# CS231n Quiz Review

# Quiz 1. Attention & Self Attention

**11.** Attention operation is permutation invariant.

0/1 POINT
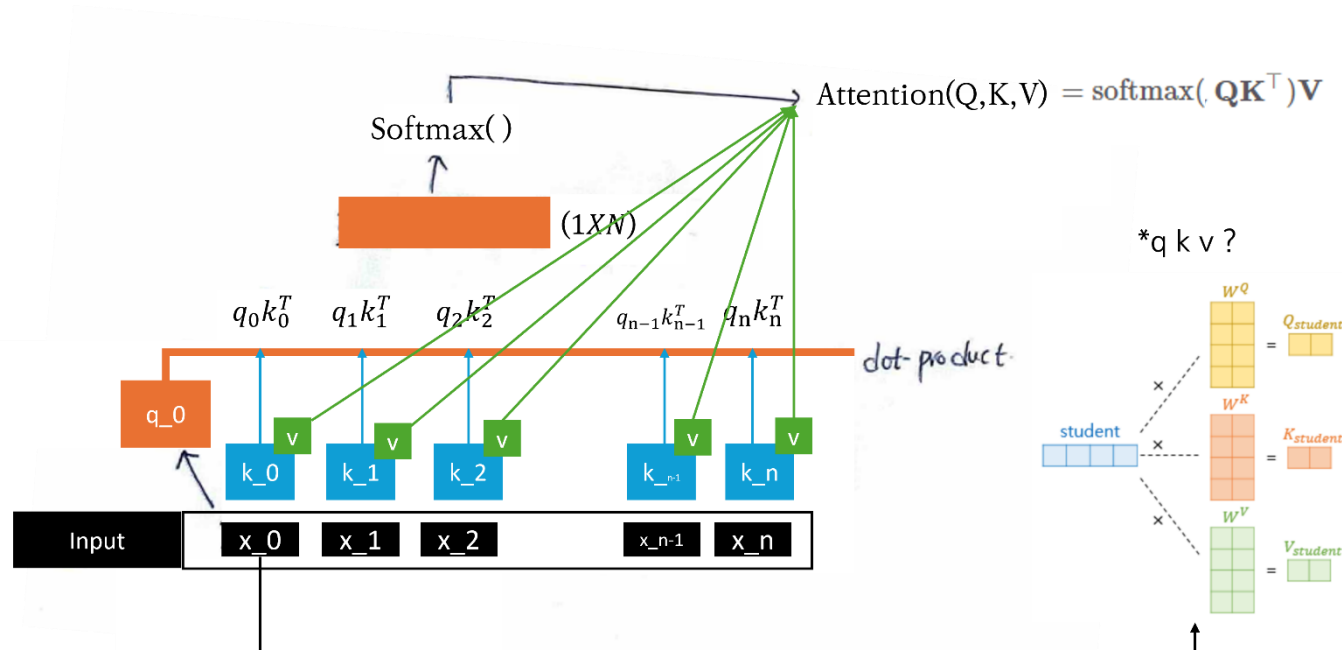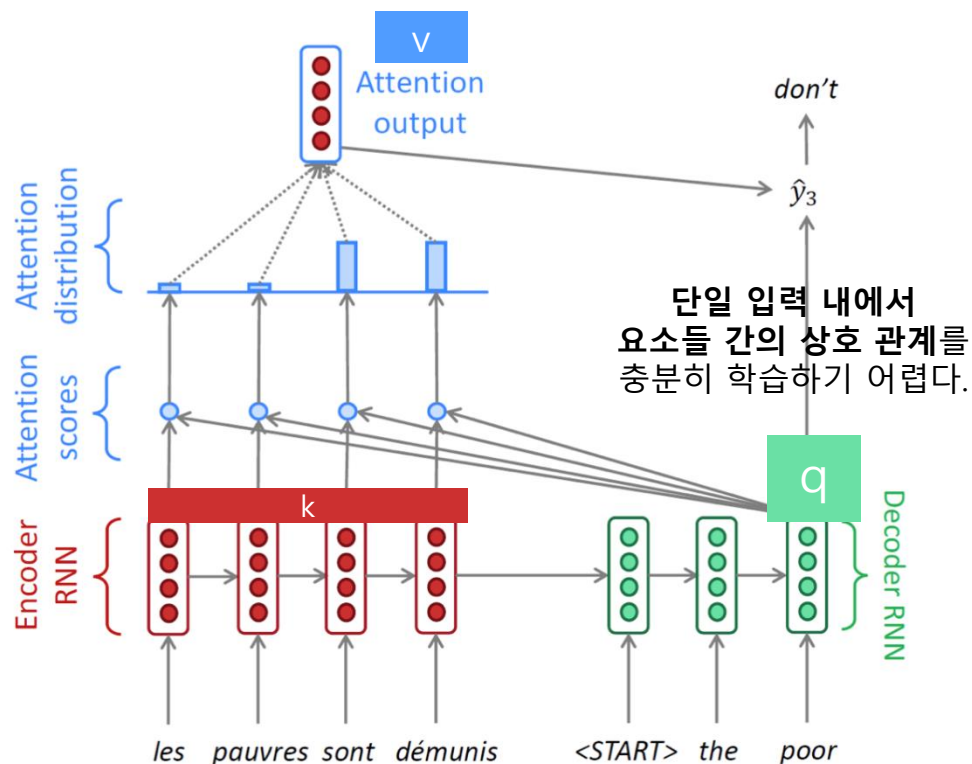
**T** True

**F** False

**Attention**
: 하나의 입력과 다른 입력과의 관계를 바탕으로 중요한 정보에 attention하는 메커니즘

**Self Attention**
: 단일 입력 내에서 모든 요소들 간의 관계를 학습



don't

$\hat{y}_3$

**단일 입력 내에서 요소들 간의 상호 관계를 충분히 학습하기 어렵다.**

Attention output

v

Attention distribution

Attention scores

k

Encoder RNN

q

Decoder RNN

les  pauvres  sont  démunis  <START>  the  poor

Softmax( )

$(1 X N)$

$q_0 k_0^T$  $q_1 k_1^T$  $q_2 k_2^T$  $q_{n-1} k_{n-1}^T$  $q_n k_n^T$

dot-product

Attention(Q,K,V) = softmax( $\mathbf{Q K}^\top$ ) $\mathbf{V}$

*q k v ?

$W^Q$  $Q_{student}$

$W^K$  $K_{student}$

$W^V$  $V_{student}$

q_0  v  v  v  v  v

k_0  k_1  k_2  k_n-1  k_n

student

Input  x_0  x_1  x_2  x_n-1  x_n

# Quiz 2. ResNet

**58.** Residual networks (ResNets) always help mitigate the degradation problem in deep neural networks by allowing training of substantially deeper models.
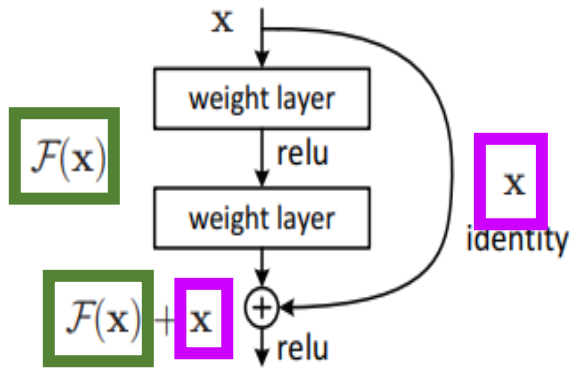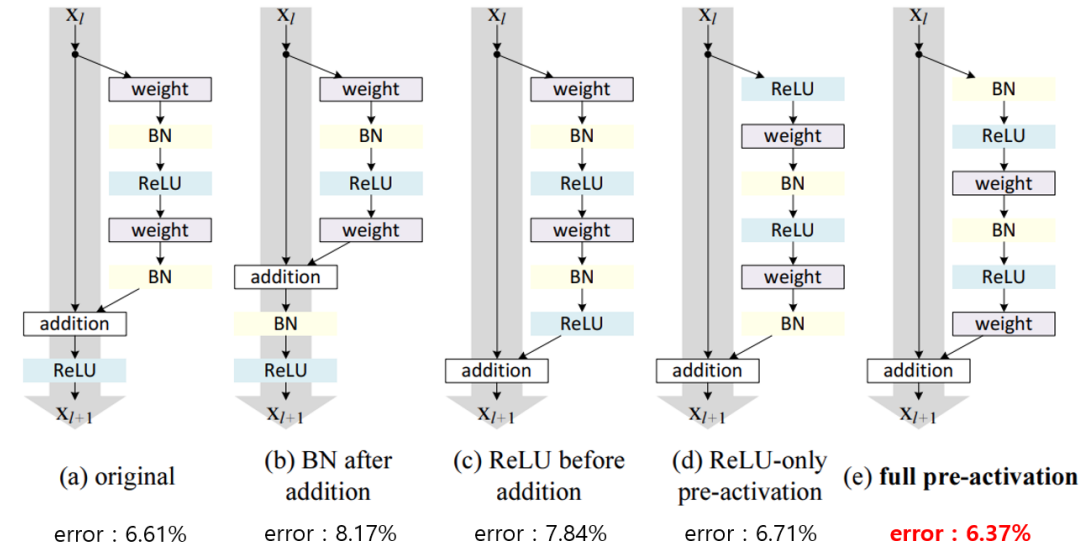
0/1 POINT

**T** True

**F** False

## 1. Residual Block
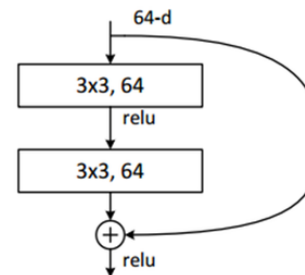


Figure 2. Residual learning: a building block.

$$F(x) + x$$

입력을 출력에 더해줌으로써 기울기가 직접 전달될 수 있게 하여
Gradient Vanishing 문제 완화
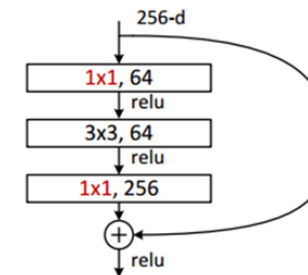
## 2. Pre-Activation

Batch Normalization과
ReLU 위치에 따른 성능 평가 지표



| (a) original | (b) BN after addition | (c) ReLU before addition | (d) ReLU-only pre-activation | (e) **full pre-activation** |
| error : 6.61% | error : 8.17% | error : 7.84% | error : 6.71% | **error : 6.37%** |

## 3. Bottleneck



dimension을 줄이기 위한 목적.
-> 파라미터 감소

ResNet-50 이상에서 사용

# Quiz 3. Batch Normalization

✗ **7.** Batch normalization can have an implicit regularizing effect, especially with smaller minibatches.

0/1 POINT

**T** True

**F** False

# Layer Normalization

Lei Ba, Jimmy, Jamie Ryan Kiros, and Geoffrey E. Hinton. "Layer normalization." *ArXiv e-prints* (2016): arXiv-1607.

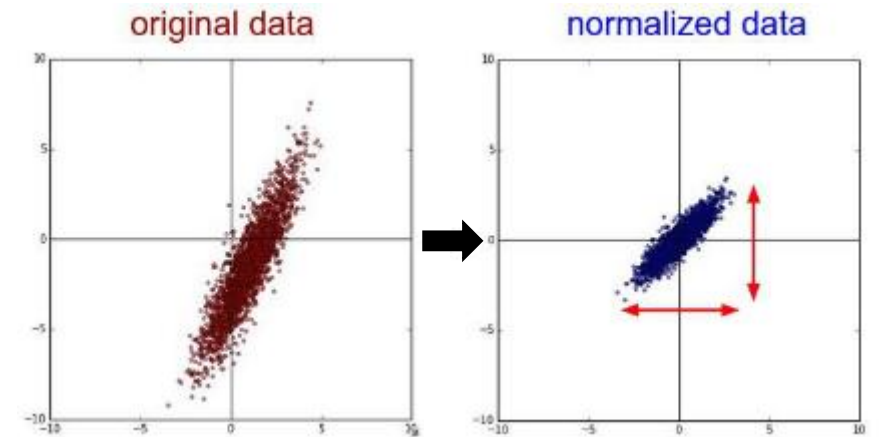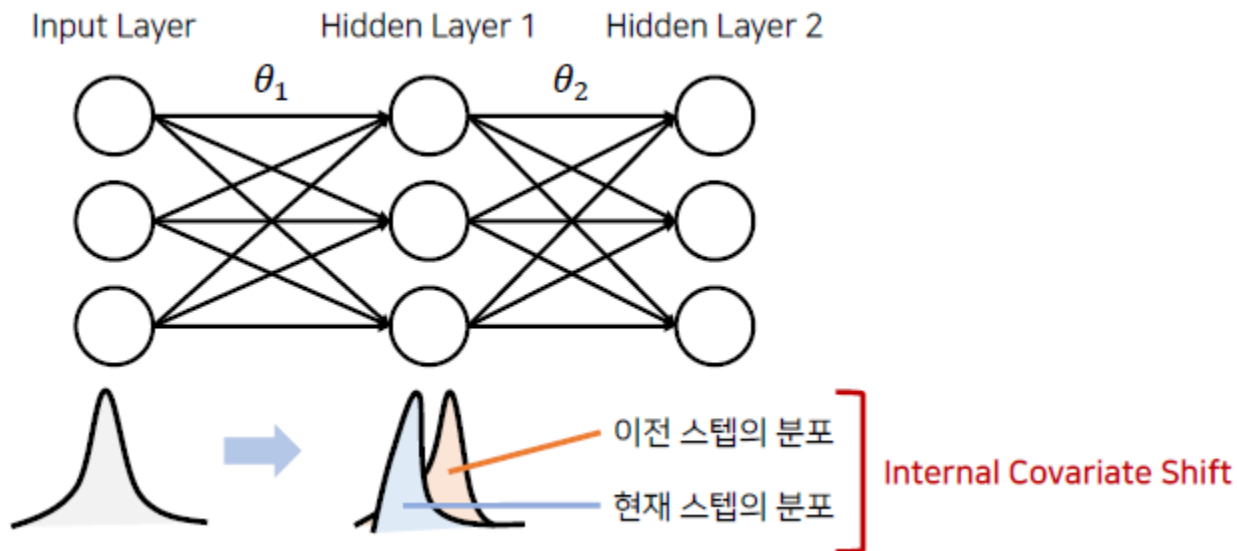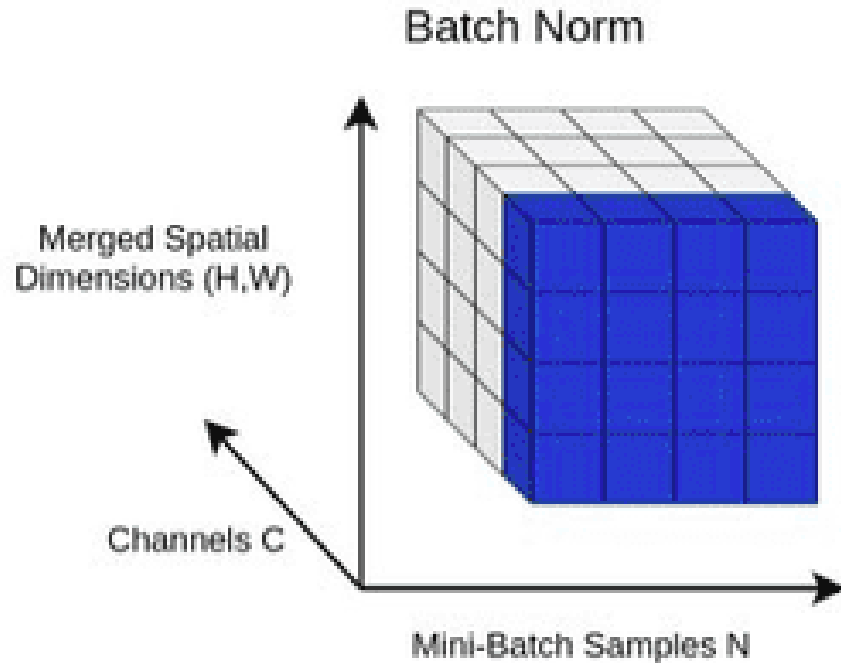# Introduction

긴 학습 시간 소요

**Normalization**

**Covariate Shift**



Input Layer    Hidden Layer 1    Hidden Layer 2

$\theta_1$     $\theta_2$

이전 스텝의 분포

현재 스텝의 분포

Internal Covariate Shift



original data     normalized data

# Batch Normalization

### Batch Norm



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

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

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

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

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad \text{// normalize}$$
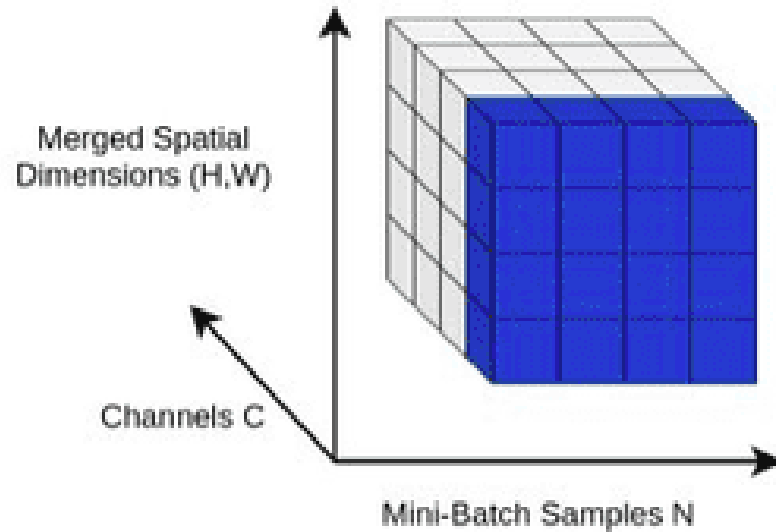
$$y_i \leftarrow \gamma\widehat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \qquad \text{// scale and shift}$$

**Algorithm 1:** Batch Normalizing Transform, applied to activation $x$ over a mini-batch.
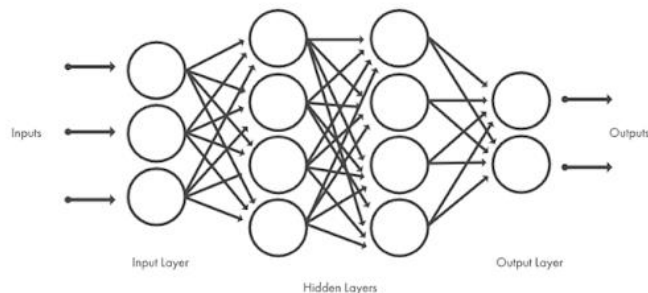
1. Sequence Length Variability in Recurrent Neural Networks (RNNs)

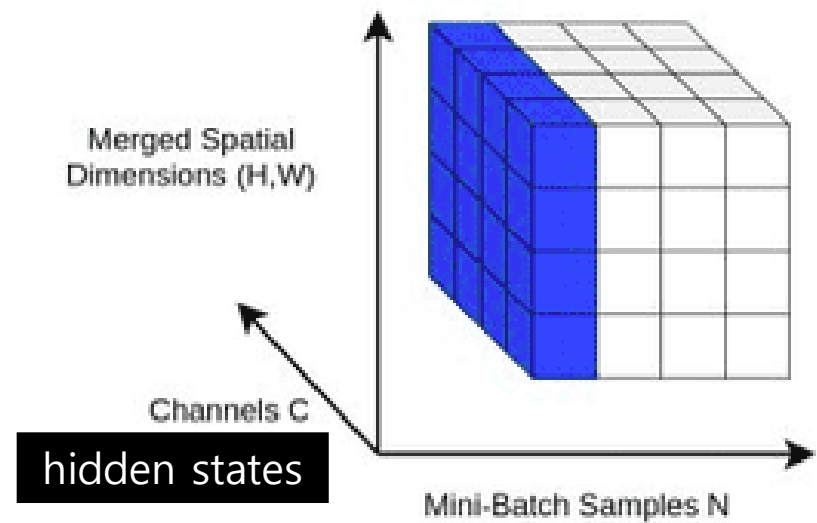2. Limitations in Online Learning and Large Distributed Models

Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. ICML, 2015.

# Layer Normalization

## Batch Norm

Merged Spatial
Dimensions (H,W)

Channels C

Mini-Batch Samples N

$$\bar{a}_i^l = \frac{g_i^l}{\sigma_i^l} \left( a_i^l - \mu_i^l \right)$$

Inputs

Outputs

Input Layer

Hidden Layers

Output Layer

Feed-forward neural network

## Layer Norm

Merged Spatial
Dimensions (H,W)

Channels C

hidden states

Mini-Batch Samples N

$$\mathbf{h}^t = f \left[ \frac{\mathbf{g}}{\sigma^t} \odot \left( \mathbf{a}^t - \mu^t \right) + \mathbf{b} \right]$$

RNN

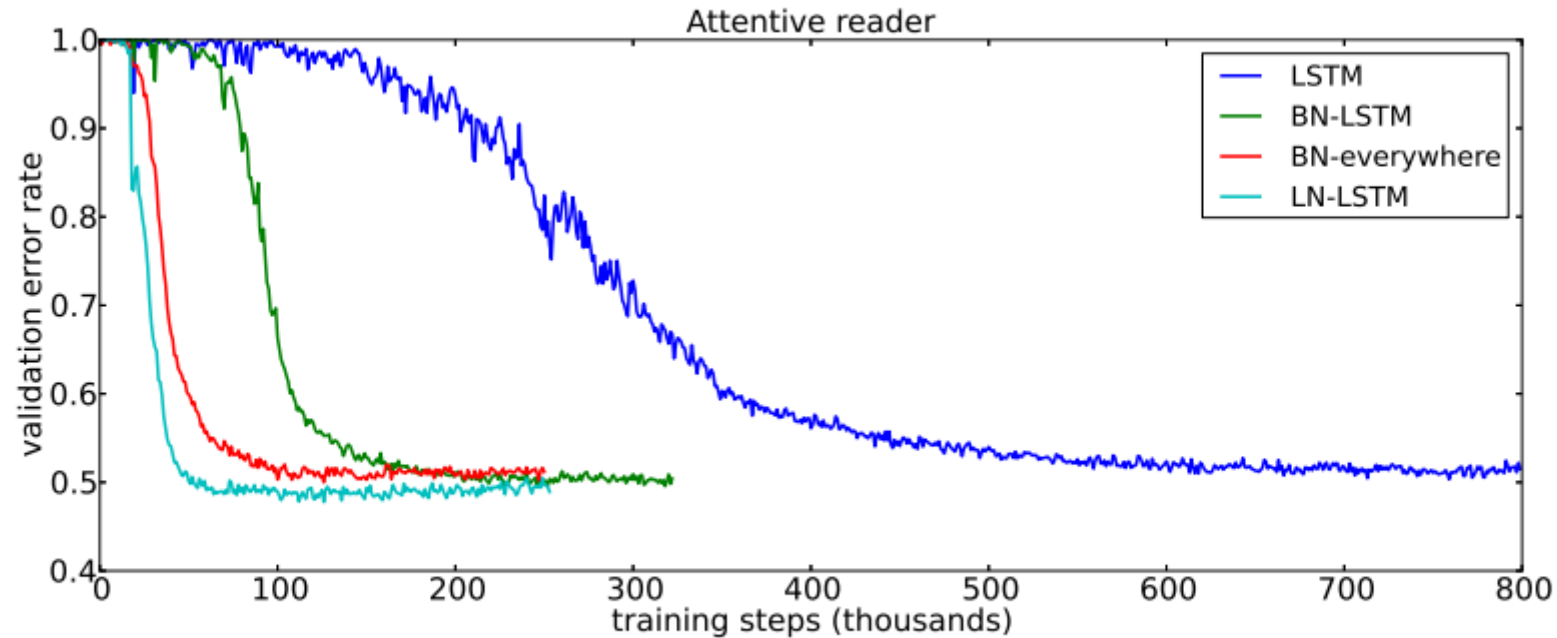**6.2 Teaching machines to read and comprehend**



Figure 2: Validation curves for the attentive reader model. BN results are taken from [Cooijmans et al., 2016].

**6.7 Convolutional Networks**

# Quiz 3. Normalization

7. Batch normalization can have an implicit regularizing effect, especially with smaller minibatches.
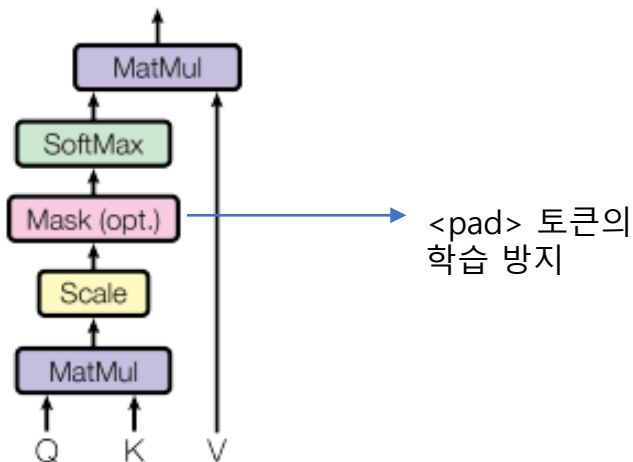
0/1 POINT

True

False

"regularizing effect"
- batch마다 달라지는 통곗값
- shift 파라미터로 인한 정규화 값의 위치 조정

\* Scaled Dot-Product Attention

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

**Scaled Dot-Product Attention**



<pad> 토큰의
학습 방지

Query, Key 벡터의 차원을 $d_k$,
$Q, K$의 요소들이 평균 0, 분산 1인 독립적인 랜덤변수라고 가정

$$Q \cdot K = \sum_{i=1}^{d_k} g_i k_i$$

\* 기댓값

$$E[g_i k_i] = E[g_i] E[k_i] = 0$$

$$E[Q \cdot K] = \sum_{i=1}^{d_k} E[g_i k_i] = 0$$

\* 분산

$$Var(g_i k_i) = Var(g_i) \cdot Var(k_i) = 1$$

$$Var(Q \cdot K) = Var(\sum_{i=1}^{d_k} g_i k_i) = \sum_{i=1}^{d_k} Var(g_i k_i) = d_k$$

분산 $= d_k$ → 차원 $d_k$가 커질수록 분산이
비례하여 증가함

→ 내적값의 분산을 1로 맞추기 위해 $\sqrt{d_k}$로 scaling

$$\frac{Q \cdot K}{\sqrt{d_k}} = \frac{\sum_{i=1}^{d_k} g_i k_i}{\sqrt{d_k}}$$

분산의 성질에 의하여

$$Var(\frac{Q \cdot K}{\sqrt{d_k}}) = \frac{Var(Q \cdot K)}{d_k} = \frac{d_k}{d_k} = 1$$

Scaling을 수행하면 내적 결과의 분산은 1이 됨

차원이 증가하더라도 내적값의 크기가 적절한 범위로 유지되므로,
Attention weight가 특정 위치로 치우치는 현상을 방지할 수 있다.