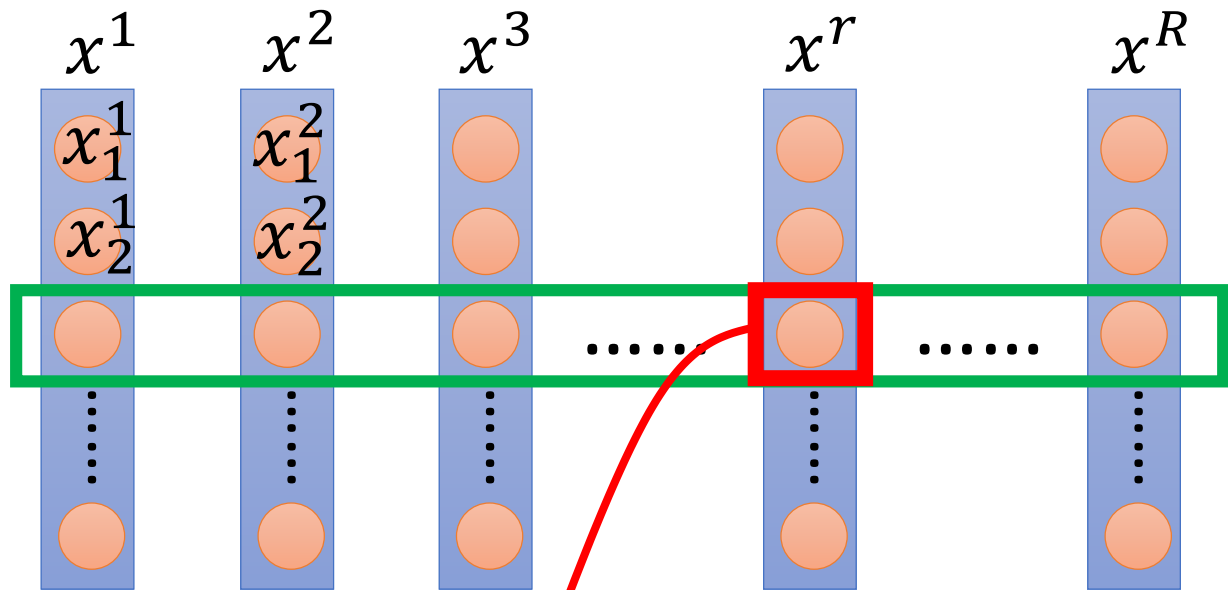


Batch Normalization

Feature Scaling



For each dimension i :
mean: m_i
standard deviation: σ_i

$$x_i^r \leftarrow \frac{x_i^r - m_i}{\sigma_i}$$

The means of all dimensions are 0, and the variances are all 1

In general, gradient descent converges much faster with feature scaling than without it.

How about Hidden Layer?

Feature Scaling



x^1



Layer 1



a^1



Feature Scaling ?



a^1



Layer 2



a^2

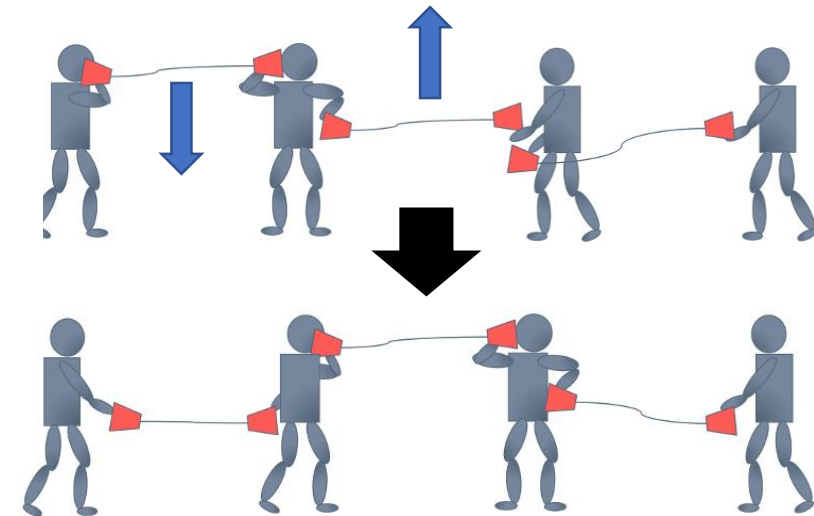
.....

Feature Scaling ?



a^2

.....



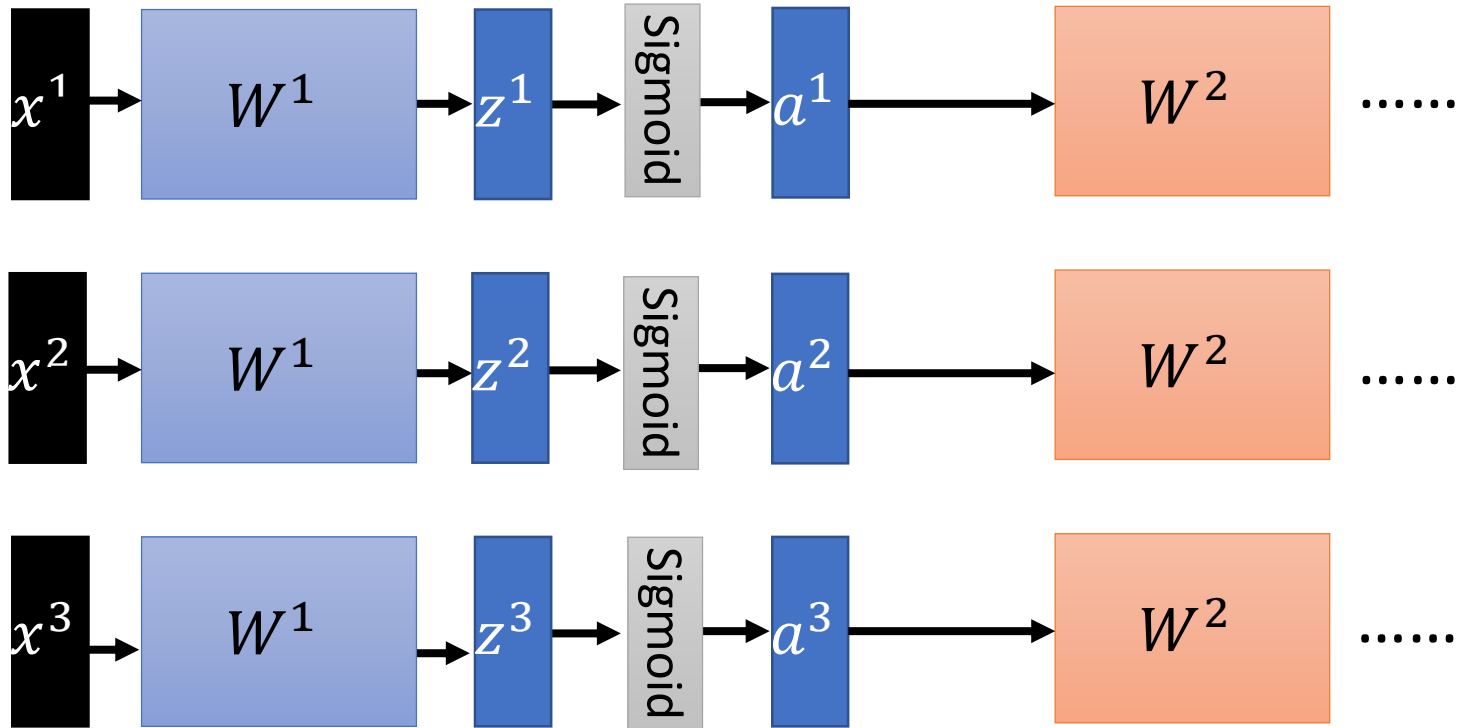
Internal Covariate Shift

Difficulty: their statistics change during the training ...

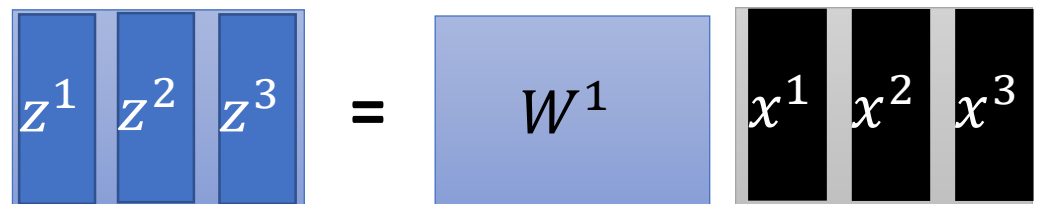
➡ **Batch normalization**

Smaller learning rate can be helpful, but the training would be slower.

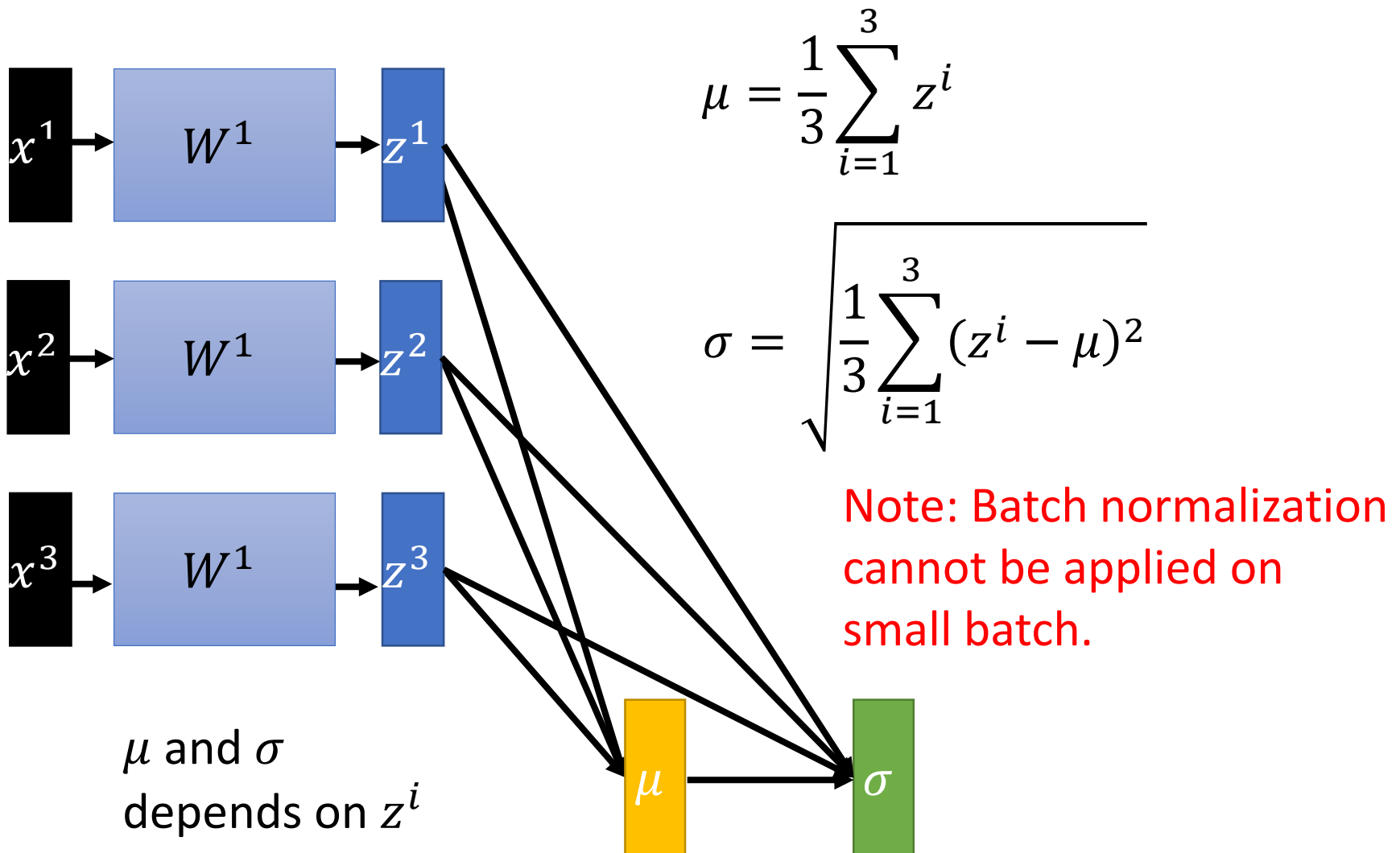
Batch



Batch

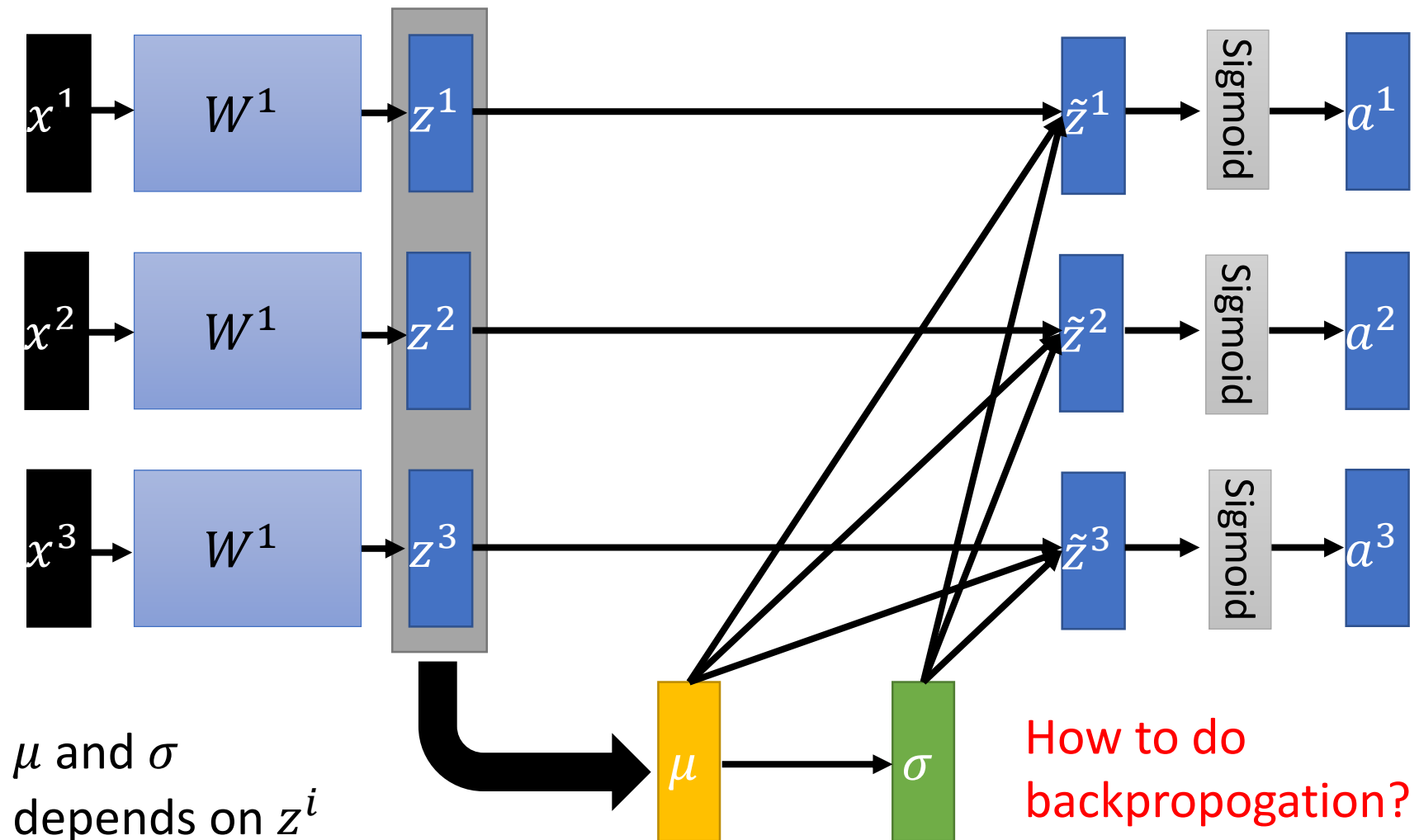


Batch normalization



Batch normalization

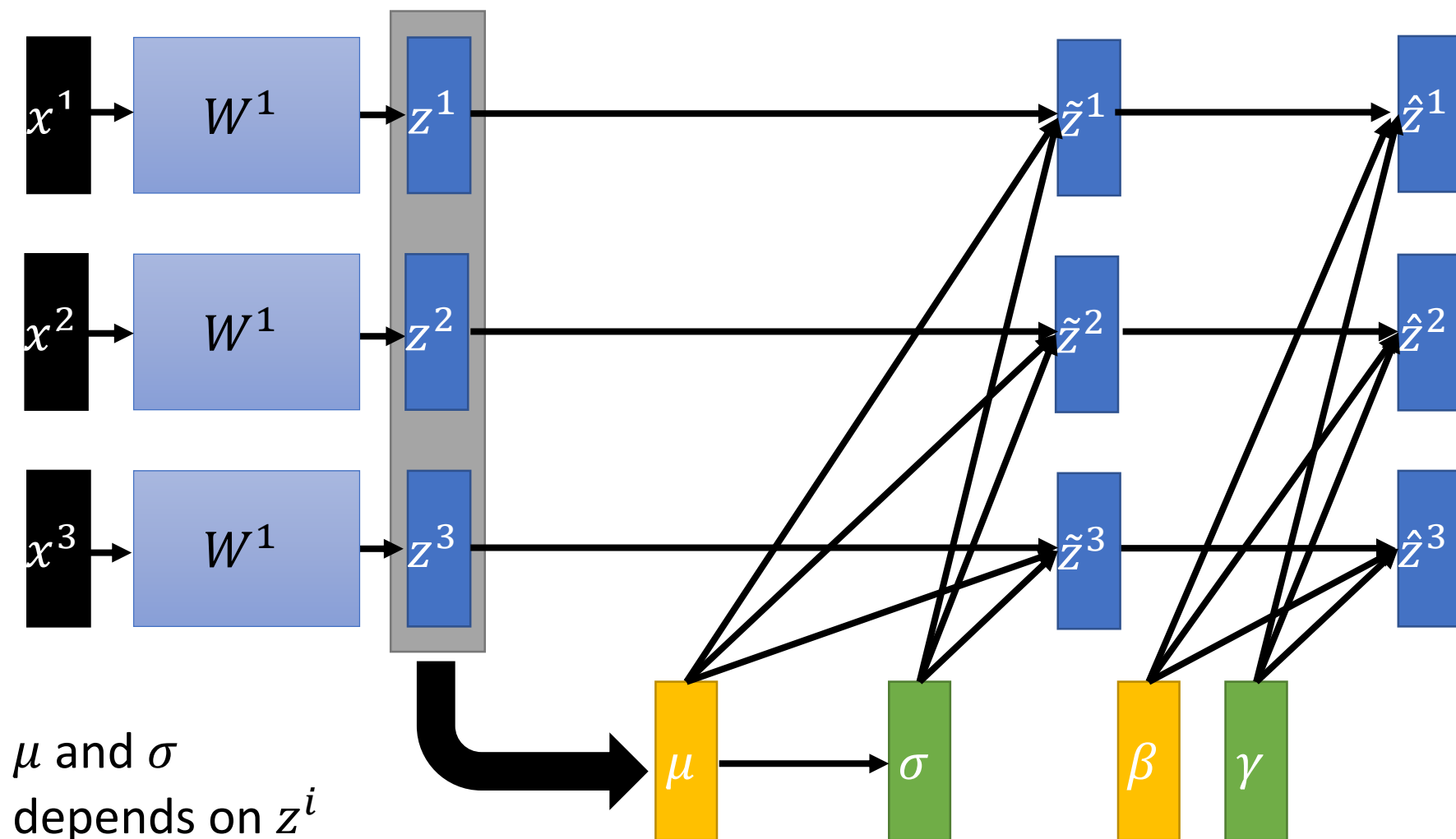
$$\tilde{z}^i = \frac{z^i - \mu}{\sigma}$$



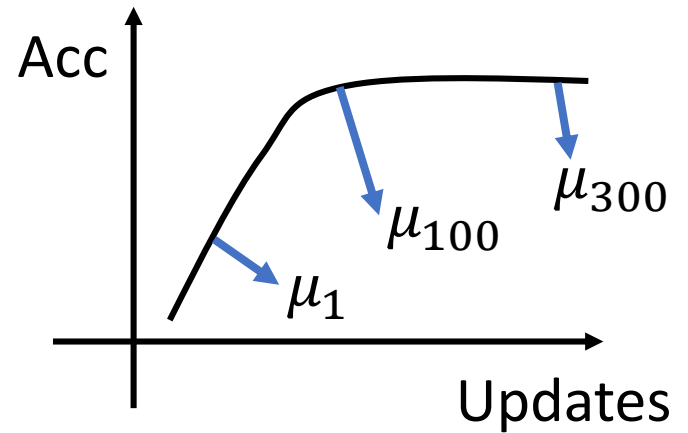
Batch normalization

$$\tilde{z}^i = \frac{z^i - \mu}{\sigma}$$

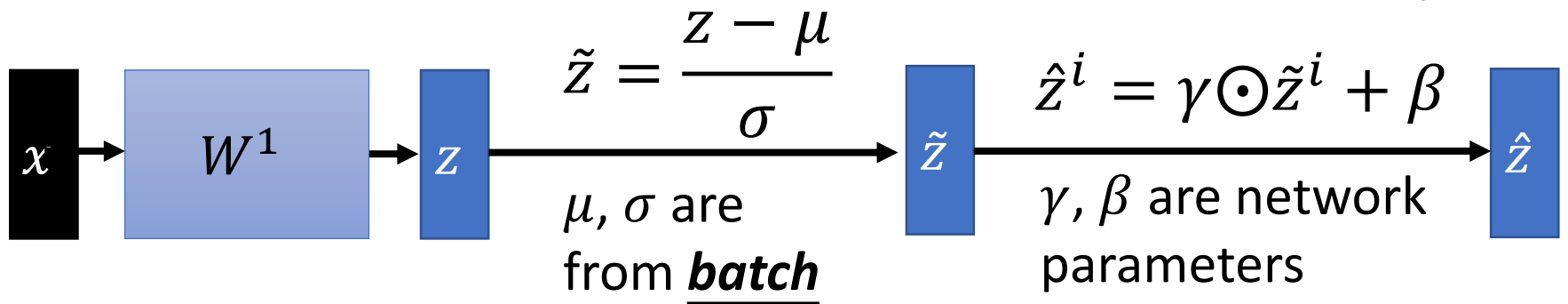
$$\hat{z}^i = \gamma \odot \tilde{z}^i + \beta$$



Batch normalization



- At testing stage:



We do not have batch at testing stage.

Ideal solution:

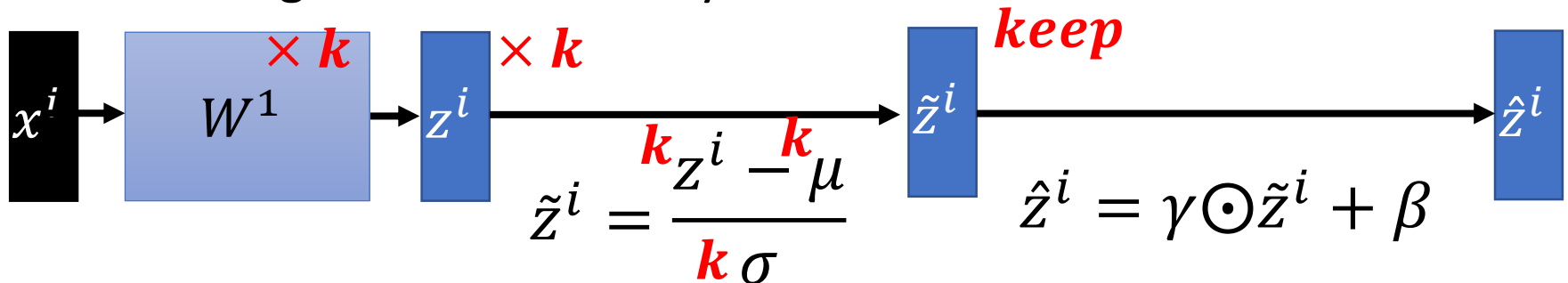
Computing μ and σ using the whole training dataset.

Practical solution:

Computing the moving average of μ and σ of the batches during training.

Batch normalization - Benefit

- BN reduces training times, and make very deep net trainable.
 - Because of less Covariate Shift, we can use larger learning rates.
 - Less exploding/vanishing gradients
 - Especially effective for sigmoid, tanh, etc.
- Learning is less affected by initialization.



- BN reduces the demand for regularization.

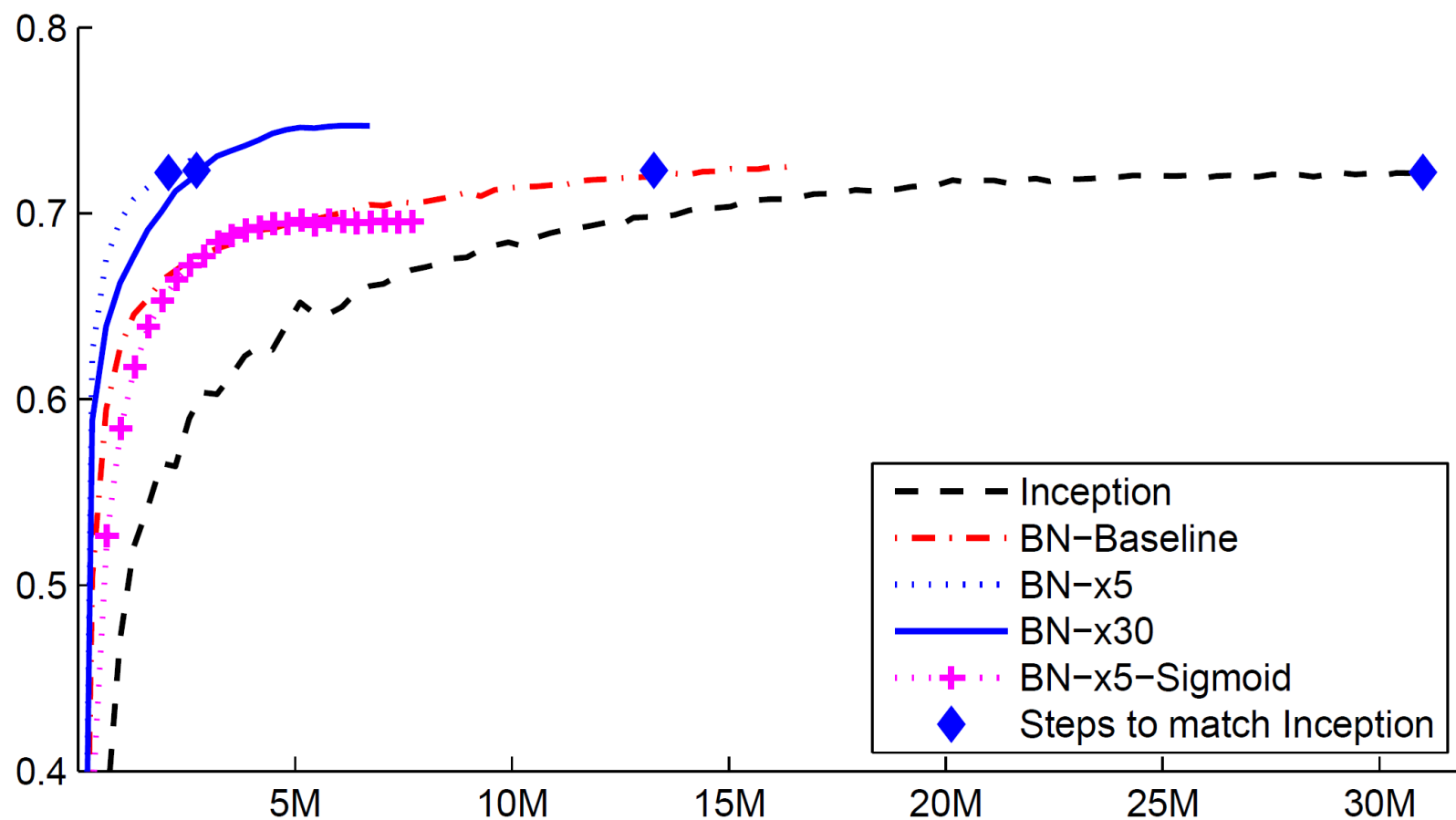


Figure 2: *Single crop validation accuracy of Inception and its batch-normalized variants, vs. the number of training steps.*