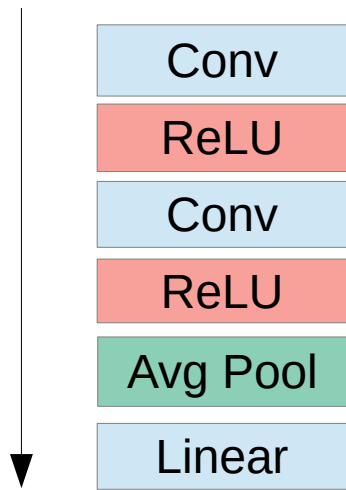


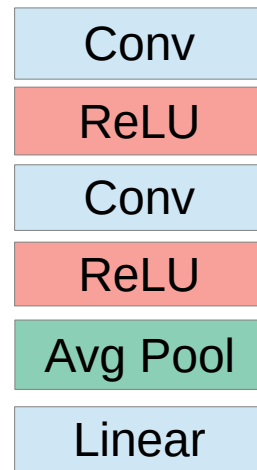
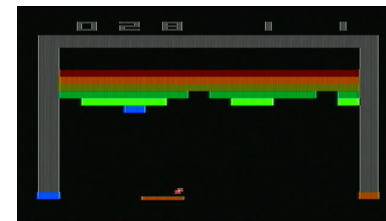


Sequence Models

Feed-Forward Networks

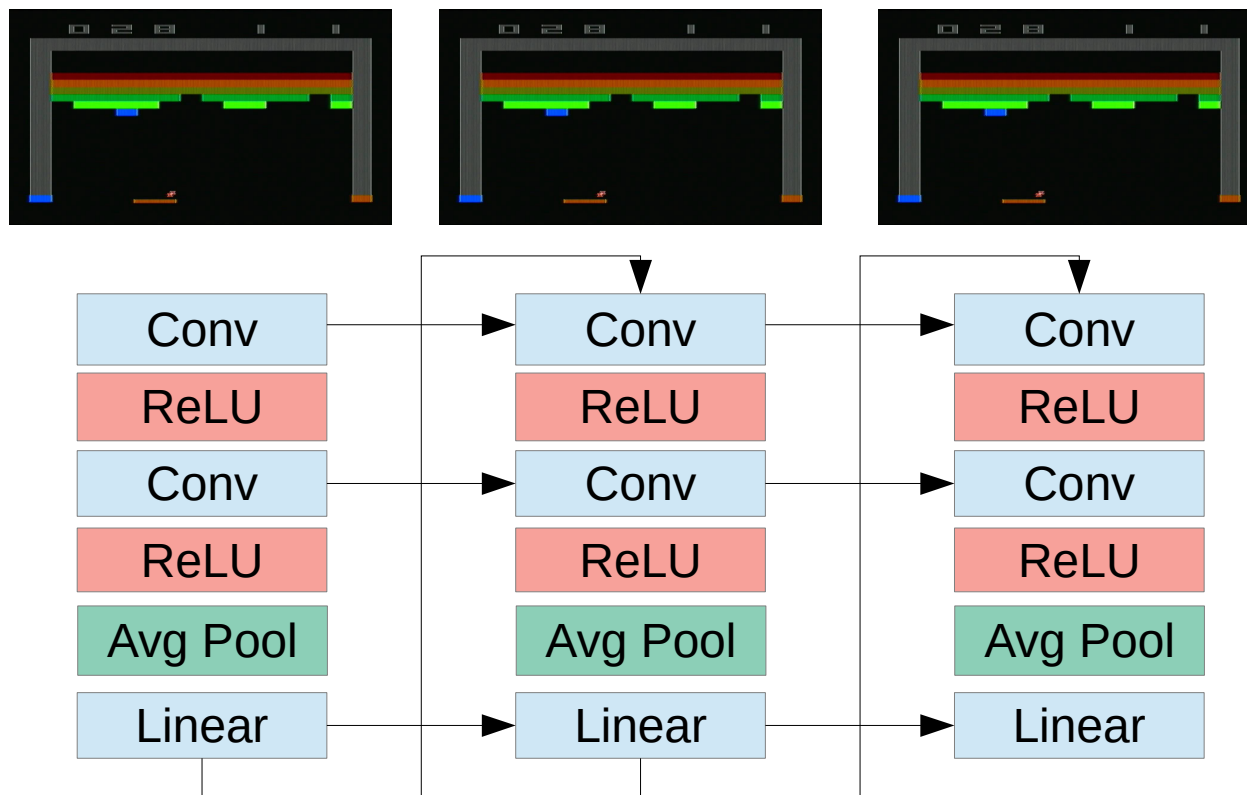


Information never flows back to previous layers

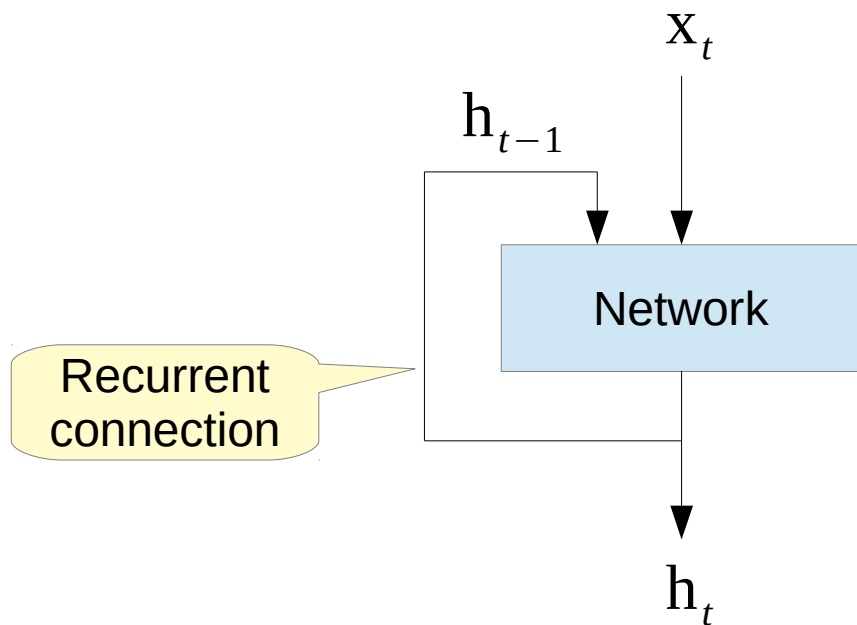


Acyclic computation graph

Memory



Recurrent Neural Network (RNN)



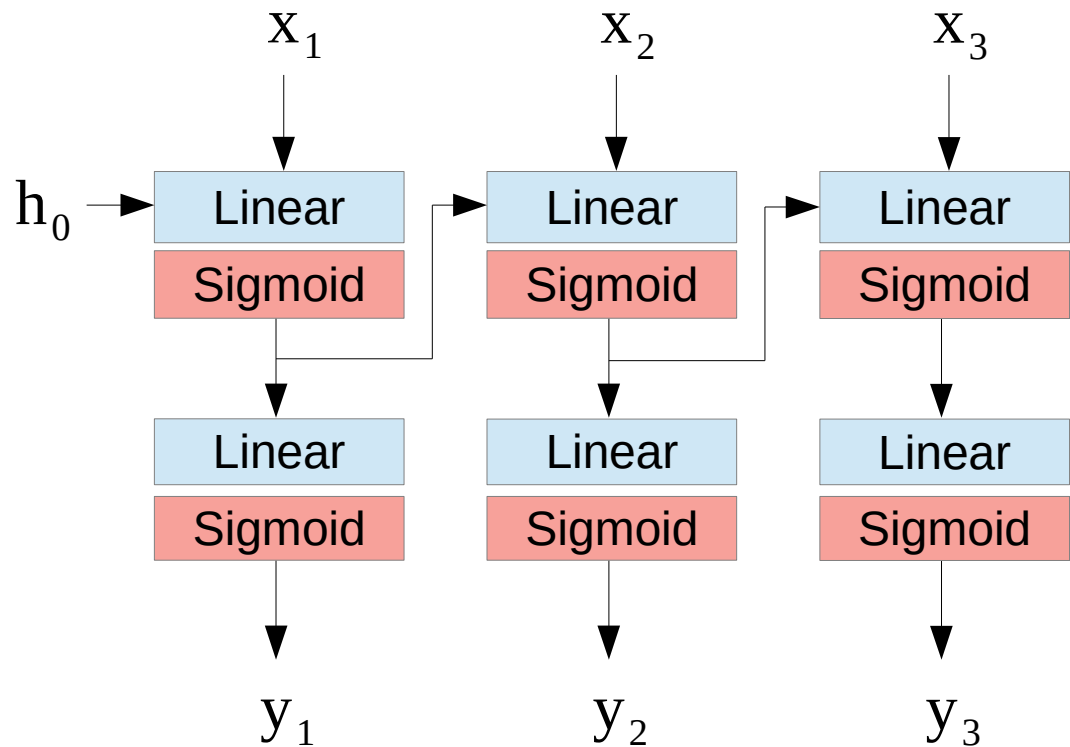
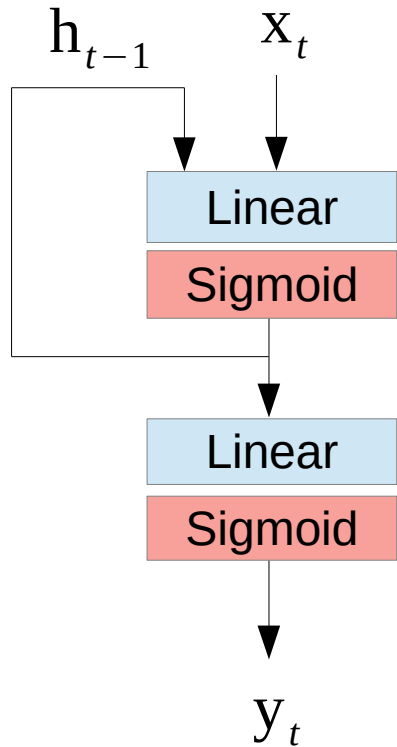
$$h_t = f_h(x_t, h_{t-1}, \theta_h)$$

Initial state: h_0

Zero

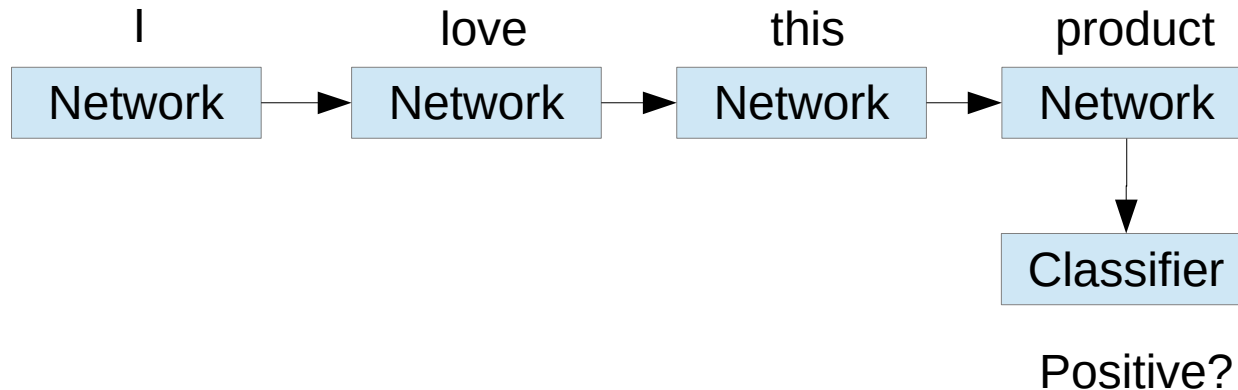
Learned

Unrolling RNNs

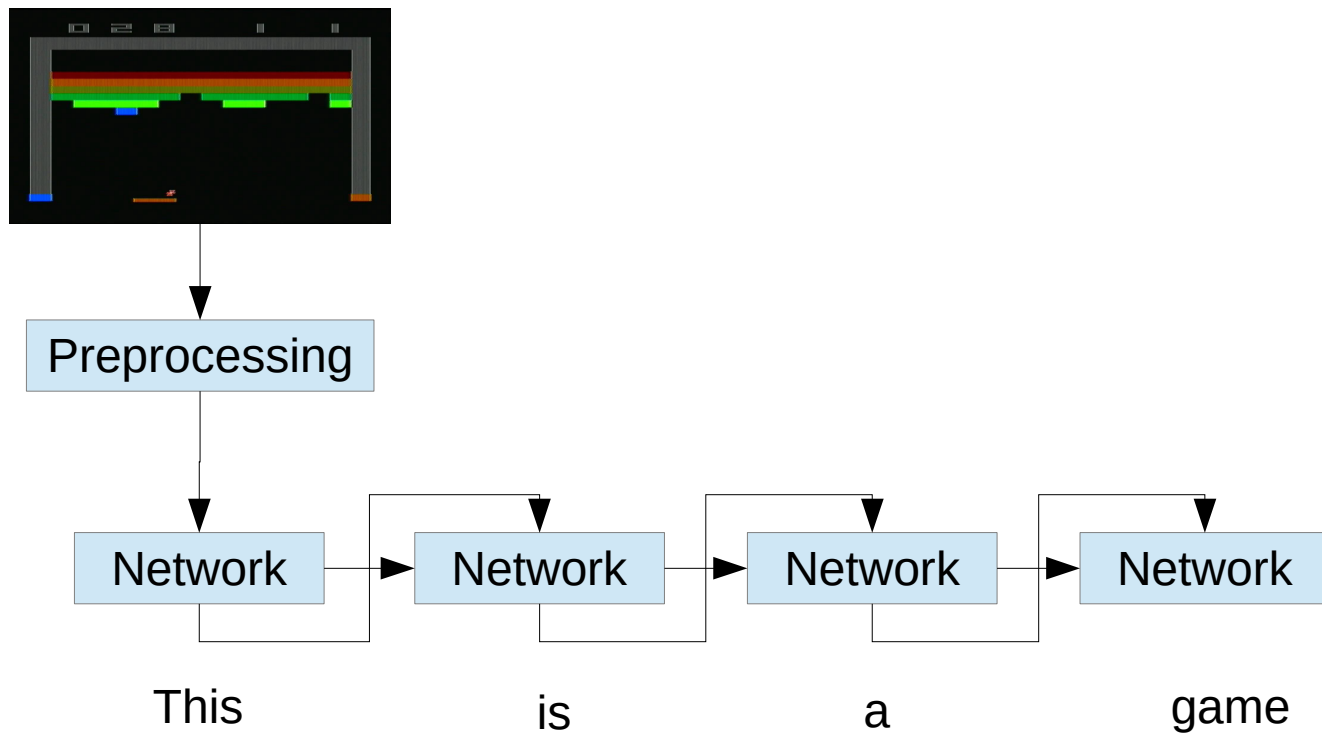


Applications

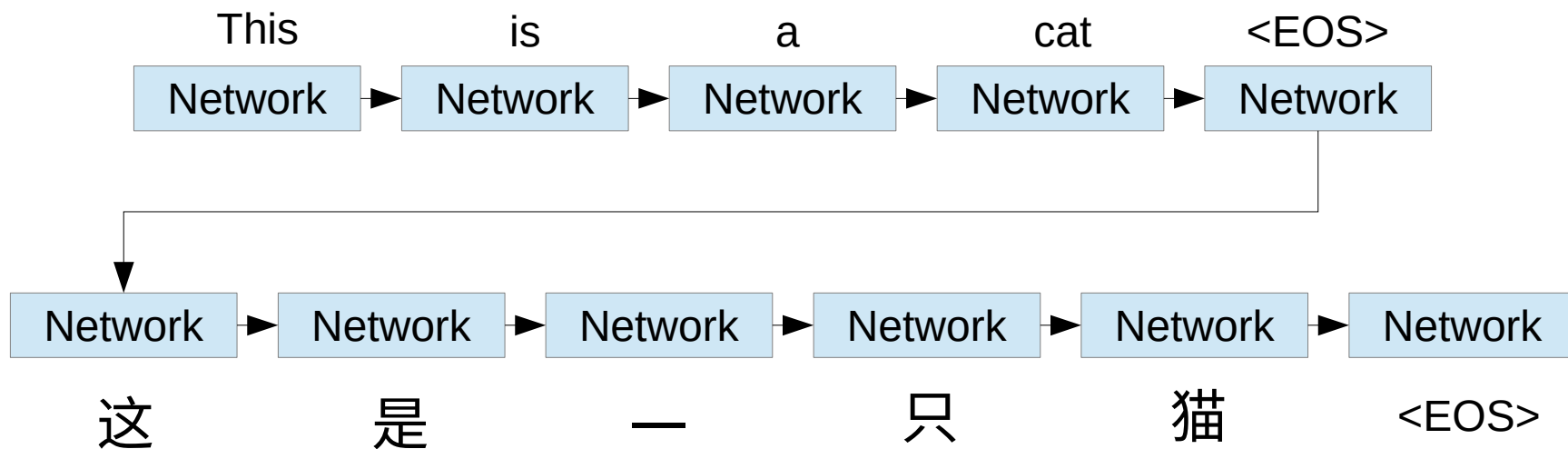
- Natural language processing
 - Either the input or output (or both) is a sequence of words
- Example: Is a review positive or negative?



Example: Language Generation

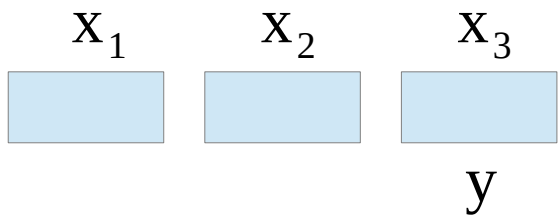


Example: Translation

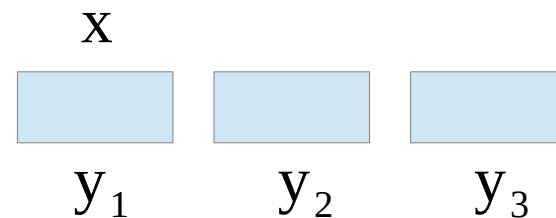


Kinds of RNN

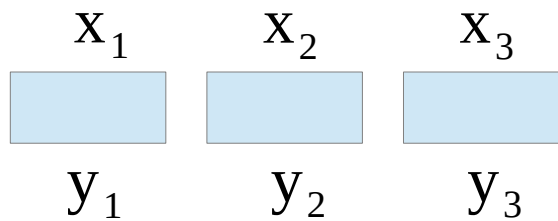
many-to-1



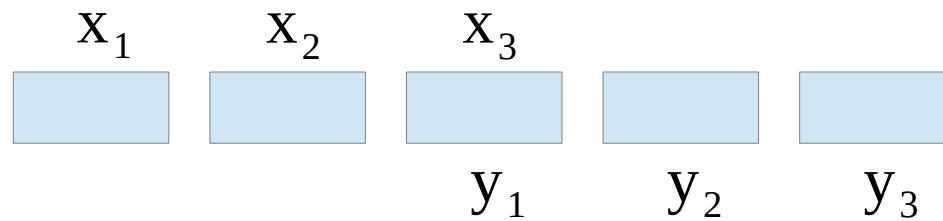
1-to-many



many-to-many

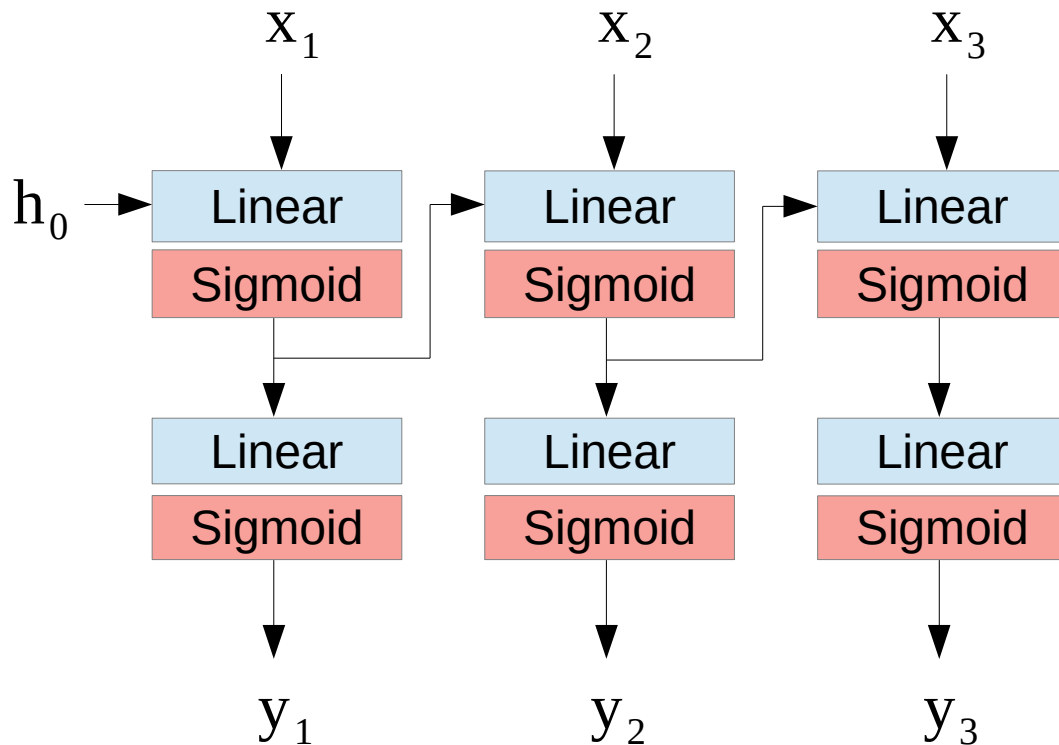


many-to-many

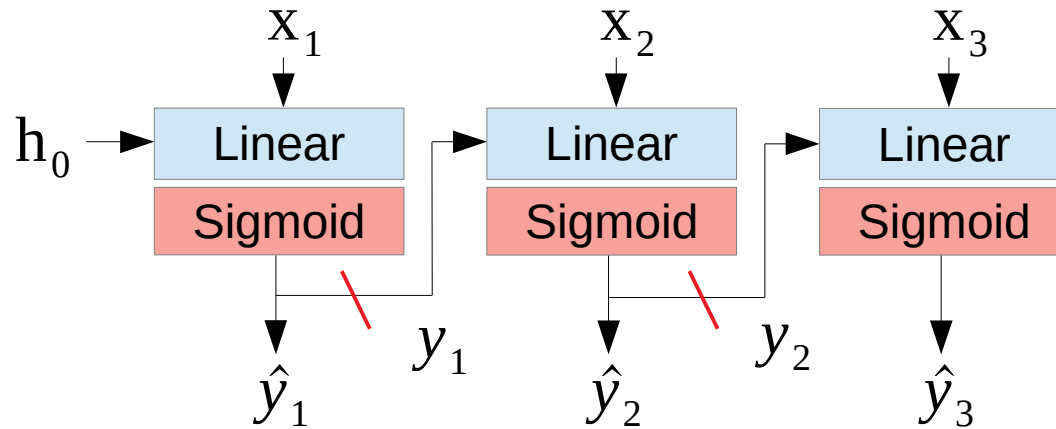


Training RNNs

- Unrolled RNNs have a feed-forward computation graph
- Regular backprop handles shared weights
- Long sequences leads to vanishing/exploding gradients



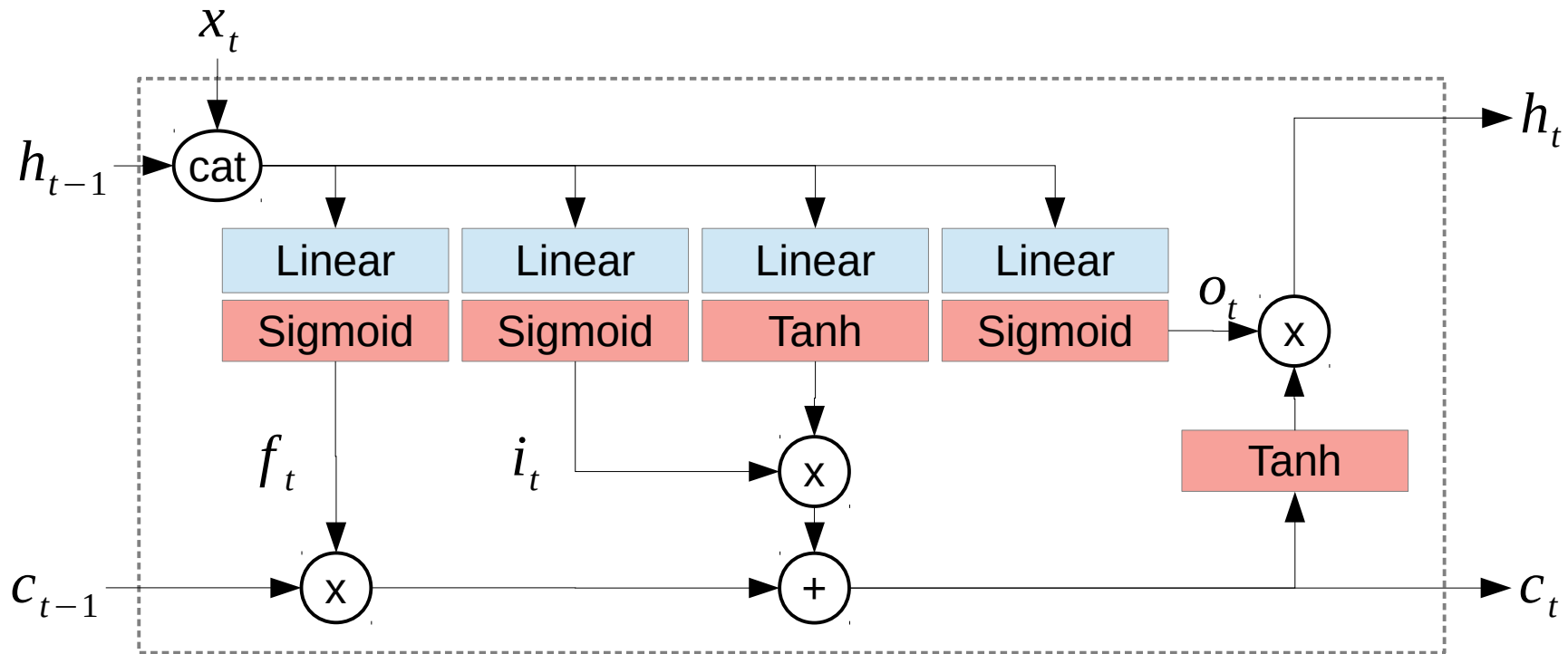
Vanishing/Exploding Gradients



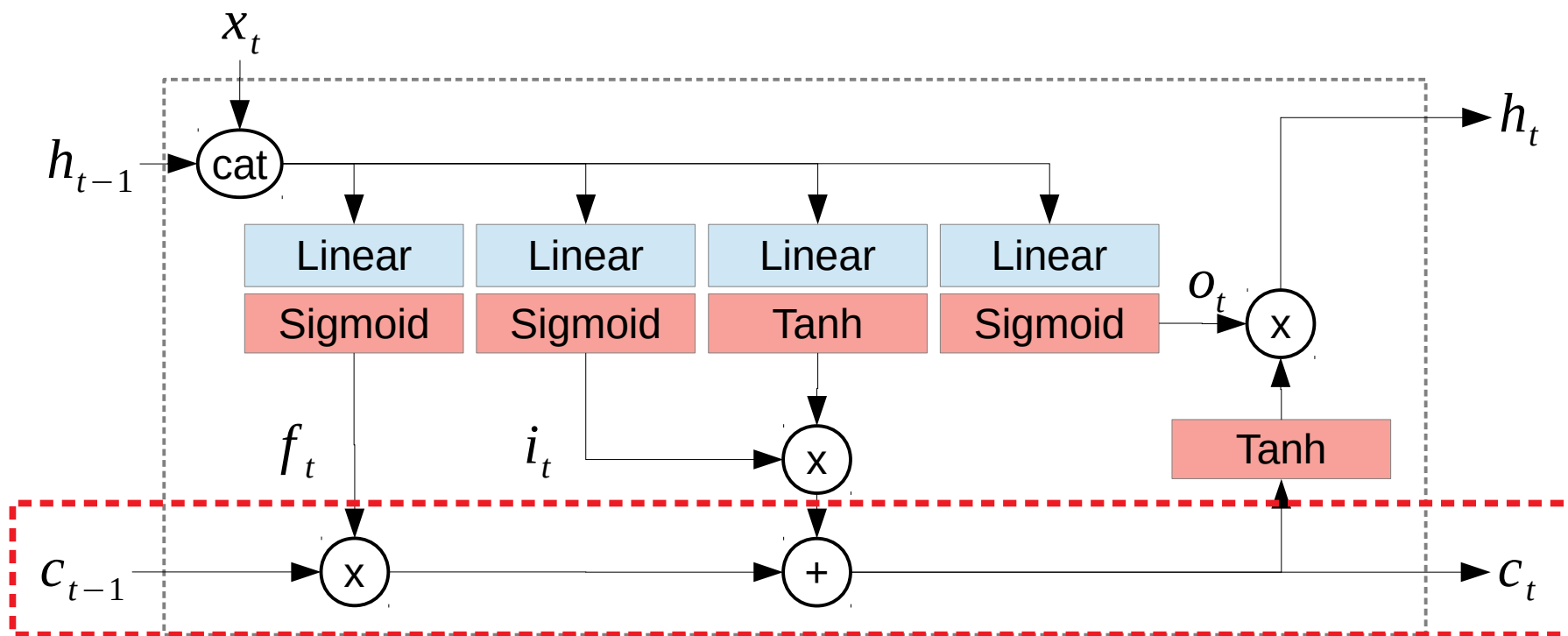
Exploding gradients:
$$\nabla L(\theta) = \min\left(1, \frac{\epsilon}{|\nabla L(\theta)|}\right) \nabla L(\theta)$$

Vanishing gradients: must rely on network structure

Long Short-Term Memory (LSTM)



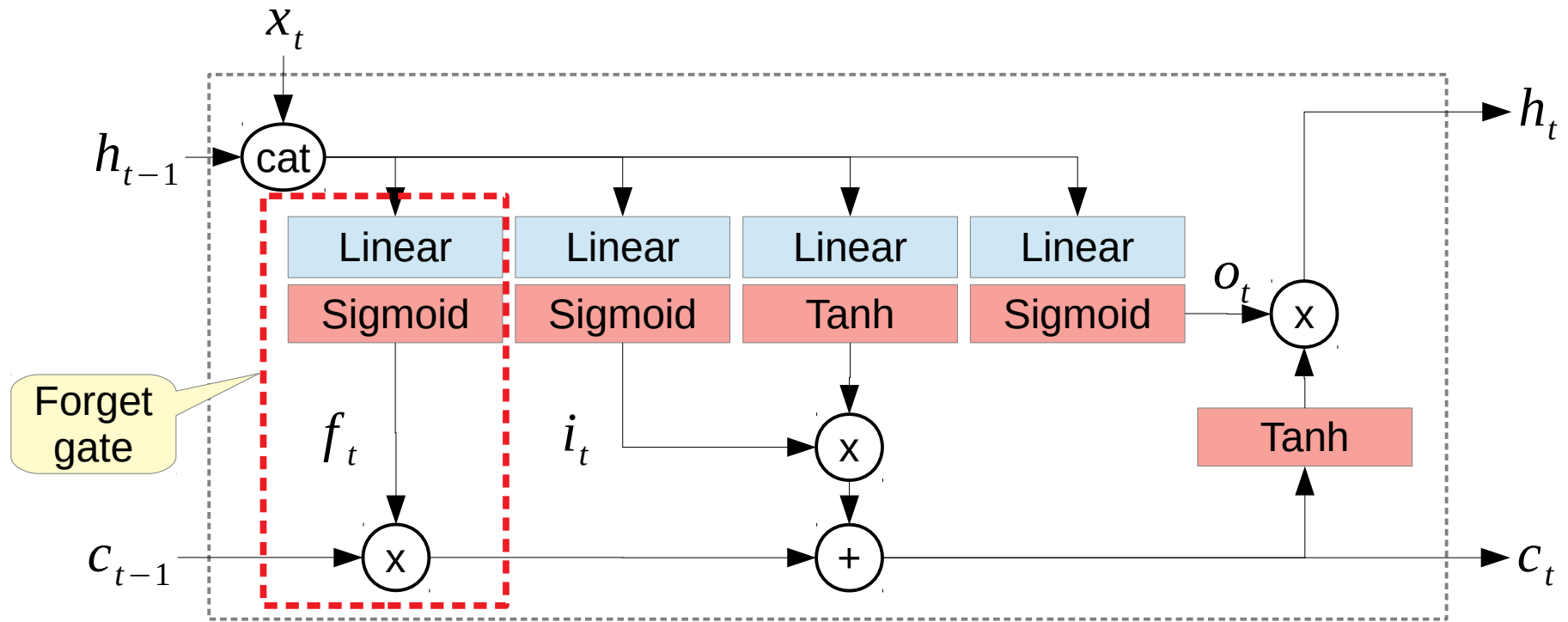
Long Short-Term Memory (LSTM)



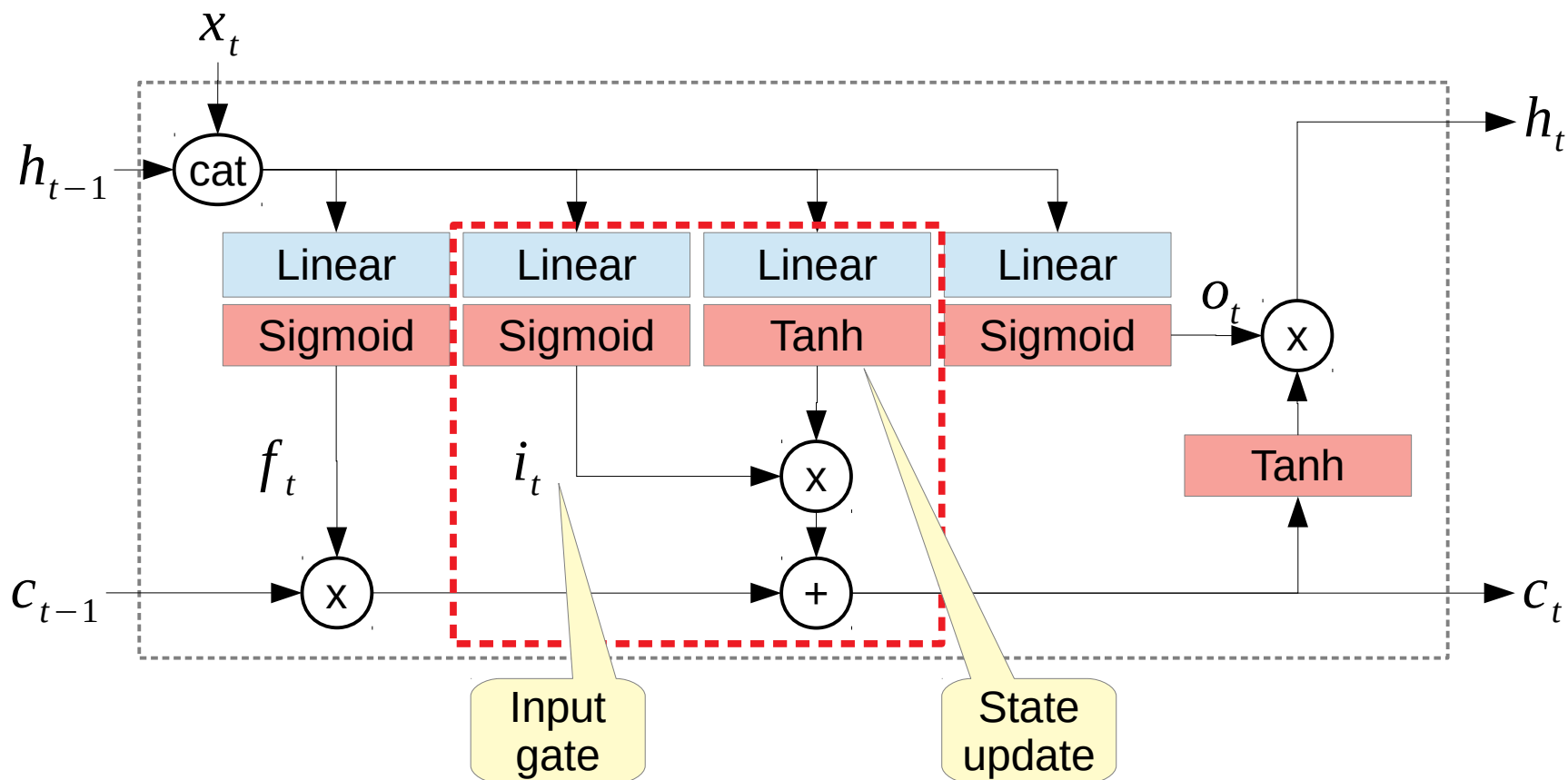
Long-term memory

Similar to residual connections

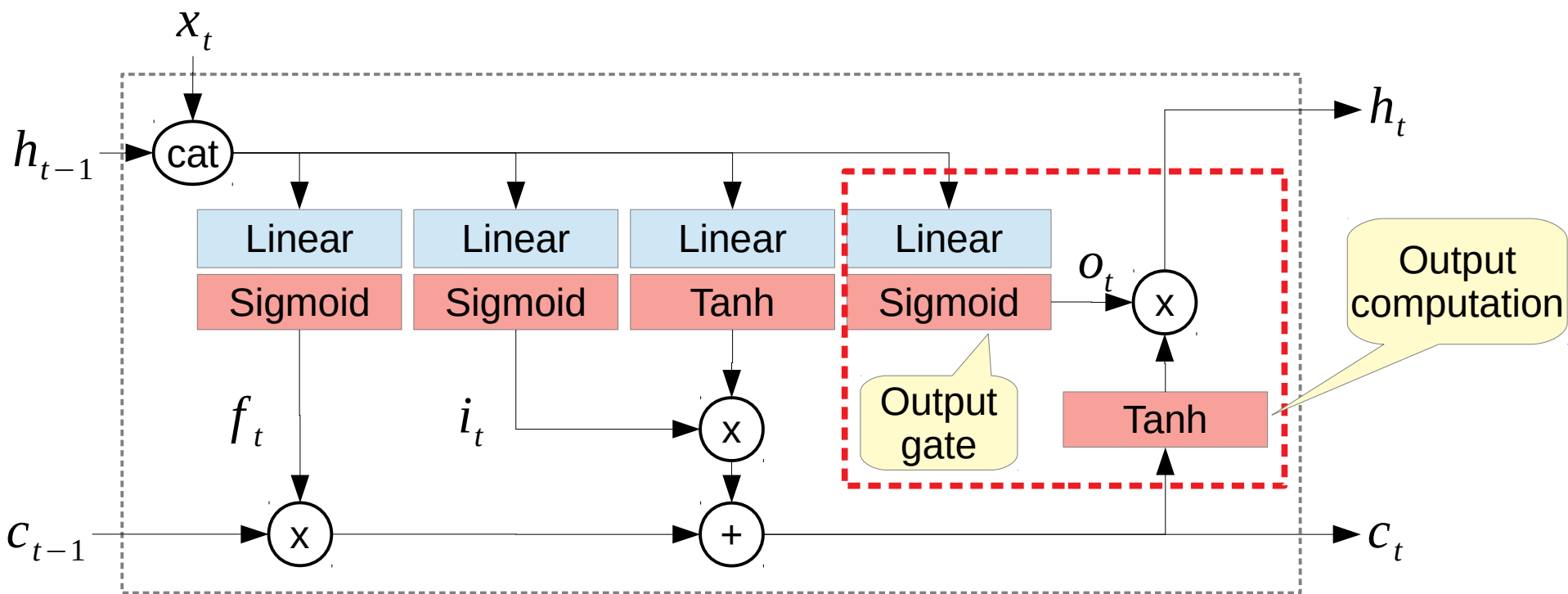
Long Short-Term Memory (LSTM)



Long Short-Term Memory (LSTM)

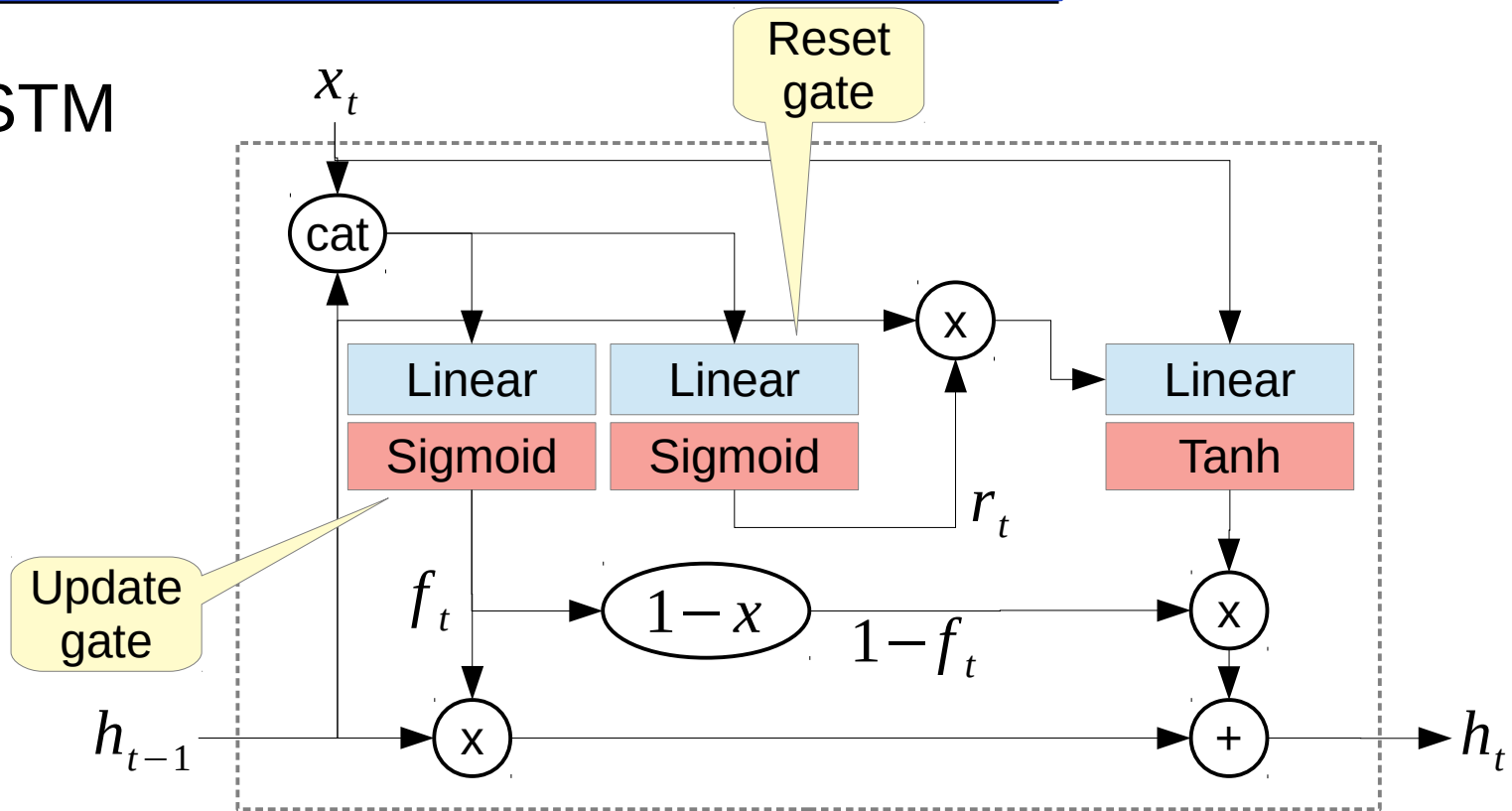


Long Short-Term Memory (LSTM)



Gated Recurrent Unit (GRU)

- Simplified LSTM
- Single state
- Fewer gates
- Similar performance



LSTM/GRU Networks

