

Computer Vision: Fall 2022 — Lecture 16

Dr. Karthik Mohan

Univ. of Washington, Seattle

November 30, 2022

References

Generic ML/DL

- ① [Good Book for Machine Learning Concepts](#)
- ② [Deep Learning Reference](#)

CNN

- ① [Convolutional Neural Networks for Visual Recognition](#)
- ② [Convolutional Neural Net Tutorial](#)
- ③ [CNN Transfer Learning](#)
- ④ [PyTorch Transfer Learning Tutorial](#)

CNN Publication References

CNN surveys

- ① Convolutional Neural Networks: A comprehensive survey, 2019
- ② A survey of Convolutional Neural Networks: Analysis, Applications, and Prospects, 2021

CNN Archs

- ① GoogLeNet
- ② Top models on ImageNet
- ③ ResNet ILSVRC paper

Object Detection and Image Segmentation References

Object Detection

- 1 A survey of modern deep learning based object detection methods
- 2 YOLO Survey
- 3 YOLO Original Paper

Image Captioning

- 1 From Show to Tell: A survey on Deep Learning-based Image Captioning
- 2 A survey of image captioning models

↳ Lightweight

Today

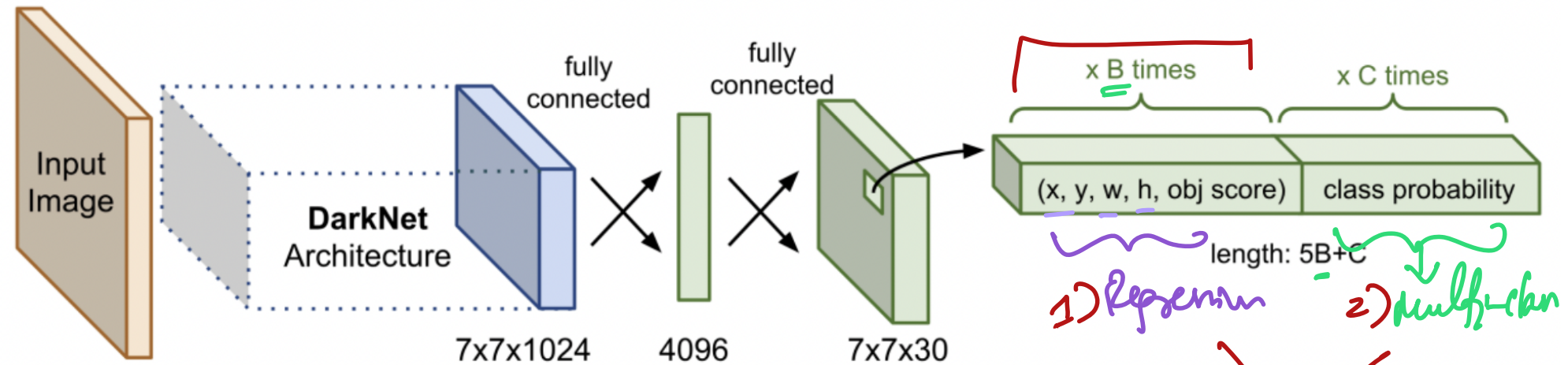
- 1 Recap of YOLO for object detection
- 2 Image Captioning intro and models

YOLO - Single Stage Detection

You only look once!

Image

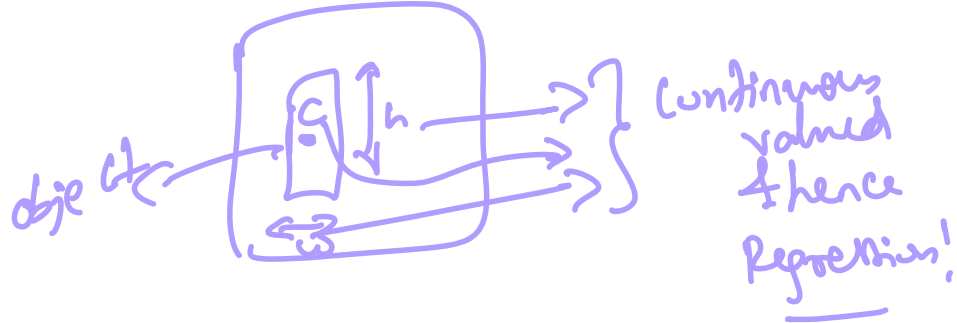
1) multiple Bounding Boxes for each grid cell
2) class



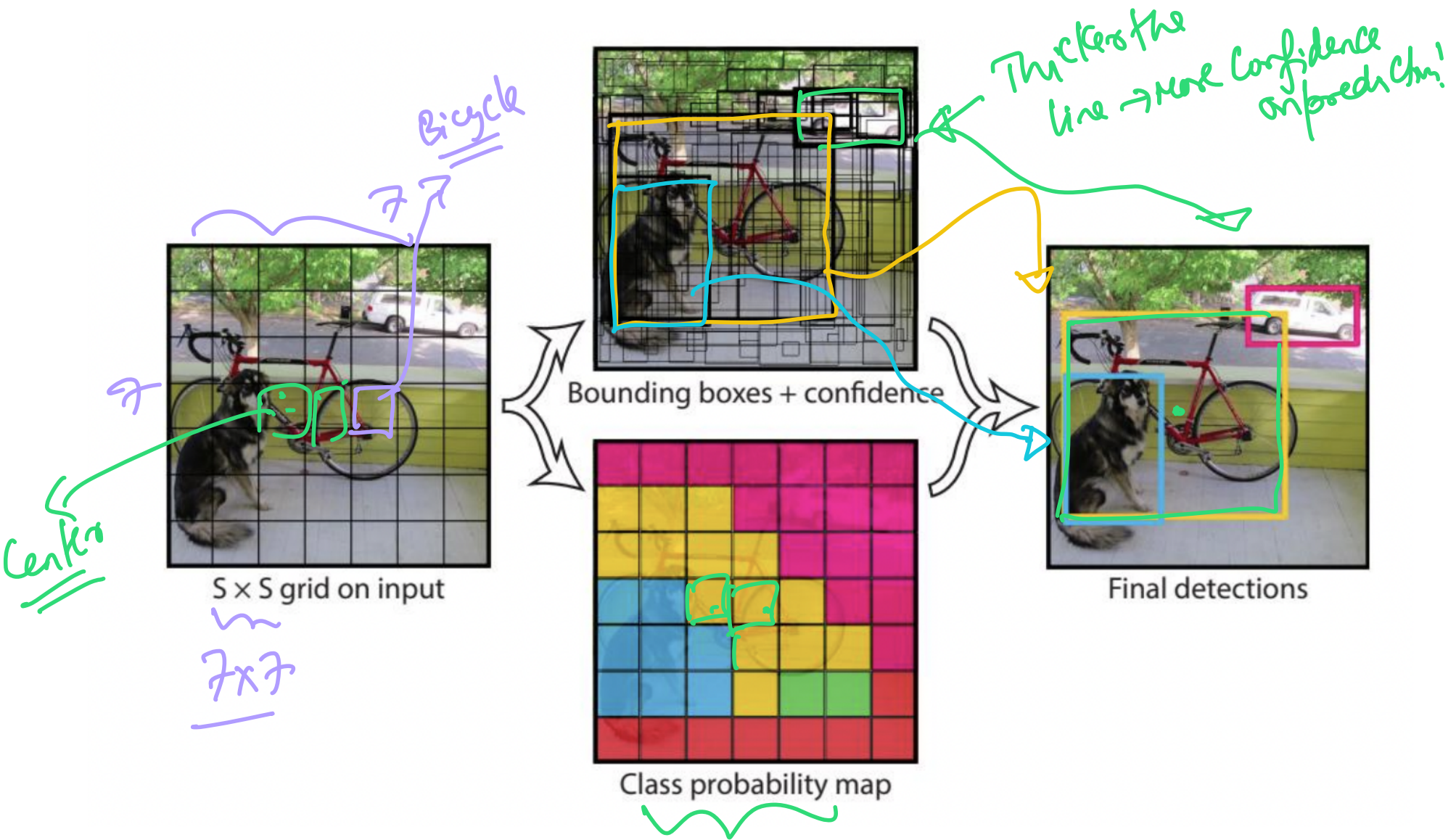
1) Regression *2) multi-class*

Combine the two

Yolo v1
v2
v3
v4
v5
Improvement



YOLO Breakdown



YOLO benefits

- 1 **Single pass**


YOLO benefits

- ① **Single pass**
- ② **Fast**

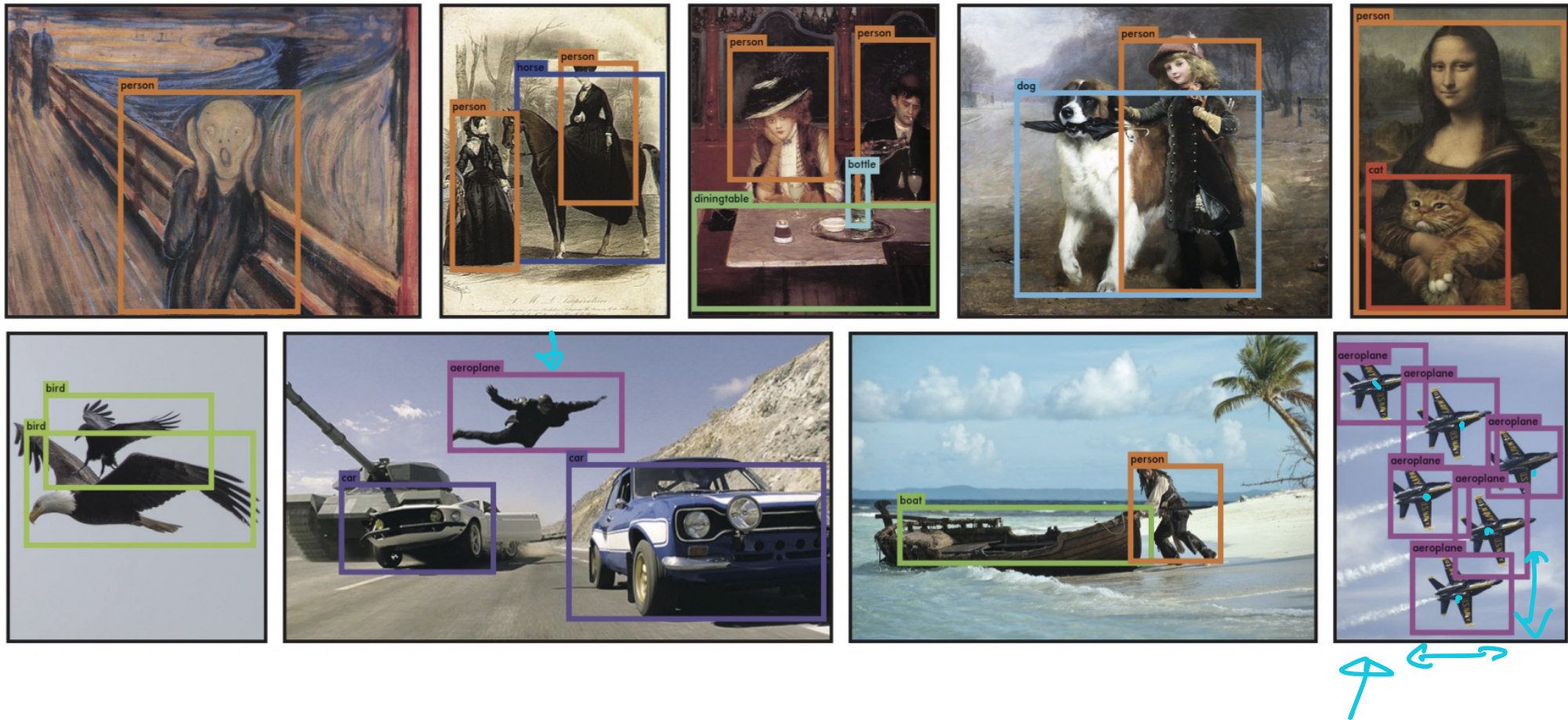
YOLO benefits

- ① **Single pass**
- ② **Fast**
- ③ **Global reasoning**

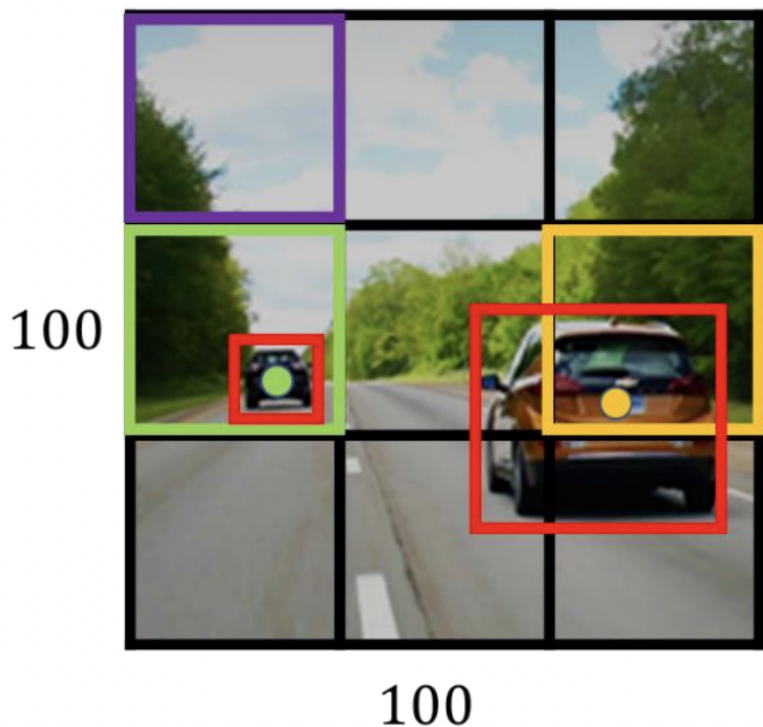
YOLO benefits

- ① **Single pass**
 - ② **Fast**
 - ③ **Global reasoning**
 - ④ **More generalized representations**
- 

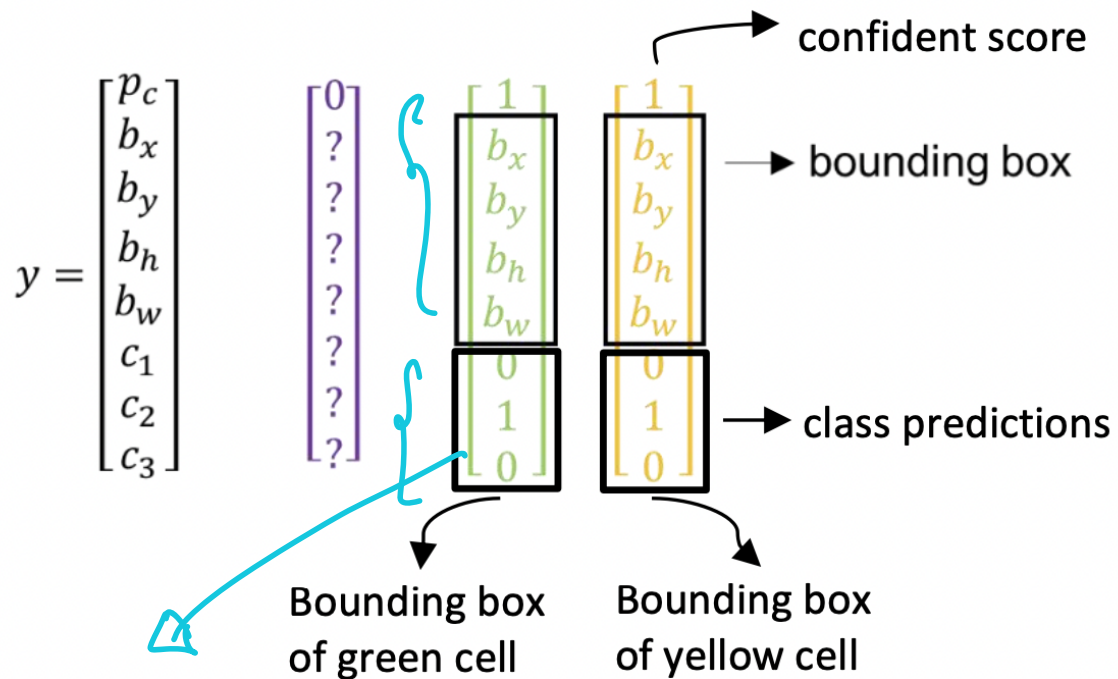
YOLO examples



YOLO labels



Labels for training for each grid cell:

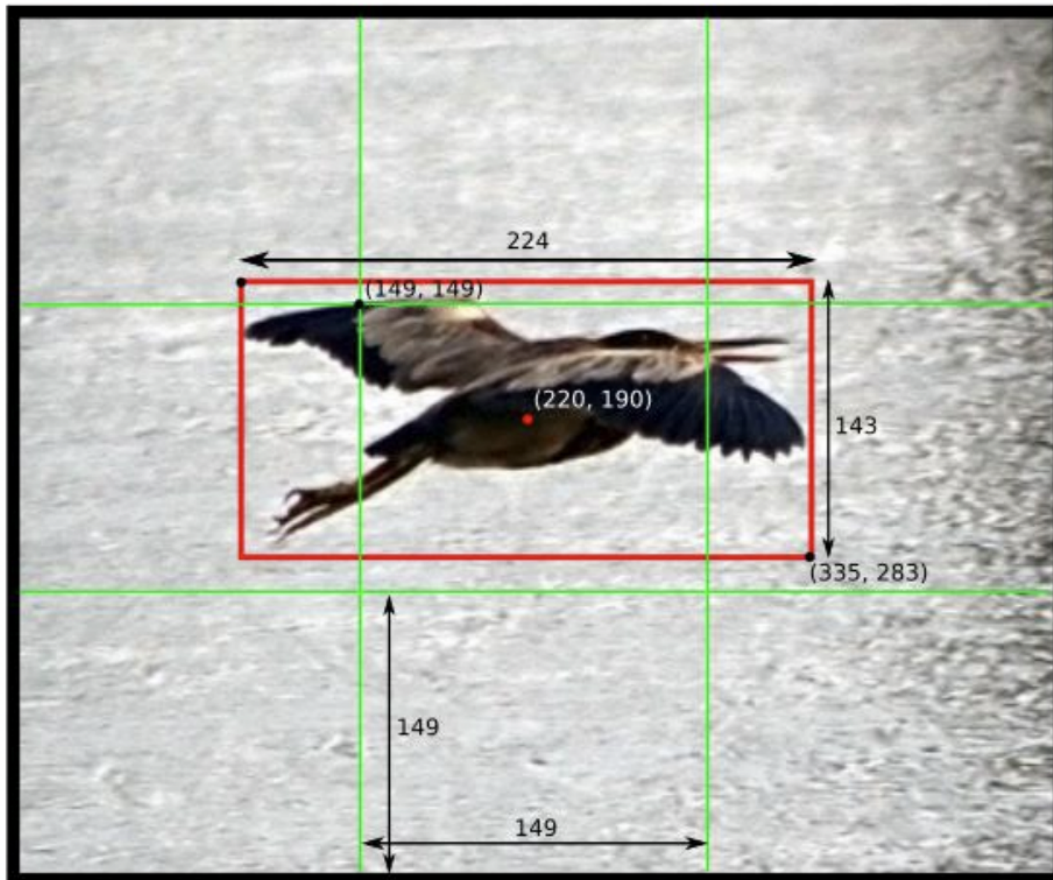


one-hot
Encoding

$S \times S \times (S \times B + C)$
Labels for an image

YOLO Bounding Box

(0, 0)



$$x = (220 - 149) / 149 = 0.48$$

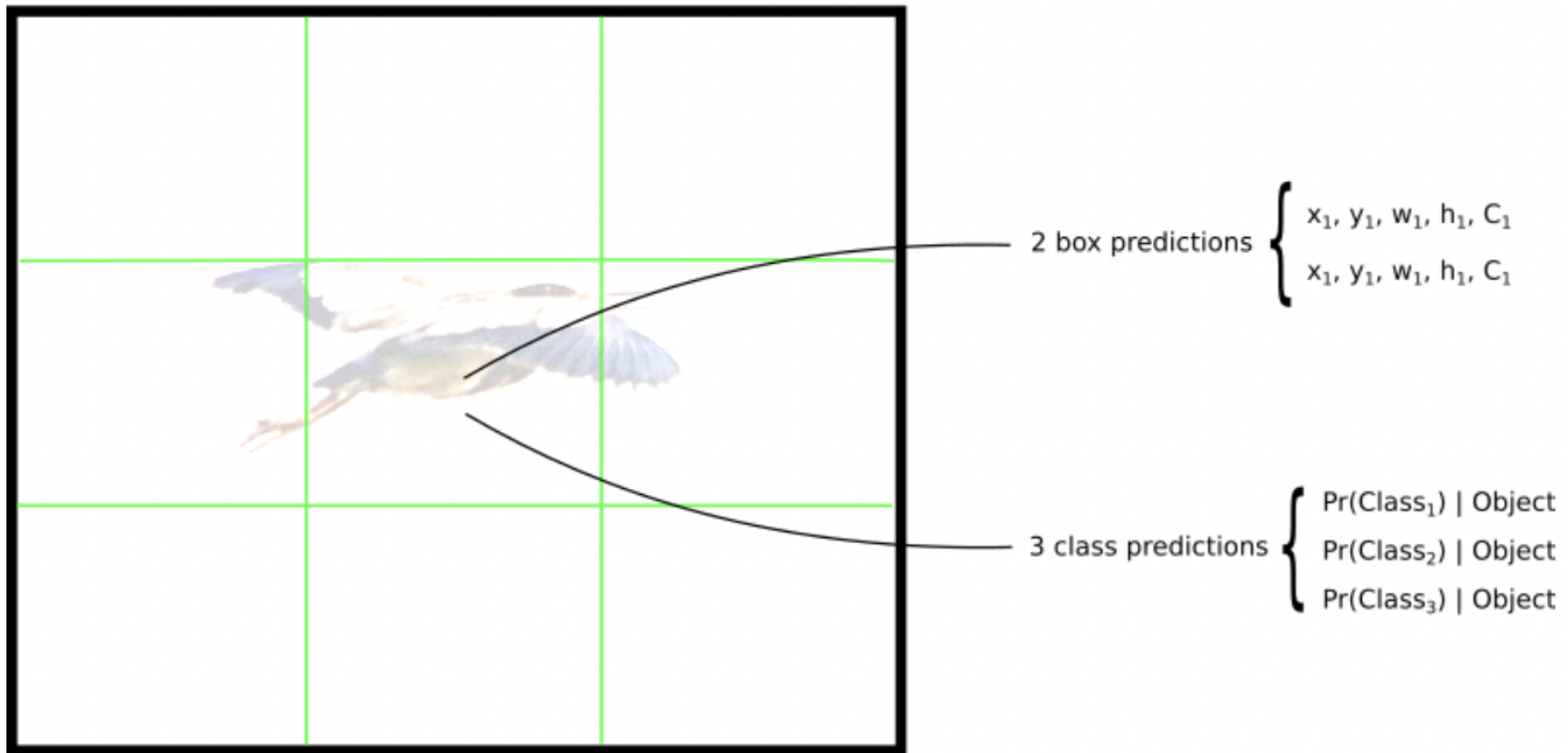
$$y = (190 - 149) / 149 = 0.28$$

$$w = 224 / 448 = 0.50$$

$$h = 143 / 448 = 0.32$$

(447, 447)

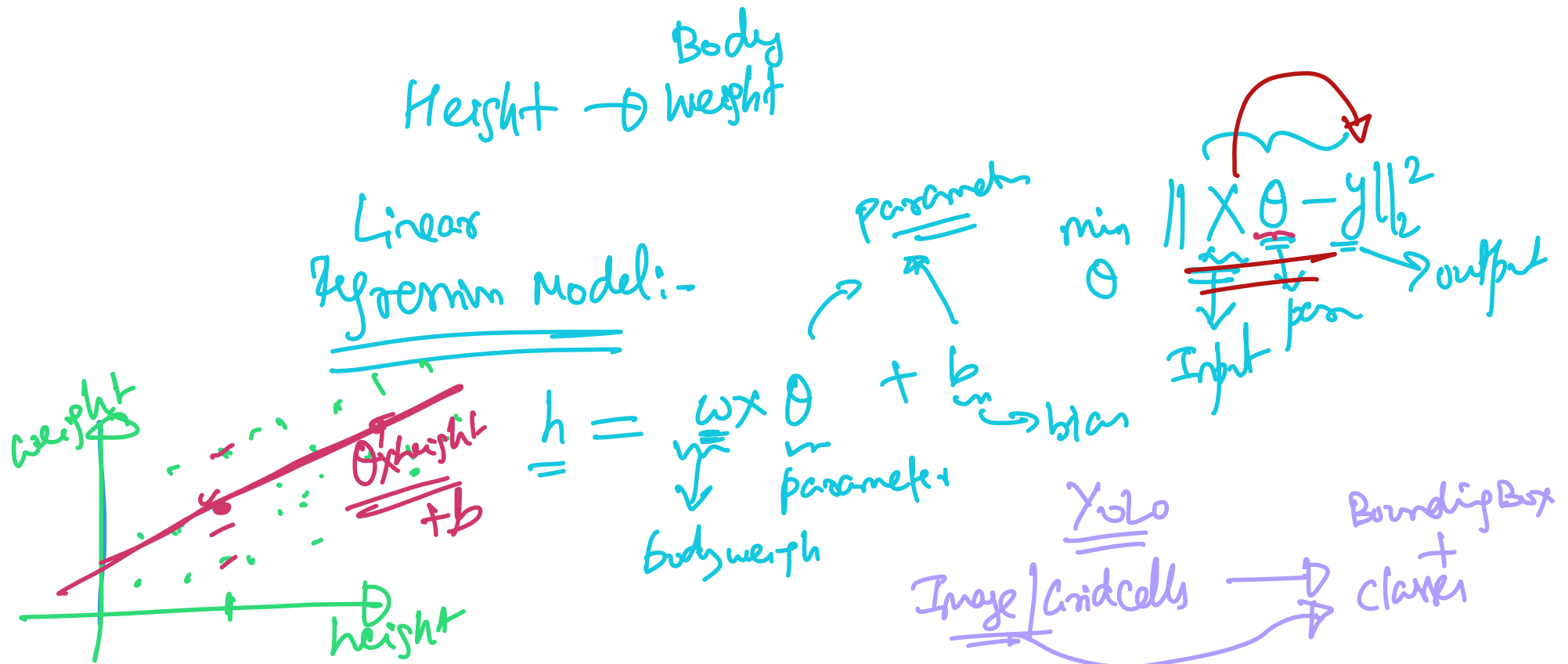
YOLO Bounding Box 2



YOLO and Regression

YOLO Loss Function - Regression!

YOLO loss function turns out to be just like a Regression Loss! Why Regression?

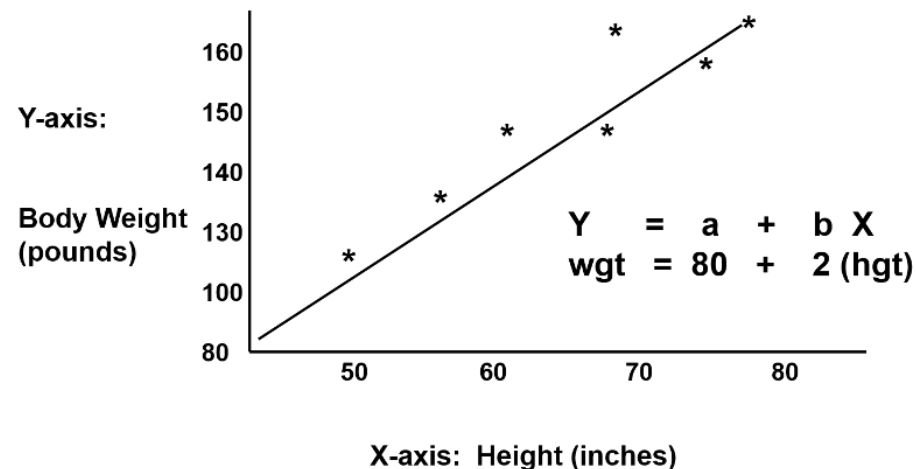


YOLO and Regression

YOLO Loss Function - Regression!

YOLO loss function turns out to be just like a Regression Loss! Why Regression?

Linear Regression Classic Example



ICE #1

$e_1 = +0.1$ $e_2 = -0.2$ $e_3 = +0.1$
Sharpness (pred, truth)
(0.5, 0.4), (0.6, 0.8), (0.3, 0.2)

Regression (2 mins)

You want to predict the 'sharpness' of an image when the input is an image. Sharpness for this exercise is defined on a continuous scale between 0 and 1. The training data looks like $\{Image, Sharpness\}$ where Image is the input and Sharpness (on a continuous scale) is the output. You devise an ingenious loss function as follows: Take the prediction \hat{y}_i of the sharpness, subtracts it from the ground truth sharpness y_i , and obtain the error, e_i . Define the loss, $L = \frac{1}{N} \sum_i e_i$. You then minimize the loss as you hope a good model for sharpness would give zero errors and hence a close to zero loss. Optimizing the loss function:

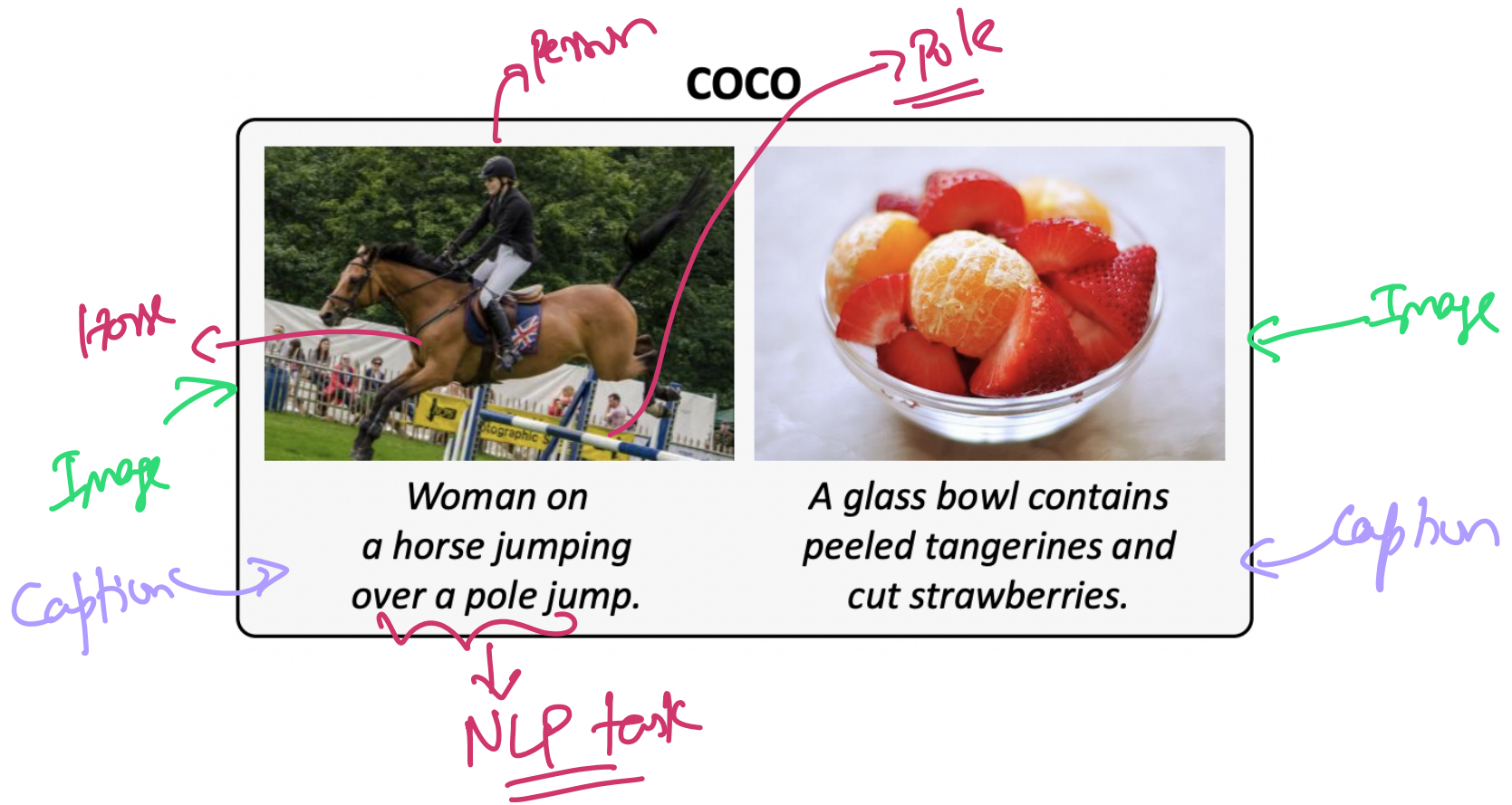
- 1 Will help you train a good model for sharpness
- 2 Is a good idea but may have to watch out for overfitting
- 3 Would not be a good idea
- 4 Could result in a model with overall zero error but poor individual predictions

n - # examples
Better!
 $L = \frac{1}{N} \sum_i e_i^2$

Next Topic: Image Captioning Models



COCO Data Set

MS-COCO



COCO Data Set

COCO

	
<p>Woman on a horse <u>jumping</u> over a pole jump.</p>	<p>A glass bowl contains peeled tangerines and cut strawberries.</p>

COCO



Why Image Captioning?

① Virtual Assistants

Why Image Captioning?

- ① Virtual Assistants
- ② Visually impaired assistance

Why Image Captioning?

- ① Virtual Assistants
- ② Visually impaired assistance
- ③ Robotics

Why Image Captioning?

- ① Virtual Assistants
- ② Visually impaired assistance
- ③ Robotics
- ④ Any other use case?

CUB-200 Data Set

CUB-200



This bird is blue with white on its chest and has a very short beak.

}

1) What's the right level of detail to describe an image?

2) Style(Net)

↳ style in which caption generator
- Humorous . . .
- Serious

CUB-200 Data Set

CUB-200



This bird is blue with white on its chest and has a very short beak.

CUB-200



Fashion Data Set

(Product)

Fashion Captioning



A decorative leather padlock on a compact bag with croc embossed leather.

Fashion Data Set

Fashion Captioning



Fashion Captioning



Text Caps Data Set

TextCaps



*The billboard displays
'Welcome to Yakima The
Palm Springs of Washington'.*

Text Caps Data Set

TextCaps



*The billboard displays
'Welcome to Yakima The
Palm Springs of Washington'.*

TextCaps



Encoder-Decoder Model for Image Captioning

Encoder-Decoder Model for Image Captioning



Image

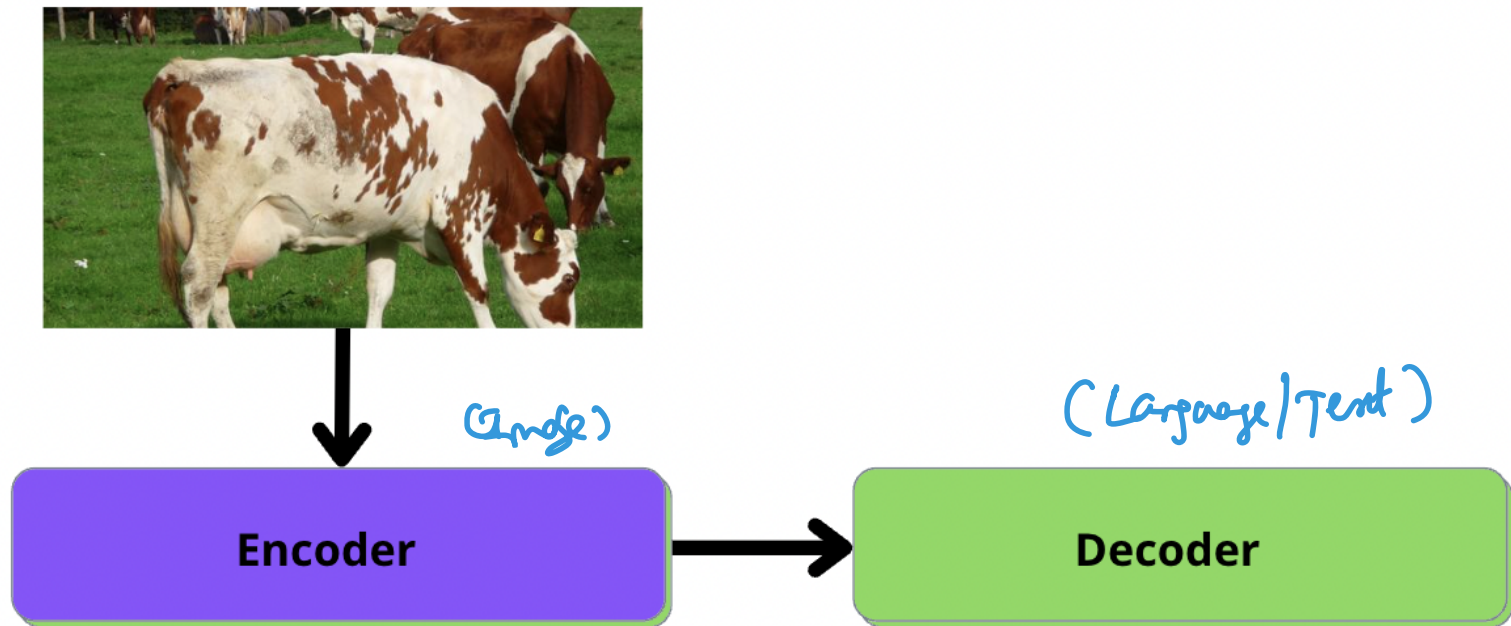
Encoder-Decoder Model for Image Captioning



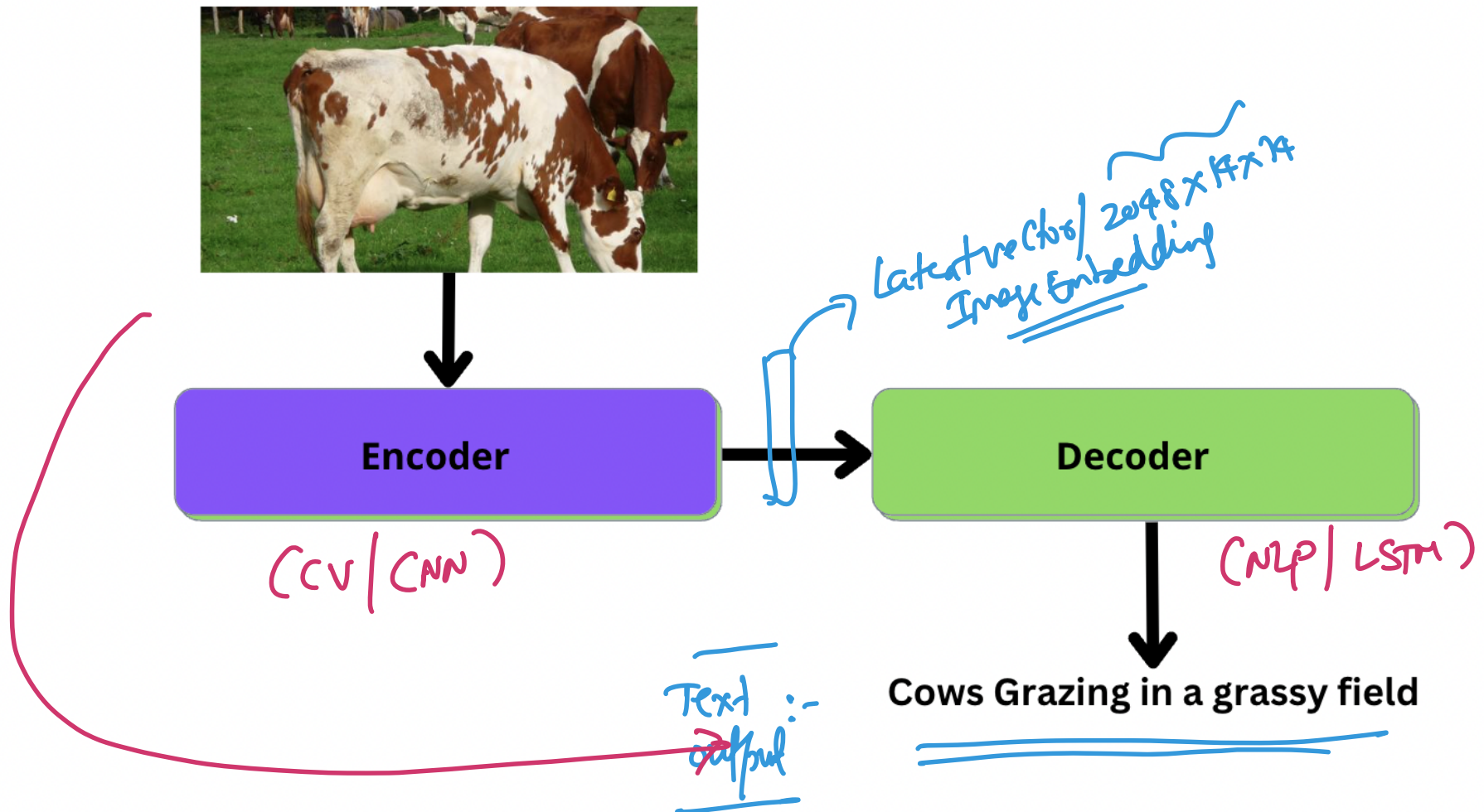
(Image Encoding)

Encoder

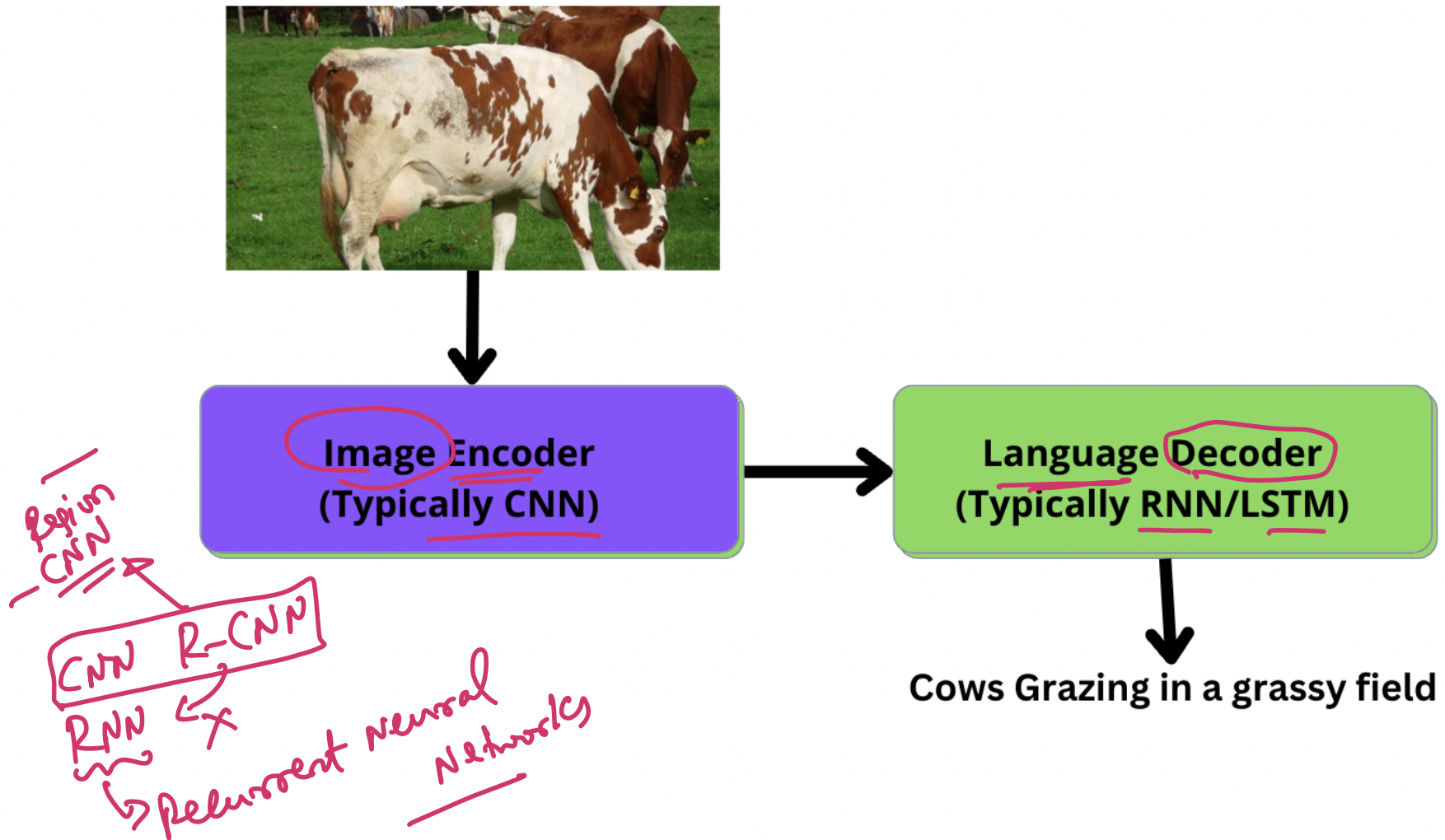
Encoder-Decoder Model for Image Captioning



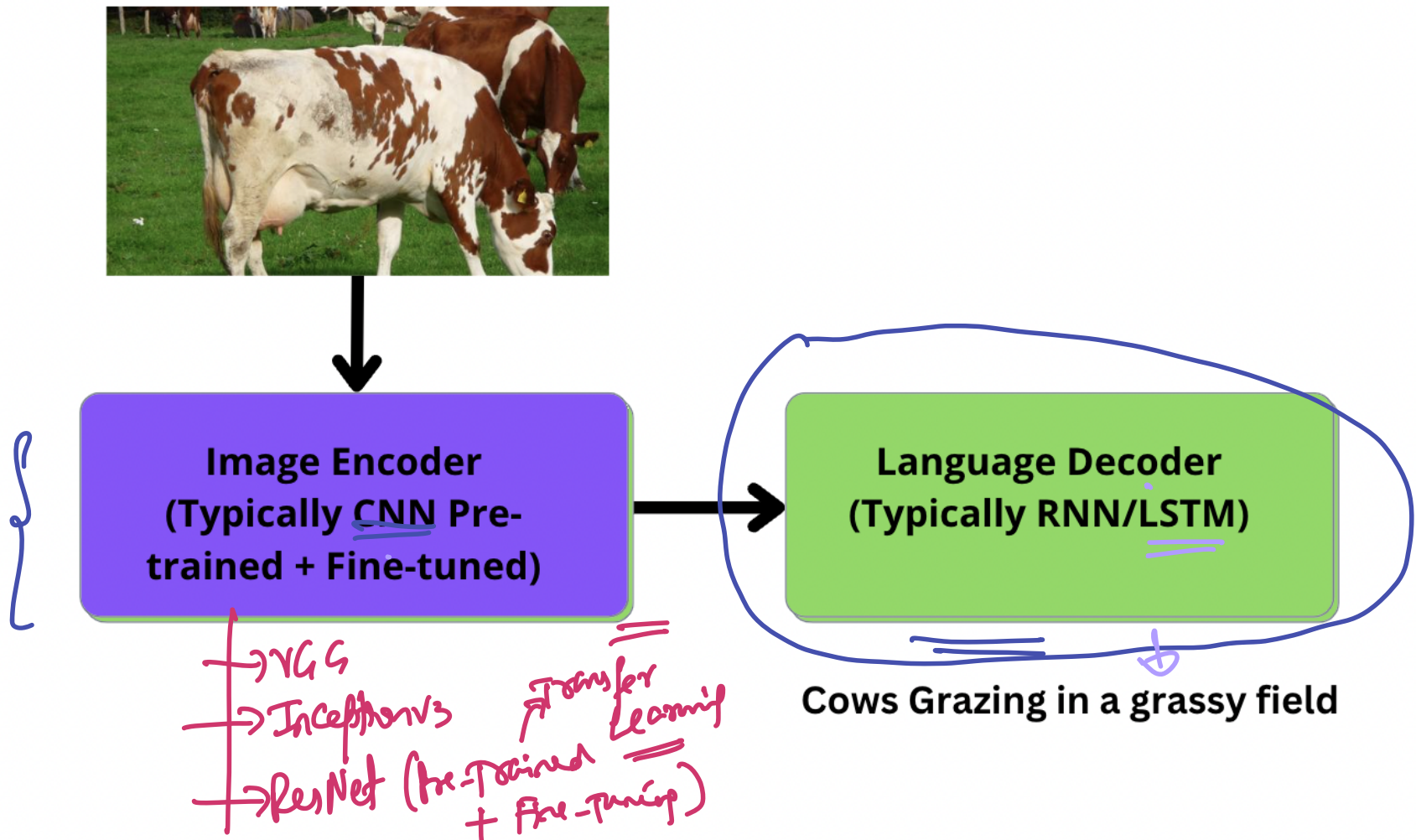
Encoder-Decoder Model for Image Captioning



Encoder-Decoder Model for Image Captioning



Encoder-Decoder Model for Image Captioning



Deep Learning Zoo

Neural Networks Zoo

- Input Cell
- Backfed Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Capsule Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Gated Memory Cell
- Kernel
- Convolution or Pool

A mostly complete chart of Neural Networks

©2019 Fjodor van Veen & Stefan Lejnen asimovinstitute.org

(Logistic Regression)

Perceptron (P)

Feed Forward (FF)

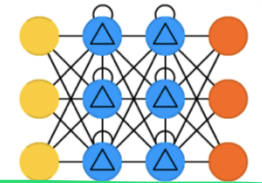
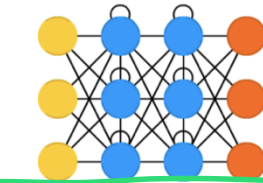
Radial Basis Network (RBF)

Deep Feed Forward (DFF)

Recurrent Neural Network (RNN)

Long / Short Term Memory (LSTM)

Gated Recurrent Unit (GRU)



(LSTM/NN)

→ NLP

Auto Encoder (AE)

Variational AE (VAE)

Denosing AE (DAE)

Sparse AE (SAE)



Winker '23 Course

Markov Chain (MC)

Hopfield Network (HN)

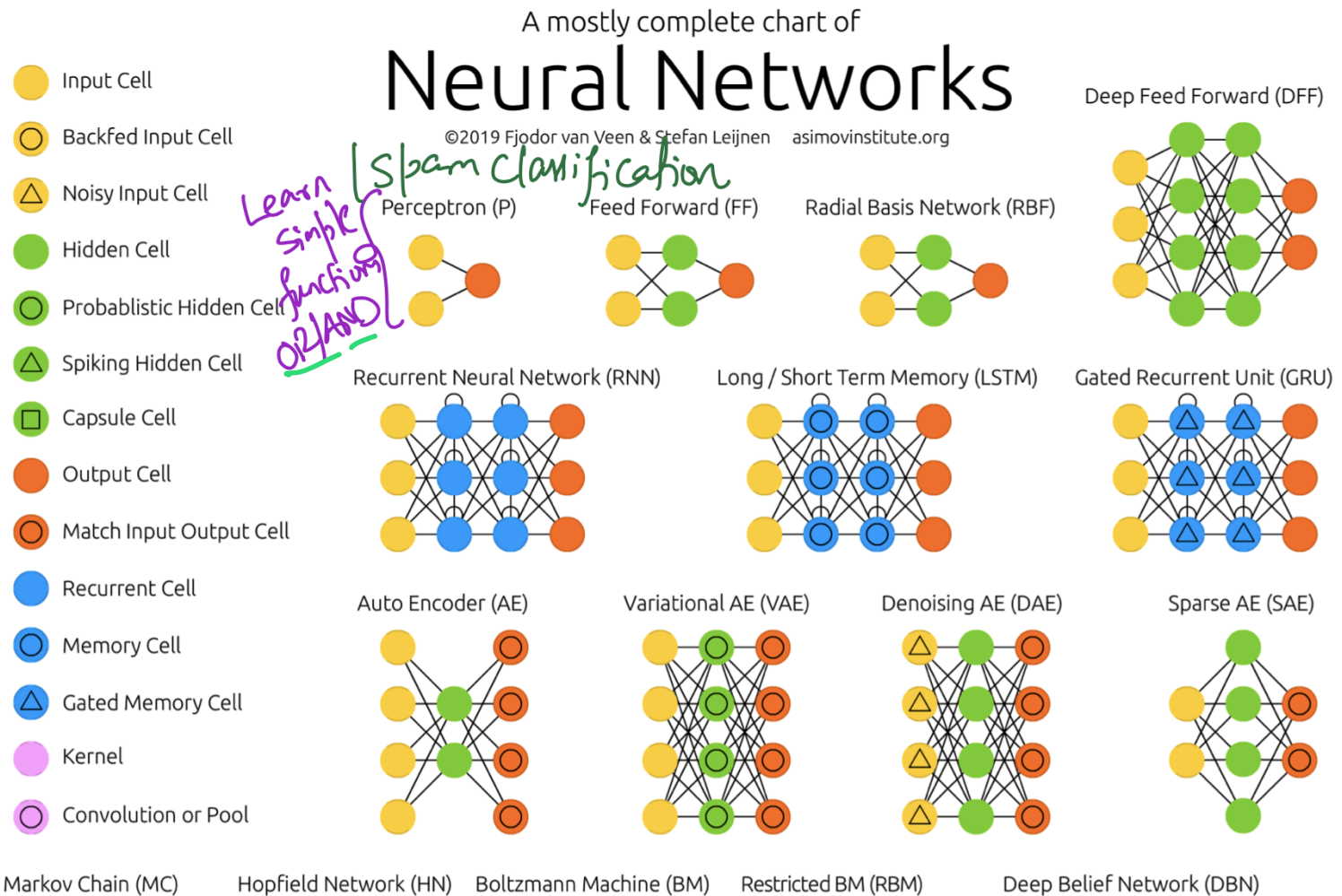
Boltzmann Machine (BM)

Restricted BM (RBM)

Deep Belief Network (DBN)














Deep Learning Zoo

Neural Networks Zoo



Deep Learning Zoo

Neural Networks Zoo

-  Input Cell
-  Backfed Input Cell
-  Noisy Input Cell
-  Hidden Cell
-  Probabilistic Hidden Cell
-  Spiking Hidden Cell
-  Capsule Cell
-  Output Cell
-  Match Input Output Cell
-  Recurrent Cell
-  Memory Cell
-  Gated Memory Cell
-  Kernel
-  Convolution or Pool

A mostly complete chart of Neural Networks

©2019 Fjodor van Veen & Stefan Leijnen asimovinstitute.org

Perceptron (P)



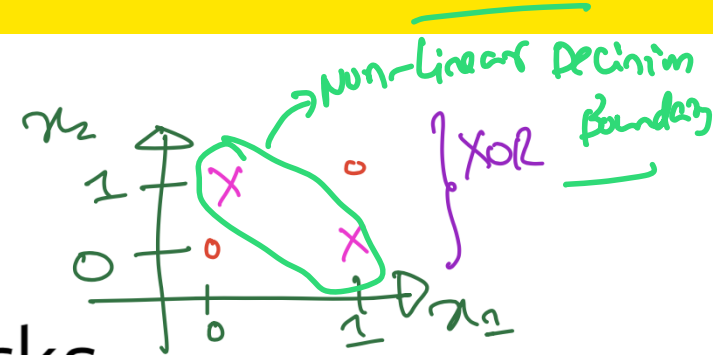
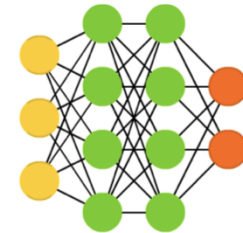
Learn XOR / Risk Assessment



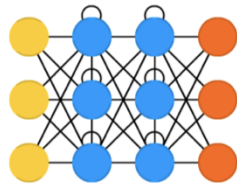
Radial Basis Network (RBF)



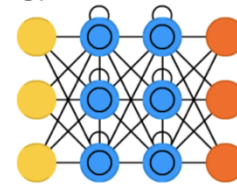
Deep Feed Forward (DFF)



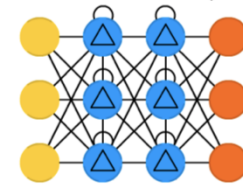
Recurrent Neural Network (RNN)



Long / Short Term Memory (LSTM)



Gated Recurrent Unit (GRU)



Auto Encoder (AE)



Variational AE (VAE)



Denosing AE (DAE)



Sparse AE (SAE)



Markov Chain (MC)

Hopfield Network (HN)







Boltzmann Machine (BM)

Restricted BM (RBM)

Deep Belief Network (DBN)

Deep Learning Zoo

Neural Networks Zoo

-  Input Cell
-  Backfed Input Cell
-  Noisy Input Cell
-  Hidden Cell
-  Probablistic Hidden Cell
-  Spiking Hidden Cell
-  Capsule Cell
-  Output Cell
-  Match Input Output Cell
-  Recurrent Cell
-  Memory Cell
-  Gated Memory Cell
-  Kernel
-  Convolution or Pool

A mostly complete chart of Neural Networks

©2019 Fjodor van Veen & Stefan Leijnen asimovinstitute.org

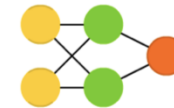
Perceptron (P)



Feed Forward (FF)

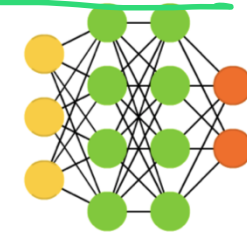


Radial Basis Network (RBF)

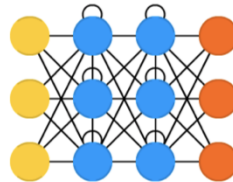


Recommender System / Classification

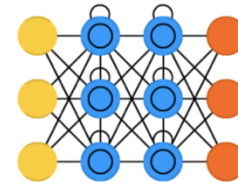
Deep Feed Forward (DFF)



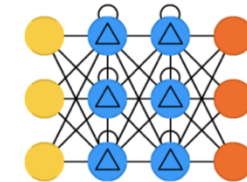
Recurrent Neural Network (RNN)



Long / Short Term Memory (LSTM)



Gated Recurrent Unit (GRU)



Auto Encoder (AE)



Variational AE (VAE)



Denosing AE (DAE)



Sparse AE (SAE)



Markov Chain (MC)

Hopfield Network (HN)

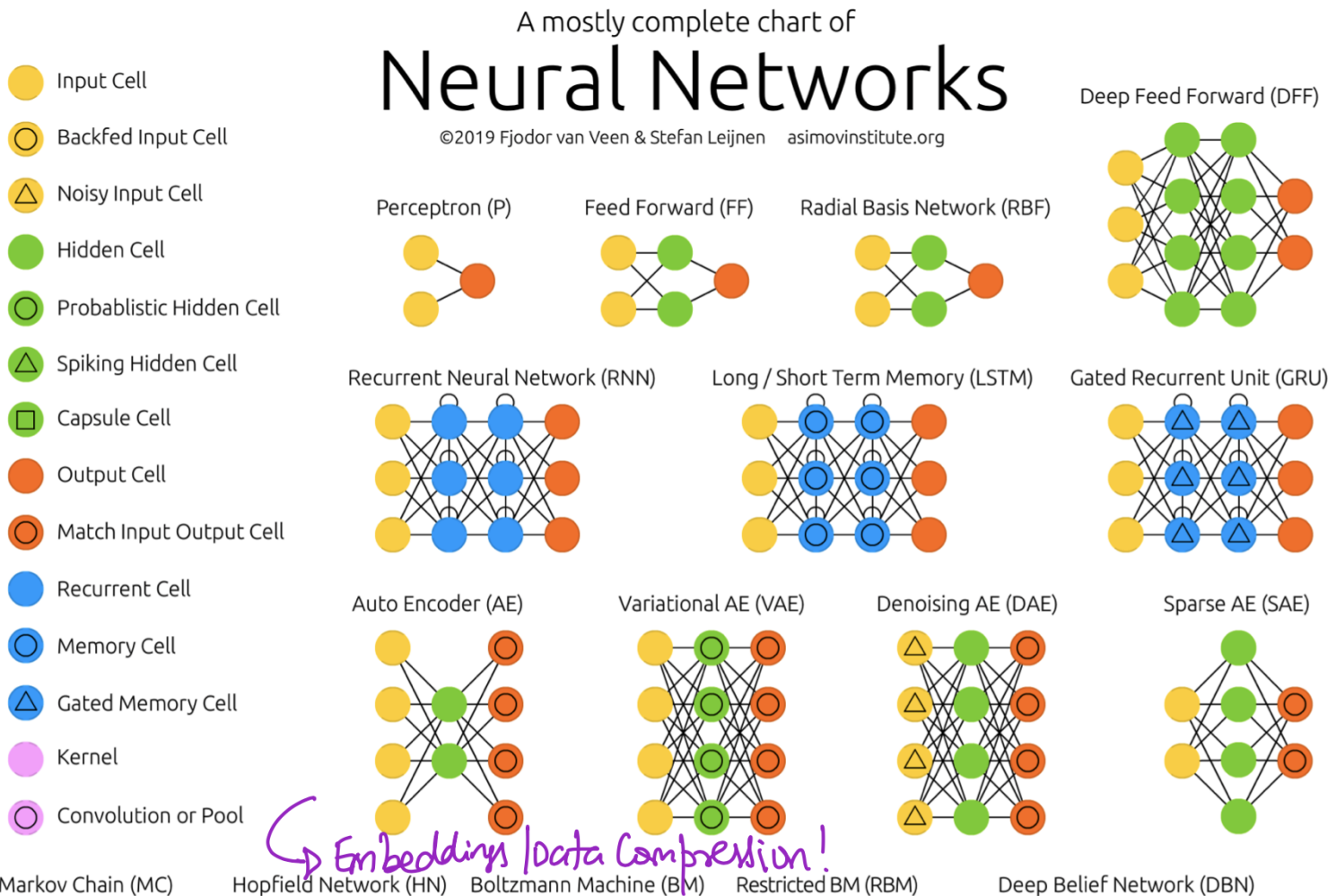
Boltzmann Machine (BM)

Restricted BM (RBM)

Deep Belief Network (DBN)

Deep Learning Zoo

Neural Networks Zoo



Deep Learning Zoo

Neural Networks Zoo

- Input Cell
- Backfed Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Capsule Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Gated Memory Cell
- Kernel
- Convolution or Pool

A mostly complete chart of Neural Networks

©2019 Fjodor van Veen & Stefan Leijnen asimovinstitute.org

Perceptron (P)



Feed Forward (FF)



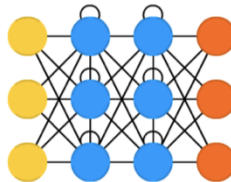
Radial Basis Network (RBF)



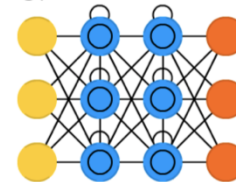
Deep Feed Forward (DFF)



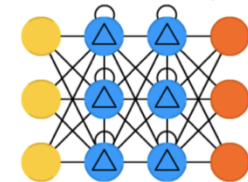
Recurrent Neural Network (RNN)



Long / Short Term Memory (LSTM)



Gated Recurrent Unit (GRU)



Auto Encoder (AE)



Variational AE (VAE)



Denoising AE (DAE)



Sparse AE (SAE)



Markov Chain (MC)

Hopfield Network (HN)

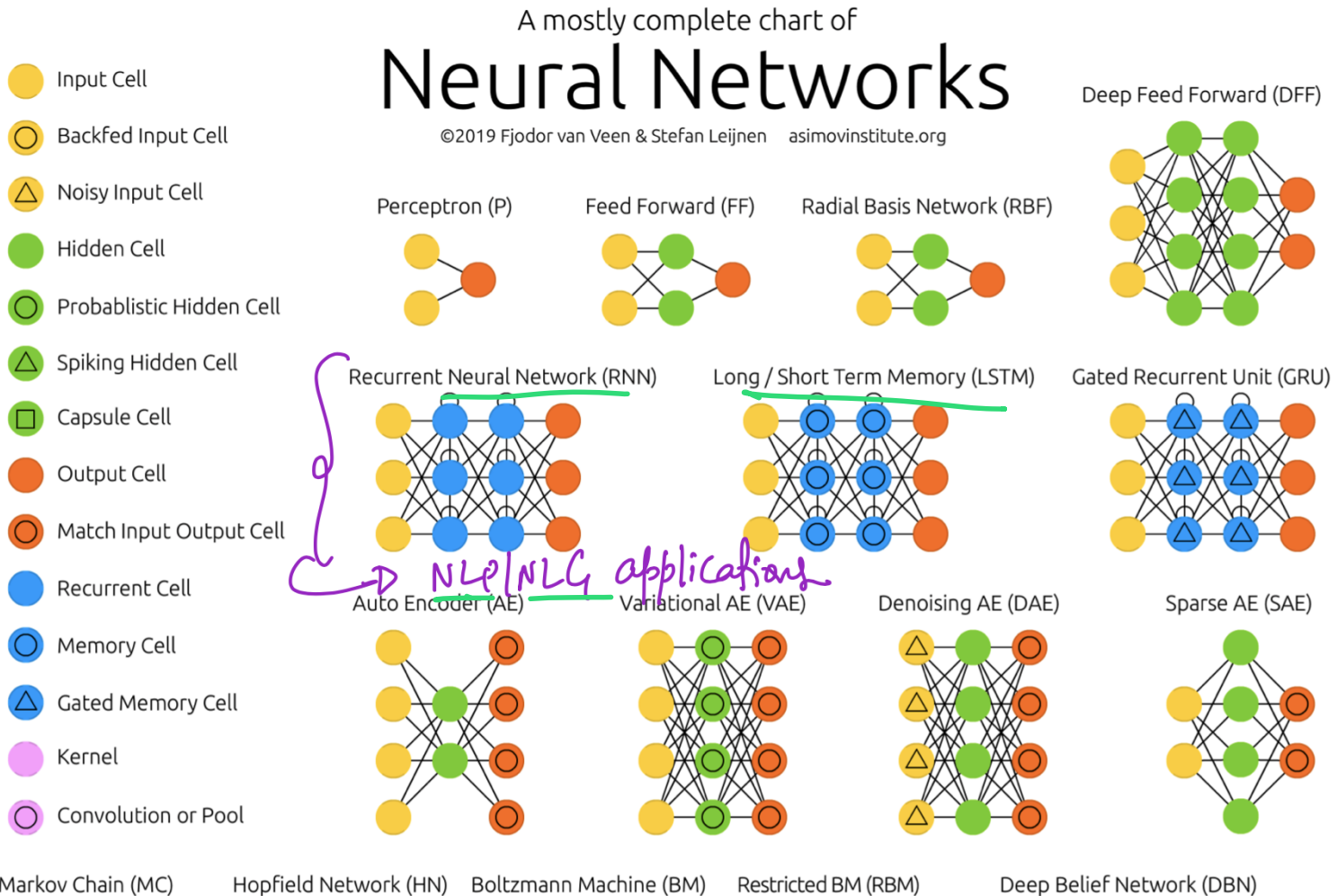
Boltzmann Machine (BM)

Restricted BM (RBM)

Deep Belief Network (DBN)





Deep Learning Zoo

Neural Networks Zoo



Deep Learning Zoo

Neural Networks Zoo

-  Input Cell
-  Backfed Input Cell
-  Noisy Input Cell
-  Hidden Cell
-  Probabilistic Hidden Cell
-  Spiking Hidden Cell
-  Capsule Cell
-  Output Cell
-  Match Input Output Cell
-  Recurrent Cell
-  Memory Cell
-  Gated Memory Cell
-  Kernel
-  Convolution or Pool

TODAY! I & II

A mostly complete chart of Neural Networks

©2019 Fjodor van Veen & Stefan Leijnen asimovinstitute.org

Perceptron (P)



Feed Forward (FF)



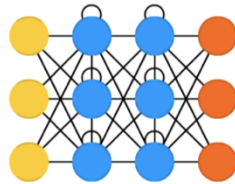
Radial Basis Network (RBF)



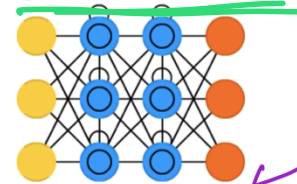
Deep Feed Forward (DFF)



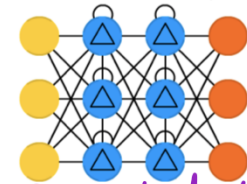
Recurrent Neural Network (RNN)



Long / Short Term Memory (LSTM)



Gated Recurrent Unit (GRU)



Auto Encoder (AE)



Variational AE (VAE)



Denosing AE (DAE)



Sparse AE (SAE)



Markov Chain (MC)

Hopfield Network (HN)

Boltzmann Machine (BM)

Restricted BM (RBM)

Deep Belief Network (DBN)

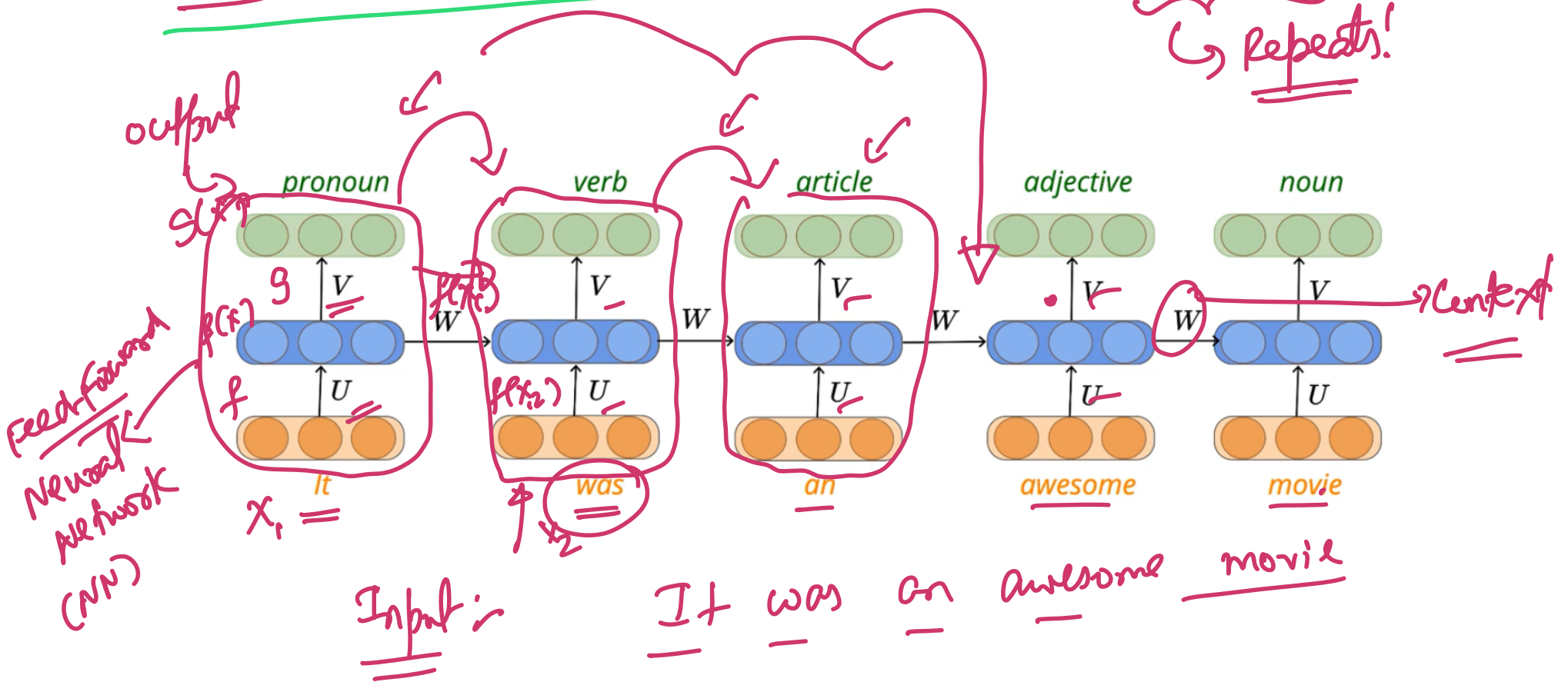
I. Recommender Systems Case-Study

II. NLG Case Study!

LSTM

↳ Long Short Term Memory Network

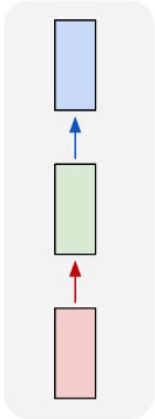
RNN: Re Current Neural network
↳ Repeats!



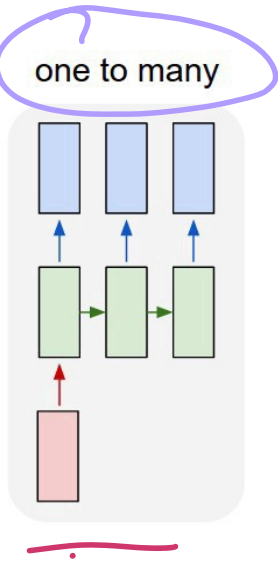
LSTM

input output

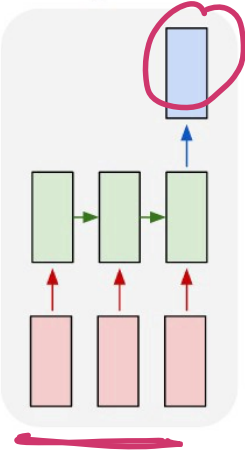
one to one



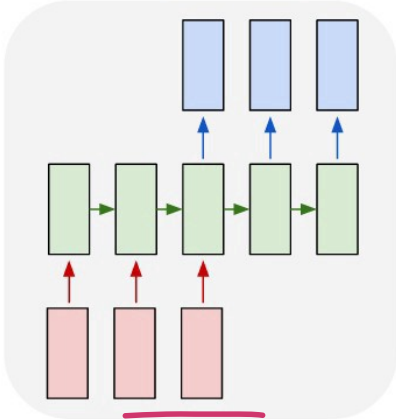
one to many



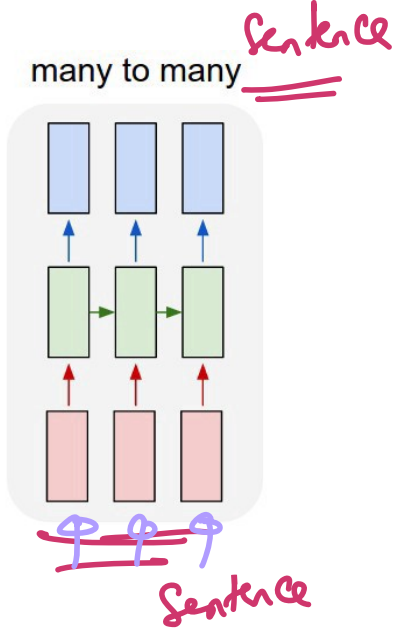
many to one



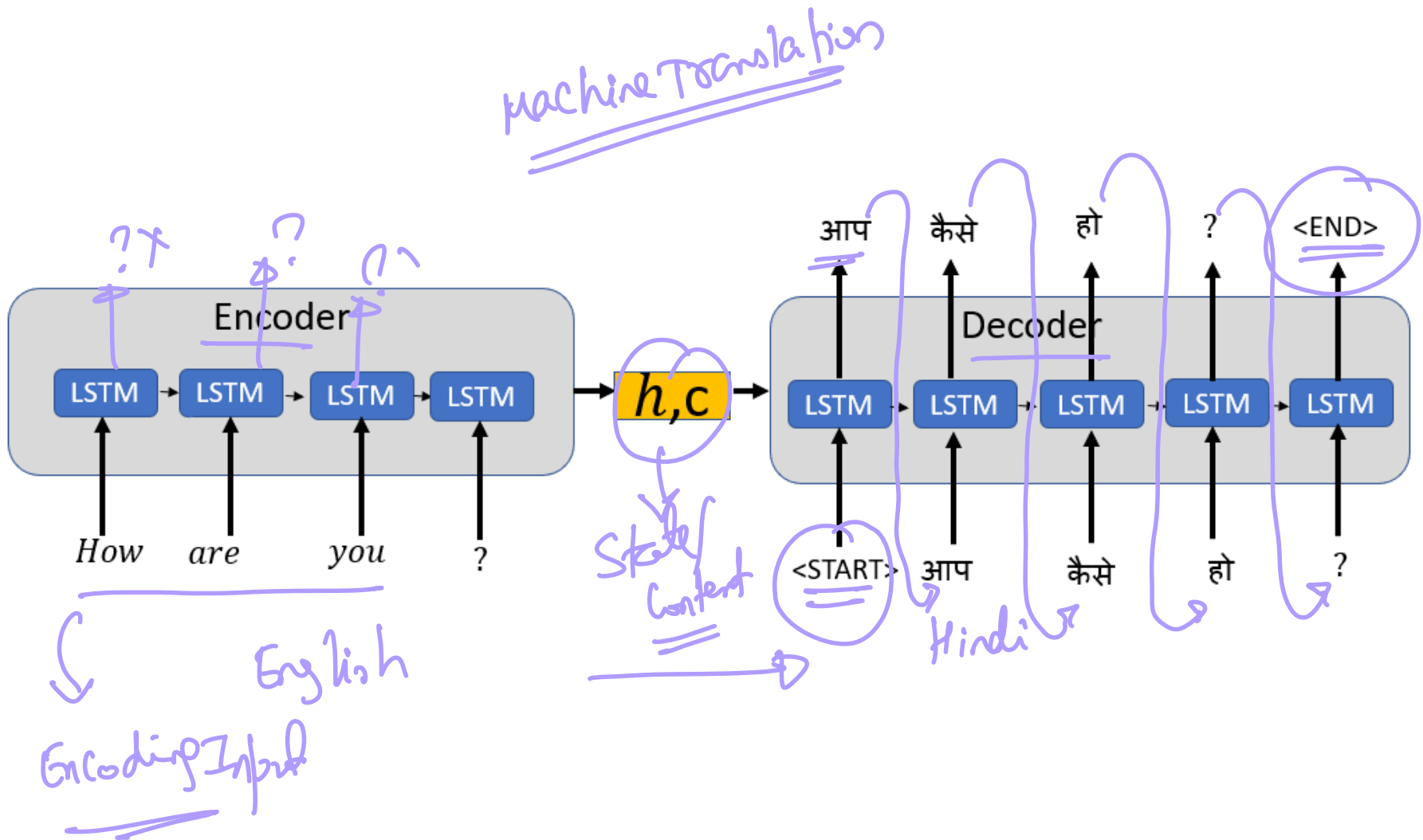
many to many



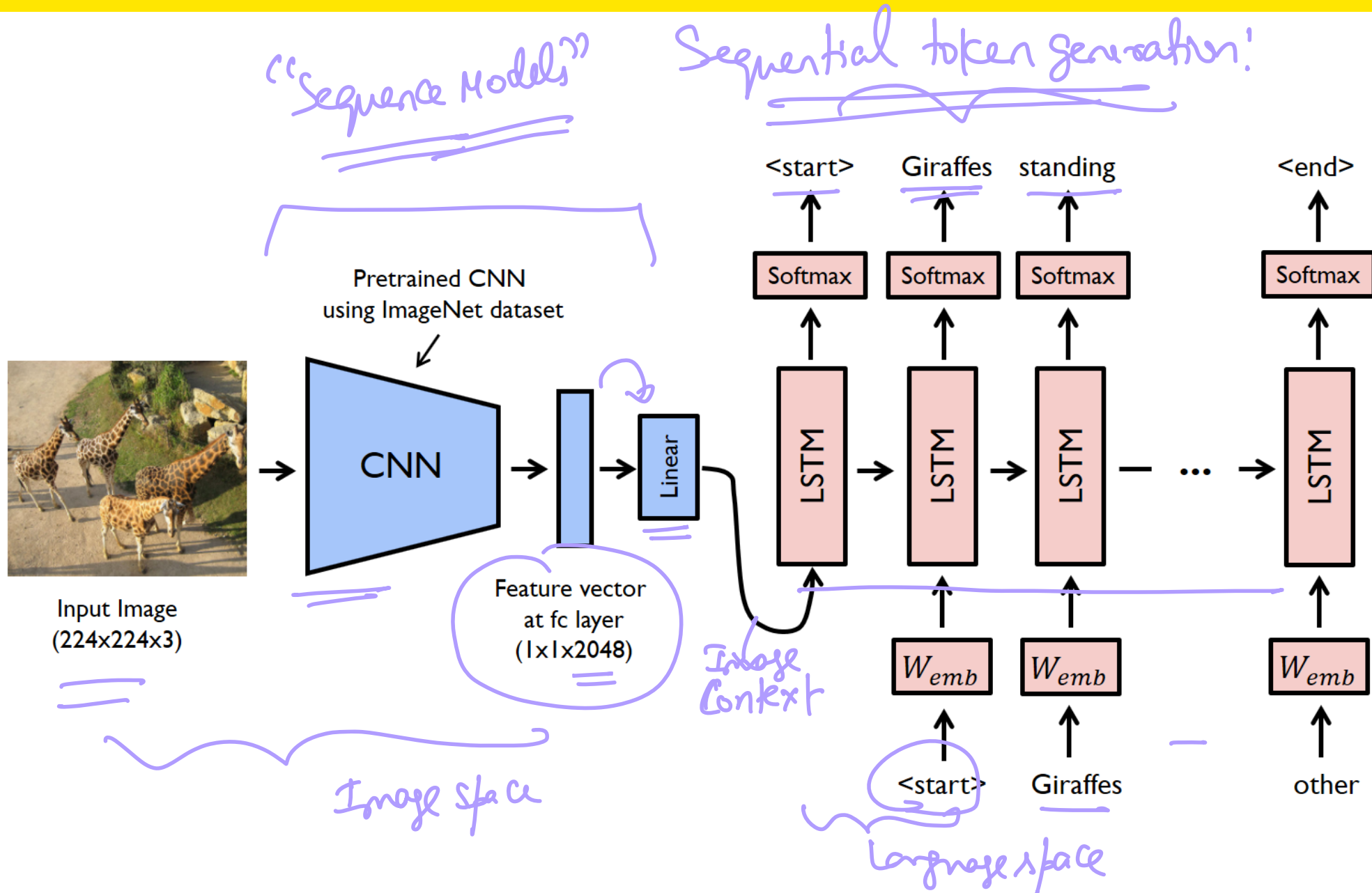
many to many



LSTM



LSTM



Metrics for Image Captioning

- ① BLEU (Bilingual evaluation understudy)

Metrics for Image Captioning

- 1 BLEU (Bilingual evaluation understudy)
- 2 **METEOR** (Metric for Evaluation of Translation with Explicit Ordering)

Improvement on BLEU score!

Example:-

Poed:-

Cows grazing on grass

Tenth:-

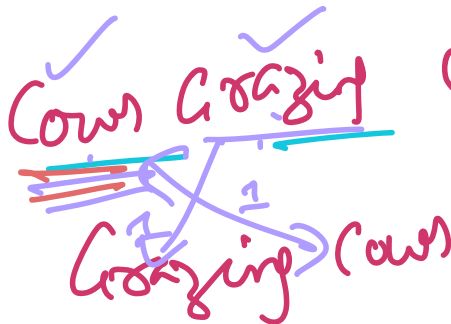
Grazing cows on grassy field

unigrams

METEOR Metric

Pred: -

Truth: -



on Grass

on grassy field $c=3$
 $\mu=4$

1) Compute Precision - (P) $3/4$

2) Compute Recall - (R) $3/5$

3) F-Score :- $\frac{10PR}{R+9P}$
 ↳ Biased towards Recall!

4) Fragmentation -

$$\text{frag} = \frac{c}{\mu}$$

↳ # chunks

tokens ident

5) Penalty
 $p = 0.5 \times (\text{frag})^2$

6) METEOR
 $\text{f-score} (1-p)$

METEOR Example

Taken from webpage on METEOR

Round →

	audatex	invests	90	million	euros	a	year	in	developing	these	databases	.
audatex	•											
invested		o										
each												
year												
,												
90			•									
million				•								
to												
develop												
these												
databases												
.												

Matrix

Include year

Don't include year

6/9

6/9

Segment 2437

P: 0.807
 R: 0.764
 Frag: 0.490
 Score: 0.393

$$F = \frac{2 \times P \times R}{P + R}$$

$$F = \frac{2 \times 0.807 \times 0.764}{0.807 + 0.764} = 0.58$$

$$F_{\text{score}} \times (1 - P) = 0.58 \times 0.490 = 0.287$$

$$F_{\text{score}} = \frac{4}{7} = 0.58$$

$$\frac{3}{6} = 0.5$$

$$\frac{6}{9}$$

$$\frac{6}{9}$$

$$\frac{3}{6} = 0.5$$

$$\frac{4}{7} = 0.58$$

$$\frac{3}{6} = 0.5$$

$$F_{\text{score}} \times (1 - P) = 0.58 \times 0.490 = 0.287$$

ICE #2

$$\begin{aligned} p &= \frac{5}{6} \\ r &= \frac{5}{6} \\ \text{F-score} &= \frac{5}{6} \\ \text{frag} &= \frac{3}{5} = 0.6 \quad \checkmark \\ p &= 0.5 \times 0.6^3 = 0.108 \\ \text{METEOR} &= \frac{5}{6} \times (1 - p) = 0.74 \end{aligned}$$

Cows Grazing

Consider that your image captioning model generated the sentence:

“Found Grazing cows on the grass” and the true caption was “Grazing cows found over the grass”. What is the fragmentation value and METEOR score in this case?

pred → Truth

$$\begin{aligned} \text{METEOR} &= \text{F-Score} \times (1 - p) \quad \text{penalty} \\ &= \frac{1 \times \frac{5}{6} \times \frac{5}{6}}{\frac{5}{6} + 1 \times \frac{5}{6}} \times (1 - 0.108) \\ &= \frac{5}{6} \times (1 - 0.108) \\ &= \frac{5}{6} \times 0.892 \\ &= 0.74 \end{aligned}$$

$\text{frag} = \frac{C}{\mu} \rightarrow \# \text{Common words}$ (chunks)

Image Captioning Models

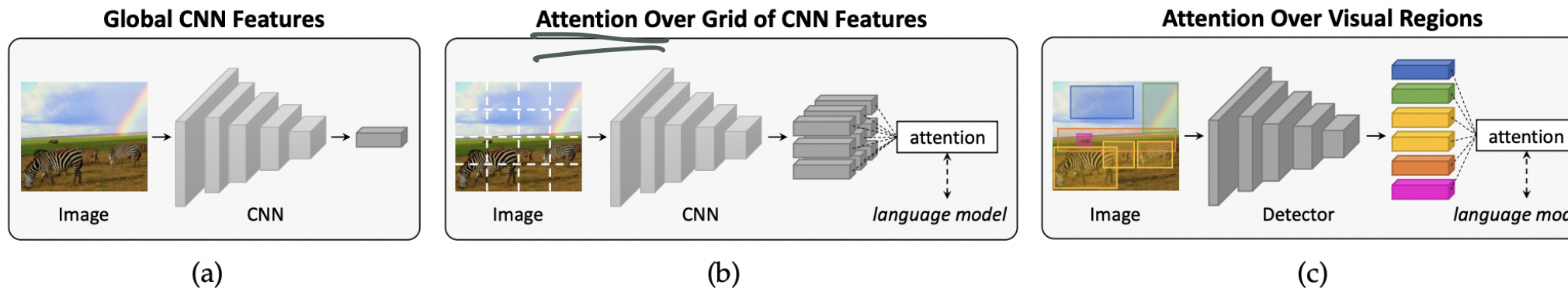


Fig. 2: Three of the most relevant visual encoding strategies for image captioning: **(a)** global CNN features; **(b)** fine-grained features extracted from the activation of a convolutional layer, together with an attention mechanism guided by the language model; **(c)** image region features coming from a detector, together with an attention mechanism.

Show & Tell

Breakout for Takeaways!

Discuss Takeaways (5 mins)

From today's lecture in your zoom group

Next Lecture

- ① More on Image Captioning Models
- ② Show and Tell Image Captioning Models