# MLP&CNN

Jinna

School of information science and engineering

2016.9.26

计算机视觉实验室

# CONTENT

计算机视觉实验室

**Header File Here!**

**Code Here!**

Training and testing code of MLP

**error**

Training and testing code of MLP

**training_accuarcy**



epoch_number

Training and testing code of MLP

**testing result**

**What I learned from this process**

C

1. Remember initializing variables
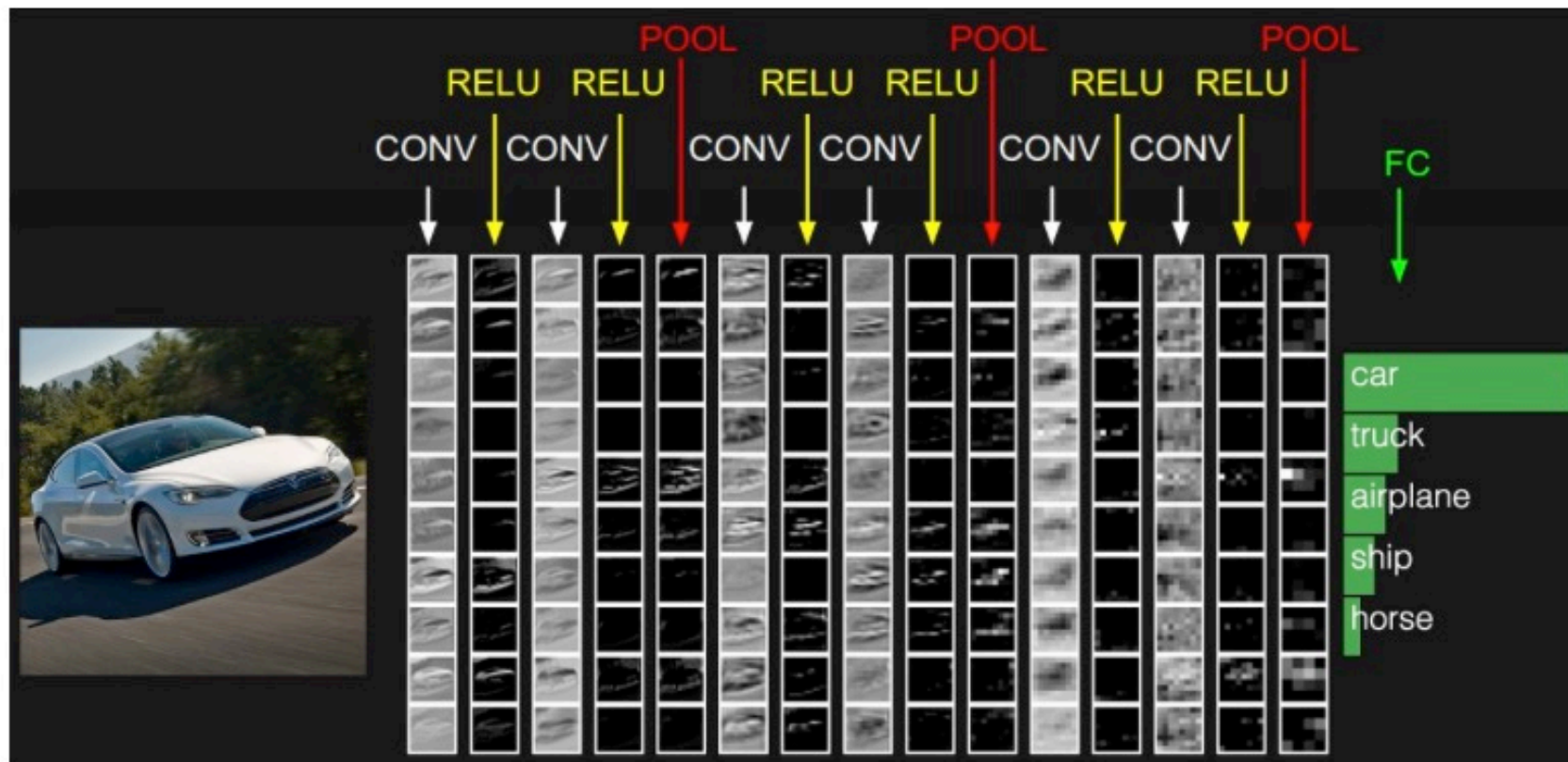
2. Organize your logic

3. Recommend using class, header File

MLP

4. Reset values

5. BP & gradient descent

6. Choose right activation function

7. Batch_size

8. Coding is difficult but interesting

# CNN：Convolutional Neural Network

# CNN Net:

**LeNet**
**AlexNet**
**VGG Net**
**GoogLeNet**

http://cs231n.github.io/convolutional-networks/

# LeNet:

## LeNet:

Application of CNN

计算机视觉实验室

**Thanks for your attention！**