



# Artificial Neural Networks

[2500WETANN]

José Oramas

# Course Announcements

- **Papers for the Research Assignment have been distributed**

Deadline: 05 / 06 / 2022

Let me know in case there are issues with the assigned paper

- **Time and Place to be Confirmed**

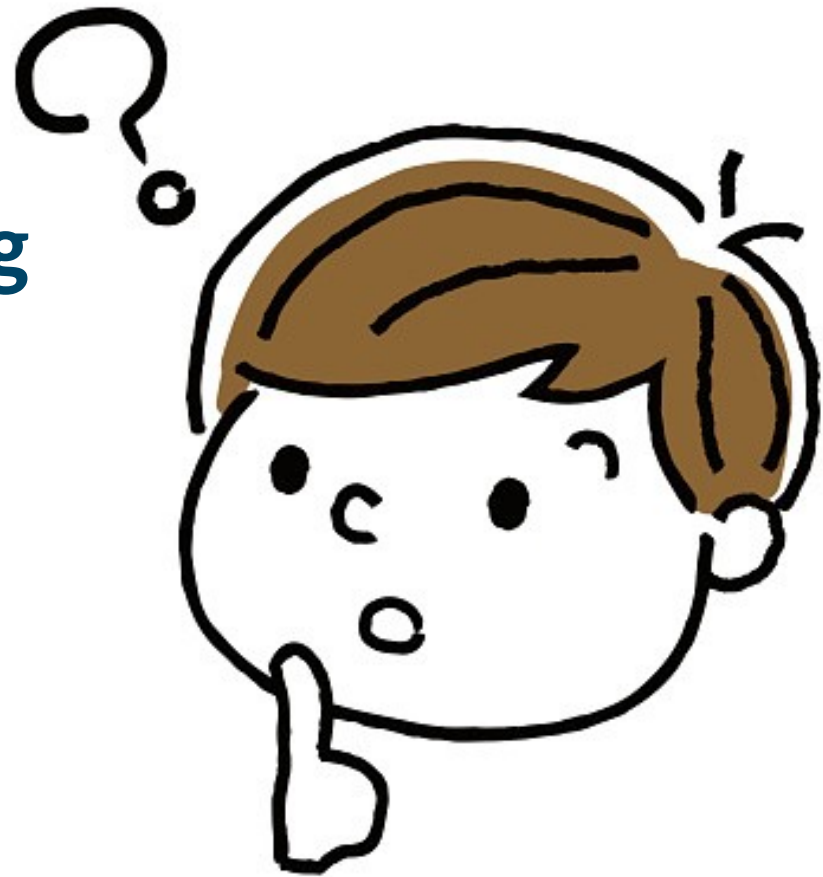


# Transfer Learning

[ Reusing Existing Models ]

José Oramas

Based on your experience...  
**What problems did you  
encounter when training  
models?**



Based on your experience...  
**What problems did you  
encounter when training  
models?**

**Today:** Existing models → New models



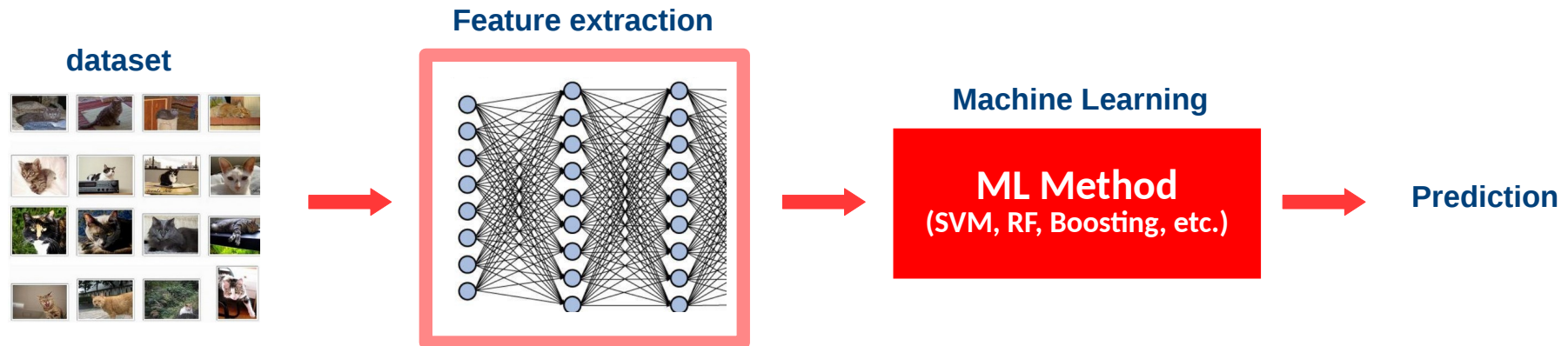
# Models as Feature Extractors

[ while we wait for GPUs to get better ]

# Transfer Learning

## 1- Pre-trained Models as Feature Extractors

- Push data through the model
- Collect activations in a given layer
- Use activations as input features in a classical ML method

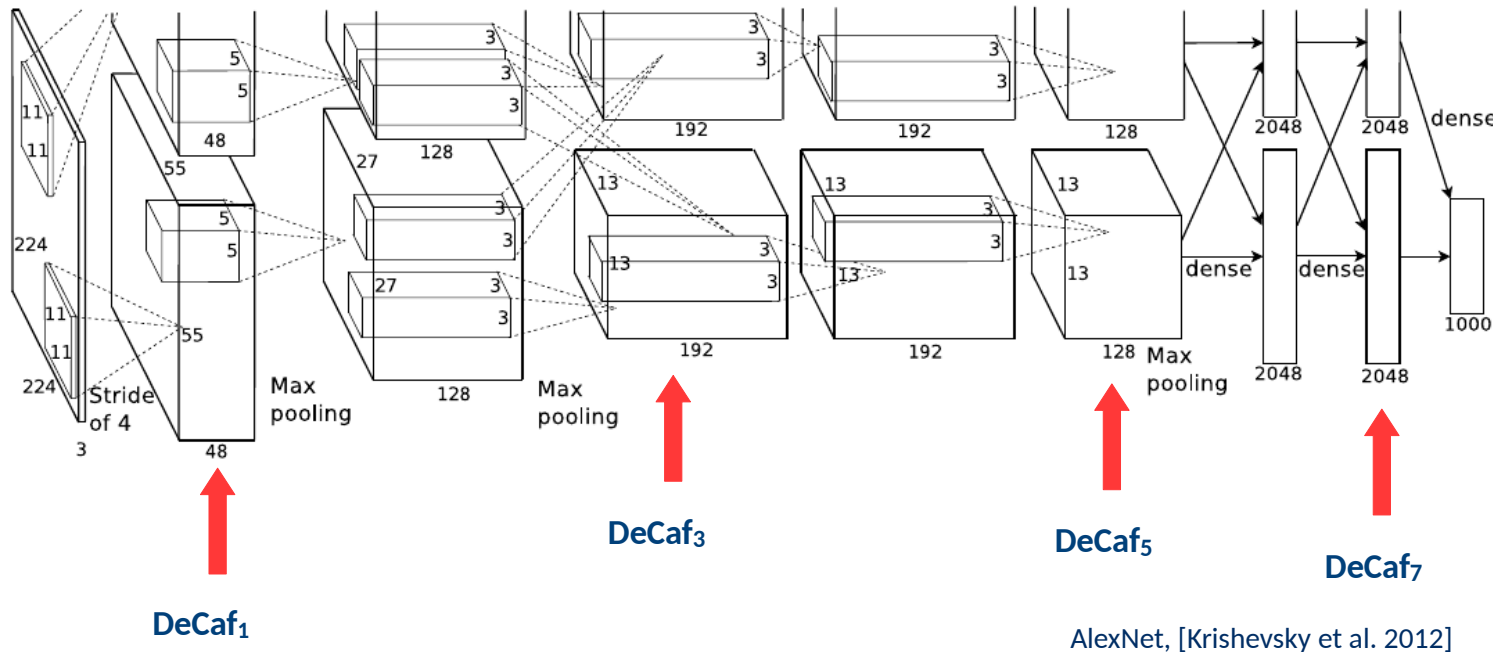




# Transfer Learning

## 1- Pre-trained Models as Feature Extractors

- **Example: DeCaf** [Donahue et al., ICML'14]



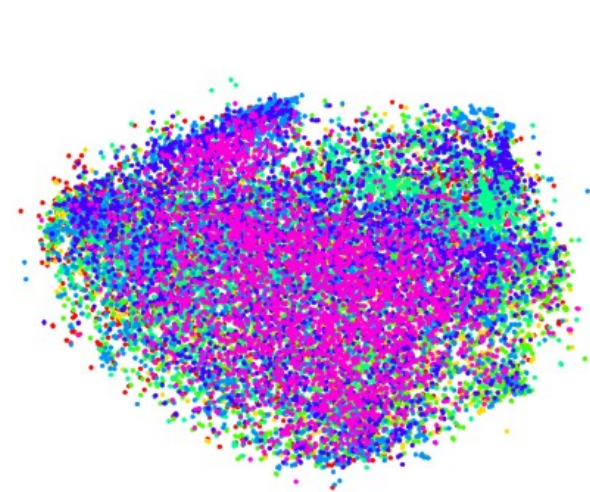
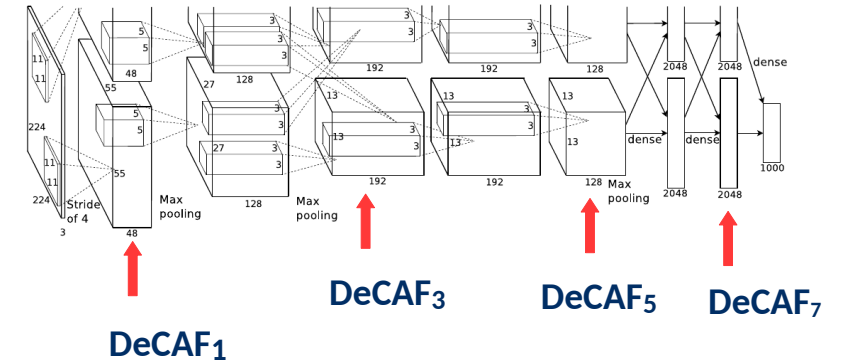
- Take a pre-trained model [AlexNet trained on ImageNet]
- Extract activations from given parts of the network. [DeCaf<sub>i</sub>]
- Train a standard ML methods based on the extracted activations.



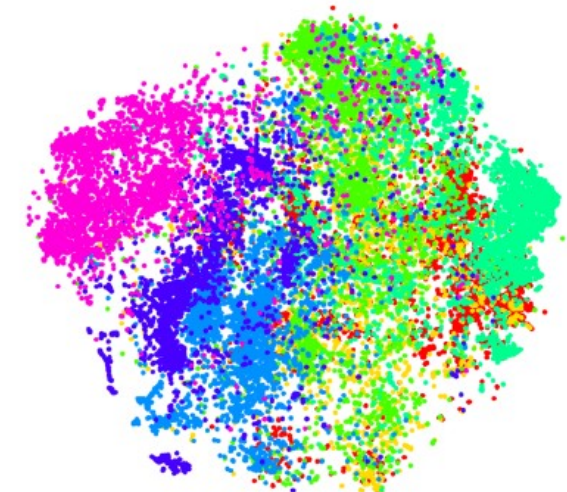
# Transfer Learning

## 1- Pre-trained Models as Feature Extractors

- DeCaf – Feature Visualization



(c) DeCAF<sub>1</sub>



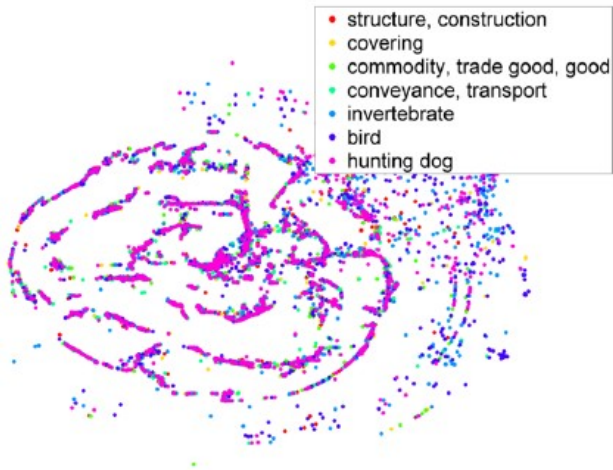
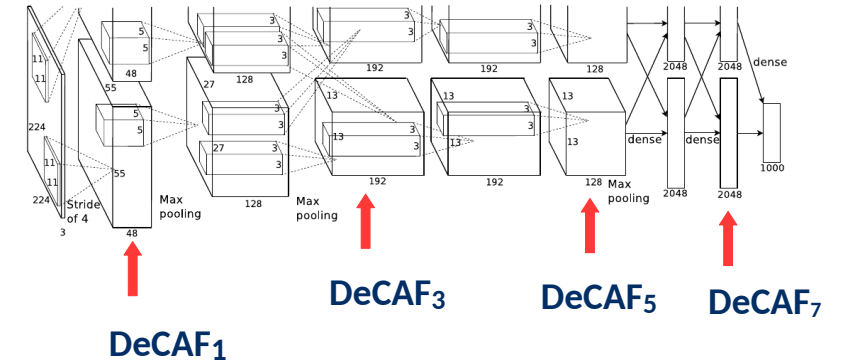
(d) DeCAF<sub>6</sub>

- Visualizing the representation via t-SNE [van der Maaten et al., 2008]

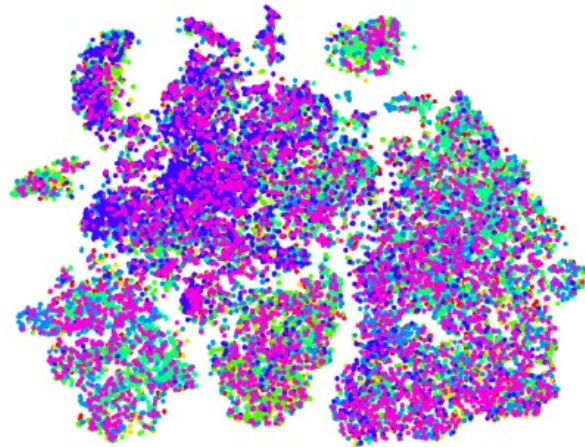
# Transfer Learning

## 1- Pre-trained Models as Feature Extractors

### – DeCaf – Feature Visualization



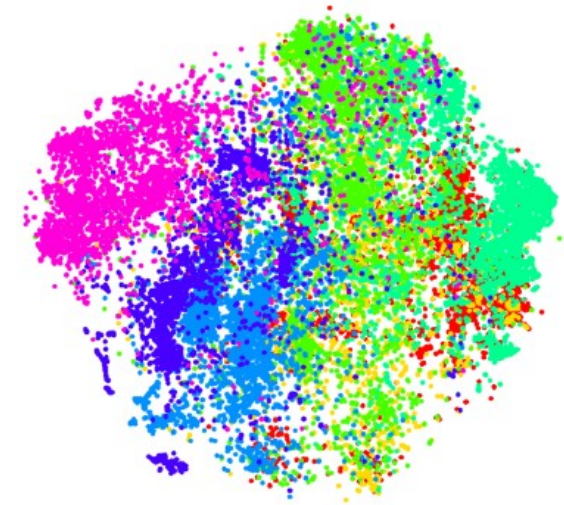
(a) LLC



(b) GIST



(c) DeCAF<sub>1</sub>



(d) DeCAF<sub>6</sub>

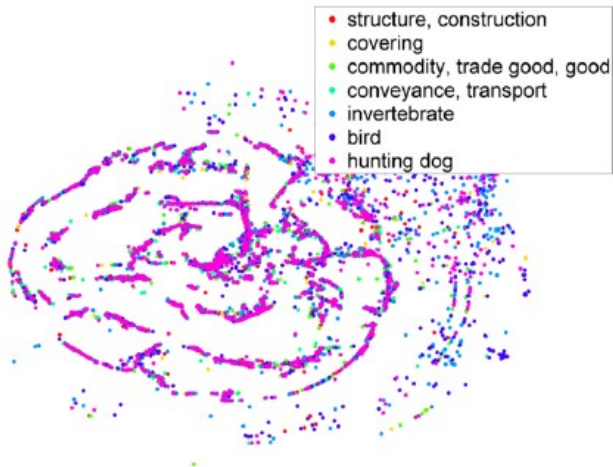
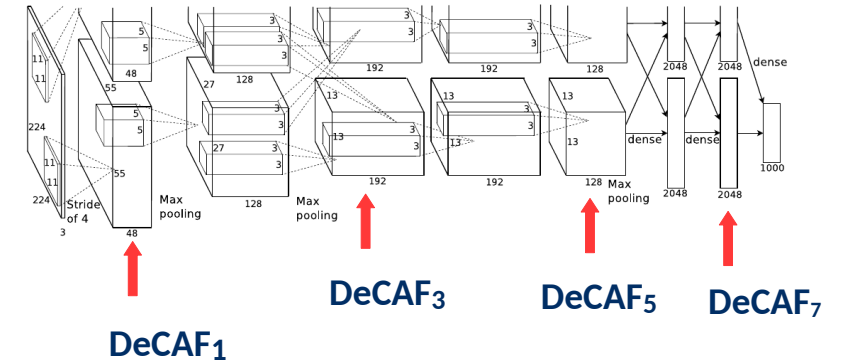
- Visualizing the representation via t-SNE [van der Maaten et al., 2008]



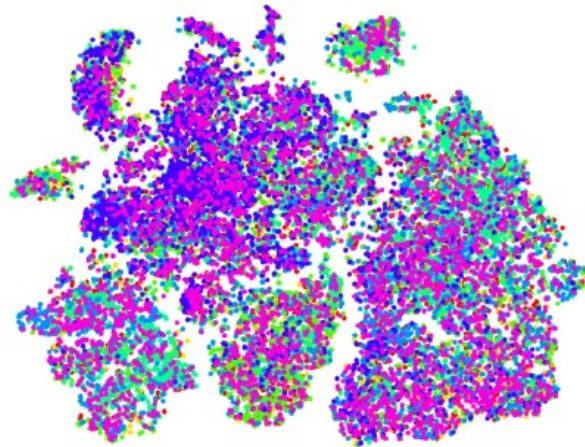
# Transfer Learning

## 1- Pre-trained Models as Feature Extractors

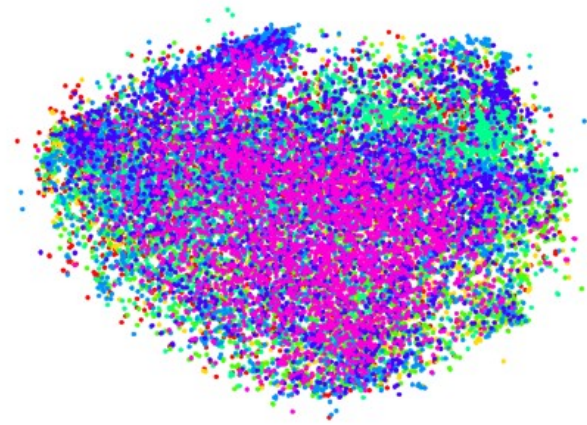
### – DeCaf – Feature Visualization



(a) LLC



(b) GIST



(c) DeCAF<sub>1</sub>



(d) DeCAF<sub>6</sub>

- Visualizing the representation via t-SNE [van der Maaten et al., 2008]

**Q: Why DeCAF<sub>1</sub> and DeCAF<sub>6</sub> look so different?**

# Transfer Learning

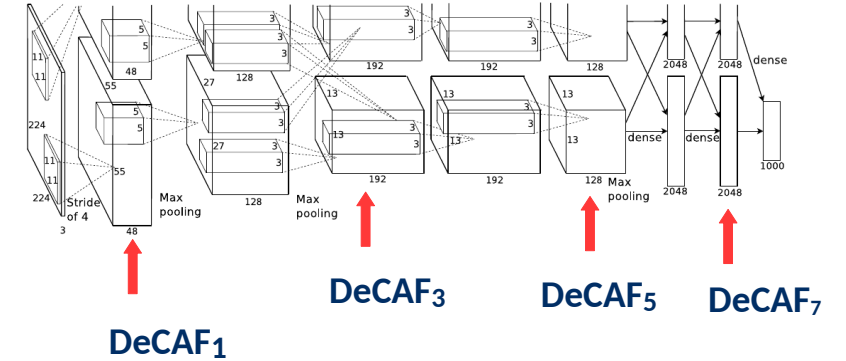
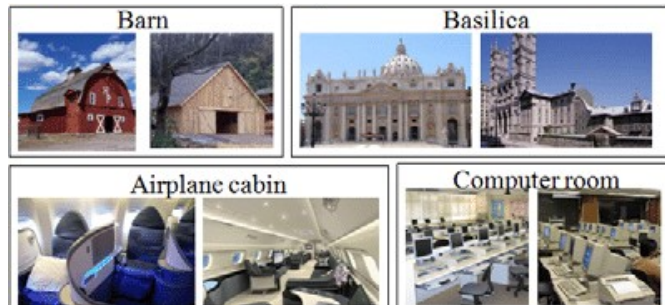
## 1- Pre-trained Models as Feature Extractors

### – DeCaf – Feature Generalization

ILSVRC'12



SUN-397

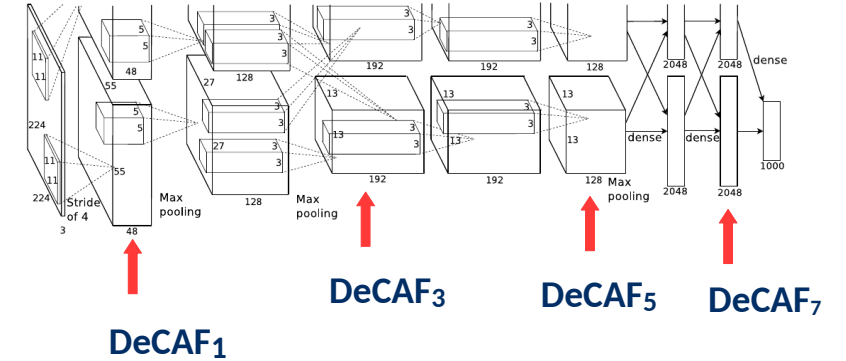




# Transfer Learning

## 1- Pre-trained Models as Feature Extractors

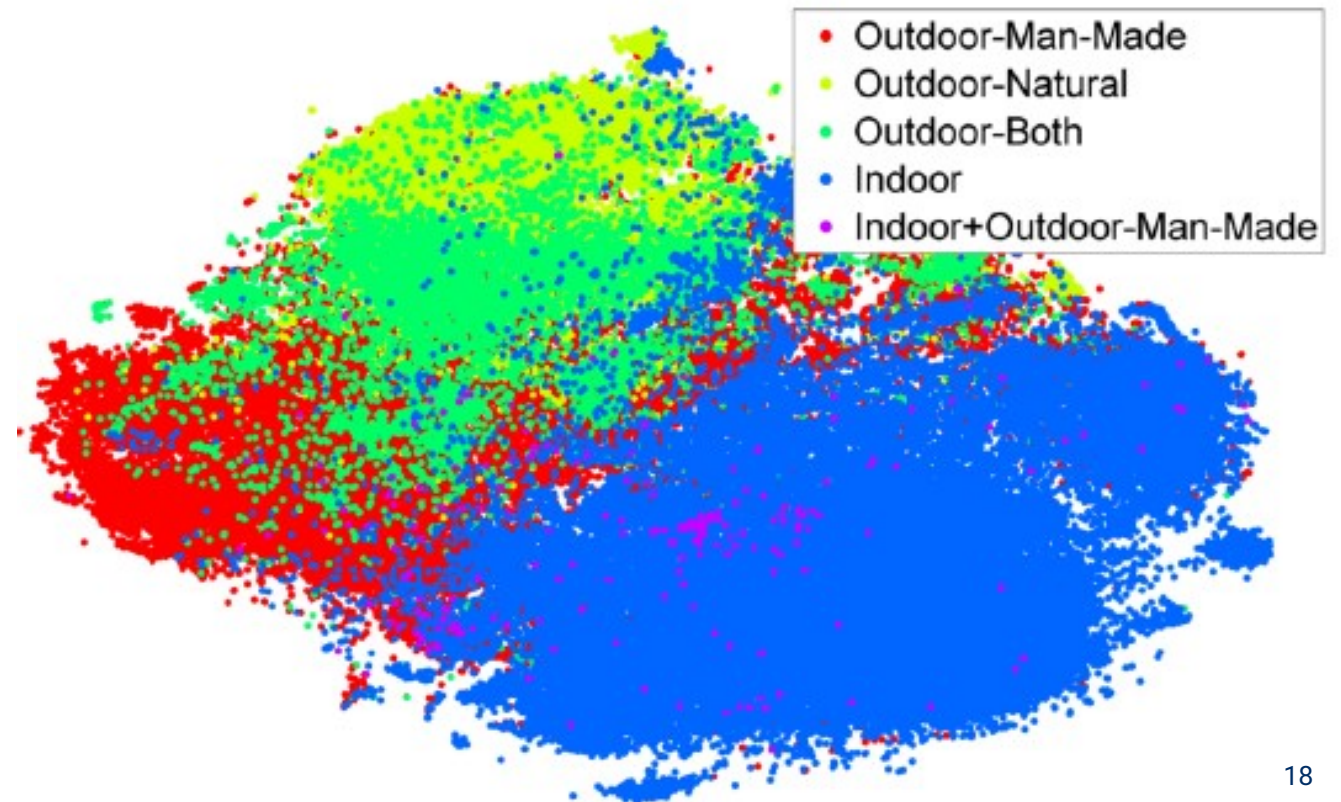
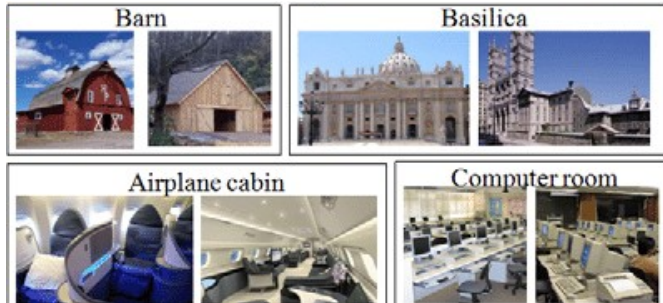
### – DeCaf – Feature Generalization



ILSVRC'12



SUN-397

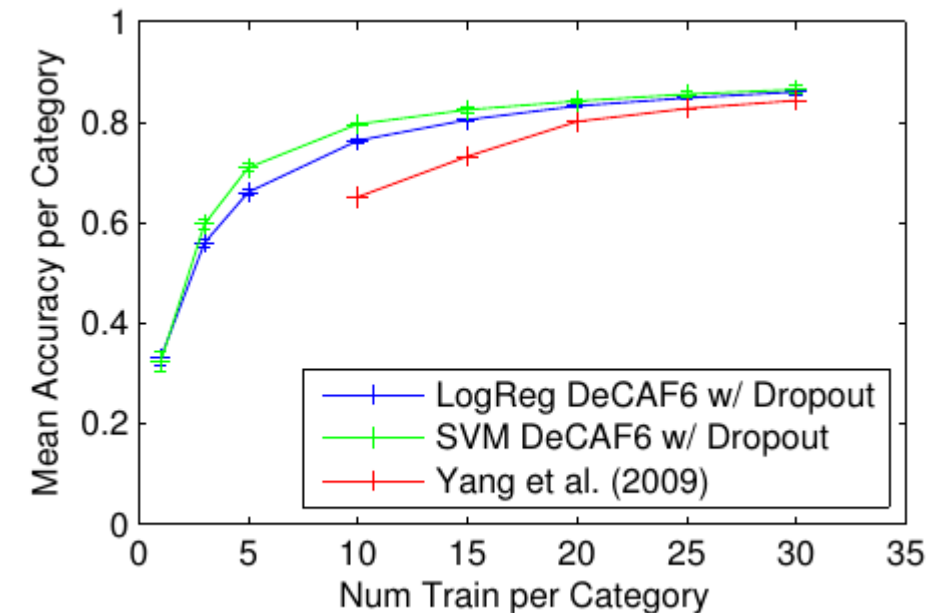
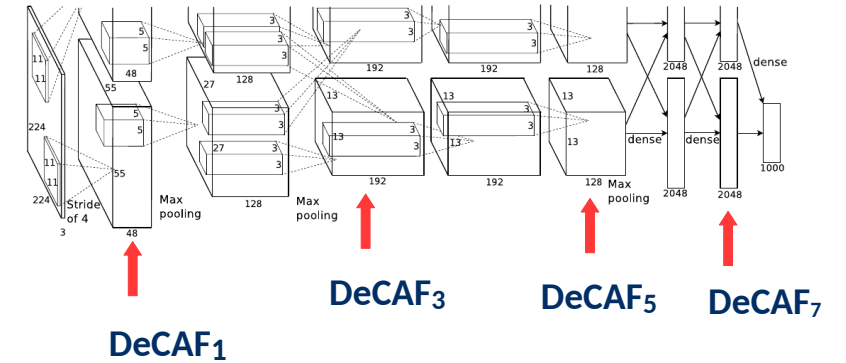


# Transfer Learning

## 1- Pre-trained Models as Feature Extractors

### – DeCaf – Quantitative Analysis

	DeCAF <sub>5</sub>	DeCAF <sub>6</sub>	DeCAF <sub>7</sub>
LogReg	$63.29 \pm 6.6$	$84.30 \pm 1.6$	$84.87 \pm 0.6$
LogReg with Dropout	-	$86.08 \pm 0.8$	$85.68 \pm 0.6$
SVM	$77.12 \pm 1.1$	$84.77 \pm 1.2$	$83.24 \pm 1.2$
SVM with Dropout	-	<b><math>86.91 \pm 0.7</math></b>	$85.51 \pm 0.9$
Yang et al. (2009)		84.3	
Jarrett et al. (2009)		65.5	

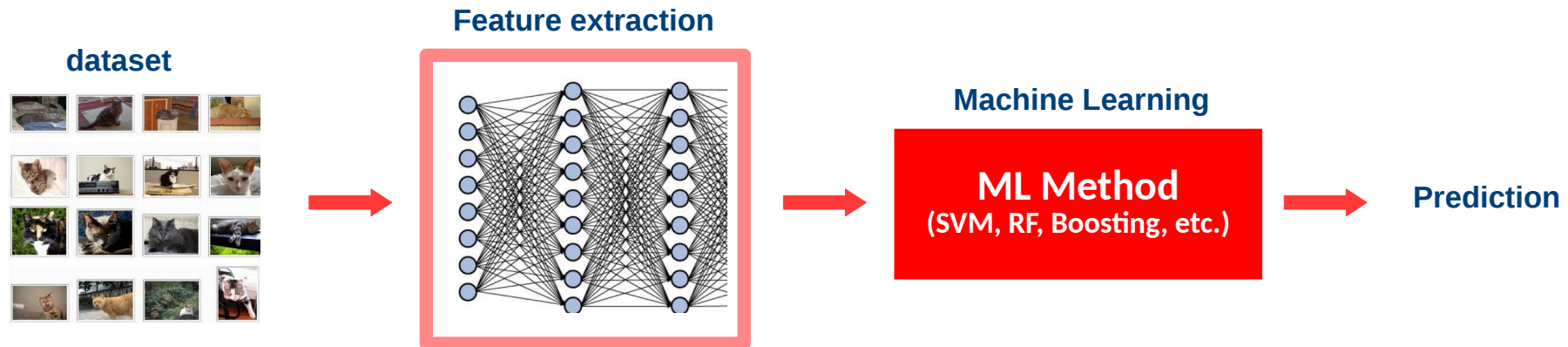


- Classification performance on the Caltech-101 dataset [30 training images/class]

# Transfer Learning

## 1- Pre-trained Models as Feature Extractors

- Push data through the model
- Collect activations in a given layer
- Use activations as input features in a classical ML method





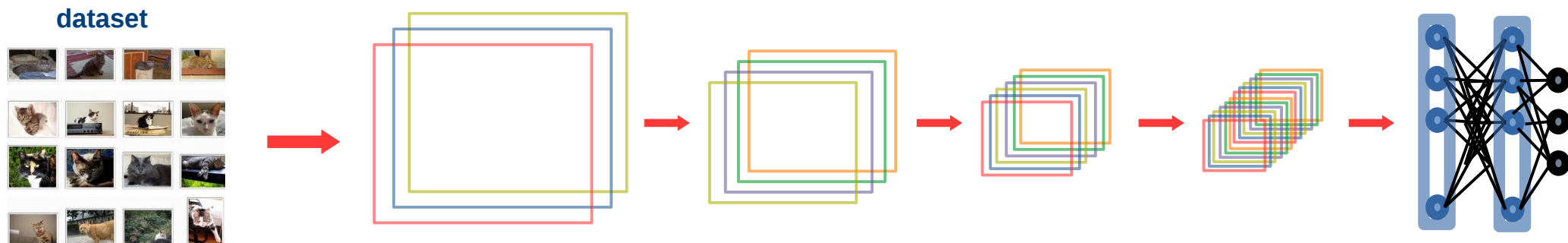
# Adapt an Existing Model

[ the standard practice ]

# Transfer Learning

## 2- Adapting a Pre-trained Model (aka “Fine-Tuning”)

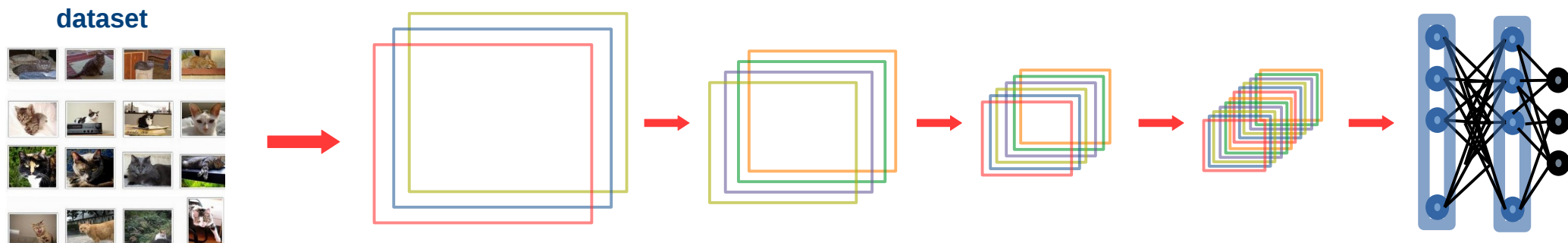
- Adjust Final Layer
- Update the weights of some layers to adapt to new tasks
- Freeze some weights, retrain others



# Transfer Learning

## 2- Adapting a Pre-trained Model (aka “Fine-Tuning”)

- Adjust Final Layer
- Update the weights of some layers to adapt to new tasks
- Freeze some weights, retrain others



# Break

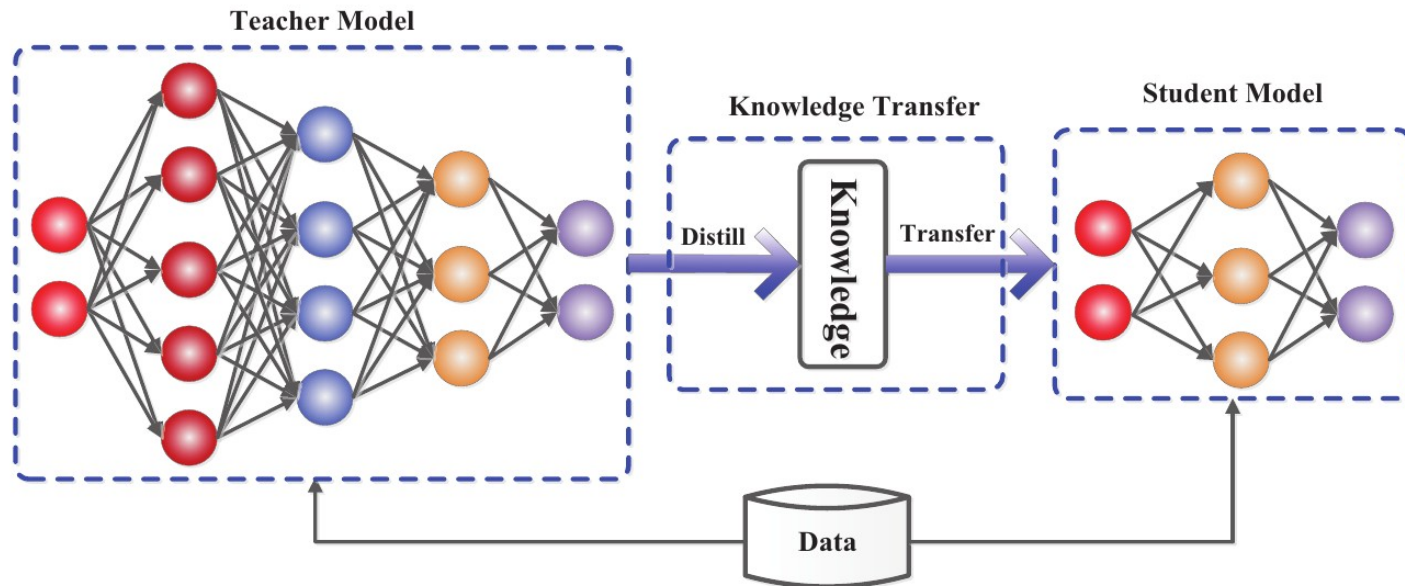
# Extracting Relevant Information

[ so my model can run without a GPU cluster :D ]

# Transfer Learning

## 3- Extracting Relevant Information (aka. “Distillation”)

- Move knowledge from a large model to an optimized one
- Teacher-student Architecture.

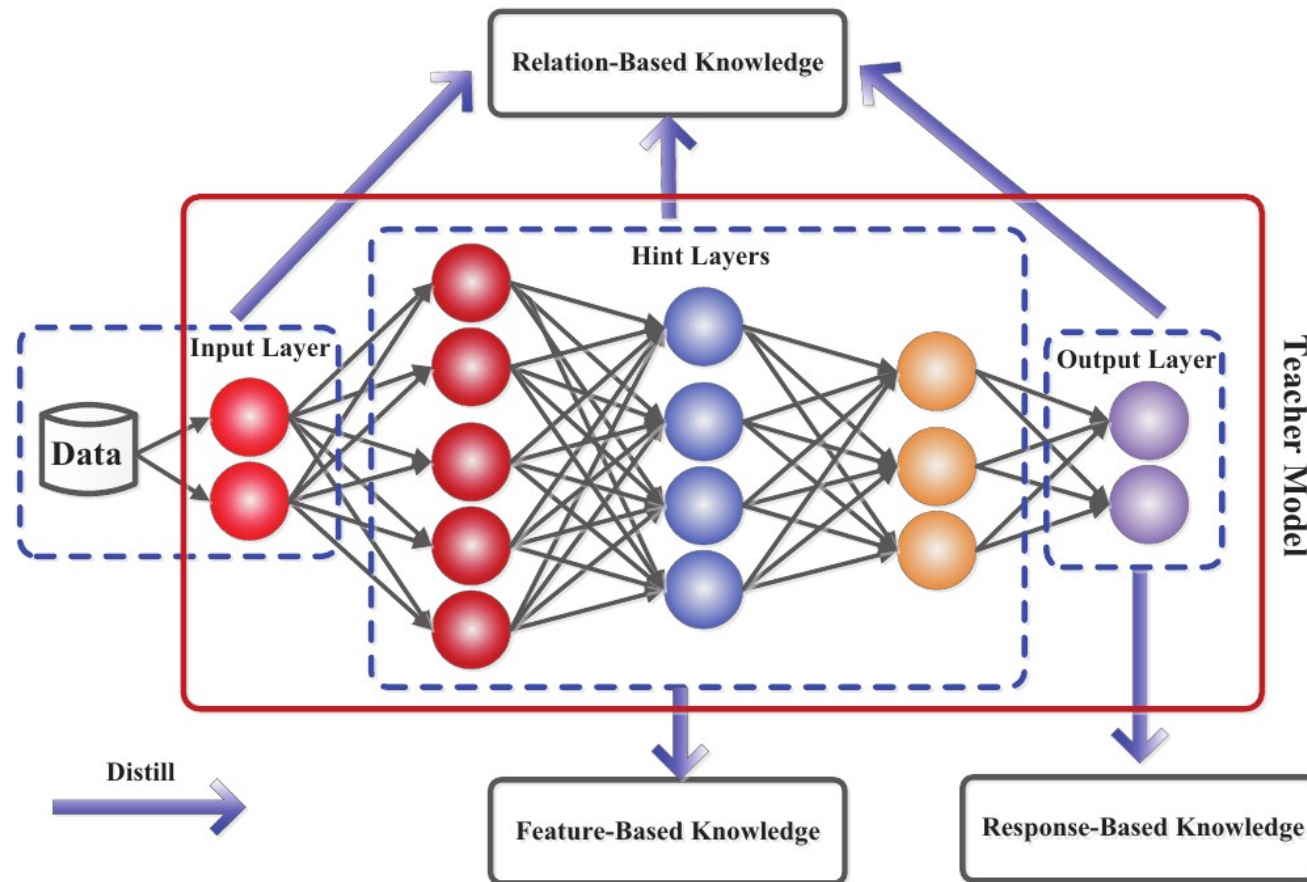


## Directions

- Parameter Pruning and Sharing [Network Quantization]
- Identify Redundant Parameters [ low-rank factorization ]
- Compression of Conv. filters
- Knowledge Distillation

# Transfer Learning

## 3- Extracting Relevant Information (aka. “Distillation”)



### Knowledge Variants

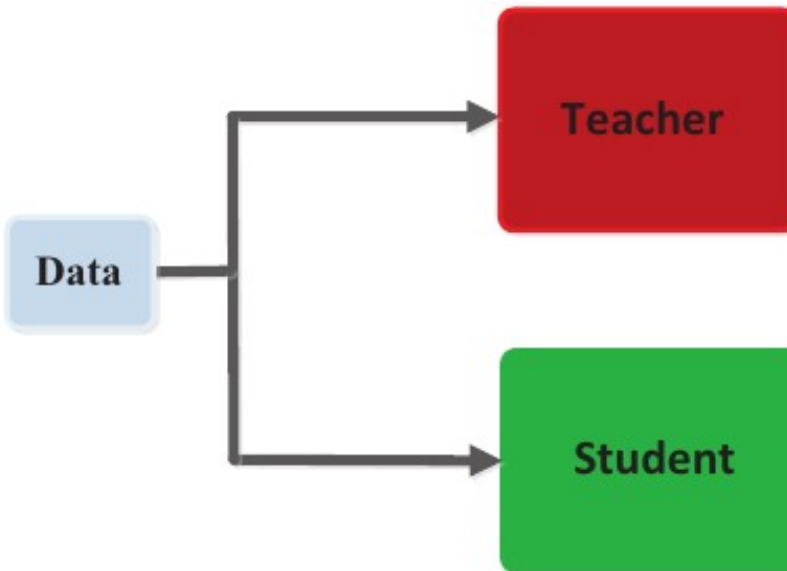
- Responses
  - Mimic the outputs
- Features
  - Mimic outputs and intermediate states
- Relations
  - Model internal relationships



# Transfer Learning

## 3- Extracting Relevant Information (aka. “Distillation”)

- Response-based Knowledge  
→ Mimic the outputs

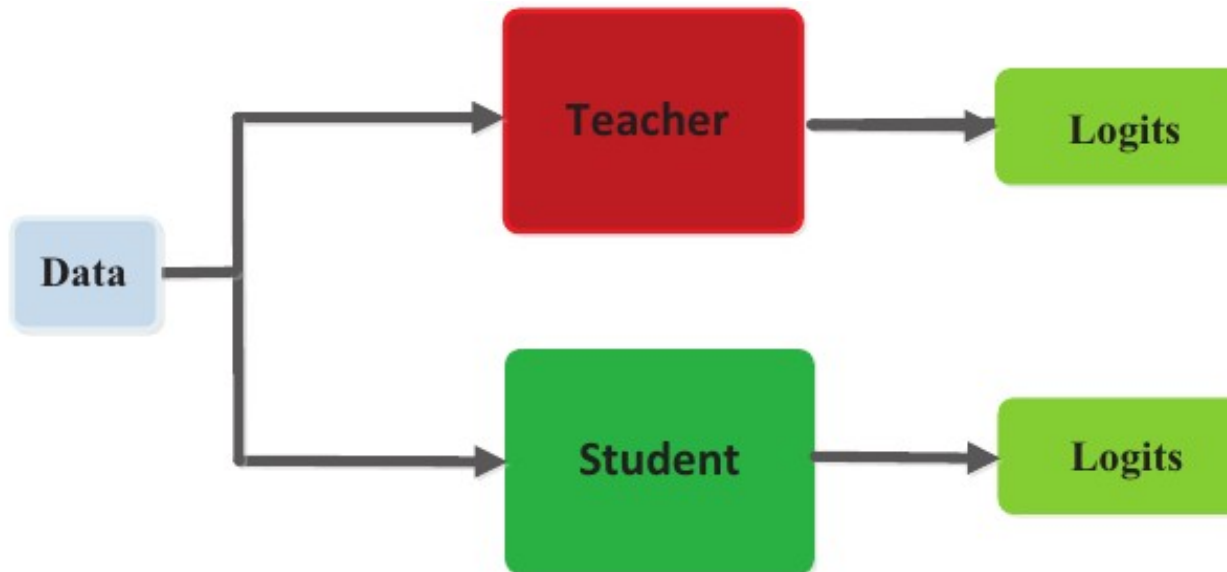


# Transfer Learning

## 3- Extracting Relevant Information (aka. “Distillation”)

- Response-based Knowledge

→ Mimic the outputs

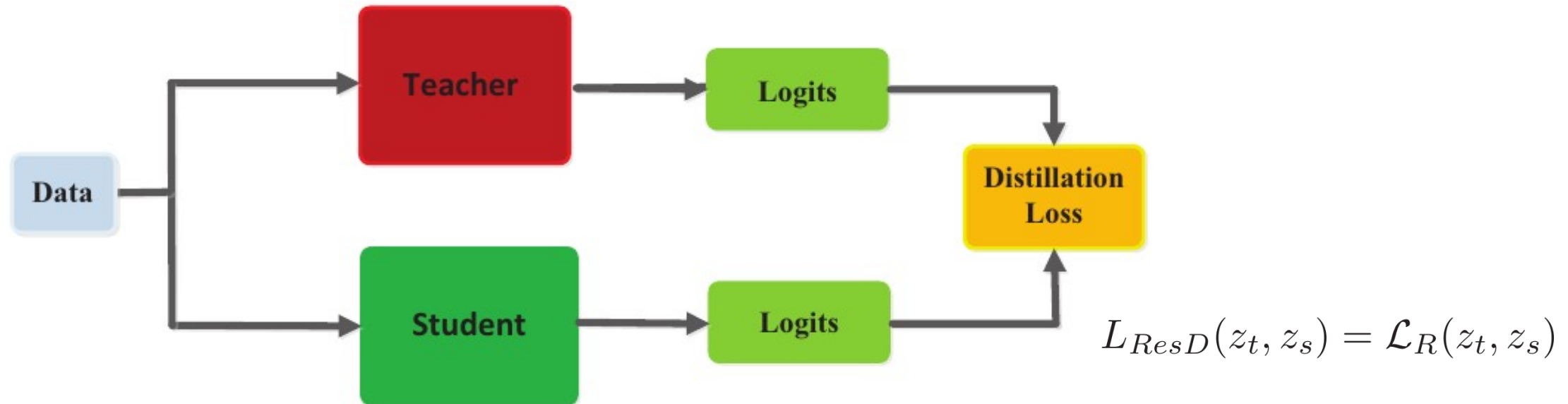


# Transfer Learning

## 3- Extracting Relevant Information (aka. “Distillation”)

- Response-based Knowledge

→ Mimic the outputs

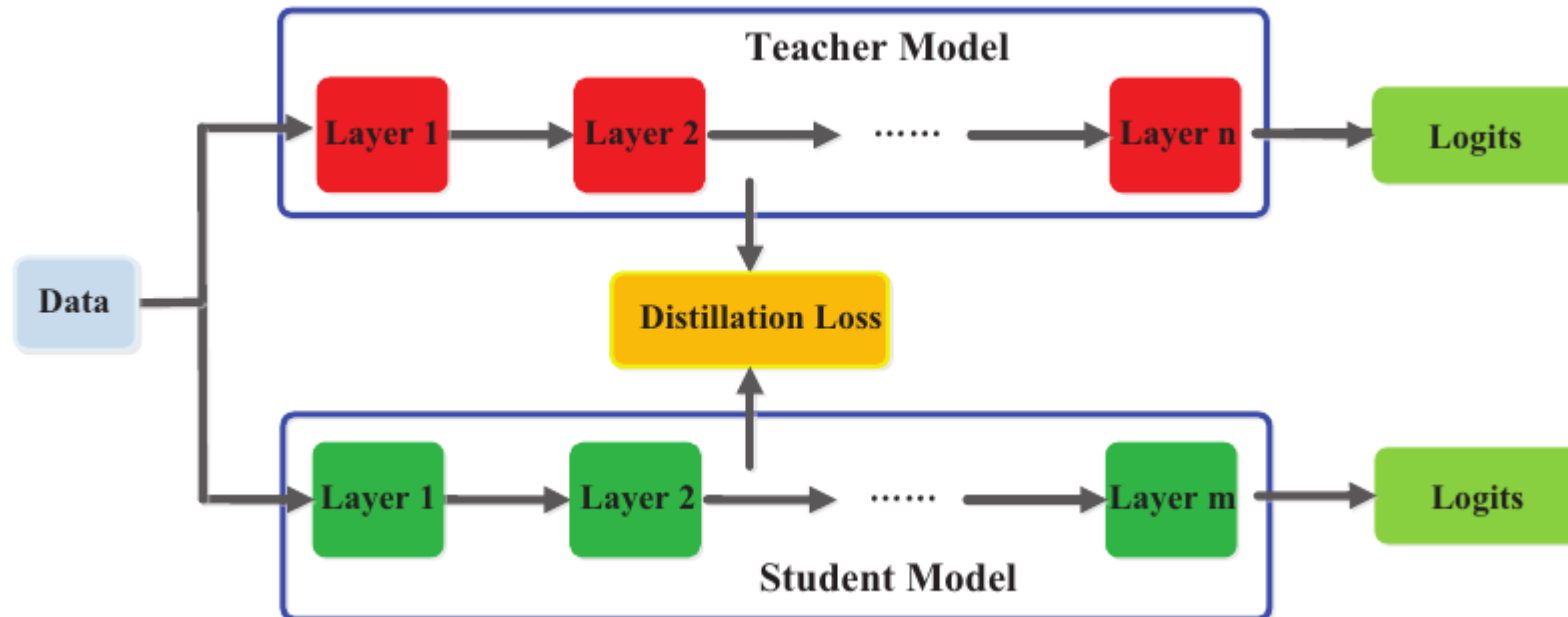


# Transfer Learning

## 3- Extracting Relevant Information (aka. “Distillation”)

- Feature-based Knowledge

→ Mimic the outputs and intermediate states



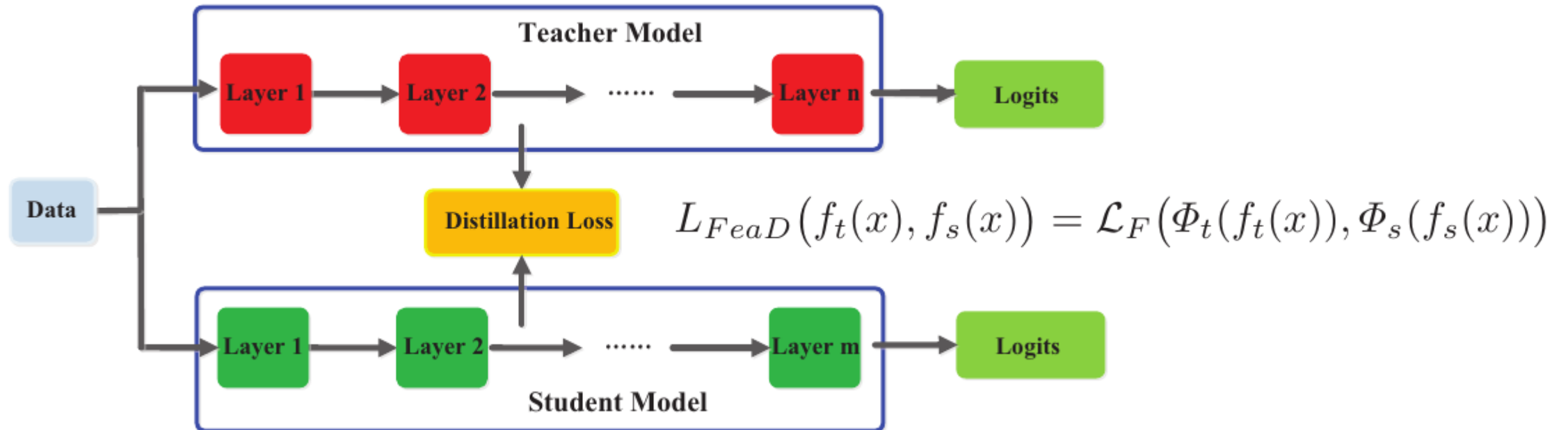
[Gou et al., 2021]

# Transfer Learning

## 3- Extracting Relevant Information (aka. “Distillation”)

- Feature-based Knowledge

→ Mimic the outputs and intermediate states

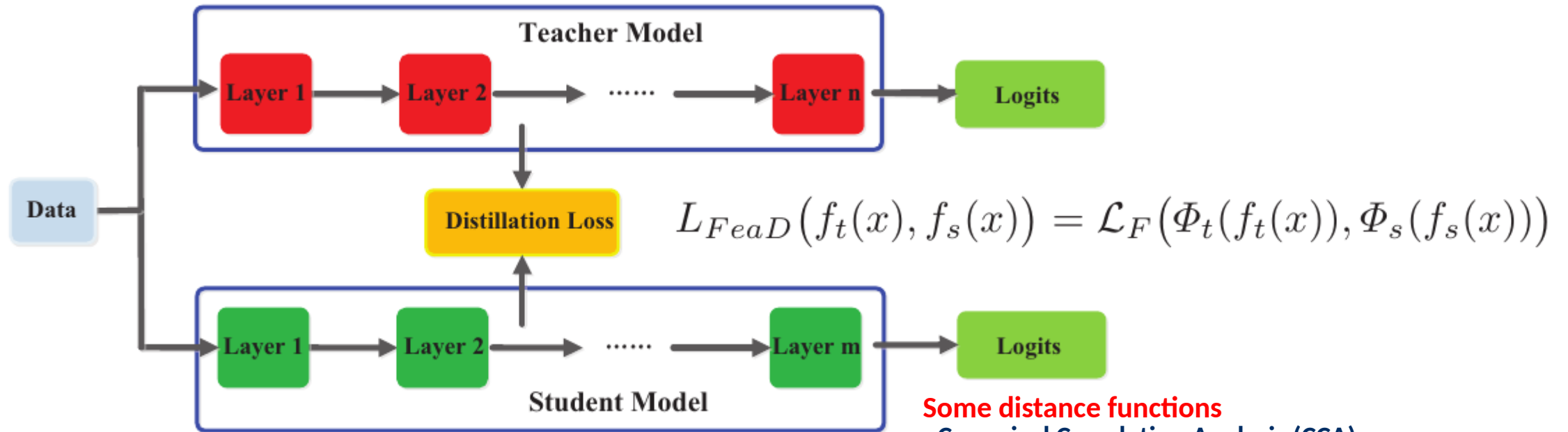


# Transfer Learning

## 3- Extracting Relevant Information (aka. “Distillation”)

### – Feature-based Knowledge

→ Mimic the outputs and intermediate states



[Gou et al., 2021]

### Some distance functions

- Canonical Correlation Analysis (CCA)
- Centered Kernel Alignment (CKA), [Kornblith et al. 2019]
- Orthogonal Procrustes, [Ding et al. 2021]

# Transfer Learning

## 3- Extracting Relevant Information (aka. “Distillation”)

- Relation-based Knowledge

→ Exploit relationships between feature maps or data samples

The diagram illustrates the Relation-based Knowledge Distillation loss function. It features the equation 
$$L_{RelD}(f_t, f_s) = \mathcal{L}_{R^1}(\Psi_t(\hat{f}_t, \check{f}_t), \Psi_s(\hat{f}_s, \check{f}_s))$$
 in the center. Above the equation, the text "Similarity Function" is written in red, with two red arrows pointing down to the  $\Psi_t$  and  $\Psi_s$  terms. Below the equation, the text "Features" is written in red, with four red arrows pointing up to the  $\hat{f}_t$ ,  $\check{f}_t$ ,  $\hat{f}_s$ , and  $\check{f}_s$  terms.

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

Similarity Function

Features

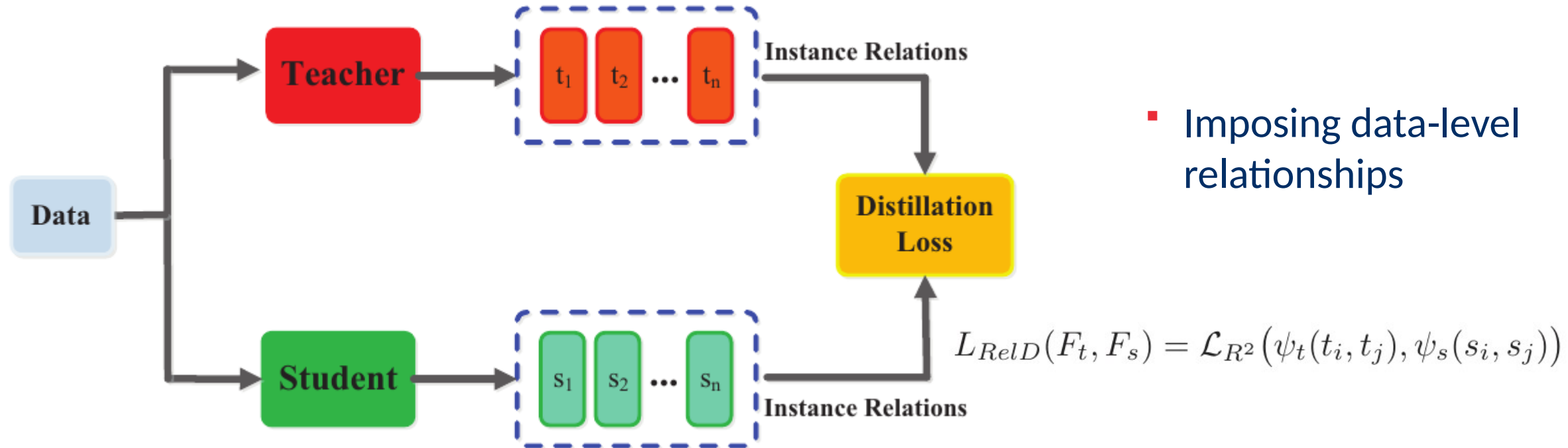


# Transfer Learning

## 3- Extracting Relevant Information (aka. “Distillation”)

- Relation-based Knowledge

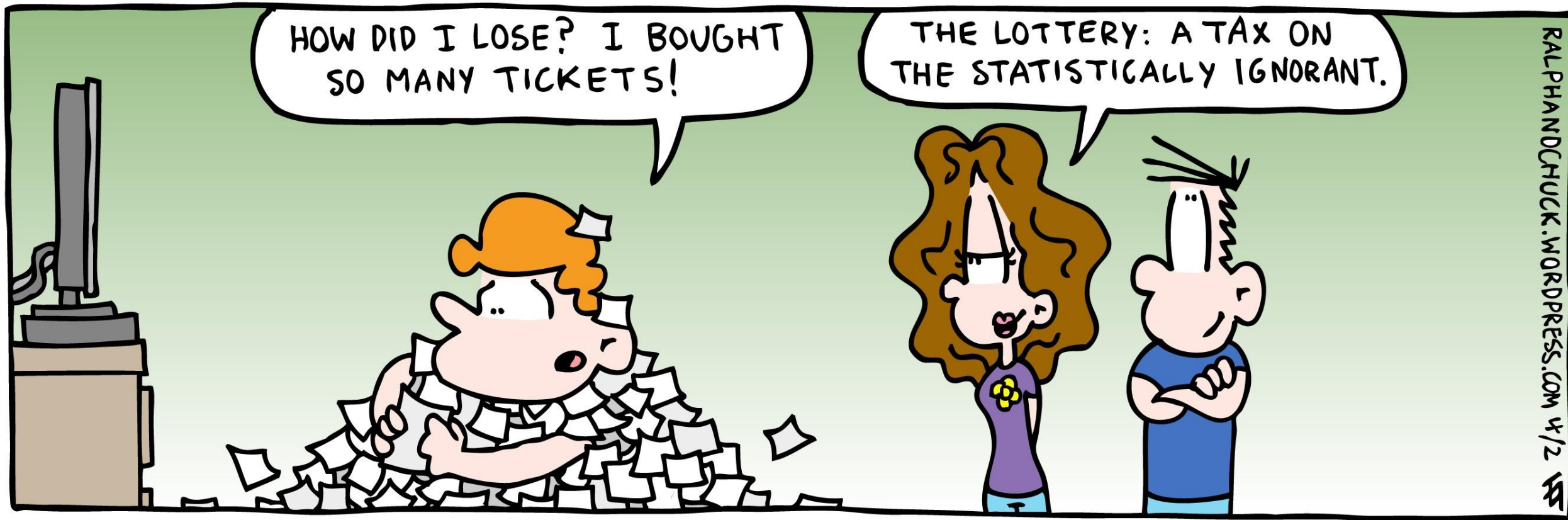
→ Exploit relationships between feature maps or data samples



Yes nice, but...  
**Is it always possible to  
reach a smaller model?**



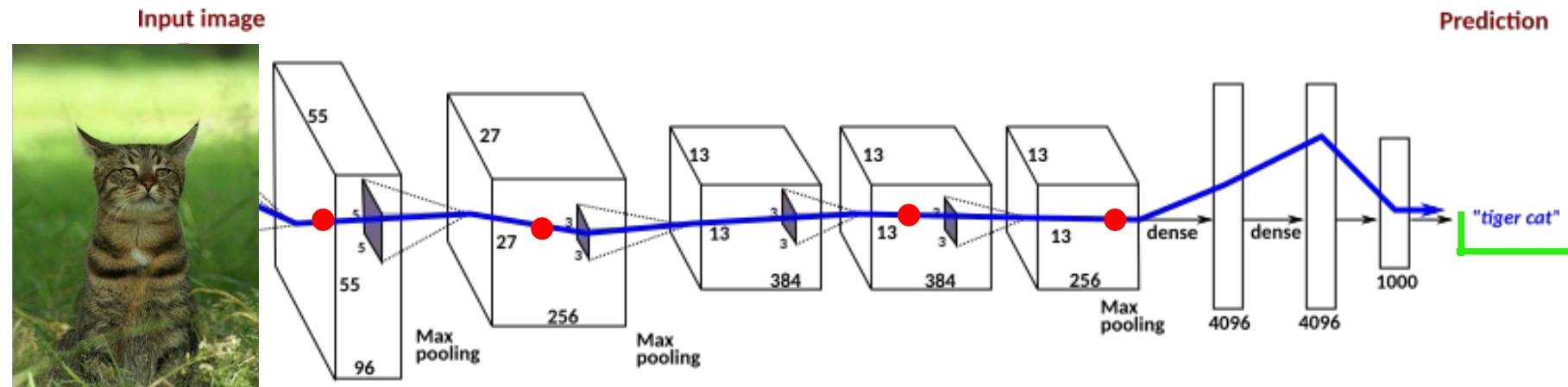
# Lottery Ticket Hypothesis



**The Lottery Ticket Hypothesis.** *Training succeeds for a given network if one of its subnetworks (a “winning ticket”) has been randomly initialized such that it can be trained in isolation to high accuracy in at most the number of iterations necessary to train the original network.*

# Lottery Ticket Hypothesis

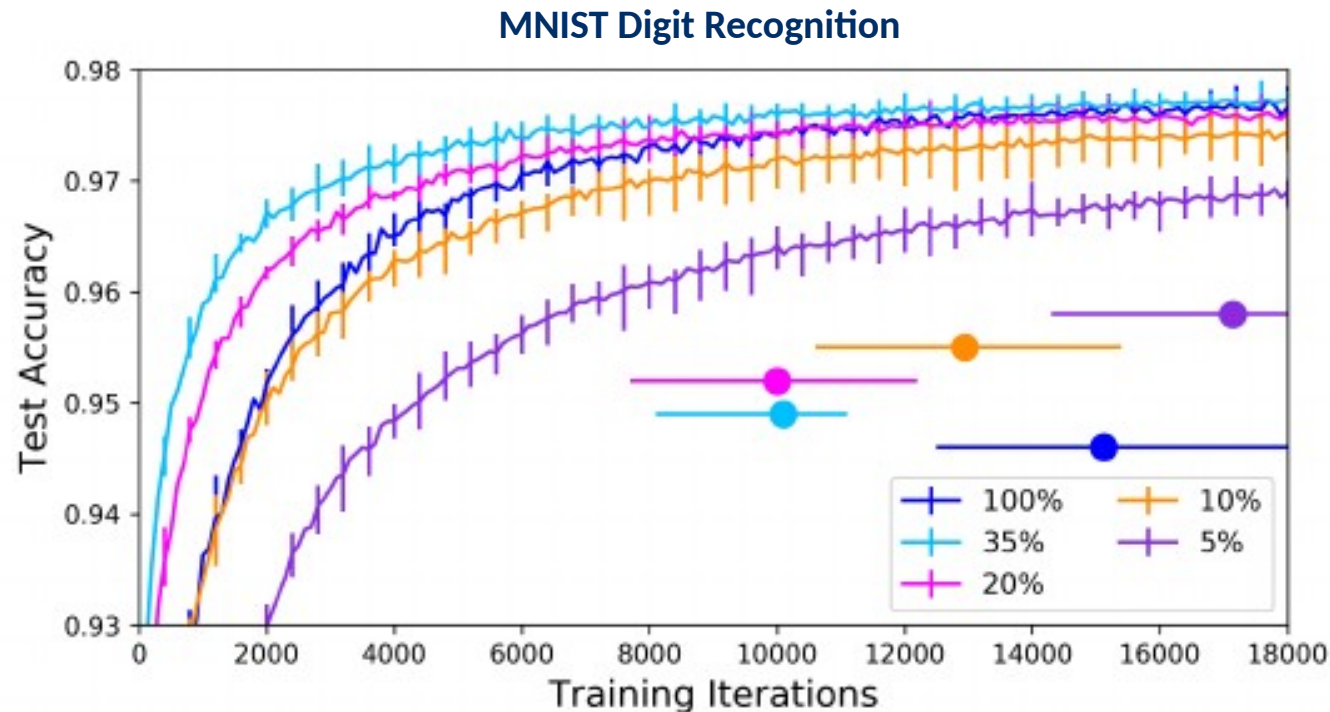
## Fingerprinting Internal Network Activations [Oramas et al., 2019]



- Given a large amount of possible network initializations
- The network successfully trains because a sub-network was initialized properly
- This sub-network can reach similar performance to that of the complete model

# Lottery Ticket Hypothesis

- Given a large amount of possible network initializations
- The network successfully trains because a sub-network was initialized properly
- This sub-network can reach similar performance to that of the complete model



# Lottery Ticket Hypothesis

- Given a large amount of possible network initializations
- The network successfully trains because a sub-network was initialized properly
- This sub-network can reach similar performance to that of the complete model

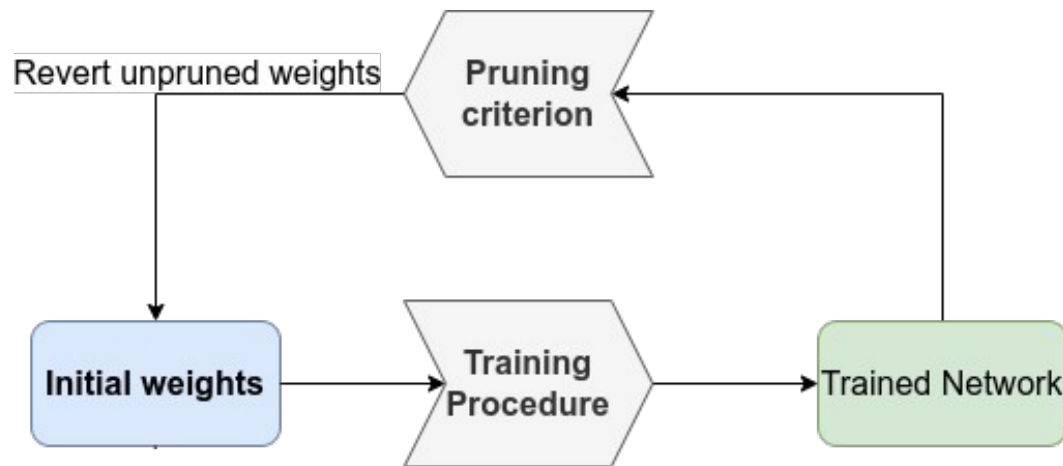


## Finding the Winning Ticket

1. Randomly initialize a neural network.
2. Train the network until it converges.
3. Prune a fraction of the network.
4. Reset the weights of the remaining portion of the network to their values from (1)  
(i.e., the initializations they received before training began).

# Lottery Ticket Hypothesis

- Given a large amount of possible network initializations
- The network successfully trains because a sub-network was initialized properly
- This sub-network can reach similar performance to that of the complete model



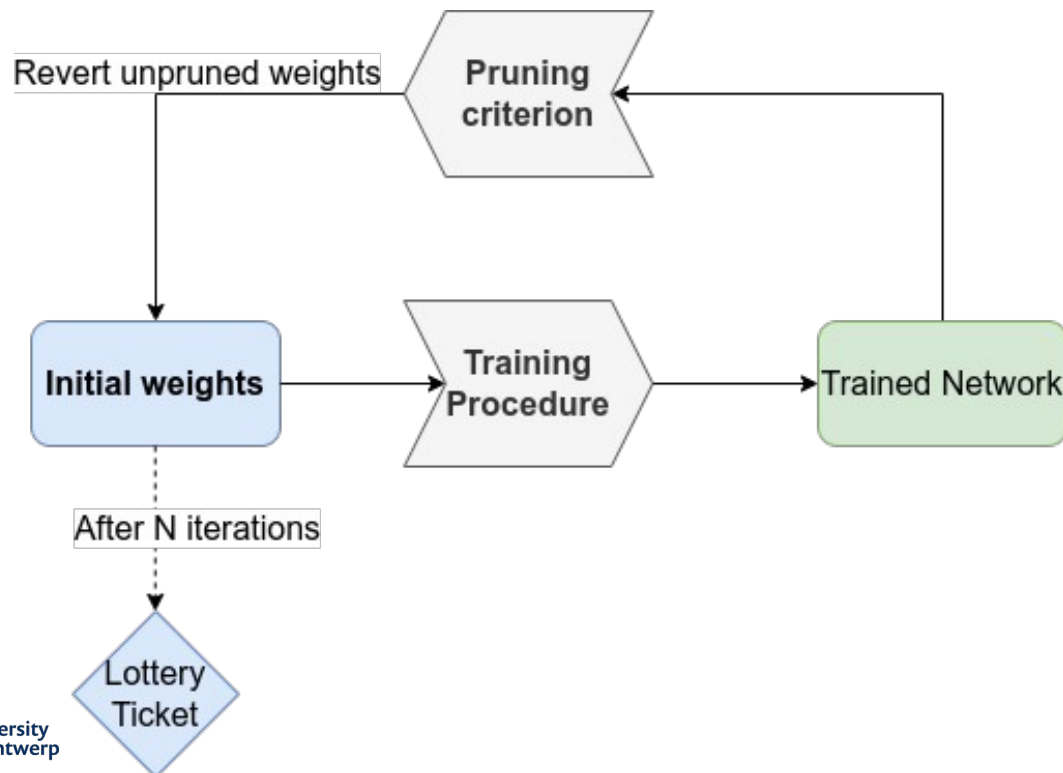
## Finding the Winning Ticket

1. Randomly initialize a neural network.
2. Train the network until it converges.
3. Prune a fraction of the network.
4. Reset the weights of the remaining portion of the network to their values from (1)  
(i.e., the initializations they received before training began).



# Lottery Ticket Hypothesis

- Given a large amount of possible network initializations
- The network successfully trains because a sub-network was initialized properly
- This sub-network can reach similar performance to that of the complete model



## Finding the Winning Ticket

1. Randomly initialize a neural network.
2. Train the network until it converges.
3. Prune a fraction of the network.
4. Reset the weights of the remaining portion of the network to their values from (1)  
(i.e., the initializations they received before training began).

# Summarizing

[ Finally :D ]

# Summarizing

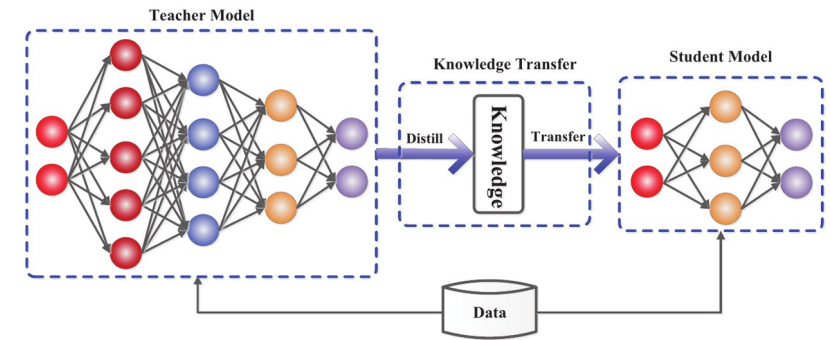
- **Multiple ways to reuse pre-trained models**

With different pros and cons

Reuse features

Reuse architecture

Optimize the architecture on a given principle



# Summarizing

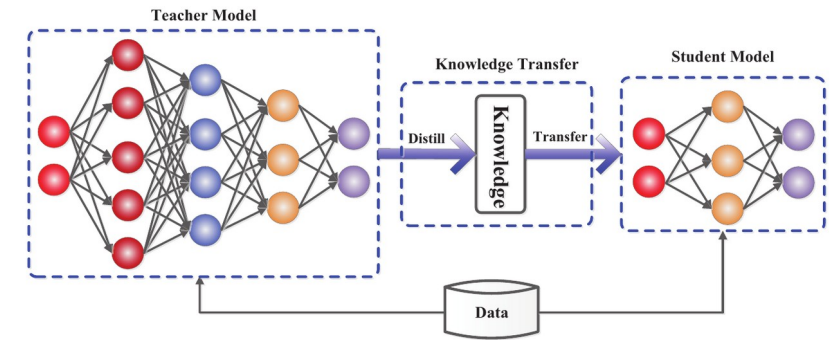
## ■ Multiple ways to reuse pre-trained models

With different pros and cons

Reuse features

Reuse architecture

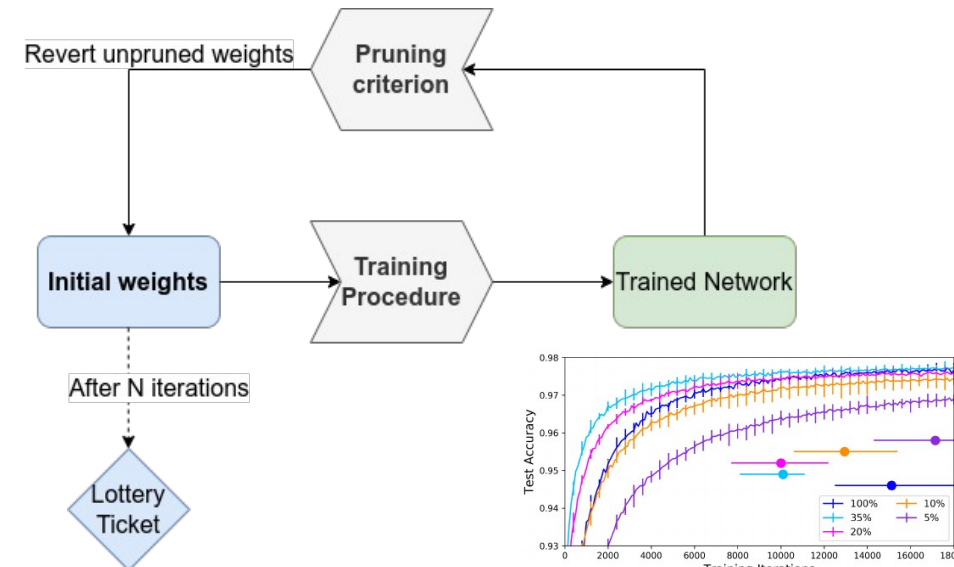
Optimize the architecture on a given principle



## ■ Winning Ticket Representation

Sparse → leads to lighter models

With generalization capabilities



# Questions?

# References

- A. Csiszárík, P. Kőrösi-Szabó, Á. K. Matszangosz, G Papp†, D. Varga. **Similarity and Matching of Neural Network Representations**. NeurIPS 2021.
- F. Ding, J. Denain, J. Steinhardt, **Grounding Representation Similarity Through Statistical Testing**, NeurIPS 2021.
- J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, T. Darrell, **DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition**, ICML 2014
- J Frankle, M. Carbin, **The Lottery Ticket Hypothesis: Training Pruned Neural Networks**, ICLR 2019
- J. Gou, B. Yu, S. Maybank, D. Tao, **Knowledge Distillation: A Survey**, IJCV 2021
- G. Hinton, O. Vinyals, & J. Dean. **Distilling the knowledge in a neural network**. ArXiv:1503.02531, 2015.
- S. Kornblith, M. Norouzi, H. Lee, G. Hinton. **Similarity of Neural Network Representations Revisited**, ICML 2019.
- L.J.P. van der Maaten and G.E. Hinton. **Visualizing High-Dimensional Data Using t-SNE**. Journal of Machine Learning Research 9, 2008



# Artificial Neural Networks

[2500WETANN]

José Oramas