

# A Robust Sampling for Optimization in Deep Learning

## Master thesis

Aurora Cobo Aguilera

Tutor: Antonio Artés Rodríguez  
University Master in Multimedia and Communications  
Signal Theory and Communications Department  
Universidad Carlos III de Madrid

Leganés, 26<sup>th</sup> July 2018

## 1. Introduction

# Introduction

**A Robust Sampling for Optimization in Deep Learning**

## 1. Introduction

# Introduction

**A Robust Sampling for Optimization in Deep Learning**



## 1. Introduction

# Introduction

**A Robust Sampling for Optimization in Deep Learning**

## 1. Introduction

# Introduction

**A Robust Sampling for Optimization in Deep Learning**

**(Namkoong & Duchi, 2017)**



**Regularized estimation  
based on minimizing the  
variance**

## 1. Introduction

# Introduction

A Robust Sampling for **Optimization** in Deep Learning

## 1. Introduction

# Introduction

A Robust Sampling for **Optimization** in Deep Learning

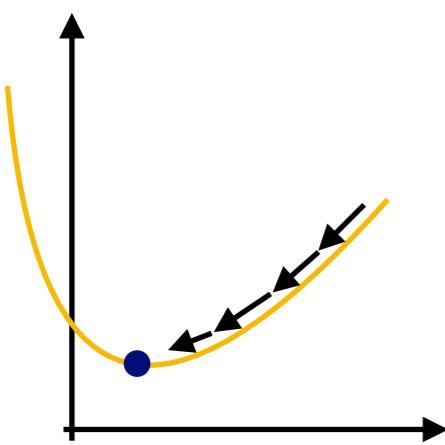
**Empirical Risk Minimization**

$$\min_{\theta}$$

## 1. Introduction

# Introduction

A Robust Sampling for **Optimization** in Deep Learning



## Empirical Risk Minimization

$$\min_{\theta} R_{\text{emp}}(\theta) = \frac{1}{n} \sum_{i=1}^n Q(x_i, \theta)$$

## 1. Introduction

# Introduction

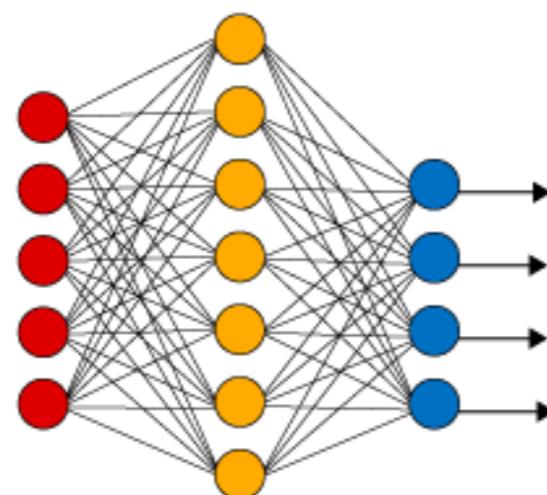
**A Robust Sampling for Optimization in Deep Learning**

## 1. Introduction

# Introduction

## A Robust Sampling for Optimization in Deep Learning

**Simple Neural Network**

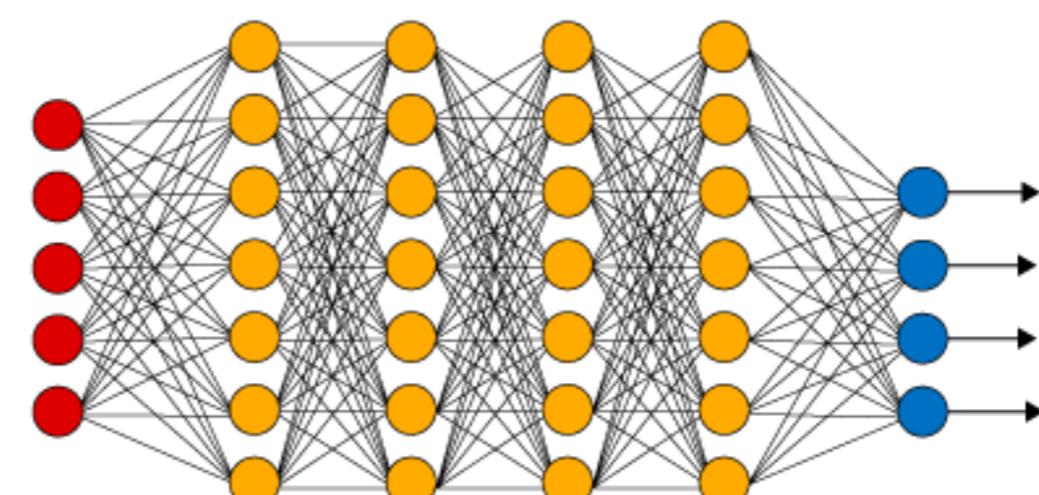


● Input Layer

○ Hidden Layer

● Output Layer

**Deep Learning Neural Network**



## 1. Introduction

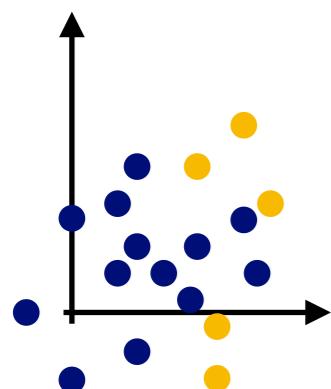
# Introduction

A Robust **Sampling** for Optimization in Deep Learning

## 1. Introduction

# Introduction

A Robust **Sampling** for Optimization in Deep Learning

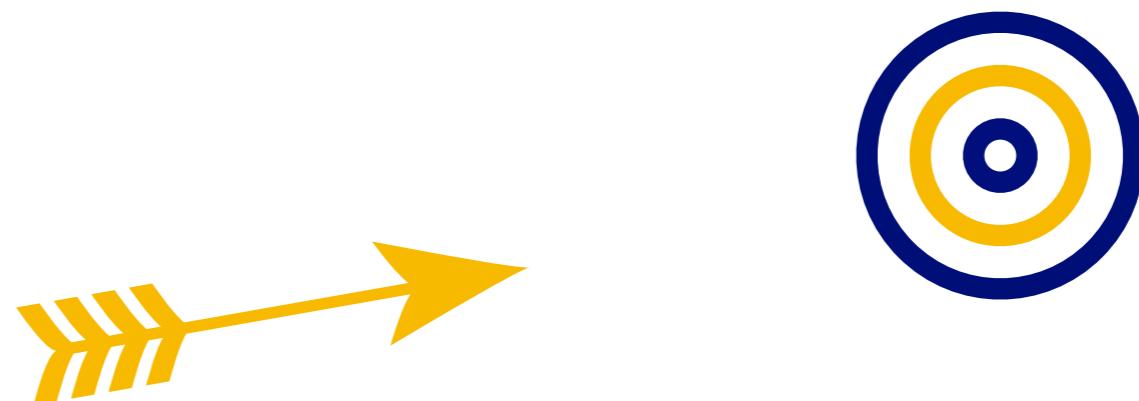


Select the samples of the mini-batch

## 1. Introduction

# Objectives

- 1. Improve the performance**
- 2. A faster convergence**
- 3. Better explanation of outliers**



## 1. Introduction

# Objectives

- 1. Improve the performance**
- 2. A faster convergence**
- 3. Better explanation of outliers**

Supply an application to Deep Learning  
Supply the numerical evidence



## 2. Variance-based robust regularization

# Empirical Risk

What we would like to have

What we have



## 2. Variance-based robust regularization

# Empirical Risk

## What we would like to have

- Real distribution of the data

## What we have



## 2. Variance-based robust regularization

# Empirical Risk

## What we would like to have

- Real distribution of the data
- True risk

## What we have



## 2. Variance-based robust regularization

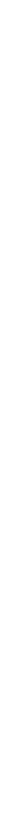
# Empirical Risk

## What we would like to have

- Real distribution of the data
- True risk

$$R(f) = \mathbb{E} \{Q(f(x))\}$$

## What we have



## 2. Variance-based robust regularization

# Empirical Risk

## What we would like to have

- Real distribution of the data
- True risk

$$R(f) = \mathbb{E} \{Q(f(x))\}$$

## What we have

- A set of samples  $\{x_i, y_i\}$



## 2. Variance-based robust regularization

# Empirical Risk

## What we would like to have

- Real distribution of the data
- True risk

$$R(f) = \mathbb{E} \{Q(f(x))\}$$

## What we have

- A set of samples  $\{x_i, y_i\}$
- Empirical risk

## 2. Variance-based robust regularization

# Empirical Risk

## What we would like to have

- Real distribution of the data
- True risk

$$R(f) = \mathbb{E} \{Q(f(x))\}$$

## What we have

- A set of samples  $\{x_i, y_i\}$
- Empirical risk

$$R_{emp}(f) = \frac{1}{n} \sum_{i=1}^n Q(f(x_i))$$

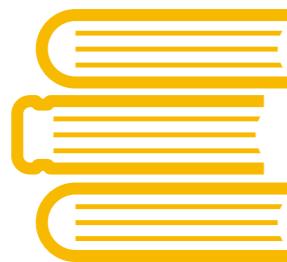
## 2. Variance-based robust regularization

# Empirical Risk Minimization

$$\min_f R_{emp}(f) = \frac{1}{n} \sum_{i=1}^n Q(f(x_i))$$

## 2. Variance-based robust regularization

# Robust Risk Minimization



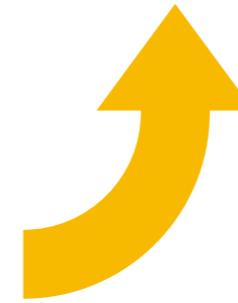
$$R(f) \leq R_{emp}(f) + \dots$$

## 2. Variance-based robust regularization

# Robust Risk Minimization



$$R(f) \leq R_{emp}(f) + \dots$$



## 2. Variance-based robust regularization

# Robust Risk Minimization



$$R(f) \leq R_{emp}(f) + \dots$$

$$\text{Var}[Q(f)]$$



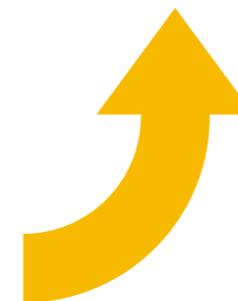
## 2. Variance-based robust regularization

# Robust Risk Minimization



$$R(f) \leq R_{emp}(f) + \dots$$

$$\text{Var}[Q(f)]$$



$$\min_f \left\{ \frac{1}{n} \sum_{i=1}^n Q(f(x_i)) \right\}$$

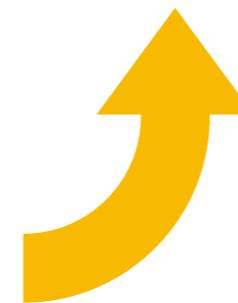
## 2. Variance-based robust regularization

# Robust Risk Minimization



$$R(f) \leq R_{emp}(f) + \dots$$

$$\text{Var}[Q(f)]$$



$$\min_f \left\{ \frac{1}{n} \sum_{i=1}^n Q(f(x_i)) + g(\text{Var}[Q(f(x_i))]) \right\}$$

2. Variance-based robust regularization

# Variance-based robust regularization



**(Namkoong & Duchi, 2017) proposed...**

## 2. Variance-based robust regularization

# Variance-based robust regularization



**(Namkoong & Duchi, 2017) proposed...**

- ◆ Risk minimization alternative

## 2. Variance-based robust regularization

# Variance-based robust regularization



**(Namkoong & Duchi, 2017) proposed...**

- ◆ Risk minimization alternative
- ◆ Optimal and computationally efficient solution

## 2. Variance-based robust regularization

# Variance-based robust regularization



**(Namkoong & Duchi, 2017) proposed...**

- ◆ Risk minimization alternative
- ◆ Optimal and computationally efficient solution
- ◆ Tractable convex formulation

## 2. Variance-based robust regularization

# Variance-based robust regularization



**(Namkoong & Duchi, 2017) proposed...**

- ◆ Risk minimization alternative
- ◆ Optimal and computationally efficient solution
- ◆ Tractable convex formulation
- ◆ Regularization of empirical risk by adding a variance term

## 2. Variance-based robust regularization

# Variance-based robust regularization



**(Namkoong & Duchi, 2017) proposed...**

- ◆ Risk minimization alternative
- ◆ Optimal and computationally efficient solution
- ◆ Tractable convex formulation
- ◆ Regularization of empirical risk by adding a variance term

Bias  $\longleftrightarrow$  Variance



## 2. Variance-based robust regularization

# Alternative as a *min-max problem*

## 1. Minimization of the empirical risk

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n p_i \cdot q_i(\theta, x_i)$$

## 2. Variance-based robust regularization

# Alternative as a *min-max problem*

## 1. Minimization of the empirical risk

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n p_i \cdot q_i(\theta, x_i)$$

$$\max_p \sum_{i=1}^n p_i \cdot q_i$$

## 2. Maximization of the robust objective

## 2. Variance-based robust regularization

# Alternative as a *min-max problem*

**1. Minimization of the empirical risk**

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n p_i \cdot q_i(\theta, x_i)$$



$$\max_p \sum_{i=1}^n p_i \cdot q_i$$

**2. Maximization of the robust objective**



## 2. Variance-based robust regularization

# Alternative as a *min-max problem*

## Constraints

## 2. Variance-based robust regularization

# Alternative as a *min-max problem*

## Constraints

$$p \in \mathcal{P}_n = \left\{ p \in \mathbb{R}_+^n : \frac{1}{2} \|np - 1\|_2^2 \leq \rho, \langle \mathbf{1}, p \rangle = 1 \right\}$$

## 2. Variance-based robust regularization

# Alternative as a *min-max problem*

## Results

## 2. Variance-based robust regularization

# Alternative as a *min-max problem*

## Results

- ❖ **Faster rates of convergence**

## 2. Variance-based robust regularization

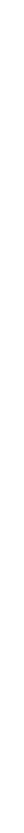
# Alternative as a *min-max problem*

## Results

- ❖ Faster rates of convergence
- ❖ Improvement of out-of-sample test performance

### 3. Application on Deep Learning

# Application on Deep Learning



### 3. Application on Deep Learning

# Application on Deep Learning

**Alternative to risk minimization**



### 3. Application on Deep Learning

# Application on Deep Learning

**Alternative to risk minimization**

**Deep Learning method**

### 3. Application on Deep Learning

# Application on Deep Learning

**Alternative to risk minimization**

**Deep Learning method**



### 3. Application on Deep Learning

# Application on Deep Learning

Alternative to risk minimization



Deep Learning method



**Cost function**

## 3. Application on Deep Learning

# Application on Deep Learning

Alternative to risk minimization



Deep Learning method



**Cost function**



## 3. Application on Deep Learning

# Application on Deep Learning

**Alternative to risk minimization**



**Deep Learning method**



**Cost function**



**Optimization procedure**

## 3. Application on Deep Learning

# Application on Deep Learning

**Alternative to risk minimization**



**Deep Learning method**



**Cost function**



**Optimization procedure**



## 3. Application on Deep Learning

# Application on Deep Learning

**Alternative to risk minimization**



**Deep Learning method**

**Cost function**



**Optimization procedure**



**Model**

### 3. Application on Deep Learning

# Convolutional Neural Network (CNN)

Image classification task

### 3. Application on Deep Learning

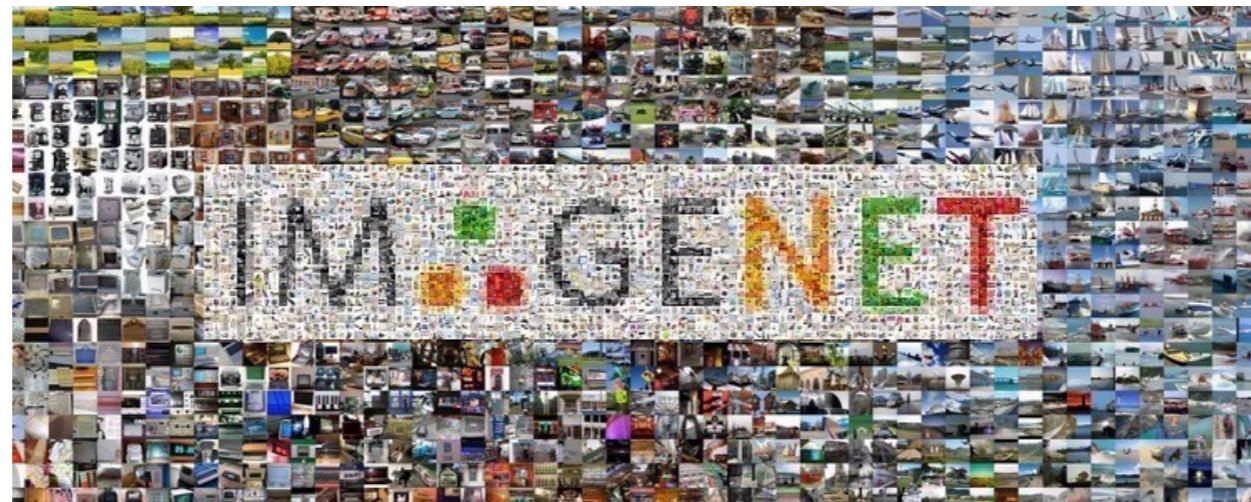
# Convolutional Neural Network (CNN)

Datasets

Image classification task

### 3. Application on Deep Learning

# Convolutional Neural Network (CNN)



Datasets

Image classification task

## 3. Application on Deep Learning

# Convolutional Neural Network (CNN)

**Datasets**

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9

**Image classification task**

### 3. Application on Deep Learning

# Convolutional Neural Network (CNN)



# Datasets

The image displays a 10x10 grid of handwritten digits from the MNIST dataset. Each digit is rendered in a different font style. In the bottom-left corner, there is a smaller inset image showing three digits (1, 6, 2) with blue boxes highlighting specific pixels, likely illustrating a neural network's receptive field or feature extraction process.

## 3. Application on Deep Learning

# Convolutional Neural Network (CNN)



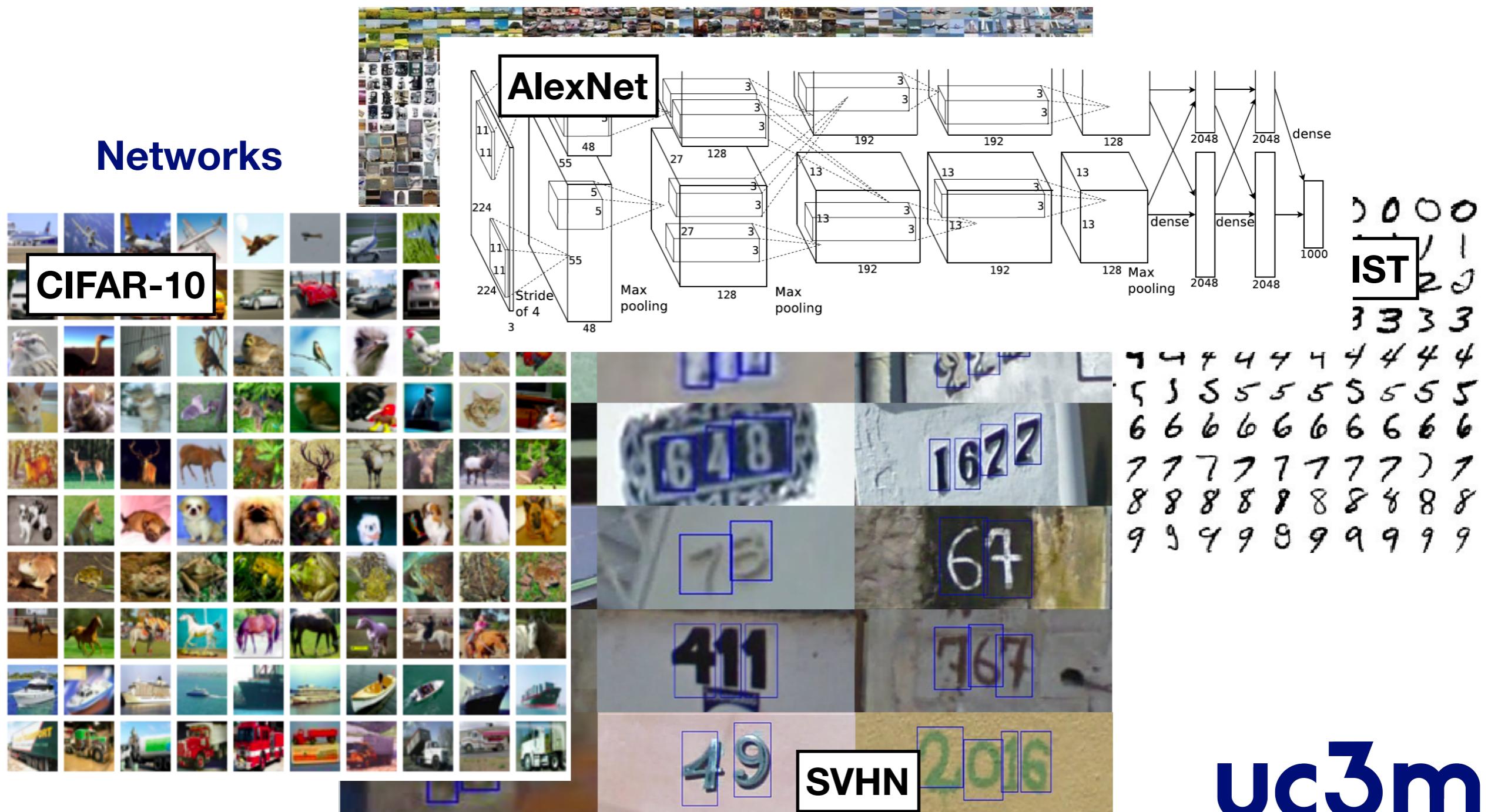
## 3. Application on Deep Learning

# Convolutional Neural Network (CNN)



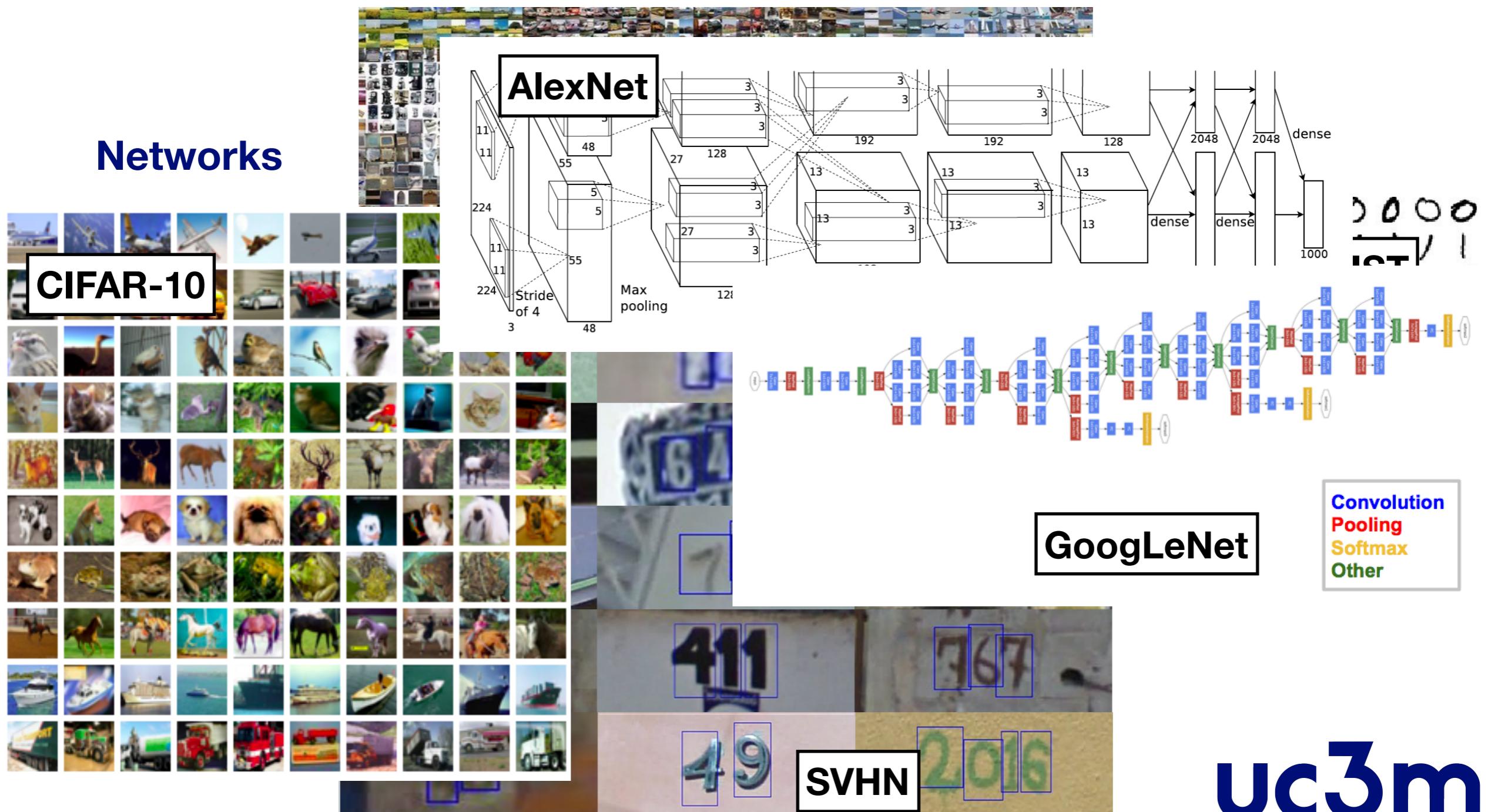
## 3. Application on Deep Learning

# Convolutional Neural Network (CNN)



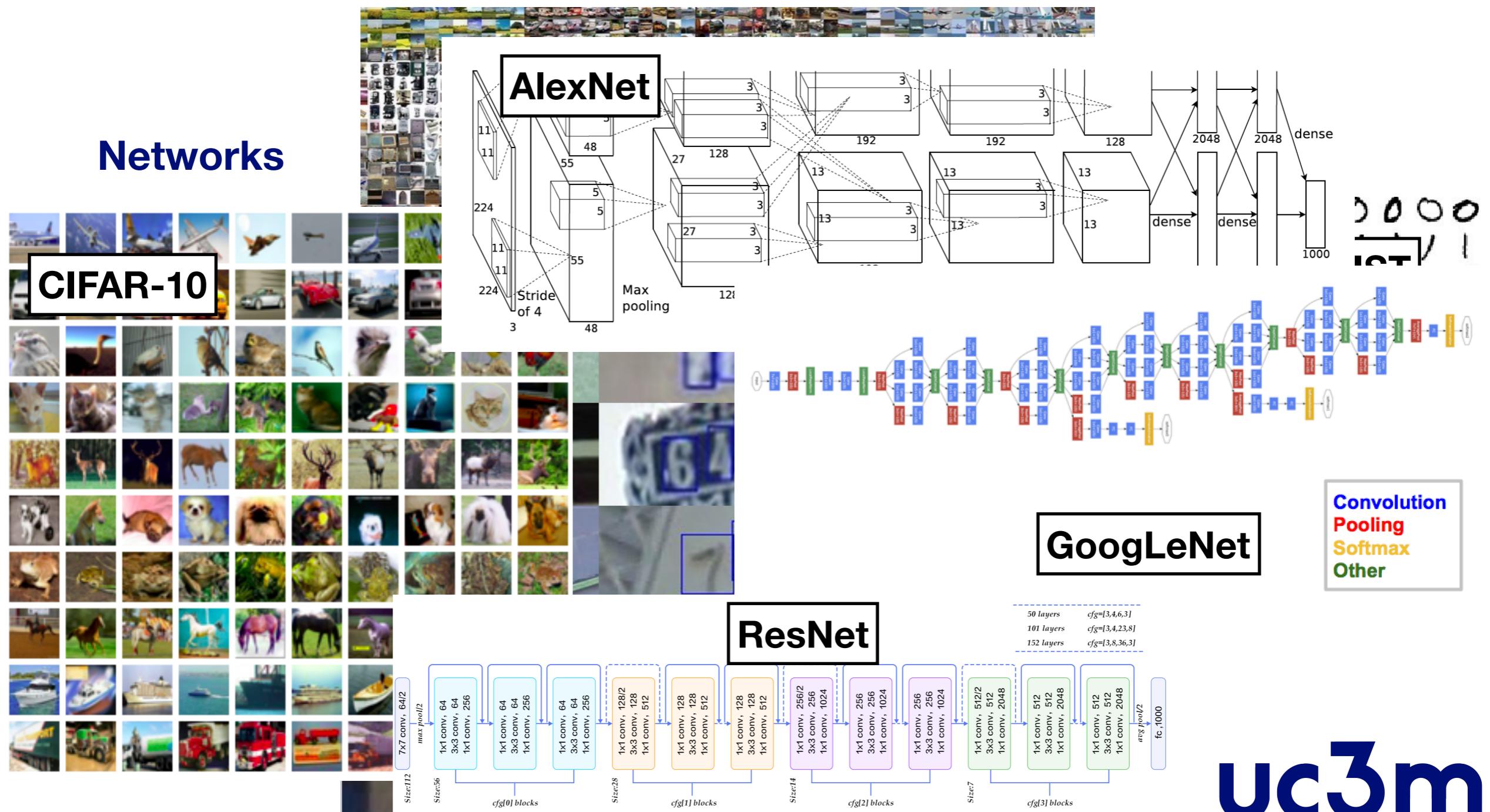
## 3. Application on Deep Learning

# Convolutional Neural Network (CNN)



## 3. Application on Deep Learning

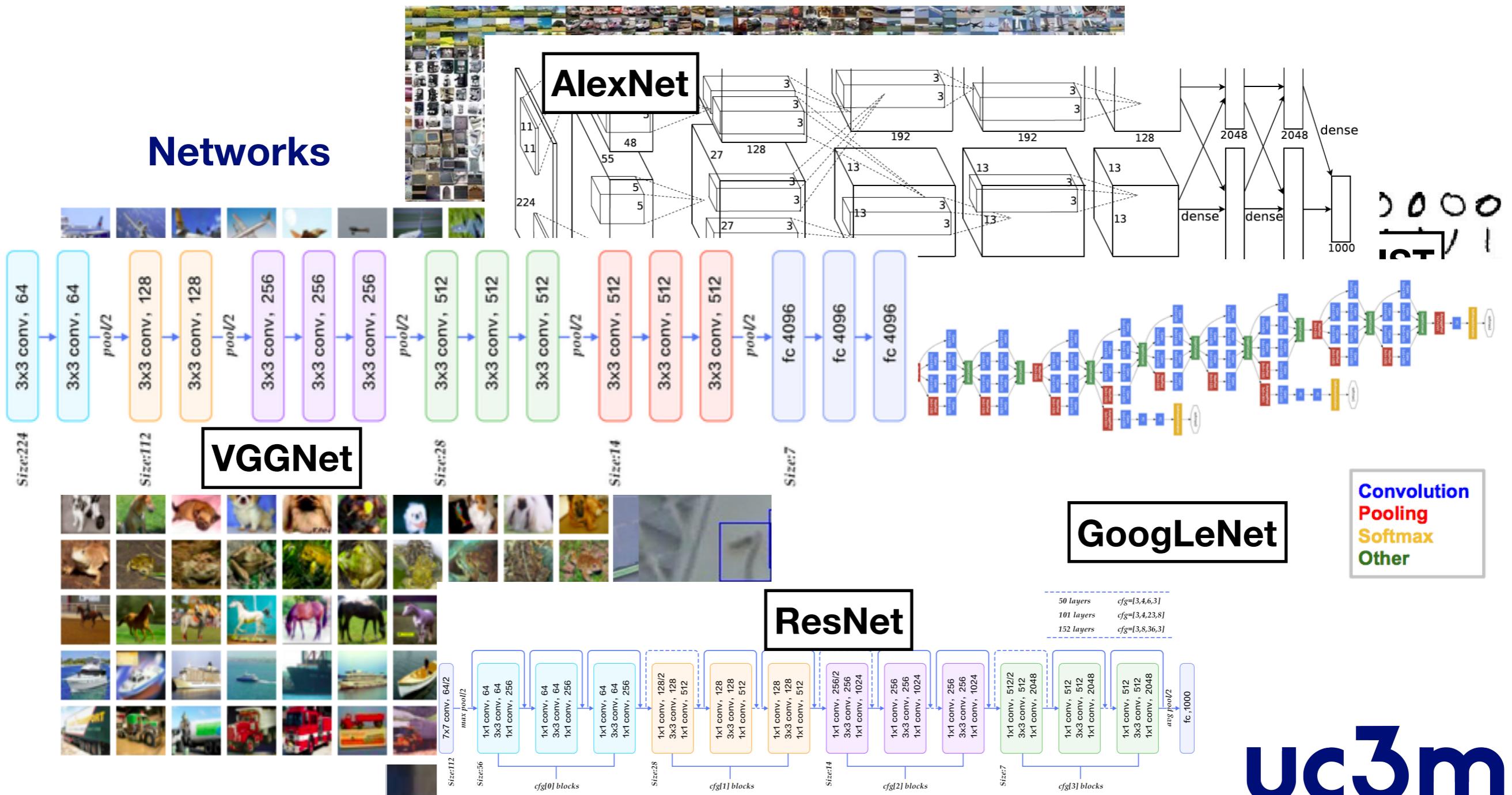
# Convolutional Neural Network (CNN)



## 3. Application on Deep Learning

# Convolutional Neural Network (CNN)

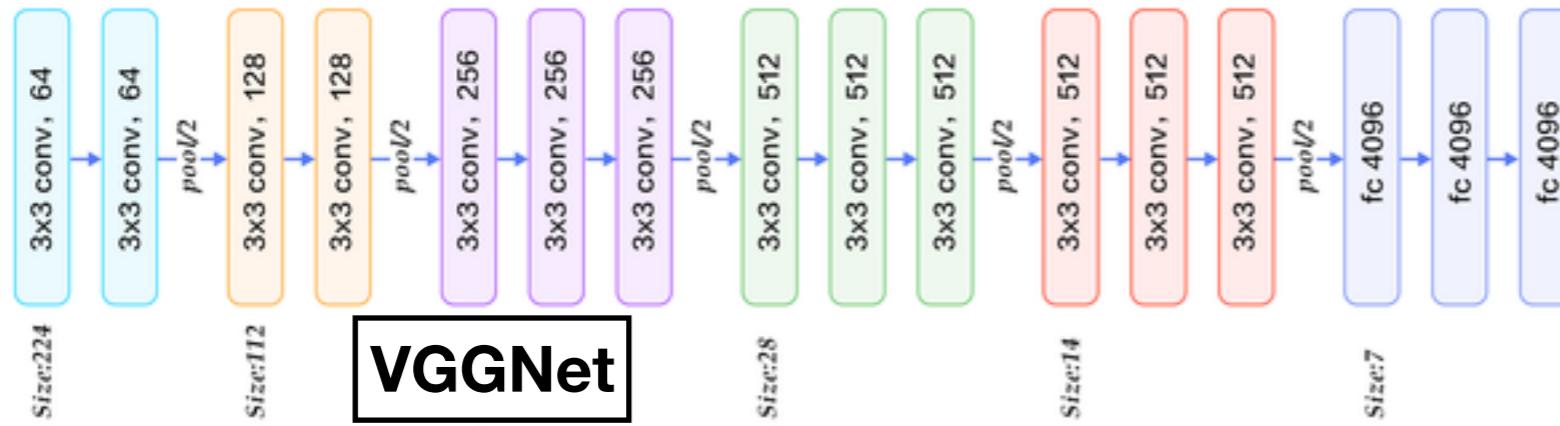
## Networks



## 3. Application on Deep Learning

# Convolutional Neural Network (CNN)

## Networks



## Datasets

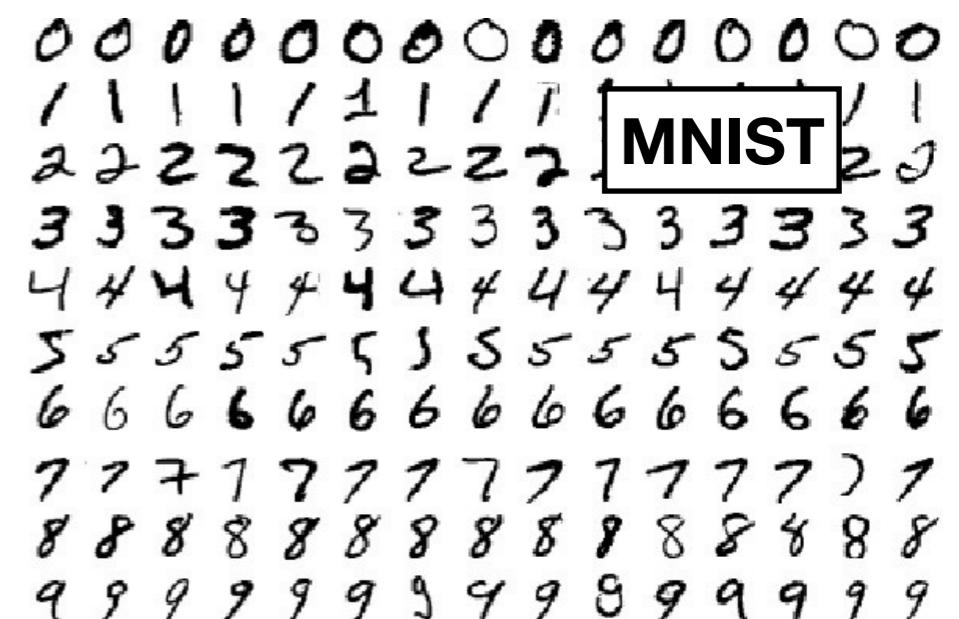


Image classification task

## 3. Application on Deep Learning

# Application on Deep Learning

**Alternative to risk minimization**



**Deep Learning method**

**Cost function**



**Optimization procedure**

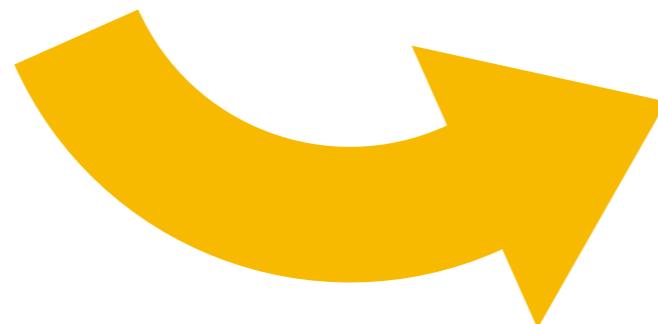


**Model**

### 3. Application on Deep Learning

# Application on Deep Learning

**Alternative to risk minimization**



**Deep Learning method**

=  
**Cost function**

+  
**Optimization procedure**

+  
**Model**

## 3. Application on Deep Learning

# Application on Deep Learning

## Alternative to risk minimization



**10<sup>6</sup> - 10<sup>9</sup> parameters**  
**10<sup>5</sup> - 10<sup>6</sup> samples**

**Deep Learning method**

=  
**Cost function**

+  
**Optimization procedure**

+  
**Model**

## 3. Application on Deep Learning

# Application on Deep Learning

**Alternative to risk minimization**



**Computational cost**

**Deep Learning method**

=  
**Cost function**

+  
**Optimization procedure**

+  
**Model**

### 3. Application on Deep Learning

# Variance reducer regularization



Applied to Mini-batch...



### 3. Application on Deep Learning

# Variance reducer regularization

Applied to Mini-batch...

Traditional methods



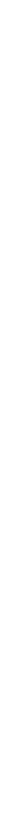
## 3. Application on Deep Learning

# Variance reducer regularization

Applied to Mini-batch...

Traditional methods

\*Sample uniformly



## 3. Application on Deep Learning

# Variance reducer regularization

Applied to Mini-batch...

## Traditional methods

- ❖ Sample uniformly
- ❖ Unbiased estimation



## 3. Application on Deep Learning

# Variance reducer regularization

Applied to Mini-batch...

## Traditional methods

- ❖ Sample uniformly
- ❖ Unbiased estimation
- ❖ High variance



## 3. Application on Deep Learning

# Variance reducer regularization

Applied to Mini-batch...

## Traditional methods

- ❖ Sample uniformly
- ❖ Unbiased estimation
- ❖ High variance

## Robust solution



## 3. Application on Deep Learning

# Variance reducer regularization

Applied to Mini-batch...

## Traditional methods

- ❖ Sample uniformly
- ❖ Unbiased estimation
- ❖ High variance

## Robust solution

- ❖ Non-uniform sampling

## 3. Application on Deep Learning

# Variance reducer regularization

Applied to Mini-batch...

## Traditional methods

- ❖ Sample uniformly
- ❖ Unbiased estimation
- ❖ High variance

## Robust solution

- ❖ Non-uniform sampling
- ❖ Balance bias-variance

## 3. Application on Deep Learning

# Variance reducer regularization

Applied to Mini-batch...

## Traditional methods

- ❖ Sample uniformly
- ❖ Unbiased estimation
- ❖ High variance

## Robust solution

- ❖ Non-uniform sampling
- ❖ Balance bias-variance
- ❖ Reduce variance

### 3. Methodology

# Methodology

Applied to Mini-batch...

### 3. Methodology

# Methodology

Applied to Mini-batch...

- 1. VR-M: Repeat a percentage of the worst performed samples from one mini-batch to the next one.**

## 3. Methodology

# Methodology

Applied to Mini-batch...

- 1. VR-M: Repeat a percentage of the worst performed samples from one mini-batch to the next one.**
  
- 2. VR-E: Repeat a percentage of the worst performed samples from one epoch to the next one.**

### 3. Methodology

# VR-M vs VR-E

---

**Algorithm 1** Variance Reducer per Mini-batch (VR-M)

---

**Input:** Datasets  $\mathcal{D}_{\text{train}}$  and  $\mathcal{D}_{\text{test}}$ .

**Output:** Test accuracy  $\eta$ .

```

1: Initialize parameters  $\theta$ , number of epochs  $E$  and repetition rate  $\epsilon$ ;
2: for  $e = 1 \dots E$  do
3:   Divide dataset  $\mathcal{D}_{\text{train}}$  in  $M$  mini-batches;
4:   for  $m = 1 \dots M$  do
5:      $\{\mathbf{x}_m, \mathbf{y}_m\} \leftarrow$  Obtain next mini-batch  $m$ ;
6:      $\{\mathbf{x}_m, \mathbf{y}_m\} \leftarrow$  Substitute  $\epsilon \cdot M$  samples with the  $\{\mathbf{x}_{m-1}, \mathbf{y}_{m-1}\}$  of highest
       $\ell_{m-1}$ ;
7:      $\ell_m \leftarrow$  Evaluate cross-entropy in mini-batch  $m$ ;
8:      $\theta \leftarrow$  Update parameters with stochastic gradient descent (SGD);
9:   end for
10:   $\eta_e \leftarrow$  Compute test accuracy on  $\mathcal{D}_{\text{test}}$ ;
11:  Shuffle  $\mathcal{D}_{\text{train}}$ ;
12: end for
```

---



---

**Algorithm 2** Variance Reducer per Epoch (VR-E)

---

**Input:** Datasets  $\mathcal{D}_{\text{train}}$  and  $\mathcal{D}_{\text{test}}$ .

**Output:** Test accuracy  $\eta$ .

```

1: Initialize parameters  $\theta$ , number of epochs  $E$  and repetition rate  $\epsilon$ ;
2: Initialize  $\mathcal{D}_{\text{train}}^{(1)} = \mathcal{D}_{\text{train}}$ ;
3: for  $e = 1 \dots E$  do
4:   Divide dataset  $\mathcal{D}_{\text{train}}^{(e)}$  in  $M$  mini-batches;
5:   for  $m = 1 \dots M$  do
6:      $\{\mathbf{x}_m, \mathbf{y}_m\} \leftarrow$  Obtain next mini-batch  $m$ ;
7:      $\ell_m^{(e)} \leftarrow$  Evaluate cross-entropy in mini-batch  $m$ ;
8:      $\theta \leftarrow$  Update parameters with stochastic gradient descent (SGD);
9:   end for
10:   $\eta_e \leftarrow$  Compute test accuracy on  $\mathcal{D}_{\text{test}}$ ;
11:  Shuffle  $\mathcal{D}_{\text{train}}$ ;
12:   $\mathcal{D}_{\text{train}}^{(e+1)} \leftarrow \mathcal{D}_{\text{train}}$ ;
13:   $\mathcal{D}_{\text{train}}^{(e+1)} \leftarrow$  Substitute  $\epsilon \cdot E$  samples with  $\{x_i, y_i\} \in \mathcal{D}_{\text{train}}^{(e)}$  of highest  $\ell^{(e)}$ ;
14: end for
```

---

## 3. Methodology

# Methodology

Applied to Mini-batch...

- 1. VR-M: Repeat a percentage of the worst performed samples from one mini-batch to the next one.**
- 2. VR-E: Repeat a percentage of the worst performed samples from one epoch to the next one.**
- 3. PVR-M/E: Resample half of the data points from the worst performed ones each mini-batch or epoch.**

## 5. Experiments

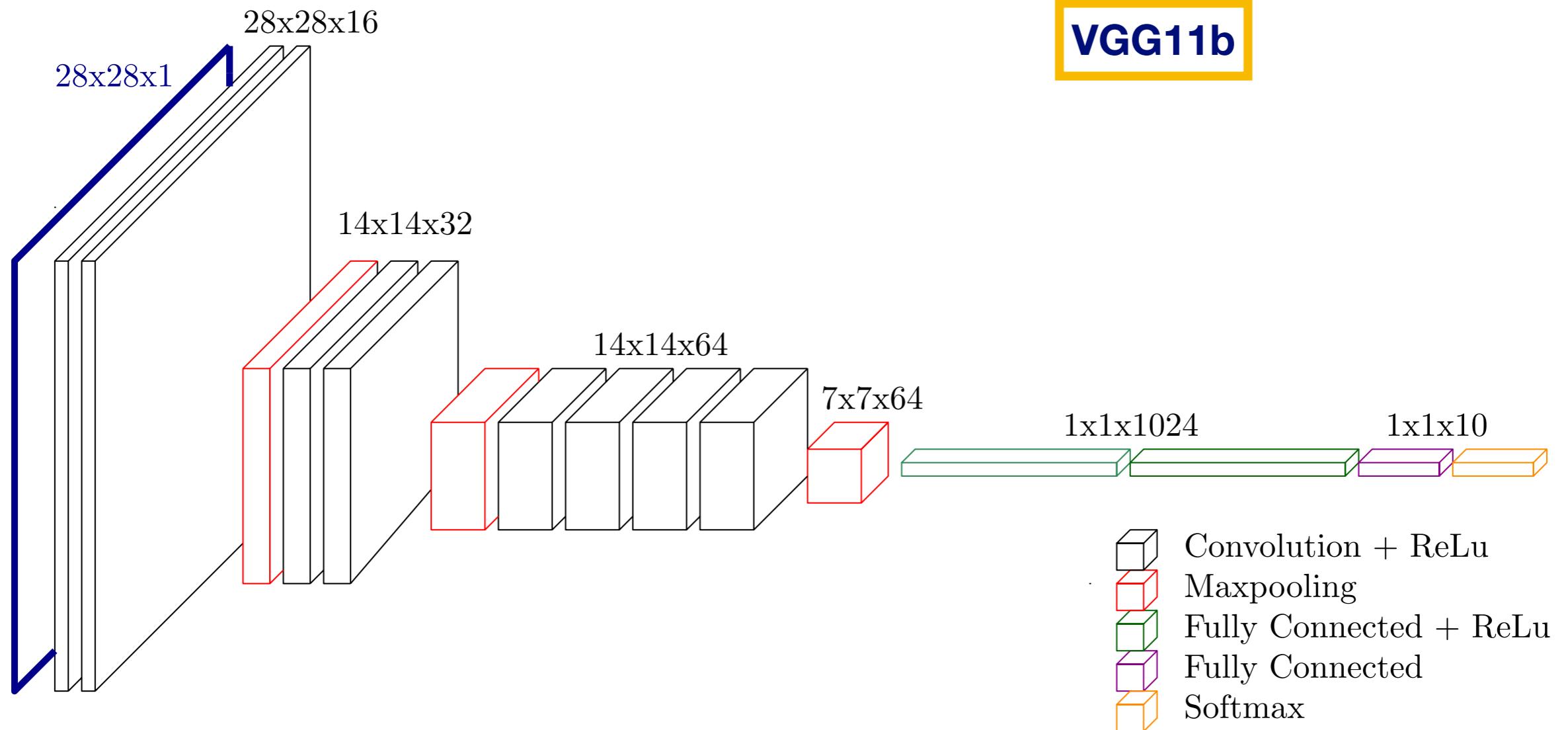
# MNIST Dataset



- **60000 training images**
- **10000 test images**
- **28x28 pixels**
- **Gray scale**

## 5. Experiments

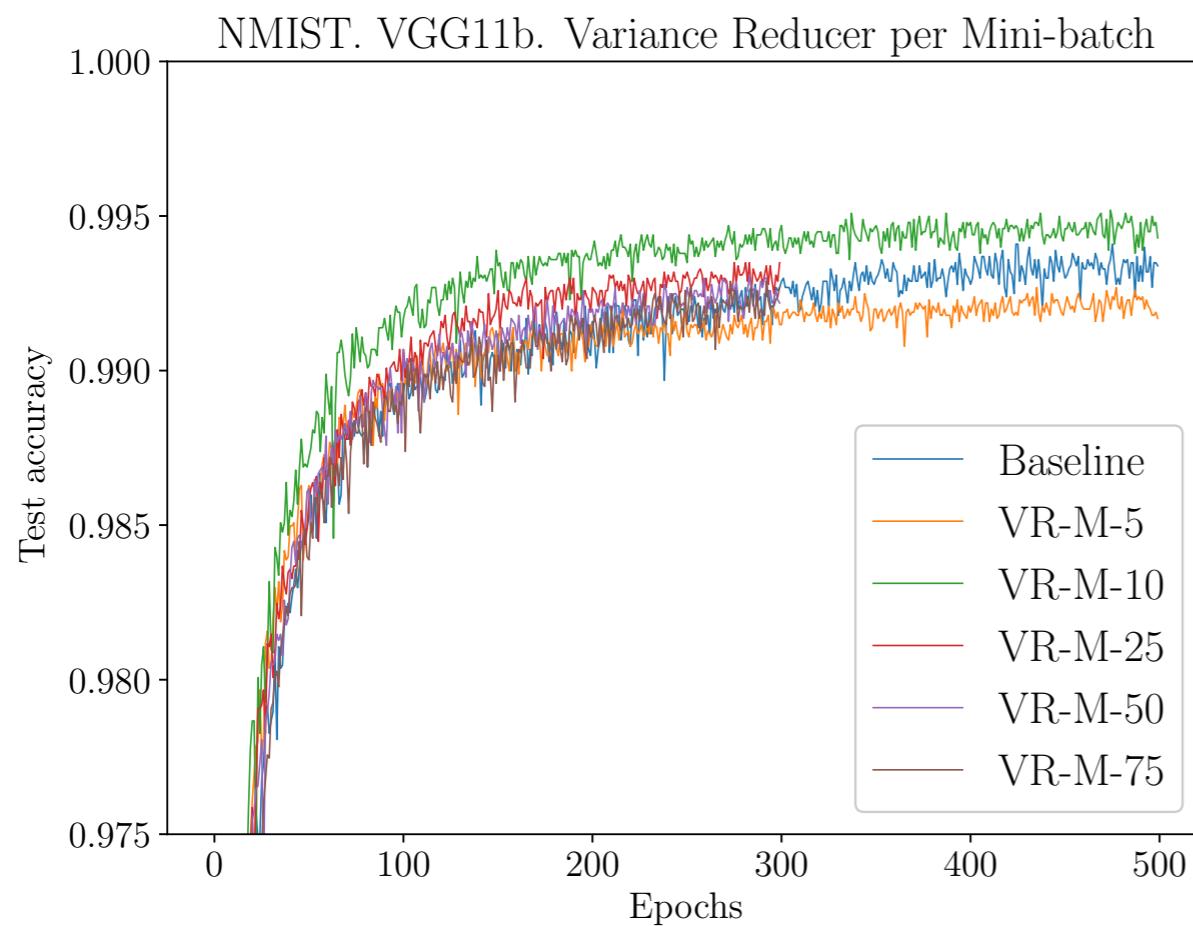
# Very deep CNN



Based on (Simonyan & Zisserman, 2015)

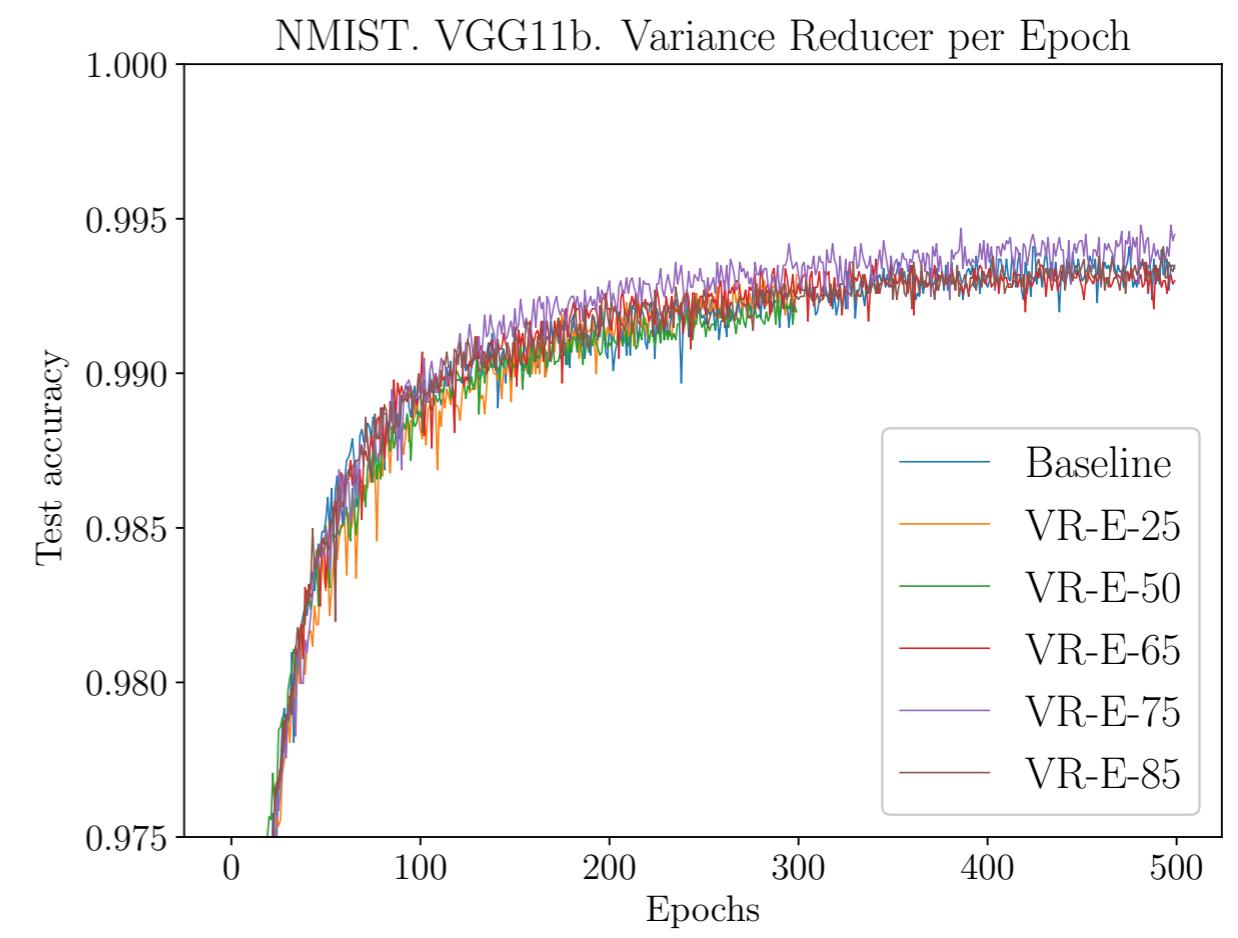
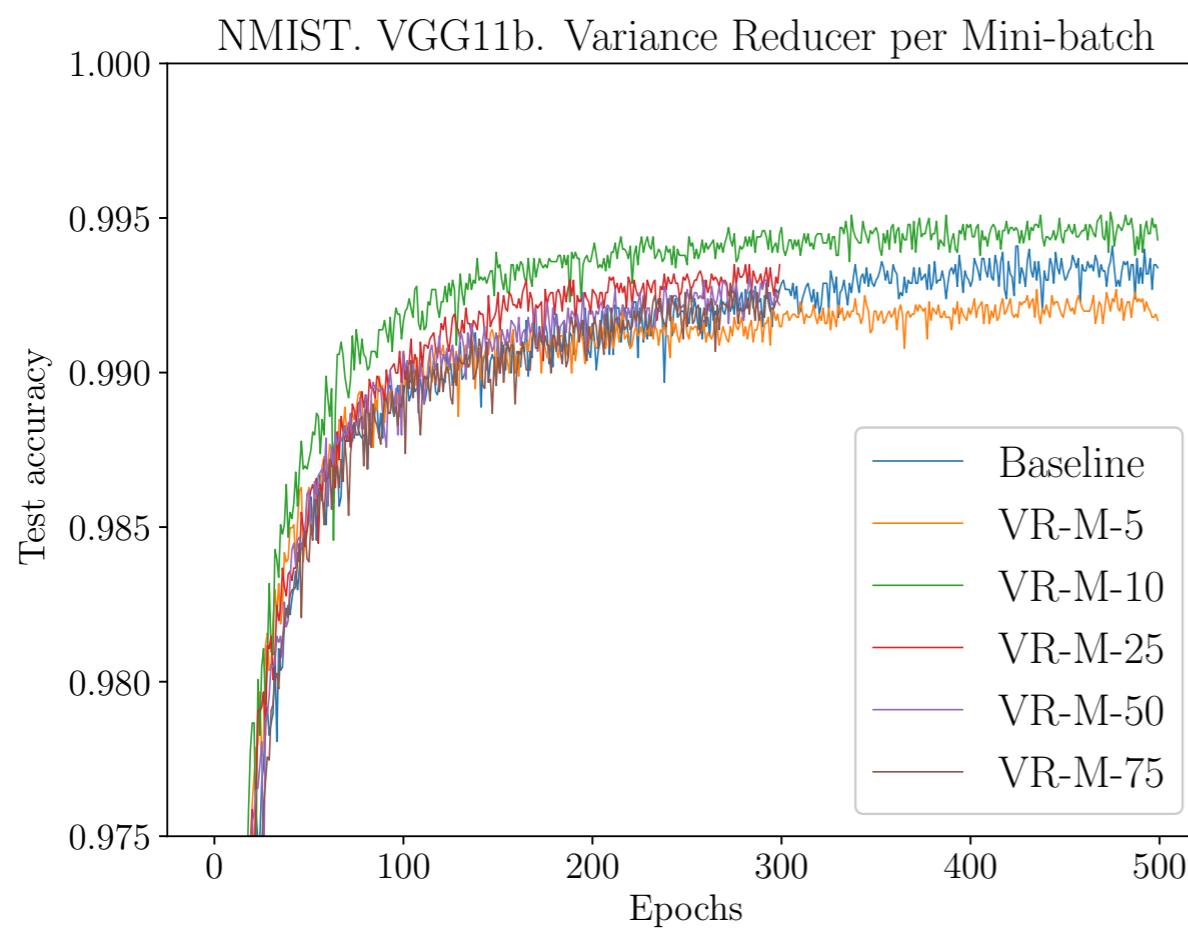
## 5. Experiments

# Results



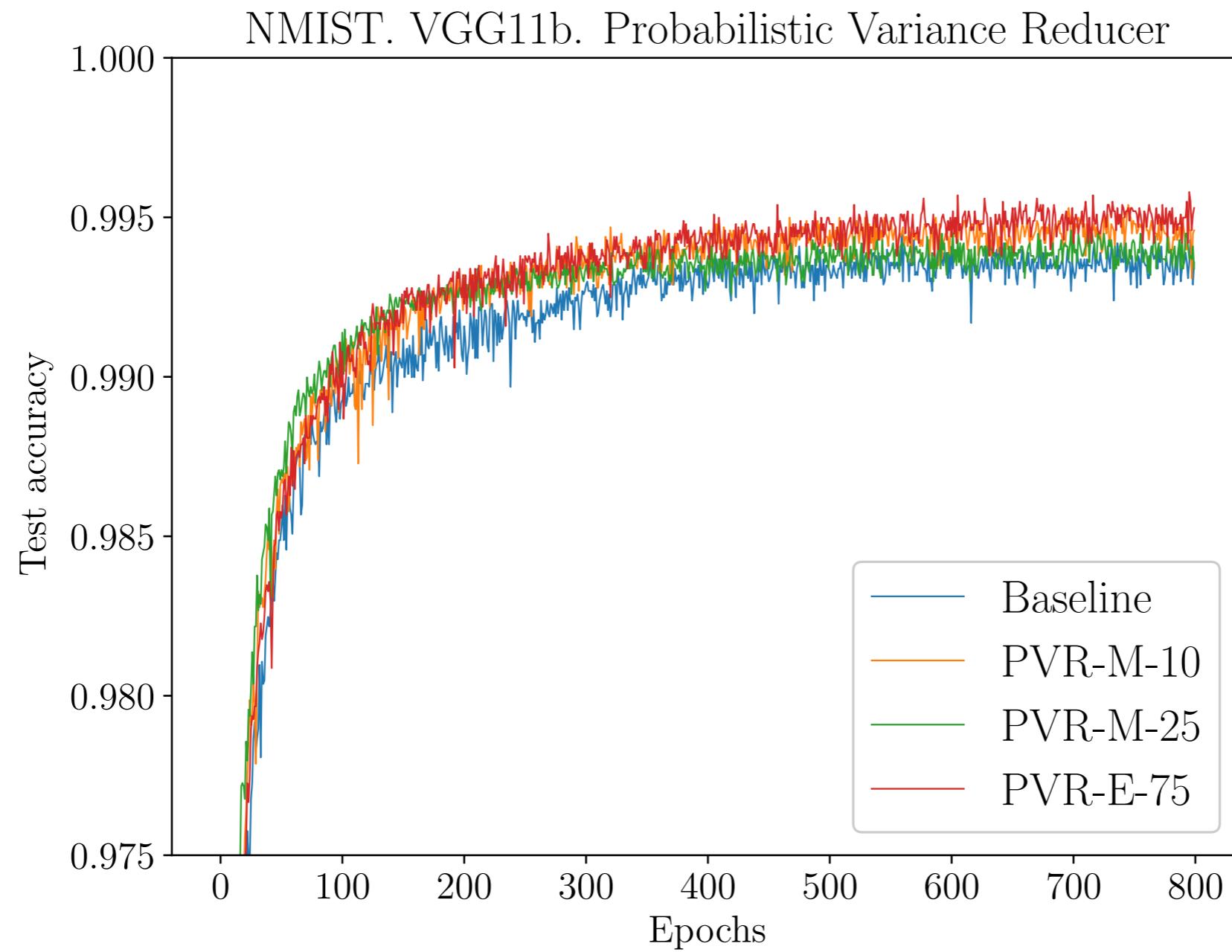
## 5. Experiments

## Results



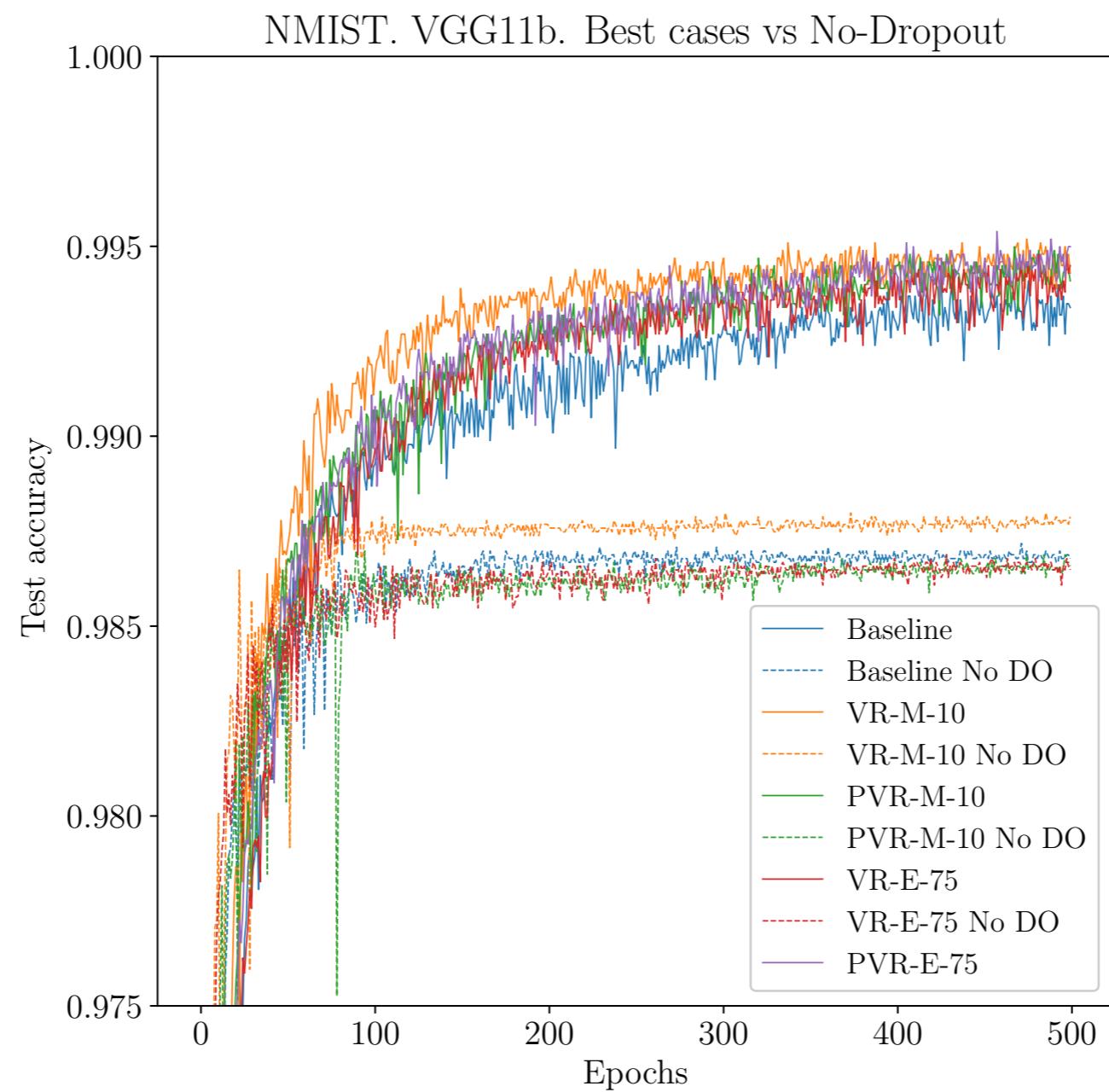
## 5. Experiments

# Results



## 5. Experiments

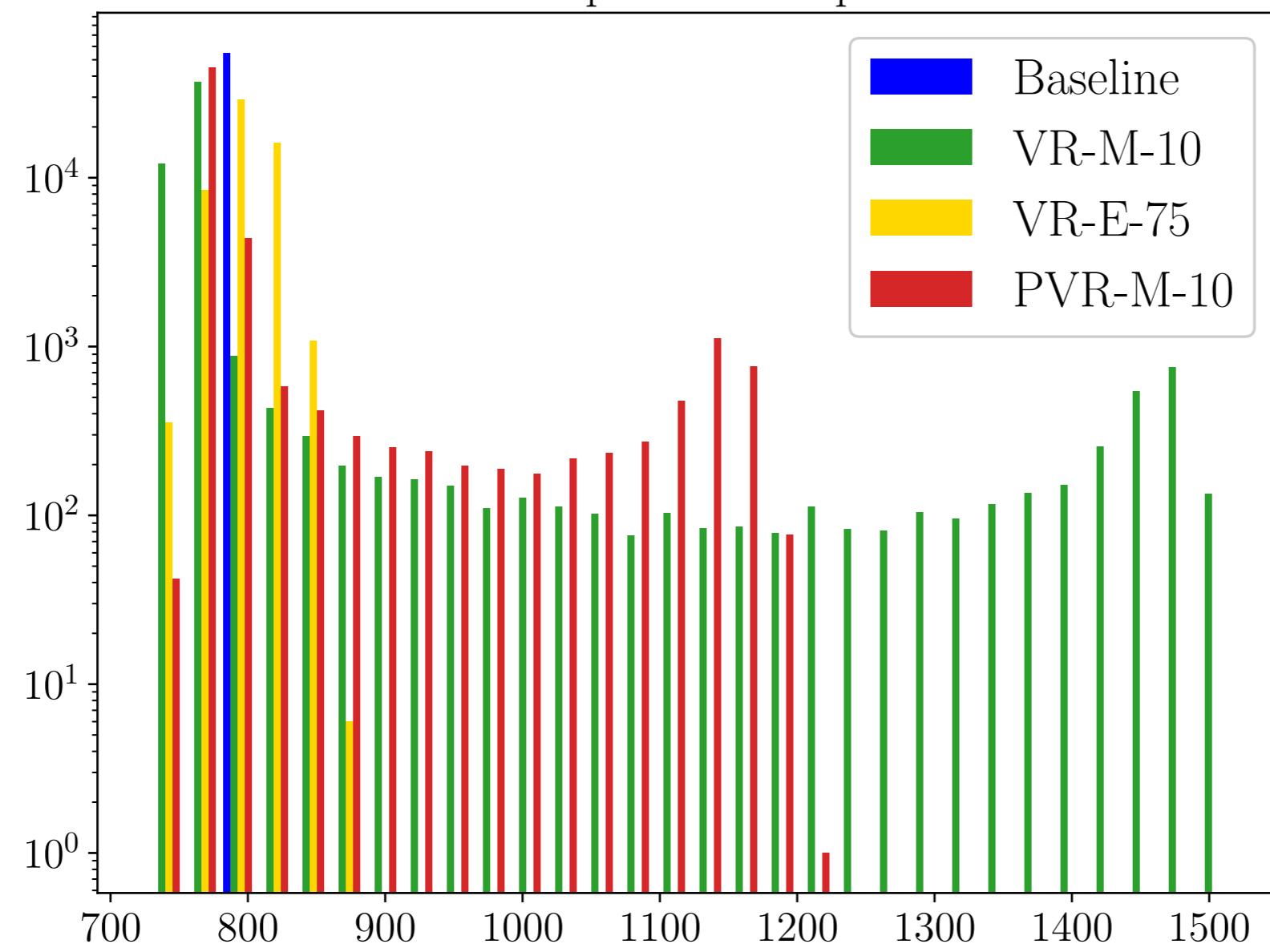
# Results



## 5. Experiments

# Results

Count of samples in the optimization



## 6. Conclusions

# Conclusions

- ❖ Study of a **novel robust estimation based on a reduction of the variance.**

## 6. Conclusions

# Conclusions

- ❖ Study of a novel robust estimation based on a reduction of the variance.
- ❖ Prove theoretical evidences with experiments in CNN.

## 6. Conclusions

# Conclusions

- ❖ Study of a **novel robust estimation based on a reduction of the variance.**
- ❖ Prove theoretical evidences with experiments in **CNN**.
- ❖ Obtain methods with different characteristics:
  - **VR-M with faster convergence.**
  - **VR-E with better performance.**
  - **PVR to deal with misclassified samples.**

## 6. Conclusions

# Conclusions

Epoch	10	100	300	500
<b>Baseline</b>	95,49%	98,91%	99,26%	99,33%
<b>VR-M-10</b>	<b>96,03%</b>	<b>99,10%</b>	<b>99,44%</b>	99,43%
<b>PVR-M-10</b>	95,10%	99,06%	99,37%	99,41%
<b>VR-E-75</b>	95,56%	98,93%	99,35%	99,45%
<b>PVR-E-75</b>	95,43%	98,96%	99,36%	<b>99,50%</b>

## 6. Conclusions

# Conclusions

- ❖ Study of a **novel robust estimation based on a reduction of the variance.**
- ❖ Prove theoretical evidences with experiments in **CNN**.
- ❖ Obtain methods with different characteristics:
  - **VR-M with faster convergence.**
  - **VR-E with better performance.**
  - **PVR to deal with misclassified samples.**
- ❖ Open work with the study of other scenarios and larger datasets.

## 7. References

# References

- ✓ H. Namkoong and J. C. Duchi. **Variance-based regularization with convex objectives**. NIPS 2017
- ✓ Z. Borsos, A. Krause and K. Levy. **Online variance reduction for stochastic optimization**. arXiv: 1802.04715, 2018
- ✓ K. Simonyan and A. Zisserman. **Very deep convolutional networks for large-scale image recognition**. In *In Proc. ICLR*, 2015.

**Thank you!**