

Sequence - to - sequence

Language modeling

$$P(x_t | x_{t-1}, \dots, x_1)$$

How to train? Collect x_1, x_2, \dots, x_T

"Teacher forcing": Feed in sequence up to x_t ,
train model to predict x_t

Sequence classification

$$P(y | x_1, \dots, x_T)$$

$$h_t = f(\underbrace{x_t, h_{t-1}}_{RNN})$$

$$o = W_o h_T + b_o$$

\uparrow
scores

Sequence - to - sequence

$$p(y_t | \underbrace{y_{t-1}, y_{t-2}, \dots, y_1}_{\text{current output}}, \underbrace{x_T, x_{T-1}, \dots, x_1}_{\text{entire input}})$$

ex Machine translation, x = English sentence
 y = French translation of x

How to build a seq2seq model with RNNs:

1. $h_t = f(x_t, h_{t-1})$ "encoder" RNN

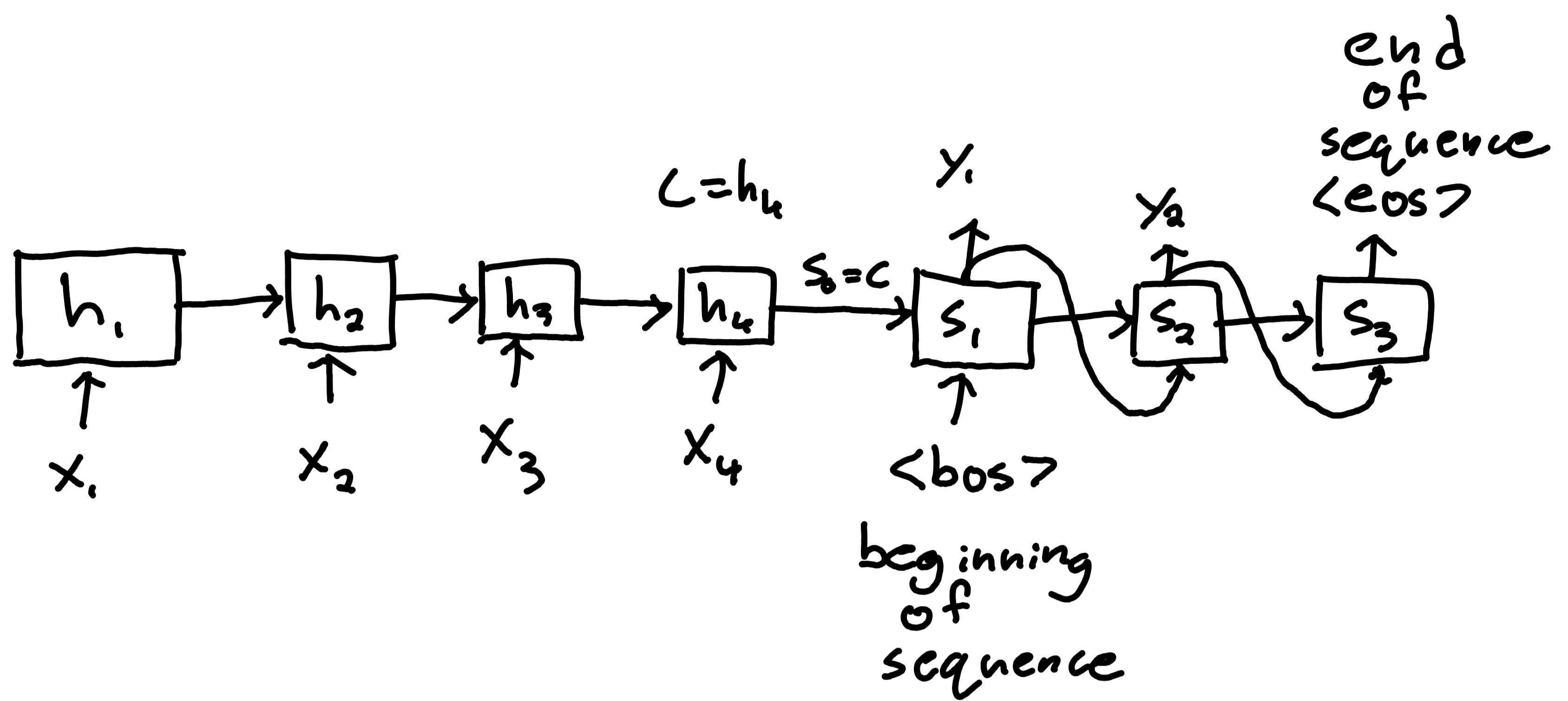
2. $c = h_T$ "context"

3. $s_t = g(y_{t-1}, s_{t-1}, c)$ "decoder" RNN

↳ used to predict y_t

How to make g a function of c ?

e.g. $s_0 = c$ or $\text{concat}(\text{embed}(y_{t-1}), c)$



Padding

Batch of sequences has shape $B \times T$
→ all sequences must be length-T!

Use padding:

B	2	3	9	4	5	3
	1	5	3	0	0	0
T						

padding

① is special token (like <bos>)

For encoder, padding token → ignore input, set $h_t = h_{t-1}$

For decoder, padding token → don't apply loss

How to generate outputs?

1. Auto regressive sampling

$$y_t \sim p(y_t | y_{t-1}, \dots, y_1, x_T, \dots, x_1)$$

$$y_{t+1} \sim p(y_{t+1} | \underbrace{y_t, \dots, y_1, x_T, \dots, x_1}_{\text{from previous step}})$$

2. Greedy sampling:

$$y_t = \operatorname{argmax}[p(y_t | y_{t-1}, \dots, y_1, x_T, \dots, x_1)]$$

3. Top-k sampling:

$$y_t \sim \operatorname{top-k}[p(y_t | y_{t-1}, \dots, y_1, x_T, \dots, x_1)]$$

4. Top-p sampling (Nucleus sampling)

Only sample the tokens that make up p% of the total probability mass

5. Beam search: Consider the k most probable prefixes at each timestep.

Attention

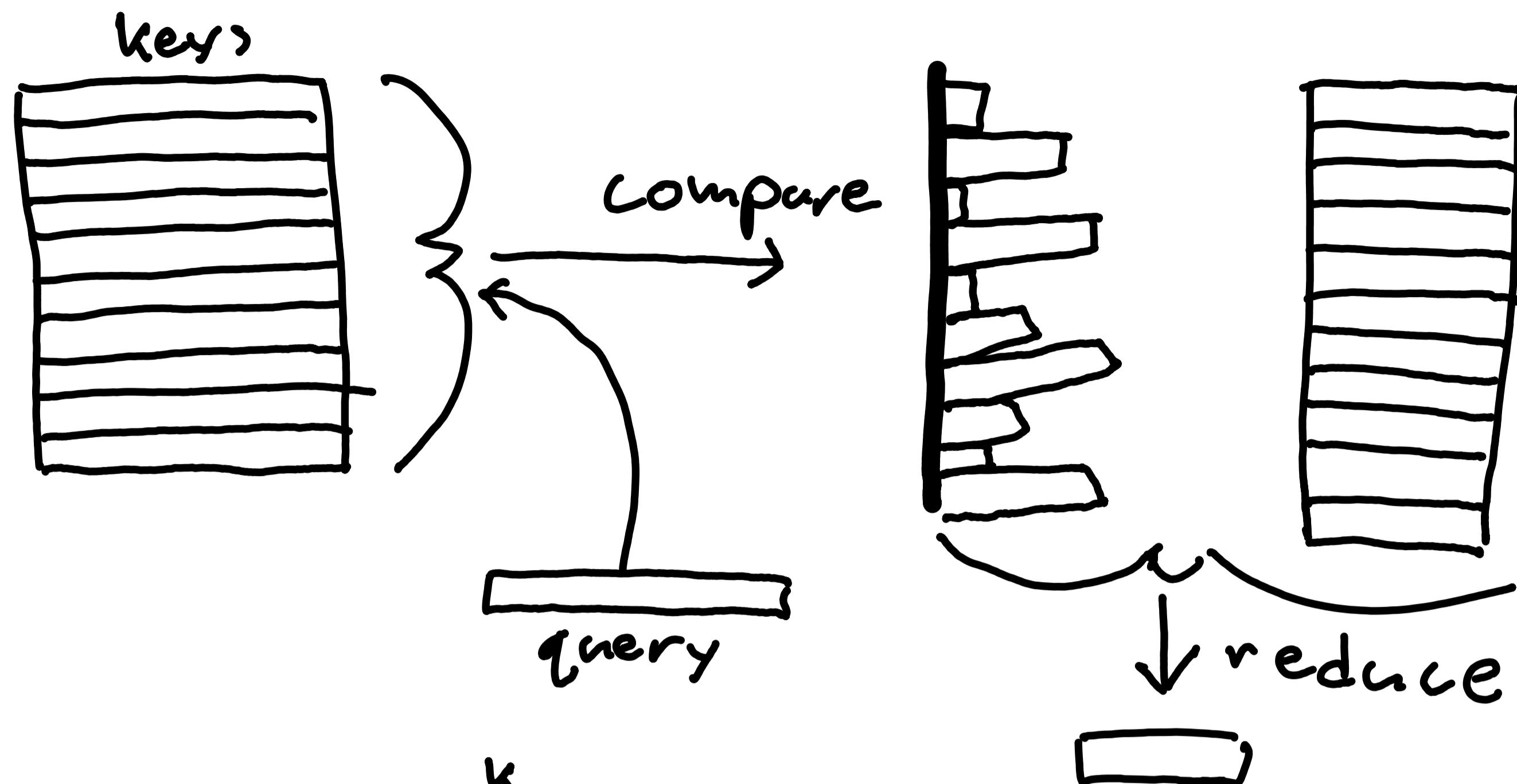
Read-only associative array / key-value store

- Values: set of things we're storing
- Keys: identifiers, one for each value
- Query: identifier of the thing to retrieve

Implements two functions:

- Comparison function: determines how well the query matches a given key $\alpha(k_i, q)$
 - Redaction function: Given results from the comparison function, return something from values
- ~~ex dict.~~ Comparison function tests for equality, reduction function returns the value for matching key.

"Soft" (differentiable) associative memory
aka attention mechanism



Keys: $K_i \in \mathbb{R}^k$, $i \in \{1, \dots, N\}$

query: $q \in \mathbb{R}^k$

values: $V_i \in \mathbb{R}^v$, $i \in \{1, \dots, N\}$

comparison: $a_i = \alpha(q, K_i) \in \mathbb{R}$, $a \in \mathbb{R}^N$

Reduction: $f(\{V_1, \dots, V_N\}, a) \in \mathbb{R}^v$

~~ex~~ Soft 1-nearest neighbor classifier
query: data point we want to classify
keys: datapoints from the training set
values: labels of datapoints from the training set

$$\alpha(q, k_i) = -\|q - k_i\| = a;$$

$$f(\{v_1, \dots, v_N\}, a) = \text{softmax}(a) \cdot v$$

$$V = \begin{matrix} R^N \\ \hline \hline \hline \hline \hline \end{matrix}$$

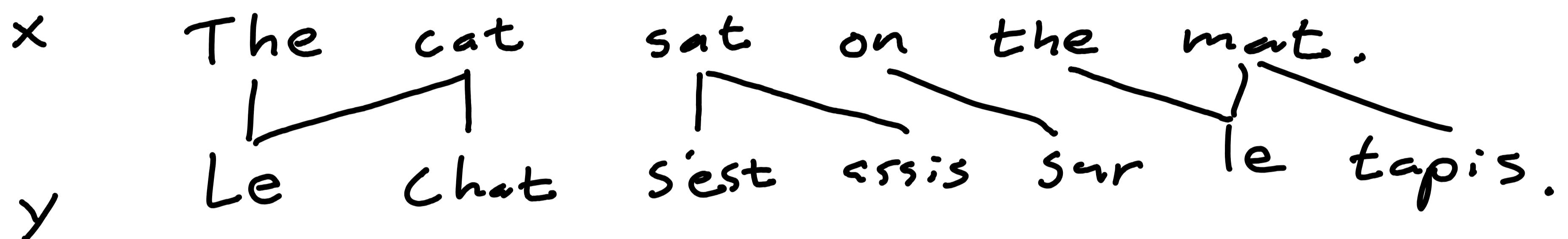
\downarrow
dot
product

if we replace softmax with arg max, just a normal
1-nn classifier

Sequence - to - sequence with attention

If $c = h_T$, then h_T has to encode all relevant information about x .

Decoder output at some time step may depend on some specific inputs.



Rather than having a single c , let's have a different c_t at each output timestep.

query: s_{t-1} (decoder state from previous timestep)

values: h_1, h_2, \dots, h_T (encoder states)

keys: h_1, h_2, \dots, h_T

$$d(q, k_i) = w_v^T \tanh(w_q q + w_k k_i) = q_i$$

learnable parameters

$$c_t' = \text{softmax}(q) \cdot [h_1, h_2, \dots, h_T]$$

use c_t' as an additional input to
the decoder RNN.

$$s_t' = g(s_{t-1}, y_{t-1}, c_t')$$

Self-attention

In self-attention, the queries, keys, and values all come from the same sequence.

$$x_1, x_2, \dots, x_N, x_i \in \mathbb{R}^m$$

$$\text{queries: } W_Q x_1, W_Q x_2, \dots, W_Q x_N = q_1, \dots, q_N \in \mathbb{R}^d$$

$$\text{keys: } W_K x_1, W_K x_2, \dots, W_K x_N = k_1, \dots, k_N \in \mathbb{R}^d$$

$$\text{values: } W_V x_1, W_V x_2, \dots, W_V x_N = v_1, \dots, v_N \in \mathbb{R}^v$$

$$\alpha(q_i, k_j) = \frac{q_i^T k_j}{\sqrt{d}} \quad (\text{if } q_i \text{ and } k_j \text{ are } \mathcal{N}(0, I))$$

then $q_i^T k_j / \sqrt{d}$ will be too

$$f(\{v_1, \dots, v_N\}, a) = \text{softmax}(a) \cdot v$$

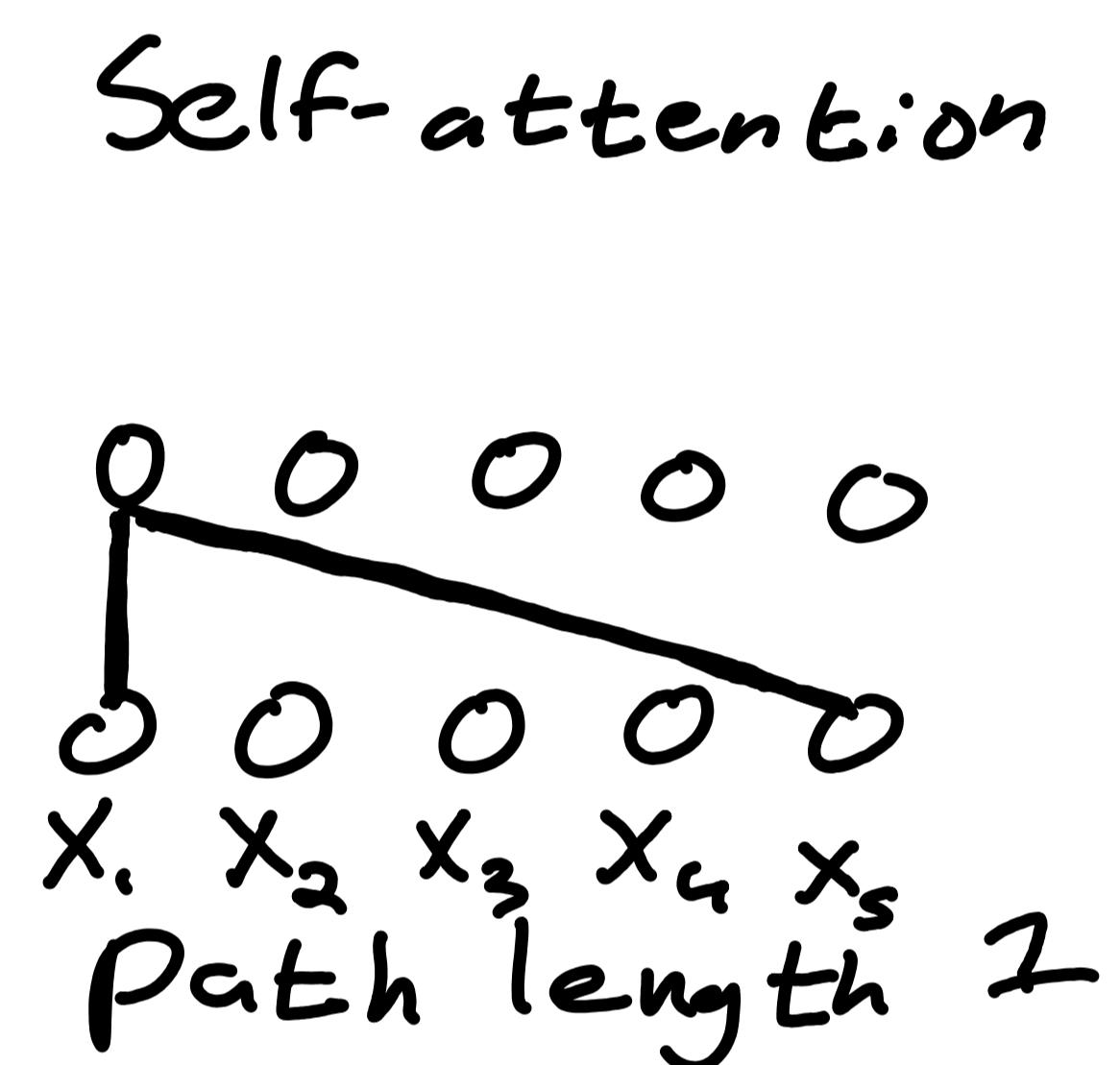
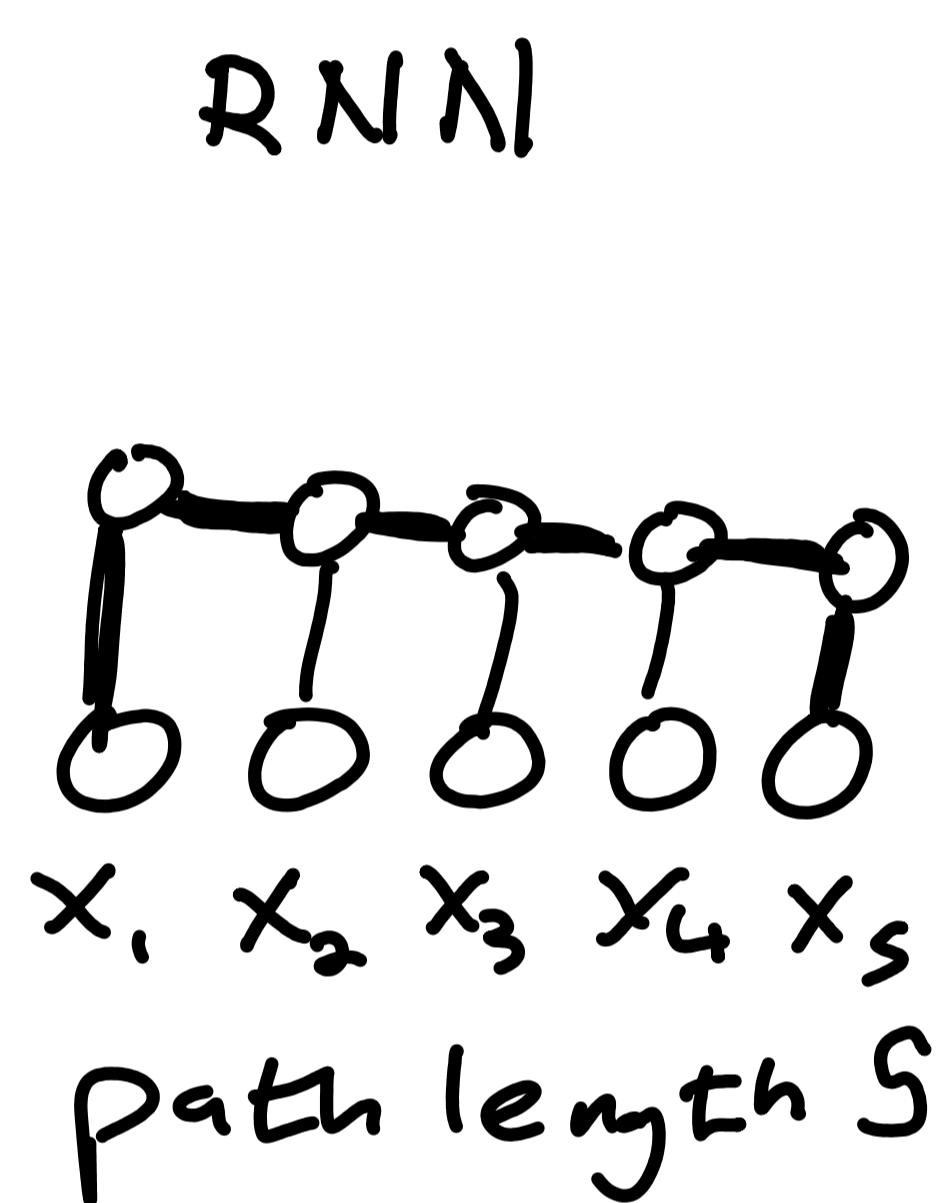
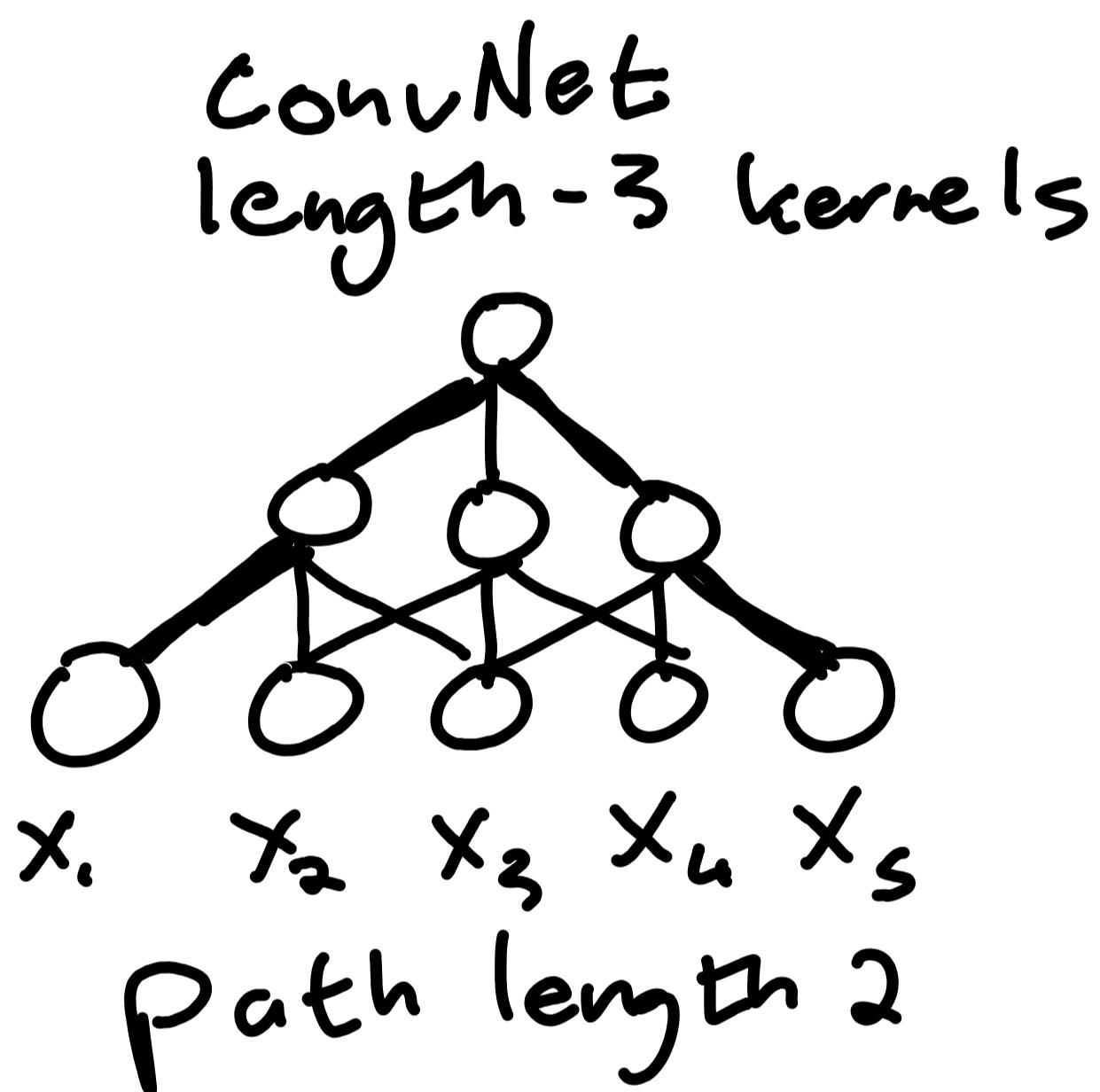
length- N sequence \rightarrow length- N sequence
of weighted averages
of the (projected) input

"Path length"

How many operations are required to "mix" information between two points in the input?

$$x = [x_1, x_2, \dots, x_5]$$

For different models, what is the path length between x_1 and x_5



•