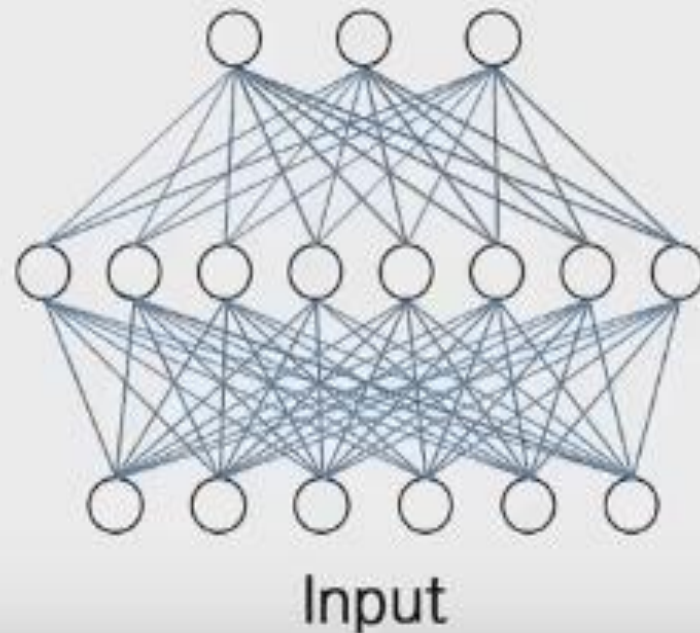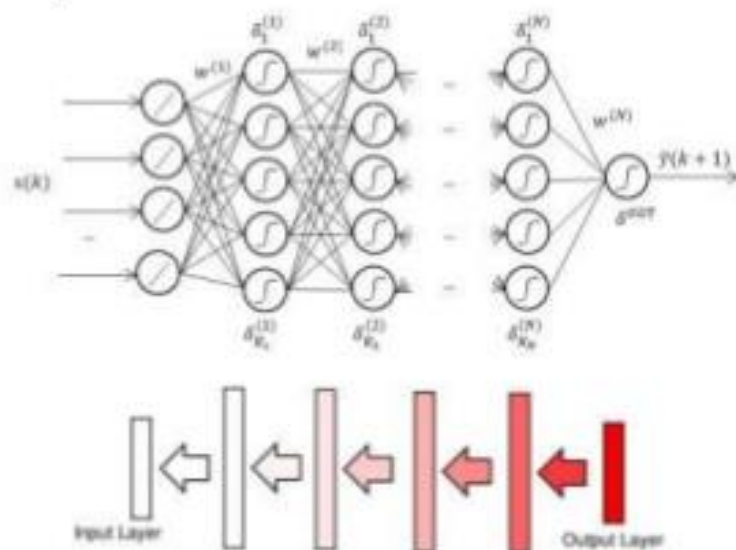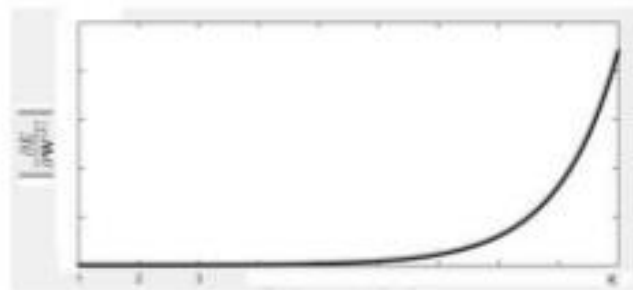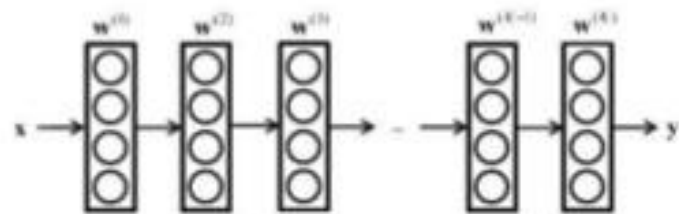# Learning Problem in DNN

Training DNN is more difficult

- Due to many parameters



Input

# Bad effect of vanishing (exploding) gradients: a problem



$$\delta_j^{(m)} = f_j^{(m-1)'} \sum_i w_{ij}^{(m)} \delta_i^{(m+1)}, \qquad => \qquad \frac{\partial E(k)}{\partial w_{ji}^{(m)}} \to 0 \ \ for \, m \to 1$$
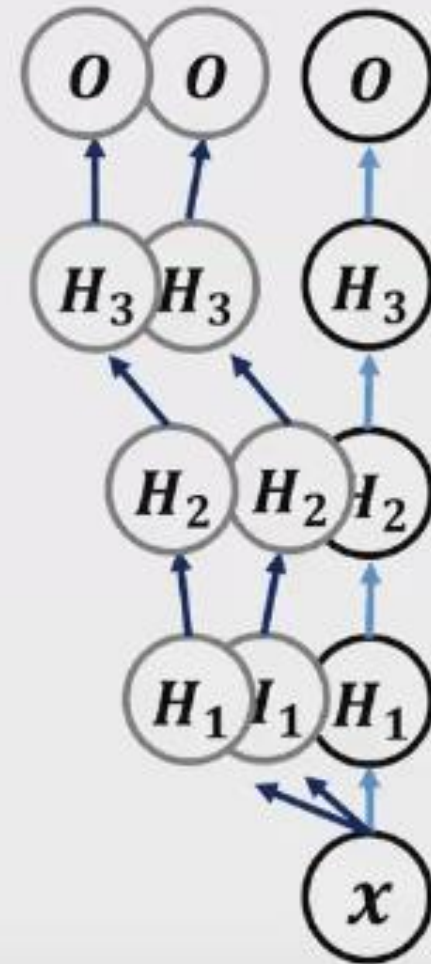
$$\frac{\partial E(k)}{\partial w_{ji}^{(m)}} = \delta_j^{(m)} z_i^{(m-1)},$$
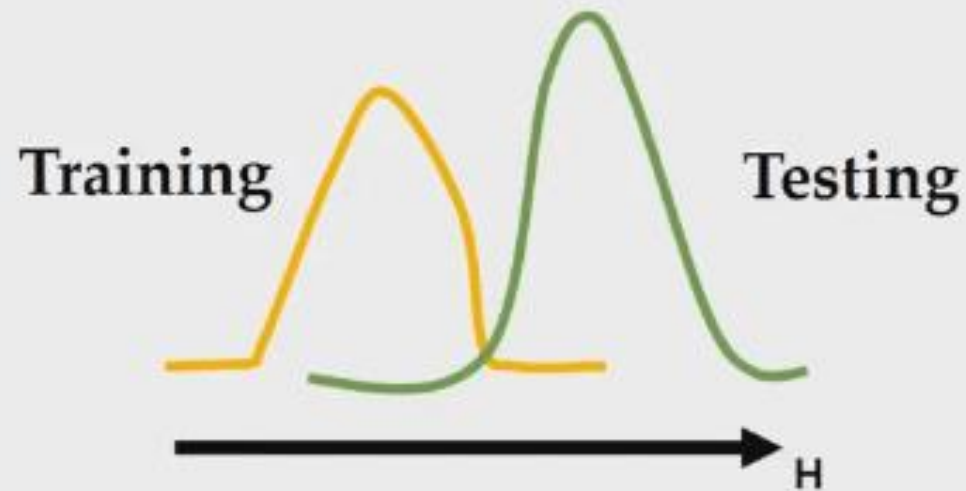
# Learning Problem in DNN

Training DNN is more difficult

- Due to many parameters
- Small change in all weights could make vary different value in upper layer
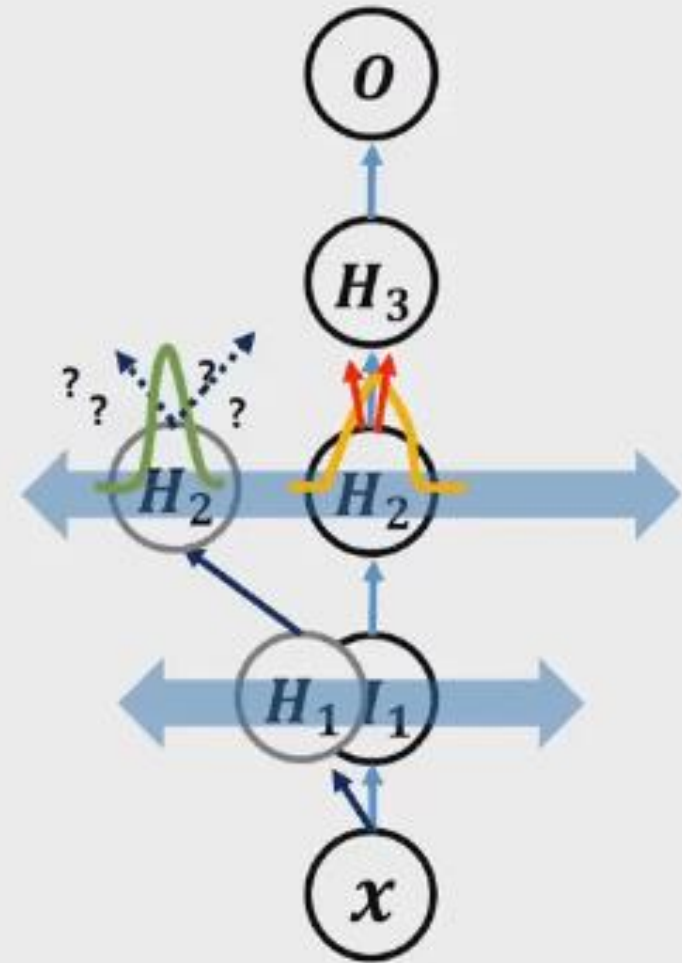
# Learning Problem in DNN

This variance is called 'Internal Covariate Shift'

Training         Testing

$H$

It is similar to the problem where dist. of training and testing are different

$O$

$H_3$

$H_2$    $H_2$

$H_1 I_1$

$x$

# Learning Problem in DNN
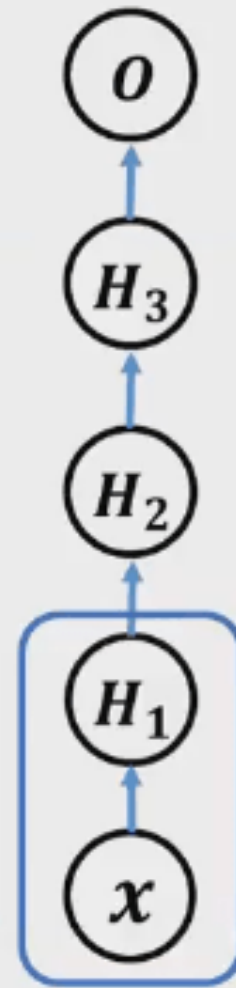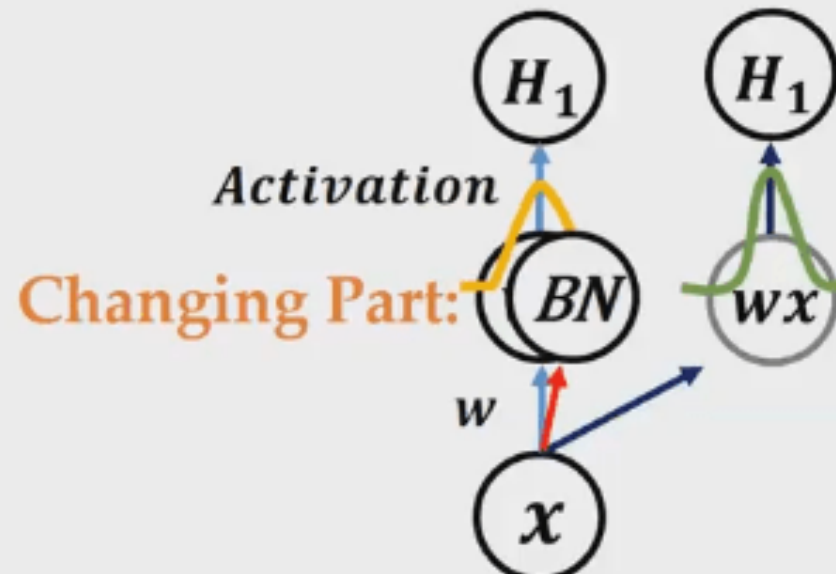
To address this Internal Covariate Shift problem, previous studies used

- careful initialization     Difficult

- small learning rate     Slow

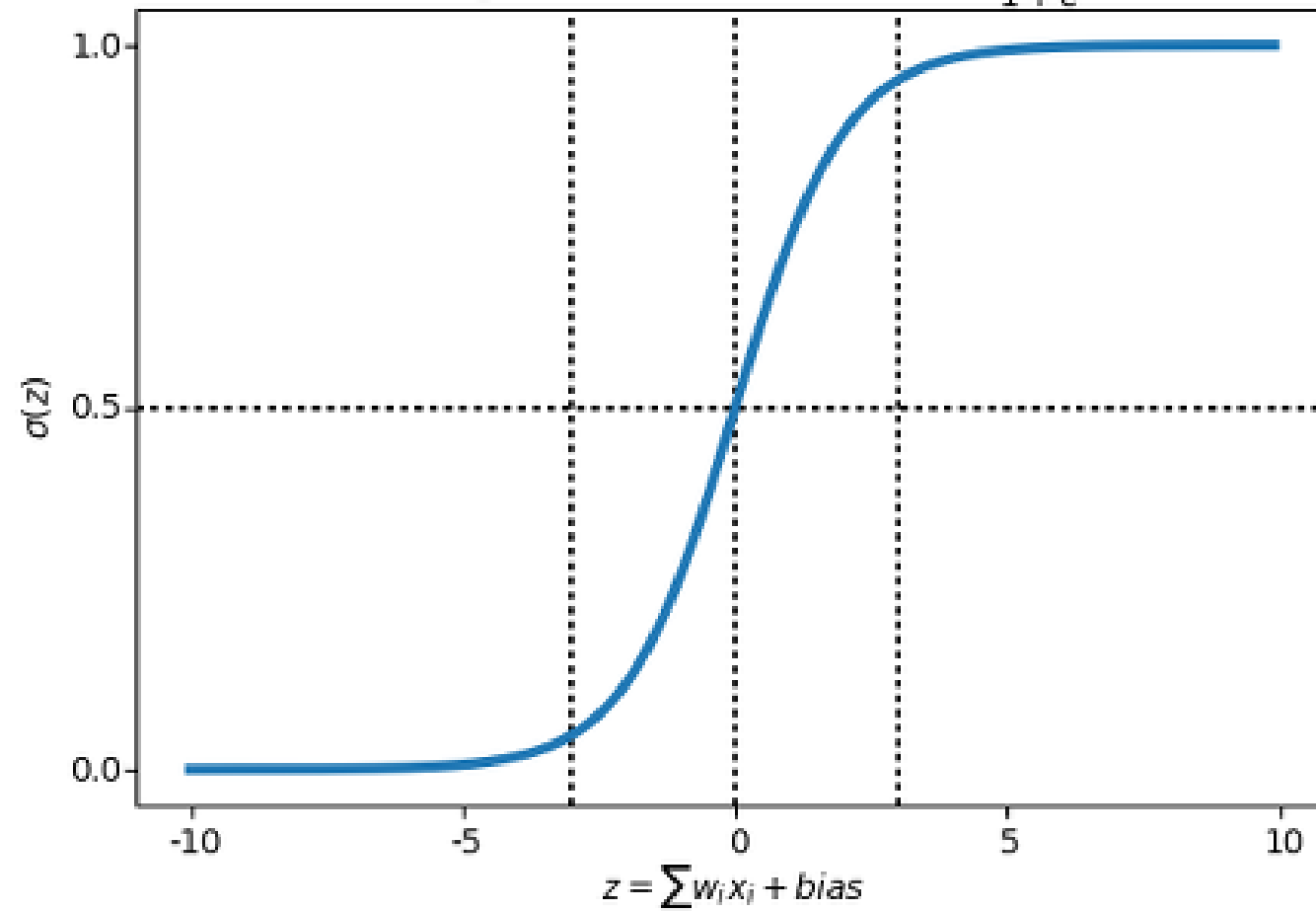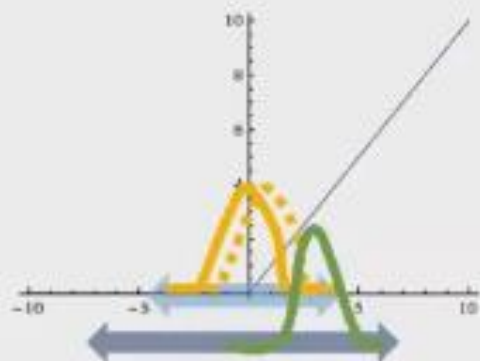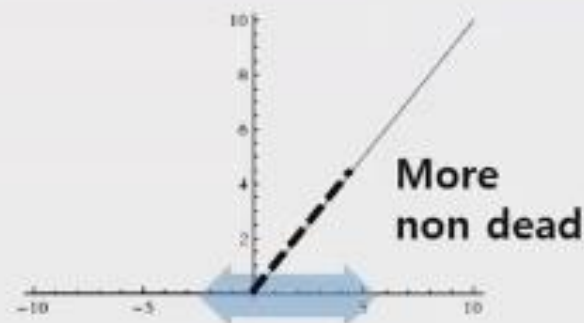# Batch Normalization

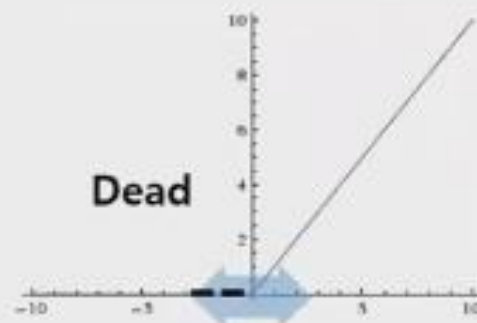Batch Norm. want to restrict change of $wx$ $(wh)$

Sigmoid Function $\sigma(z) = \frac{1}{1+e^{-z}}$

**0 mean and 1 variance is preferred
But it is not best**

Dead

More non dead

$$y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \mathrm{BN}_{\gamma,\beta}(x_i)$$

// scale and shift

**Input:** Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;
Parameters to be learned: $\gamma$, $\beta$

**Output:** $\{y_i = \text{BN}_{\gamma,\beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m}\sum_{i=1}^{m} x_i \qquad\qquad \text{// mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m}\sum_{i=1}^{m}(x_i - \mu_{\mathcal{B}})^2 \qquad\qquad \text{// mini-batch variance}$$

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad\qquad \text{// normalize}$$

$$y_i \leftarrow \gamma\widehat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \qquad\qquad \text{// scale and shift}$$

NN without BN

Output

Hidden
Layer

Input

$W_2$

$W_1$

NN without BN

Output

Hidden
Layer

BN Layer

Input

$W_2$

$\gamma, \beta$

$W_1$

**Input:** Network $N$ with trainable parameters $\Theta$;
        subset of activations $\{x^{(k)}\}_{k=1}^{K}$

**Output:** Batch-normalized network for inference, $N_{\text{BN}}^{\text{inf}}$

1: $N_{\text{BN}}^{\text{tr}} \leftarrow N$    // Training BN network
2: **for** $k = 1 \ldots K$ **do**
3:      Add transformation $y^{(k)} = \text{BN}_{\gamma^{(k)}, \beta^{(k)}}(x^{(k)})$ to $N_{\text{BN}}^{\text{tr}}$ (Alg. 1)
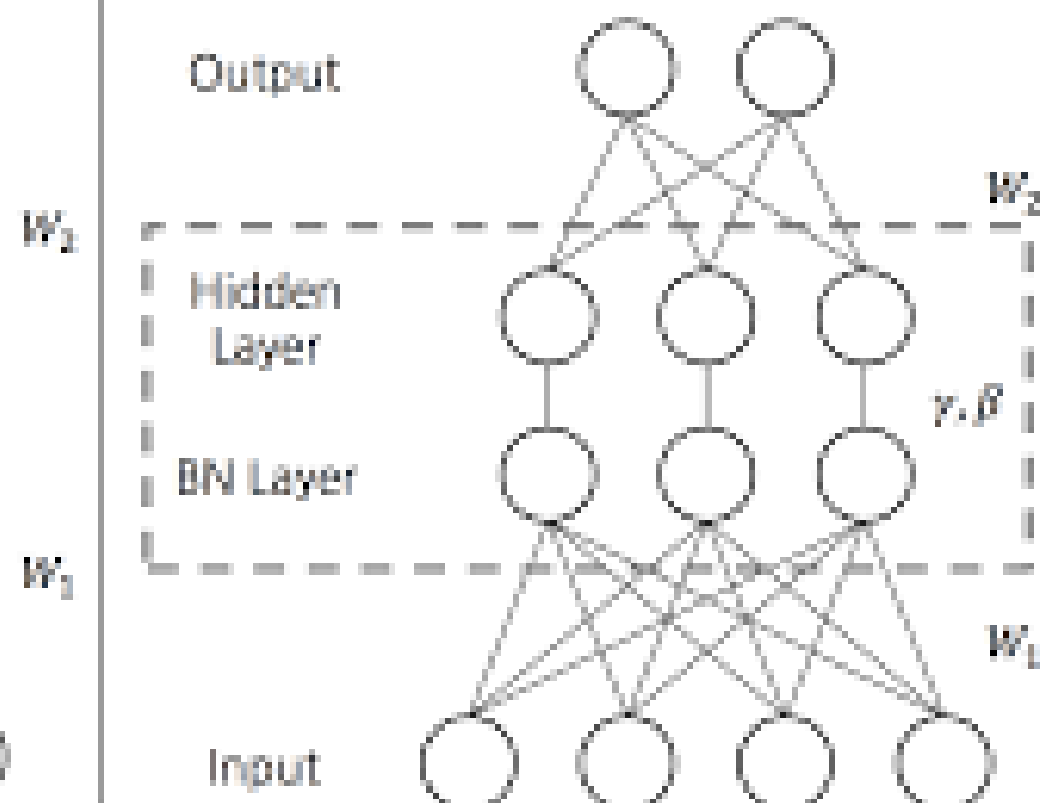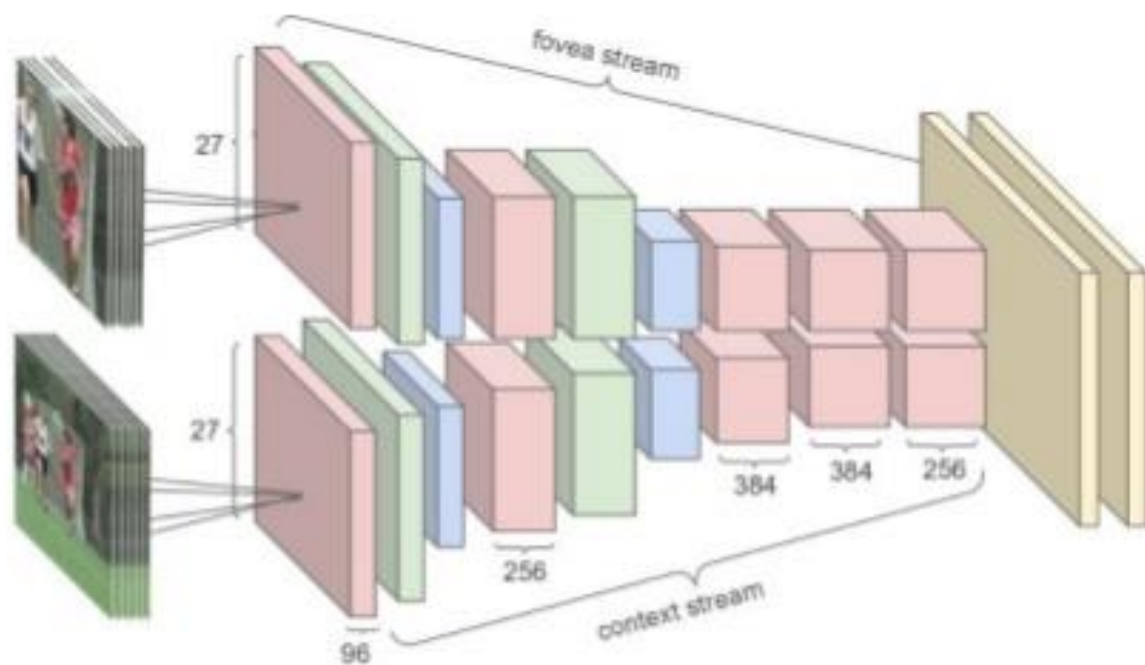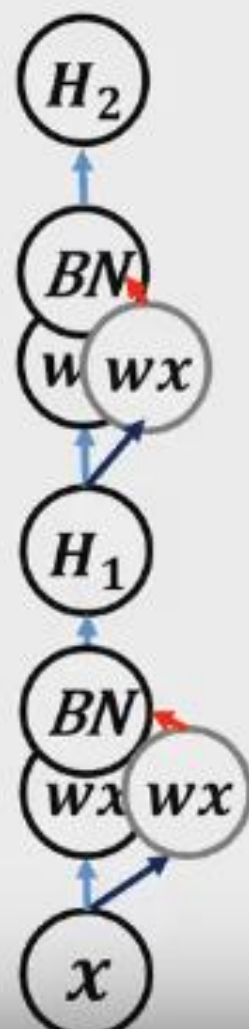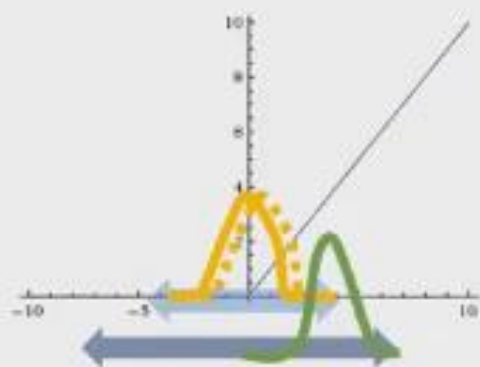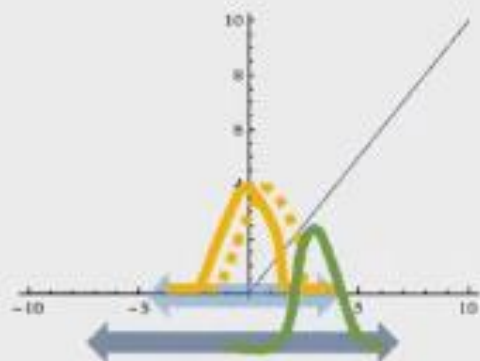4:      Modify each layer in $N_{\text{BN}}^{\text{tr}}$ with input $x^{(k)}$ to take $y^{(k)}$ instead
5: **end for**
6: Train $N_{\text{BN}}^{\text{tr}}$ to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^{K}$
7: $N_{\text{BN}}^{\text{inf}} \leftarrow N_{\text{BN}}^{\text{tr}}$    // Inference BN network with frozen
                       // parameters

8: **for** $k = 1 \ldots K$ **do**
9:      // For clarity, $x \equiv x^{(k)}, \gamma \equiv \gamma^{(k)}, \mu_{\mathcal{B}} \equiv \mu_{\mathcal{B}}^{(k)}$, etc.
10:      Process multiple training mini-batches $\mathcal{B}$, each of size $m$, and average over them:
$$\text{E}[x] \leftarrow \text{E}_{\mathcal{B}}[\mu_{\mathcal{B}}]$$
$$\text{Var}[x] \leftarrow \frac{m}{m-1}\text{E}_{\mathcal{B}}[\sigma_{\mathcal{B}}^2]$$

11:      In $N_{\text{BN}}^{\text{inf}}$, replace the transform $y = \text{BN}_{\gamma,\beta}(x)$ with
$$y = \frac{\gamma}{\sqrt{\text{Var}[x]+\epsilon}} \cdot x + \left(\beta - \frac{\gamma\,\text{E}[x]}{\sqrt{\text{Var}[x]+\epsilon}}\right)$$
12: **end for**

# Recognition: DeepVideo: Multiscale

Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R., & Fei-Fei, L. (2014, June). Large-scale video classification with convolutional neural networks. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on* (pp. 1725-1732). IEEE.

14

# Advantage of Batch Norm.

- **Regularization Effect**
  **(So, Dropout is not necessary)**

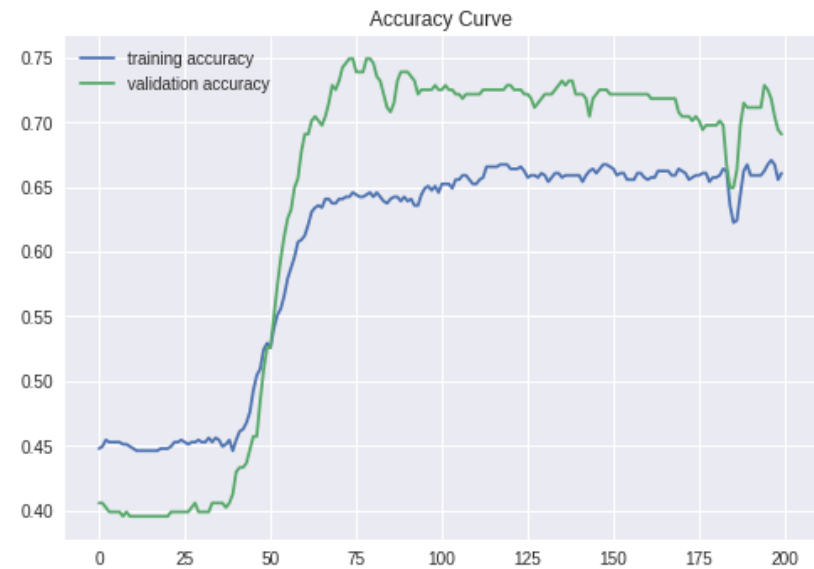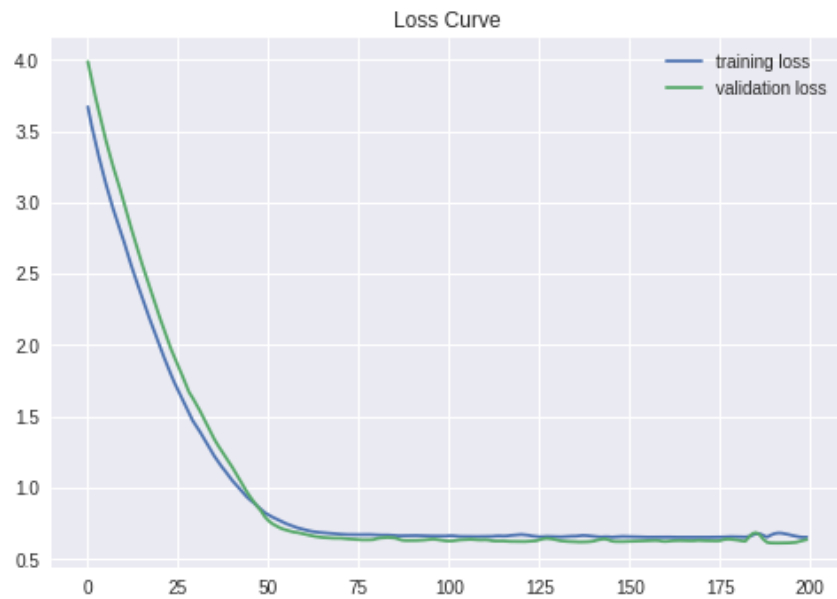**Because Mini-batch statistics $x$, E, Var is not deterministic but stochastic**

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^{m} x_i \qquad \text{// mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_{\mathcal{B}})^2 \qquad \text{// mini-batch variance}$$

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad \text{// normalize}$$

$$y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \mathrm{BN}_{\gamma,\beta}(x_i) \qquad \text{// scale and shift}$$

# 배치 정규화를 적요하지 않은 결과



# 배치 정규화를 적용한 결과