



Northeastern

CS 7140: ADVANCED MACHINE LEARNING

Langevin Monte Carlo

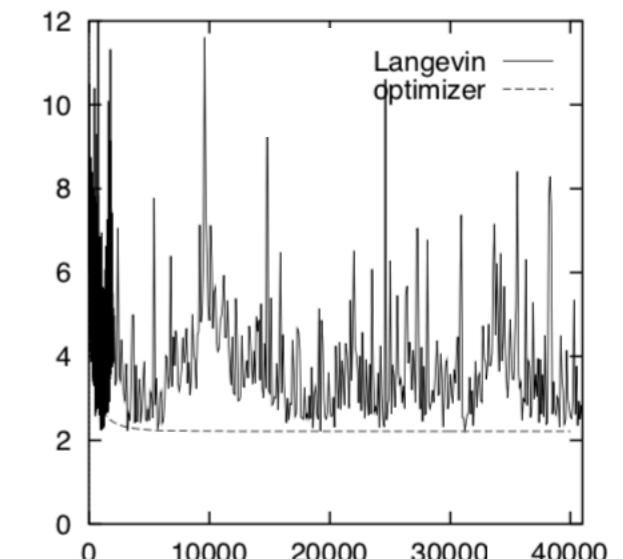
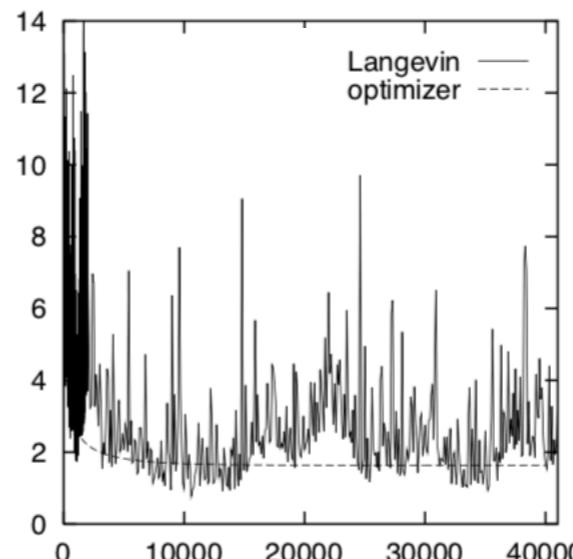
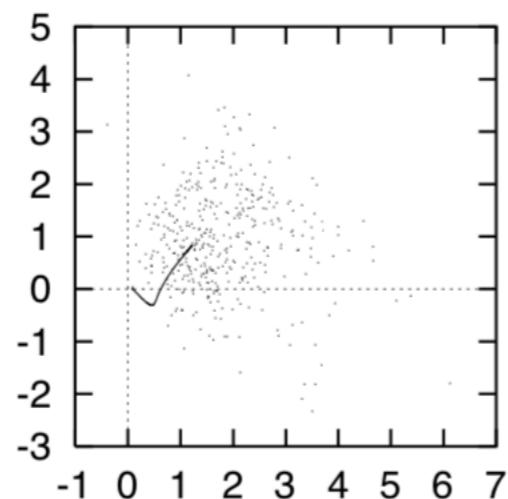
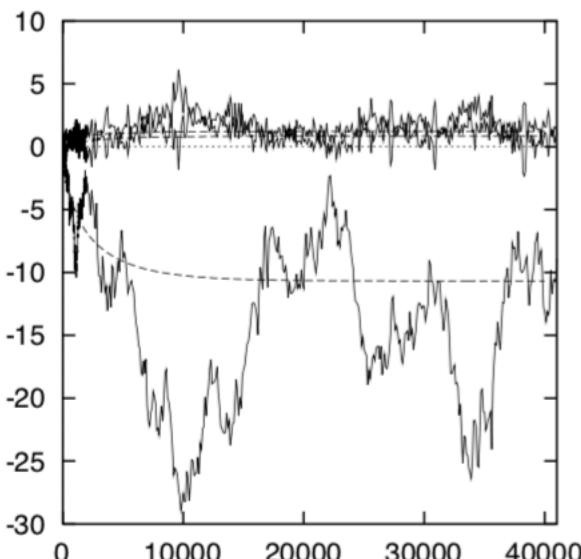
- Monte Carlo: evaluate the integral of a distribution

$$\mathbb{E}[f(w)] \approx \frac{1}{R} \sum_r f(w^{(r)})$$

- Langevin: gradient descent with added noise $p \sim \mathcal{N}(0,1)$

$\Delta w = \frac{1}{2}\epsilon^2 g + \epsilon p$ accept or reject depending on $J(w)$, g is the gradient

- ϵ controls the step size, too *large*, moves may be rejected; too *small*, progress around the state space is slow



Gaussian Approximation

- Taylor expansion $J(w) \approx J(w_0) + \frac{1}{2}(w - w_0)^\top A(w - w_0) + \dots$

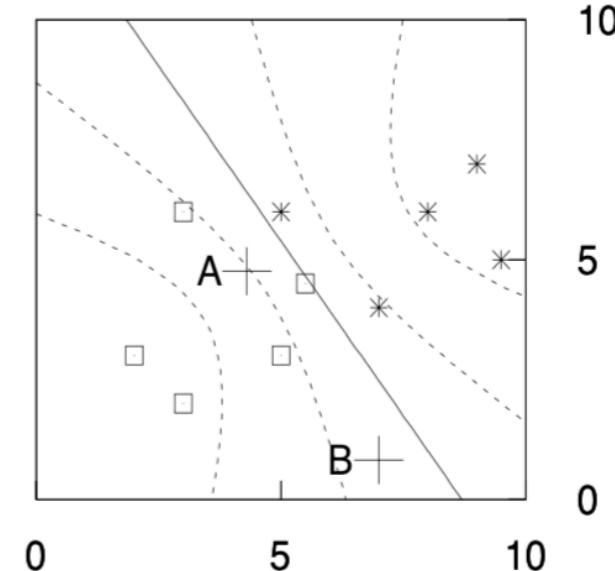
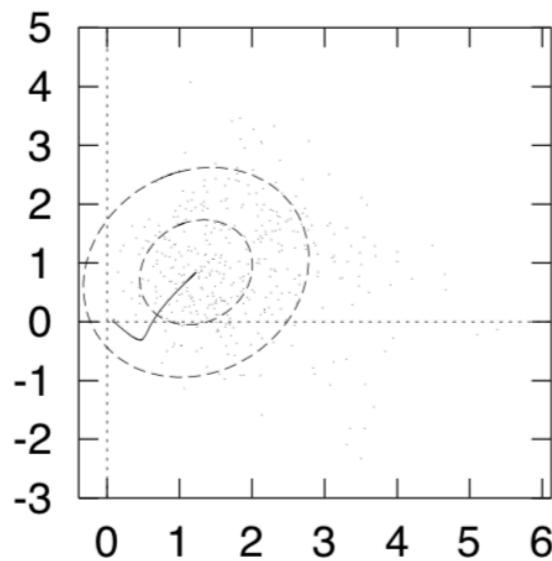
where $A_{ij} \equiv \frac{\partial^2}{\partial w_i \partial w_j} J(w)$ is the *Hessian* matrix

- Approximate the posterior with a Gaussian

$$Q(w; w_0, A) = [\det A / 2\pi]^{1/2} \exp \left[-\frac{1}{2}(w - w_0)^\top A(w - w_0) \right]$$

- Compute the integral using Gaussian posterior

$$\Delta w = w - w_0 \quad P(w | D, \alpha) \approx (1/Z_Q) \exp(-\frac{1}{2}\Delta w^\top A \Delta w)$$



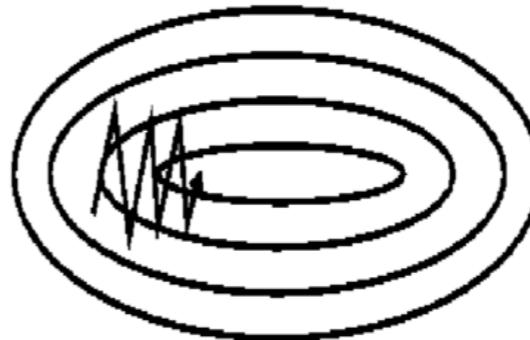
Accelerating SGD

- Momentum

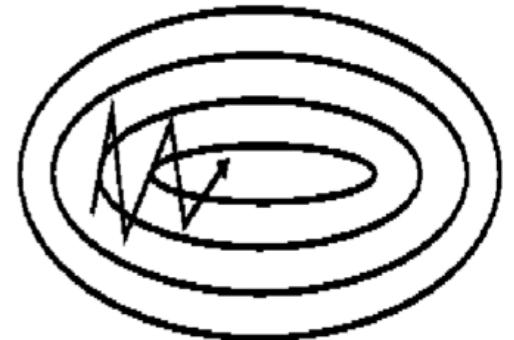
physical interpretation: update velocity

$$w = w - \alpha \tilde{\nabla}_w \mathcal{L}$$

$$\begin{aligned} v_t &= \alpha v_{t-1} - \epsilon \tilde{\nabla}_w L(w_t) \\ w_{t+1} &= w_t + v_t \end{aligned}$$



SGD without momentum



SGD with momentum

- Nesterov Momentum

stronger theoretical guarantee

$$\begin{aligned} v_t &= \alpha v_{t-1} - \epsilon \tilde{\nabla}_w L(w_t + \alpha v_{t-1}) \\ w_{t+1} &= w_t + v_t \end{aligned}$$

Adaptive Learning Rate

Cost is sensitive to learning rate only in some directions in parameter space

Use a separate learning rate for each parameter

- AdaGrad

gradient

approximate Hessian

update

- RMSProp

gradient

approximate Hessian

update

$$g = \tilde{\nabla}_w L(w_t)$$

$$r_t = r_{t-1} + g \odot g$$

$$w_{t+1} = w_t - \frac{\epsilon}{\delta + \sqrt{r}} \odot g$$

$$g = \tilde{\nabla}_w L(w_t)$$

$$r_t = \rho r_{t-1} + (1 - \rho) g \odot g$$

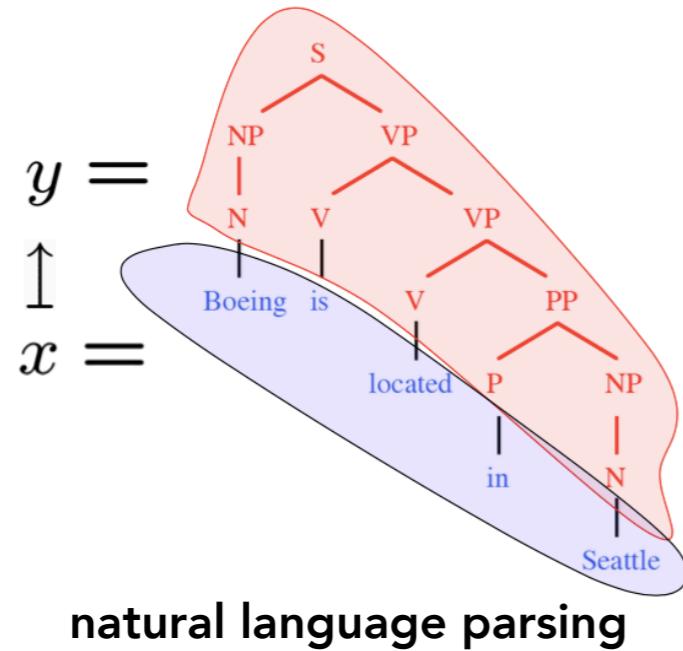
$$w_{t+1} = w_t - \frac{\epsilon}{\sqrt{\delta + r}} \odot g$$

STRUCTURED PREDICTION

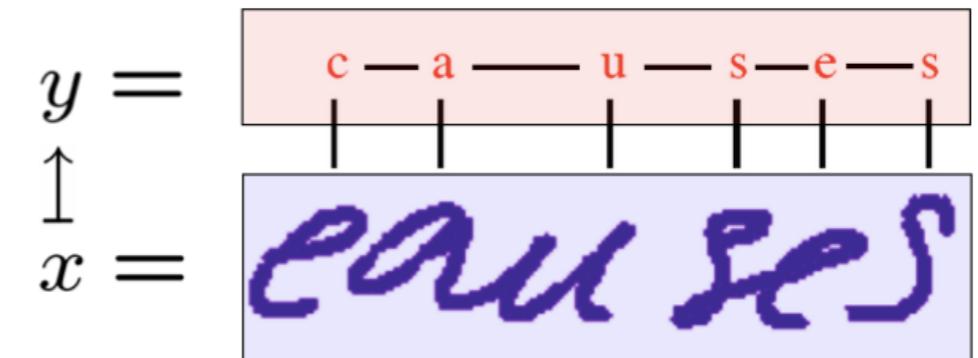
Structured Prediction

- A prediction function $f: \mathcal{X} \longrightarrow \mathcal{Y}$, where the output \mathcal{Y} domain is structured
 $y^* = f(x) = \operatorname{argmax}_y g(x, y)$ MAP inference
- $g(x, y) = p(y | x)$ when the model is probabilistic
- Many MAP inference problems are NP Hard:
combinatorial number of possibilities

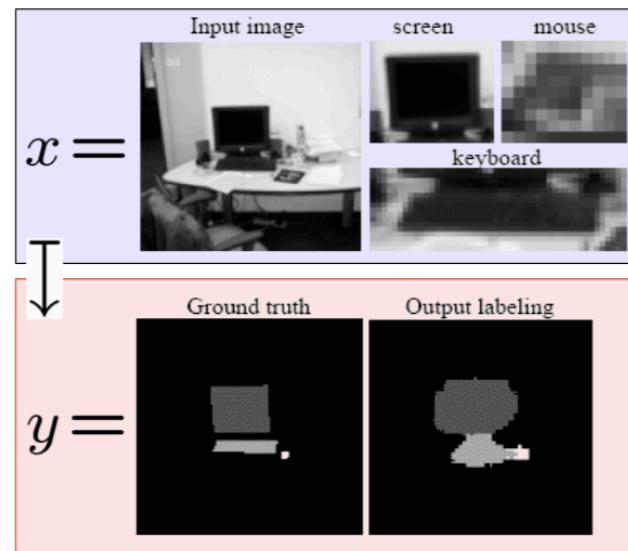
Examples Applications



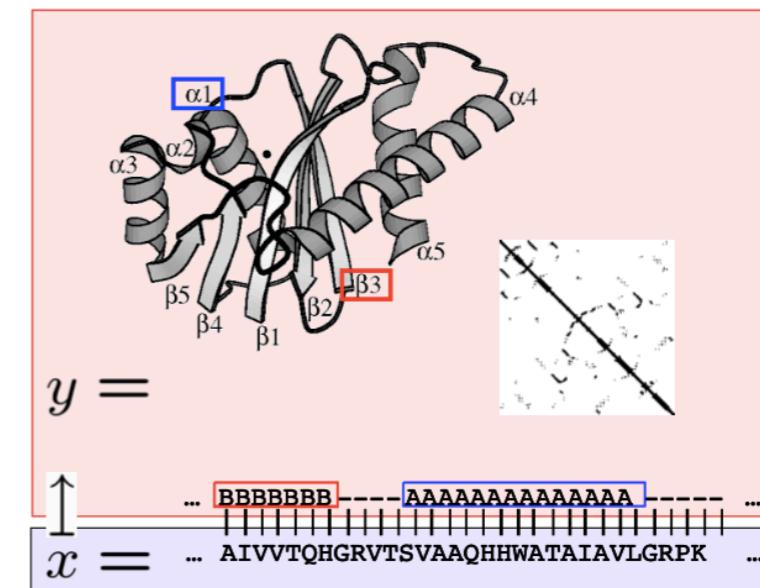
natural language parsing



hand-written text recognition



visual scene understanding



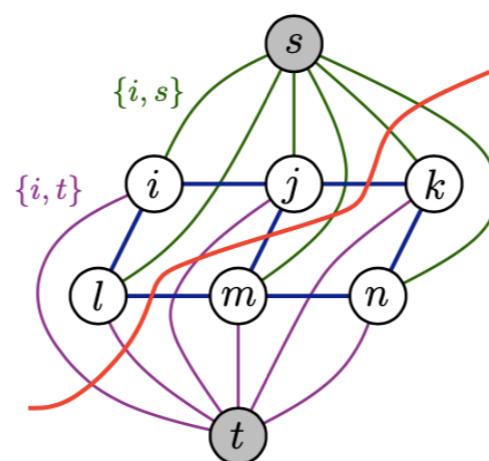
protein structure prediction

Graph Cuts

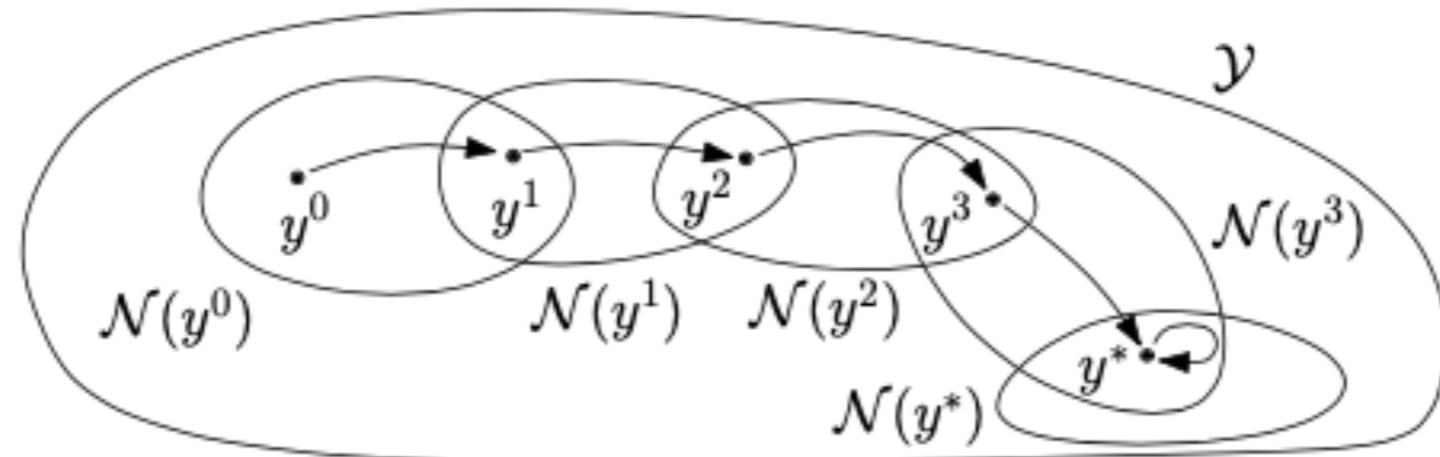
- Image foreground segmentation:
 - foreground 1, background 0, $C(x_i, x_j)$ intensity contrast
 - $$g(x, y; w) = \sum_{i \in V} \log p(y_i | x_i) + w \sum_{(i,j) \in E} C(x_i, x_j)[y_i \neq y_j]$$
 - find a global maximizer with graph cuts method

- Giving up Optimality

- computational hard
- model uncertainty



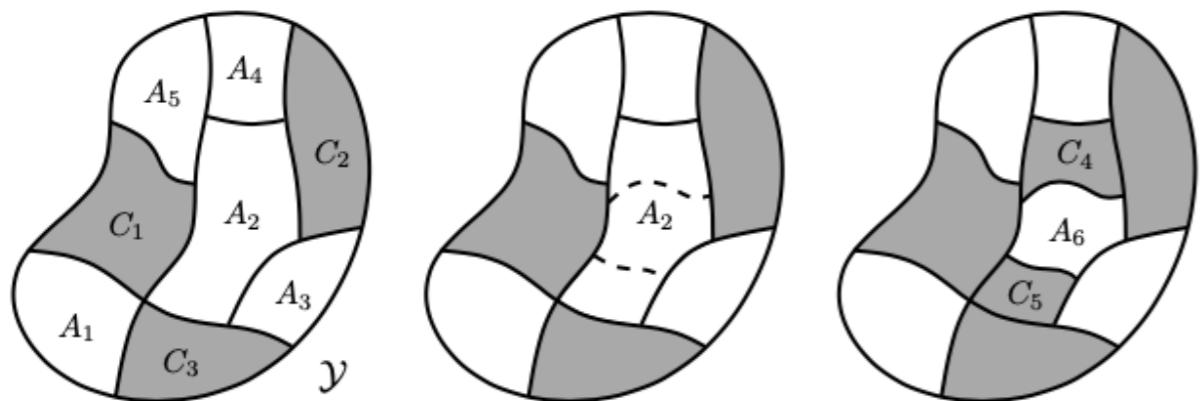
Local Search



- Improve the solution by iteratively solving $\operatorname{argmax}_{y \in N_s(y^t)} g(x, y)$ within the neighborhood of current solution
- Terminate when no further improvement can be made
- Optimality cannot be guaranteed. Approximation algorithms offer worst-case a priori guarantees.

Branch and Bound

- Worst Case: Implicit enumeration of \mathcal{Y} and is guaranteed to find the optimal solution
- Divide-and-conquer algorithm
 - Partition y into active and closed nodes.
 - Closed nodes do not contain a better solution
 - Move nodes from active to closed



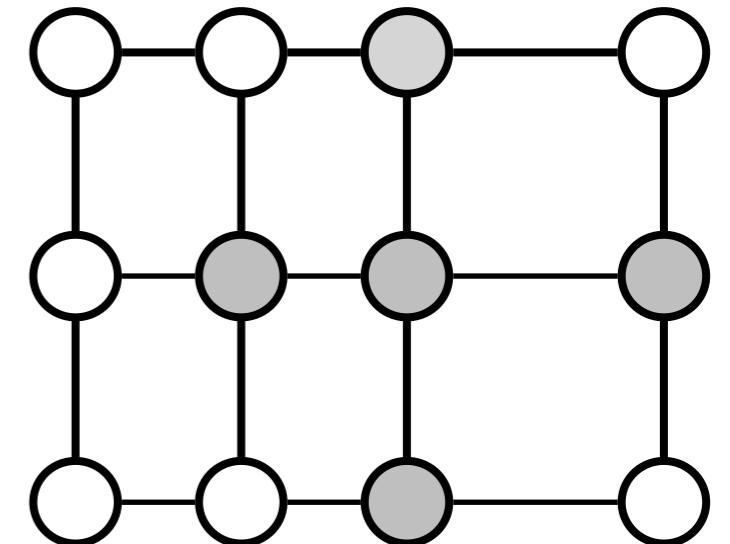
Learning to Plan



- Learning the cost function with the intended property, e.g. preferring roads and avoiding vegetation
- Update the cost function as search continues
- Use A* for prediction and loss-augmented prediction

Gradient-Based Optimization

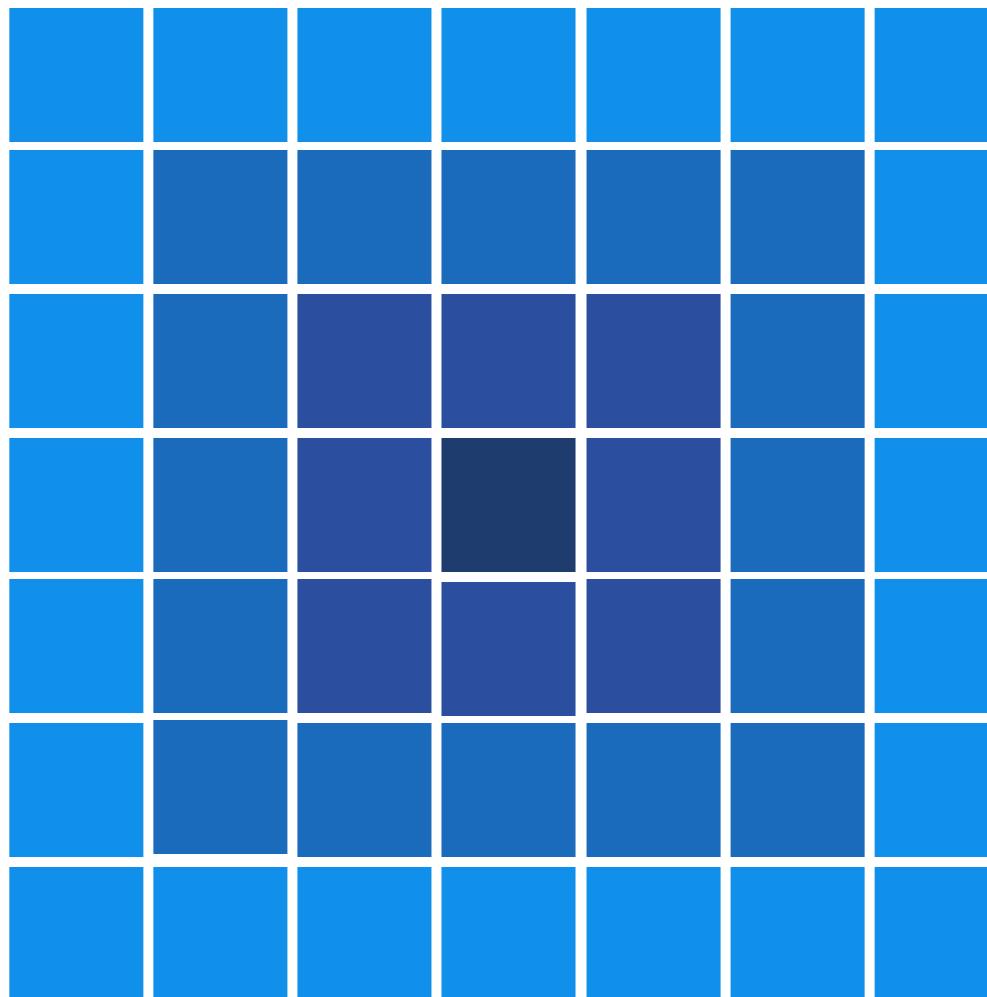
- Lagrangian Relaxation
 - Objective $J(x)$, equality constraint $x = c$
 - Lagrangian: $J(x) + \lambda(x - c)$
- MAP-MRF:
 - $E(y; \theta, \mu) = \sum_{m,n} \theta_{mn} \mu(y_m, y_n) + \sum_n \theta_n \mu(y_n)$
 - constraints: $\sum_n \mu_n(y_n) = 1, \sum_n \mu_{m,n}(y_m, y_n) = \mu_m(y_m)$



DEEP STRUCTURED PREDICTION

Spatial Structure

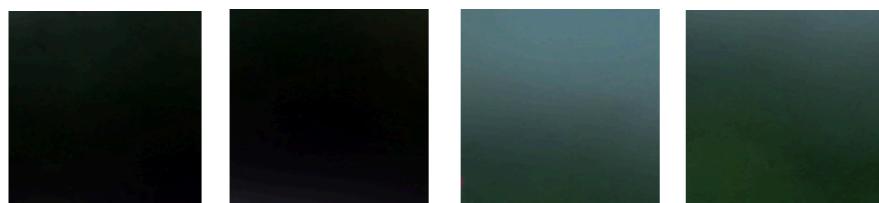
- to reduce the amount of learning, we can exploit the ***spatial structure*** of image data



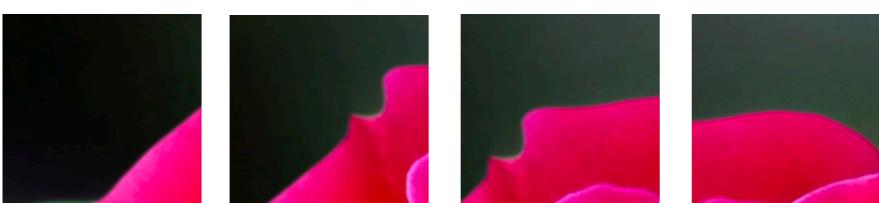
Locality

nearby areas tend to contain stronger patterns

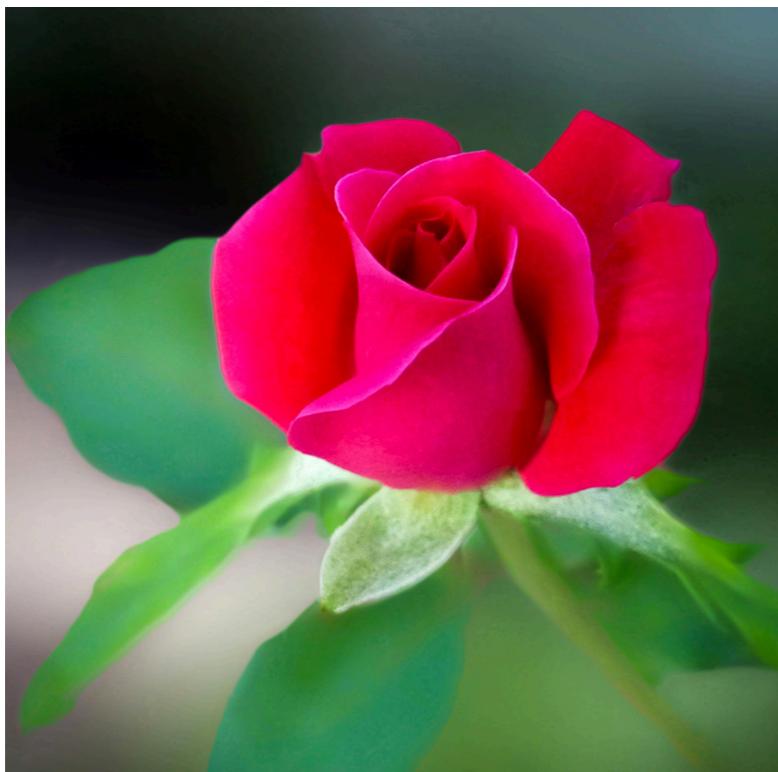
nearby **pixels** tend to be similar and vary
in particular ways



nearby **patches** tend to share characteristics
and are combined in particular ways

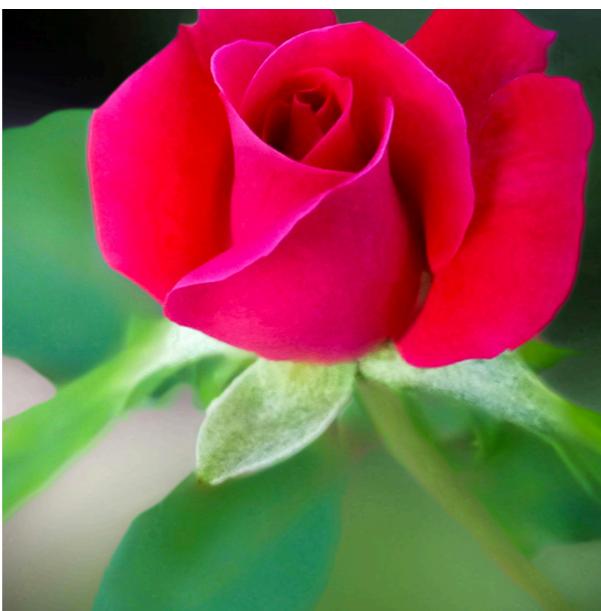
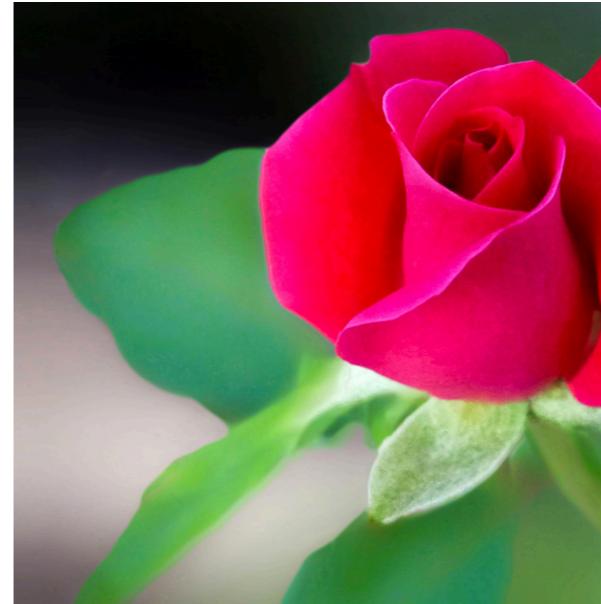


nearby **regions** tend to be found
in particular arrangements



Translation Invariance

relative (rather than absolute) positions are relevant



Rose is independent of absolute location of its pixels

Inductive Biases

locality

*nearby areas tend
to contain stronger
patterns*

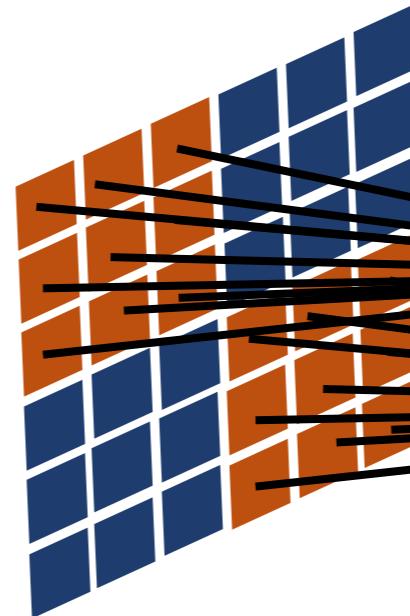


inputs can be
restricted to regions

maintain spatial ordering

translation invariance

*relative positions
are relevant*



same filters can be applied
throughout the input

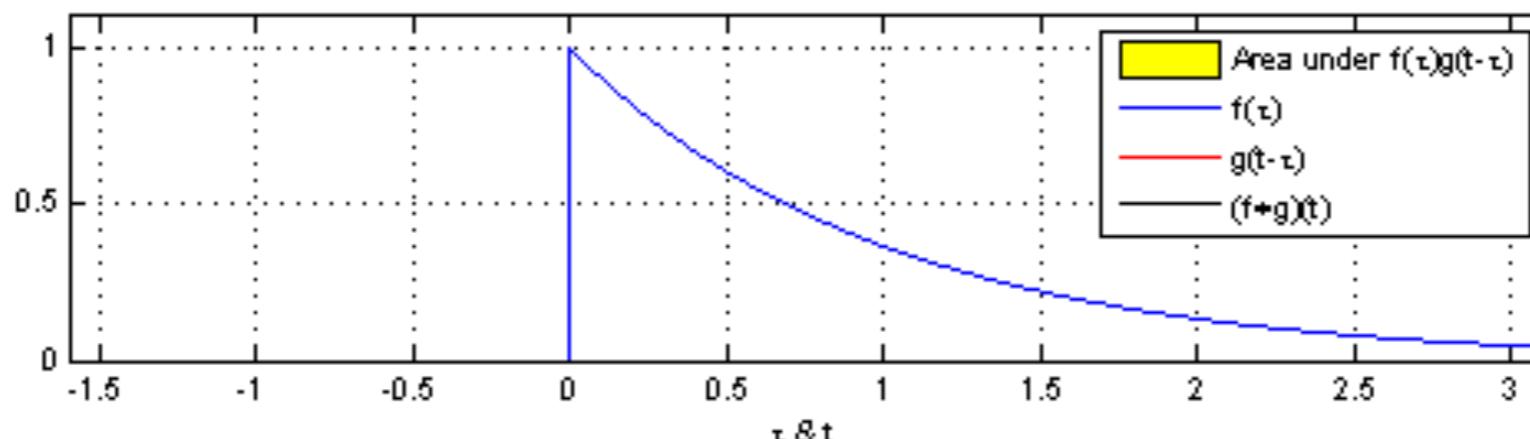
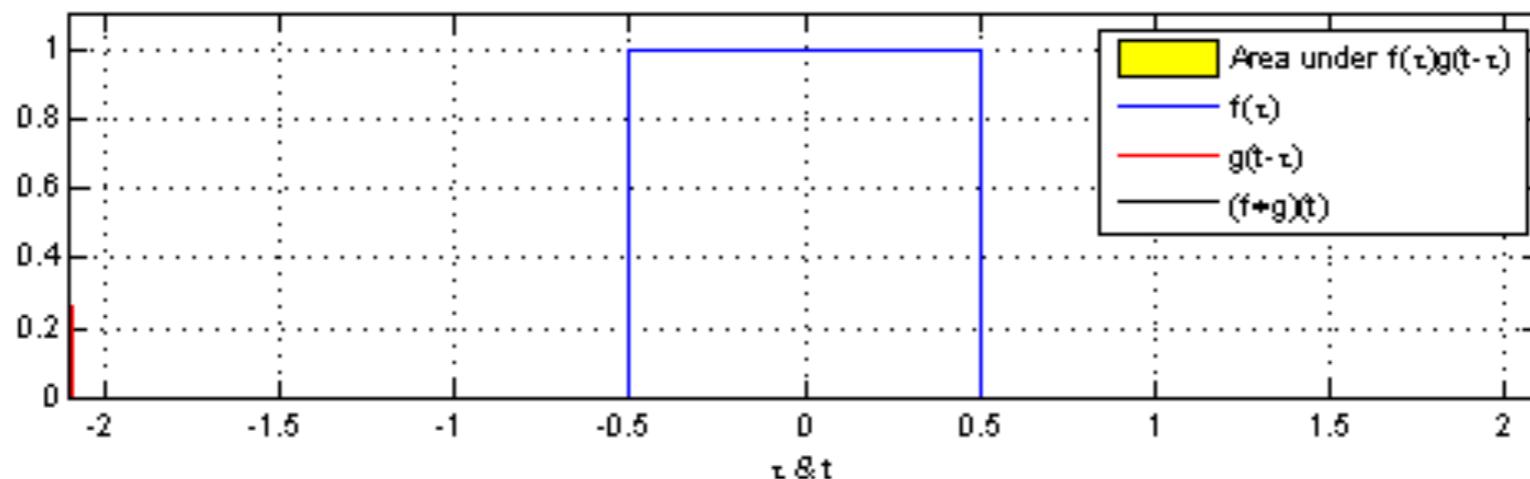
same weights

Convolution

- A integral: the amount of overlap of one function as it is shifted over another function

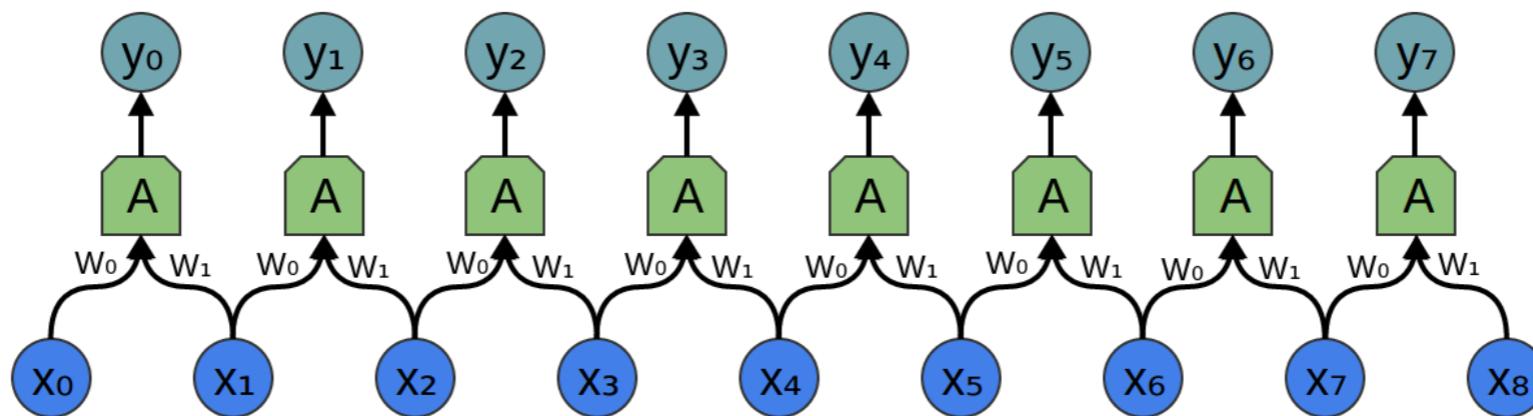
$$(f * g)(t) = \int_{-\infty}^{\infty} f(\tau)g(t - \tau)d\tau = (g * f)(t)$$

filter



Convolution

- 1D discrete
$$(f * g)(t) = \sum_i f(i)g(t - i)$$
- 2D discrete
$$(f * g)(x, y) = \sum_i \sum_j f(i, j)g(x - i, y - j)$$
- 2D cross correlation
$$(f * g)(x, y) = \sum_i \sum_j f(i, j)g(x + i, y + j)$$



allow weights sharing, reduce parameter size

Convolutional Filter

filter weights: $\begin{pmatrix} 0 & 1 & 2 \\ 2 & 2 & 0 \\ 0 & 1 & 2 \end{pmatrix}$

3 ₀	3 ₁	2 ₂	1	0
0 ₂	0 ₂	1 ₀	3	1
3 ₀	1 ₁	2 ₂	2	3
2	0	0	2	2
2	0	0	0	1

12	12	17
10	17	19
9	6	14

$$(f * g)(x, y) = \sum_i \sum_j f(i, j)g(x - i, y - j)$$

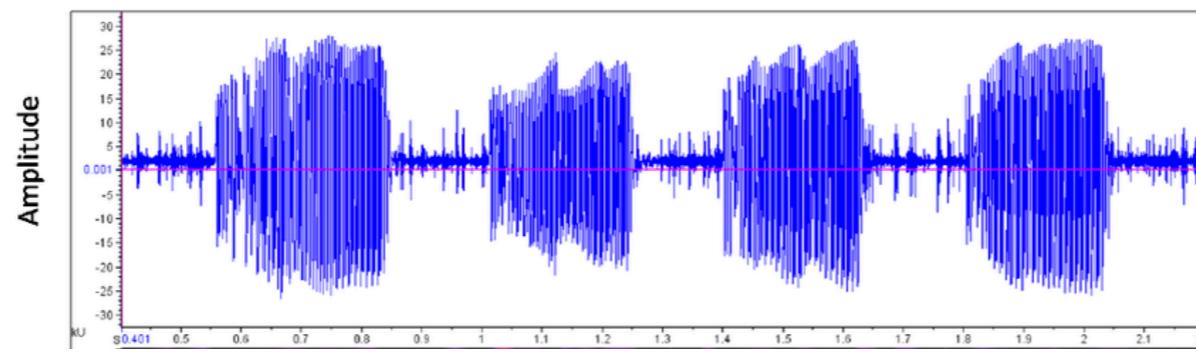
take inner (dot) product of filter and each input location

measures degree of filter feature at input location

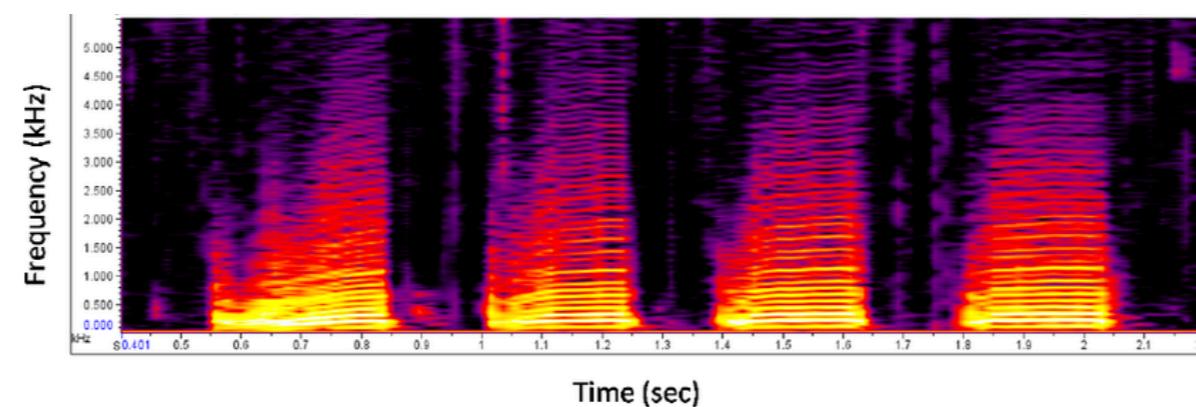
→ *feature map*

Convolution Applications

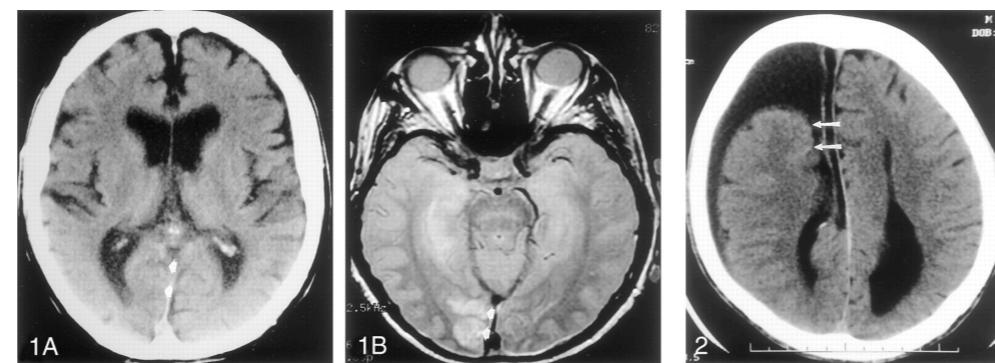
- 1 D Conv



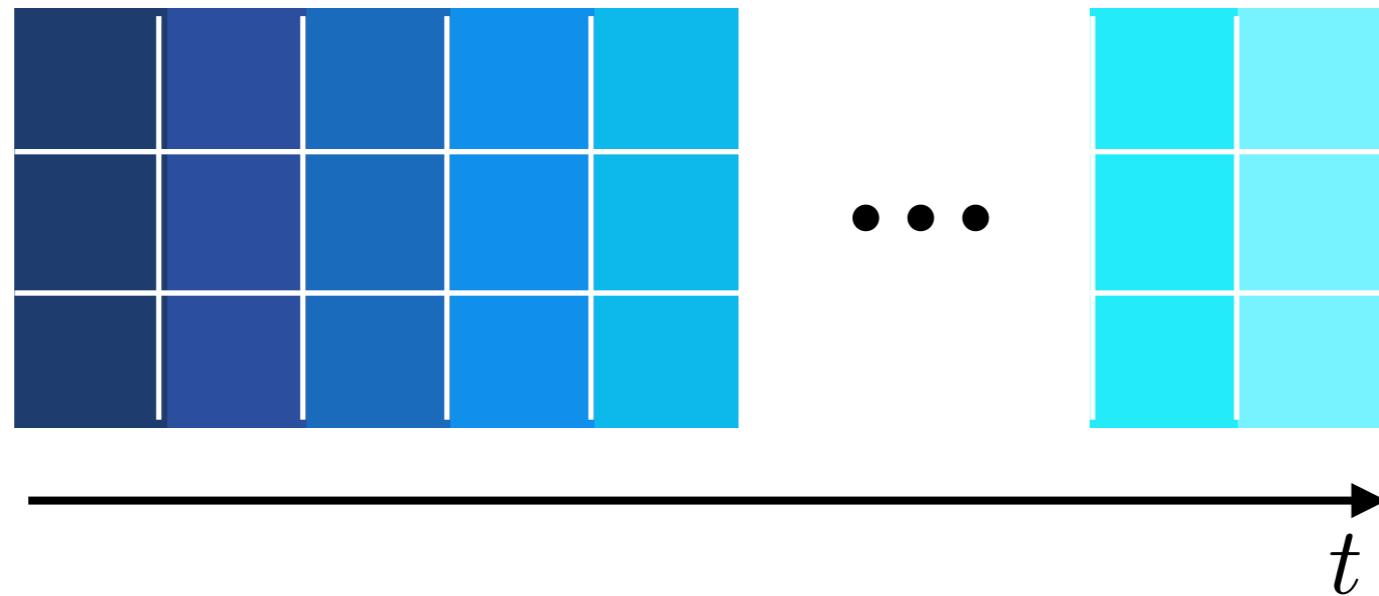
- 2 D Conv



- 3 D Conv



Temporal Structure



sequence data also tends to obey

locality: nearby regions tend to form stronger patterns

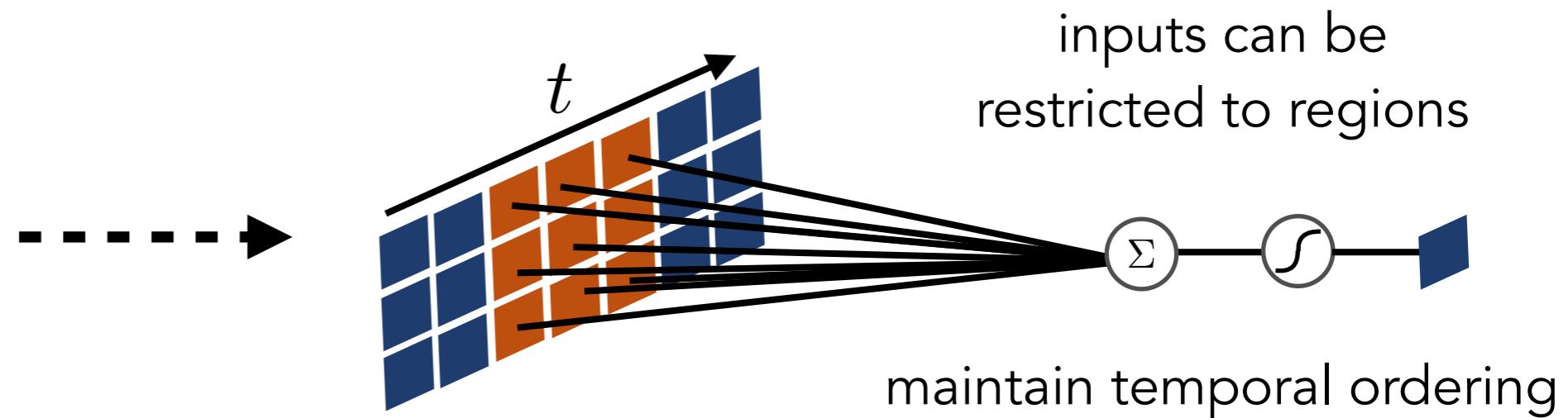
translation invariance: patterns are relative rather than absolute

but has a single axis on which extended patterns occur

let's translate **locality** and **translation invariance** into *inductive biases*

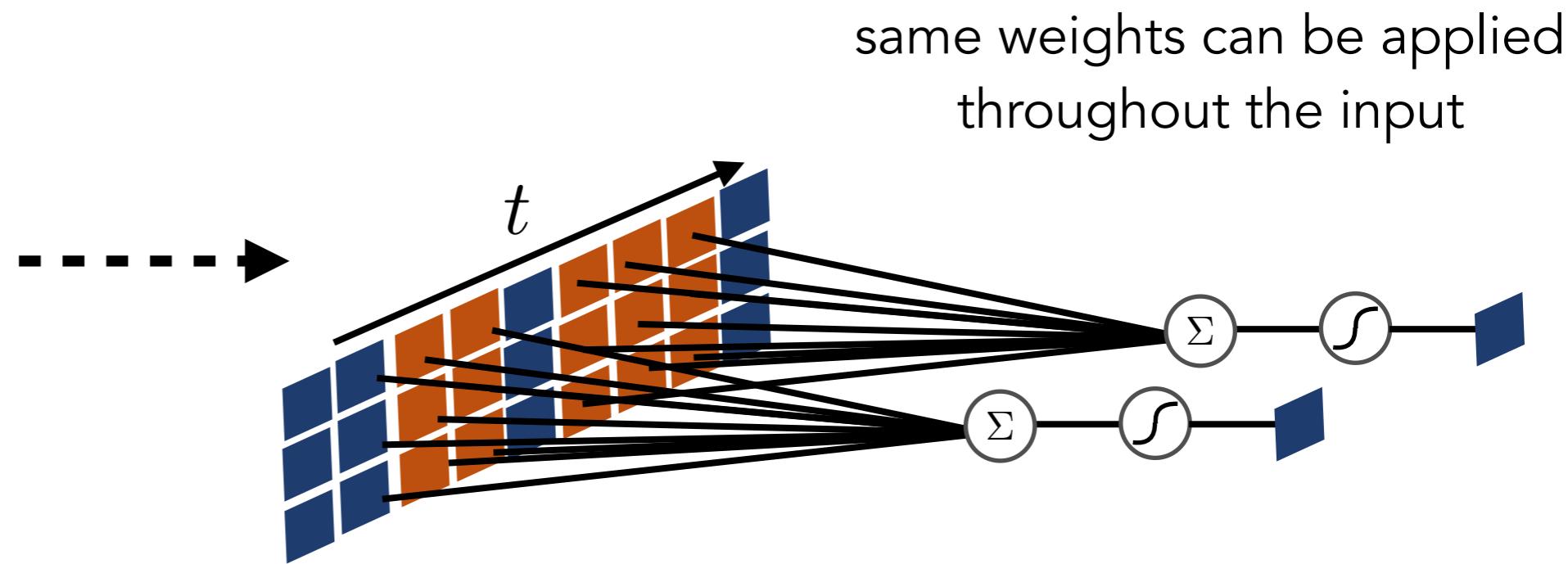
locality

nearby areas tend
to contain stronger
patterns

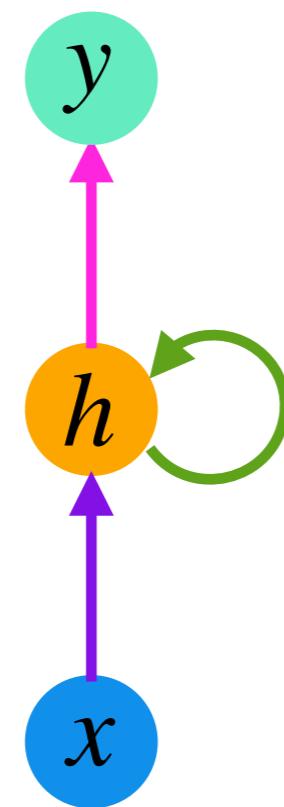


translation invariance

relative positions
are relevant



a **recurrent neural network (RNN)** can be expressed as



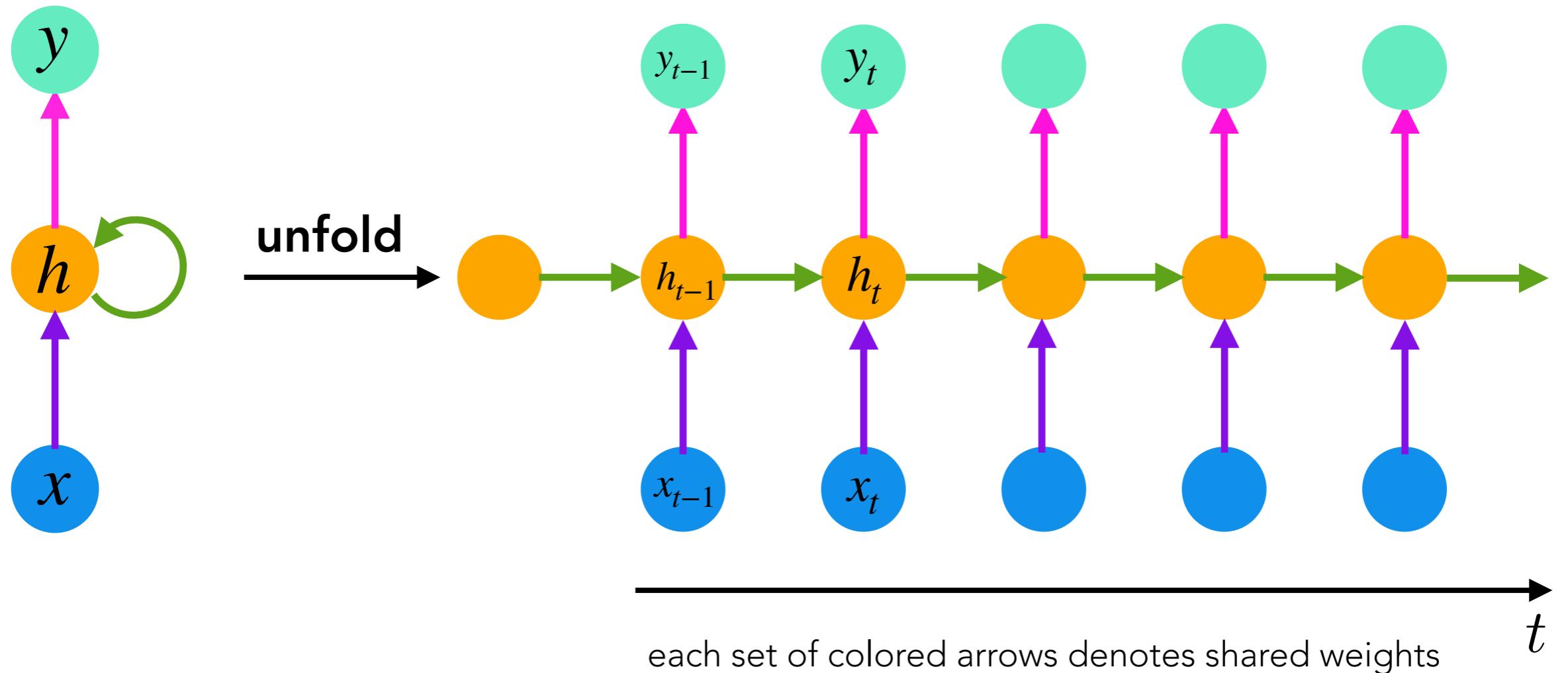
Hidden State

$$\mathbf{h}_t = \sigma(\mathbf{W}_h^\top [\mathbf{h}_{t-1}, \mathbf{x}_t])$$

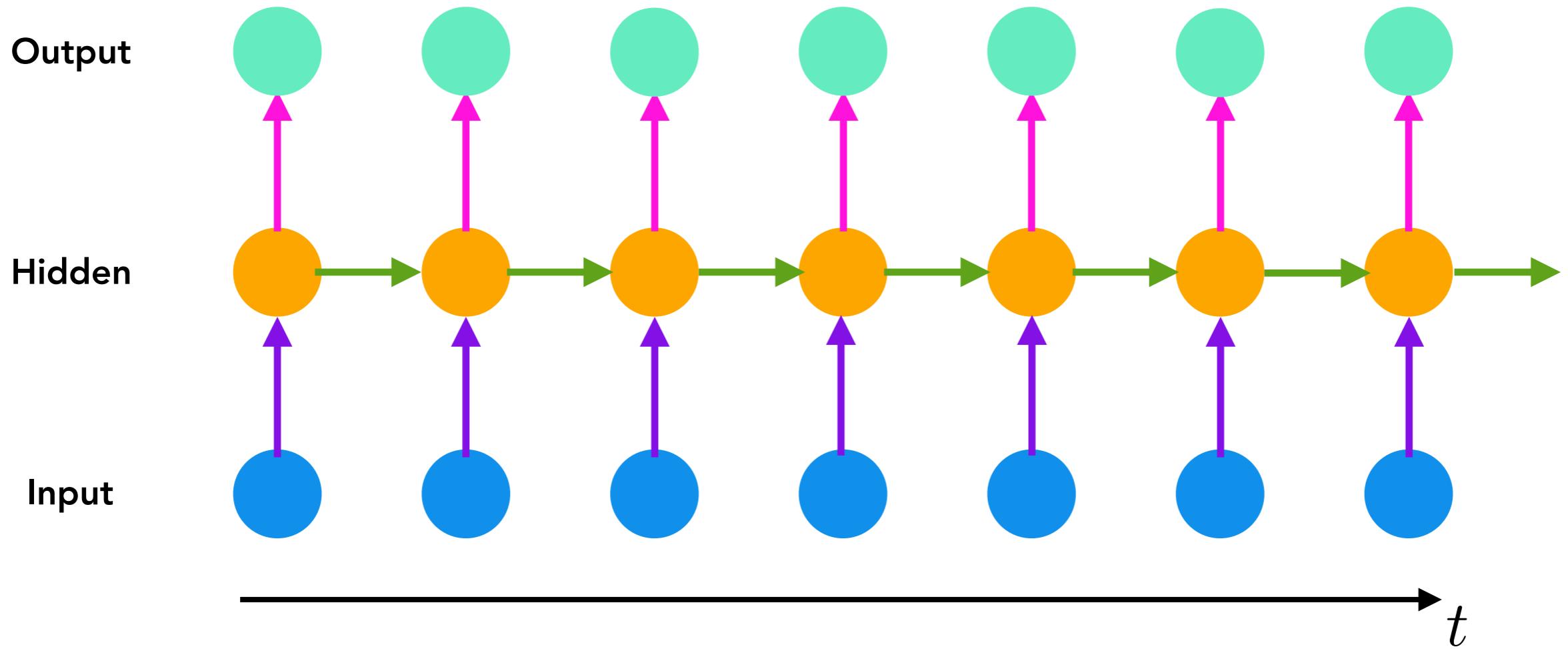
Output

$$\mathbf{y}_t = \sigma(\mathbf{W}_y^\top \mathbf{h}_t)$$

unfold a recurrent neural network



mirror the sequential structure of the data,
we can process the data sequentially

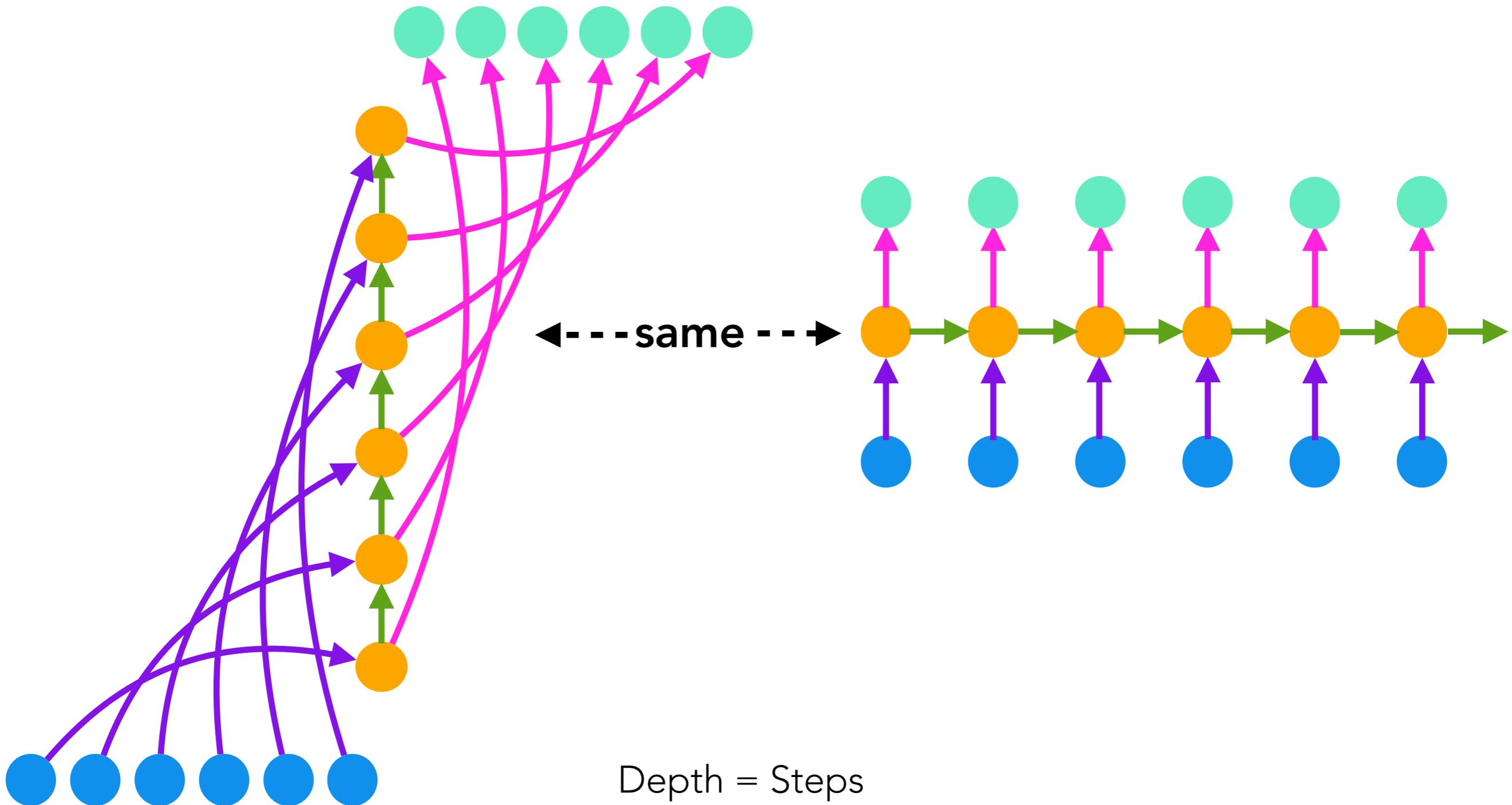


maintain an *internal representation* during processing

→ potentially *infinite* effective input window

→ fixed number of parameters

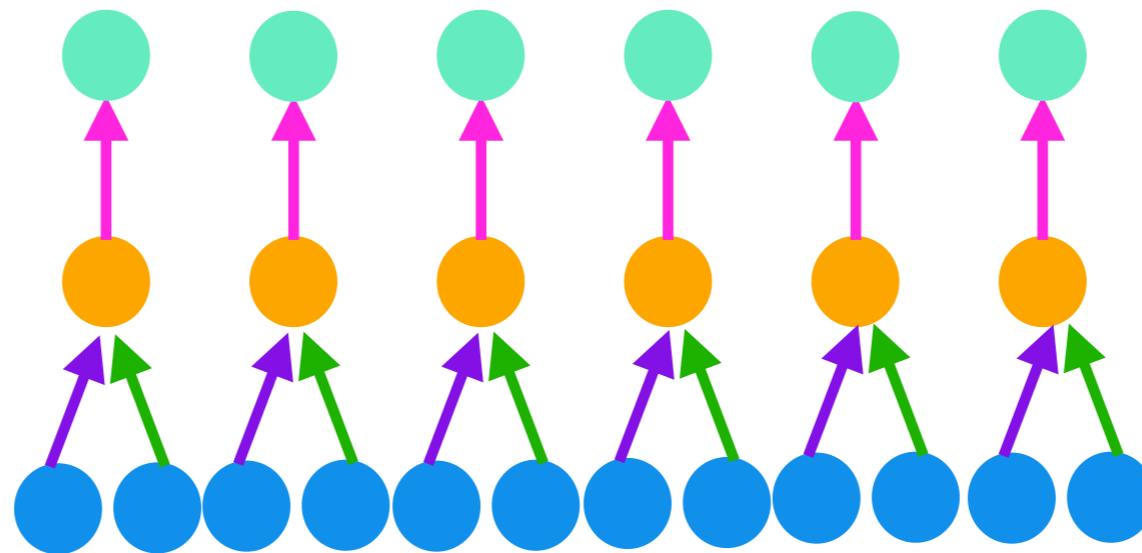
basic recurrent networks are also a **special case**
of standard neural networks with skip connections and shared weights



Weight Sharing

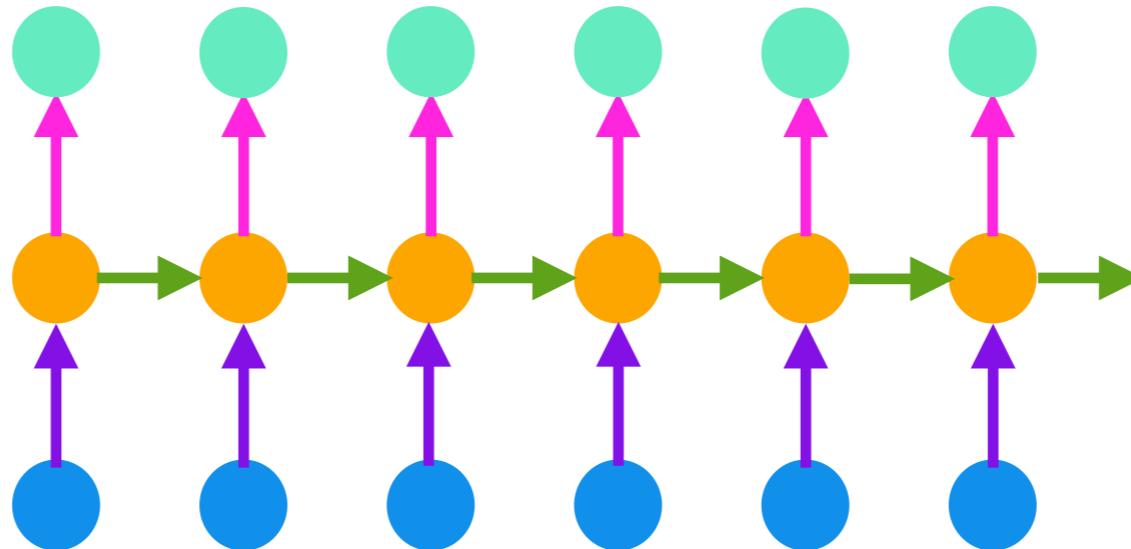
CNN

*sharing weights
across filters*



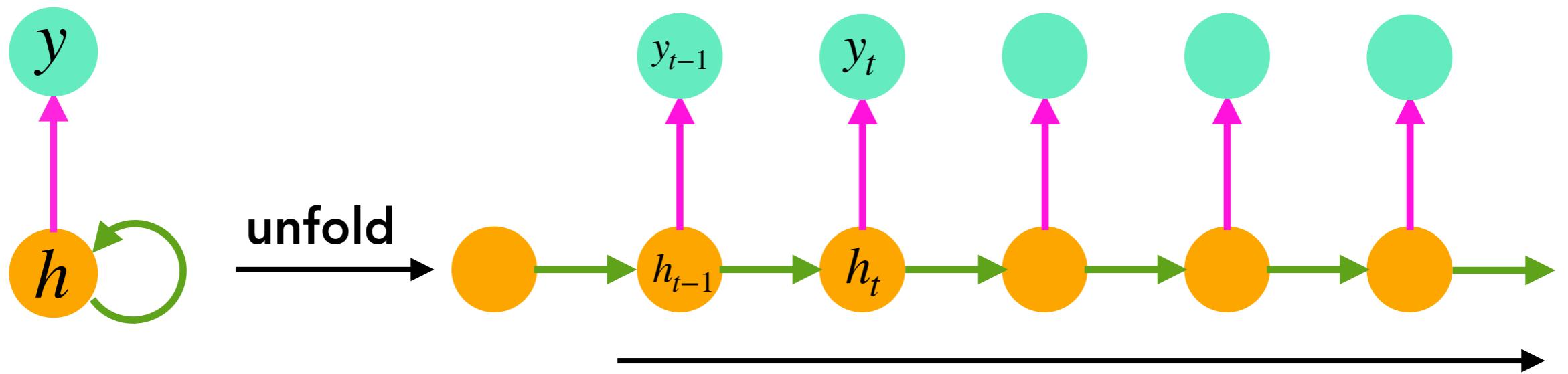
RNN

*sharing weights
across transitions*



Probabilistic Interpretation

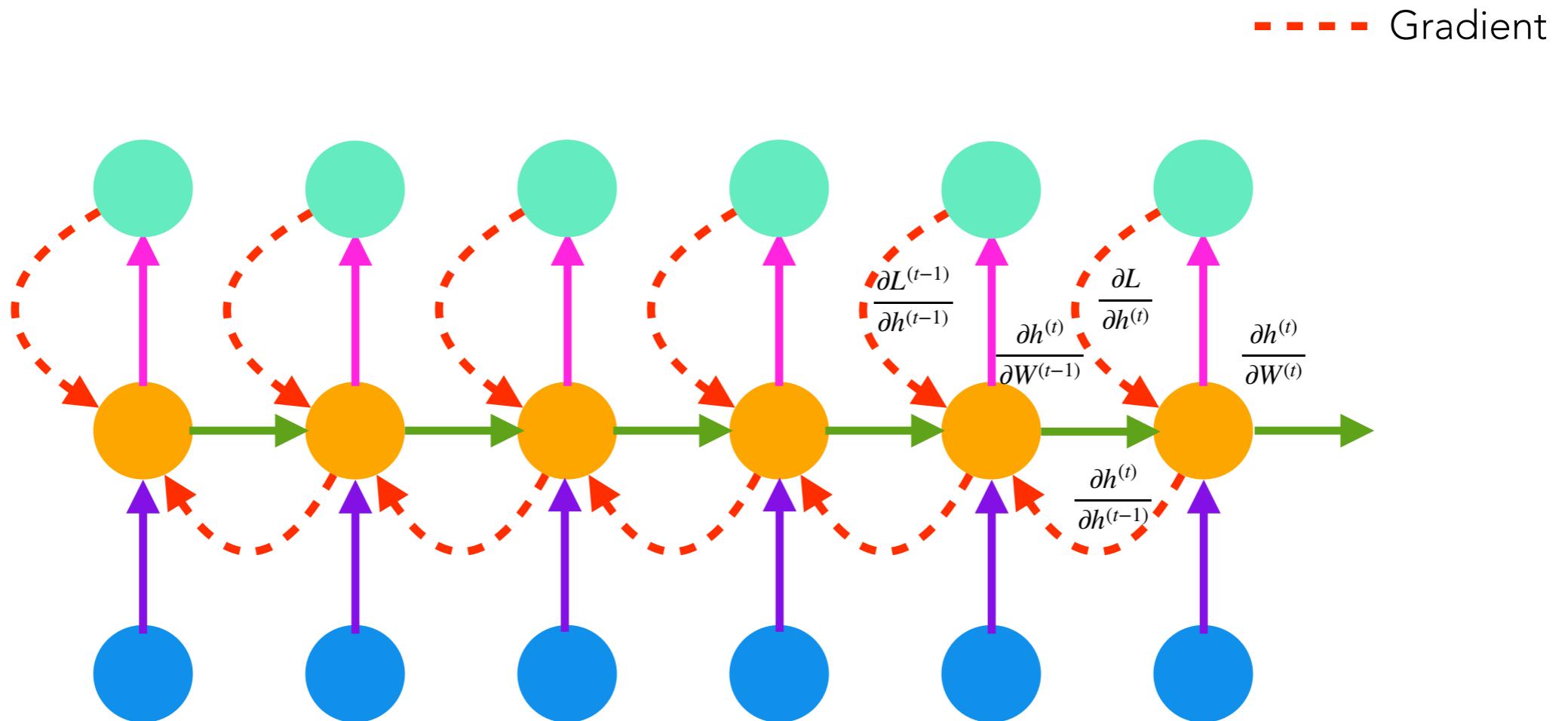
- Consider a simple RNN without inputs



$$p(y_1, y_2, \dots, y_T) = \prod_{t=1}^T p(y_t | y_{t-1}, \dots, y_1)$$

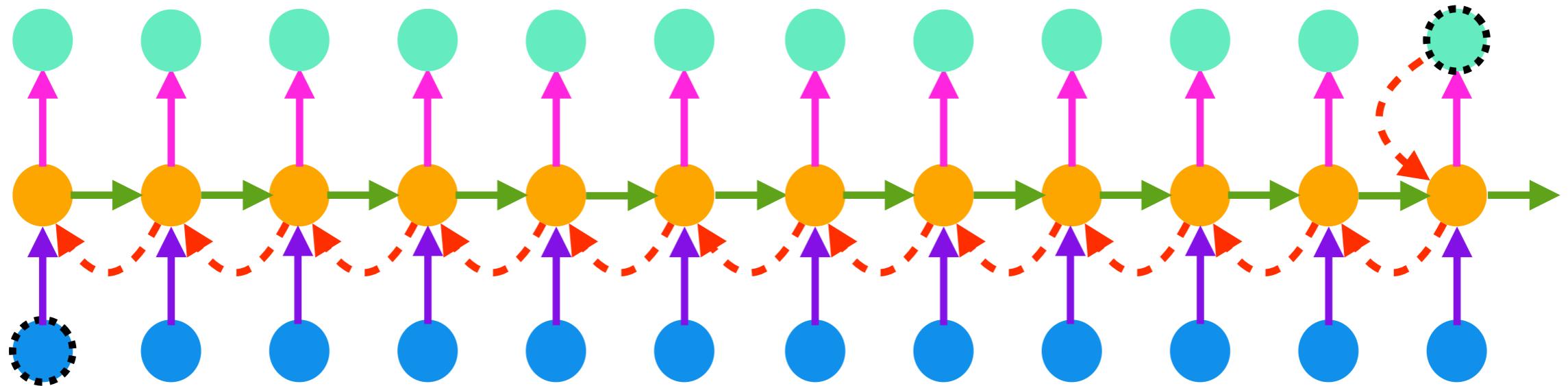
$$p(y_t | y_{t-1}, \dots, y_1) = \text{RNN}(h_t) \quad \text{deterministic hidden states}$$

backpropagation
resulting in ***backpropagation through time (BPTT)***



Introducing dummy variable $W^{(t)}$, which are copies of W but only uses at time step t.

primary difficulty of training RNNs involves propagating information over long horizons

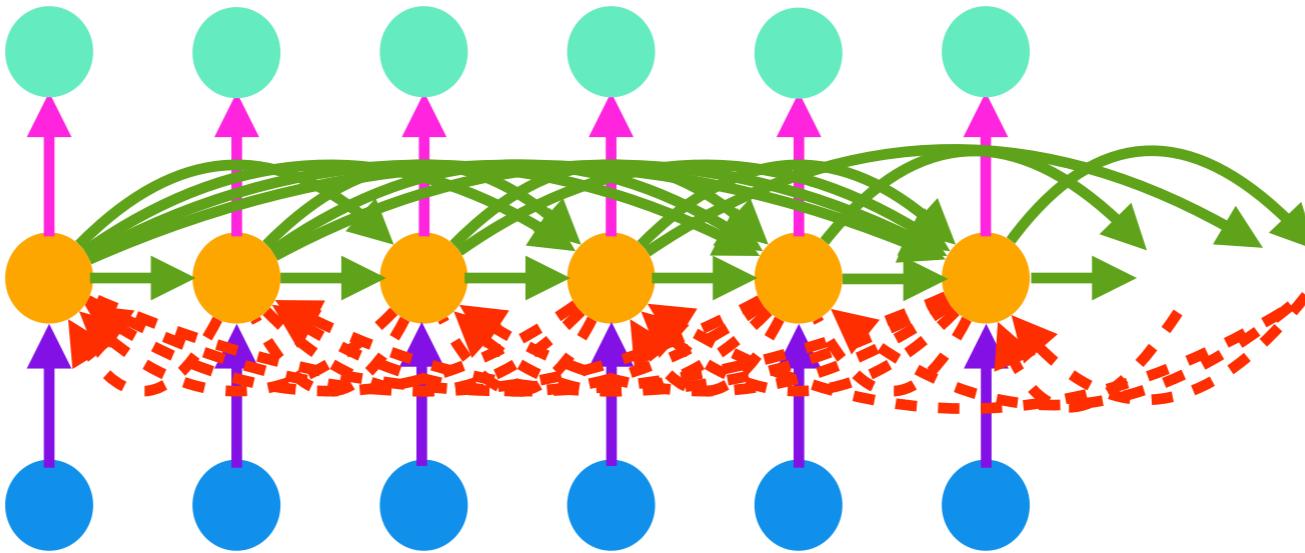


e.g. input at one step is predictive of output at much *later* step

learning extended sequential dependencies
requires propagating gradients over long horizons

- vanishing / exploding gradients
- large memory/computational footprint

naïve attempt to fix information propagation issue



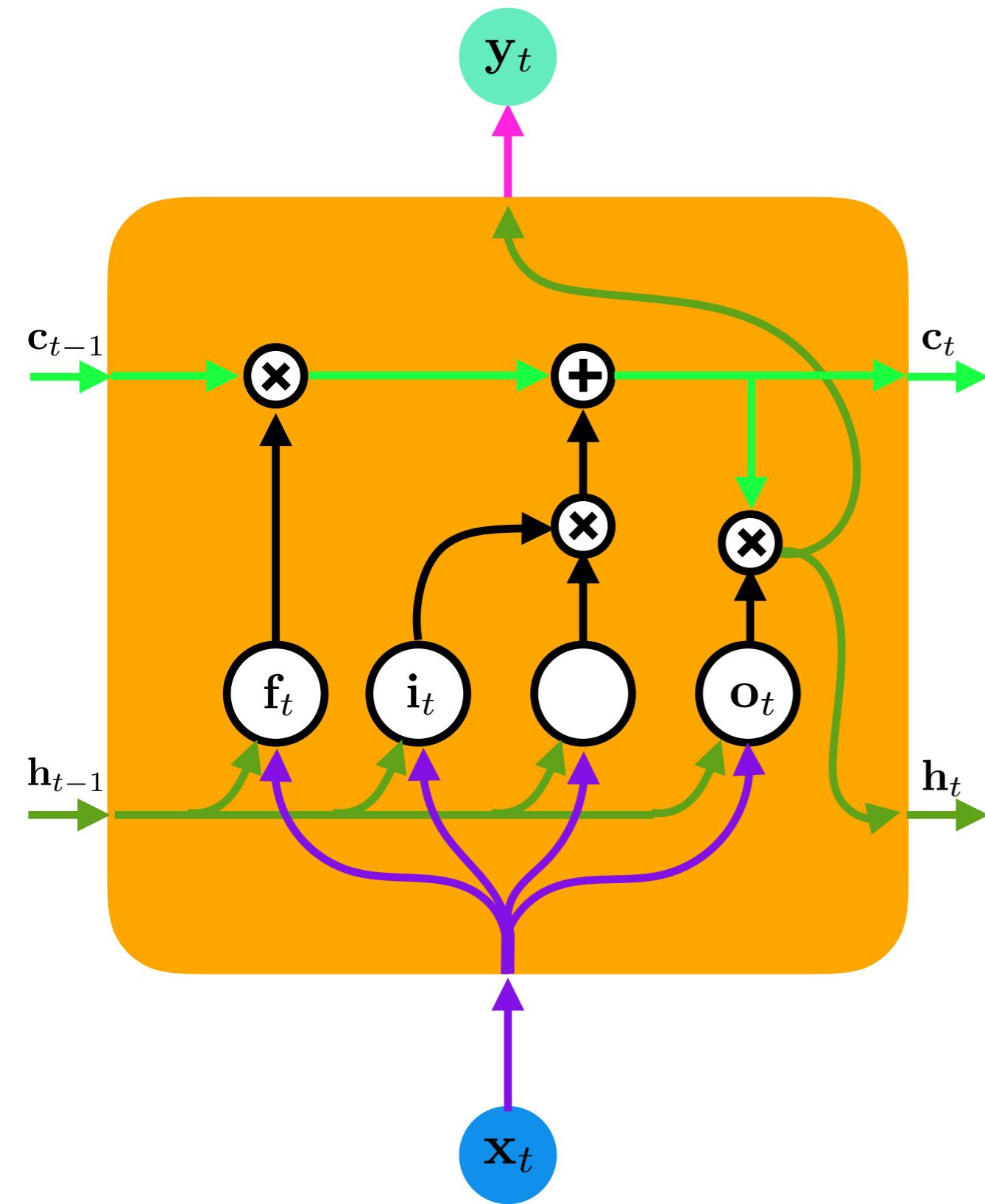
add skip connections across steps

information, gradients can propagate more easily

but...

- increases computation
- must set limit on window size

add trainable ***memory*** to the network
read from and write to “**cell**” state



Long Short-Term Memory (LSTM)

Forget Gate

$$\mathbf{f}_t = \sigma(\mathbf{W}_f^\top [\mathbf{h}_{t-1}, \mathbf{x}_t])$$

Input Gate

$$\mathbf{i}_t = \sigma(\mathbf{W}_i^\top [\mathbf{h}_{t-1}, \mathbf{x}_t])$$

Cell State

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tanh(\mathbf{W}_c^\top [\mathbf{h}_{t-1}, \mathbf{x}_t])$$

Output Gate

$$\mathbf{o}_t = \sigma(\mathbf{W}_o^\top [\mathbf{h}_{t-1}, \mathbf{x}_t])$$

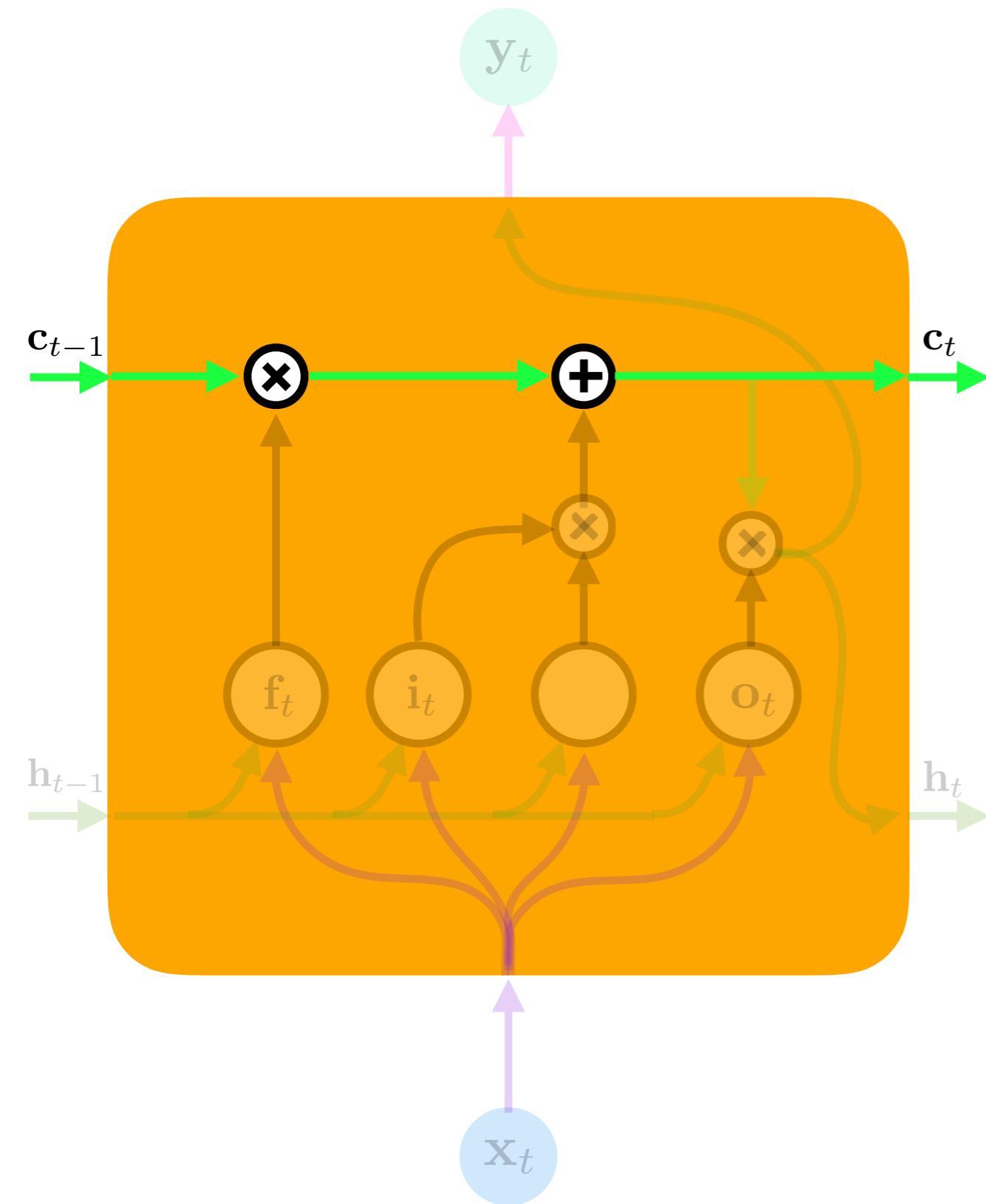
Hidden State

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$

Output

$$\mathbf{y}_t = \sigma(\mathbf{W}_y^\top \mathbf{h}_t)$$

add trainable **memory** to the network
read from and write to “**cell**” state



Long Short-Term Memory (LSTM)

Forget Gate

$$f_t = \sigma(W_f^\top [h_{t-1}, x_t])$$

Input Gate

$$i_t = \sigma(W_i^\top [h_{t-1}, x_t])$$

Cell State

$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh(W_c^\top [h_{t-1}, x_t])$$

Output Gate

$$o_t = \sigma(W_o^\top [h_{t-1}, x_t])$$

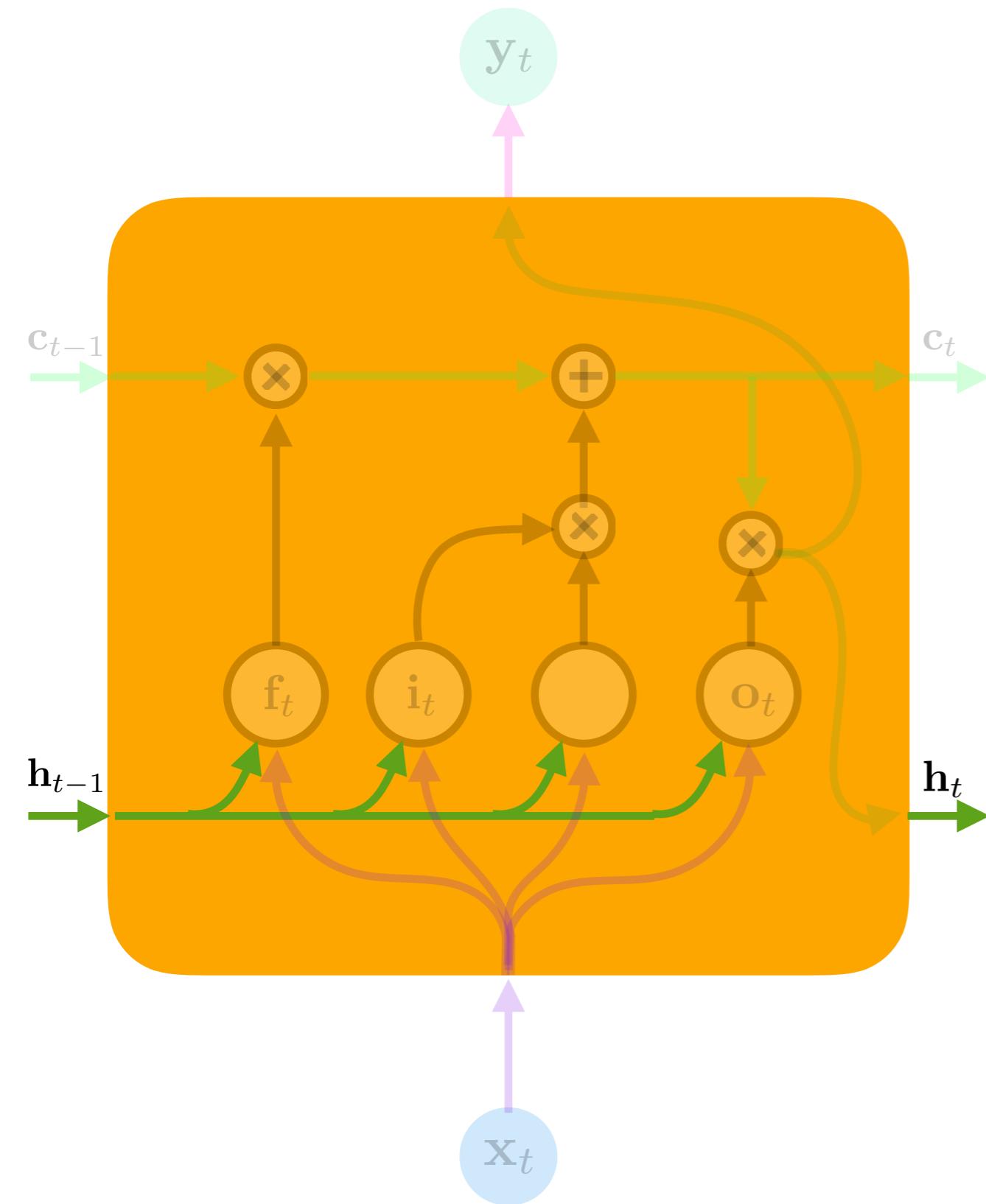
Hidden State

$$h_t = o_t \odot \tanh(c_t)$$

Output

$$y_t = \sigma(W_y^\top h_t)$$

add trainable ***memory*** to the network
read from and write to “***cell***” state



Long Short-Term Memory (LSTM)

Forget Gate

$$f_t = \sigma(W_f^\top [h_{t-1}, x_t])$$

Input Gate

$$i_t = \sigma(W_i^\top [h_{t-1}, x_t])$$

Cell State

$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh(W_c^\top [h_{t-1}, x_t])$$

Output Gate

$$o_t = \sigma(W_o^\top [h_{t-1}, x_t])$$

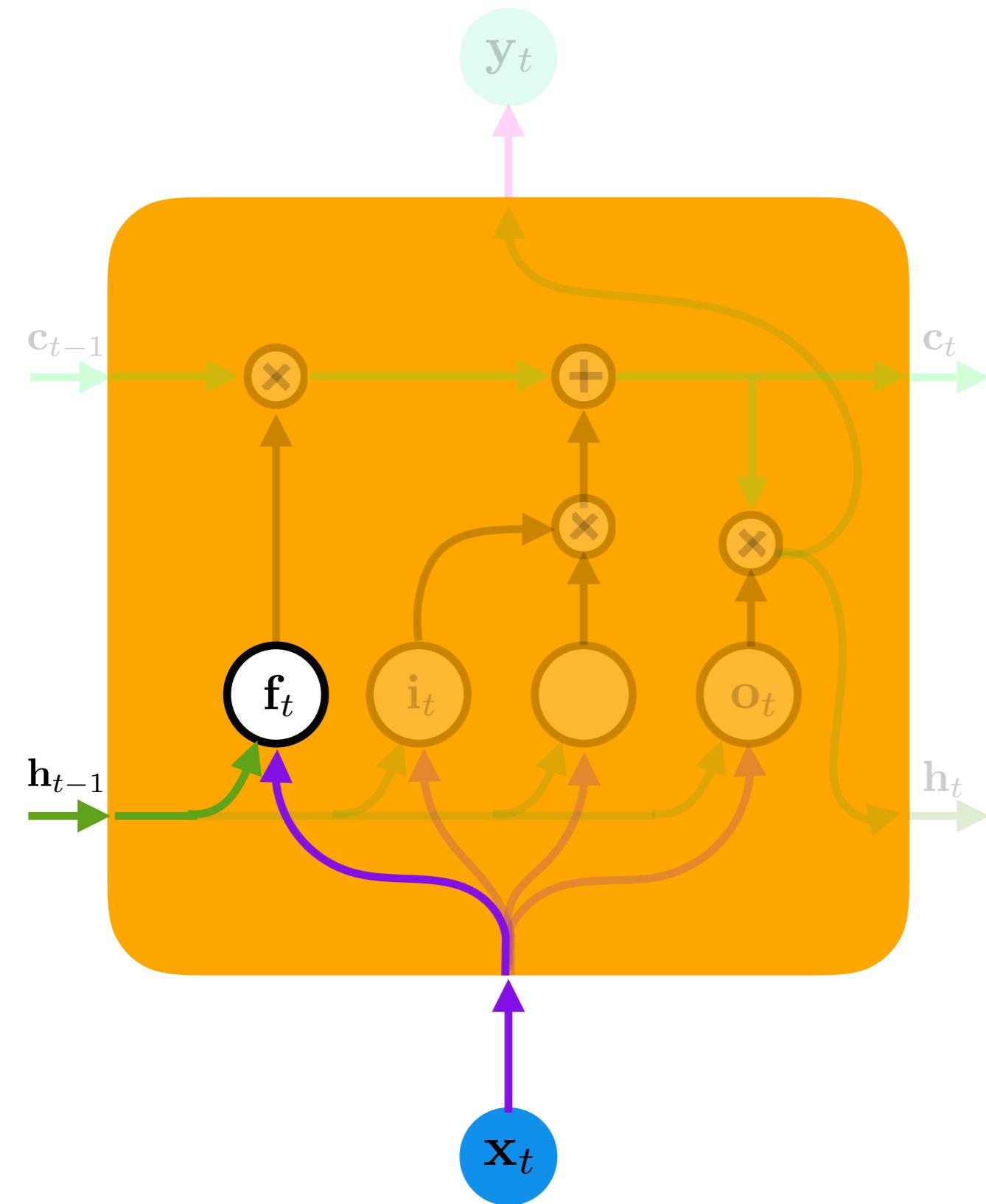
Hidden State

$$h_t = o_t \odot \tanh(c_t)$$

Output

$$y_t = \sigma(W_y^\top h_t)$$

add trainable **memory** to the network
read from and write to “**cell**” state



Long Short-Term Memory (LSTM)

Forget Gate

$$f_t = \sigma(W_f^\top [h_{t-1}, x_t])$$

Input Gate

$$i_t = \sigma(W_i^\top [h_{t-1}, x_t])$$

Cell State

$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh(W_c^\top [h_{t-1}, x_t])$$

Output Gate

$$o_t = \sigma(W_o^\top [h_{t-1}, x_t])$$

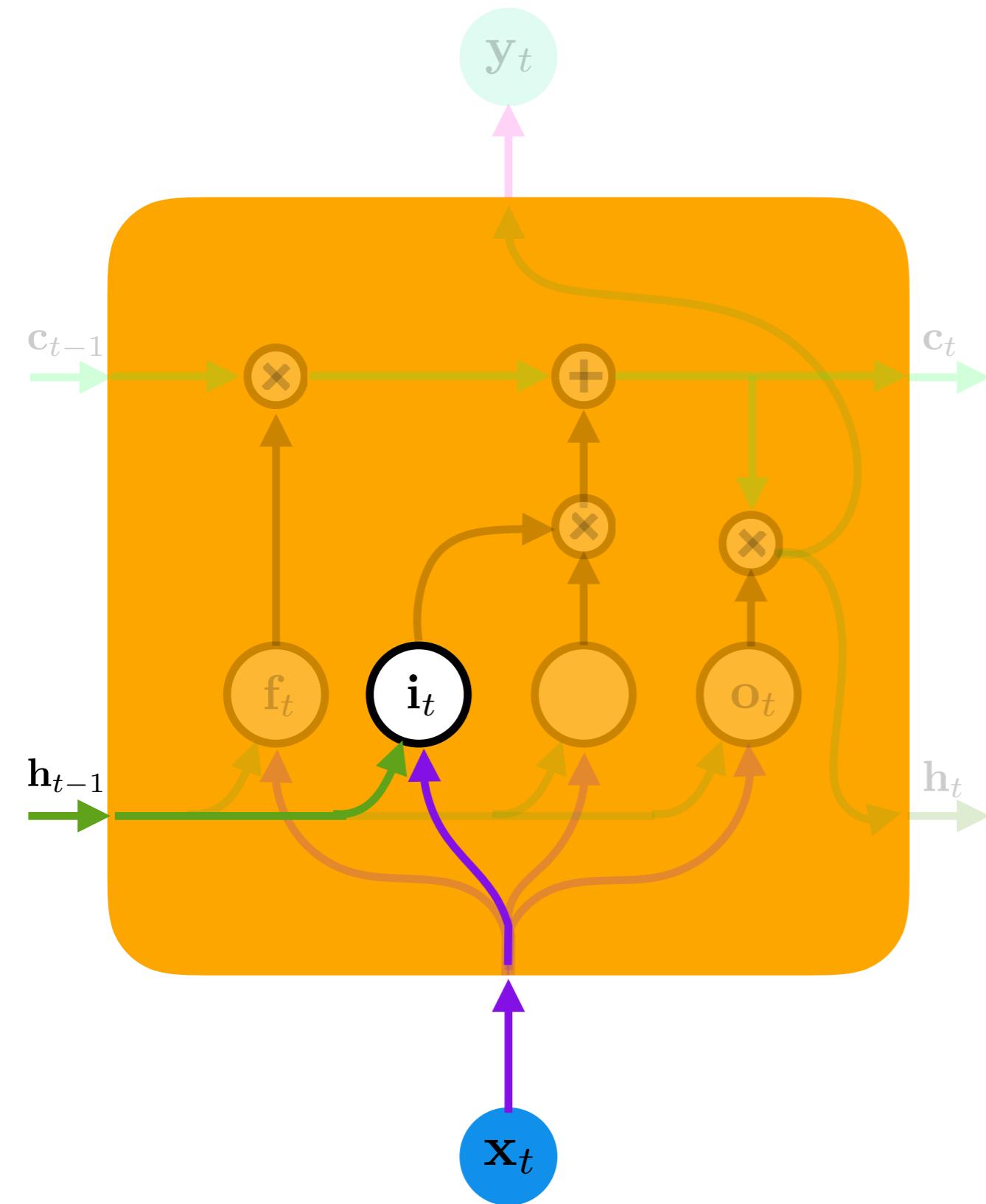
Hidden State

$$h_t = o_t \odot \tanh(c_t)$$

Output

$$y_t = \sigma(W_y^\top h_t)$$

add trainable ***memory*** to the network
read from and write to “***cell***” state



Long Short-Term Memory (LSTM)

Forget Gate

$$f_t = \sigma(\mathbf{W}_f^\top [\mathbf{h}_{t-1}, \mathbf{x}_t])$$

Input Gate

$$i_t = \sigma(\mathbf{W}_i^\top [\mathbf{h}_{t-1}, \mathbf{x}_t])$$

Cell State

$$\mathbf{c}_t = f_t \odot \mathbf{c}_{t-1} + i_t \odot \tanh(\mathbf{W}_c^\top [\mathbf{h}_{t-1}, \mathbf{x}_t])$$

Output Gate

$$o_t = \sigma(\mathbf{W}_o^\top [\mathbf{h}_{t-1}, \mathbf{x}_t])$$

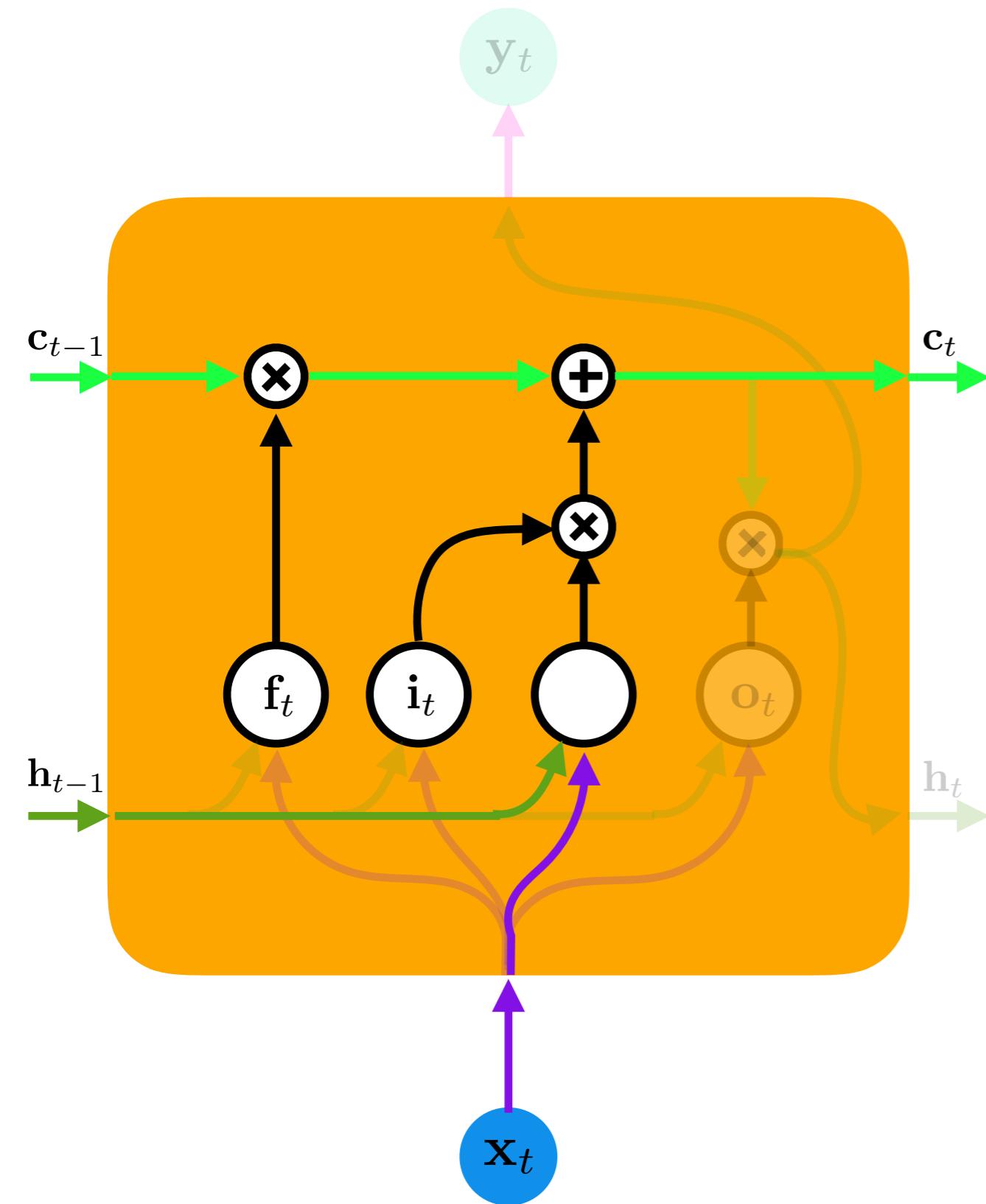
Hidden State

$$\mathbf{h}_t = o_t \odot \tanh(\mathbf{c}_t)$$

Output

$$\mathbf{y}_t = \sigma(\mathbf{W}_y^\top \mathbf{h}_t)$$

add trainable ***memory*** to the network
read from and write to “***cell***” state



Long Short-Term Memory (LSTM)

Forget Gate

$$f_t = \sigma(W_f^T[h_{t-1}, x_t])$$

Input Gate

$$i_t = \sigma(W_i^T[h_{t-1}, x_t])$$

Cell State

$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh(W_c^T[h_{t-1}, x_t])$$

Output Gate

$$o_t = \sigma(W_o^T[h_{t-1}, x_t])$$

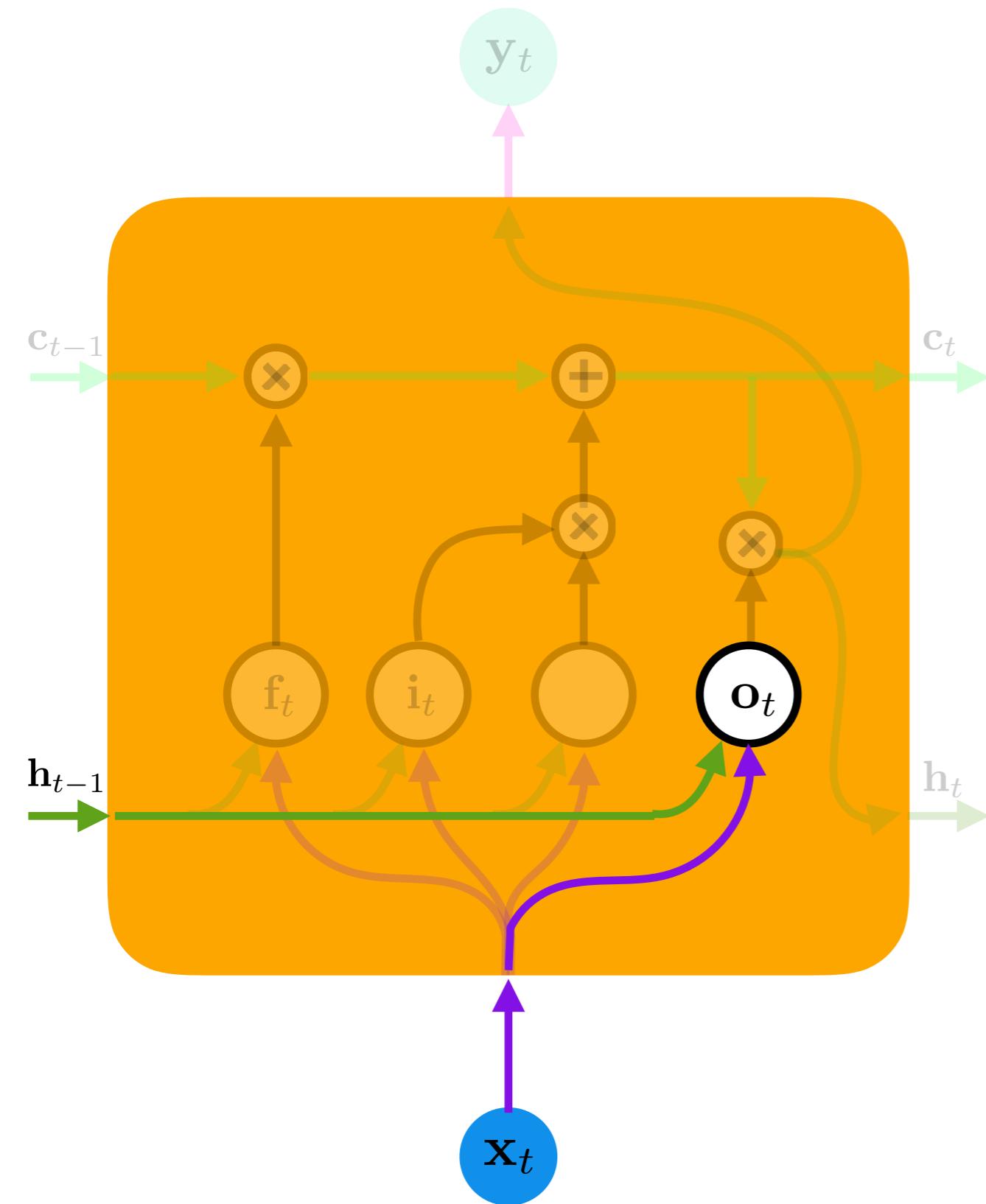
Hidden State

$$h_t = o_t \odot \tanh(c_t)$$

Output

$$y_t = \sigma(W_y^T h_t)$$

add trainable ***memory*** to the network
read from and write to “***cell***” state



Long Short-Term Memory (LSTM)

Forget Gate

$$f_t = \sigma(W_f^T [h_{t-1}, x_t])$$

Input Gate

$$i_t = \sigma(W_i^T [h_{t-1}, x_t])$$

Cell State

$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh(W_c^T [h_{t-1}, x_t])$$

Output Gate

$$o_t = \sigma(W_o^T [h_{t-1}, x_t])$$

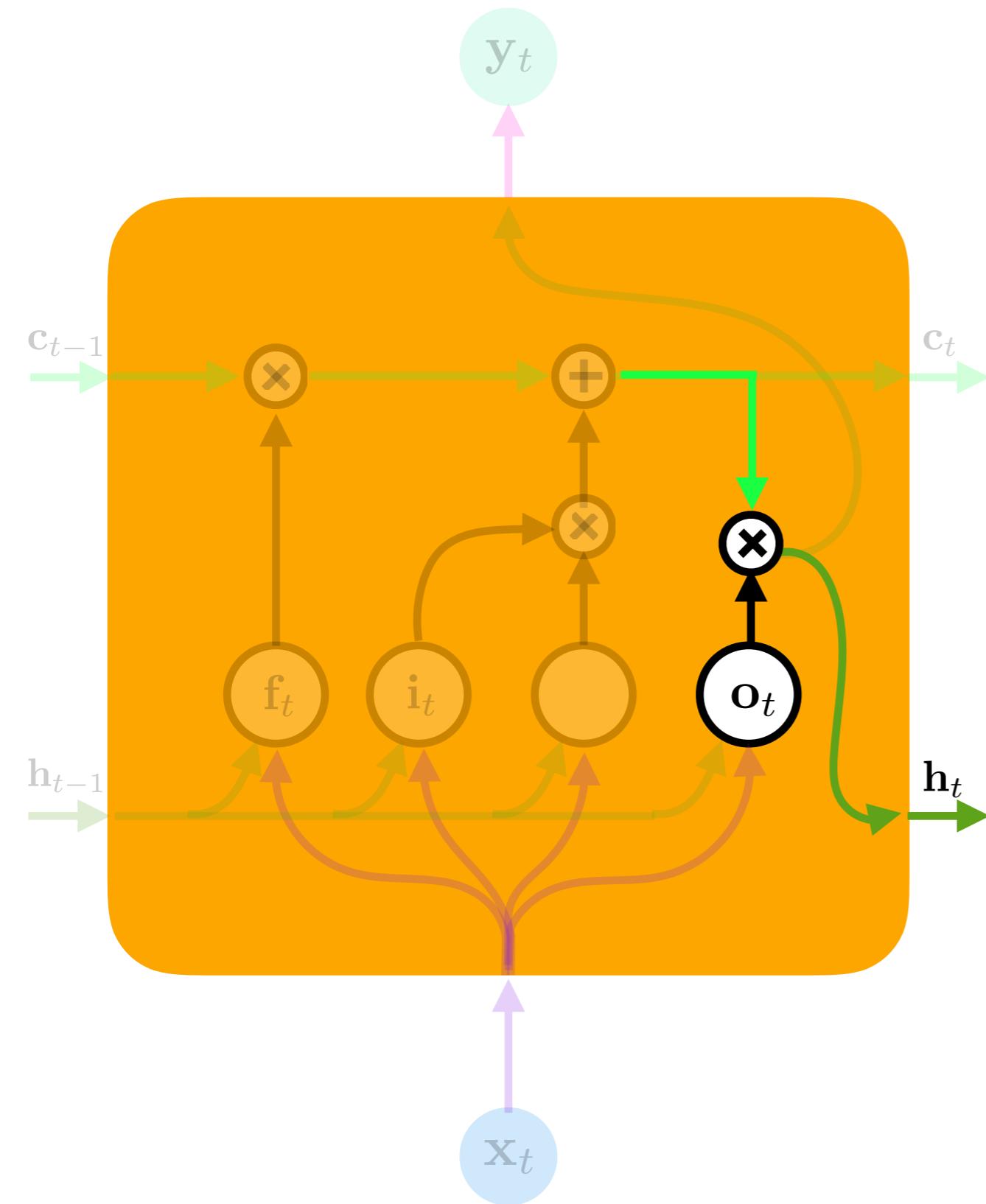
Hidden State

$$h_t = o_t \odot \tanh(c_t)$$

Output

$$y_t = \sigma(W_y^T h_t)$$

add trainable ***memory*** to the network
read from and write to “***cell***” state



Long Short-Term Memory (LSTM)

Forget Gate

$$f_t = \sigma(W_f^\top [h_{t-1}, x_t])$$

Input Gate

$$i_t = \sigma(W_i^\top [h_{t-1}, x_t])$$

Cell State

$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh(W_c^\top [h_{t-1}, x_t])$$

Output Gate

$$o_t = \sigma(W_o^\top [h_{t-1}, x_t])$$

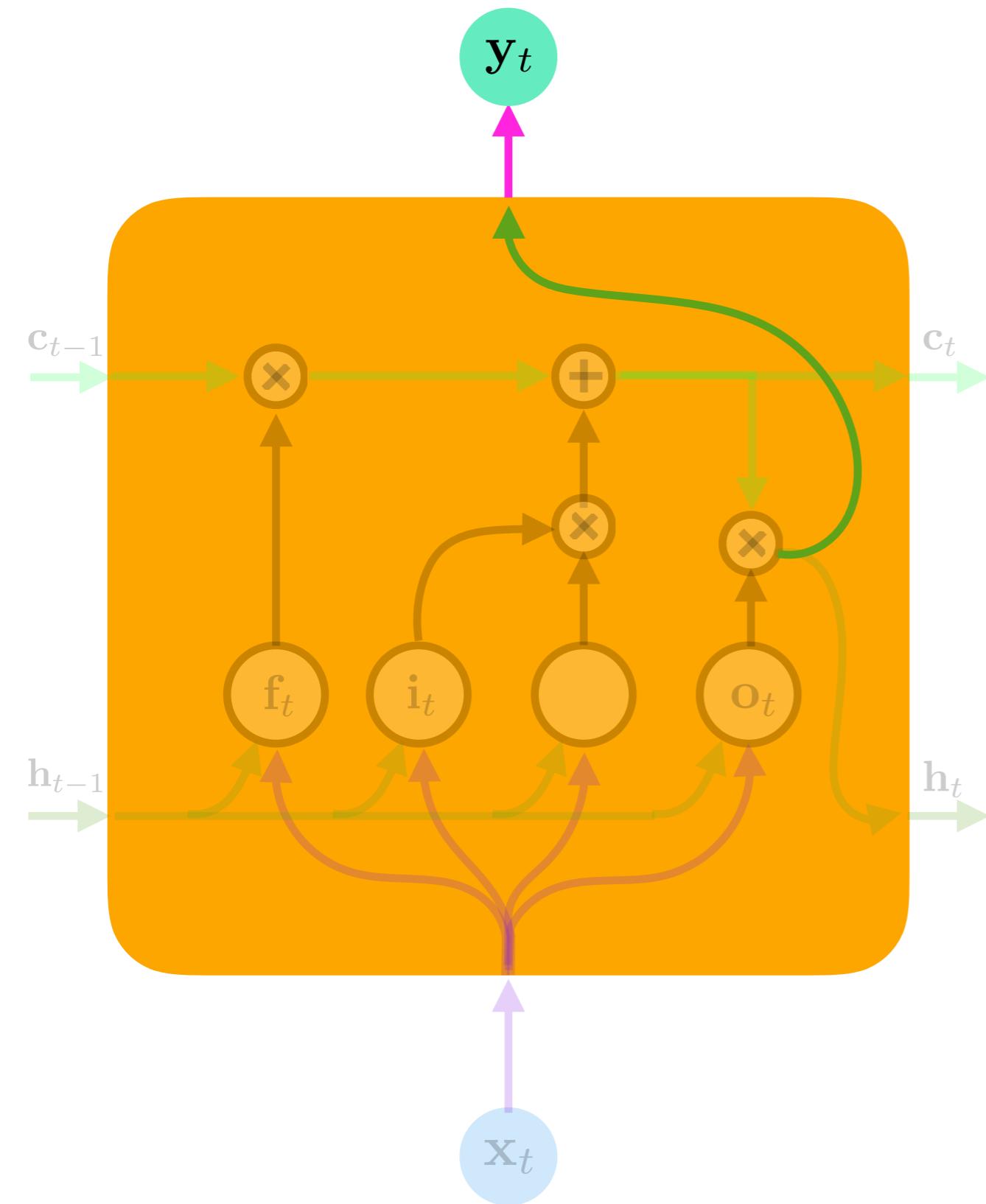
Hidden State

$$h_t = o_t \odot \tanh(c_t)$$

Output

$$y_t = \sigma(W_y^\top h_t)$$

add trainable ***memory*** to the network
read from and write to “***cell***” state



Long Short-Term Memory (LSTM)

Forget Gate

$$f_t = \sigma(W_f^\top [h_{t-1}, x_t])$$

Input Gate

$$i_t = \sigma(W_i^\top [h_{t-1}, x_t])$$

Cell State

$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh(W_c^\top [h_{t-1}, x_t])$$

Output Gate

$$o_t = \sigma(W_o^\top [h_{t-1}, x_t])$$

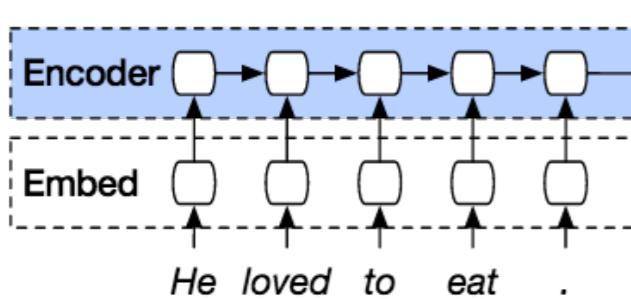
Hidden State

$$h_t = o_t \odot \tanh(c_t)$$

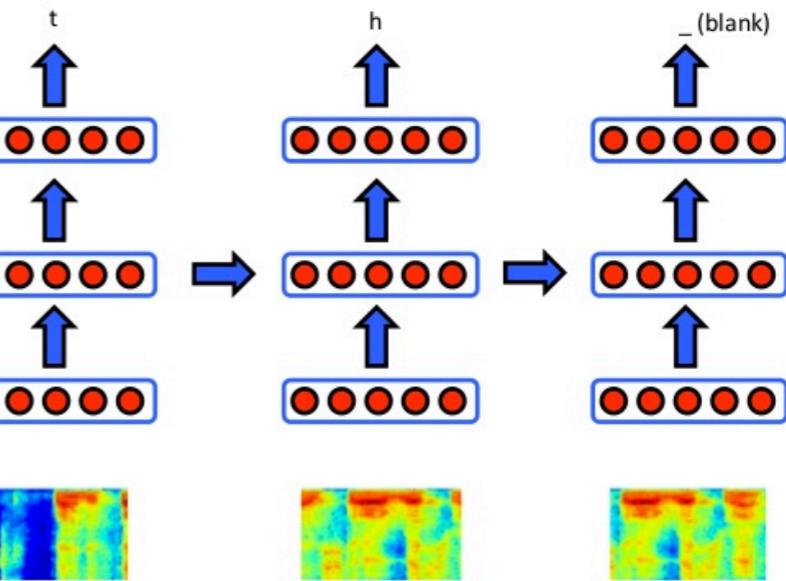
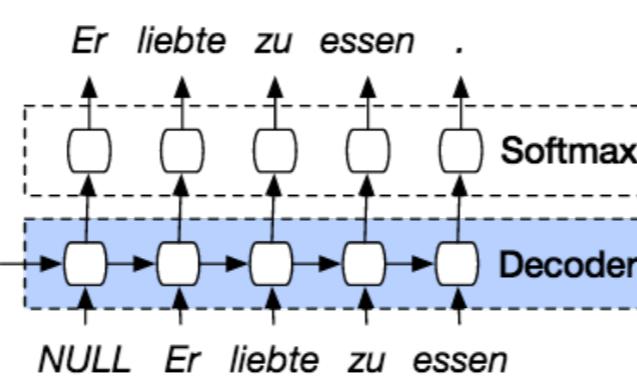
Output

$$y_t = \sigma(W_y^\top h_t)$$

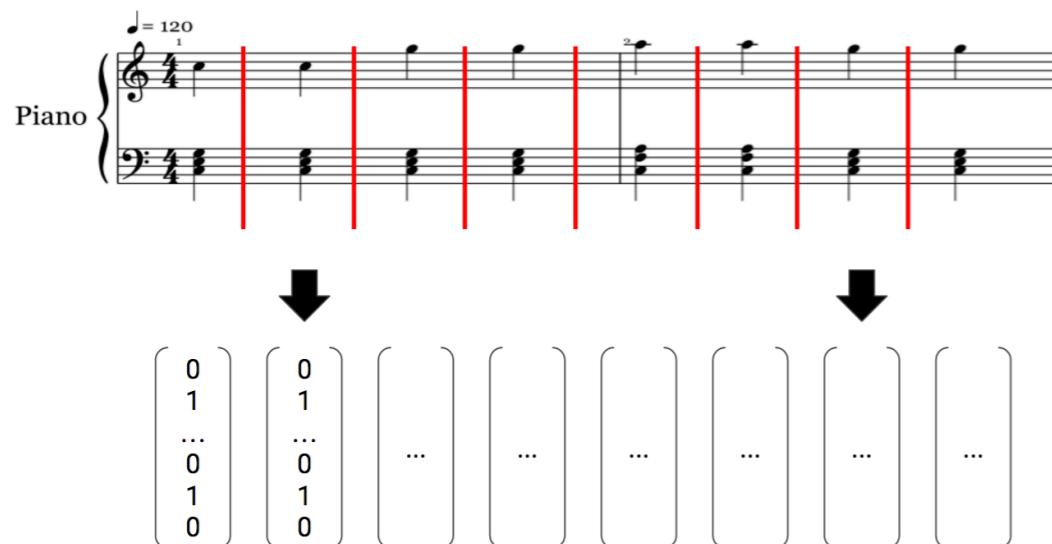
RNN Applications



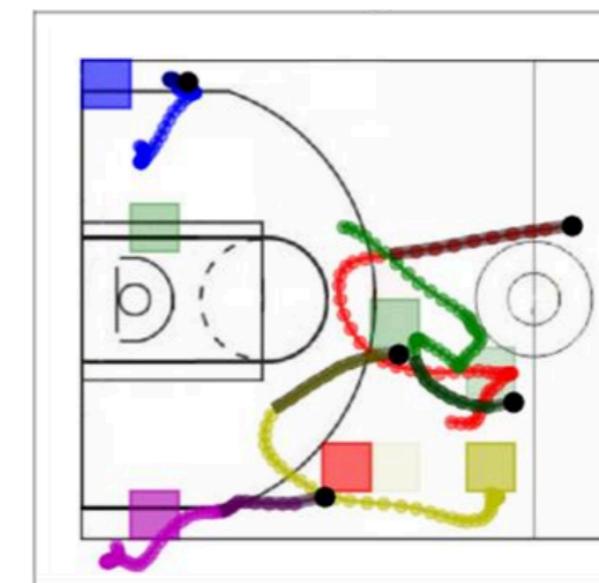
machine translation



speech recognition



music composition



trajectory generation