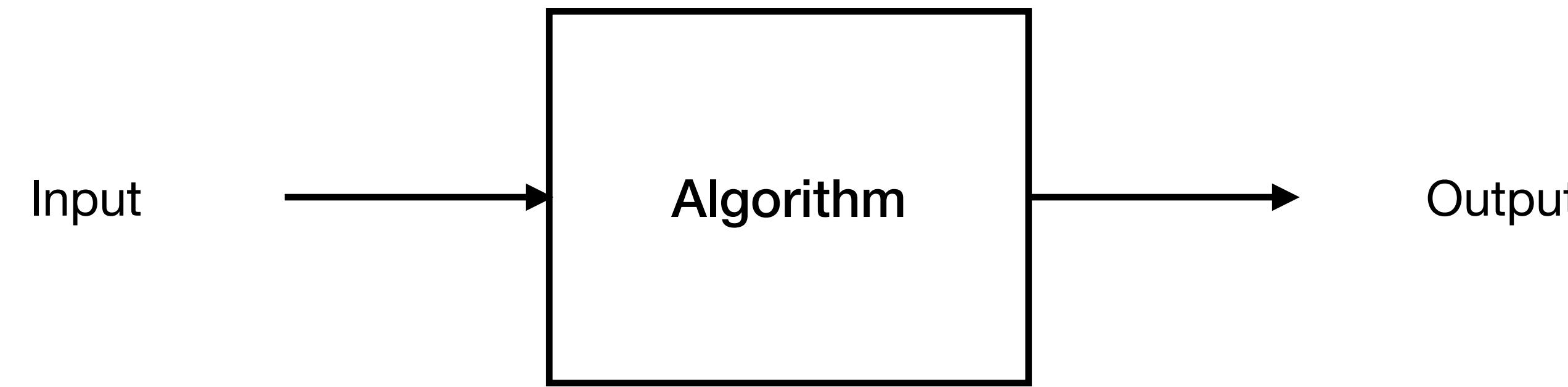


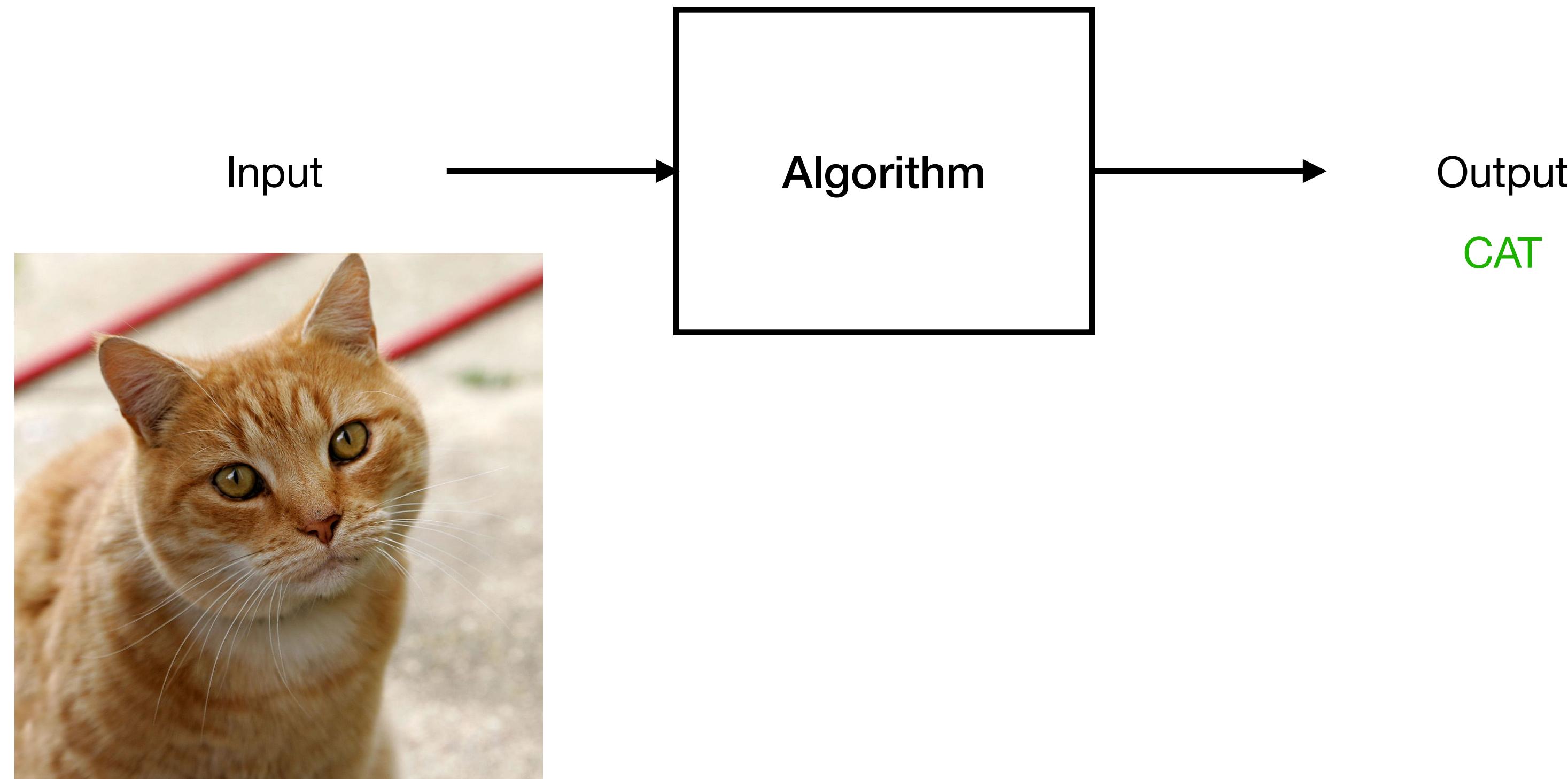
Sequence Prediction

One Step Prediction



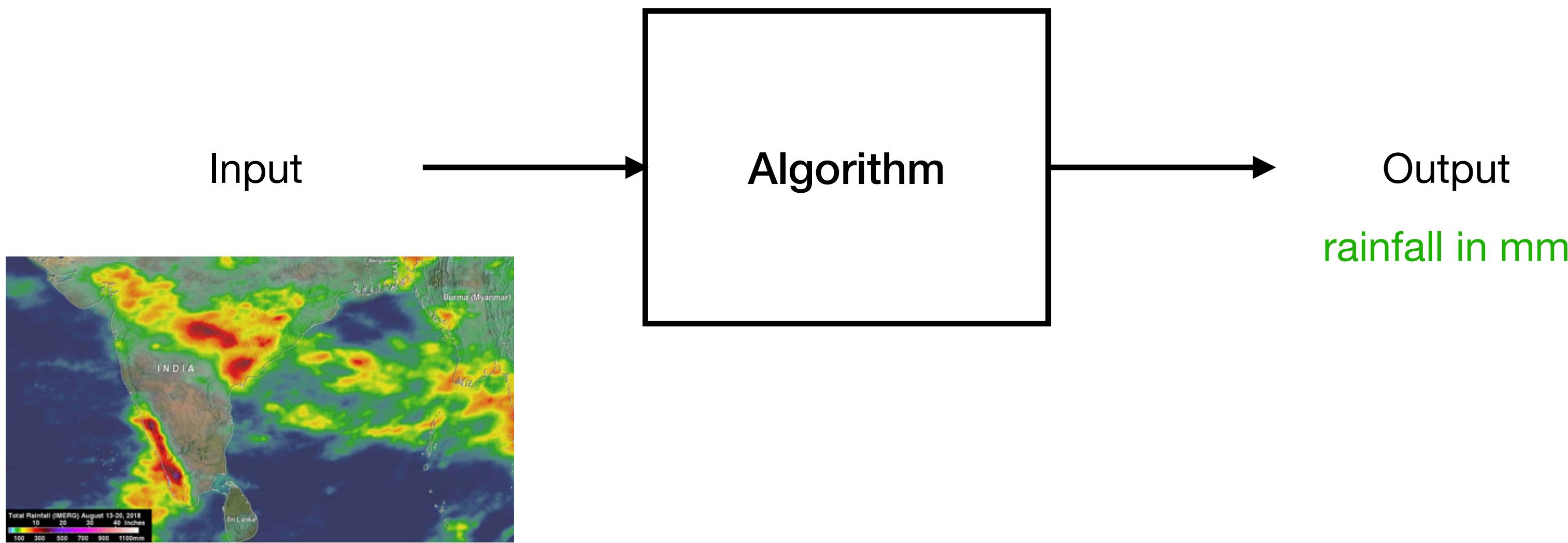
One Step Prediction

Image Classification



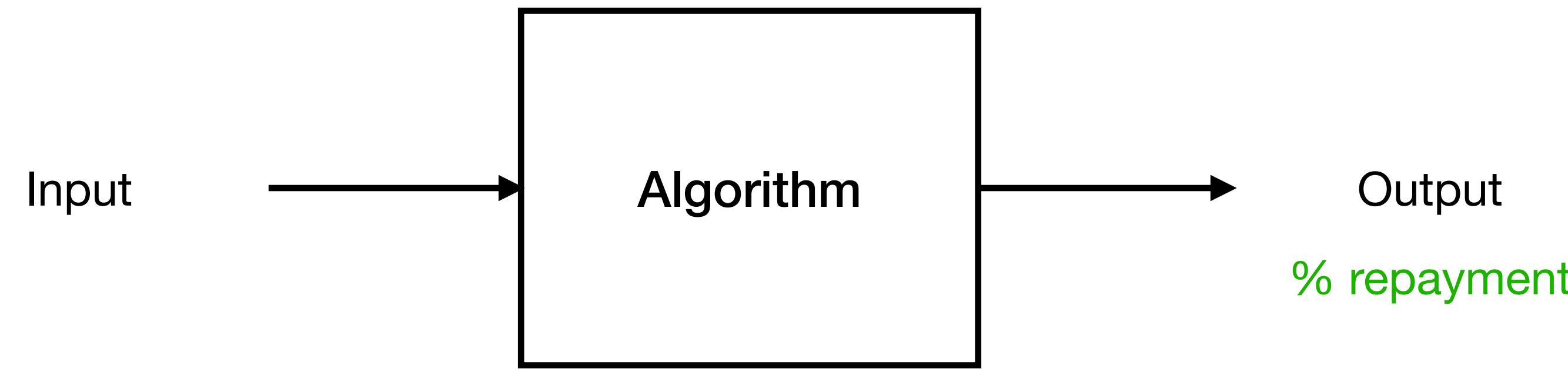
One Step Prediction

Rainfall Prediction



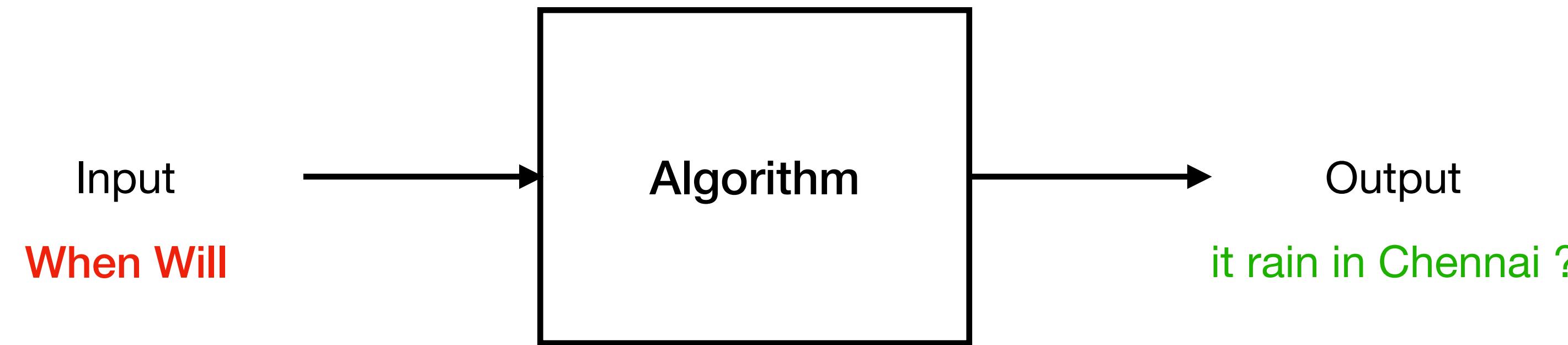
One Step Prediction

Loan Repayment



Sequence Prediction

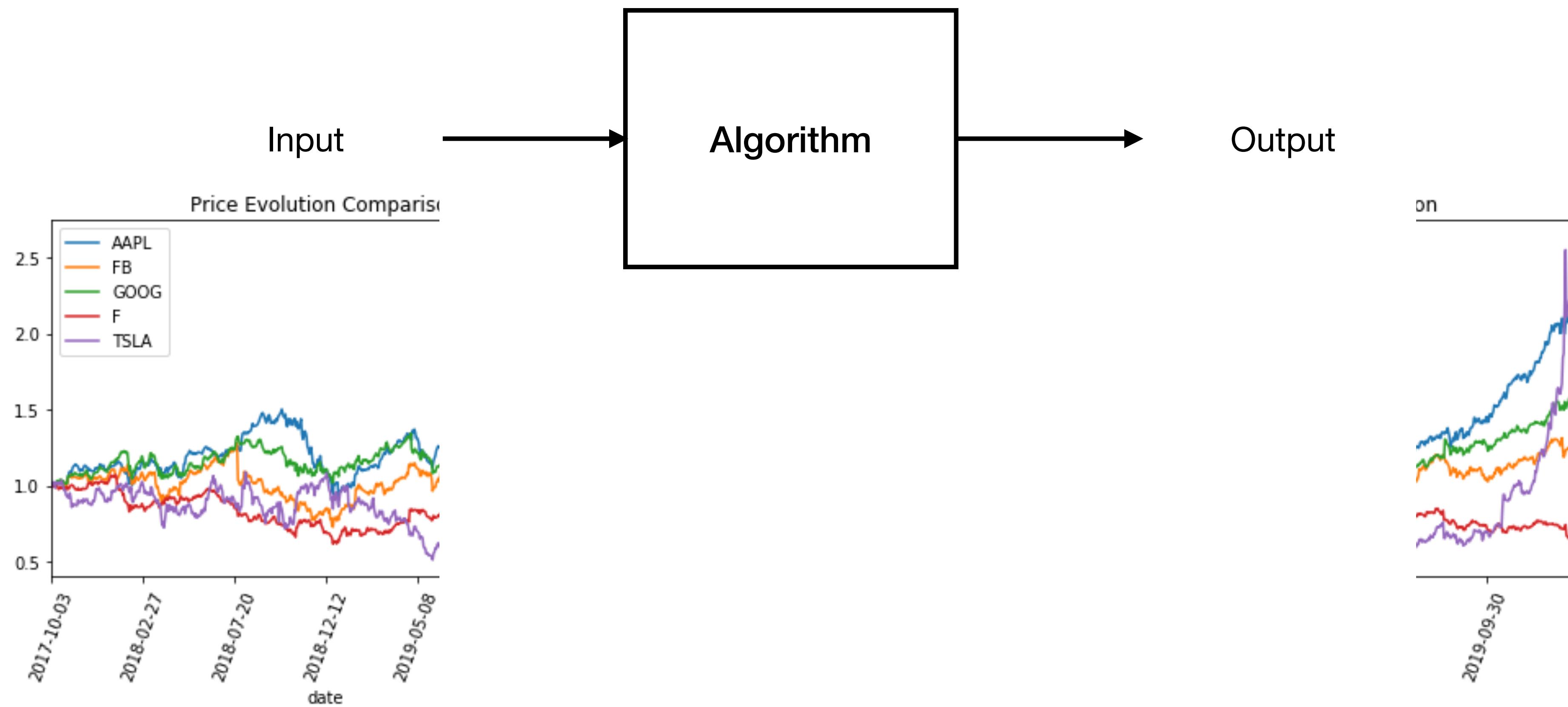
Sentence Completion



When will it rain in Chennai ?

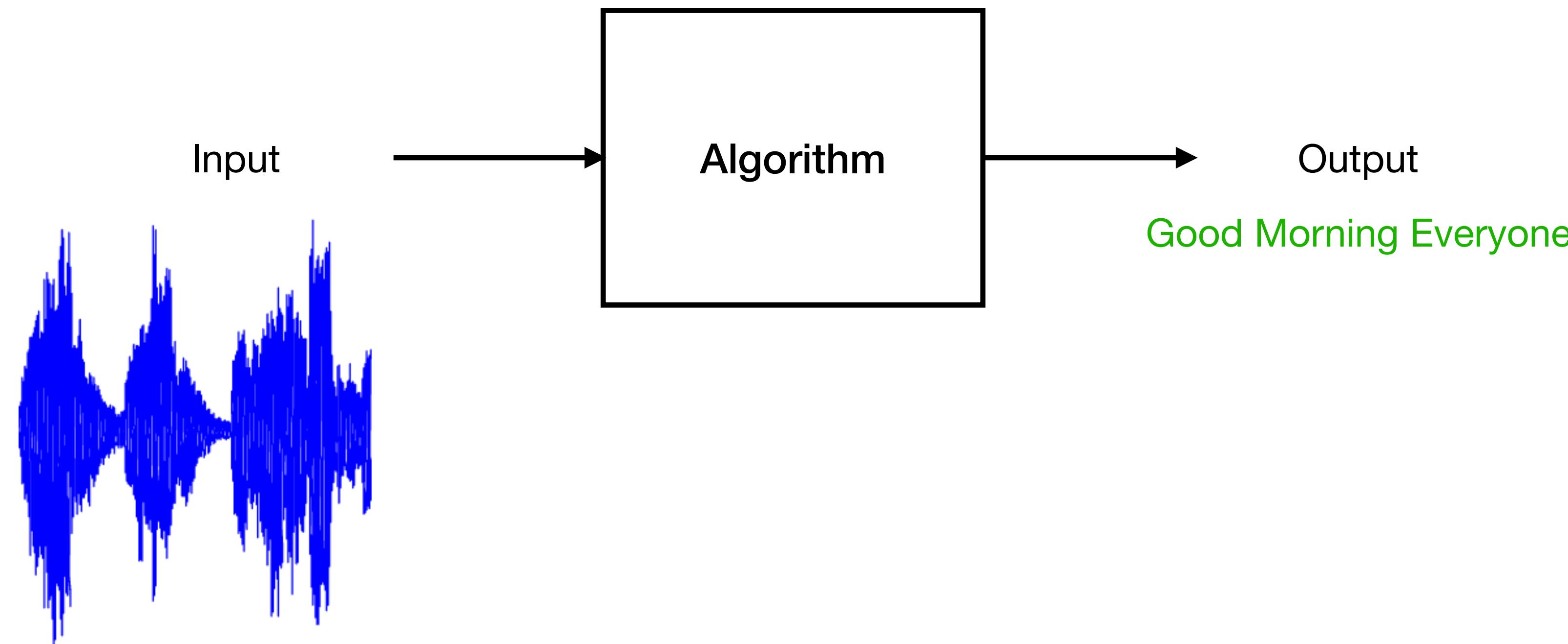
Sequence Prediction

Stock Prediction



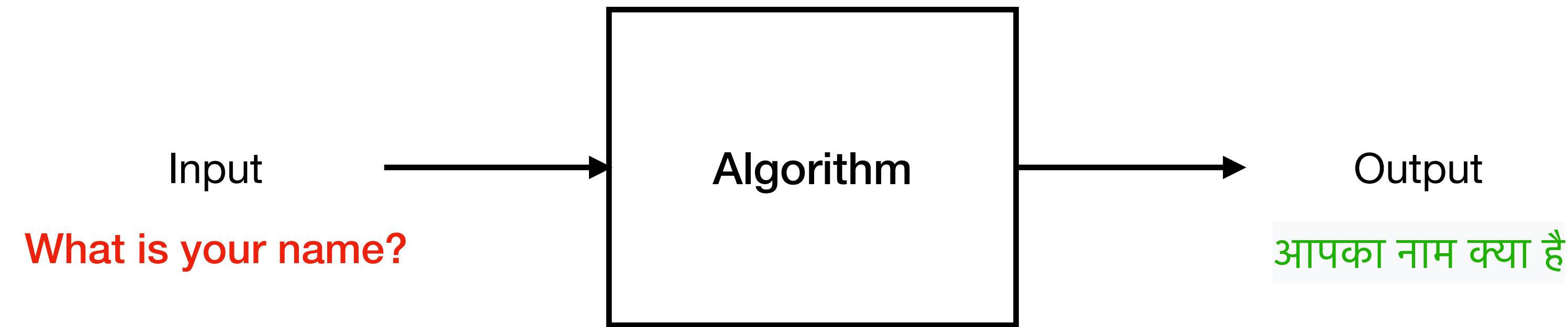
Sequence Prediction

Speech Recognition



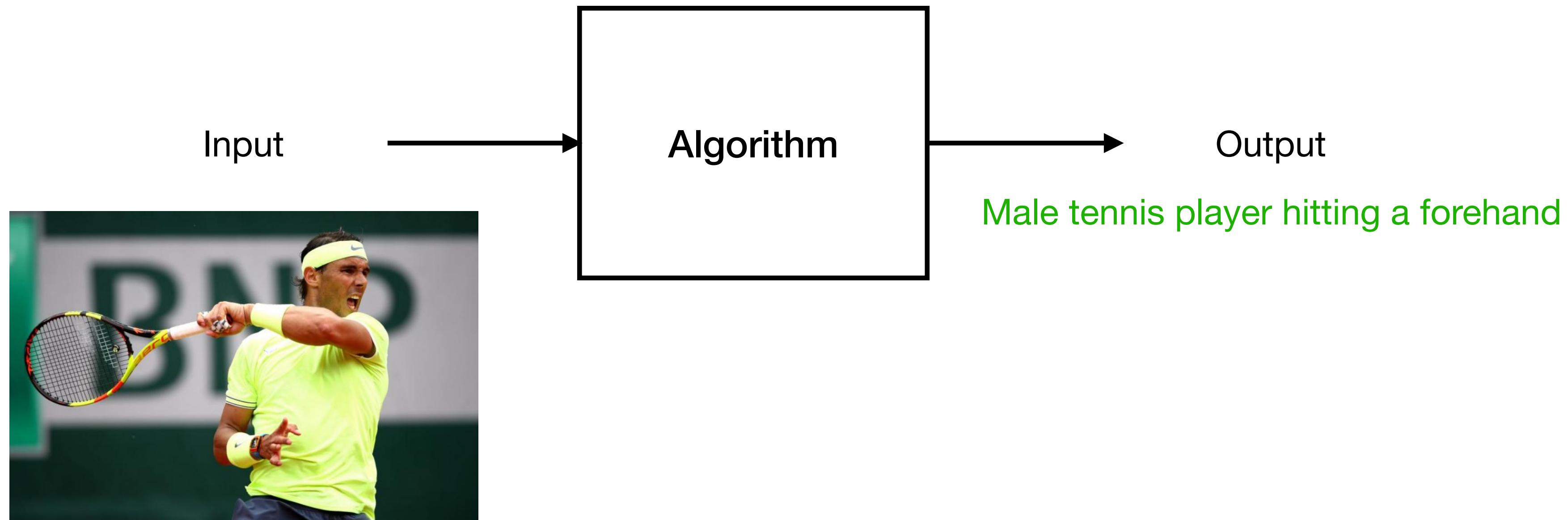
Sequence Prediction

Machine Translation



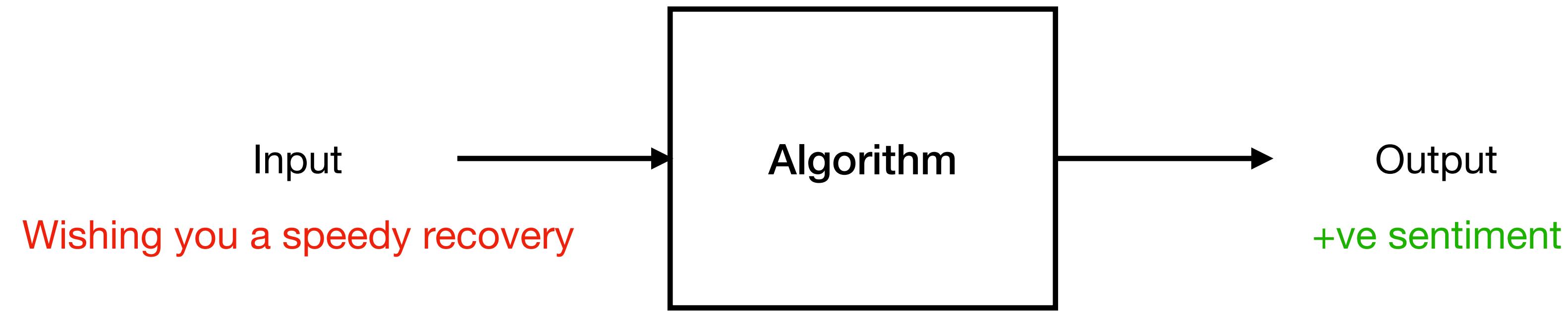
Sequence Prediction

Image Captioning



Sequence Prediction

Sentiment Analysis



- Different types of sequence prediction problems
- Place RNNs, LSTMs and Transformers amongst various
(these are architectures for sequence prediction)
other deep neural network architectures
- List down aspects of sequence prediction
(gives an idea of how to derive architecture for
RNN, LSTM and Transformers)
- Build RNNs from scratch
- Recurrences (linear)
- Backpropagation Through Time
- Issues with RNN
- LSTMs
- Transformers

One Shot Prediction

think about the algorithm that tags images

- 1) snapchat
- 2) Instagram
- 3) Facebook

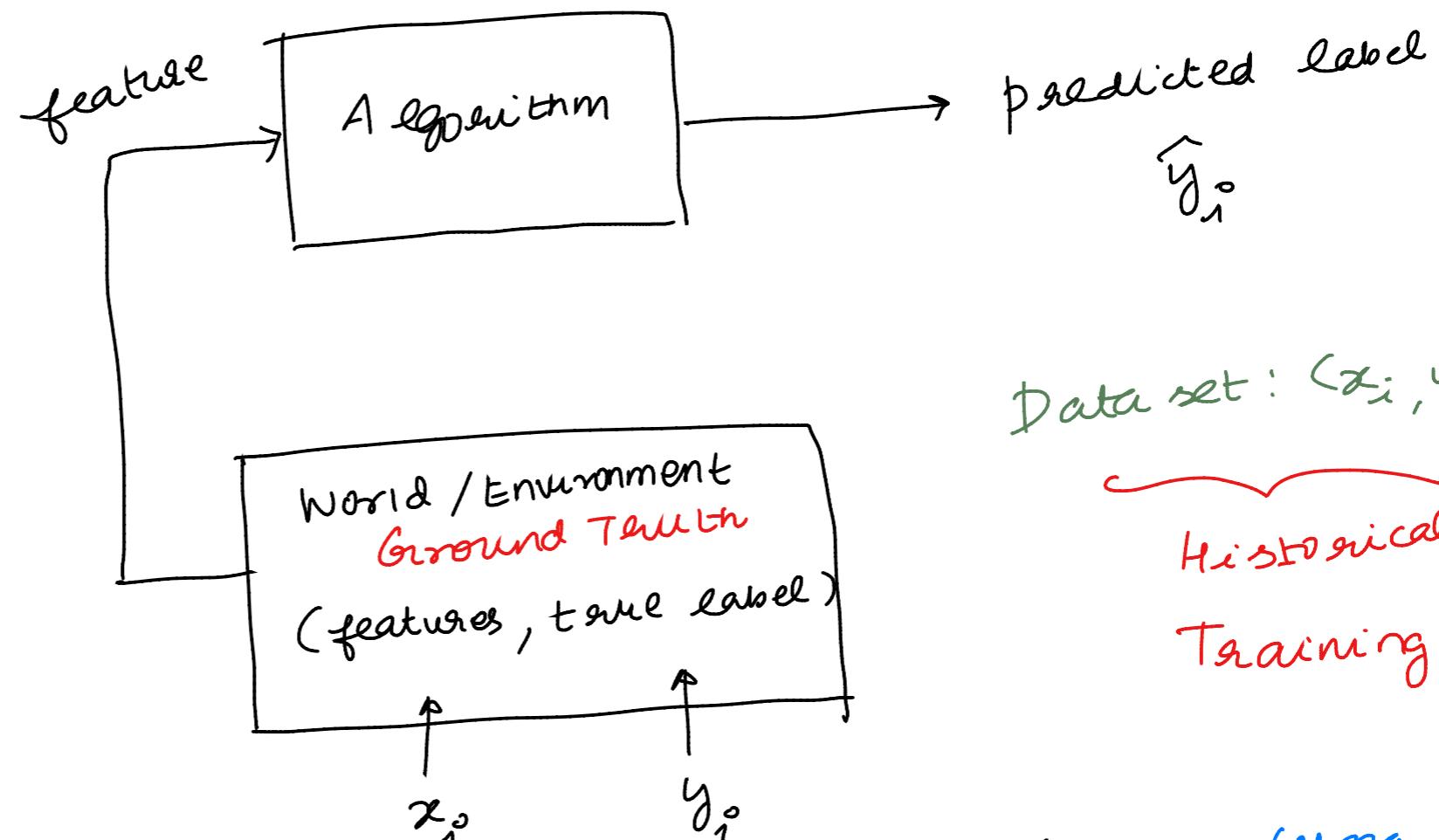
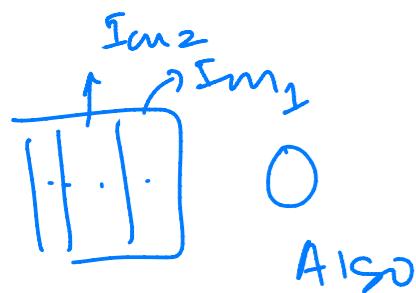


Image Recognition:

$$x_i \in \text{image} \subseteq \mathbb{R}^d$$

$y_i \in \mathcal{Y} = \{\text{cat, dog, human, house, ..., train, bus}\}$

Imagenet : 1000 categories classes.

Dataset: $(x_i, y_i)_{i=1}^n$

Historical Data

Training Data

(Mega Pixel Camera)

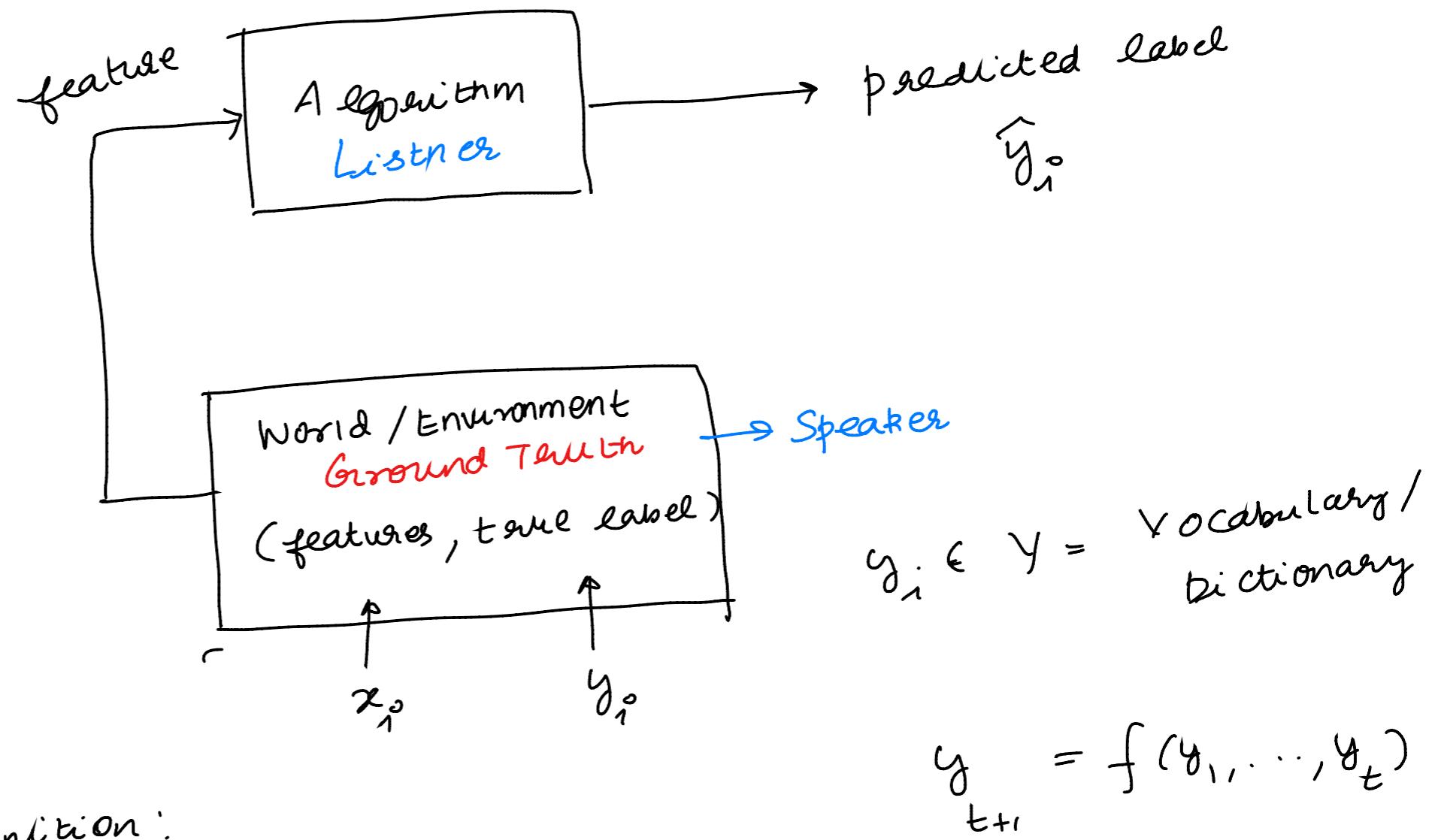
 $x_i \in \mathbb{R}^{3 \times 10^6}$
 $\mathbb{R}^{3 \times 1000 \times 1000}$

Sequence Prediction

Dataset

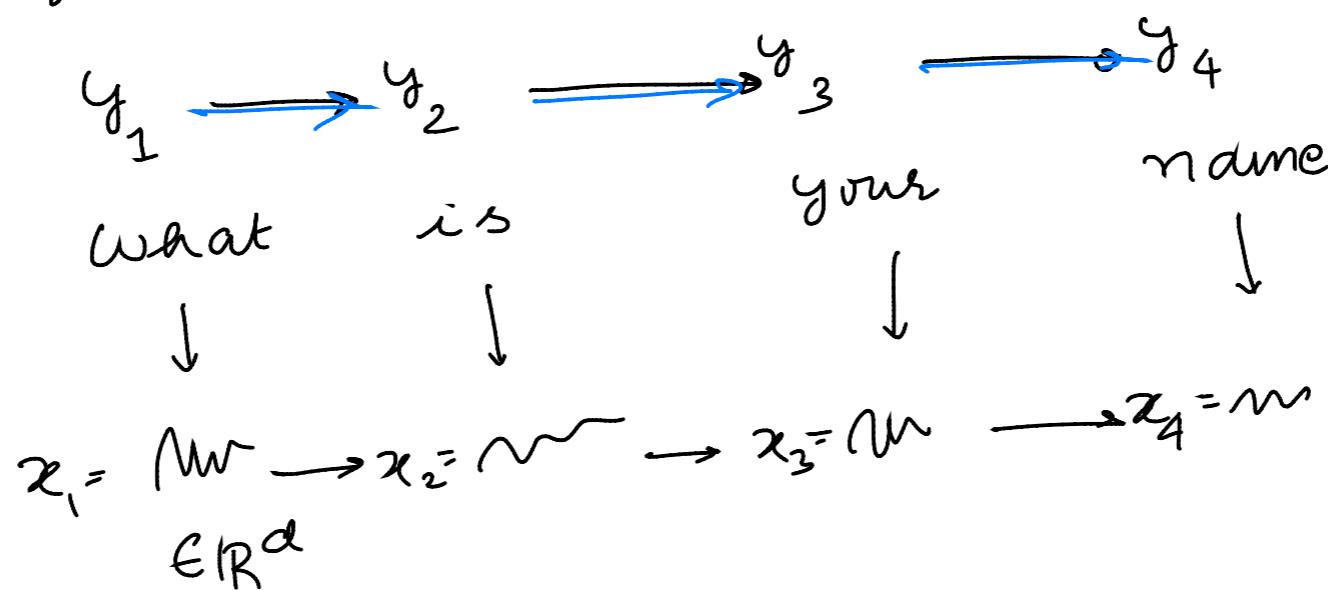
$$(x_i^o, y_i^o)_{i=1}^n$$

$\sim \sim$
Historical
Training Data

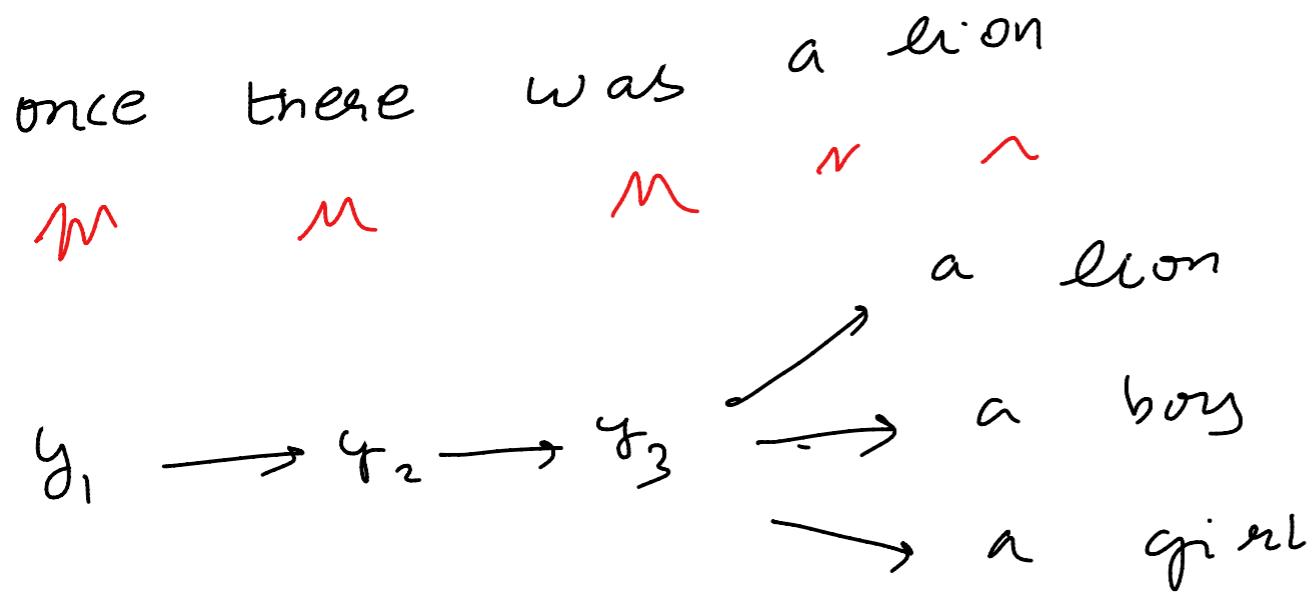


$$y_t = f(y_1, \dots, y_t)$$

Speech recognition:



NOT Speech Recognition:



example shows
that text is
temporally connected

Stock Prediction:

All the
information
that dictates
dynamics of finance

$$y_1 \rightarrow y_2 \rightarrow \dots \rightarrow y_t$$

$$x_1 \quad x_2 \quad \dots \quad x_t \\ \in \mathbb{R}^{100}$$

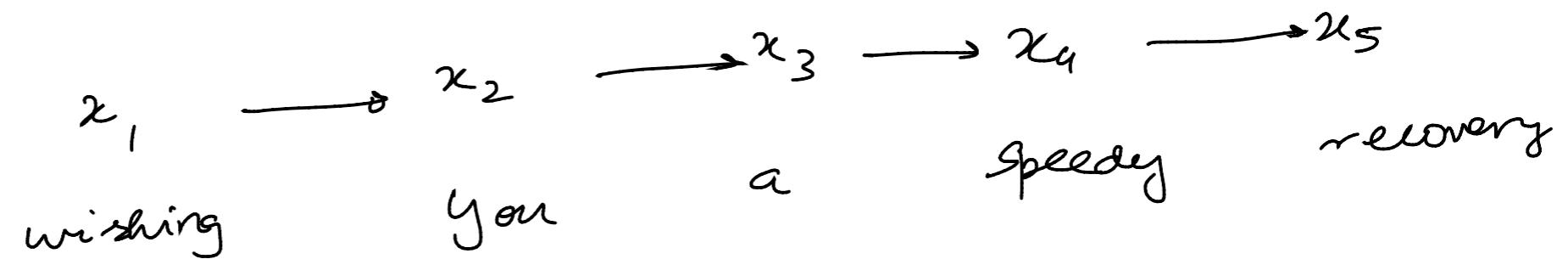
Price of 100 Stock :

Observed / watched

Predicted :

$$\hat{y}_{t+1} = \hat{x}_{t+1}, \dots, \hat{y}_{t+k} = \hat{x}_{t+k}$$

Sentiment Analysis:



$$\hat{y} = +ve / -ve$$

Dataset: $x_1 \dots x_5$ + ve / -ve
 $x_1 \dots x_{100}$ +ve / -ve

$ANN = MLP$ (Generic Word)

Depth

Deep Neural Network

(several variants based on the problem they solve)

One step prediction

Feed Forward

Sequence Prediction

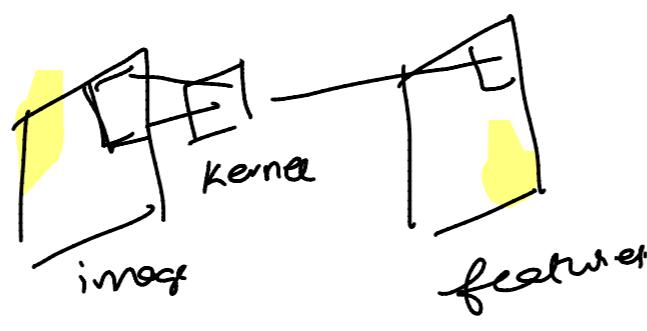
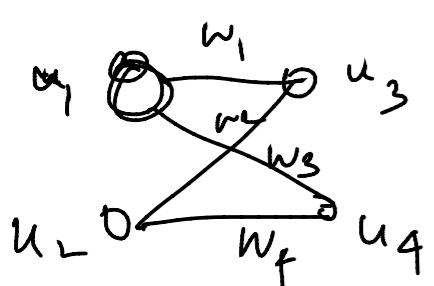
No feedback

Transformers

Feedback

RNN
LSTM
GRU

Dense Fully connected with weight sharing (for images)
Sparse weights + uses convolutions to extract local features from images.



One step prediction: 1 input, 1 output

Sequence Prediction

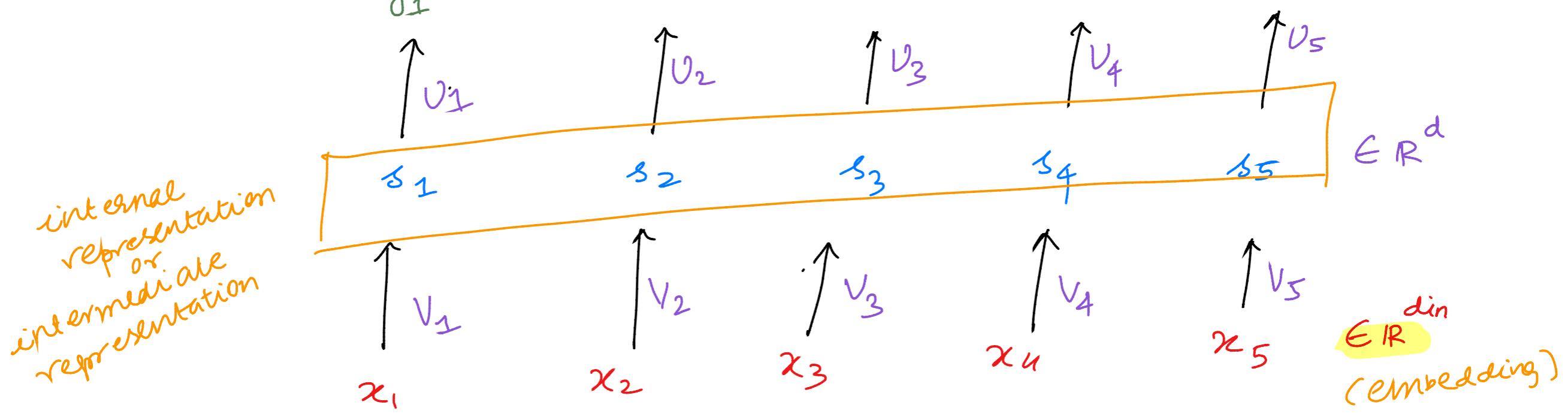
- Both input and output sequences are of different length
- many ip, many op : machine translation,
Speech recognition
- many-ip; one op : sentiment analysis
- one-ip; many-op : image captioning
- Sequence prediction is always not a "time series"
- Inputs and outputs are temporally related

Part of speech tagging

(domain = tag)

Output :

noun	Verb	article	adjective	noun
y_1	y_2	y_3	y_4	$y_5 \in \mathbb{R}^{d_{out}}$



Input :

man is a social animal

(domain = vocabulary)

$$U_i \in \mathbb{R}^{d \times d_{in}}$$

$$V_i \in \mathbb{R}^{d_{out} \times d}$$

$$s_i^o = U_i x_i \quad (\text{input to internal})$$

$$y_i^o = V_i s_i^o \quad (\text{internal to output})$$

- U_i, V_i are different for each $i^o \Rightarrow$ no temporal connection yet

$$s_i^o = U x_i$$

$$y_i^o = V s_i^o$$

⇒ one step prediction applied at different times

A step back: How are sequences generated?

$y_1 \ y_2 \ y_3 \ y_4 \ y_5 \ y_6$
0, 0, 1, 0, 0, 1, 0, 0, 1, - - -
Matrix Equivalent

$$S = [0, 0, 1]$$

i = 0
while (true)

{
print S[i:i+3]

i = i + 1
}

$$s_0 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$s_{t+1} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} s_t$$

$$s_1 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \quad s_2 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad s_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$y_t = [0 \ 0 \ 1] s_t$$

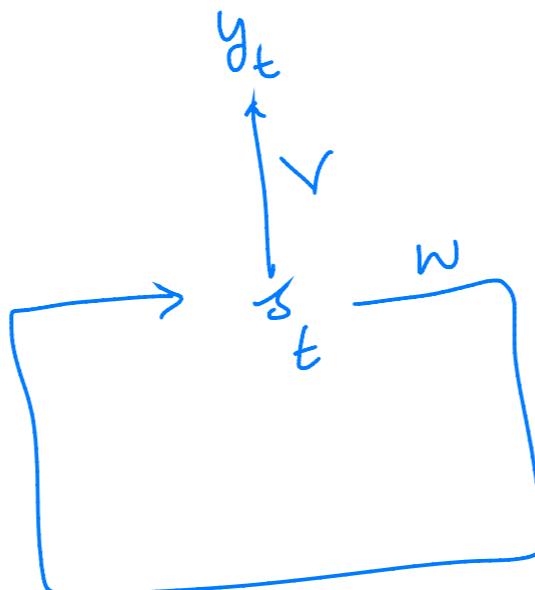
$$y_1 = [0 \ 0 \ 1] \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = 0$$

$$y_2 = \begin{bmatrix} 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = 0, \quad y_3 = \begin{bmatrix} 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = 1$$

0, 0, 1, 0, 0, 1 - - -

$$W = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

$$\begin{aligned} s_{t+1} &= W s_t \\ y_t &= \sqrt{s_t} \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Linear Recurrence without input}$$

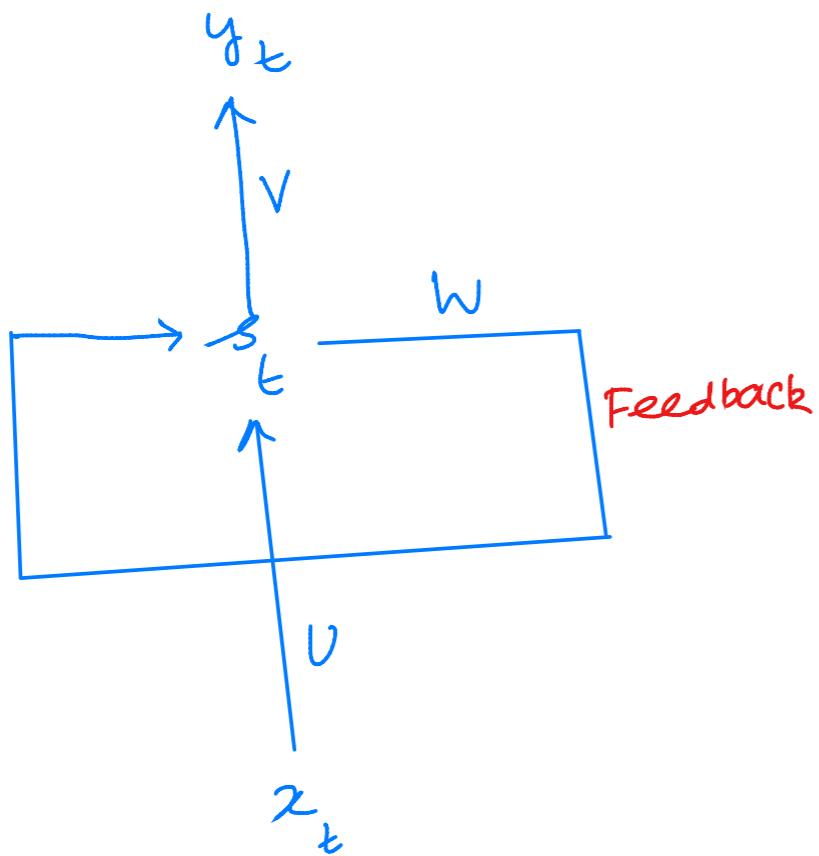


Feedback

$$s_{t+1} = W s_t + U x_t \quad \text{input}$$

output $\rightarrow y_t = \sqrt{s_t}$

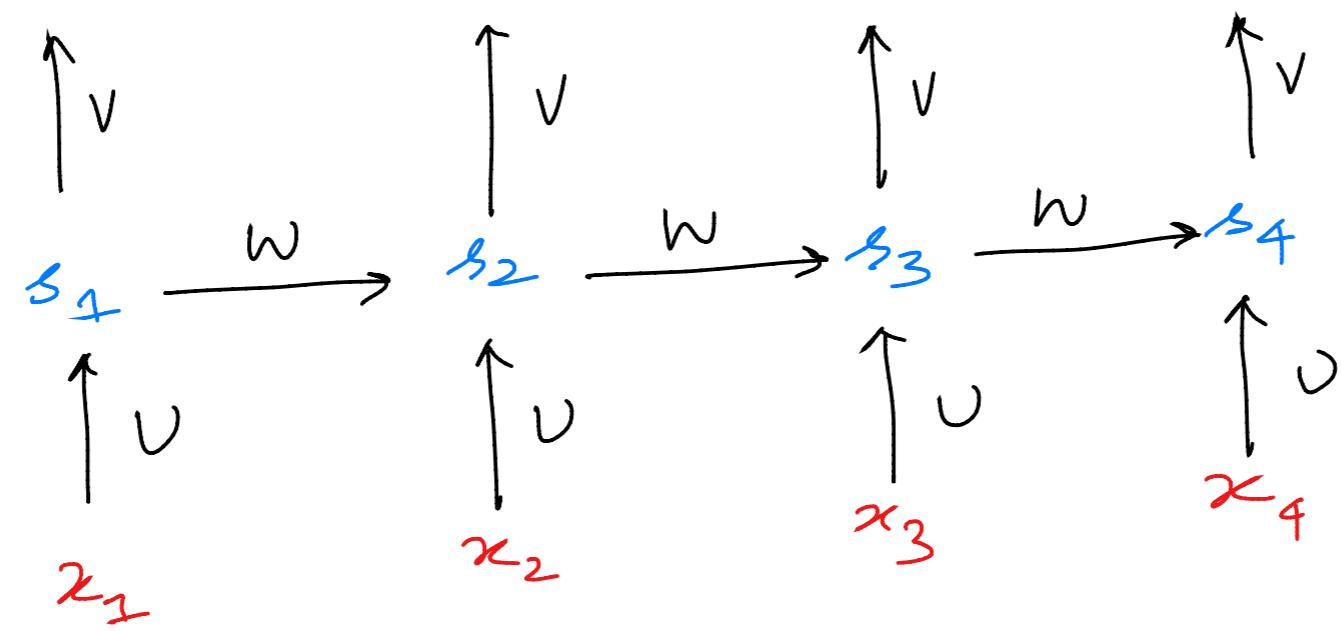
Linear Recurrence
with input



unfurling the recurrence

$$\begin{aligned}
 s_{t+1} &= W s_t + U x_t \\
 &= W (W s_{t-1} + U x_{t-1}) + U x_t \\
 &= W^2 s_{t-1} + W U x_{t-1} + U x_t \\
 &= W^2 (W s_{t-2} + U x_{t-2}) + W U x_{t-1} + U x_t \\
 &\vdots \\
 &= \underbrace{W^t s_1}_{\text{how much of initial memory } s_1 \text{ is remaining}} + \underbrace{W^{t-1} U x_1}_{\text{summary input at } t=1, \text{i.e. } x_1} + \underbrace{W^{t-2} U x_2}_{\text{summary }} + \dots + U x_t
 \end{aligned}$$

summary of x_2



Recurrent Neural Network

$$s_t = \sigma_s (W s_{t-1} + V x_t + \text{bias})$$

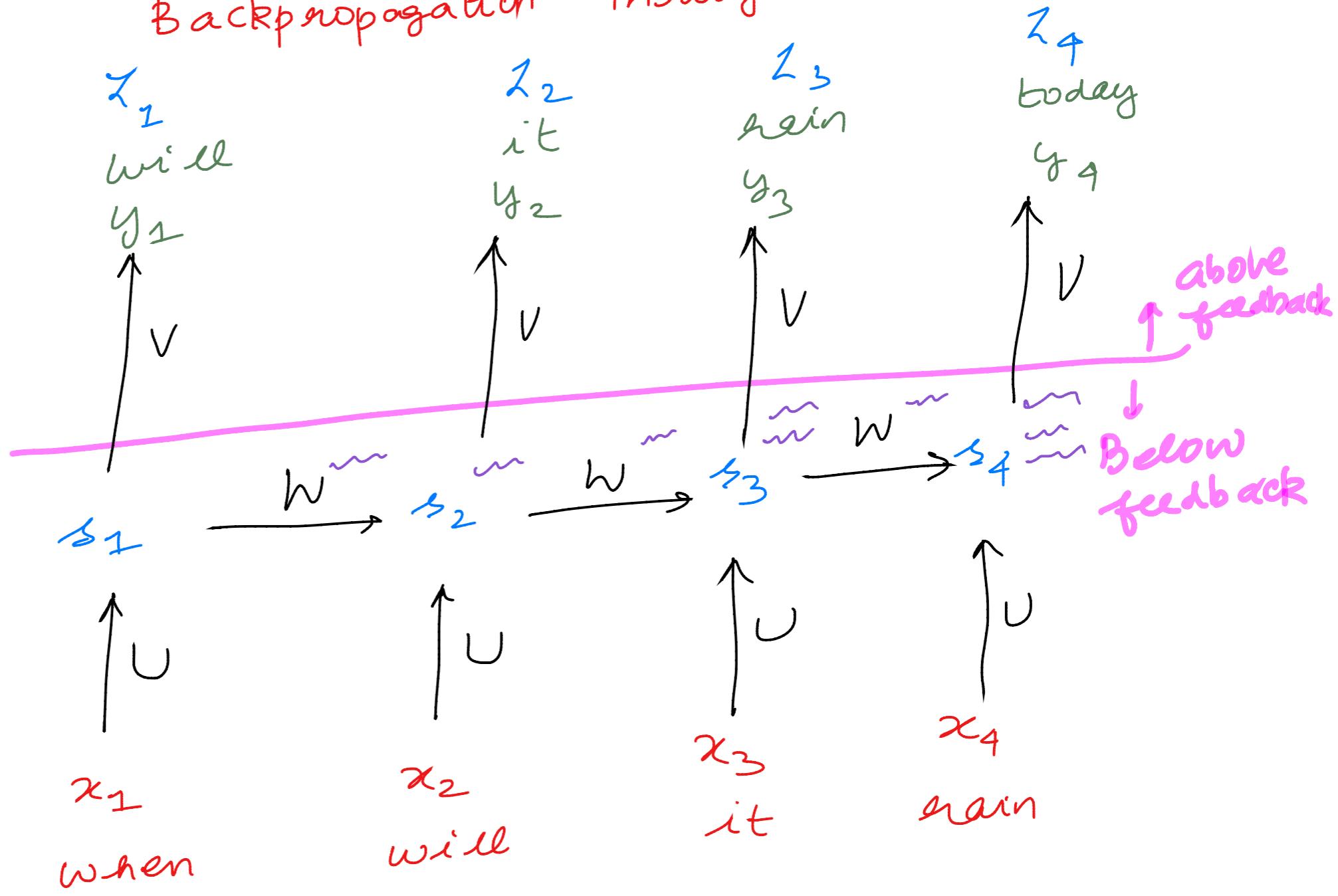
$$y_t = \sigma_y (\nabla s_t + \text{bias})$$

- $W \in \mathbb{R}^{d \times d}$
- $V \in \mathbb{R}^{d \times \text{din}}$
- $\nabla \in \mathbb{R}^{\text{dout} \times d}$

RNN Training

Backpropagation Through Time

Output :



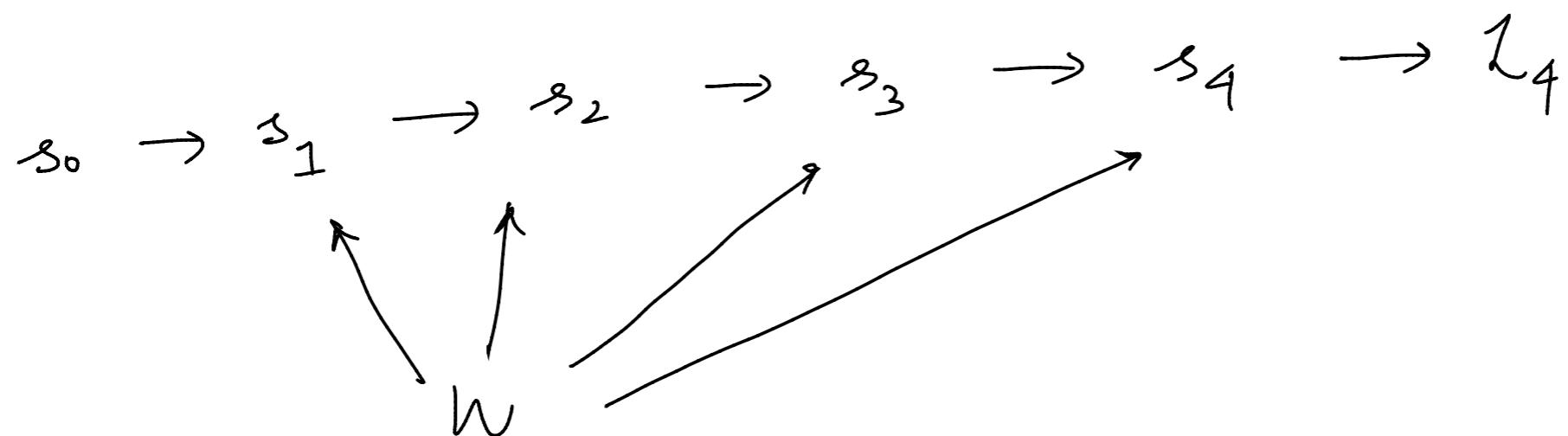
Input :

y_1 can be the probability over vocabulary
 z_1 : can be cross entropy loss

$$\mathcal{L} = \mathcal{L}_1 + \mathcal{L}_2 + \mathcal{L}_3 + \mathcal{L}_4$$

$\frac{\partial \mathcal{L}}{\partial v}$ calculation is standard back propagation
because v is above the feedback

$\frac{\partial L}{\partial w}$ involves time



Moral: change in w causes a ripple effect
and capturing this is called
Backpropagation Through Time

$$\frac{\partial L_4}{\partial w} = \frac{\partial L_4}{\partial s_4} \cdot \frac{\partial s_4}{\partial w}$$

To compute

$$\frac{\partial s_4}{\partial w}$$

• direct part :

$$\frac{\partial^+ s_4}{\partial w} = \text{keeps } s_3 \text{ a constant}$$

• indirect part :

effect of changing w
on s_3

$$\frac{\partial s_4}{\partial w} = \underbrace{\frac{\partial^+ s_4}{\partial w}}_{\text{direct}} + \underbrace{\frac{\partial s_4}{\partial s_3} \frac{\partial s_3}{\partial w}}_{\text{indirect}}$$

$$= \frac{\partial^{+} s_4}{\partial w} + \frac{\partial s_4}{\partial s_3} \left(\frac{\partial^{+} s_3}{\partial w} + \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial w} \right)$$

$$\frac{\partial s_4}{\partial w} = \frac{\partial^{+} s_4}{\partial w} + \frac{\partial s_4}{\partial s_3} \frac{\partial^{+} s_3}{\partial w} + \frac{\partial s_4}{\partial s_3} \cancel{\frac{\partial s_3}{\partial s_2}} \frac{\partial s_2}{\partial w} + \frac{\partial s_4}{\partial s_3} \cdot \cancel{\frac{\partial s_3}{\partial s_2}} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial w}$$

$$= \frac{\partial^{+} s_4}{\partial w} + \frac{\partial s_4}{\partial s_3} \frac{\partial^{+} s_3}{\partial w} + \frac{\partial s_4}{\partial s_2} \frac{\partial^{+} s_2}{\partial w} + \frac{\partial s_4}{\partial s_1} \frac{\partial^{+} s_1}{\partial w}$$

$$\frac{\partial s_4}{\partial w} = \sum_{k=1}^1 \frac{\partial s_4}{\partial s_k} \frac{\partial^{+} s_k}{\partial w}$$

Backpropagation Through Time

$$\frac{\partial L_4}{\partial w} = \frac{\partial L}{\partial s_4} \left(\sum_{k=1}^1 \frac{\partial s_4}{\partial s_k} \frac{\partial^{+} s_k}{\partial w} \right)$$

Taken care under the hood by Tensorflow / PyTorch / Caffe

Issues with RNNs

$$s_t = \sigma_s (W s_{t-1} + U x_t + b_s)$$

$$y_t = \sigma_y (V s_t + b_y)$$

Simplify by removing σ_s , σ_y , b_s , b_y

$$s_t = W s_{t-1} + U x_t$$

$$y_t = V s_t$$

Further simplify by making dimension = 1

$$s_t = W s_{t-1} + U x_t$$

$$y_t = V s_t$$

Further simplify, let $U = 0$

$$s_t = w s_{t-1}$$

$$y_t = v s_t$$

$$s_1 = w s_0, \quad s_2 = w s_1 = w^2 s_0, \quad s_3 = w^3 s_0, \dots, \quad s_t = w^t s_0$$

$|w| > 1$: for $s_0 \neq 0$, $s_t \rightarrow \infty$ (exploding memory)

$|w| < 1$: $s_t \rightarrow 0$ (forgetting/fading memory)

$w = 1$: s_t remains the same

$w = -1$: s_t flips sign

Moral: Because of feedback memory either explodes or fades

Long - Short Term Memory Networks

explicit control of head, writing and forgetting

We will use the white board analogy by Prof. Mitesh

Idea: Think of b_t to be whiteboard

computation of $ac(bd+a) + ad$

toy example:

constraint: we can only have a summary of
3 statements at a time

$$\bullet \quad ac = 5$$

$$\bullet \quad bd = 33$$

$$\bullet \quad bd+a = 34$$

selective writings

$ac, bd, ad, a, bd+a$

selective reading of input
 $ad \vee (bd+a)$

- $ac = 5$
 - $ac(bd+a) = 170$
 - $bd+a = 34$
- selectively forgetting
bd

- $ac = 5$
 - $ac(bd+a) = 170$
 - $ad = 11$
 - $ad + ac(bd+a) = 181$
 - $ac(bd+a) = 170$
 - $ad = 11$
- select $bd+a = 34$
for forgetting
- push $ac = 5$
out

Moral : Selective Reading, Write, Forget

LSTM

Selective Write

$$\begin{bmatrix} 1.2 \\ -3 \\ \vdots \\ \vdots \\ 5 \end{bmatrix} \odot \begin{bmatrix} 0.2 \\ 0.5 \\ \vdots \\ \vdots \\ 0.3 \end{bmatrix} = \begin{bmatrix} 1.2 \times 0.2 \\ -3 \times 0.5 \\ \vdots \\ \vdots \\ 5 \times 0.3 \end{bmatrix}$$

s_{t-1} element wise product O_{t-1} h_{t-1}

$$h_{t-1} = s_{t-1} \odot o_{t-1}$$

Output Gate

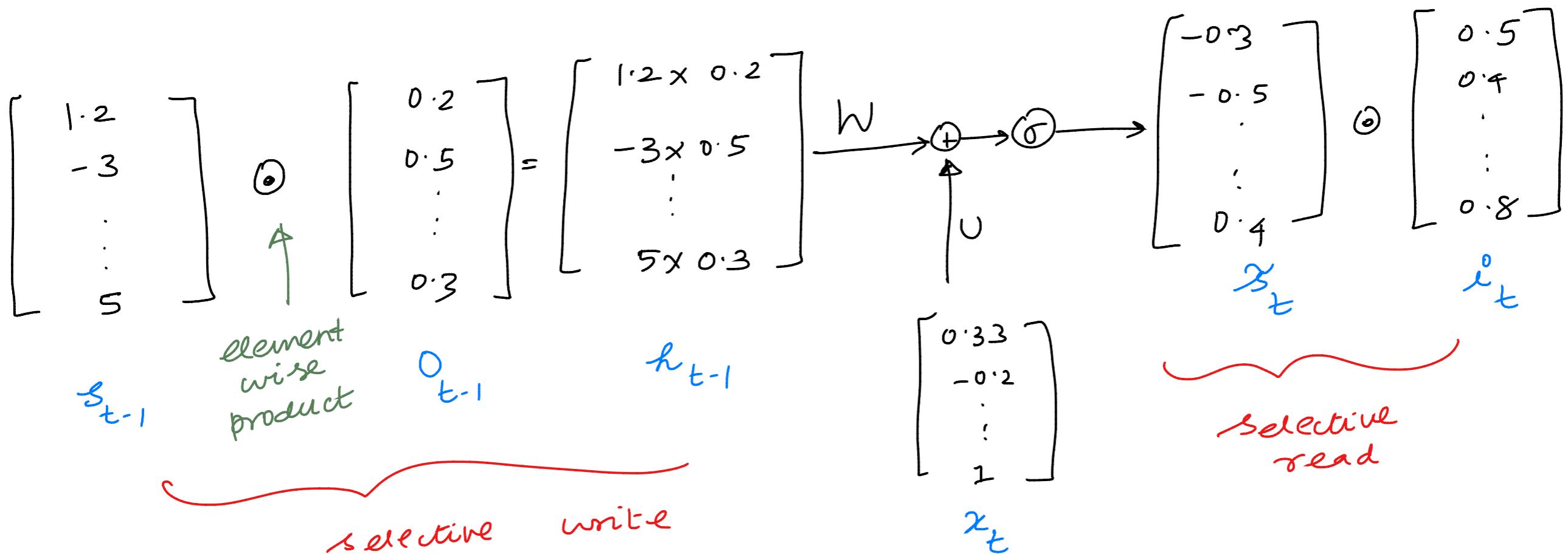
$$o_{t-1} = \sigma (W_o h_{t-2} + U_o x_{t-1} + b_o)$$

$$\tilde{s}_t = \sigma (W h_{t-1} + U x_t + b)$$

$$i_t = \sigma (W_i h_{t-1} + U_i x_t + b_i)$$

Input Gate

$$s_t = \tilde{s}_t \odot i_t$$



Selective Forget

$$f_t = \sigma(W_f h_{t-1} + v_f x_t + b_f)$$

$$s_t = f_t \underbrace{\odot s_t}_{\begin{matrix} \uparrow \\ \text{forget} \\ \text{summary} \end{matrix}} + \underbrace{i_t \odot \tilde{s}_t}_{\begin{matrix} \text{selective} \\ \text{reading} \end{matrix}}$$

LSTM

States

$$o_t = \sigma(W_o h_{t-1} + v_o x_t + b_o)$$

$$i_t = \sigma(W_i h_{t-1} + v_i x_t + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + v_f x_t + b_f)$$

States

$$\tilde{s}_t = \sigma(W h_{t-1} + v x_t + b)$$

$$s_t = f_t \odot s_{t-1} + i_t \odot \tilde{s}_t$$

$$h_t = \frac{o}{t} \odot \sigma(s_t)$$