



Week 8 | Lecture 09

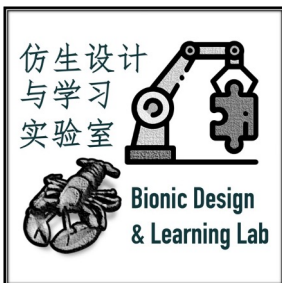
Convolutional Networks

Wan Fang

Southern University of Science and Technology

- **Convolutional Networks**
 - A Design Challenge in 3D Volumes
 - Convolutional Operation
 - The Design of a Convolutional Layer
- **Layers in ConvNets**
 - Three Stages of a Convolutional Layer
 - Common Architectures of ConvNets
 - CONV, POOL, FC Layers
- **A Universal Workflow**
 - Define the Problem & Dataset Assembly
 - Success Metric & Evaluation Protocol
 - Scale up the Model with Regularization
- **Learning a Transition Model**

Convolutional Networks



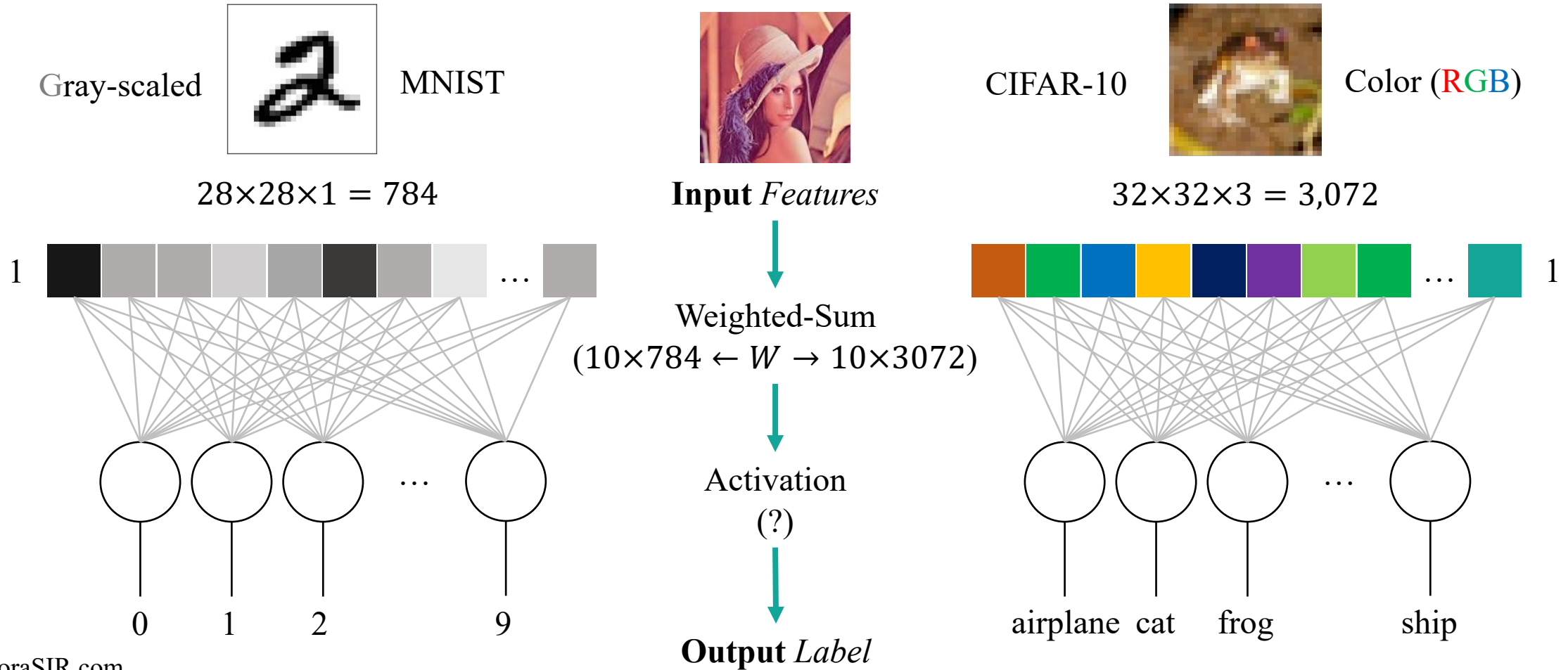
AncoraSIR.com



A Design Challenge with Increasing Dimensions

Regular Neural Nets don't scale well to full images

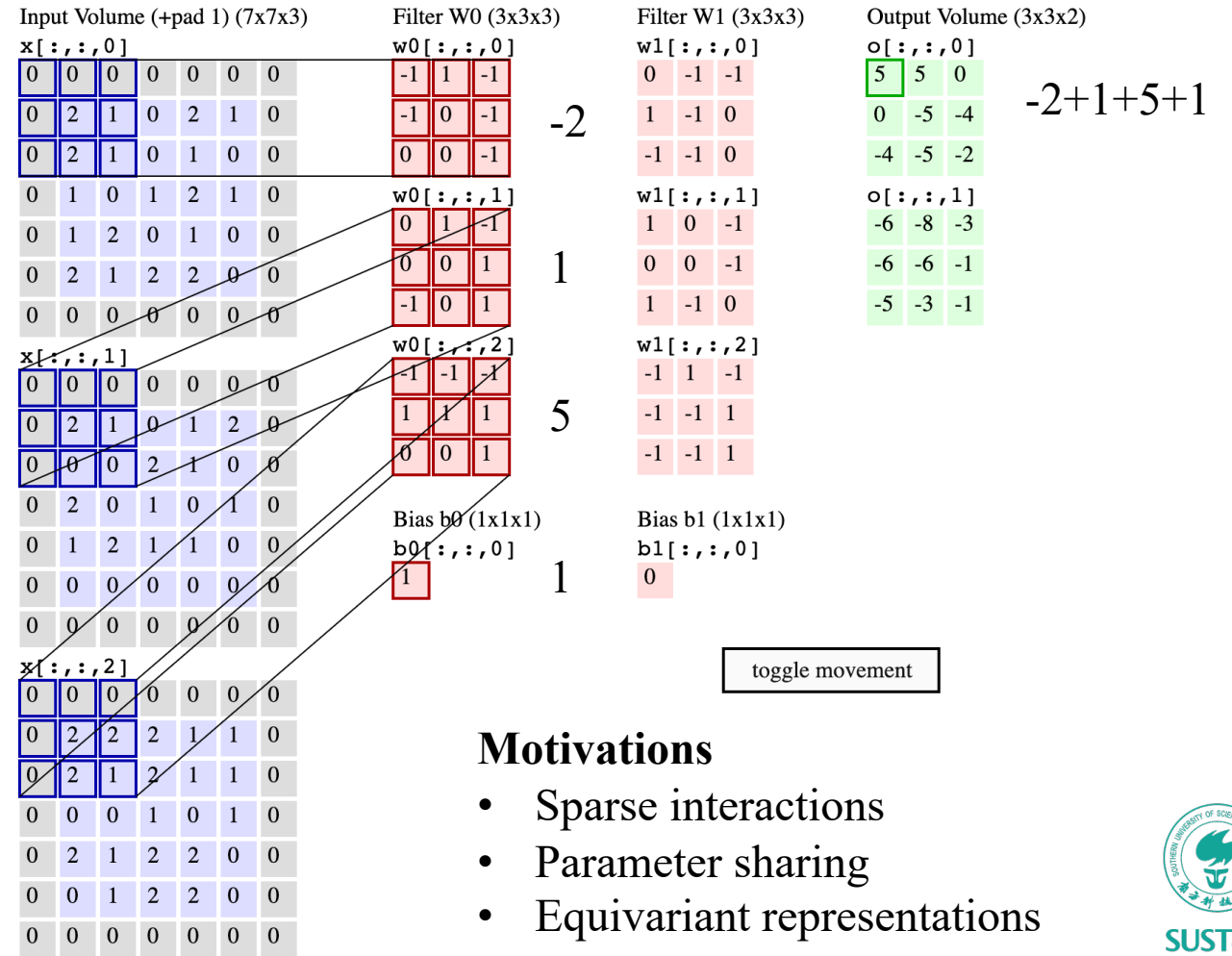
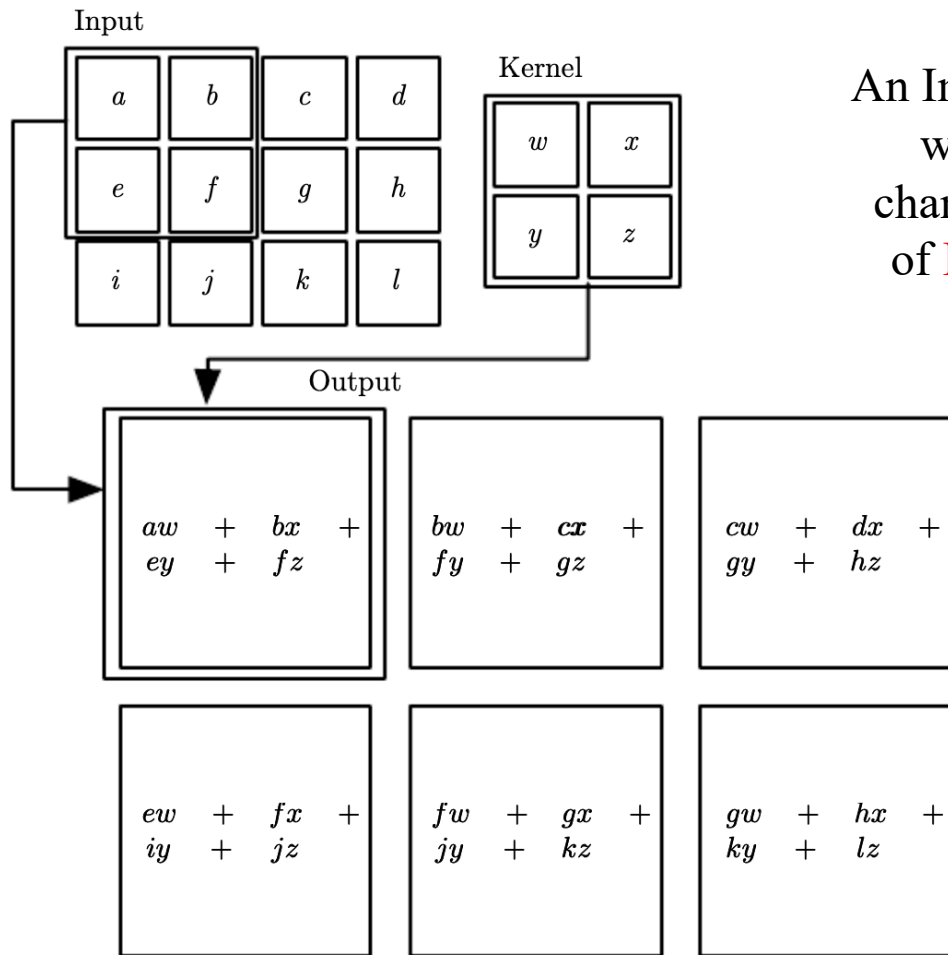
$$512 \times 512 \times 3 = 765,432$$



Convolutional Operation

$$s(t) = \int x(a)w(t - a)da = (x * w)(t)$$

An Image with 3 channels of RGB

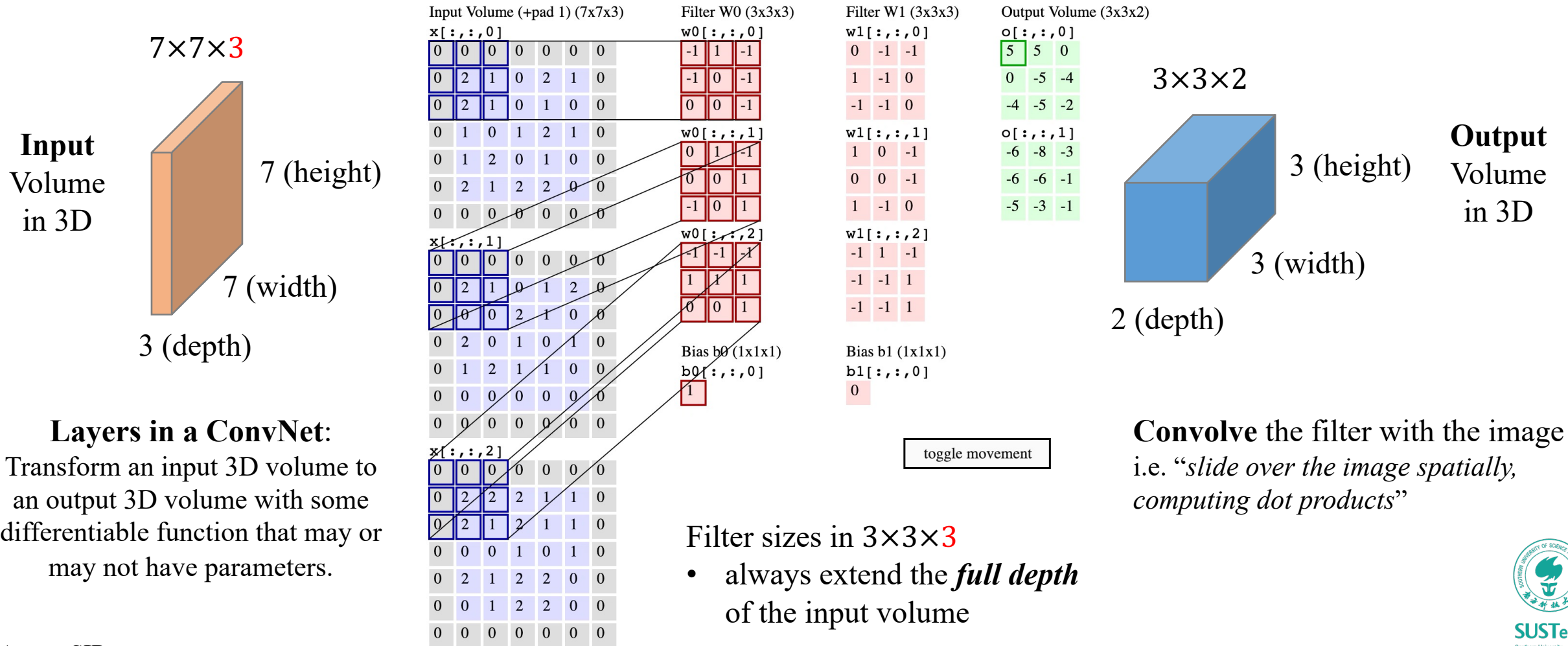


Motivations

- Sparse interactions
- Parameter sharing
- Equivariant representations

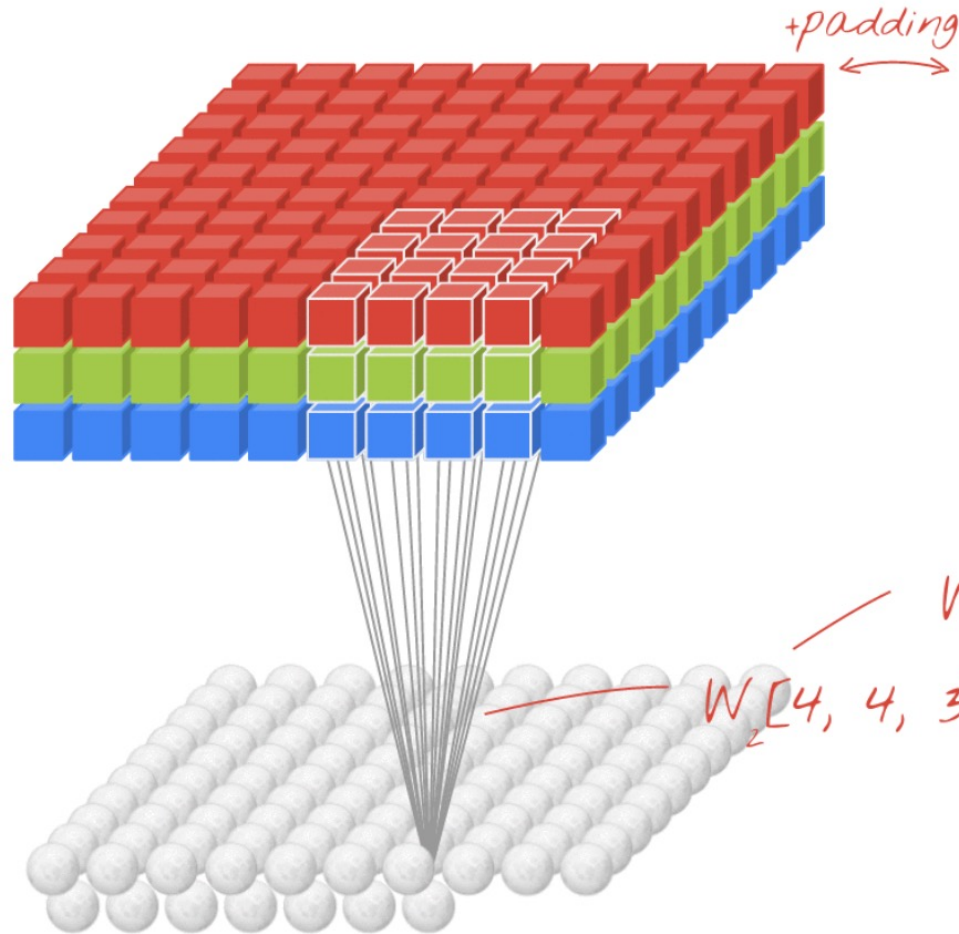
Convolution in 3D Volumes

Preserved spatial structure between the input and output volumes in width, height, number of channels



The Design of a Convolutional Layer

Defined by the filter (or kernel) size, the number of filters applied and the stride

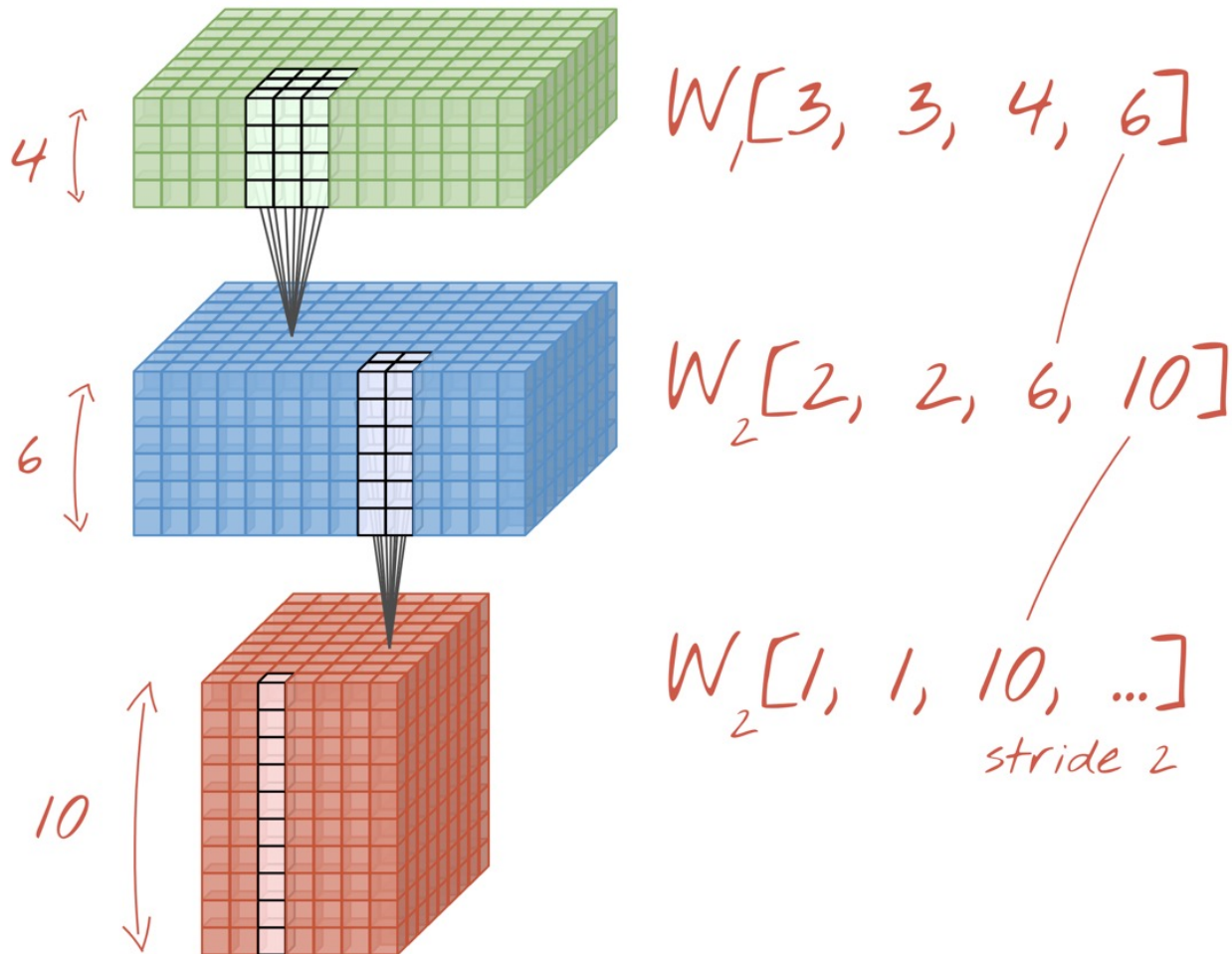


$W_1[4, 4, 3]$
 $W_2[4, 4, 3]$ | $W[4, 4, 3, 2]$

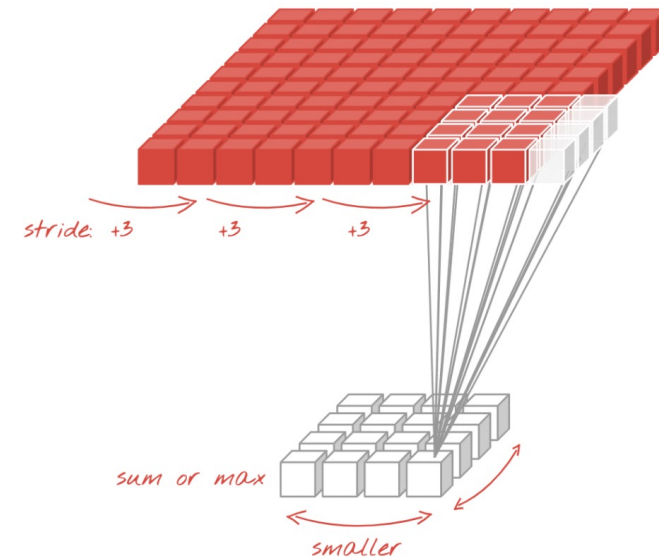
filter size input channels output channels

Output Volume Size

Defined by the filter (or kernel) size, the number of filters applied and the stride



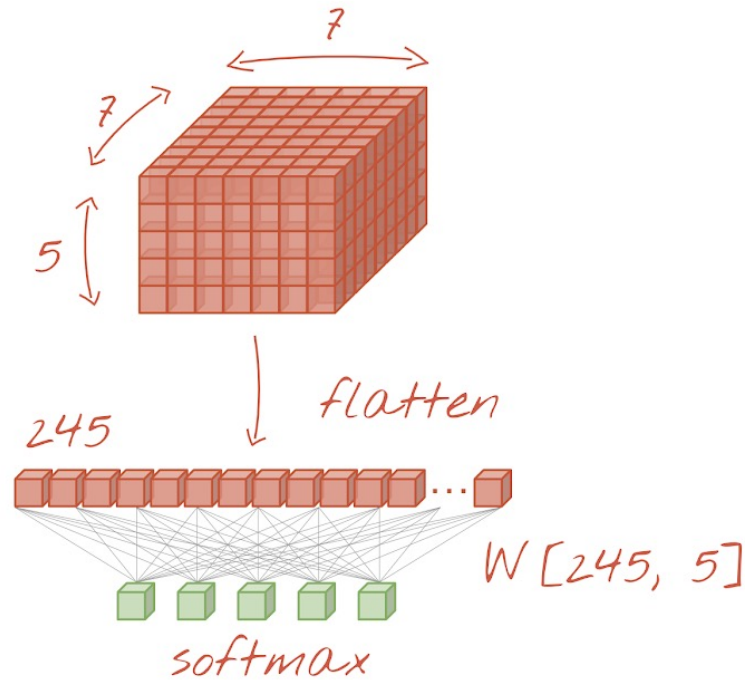
- Depth (number of channels):
 - adjusted by using more or fewer filters
- Width & Height:
 - adjusted by using a stride > 1
 - (or with a max-pooling operation)



The Last Layer

From a Cubic Volume in 3D to predicted labels

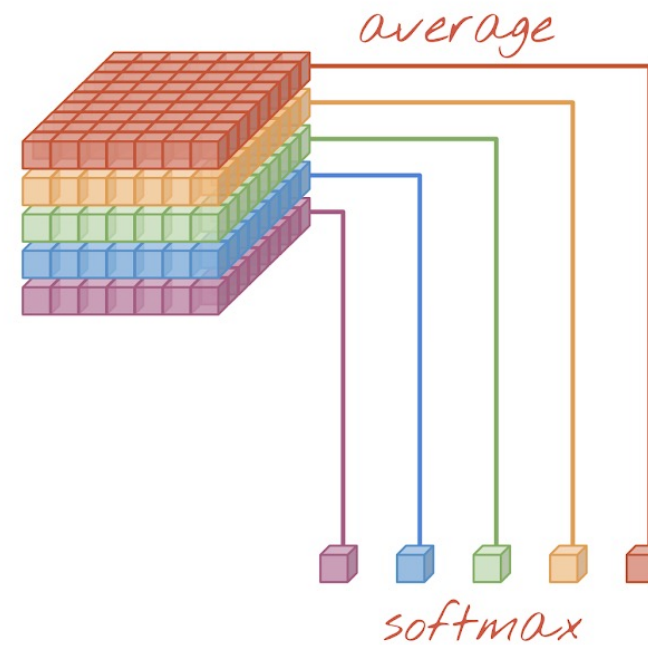
Fully connected layer



1225 weights



Global average pooling



0 weights

Similar like a normal neural network

Expensive in #weights

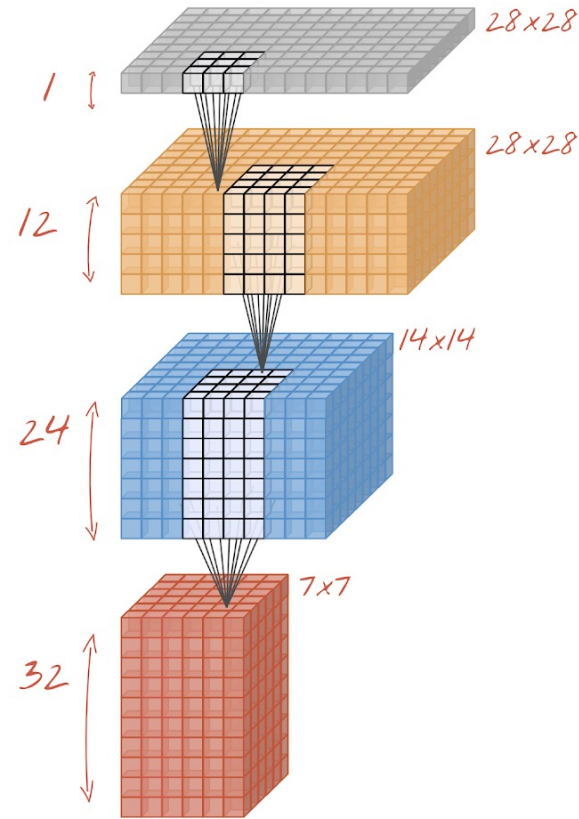
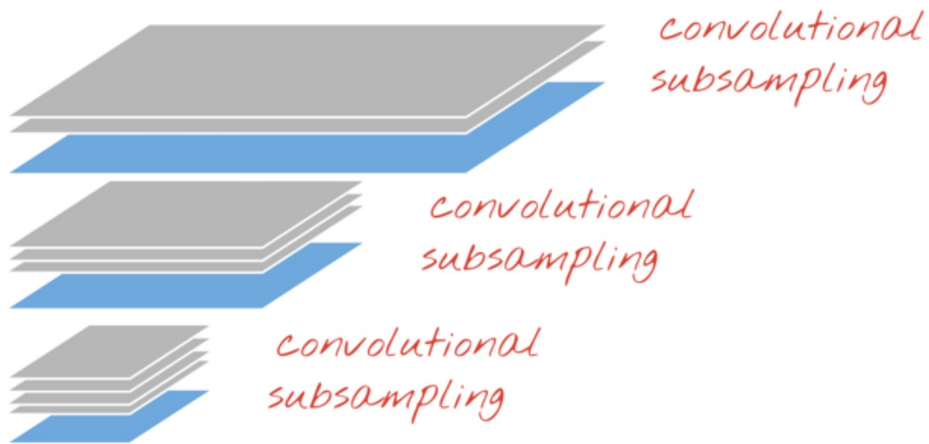
But preserves the location data (x, y)

Much lighter in calculation

The average pooling explicitly discards all location data

Stacking Up a ConvNet

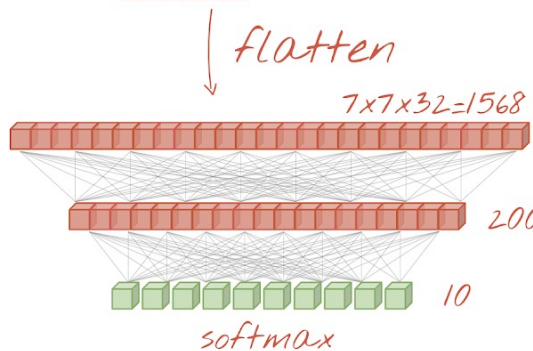
Layer-by-layer



Convolutional 3x3 filters=12
 $W_1[3, 3, 1, 12]$

Convolutional 6x6 filters=24
 $W_2[6, 6, 12, 24]$ stride 2

Convolutional 6x6 filters=32
 $W_3[6, 6, 24, 32]$ stride 2



Dense layer
 $W_4[1568, 200]$

Softmax dense layer
 $W_5[200, 10]$

Layers in ConvNets

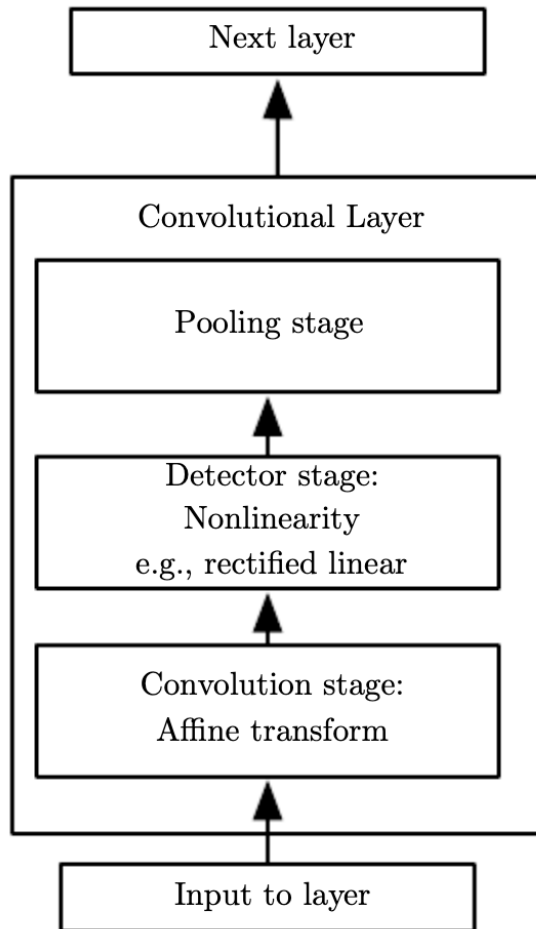


AncoraSIR.com



The Three Stages of a Typical ConvNet Layer

The Convolution, Detector and Pooling Stages



- The maximum output within a rectangular neighborhood (max-pooling)
- The average of a rectangular neighborhood
- The L2 norm of a rectangular neighborhood
- A weighted average based on the distance from the central pixel

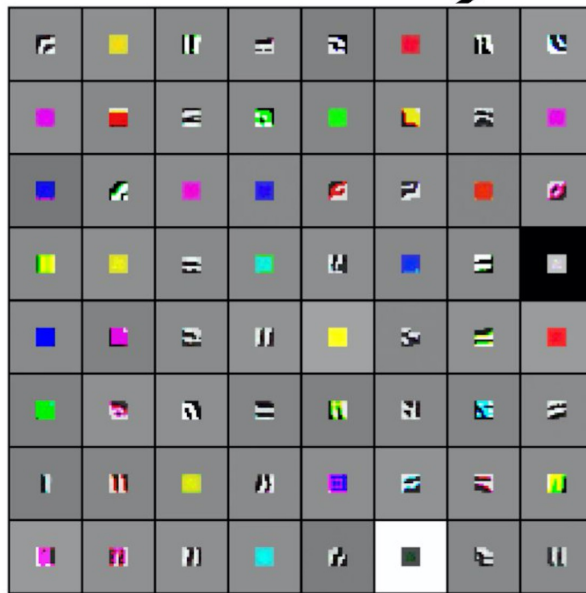
Replace the output of the net at a certain location with a summary statistic of the nearby outputs (can be viewed as a further abstraction of the learned features)

Each linear activation is run through a nonlinear activation function, such as ReLU (can be viewed as activation function)

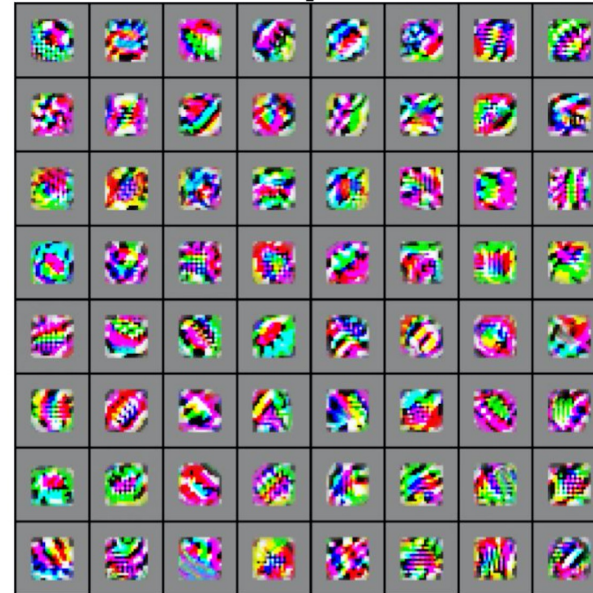
Performs several convolutions in parallel to produce a set of linear activations (can be viewed as weighted-sum)

A Visualized Understanding of ConvNet

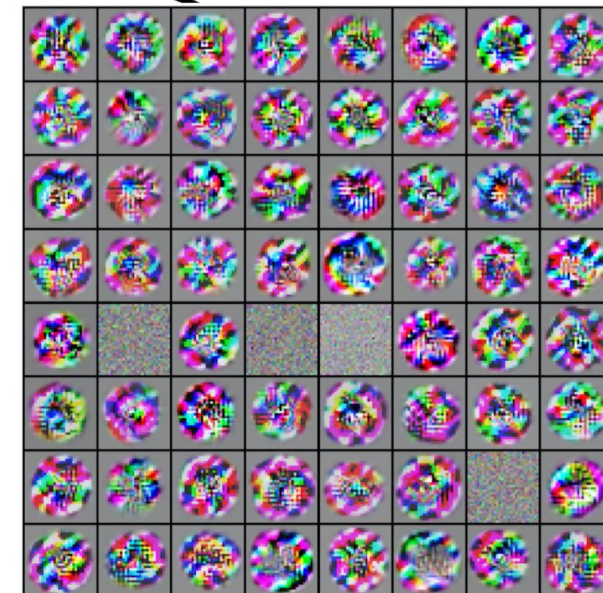
Multi-layered abstraction of 3D features towards a linearly separable classification



VGG-16 Conv1_1



VGG-16 Conv3_2



VGG-16 Conv5_3

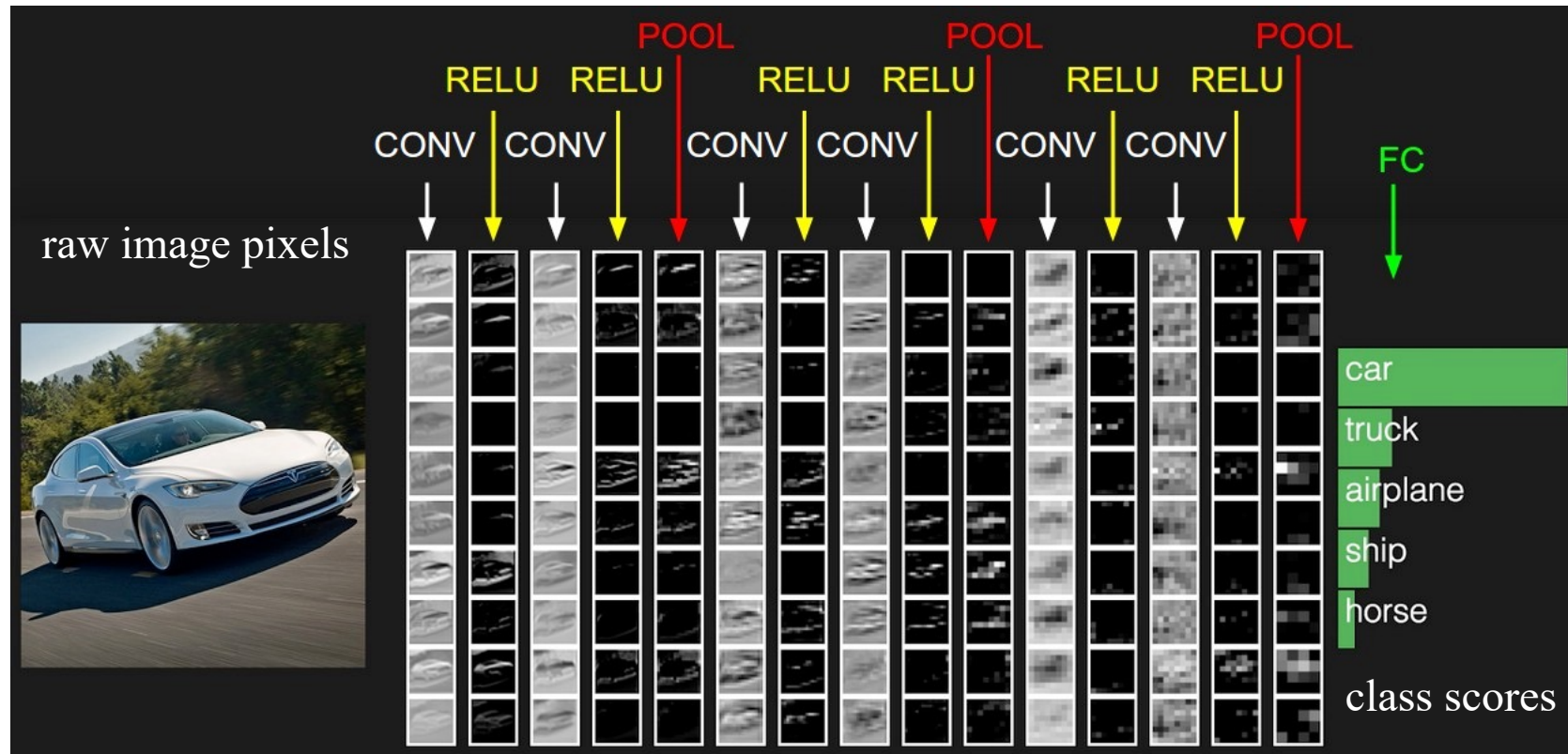
A Simple ConvNet for CIFAR-10 Classification

[INPUT - CONV - RELU - POOL - FC]

CONV layer compute the output of neurons that are connected to local regions in the input, i.e. $[32 \times 32 \times 12]$ with 12 filters.

RELU layer will apply an elementwise activation function, such as the $\max(0, x)$ thresholding at zero. This leaves the size of the volume unchanged ($[32 \times 32 \times 12]$).

POOL layer will perform a downsampling operation along the spatial dimensions (width, height), resulting in volume such as $[16 \times 16 \times 12]$.



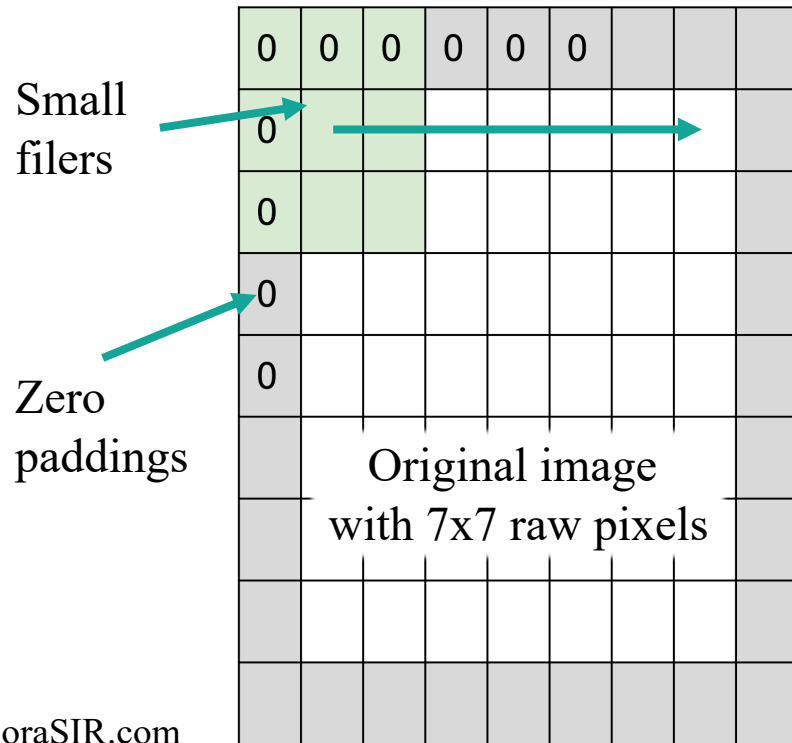
INPUT layer $[32 \times 32 \times 3]$ will hold the raw pixel values of the image

FC (i.e. fully-connected) layer will compute the class scores, resulting in volume of size $[1 \times 1 \times 10]$, where each of the 10 numbers correspond to a class score

Convolutional Layer

Small filters that slide across the input volume

- Small-size filers
 - e.g. 3x3 or at most 5x5, using a stride of $S=1$,
 - Padding the input volume with zeros to avoid altering the spatial dimensions of the input.



INPUT features: 7x7

Filer size: 3x3

Stride: 1 (move step-by-step)

Padding: 1 pixel of 0 on all borders

OUTPUT features: 7x7

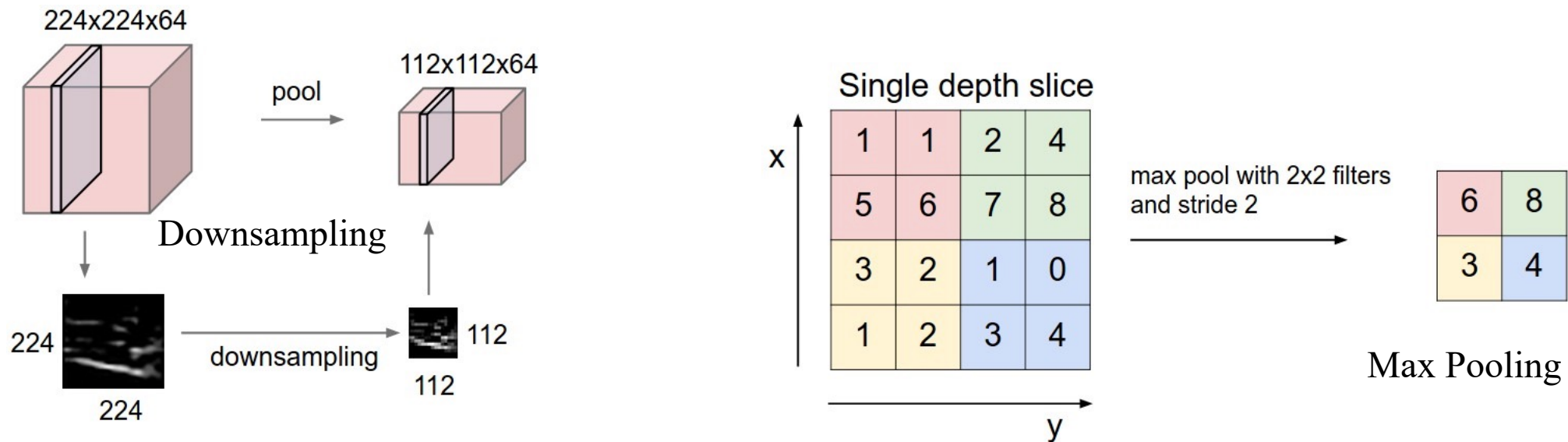
What if without paddings on the border?

- *The spatial dimensions of the input will be changed, causing information loss on the border*

Pooling Layer

Downsampling the spatial dimensions of the input volume

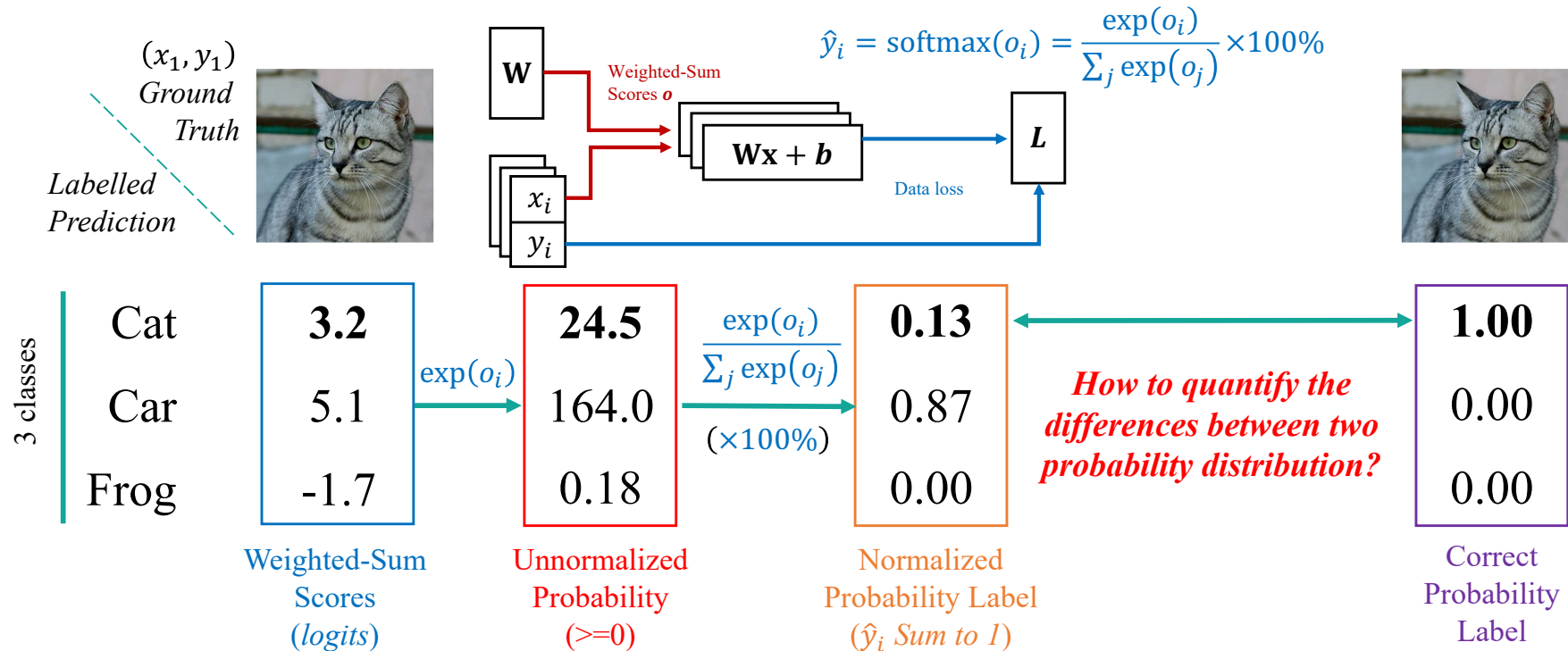
- A network-wise regularization
 - Progressively reduce the spatial size of the representation to reduce the amount of parameters and computation in the network, and hence to also control overfitting
 - Operates over each activation map independently
 - Usually, no need to zero padding (no convolutional operations)



Fully-Connected Layer

Full connections to all activations in the previous layer, as seen in regular Neural Networks

- Contains neurons that connect to the entire input volume
- Softmax is a common choice



ConvNet Architectures

Common choice of hyperparameters of ConvNet designs

- **INPUT** \rightarrow $[[\text{CONV} \rightarrow \text{RELU}] * N \rightarrow \text{POOL?}] * M \rightarrow [\text{FC} - \text{RELU}] * K \rightarrow \text{FC}$
 - the * indicates repetition,
 - the POOL? indicates an optional pooling layer.
 - $N \geq 0$ (and usually $N \leq 3$), $M \geq 0$, $K \geq 0$ (and usually $K < 3$)
- **INPUT** (that contains the image) should be divisible by 2 many times
 - 32 (e.g. CIFAR-10), 64, 96 (e.g. STL-10), or 224 (e.g. ImageNet), 384, and 512
- **CONV** should be using small filters using a stride of $S=1$
 - 3x3 or at most 5x5 with zero padding of the input volume
- **POOL** downsamples the spatial dimensions of the input
 - Common setting is to use max-pooling with 2x2 receptive fields with a stride of 2

A Universal Workflow



AncoraSIR.com



Defining the Problem

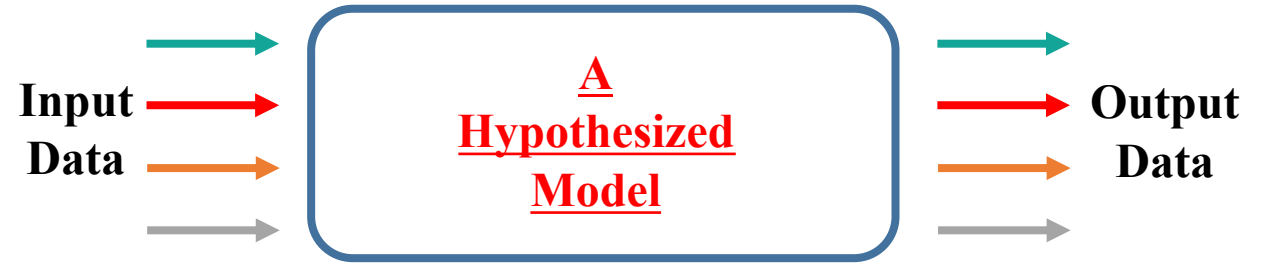
Assembling a Dataset

- **Data availability**
 - Usually the limiting factor at this stage.
- **Identifying the problem type**
 - Guide your choice of model architecture, loss function, etc.
- **A Hypothesis of Modeling**
 - *Can you build a model using price history to predict stock market? ☹️*
- **Non-stationary Problem?**
 - *Will your data change over time?*

What will your input data be?

What are you trying to predict?

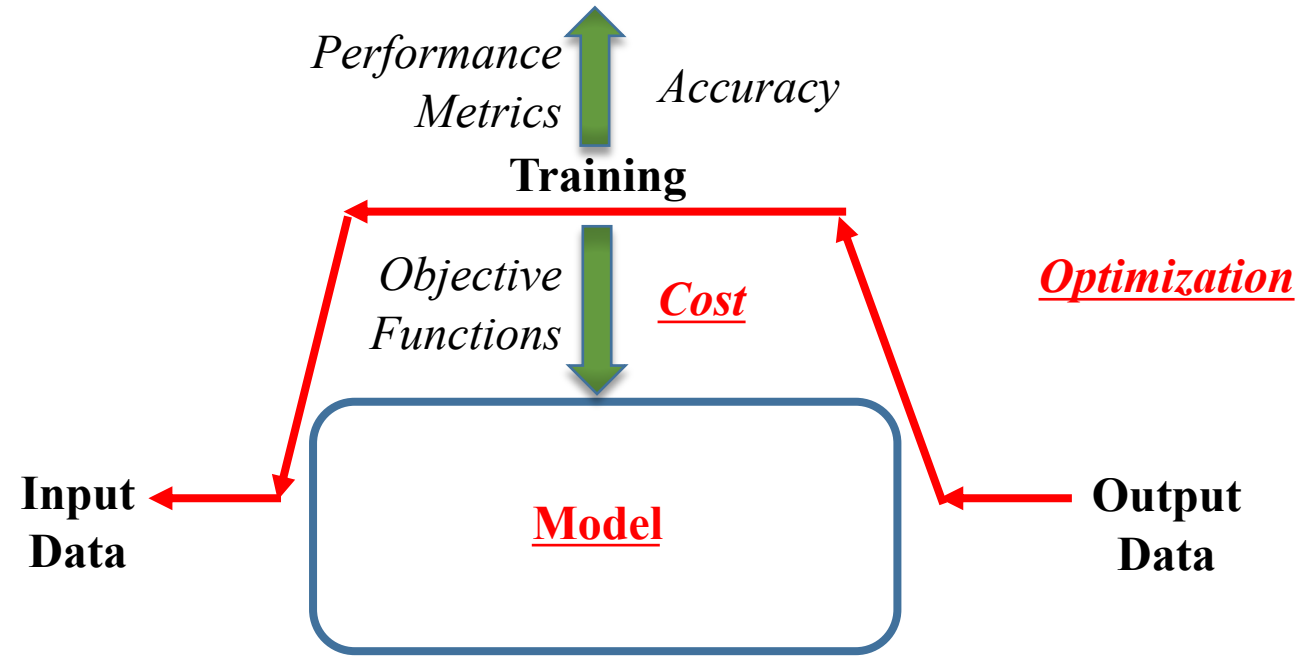
- Can only learn to predict something **if you have available training data**
- Can only be used to **learn patterns** that are present in training data
 - *A fundamental assumption that the future will behave like the past*



- *What type of problem are you facing?*
- *Is it binary classification?*
- *Multiclass classification?*
- *Scalar regression?*
- *Vector regression?*
- *Multiclass, multilabel classification?*
- *Something else, like clustering, generation, or reinforcement learning?*

Choosing a Measure of Success

You need to be able to observe it to control it



Choosing a Measure of Success

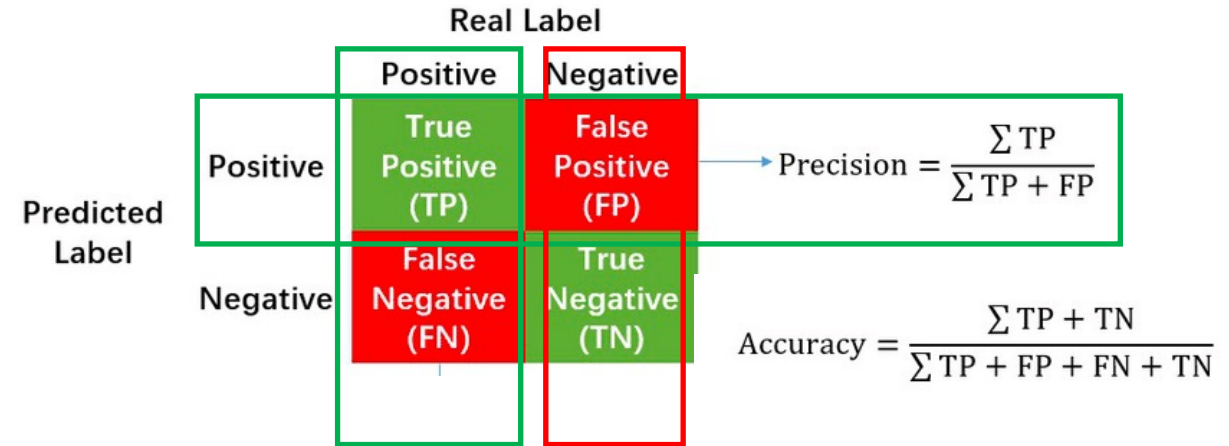
You need to be able to observe it to control it

Your metric for success will guide the choice of a loss function

- Accuracy? Precision and recall? **What your model optimizes?**
- For balanced-classification problems (every class is equally likely): *accuracy and area under the receiver operating characteristic curve (ROC AUC)* are common metrics
- For class-imbalanced problems, you can use precision and recall

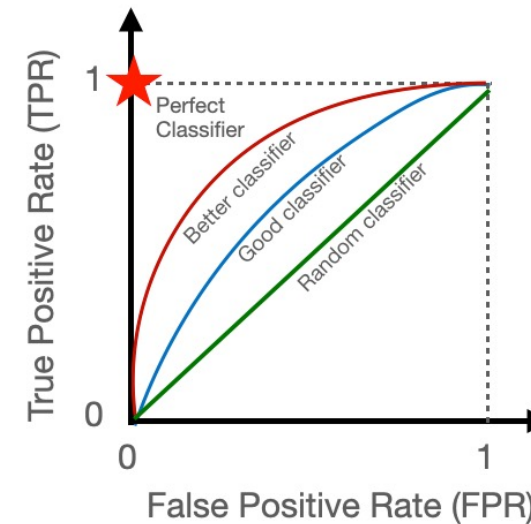
TP=63	FN=37	100
FP=28	TN=72	100
91	109	200

AncoraSIR.com



True Positive Rate (TPR)
 $= TP / (TP + FN)$
 $= \text{efficiency to identify the signal (Recall)}$

False Positive Rate (FPR)
 $= FP / (FP + TN)$
 $= \text{inefficiency to reject background}$



Source



Choosing a Measure of Success

You need to be able to observe it to control it

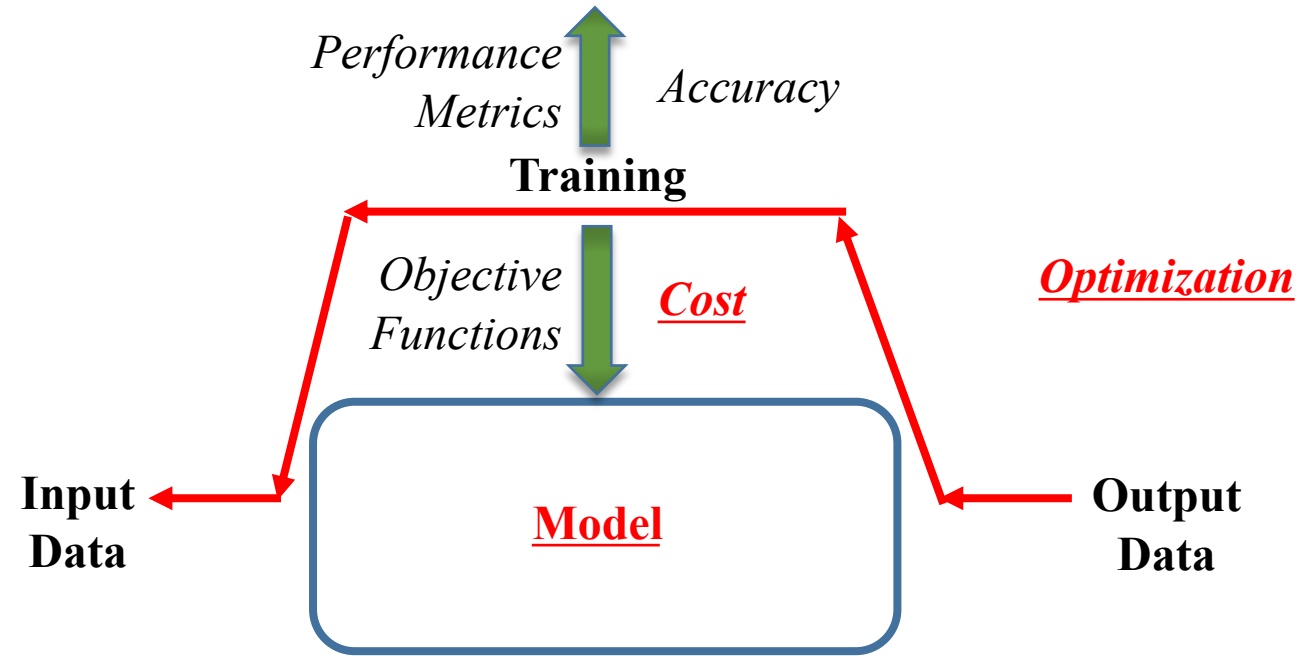
Your metric for success will guide the choice of a loss function

- Accuracy? Precision and recall? **What your model optimizes?**
- For balanced-classification problems (every class is equally likely): *accuracy and area under the receiver operating characteristic curve (ROC AUC)* are common metrics
- For class-imbalanced problems, you can use precision and recall

It isn't uncommon to have to **define your own custom metric**.

- Common in Robotics
- Read more papers
- Be clear and direct about your modeling goal

AncoraSIR.com



		Real Label		
		Positive	Negative	
Predicted Label	Positive	True Positive (TP)	False Positive (FP)	Precision = $\frac{\sum TP}{\sum TP + FP}$
	Negative	False Negative (FN)	True Negative (TN)	

Recall = $\frac{\sum TP}{\sum TP + FN}$

Accuracy = $\frac{\sum TP + TN}{\sum TP + FP + FN + TN}$

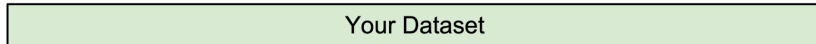
Deciding on an Evaluation Protocol

How you'll measure your current progress

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the data

BAD: $K = 1$ always works perfectly on training data



Maintaining a hold-out validation set

- The way to go when you have plenty of data

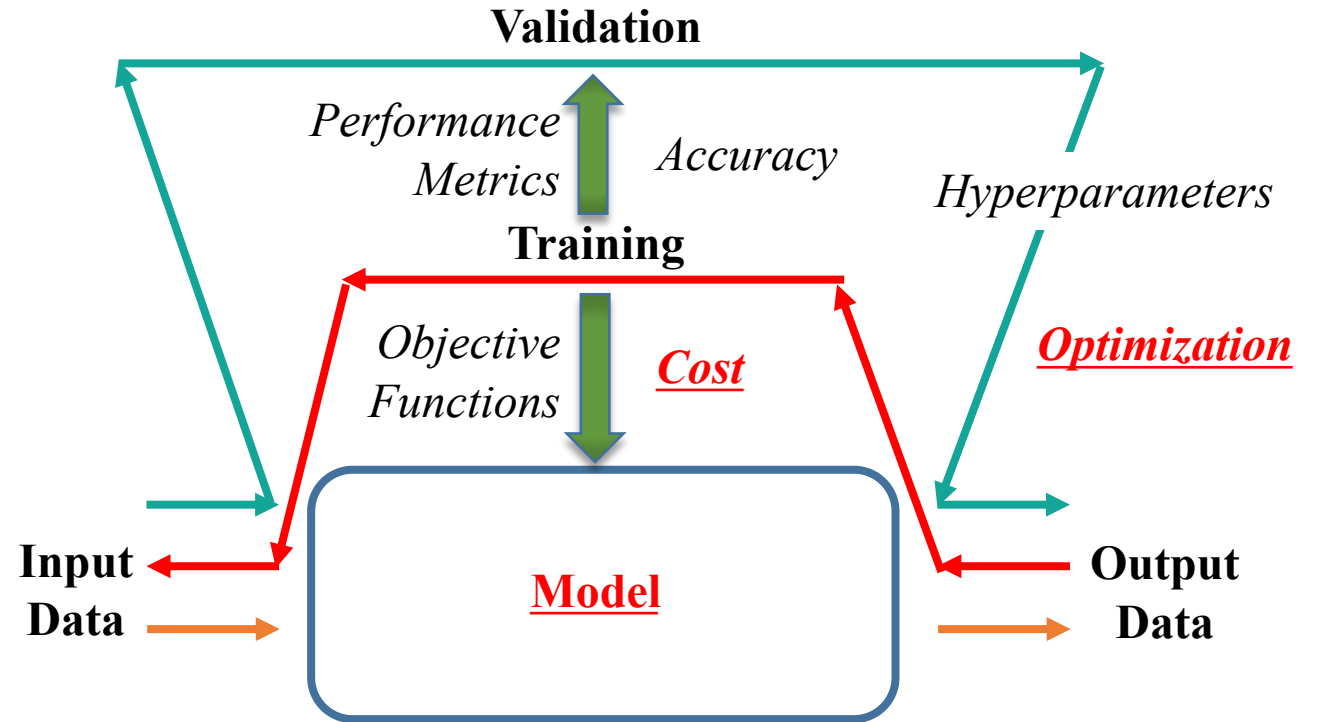
Doing K-fold cross-validation

- The right choice when you have too few samples for hold-out validation to be reliable

Doing iterated K-fold validation

- For performing highly accurate model evaluation when little data is available

AncoraSIR.com



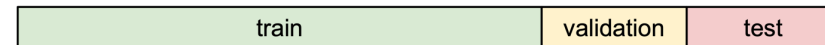
Idea #2: Split data into **train** and **test**, choose hyperparameters that work best on test data

BAD: No idea how algorithm will perform on new data

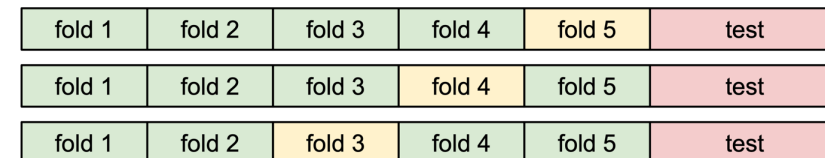


Idea #3: Split data into **train**, **val**, and **test**; choose hyperparameters on val and evaluate on test

Better!



Idea #4: Cross-Validation: Split data into **folds**, try each fold as validation and average the results



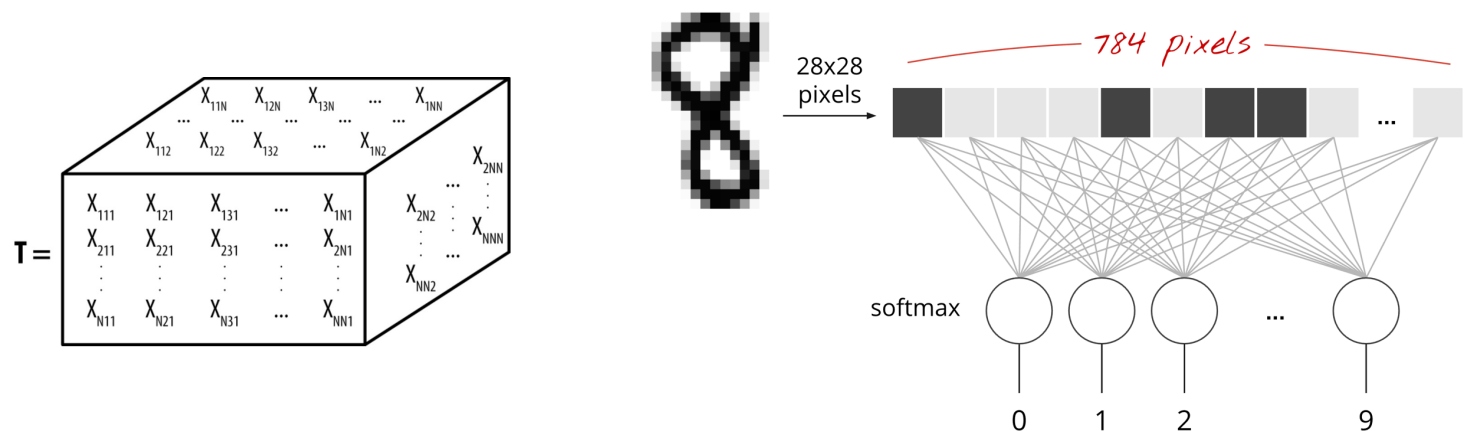
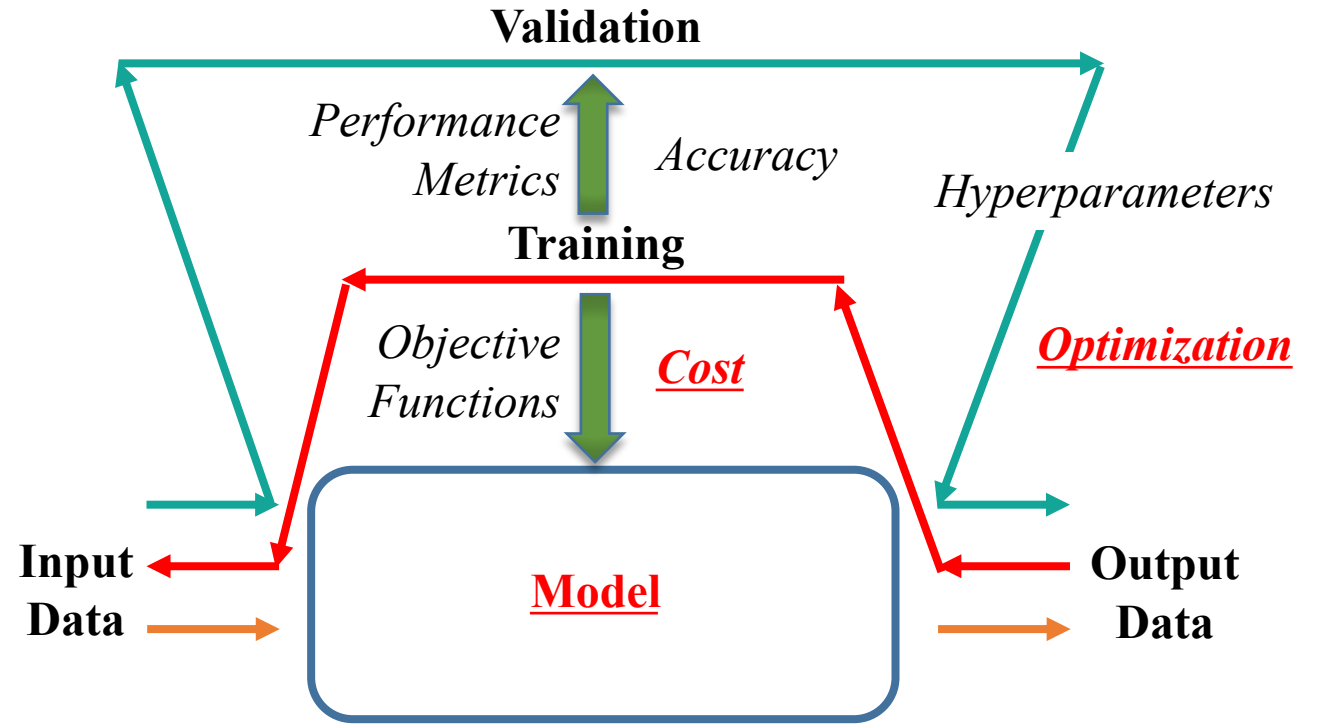
Useful for small datasets, but not used too frequently in deep learning



Preparing Your Data

Format your data in a way that can be fed into a model of deep neural network

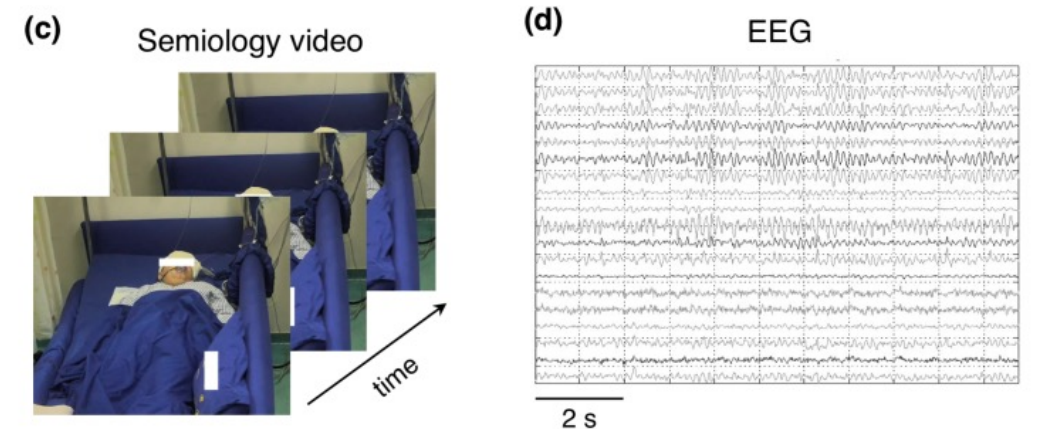
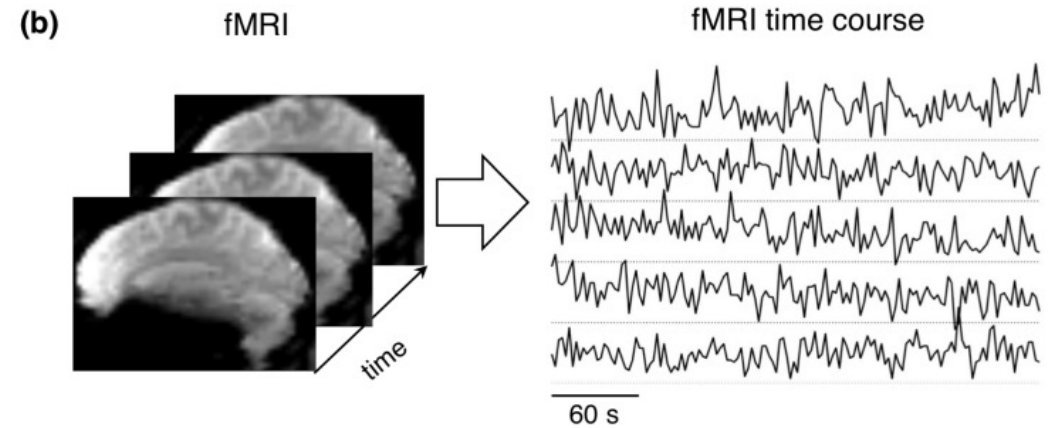
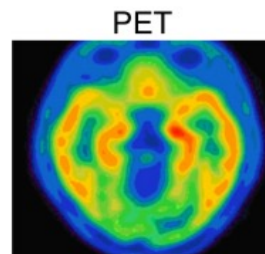
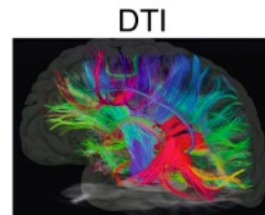
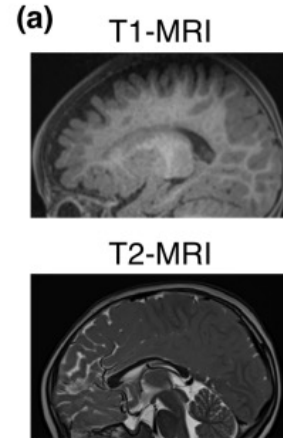
- Your data should be formatted as **tensors**
- The values taken by these tensors should usually be **scaled** to small values
 - For example, in the $[-1, 1]$ range or $[0, 1]$ range
- **Normalize** the data if different features take values in different ranges



Preparing Your Data

Format your data in a way that can be fed into a model of deep neural network

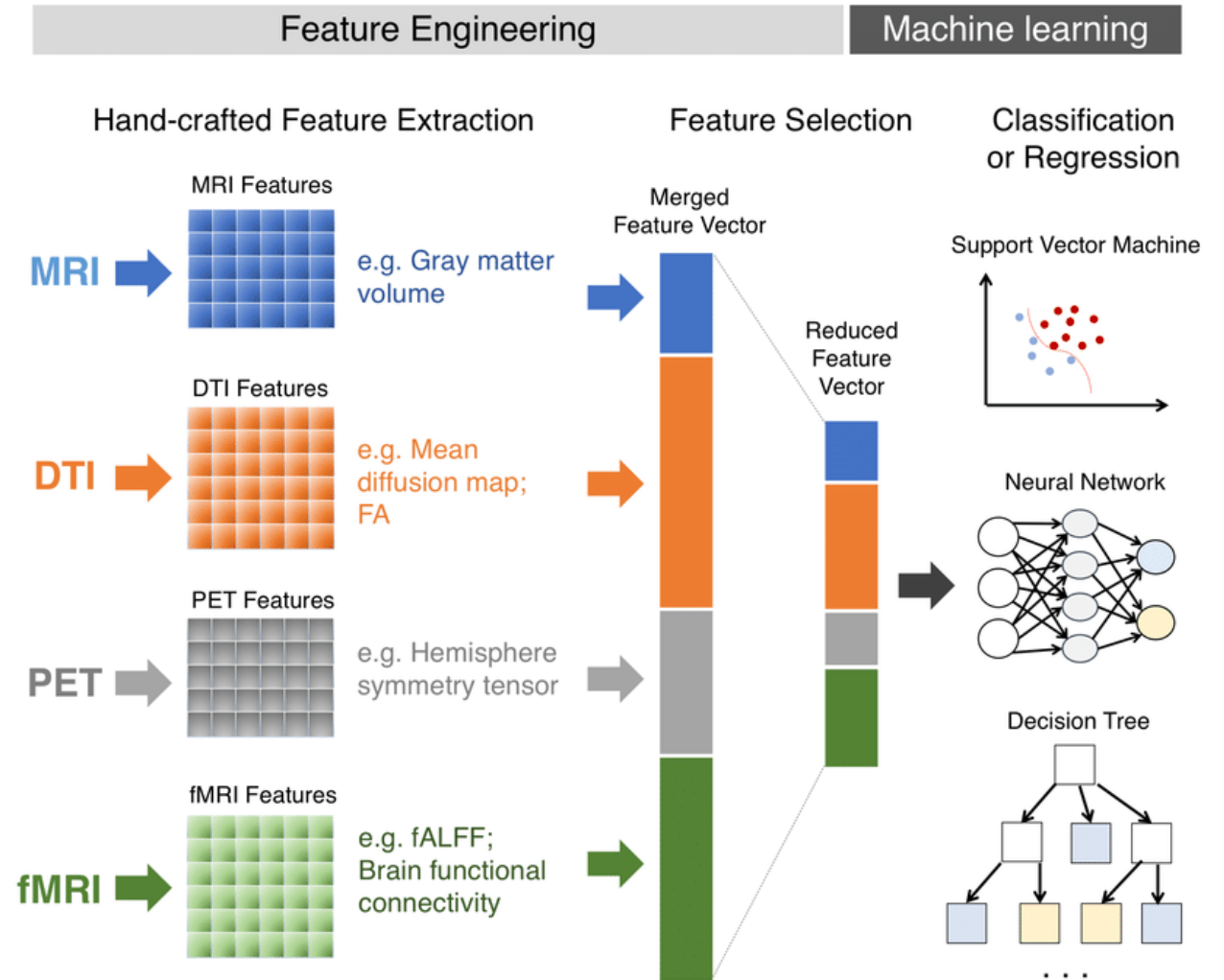
- Your data should be formatted as **tensors**
- The values taken by these tensors should usually be **scaled** to small values
 - For example, in the $[-1, 1]$ range or $[0, 1]$ range
- **Normalize** the data if different features take values in different ranges
- **Feature engineering** may be required, especially for small-data problems



Preparing Your Data

Format your data in a way that can be fed into a model of deep neural network

- Your data should be formatted as **tensors**
- The values taken by these tensors should usually be **scaled** to small values
 - For example, in the $[-1, 1]$ range or $[0, 1]$ range
- **Normalize** the data if different features take values in different ranges
- **Feature engineering** may be required, especially for small-data problems



Developing a Model Better than a Baseline

To achieve statistical power by building a small model that beats a dumb baseline

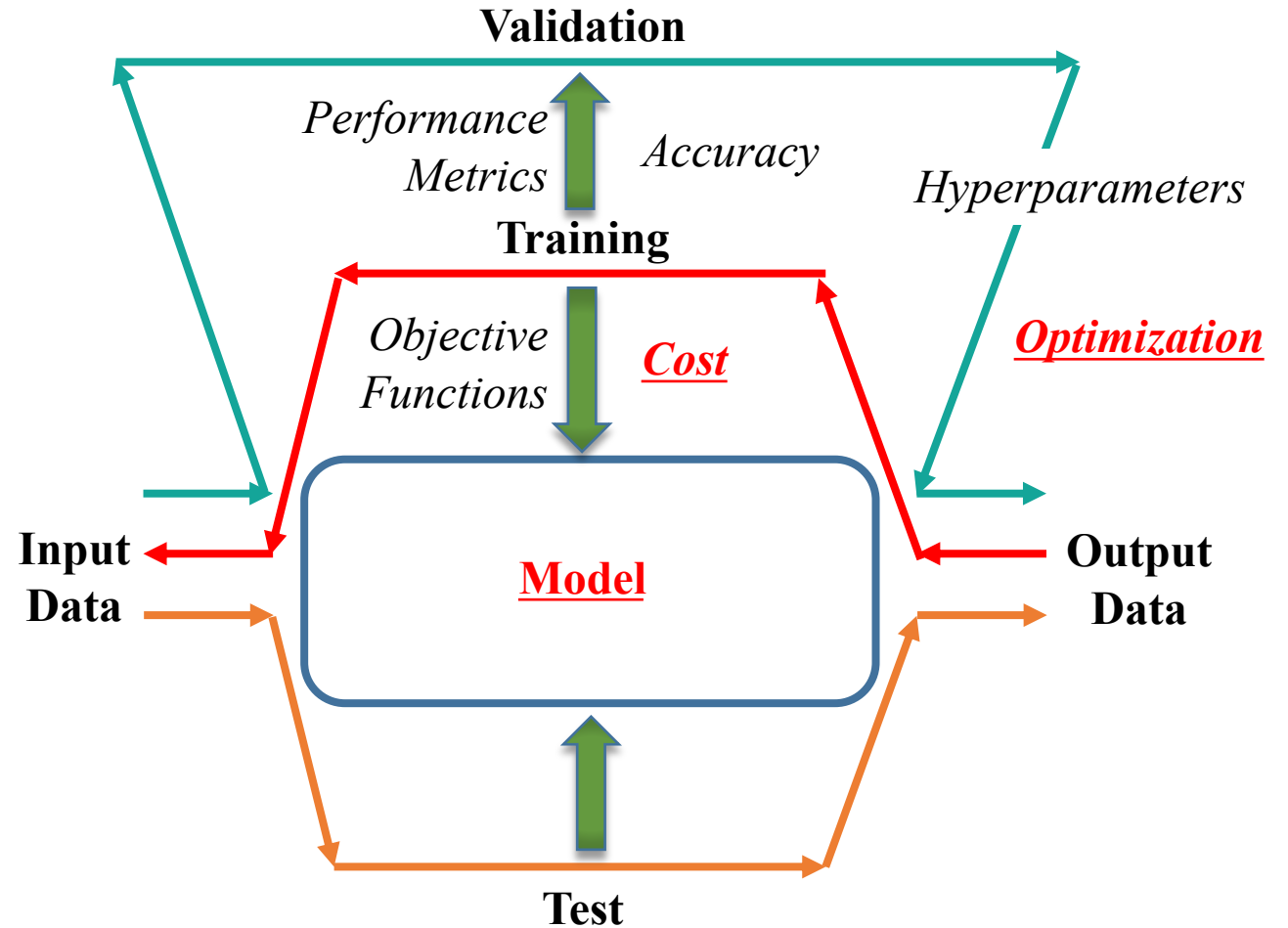
Check your assumptions against a baseline model

- Your outputs can be predicted given your inputs.
- The available data is sufficiently informative to learn the relationship between inputs and outputs.
- (*asking the right questions to your data*)

Key choices to build a model

- *Last-layer activation*—This establishes useful constraints on the network’s output.
- *Loss function*—This should match the type of problem you’re trying to solve.
- *Optimization configuration*—What optimizer will you use? What will its learning rate be?

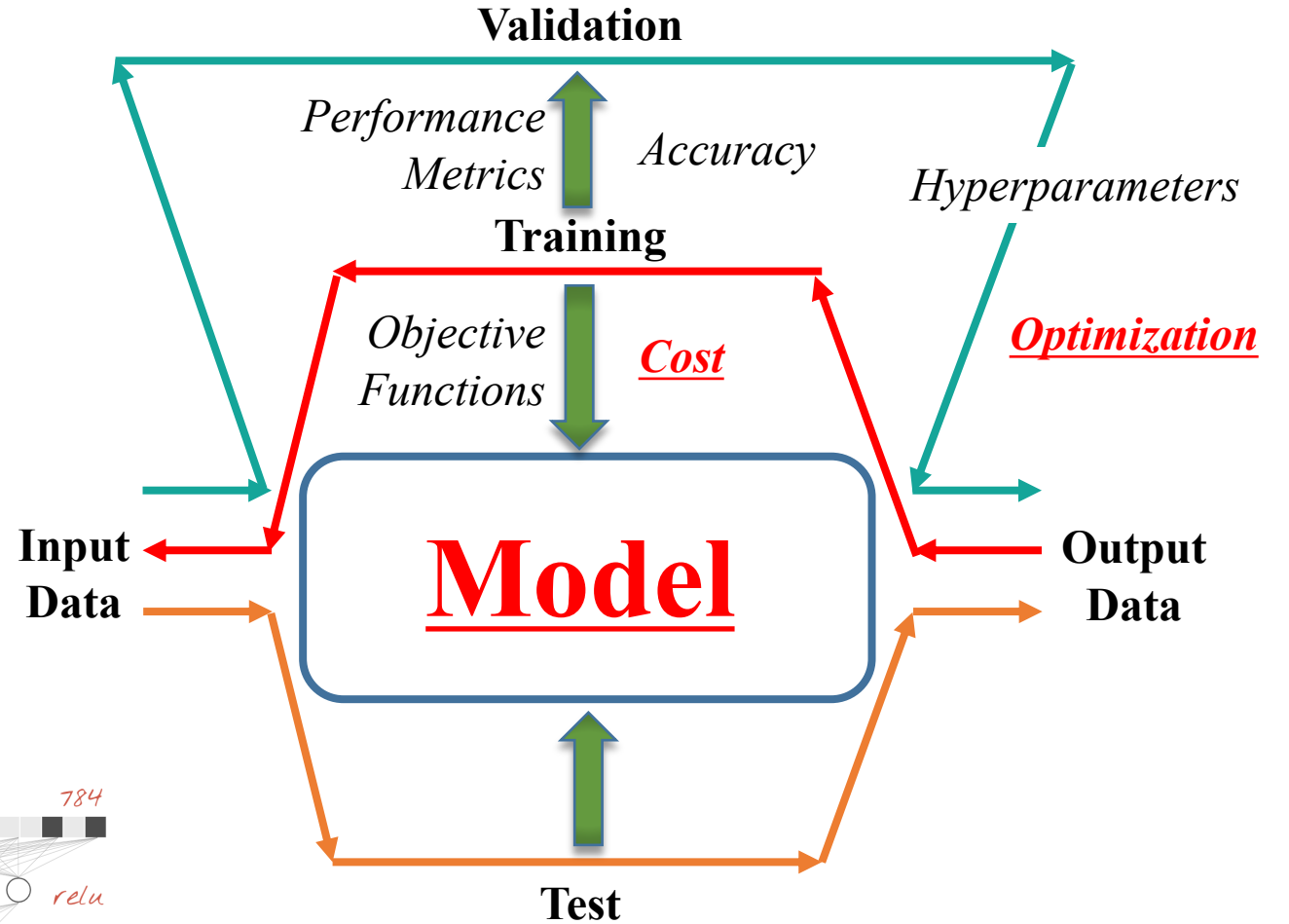
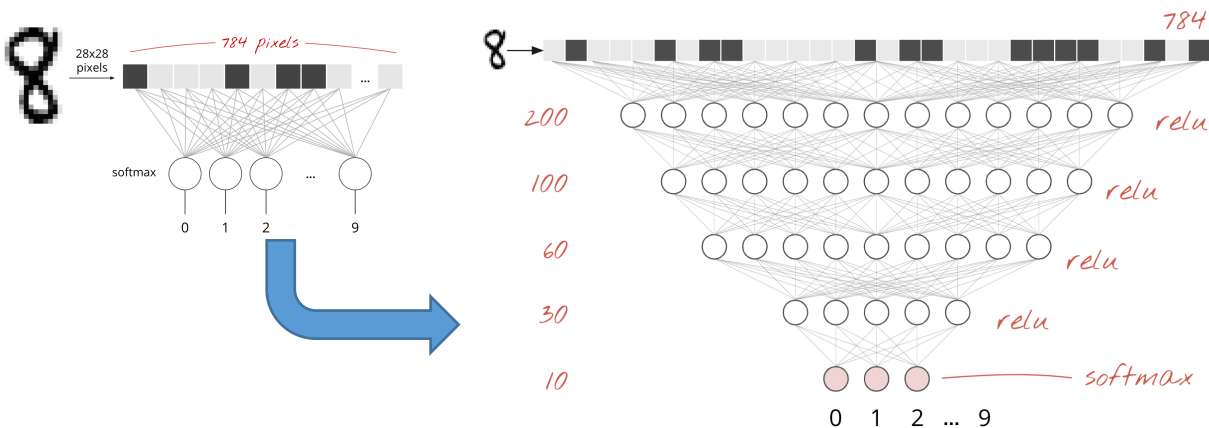
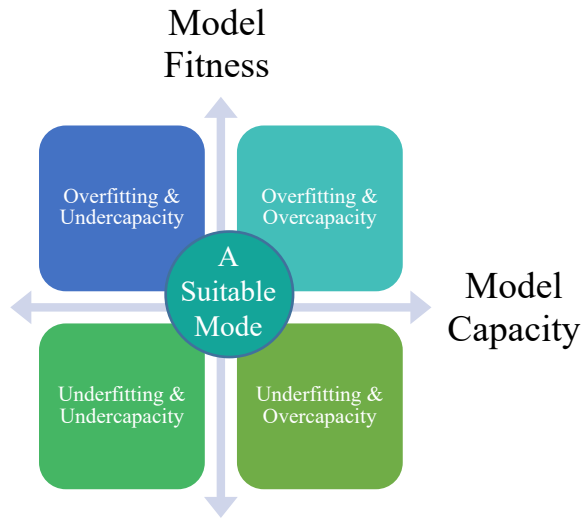
AncoraSIR.com



Problem type	Last-layer activation	Loss function
Binary classification	sigmoid	binary_crossentropy
Multiclass, single-label classification	softmax	categorical_crossentropy
Multiclass, multilabel classification	sigmoid	binary_crossentropy
Regression to arbitrary values	None	mse
Regression to values between 0 and 1	sigmoid	mse or binary_crossentropy

Scaling Up: Developing a Model that Overfits

1. Add layers.
2. Make the layers bigger.
3. Train for more epochs.



- *Is your model sufficiently powerful?*
- *Does it have enough layers and parameters to properly model the problem at hand?*

Regularizing Your Model

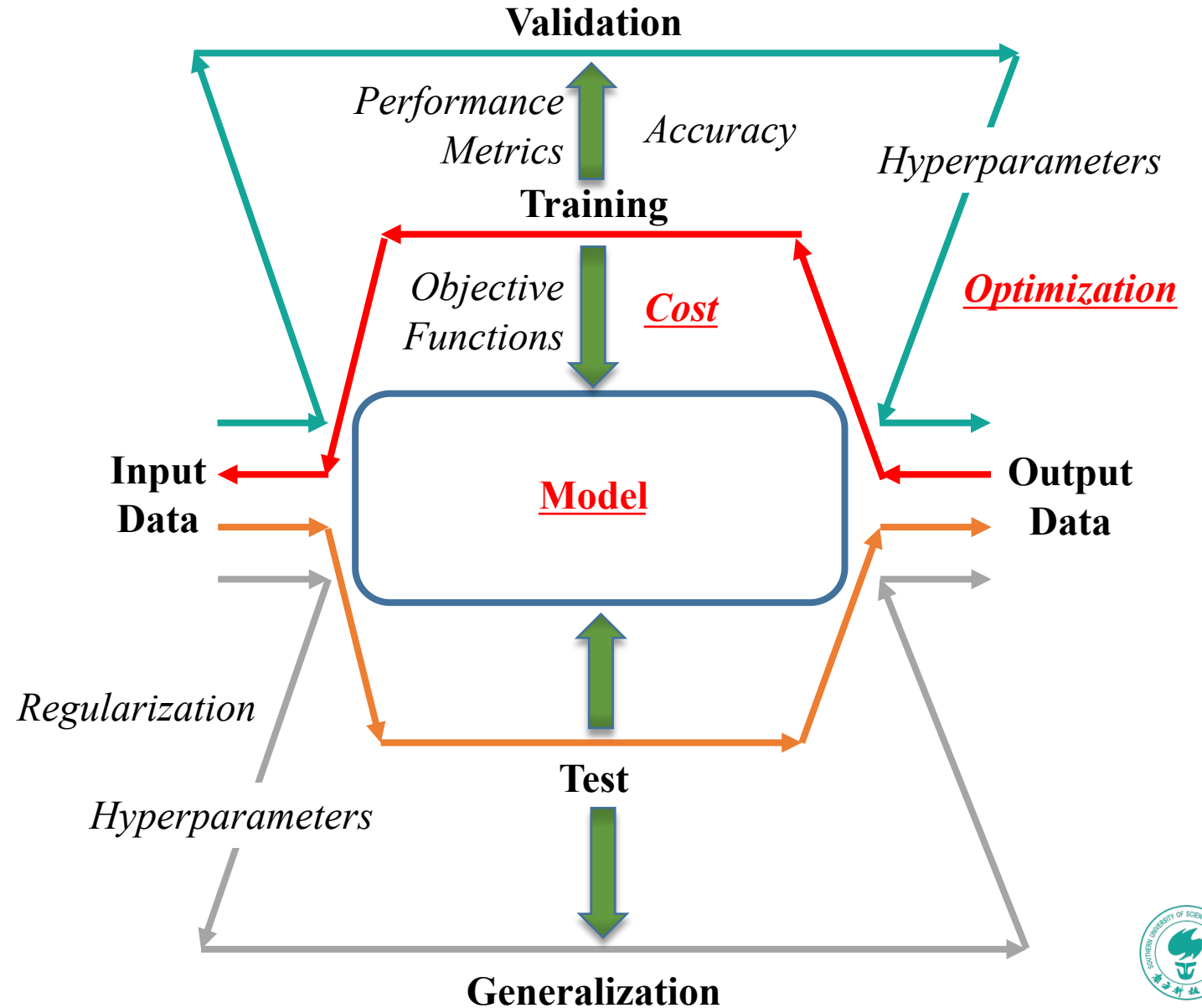
The most time-consuming step ...

- *Modify your model, train it, evaluate on your validation data (not the test data, at this point),*
- *Modify it again, and repeat, until the model is as good as it can get*

Things you can try

- Add dropout
- Try different architectures: add or remove layers
- Add L1 and/or L2 regularization
- Try different hyperparameters
 - such as the number of units per layer, or
 - the learning rate of the optimizer
- Optionally, iterate on feature engineering
 - add new features, or
 - remove features that don't seem to be informative

AncoraSIR.com



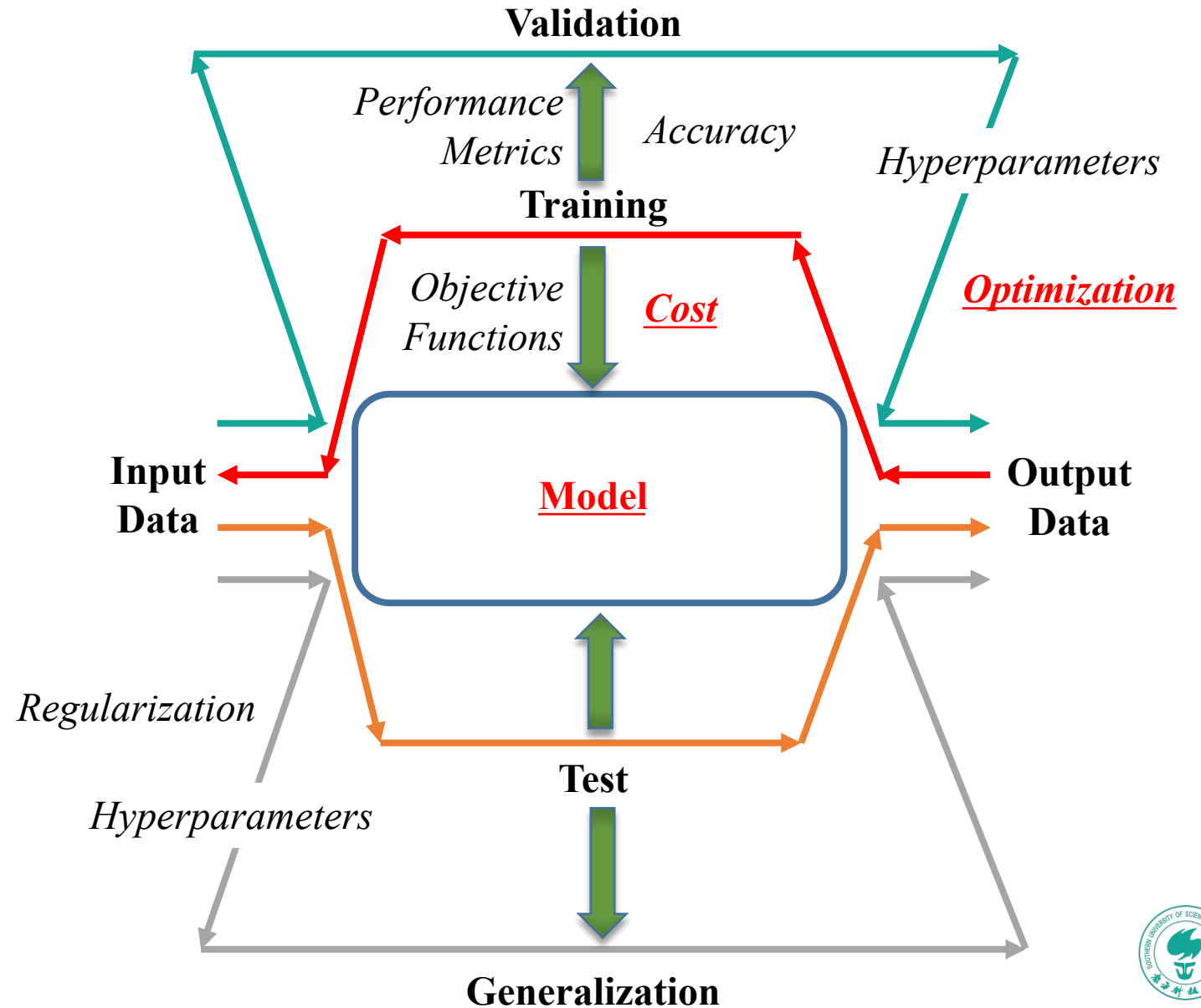
Tuning Your Hyperparameters

Information Leak => Model Overfit

- Every time you use feedback from your validation process to tune your model, you leak information about the validation process into the model.
- This makes the evaluation process less reliable.

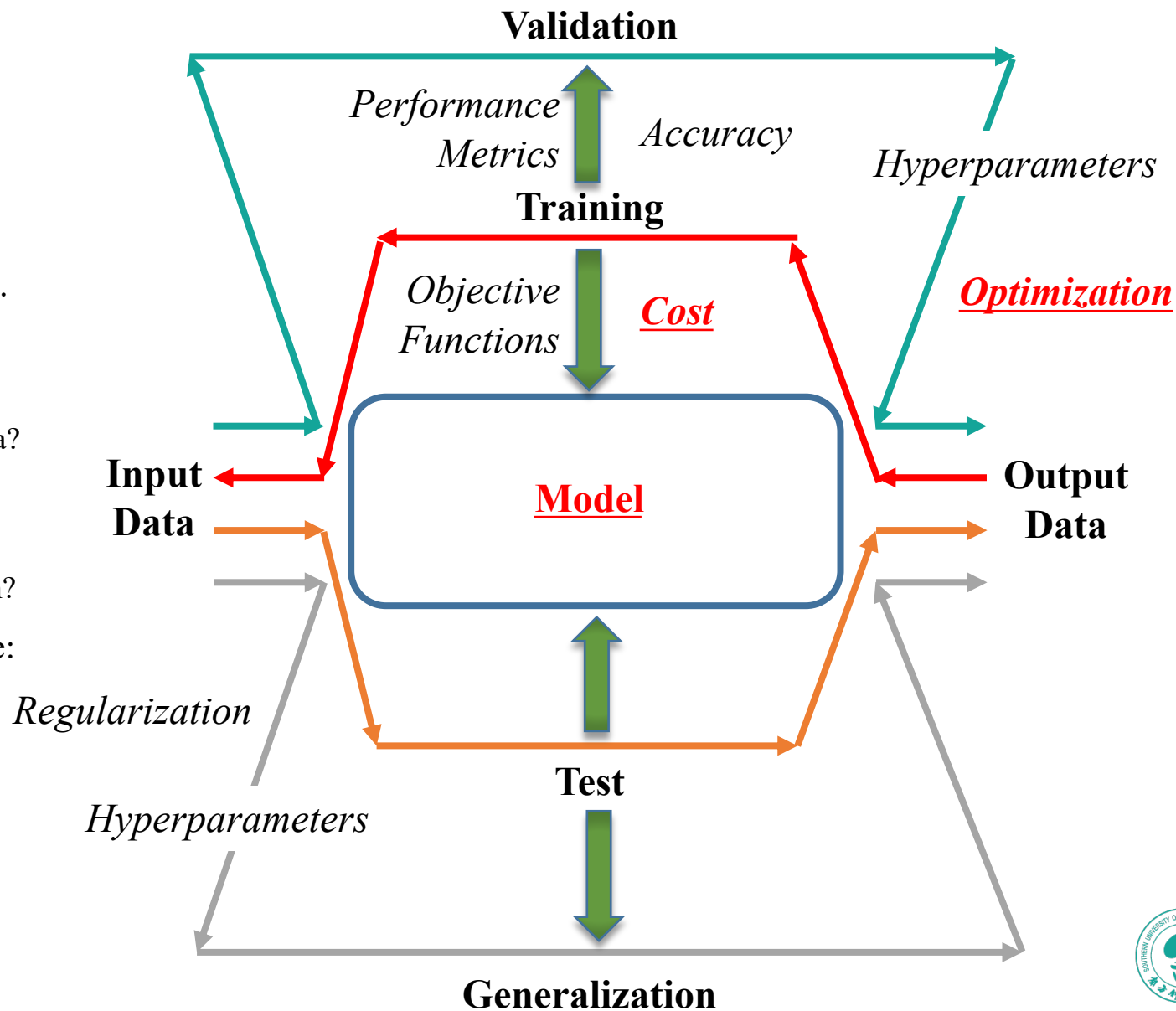
Once you've developed a satisfactory model configuration,

- you can train your final production model on all the available data (training and validation) and
- evaluate it one last time on the test set
- If good, then proceed
- Or if worse, you may want to switch to a more reliable evaluation protocol (may be caused by overfitting)

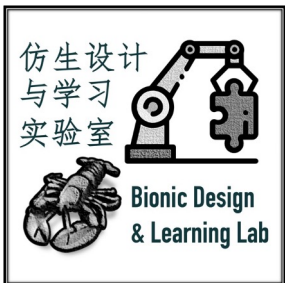


A Universal Workflow

- Define the problem and the data on which you'll train.
 - Collect this data or annotate it with labels if need be.
- Choose how you'll measure success on your problem.
 - Which metrics will you monitor on your validation data?
- Determine your evaluation protocol:
 - Hold-out validation? K-fold validation?
 - Which portion of the data should you use for validation?
- Develop a model that does better than a basic baseline:
 - A model with statistical power.
- Develop a model that overfits.
- Regularize your model and tune its hyperparameters,
 - Based on performance on the validation data.
 - A lot of machine-learning research tends to focus only on this step—but keep the big picture in mind.



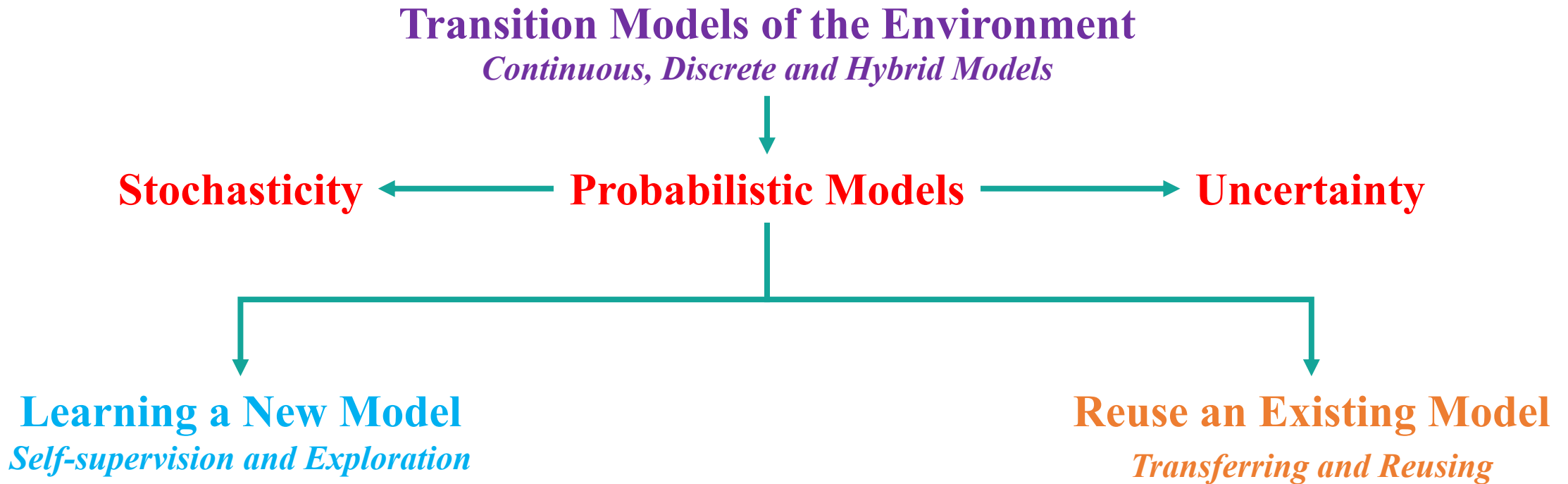
Learning a Transition Model



AncoraSIR.com

Learning a Transition Model of the Environment

The robot must learn a model of how its actions affect the task state, and the resulting background cost, for use in planning



Continuous Transition Models

*Regression methods can be used to learn **low-level** transition models for predicting the next state as a set of continuous values, even when the set of actions is discrete*

Scoop & Dump Task $l(h_0, a_0, \dots, a_T, h_g) = \|\mathcal{F}(h_0, a_0, \dots, a_T) - h_g\|_1$

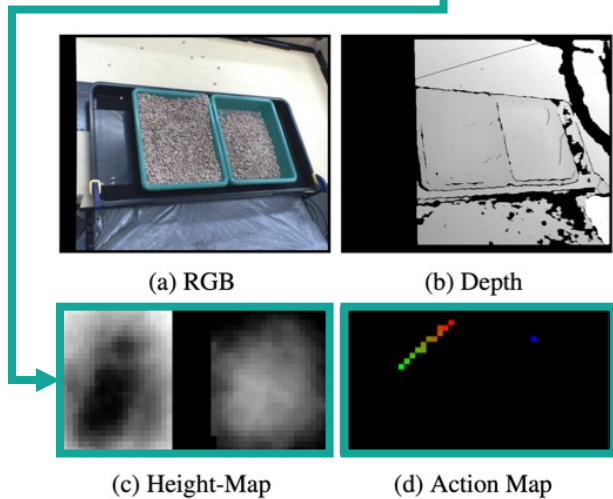
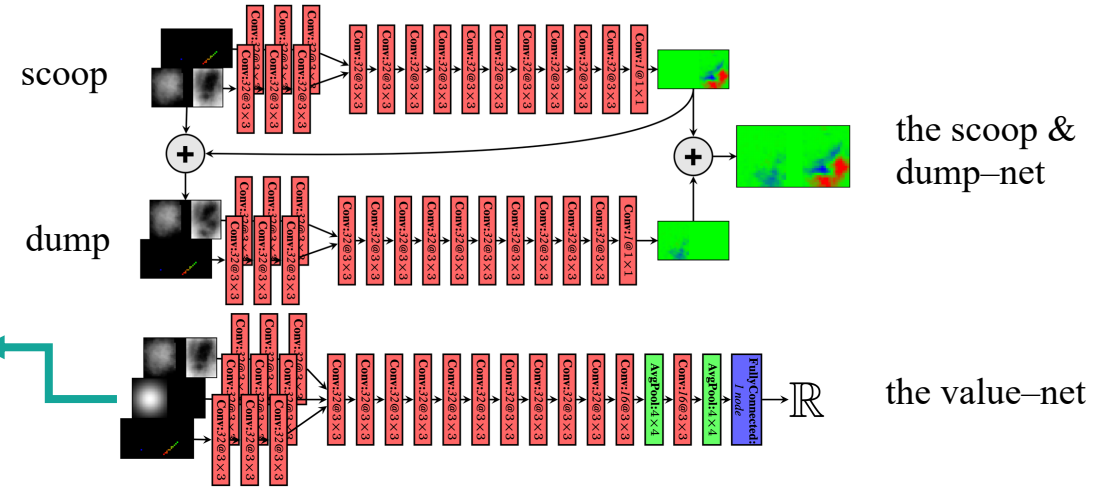
h_0 : a given *initial state* of the environment

h_g : a given *goal state* of the environment

a_0, \dots, a_T : a series of *robot actions*

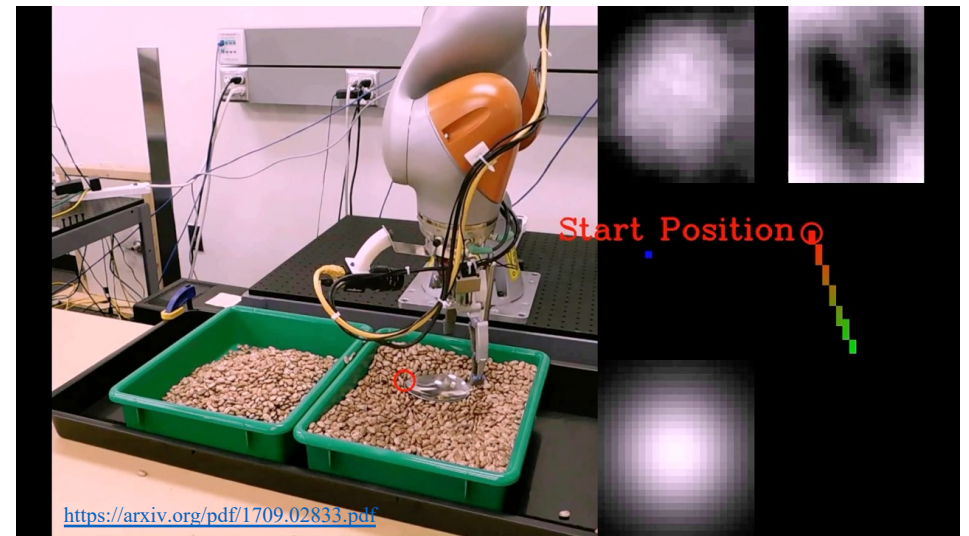
$\|Norm\|_1$: L1 norm to be minimized as distance to the goal state

$\mathcal{F}(h_0, a_0, \dots, a_T) = h_{T+1}$ applies actions sequentially to reach h_{T+1}



Robot action (a 9D vector)

- *Scoop action*
 - the start location (2D)
 - the start angle (1D)
 - the end location (2D)
 - the end angle (1D)
 - the roll angle (1D)
- *Dump action*
 - the dump location (2D)





DES 5002: Designing Robots for Social Good

Autumn 2022

Thank you~

Wan Fang

Southern University of Science and Technology