

AutoFHE: **Automated** Adaption of CNNs for Efficient Evaluation over **FHE**



WEI AO



Vishnu Boddeti



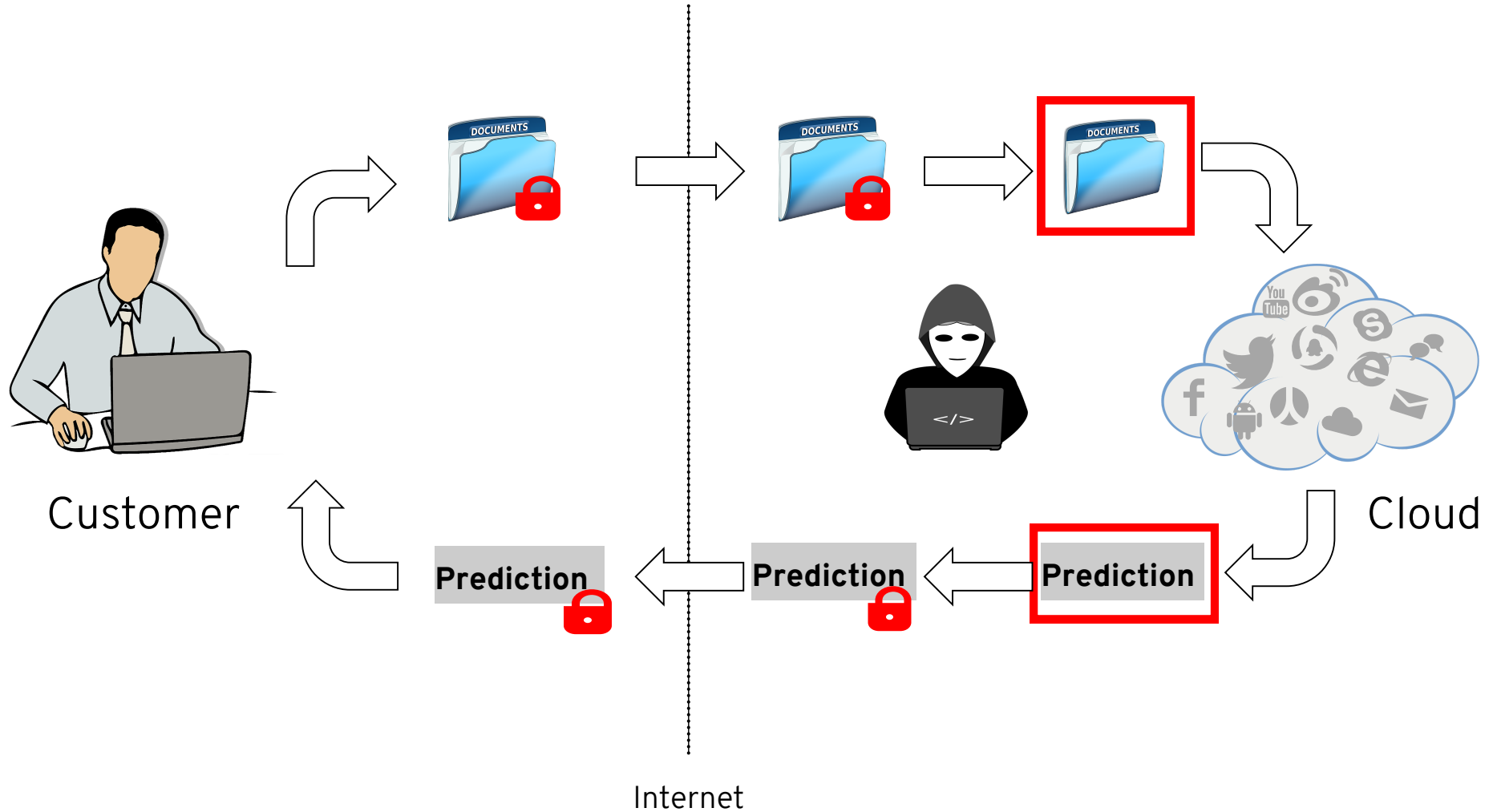
wei-ao.github.io

Michigan State University

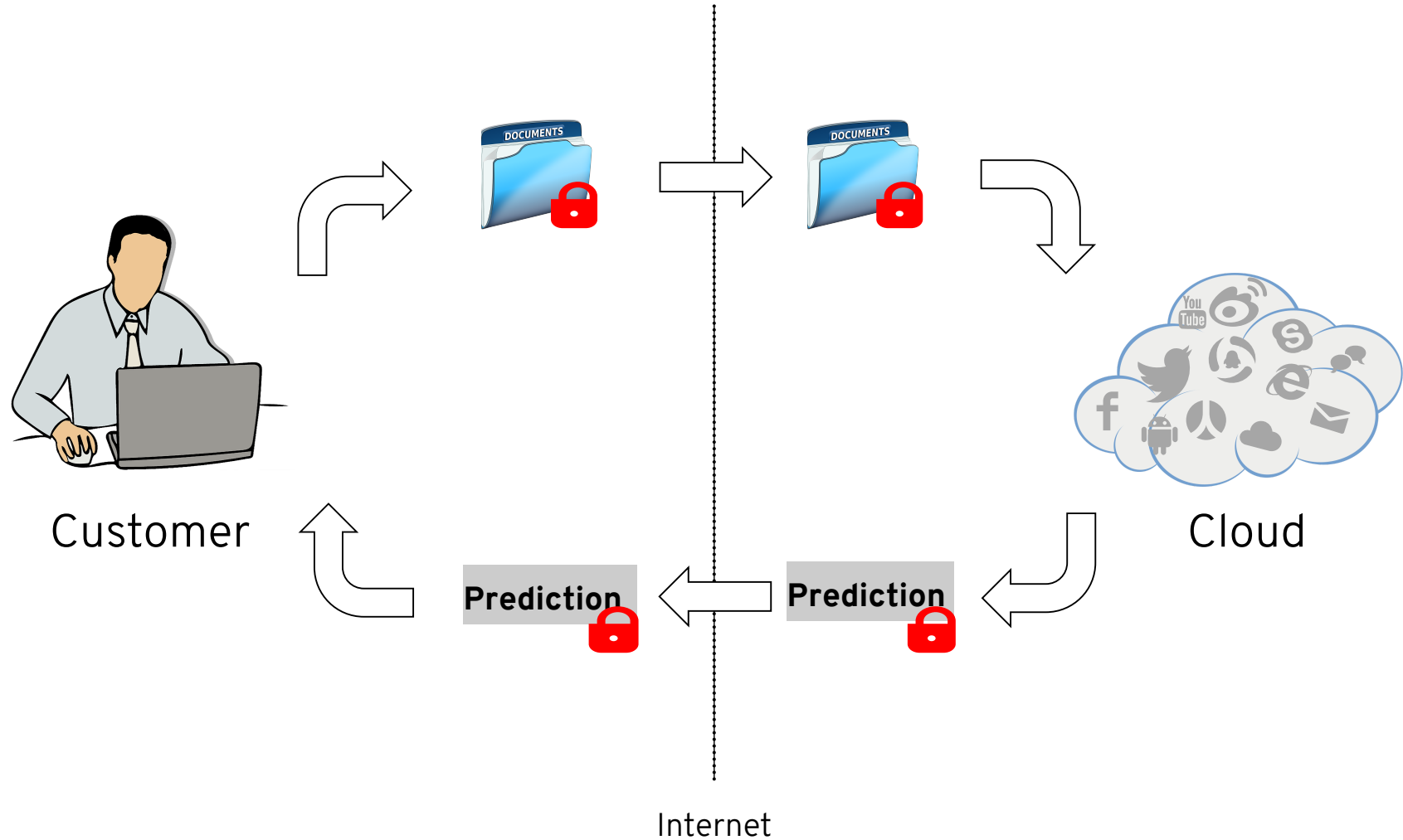
Philadelphia, PA, USA 2024

Secure deep learning under fully homomorphhic encryption

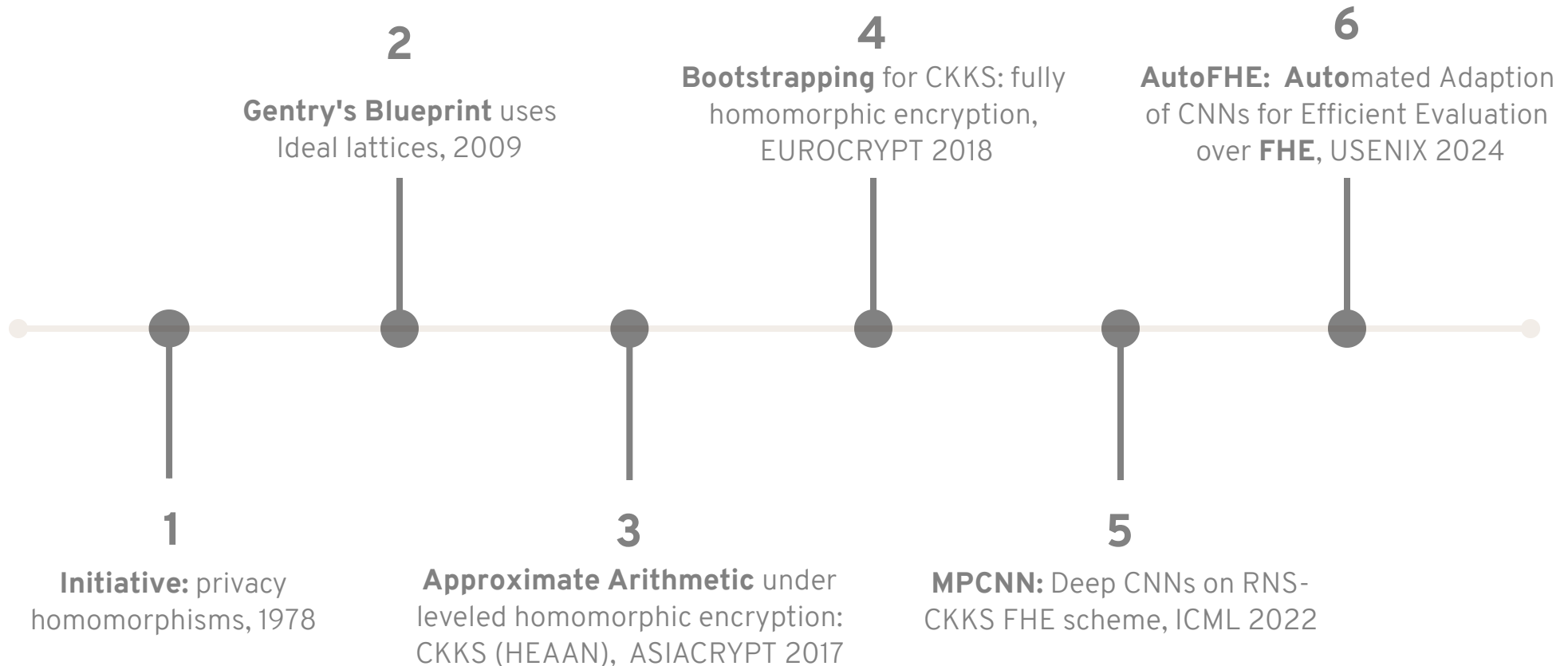
Deep Learning as a Service (DLaaS)



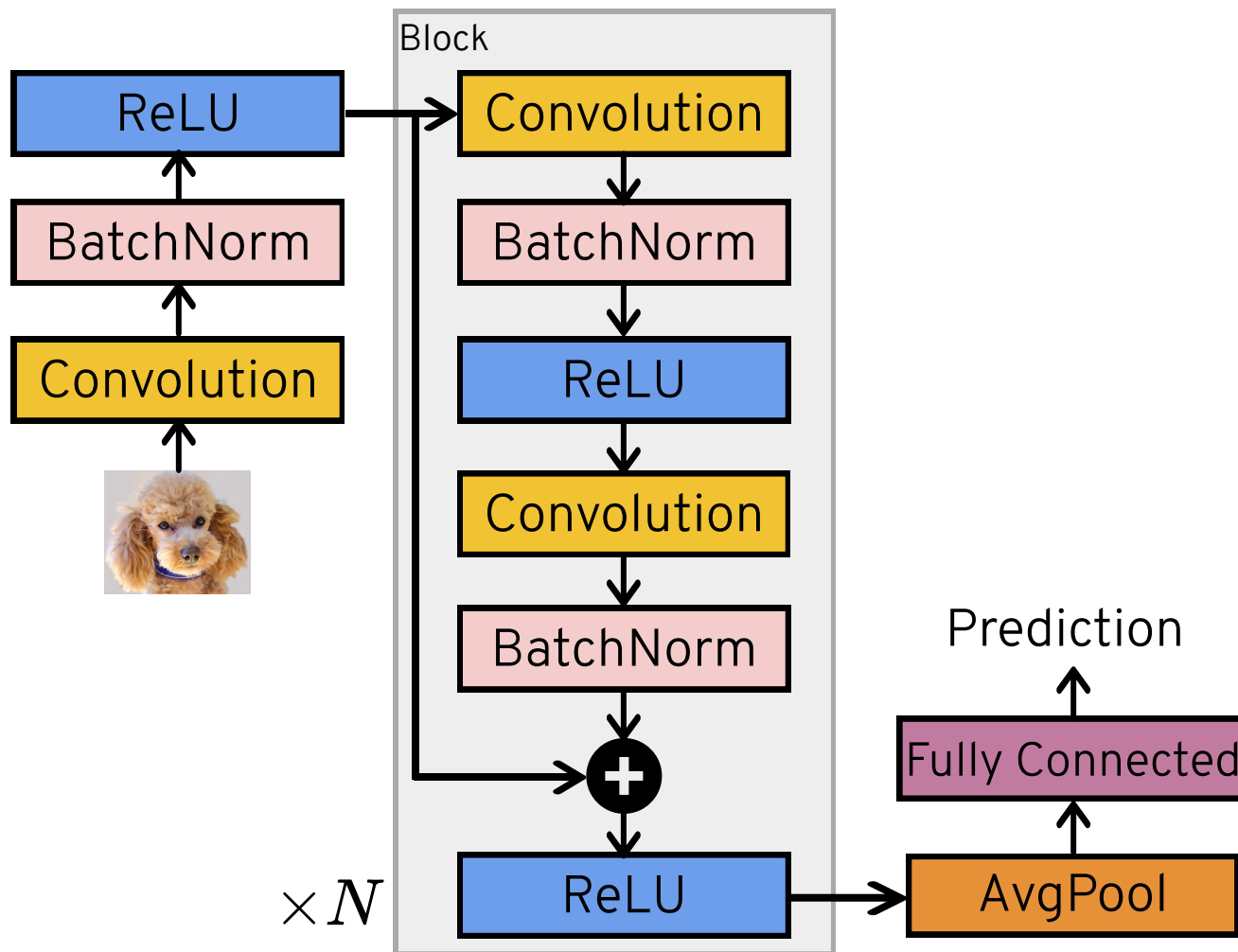
Secure DLaaS under Fully Homomorphic Encryption (FHE)



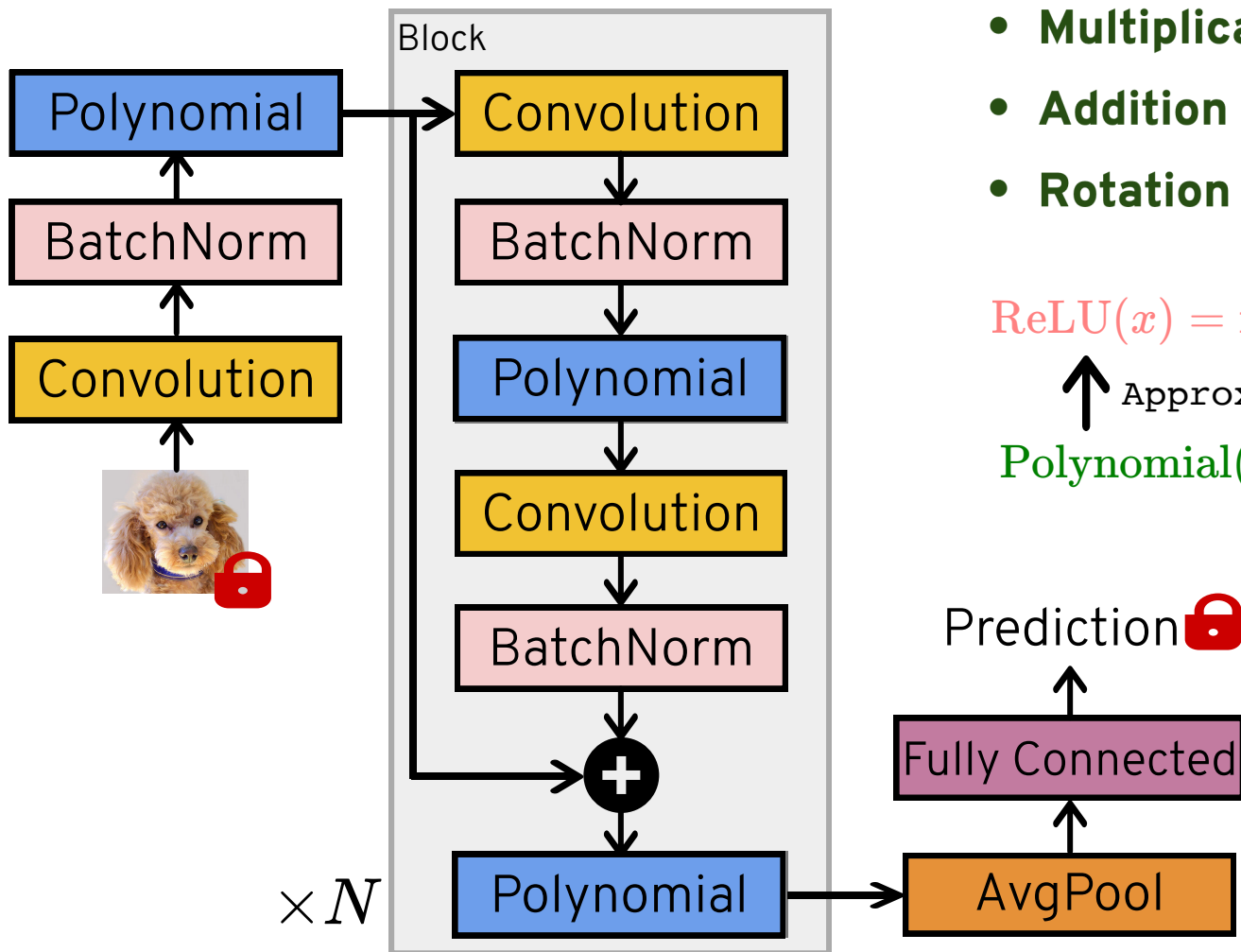
From **Secure Computation** to **Secure Deep Learning**



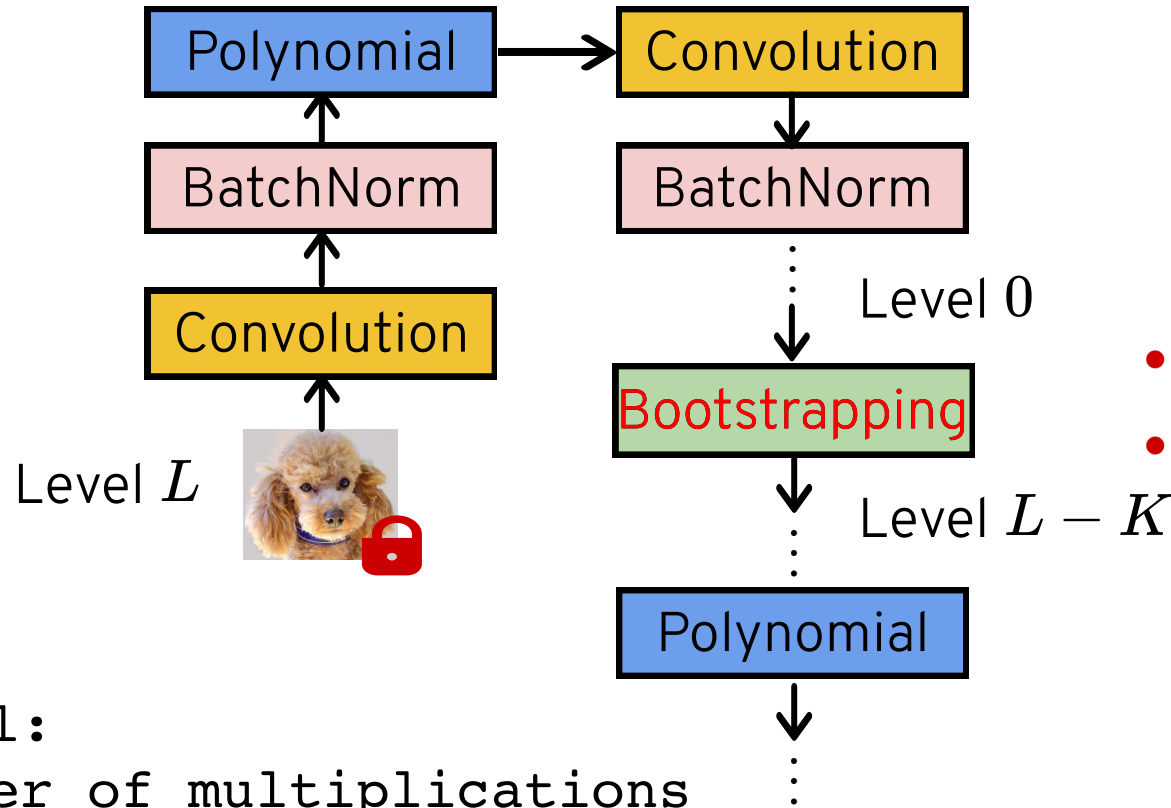
Convolutional Neural Networks (CNNs)



CNNs under Homomorphic Encryption (HE)



Deep CNNs under Fully Homomorphic Encryption (FHE)



- **High Latency**
- **High Memory Footprint**

Level:
number of multiplications
allowed to evaluate

Deep CNNs under Fully Homomorphic Encryption (FHE)

- Security Requirement
- Inference Latency
- Prediction Accuracy

Deep CNNs under Fully Homomorphic Encryption (FHE)

- Security Requirement
- Inference Latency
- Prediction Accuracy

Cryptographic Parameters

Cyclotomic polynomial degree N

Level L

Modulus $Q_\ell = \prod_{i=0}^{\ell} q_i, 0 \leq \ell \leq L$

Bootstrapping depth K

Hamming weight h

Deep CNNs under Fully Homomorphic Encryption (FHE)

- Security Requirement

- Inference Latency

- Prediction Accuracy

Cryptographic Parameters

Cyclotomic polynomial degree N

Level L

Modulus $Q_\ell = \prod_{i=0}^{\ell} q_i, 0 \leq \ell \leq L$

Bootstrapping depth K

Hamming weight h

Polynomial CNNs

Conv, BN, pooling, FC layers: packing

Polynomials: degree -> depth

Number of layers: ResNet20, ResNet32

Input image resolution

Channels/kernels

Deep CNNs under Fully Homomorphic Encryption (FHE)

- Security Requirement

- Inference Latency

- Prediction Accuracy

Cryptographic Parameters

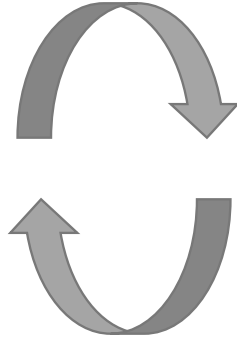
Cyclotomic polynomial degree N

Level L

Modulus $Q_\ell = \prod_{i=0}^{\ell} q_i, 0 \leq \ell \leq L$

Bootstrapping depth K

Hamming weight h



Polynomial CNNs

Conv, BN, pooling, FC layers: packing

Polynomials: degree -> depth

Number of layers: ResNet20, ResNet32

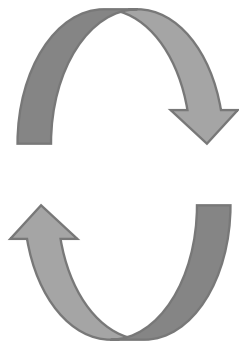
Input image resolution

Channels/kernels

Hand-crafted Design of Polynomial for CNNs under FHE

Cryptographic Parameters

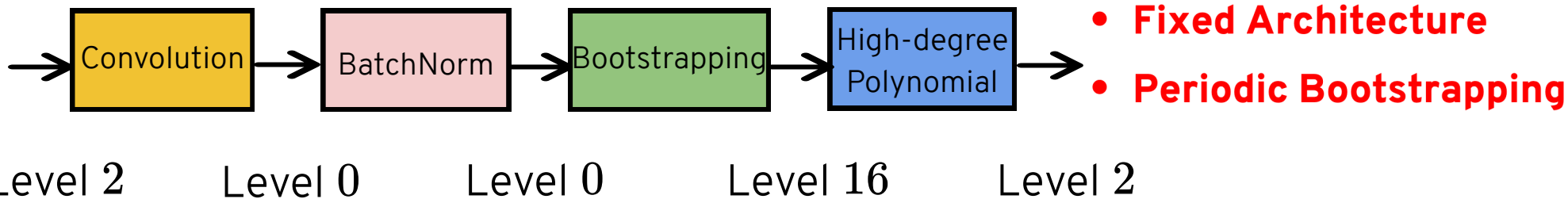
$$N, L, Q_\ell = \prod_{i=0}^{\ell} q_i (0 \leq \ell \leq L), K, h$$



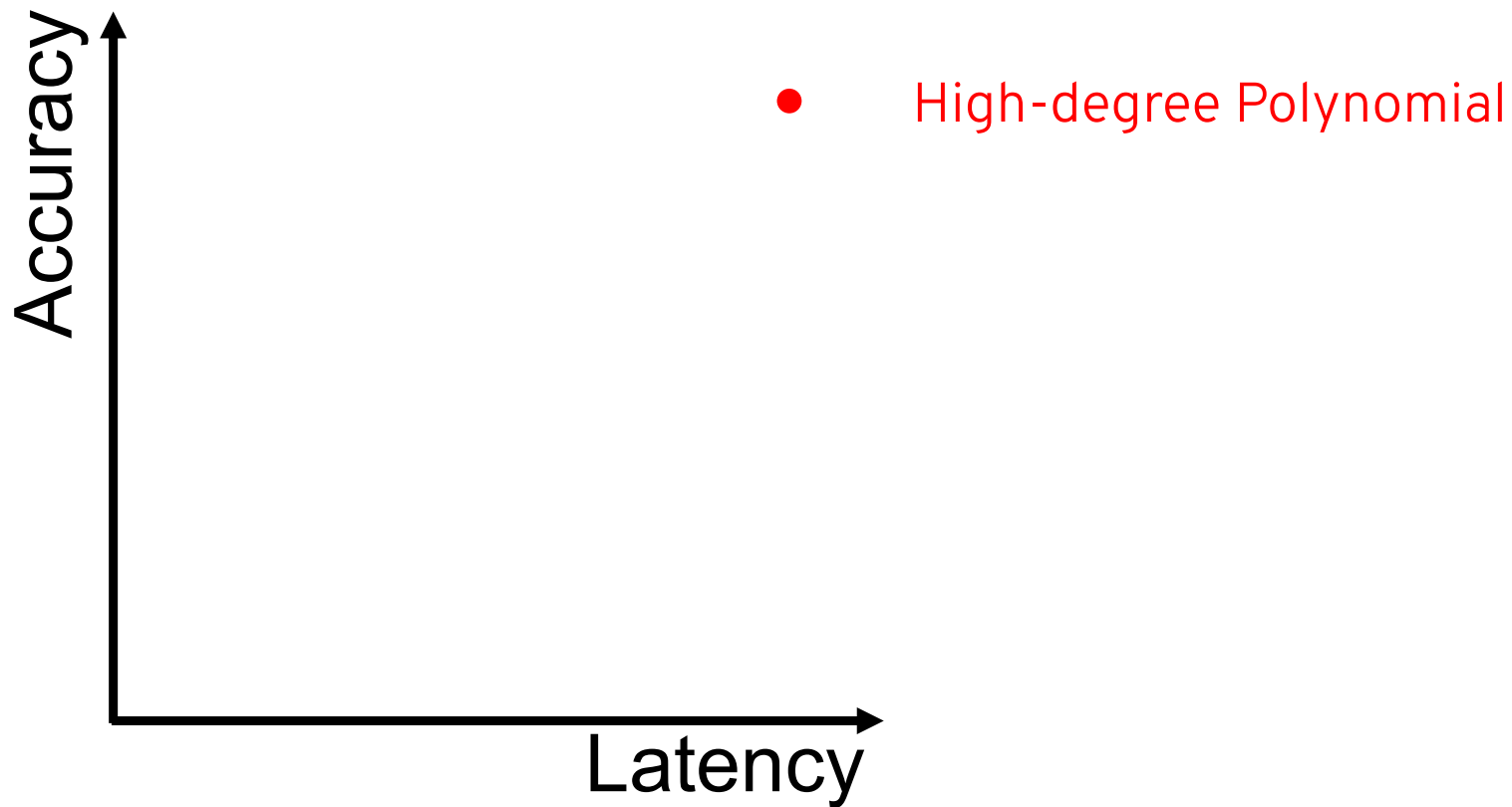
Polynomial CNNs

Polynomials: degree \rightarrow depth

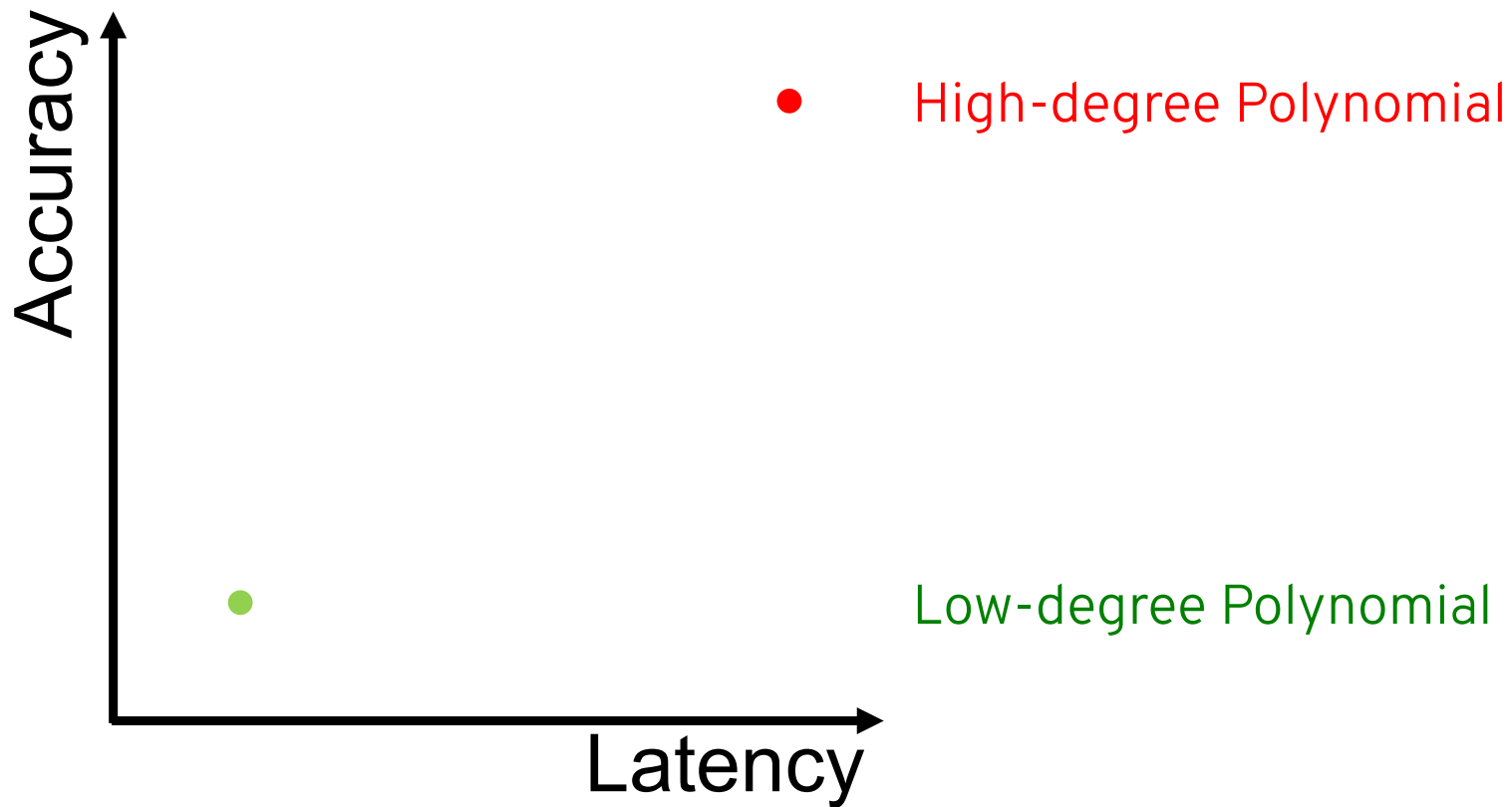
MPCNN [1]:



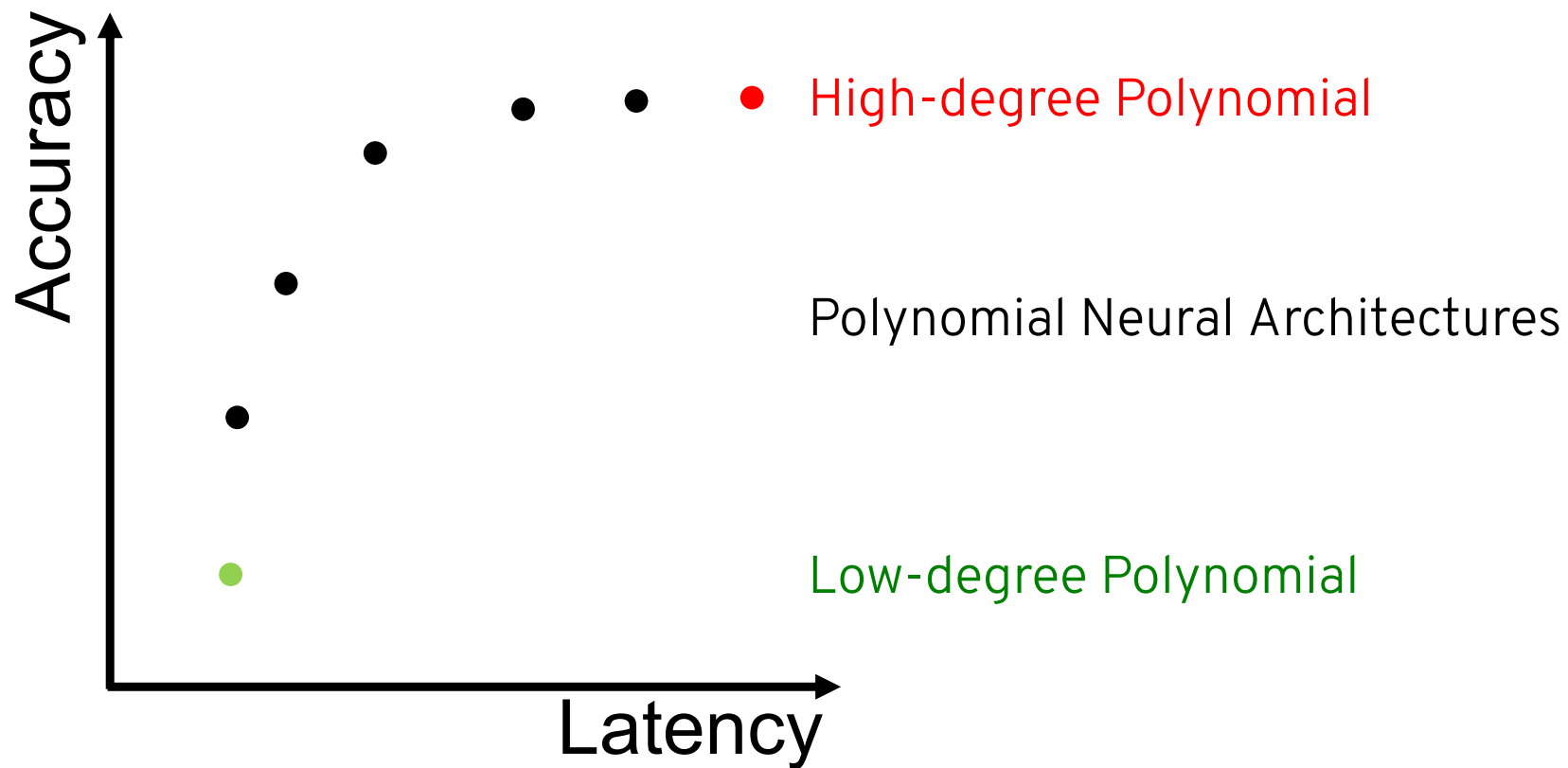
Hand-crafted Design of Polynomial for CNNs under FHE



Hand-crafted Design of Polynomial for CNNs under FHE



Hand-crafted Design of Polynomial for CNNs under FHE

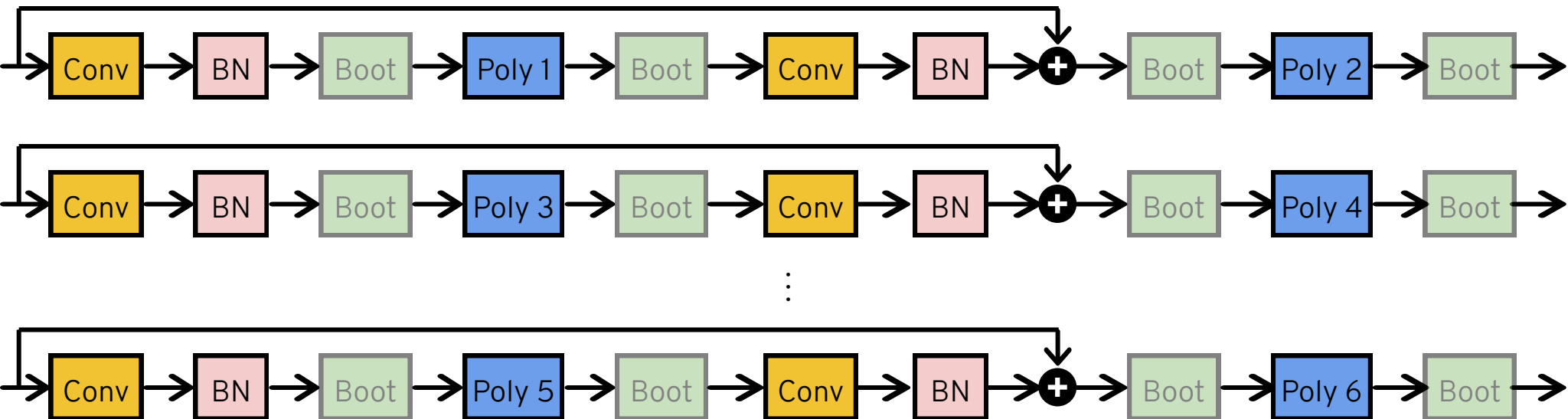


How to obtain all possible polynomial neural architectures?

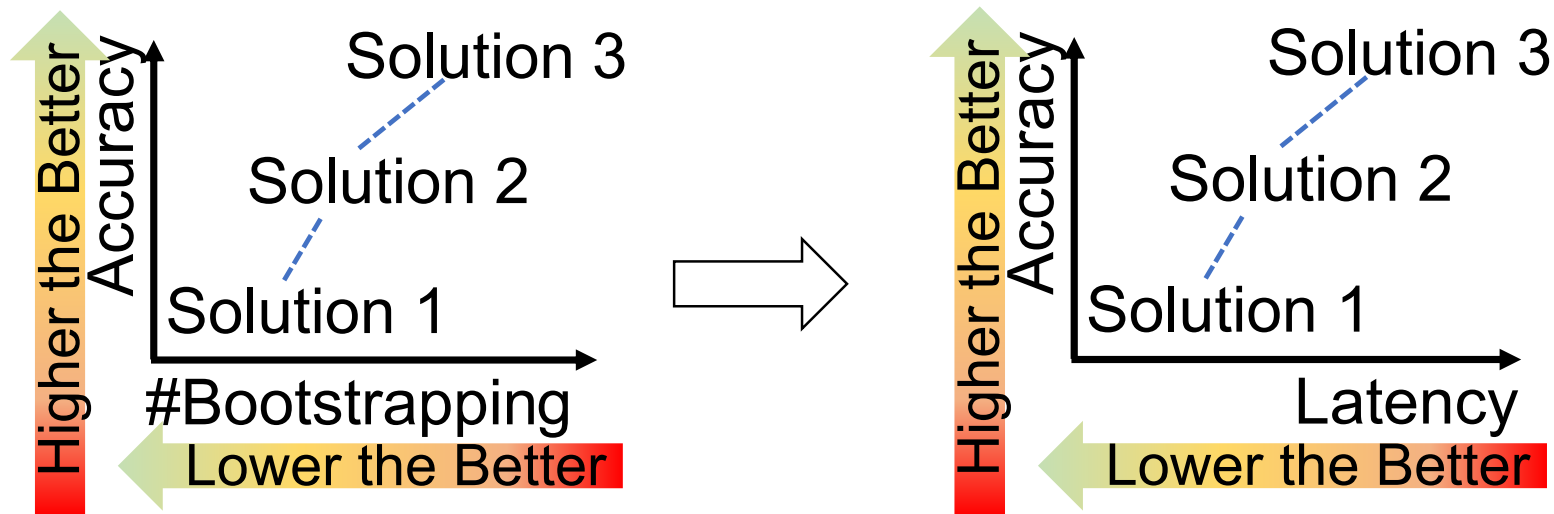
Key Insight

Optimize the
end-to-end polynomial neural **architecture**
instead of the polynomial function

Optimization of End-to-End Polynomial Neural Architecture



Optimization of End-to-End Polynomial Neural Architecture

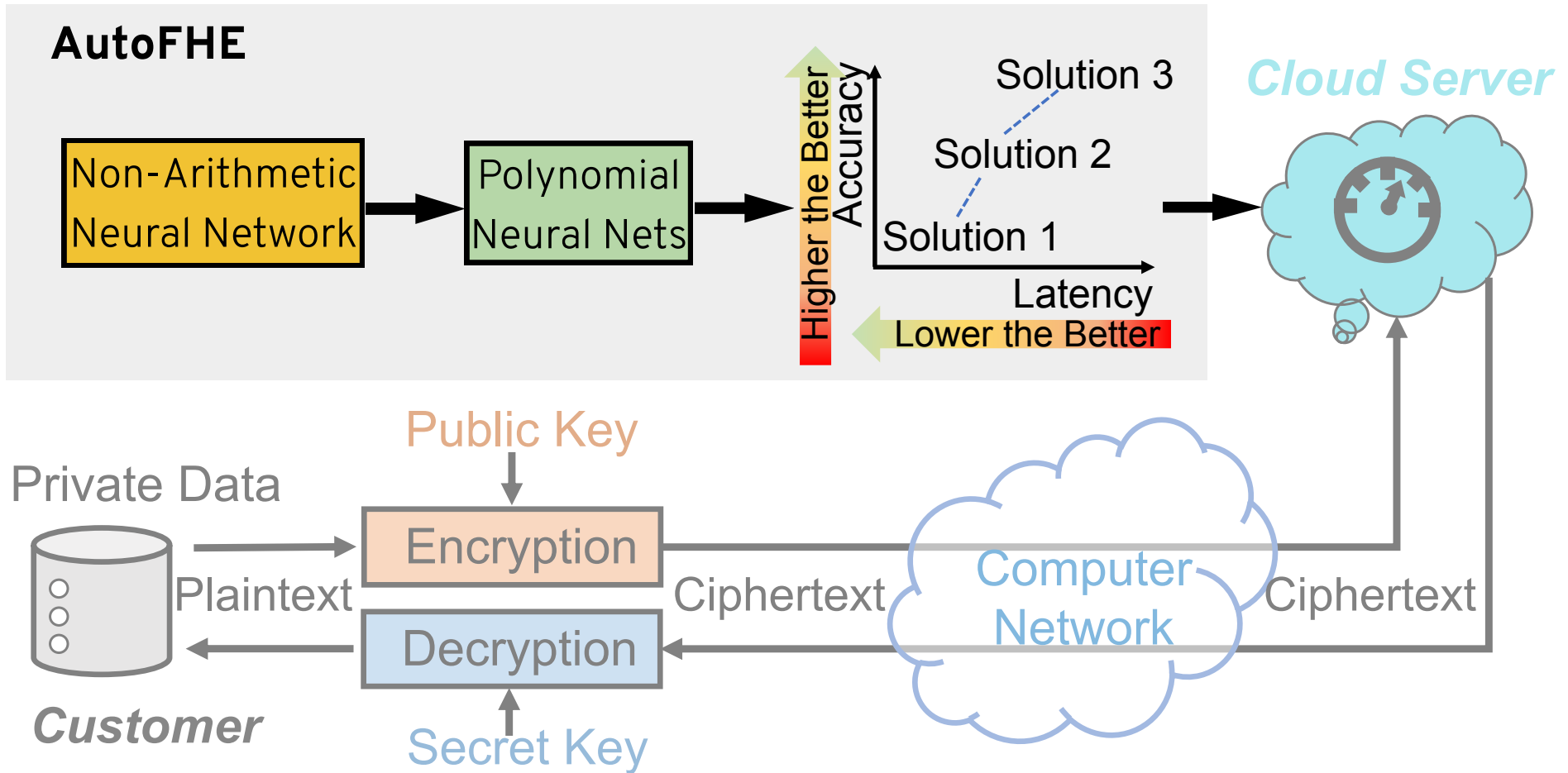


To meet different requirements in real world



- I want a faster response
- I can wait for an accurate result

AutoFHE: Automated Adaption of CNNs under FHE



EvoReLU: Evolutionary Mixed-Degree Polynomial Approximation of ReLU

Forward Propagation

$$\text{EvoReLU}(x) = \begin{cases} x, & d = 1 \\ \alpha_2 x^2 + \alpha_1 x + \alpha_0, & d = 2 \\ x \cdot (\mathcal{F}(x) + 0.5), & d > 2 \end{cases}$$

High-degree composite polynomial [2]:

$$\mathcal{F}(x) = (f_K^{d_K} \circ \dots \circ f_k^{d_k} \circ \dots \circ f_1^{d_1})(x), 1 \leq k \leq K$$

- **Pruning:** DeepReDuce, SAFENet, Delphi
- **Quadratic:** LoLa, CryptoNets, HEMET
- **High-degree approximation:** MPCNN

 Differentiable Evolution

[2] Lee, Eunsang, Joon-Woo Lee, Jong-Seon No, and Young-Sik Kim. "Minimax approximation of sign function by composite polynomial for homomorphic comparison."

EvoReLU: Evolutionary Mixed-Degree Polynomial Approximation of ReLU

Backward Propagation

$$\frac{\partial \text{EvoReLU}(x)}{\partial x} = \begin{cases} 1, & d = 1 \\ 2\alpha_2 x + \alpha_1, & d = 2 \\ \partial \text{ReLU}(x) / \partial x, & d > 2 \end{cases}$$

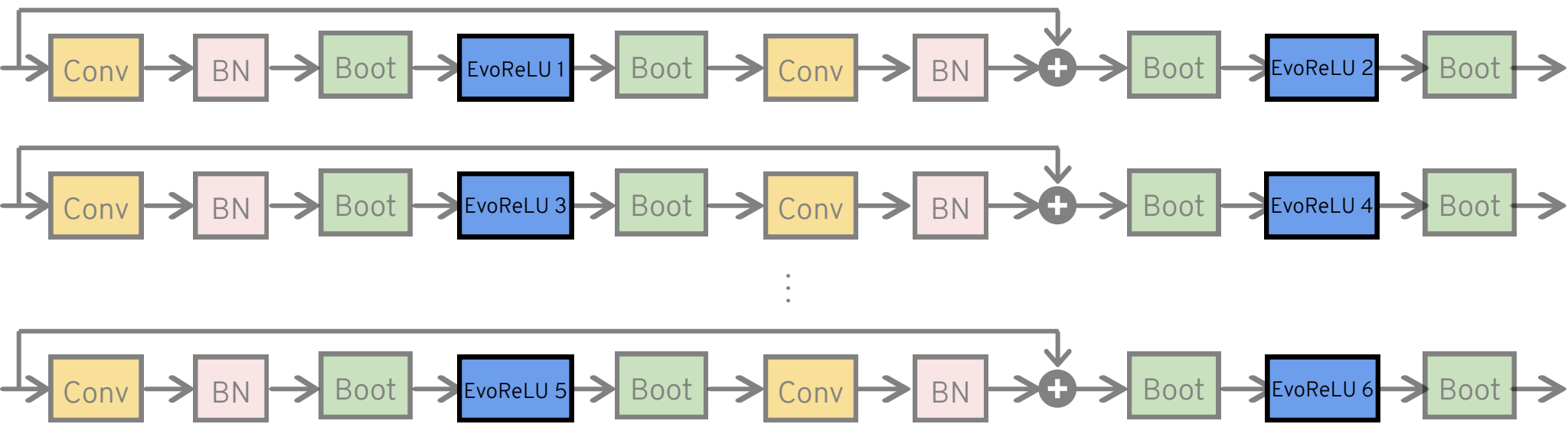
- Gradient
- Gradient
- Straight-through estimated gradient

- **Make training more stable**

How to **optimize** end-to-end
polynomial neural architecture?

Multi-Objective evolutionary optimization

Joint Search for **Layerwise EvoReLU** and **Bootstrapping Operations**



Joint search
problem



Multi-objective
optimization



- **Flexible Architecture**
- **On-demand Bootstrapping**

Multi-Objective Optimization

Single Objective

- Accuracy
- Latency

Scalarization of Multiple Objectives

$$\alpha \cdot \text{Accuracy} + \beta \cdot \text{Latency}$$

- Only generate a single solution
- Hard to tune balancing weights
- Not Pareto optimal

Multi-Objective Optimization

$$\min \{1 - \text{Accuracy}, \# \text{Bootstrapping}\}$$

- Multiple solutions on the Pareto front
- No need to tune weights
- Pareto optimal

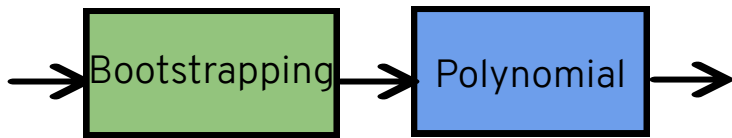
Multi-Objective Optimization

Multi-Objective Optimization

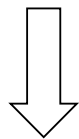
$\min \{1 - \text{Accuracy}, \text{Depth of polys}\}$

Level 4

Depth 9



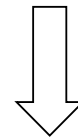
Drop 4 Levels



- Not necessarily reduce bootstrapping operations

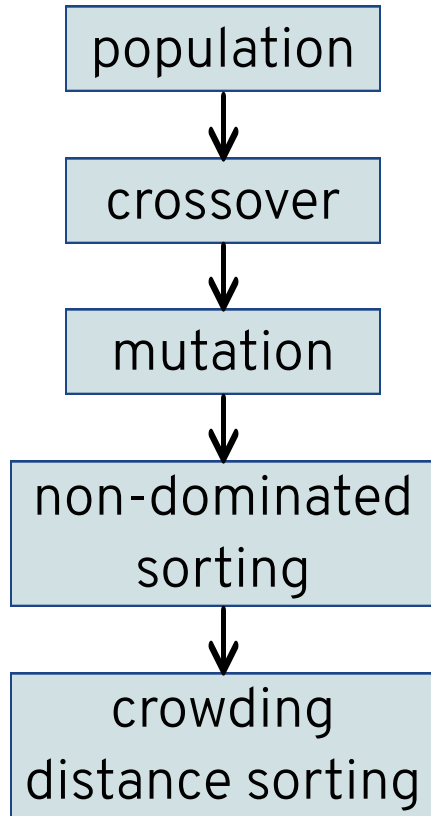
Multi-Objective Optimization

$\min \{1 - \text{Accuracy}, \#\text{Bootstrapping}\}$

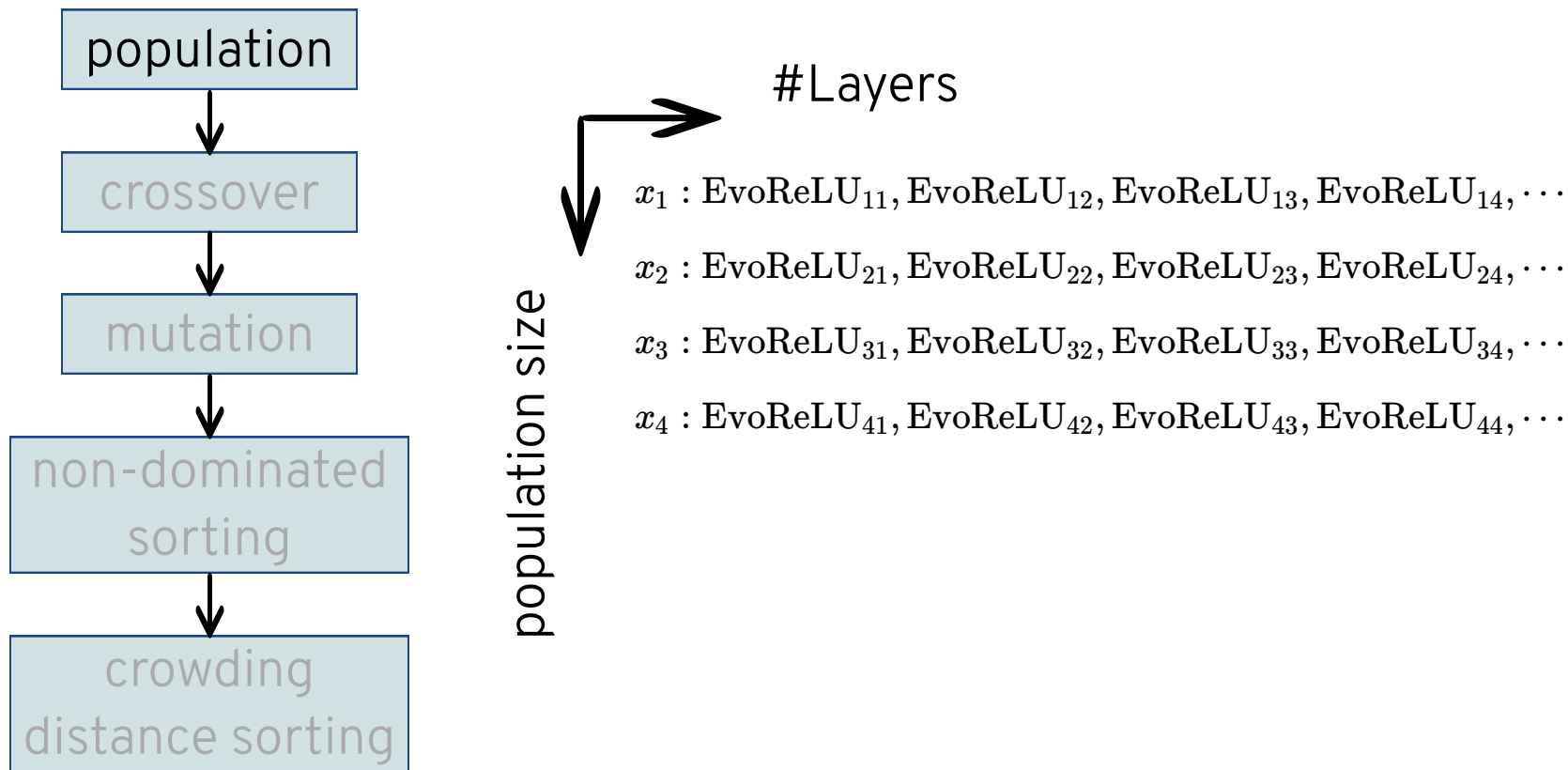


- Directly reduce bootstrapping operations

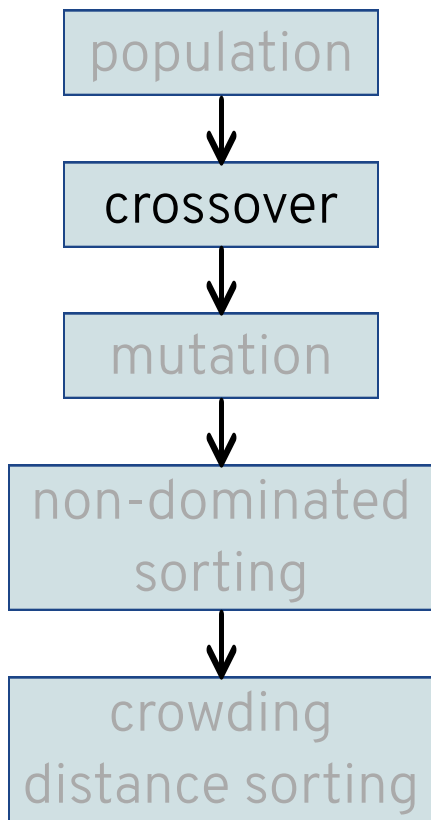
Evolutionary Multi-Objective Optimization



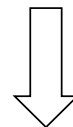
Evolutionary Multi-Objective Optimization



Evolutionary Multi-Objective Optimization

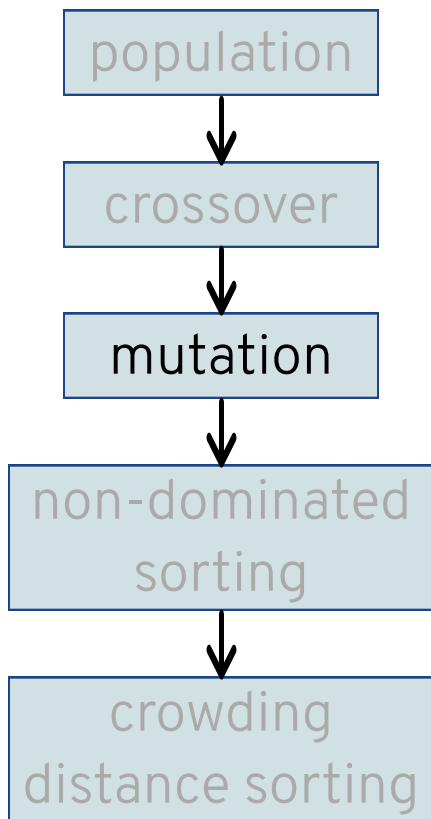


x_1 : EvoReLU₁₁, EvoReLU₁₂, EvoReLU₁₃, EvoReLU₁₄, ...
 x_2 : EvoReLU₂₁, EvoReLU₂₂, EvoReLU₂₃, EvoReLU₂₄, ...



x'_1 : EvoReLU₂₁, EvoReLU₁₂, EvoReLU₂₃, EvoReLU₁₄, ...
 x'_2 : EvoReLU₁₁, EvoReLU₂₂, EvoReLU₁₃, EvoReLU₂₄, ...

Evolutionary Multi-Objective Optimization



$x_1 : \text{EvoReLU}_{11}, \text{EvoReLU}_{12}, \text{EvoReLU}_{13}, \text{EvoReLU}_{14}, \dots$

$x_2 : \text{EvoReLU}_{21}, \text{EvoReLU}_{22}, \text{EvoReLU}_{23}, \text{EvoReLU}_{24}, \dots$

$x_3 : \text{EvoReLU}_{31}, \text{EvoReLU}_{32}, \text{EvoReLU}_{33}, \text{EvoReLU}_{34}, \dots$

$x_4 : \text{EvoReLU}_{41}, \text{EvoReLU}_{42}, \text{EvoReLU}_{43}, \text{EvoReLU}_{44}, \dots$



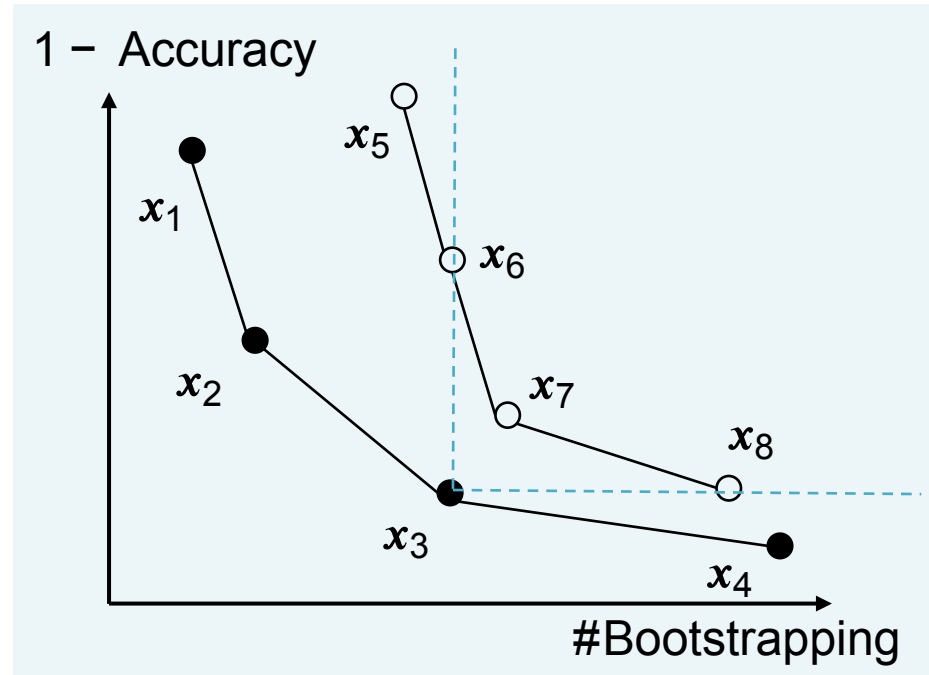
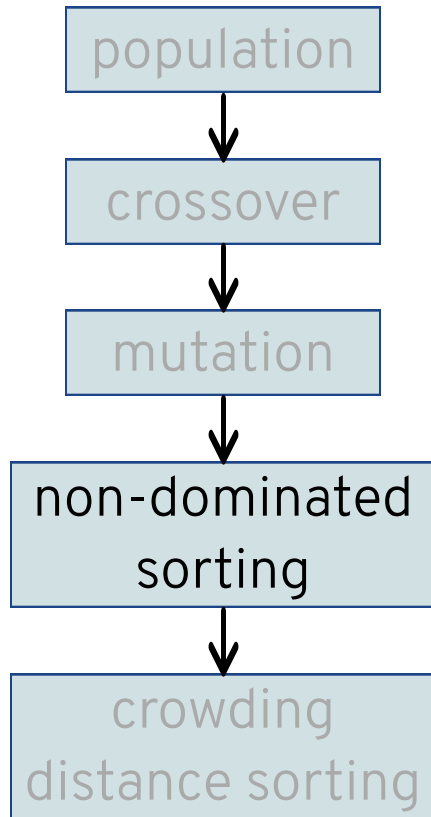
$x'_1 : \text{EvoReLU}_{11}, \text{EvoReLU}_{12}, \text{EvoReLU}'_{13}, \text{EvoReLU}_{14}, \dots$

$x'_2 : \text{EvoReLU}_{21}, \text{EvoReLU}_{22}, \text{EvoReLU}'_{23}, \text{EvoReLU}_{24}, \dots$

$x'_3 : \text{EvoReLU}'_{31}, \text{EvoReLU}_{32}, \text{EvoReLU}_{33}, \text{EvoReLU}_{34}, \dots$

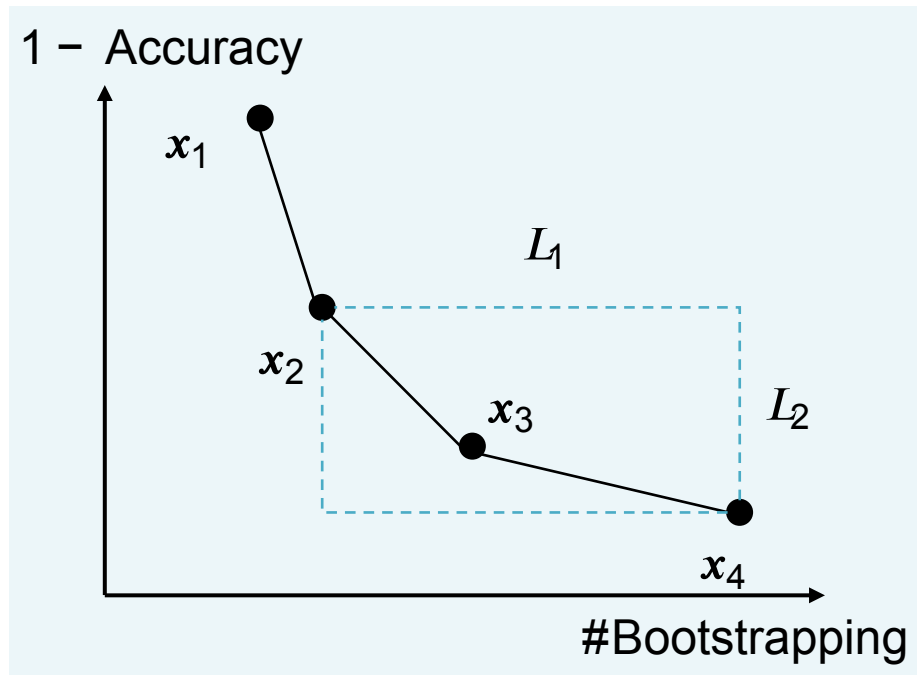
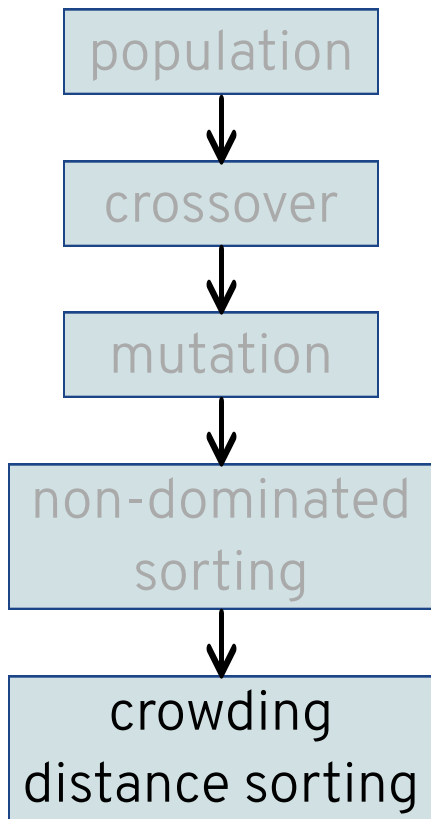
$x'_4 : \text{EvoReLU}_{41}, \text{EvoReLU}'_{42}, \text{EvoReLU}_{43}, \text{EvoReLU}_{44}, \dots$

Evolutionary Multi-Objective Optimization



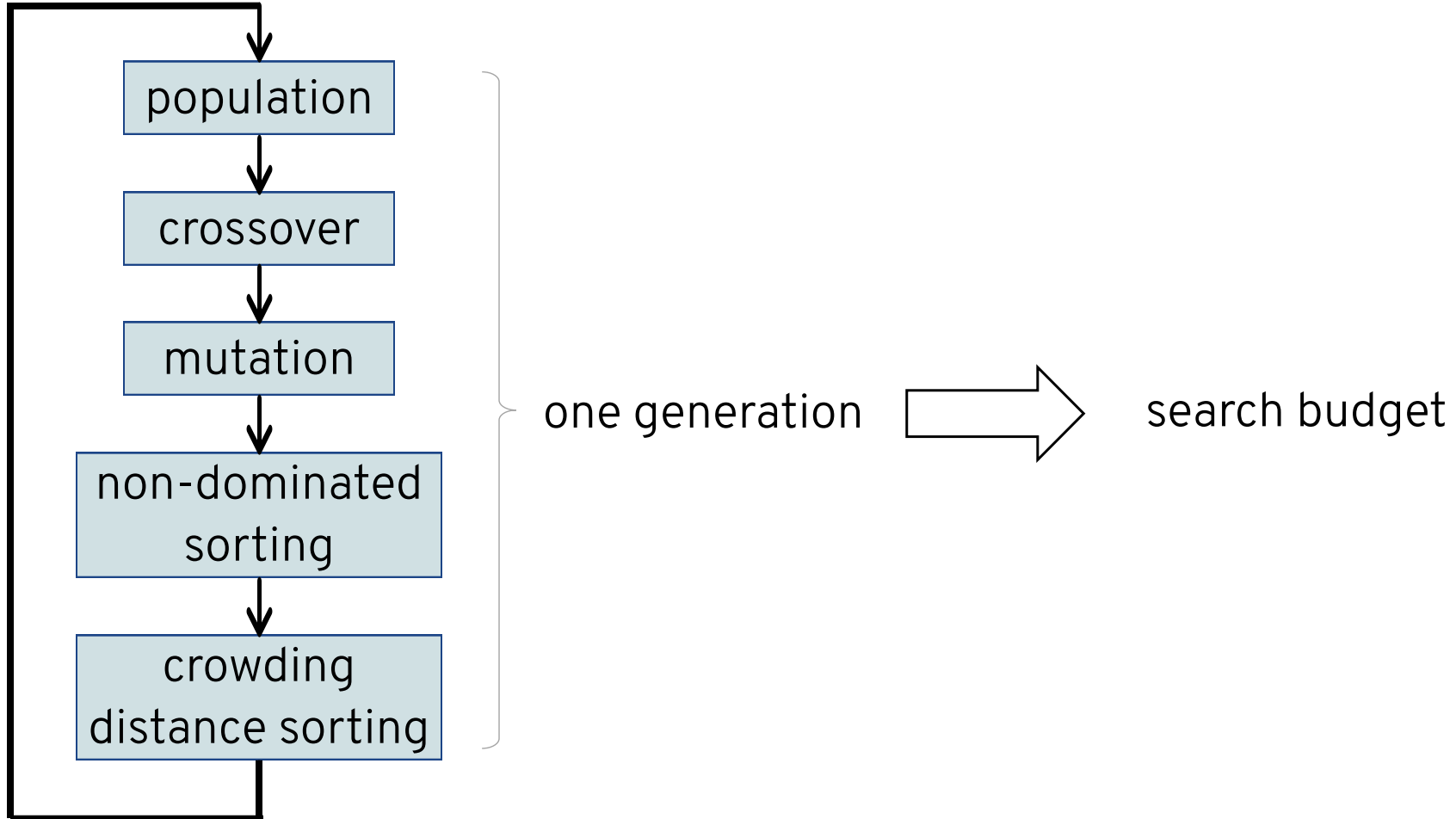
x_3 dominates x_6 , x_7 , and x_8
i.e. x_3 is better than x_6 , x_7 , and x_8

Evolutionary Multi-Objective Optimization



crowding distance of x_3 is $\frac{L_1 + L_2}{2}$

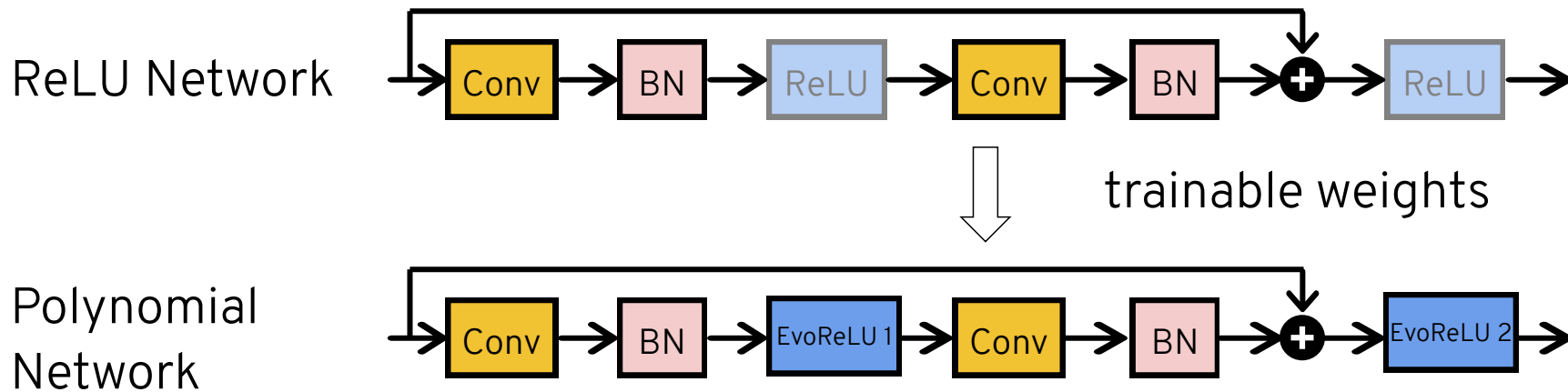
Evolutionary Multi-Objective Optimization



How to **fine-tune** polynomial CNNs?

Neural network adaption

Trainable Weight Adaption and Knowledge Transferring



Fine-tuning objective

$$\mathcal{L}_{train} = (1 - \tau)\mathcal{L}_{CE} + \tau\mathcal{L}_{KL}$$

- Inherit representation learning ability
- Adapt trainable weights to EvoReLU

Experiments on encrypted CIFAR10 dataset under FHE

Dataset: CIFAR10

50,000 training images

10,000 test images

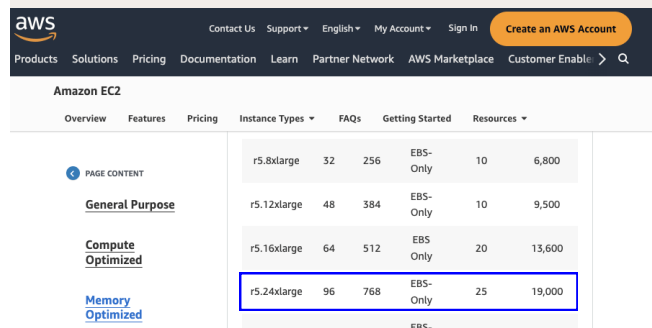
32x32 resolution, 10 classes

Hardware & Software

Amazon AWS, r5.24xlarge

96 CPUs, 768 GB RAM

Microsoft SEAL, 3.6



Amazon EC2

Instance Type	Number of Instances	Price per Hour	On-Demand Price per Hour	On-Demand Price per Month	
r5.xlarge	32	256	EBS-Only	10	6,800
r5.12xlarge	48	384	EBS-Only	10	9,500
r5.16xlarge	64	512	EBS-Only	20	13,600
r5.24xlarge	96	768	EBS-Only	25	19,000

Experimental Setup

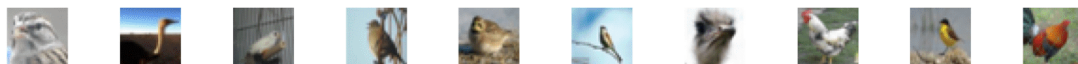
plane



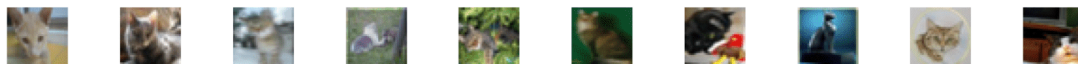
auto



bird



cat



deer



dog



frog



horse



ship



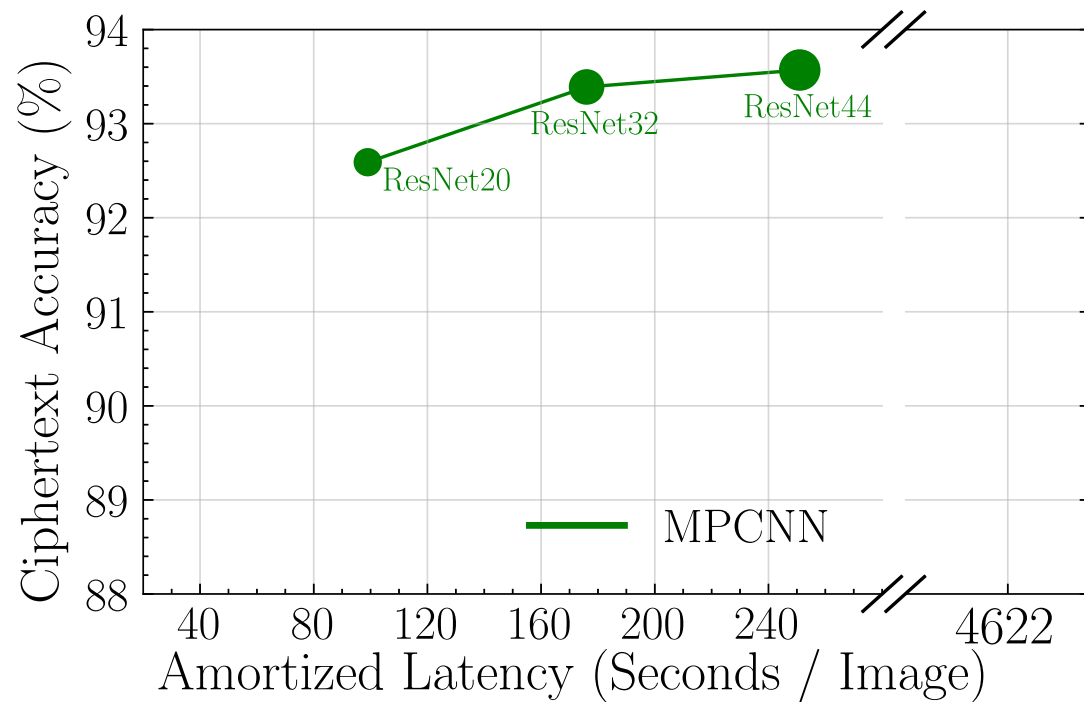
truck



[3]Alex Krizhevsky. CIFAR example images (online).

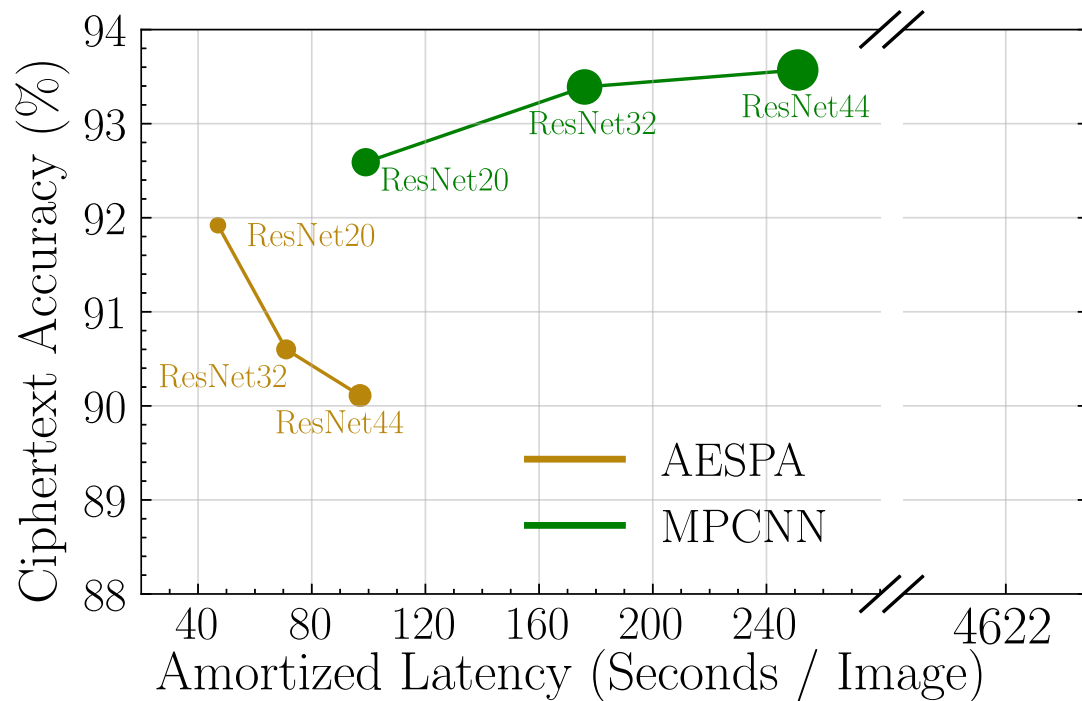
Latency and Accuracy Trade-offs under FHE

Approach	MPCNN
Venue	ICML22
Scheme	CKKS
Polynomial	high-degree
Layerwise	no
Strategy	approx
Arch	manual



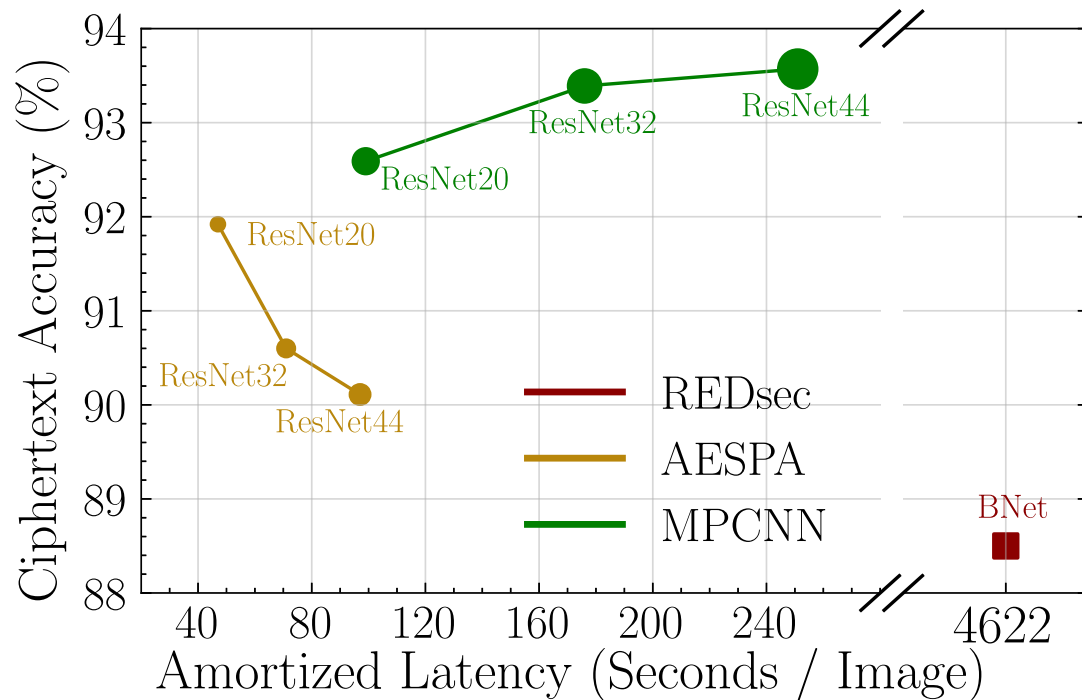
Latency and Accuracy Trade-offs under FHE

Approach	MPCNN	AESPA
Venue	ICML22	arXiv22
Scheme	CKKS	CKKS
Polynomial	high-degree	low-degree
Layerwise	no	no
Strategy	approx	train
Arch	manual	manual



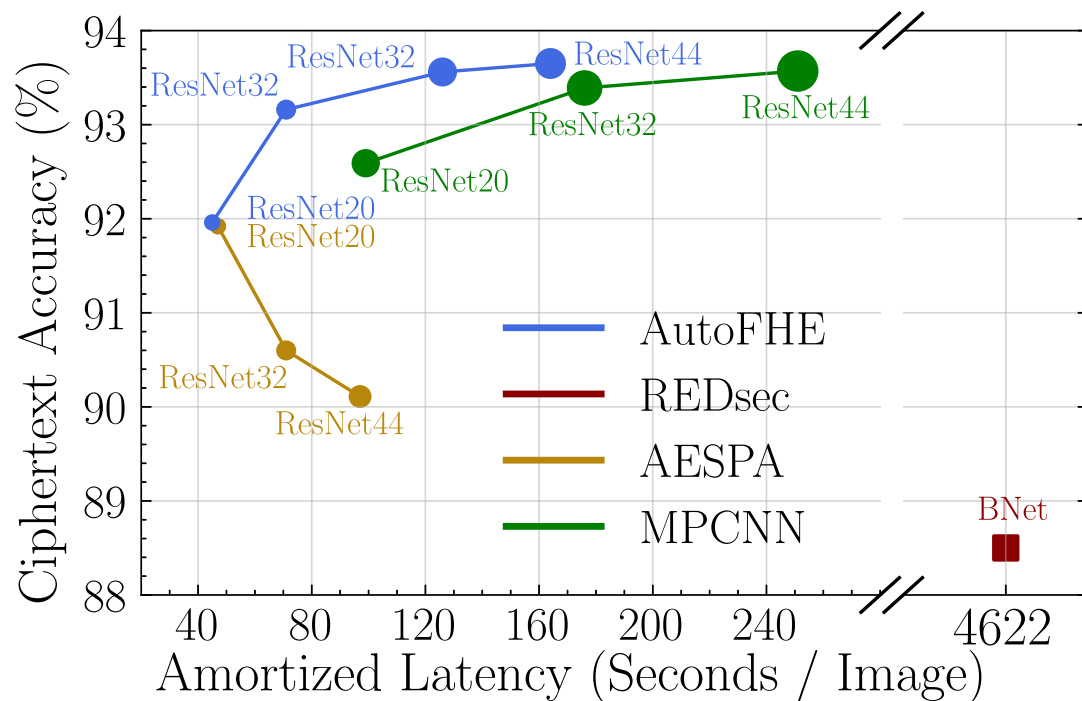
Latency and Accuracy Trade-offs under FHE

Approach	MPCNN	AESPA	REDsec
Venue	ICML22	arXiv22	NDSS23
Scheme	CKKS	CKKS	TFHE
Polynomial	high-degree	low-degree	n/a
Layerwise	no	no	n/a
Strategy	approx	train	train
Arch	manual	manual	manual

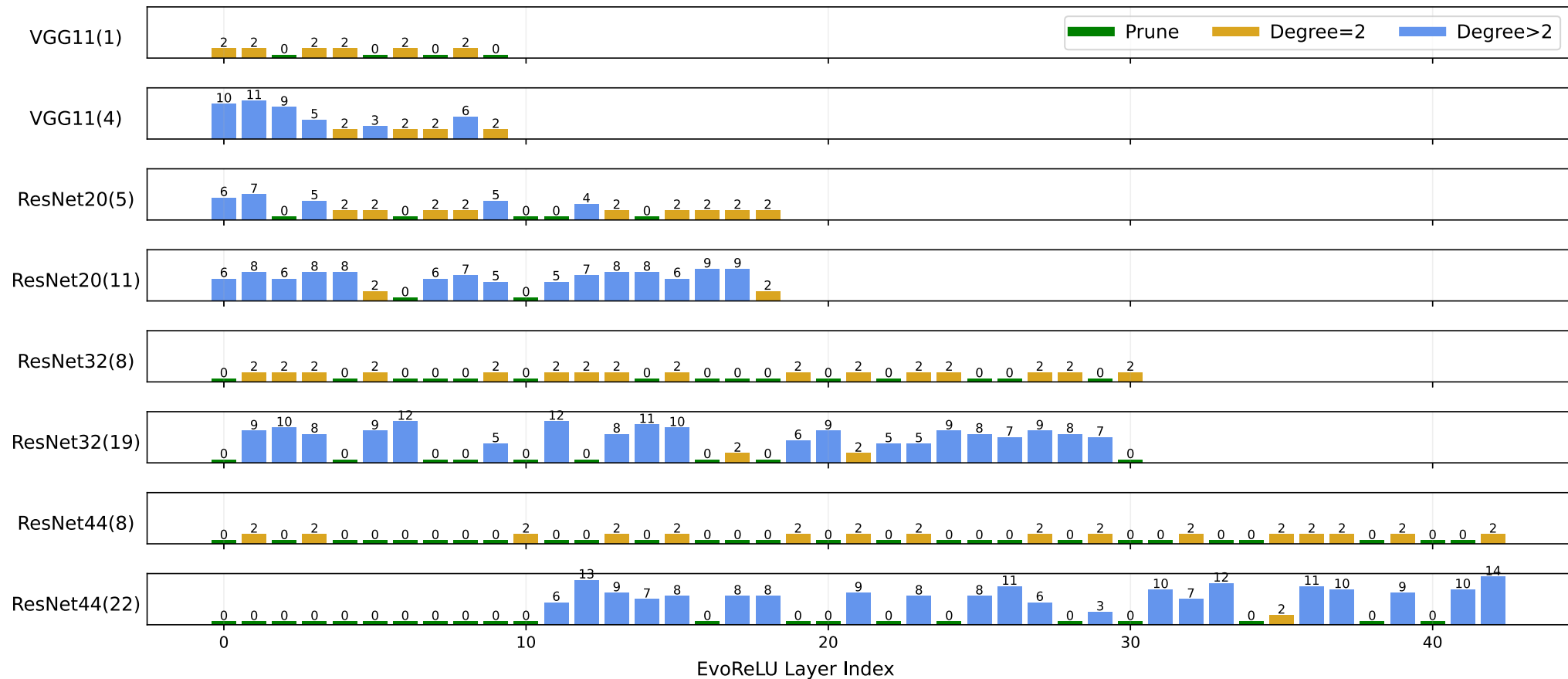


Latency and Accuracy Trade-offs under FHE

Approach	MPCNN	AESPA	REDsec	AutoFHE
Venue	ICML22	arXiv22	NDSS23	USENIX24
Scheme	CKKS	CKKS	TFHE	CKKS
Polynomial	high-degree	low-degree	n/a	mixed
Layerwise	no	no	n/a	yes
Strategy	approx	train	train	adapt
Arch	manual	manual	manual	search

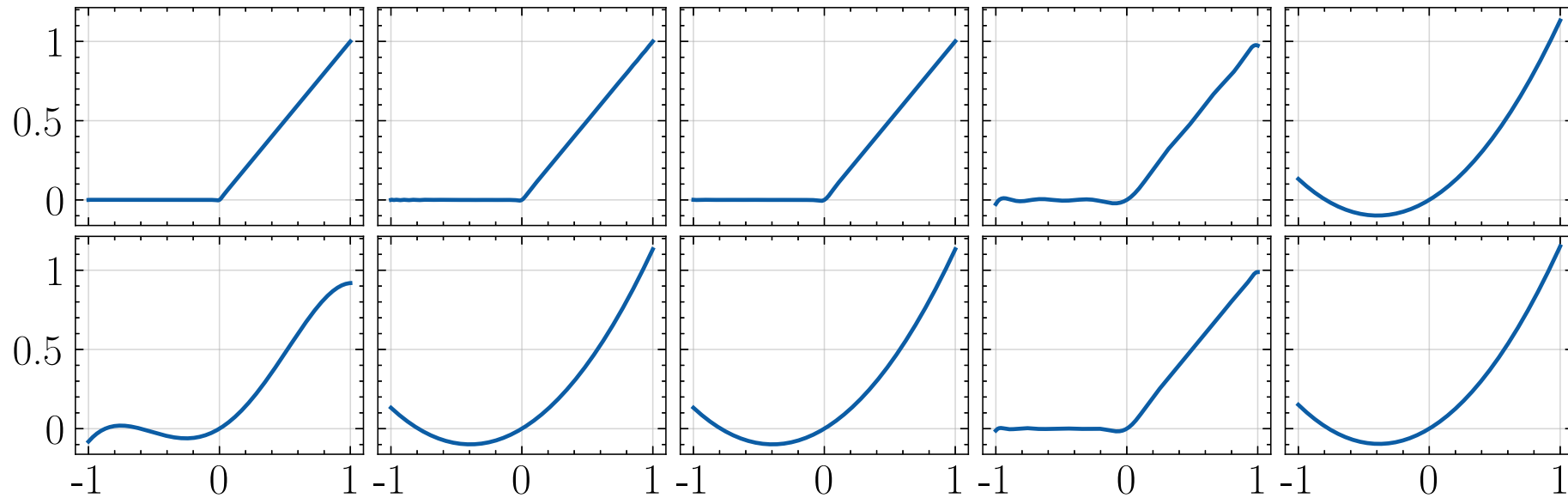


Multiplicative Depth of Layerwise EvoReLU

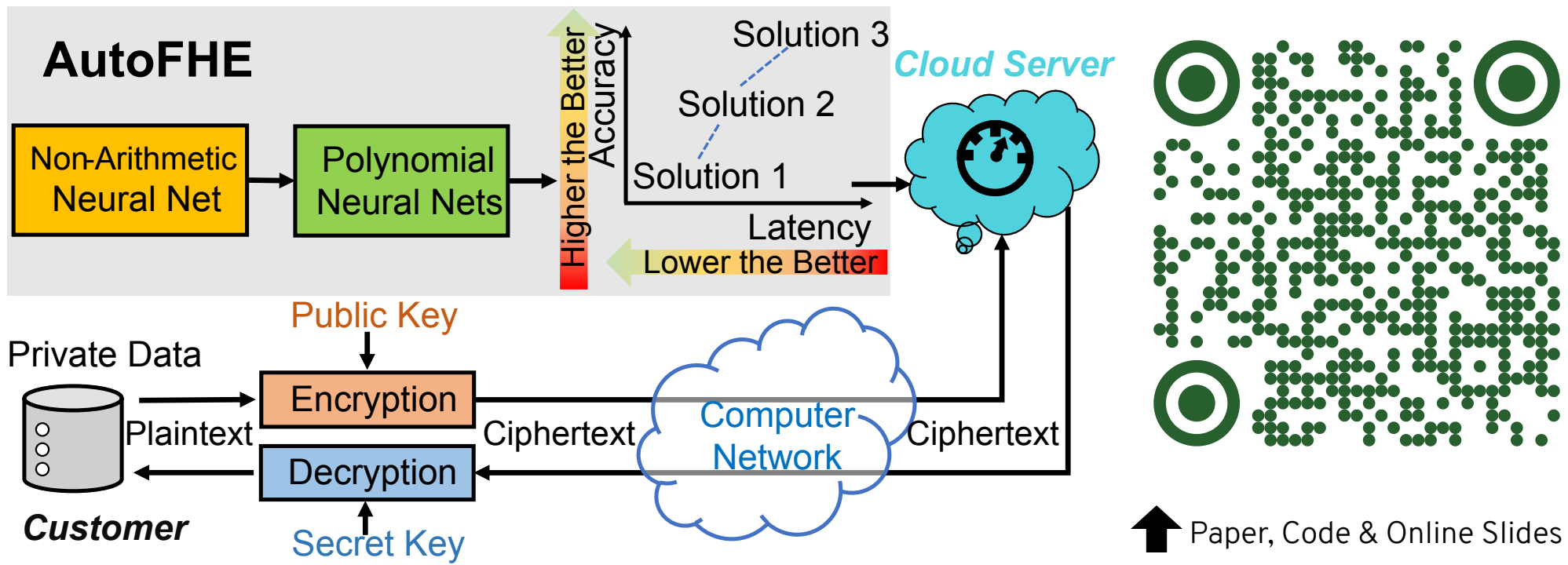


Layerwise EvoReLU

VGG11(4)



Conclusion



AutoFHE optimizes end-to-end polynomial neural architecture

- Multi-objective optimization generates Pareto-effective solutions to meet different requirements
- Joint optimization of layerwise EvoReLU and bootstrapping results in optimal polynomial neural architectures
- Adapted neural networks can inherit representation learning ability from ReLU networks