

# MATLAB EXPO

엣지 AI기술을 위한 딥러닝 모델의 1-bit 양자화  
김정훈, 네패스



# 목차

1. 회사 및 발표자 약력 (Introduction to Organization and Business)
  1. 프로젝트 개요 (Project Overview)
  1. 개발 배경 (Backgrounds)
  1. 기술적인 해결과제 (Project Goals and Challenges)
  1. MathWorks 솔루션을 통한 해결 방안 및 결과 (How did we get there and leverage MathWorks)
  1. 결과 및 정리 (Achievements and Outlook)
  1. 결론 (Concluding Remarks)

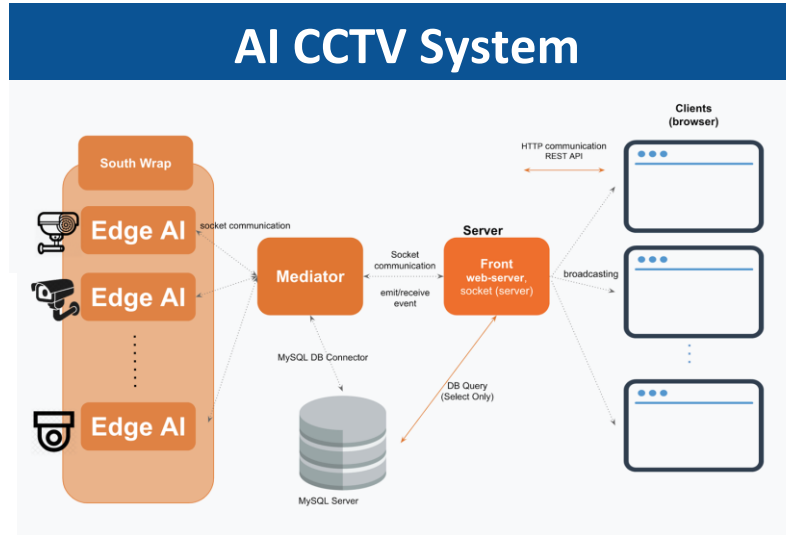
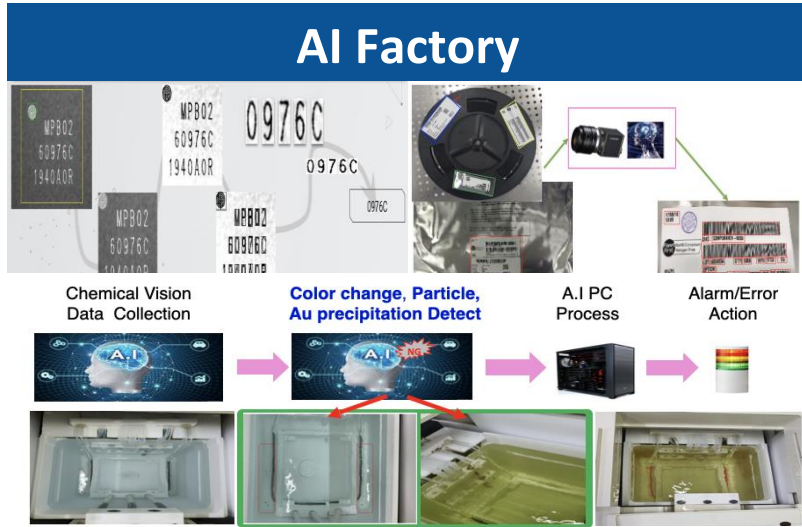
# 발표자 김정훈

- 고려대학교 전기전자공학 석사
- (주)네패스 딥러닝 엔지니어 전문연구요원
- Contact: jeonghoon.samuel@gmail.com
- 관심 키워드: Neural Network Quantization, Compact Network Design, Robotics Perceptions, State Estimation, ...
- 2019 활동 사항:
  - Google Developer Group Gwangju DevFest 2019 발표자
  - 한국 기술 교육 대학교 온라인 평생 교육원 자문
  - 삼성 오픈 소스 컨퍼런스 (SOSCON) 2019 심사위원
  - Mathworks Advisory Board 2019
  - 인공지능 로보틱스 커뮤니티 운영자
  - Neural Network Quantization & Compact Network Development 스테디 리더

함께 연구하고, 나누며, 성장하는데에 큰 보람을 느끼고 있습니다!

오늘의 발표 내용: **Binarized Neural Networks on MATLAB**





Device	Solution
Hardware	Software

**Real Applications**  
 &  
**Core Value Research**

# 프로젝트 개요: 매트랩을 통한 이진화 신경망 학습기 개발

---

## Binarized Neural Networks: Training Neural Networks with Weights and Activations Constrained to +1 or -1

---

Matthieu Courbariaux\*<sup>1</sup>

Itay Hubara\*<sup>2</sup>

Daniel Soudry<sup>3</sup>

Ran El-Yaniv<sup>2</sup>

Yoshua Bengio<sup>1,4</sup>

<sup>1</sup>Université de Montréal

<sup>2</sup>Technion - Israel Institute of Technology

<sup>3</sup>Columbia University

<sup>4</sup>CIFAR Senior Fellow

\*Indicates equal contribution. Ordering determined by coin flip.

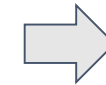
MATTHIEU.COURBARIAUX@GMAIL.COM

ITAYHUBARA@GMAIL.COM

DANIEL.SOUDRY@GMAIL.COM

RANI@CS.TECHNION.AC.IL

YOSHUA.UMONTREAL@GMAIL.COM

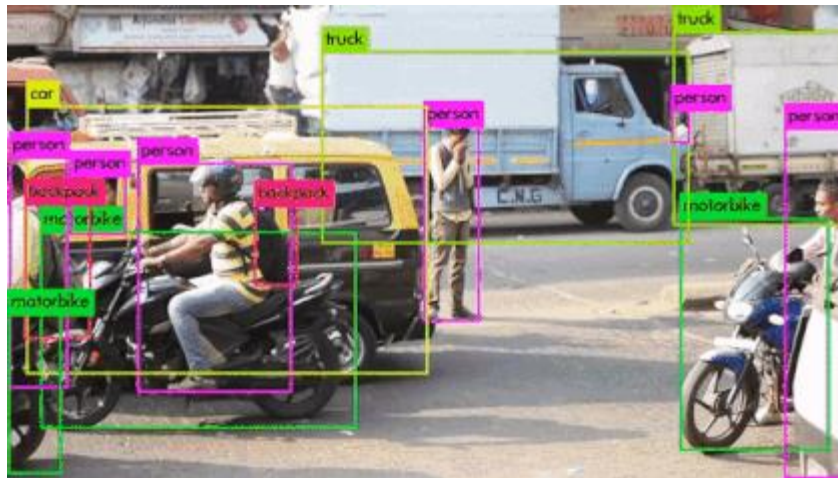


**MATLAB**

MATLAB은 제공해주는 기능만 사용 할 수 있는 폐쇄적 개발환경인가?  
이 논문을 MATLAB으로 어떻게 구현할 수 있을까?

# 개발 배경

훌륭한 성능을 내는 신경망 기반 어플리케이션들

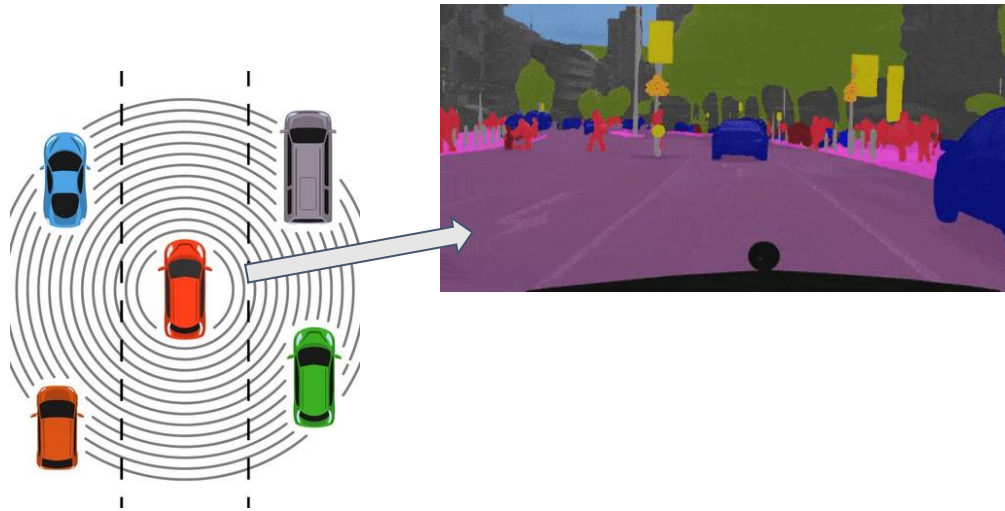


Object Detection



Semantic Segmentation

# 개발 배경

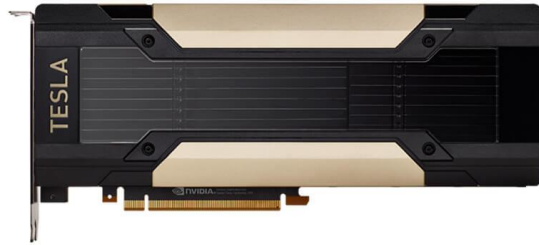


Real Time Performance



Customer

# 개발 배경



x4

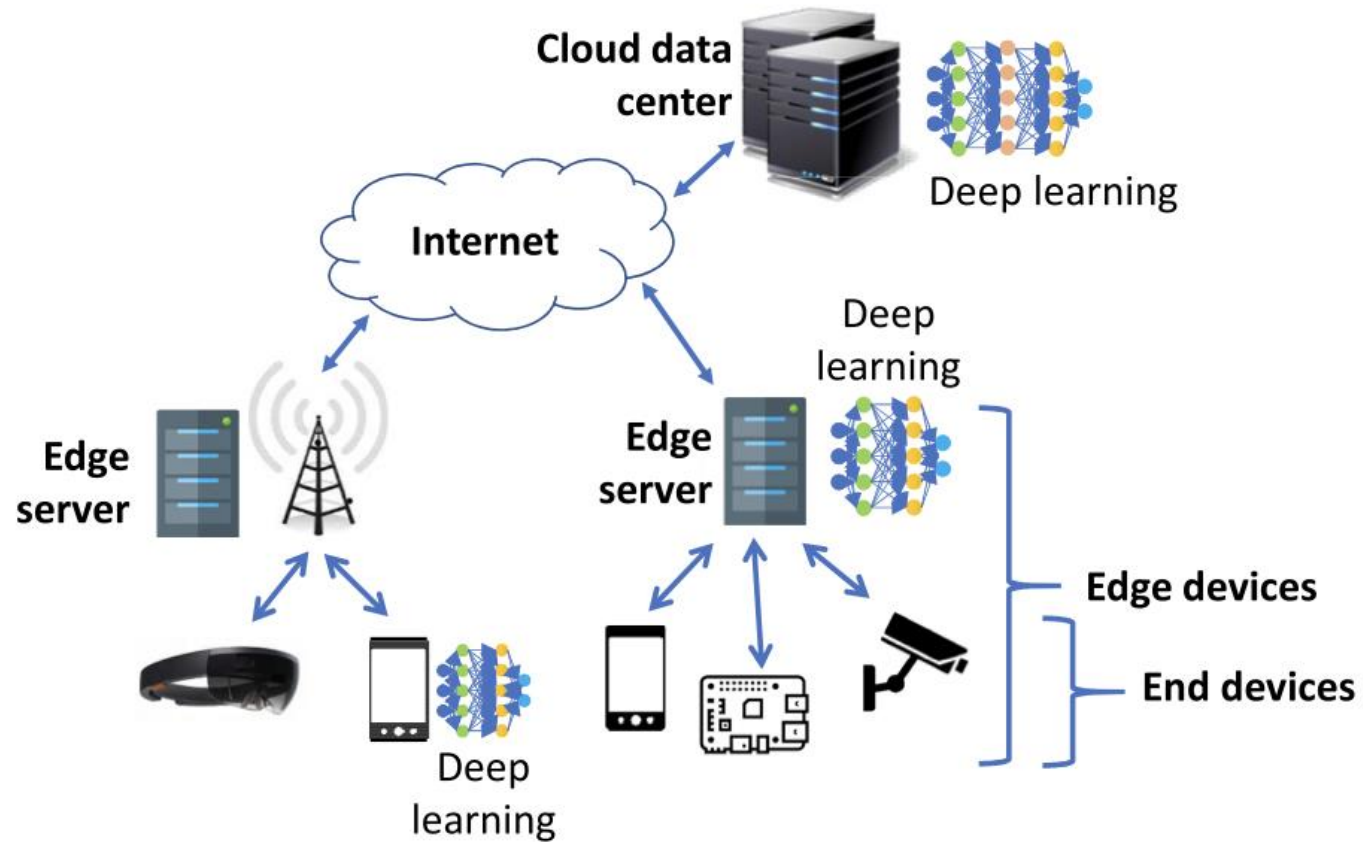
On a High-end GPU



Customer



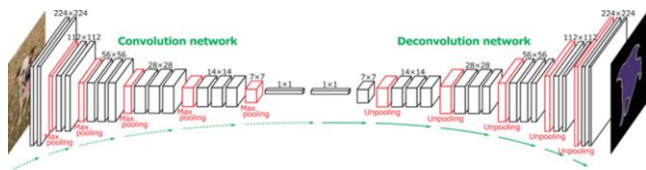
# 개발 배경



Various **Inference Environments** of Deep Learning Network

Chen, Jiasi, and Xukan Ran. "Deep learning with edge computing: A review." Proceedings of the IEEE 107.8 (2019): 1655-1674.

# 개발 배경



Neural Network

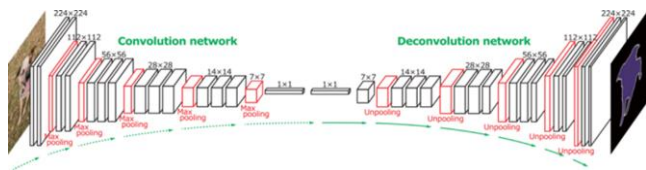


?



Target Environment

# 개발 배경



Neural Network



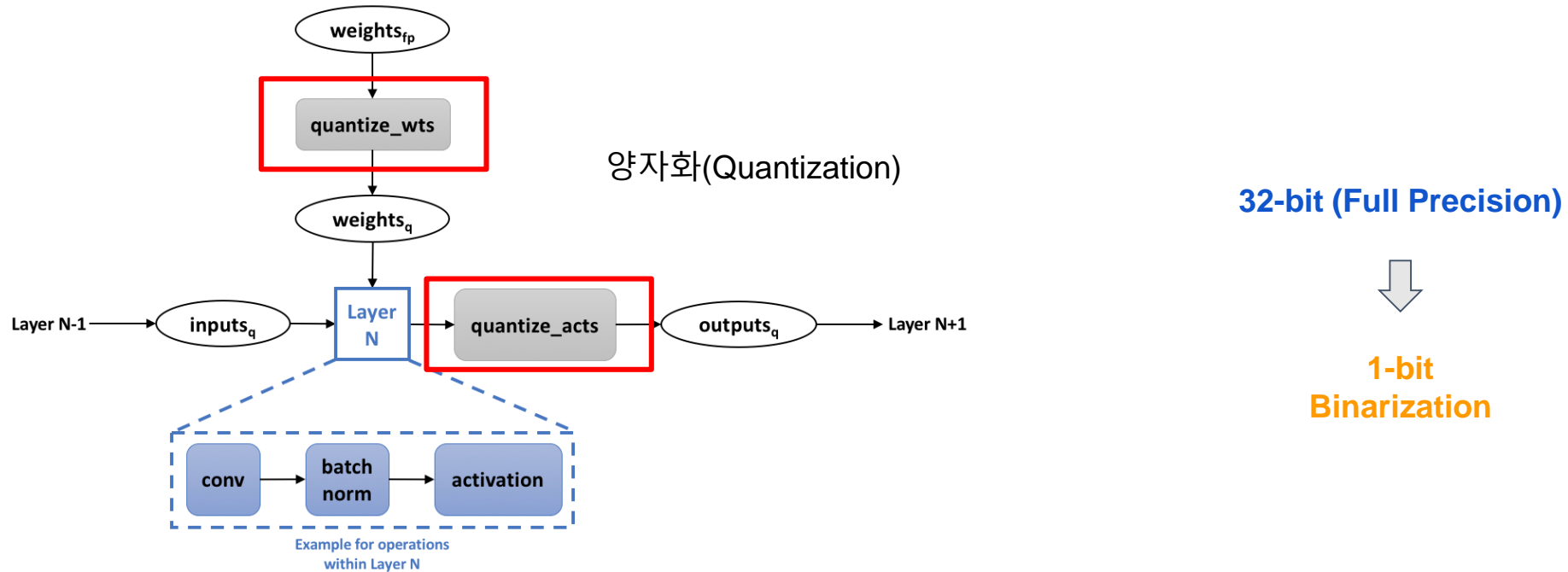
Neural Network  
Quantization



Target Environment

# 기술적인 해결과제

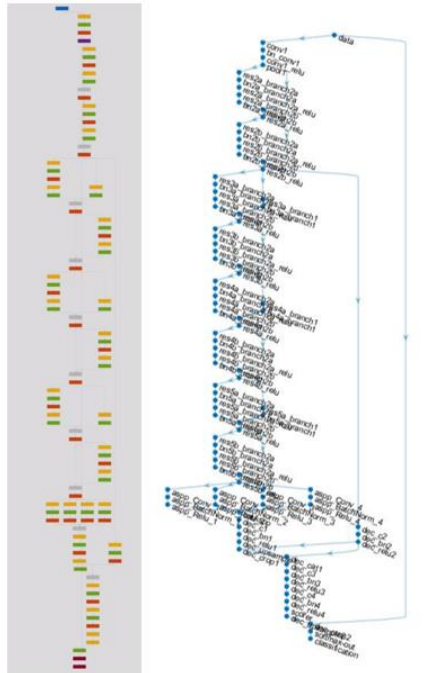
네트워크를 구성하는 weight와 activation의 bit를 낮추어 네트워크를 **경량화** 하거나 연산을 **가속**



MATLAB에서 **1-bit quantization**에 대하여 구현해보기!

(현재 MATLAB 2020a에서 8-bit quantization 기능을 일부 지원하고 있습니다. 오늘 바로 사용해 보세요!)

# 기술적인 해결과제



DeepLabV3+ (Baseline: ResNet18)

Width	Height	Input Channel	Output Channel	Total
7	7	3	64	9408
3	3	64	64	36864
3	3	64	64	36864
3	3	64	64	36864
3	3	64	64	36864
3	3	64	64	36864
1	1	64	48	3072
1	1	64	128	8192
3	3	64	128	73728
3	3	128	128	147456
3	3	128	128	147456
3	3	128	128	147456
3	3	128	256	294912
1	1	128	256	32768
3	3	256	256	589824
3	3	256	256	589824
3	3	256	256	589824
1	1	256	512	131072
3	3	256	512	1179648
3	3	512	512	2359296
3	3	512	512	2359296
3	3	512	512	2359296
3	3	512	256	1179648
3	3	512	256	1179648
3	3	512	256	1179648
1	1	512	256	131072
1	1	1024	256	262144
8	8	256	256	4194304
3	3	304	256	700416
3	3	256	256	589824
1	1	256	11	2816
8	8	11	11	7744

Architecture Weight Count

만약 DeepLabV3+ 를 1-bit Quantization 한다면?

## FP32 vs 1 Bit

Total Weight Count

	bit	byte	
FP32	659111936	82388992	82.388992MB
1bit	20597248	2574656	2.574656 MB

Maximum Activation (Input Size: 360 x 480 x 3)

	bit	byte	
FP32	105,062,400	13132800	13.1328 MB
1bit	3,283,200	172800	0.1728 MB

당연한 결과지만, 32-bit와 1-bit는 Volume에서 32배 차이!

# 기술적인 해결과제

---

## Binarized Neural Networks: Training Neural Networks with Weights and Activations Constrained to +1 or -1

---

Matthieu Courbariaux\*<sup>1</sup>

Itay Hubara\*<sup>2</sup>

Daniel Soudry<sup>3</sup>

Ran El-Yaniv<sup>2</sup>

Yoshua Bengio<sup>1,4</sup>

<sup>1</sup>Université de Montréal

<sup>2</sup>Technion - Israel Institute of Technology

<sup>3</sup>Columbia University

<sup>4</sup>CIFAR Senior Fellow

\*Indicates equal contribution. Ordering determined by coin flip.

MATTHIEU.COURBARIAUX@GMAIL.COM

ITAYHUBARA@GMAIL.COM

DANIEL.SOUDRY@GMAIL.COM

RANI@CS.TECHNION.AC.IL

YOSHUA.UMONTREAL@GMAIL.COM

### Check point:

1. Binarization function
2. Backpropagation of BNNs (Algorithm1)
3. Computing gradients and updates

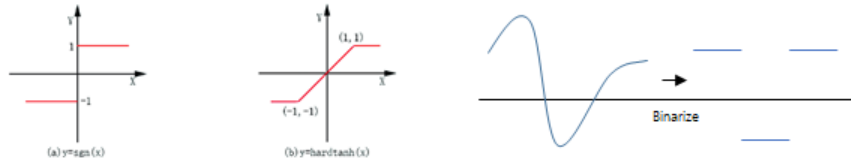
# 기술적인 해결과제 : BNN 구현

## Binarization

### Deterministic Binarization

$$x^b = \text{Sign}(x) = \begin{cases} +1 & \text{if } x \geq 0, \\ -1 & \text{otherwise,} \end{cases} \quad (1)$$

where  $x^b$  is the binarized variable (weight or activation)  
and  $x$  the real-valued variable



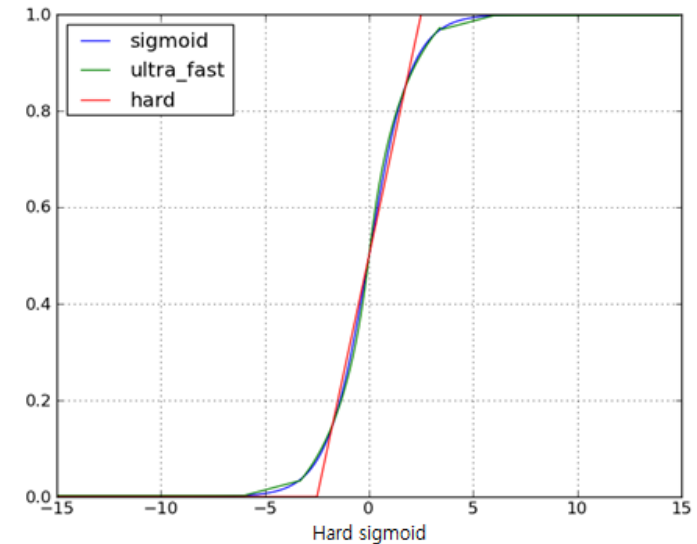
1.12	3.42	-1.5	-12	1	1	-1	-1
32	-1	-5	15	1	-1	-1	1
24	0.55	-54	0.24	1	1	-1	1
-0.1	0.1	-0.2	2	-1	1	-1	1

### Stochastic Binarization

$$x^b = \begin{cases} +1 & \text{with probability } p = \sigma(x), \\ -1 & \text{with probability } 1 - p, \end{cases} \quad (2)$$

where

$$\sigma(x) = \text{clip}\left(\frac{x+1}{2}, 0, 1\right) = \max\left(0, \min\left(1, \frac{x+1}{2}\right)\right). \quad (3)$$



# 기술적인 해결과제 : BNN 구현

## Straight-through estimator?

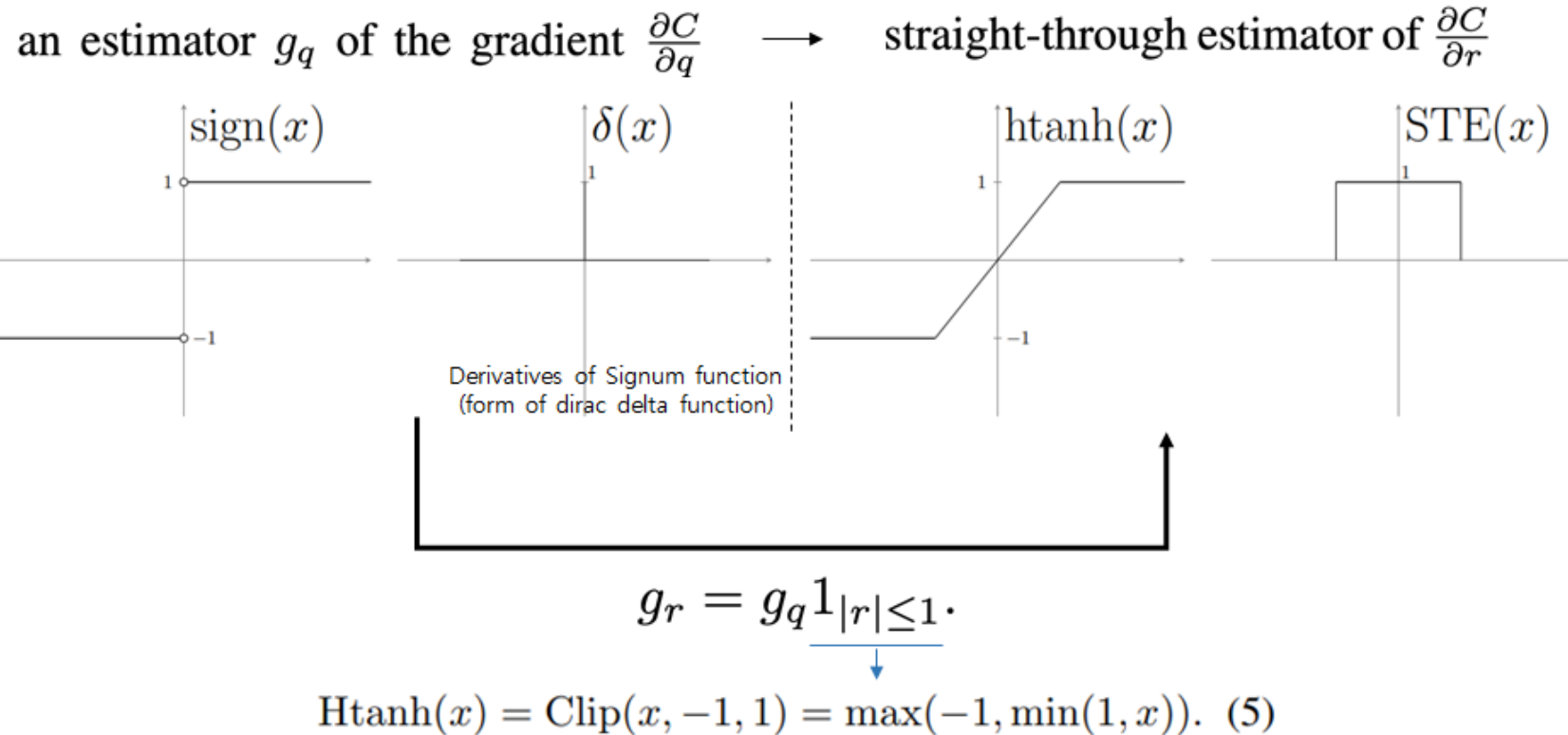


Figure from Darabi, Sajad, et al. "BNN+: Improved binary network training." arXiv preprint arXiv:1812.11800 (2018).



# 기술적인 해결과제 : BNN 구현

## Algorithm 1\*

**Algorithm 1** Training a BNN.  $C$  is the cost function for minibatch,  $\lambda$  - the learning rate decay factor and  $L$  the number of layers.  $\circ$  indicates element-wise multiplication. The function `Binarize()` specifies how to (stochastically or deterministically) binarize the activations and weights, and `Clip()`, how to clip the weights. `BatchNorm()` specifies how to batch-normalize the activations, using either batch normalization (Ioffe & Szegedy, 2015) or its shift-based variant we describe in Algorithm 3. `BackBatchNorm()` specifies how to backpropagate through the normalization. `Update()` specifies how to update the parameters when their gradients are known, using either ADAM (Kingma & Ba, 2014) or the shift-based AdaMax we describe in Algorithm 4.

**Require:** a minibatch of inputs and targets  $(a_0, a^*)$ , previous weights  $W$ , previous BatchNorm parameters  $\theta$ , weights initialization coefficients from (Glorot & Bengio, 2010)  $\gamma$ , and previous learning rate  $\eta$ .

**Ensure:** updated weights  $W^{t+1}$ , updated BatchNorm parameters  $\theta^{t+1}$  and updated learning rate  $\eta^{t+1}$ .

```
{ 1. Computing the parameters gradients:}
{ 1.1. Forward propagation:}
for  $k = 1$  to  $L$  do
   $W_k^b \leftarrow \text{Binarize}(W_k)$ 
   $s_k \leftarrow a_{k-1}^b W_k^b$ 
   $a_k \leftarrow \text{BatchNorm}(s_k, \theta_k)$ 
  if  $k < L$  then
     $a_k^b \leftarrow \text{Binarize}(a_k)$ 
  end if
end for
{ 1.2. Backward propagation:}
{ Please note that the gradients are not binary.}
Compute  $g_{a_L} = \frac{\partial C}{\partial a_L}$  knowing  $a_L$  and  $a^*$ 
for  $k = L$  to  $1$  do
  if  $k < L$  then
     $g_{a_k} \leftarrow g_{a_k^b} \circ 1_{|a_k| \leq 1}$ 
  end if
   $(g_{s_k}, g_{\theta_k}) \leftarrow \text{BackBatchNorm}(g_{a_k}, s_k, \theta_k)$ 
   $g_{a_{k-1}^b} \leftarrow g_{s_k} W_k^b$ 
   $g_{W_k^b} \leftarrow g_{s_k}^{\top} a_{k-1}^b$ 
end for
{ 2. Accumulating the parameters gradients:}
for  $k = 1$  to  $L$  do
   $\theta_k^{t+1} \leftarrow \text{Update}(\theta_k, \eta, g_{\theta_k})$ 
   $W_k^{t+1} \leftarrow \text{Clip}(\text{Update}(W_k, \gamma_k \eta, g_{W_k^b}), -1, 1)$ 
   $\eta^{t+1} \leftarrow \lambda \eta$ 
end for
```

\*Hubara, Itay, et al. "Binarized neural networks." Advances in neural information processing systems. 2016.

# 기술적인 해결과제 : BNN 구현

## Algorithm 1: Training a BNN

**Ensure:** updated weights  $W^{t+1}$ , updated BatchNorm parameters  $\theta^{t+1}$  and updated learning rate  $\eta^{t+1}$ .

{1. Computing the parameters gradients:}

{1.1. Forward propagation:}

**for**  $k = 1$  to  $L$  **do**

$W_k^b \leftarrow \text{Binarize}(W_k)$

$s_k \leftarrow a_{k-1}^b W_k^b$

$a_k \leftarrow \text{BatchNorm}(s_k, \theta_k)$

**if**  $k < L$  **then**

$a_k^b \leftarrow \text{Binarize}(a_k)$

**end if**

**end for**

{1.2. Backward propagation:}

{Please note that the gradients are not binary.}

Compute  $g_{a_L} = \frac{\partial C}{\partial a_L}$  knowing  $a_L$  and  $a^*$

**for**  $k = L$  to 1 **do**

**if**  $k < L$  **then**

$g_{a_k} \leftarrow g_{a_k^b} \circ 1_{|a_k| \leq 1}$

**end if**

$(g_{s_k}, g_{\theta_k}) \leftarrow \text{BackBatchNorm}(g_{a_k}, s_k, \theta_k)$

$g_{a_{k-1}^b} \leftarrow g_{s_k} W_k^b$

$g_{W_k^b} \leftarrow g_{s_k} a_{k-1}^b$

**end for**

{2. Accumulating the parameters gradients:}

**for**  $k = 1$  to  $L$  **do**

$\theta_k^{t+1} \leftarrow \text{Update}(\theta_k, \eta, g_{\theta_k})$

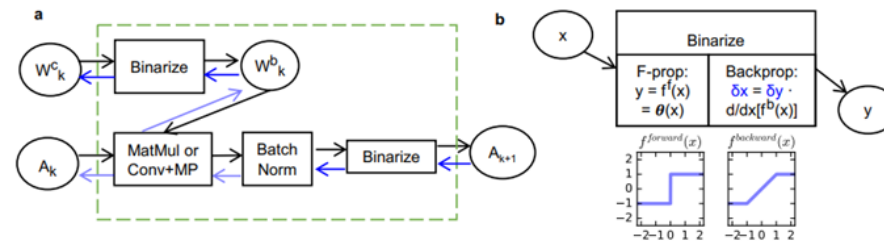
$W_k^{t+1} \leftarrow \text{Clip}(\text{Update}(W_k, \gamma_k \eta, g_{W_k^b}), -1, 1)$

$\eta^{t+1} \leftarrow \lambda \eta$

**end for**

- Binary Activation : Signum Activation Layer (backward: using STE)
  - Binary Weight : When using a weight  $w^r$ , quantize it using  $w^b = \text{Sign}(w^r)$
- Constrain each real-valued weight between -1 and 1.  
 → Clip weights during training
- The real-valued weights would otherwise grow very large

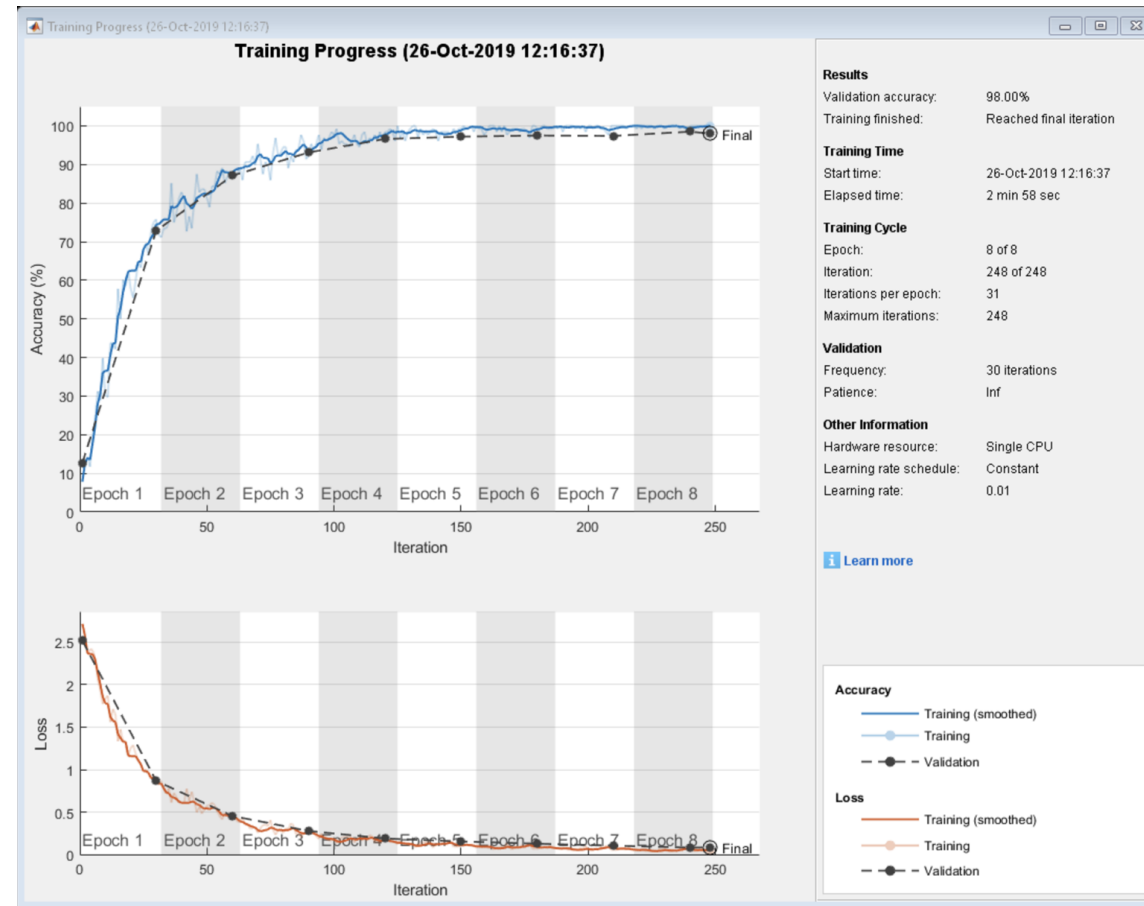
without any impact on the binary weights.



# MathWorks 솔루션을 통한 해결 방안 및 결과

## 기존 함수를 활용한 Model Training

```
layers = [  
    imageInputLayer([28 28 1])  
  
    convolution2dLayer(3,8,'Padding','same')  
    batchNormalizationLayer  
    reluLayer  
  
    maxPooling2dLayer(2,'Stride',2)  
  
    convolution2dLayer(3,16,'Padding','same')  
    batchNormalizationLayer  
    reluLayer  
  
    maxPooling2dLayer(2,'Stride',2)  
  
    convolution2dLayer(3,32,'Padding','same')  
    batchNormalizationLayer  
    reluLayer  
  
    fullyConnectedLayer(10)  
    softmaxLayer  
    classificationLayer];  
  
options = trainingOptions('sgdm', ...  
    'MaxEpochs',8, ...  
    'ValidationData',{XValidation,YValidation}, ...  
    'ValidationFrequency',30, ...  
    'Verbose',false, ...  
    'Plots','training-progress');  
  
net = trainNetwork(XTrain,YTrain,layers,options);
```



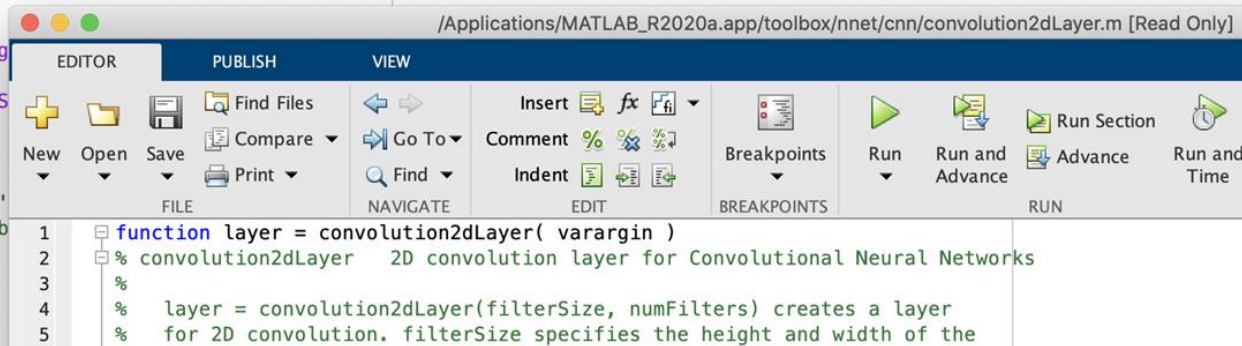
<https://kr.mathworks.com/help/deeplearning/ref/trainingoptions.html>

# MathWorks 솔루션을 통한 해결 방안 및 결과

내부 구조는 어떻게 생겼나?

```
layers = [  
    imageInputLayer([imgsize 3], 'Name', 'Input', 'AverageImage', img_mean) %% Should have Mean  
  
    convolution2dLayer(5,64, 'Padding', 2, 'Stride', 1, 'Name', 'Conv1')  
    batchNormalizationLayer('Name', 'bn1')  
    reluLayer('Name', 'r1')  
    maxPooling2dLayer(3, 'Stride', 2, 'Padding', 1, 'Name', 'Pool1')  
  
    convolution2dLayer(5,64, 'Padding', 2, 'Stride', 1, 'Name', 'Conv2')  
    batchNormalizationLayer('Name', 'bn2')  
    reluLayer('Name', 'r2')  
    maxPooling2dLayer(3, 'Stride', 2, 'Padding', 1, 'Name', 'Pool2')  
  
    convolution2dLayer(5,128, 'Padding', 2, 'Stride', 1, 'Name', 'Conv3')  
    batchNormalizationLayer('Name', 'bn3')  
    reluLayer('Name', 'r3')  
  
    fullyConnectedLayer(3072, 'Name', 'fc1')  
    batchNormalizationLayer('Name', 'fc1-bn')  
    reluLayer('Name', 'fc1-r')  
    dropoutLayer(0.5, 'Name', 'drop1')  
];
```

ctrl+D (Mac은 command+shift+D)를 통해 라이브러리 함수 열람



# MathWorks 솔루션을 통한 해결 방안 및 결과

## Convolution Layers를 살펴보면!

convolution2dLayer -> Convolution2D -> Executionstrategy  
-> Convolution2DLayer

```
1 function layer = convolution2dLayer( varargin )
2 % convolution2dLayer 2D convolution layer for Convolutional Neural Network
...
128 internalLayer = nnet.internal.cnn.layer.Convolution2D( args.Name, ...
129 args.FilterSize, ...
130 args.NumChannels, ...
131 args.NumFilters, ...
132 args.Stride, ...
133 args.DilationFactor, ...
134 args.PaddingMode, ...
135 args.PaddingSize);
...
145 layer = nnet.cnn.layer.Convolution2DLayer( internalLayer );
146 layer.WeightsInitializer = args.WeightsInitializer;
147 layer.BiasInitializer = args.BiasInitializer;
148 layer.Weights = args.Weights;
149 layer.Bias = args.Bias;
150 end
```

```
1 classdef Convolution2D < nnet.internal.cnn.layer.FunctionalLayer ...
2 & nnet.internal.cnn.layer.CPUUsableLayer
3 % Convolution2D Implementation of the 2D convolution layer
4
5 % Copyright 2015-2019 The MathWorks, Inc.
...
115 function this = Convolution2D( ...
116 name, filterSize, numChannels, numFilters, stride, dilationFactor, paddingMode, paddingSize)
...
169 function Z = predict( this, X )
170 % predict Forward input data through the layer and output the
171 if( this.usingFilterGroups() )
172 Z = this.predictTwoFilterGroupsWithCaching( X );
173 else
174 Z = this.predictNormal( X );
175 end
176 end
...
178 function [Z, memory] = forward( this, X )
179 % forward Forward propagate data during training
180 memory = [];
181 if( this.usingFilterGroups() )
182 Z = this.forwardTwoFilterGroups( X );
183 else
184 Z = this.forwardNormal( X );
185 end
186 end
...
188 function varargout = backward( this, X, ~, dZ, ~ )
189 % backward Back propagate the derivative of the loss function
190 % through the layer
191 if( this.usingFilterGroups() )
192 [varargout{1:nargout}] = this.backwardTwoFilterGroups( X, [], dZ );
193 else
194 [varargout{1:nargout}] = this.backwardNormal( X, [], dZ, [] );
195 end
196 end
```

```
545 function this = setHostStrategy( this )
546 % setHostStrategy Use MklDnn only if MklDnn is featured on and no dilation
547 noDilation = isequal( this.DilationFactor, [1 1] );
548 if nnet.internal.cnn.host.useMKLDNN && noDilation
549 this.ExecutionStrategy = nnet.internal.cnn.layer.util.Convolution2DHostMklDnnStrategy();
550 else
551 this.ExecutionStrategy = nnet.internal.cnn.layer.util.Convolution2DHostStridedConvStrategy();
552 end
553 end
...
555 function this = setGPUStrategy( this )
556 this.ExecutionStrategy = nnet.internal.cnn.layer.util.Convolution2DGPUStrategy();
557 end
558 end
...
560 methods( Access=protected )
561 function this = setFunctionalStrategy( this )
562 this.ExecutionStrategy = ...
563 nnet.internal.cnn.layer.util.Convolution2DFunctionalStrategy();
564 end
565 end
```

# MathWorks 솔루션을 통한 해결 방안 및 결과

## Convolution Layers를 살펴보면!

convolution2dLayer -> Convolution2D -> Executionstrategy  
-> Convolution2DLayer

```
1 function layer = convolution2dLayer( varargin )
2 % convolution2dLayer 2-D convolution layer for Convolutional Neural Network
...
128 internalLayer = nnet.internal.cnn.layer.Convolution2D( args.Name, ...
129 args.FilterSize, ...
130 args.NumChannels, ...
131 args.NumFilters, ...
132 args.Stride, ...
133 args.DilationFactor, ...
134 args.PaddingMode, ...
135 args.PaddingSize);
...
145 layer = nnet.cnn.layer.Convolution2DLayer( internalLayer);
146 layer.WeightsInitializer = args.WeightsInitializer;
147 layer.BiasInitializer = args.BiasInitializer;
148 layer.Weights = args.Weights;
149 layer.Bias = args.Bias;
150 end
```

```
1 classdef Convolution2D < nnet.internal.cnn.layer.FunctionalLayer ...
2 & nnet.internal.cnn.layer.CPUUsableLayer
3 % Convolution2D Implementation of the 2-D convolution layer
4
5 % Copyright 2015-2019 The MathWorks, Inc.
...
115 function this = Convolution2D( ...
116 name, filterSize, numChannels, numFilters, stride, dilationFactor, paddingMode, paddingSize)
...
169 function Z = predict( this, X )
170 % predict Forward input data through the layer and output
171 if( this.usingFilterGroups() )
172 Z = this.predictTwoFilterGroupsWithCaching( X );
173 else
174 Z = this.predictNormal( X );
175 end
176 end
...
178 function [Z, memory] = forward( this, X )
179 % forward Forward propagate data during training
180 memory = [];
181 if( this.usingFilterGroups() )
182 Z = this.forwardTwoFilterGroups( X );
183 else
184 Z = this.forwardNormal( X );
185 end
186 end
...
188 function varargout = backward( this, X, ~, dZ, ~ )
189 % backward Back propagate the derivative of the loss function
190 % through the layer
191 if( this.usingFilterGroups() )
192 [varargout{1:nargout}] = this.backwardTwoFilterGroups( X, [], dZ );
193 else
194 [varargout{1:nargout}] = this.backwardNormal( X, [], dZ, [] );
195 end
196 end
```

```
545 function this = setHostStrategy( this )
546 % setHostStrategy Use MklDnn only if MklDnn is featured on and no dilation
547 noDilation = isequal( this.DilationFactor, [1 1] );
548 if nnet.internal.cnn.host.UseMKLDNN && noDilation
549 this.ExecutionStrategy = nnet.internal.cnn.layer.util.Convolution2DHostMklDnnStrategy();
550 else
551 this.ExecutionStrategy = nnet.internal.cnn.layer.util.Convolution2DHostStridedConvStrategy();
552 end
553 end
...
555 function this = setGPUstrategy( this )
556 this.ExecutionStrategy = nnet.internal.cnn.layer.util.Convolution2DGPUstrategy();
557 end
558 end
...
560 methods( Access=protected )
561 function this = setFunctionalStrategy( this )
562 this.ExecutionStrategy = ...
563 nnet.internal.cnn.layer.util.Convolution2DFunctionalStrategy();
564 end
565 end
```

# MathWorks 솔루션을 통한 해결 방안 및 결과

## Convolution Layers를 살펴보면!

### Execution strategy

```
1 classdef Convolution2DGPUstrategy < nnet.internal.cnn.layer.util.ExecutionSt
2     % Convolution2DGPUstrategy Execution strategy for running the
3     % convolution on the GPU
4
5
6
7     methods
8         function [Z, memory] = forward(~, X, ...
9             weights, bias, ...
10            topPad, leftPad, ...
11            bottomPad, rightPad, ...
12            verticalStride, horizontalStride, ...
13            verticalDilation, horizontalDilation)
14            paddingSize = [topPad bottomPad leftPad rightPad];
15            if iPaddingIsSymmetric(paddingSize)
16                Z = nnet.internal.cnn.gpu.convolveForward2D( ...
17                    X, weights, ...
18                    topPad, leftPad, ...
19                    bottomPad, rightPad, ...
20                    verticalStride, horizontalStride, ...
21                    verticalDilation, horizontalDilation) + bias;
22            else
23                X = iPadArray(X, paddingSize);
24                Z = nnet.internal.cnn.gpu.convolveForward2D( ...
25                    X, weights, ...
26                    0, 0, ...
27                    0, 0, ...
28                    verticalStride, horizontalStride, ...
29                    verticalDilation, horizontalDilation) + bias;
30            end
31            memory = [];
32        end
33    end
34
```

```
34     function [dX, dW] = backward(~, ...
35         X, weights, dZ, ...
36         topPad, leftPad, ...
37         bottomPad, rightPad, ...
38         strideHeight, strideWidth, ...
39         verticalDilation, horizontalDilation)
40     paddingSize = [topPad bottomPad leftPad rightPad];
41     needsWeightGradients = nargin > 1;
42     if iPaddingIsSymmetric(paddingSize)
43         dX = nnet.internal.cnn.gpu.convolveBackwardData2D( ...
44             X, weights, dZ, ...
45             topPad, leftPad, ...
46             bottomPad, rightPad, ...
47             strideHeight, strideWidth, ...
48             verticalDilation, horizontalDilation);
49     else if needsWeightGradients
50         dW{1} = nnet.internal.cnn.gpu.convolveBackwardFilter2D( ...
51             X, weights, dZ, ...
52             topPad, leftPad, ...
53             bottomPad, rightPad, ...
54             strideHeight, strideWidth, ...
55             verticalDilation, horizontalDilation);
56     end
57     else
58         X = iPadArray(X, paddingSize);
59         dX = nnet.internal.cnn.gpu.convolveBackwardData2D( ...
60             X, weights, dZ, ...
61             0, 0, ...
62             0, 0, ...
63             strideHeight, strideWidth, ...
64             verticalDilation, horizontalDilation);
65         dX = iUnpadArray(dX, paddingSize);
66     end
67 end
```

# MathWorks 솔루션을 통한 해결 방안 및 결과

## Binarized Convolutional Layers

Ensure: updated weights  $W^{t+1}$ , updated BatchNorm parameters  $\theta^{t+1}$  and updated learning rate  $\eta^{t+1}$ .

{1. Computing the parameters gradients:}

{1.1. Forward propagation:}

for  $k = 1$  to  $L$  do

$W_k^b \leftarrow \text{Binarize}(W_k)$

$s_k \leftarrow a_{k-1}^b W_k^b$

$a_k \leftarrow \text{BatchNorm}(s_k, \theta_k)$

if  $k < L$  then

$a_k^b \leftarrow \text{Binarize}(a_k)$

end if

end for

{1.2. Backward propagation:}

{Please note that the gradients are not binary.}

Compute  $g_{a_L} = \frac{\partial C}{\partial a_L}$  knowing  $a_L$  and  $a^*$

for  $k = L$  to 1 do

if  $k < L$  then

$g_{a_k} \leftarrow g_{a_k}^b \circ 1_{|a_k| \leq 1}$

end if

$(g_{s_k}, g_{\theta_k}) \leftarrow \text{BackBatchNorm}(g_{a_k}, s_k, \theta_k)$

$g_{a_{k-1}^b} \leftarrow g_{s_k} W_k^b$

$g_{W_k^b} \leftarrow g_{s_k} a_{k-1}^b$

end for

{2. Accumulating the parameters gradients:}

for  $k = 1$  to  $L$  do

$\theta_k^{t+1} \leftarrow \text{Update}(\theta_k, \eta, g_{\theta_k})$

$W_k^{t+1} \leftarrow \text{Clip}(\text{Update}(W_k, \gamma_k \eta, g_{W_k^b}), -1, 1)$

$\eta^{t+1} \leftarrow \lambda \eta$

end for

다른 이름으로 저장 후 수정!

```

8         function [Z, memory] = forward(~, X, ...
9             weights, bias, ...
10            topPad, leftPad, ...
11            bottomPad, rightPad, ...
12            verticalStride, horizontalStride, ...
13            verticalDilation, horizontalDilation)
14 -         paddingSize = [topPad bottomPad leftPad rightPad];
15
16 -         weights= sign(weights); %%%%%%%%%%%%%%%
17 -         weights(weights==0)=1; %%%%%%%%%%%%%%%
    
```

```

38        function [dX,dW] = backward(~, ...
39            X, weights, dZ, ...
40            topPad, leftPad, ...
41            bottomPad, rightPad, ...
42            strideHeight, strideWidth, ...
43            verticalDilation, horizontalDilation)
44 -         paddingSize = [topPad bottomPad leftPad rightPad];
45 -         needsWeightGradients = nargin > 1;
46
47 -         weightsFP=weights; %%%%%%%%%%%%%%%
48 -         weights= sign(weights); %%%%%%%%%%%%%%%
49 -         weights(weights==0)=1; %%%%%%%%%%%%%%%
    
```

이진화 컨벌루션 함수 생성!

- Binarizedconvolution2dLayer.m
- BinarizedConvolution2D.m
- BinarizedConvolution2DFunctionalStrategy.m
- BinarizedConvolution2DGPUStrategy.asv
- BinarizedConvolution2DGPUStrategy.m
- BinarizedConvolution2DHostMklDnnStrategy.m
- BinarizedConvolution2DHostStridedConvStrategy.m



# MathWorks 솔루션을 통한 해결 방안 및 결과

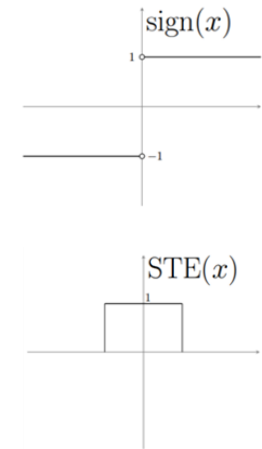
Custom Layer를 활용하여 필요한 활성 함수 만들기!

```
classdef myLayer < nnet.layer.Layer
    properties
        % (Optional) Layer properties.
        % Layer properties go here.
    end
    properties (Learnable)
        % (Optional) Layer learnable parameters.
        % Layer learnable parameters go here.
    end
    methods
        function layer = myLayer()
            % (Optional) Create a myLayer.
            % This function must have the same name as the class.
            % Layer constructor function goes here.
        end
        function [Z1, ..., Zm] = predict(layer, X1, ..., Xn)
            % Forward input data through the layer at prediction time and
            % output the result.
            %
            % Inputs:
            %     layer - Layer to forward propagate through
            %     X1, ..., Xn - Input data
            % Outputs:
            %     Z1, ..., Zm - Outputs of layer forward function
            % Layer forward function for prediction goes here.
        end
    end
end
```



Sign Activation Layer → Binary Activation!

```
1 classdef SignumActivation < nnet.layer.Layer
2     methods
3         function layer = SignumActivation(name)
4             layer.Name = name;
5         end
6
7         function [Z,memory] = forward(~, X)
8             Z=sign(X);
9             Z(Z==0)=1;
10            memory=[];
11        end
12
13        function Z = predict(~, X)
14            Z=sign(X);
15            Z(Z==0)=1;
16        end
17
18        function dLdX = backward(~, X, ~, dLdZ, ~)
19            dLdX = (dLdZ) .* (X>-1 & X<1) ;
20        end
21    end
22 end
```



SignumActivation.m

<https://kr.mathworks.com/help/deeplearning/ug/define-custom-deep-learning-layers.html>

# MathWorks 솔루션을 통한 해결 방안 및 결과

## New Functions for BNN Training

**Ensure:** updated weights  $W^{t+1}$ , updated BatchNorm parameters  $\theta^{t+1}$  and updated learning rate  $\eta^{t+1}$ .

{1. Computing the parameters gradients:}

{1.1. Forward propagation:}

for  $k = 1$  to  $L$  do

$W_k^b \leftarrow \text{Binarize}(W_k)$

$s_k \leftarrow a_{k-1}^b W_k^b$

$a_k \leftarrow \text{BatchNorm}(s_k, \theta_k)$

if  $k < L$  then

$a_k^b \leftarrow \text{Binarize}(a_k)$

end if

end for

{1.2. Backward propagation:}

{Please note that the gradients are not binary.}

Compute  $g_{a_L} = \frac{\partial C}{\partial a_L}$  knowing  $a_L$  and  $a^*$

for  $k = L$  to 1 do

if  $k < L$  then

$g_{a_k} \leftarrow g_{a_k^b} \circ 1_{|a_k| \leq 1}$

end if

$(g_{s_k}, g_{\theta_k}) \leftarrow \text{BackBatchNorm}(g_{a_k}, s_k, \theta_k)$

$g_{a_{k-1}^b} \leftarrow g_{s_k} W_k^b$

$g_{W_k^b} \leftarrow g_{s_k}^T a_{k-1}^b$

end for

{2. Accumulating the parameters gradients:}

for  $k = 1$  to  $L$  do

$\theta_k^{t+1} \leftarrow \text{Update}(\theta_k, \eta, g_{\theta_k})$

$W_k^{t+1} \leftarrow \text{Clip}(\text{Update}(W_k, \gamma_k \eta, g_{W_k^b}), -1, 1)$

$\eta^{t+1} \leftarrow \lambda \eta$

end for

### CNN의 1-bit quantization을 위한 함수들

- Binarizedconvolution2dLayer.m
- BinarizedConvolution2D.m
- BinarizedConvolution2DFunctionalStrategy.m
- BinarizedConvolution2DGPUStrategy.asv
- BinarizedConvolution2DGPUStrategy.m
- BinarizedConvolution2DHostMklDnnStrategy.m
- BinarizedConvolution2DHostStridedConvStrategy.m
- SignumActivation.m

```
layers = [  
    imageInputLayer([32 32 3], 'Name', 'input', 'Normalization', 'none')  
  
    Binarizedconvolution2dLayer(3,64, 'Padding', 'same', 'Stride', 1, 'BiasLearnRateFactor', 0, 'Name', 'binConv1')  
    batchNormalizationLayer('Name', 'BatchNorm1')  
    SignumActivation('Sign1')  
    Binarizedconvolution2dLayer(3,64, 'Padding', 'same', 'Stride', 1, 'BiasLearnRateFactor', 0, 'Name', 'binConv2')  
    batchNormalizationLayer('Name', 'BatchNorm2')  
    SignumActivation('Sign2')  
  
    maxPooling2dLayer(2, 'Stride', 2, 'Name', 'MaxPool1')  
  
    Binarizedconvolution2dLayer(3,128, 'Padding', 'same', 'Stride', 1, 'BiasLearnRateFactor', 0, 'Name', 'binConv3')  
    batchNormalizationLayer('Name', 'BatchNorm3')  
    SignumActivation('Sign3')  
    Binarizedconvolution2dLayer(3,128, 'Padding', 'same', 'Stride', 1, 'BiasLearnRateFactor', 0, 'Name', 'binConv4')  
    batchNormalizationLayer('Name', 'BatchNorm4')  
    SignumActivation('Sign4')  
  
    maxPooling2dLayer(2, 'Stride', 2, 'Name', 'MaxPool2')  
  
    Binarizedconvolution2dLayer(3,256, 'Padding', 'same', 'Stride', 1, 'BiasLearnRateFactor', 0, 'Name', 'binConv5')  
    batchNormalizationLayer('Name', 'BatchNorm5')  
    SignumActivation('Sign5')  
    Binarizedconvolution2dLayer(3,256, 'Padding', 'same', 'Stride', 1, 'BiasLearnRateFactor', 0, 'Name', 'binConv6')  
    batchNormalizationLayer('Name', 'BatchNorm6')  
    SignumActivation('Sign6')  
  
    maxPooling2dLayer(2, 'Stride', 2, 'Name', 'MaxPool3')  
  
    averagePooling2dLayer(4, 'Name', 'avePool1')  
    SignumActivation('SignAve')  
    %===== :Classifier: =====%  
    Binarizedconvolution2dLayer(1,10, 'Padding', 'same', 'Stride', 1, 'BiasLearnRateFactor', 0, 'Name', 'binConv7')  
    batchNormalizationLayer('Name', 'BatchNorm7')  
    softmaxLayer('Name', 'softmax')  
    classificationLayer('Name', 'classoutput')  
];
```

BNN(VGG7 model)

# MathWorks 솔루션을 통한 해결 방안 및 결과

## CIFAR-10 Training Results

**Validation Data Confusion Matrix**

Output Class	Validation Data Confusion Matrix										Accuracy
	airplane	automobile	bird	cat	deer	dog	frog	horse	ship	truck	
airplane	865 8.6%	9 0.1%	45 0.4%	23 0.2%	22 0.2%	13 0.1%	2 0.0%	13 0.1%	62 0.6%	20 0.2%	80.5%
automobile	16 0.2%	931 9.3%	1 0.0%	8 0.1%	2 0.0%	2 0.0%	1 0.0%	3 0.0%	15 0.1%	45 0.4%	90.9%
bird	34 0.3%	2 0.0%	750 7.5%	41 0.4%	42 0.4%	25 0.3%	31 0.3%	12 0.1%	9 0.1%	4 0.0%	78.9%
cat	14 0.1%	1 0.0%	25 0.3%	630 6.3%	25 0.3%	99 1.0%	26 0.3%	14 0.1%	1 0.0%	4 0.0%	75.1%
deer	9 0.1%	0 0.0%	44 0.4%	56 0.6%	805 8.1%	30 0.3%	12 0.1%	32 0.3%	4 0.0%	3 0.0%	80.9%
dog	4 0.0%	0 0.0%	39 0.4%	122 1.2%	18 0.2%	750 7.5%	8 0.1%	25 0.3%	0 0.0%	1 0.0%	77.6%
frog	5 0.1%	1 0.0%	56 0.6%	69 0.7%	49 0.5%	35 0.4%	909 9.1%	8 0.1%	7 0.1%	2 0.0%	79.7%
horse	8 0.1%	2 0.0%	27 0.3%	25 0.3%	32 0.3%	35 0.4%	4 0.0%	885 8.8%	2 0.0%	2 0.0%	86.6%
ship	33 0.3%	9 0.1%	5 0.1%	8 0.1%	3 0.0%	5 0.1%	4 0.0%	1 0.0%	884 8.8%	11 0.1%	91.8%
truck	12 0.1%	45 0.4%	8 0.1%	18 0.2%	2 0.0%	6 0.1%	3 0.0%	7 0.1%	16 0.2%	908 9.1%	88.6%
	86.5%	93.1%	75.0%	63.0%	80.5%	75.0%	90.9%	88.5%	88.4%	90.8%	83.2%
	13.5%	6.9%	25.0%	37.0%	19.5%	25.0%	9.1%	11.5%	11.6%	9.2%	16.8%
	airplane	automobile	bird	cat	deer	dog	frog	horse	ship	truck	
											<b>Target Class</b>

TABLE III. COMPARISON OF BINARIZED VGG-11 CNNs (NOTE THAT, THE INTEGER CONVOLUTIONAL LAYER (ICONV1) USES 1 BIT WEIGHT AND 8 BIT INPUT.).

Layer	Baseline				Neuron Pruning				Proposed			
	Output Dim.	Input # Fmaps	Output Fmaps	Weight [bits]	Output Dim.	Input # Fmaps	Output Fmaps	Weight [bits]	Output Dim.	Input # Fmaps	Output Fmaps	Weight [bits]
ICov1	32x 32	3	64	1.7K	32x 32	3	64	1.7K	32x 32	3	64	1.7K
BCov2	32x 32	64	64	36.8K	32x 32	64	64	36.8K	32x 32	64	64	36.8K
Max Pool	16x 16	64	64		16x 16	64	64		16x 16	64	64	
BCov3	16x 16	64	128	73.7K	16x 16	64	128	73.7K	16x 16	64	128	73.7K
BCov4	16x 16	128	128	147.4K	16x 16	128	128	147.4K	16x 16	128	128	147.4K
Max Pool	8x 8	128	128		8x 8	128	128		8x 8	128	128	
BCov5	8x 8	128	256	294.9K	8x 8	128	256	294.9K	8x 8	128	256	294.9K
BCov6	8x 8	256	256	589.8K	8x 8	256	256	589.8K	8x 8	256	256	589.8K
Max Pool	4x 4	256	256		4x 4	256	256		4x 4	256	256	
BFC1	1x 1	4096	4096	16.7M	32x 32	4096	354	14.4M	(Ave Pool)			
BFC2	1x 1	4096	4096	16.7M	32x 32	354	31	10.9K				
BFC3	1x 1	4096	10	40.9K	32x 32	31	10	310	1x 1	256	10	2.5K
(fc total)				(33.6M)				(14.6M)				(2.5K)
Total				34.7M				26.0M				11.5M
Error Rate				18.6%				19.1%				18.2%

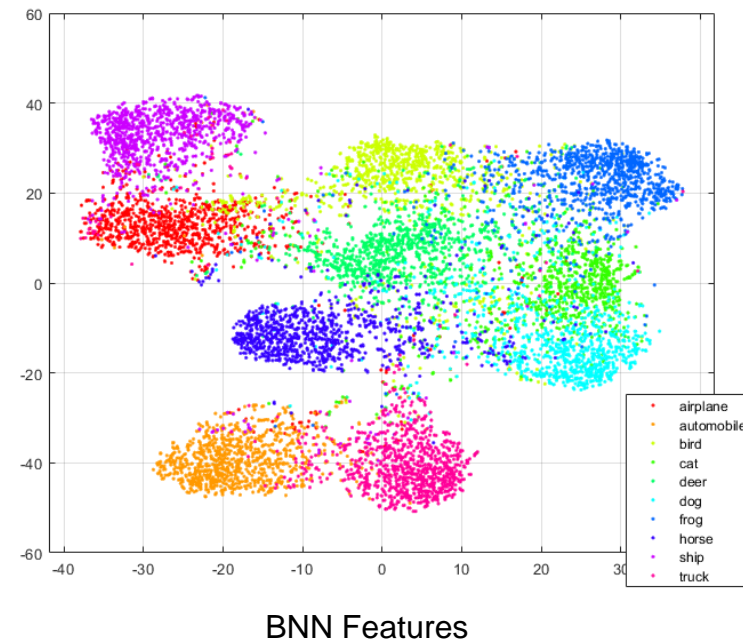
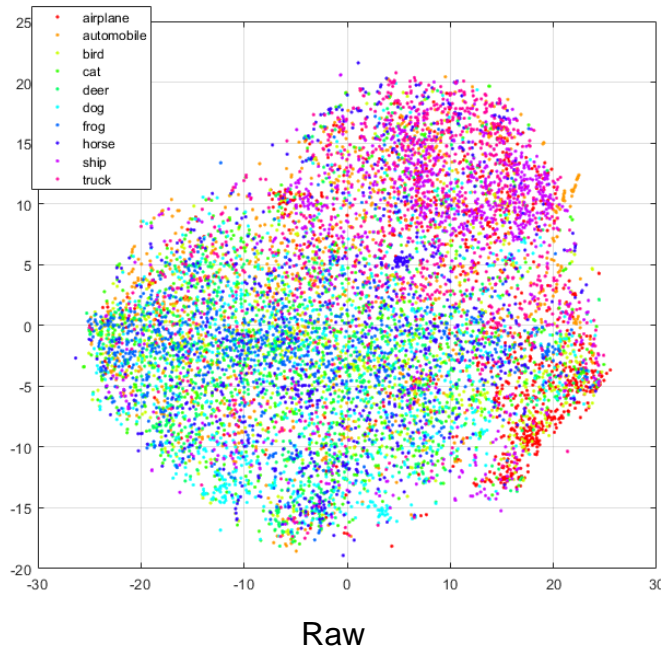
Error Rate = 18.2%

Paper Reference\*

\*Nakahara, Hiroki, Tomoya Fujii, and Shimpei Sato. "A fully connected layer elimination for a binarized convolutional neural network on an FPGA." 2017 27th International Conference on Field Programmable Logic and Applications (FPL). IEEE, 2017.

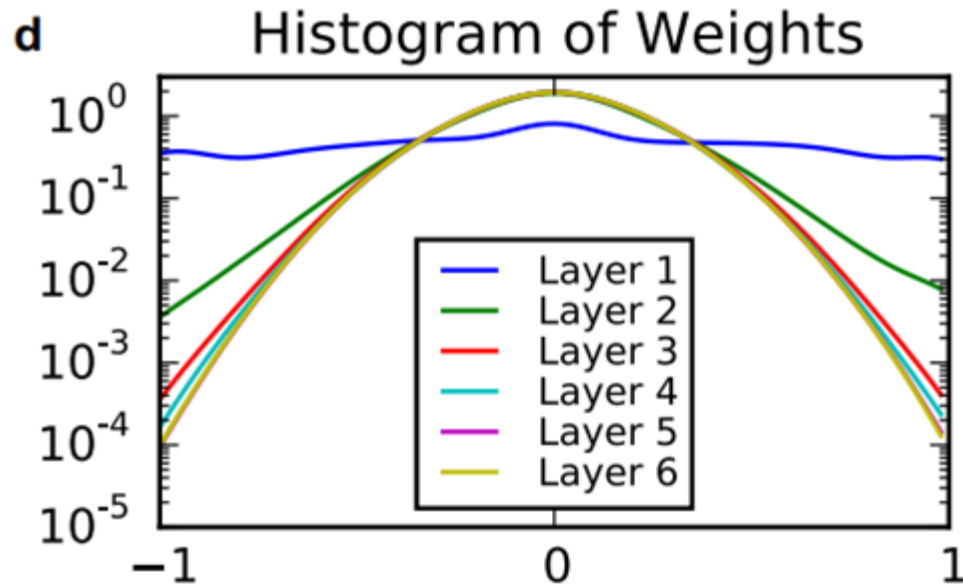
# MathWorks 솔루션을 통한 해결 방안 및 결과

## Latent Space Visualization(t-SNE)



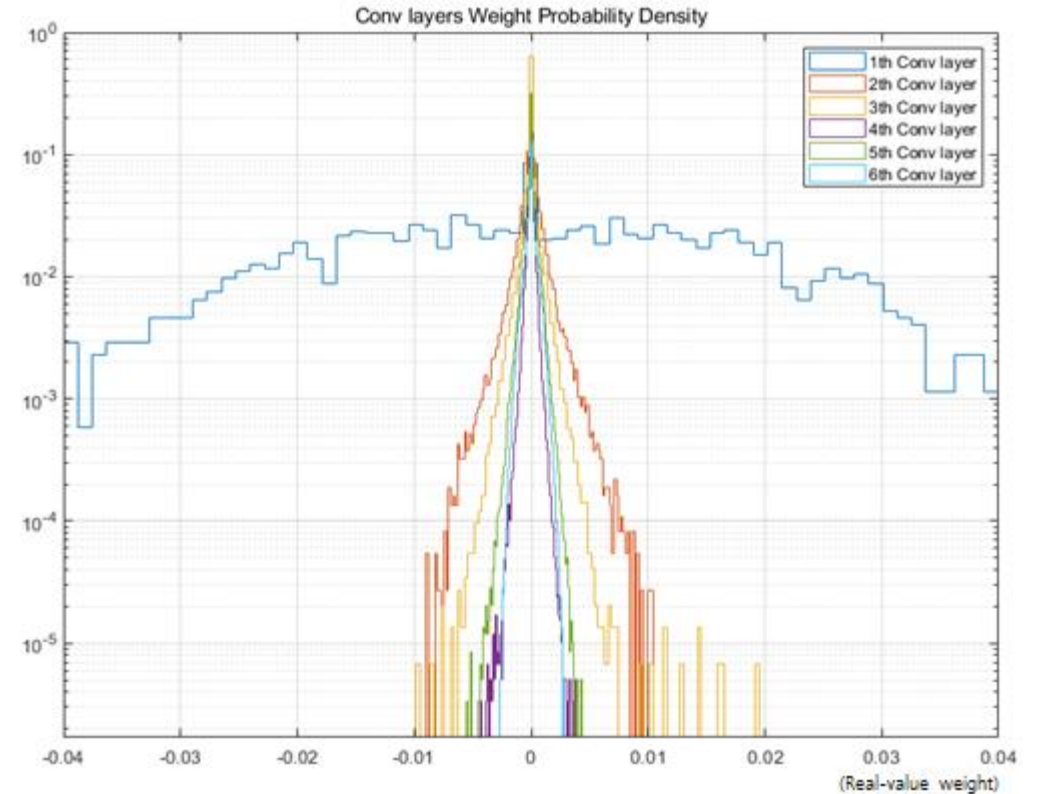
# MathWorks 솔루션을 통한 해결 방안 및 결과

## BNN Training Results



Anderson, Alexander G., and Cory P. Berg. "The high-dimensional geometry of binary neural networks." ICLR2018.

What paper said



What I've got

# 결과 및 정리 (Achievements and Outlook)

- BNN은 신경망 양자화(Quantization)기법 중 하나로 모델의 가중치와 활성화 값을 1-bit으로 양자화
- 1-bit 양자화를 통해 얻을 수 있는 이득으로는..
  1. 모델 사이즈 감소
  2. 메모리 소비 감소
  3. 연산 가속 효과 기대
  4. 하드웨어 구성 시 집적도 향상
- MATLAB을 통한 BNN 논문 구현 방법을 소개하고 결과물을 기존 MATLAB 함수들을 사용하여 평가

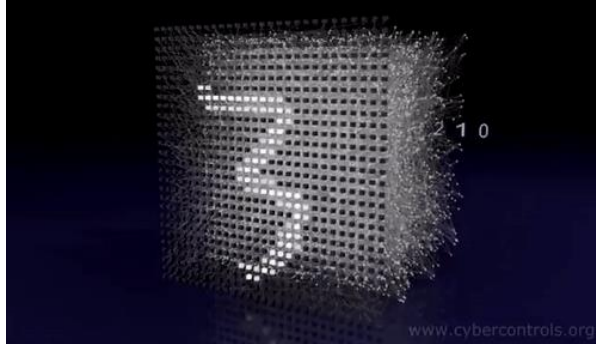
MATLAB으로 논문을 구현하며 좋았던 장점은...

- 레이어 별 함수 내부로 접근이 가능하여 구조를 살펴볼 수 있음
- 필요한 기능을 구현하기 위해 해당 함수 및 클래스에 직접 찾아가 모듈별로 직관적으로 구현 용이

MATLAB으로 구현하며 겪은 단점은..

- 모든 함수에 접근이 가능한 것은 아님
- 주석이나 설명이 불친절 한 경우나 함수가 복잡하게 얽혀 있는 경우가 많음

# 결론 (Concluding Remarks)



Jetson Nano



Raspberry Pi



PC/CPU/GPU



ASIC



FPGA

- 딥러닝 네트워크 모델은 최종 단계에서 다양한 환경에 배포됨
- 엣지 환경에서 클라우드 환경까지 다양한 환경에는 다양한 제약 조건이 존재
- 모델 배포 단계에서 타겟 환경을 고려한 경량화 및 압축 기술은 필수적
- 이를 위해 신경망의 양자화 및 모델 경량화 기술이 중요
- 오늘은 이러한 신경망의 양자화 중 1-bit(이진화) 신경망에 대하여 소개
- 양자화된 네트워크 모델은 추론 시에 가속기와 함께 할 때 효과가 극대화
- 각 환경에는 그 환경에 맞는 다양한 형태의 가속, 연산기가 존재
- 신경망 양자화는 모델 사이즈의 감소와 더불어 연산기와 함께 할 때 연산 가속, 하드웨어의 집적도, 메모리 소모 감소 등 다양한 이득을 얻을 수 있음
- Google Coral, Xilinx FPGA, Arm의 Ethos, Mali, Arm nn, NVIDIA의 TensorRT 등