# Lecture 9: Recurrent Neural Networks (RNNs)
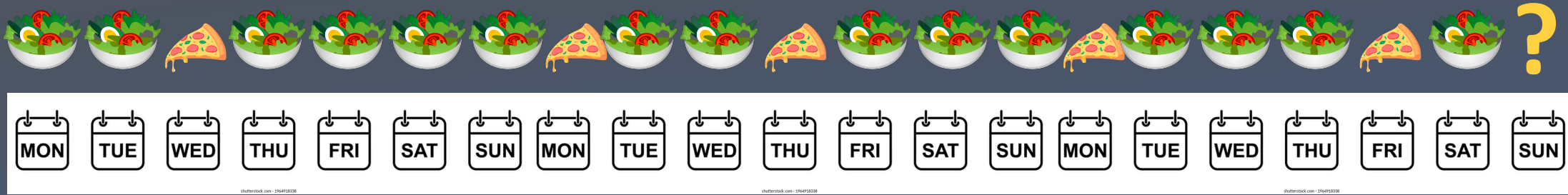
COMPSCI/DATA 182: Deep Learning

09/26/2024

# In **Sequences**, we Trust

- *What is the entree for today ?*
  - What do I factor from history ?
    - How much history ?
    - What tidbits do I need to retain ?
    - What tidbits can I forget about ?
    - Can knowledge of the future help me, as well ?

# Latent Autoregressive Models

$$P(x_1, x_2, \ldots, x_T),$$

$$P(x_1, \ldots, x_T) = P(x_1) \prod_{t=2}^{T} P(x_t \mid x_{t-1}).$$

$$P(x_1, \ldots, x_T) = P(x_1) \prod_{t=2}^{T} P(x_t \mid x_{t-1}, \ldots, x_1).$$

- Leveraging sequences

# TEXT: The Goldmine

- Text is among the most common forms of sequence data encountered in deep learning
  - Common choices for token: characters, words, and word pieces.
- Language models estimate the joint probability of a text sequence
  - For long sequences, ngrams provide a convenient model by truncating the dependence
- Language models can be scaled up with increased data size, model size, and amount in training compute
- PERPLEXITY

$$\text{Perplexity} = \exp(\text{Cross-Entropy Loss})$$

$$\text{Perplexity} = \exp\left(-\frac{1}{N}\sum_{i=1}^{N}\log P(w_i|w_{1:i-1})\right)$$

- Large language models can perform desired tasks by predicting output text given input text instructions.

| Input sequences: | th | e tim | e mac | hine | by h | g wel | ls |
|---|---|---|---|---|---|---|---|

| Target sequences: | the | time | mach | ine b | y h g | well | s |
|---|---|---|---|---|---|---|---|

# Sequences, in Deep Learning

$$\mathbf{H} = \phi(\mathbf{X}\mathbf{W}_{xh} + \mathbf{b}_h).$$

$$\mathbf{H}_t = \phi(\mathbf{X}_t\mathbf{W}_{xh} + \mathbf{H}_{t-1}\mathbf{W}_{hh} + \mathbf{b}_h).$$

H = Hidden layer output
X = Input
W = Hidden layer parameters
$\phi$ = activation function
b = bias parameter
O = output
Batch size n, and d inputs
t = time (stamp)

$$\mathbf{O} = \mathbf{H}\mathbf{W}_{hq} + \mathbf{b}_q,$$

$$\mathbf{O}_t = \mathbf{H}_t\mathbf{W}_{hq} + \mathbf{b}_q.$$
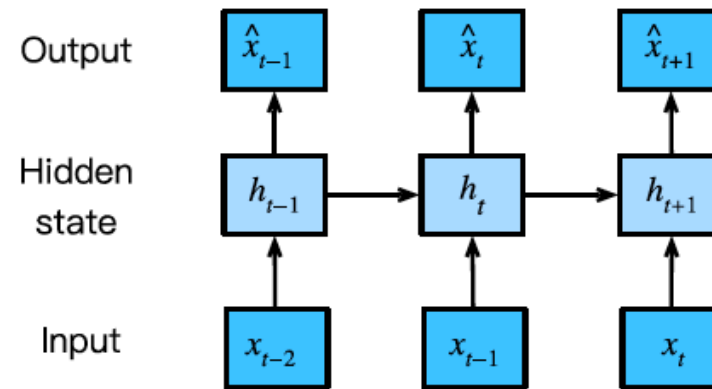
Hidden layer

Output layer

- **Hidden Units**

# Hidden Units

$$\mathbf{H} = \phi(\mathbf{X}\mathbf{W}_{xh} + \mathbf{b}_h).$$

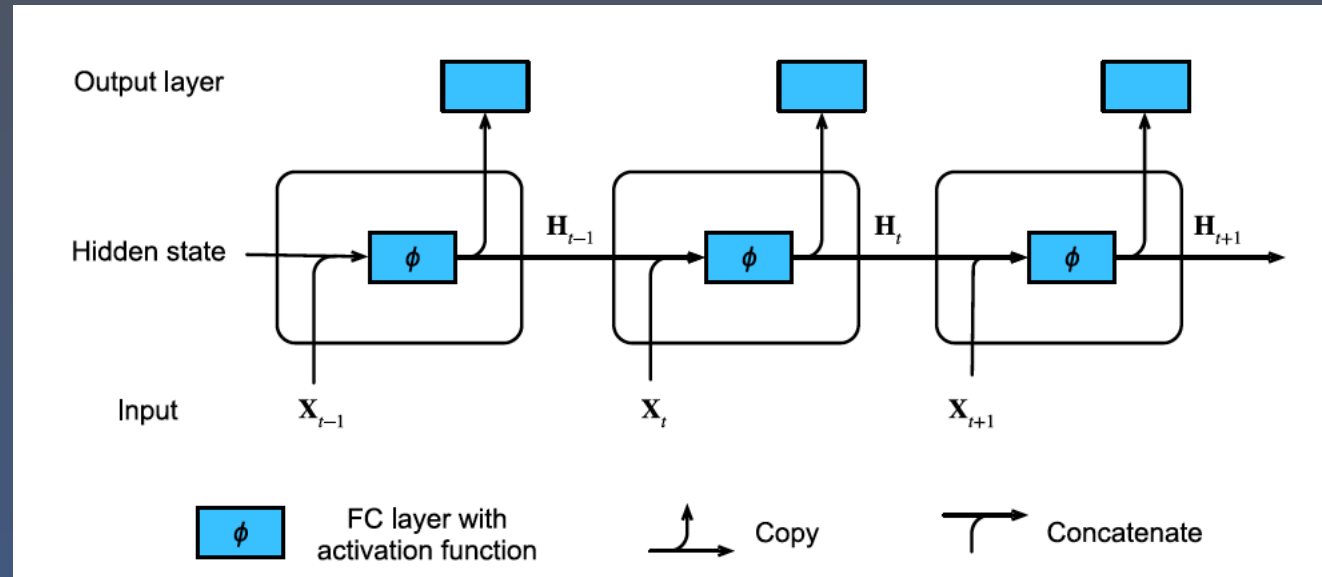$$\mathbf{O} = \mathbf{H}\mathbf{W}_{hq} + \mathbf{b}_q,$$

$$\mathbf{H}_t = \phi(\mathbf{X}_t\mathbf{W}_{xh} + \mathbf{H}_{t-1}\mathbf{W}_{hh} + \mathbf{b}_h).$$

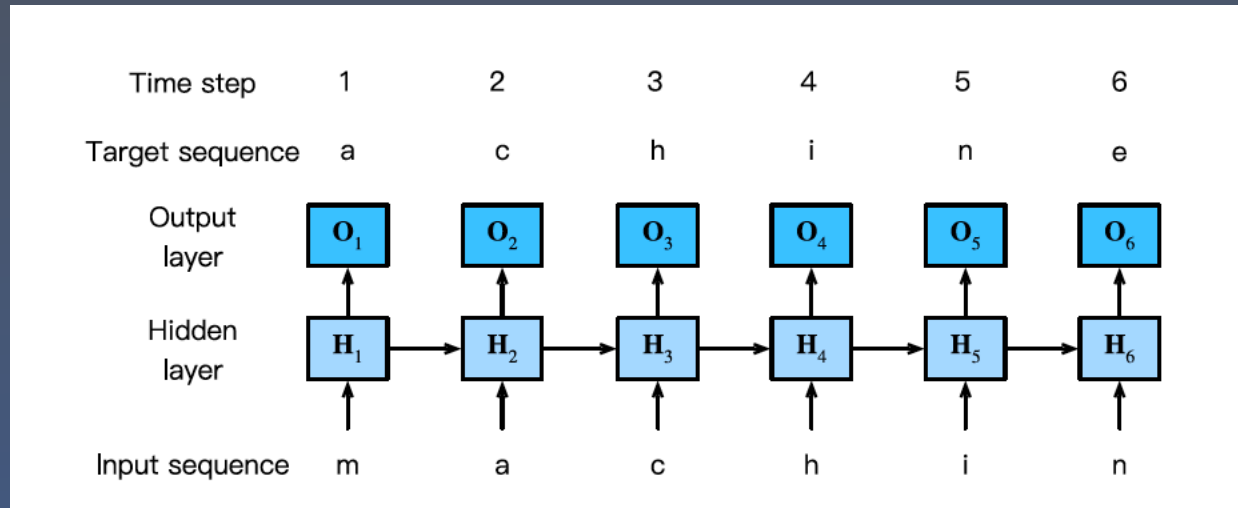$$\mathbf{O}_t = \mathbf{H}_t\mathbf{W}_{hq} + \mathbf{b}_q.$$

# Recurrence



- **Recurrent** computation
- Even at *different* time steps, RNNs always use these *same* model parameters
  - the parametrization cost of an RNN does **not** grow as the number of time steps increases

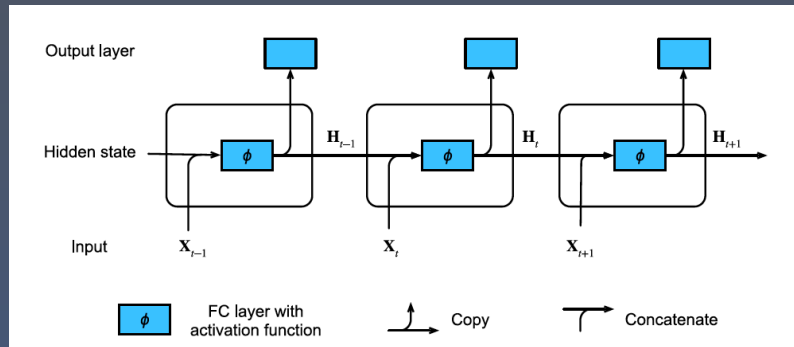# Character sequence (model)

# Gradient Clipping

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L\|\mathbf{x} - \mathbf{y}\|$$

$$|f(\mathbf{x}) - f(\mathbf{x} - \eta\mathbf{g})| \leq L\eta\|\mathbf{g}\|.$$

$$\mathbf{g} \leftarrow \min\left(1, \frac{\theta}{\|\mathbf{g}\|}\right)\mathbf{g}.$$

- Time steps = hidden layers
- More time steps —> more (hidden) layers
  - Vanishing & Exploding gradients

# Loss, in RNNs



$$h_t = f(x_t, h_{t-1}, w_{\mathrm{h}}),$$
$$o_t = g(h_t, w_{\mathrm{o}}),$$

$$L(x_1, \ldots, x_T, y_1, \ldots, y_T, w_{\mathrm{h}}, w_{\mathrm{o}}) = \frac{1}{T} \sum_{t=1}^{T} l(y_t, o_t).$$

- Unrolling across time steps
- Sum the gradients in unrolled (can be rather long)

# Backpropagation through time

$$L(x_1, \ldots, x_T, y_1, \ldots, y_T, w_h, w_o) = \frac{1}{T} \sum_{t=1}^{T} l(y_t, o_t).$$

$$h_t = f(x_t, h_{t-1}, w_h),$$
$$o_t = g(h_t, w_o),$$

$$\frac{d}{dx} f(g(x)) = \frac{df}{dg} \cdot \frac{dg}{dx}$$

$$\frac{\partial L}{\partial w_h} = \frac{1}{T} \sum_{t=1}^{T} \frac{\partial l(y_t, o_t)}{\partial w_h}$$
$$= \frac{1}{T} \sum_{t=1}^{T} \frac{\partial l(y_t, o_t)}{\partial o_t} \frac{\partial g(h_t, w_o)}{\partial h_t} \frac{\partial h_t}{\partial w_h}.$$

$$a_0 = 0 \quad a_t = b_t + c_t a_{t-1}$$

$$a_t = b_t + \sum_{i=1}^{t-1} \left( \prod_{j=i+1}^{t} c_j \right) b_i$$

$$\frac{\partial h_t}{\partial w_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial w_h}$$

$$\frac{\partial h_t}{\partial w_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \sum_{i=1}^{t-1} \left( \prod_{j=i+1}^{t} \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h}$$
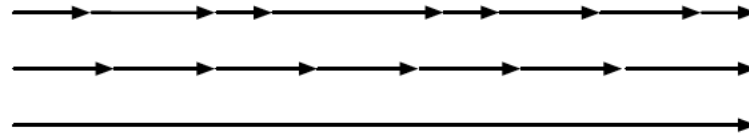
- Gradient can blow up

# Truncation

$$\partial h_{t-\tau} / \partial w_{\text{h}}$$

- Terminate the sequence
  - Fixed, Random
- *Short term* dependencies
  - A good thing !
- Truncation (regular or randomized):
  - Computational convenience
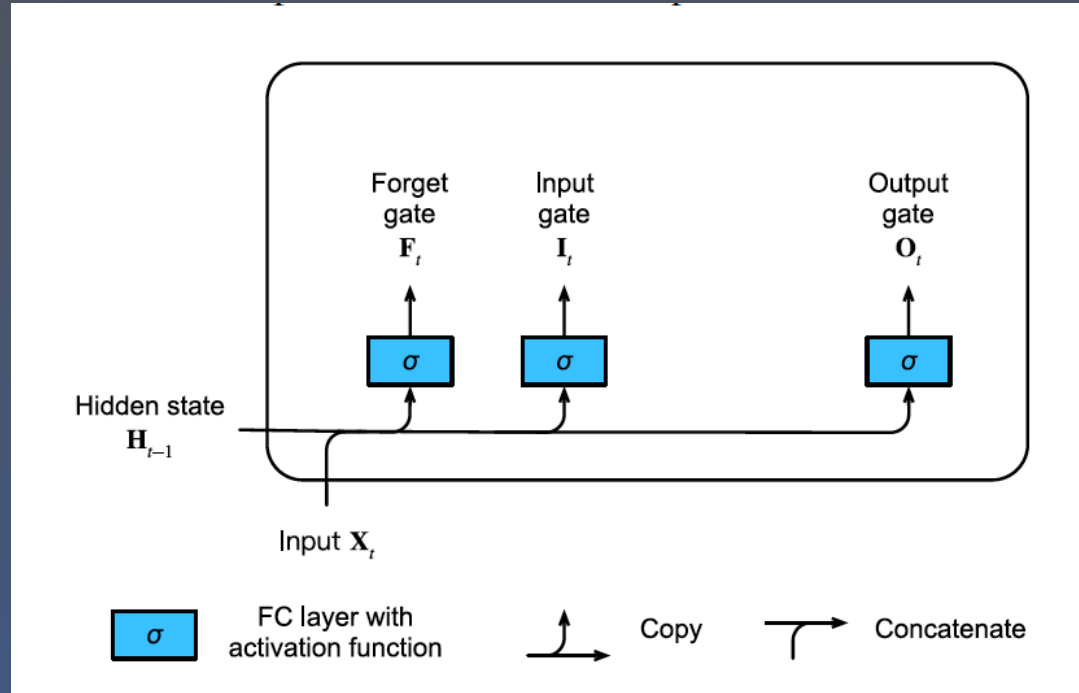  - Numerical stability

the time machine by h g wells

What about **vanishing** gradients ?

# Long Short Term Memory: **LSTM**

- Regular recurrent node **—> Memory Cell**
- **Gates**:
  - Input gate
  - Forget gate
  - Output gate

# Gates



$$I_t = \sigma(X_t W_{xi} + H_{t-1} W_{hi} + b_i),$$
$$F_t = \sigma(X_t W_{xf} + H_{t-1} W_{hf} + b_f),$$
$$O_t = \sigma(X_t W_{xo} + H_{t-1} W_{ho} + b_o),$$

# Memory Cell Update

$$\tilde{\mathbf{C}}_t = \tanh(\mathbf{X}_t \mathbf{W}_{xc} + \mathbf{H}_{t-1} \mathbf{W}_{hc} + \mathbf{b}_c),$$
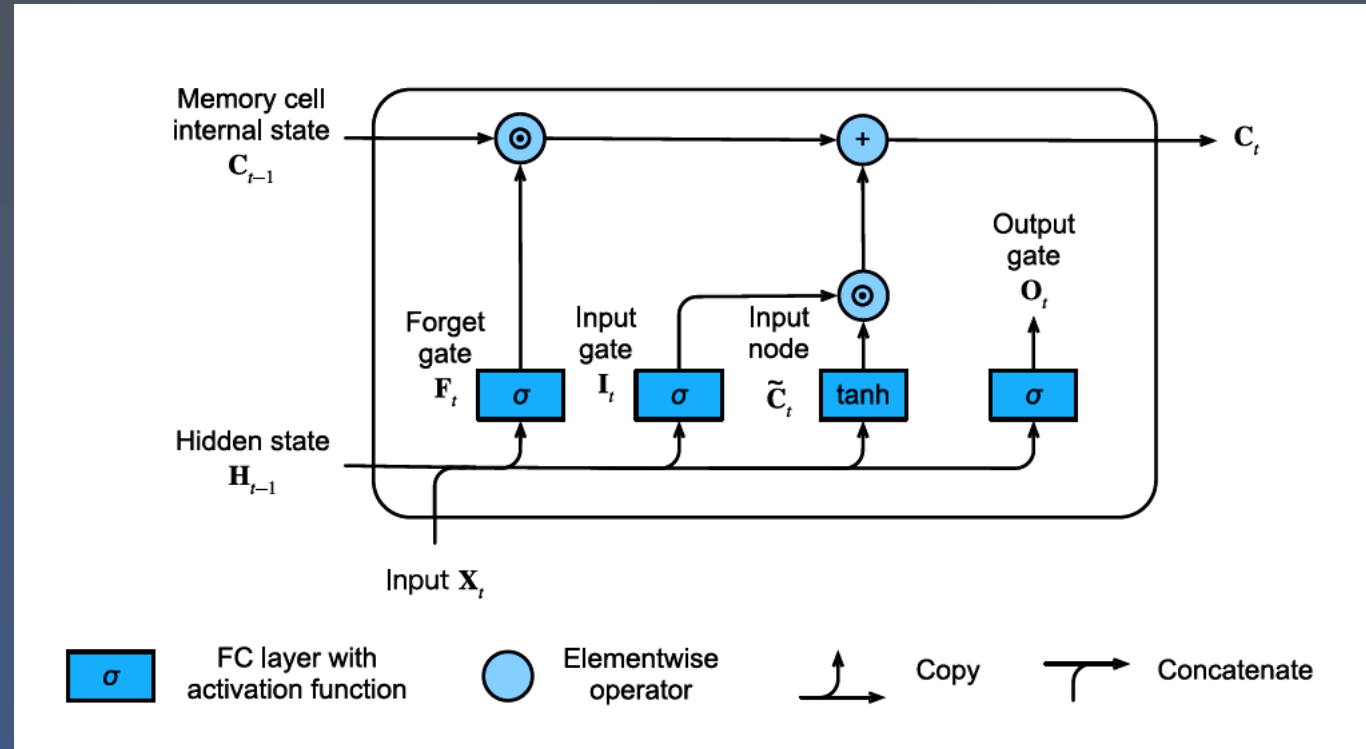
- Input node

$$\mathbf{C}_t = \mathbf{F}_t \odot \mathbf{C}_{t-1} + \mathbf{I}_t \odot \tilde{\mathbf{C}}_t.$$
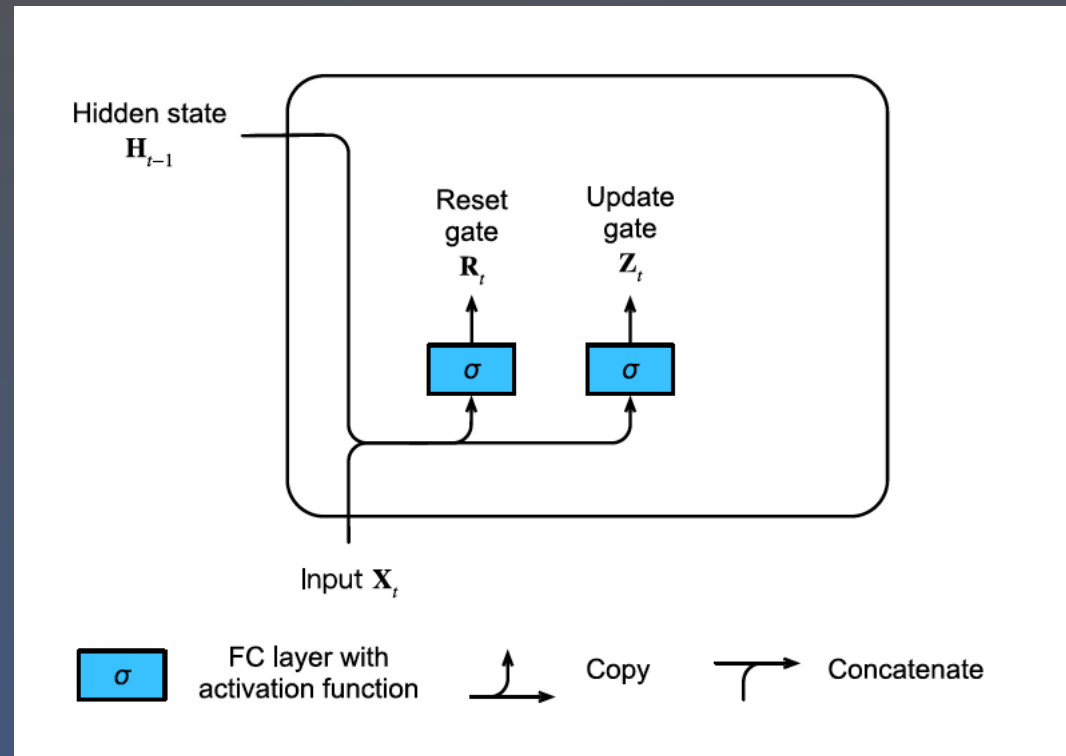
- Memory cell state, update

# Hidden State Update



$$\mathbf{H}_t = \mathbf{O}_t \odot \tanh(\mathbf{C}_t).$$

- Output gate: 1 vs 0 ?

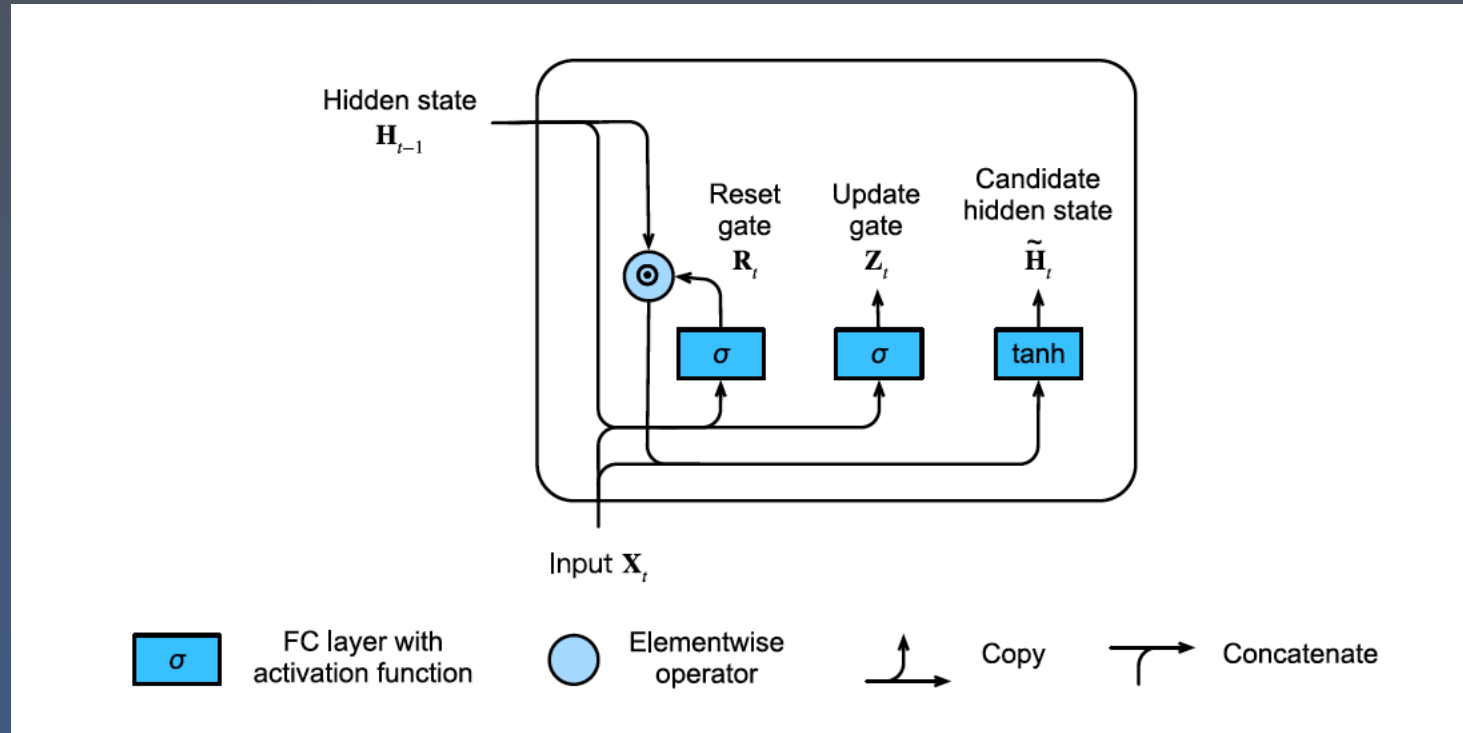# Gated Recurrence Unit GRU



$$\mathbf{R}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xr} + \mathbf{H}_{t-1} \mathbf{W}_{hr} + \mathbf{b}_r),$$
$$\mathbf{Z}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xz} + \mathbf{H}_{t-1} \mathbf{W}_{hz} + \mathbf{b}_z),$$
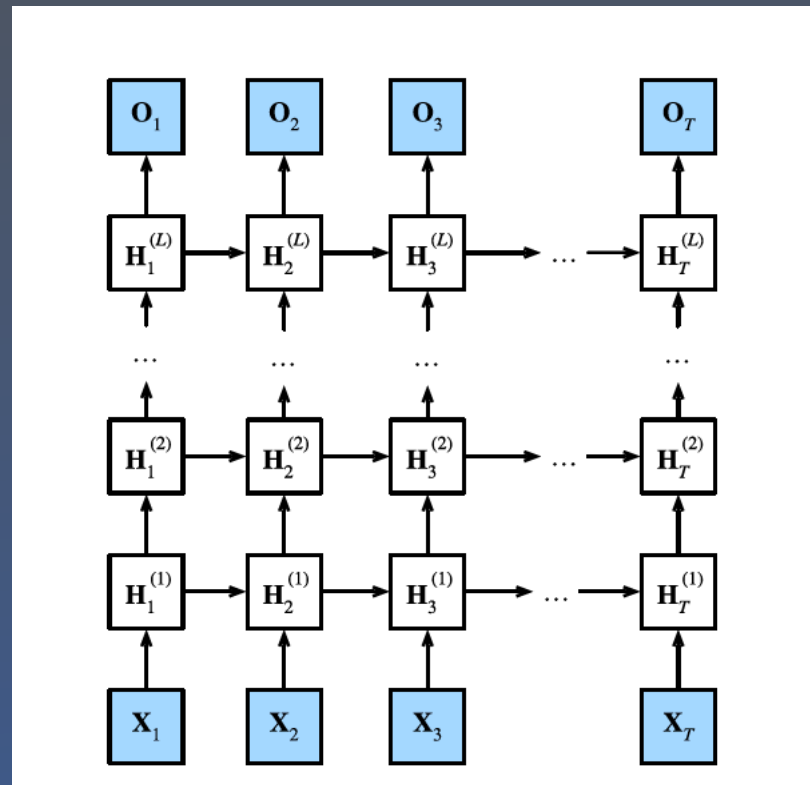
- Update gate : Input + Forget gates (of LSTM) combined
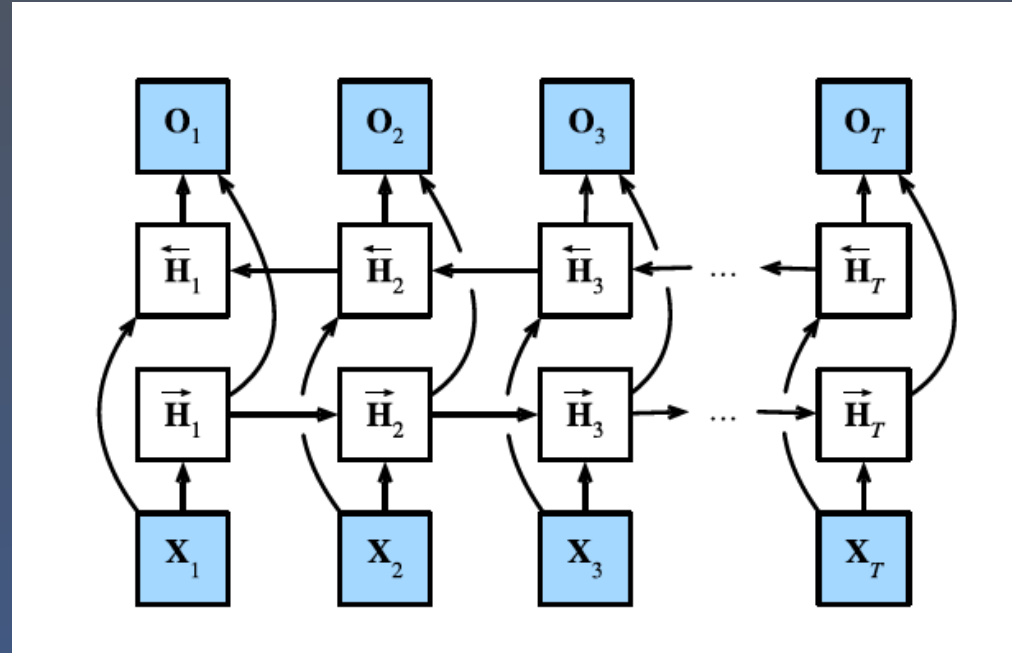- Reset gate

# GRU : Hidden State



$$\tilde{\mathbf{H}}_t = \tanh(\mathbf{X}_t \mathbf{W}_{\mathrm{xh}} + (\mathbf{R}_t \odot \mathbf{H}_{t-1}) \mathbf{W}_{\mathrm{hh}} + \mathbf{b}_{\mathrm{h}}),$$
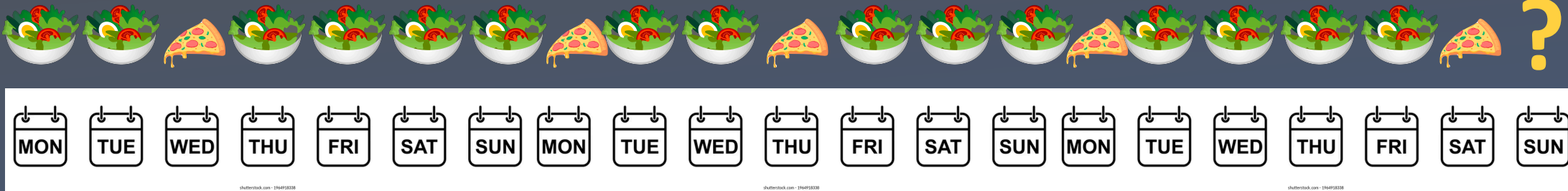
# Deep RNNs

# Bidirectional RNNs



- Go both forward and backward in the sequence

# Summary

The prediction, is based on understanding *"the psychology of the individual" (- Jeeves)*, which in turn is best facilitated by **Recurrent Neural Networks (RNNs)**. We maintain a detailed, day-by-day eating record of the individual, stored as a **Sequence** of time-stamped entries in **Hidden Units**. To keep things manageable, we periodically **Truncate** this sequence, either at fixed intervals or randomly. The individual's clear pattern of eating habits is stored in **Long Short-Term Memory (LSTM)**, where the **Input Gate** carefully preserves the meals from the last 2 days, and the **Forget Gate** ensures nothing beyond those two days is remembered. This recent history is used to predict today's lunch through the Output Gate. A **GRU** simplifies the process, reducing memory to "just the last two days." The **Bidirectional RNN** anticipates tomorrow's pizza, accounting for both past meals and future plans.