Let's revisit the task of autocompletion.



vocab
$= \{ I, at, today, home, am \}$

| j | | |
|---|---|---|
| j=0 | 0.1 | I |
| j=1 | 0.1 | at |
| j=2 | 0.2 | today |
| j=3 | 0.1 | home |
| j=4 | 0.5 | am |

At any time step $t$, we want to compute

$$\underset{j \in V}{\arg\max} \; P\left( y_t = j \mid y_{t-1}, y_{t-2}, \ldots, y_1 \right)$$

word ⟵ Vocabulary

eg.  $P\left( y_4 = home \mid at, am, I \right)$
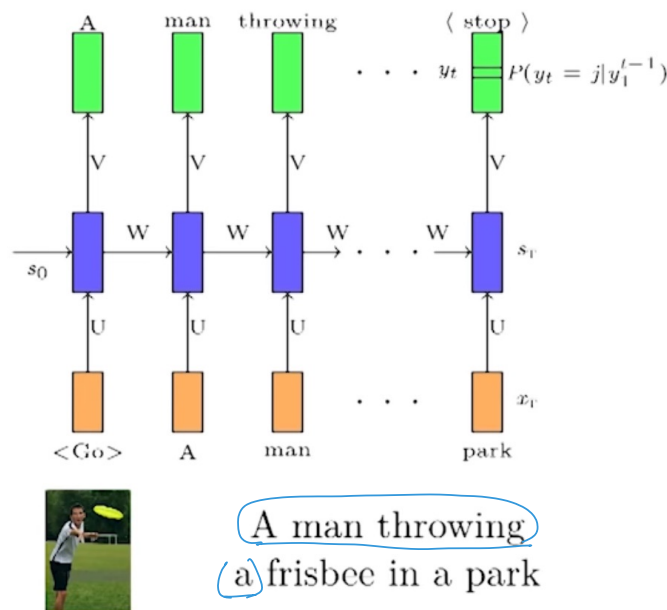
Using RNN,

$$P\left( y_t = j \mid y_{t-1}, y_{t-2}, \ldots, y_1 \right) = \left( softmax\left( V s_t + c \right) \right)_j$$

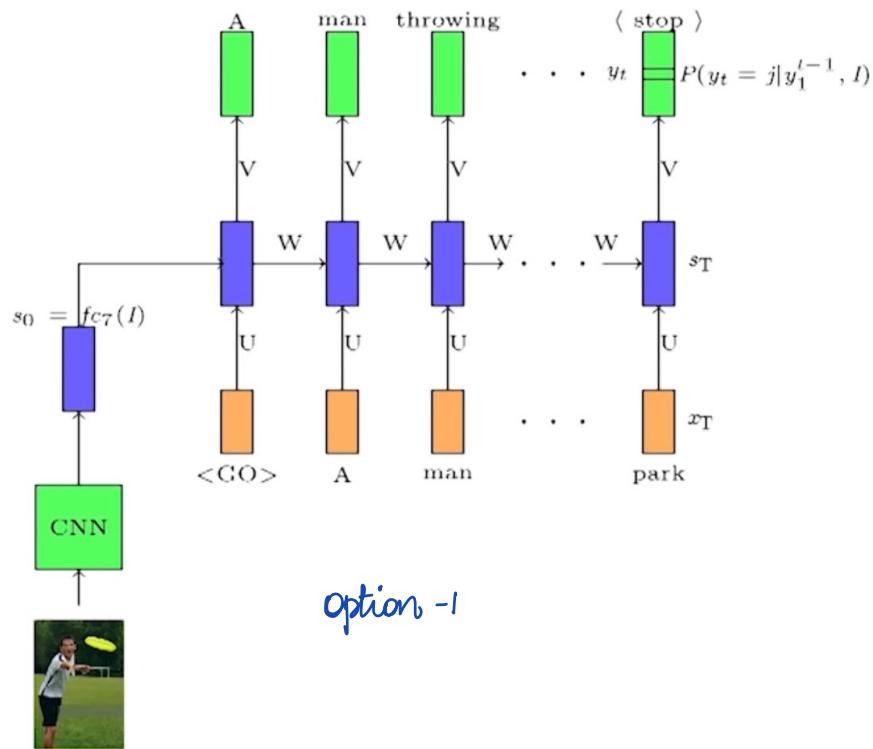$s_t$ captures all the information until time step $t$.

At test time,



o/p: I    am    at    home   today   &lt;stop&gt;

&lt;GO&gt;

$\left.\begin{array}{c}\end{array}\right\}$ one hot encoding of the words.

<span style="color:red">**Image Captioning**</span>



A   man  throwing     ⟨ stop ⟩

$\cdots \quad y_t \quad P(y_t = j|y_1^{t-1})$

9:05 AM

$s_0$   W   W   W $\cdots$ W   $s_\tau$

U   U   U   U

&lt;Go&gt;    A    man       park

A man throwing a frisbee in a park

We want to generate a sentence given an image. → Image.
We are interested in $P\left(y_t = j \mid y_{t-1}, \ldots, y_1, I\right)$

Since CNN architectures are good for images, we use it to learn important features from the image and then pass it on to the RNN model.
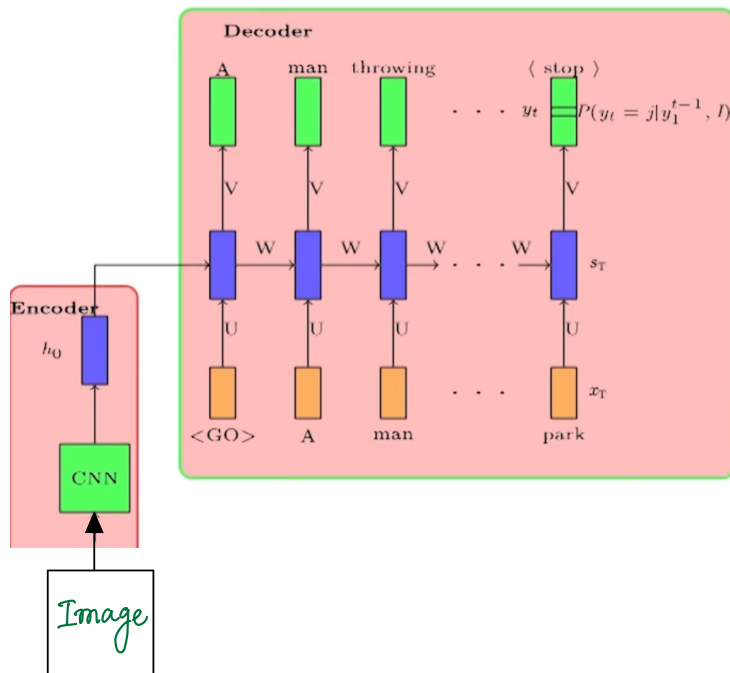


A    man   throwing    ⟨ stop ⟩

$$y_t \quad P(y_t = j | y_1^{t-1}, I)$$

$$s_0 = fc_7(I)$$

$$s_T$$

$$x_T$$

<GO>   A   man    park

CNN

Option -1

A    man    throwing    ⟨ stop ⟩

$y_t$  $P(y_t = j | y_1^{t-1}, I)$

V    V    V        V

W    W    W . . . W    $s_T$

$fc_7(I)$

U    U    U        U

CNN        $x_T$

⟨GO⟩    A    man        park

option-2

## Encoder - Decoder

**Decoder**

A    man    throwing    ⟨ stop ⟩

$y_t$  $P(y_t = j | y_1^{t-1}, I)$

V    V    V        V

W    W    W . . . W    $s_T$

**Encoder**

$h_0$

U    U    U        U

$x_T$

CNN

⟨GO⟩    A    man        park

Image

In image captioning,
A CNN is used to "encode" the image
        − Learn good feature representation of
        the input.

An RNN is used to "decode" a sentence from this encoding.

**Task:** Image Captioning

**Data:** $\{ x_i = \text{image} , \ y_i = \text{caption} \}_{i=1}^{N}$

**Model:**
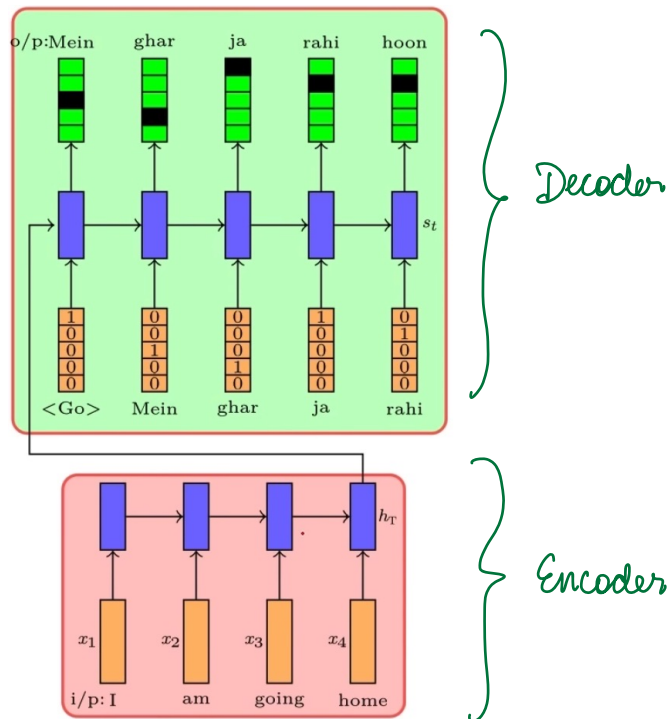
Encoder

$$s_0 = \text{CNN}(x_i)$$

Decoder

$$s_t = \text{RNN}\left( s_{t-1} , \ y_{t-1} \right)$$

**Parameters:** $U, V, W, b, c$ and all weights and biases of CNN.

Training all parameters of encoder and decoder together.

— End to end model

# Machine Translation



o/p:Mein · ghar · ja · rahi · hoon

Decoder

$s_t$

$\begin{smallmatrix}1\\0\\0\\0\\0\end{smallmatrix}$ $\begin{smallmatrix}0\\0\\0\\1\\0\end{smallmatrix}$ $\begin{smallmatrix}0\\0\\0\\0\\1\end{smallmatrix}$ $\begin{smallmatrix}1\\0\\0\\0\\0\end{smallmatrix}$ $\begin{smallmatrix}0\\1\\0\\0\\0\end{smallmatrix}$

<Go> · Mein · ghar · ja · rahi

$h_T$

Encoder

$x_1$ · $x_2$ · $x_3$ · $x_4$

i/p: I · am · going · home

# Video Captioning



A · man · walking · on · a · rope

Decoder

Encoder

$s_1$ · $s_2$

CNN · CNN · CNN