

How to AI (Almost) Anything

Lecture 3 – Common model architectures

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

For project:

- 4-page project proposal due tonight (2/25), email to me.
- Meet with me after class 2-3pm if need feedback about proposal ideas.

Reading assignment due tomorrow Wednesday (2/26).

This Thursday (2/27): first reading discussion on **data and learning**.

Bitter lesson

Grokking/double descent

Logistics – Reading Assignments

Roles and Grading

[Role assignments for every reading are linked here](#)

Scientific Peer Reviewer

Task: Complete a full review of the paper, recommending acceptance or rejection, and address all prompts in the review form (e.g., technical soundness, clarity, originality, significance). See [an example of review instructions here](#) (navigate down to 'review form').

Grading

(Unacceptable)

- Review is incomplete or missing major components.
- Minimal effort with vague or unsubstantiated comments.
- No clear recommendation or justification provided.

(Needs Improvement)

- Review touches on a few relevant points but lacks depth and critical evaluation.
- Feedback is superficial and misses key aspects of the paper.

(Adequate)

- Provides a basic review covering most areas.
- Offers some useful feedback but misses deeper analysis or constructive suggestions.
- The recommendation is justified but could be stronger.

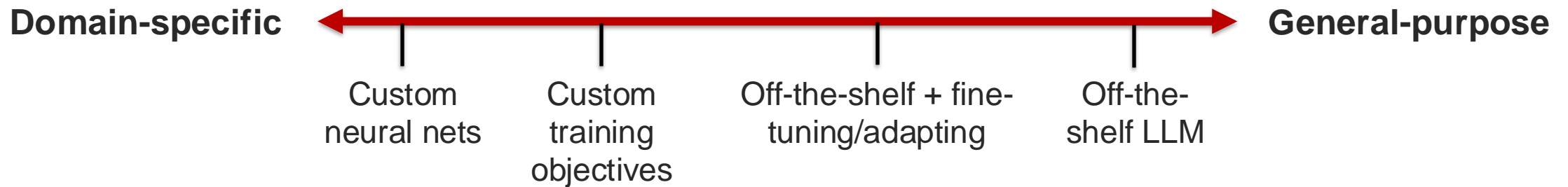
(Good)

- Thoughtful, thorough review with balanced criticism and praise.
- Clearly explains the paper's strengths and weaknesses.
- Provides helpful suggestions for improvement.

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

Lecture Topics *(subject to change, based on student interests and course discussions)*

Domain-specific/custom models

Week 4 (2/25): Common model architectures

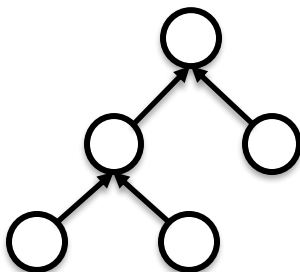
Week 5 (3/4): Multimodal connections and alignment

Week 6 (3/11): Multimodal interactions and fusion

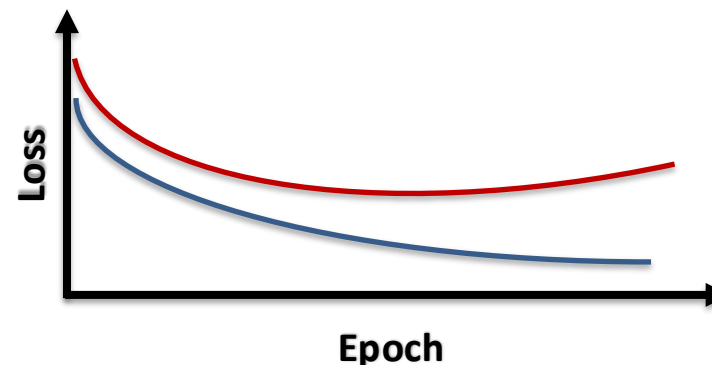
Week 7 (3/18): Cross-modal transfer



Spatial



Hierarchical



Lecture Topics *(subject to change, based on student interests and course discussions)*

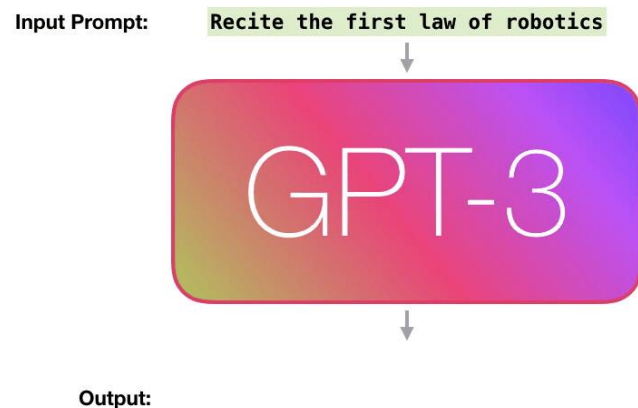
General architectures and adapting pre-trained models

Week 9 (4/1): Pre-training, scaling, fine-tuning LLMs

Week 10 – No class, member's week

Week 11 (4/15): Large multimodal models

Week 12 (4/22): Modern generative AI



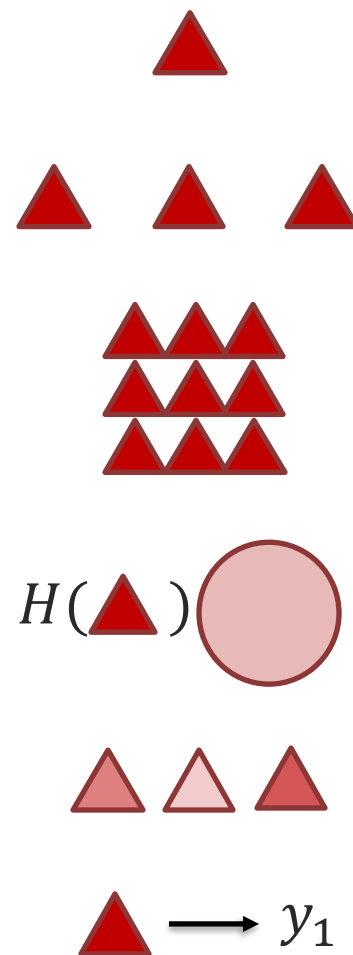
*An armchair in
the shape of an
avocado*



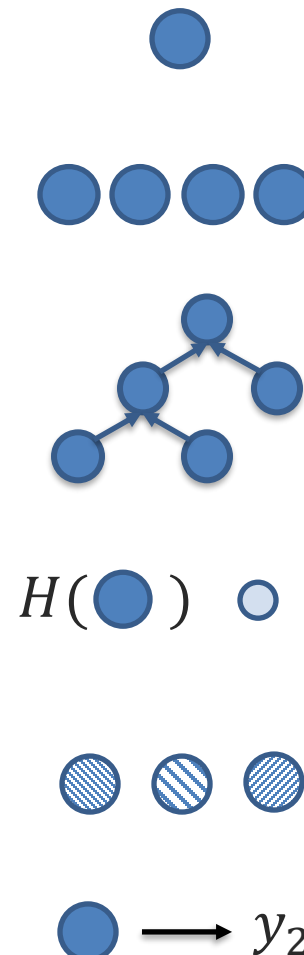
Modality Profile

- 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

Modality A



Modality B



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



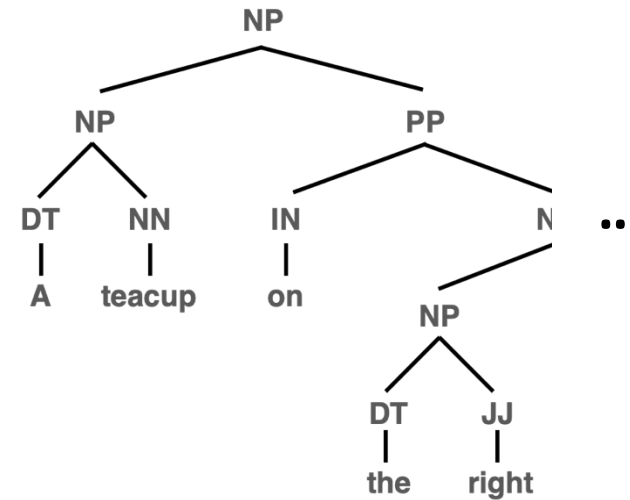
objects per image



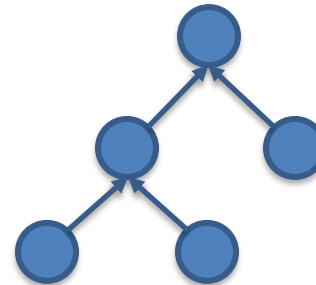
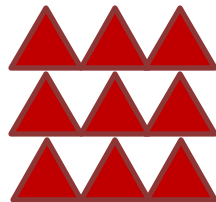
words per minute

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.

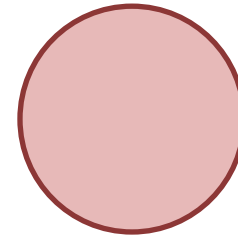


*A teacup on the right of a laptop
in a clean room.*

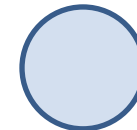
4

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

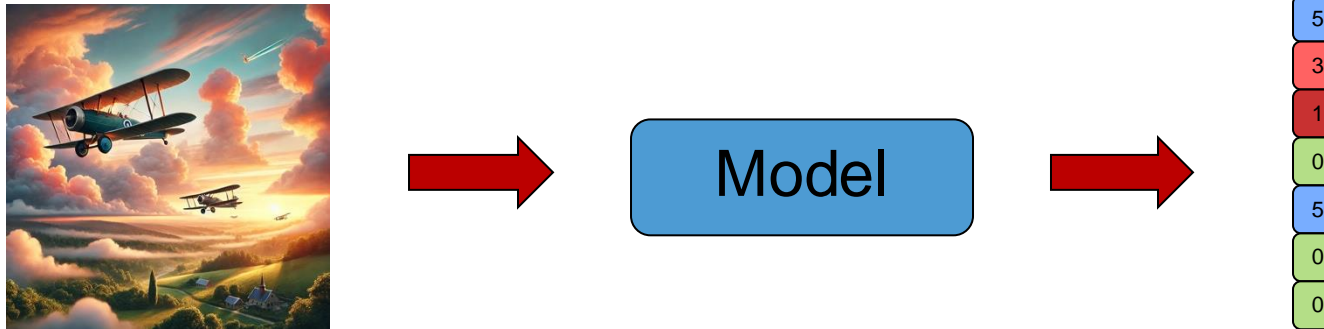


teacup → **teacip**

right → **rihjt**

Unified View of Deep Learning Models

1. Learning representations



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

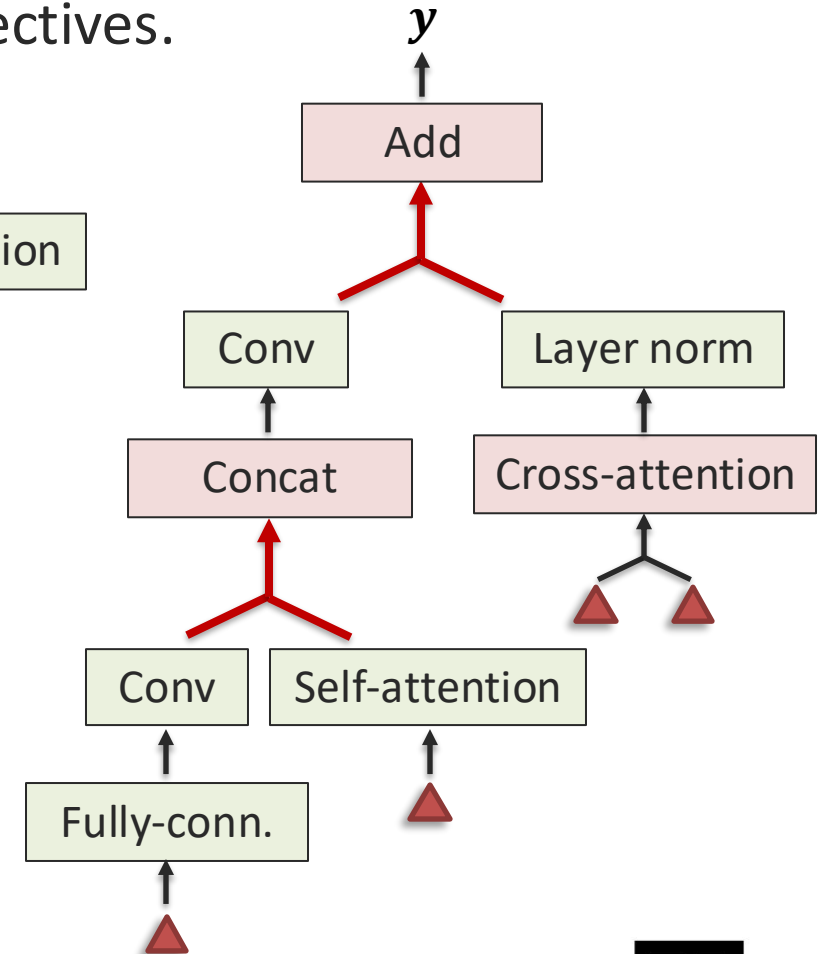


2. Basic information aggregation blocks



3. Compute loss function

4. Take gradients, update with stochastic gradient descent

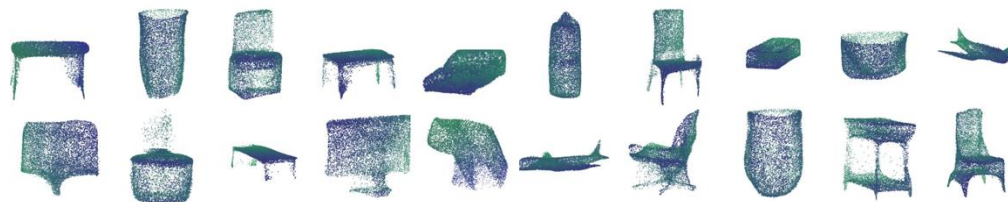


A Simple Classification Example

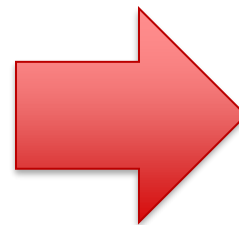
Sets and point clouds



Sets



Point clouds

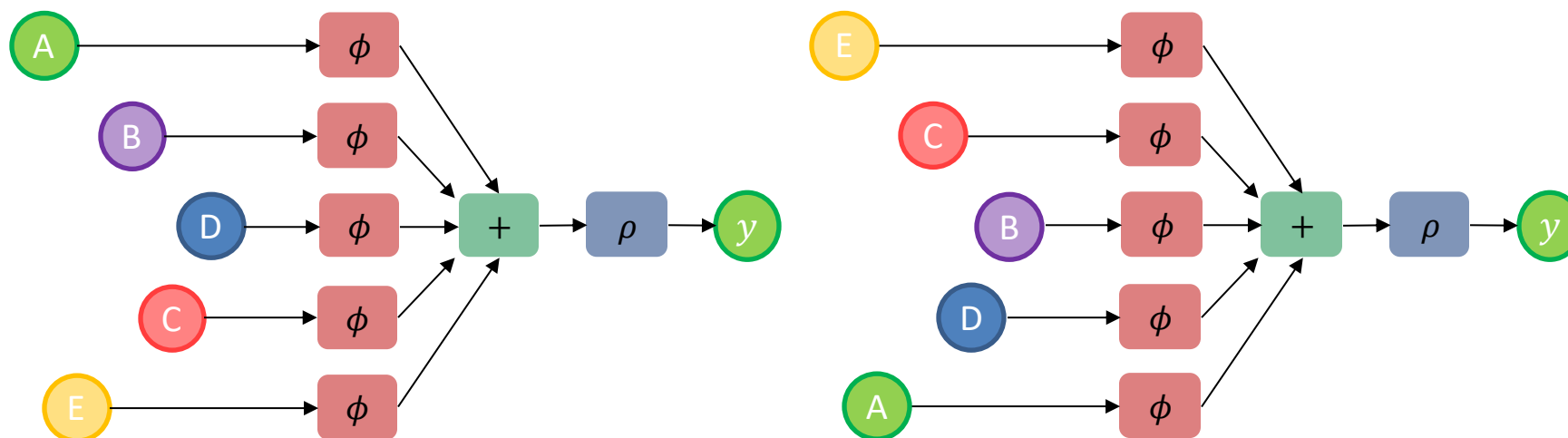


Set anomaly detection
Set expansion
Set completion
Point cloud classification
Point cloud generation

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

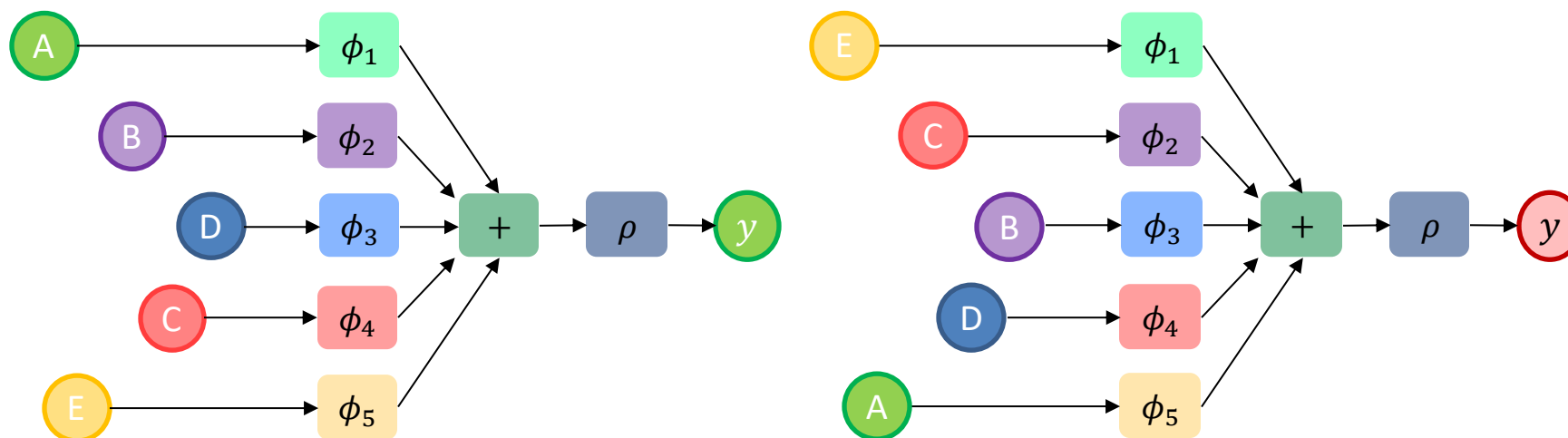


Giving ABDCE also gives ECBDA, BCAED etc...

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

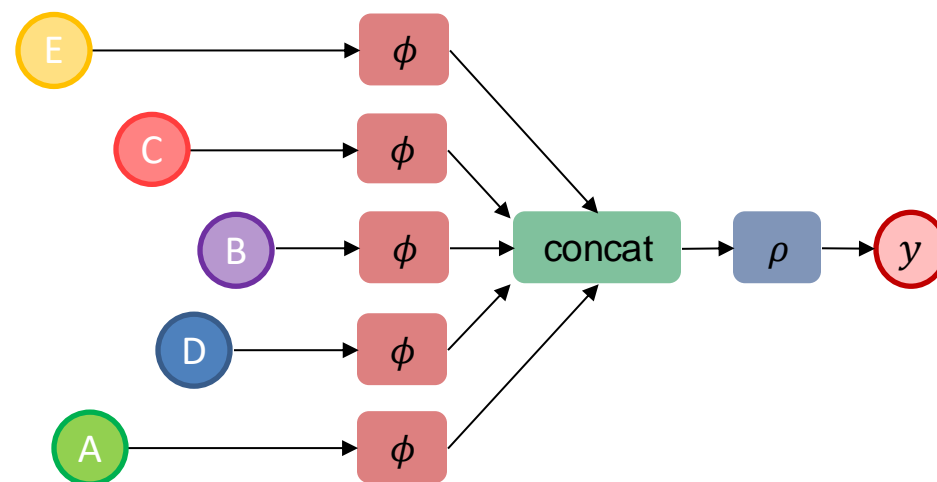
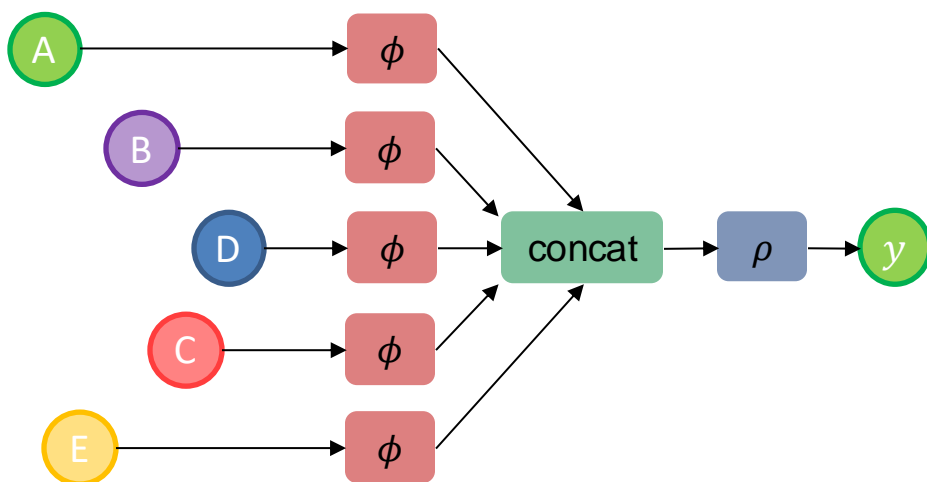


Need to give ABDCE, then ECBDA, and more.... 5! more samples needed

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

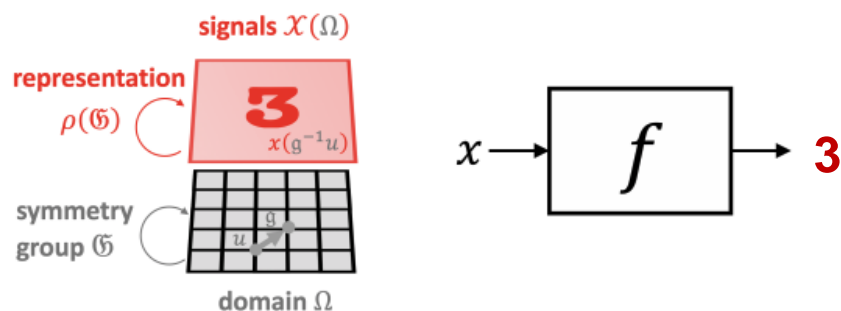


Need to give ABDCE, then ECBDA, and more.... 5! more samples needed

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.

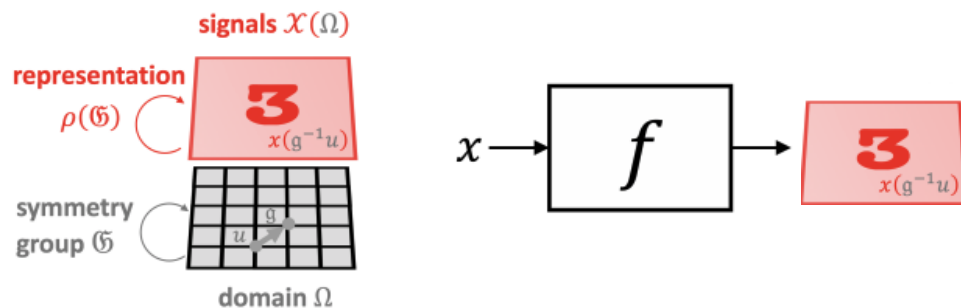


sunset

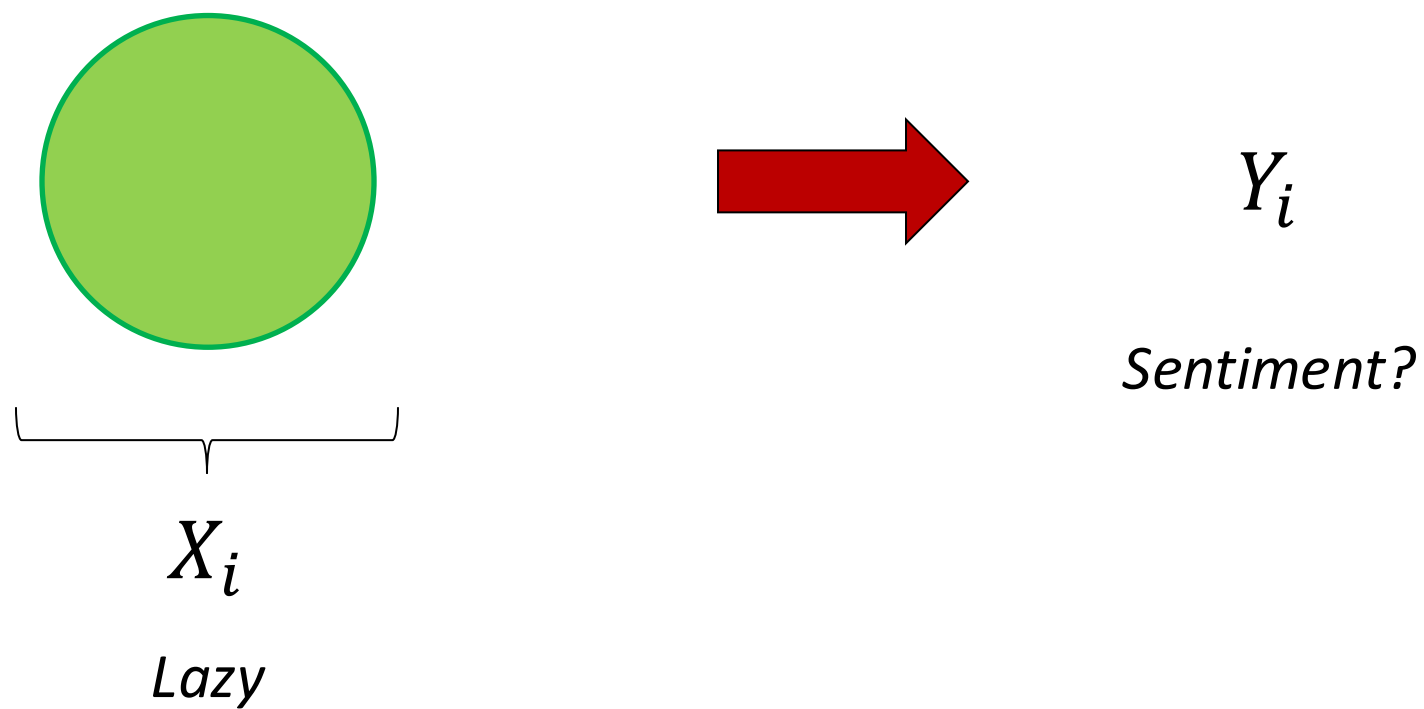
Structure

Data equivariances – example of image segmentation

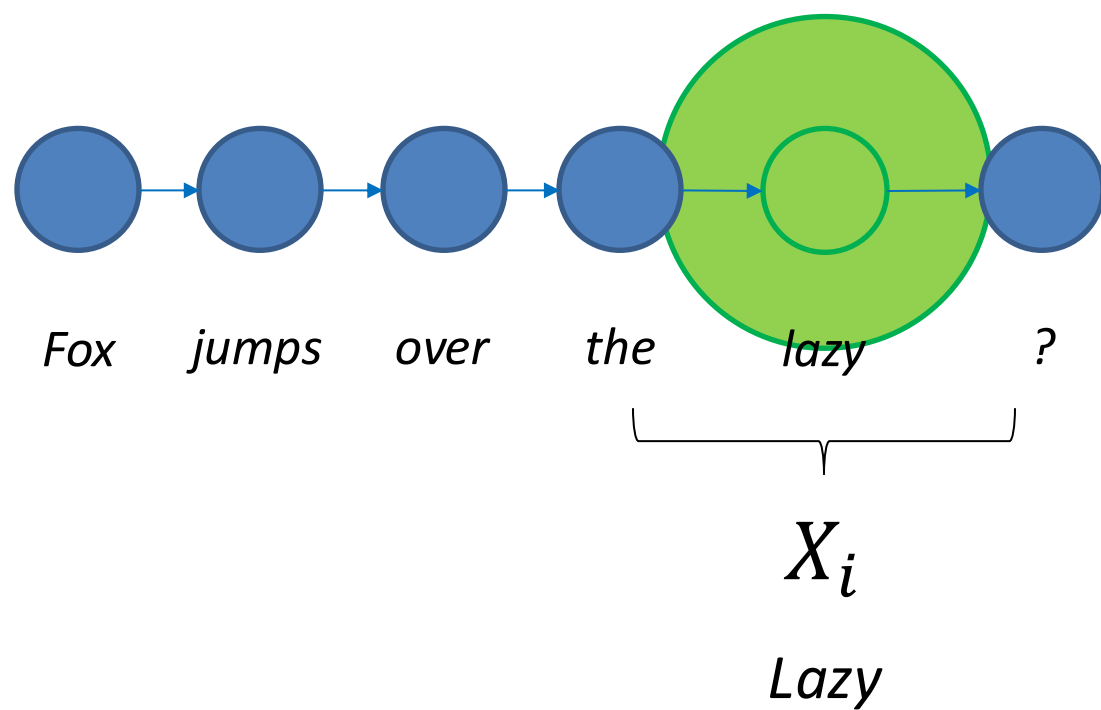
A function $f : \mathcal{X}(\Omega) \rightarrow \mathcal{X}(\Omega)$ is \mathfrak{G} -equivariant if $f(\rho(\mathfrak{g})x) = \rho(\mathfrak{g})f(x)$ for all $\mathfrak{g} \in \mathfrak{G}$, i.e., group action on the input affects the output in the same way.



Elements

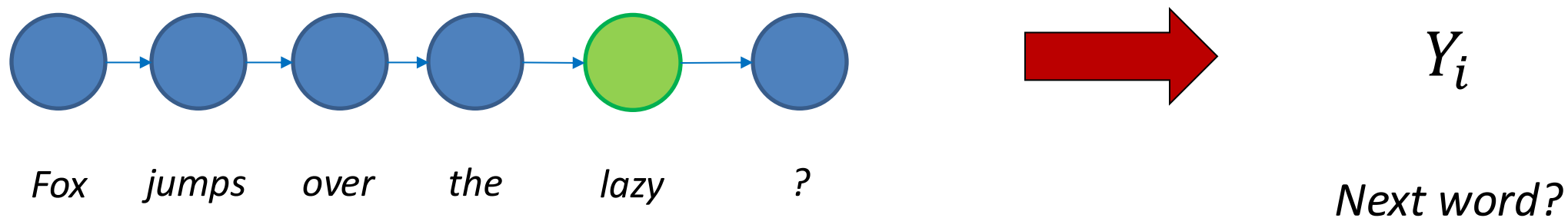


Sequences

 Y_i

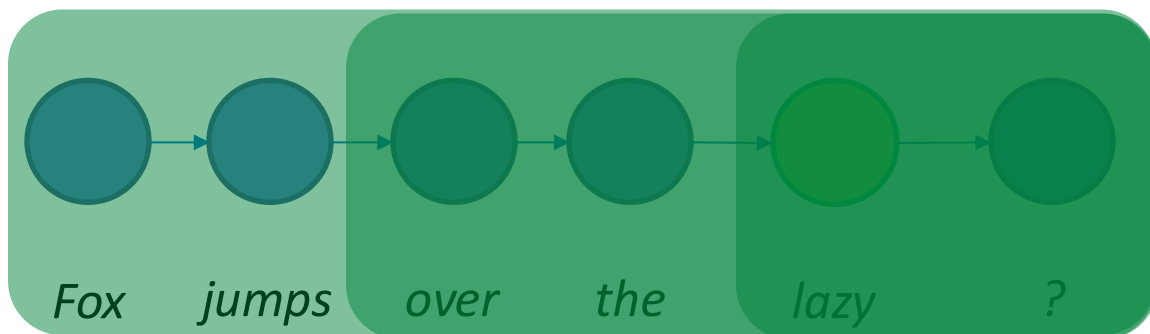
Next word?

Sequences



How do we aggregate information?

Sequences

 Y_i

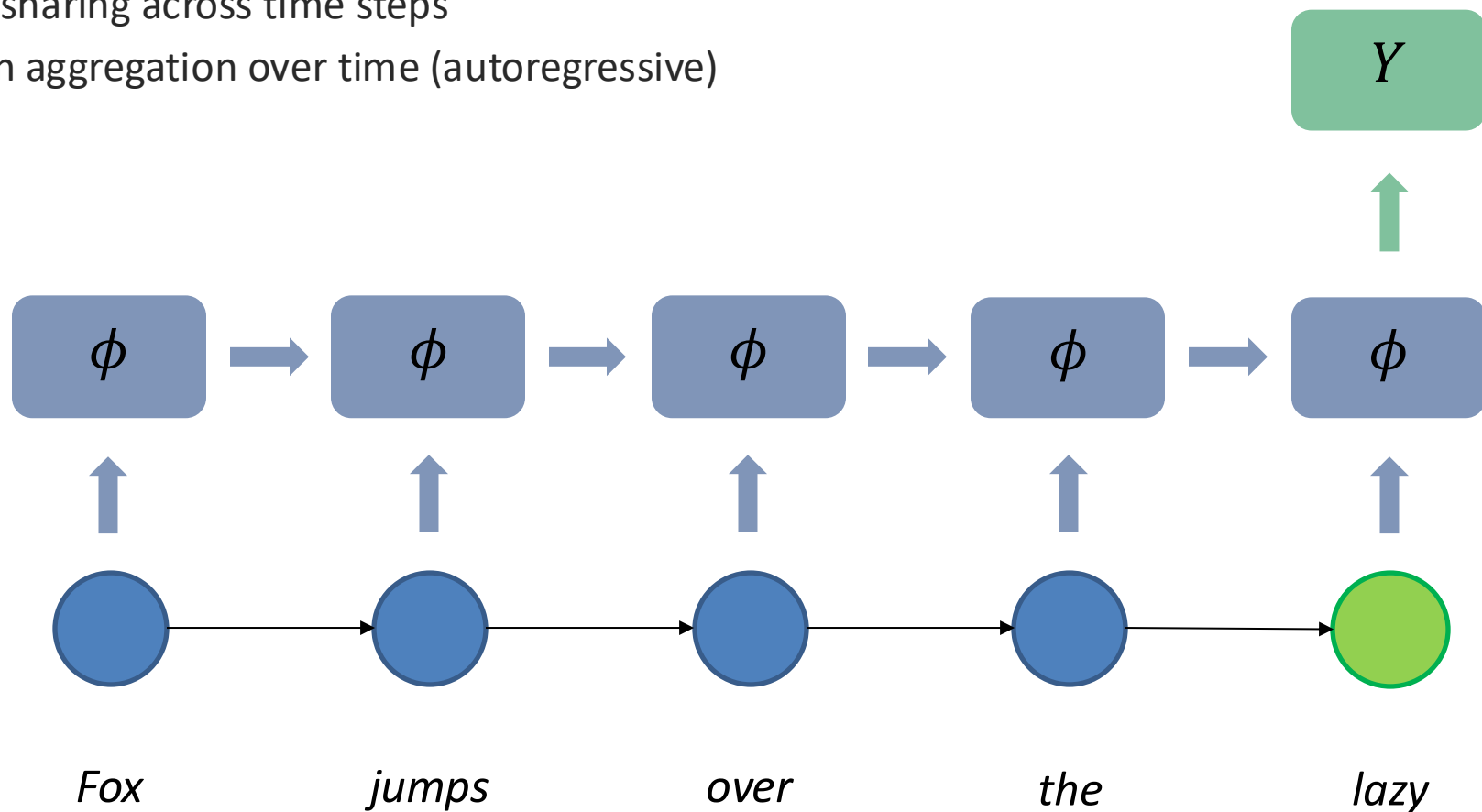
Next word?

How do we aggregate information?

Sequence Classification

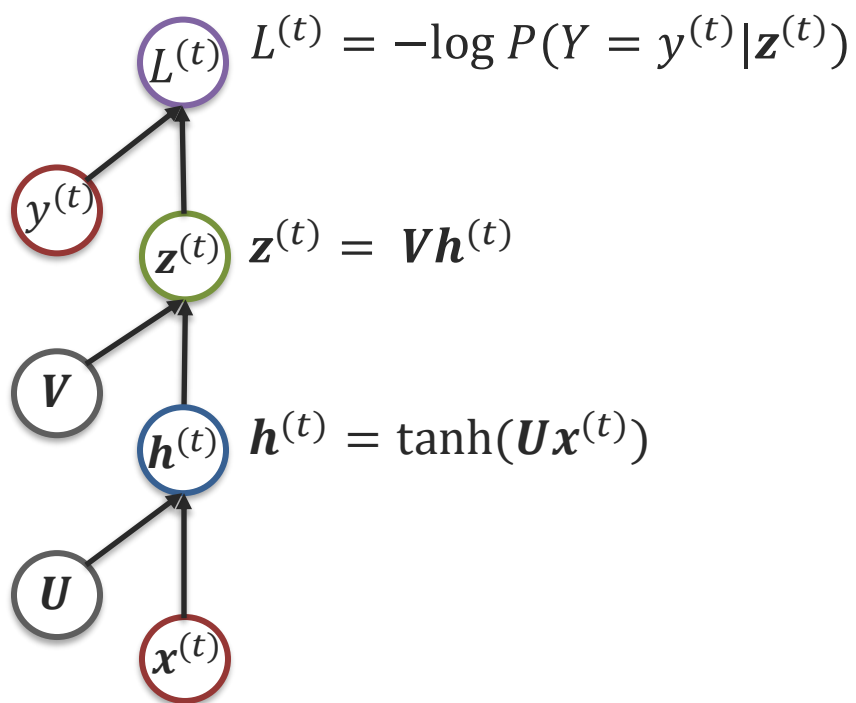
Models for sequential data must be invariant to time, but equivariant to word order.

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

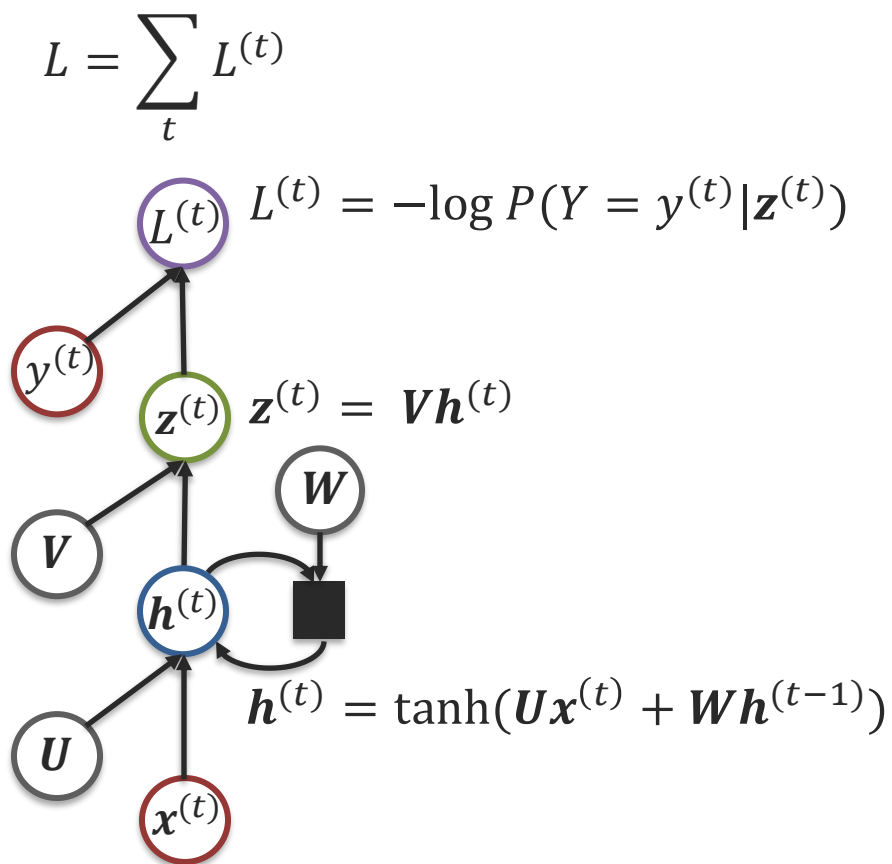


Sequence Classification

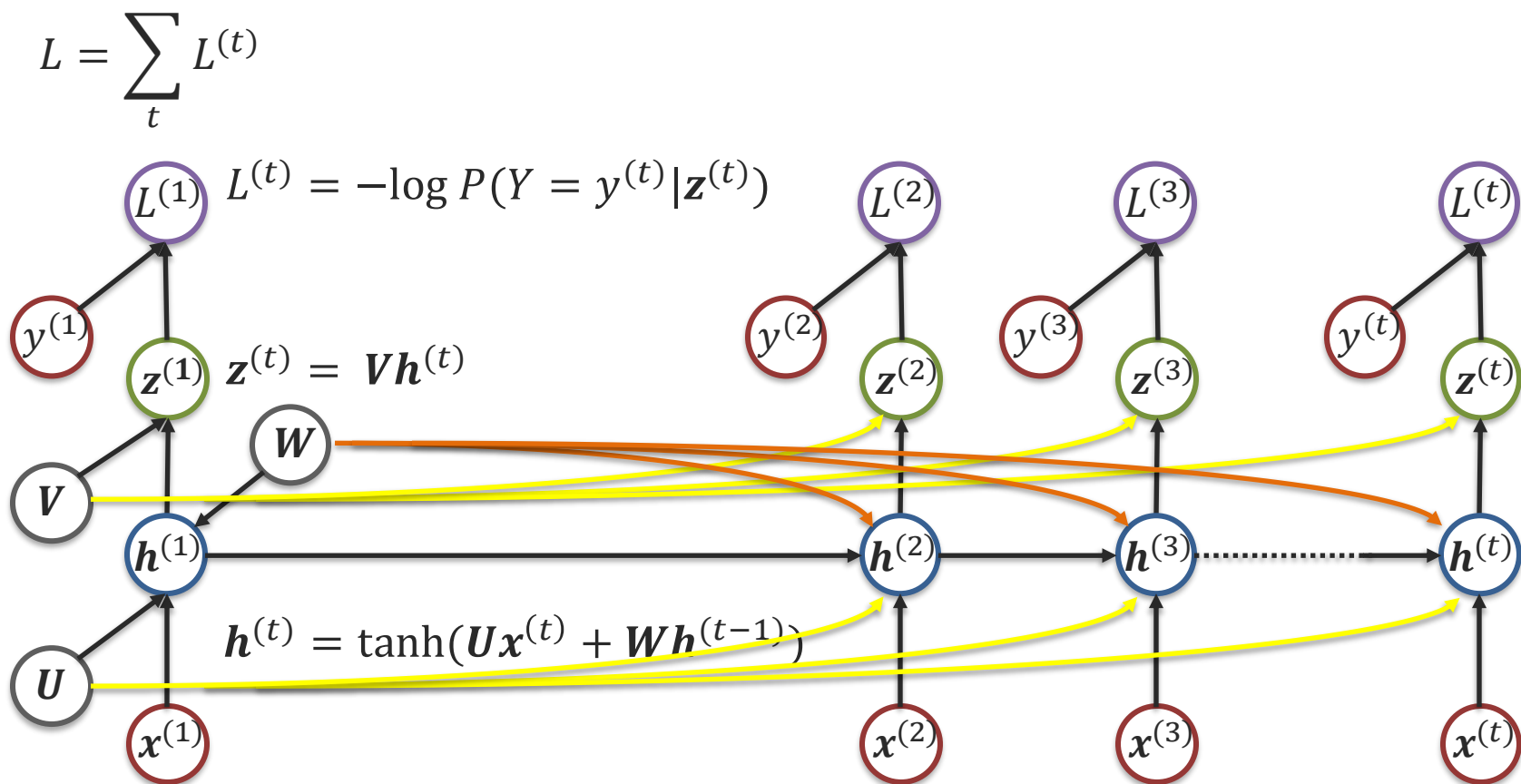
Feedforward Neural Network



Sequence Classification

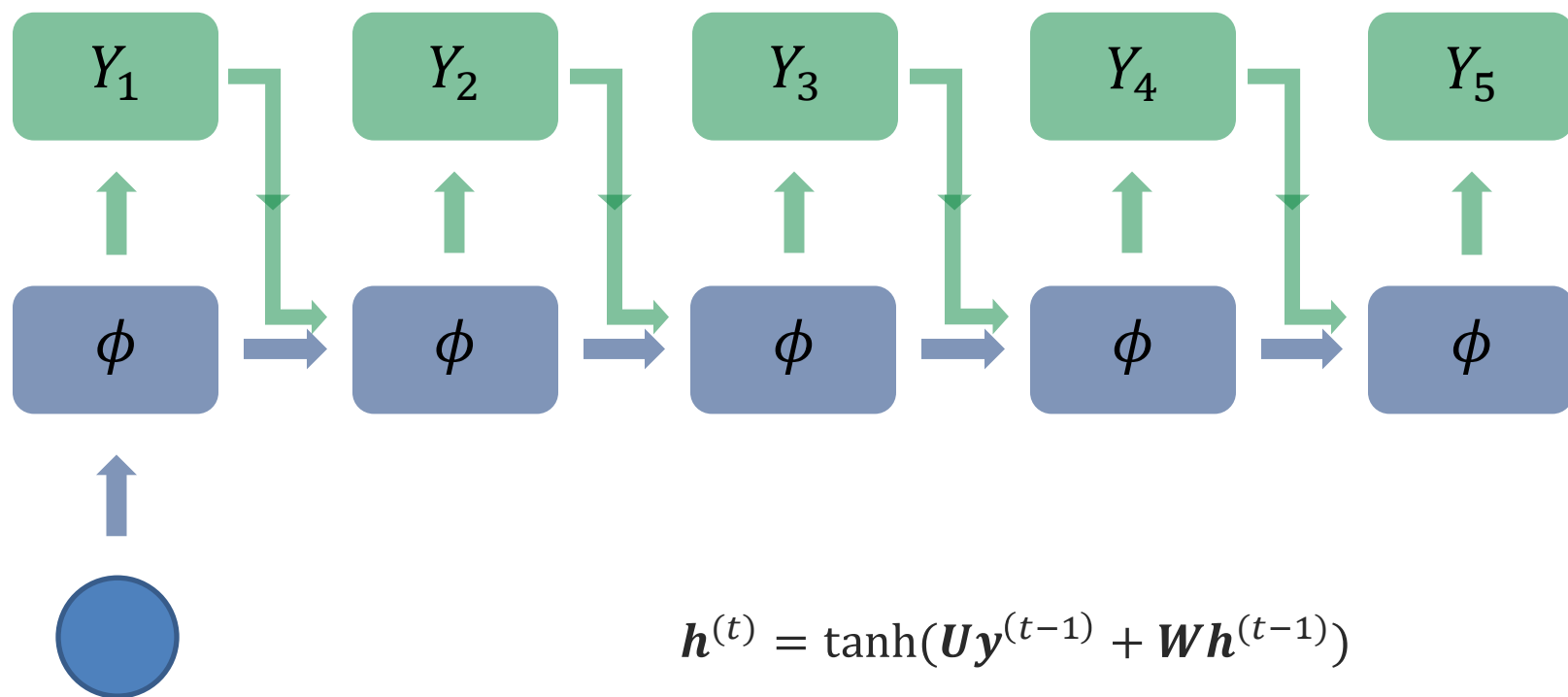


Sequence Classification



Same model parameters are used for all time steps.

Sequence Generation



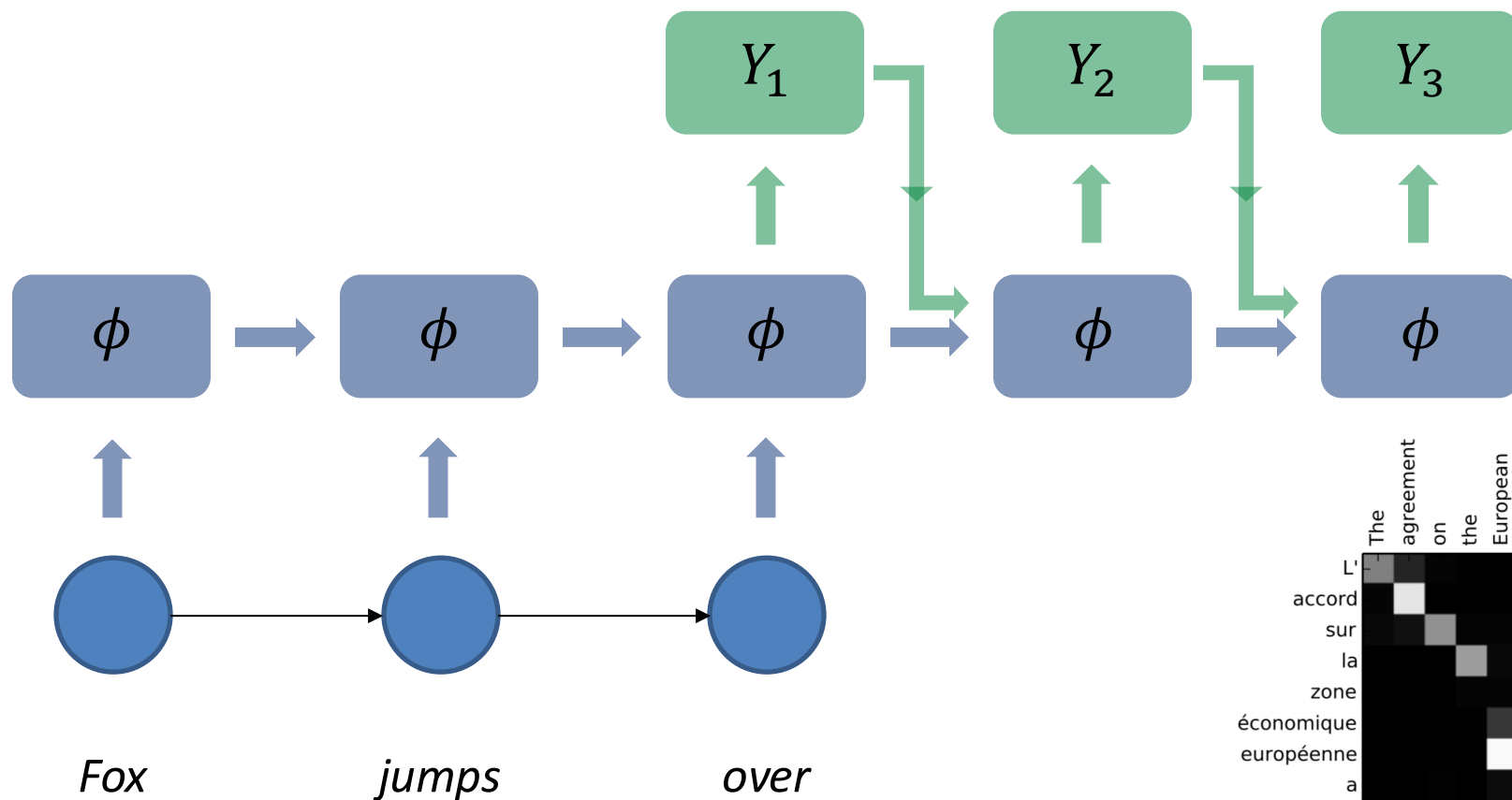
$$\mathbf{h}^{(t)} = \tanh(\mathbf{U}\mathbf{y}^{(t-1)} + \mathbf{W}\mathbf{h}^{(t-1)})$$

Fox

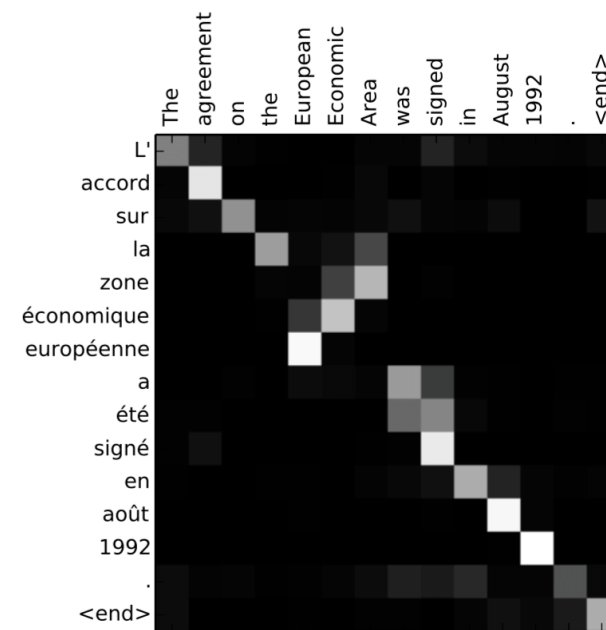
e.g, text or music generation

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

Sequence-to-Sequence Models



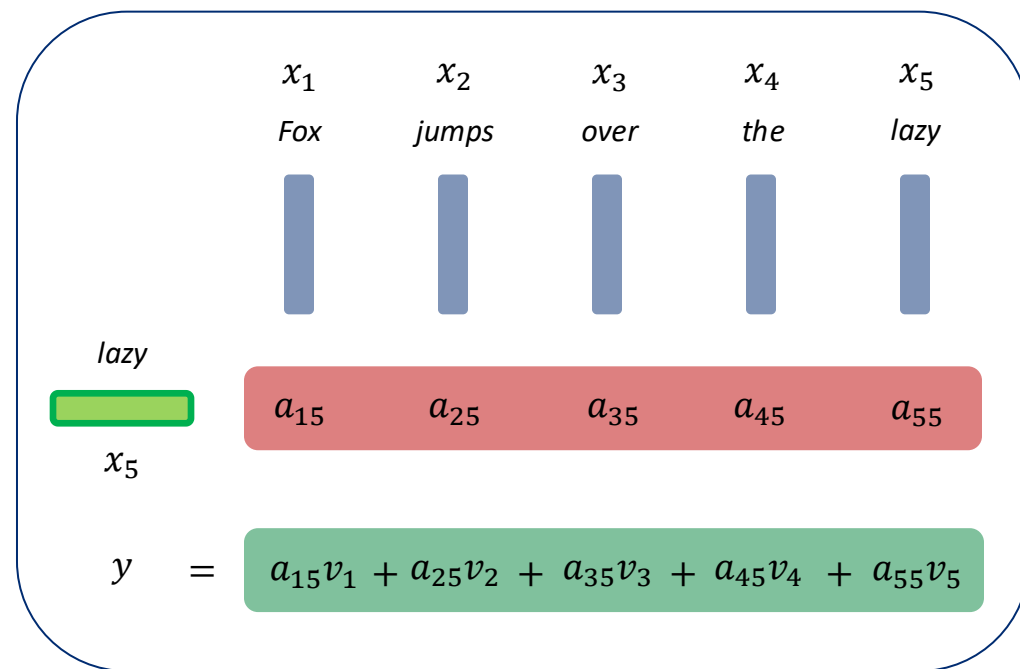
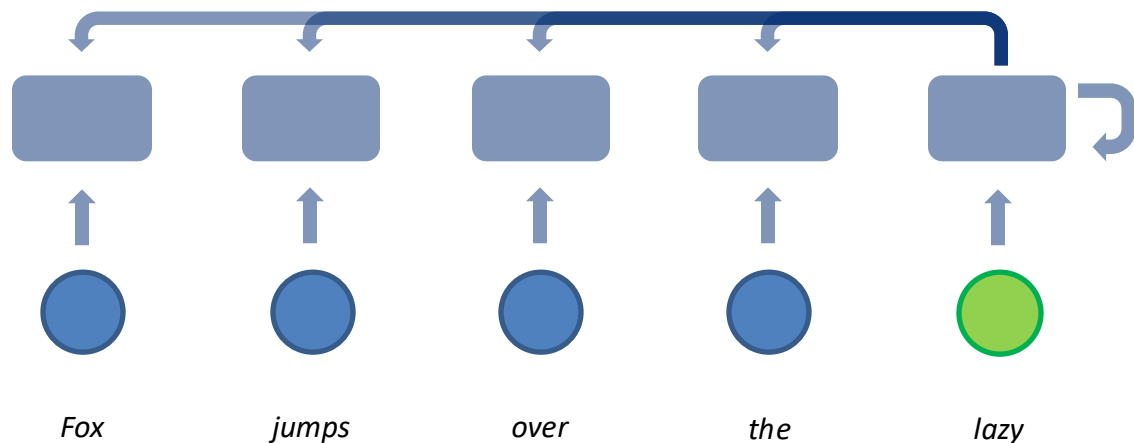
e.g, machine translation -> birth of attention-based models



Modern Sequence Models

Birth of attention-based models

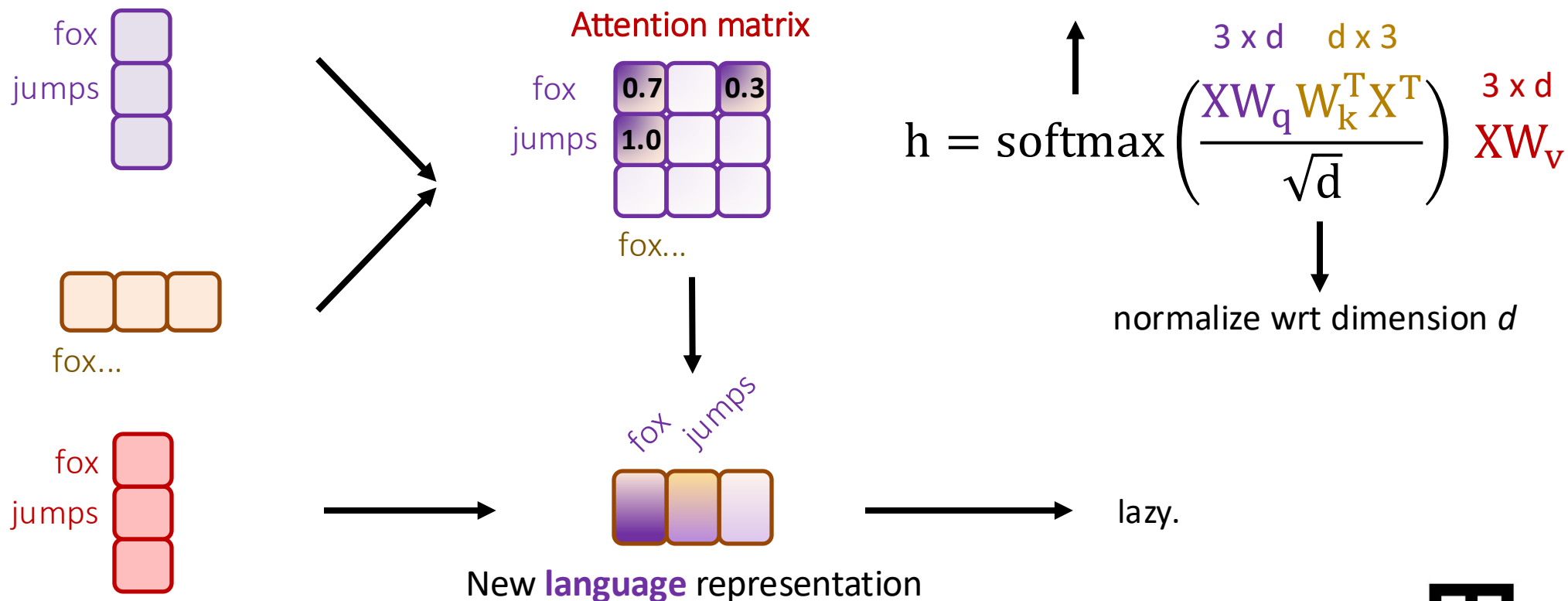
- Dynamic weights for different elements



Modern Sequence Models

Birth of attention-based models

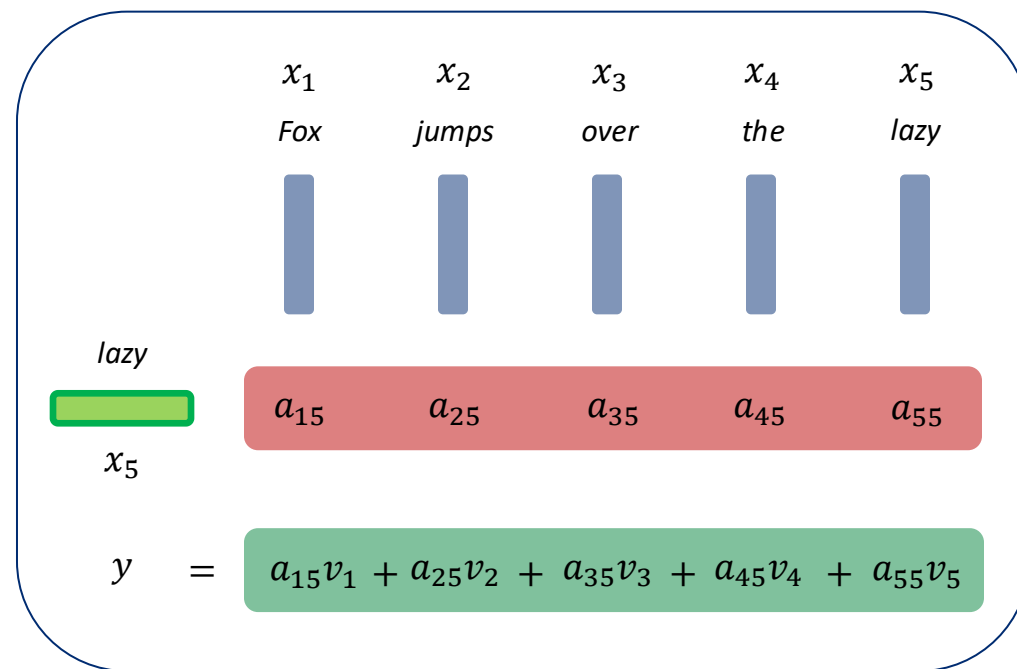
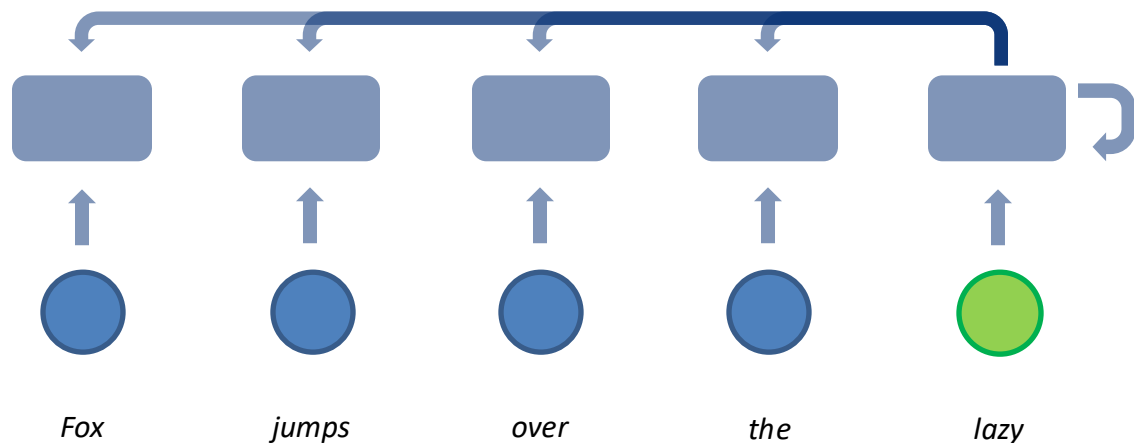
- Dynamic weights for different elements



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)



$$h = \text{softmax} \left(\frac{\overset{T \times d}{X} \overset{d \times T}{W_q} \overset{T \times T}{X^T}}{\sqrt{d}} \right) \overset{T \times d}{XW_v}$$

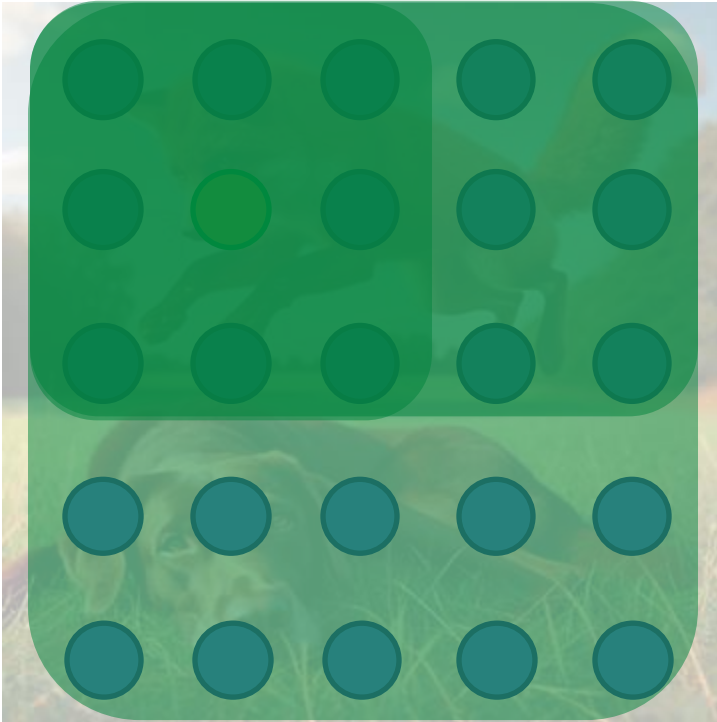
Spatial Data

 Y_i

Is there a fox?

How do we aggregate information?

Spatial Data

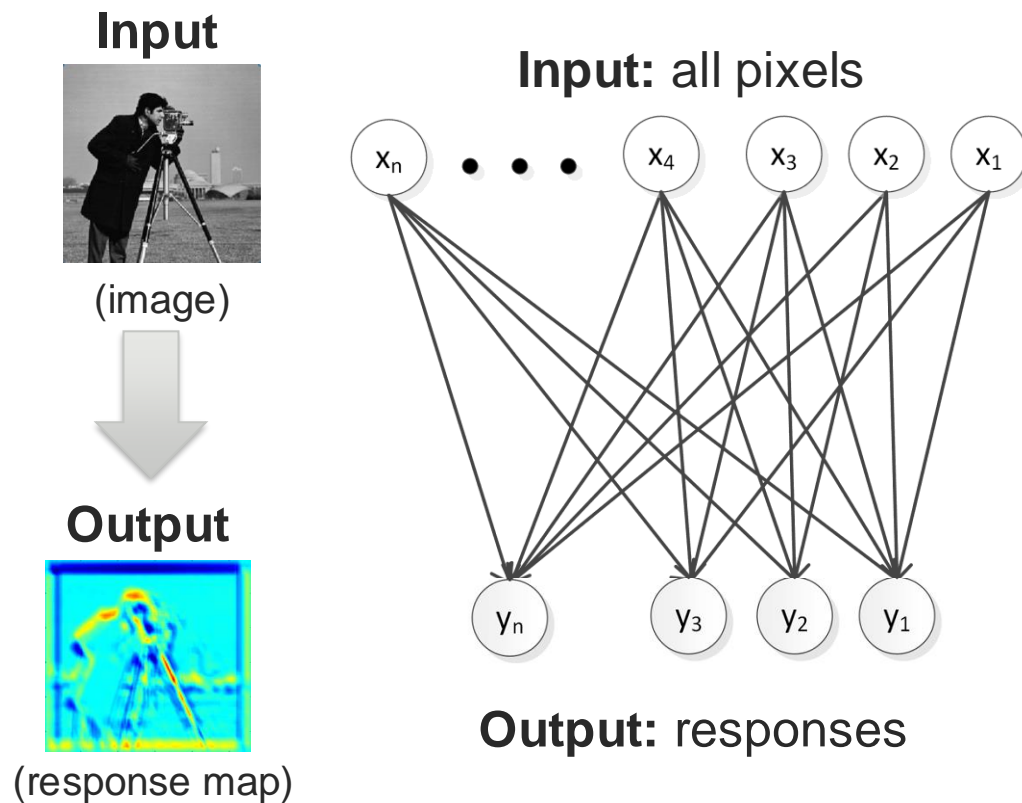
 Y_i

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!

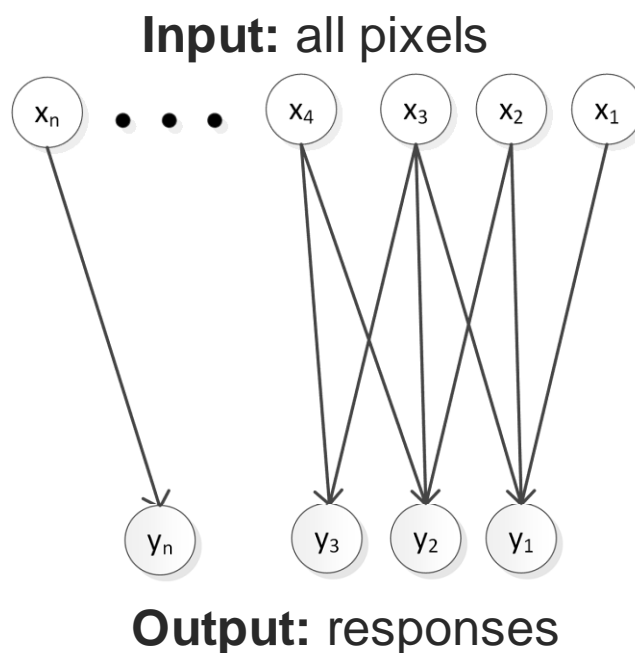
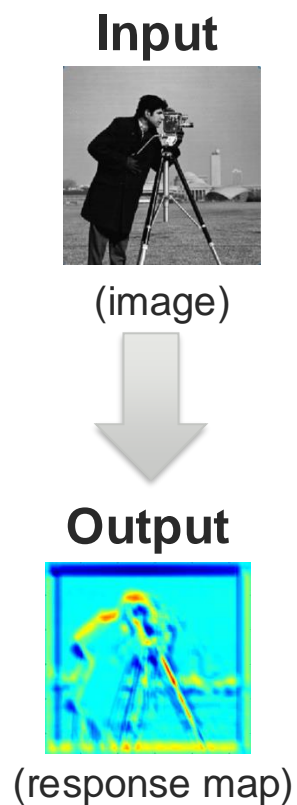
200 × 200 image
requires
40,000 × n
(where n is size of output)
parameters

**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

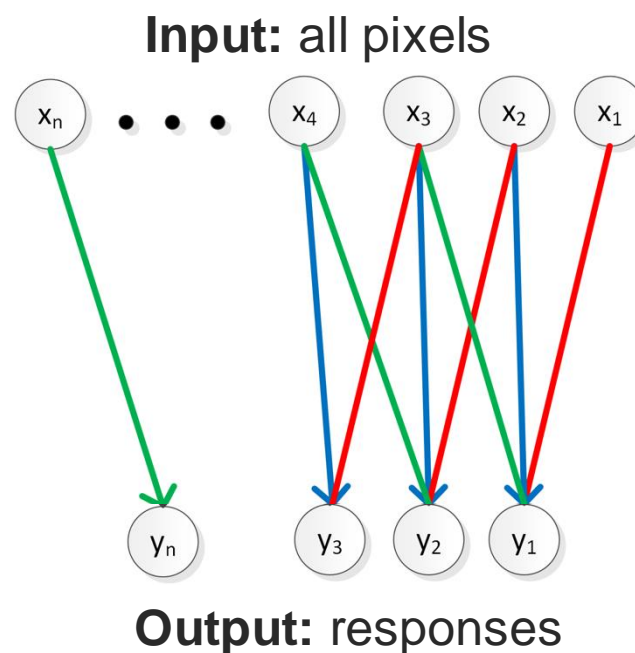
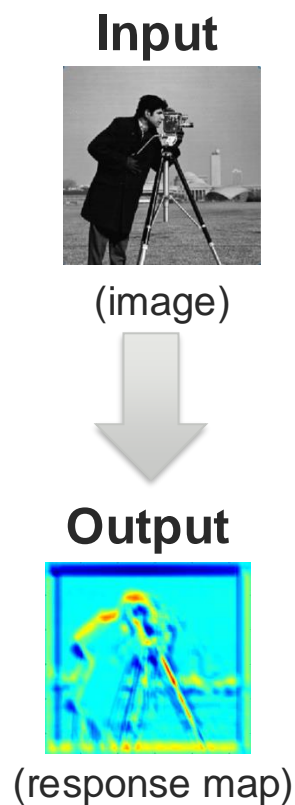


Example with
1D filter:

w_1	w_2	w_3
-------	-------	-------

Convolutional Neural Networks

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



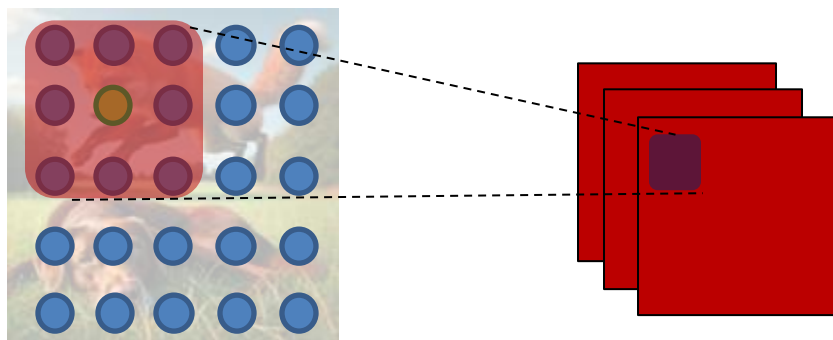
Example with
1D filter:



Convolutional Neural Networks

Models for spatial data need to be invariant to spatial translation

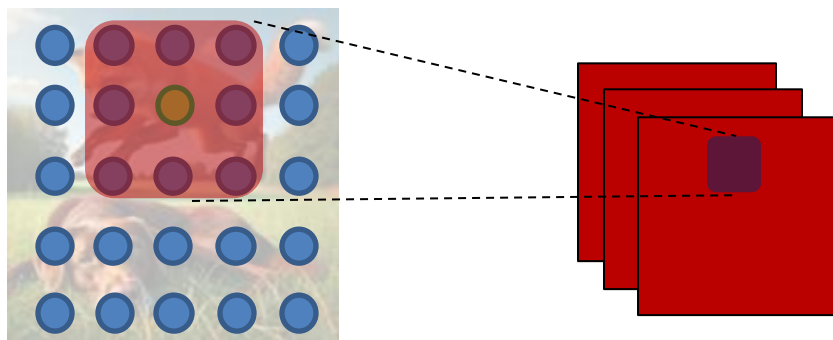
1. Parameter sharing across $k \times k$ convolutional filter



Convolutional Neural Networks

Models for spatial data need to be invariant to spatial translation

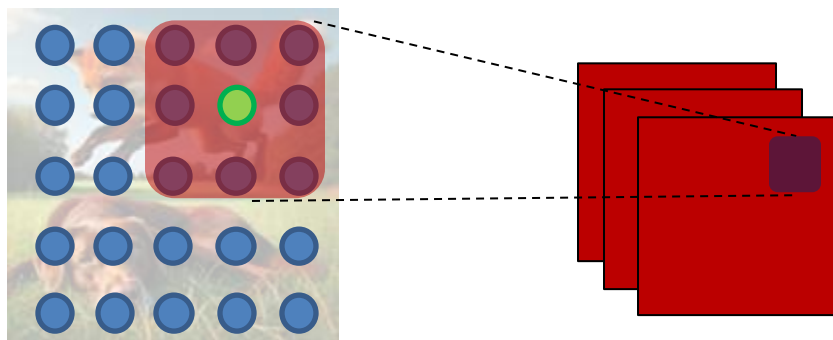
1. Parameter sharing across $k \times k$ convolutional filter



Convolutional Neural Networks

Models for spatial data need to be invariant to spatial translation

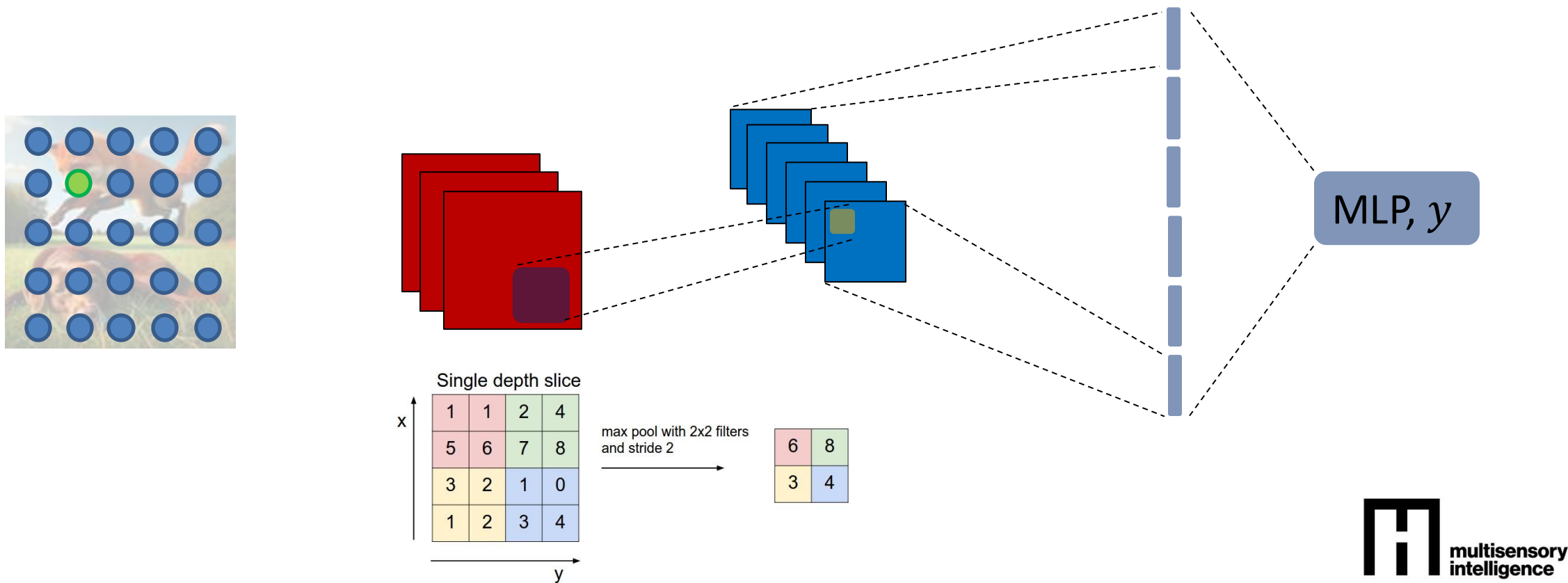
1. Parameter sharing across $k \times k$ convolutional filter



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



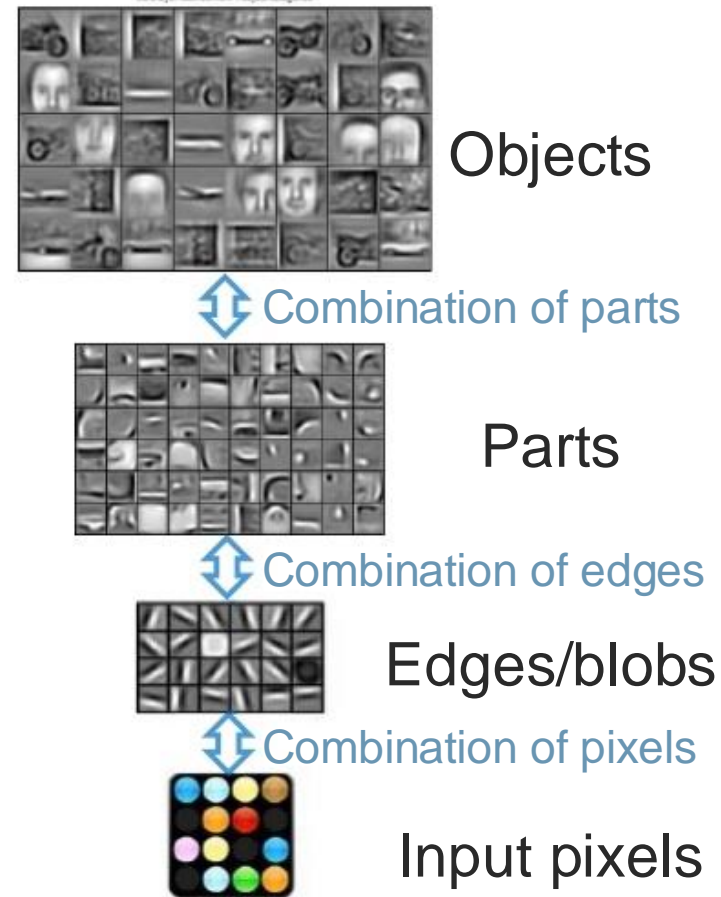
Convolutional Neural Networks

Multiple convolutional layers

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

Multiple pooling layers

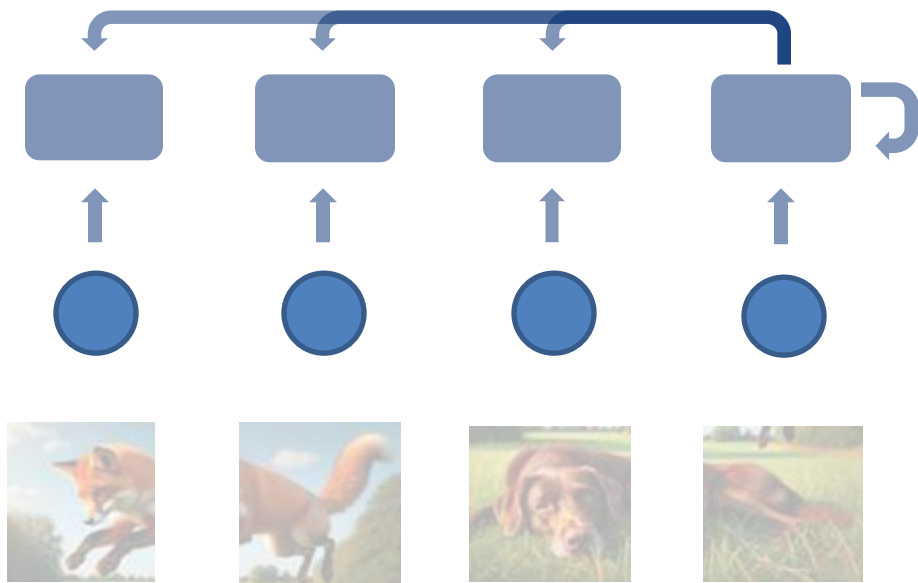
→ Allows the network to learn increasingly abstract & summarized information



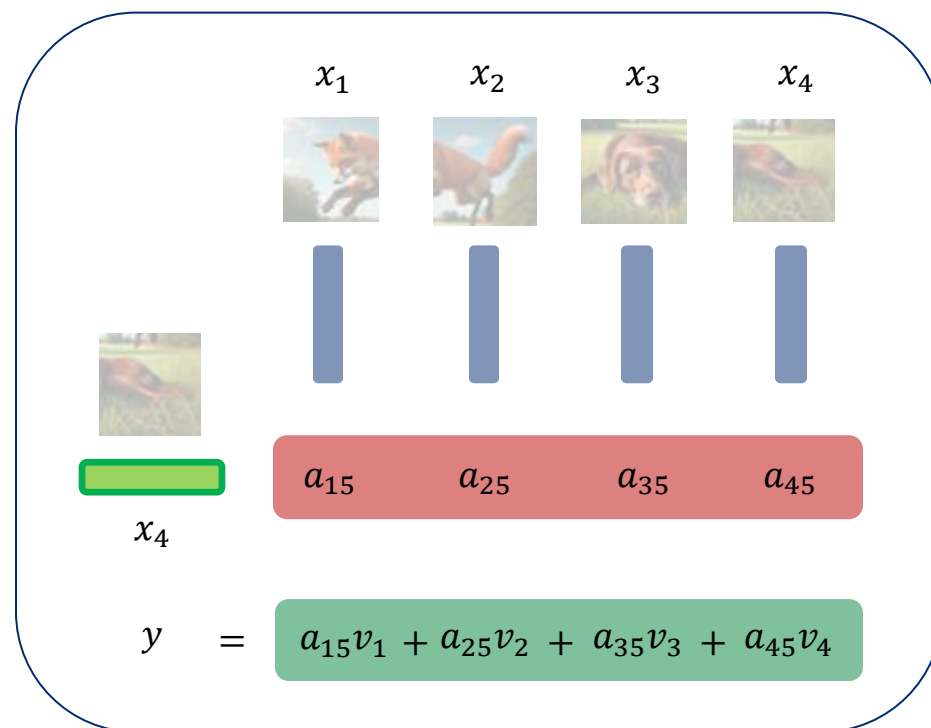
Vision Transformer

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



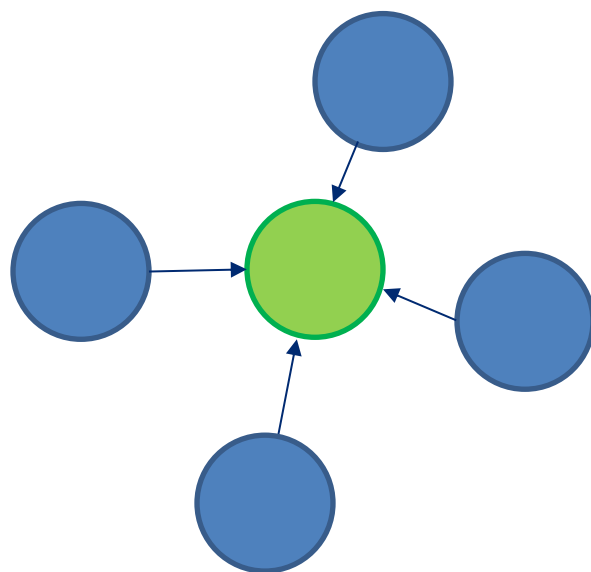
$$h = \text{softmax} \left(\frac{\overset{T \times d}{XW_q} \overset{d \times T}{W_k^T X^T}}{\sqrt{d}} \right) \overset{T \times d}{XW_v}$$



Vision Transformer



Graphs

 Y_i

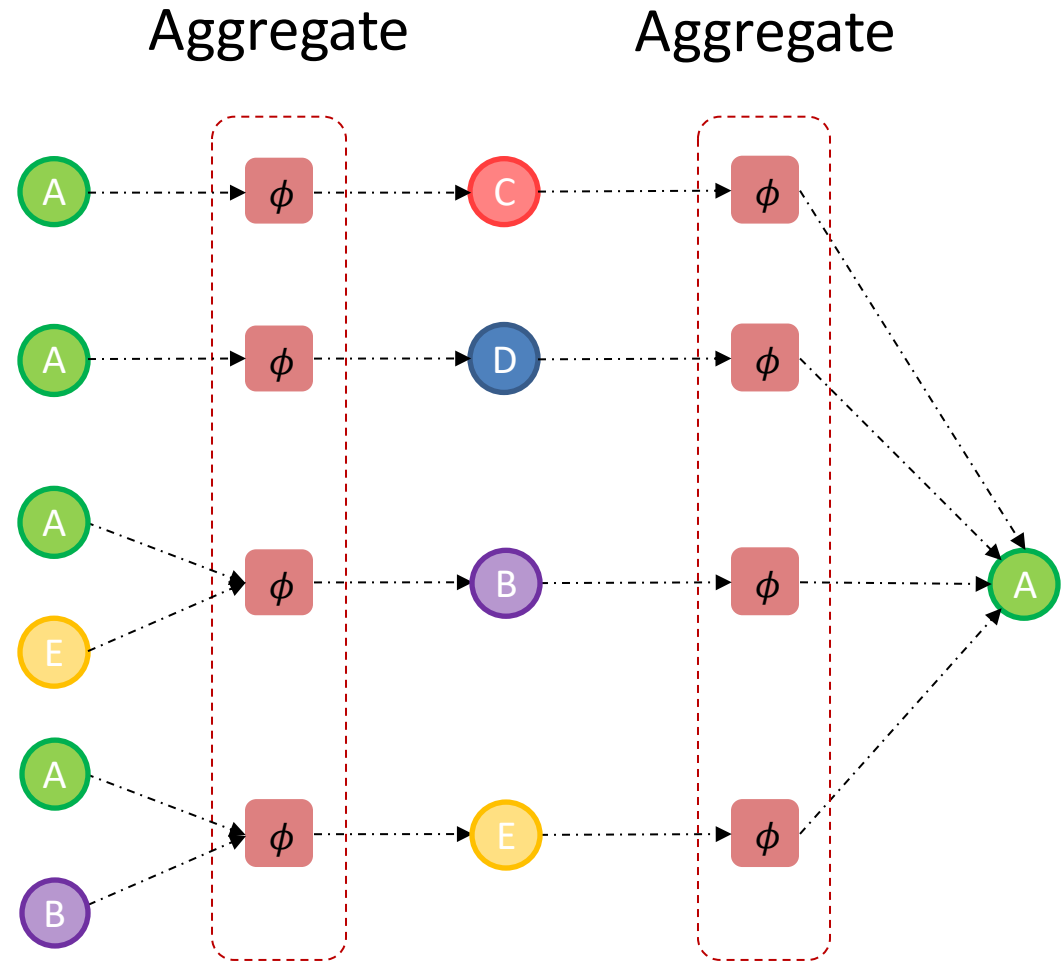
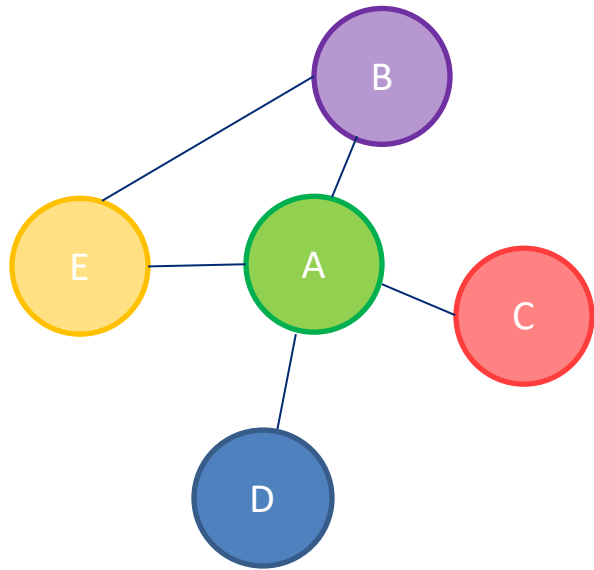
What molecule is this?

How do we aggregate information?

Graph Neural Networks

Models for graph data:

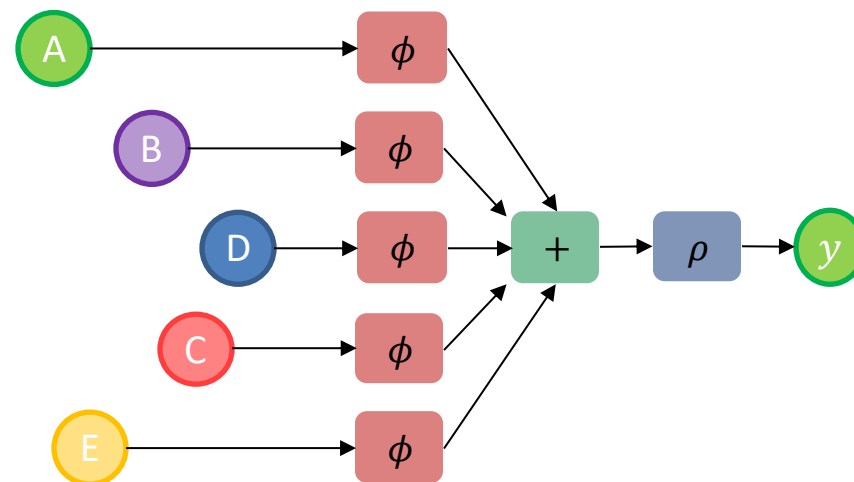
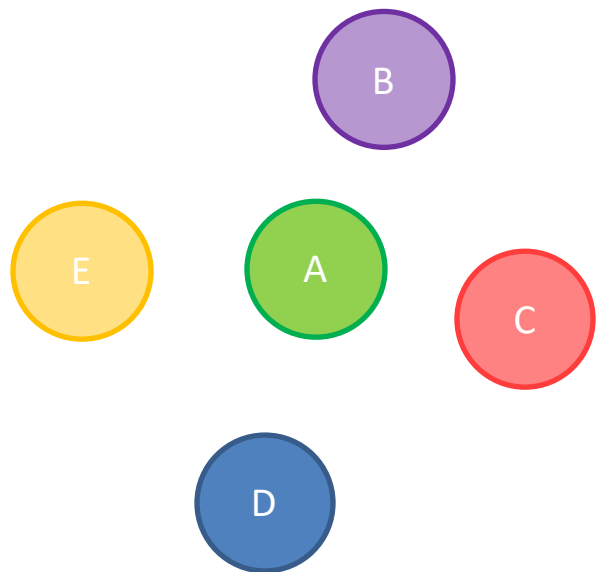
1. Parameter sharing across nodes
2. Information aggregation over neighbors (edges)



Graph Recover Sets

Sets are graphs with only nodes, no edges

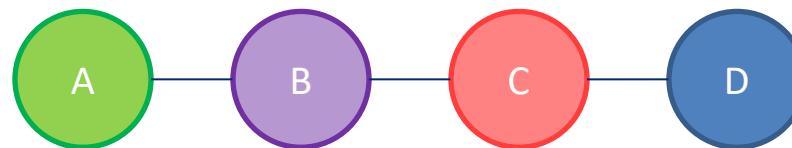
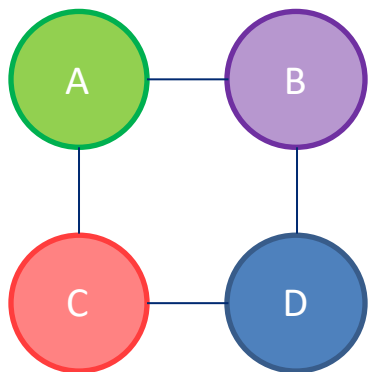
1. Parameter sharing across nodes -> set elements
2. Information aggregation over neighbors -> no neighbors



Graph Recover Spatial and Temporal Data

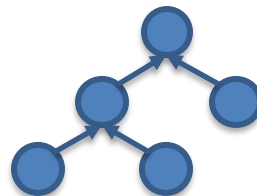
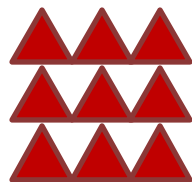
Spatial data and sequential data

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

Reading assignment due tomorrow Wednesday (2/26).

For project:

- Project proposal due tonight (2/25). Email to me.
- Meet with me 2-3pm if need feedback about proposal ideas.

This Thursday (2/27): first reading discussion on **data and learning**.