

Multimodal AI

Lecture 2.2 – Data & Heterogeneity

Paul Liang

Assistant Professor

MIT Media Lab & MIT EECS



<https://pliang279.github.io>

ppliang@mit.edu

 [@pliang279](https://twitter.com/pliang279)



Assignments for This Coming Week

HW1: reading assignment + homework due next Tuesday 2/17

For project:

- Project proposal instructions released, due Tuesday (2/24). Submit on canvas
- Meet with me 4-5pm if need feedback about proposal ideas.

Clarification: 2-4 students per group. Listeners can help out on projects and will not be counted in the group size.

No 1 student groups, unless exceptional circumstances.

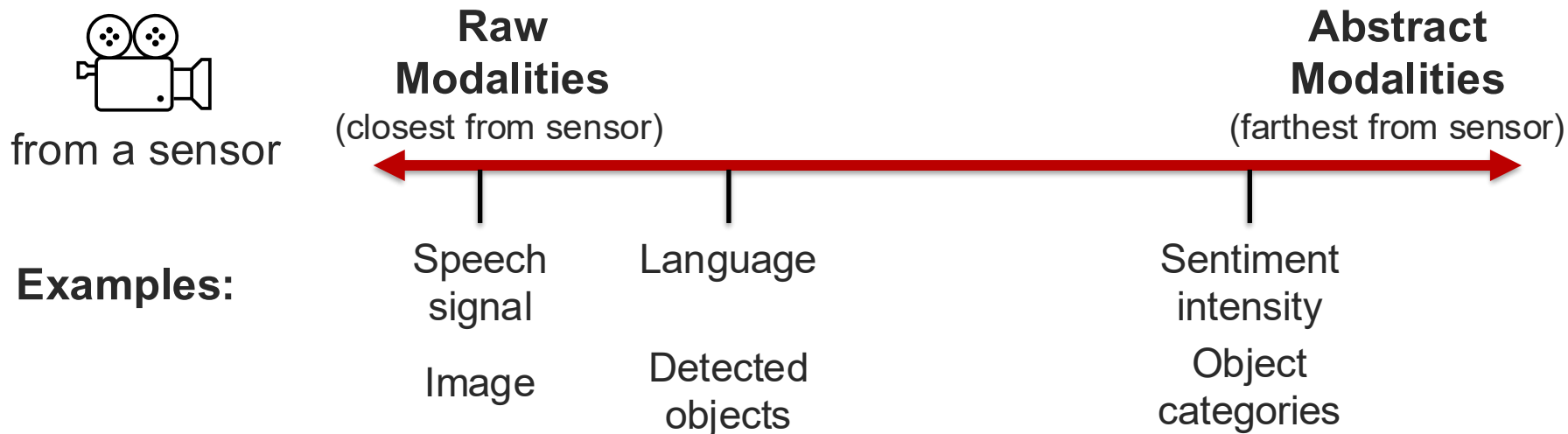
Lecture Outline

- 1 Vision, language, audio, sensing, set, graph modalities
- 2 Modality profile
- 3 Types of data and labels
- 4 Common learning objectives and generalization

What is a Sensory Modality?

Sensory modality

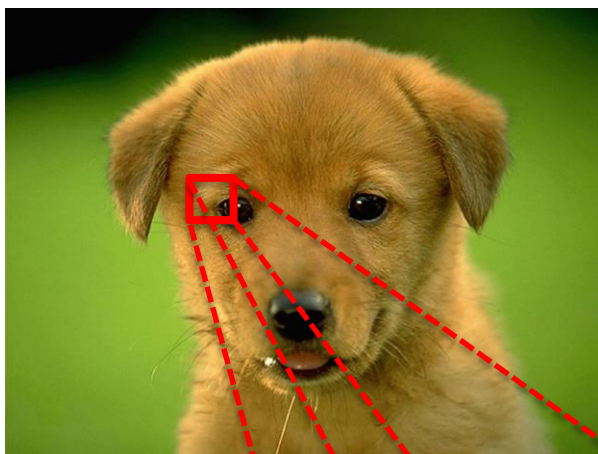
Modality refers to the way in which something expressed or perceived.



Most of AI is about learning abstractions, or representations, from data.

Visual Modality

Color image



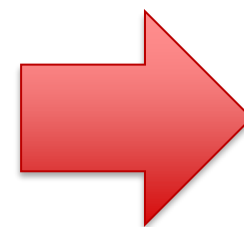
Each pixel is represented in \mathcal{R}^d , d is the number of colors ($d=3$ for RGB)

88	82	84	88	85	83	80	93	102
88	80	78	80	80	78	73	94	100
85	79	80	78	77	74	65	91	99
38	35	40	35	39	74	77	70	65
20	25	23	28	37	69	64	60	57
22	26	22	28	40	65	64	59	34
24	28	24	30	37	60	58	56	66
21	22	23	27	38	60	67	65	67
23	22	22	25	38	59	64	67	66

Input observation x_i

88
88
85
38
20
22
24
21
23
82
80
79
35
25
26
28
22
22
84
78
80
⋮

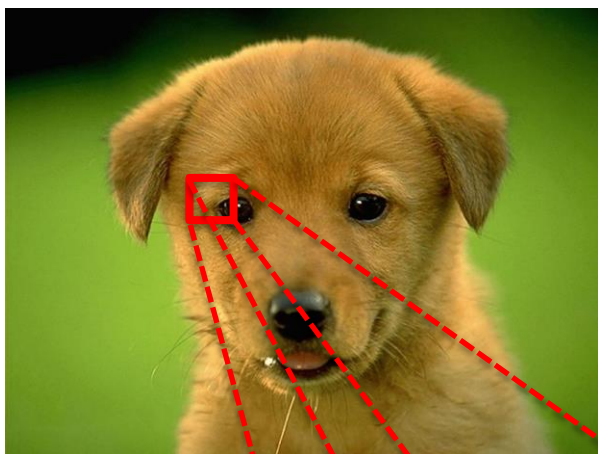
Binary classification problem



Dog ?

label $y_i \in \mathcal{Y} = \{0,1\}$

Visual Modality



Each pixel is represented in \mathcal{R}^d , d is the number of colors ($d=3$ for RGB)

88	82	84	88	85	83	80	93	102
88	80	78	80	80	78	73	94	100
85	79	80	78	77	74	65	91	99
38	35	40	35	39	74	77	70	65
20	25	23	28	37	69	64	60	57
22	26	22	28	40	65	64	59	34
24	28	24	30	37	60	58	56	66
21	22	23	27	38	60	67	65	67
23	22	22	25	38	59	64	67	66

Input observation x_i

88
88
85
38
20
22
24
21
23
82
80
79
35
25
26
28
22
22
84
78
80
⋮

Multi-class classification problem

Duck

-or-

Cat

-or-

Dog

-or-

Pig

-or-

Bird ?

label $y_i \in \mathcal{Y} = \{0,1,2,3, \dots\}$

Language Modality

Written language

★★★★★ **Masterful!**

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humourous manner.

0 of 4 people found this review helpful

Spoken language

MARTHA (CON'T)

Look around you. Look at all the great things you've done and the people you've helped.

CLARK

But you've only put up the good things they say about me.

MARTHA

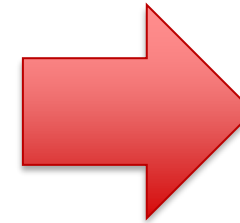
Clark, honey. If I were to use the bad things they say I could cover the barn, the house and the outhouse.

Input observation x_i

0
1
0
0
1
0
1
0
0
0
0
1
0
0
0
0
1
0
0
0
⋮

“bag-of-words” vector

$|x_i|$ = number of words in dictionary



Document-level
classification

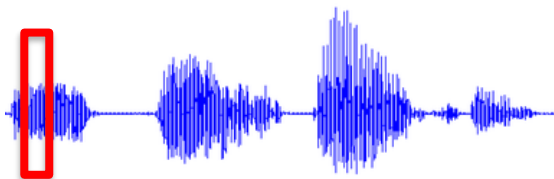
Sentiment ?
(positive or negative)

Response?

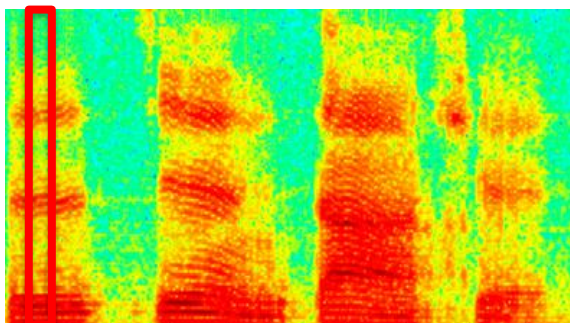
What happens with word ordering?

Acoustic Modality

Digitalized acoustic signal



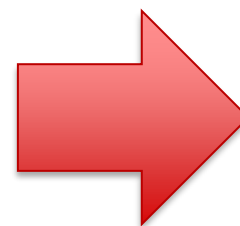
- Sampling rates: 8~96kHz
- Bit depth: 8, 16 or 24 bits
- Time window size: 20ms
 - Offset: 10ms



Spectrogram

Input observation x_i

0.21
0.14
0.56
0.45
0.9
0.98
0.75
0.34
0.24
0.11
0.02



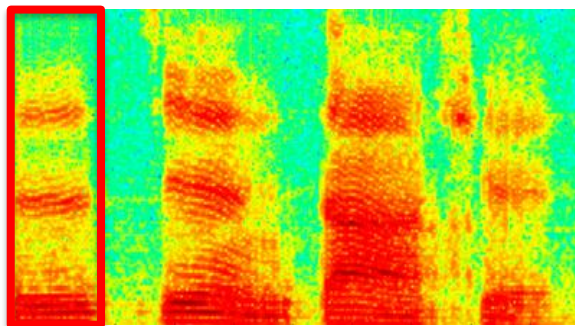
Spoken word ?

Acoustic Modality

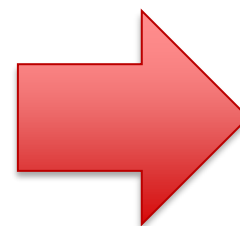
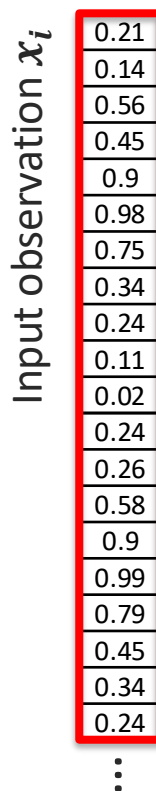
Digitalized acoustic signal



- Sampling rates: 8~96kHz
- Bit depth: 8, 16 or 24 bits
- Time window size: 20ms
 - Offset: 10ms



Spectrogram

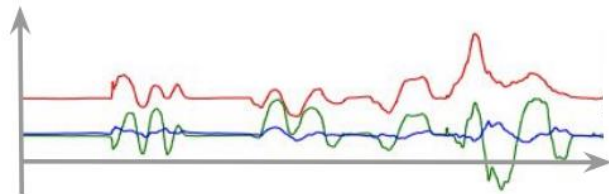


Emotion ?

Spoken word ?

Voice quality ?

Sensor Modality



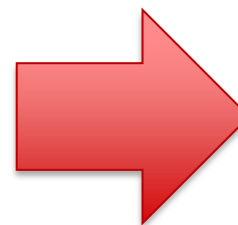
Time series data across six-axis Force-Torque sensor:
 $T \times 6$ signal.

Force-Torque Sensor



Proprioception

Measure values internal to the system (robot);
e.g. motor speed, wheel load, **robot arm joint angles**, battery voltage.

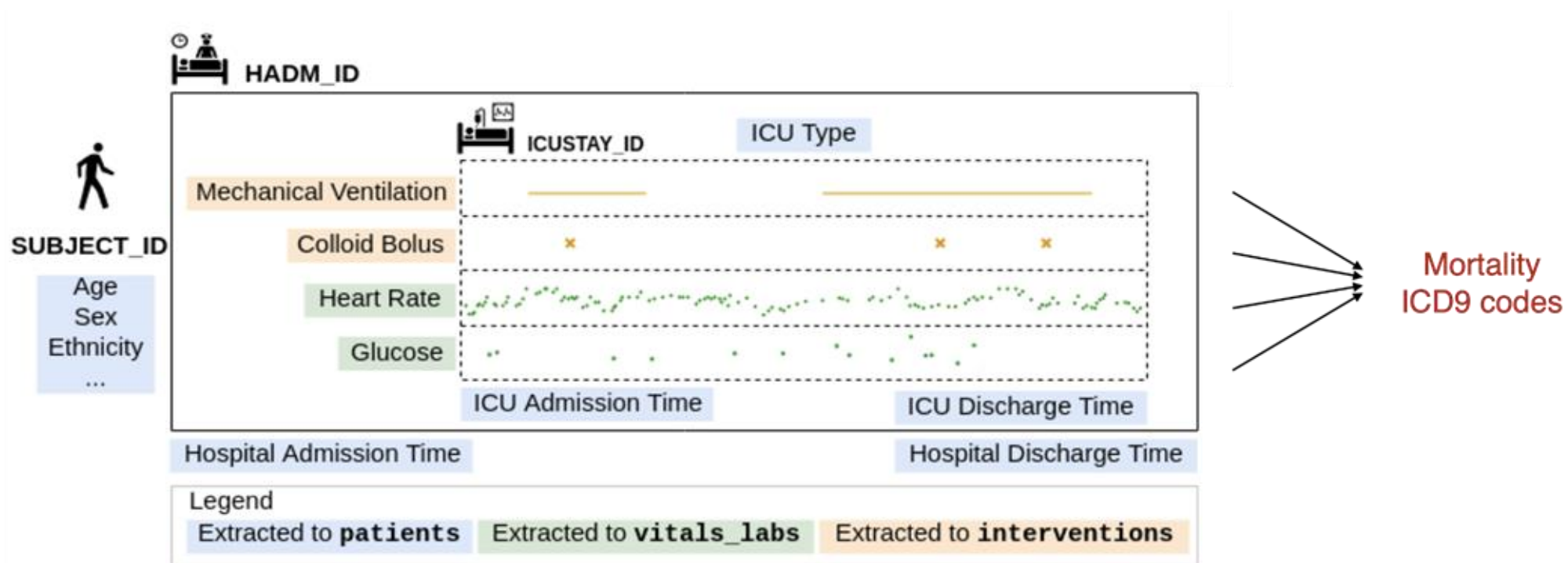


Time series data across current position and velocity of the end-effector:
 $T \times 2$ signal.



Object property
Next action

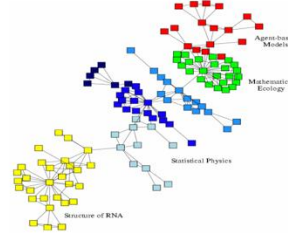
Tabular Modality



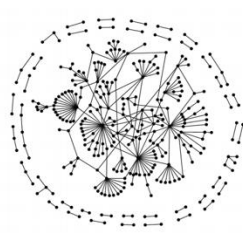
Graph Modality



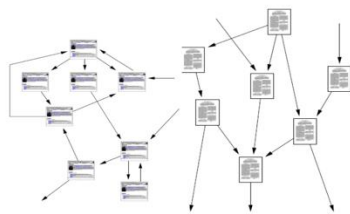
Social networks



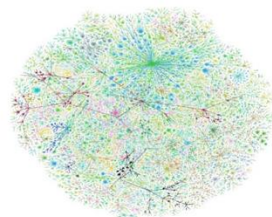
Economic networks



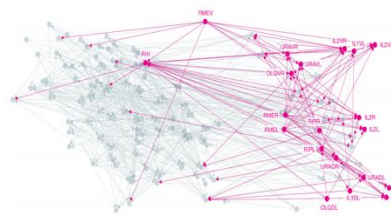
Biomedical networks



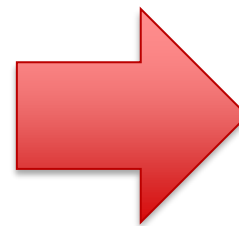
Information networks:
Web & citations



Internet



Networks of neurons



Tasks on graphs:
Node classification
Link prediction

...

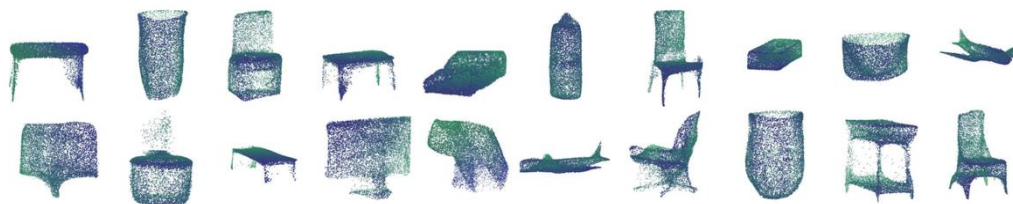
Using graphs:
Knowledge graphs for
QA
Social network for
sentiment analysis

...

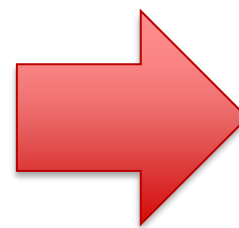
Set Modality



Sets



Point clouds



- Set anomaly detection
- Set expansion
- Set completion
- Point cloud classification
- Point cloud generation

Modality Profile

The qualities and structures that are unique to a data modality.



*A teacup on the right of a laptop
in a clean room.*

Modality Profile

The distribution of individual elements within that modality.



A *teacup* on the *right* of a *laptop*
in a *clean room*.

① **Distribution:** discrete or continuous, support



● *{teacup, right, laptop, clean, room}*

Modality Profile

The frequency at which elements appear or are sampled.



*A teacup on the right of a laptop
in a clean room.*

2 **Granularity:** sampling rate and frequency



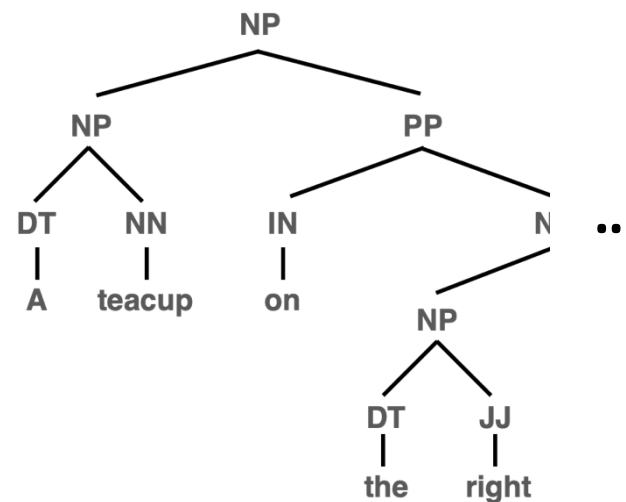
objects per image



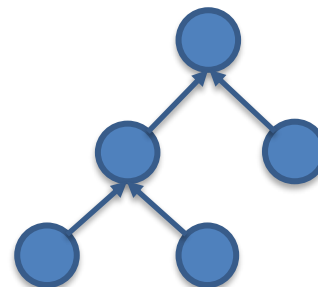
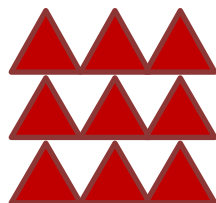
words per minute

Modality Profile

The way elements compose with each other to form entire data.



3 **Structure:** static, temporal, spatial, hierarchical



Modality Profile

The total information contained in the elements and their composition.

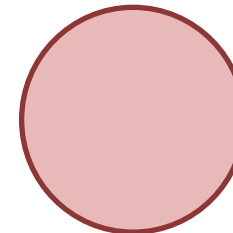


*A teacup on the right of a laptop
in a clean room.*

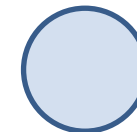
4

Information: entropy and density

$H(\blacktriangle)$



$H(\bullet)$



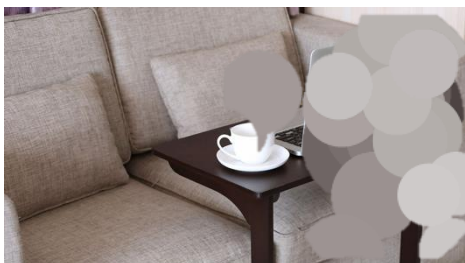
Modality Profile

The natural imperfections in the data modality.



*A teacup on the right of a laptop
in a clean room.*

5 **Noise:** uncertainty, signal-to-noise ratio, missing data

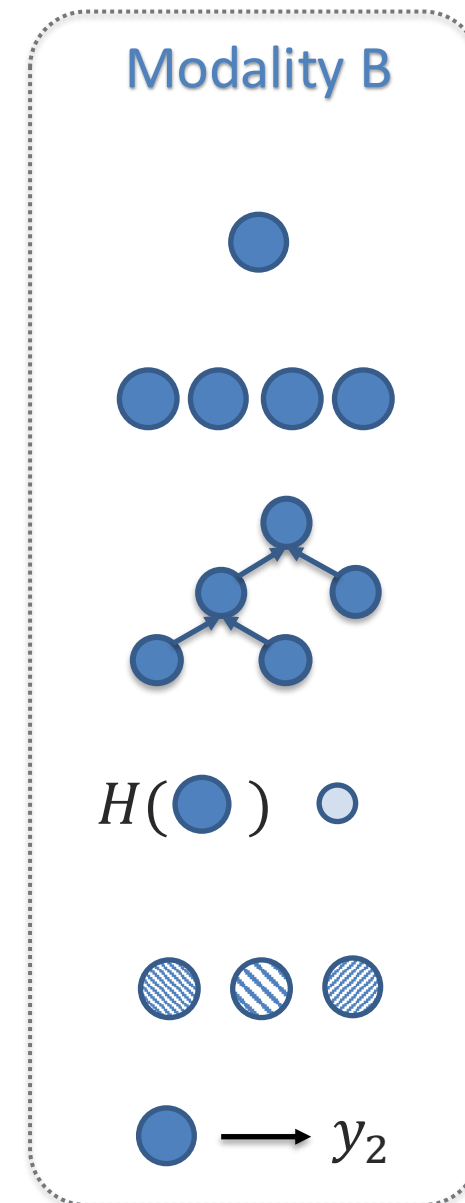
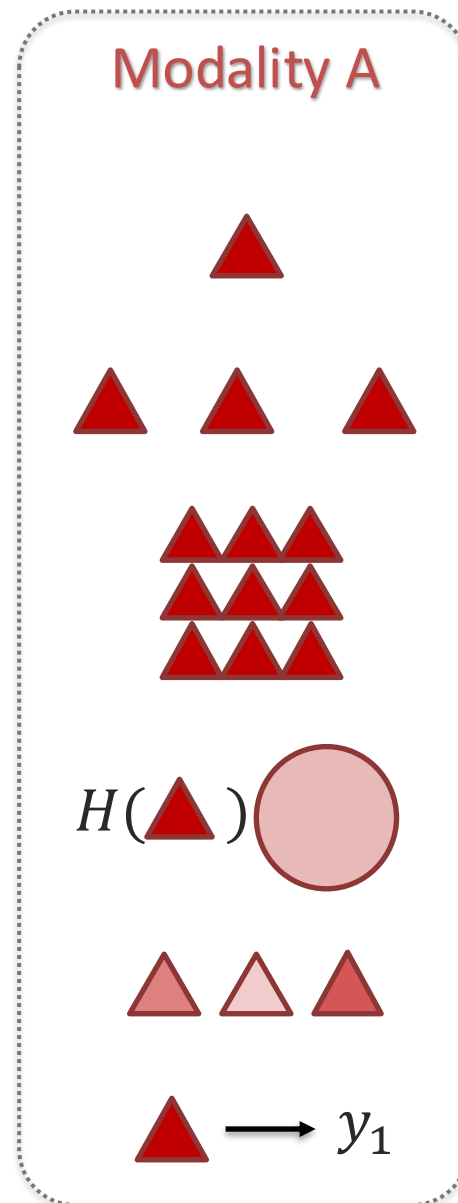


teacup → **teacip**

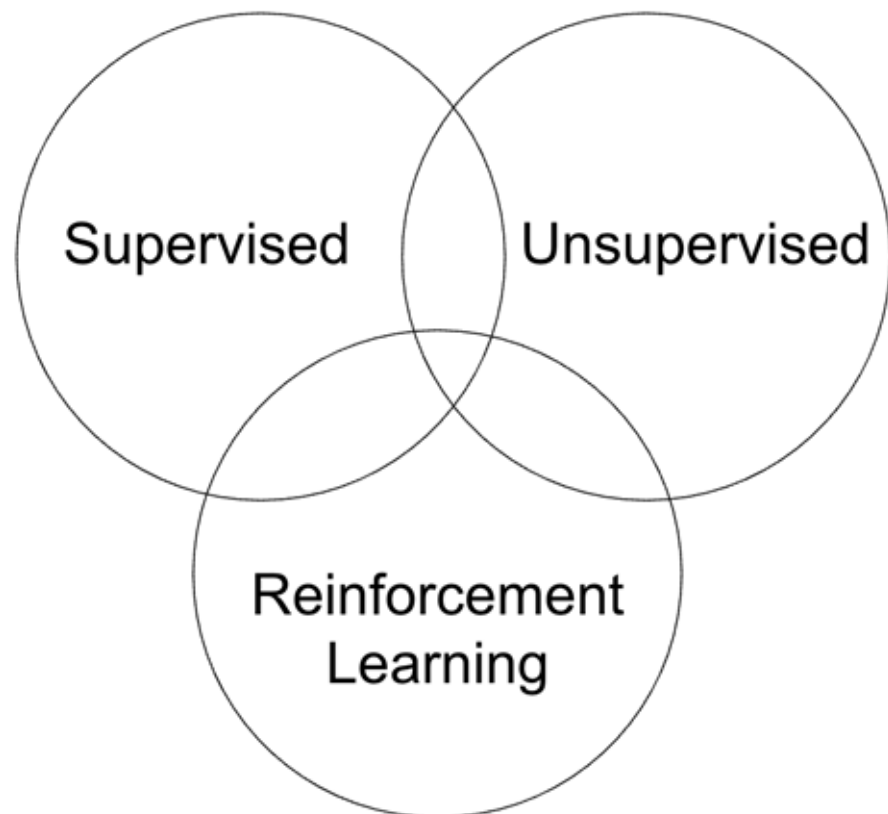
right → **rihjt**

Modality Profile

- 1 **Element representations:**
Discrete, continuous, granularity
- 2 **Element distributions:**
Density, frequency
- 3 **Structure:**
Temporal, spatial, latent, explicit
- 4 **Information:**
Abstraction, entropy
- 5 **Noise:**
Uncertainty, noise, missing data
- 6 **Relevance:**
Task, context dependence



Types of Learning Paradigms



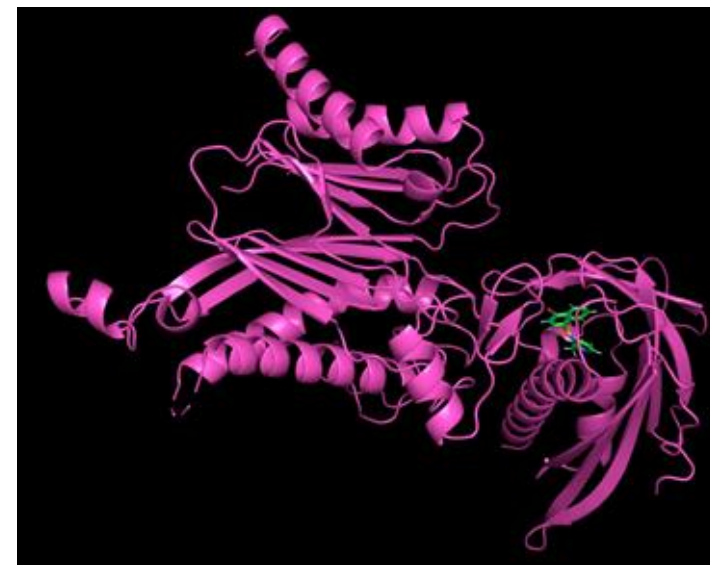
(the categorization can be refined, e.g. there are active learning, semi-supervised, selective, contrastive, few-shot, inverse reinforcement learning...)

Supervised Learning

Goal: correctly classify so far unseen test images



Goal: predict to what degree a drug candidate binds to the intended target protein (based on a dataset of already-screened molecules against the target)



- Learning a machine translation system from pairs of sentences

Spanish (input)

Aquí tienes un bolígrafo

Las conferencias de ML son divertidas

Todo el mundo debería estudiar AI

English (output)

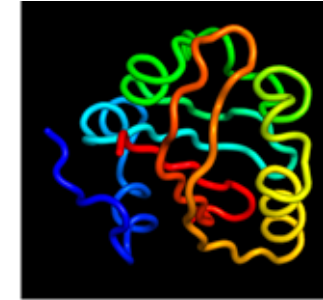
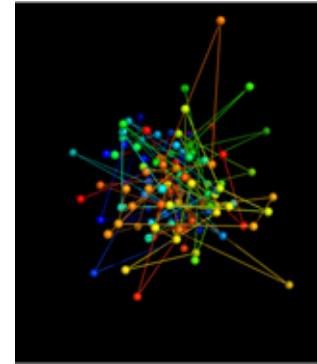
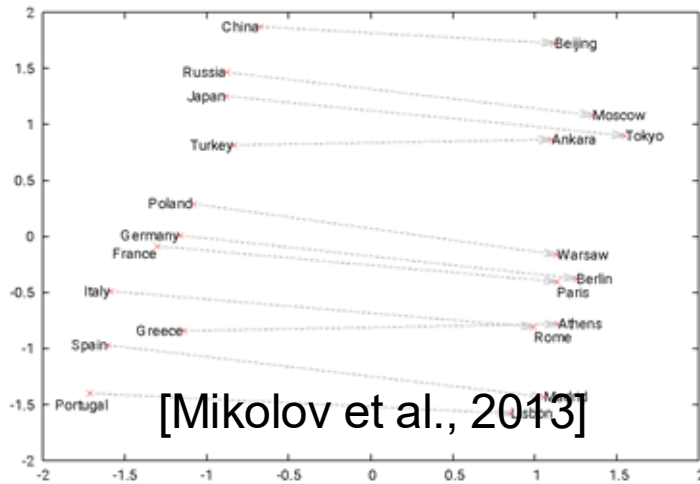
Here's a pen

ML conferences are fun

Everyone should study AI

Unsupervised Learning

dimensionality reduction, embedding



[courtesy of Jason Yim]

Over 3D protein structures, etc.

**+Self-Supervised
paradigm**

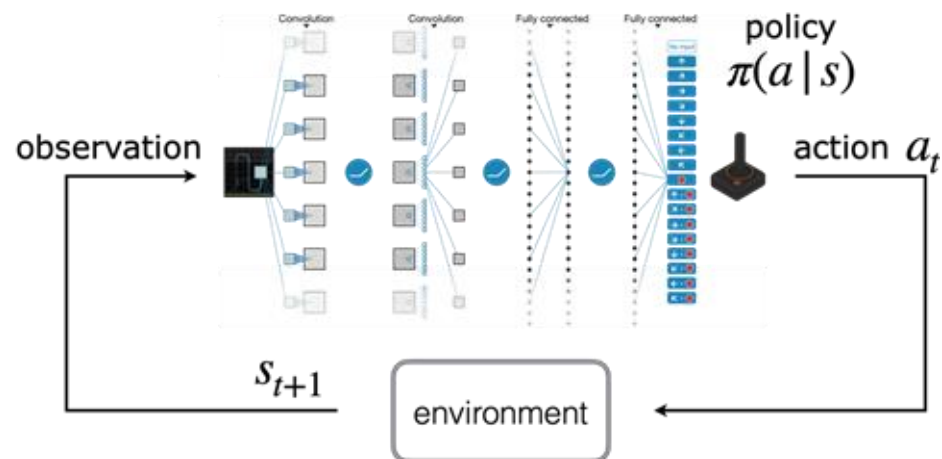
de-noising diffusion models over images



[image from
Rissanen et al 2022]

[Slides adapted from 6.790]

Reinforcement Learning



Step 1

Collect demonstration data and train a supervised policy.

A prompt is sampled from our prompt dataset.

Explain reinforcement learning to a 6 year old.

A labeler demonstrates the desired output behavior.

We give treats and punishments to teach...

This data is used to fine-tune GPT-3.5 with supervised learning.

SFT

Step 2

Collect comparison data and train a reward model.

A prompt and several model outputs are sampled.

Explain reinforcement learning to a 6 year old.

A In reinforcement learning, the agent is...
B Explain rewards...
C In machine learning...
D We give treats and punishments to teach...

A labeler ranks the outputs from best to worst.

D > C > A > B

This data is used to train our reward model.

RM
D > C > A > B

ChatGPT

Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

A new prompt is sampled from the dataset.

Write a story about otters.

The PPO model is initialized from the supervised policy.

PPO

The policy generates an output.

Once upon a time...

The reward model calculates a reward for the output.

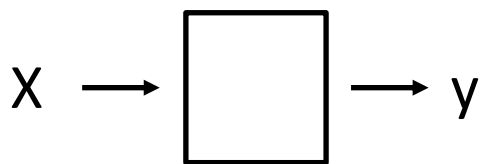
RM

The reward is used to update the policy using PPO.

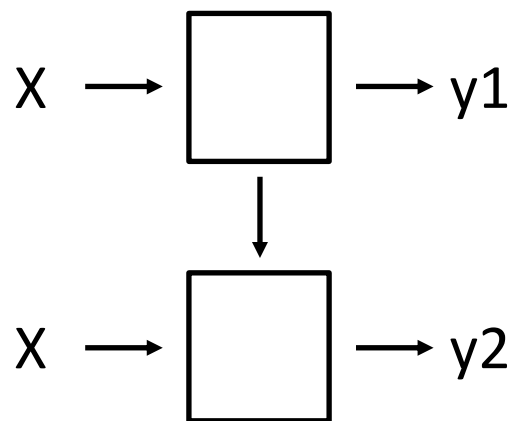
r_k

More Learning Paradigms

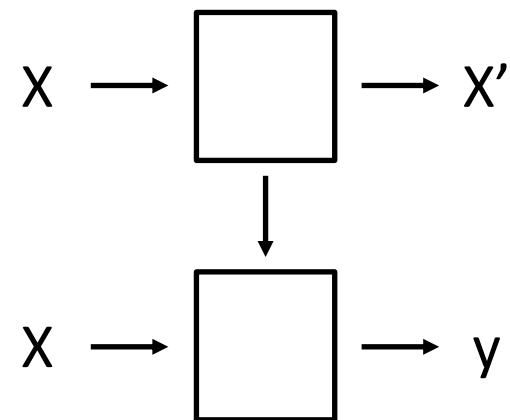
Supervised learning



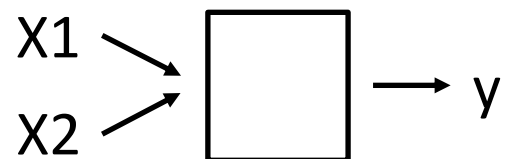
Transfer learning



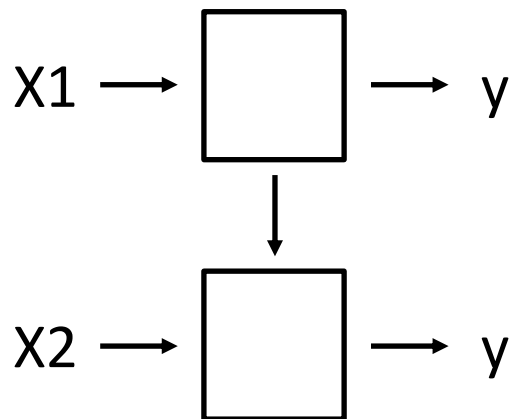
Unsupervised/self-supervised pre-training



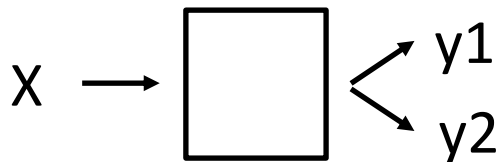
Multimodal (supervised) learning



Cross-modal learning

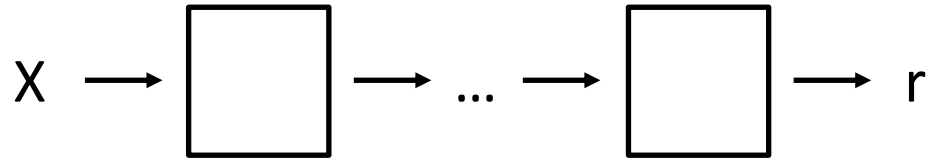


Multitask learning

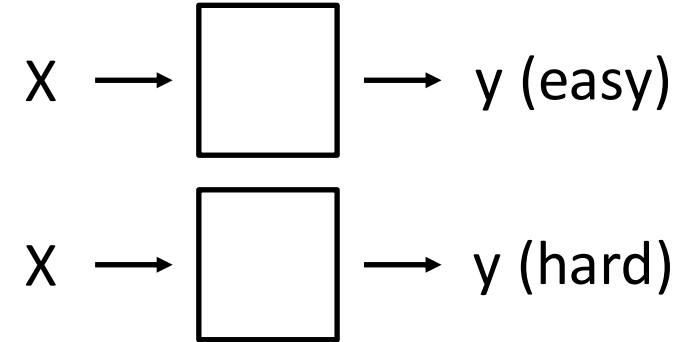


More Interactive Learning Paradigms

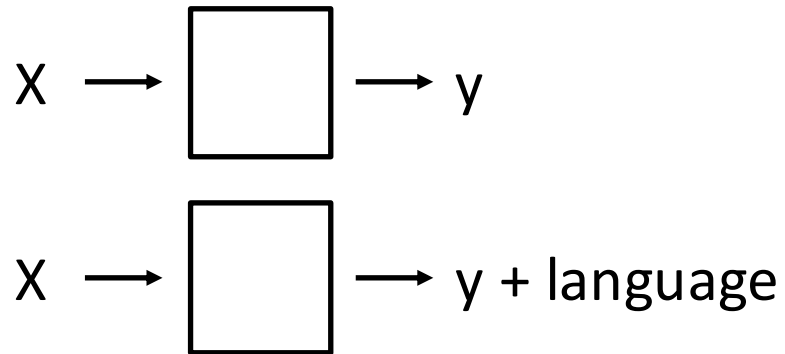
Reinforcement learning



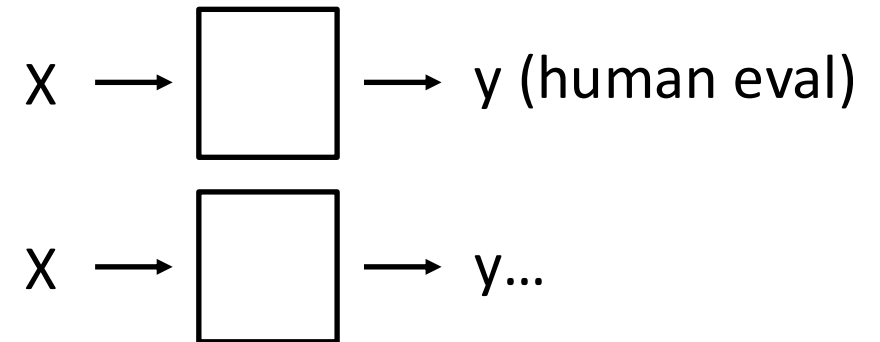
Curriculum/active learning



LLM adaptation



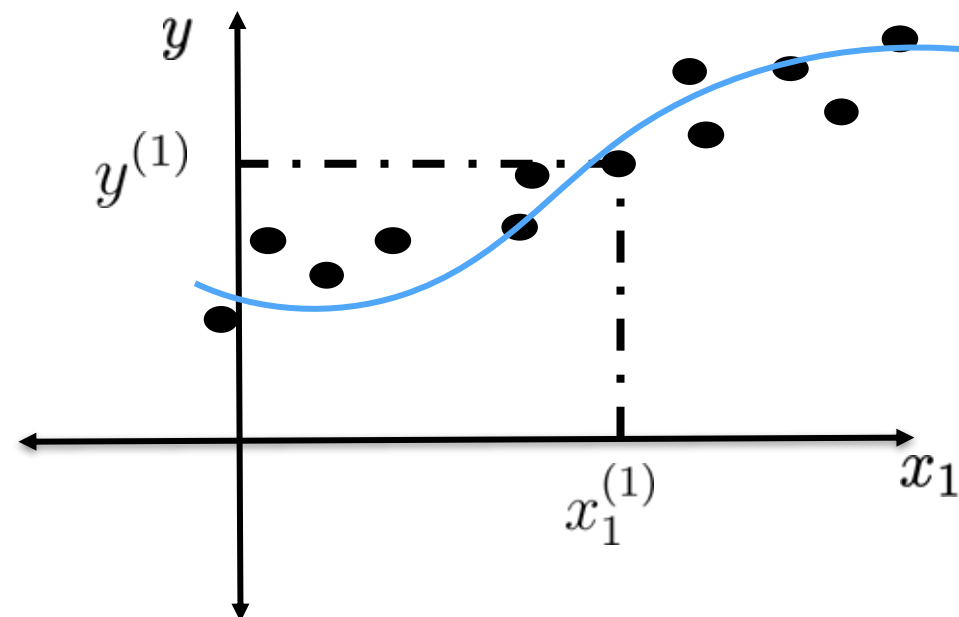
Human-in-the-loop learning



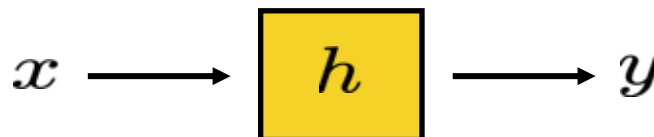
Learning Process

We want a “good” way to label our data

- How to label? Learn a model
- We typically consider a class of possible models



Input:
Data



Output:
Label

how well our model labels new data depends largely on how good the chosen model class is

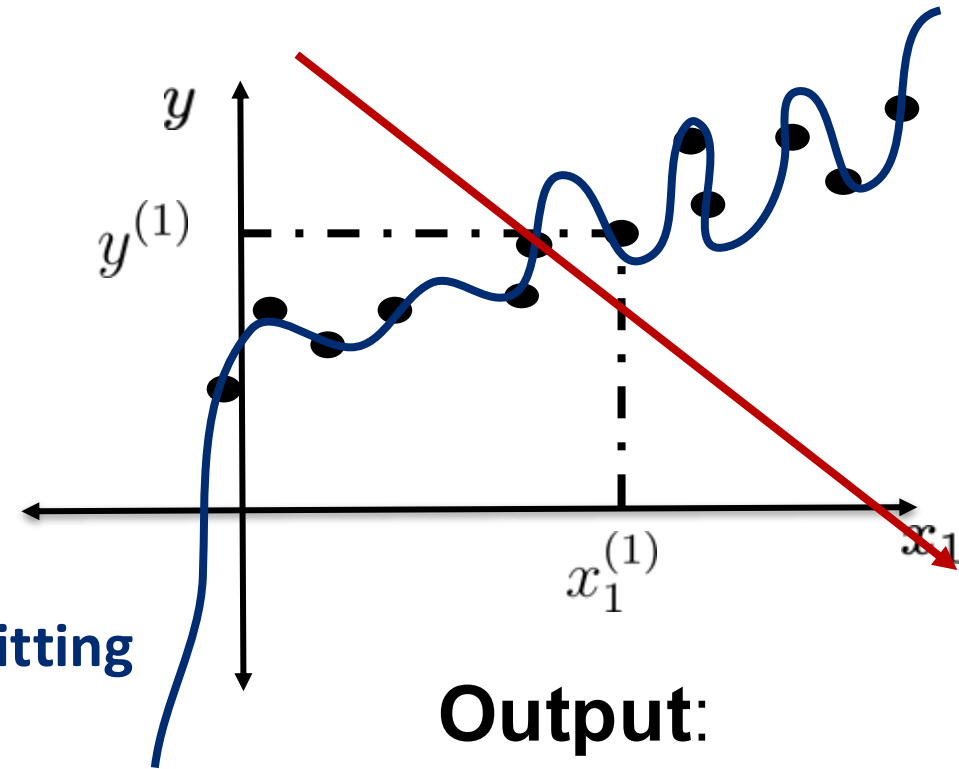
Overfitting vs Generalization

What we really want is to generalize to **future data!**

What we don't want:

- Model does not capture the input-output relationship → **Underfitting**
- Model too specialized to training data → **Overfitting**

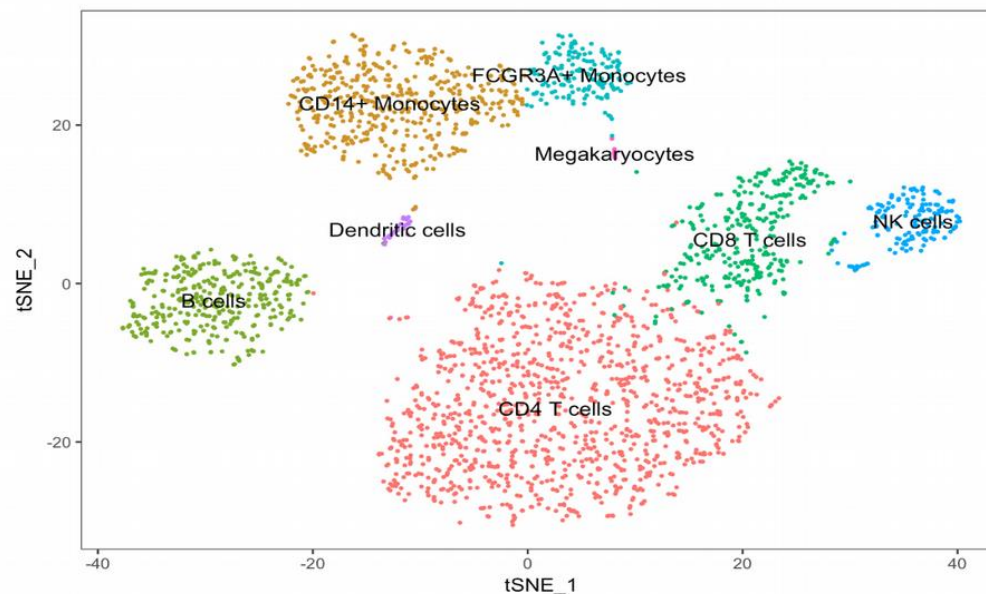
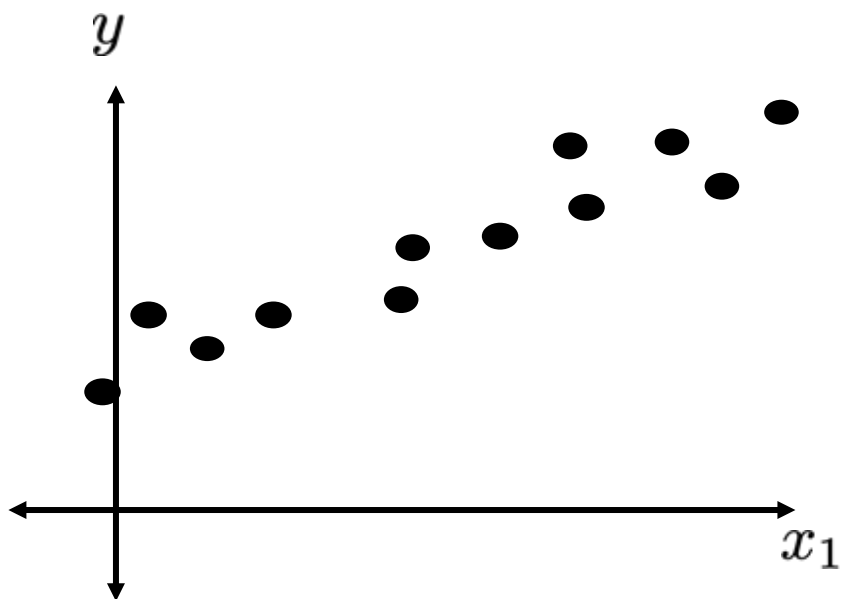
Split collected data into training, validation, and testing.
Critical to make sure test data conditions match real-time deployment conditions.



Output:
Label

Summary: How To Data

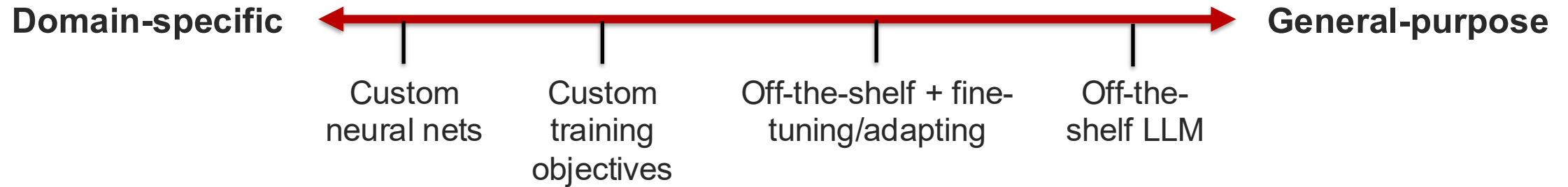
1. Decide how much data to collect, and how much to label (costs and time)
2. Clean data: normalize/standardize, find noisy data, anomaly/outlier detection
3. Visualize data: plot, dimensionality reduction (PCA, t-sne), cluster analysis
4. Decide on evaluation metric (proxy + real, quantitative and qualitative)
5. Choose model class and learning algorithm (more coming up)



Lecture Outline

- 1 A unifying paradigm of model architectures
- 2 Temporal sequence models
- 3 Spatial convolution models
- 4 Models for sets and graphs

Two General Modeling Paradigms



Your decision will depend on many factors.

Designing Models for Data

What is a good model?

One that captures the:

- right semantic information
- at the right granularity
- using an appropriate amount of data and labels
- with the right resource constraints
- with the right level of usability (explainability, accessibility, etc.)
- and more...

Domain-specific



General-purpose

Unified View of Deep Learning Models

1. Learning representations



Model



2. Combining representations (information aggregation)

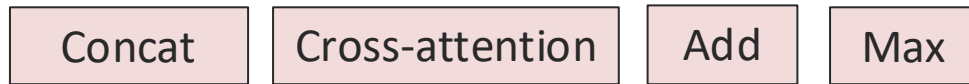
Unified View of Deep Learning Models

Composing differentiable functions and training objectives.

1. Basic representation building blocks for each element

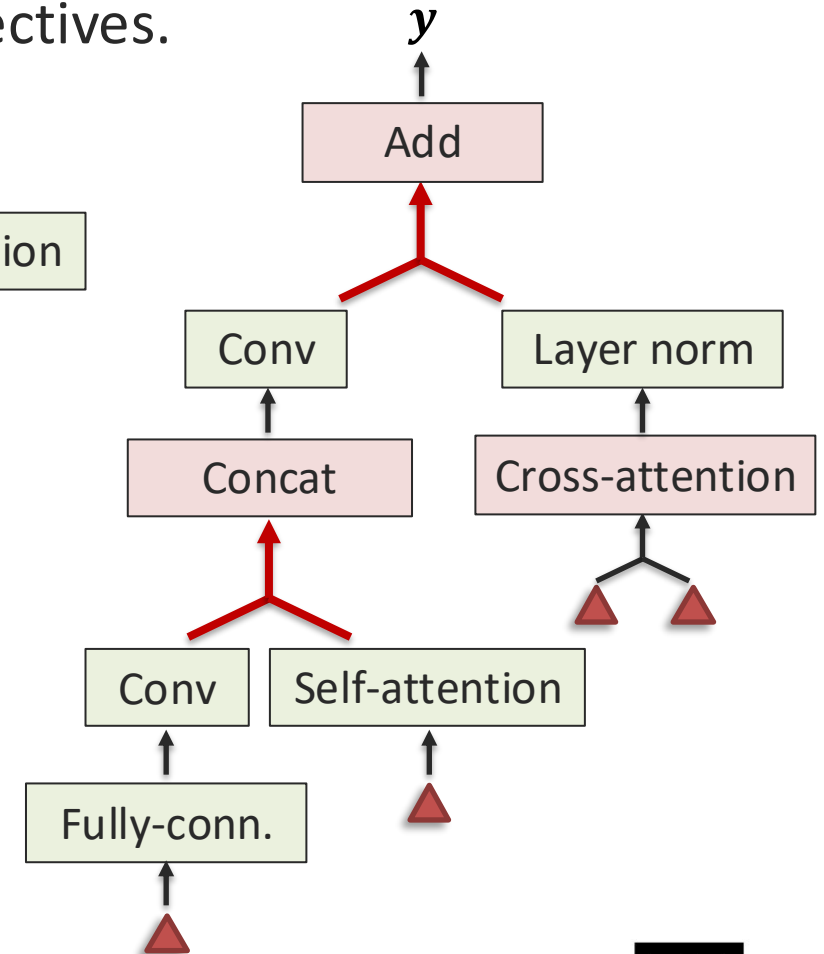


2. Basic information aggregation blocks



3. Compute loss function

4. Take gradients, update with stochastic gradient descent

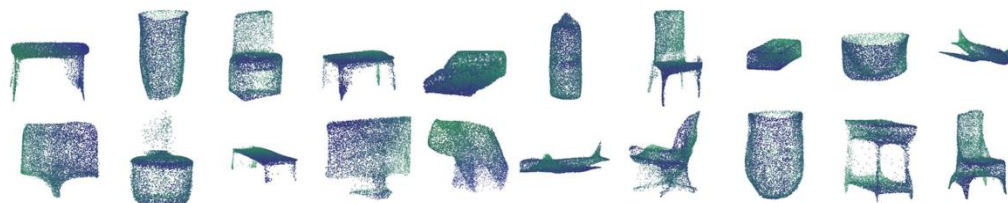


A Simple Classification Example

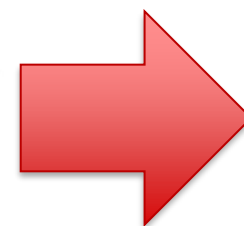
Sets and point clouds



Sets



Point clouds

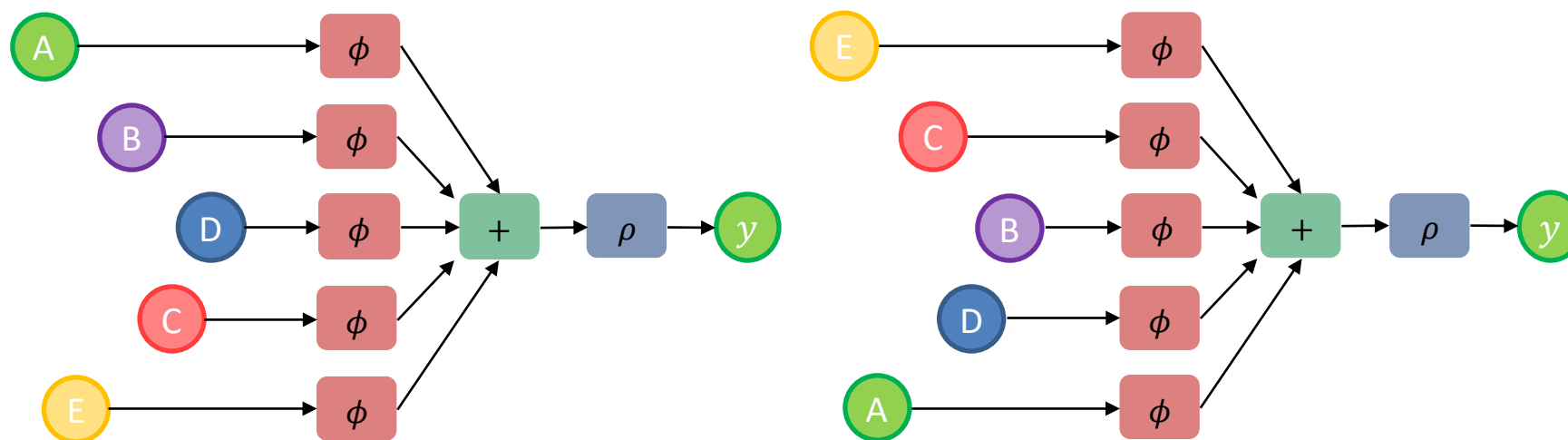


Set anomaly detection
Set expansion
Set completion
Point cloud classification
Point cloud generation

A Simple Classification Example

Models for set-based data must be invariant to element order.

1. Parameter sharing for each set element
2. Permutation invariant aggregation function

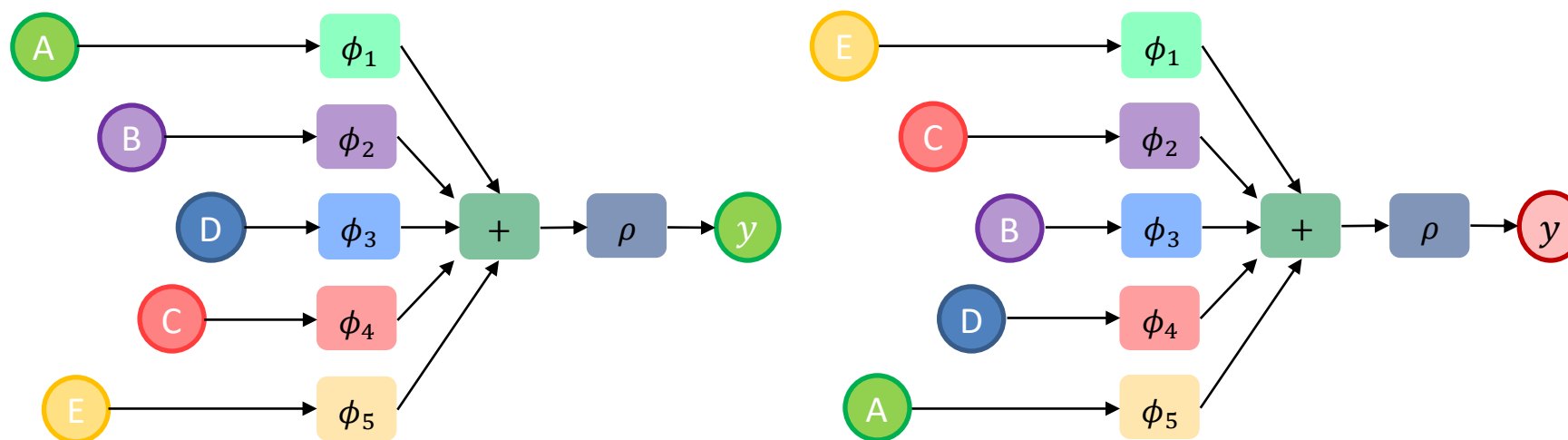


Giving ABDCE also gives ECBDA, BCAED etc...

A Simple Classification Example

Models for set-based data must be invariant to element order.

1. No parameter sharing for each set element
2. Permutation invariant aggregation function

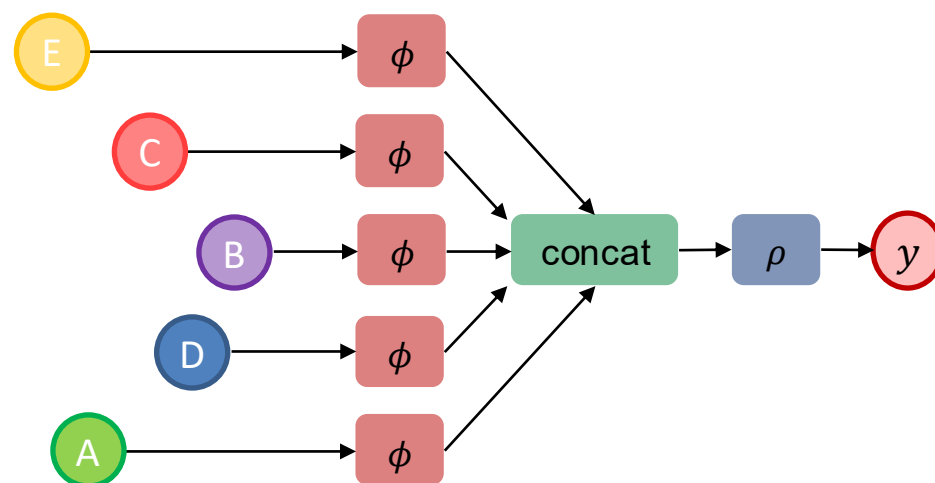
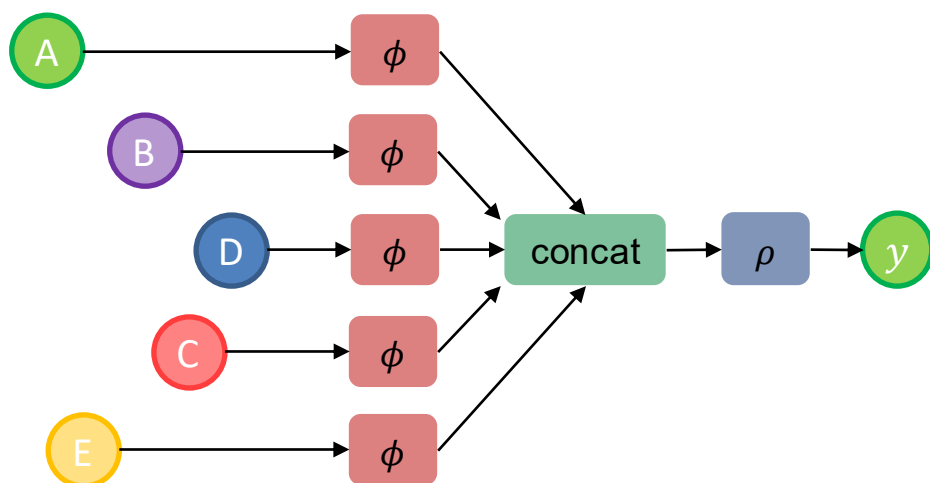


Need to give ABDCE, then ECBDA, and more.... 5! more samples needed

A Simple Classification Example

Models for set-based data must be invariant to element order.

1. Parameter sharing for each set element
2. Not permutation invariant aggregation function

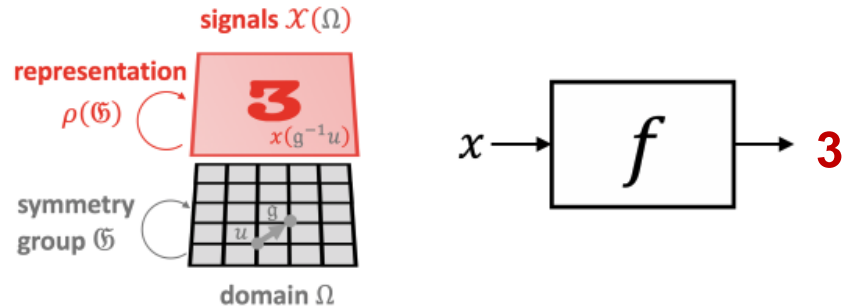


Need to give ABDCE, then ECBDA, and more.... 5! more samples needed

Structure

Data invariances – example of image classification

A function $f : \mathcal{X}(\Omega) \rightarrow \mathcal{Y}$ is \mathfrak{G} -invariant if $f(\rho(\mathfrak{g})x) = f(x)$ for all $\mathfrak{g} \in \mathfrak{G}$ and $x \in \mathcal{X}(\Omega)$, i.e., its output is unaffected by the group action on the input.

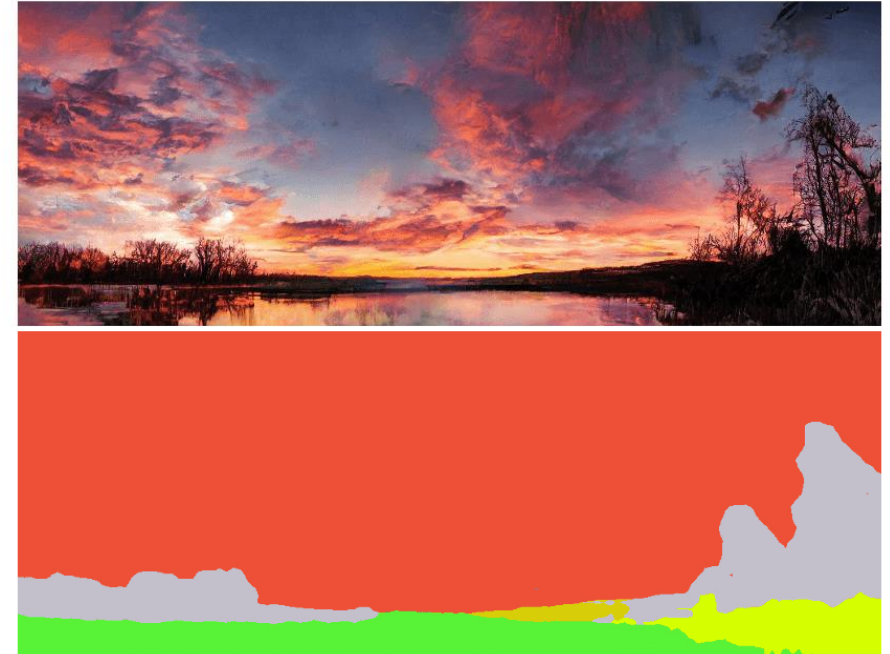
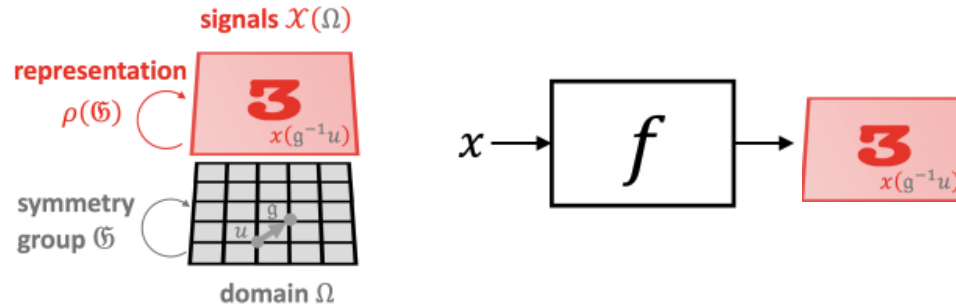


sunset

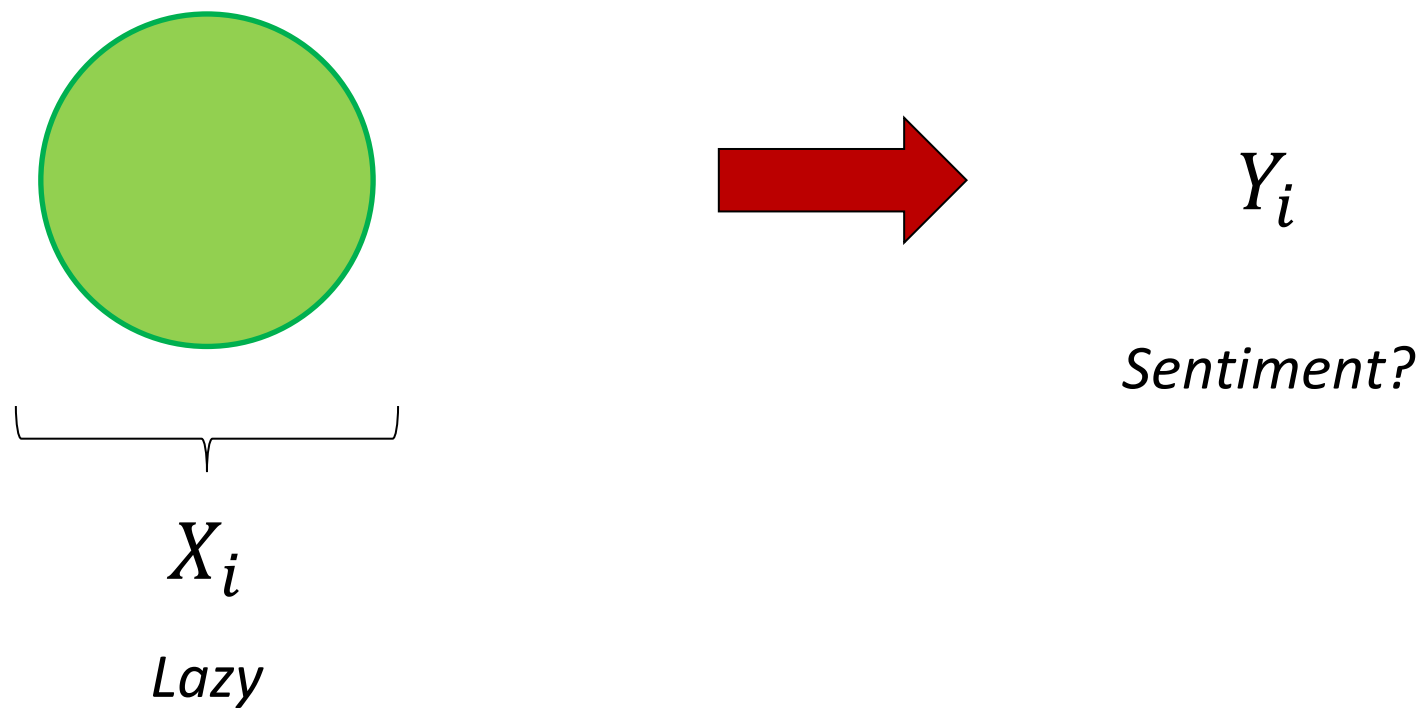
Structure

Data equivariances – example of image segmentation

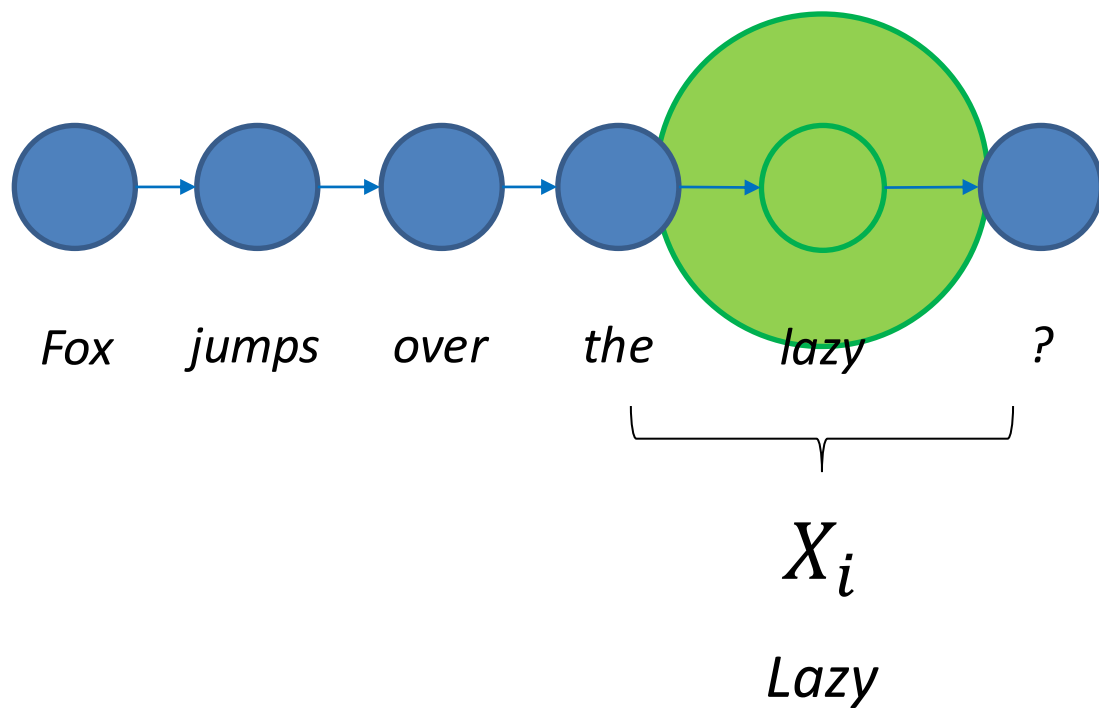
A function $f : \mathcal{X}(\Omega) \rightarrow \mathcal{X}(\Omega)$ is \mathfrak{G} -equivariant if $f(\rho(\mathfrak{g})x) = \rho(\mathfrak{g})f(x)$ for all $\mathfrak{g} \in \mathfrak{G}$, i.e., group action on the input affects the output in the same way.



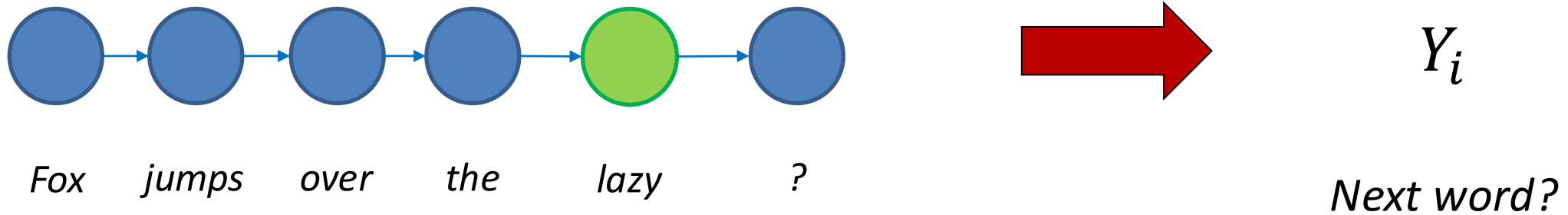
Elements



From Token to Sequences

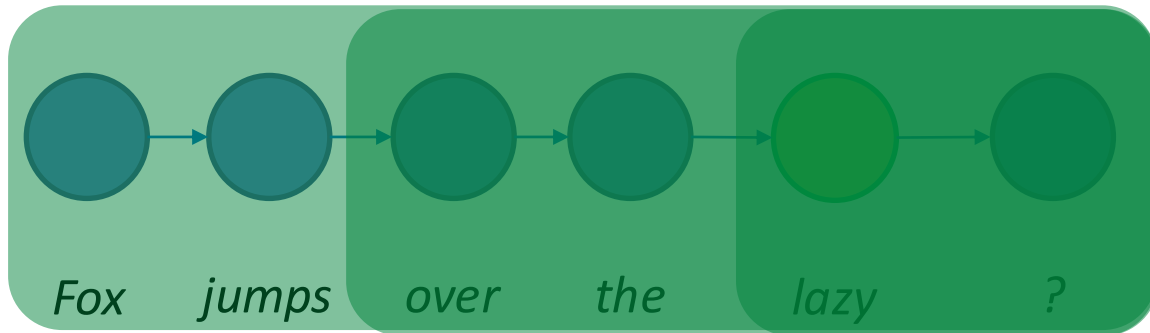


From Token to Sequences



How do we aggregate information?

From Token to Sequences

 Y_i

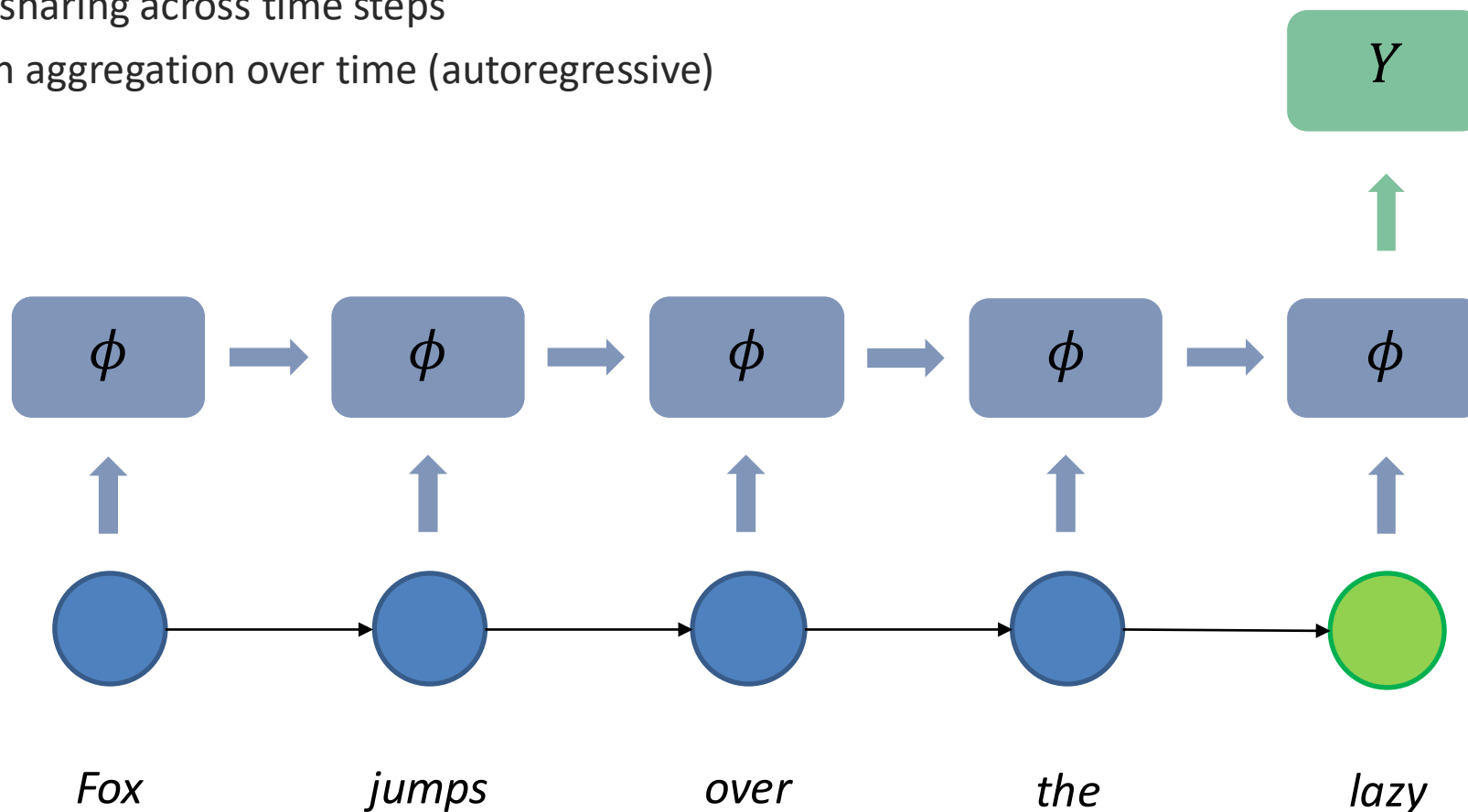
Next word?

How do we aggregate information?

Sequence Classification

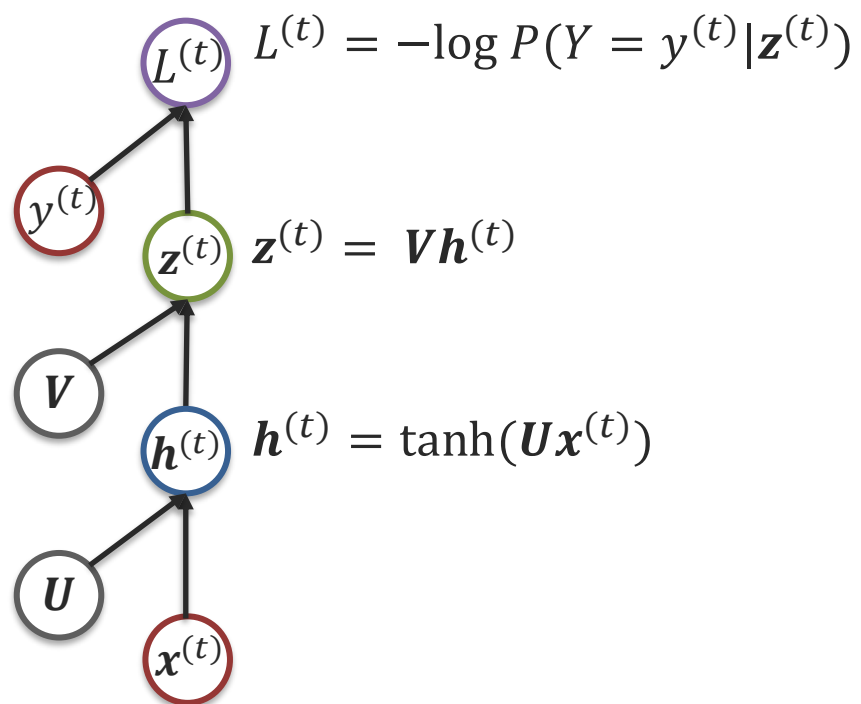
Models for sequential data must be invariant to time, but equivariant to word order.

1. Parameter sharing across time steps
2. Information aggregation over time (autoregressive)

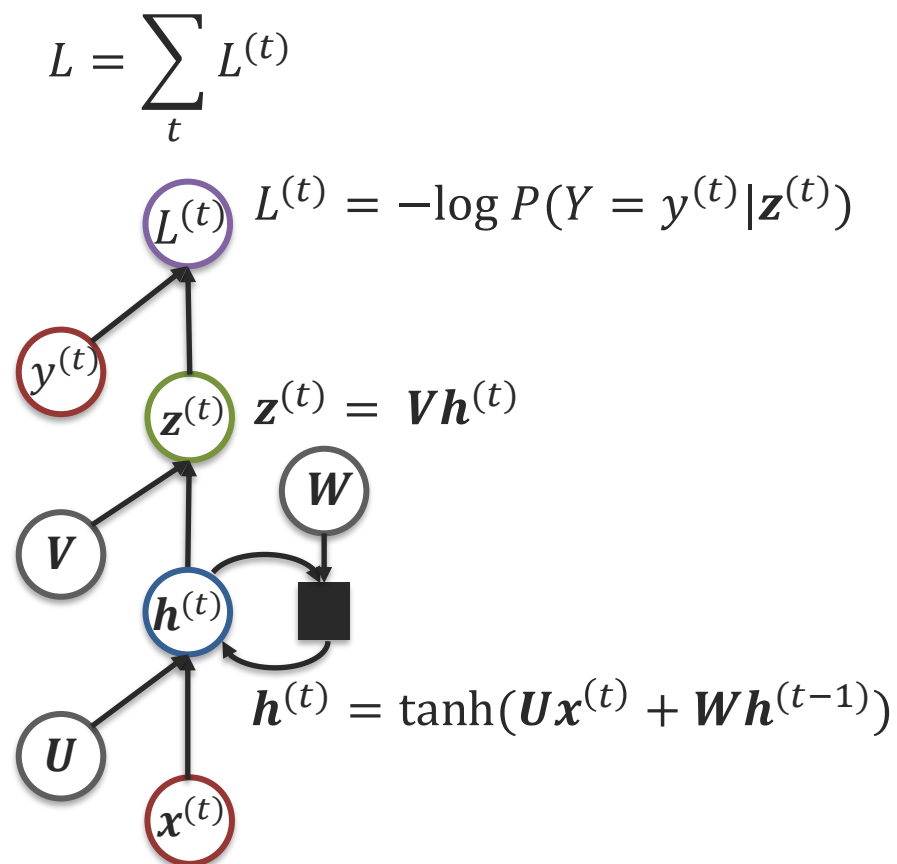


Sequence Classification

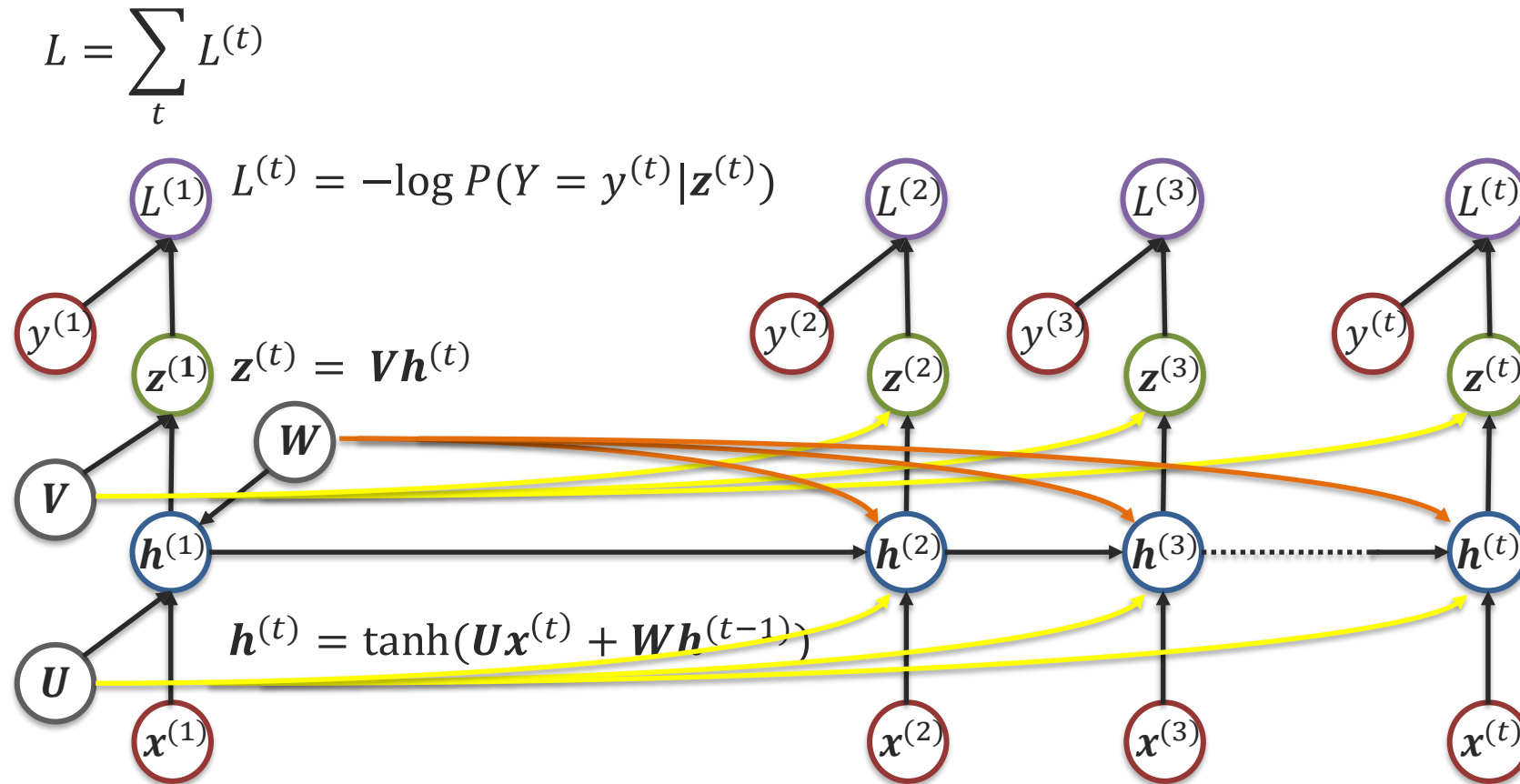
Feedforward Neural Network



Sequence Classification

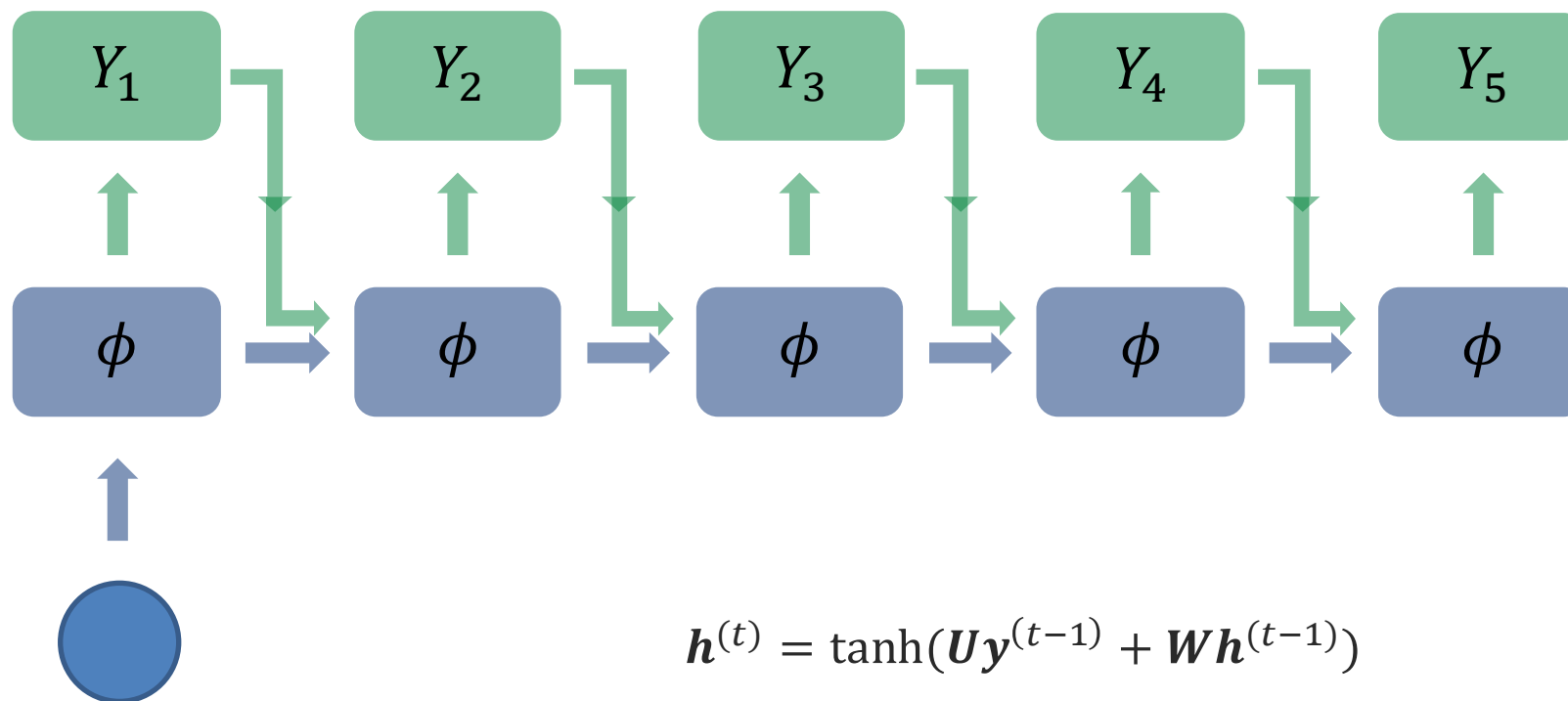


Sequence Classification



Same model parameters are used for all time steps.

Sequence Generation



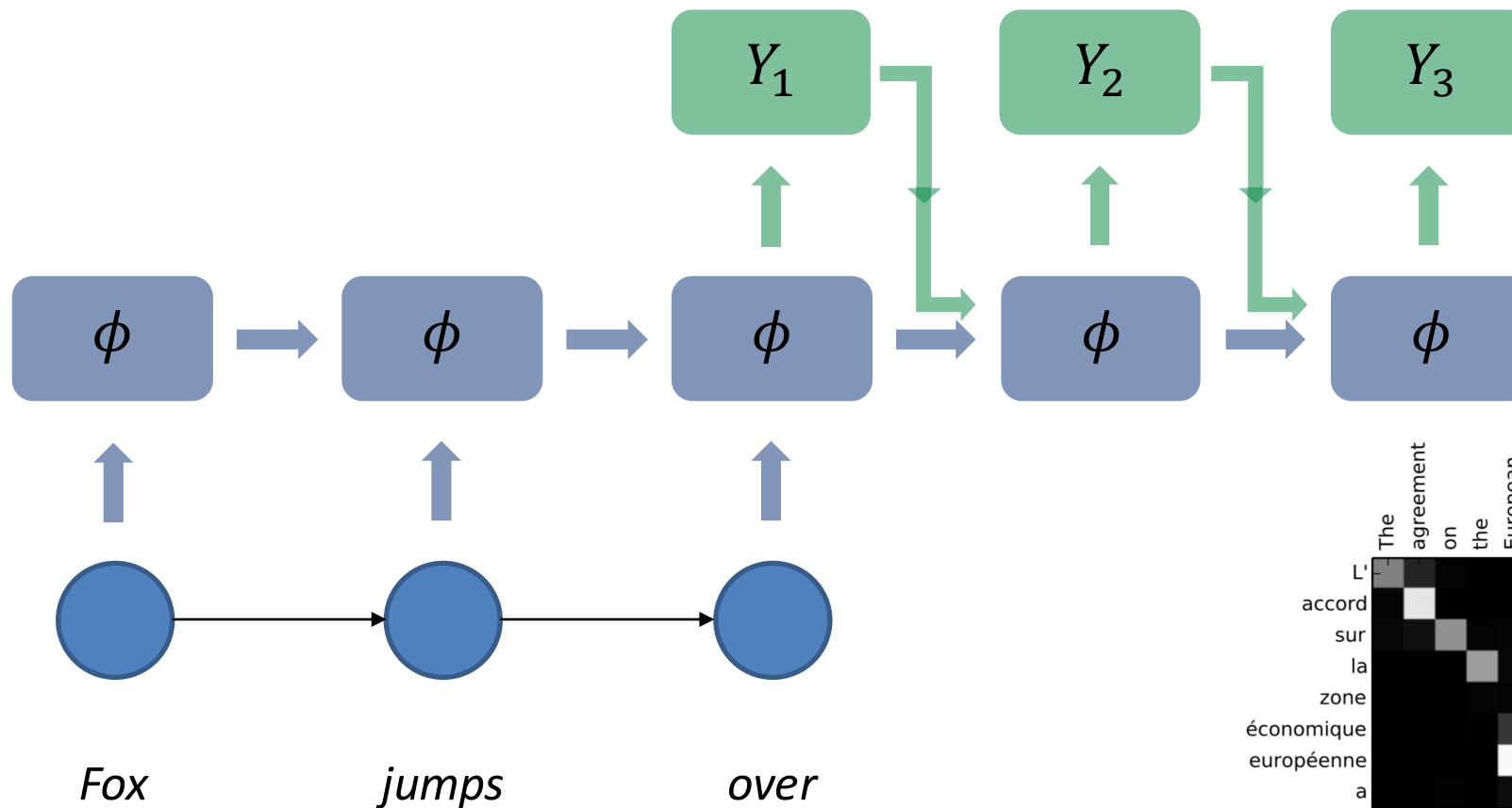
$$\mathbf{h}^{(t)} = \tanh(\mathbf{U}\mathbf{y}^{(t-1)} + \mathbf{W}\mathbf{h}^{(t-1)})$$

Fox

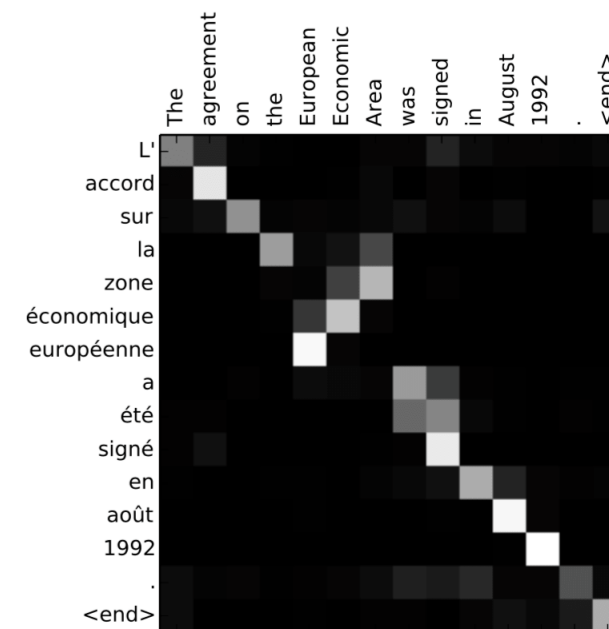
e.g, text or music generation

Modern versions: RNN -> LSTM -> TCN -> State space models

Sequence-to-Sequence Models



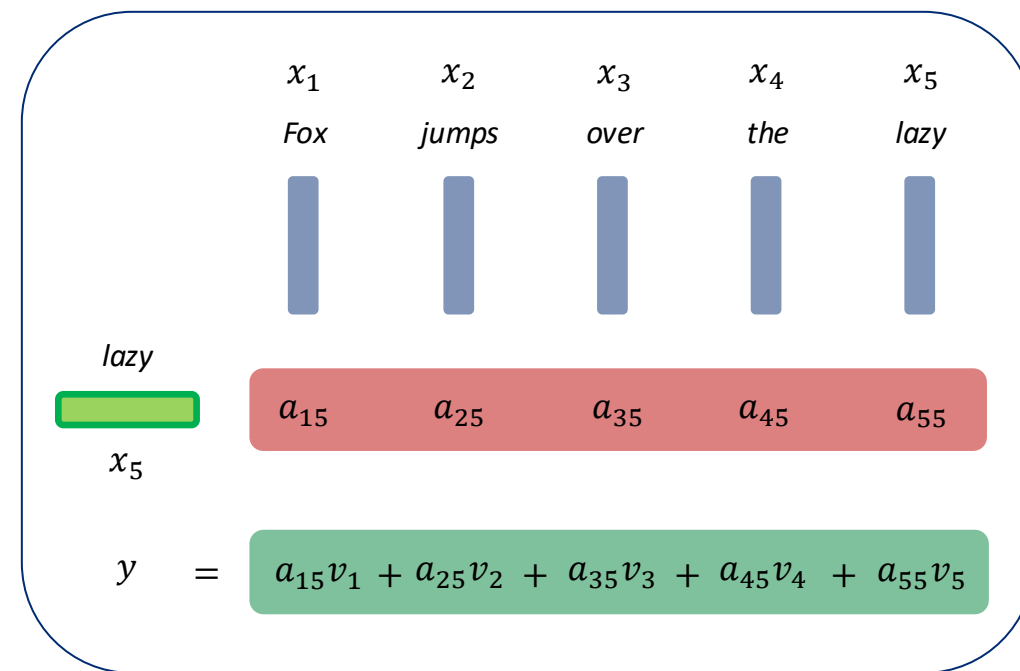
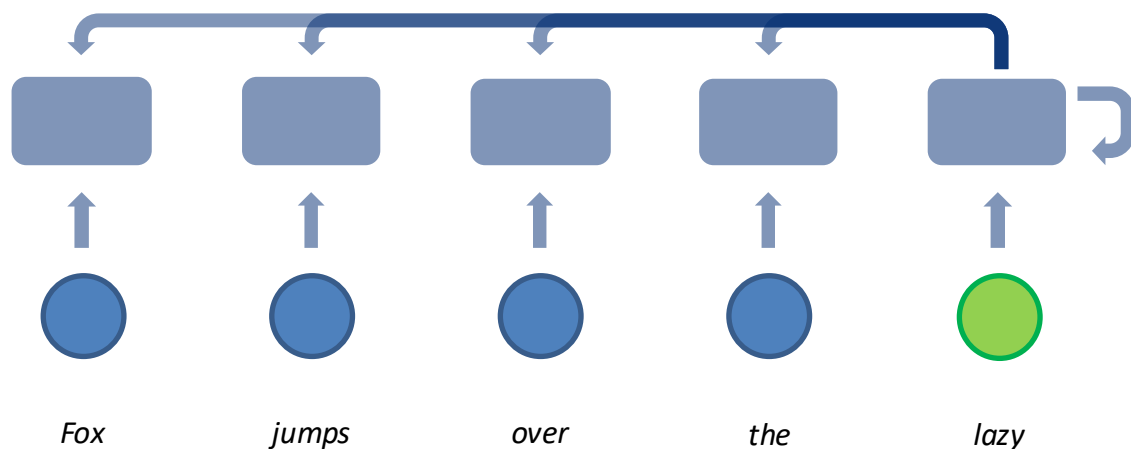
e.g, machine translation -> birth of attention-based models



Modern Sequence Models

Birth of attention-based models

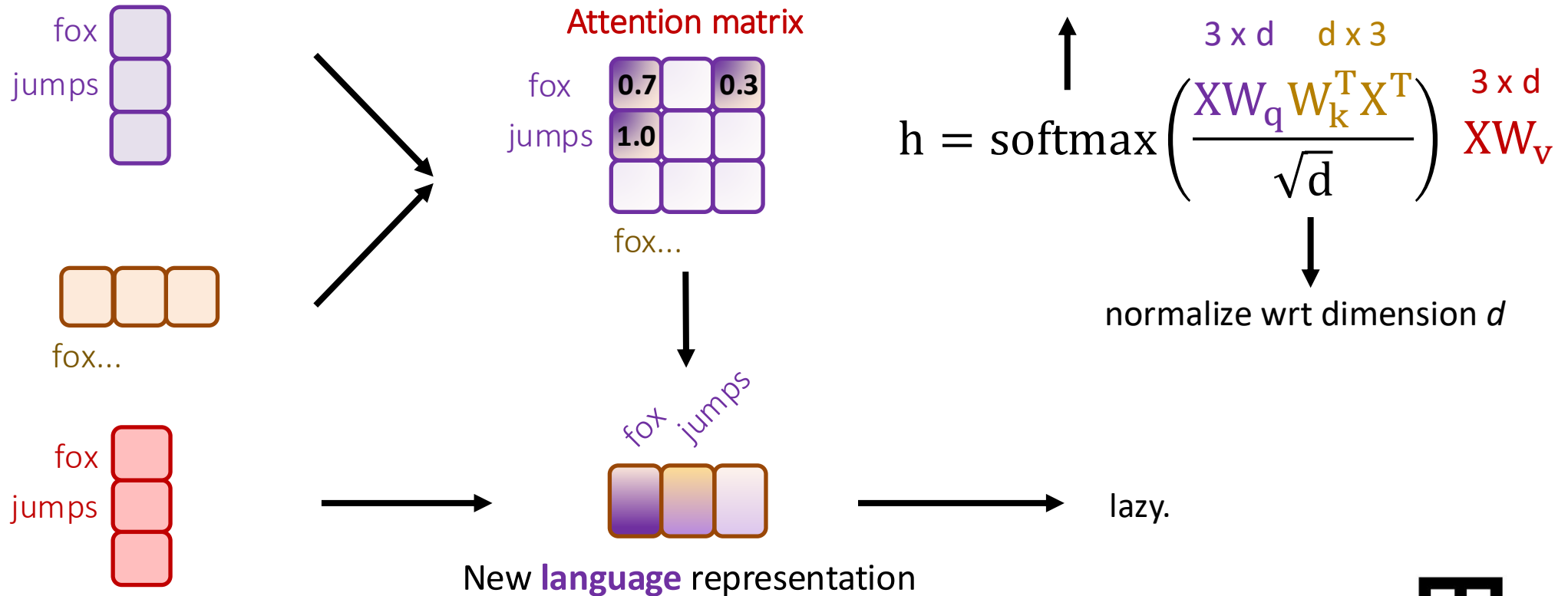
– Dynamic weights for different elements



Modern Sequence Models

Birth of attention-based models

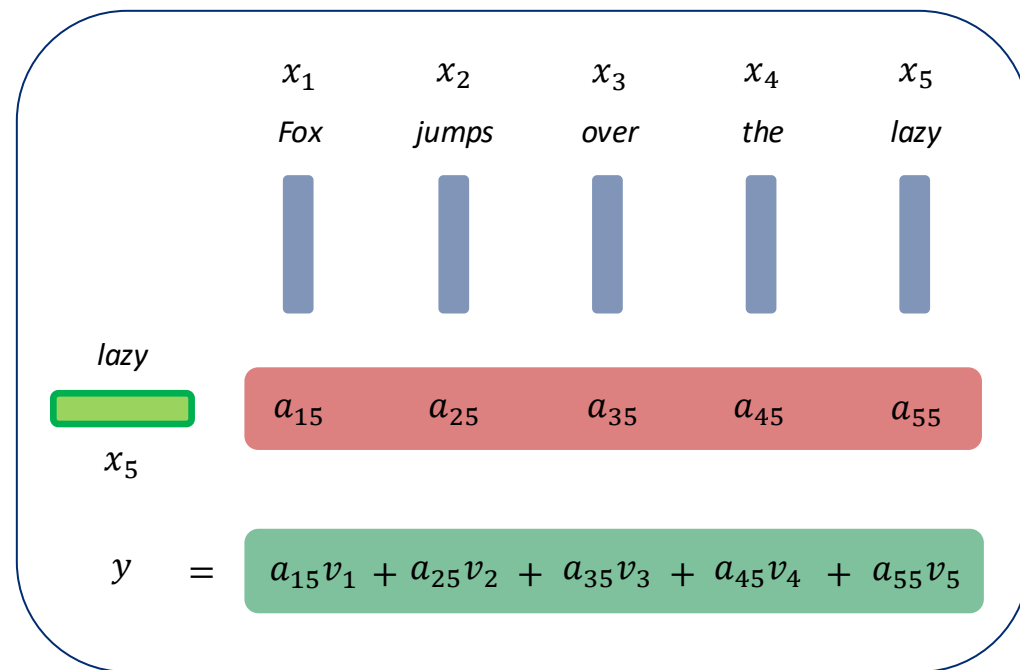
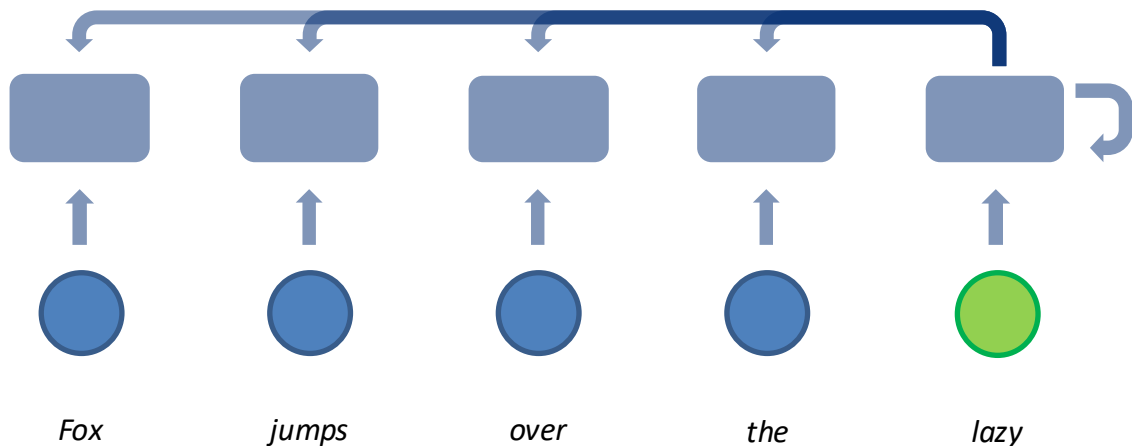
– Dynamic weights for different elements



Modern Sequence Models

Models for sequential data must be invariant to time, but equivariant to word order.

1. Parameter sharing across time steps
2. Information aggregation over time (in parallel)



$$h = \text{softmax} \left(\frac{XW_q W_k^T X^T}{\sqrt{d}} \right) XW_v$$

$T \times d$ $d \times T$ $T \times d$
 X W_q W_k^T X^T X W_v

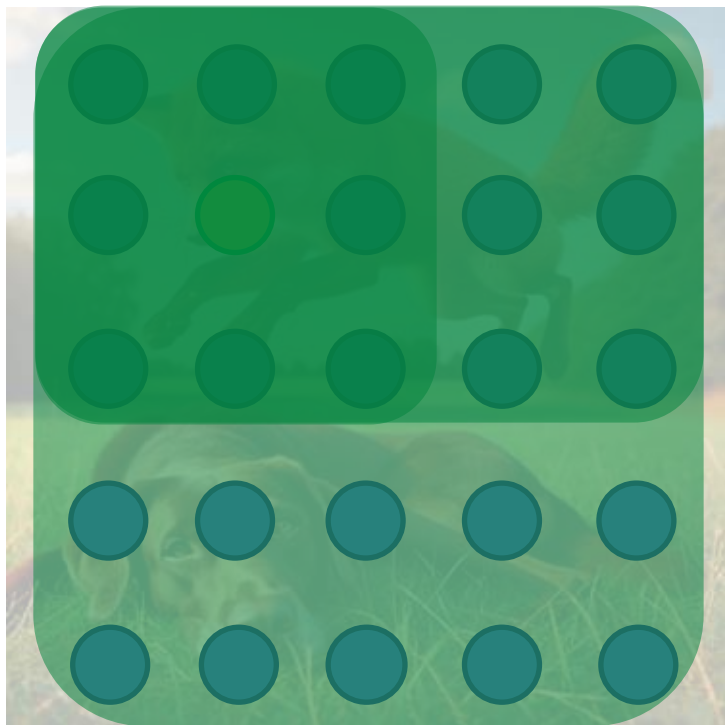
Spatial Data

 Y_i

Is there a fox?

How do we aggregate information?

Spatial Data

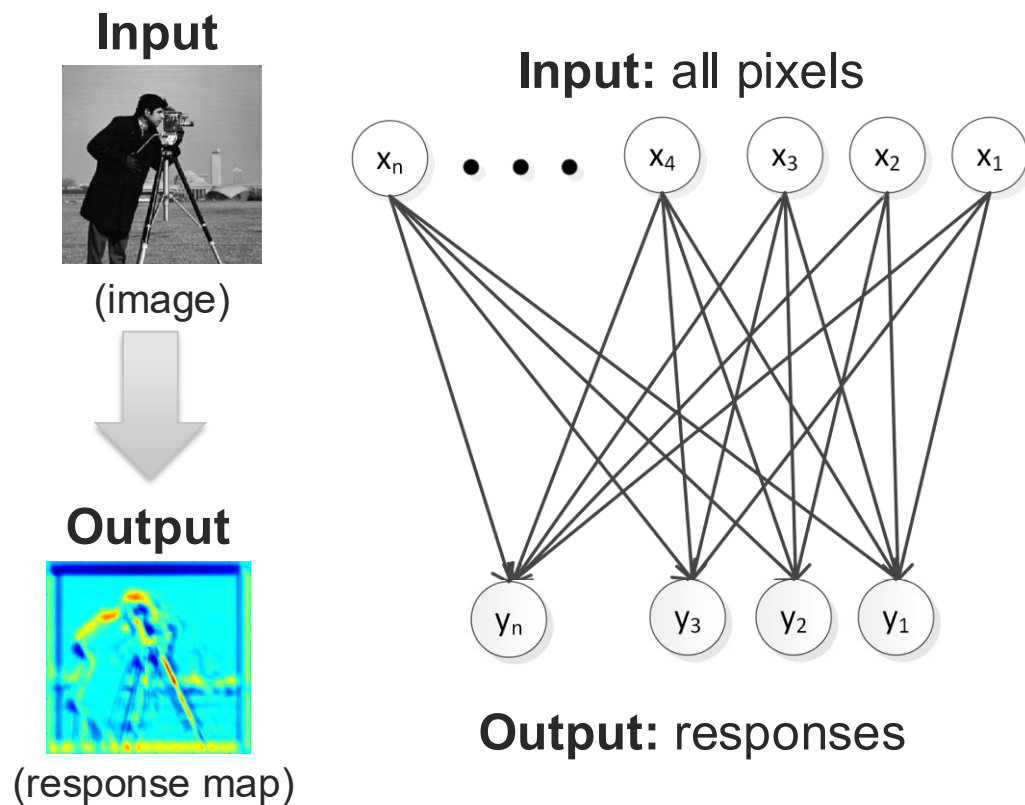
 Y_i

Is there a fox?

How do we aggregate information?

Convolutional Neural Networks

Models for spatial data need to be invariant to spatial translations.



Not efficient!

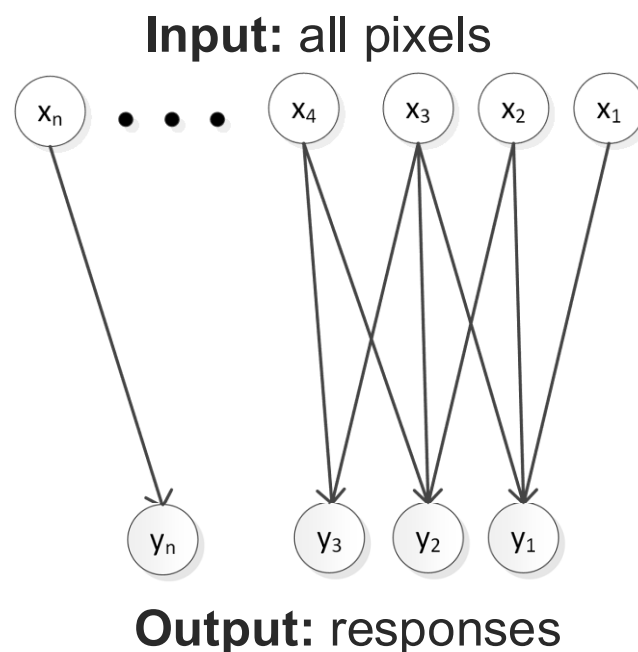
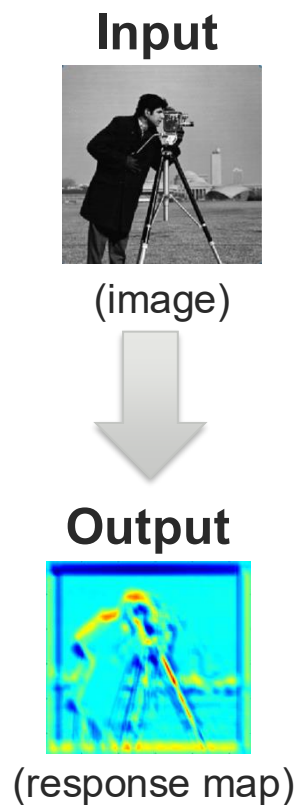
200 × 200 image
requires
40,000 × n
parameters
(where n is size of output)

**And it may learn different outputs
for different pixel positions**

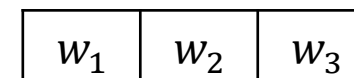
➔ Not spatial invariant

Convolutional Neural Networks

Modification 1: Only apply the filter to a small sliding window
 -> for efficiency and locality

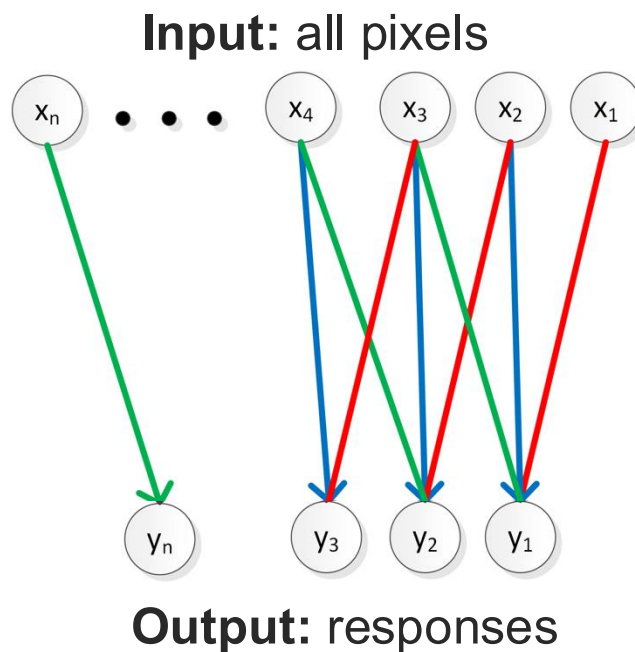
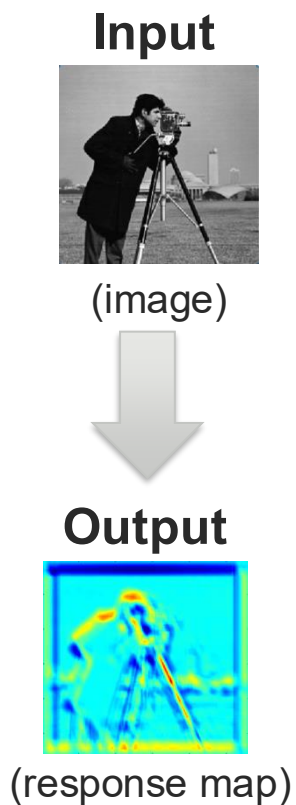


Example with
1D filter:



Convolutional Neural Networks

Modification 2: Same filter applied to all sliding windows
-> for spatial invariance



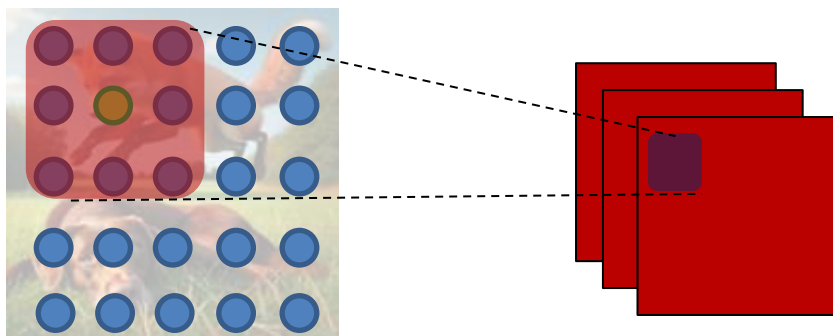
Example with
1D filter:



Convolutional Neural Networks

Models for spatial data need to be invariant to spatial translation

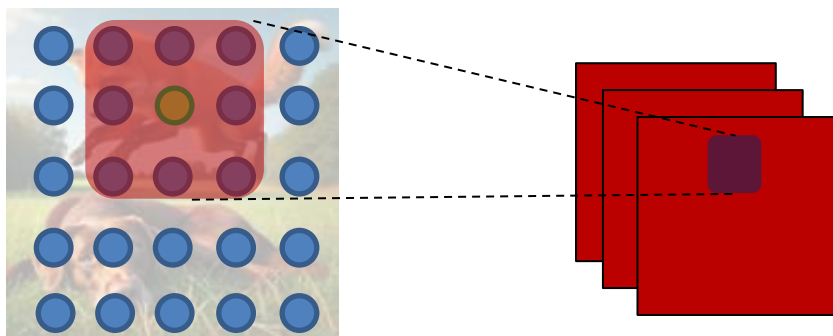
1. Parameter sharing across $k \times k$ convolutional filter



Convolutional Neural Networks

Models for spatial data need to be invariant to spatial translation

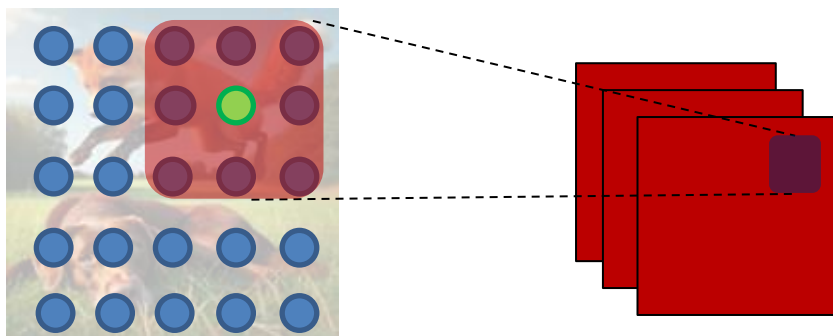
1. Parameter sharing across $k \times k$ convolutional filter



Convolutional Neural Networks

Models for spatial data need to be invariant to spatial translation

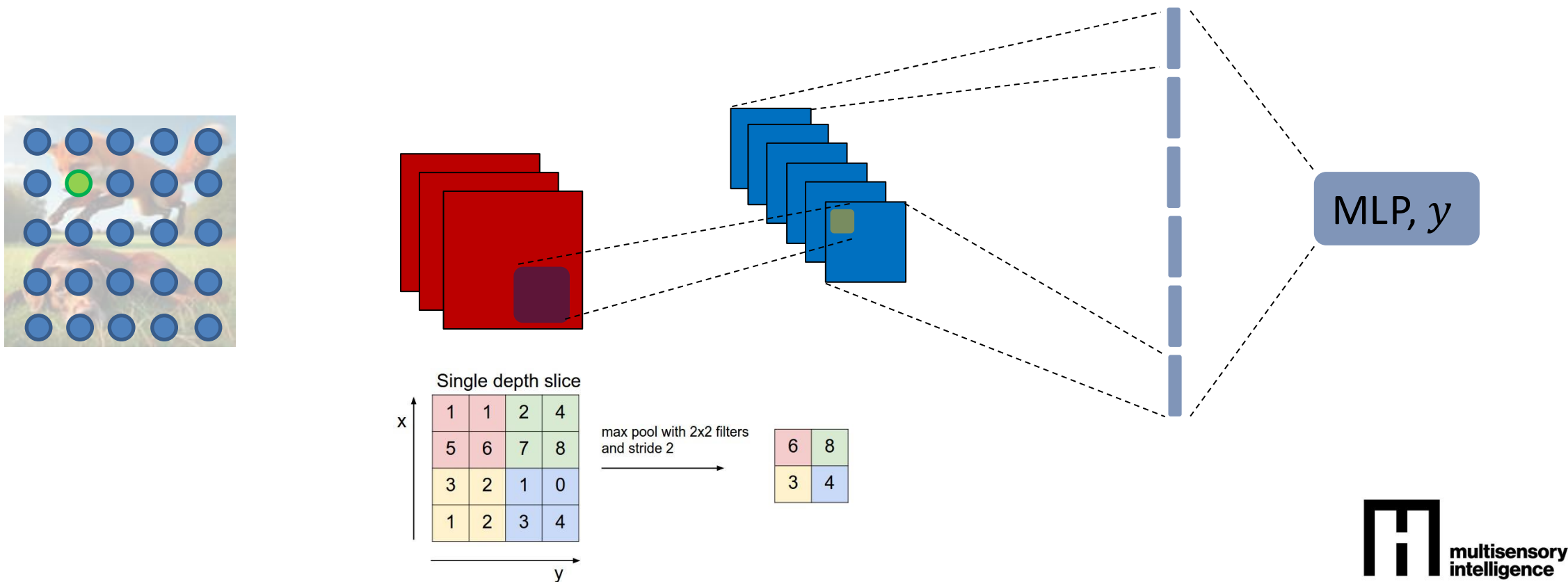
1. Parameter sharing across $k \times k$ convolutional filter



Convolutional Neural Networks

Models for spatial data need to be invariant to spatial translation

1. Parameter sharing across $k \times k$ convolutional filter
2. Information aggregation over $k \times k$ pooling region



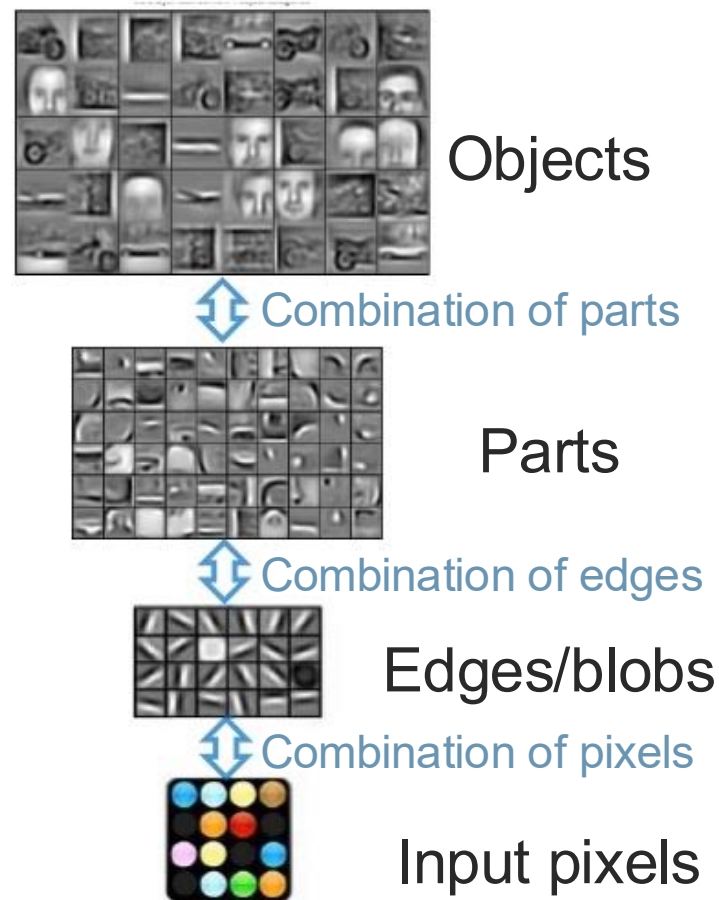
Convolutional Neural Networks

Multiple convolutional layers

→ Allows the network to learn combinations of sub-parts, to increase complexity

Multiple pooling layers

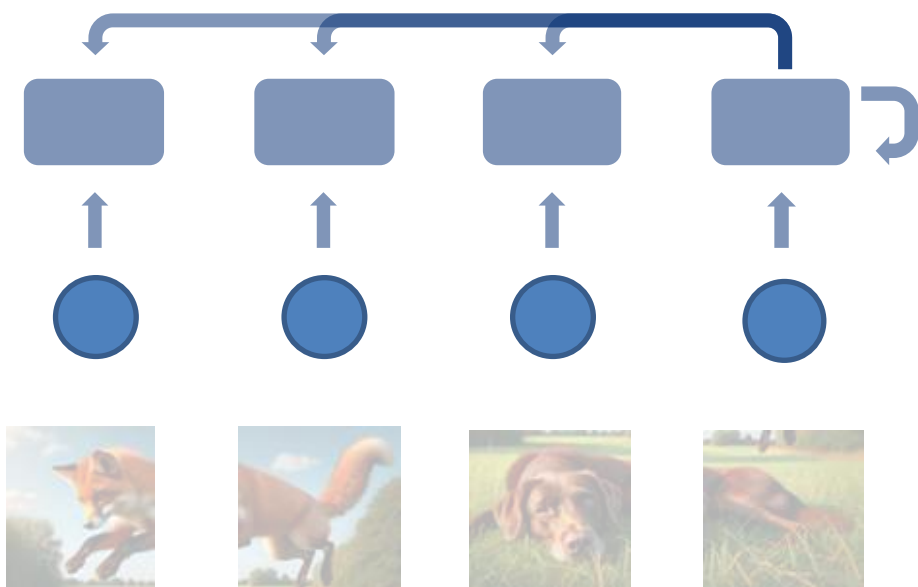
→ Allows the network to learn increasingly abstract & summarized information



Vision Transformer

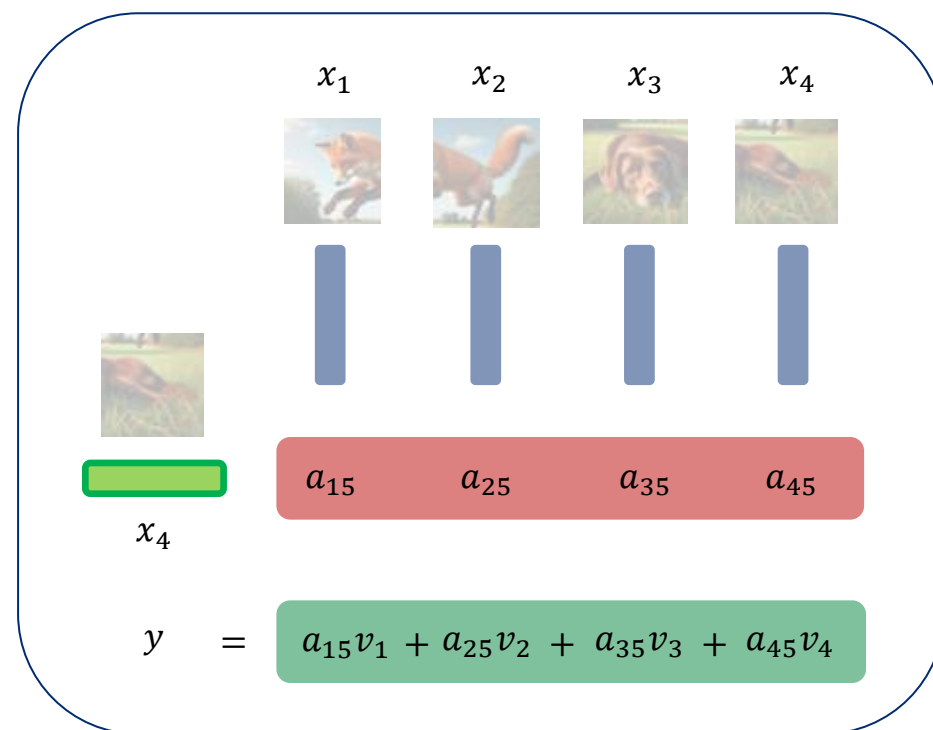
Models for spatial data need to be invariant to spatial translation

1. Parameter sharing across $k \times k$ self-attention region
2. Information aggregation over $k \times k$ patch region



$$h = \text{softmax} \left(\frac{XW_q W_k^T X^T}{\sqrt{d}} \right) XW_v$$

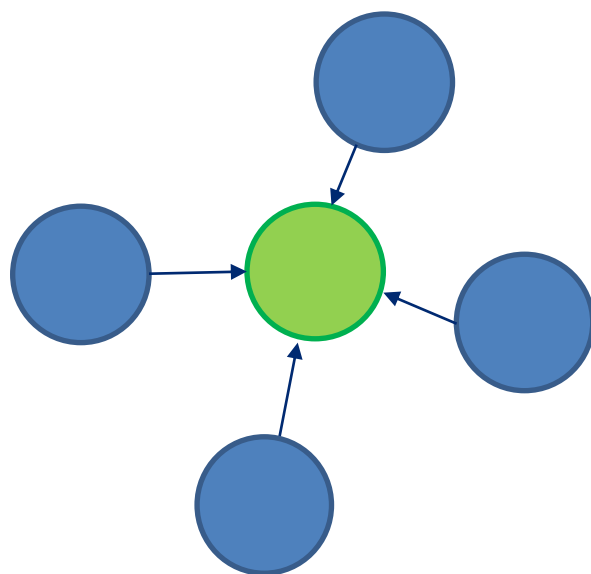
$T \times d$ $d \times T$
 $T \times d$
 XW_v



Vision Transformer



Graphs

 Y_i

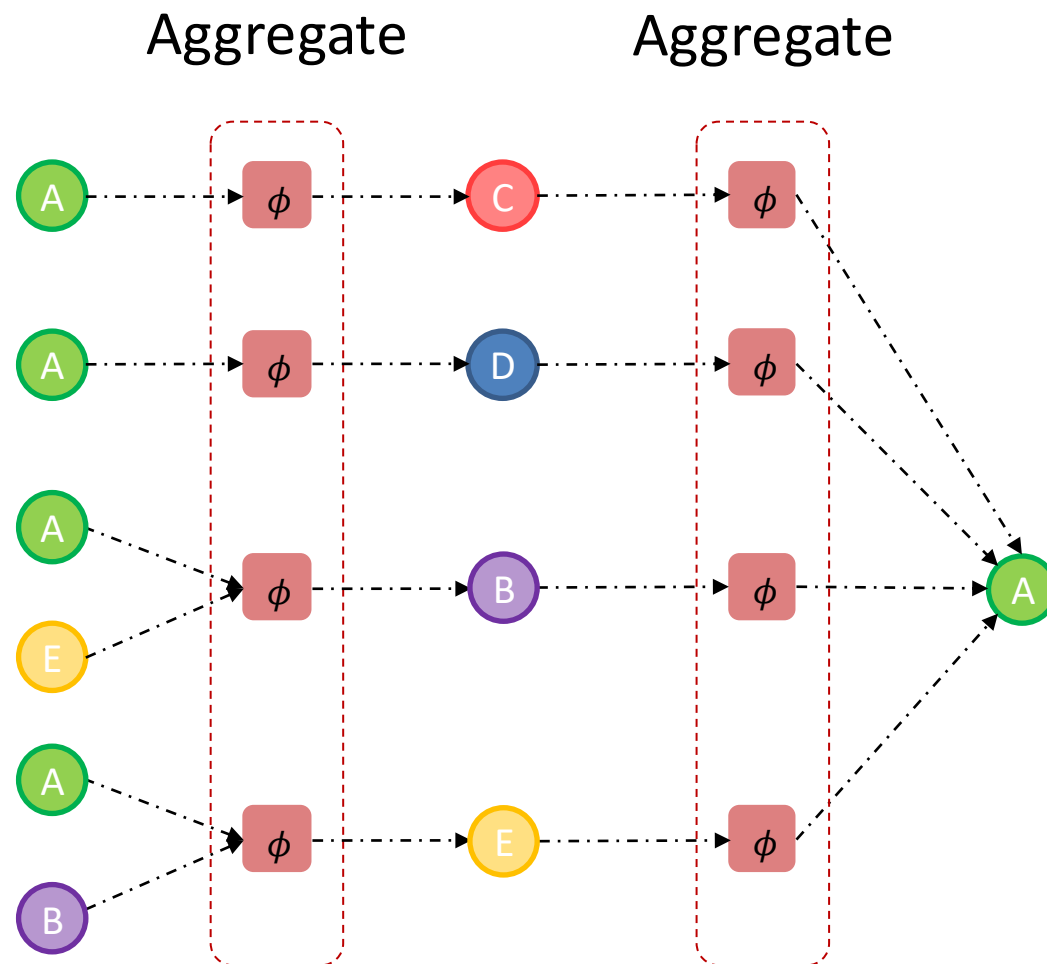
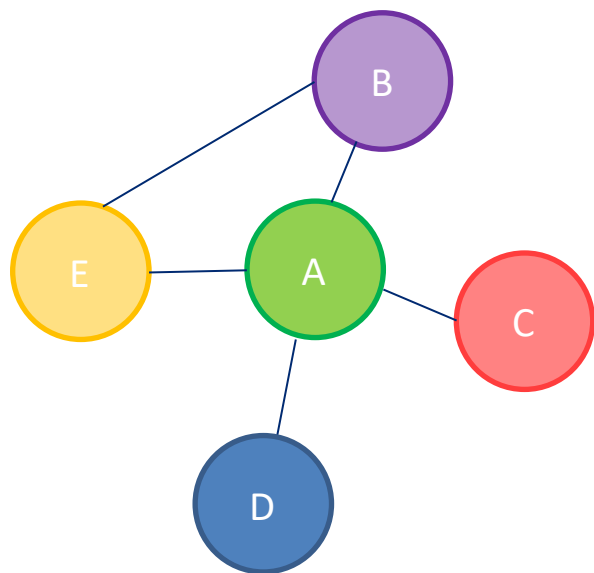
What molecule is this?

How do we aggregate information?

Graph Neural Networks

Models for graph data:

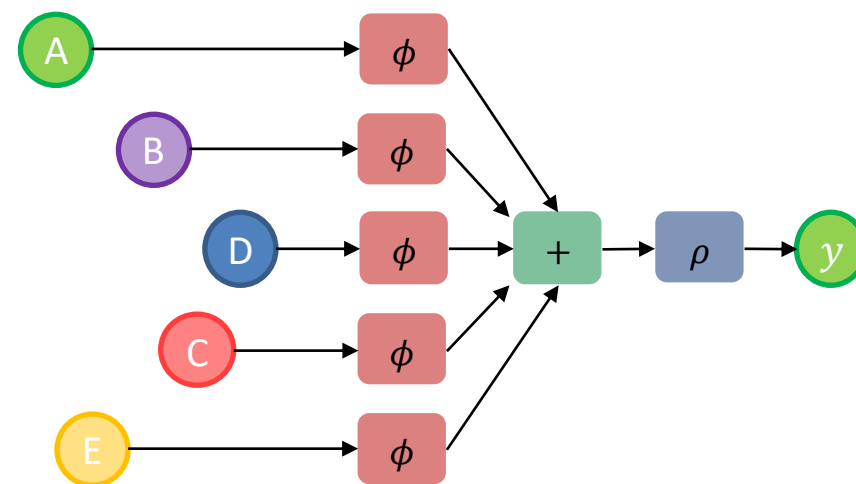
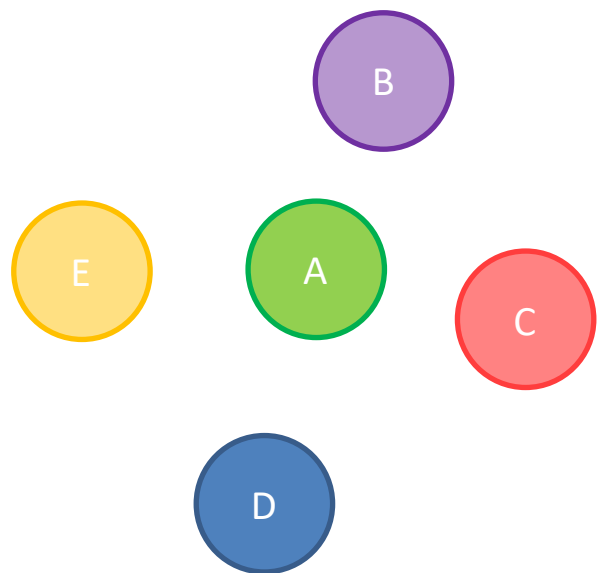
1. Parameter sharing across nodes
2. Information aggregation over neighbors (edges)



Graph Recover Sets

Sets are graphs with only nodes, no edges

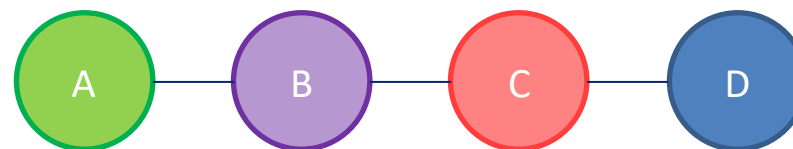
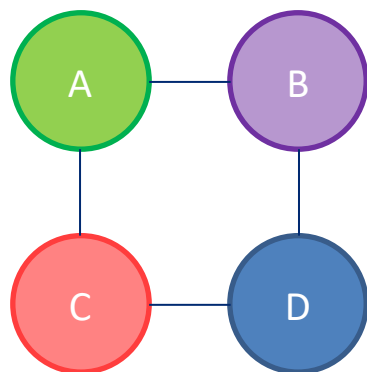
1. Parameter sharing across nodes \rightarrow set elements
2. Information aggregation over neighbors \rightarrow no neighbors



Graph Recover Spatial and Temporal Data

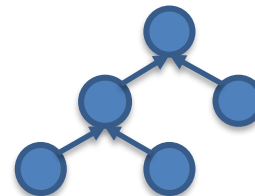
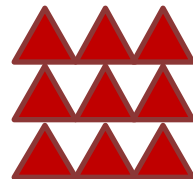
Spatial data and sequential data

1. Parameter sharing across nodes
2. Information aggregation over neighbors



Summary: How To Model

1. Decide how much data to collect, and how much to label (costs and time)
2. Clean data: normalize/standardize, find noisy data, anomaly/outlier detection
3. Visualize data: plot, dimensionality reduction (PCA, t-sne), cluster analysis
4. Decide on evaluation metric (proxy + real, quantitative and qualitative)
5. Choose modeling paradigm - domain-specific vs general-purpose
6. Figure out base elements and their representation
7. Figure out data invariances & equivariances (+other parts of modality profile)
8. Iterate between data collection, model design, model training, hyperparameter tuning etc. until satisfied.



Lecture Summary

- 1 A unifying paradigm of model architectures
- 2 Temporal sequence models
- 3 Spatial convolution models
- 4 Models for sets and graphs

Assignments for This Coming Week

HW1: reading assignment + homework due next Tuesday 2/17

For project:

- Project proposal instructions released, due Tuesday (2/24). Submit on canvas
- Meet with me 4-5pm if need feedback about proposal ideas.