - k^{d_c})) \cdot t^{l_{new}^2-1}]}{(s-k)}$$

$$= \sum_{j=0}^{l_{new}^2-1} \sum_{t=0}^{d_c-1} \sum_{i=0}^{d_c-t} l_{i+t,j} k^{i-1} s^i t^j.$$

Note that $C_T$ correctly commits to $\mathbf{T}(s,t)$, then given $\pi^{S1} = (C_T, e, \sigma, \tau)$, a verifier can compute $U$ which is equal to $e(h^a \bar{g}^b, g^*)$ (see line 28 in Fig. 8) by

$$U = e(h^\sigma \bar{g}^\tau, g^*) \cdot e(C_T^1, g_1^*/g^{*k})^e \cdot e(C_L^1/C_{L,k}^1, g^*)^{-e}$$

$$= e(h^\sigma \bar{g}^\tau, g^*) \cdot e(h^{r_T} \bar{g}^{\mathsf{T}(s,t)}, g^{*(s-k)})^e \cdot e(\frac{h^{r_T}}{h^{r_{T,k}}} \frac{g^{\mathsf{L}_1(s,t)}}{g^{\mathsf{L}_1(k,t)}}, g^*)^{-e}$$

$$= e(h^\sigma \bar{g}^\tau, g^*) \cdot e(h, g^*)^{e r_T(s-k)} \cdot e(g, g^*)^{e \mathsf{T}(s,t)(s-k)} \cdot$$
$$e(h, g^*)^{-e(r_T - r_{T,k})} \cdot e(g, g^*)^{-e(\mathsf{L}_1(s,t) - \mathsf{L}_1(k,t))}$$

$$= e(h, g^*)^\sigma \cdot e(h, g^*)^{e(r_{T,k} - r_T)} \cdot e(\bar{g}, g^*)^{\tau + e r_T} \cdot$$
$$e(g, g^*)^{e \mathsf{T}(s,t)(s-k)} \cdot e(g, g^*)^{-e(\mathsf{L}_1(s,t) - \mathsf{L}_1(k,t))}$$

$$= e(h, g^*)^a \cdot e(\bar{g}, g^*)^b \cdot 1$$

$$= e(h^a \bar{g}^b, g^*).$$

Based on *Theorem H.1.* of scheme (*b*) [21], we proceed to derive the correctness of our QMP-based zk-SNARK proofs for relation $R_{PN}^{S2}[L_1]$ (*i.e.*, scheme (*c*)) by verifying $\pi_{L_1}^{S2} := (A, B, C, D, C_{L_1,k}')$, where $C_{L_1,k}'$ will be used in $CP_{Link}$. Specifically, a verifier needs to verify that $(\star)$ $\mathsf{tr}\{e(A,B)\} \stackrel{?}{==} \mathsf{tr}\{e(g^\alpha, h^\beta) \cdot e(D, h^\gamma) \cdot e(C, h^\delta)\}$, where

$$\mathsf{tr}\{e(A,B)\} = \mathsf{tr}\{e(g^\alpha \cdot g^{Z^T L(z)} \cdot g^{\delta t}, h^\beta \cdot h^{R(z)} \cdot h^{\delta s})\}$$

$$= e(g,h)^{\mathsf{tr}\{(\alpha + Z^T L(z) + \delta t)(\beta + R(z) + \delta s)\}}$$

$$= e(g,h)^{(\alpha + \delta t)(\beta + R(z) + \delta s) + \mathsf{tr}\{Z^T L(z)(\beta + R(z) + \delta s)\}}$$

$$= e(g,h)^{(\alpha + \delta t)(\beta + R(z) + \delta s) + \beta Z^T L(z) + \mathsf{tr}\{Z^T L(z) R(z)\} + \delta s Z^T L(z)},$$

$$\mathsf{tr}\{e(g^\alpha, h^\beta) \cdot e(D, h^\gamma) \cdot e(C, h^\delta)\}$$

$$= \mathsf{tr}\{e(g,h)^{\alpha\beta} \cdot e(g^{\frac{\beta Z^T L(z) + \alpha R(z) + Z^T O(z) + v\eta}{\gamma}}, h^\gamma)$$
$$\cdot e(g^{h(z,Z)t(z)/\delta + (\alpha + Z^T L(z) + \delta t)s + (\beta + R(z) + \delta s)t - ts\delta - v\eta/\delta}, h^\delta)\}$$

$$= e(g,h)^{\alpha\beta + \beta Z^T L(z) + \alpha R(z) + \mathsf{tr}\{Z^T O(z)\} + v\eta}$$
$$\cdot e(g,h)^{h(z,Z)t(z) + (\alpha + Z^T L(z) + \delta t)s\gamma + (\beta + R(z) + \delta s)t\gamma - ts\gamma^2 - v\eta}.$$

Due to $p(x,Z) = \mathsf{tr}\{Z^T \cdot L(x) \cdot R(x)\} - \mathsf{tr}\{Z^T \cdot O(x)\} = h(x,Z)t(x)$, formula $(\star)$ holds.

The presentation of the correctness of schemes (*b*) and (*d*) is omitted, as they can be directly found in [21].

The security property relies on the correctness of L-FHE, and the knowledge soundness of CaP zk-SNARK we leveraged. The correctness of L-FHE ensures that $Y_{et}^{r(3)} = f^{avg}(f^{act}(f^{conv}(X_{et}^{r(1)})))$ decrypts to $Y^{r(3)} = f^{avg}(f^{act}(f^{conv}(X^{r(1)}))$. For the knowledge soundness of our

two-step zk-SNARK proofs (Fig. 8 and Fig. 9), we rely on the knowledge soundness of the two schemes (*a*) and (*c*), and assume an adversary can black-box query a random oracle HASH. We next analyze it specifically for the first layer $L_1$ from two aspects.

On the one hand, based on the knowledge soundness of schemes (*a*) and (*c*), if any adversary $\mathcal{A}$ can provide valid proofs with regard to layer-by-layer computations, *e.g.,* $\pi^{S1}$ and $\pi^{S2}$ for $L_1$, there exists an extractor who is able to output the witnesses satisfying the corresponding defined relations $R_{PN}^{S_1}, R_{PN}^{S_2}[L_1]$, with all but negligible probability. Particularly, the knowledge soundness of scheme (*a*) gives us that $\mathbf{T}(s,t) = \sum_{j=0}^{l_{new}^2-1} \sum_{t=0}^{d_c-1} \sum_{i=0}^{d_c-t} l_{i+t,j} k^{i-1} s^i t^j$ is correct evaluation value in the random point $k$ of the polynomial $L_1(x,t)$; the knowledge soundness of scheme (*c*) gives us that $\text{tr}\{Z^T \cdot L(z) \cdot R(z)\} - \text{tr}\{Z^T \cdot O(z)\} = h(z,Z)t(z)$.

On the other hand, the remaining probability $\mathcal{A}$ can cheat is that $\mathbf{W}(x,t) \odot \mathbf{X}(x,t) - \mathbf{Y}^*(x,t)$ is a non-zero polynomial while $\mathbf{W}(k,t) \odot \mathbf{X}(k,t) - \mathbf{Y}^*(k,t) = 0$ (meaning $\mathbf{Y}_{et}^{*r(1)} = \mathbf{W}_{ed}^{r(1)} \odot \mathbf{X}_{et}^{r(1)}$), which is negligible in $L_1$. This depends on that we ensure parameters $q \gg d_c$ ($q = 2^{109}$ and $d_c = 4096$) and the point $k$ is randomly generated by the random oracle HASH. The probability that $k$ is the root of the non-zero polynomial $\mathbf{W}(x,t) \odot \mathbf{X}(x,t) - \mathbf{Y}^*(x,t)$ thus is negligible, *i.e.,* $d_c/q = 4096/2^{109}$. Last, to analyze the knowledge soundness of the proofs for other layers of PriorNet is similar to the above two-aspect analysis for $L_1$, due to the nature of layer-wise computation.

The privacy property relies on the semantic security of L-FHE, the hiding property of comments and the zero-knowledge of the leveraged CaP zk-SNARK schemes. In terms of semantic security, we derive from the security of a previously proposed L-FHE (applied to PriorNet), such that for any PPT adversary $\mathcal{A}$ who has acess to the resulting proofs, the probability of the following experiment $\text{Exp}_{\mathcal{A}}^{Privacy}[\text{L-FHE}, f^{avg}(f^{act}(f^{conv}(\cdot), \lambda]$ outputting 1 is not larger than $1/2 + \text{negl}(\lambda)$:

$$\text{Exp}_{\mathcal{A}}^{Privacy}[\text{L-FHE}, f^{avg}(f^{act}(f^{conv}(\cdot), \lambda] :$$
$$b \leftarrow 0, 1;$$
$$(pk_u, sk_u) \leftarrow \text{L-FHE.KeyGen}(1^\lambda);$$
$$(\mathbf{X}_0^{r(1)}, \mathbf{X}_1^{r(1)}) \leftarrow \mathcal{A}(pk_u);$$
$$(\mathbf{X}_{et,b}^{r(1)}) \leftarrow \text{L-FHE.Enc}(\mathbf{X}_b^{r(1)});$$
$$b^* \leftarrow \mathcal{A}(pk_u, \mathbf{X}_{et,b}^{r(1)}, \pi^{S1}, \pi^{S2}),$$
$$\text{If } b^* = b, \text{output 1, else 0.}$$

For commitments, we leverage the MPoly.Com scheme which is perfect hiding as proven in *Theorem* 8 of [22].

We now analyse zero-knowledge based on the zero-knowledge simulators of the underlying zk-SNARKs, including $(\text{Sim}^{\text{MUniEv}-\Pi.\text{Setup}}, \text{Sim}^{\text{MUniEv}-\Pi.\text{Prove}})$ and $(\text{Sim}^{\text{AC}-\Pi.\text{Setup}}, \text{Sim}^{\text{AC}-\Pi.\text{Prove}})$, respectively. $\text{Sim}^{\text{MUniEv}-\Pi.\text{Setup}}$ and $\text{Sim}^{\text{AC}-\Pi.\text{Setup}}$ are the same as the algorithms of $\text{MUniEv} - \Pi.\text{Setup}$ and $\text{AC} - \Pi.\text{Setup}$ respectively, and generate specific common random strings for the two zk-SNARKs by running $\text{Sim}^{\text{MUniEv}-\Pi.\text{Setup}}(1^\lambda) \rightarrow (\text{crs}^{S1}, \text{td}^{S1})$ and $\text{Sim}^{\text{AC}-\Pi.\text{Setup}}(1^\lambda) \rightarrow (\text{crs}^{S2}, \text{td}^{S2})$. Then, we use the simulations $\text{Sim}^{\text{MUniEv}-\Pi.\text{Prove}}$ and $\text{Sim}^{\text{AC}-\Pi.\text{Prove}}$ to generate simulated proofs with the crs for the corresponding statements. With the simulators, we proceed with the following games, in which any PPT distinguisher $\mathcal{D}$ successfully distinguishes the simulation with a negligible probability.

**Hybrid 0.** This starts with the real algorithms in Fig. 8 and Fig. 9, where the proofs $\pi^{S1}$ and $\pi^{S2}$ are generated by $\text{MUniEv} - \Pi.\text{Prove}$ and $\text{AC} - \Pi.\text{Prove}$, respectively.

**Hybrid 1.** This remains unchanged in the usage of witnesses to be proven and we generate commitments as the line 11-12, 19-23 and 27 presented of Fig. 8. But we adopt the zero-knowledge simulators including $(\text{Sim}^{\text{MUniEv}-\Pi.\text{Setup}}, \text{Sim}^{\text{MUniEv}-\Pi.\text{Prove}})$ and $(\text{Sim}^{\text{AC}-\Pi.\text{Setup}}, \text{Sim}^{\text{AC}-\Pi.\text{Prove}})$ to generate proofs. Based on that the real $\text{MUniEv} - \Pi.\text{Prove}$ and $\text{AC} - \Pi.\text{Prove}$ are zero-knowledge algorithms, Hybrid 0 and Hybrid 1 are indistinguishable.

**Hybrid 2.** This runs the same algorithms as Hybrid 1 except that we replace the values in the above commitments with zeros. Based on the hiding property of the leveraged MPoly.Com scheme, Hybrid 1 and Hybrid 2 are indistinguishable.

## VII. EXPERIMENTS

**Implementation.** The main components of our implementation include zk-SNARK systems, FHE and polynomial commitment. We firstly use the libsnark library in C++ to implement proving/verification based on Groth16 of QAP and QMP-based circuits during the process of CNN prediction. We then use the Microsoft SEAL library to encrypt test inputs and evaluate PriorNet on ciphered inputs, with necessary parameters setting (*i.e.,* $d_c = 4096, q = 2^{109}, t = 1032193 \approx 2^{20}$) as recommended. We proceed to implement the component of polynomial commitment using the libff library of libsnark, for generating commitments to encrypted/unencrypted data.

In the aspect of machine learning models, we use well-trained CNN models to run the prediction process over the single-channel MNIST and three-channel CIFAR-10 datasets. Specifically, we start from a toy CNN with a convolution layer, a ReLU layer plus an average pooling layer, and generate proofs w.r.t the toy CNN running on the MNIST dataset. We then extend it to LeNet-5 over the MNIST dataset. The LeNet-5 architecture is $32 \times 32 \times 1 \xrightarrow[\text{filter size}=5\times5\times6]{conv_1 \text{ and } act_1} 28 \times 28 \times 6 \xrightarrow[\text{size}=2\times2]{pool_1} 14 \times 14 \times 6 \xrightarrow[\text{filter size}=5\times5\times16]{conv_2 \text{ and } act_2} 10 \times 10 \times 16 \xrightarrow[\text{size}=2\times2]{pool_2} 5 \times 5 \times 16 \xrightarrow{fc_3 \text{ and } act_3} 120 \times 1 \times 1 \xrightarrow{fc_4 \text{ and } act_4} 84 \times 1 \times 1 \xrightarrow{fc_o} 10 \times 1 \times 1$. We additionally run the model on the CIFAR-10 dataset. In this part, we spend major efforts on approximately handling convolution operations over the two datasets with different-scale matrix multiplication, before generating proofs. Due to the parameters and pixel values of data are floating-point numbers, not adapting to zk-SNARK systems, we use a generic 8-bit unsigned quantization technique to transform them into integers in $[0, 255]$. We conduct our MNIST experiments on

a Ubuntu 18.04.6 server with a 6-core Ryzen5 processor, 7.8 GB RAM and 97.2 GB Disk space, and execute CIFAR-10 experiments on a docker container with Ubuntu18.04, Quad-core Intel i5-7500 processor running at 3.4 GHz and 48 GB RAM.

**Evaluation.** We evaluate the running time and the storage overhead of performing PriorNet over a single encrypted image, by taking LeNet-5 as an example and selecting different split points. The evaluation results are shown in TABLE II. We also evaluate the overhead with an increasing number of encrypted images, when a split point is selected at the activation layer after the first full connection layer, see TABLE III.

TABLE II: Running time (s) and storage overhead (MB) of performing PriorNet. If the split point is at $fc_3$, the PriorNet $F_{\mathbf{M}_1}$ can be represented by $f^{fc_3}(f^{pool_2}(f^{act_2}(f^{conv_2}(f^{pool_1}(f^{act_1}(f^{conv_1}(\cdot))))))))$.

| Split point | $conv_1$ | $conv_2$ | $fc_3$ | $fc_4$ | $fc_5$ |
|---|---|---|---|---|---|
| **Running time** | 12.95 | 46.08 | 53.70 | 55.50 | 56.26 |
| **Storage** | 5.11 | 84.31 | 1702.23 | 2035.23 | 2063.13 |

TABLE III: Running time (s) with an increasing number of encrypted images.

| Dataset size | 10 | 30 | 50 | 80 | 100 |
|---|---|---|---|---|---|
| **Running time** | 1272.68 | 3807.05 | 6343.02 | 10261.01 | 12687.47 |

Next, we compare QMP-based zk-SNARK and QAP-based zk-SNARK for matrix multiplication in terms of the proving time, Setup time and CRS size, as shown Fig.12 and TABLE IV. Concretely, we simulate the operations of matrix multiplication in increasing dimensions up to $200 \times 200$ padding with random integers in $[0, 10]$. Then, we generate proofs for the matrix multiplication operations using the QAP-based and QMP-based zk-SNARK. Fig.12 shows that the QMP-based zk-SNARK is efficient than the QAP-based one in generating CRS and proof, as the matrix dimension increases. For the $200 \times 200$ matrix multiplication, the QMP-based zk-SNARK is $17.6\times$ and $13.9\times$ faster than the QAP-based one in Setup time and proving time, respectively. Besides, the QMP-based zk-SNARK obviously produces smaller CRS size than the QAP-based zk-SNARK, see TABLE IV.
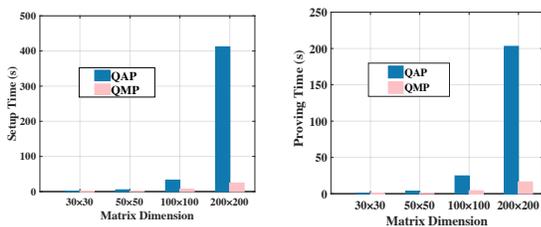


Fig. 12: Comparison on Setup time and proving time.

TABLE IV: CRS size comparison (KB).

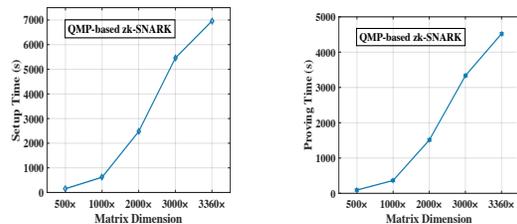| Matrix dimension | $30 \times 30$ | $50 \times 50$ | $100 \times 100$ | $200 \times 200$ |
|---|---|---|---|---|
| **QAP** | 2,911.22 | 12,445.73 | 97,230.09 | 768,503.09 |
| **QMP** | 84.42 | 233.83 | 934.21 | 3,735.72 |
| **Matrix dimension** | $500\times$ | $1000\times$ | $2000\times$ | $3000\times$ |
| **QMP** | 23,346.32 | 93,384.16 | 373,535.53 | 840,454.47 |



Fig. 13: Performance of our QMP-based zk-SNARK handling matrix multiplication in dimensions greater than $220 \times 220$.
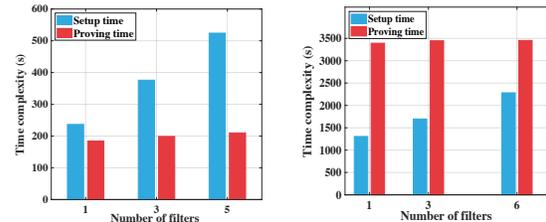


Fig. 14: Setup time and proving time for a conv. layer with 1000 images from MNIST and CIFAR-10.

We can see in Fig. 13 that the QMP-based zk-SNARK can handle the matrix multiplication in increasing dimensions up to $3360 \times 3360$, which is its merit, compared to the QAP-based zk-SNARK supporting the maximum number of multiplication bound $10^7$ [14]. Also, the proving time almost increases linearly by the dimension of matrices.

We proceed to apply the QMP-based zk-SNARK in a convolution layer with stride $(1, 1)$, taking as inputs 1000 single-channel images of $28 \times 28$ from MNIST and 1000 three-channel images of $32 \times 32$ from CIFAR-10. For the former dataset, we use 5 filters of $5 \times 5$ while for the latter one, we use 6 filters. We transform the convolution operations into a matrix multiplication between a weight matrix and an input matrix both in dimension $3360 \times 3360$ for MNIST (similarly, $4704 \times 4704$ for CIFAR-10). We note $3360 = 28 \times (28 - 5 + 1) \times 5 \approx mn^2$, where $m, n$ mean the dimension of a filter and an image, respectively. Here, $m \geq M$ which is the number of filters. Different from the aforementioned experiments where the values of matrices are random integers, the weight matrix of $3360 \times 3360$ here is strategically assigned with the weights of filters plus some zero elements as padding, for ensuring computation correctness (recall it in Fig.5), and similarly, the input matrix of $3360 \times 3360$ is assigned with the pixel values of 1000 images, padding with $3360 \times (3360 - 1000)$ zeros. As a result, the average performance results of 5 runs are shown in Fig. 14. We discover that compared to the aforementioned experiments as presented in Fig.13, the Setup and proving time here are relatively smaller. The main reason can be the

TABLE V: Other performance metrics.

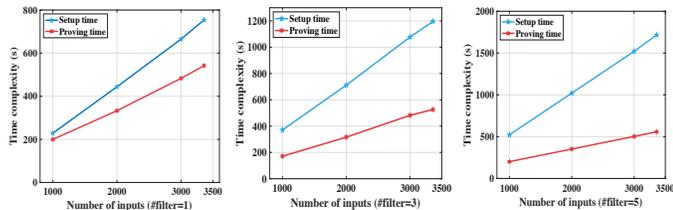| Datasets | #Filter | Verification time (s) | Size (KB) | |
|---|---|---|---|---|
| | | | **CRS** | **Proof** |
| **MNIST** | 1 | 54559 | 1,054,266.00 | 351,421.97 |
| | 3 | 54002 | 1,054,266.00 | 351,421.97 |
| | 5 | 56011 | 1,054,266.00 | 351,421.97 |
| **CIFAR-10** | 1 | 375942 | 15353859056 | 351,421.97 |
| | 3 | 376144 | 15353859056 | 351,421.97 |
| | 6 | 408911 | 16927629296 | 351,421.97 |

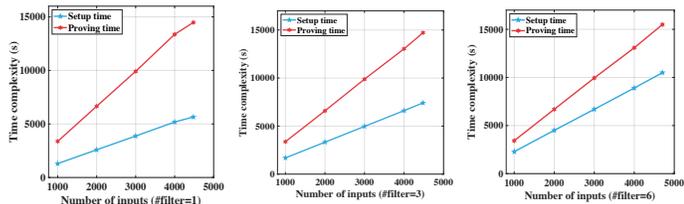Fig. 15: Time complexity on MNIST dataset.



Fig. 16: Time complexity on CIFAR-10 dataset.

padding zero elements inside the weight and input matrices cancel a lot of multiplications. Besides the Setup and proving time, TABLE V demonstrates the corresponding verification time, as well as constant CRS size and proof size.

We further conduct additional experiments (named Exp ($\star$)) on different number of filters and an increasing number of inputs from MNIST and CIFAR-10 datasets, as elaborated in Fig. 15 and Fig. 16. We can see from the figures that the proving time basically stays stable regardless of the number of filters, but the Setup time increases linearly by the number of filters $M$. We note that the complexity of the proving time is $O(mn^2 \cdot mn^2), m \geq M$ (resp. $O(Mn^2 \cdot Mn^2)$ if $m < M$), where $mn^2$ (resp. $Mn^2$) is the number of inputs, and the proving time is independent of the number of filters $M$.

We proceed to generate QAP-based proofs for the ReLU and average pooling operations on the $3360 \times 3360$ matrix, named Exp ($\star\star$). Note that a ReLU operation needs 20 constraints and an average pooling operation needs 144 constraints. The evaluated performance is elaborated in TABLE VI.

TABLE VI: Performance for Exp ($\star\star$) (s).

| Layer | Setup Time | Proving Time | Verification Time |
|---|---|---|---|
| ReLU | 5520.61 | 1448.83 | 14.78 |
| Pooling | 196.43 | 49.93 | 14.78 |

## VIII. Conclusion

The paper discusses using zk-SNARK systems for verifiable CNN testing on encrypted test data. The authors optimize matrix multiplication relations by representing convolution operations with a single MM computation and using a new QMP. This reduces the multiplication gate and proof generation overhead. They also aggregate multiple proofs into a single proof for the same CNN but different test datasets. They provide a proof-of-concept implementation and share their implementation code publicly.

## Acknowledgment

## References

[1] Y. LeCun, B. Boser *et al.*, "Backpropagation applied to handwritten zip code recognition," *Neural computation*, vol. 1, no. 4, pp. 541–551, 1989.
[2] R. Das, E. Piciucco *et al.*, "Convolutional neural network for finger-vein-based biometric identification," *IEEE Transactions on Information Forensics and Security*, vol. 14, no. 2, pp. 360–373, 2018.
[3] G. Accident, "A google self-driving car caused a crash for the first time." https://www.theverge.com/2016/2/29/11134344/google-self-driving-car-crash-report, 2016.
[4] "The two-year fight to stop amazon from selling face recognition to the police." https://www.technologyreview.com/2020/06/12/1003482/amazon-stopped-selling-police-face-recognition-fight/, 2020.
[5] M. Wicker, X. Huang *et al.*, "Feature-guided black-box safety testing of deep neural networks," in *Proc. of TACAS*, 2018.
[6] A. Aggarwal, S. Shaikh, S. Hans, S. Haldar, R. Ananthanarayanan, and D. Saha, "Testing framework for black-box ai models," in *Proc. of IEEE/ACM ICSE-Companion*, 2021.
[7] L. Ma, F. Juefei-Xu *et al.*, "Deepgauge: Multi-granularity testing criteria for deep learning systems," in *Proc. of ACM/IEEE ICASE*, 2018.
[8] D. Hendrycks and T. Dietterich, "Benchmarking neural network robustness to common corruptions and perturbations," in *ICLR*, 2018.
[9] L. Aroyo and P. Paritosh, "Uncovering unknown unknowns in machine learning," https://ai.googleblog.com/2021/02/uncovering-unknown-unknowns-in-machine.html, 2021.
[10] Z. Ghodsi, T. Gu, and S. Garg, "Safetynets: Verifiable execution of deep neural networks on an untrusted cloud," in *Proc. of NIPS*, 2017.
[11] M. I. M. Collantes *et al.*, "Safetpu: A verifiably secure hardware accelerator for deep neural networks," in *Proc. Of IEEE VTS*, 2020.
[12] S. Lee, H. Ko, J. Kim, and H. Oh, "vcnn: Verifiable convolutional neural network." https://eprint.iacr.org/2020/584.pdf, 2020.
[13] B. Feng, L. Qin, Z. Zhang, Y. Ding, and S. Chu, "Zen: An optimizing compiler for verifiable, zero-knowledge neural network inferences," https://eprint.iacr.org/2021/087.pdf, 2021.
[14] J. Keuffer, R. Molva, and H. Chabanne, "Efficient proof composition for verifiable computation," in *Proc. of ESORICS*, 2018.
[15] L. Zhao, Q. Wang *et al.*, "Veriml: Enabling integrity assurances and fair payments for machine learning as a service," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 10, pp. 2524–2540, 2021.
[16] A. Madi, R. Sirdey *et al.*, "Computing neural networks with homomorphic encryption and verifiable computing," in *ACNS Workshops*, 2020.
[17] T. Liu, X. Xie *et al.*, "zkcnn: Zero knowledge proofs for convolutional neural network predictions and accuracy," in *Proc. of ACM CCS*, 2021.
[18] J. Thaler, "Time-optimal interactive proofs for circuit evaluation," in *Annual Cryptology Conference*, 2013, pp. 71–89.
[19] J. Groth, "On the size of pairing-based non-interactive arguments," in *EUROCRYPT*, 2016, pp. 305–326.
[20] S. Agrawal, C. Ganesh *et al.*, "Non-interactive zero-knowledge proofs for composite statements," in *CRYPTO*, 2018, pp. 643–673.
[21] M. Campanelli *et al.*, "Legosnark: Modular design and composition of succinct zero-knowledge proofs," in *Proc. of ACM CCS*, 2019.

[22] D. Fiore, A. Nitulescu, and D. Pointcheval, "Boosting verifiable computation on encrypted data," in *PKC*, no. 12111, 2020, pp. 124–154.

[23] D. Mouris and N. G. Tsoutsos, "Zilch: A framework for deploying transparent zero-knowledge proofs," *IEEE Transactions on Information Forensics and Security*, 2021.

[24] "Bvlc/caffe," https://github.com/BVLC/caffe/wiki/Model-Zoo, 2019.

[25] Kaggle, "Data science competition platform." https://www.kaggle.com/.

[26] A. Marketplace, "Machine learning solutions." https://aws.amazon.com/marketplace/solutions/machine-learning.

[27] G. Xu *et al.*, "Verifynet: Secure and verifiable federated learning," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 911–926, 2019.

[28] Z. He, T. Zhang, and R. B. Lee, "Model inversion attacks against collaborative inference," in *Proc. of ACM ACSAC*, 2019.

[29] T. Ryffel, E. Dufour-Sans *et al.*, "Partially encrypted machine learning using functional encryption," in *NIPS*, 2019.

[30] Y. Aono, T. Hayashi *et al.*, "Privacy-preserving deep learning via additively homomorphic encryption," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 5, pp. 1333–1345, 2017.

[31] S. G. Francesco Alesiani, "Method for verifying information," https://patentimages.storage.googleapis.com/53/3c/62/0ed0b3f9bb163f/US20210091953A1.pdf, 2021.

[32] N. Gailly, M. Maller, and A. Nitulescu, "Snarkpack: Practical snark aggregation." in *Proc. of RWC*, 2022.

[33] J. Zhang, T. Liu *et al.*, "Doubly efficient interactive proofs for general arithmetic circuits with linear prover time," in *Proc. of ACM CCS*, 2021.

[34] S. Goldwasser *et al.*, "Delegating computation: interactive proofs for muggles," *Journal of the ACM*, vol. 62, no. 4, pp. 1–64, 2015.

[35] C. Niu, F. Wu, S. Tang, S. Ma, and G. Chen, "Toward verifiable and privacy preserving machine learning prediction," *IEEE Transactions on Dependable and Secure Computing*, 2020.

[36] J. Zhang, Z. Fang *et al.*, "Zero knowledge proofs for decision tree predictions and accuracy," in *Proc. of ACM CCS*, 2020.

[37] C. Weng, K. Yang, X. Xie, J. Katz, and X. Wang, "Mystique: Efficient conversions for zero-knowledge proofs with applications to machine learning," in *Proc. of USENIX Security*, 2021, pp. 501–518.

[38] C. Dong, J. Weng *et al.*, "Fusion: Efficient and secure inference resilient to malicious server and curious clients," in *Proc. of NDSS*, 2023.

[39] R. Gilad-Bachrach, N. Dowlin *et al.*, "Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy," in *ICML*, 2016, pp. 201–210.

[40] P. Mohassel and Y. Zhang, "Secureml: A system for scalable privacy-preserving machine learning," in *Proc. Of IEEE S&P*, 2017.

[41] P. Mishra, R. Lehmkuhl *et al.*, "Delphi: A cryptographic inference service for neural networks," in *USENIX Security*, 2020, pp. 2505–2522.

[42] C. Juvekar, V. Vaikuntanathan, and A. Chandrakasan, "{GAZELLE}: A low latency framework for secure neural network inference," in *Proc. of USENIX Security*, 2018.

[43] N. Kumar, M. Rathee *et al.*, "Cryptflow: Secure tensorflow inference," in *IEEE S&P*, 2020, pp. 336–353.

[44] K. Huang, X. Liu *et al.*, "A lightweight privacy-preserving cnn feature extraction framework for mobile sensing," *IEEE Transactions on Dependable and Secure Computing*, vol. 18, no. 3, pp. 1441–1455, 2019.

[45] E. Hesamifard, H. Takabi, and M. Ghasemi, "Cryptodl: Deep neural networks over encrypted data," https://arxiv.org/abs/1711.05189, 2017.

[46] X. Liu, R. H. Deng *et al.*, "An efficient privacy-preserving outsourced calculation toolkit with multiple keys," *IEEE Transactions on Information Forensics and Security*, vol. 11, no. 11, pp. 2401–2414, 2016.

[47] X. Liu *et al.*, "Efficient and privacy-preserving outsourced calculation of rational numbers," *IEEE Transactions on Dependable and Secure Computing*, vol. 15, no. 1, pp. 27–39, 2016.

[48] X. Liu, R. H. Deng *et al.*, "Privacy-preserving outsourced calculation toolkit in the cloud," *IEEE Transactions on Dependable and Secure Computing*, vol. 17, no. 5, pp. 898–911, 2018.

[49] R. Gennaro, C. Gentry *et al.*, "Quadratic span programs and succinct nizks without pcps," https://eprint.iacr.org/2012/215.pdf, 2012.

[50] J. Fan and F. Vercauteren, "Somewhat practical fully homomorphic encryption," https://eprint.iacr.org/2012/144.pdf, 2012.

[51] X. Chen, C. Liu *et al.*, "Targeted backdoor attacks on deep learning systems using data poisoning," https://arxiv.org/abs/1712.05526, 2017.

[52] G. Xu, H. Li *et al.*, "Secure and verifiable inference in deep neural networks," in *Proc. of ACM ACSAC*, 2020.

[53] Z. He, T. Zhang *et al.*, "Sensitive-sample fingerprinting of deep neural networks," in *CVPR*, 2019.

[54] Y. Zhang *et al.*, "Stealing neural network structure through remote fpga side-channel analysis," *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 4377–4388, 2021.

[55] M. Yan *et al.*, "Cache telepathy: Leveraging shared resource attacks to learn dnn architectures," in *Proc. of USENIX Security*, 2020.

[56] G. Saileshwar *et al.*, "Mirage: Mitigating conflict-based cache attacks with a practical fully-associative design." in *Proc. of USENIX Security*, 2021.

[57] A. E. Kosba, D. Papadopoulos *et al.*, "Trueset: Faster verifiable set computations," in *Proc. of USENIX Security*, 2014.

## A. Discussion

In terms of verification time and proof size, the QMP-based zk-SNARK has a higher overhead than the QAP-zk-SNARK. Its verification time grows faster than that of QAP-based zk-SNARK, see TABLE VIII. The proof size also becomes lager when the matrix dimension turns lager as shown in TABLE VII, while the proof size in the QAP-based zk-SNARK keeps 1019 bits regardless of the matrix dimension. We next see how the proof size becomes longer with the matrix dimension increasing. We note that the bounded number of multiplications the QAP zk-SNARK can handle is $220 * 220 * 220 = 10,648,000$, and then the QAP zk-SNARK would be called multiple times when handling the multiplication operations more than $10,648,000$, which results in multiple 1019-bit proofs. Suppose that the QAP-based zk-SNARK proofs for the $3000 \times 3000$ matrix multiplication are totally $\frac{3000*3000*3000}{220*220*220} \times 1019$ bits. Also, the QMP-based zk-SNARK proofs for the same matrix multiplication are $2,295,000,764$ bits. The proof size is nearly $888$ times larger than that of the above QAP-based zk-SNARK proof. We observe that the times of magnitude become smaller as the matrix dimension increases, see TABLE VII, which may mean that the QMP-based zk-SNARK is more suitable to handle large matrix multiplication.

TABLE VII: Times of magnitude in proof size.

| Matrix dimension | 1000× | 2000× | 3000× | 3360× | 4000× |
|---|---|---|---|---|---|
| Times | 2665 | 1332 | 888 | 793 | 666 |

TABLE VIII: Verification time comparison (s).

| Matrix dimension | 30 × 30 | 50 × 50 | 100 × 100 | 200 × 200 |
|---|---|---|---|---|
| QAP | 0.005 | 0.007 | 0.014 | 0.044 |
| QMP | 4.896 | 13.580 | 54.630 | 215.800 |

In our future work, we would introduce a random sampling strategy in the phase of proof generation, aiming to reduce the verification time. A straightforward idea can be adopted by randomly sampling a bounded number of values in two matrices to be multiplied, and resetting the non-chosen values as zeros. In such a way, only the sampled values in the two matrices are multiplied and only their multiplication correctness need to be proved. But here two noteworthy issues should be considered: (1) how to generate the randomness used for sampling against a distrusted prover; (2) how to determine the bound of the number of sampled values for ensuring computational soundness in zk-SNARK.