



The
University
Of
Sheffield.

Introduction to Some Parts of Deep Learning

Jiahui Liu

Supervised by Kevin Li Sun

Department of Computer Science
Faculty of Engineering

Knowledge Distillation (Part 1)

Distilling the knowledge in a neural network

Authors: Hinton G, Vinyals O, Dean J.

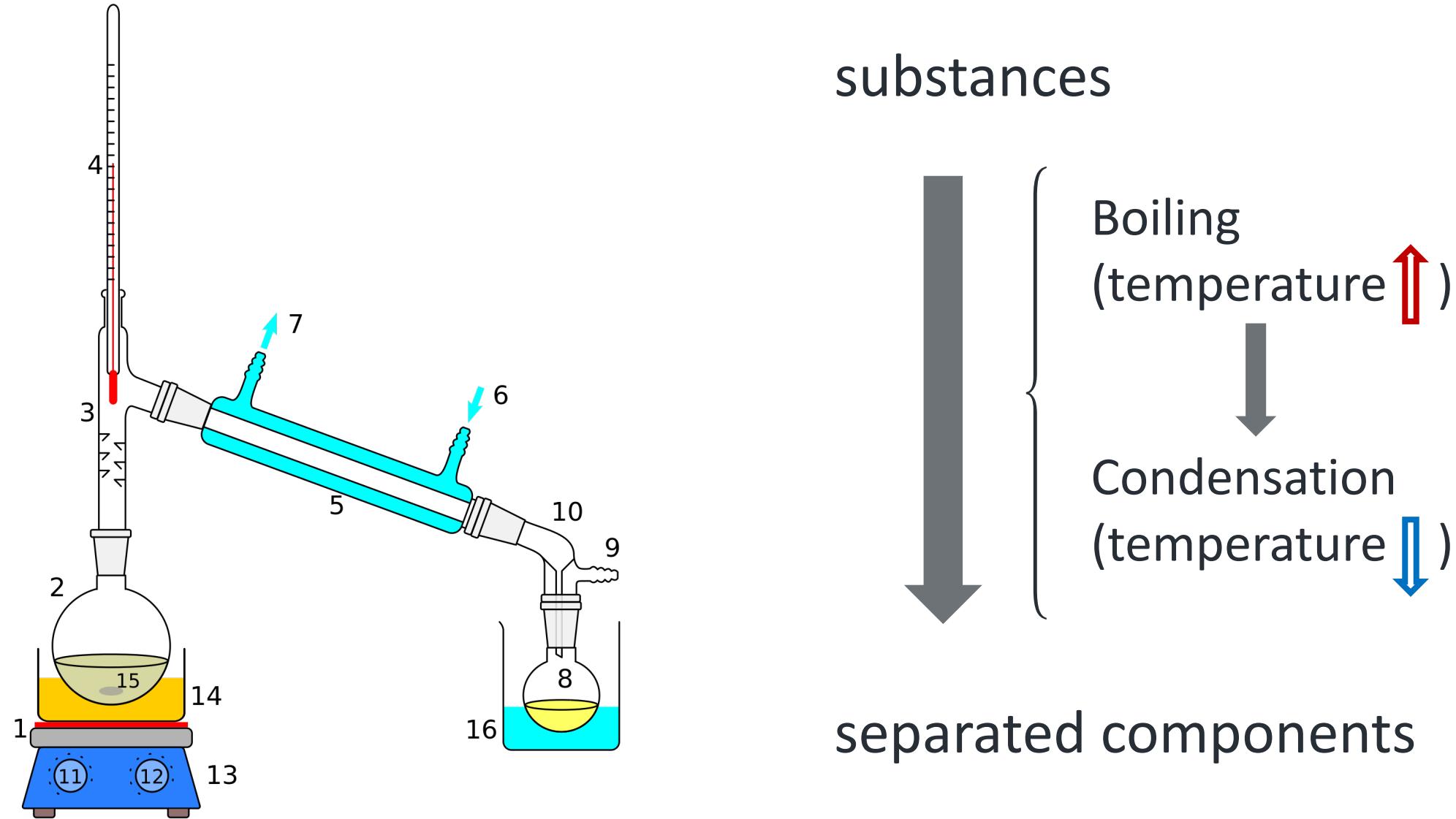
Hinton G, Vinyals O, Dean J. Distilling the knowledge in a neural network[J]. arXiv preprint arXiv:1503.02531, 2015.c

Bayesian Deep Learning (Part 2)

Bayesian Deep Learning

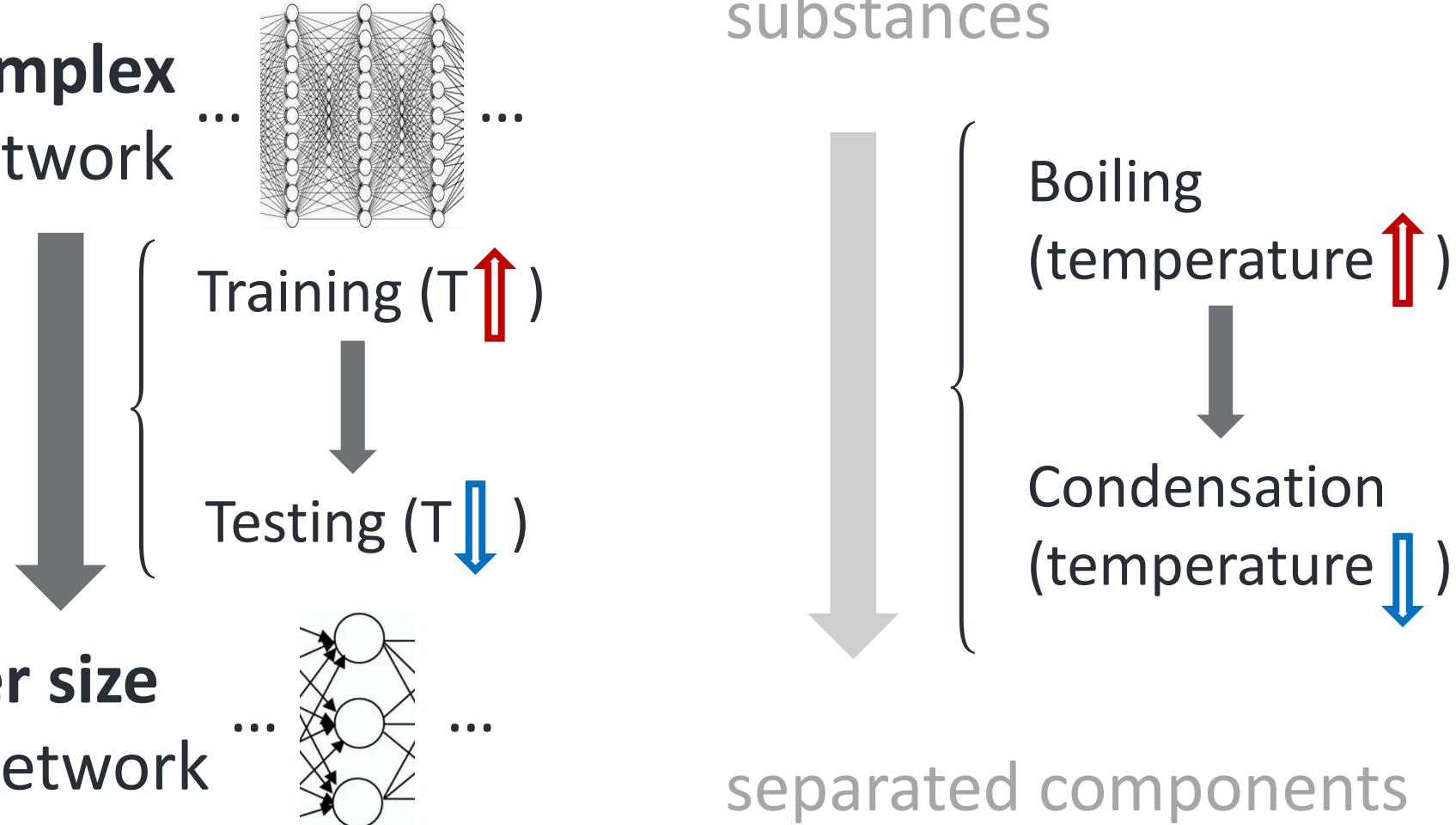
Authors: Yarin Gal

Knowledge Distillation



Knowledge Distillation

a **very complex**
neural network

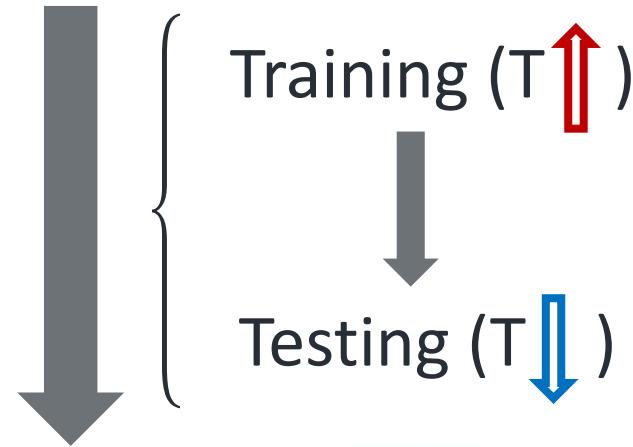
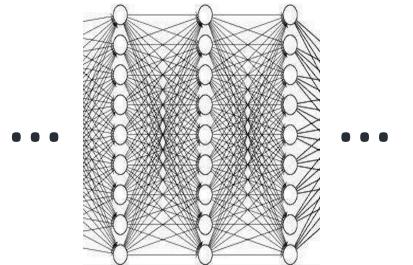


a **smaller size**
neural network

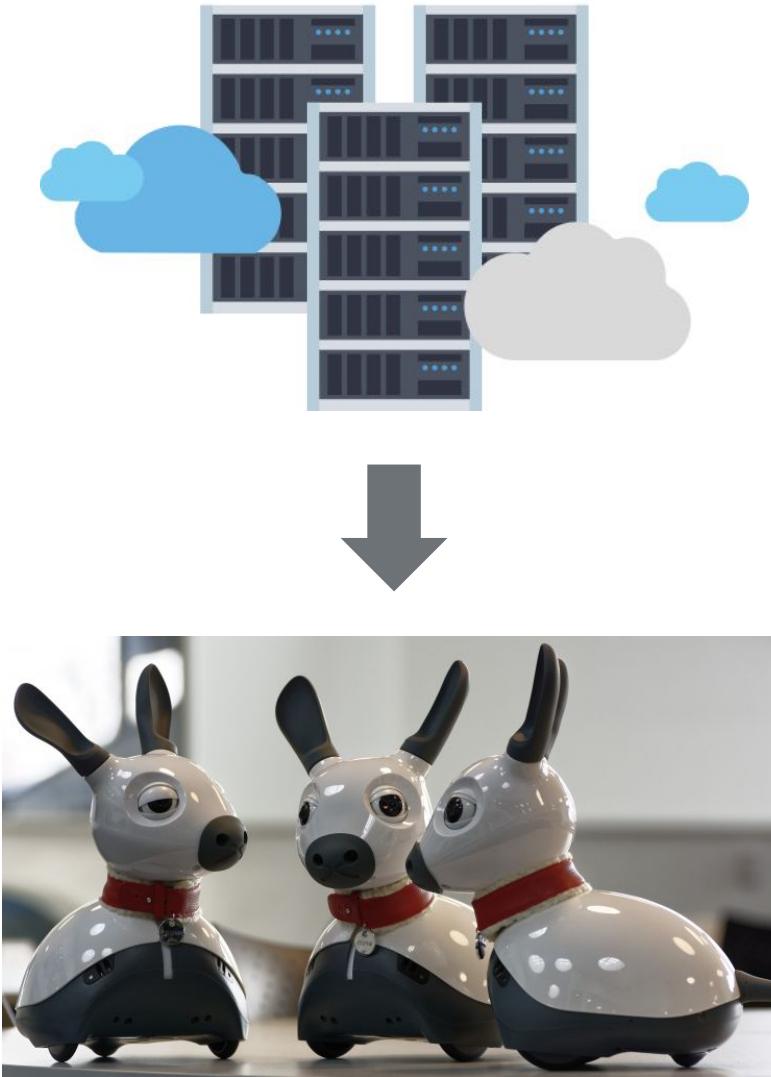
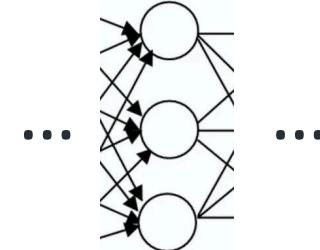
substances
separated components

Knowledge Distillation

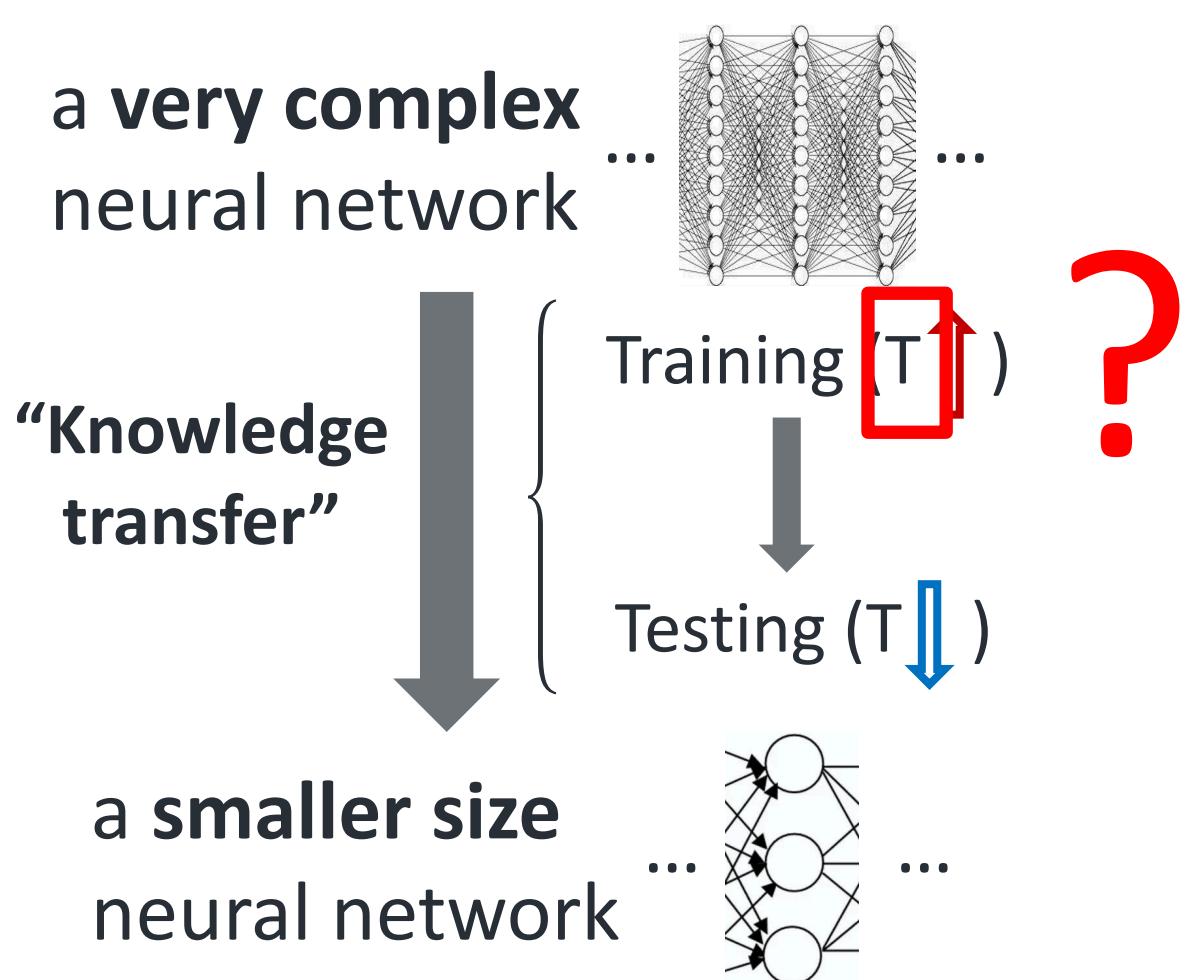
a **very complex**
neural network



a **smaller size**
neural network



Knowledge Distillation



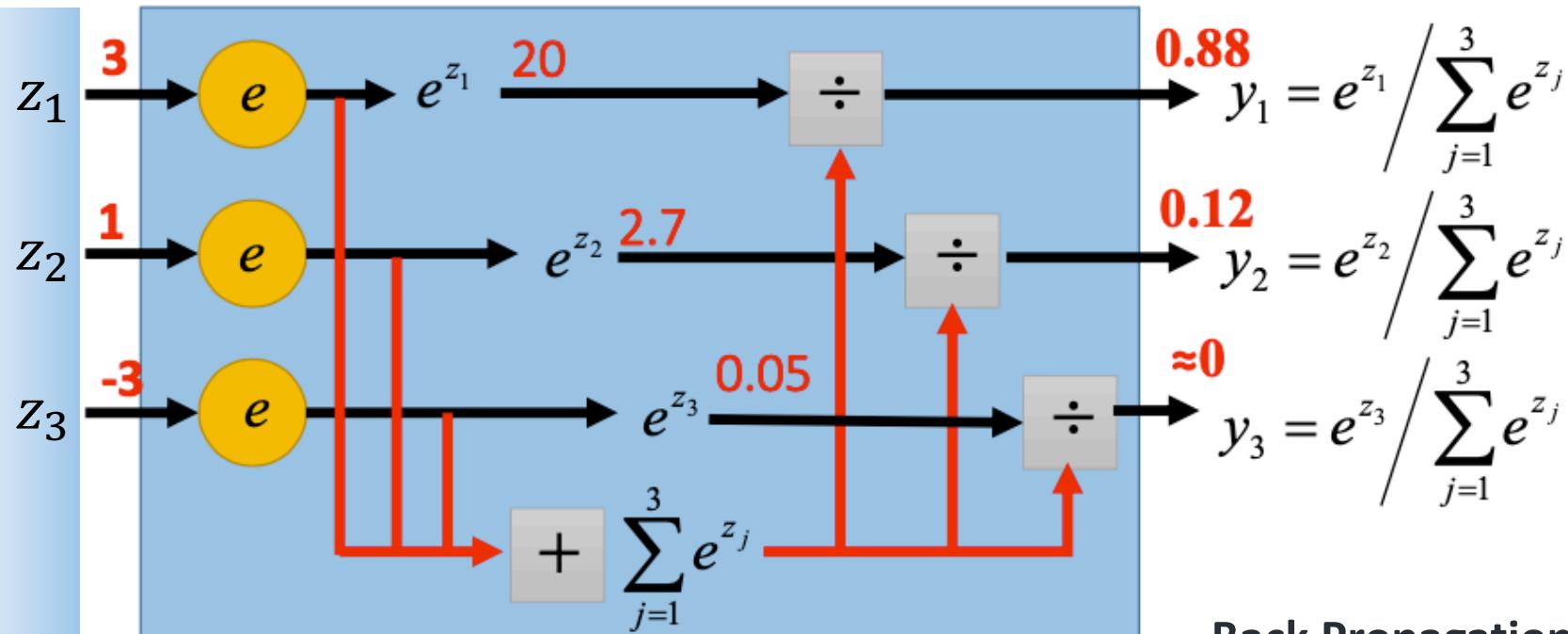
Teacher Network
(Well trained)

Student Network

Softmax

Softmax Layer

logits



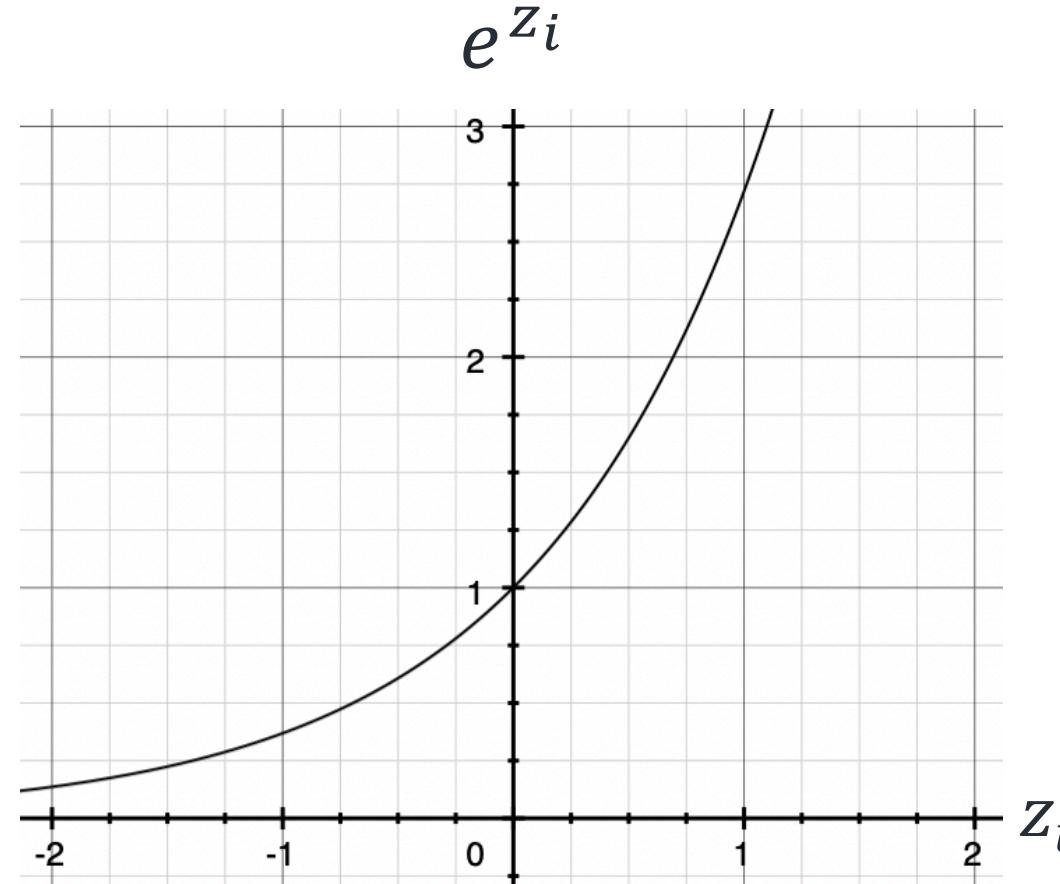
For example (training):

| Output | Target |
|--------|--------|
| 0.88 | 1 |
| 0.12 | 0 |
| 0 | 0 |

Softmax

For the Output:

$$q_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$$



| “Soft” | “Hard” |
|--------|--------|
| Output | Target |
| 0.88 | 1 |
| 0.12 | 0 |
| 0 | 0 |

Loss

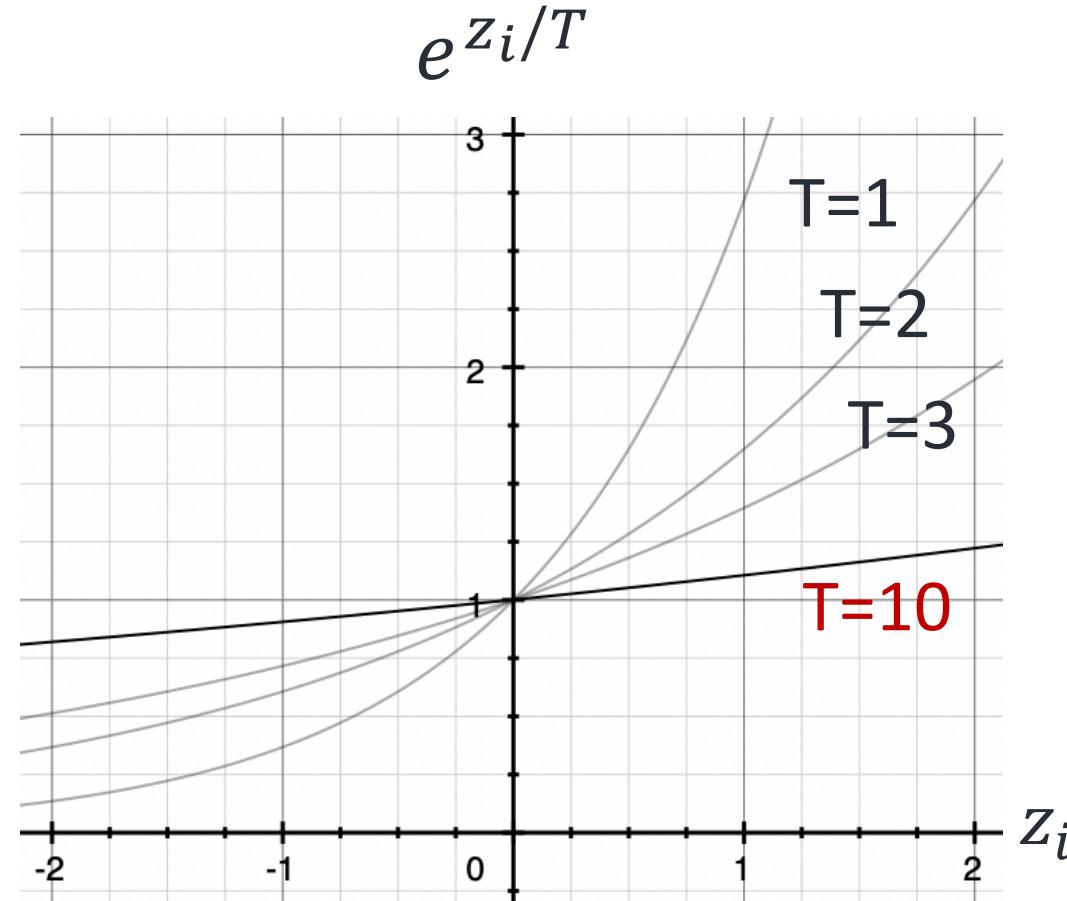
Knowledge Distillation

For the Output:

$$q_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

↓ “Softer”

$$q_i = \frac{e^{z_i/T}}{\sum_j e^{z_j/T}}$$



$T \rightarrow \text{higher}: e^{z_i/T} \rightarrow \text{more gentle } \vec{q} \rightarrow \text{softer}$

| “Soft” | “Hard” |
|--------|--------|
| Output | Target |
| 0.88 | 1 |
| 0.12 | 0 |
| 0 | 0 |

Loss

Knowledge Distillation

For the Output:

$$q_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

↓ “Softer”

$$q_i = \frac{e^{z_i/T}}{\sum_j e^{z_j/T}}$$

Well trained
teacher
network

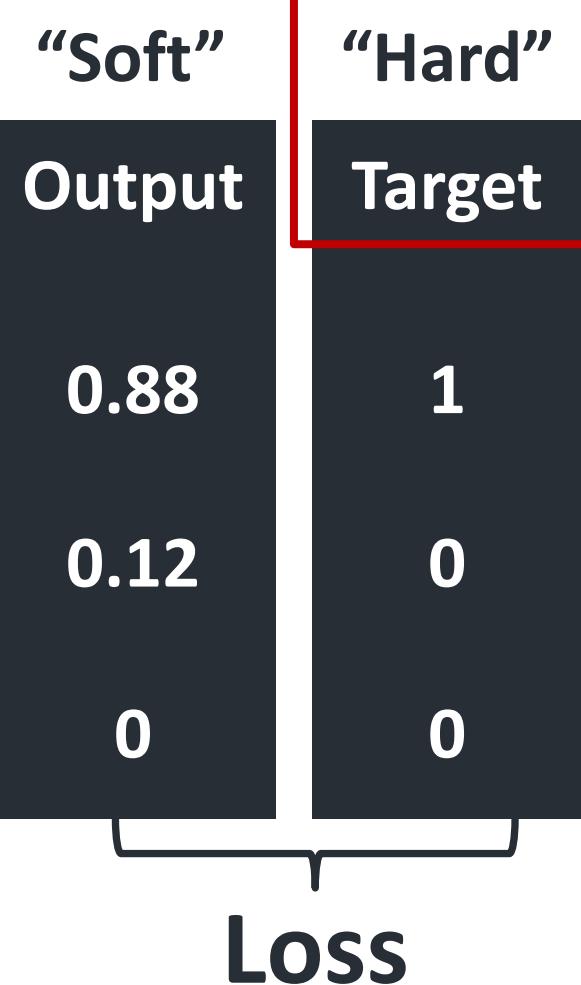
↓
 T

Softer Output
(Soft Target)

$T=3$

Softer
Output
0.61
0.31
0.08

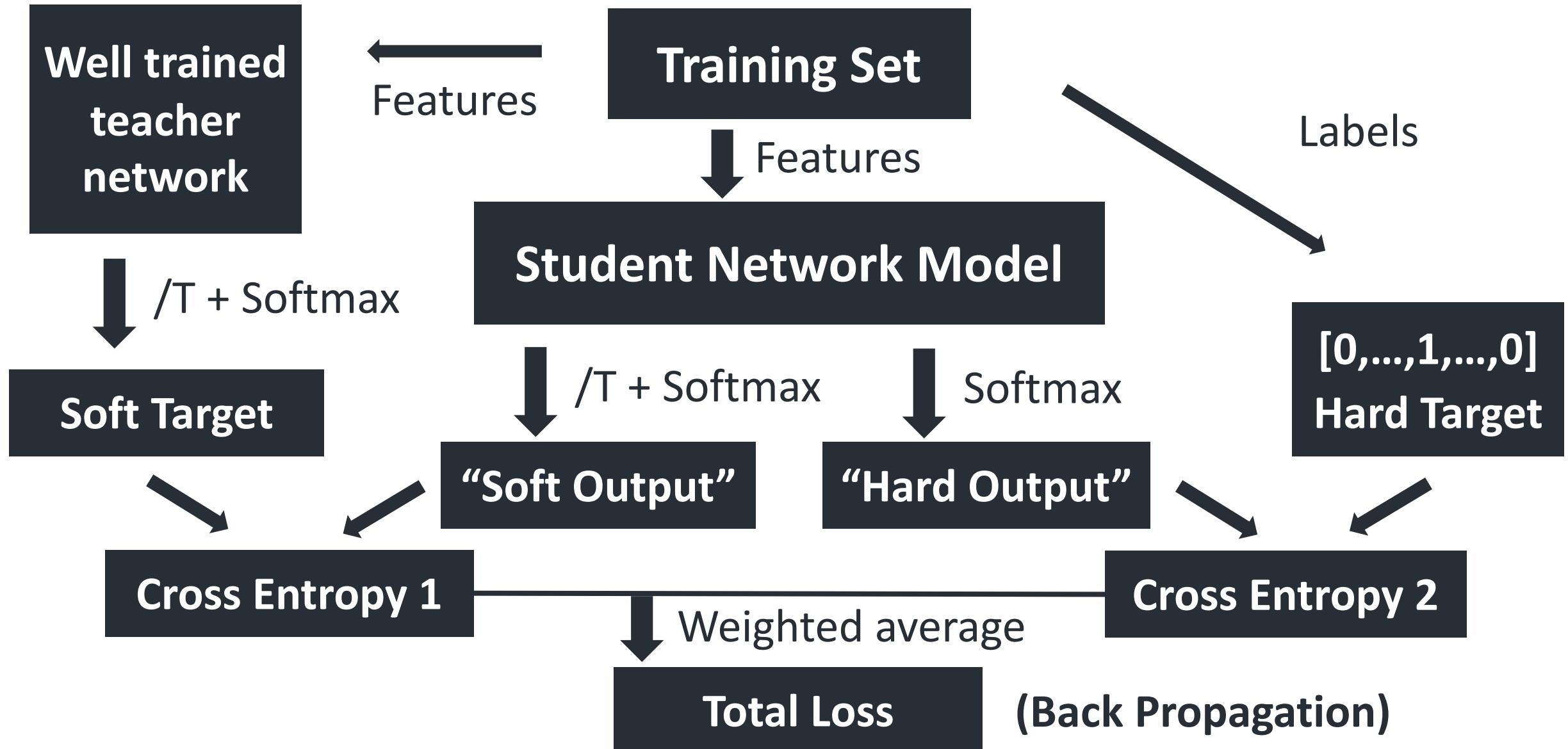
“Softer”
←



$T \rightarrow$ higher: $e^{z_i/T} \rightarrow$ more gentle $\vec{q} \rightarrow$ softer

Knowledge Distillation

(Train a student network)



Knowledge Distillation

For the Soft Target Loss:

Teacher Network logits: v_i

Student Network logits: z_i

$$C = - \sum_{j=1}^C p_j \log q_j$$

$$p_i = \frac{\exp(v_i/T)}{\sum_j \exp(v_j/T)} \quad q_i = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}$$

$$\frac{\partial C}{\partial z_i} = \frac{1}{T} (q_i - p_i) = \frac{1}{T} \left(\frac{e^{z_i/T}}{\sum_j e^{z_j/T}} - \frac{e^{v_i/T}}{\sum_j e^{v_j/T}} \right)$$

$$\frac{\partial C}{\partial z_i} \approx \frac{1}{T} \left(\frac{1 + z_i/T}{N + \sum_j z_j/T} - \frac{1 + v_i/T}{N + \sum_j v_j/T} \right)$$

$$\frac{\partial C}{\partial z_i} \approx \frac{1}{NT^2} (z_i - v_i)$$

In the high temperature limit, distillation is equivalent to minimizing MSE.

Knowledge Distillation (Part 1)

Distilling the knowledge in a neural network

Authors: Hinton G, Vinyals O, Dean J.

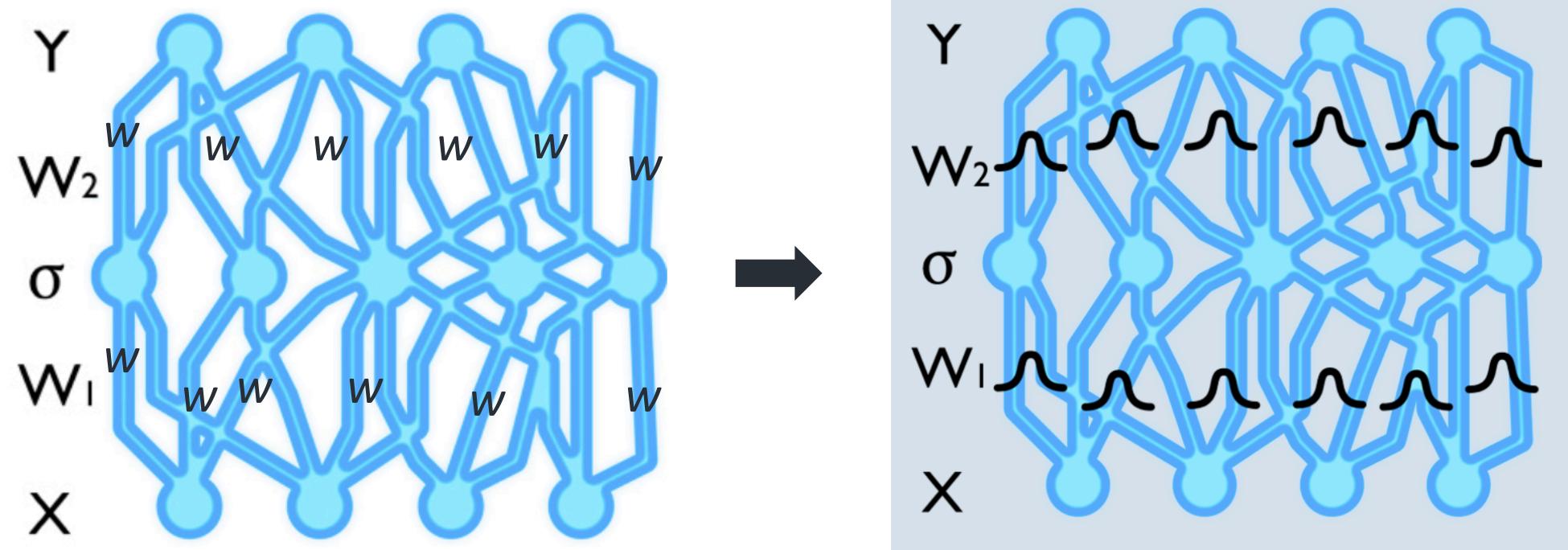
Hinton G, Vinyals O, Dean J. Distilling the knowledge in a neural network[J]. arXiv preprint arXiv:1503.02531, 2015.c

Bayesian Deep Learning (Part 2)

Bayesian Deep Learning

Author: Yarin Gal

Bayesian Deep Learning



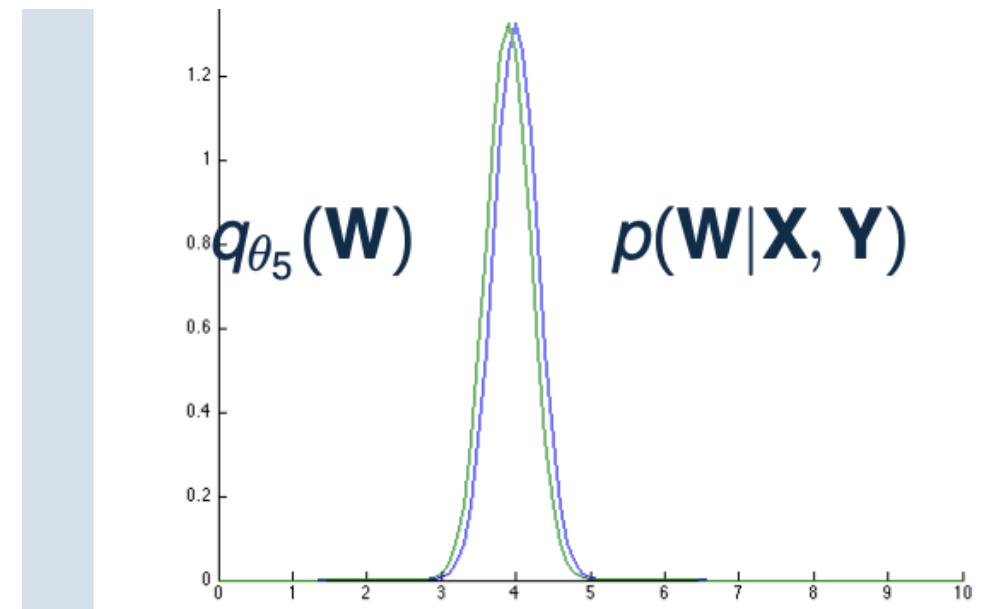
Bayesian Deep Learning

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$$

$$p(w|x, y) = \frac{p(y|x, w)p(w)}{\int p(y|x, w)p(w)dw}$$

Define simple distribution $q_M(W)$
and approximate $q_M(W) \approx p(W|X, Y)$

Variational Inference



KL Divergence (Relative entropy)

Describe the approximation effect of two distributions.

$$p(w|x, y) = \frac{p(y|x, w)p(w)}{\int p(y|x, w)p(w)dw} \quad q_{\mathbf{M}}(\mathbf{W})$$

$$\begin{aligned} KL(q||p) &= \sum_{i=1}^n q(W) \log \frac{q(W)}{p(W|X, Y)} = E_q[\log \frac{q(W)}{p(W|X, Y)}] \\ &= E_q[\log q(W) - \log p(W) - \log p(Y|X, W) + \log p(Y|x)] \\ &= E_q[\log q(W) - \log p(W)] - E_q[\log p(Y|X, W)] + const \end{aligned}$$

Bayesian Deep Learning

Describe the approximation effect of two distributions.

$$p(w|x, y) = \frac{p(y|x, w)p(w)}{\int p(y|x, w)p(w)dw} \quad q_{\mathbf{M}}(\mathbf{W})$$

$$\begin{aligned} KL(q||p) &= \sum_{i=1}^n q(W) \log \frac{q(W)}{p(W|X, Y)} = E_q[\log \frac{q(W)}{p(W|X, Y)}] \\ &= E_q[\log q(W) - \log p(W) - \log p(Y|X, W) + \log p(Y|x)] \\ &= E_q[\log q(W) - \log p(W)] - E_q[\log p(Y|X, W)] + const \end{aligned}$$

$$Loss = E_q[\log q(W) - \log p(W)] - E_q[\log p(Y|X, W)]$$

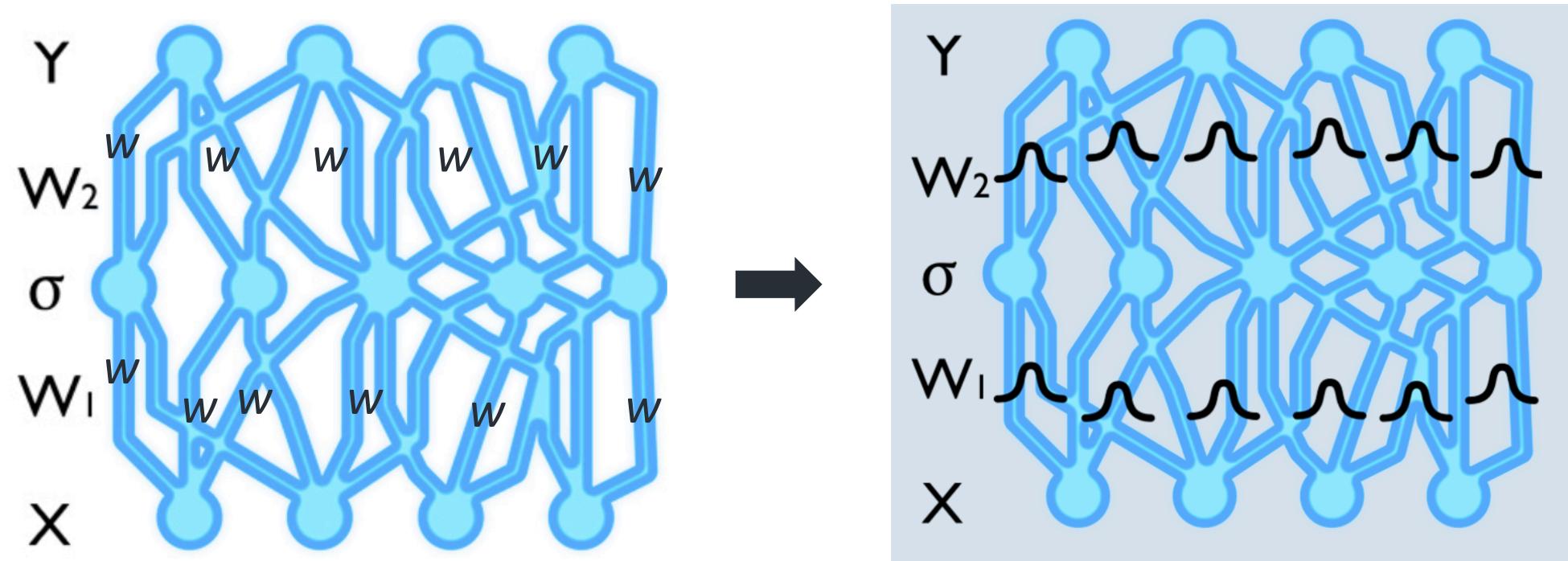
Uncertainty

Aleatoric uncertainty: noise inherent in the data

Epistemic uncertainty: model's lack of knowledge

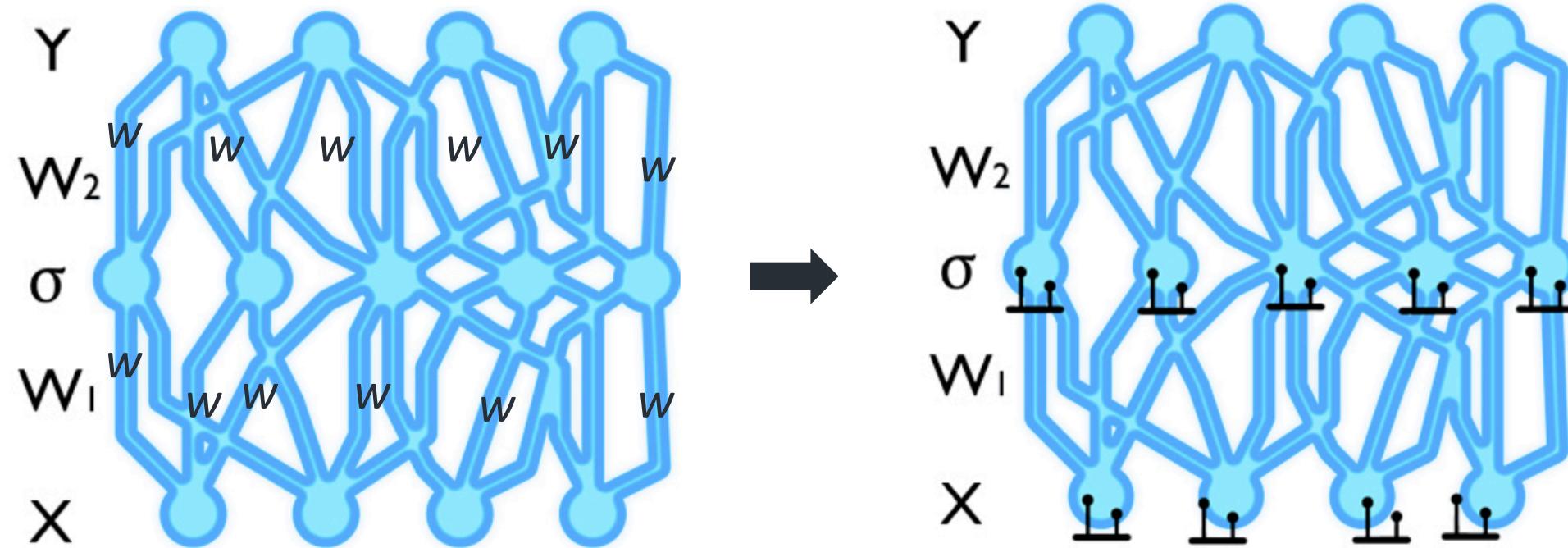
Predictive uncertainty: sum of the two

Bayesian Deep Learning



$q_M(W)$: Gaussian distribution

Bayesian Deep Learning



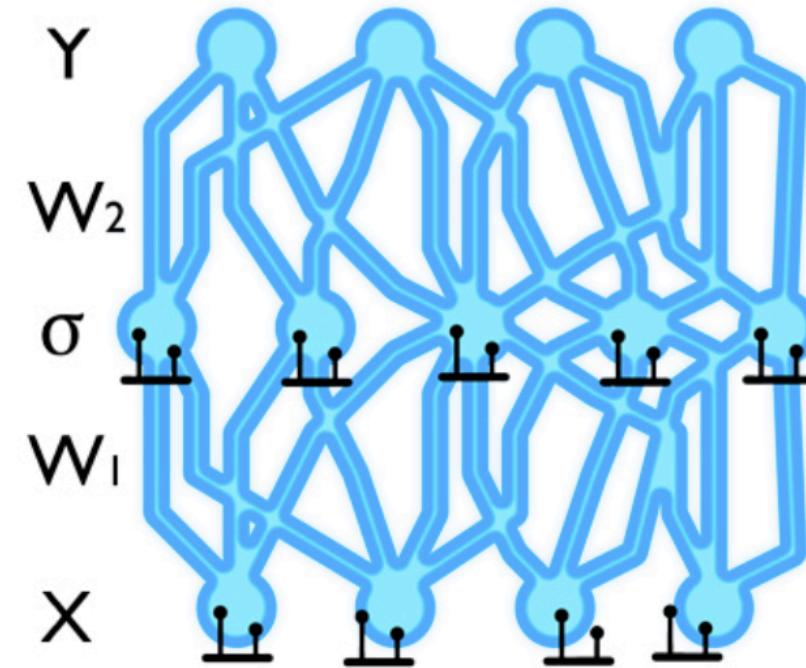
$q_M(W)$: Bernoulli distribution

Bayesian Deep Learning

Variational inference with $q_M(W)$

=

Dropout neural network



$q_M(W)$: Bernoulli distribution

Questions?



The
University
Of
Sheffield.

Thanks for listening

Jiahui Liu
Supervised by Kevin Li Sun
Department of Computer Science
Faculty of Engineering