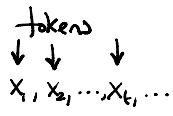


Last time: RNNs

① Language modeling.

Language as a sequence of symbols



Objective: Maximize $P(x_1, x_2, \dots, x_t, \dots)$ of data
However, just like images, this is a high-dimensional space.

Any specific sequence has exponentially small probability!

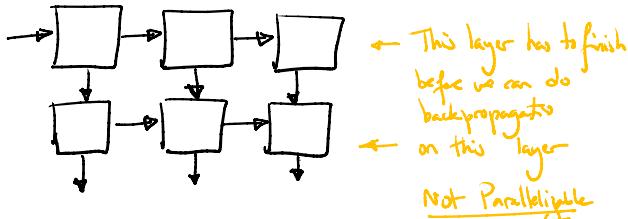
Instead: Write $P(x_1, x_2, \dots, x_t, \dots) \approx P(x_1)P(x_2|x_1)\dots P(x_t|x_{t-1}, \dots)$

Maximize $\sum_t \log P(x_t | x_1, x_2, \dots)$
this probability is not as small!

- RNNs can be used to represent these partial probabilities
How? Often, we let y_t (the output) be logits
We go from logits to probabilities via a softmax:

$$P(x_{t+1} = i | x_1, x_2, \dots) = \frac{\exp(y_{t+1}[i])}{\sum_i \exp(y_{t+1}[i])}$$

- Issue with RNN: they're best with multiple layers



② (Self) Attention

Instead of recurrence, why not directly assess what part information is relevant? (like in databases)

- Query: What the current symbol represents

- Key: The relevance of other sequence symbols to a query.

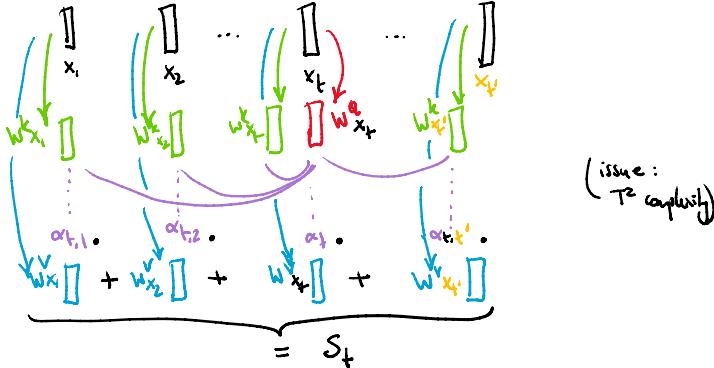
- Value: What each sequence query contributes

$$\begin{aligned} \text{weights} & \quad Q_{t,h} = W^Q x_t & \text{embedding} & \quad K_{t,h} = W^K x_t & \quad V_{t,h} = W^V x_t \\ & \quad \downarrow \text{attention weight} & & \quad \downarrow \text{scaled inner product} & \\ & \quad \alpha_{t,h} = \text{softmax} \left(\frac{1}{\sqrt{k}} Q_{t,h}^T K_{t,h} \right) & & & \end{aligned}$$

(ignore all h at first)

Then, to generate the new state variable, we combine values based on attention: $S_t = W^O \text{concat}_h \sum_h \alpha_{t,h} V_{t,h}$

We could also have multiple heads, copies w/ different weights



Notes: Each head is given by the triple (W^Q, W^K, W^V)

- The weights don't change by position.
- Position info can be placed inside x_t (positional embedding)
- Usually dimensions are maintained $\# \text{heads} \times d = k$, W^O square.

③ Transformers

Attention only tells us where the relevant information is

We still need to act on it. We do it using:

(a) Add and normalize (Residual connection)

$$S' = \text{Layer Norm}(S + X)$$

$$\begin{aligned} p_i &= \frac{1}{T} \sum_j S_{j,i} K_{j,i} \\ \alpha_i &= \frac{1}{T} \sum_j (S_{j,i} K_{j,i})^2 \\ S'_i &= \gamma S_i + X_i - \mu + \beta \end{aligned}$$

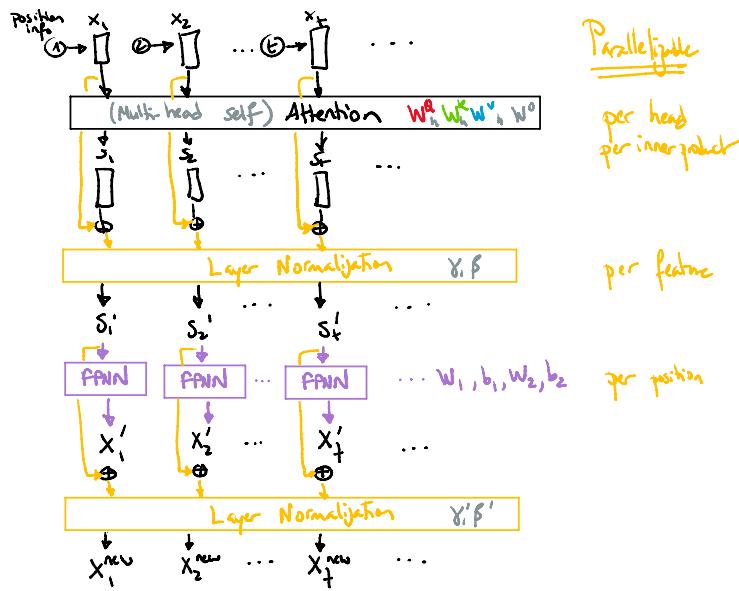
learnable learnable

(b) 2-Layer Feedforward Neural Network

$$X'_t = W^2 \text{ReLU}(W^1 S'_t + b^1) + b^2$$

(c) Add and normalize (again!)

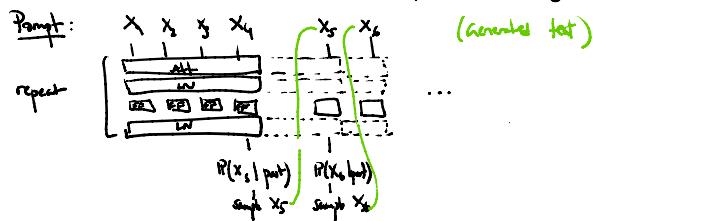
$$X'^{\text{new}} = \text{Layer Norm}(S' + X')$$



④ Using a transformer:

(a) Decoder-only architecture (ChatGPT)

- After the prompt, we put outputs as inputs.
- Attention restricted to past outputs (since we generate forward)



- Variations: Beam Search of most likely generation:

- Maintain k sequences, sample m (expand to km), prune back to top k

(b) Encoder-Decoder Architectures

- The outputs of an encoder transformer become the inputs of a decoder transformer, usually in sandwiched layer:

