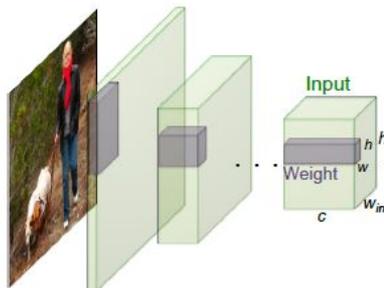# Interactive Gesture Plugin using Quantized Neural Network

Shaofan Lai

# Gesture Recognition

- Traditional methods: (Filter+Threshold)
    - Sensitive to the background.
    - Hard-coded definition of gesture
- State-of-the-art Devices: (Kinect, Leap Motion)
    - Requires extra information (e.g. depth, infrared)
    - Expensive
- Deep Neural Network (Lab setting):
    - Intensive computation
    - Works with high-resolution video
    - Millions of data samples

# Recognize Gesture with JeVois (OpenCV)

# Quantized (Binarized) Neural Network



| | Network Variations | | Operations used in Convolution | Memory Saving (Inference) | Computation Saving (Inference) | Accuracy on ImageNet (AlexNet) |
|---|---|---|---|---|---|---|
| Standard Convolution | Real-Value Inputs<br>0.11 -0.21 ... -0.34<br>-0.25 0.61 ... 0.52 | Real-Value Weights<br>0.12 -1.2 ... 0.41<br>-0.2 0.5 ... 0.68 | +, −, × | 1x | 1x | %56.7 |
| Binary Weight | Real-Value Inputs<br>0.11 -0.21 ... -0.34<br>-0.25 0.61 ... 0.52 | Binary Weights<br>1 -1 ... 1<br>-1 1 ... 1 | +, − | ~32x | ~2x | %56.8 |
| BinaryWeight Binary Input (XNOR-Net) | Binary Inputs<br>1 -1 ... -1<br>-1 1 ... 1 | Binary Weights<br>1 -1 ... 1<br>-1 1 ... 1 | XNOR, bitcount | ~32x | ~58x | %44.2 |

Fig. 1: We propose two efficient variations of convolutional neural networks. **Binary-Weight-Networks**, when the weight filters contains binary values. **XNOR-Networks**, when both weigh and input have binary values. These networks are very efficient in terms of memory and computation, while being very accurate in natural image classification. This offers the possibility of using accurate vision techniques in portable devices with limited resources.

# XNOR-Network with Darknet

- On CIFAR-10
  - Standard Network
    - top 1: 0.831600, top 5: 0.992000
  - Xnor Network
    - top 1: 0.684600, top 5: 0.976100
- Drawback
  - No time advantage
    - GPU (GTX 1080 Ti)
      - Standard~2.872s vs Bin~5.340s
    - CPU (Intel 7700K)
      - Standard~32.908s vs Bin~33.316s
  - No memory/storage advantage

```
217   if(xnor){
218       l.binary_weights = calloc(c*n*size*size, sizeof(float));
219       l.binary_input = calloc(l.inputs*l.batch, sizeof(float));
220   }
```

```
29  [maxpool]
30  size=2
31  stride=2
32
33  [convolutional]
34  batch_normalize=1
35  filters=64
36  size=3
37  stride=1
38  pad=1
39  activation=leaky
40
41  xnor=1
42
43  [maxpool]
44  size=2
45  stride=2
46
47  [convolutional]
48  batch_normalize=1
49  filters=128
50  size=3
51  stride=1
52  pad=1
53  activation=leaky
54
55  xnor=1
56
```

# My Proposal

Using <u>quantized neural network</u> to <span style="color:red">recognize the gestures</span> <u>in my workplace</u>:

- Using only JeVois Camera(s)
- Relatively consistent background
- Customized gestures
- Interacting with the system in real time

# Challenge

- Training and employment of the gesture model
    - Data (Collected in workplace? Transfered from other gesture datasets?)
    - Model (Attention-base? Detection-based? Tracking-based?)
    - Training (Unsupervised (with saliency)? Supervised (with labels)? Reinforcement Learning (with auxiliary sensors)?
- Implementation of the Quantized Neural Network
    - Memory usage improvement
    - Computation acceleration
- Hacking GNOME system plugin
    - Communication with JeVois
    - User Interface design

# References

- *Binarized Neural Networks: Training Neural Networks with Weights and Activations Constrained to +1 or -1*
- *Quantized Neural Networks: Training Neural Networks with Low Precision Weights and Activations*
- *XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks*