



CVPR 2021 Tutorial



Normalization Techniques in Deep Learning: Methods, Analyses and Applications



Lei Huang

Beihang University, Beijing, China



北京航空航天大学人工智能研究院
Institute of Artificial Intelligence, Beihang University

2021-06-19



软件开发环境国家重点实验室
State Key Laboratory of Software Development Environment

Deep Neural Networks

Computer vision

Image classification

Segmentation

3D vision

Object detection

Object tracking

Natural language processing

Sentiment Analysis

Knowledge Graph

Question answering

Machine translation

Language model

Speech and audio processing

Speech recognition

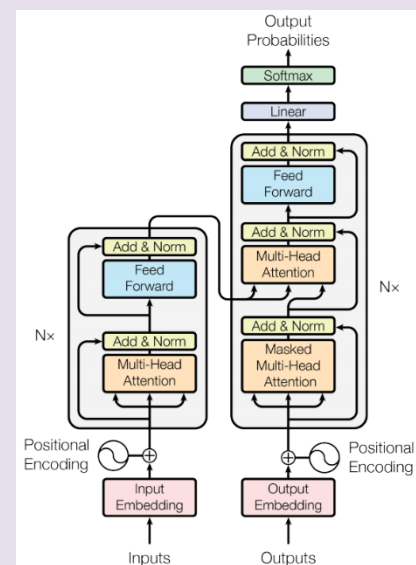
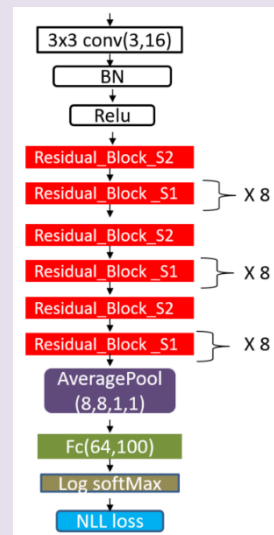
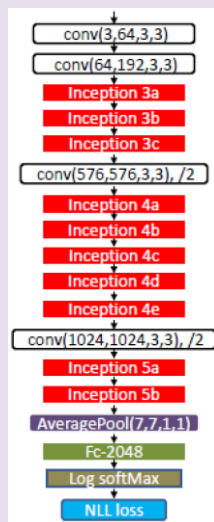
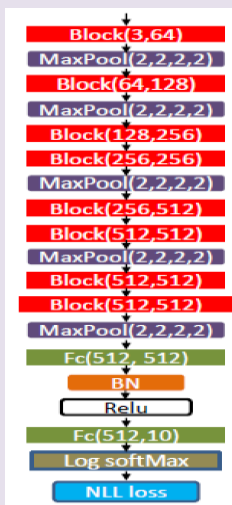
Speech coding

Speech Synthesis

Speaker identification

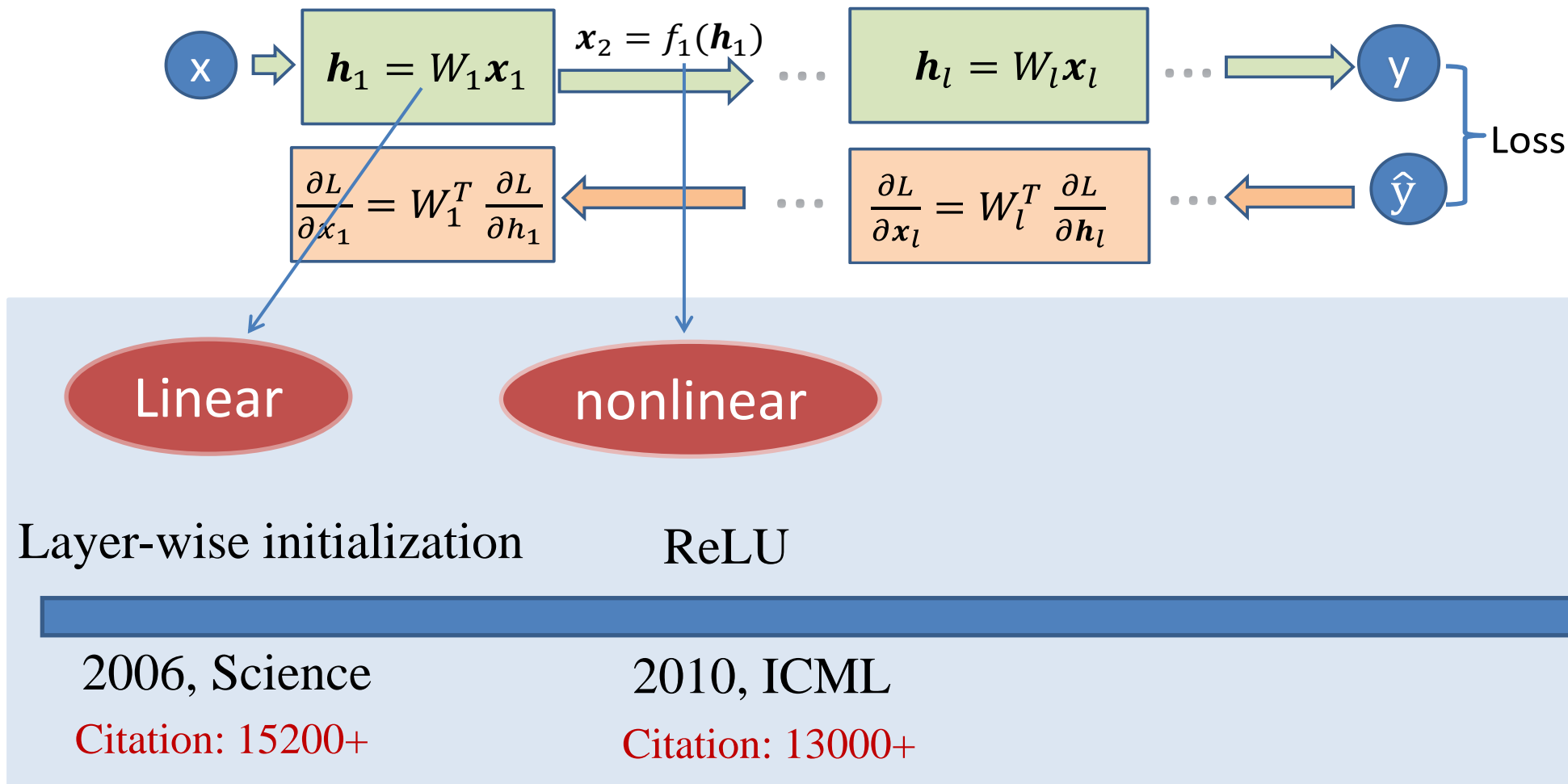
Speech Enhancement

Deep Neural Networks (DNNs)



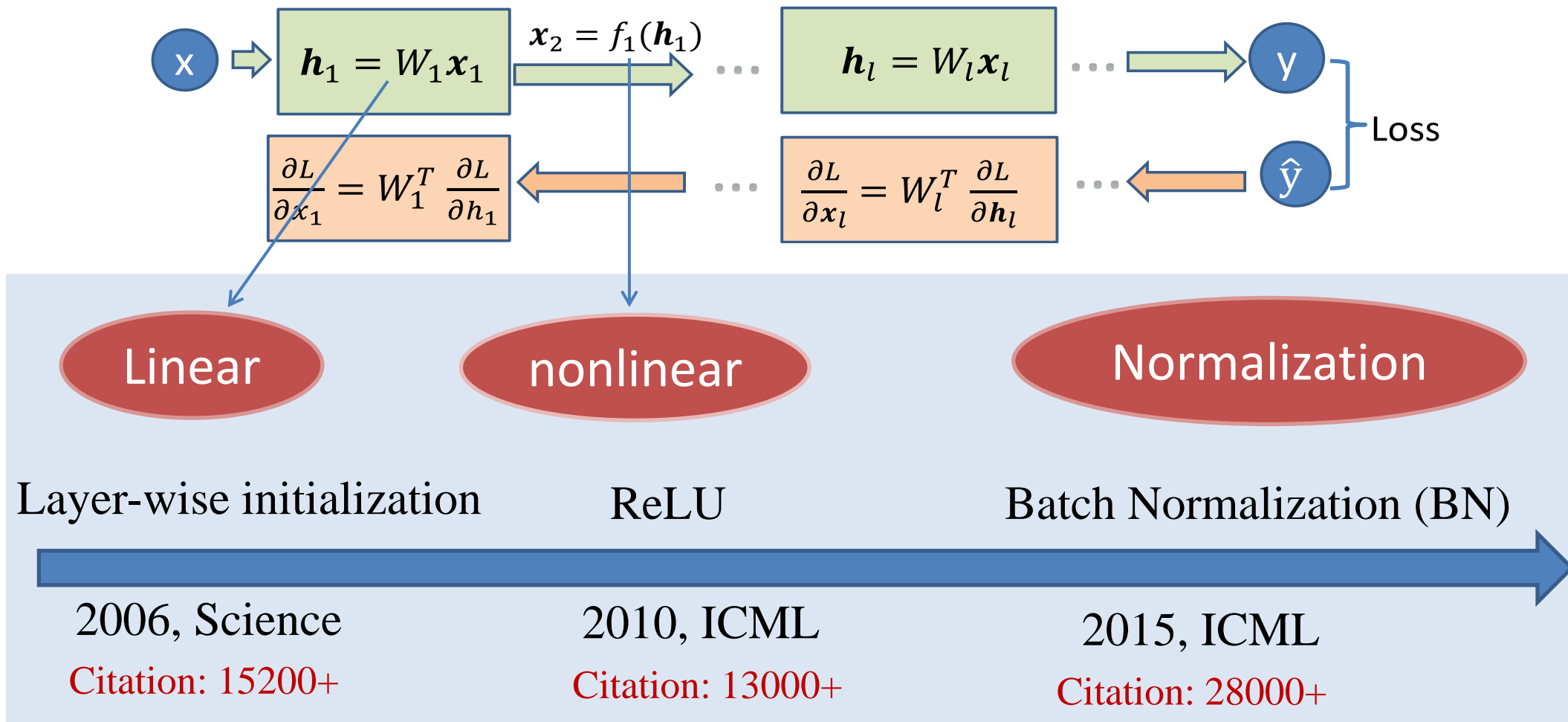
Milestones of DNNs

- The gradient explosion or vanishing problem



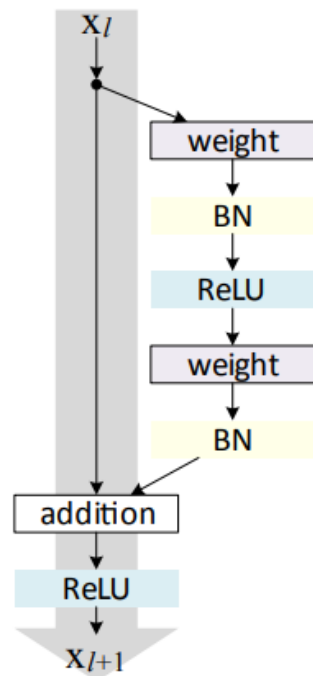
Milestones of DNNs

- The gradient explosion or vanishing problem



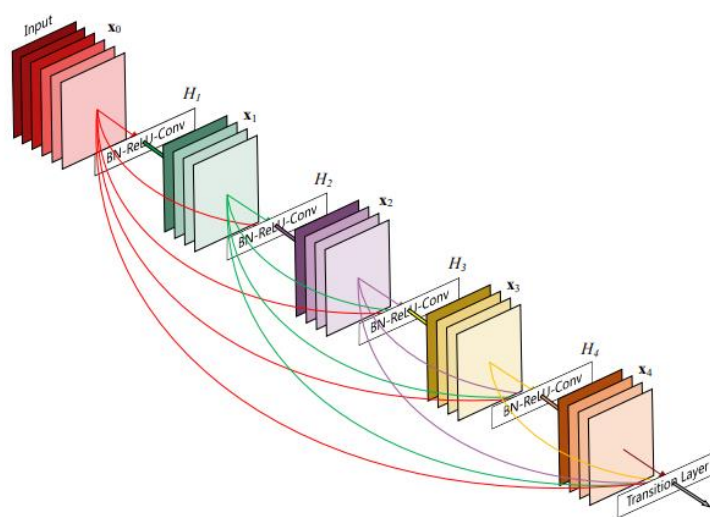
BN: Key module of DNNs

- Key module in current the state-of-the-art network architectures



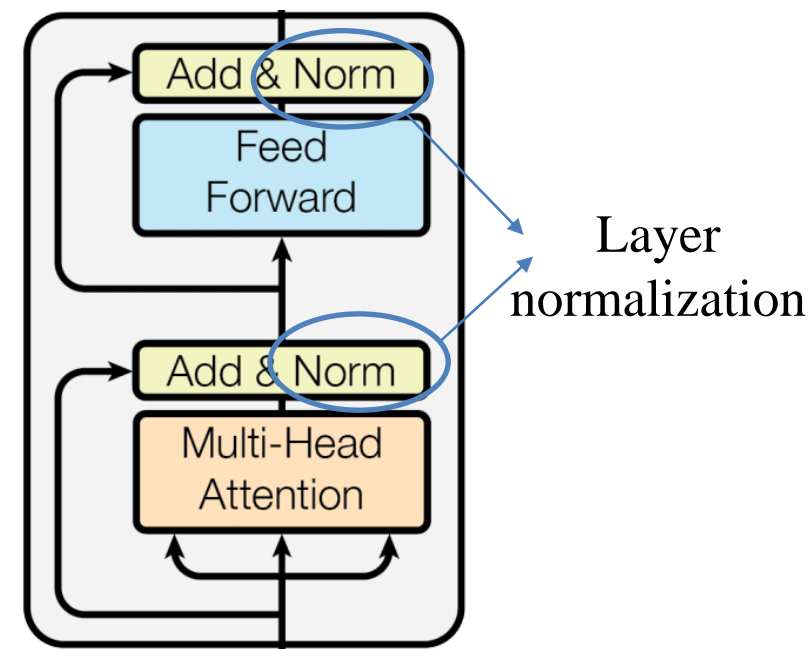
Residual Network,
CVPR 2016

Citation: 81000+



DenseNet,
CVPR 2017

Citation: 16000+



Transformer,
NeurIPS 2017

Citation: 22000+

Recently Papers Involving Normalization

- Statistics of publications on main conferences between 2020 to 2021
 - ICLR, AAAI, CVPR, ICML, ECCV, NeurIPS: 77+

	2020	2021
ICLR	4+	6+
AAAI	7+	6+
CVPR	11+	15+
ICML	5+	7+
ECCV	7+	
NeurIPS	9+	

Recently Papers Involving Normalization

- Statistics of publications on main conferences between 2020 to 2021
 - ICLR, AAAI, CVPR, ICML, ECCV, NeurIPS: 77+

	2020	2021
ICLR	4+	6+
AAAI	7+	6+
CVPR	11+	15+
ICML	5+	7+
ECCV	7+	
NeurIPS	9+	

{ ***normalization/whitening

{ 1.Rethink BN in ***
2.Towards understanding
***normalization from ***

{ 1.***normalization for ***(tasks)
2.Rethink ***normalization in ***
(tasks)



Some questions

- Why so many normalization variants? What are the main motivations behind them? And how can we present a taxonomy?
- How can we reduce the gap between empirical success of normalization techniques and our theoretical understanding of them?
- What recent advances have been made in designing/tailoring normalization techniques for different tasks, and what are the main insights behind them?



Outline

01. Motivations of
Normalization Techniques

02. Introduction of
Normalization Methods

03. Analyses of
Normalization

04. Applications of
Normalization



CVPR
VIRTUAL JUNE 19-25

Disclaimers

- Inevitably miss important related work
- Citations are only representative examples
- There are no consistent understandings of normalization, this is one

Survey paper: “**Normalization Techniques in Training DNNs: Methodology, Analysis and Application,**”
arXiv:2009.12836

Paper list related: <https://github.com/huangleiBuaa/NormalizationSurvey>



Outline

01. Motivations of
Normalization Techniques

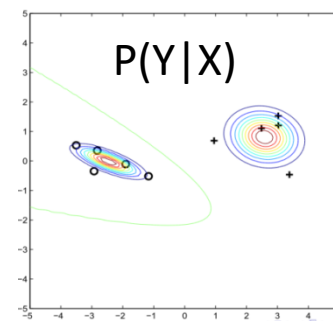
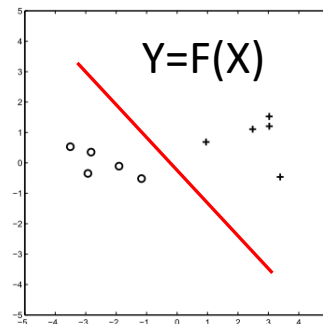
02. Introduction of
Normalization Methods

03. Analyses of
Normalization

04. Applications of
Normalization

Supervised Learning

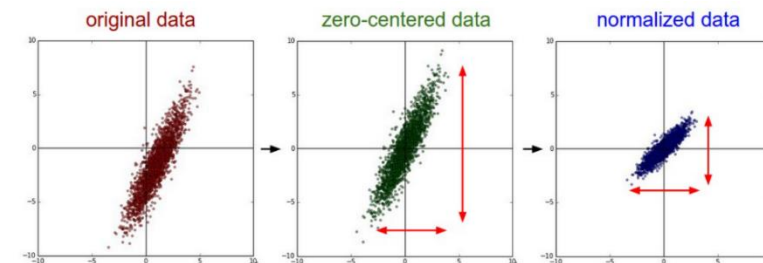
- Dataset $D=\{X, Y\}$
 - Input: X
 - Output: Y
 - Learning: $Y = F(X)$ or $P(Y|X)$
- Main types of learning models
 - Non-parametric model
 - $Y=F(X; x_1, x_2 \dots x_n)$
 - Parametric model
 - $Y=F(X; \theta)$
- Training (Fitting) : $\min_{\theta} \mathcal{L} = E_D \ell(F(X; \theta), \hat{Y})$
- Generalization



Normalization

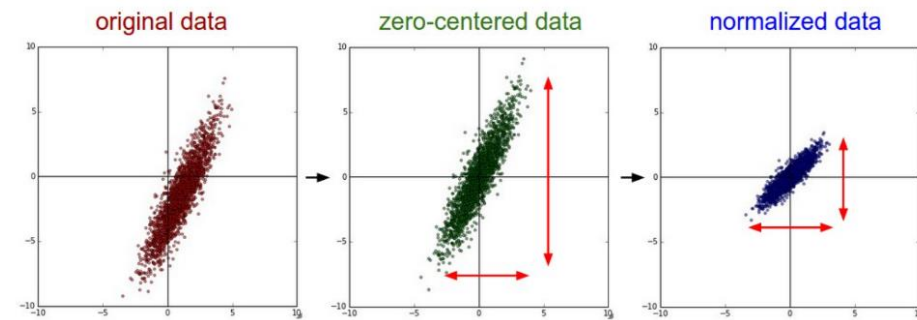
- Definition of normalization
 - In statistics: adjustments of values or distributions in statistics
 - In image processing: changing the range of pixel intensity values
 - In data processing: general reduction of data to canonical form
- Definition of normalization in this tutorial
 - Given a set of data $\mathbb{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$, the normalization operation is a function $\Phi: \mathbf{x} \mapsto \hat{\mathbf{x}}$, which ensures that the transformed data $\hat{\mathbb{D}} = \{\hat{\mathbf{x}}^{(i)}\}_{i=1}^N$ has **certain statistical properties**.

$$\hat{\mathbf{x}} = \frac{\mathbf{x} - \mu}{\sigma}$$

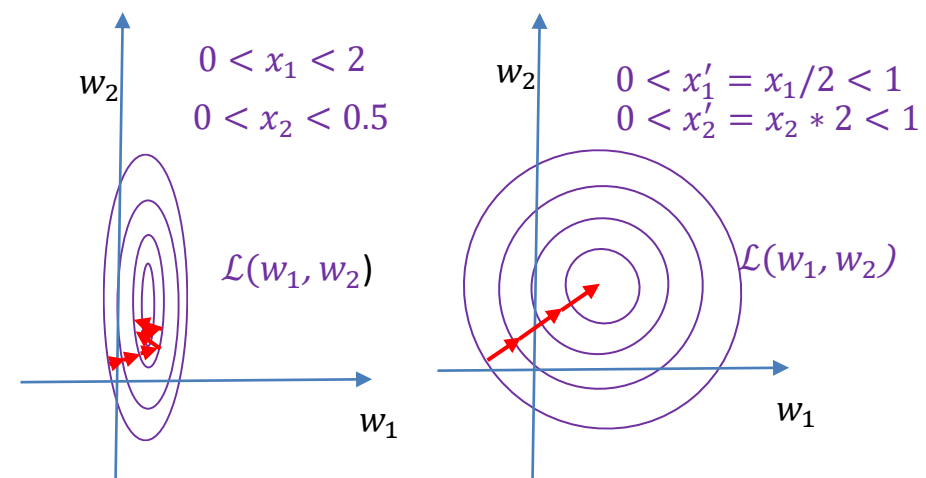


Motivation of Normalizing Input

- Improve the effects of learning
 - Non-parameter models (KNN, Kernel SVM)
 - Distance/ Similarity
- Improve optimization efficiency
 - Parametric model (logistic regression)
 - Update parameters iteratively



$$y = w_1x_1 + w_2x_2 + b, \mathcal{L} = (y - \hat{y})^2$$
$$\theta = \{w_1, w_2\}$$



$$\theta_{t+1} \leftarrow \theta_t - \eta \frac{\partial \mathcal{L}}{\partial \theta_t}$$

Normalizing Input Benefits Optimization

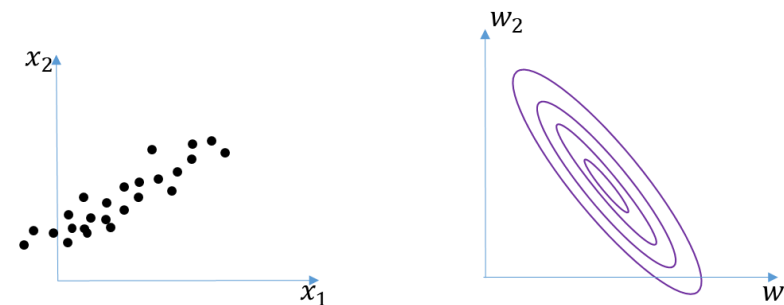
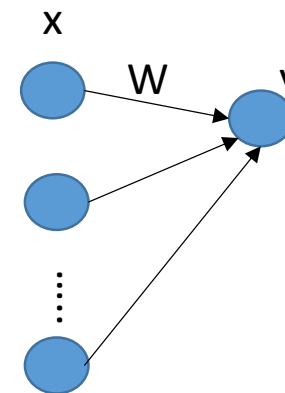
- Linear regression: $\mathcal{L}(F(\mathbf{x}), \hat{y}) = \frac{1}{2} (W\mathbf{x} - \hat{y})^2$

$$\mathcal{L}(W) = \frac{1}{2} (W^T C W - 2A^T W + b)$$

– Where $C = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^T \mathbf{x}$ is the covariance matrix

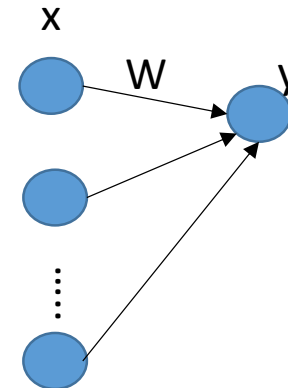
– Gradient: $\frac{\partial \mathcal{L}}{\partial W} = \sum_{i=1}^N \mathbf{x}(y - \hat{y})$

– Hessian matrix: $\mathbf{H} = \frac{\partial^2 \mathcal{L}}{\partial W \partial W} = C$

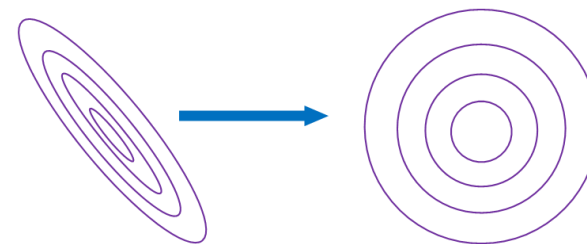


Normalizing Input Benefits Optimization

- Learning dynamics are controlled by the spectrum of curvature matrix (Hessian \mathbf{H})
 - $\lambda_{\max}(\mathbf{H})$:
 - Optimal learning rate: $\eta = \frac{1}{\lambda_{\max}(\mathbf{H})}$
 - Diverge if $\eta > \frac{2}{\lambda_{\max}(\mathbf{H})}$
 - Condition number $\kappa = \frac{\lambda_{\max}(\mathbf{H})}{\lambda_{\min}(\mathbf{H})}$ control the iterations required for convergence



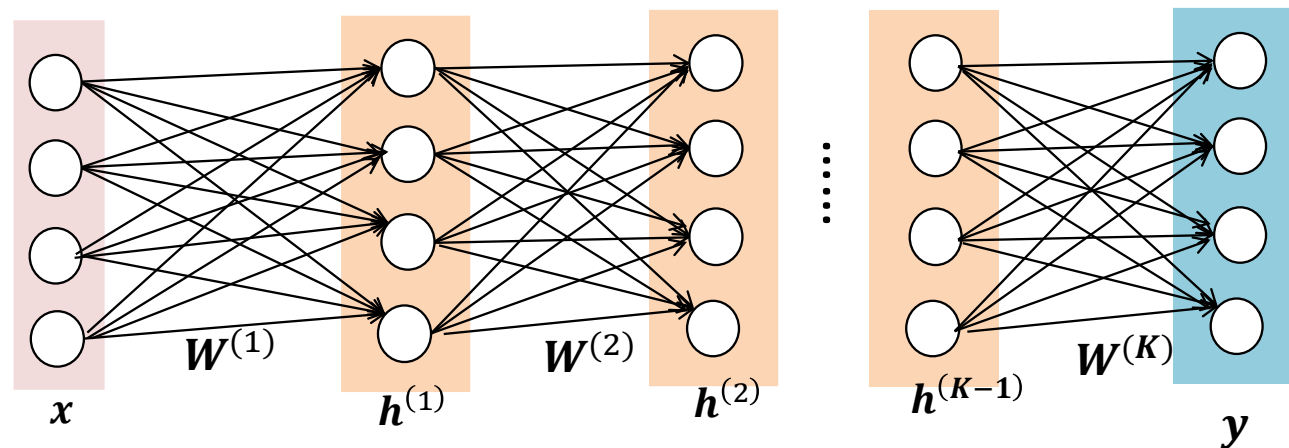
- Hessian of multiple output:
$$\mathbf{H} = \mathbb{E}_{\mathcal{D}}(\mathbf{x}\mathbf{x}^T) \otimes \mathbf{I}$$



Towards Normalizing Activations of DNNs

- Difficulty of analysis for DNNs

- Nonlinear model
- x is only linearly connected by $W^{(1)}$; Optimization is over θ , not $W^{(1)}$ only



$$\theta = \{W^{(1)}, W^{(2)}, \dots, W^{(K)}\}$$

- What we can exploit?

- Layer-wise structure
- $h^{(i)}$ is linearly connected by $W^{(i+1)}$



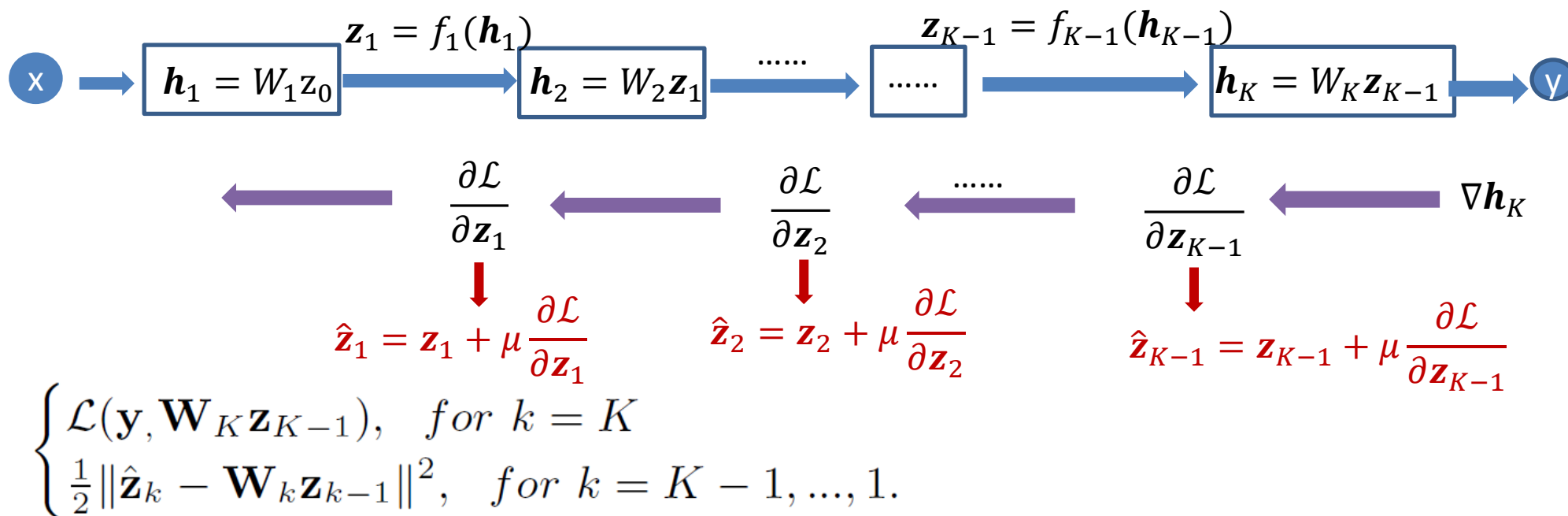
**Normalizing
Activations**

Intuitive Motivation for Normalizing Activation

- Proximal back-propagation

$$\tilde{\mathcal{L}}(\theta, \mathbf{z}) = \mathcal{L}(\mathbf{y}, f_K(\mathbf{W}_K, \mathbf{z}_{K-1})) + \sum_{k=1}^{K-1} \frac{\lambda}{2} \|\mathbf{z}_k - f_k(\mathbf{W}_k, \mathbf{z}_{k-1})\|^2$$

- Back-match propagation

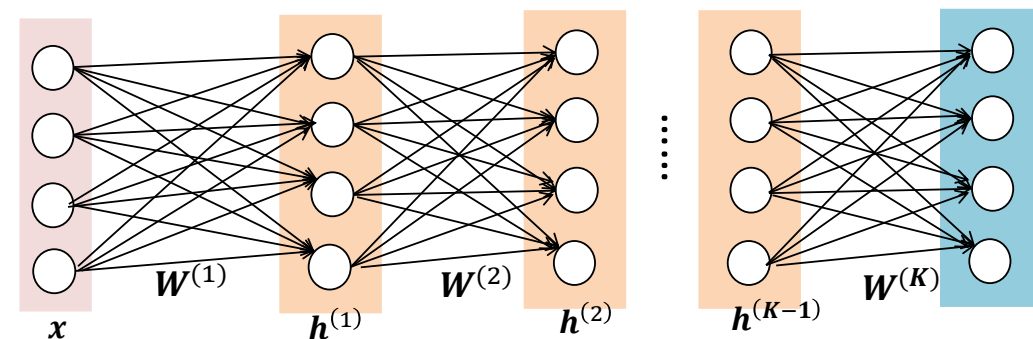


Theoretic Analysis for Normalizing Activation

- Fisher Information Matrix (FIM): $\mathbf{F} = \mathbb{E}_{p(\mathbf{x}), q(\mathbf{y}|\mathbf{x})} \left(\frac{\partial \ell^T}{\partial \theta} \frac{\partial \ell}{\partial \theta} \right)$
- Foundation: approximating FIM using the Kronecker product (K-FAC)
 - Assumption1: weight-gradients in different layers are assumed to be uncorrelated
 - Assumption2: the input and output-gradient in each layer are approximated as independent

$$\mathbf{F} = \text{diag}(F_1, \dots, F_K)$$

$\mathbb{R}^{Kd^2 \times Kd^2}$
 $\mathbb{R}^{d^2 \times d^2}$



$$F_k = \mathbb{E}_{p(\mathbf{x}), q(\mathbf{y}|\mathbf{x})} \left((\mathbf{x}_k \mathbf{x}_k^T) \otimes \left(\frac{\partial \ell}{\partial \mathbf{h}_k}^T \frac{\partial \ell}{\partial \mathbf{h}_k} \right) \right) \approx \mathbb{E}_{p(\mathbf{x})} (\mathbf{x}_k \mathbf{x}_k^T) \otimes \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left(\frac{\partial \ell}{\partial \mathbf{h}_k}^T \frac{\partial \ell}{\partial \mathbf{h}_k} \right)$$

$\mathbb{R}^{d \times d}$

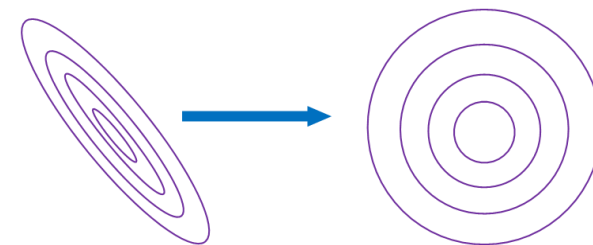
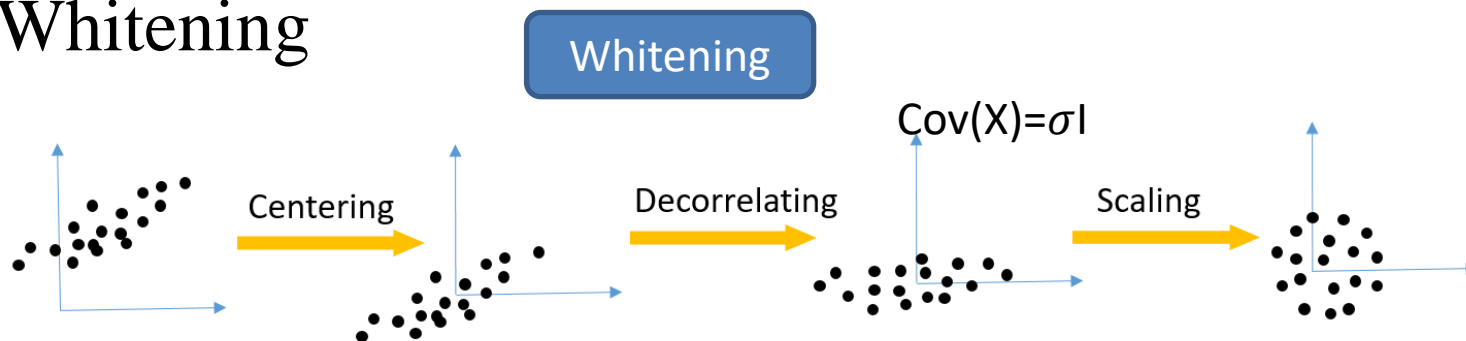
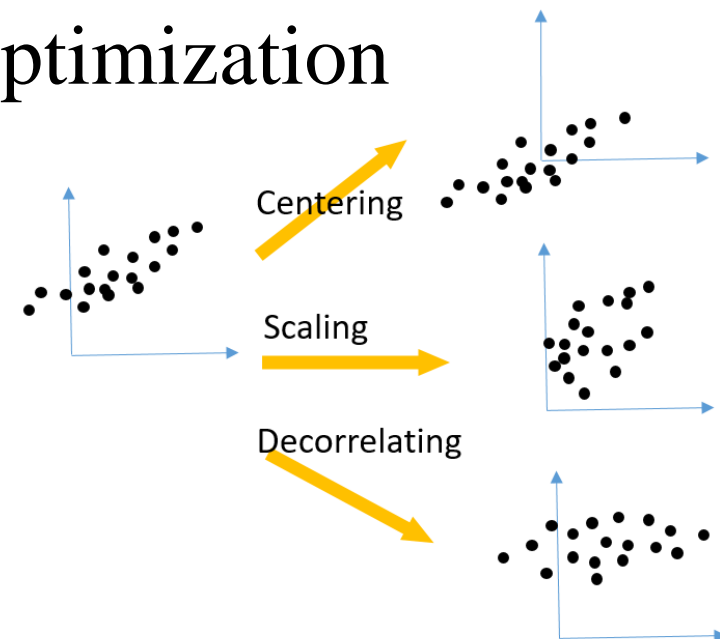
Construct Well-Conditioned Landscape

- Denoting $\Sigma_x = \mathbb{E}_{p(x)} (xx^T)$ and $\Sigma_{\nabla h} = \mathbb{E}_{p(x), q(y|x)} \left(\frac{\partial \ell^T}{\partial h} \frac{\partial \ell}{\partial h} \right)$
- Criteria
 - 1. The statistics of the layer input (e.g., Σ_x) and output-gradient ($\Sigma_{\nabla h}$) across different layers are equal (**across layer**)
 - 2. Σ_x and $\Sigma_{\nabla h}$ are well conditioned (**in layer**)
- Initialization techniques: designed to satisfy Criteria 1 and/or 2 **during initialization**
 - Arxiv-Init [Glorot and Bengio, 2010], He-Init [He et al, 2015]: for Criteria 1
 - Orthogonal Initialization [Saxe et al, 2014] : for Criteria 1 and 2
- General goals of “normalization” in DNNs: **Controlling the distribution of the activations/output-gradients during training.**

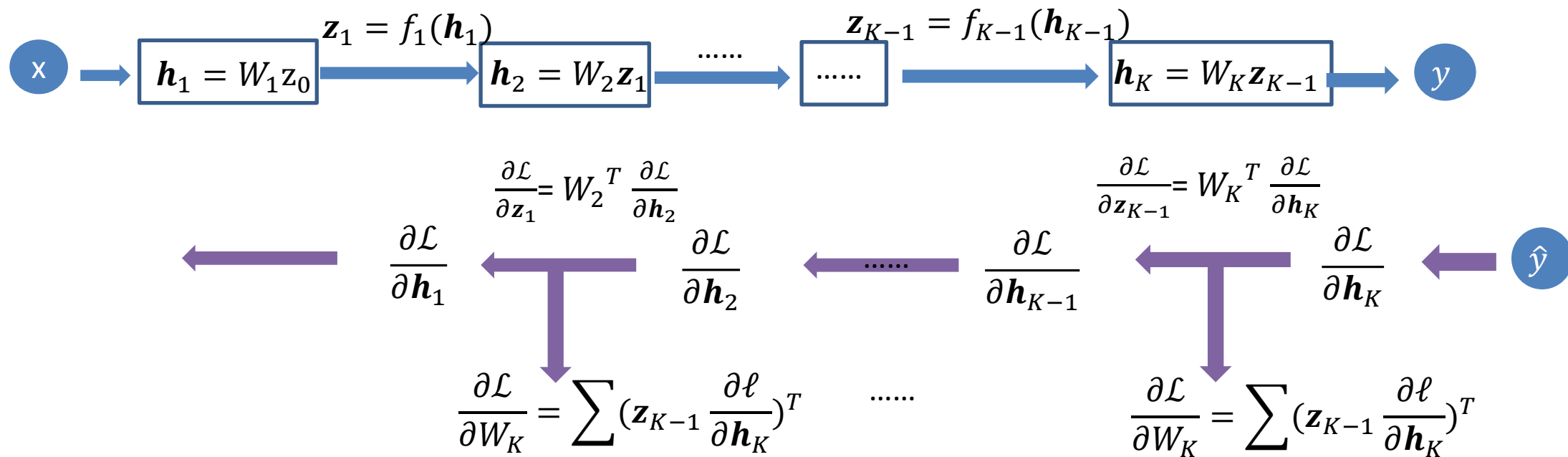
Normalization Operation

- Basic normalization operations benefits the optimization
 - Centering
 - Scaling
 - Decorrelating
- Combine above
 - Standardization
 - Whitening

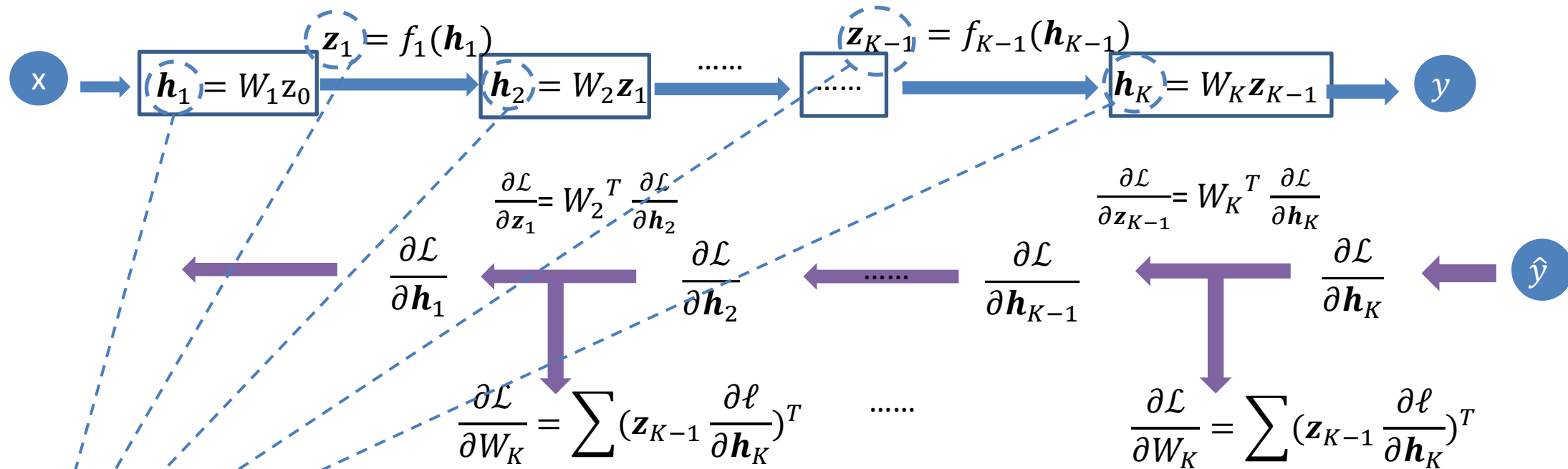
$$H = \mathbb{E}_{\mathcal{D}}(xx^T) \otimes I$$



General Picture of Normalization in DNNs

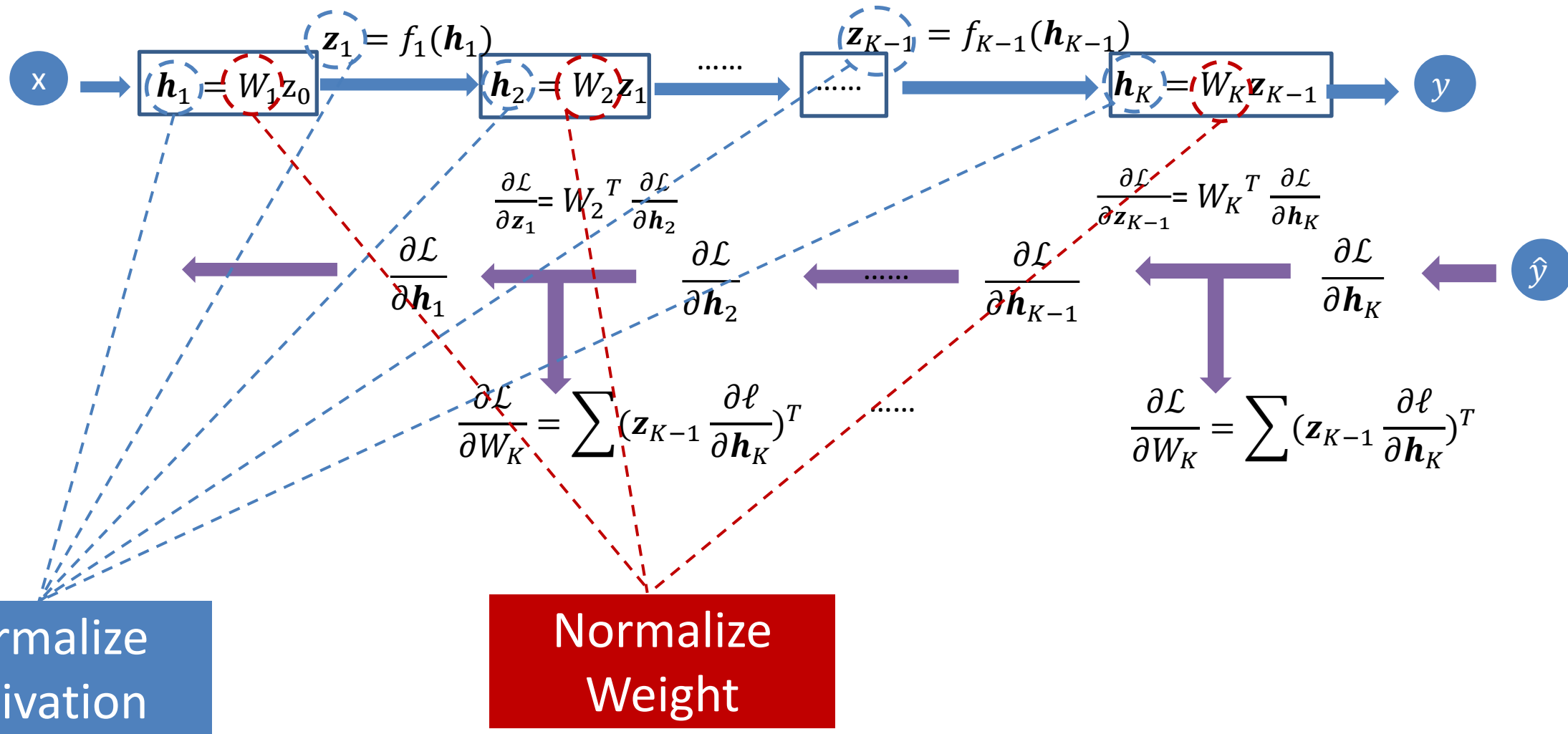


General Picture of Normalization in DNNs



Normalize
Activation

General Picture of Normalization in DNNs



General Picture of Normalization in DNNs

