

# Neural Networks

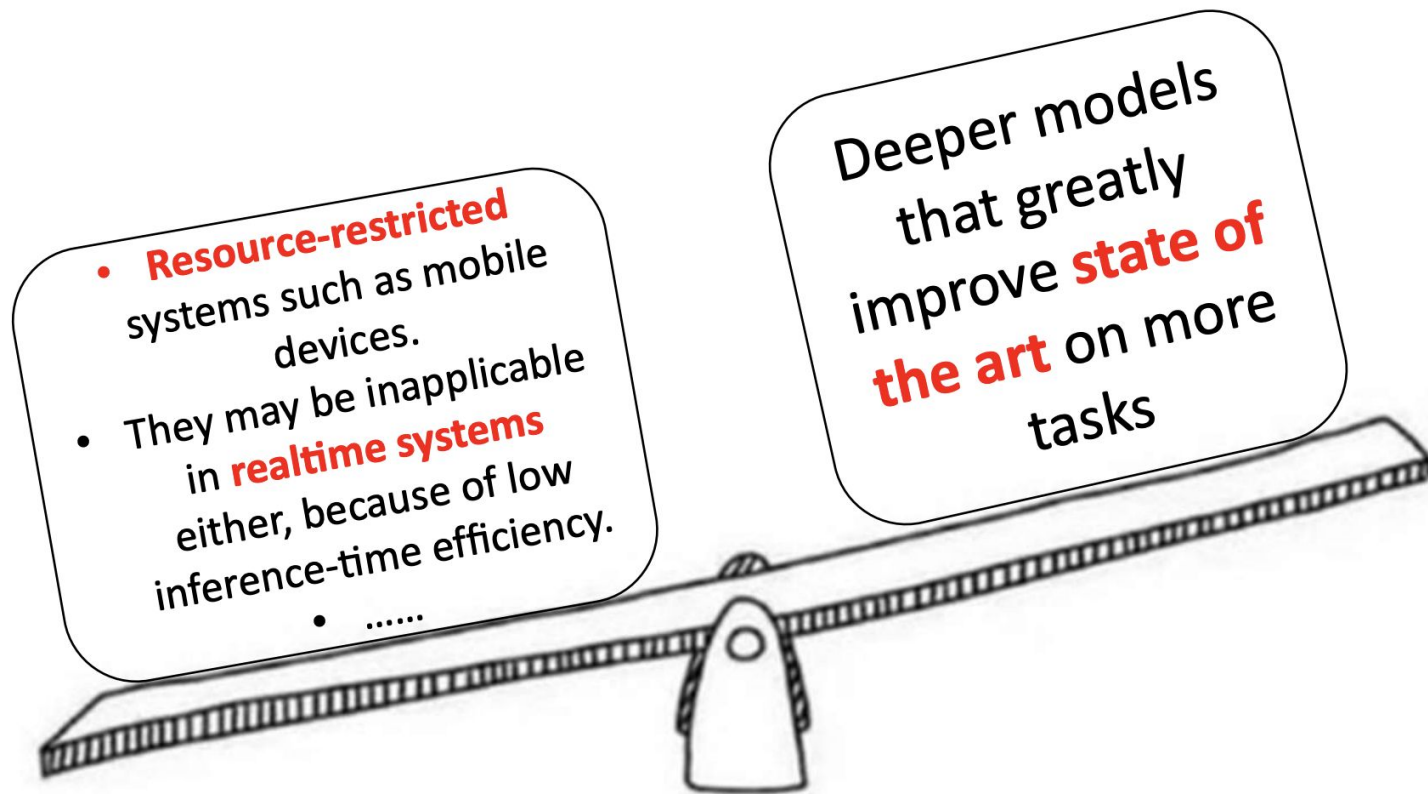
## Speed-up and Compression

**Lecture 5: Knowledge Distillation**

# Outline

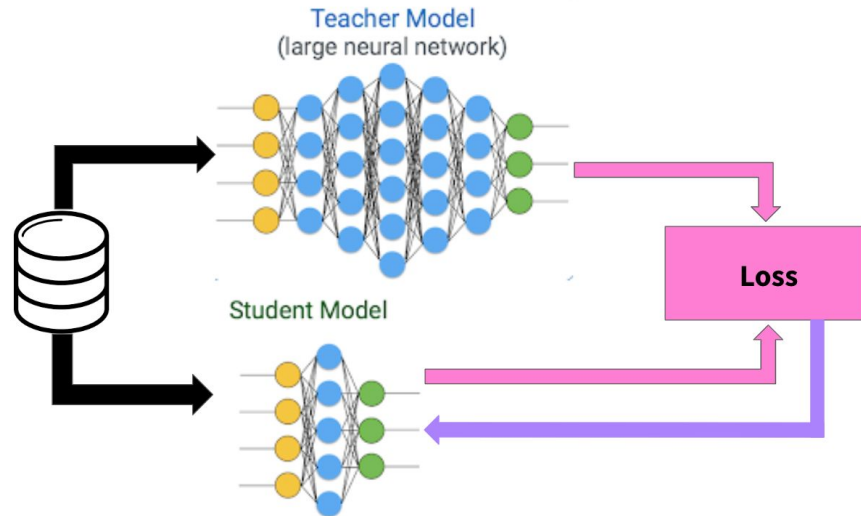
- Motivation
- Knowledge Distillation concept
- Knowledge Distillation examples
  - Using Soft targets
  - Using features (FitNets)
  - Using attention
- Other details of KD

# Motivation

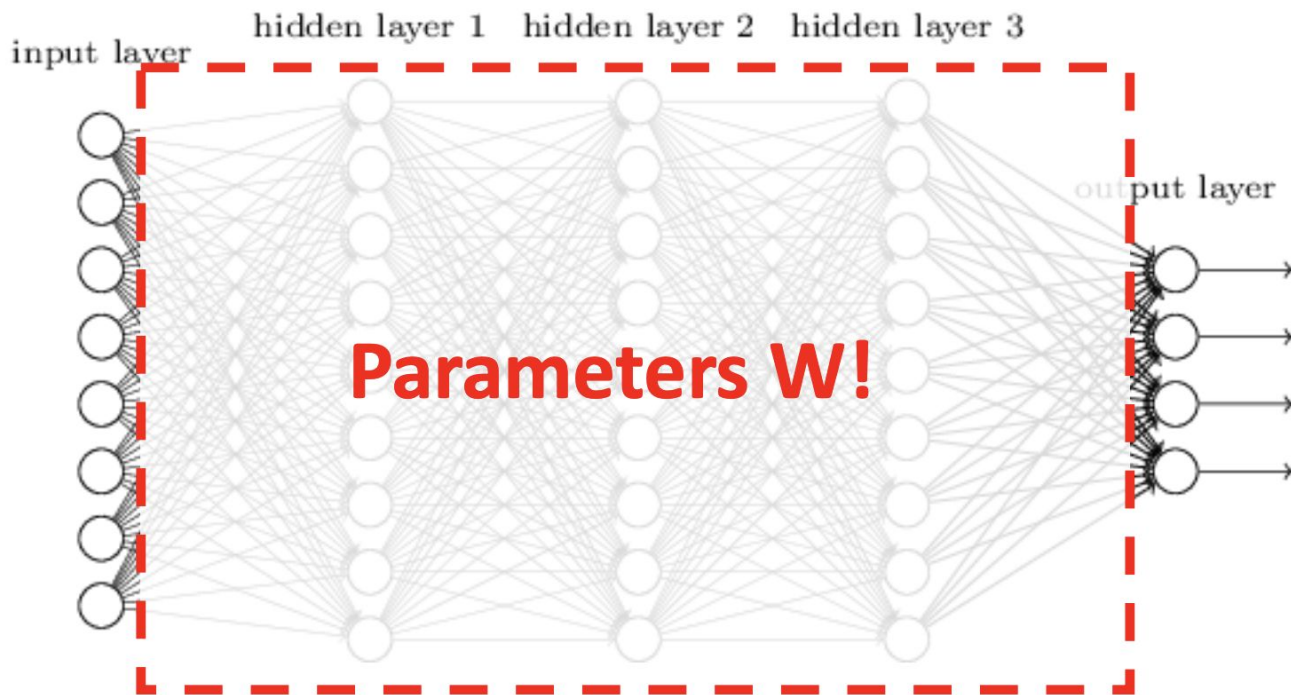


# Knowledge Distillation

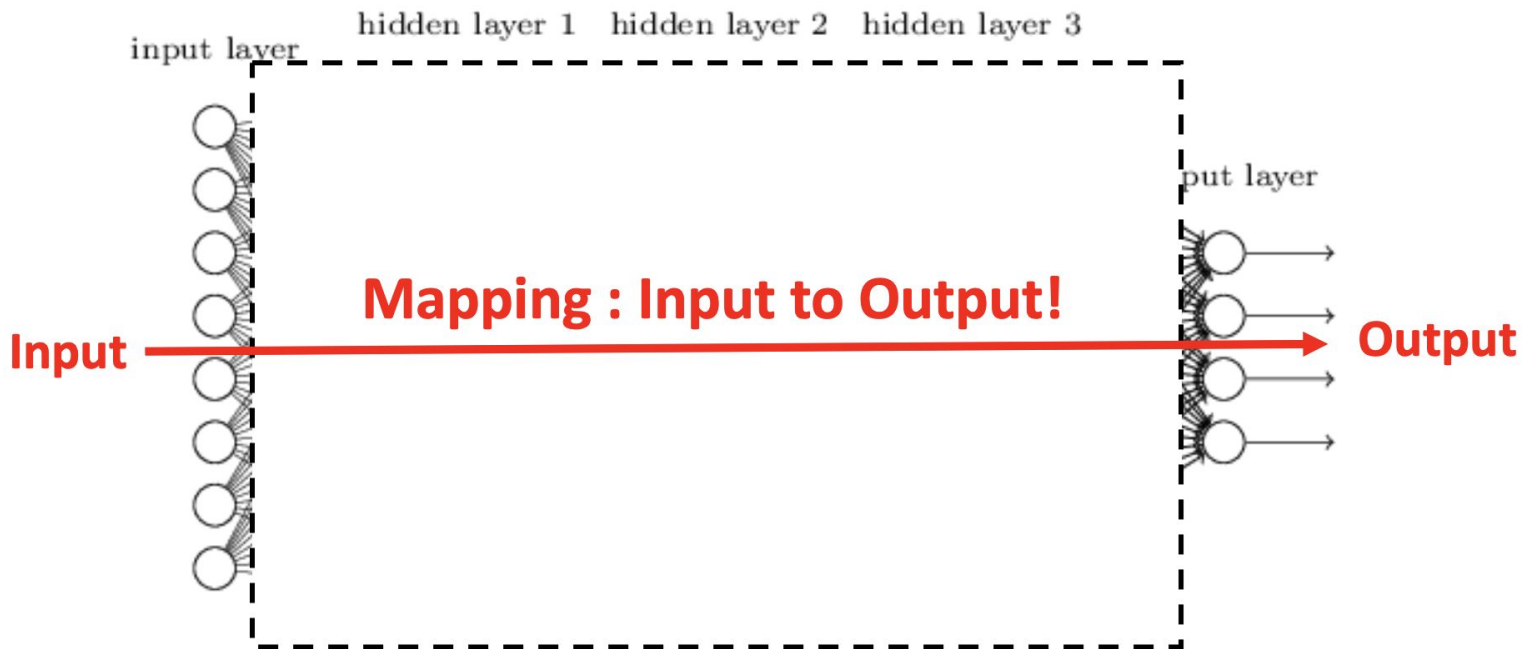
**Knowledge distillation** is a process of distilling or transferring the knowledge from a (set of) large, cumbersome model(s) to a lighter, easier-to-deploy single model



# What is knowledge?



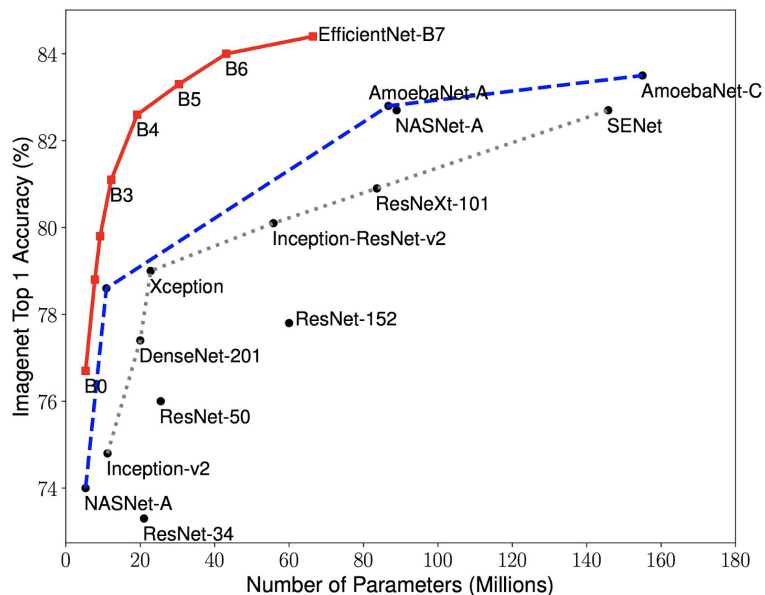
# What is knowledge?



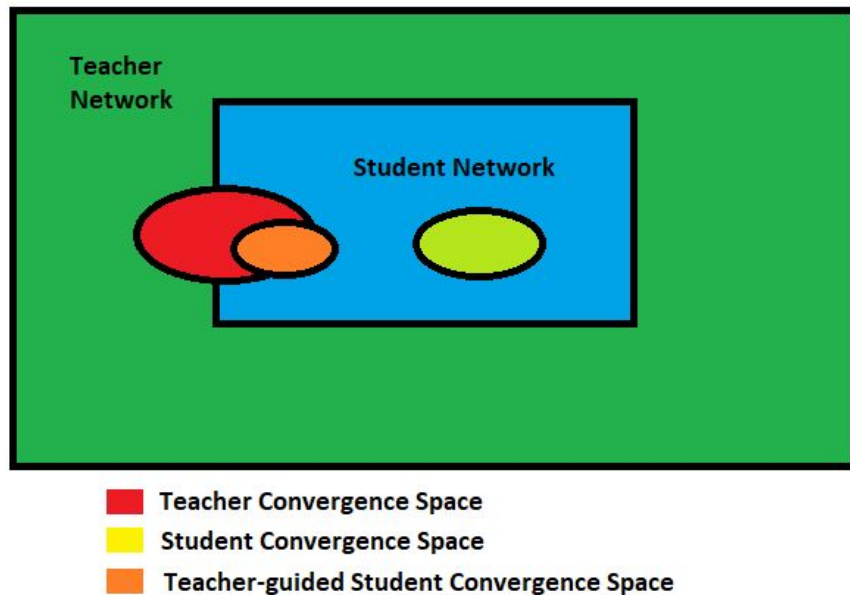
A more abstract view of the **knowledge**, that frees it from any particular instantiation, is that it is a **learned mapping from input vectors to output vectors**.

# Why Knowledge in bigger models is better?

Higher number of parameters -> better performance



We can train smaller (student) model to mimic a behaviour of bigger model (teacher)

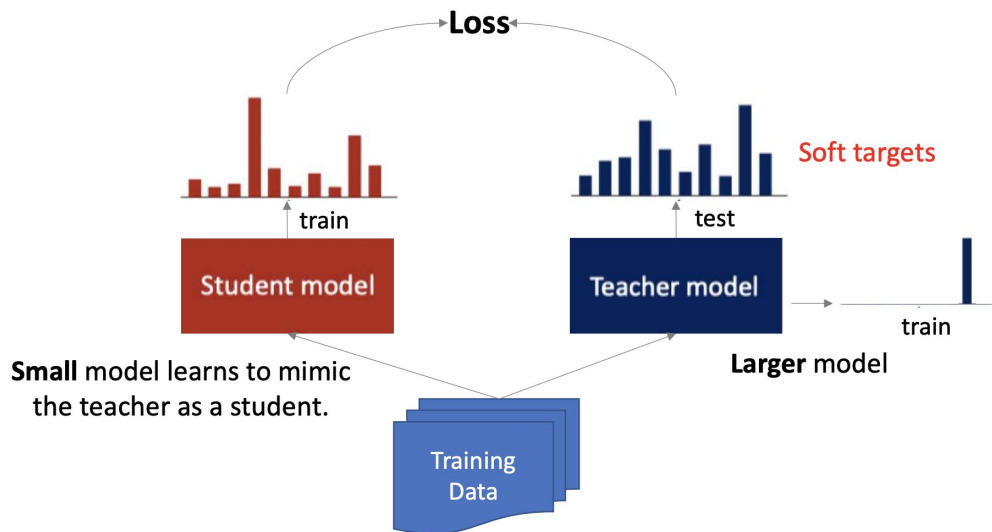
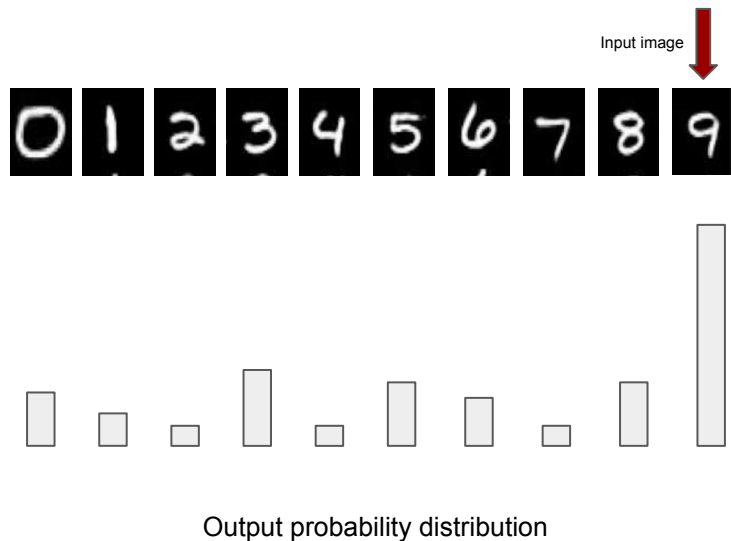


# Knowledge Distillation via Soft Targets



# Basic concept

Idea: Train less redundant model by using outputs from big models

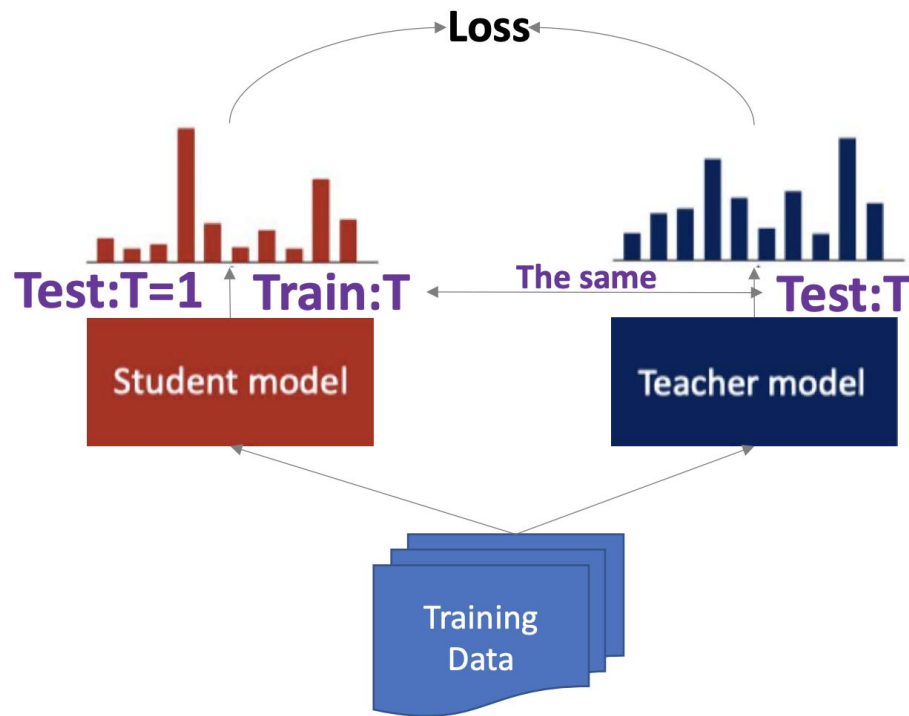
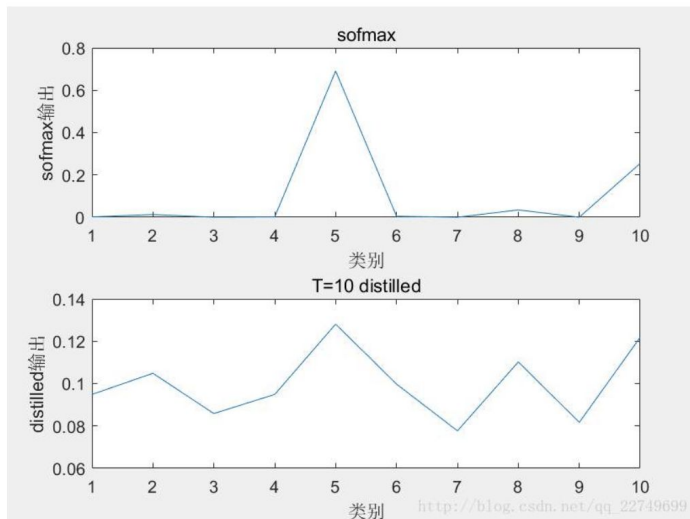


# Softmax with Temperature

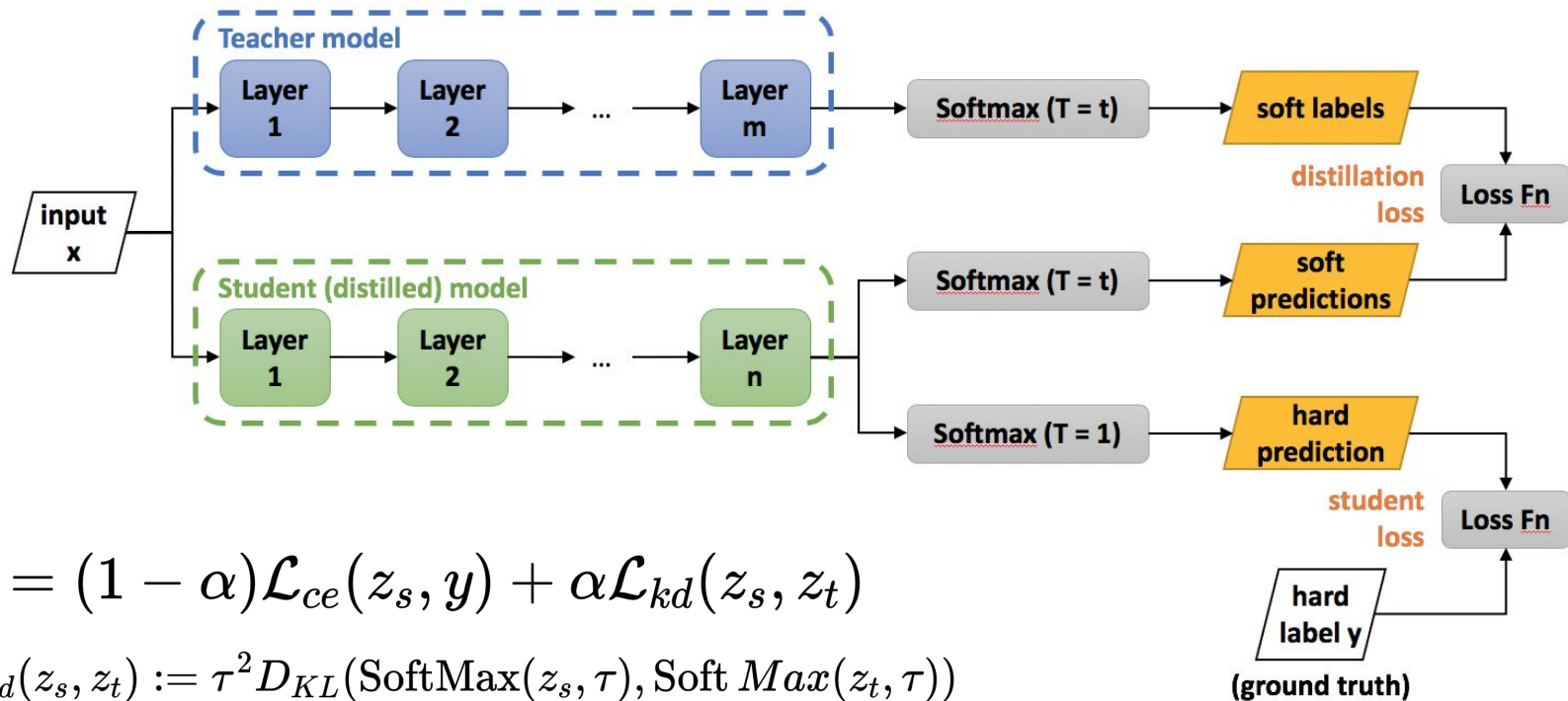
$$q_i = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}$$

Logits

Temperature



# Full Method



$$\mathcal{L} = (1 - \alpha)\mathcal{L}_{ce}(z_s, y) + \alpha\mathcal{L}_{kd}(z_s, z_t)$$

$$\mathcal{L}_{kd}(z_s, z_t) := \tau^2 D_{KL}(\text{SoftMax}(z_s, \tau), \text{SoftMax}(z_t, \tau))$$

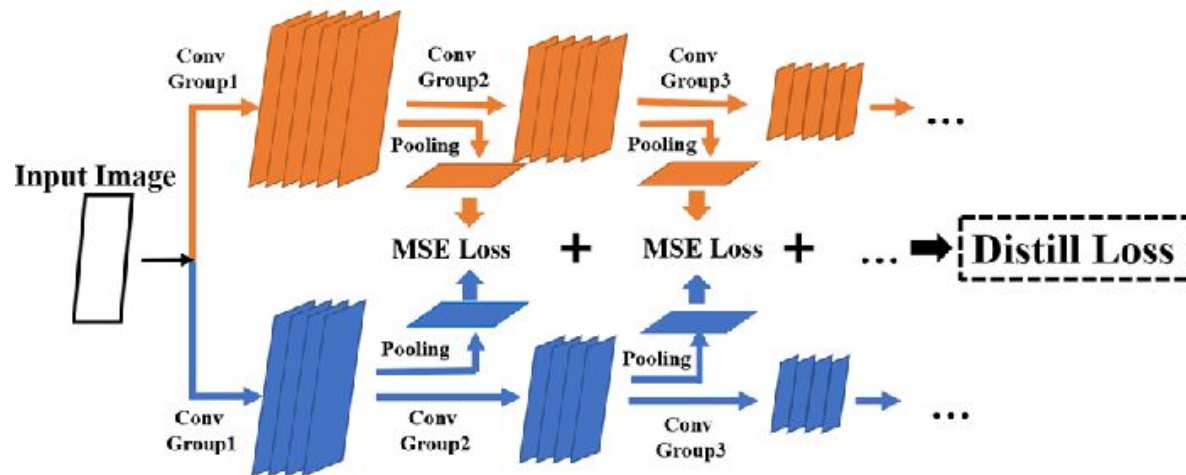
# FitNets: Hits for Thin Deep Nets

# FitNets: Concept

Idea: Use intermediate representations to guide training of student:

$$\mathcal{L}_{HT}(\mathbf{W}_{\text{Guided}}, \mathbf{W}_{\text{r}}) = \frac{1}{2} \|u_h(\mathbf{x}; \mathbf{W}_{\text{Hint}}) - r(v_g(\mathbf{x}; \mathbf{W}_{\text{Guided}}); \mathbf{W}_{\text{r}})\|^2$$

Teacher CNN (Pre-Trained)



Student CNN (Pruned From Teacher CNN)

# FitNets: Results

Algorithm	# params	Accuracy
<i>Compression</i>		
FitNet	~2.5M	<b>91.61%</b>
Teacher	~9M	90.18%
Mimic single	~54M	84.6%
Mimic single	~70M	84.9%
Mimic ensemble	~70M	85.8%
<i>State-of-the-art methods</i>		
Maxout		90.65%
Network in Network		91.2%
Deeply-Supervised Networks		<b>91.78%</b>
Deeply-Supervised Networks (19)		88.2%

Table 1: Accuracy on CIFAR-10

Algorithm	# params	Accuracy
<i>Compression</i>		
FitNet	~2.5M	<b>64.96%</b>
Teacher	~9M	63.54%
<i>State-of-the-art methods</i>		
Maxout		61.43%
Network in Network		64.32%
Deeply-Supervised Networks		<b>65.43%</b>

Table 2: Accuracy on CIFAR-100

# Attention-based Knowledge Distillation

# Attention transfer

- Attention plays a critical role in human visual experience
- It also has demonstrated important role in NNs
- We can use it to improve teacher-student training

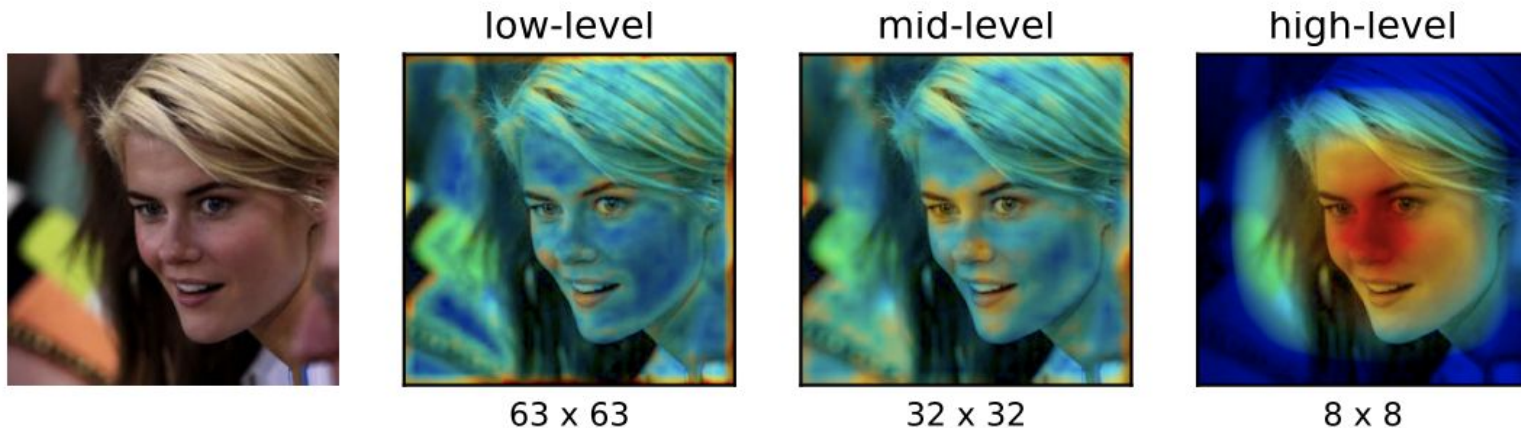


Figure: Sum of absolute values attention maps over different levels of a network trained for face recognition. Mid-level attention maps have higher activation level around eyes, nose and lips, high-level activations correspond to the whole

Credits:

Zagoruyko et al. **Paying More Attention to Attention: Improving the Performance of Convolutional Neural Networks via Attention Transfer**, ICLR 2017



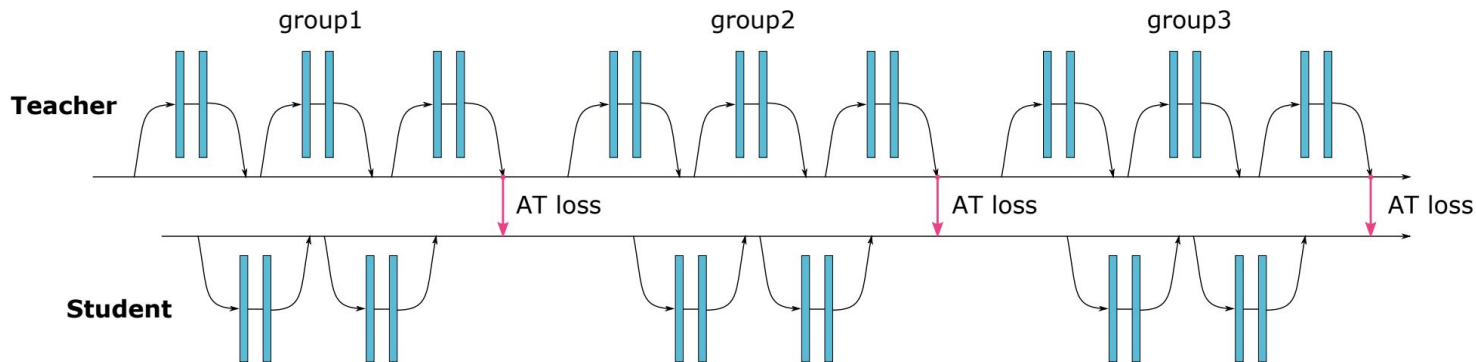
# Attention transfer

Attention loss: 
$$\mathcal{L}_{AT} = \mathcal{L}(\mathbf{W}_S, x) + \frac{\beta}{2} \sum_{j \in \mathcal{I}} \left\| \frac{Q_S^j}{\|Q_S^j\|_2} - \frac{Q_T^j}{\|Q_T^j\|_2} \right\|_p$$

where  $Q_S^j = \text{vec}(F(A_S^j))$  and  $Q_T^j = \text{vec}(F(A_T^j))$  are vectorized j-th pair of teacher - student attention maps

The work considers several types of attention:

- sum of absolute values:  $F_{\text{sum}}(A) = \sum_{i=1}^C |A_i|$
- sum of absolute values raised to the power of  $p$  (where  $p > 1$ ):  $F_{\text{sum}}^p(A) = \sum_{i=1}^C |A_i|^p$
- max of absolute values raised to the power of  $p$  (where  $p > 1$ ):  $F_{\text{max}}^p(A) = \max_{i=1, C} |A_i|^p$



# Methods comparison

KD between similar architectures (CIFAR100):

Teacher	wrn40-2	wrn40-2	resnet56	resnet32×4	vgg13
Student	wrn16-2	wrn40-1	resnet20	resnet8×4	vgg8
Teacher	76.46	76.46	73.44	79.63	75.38
Student	73.64	72.24	69.63	72.51	70.68
KD [16]	74.92	73.54	70.66	73.33	72.98
FitNet [37]	75.75	74.12	71.60	74.31	73.54
AT [46]	75.28	74.45	<b>71.78</b>	74.26	73.62

KD between different architectures:

Teacher	vgg13	ResNet50	ResNet50	resnet32×4	resnet32×4	wrn40-2
Student	MobileNetV2	MobileNetV2	vgg8	ShuffleV1	ShuffleV2	ShuffleV1
Teacher	75.38	79.10	79.10	79.63	79.63	76.46
Student	65.79	65.79	70.68	70.77	73.12	70.77
KD [16]	67.37	67.35	73.81	74.07	74.45	74.83
FitNet [37]	68.58	68.54	73.84	74.82	75.11	75.55
AT [46]	<u>69.34</u>	69.28	73.45	74.76	75.30	75.61

# Knowledge Distillation Types

- Response-Based Knowledge

$$L_{ResD}(p(z_t, T), p(z_s, T)) = \mathcal{L}_R(p(z_t, T), p(z_s, T))$$

$$p(z_i, T) = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}$$

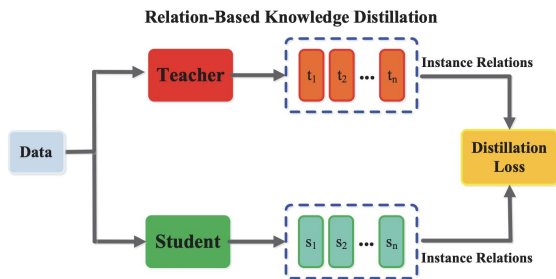
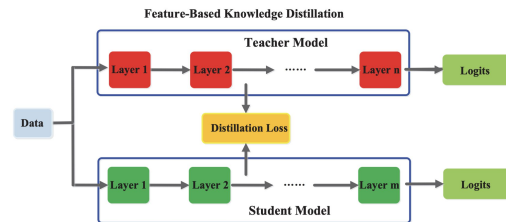
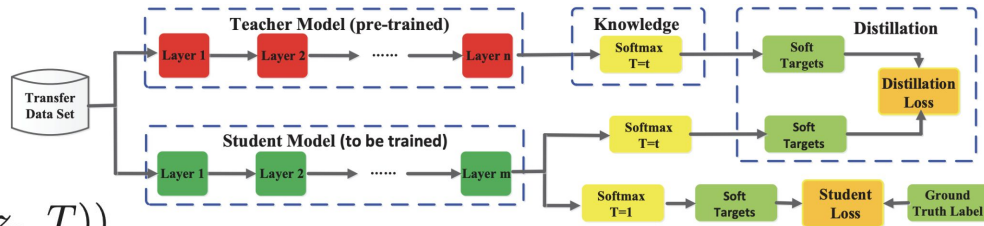
- Feature-Based Knowledge

$$L_{FeaD}(f_t(x), f_s(x)) = \mathcal{L}_F(\Phi_t(f_t(x)), \Phi_s(f_s(x)))$$

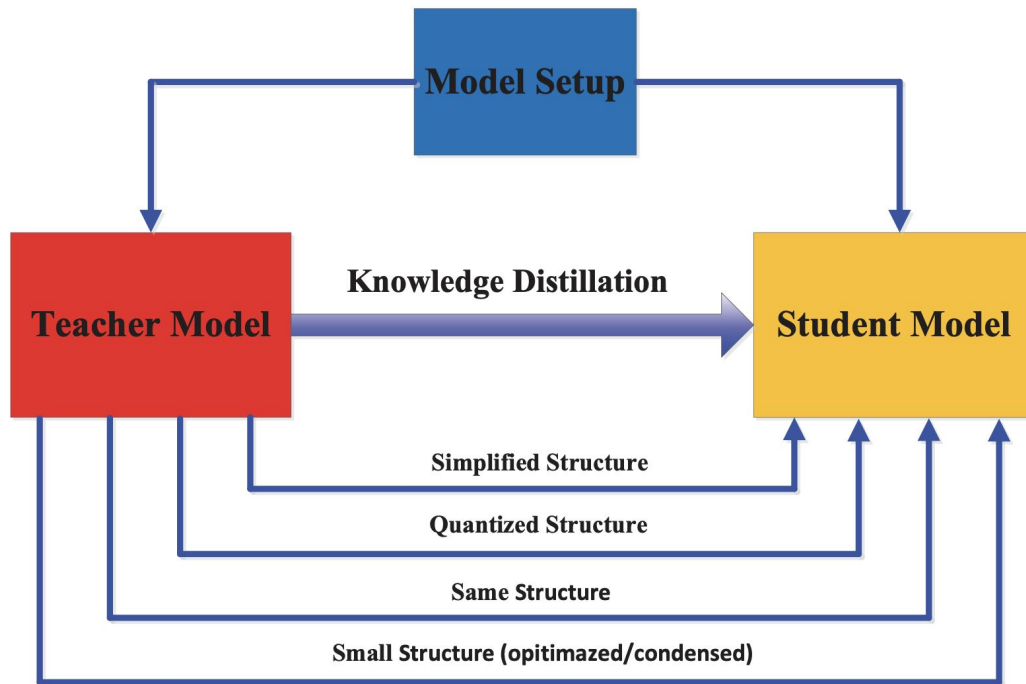
- Relation-Based Knowledge

$$L_{RelD}(f_t, f_s) = \mathcal{L}_{R^1}(\Psi_t(\hat{f}_t, \check{f}_t), \Psi_s(\hat{f}_s, \check{f}_s))$$

$$L_{RelD}(F_t, F_s) = \mathcal{L}_{R^2}(\psi_t(t_i, t_j), \psi_s(s_i, s_j))$$



# Student Setup



# Summary

**Knowledge Distillation** is a powerful technique to train small neural network using information from big pretrained network

Pros:

- Improves model performance

Cons:

- Capacity gap
- Controversy of the method

# Sources

1. Hinton et al. **Distilling the Knowledge in a Neural Network**, NIPS 2014
2. Romero et al. **FitNets: Hints for Thin Deep Nets**, ICLR 2015
3. Zagoruyko et al. **Paying More Attention to Attention: Improving the Performance of Convolutional Neural Networks via Attention Transfer**, ICLR 2017
4. <https://devopedia.org/knowledge-distillation>
5. <https://medium.com/analytics-vidhya/knowledge-distillation-in-a-deep-neural-network-c9dd59aff89b>
6. <https://neptune.ai/blog/knowledge-distillation>