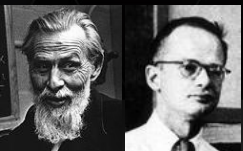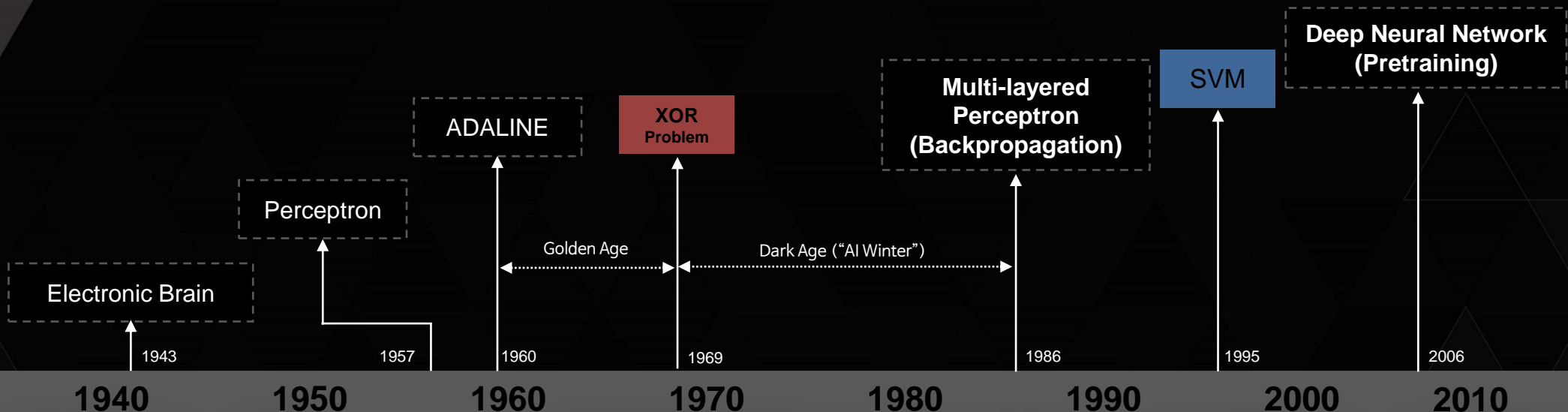# IMPLEMENTING DEEP LEARNING USING CUDNN

이예하 VUNO INC.

# CONTENTS

▸ Deep Learning Review

▸ Implementation on GPU using cuDNN

▸ Optimization Issues

▸ Introduction to VUNO-Net

# DEEP LEARNING REVIEW

# BRIEF HISTORY OF NEURAL NETWORK

Deep Neural Network (Pretraining)

SVM

Multi-layered Perceptron (Backpropagation)

ADALINE

XOR Problem

Perceptron

Golden Age

Dark Age ("AI Winter")

Electronic Brain

1943    1957    1960    1969    1986    1995    2006

**1940    1950    1960    1970    1980    1990    2000    2010**

S. McCulloch – W. Pitts

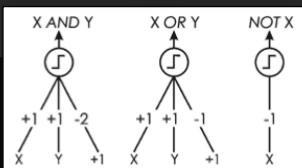F. Rosenblatt

B. Widrow – M. Hoff

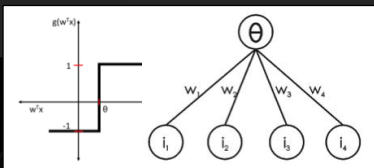M. Minsky – S. Papert

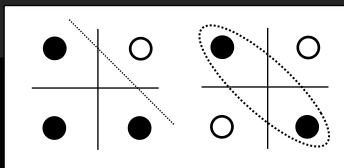D. Rumelhart – G. Hinton – R. Wiliams

V. Vapnik – C. Cortes
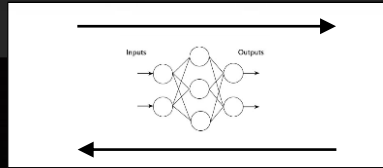
G. Hinton – S. Ruslan

- Adjustable Weights
- Weights are not Learned
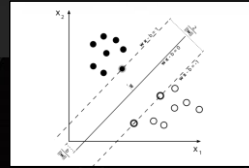
- Learnable Weights and Threshold

- XOR Problem

- Solution to nonlinearly separable problems
- Big computation, local optima and overfitting

- Limitations of learning prior knowledge
- Kernel function: Human Intervention

- Hierarchical feature Learning

# MACHINE/DEEP LEARNING IS EATING THE WORLD!



Machine Intelligence LANDSCAPE

www.shivonzilis.com/machineintelligence — Bloomberg BETA

# BUILDING BLOCKS

- ▸ Restricted Boltzmann machine
- ▸ Auto-encoder
- ▸ Deep belief Network
- ▸ Deep Boltzmann machine
- ▸ Generative stochastic networks
- ▸ Recurrent neural networks
- ▸ **Convolutional neural netwoks**

# CONVOLUTIONAL NEURAL NETWORKS

▸ LeNet-5 (Yann LeCun, 1998)

# CONVOLUTIONAL NEURAL NETWORKS

▸ Alex Net (Alex Krizhevsky et. al., 2012)



▸ GoogleNet (Szegedy et. Al., 2015)

# CONVOLUTIONAL NEURAL NETWORKS

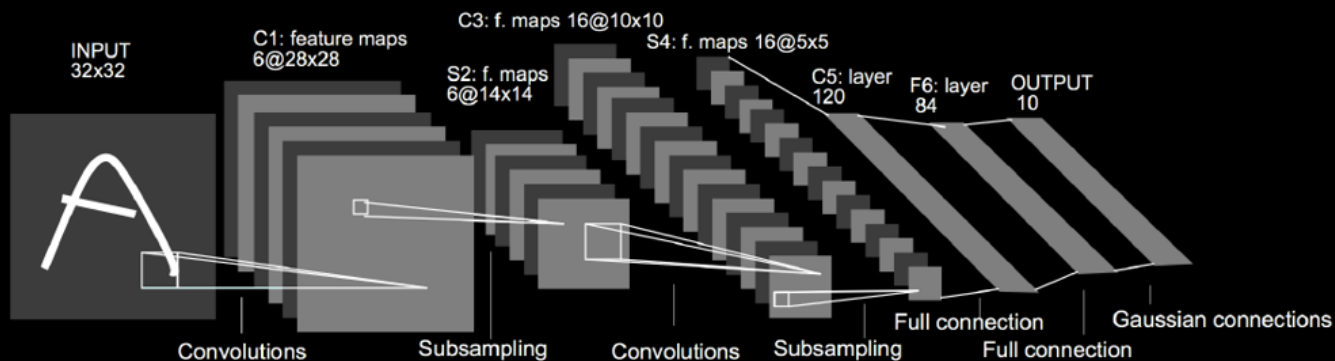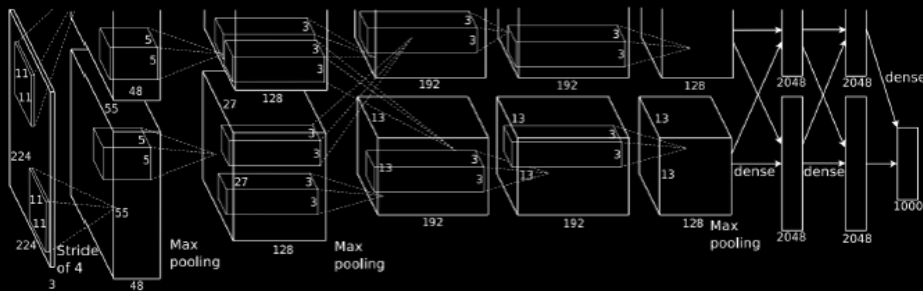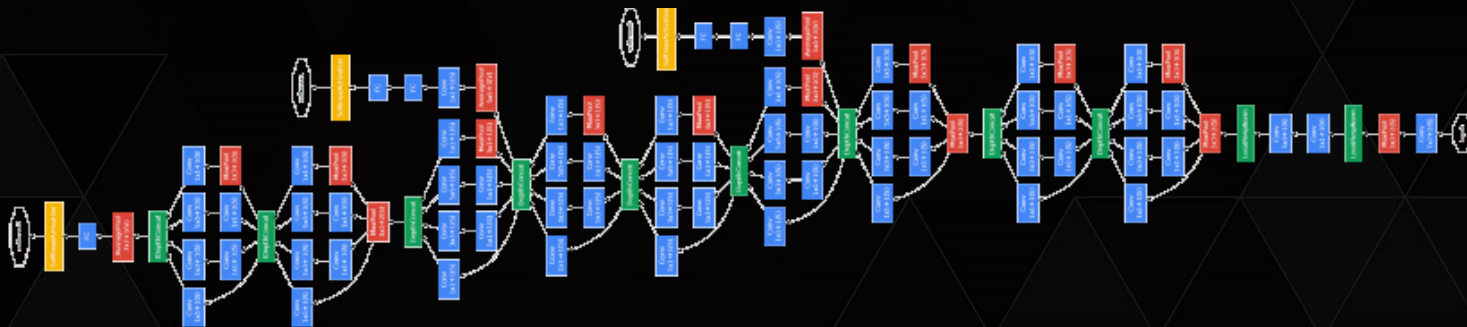| ConvNet Configuration | | | | | |
|---|---|---|---|---|---|
| A | A-LRN | B | C | D | E |
| 11 weight layers | 11 weight layers | 13 weight layers | 16 weight layers | 16 weight layers | 19 weight layers |
| input (224 × 224 RGB image) | | | | | |
| conv3-64 | conv3-64 LRN | conv3-64 conv3-64 | conv3-64 conv3-64 | conv3-64 conv3-64 | conv3-64 conv3-64 |
| maxpool | | | | | |
| conv3-128 | conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 |
| maxpool | | | | | |
| conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 conv1-256 | conv3-256 conv3-256 conv3-256 | conv3-256 conv3-256 conv3-256 conv3-256 |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 conv1-512 | conv3-512 conv3-512 conv3-512 | conv3-512 conv3-512 conv3-512 conv3-512 |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 conv1-512 | conv3-512 conv3-512 conv3-512 | conv3-512 conv3-512 conv3-512 conv3-512 |
| maxpool | | | | | |
| FC-4096 | | | | | |
| FC-4096 | | | | | |
| FC-1000 | | | | | |
| soft-max | | | | | |

Forward Pass

Backward Pass

VGG (K. Simonyan and A. Zisserman, 2015)

▸ Network
  ▸ Softmax Layer (Output)
  ▸ Fully Connected Layer
  ▸ Pooling Layer
  ▸ Convoluti   on Layer

▸ Layer
  ▸ Input / Output
  ▸ Weights
  ▸ Neuron activation

# FULLY CONNECTED LAYER - FORWARD



$$a_1 = w_{11}x_1 + w_{12}x_2 + w_{13}x_3$$

$$y_1 = f(a_1)$$

$$a_2 = w_{21}x_1 + w_{22}x_2 + w_{23}x_3$$

$$y_2 = f(a_2)$$

$$\begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

‣ Matrix calculation is very fast on GPU
   ‣ cuBLAS library

# FULLY CONNECTED LAYER - BACKWARD

**Forward pass:**
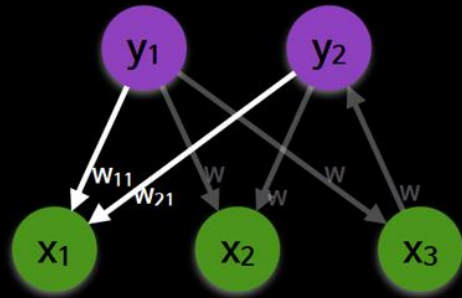$$a_1^{l+1} = w_{11} f(a_1^l) + w_{12} f(a_2^l) + w_{13} f(a_3^l)$$
$$a_2^{l+1} = w_{21} f(a_1^l) + w_{22} f(a_2^l) + w_{23} f(a_3^l)$$

**Error:**
$$\frac{\partial L}{\partial a_i^l} = \sum_{j=1}^{H} \frac{\partial L}{\partial a_j^{l+1}} \frac{\partial a_j^{l+1}}{\partial a_i^l} = \frac{\partial f(a_i^l)}{\partial a_i^l} \boxed{\sum_{j=1}^{H} \frac{\partial L}{\partial a_j^{l+1}} w_{ji}}$$

**Gradient:**
$$\frac{\partial L}{\partial w_{ji}} = \frac{\partial L}{\partial a_j^{l+1}} \frac{\partial a_j^{l+1}}{\partial w_{ji}} = \boxed{\frac{\partial L}{\partial a_j^{l+1}} f(a_i^l)}$$
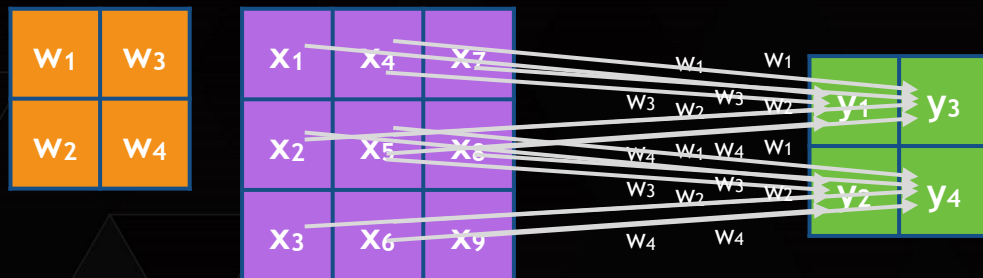
**Error:**
$$\begin{bmatrix} \frac{\partial L}{\partial a_1^l} \\ \frac{\partial L}{\partial a_2^l} \\ \frac{\partial L}{\partial a_3^l} \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix}^T \begin{bmatrix} \frac{\partial L}{\partial a_1^{l+1}} \\ \frac{\partial L}{\partial a_2^{l+1}} \end{bmatrix}$$

**Gradient:**
$$\begin{bmatrix} \frac{\partial L}{\partial w_{11}} & \frac{\partial L}{\partial w_{12}} & \frac{\partial L}{\partial w_{13}} \\ \frac{\partial L}{\partial w_{23}} & \frac{\partial L}{\partial w_{23}} & \frac{\partial L}{\partial w_{23}} \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial a_1^{l+1}} \\ \frac{\partial L}{\partial a_2^{l+1}} \end{bmatrix} \begin{bmatrix} f(a_1^l) \\ f(a_2^l) \\ f(a_3^l) \end{bmatrix}^T$$

▸ Matrix calculation is very fast on GPU

▸ Element-wise multiplication can be done efficiently using GPU thread

# CONVOLUTION LAYER - FORWARD

$$Y_i^{(l)} = B_i^{(l)} + \sum_{j=1}^{m_1^{(l-1)}} K_{i,j}^{(l)} * Y_j^{(l-1)}$$

$y_1 = f(a_1)$ $\qquad a_1 = w_1 x_1 + w_2 x_2 + w_3 x_4 + w_4 x_5$

$y_2 = f(a_2)$ $\qquad a_2 = w_1 x_2 + w_2 x_3 + w_3 x_5 + w_4 x_6$

$y_3 = f(a_3)$ $\qquad a_3 = w_1 x_4 + w_2 x_5 + w_3 x_7 + w_4 x_8$

$y_4 = f(a_4)$ $\qquad a_4 = w_1 x_5 + w_2 x_6 + w_3 x_8 + w_4 x_9$

| $w_1$ | $w_3$ |
|-------|-------|
| $w_2$ | $w_4$ |

| $x_1$ | $x_4$ | $x_7$ |
|-------|-------|-------|
| $x_2$ | $x_5$ | $x_8$ |
| $x_3$ | $x_6$ | $x_9$ |

| $y_1$ | $y_3$ |
|-------|-------|
| $y_2$ | $y_4$ |

# CONVOLUTION LAYER - BACKWARD



$$\frac{\partial L}{\partial a_1^l} = \frac{\partial f\left(a_1^l\right)}{\partial a_1^l}\left(\frac{\partial L}{\partial a_1^{l+1}}w_1\right)$$

$$\frac{\partial L}{\partial a_5^l} = \frac{\partial f\left(a_5^l\right)}{\partial a_5^l}\left(\frac{\partial L}{\partial a_1^{l+1}}w_4 + \frac{\partial L}{\partial a_2^{l+1}}w_3 + \frac{\partial L}{\partial a_3^{l+1}}w_2 + \frac{\partial L}{\partial a_4^{l+1}}w_1\right)$$

$$\frac{\partial L}{\partial a_9^l} = \frac{\partial f\left(a_9^l\right)}{\partial a_9^l}\left(\frac{\partial L}{\partial a_4^{l+1}}w_4\right)$$

# CONVOLUTION LAYER - BACKWARD

# HOW TO EVALUATE THE CONVOLUTION LAYER EFFICIENTLY?

▸ Both Forward and Backward passes can be computed with convolution scheme

▸ Lower the convolutions into a matrix multiplication (cuDNN)

   ▸ There are several ways to implement convolutions efficiently

▸ Fast Fourier Transform to compute the convolution (cuDNN_v3)

▸ Computing the convolutions directly (cuda-convnet)

# INTRODUCTION TO CUDNN

▸ cuDNN is a GPU-accelerated library of primitives for deep neural networks

▸ Convolution forward and backward

▸ Pooling forward and backward

▸ Softmax forward and backward

▸ Neuron activations forward and backward:

  ▸ Rectified linear (ReLU)

  ▸ Sigmoid

  ▸ Hyperbolic tangent (TANH)

▸ Tensor transformation functions

# INTRODUCTION TO CUDNN (VERSION 2)

▸ cuDNN's convolution routines aim for performance competitive with the fastest GEMM

▸ Lowering the convolutions into a matrix multiplication



(Sharan Chetlur et. al., 2015)

# INTRODUCTION TO CUDNN

▸ Benchmarks

**Overfeat [fast]** - Input 128x3x231x231

| Library | Class | Time (ms) | forward (ms) | backward (ms) |
|---|---|---|---|---|
| **CuDNN[R3]-fp16** | cudnn.SpatialConvolution | **313** | **107** | **206** |
| CuDNN[R3]-fp32 | cudnn.SpatialConvolution | 326 | 113 | 213 |
| fbfft | SpatialConvolutionCuFFT | 342 | 114 | 227 |
| Nervana-fp16 | ConvLayer | 355 | 112 | 242 |
| Nervana-fp32 | ConvLayer | 398 | 124 | 273 |
| cudaconvnet2* | ConvLayer | 723 | 176 | 547 |
| CuDNN[R2] * | cudnn.SpatialConvolution | 810 | 234 | 576 |
| Caffe | ConvolutionLayer | 823 | 355 | 468 |
| Torch-7 (native) | SpatialConvolutionMM | 878 | 379 | 499 |
| CL-nn (Torch) | SpatialConvolutionMM | 963 | 388 | 574 |
| Caffe-CLGreenTea | ConvolutionLayer | 2857 | 616 | 2240 |

https://github.com/soumith/convnet-benchmarks

**TRAIN MODELS UP TO 2X FASTER**



https://developer.nvidia.com/cudnn

# LEARNING VGG MODEL USING CUDNN

| ConvNet Configuration | | | | | |
|---|---|---|---|---|---|
| A | A-LRN | B | C | D | E |
| 11 weight layers | 11 weight layers | 13 weight layers | 16 weight layers | 16 weight layers | 19 weight layers |
| input (224 × 224 RGB image) | | | | | |
| conv3-64 | conv3-64 **LRN** | conv3-64 **conv3-64** | conv3-64 conv3-64 | conv3-64 conv3-64 | conv3-64 conv3-64 |
| maxpool | | | | | |
| conv3-128 | conv3-128 | conv3-128 **conv3-128** | conv3-128 conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 |
| maxpool | | | | | |
| conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 **conv1-256** | conv3-256 conv3-256 **conv3-256** | conv3-256 conv3-256 conv3-256 **conv3-256** |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 **conv1-512** | conv3-512 conv3-512 **conv3-512** | conv3-512 conv3-512 conv3-512 **conv3-512** |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 **conv1-512** | conv3-512 conv3-512 **conv3-512** | conv3-512 conv3-512 conv3-512 **conv3-512** |
| maxpool | | | | | |
| FC-4096 | | | | | |
| FC-4096 | | | | | |
| FC-1000 | | | | | |
| soft-max | | | | | |

▸ Data Layer

▸ Convolution Layer

▸ Pooling Layer

▸ Fully Connected Layer

▸ Softmax Layer

# COMMON DATA STRUCTURE FOR LAYER

▸ Device memory & tensor description for input/output data & error

 ▸ Tensor Description defines dimensions of data

```
float                  *d_input, *d_output, *d_inputDelta, *d_outputDelta
cudnnTensorDescriptor_t    inputDesc;
cudnnTensorDescriptor_t    outputDesc;
```

# DATA LAYER

create & set Tensor Descriptor

```
cudnnCreateTensorDescriptor();
cudnnSetTensor4dDescriptor();
```

```
cudnnStatus_t
cudnnSetTensor4dDescriptor( cudnnTensorDescriptor_t    tensorDesc,
                            cudnnTensorFormat_t format,
                            cudnnDataType_t dataType,
                            int n,
                            int c,
                            int h,
                            int w )
```

```
cudnnSetTensor4dDescriptor(
                outputDesc,
                CUDNN_TENSOR_NCHW,
                CUDNN_FLOAT,
                sampleCnt,
                channels,
                height,
                width
);
```

Example: 2 images (3x3x2)

| sample #1 | sample #2 |

channel #1

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

| 19 | 20 | 21 |
| 22 | 23 | 24 |
| 25 | 26 | 27 |

channel #2

| 10 | 11 | 12 |
| 13 | 14 | 15 |
| 16 | 17 | 18 |

| 28 | 29 | 30 |
| 31 | 32 | 33 |
| 34 | 35 | 36 |

# CONVOLUTION LAYER

▸ Initailization

| | |
|---|---|
| 1.1 create & set Filter Descriptor | ```cudnnCreateFilterDescriptor(&filterDesc); cudnnSetFilter4dDescriptor(…);``` |
| ↓ | |
| 1.2 create & set Conv Descriptor | ```cudnnCreateConvolutionDescriptor(&convDesc); cudnnSetConvolution2dDescriptor(…);``` |
| ↓ | |
| 1.3 create & set output Tensor Descriptor | ```cudnnGetConvolution2dForwardOutputDim(…); cudnnCreateTensorDescriptor(&dstTensorDesc); cudnnSetTensor4dDescriptor();``` |
| ↓ | |
| 1.4 Get Convolution Algorithm | ```cudnnGetConvolutionForwardAlgorithm(…); cudnnGetConvolutionForwardWorkspaceSize(…);``` |

# CONVOLUTION LAYER

## 1.1 Set Filter Descriptor

```
cudnnStatus_t
cudnnSetFilter4dDescriptor( cudnnFilterDescriptor_t filterDesc,
                            cudnnDataType_t dataType,
                            int k,
                            int c,
                            int h,
                            int w )
```

```
cudnnSetFilter4dDescriptor(
            filterDesc,
            CUDNN_FLOAT,
            filterCnt,
            input_channelCnt,
            filter_height,
            filter_width
);
```

**Example: 2 Filters (2x2x2)**

| Filter #1 | Filter #2 |
|-----------|-----------|

channel #1

| 1 | 2 | 9 | 10 |
|---|---|---|----|
| 3 | 4 | 11 | 12 |

channel #2

| 5 | 6 | 13 | 14 |
|---|---|----|----|
| 7 | 8 | 15 | 16 |

# CONVOLUTION LAYER

## 1.2 Set Convolution Descriptor

```
cudnnStatus_t
cudnnSetConvolution2dDescriptor( cudnnConvolutionDescriptor_t convDesc,
                                 int pad_h,
                                 int pad_w,
                                 int u,
                                 int v,
                                 int upscalex,
                                 int upscaley,
                                 cudnnConvolutionMode_t mode )
```

**CUDNN_CROSS_CORRELATION**

# CONVOLUTION LAYER

## 1.3 Set output Tensor Descriptor

```
cudnnStatus_t
cudnnGetConvolution2dForwardOutputDim( const cudnnConvolutionDescriptor_t
 convDesc,

                                       const cudnnTensorDescriptor_t
 inputTensorDesc,

                                       const cudnnFilterDescriptor_t filterDesc,
                                       int *n,
                                       int *c,
                                       int *h,
                                       int *w )
```

- n, c, h and w indicate output dimension
  - Tensor Description defines dimensions of data

# CONVOLUTION LAYER

## 1.4 Get Convolution Algorithm

```
cudnnStatus_t
cudnnGetConvolutionForwardAlgorithm( cudnnHandle_t              handle,
                                     const cudnnTensorDescriptor_t    srcDesc,    inputDesc
                                     const cudnnFilterDescriptor_t
filterDesc,
                                     const cudnnConvolutionDescriptor_t
convDesc,
                                     const cudnnTensorDescriptor_t              outputDesc
destDesc,
                                     cudnnConvolutionFwdPreference_t            CUDNN_CONVOLUTION_FWD_PREFER_FASTEST
preference,
                                     size_t
memoryLimitInbytes,
                                     cudnnConvolutionFwdAlgo_t        *algo
                                   )

cudnnStatus_t
cudnnGetConvolutionForwardWorkspaceSize( cudnnHandle_t    handle,
                                         const    cudnnTensorDescriptor_t
 srcDesc,
                                         const    cudnnFilterDescriptor_t
 filterDesc,
                                         const    cudnnConvolutionDescriptor_t
 convDesc,
                                         const    cudnnTensor4dDescriptor_t
 destDesc,
                                         cudnnConvolutionFwdAlgo_t
 algo,
                                         size_t
 *sizeInBytes
                                       )
```

# CONVOLUTION LAYER

## 2.1 Convolution

```
cudnnStatus_t
cudnnConvolutionForward( cudnnHandle_t                          handle,
                         const void                            *alpha,
                         const cudnnTensorDescriptor_t          srcDesc,      d_input
                         const void                            *srcData,      inputDesc
                         const cudnnFilterDescriptor_t          filterDesc,
                         const void                            *filterData,
                         const cudnnConvolutionDescriptor_t     convDesc,
                         cudnnConvolutionFwdAlgo_t              algo,
                         void                                  *workSpace,
                         size_t
workSpaceSizeInBytes,

                         const void                            *beta,

                         const cudnnTensorDescriptor_t          destDesc,     d_output
                         void                                  *destData )    outputDesc
```

# CONVOLUTION LAYER

## 2.2 Activation

```
cudnnStatus_t
cudnnActivationForward( cudnnHandle_t                          handle,
                        cudnnActivationMode_t                  mode,
                        const void                            *alpha,
                        const cudnnTensorDescriptor_t          srcDesc,
                        const void                            *srcData,
                        const void                            *beta,
                        const cudnnTensorDescriptor_t          destDesc,
                        void                                  *destData )
```

sigmoid
tanh
ReLU

outputDesc
d_output

# CONVOLUTION LAYER

## 3.1 Activation Backward

```
cudnnStatus_t
cudnnActivationBackward( cudnnHandle_t                    handle,
                         cudnnActivationMode_t            mode,
                         const void                       *alpha,
                         const cudnnTensorDescriptor_t    srcDesc,          outputDesc
                         const void                       *srcData,         d_output
                         const cudnnTensorDescriptor_t    srcDiffDesc,      outputDesc
                         const void                       *srcDiffData,     d_outputDelta
                         const cudnnTensorDescriptor_t    destDesc,         outputDesc
                         const void                       *destData,        d_output
                         const void                       *beta,
                         const cudnnTensorDescriptor_t    destDiffDesc,     outputDesc
                         void                             *destDiffData )   d_outputDelta
```

- Errors back-propagated from l+1 layer (d_outputDelta) is multiplied by $\frac{\partial f(a_i)}{\partial a_i}$
- See 22 slide (Convolution Layer – Backward)

# CONVOLUTION LAYER

## 3.2 Calculate Gradient

```
cudnnStatus_t
cudnnConvolutionBackwardFilter( cudnnHandle_t                        handle,
                                const void                           *alpha,
                                const cudnnTensorDescriptor_t        srcDesc,      | inputDesc
                                const void                           *srcData,     | d_input
                                const cudnnTensorDescriptor_t        diffDesc,     | outputDesc
                                const void                           *diffData,    | d_outputDelta
                                const cudnnConvolutionDescriptor_t   convDesc,
                                const void                           *beta,

                                const cudnnFilterDescriptor_t        gradDesc,     | filterDesc
                                void                                 *gradData )   | d_filterGradient
```

# CONVOLUTION LAYER

## 3.2 Error Backpropagation

```
cudnnStatus_t
cudnnConvolutionBackwardData( cudnnHandle_t                      handle,
                              const void                        *alpha,
                              const cudnnFilterDescriptor_t      filterDesc,
                              const void                        *filterData,
                              const cudnnTensorDescriptor_t      diffDesc,      outputDesc
                              const void                        *diffData,     d_outputDelta
                              const cudnnConvolutionDescriptor_t convDesc,
                              const void                        *beta,
                              const cudnnTensorDescriptor_t      gradDesc,      inputDesc
                              void                              *gradData      d_inputDelta
);
```

# POOLING LAYER / SOFTMAX LAYER

## Pooling Layer

| | |
|---|---|
| 1. Initilization | cudnnCreatePoolingDescriptor(&poolingDesc);<br>cudnnSetPooling2dDescriptor(…); |
| 2. Forward Pass | cudnnPoolingForward(…); |
| 3. Backward Pass | cudnnPoolingBackward(…); |

## Softmax Layer

| | |
|---|---|
| Forward Pass | cudnnSoftmaxForward(…); |

# OPTIMIZATION ISSUES

# OPTIMIZATION

Learning Very Deep ConvNet

We know the Deep ConvNet can be trained without pre-training

- weight sharing
- sparsity
- Recifier unit

But, "With fixed standard deviations very deep models have difficulties to converges" (Kaiming He et. al., 2015)

- e.g. random initialization from Gaussian dist. with 0.01 std
- >8 convolution layers

# OPTIMIZATION



| ConvNet Configuration | | | | | |
|---|---|---|---|---|---|
| A | A-LRN | B | C | D | E |
| 11 weight layers | 11 weight layers | 13 weight layers | 16 weight layers | 16 weight layers | 19 weight layers |
| input (224 × 224 RGB image) | | | | | |
| conv3-64 | conv3-64 | conv3-64 | conv3-64 | conv3-64 | conv3-64 |
| | LRN | conv3-64 | conv3-64 | conv3-64 | conv3-64 |
| maxpool | | | | | |
| conv3-128 | conv3-128 | conv3-128 | conv3-128 | conv3-128 | conv3-128 |
| | | conv3-128 | conv3-128 | conv3-128 | conv3-128 |
| maxpool | | | | | |
| conv3-256 | conv3-256 | conv3-256 | conv3-256 | conv3-256 | conv3-256 |
| conv3-256 | conv3-256 | conv3-256 | conv3-256 | conv3-256 | conv3-256 |
| | | | conv1-256 | conv3-256 | conv3-256 |
| | | | | | conv3-256 |
| maxpool | | | | | |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
| | | | conv1-512 | conv3-512 | conv3-512 |
| | | | | | conv3-512 |
| maxpool | | | | | |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
| | | | conv1-512 | conv3-512 | conv3-512 |
| | | | | | conv3-512 |
| maxpool | | | | | |
| FC-4096 | | | | | |
| FC-4096 | | | | | |
| FC-1000 | | | | | |
| soft-max | | | | | |

VGG (K. Simonyan and A. Zisserman, 2015)

GoogleNet (C. Szegedy et. al., 2014)

# OPTIMIZATION

## Initialization of Weights for Rectifier (Kaiming He et. al., 2015)

- The variance of the response in each layer

$$Var[\Delta x_2] = Var[\Delta x_{L+1}] \left( \prod_{l=2}^{L} \frac{1}{2} \hat{n}_l Var[w_l] \right)$$

- Sufficient condition that the gradient is not exponentially large/small

$$\frac{1}{2} \hat{n}_l Var[w_l] = 1, \quad \forall l$$

- Standard deviation for initialization

$$\sqrt{2/\hat{n}_l} \qquad \hat{n}_l = k_l^2 d_l$$

(spatial filter size)$^2$ x (filter Cnt)

# OPTIMIZATION

## Case study

| ConvNet Configuration | | | | | |
|---|---|---|---|---|---|
| A | A-LRN | B | C | D | E |
| 11 weight layers | 11 weight layers | 13 weight layers | 16 weight layers | 16 weight layers | 19 weight layers |
| input (224 × 224 RGB image) | | | | | |
| conv3-64 | conv3-64 LRN | conv3-64 conv3-64 | conv3-64 conv3-64 | conv3-64 conv3-64 | conv3-64 conv3-64 |
| maxpool | | | | | |
| conv3-128 | conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 |
| maxpool | | | | | |
| conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 conv1-256 | conv3-256 conv3-256 conv3-256 | conv3-256 conv3-256 conv3-256 conv3-256 |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 conv1-512 | conv3-512 conv3-512 conv3-512 | conv3-512 conv3-512 conv3-512 conv3-512 |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 conv1-512 | conv3-512 conv3-512 conv3-512 | conv3-512 conv3-512 conv3-512 conv3-512 |
| maxpool | | | | | |
| FC-4096 | | | | | |
| FC-4096 | | | | | |
| FC-1000 | | | | | |
| soft-max | | | | | |

3x3 filter

| The filter number | $\sqrt{2/\hat{n}_l}$ |
|---|---|
| 64 | 0.059 |
| 128 | 0.042 |
| 256 | 0.029 |
| 512 | 0.021 |

- When using 0.01, the std of the gradient propagated from conv10 to convey

$$1/(5.9 \times 4.2^2 \times 2.9^2 \times 2.1^4) = 1/(1.7 \times 10^4)$$

**Error vanishing**

# SPEED

## Data loading & model learning

- Reducing data loading and augmentation time
  - Data provider thread (dp_thread)
  - Model learning thread (worker_thread)

```
readData();
for(...) {
    readData();
    pthread_create(worker_thread)
    ...
    pthread_join(worker_thread)
}
```

```
readData()
{
    if(is_dp_thread_running) pthread_join(dp_thread)
    ...
    if(is_data_remain) pthread_create(dp_thread)
}
```

# SPEED

## Multi-GPU

- Data parallelization v.s. Model parallelization
  - Distribute the model, use the same data : Model Parallelism
  - Distribute the data, use the same model : Data Parallelism

- Data parallelization & Gradient Average
  - One of the easiest way to use Multi-GPU
  - The result is same with using single GPU

# PARALLELISM



## Data Parallelization

Forward     Backward     GPU cluster communication

Exchange gradients   5GB/s

2x 16GB/s

+200GB/s   Hidden layer I

+200GB/s   Hidden layer I

+200GB/s   Hidden layer II

Exchange gradients   5GB/s

2x 16GB/s

+200GB/s   Hidden layer II

The good : Easy to implement
The bad : Cost of sync increases with the number of GPU

## Model Parallelization

Forward     Backward     GPU cluster communication

+200GB/s   Hidden layer I    +200GB/s   Hidden layer I

Exchange outputs   Exchange deltas   5GB/s

2x 16GB/s    2x 16GB/s

+200GB/s   Hidden layer II    +200GB/s   Hidden layer II

Exchange outputs   Exchange deltas   5GB/s

2x 16GB/s    2x 16GB/s

The good : Larger network can be trained
The bad : Sync is necessary in all layers

# PARALLELISM



Mixing Data Parallelization and Model Parallelization

(Krizhevsky, 2014)

# PARALLELISM



Data Parallelization & Gradient Average

# VUNO-NET

# VUNO-NET

| Structure | Output | Learning |
|---|---|---|
| Convolution | Softmax | Multi-GPU Support |
| LSTM | Regression | Batch Normalization |
| MD-LSTM(2D) | Connectionist Temporal Classification | Parametric Rectifier |
| Pooling | | Initialization for Rectifier |
| Spatial Pyramid Pooling | | Stochastic Gradient Descent |
| Fully Connection | | Dropout |
| Concatenation | | Data Augmentation |

# PERFORMANCE



The state-of-the-art performance at image & speech

Speech Recognition (TIMIT*)

WR1: 82.2%  VUNO: 84.3%

Image Classification (CIFAR-100)

Kaggle Top 1 (2014) WR2: 79.3%  VUNO: 82.1%

* TIMIT is a one of most popular benchmark dataset for speech recognition task (Texas Instrument – MIT)
WR1 (World Record) – "Speech Recognitino with Deep Recurrent Neural Networks", Alex Graves, ICCASP (2013)
WR2 (World Record) – "kaggle competition: https://www.kaggle.com/c/cifar-10"

# APPLICATION



Whole Lung Quantification

Example #1

Original Image (CT) — VUNO — Golden Standard

# VISUALIZATION



SVM Features (22 features)

Histogram, Gradient, Run-length, Co-occurrence matrix, Cluster analysis, Top-hat transformation

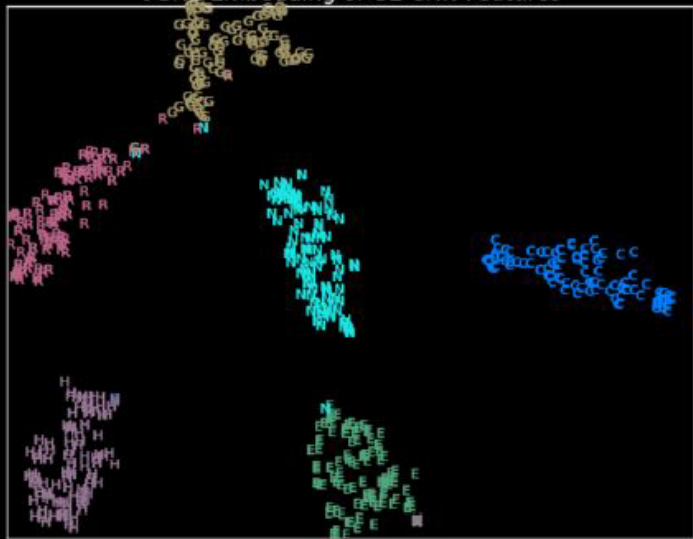t-SNE Embedding of the SVM GE Features
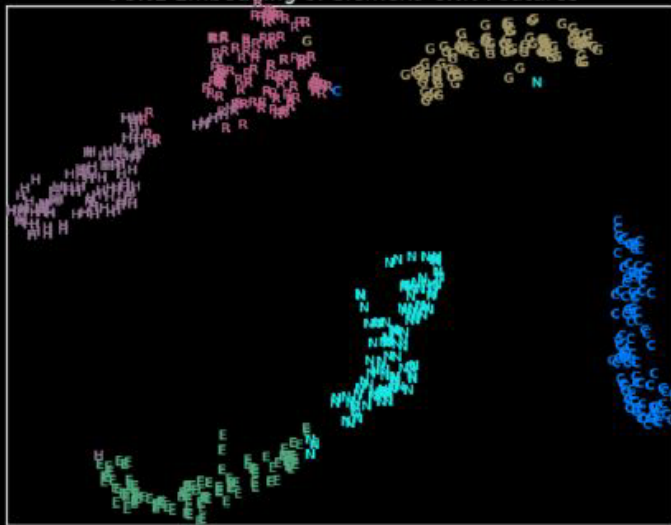
t-SNE Embedding of the SVM Siemens Features

# VISUALIZATION

Activation of top hidden layer

200 hidden nodes

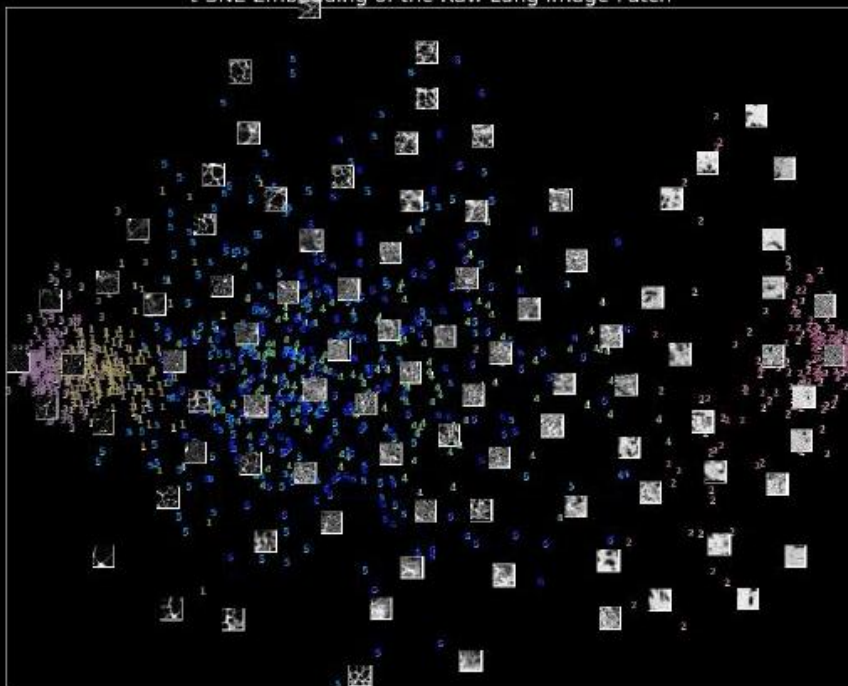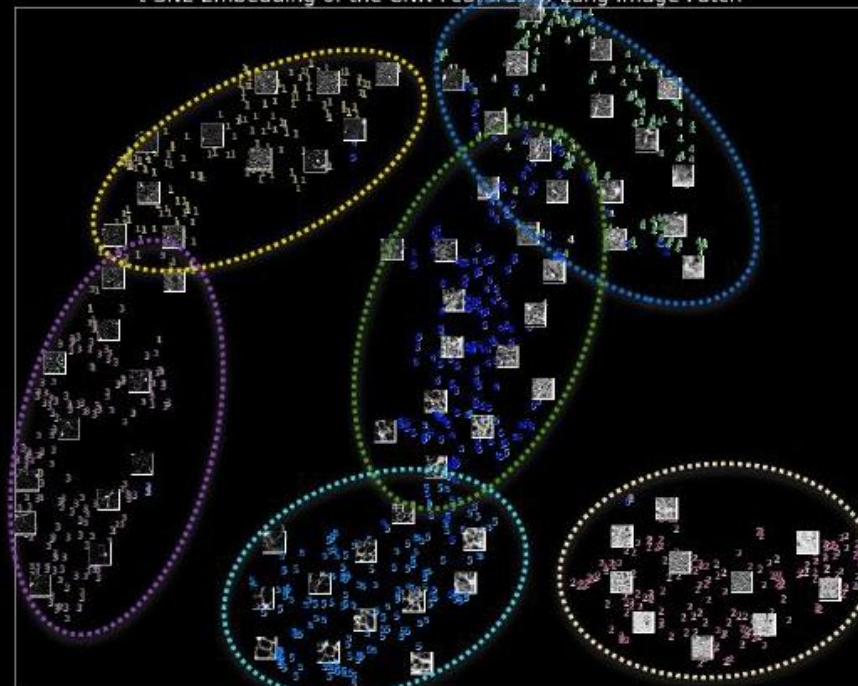# VISUALIZATION



t-SNE Embedding of the Raw Lung Image Patch

t-SNE Embedding of the CNN Features of Lung Image Patch