

The background of the slide features a large, faint watermark of the University of Pisa crest. The crest is circular and contains a classical face, likely Minerva, surrounded by the Latin motto 'IN ARTE VERITAS' and the year '1543'.

Gated Recurrent Networks

INTELLIGENT SYSTEMS FOR PATTERN RECOGNITION (ISPR)

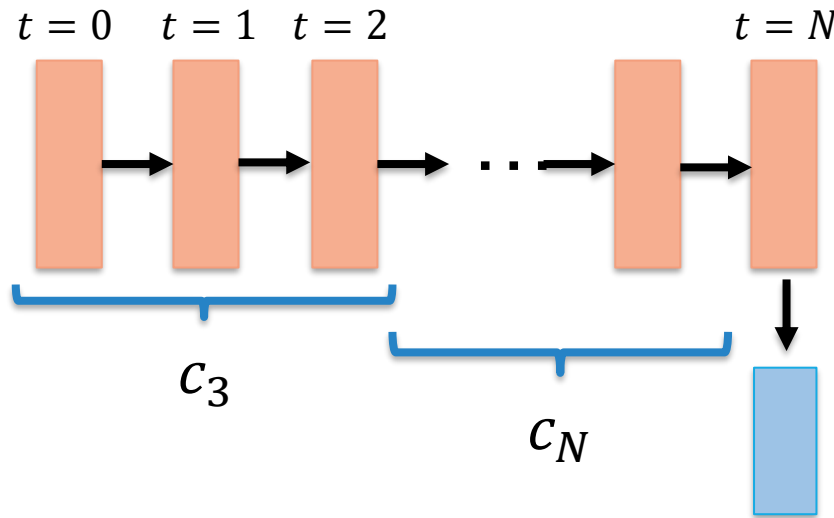
DAVIDE BACCIU – DIPARTIMENTO DI INFORMATICA - UNIVERSITA' DI PISA

davide.bacciu@unipi.it

Lecture Outline

- RNN Repetita
- Motivations
 - Learning long-term dependencies is difficult
 - Gradient issues (EVGP)
- Solving the EVGP
 - Constant error propagation
 - Adaptive forgetting
- Gated RNN
 - Long-Short Term Memories (LSTM)
 - Gated Recurrent Units (GRU)
- Advanced topics
 - Understanding and exploiting memory encoding
 - Applications

Dealing with Sequences in NN



Variable size data describing
sequentially dependent
information

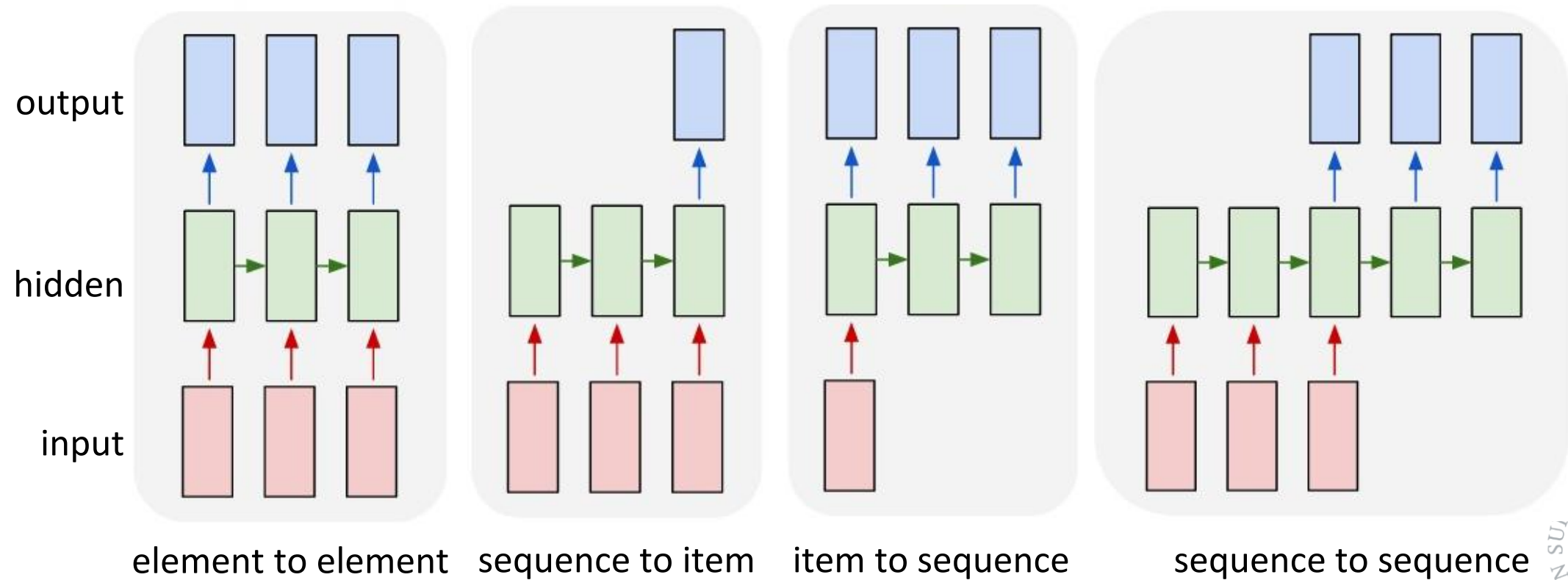
Neural models need to capture
dynamic context c_t to perform
predictions

- Recurrent Neural Network
 - Fully adaptive (Elman, SRN, ...)
 - Randomized approaches (Reservoir Computing)
- Introduce (deep) **gated** recurrent networks

RNN Design

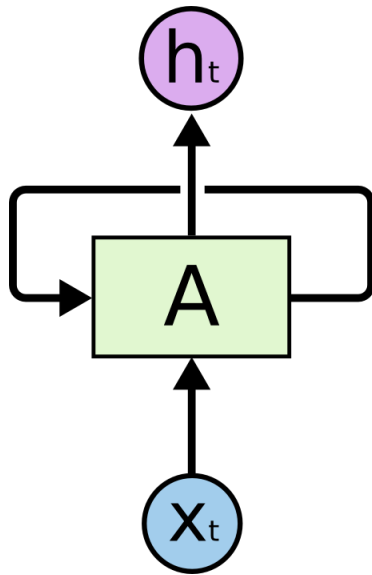
- **Inductive bias/Expressiveness:** the network structure influences the sequential data processing.
- **Training:** the network should be easy to train. Depends on the architecture, initialization, and learning algorithm.
- **Computational Efficiency:** the network should be efficient
 - training or at inference time.
 - different hardware: GPU, CPU, embedded devices.

Supervised Recurrent Tasks



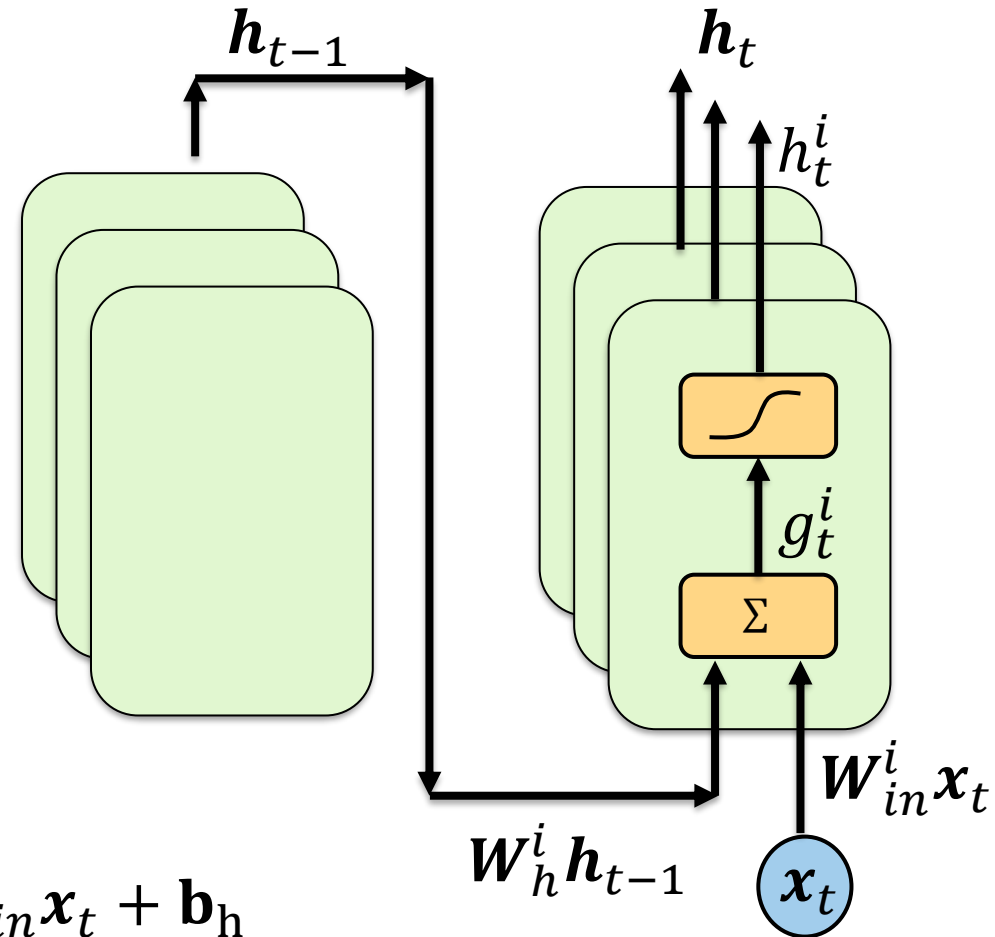
A Non-Gated RNN (a.k.a. Vanilla)

$$y_t = f(W_{out}h_t + b_{out})$$



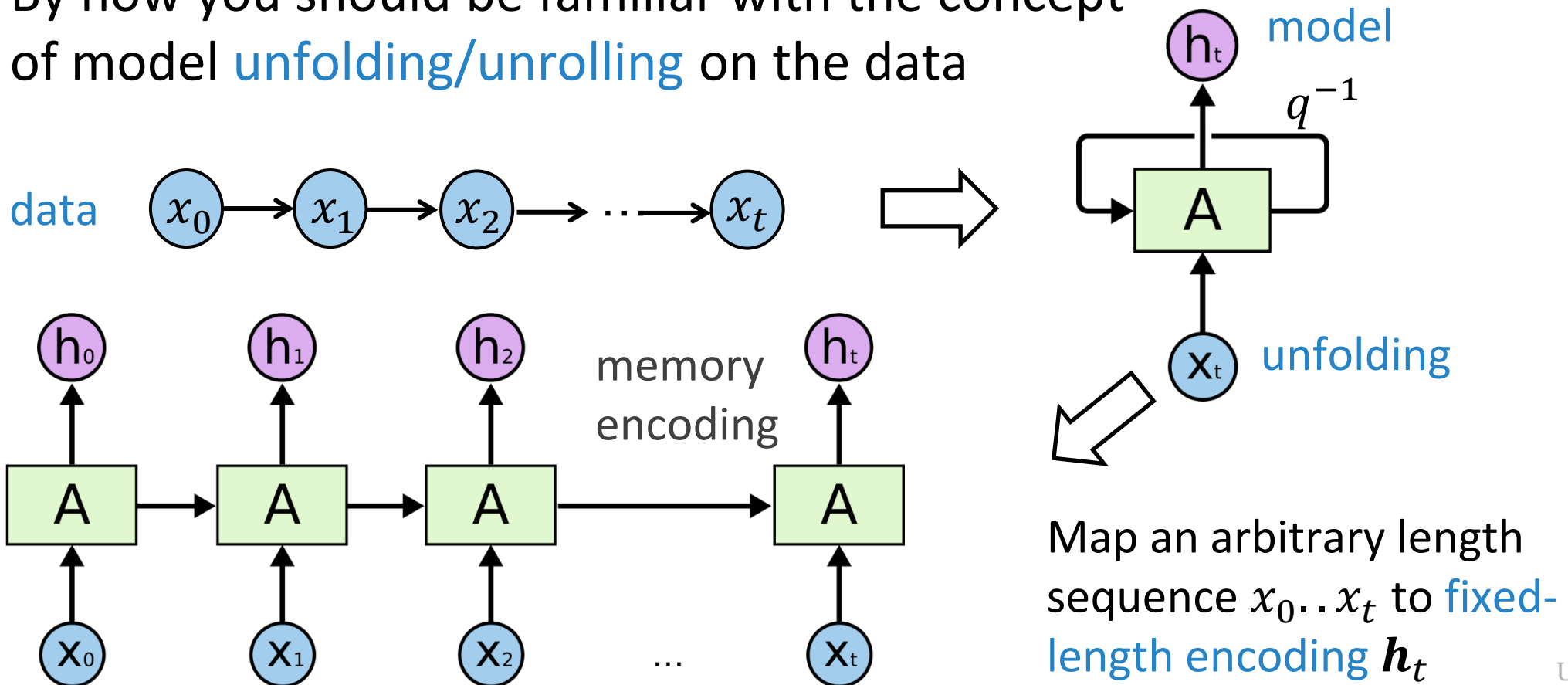
$$h_t = \tanh(g_t)$$

$$g_t(h_{t-1}, x_t) = W_h h_{t-1} + W_{in} x_t + b_h$$

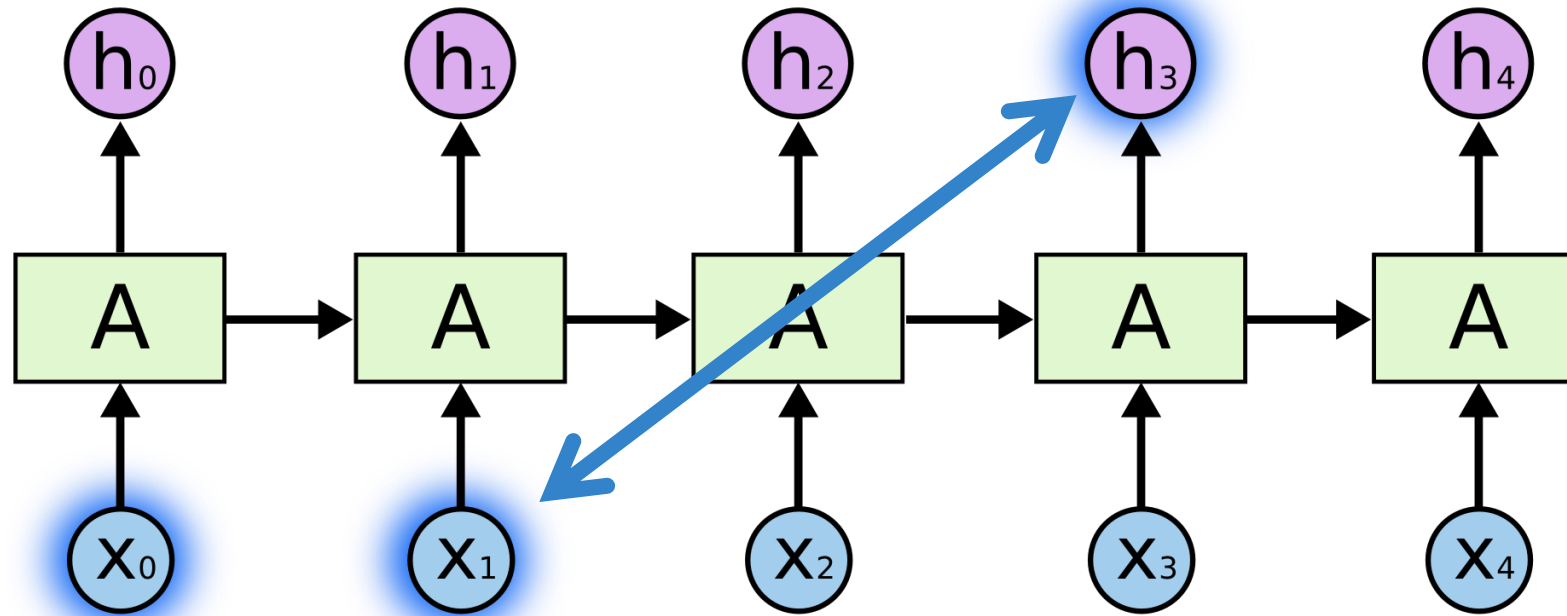


Unfolding RNN (Forward Pass)

By now you should be familiar with the concept of model **unfolding/unrolling** on the data



Learning to Encode Input History



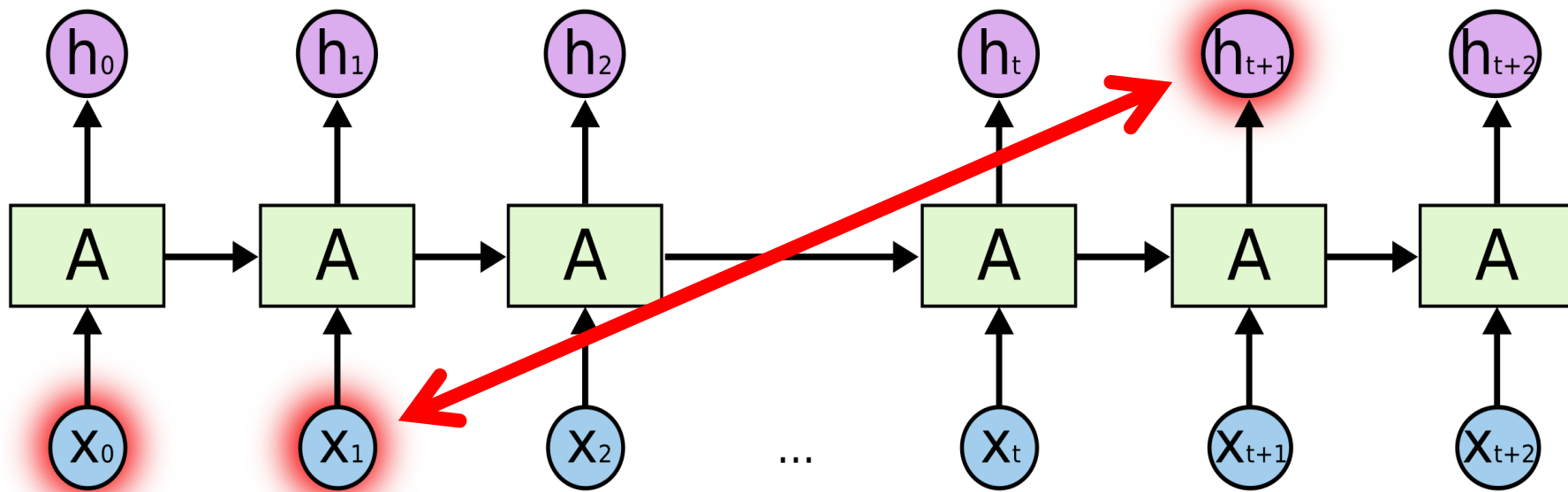
Hidden state h_t summarizes information on the history of the input signal up to time t



UNIVERSITÀ DI PISA

Learning Long-Term Dependencies is Difficult

When the time gap between the observation and the state grows there is little residual information of the input inside of the memory



What is the **cause**?

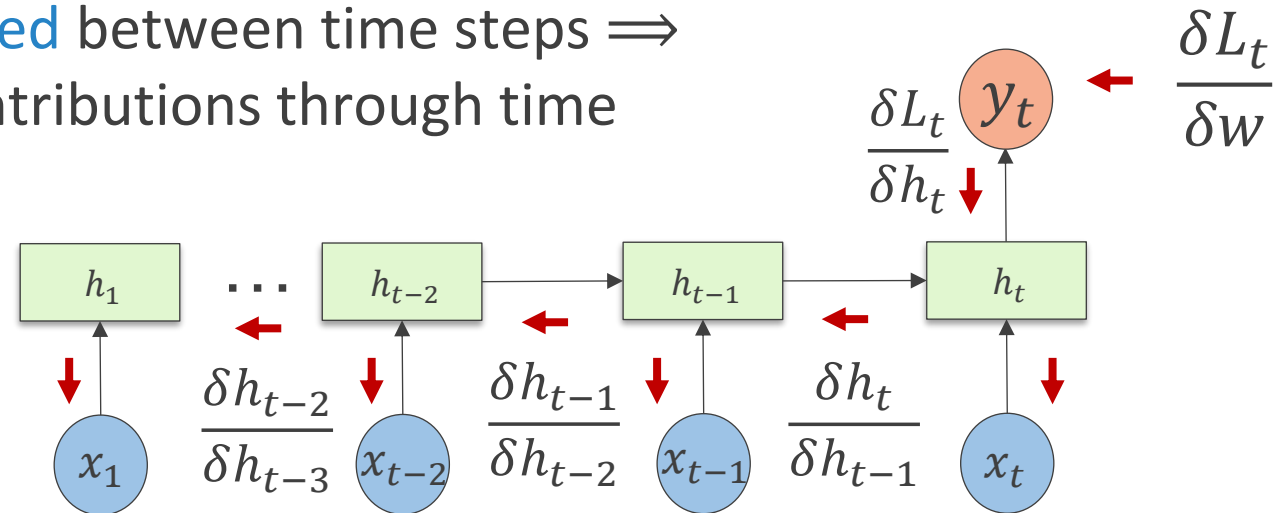
J. Hochreiter. Untersuchungen zu dynamischen neuronalen Netzen, TUM, 1991

Exploding/Vanishing Gradient

Short story: Gradients propagated over many stages tend to

- Vanish (often) \Rightarrow No learning
- Explode (rarely) \Rightarrow Instability and oscillations

Weights are shared between time steps \Rightarrow
sum gradient contributions through time

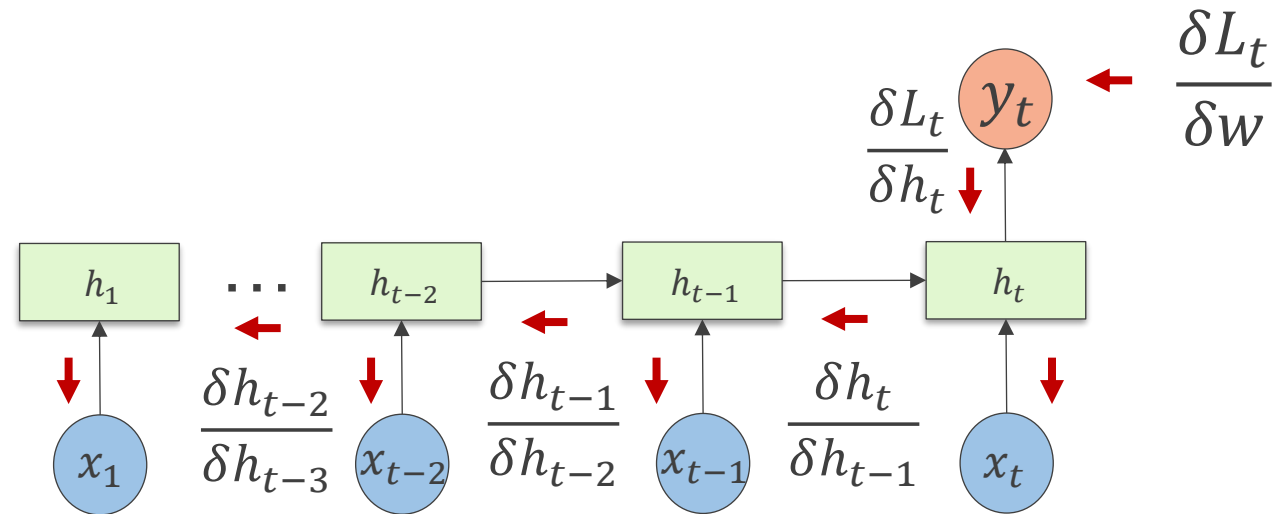


Bengio, Simard and Frasconi, Learning long-term dependencies with gradient descent is difficult. TNN, 1994



UNIVERSITÀ DI PISA

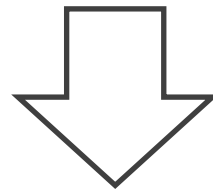
Backward propagation



A Closer Look at the Gradient

$$\frac{\delta L_t}{\delta \mathbf{W}} = \sum_{k=1}^t \frac{\delta L_t}{\delta h_t} \boxed{\frac{\delta h_t}{\delta h_k}} \frac{\delta h_k}{\delta \mathbf{W}}$$

→ This is a parameter **matrix**
 ⇒ we have a **Jacobian**



→ Inside here you have **chain rule**

$$\frac{\delta h_t}{\delta h_k} = \frac{\delta h_t}{\delta h_{t-1}} \times \frac{\delta h_{t-1}}{\delta h_{t-2}} \times \dots \times \frac{\delta h_{k+1}}{\delta h_k}$$

$$\frac{\delta L_t}{\delta \mathbf{W}} = \sum_{k=1}^t \frac{\delta L_t}{\delta h_t} \boxed{\left(\prod_{l=k}^{t-1} \frac{\delta h_{l+1}}{\delta h_l} \right)} \frac{\delta h_k}{\delta \mathbf{W}}$$

The gradient is a **recursive product of hidden activation gradients** (Jacobian)



UNIVERSITÀ DI PISA

Bounding the Gradient (I)

Given $\mathbf{h}_l = \tanh(\mathbf{W}_{hl}\mathbf{h}_{l-1} + \mathbf{W}_{in}\mathbf{x}_l)$ then $\frac{\delta \mathbf{h}_{l+1}}{\delta \mathbf{h}_l} = \mathbf{D}_{l+1}\mathbf{W}_{hl}^T$ where the **activation Jacobian** is

$$\mathbf{D}_{l+1} = \text{diag}(1 - \tanh^2(\mathbf{W}_{hl}\mathbf{h}_l + \mathbf{W}_{in}\mathbf{x}_{l+1}))$$

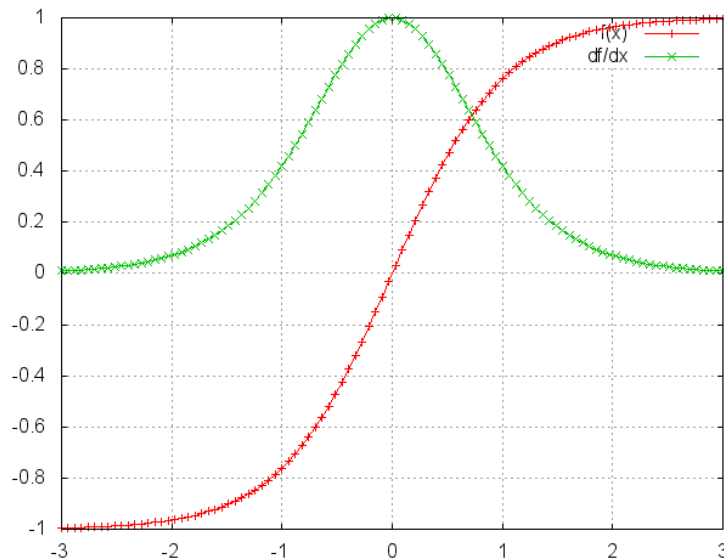
$$\frac{\delta L_t}{\delta \mathbf{h}_k} = \frac{\delta L_t}{\delta \mathbf{h}_t} \left(\prod_{l=k}^{t-1} \frac{\delta \mathbf{h}_{l+1}}{\delta \mathbf{h}_l} \right) = \frac{\delta L_t}{\delta \mathbf{h}_t} \prod_{l=k}^{t-1} \mathbf{D}_{l+1}\mathbf{W}_{hl}^T$$

We are interested in the gradient magnitude $\left\| \frac{\delta L_t}{\delta \mathbf{h}_k} \right\|$

Bounding the Gradient (II)

$$\left\| \frac{\delta L_t}{\delta \mathbf{h}_k} \right\| = \left\| \frac{\delta L_t}{\delta \mathbf{h}_t} \prod_{l=k}^{t-1} \mathbf{D}_{l+1} \mathbf{W}_{hl}^T \right\| \leq \left\| \frac{\delta L_t}{\delta \mathbf{h}_t} \right\| \prod_{l=k}^{t-1} \|\mathbf{D}_{l+1}\| \|\mathbf{W}_{hl}^T\| \approx \left\| \frac{\delta L_t}{\delta \mathbf{h}_t} \right\| \|\mathbf{D}\|^{k-1} \|\mathbf{W}_h^T\|^{k-1} \approx \left\| \frac{\delta L_t}{\delta \mathbf{h}_t} \right\| \sigma(\mathbf{D})^{k-1} \sigma(\mathbf{W}_h^T)^{k-1}$$

Bounded by the spectral radius σ for some norm and k large enough



Can shrink to zero or increase exponentially depending on the spectral properties

- $\sigma < 1 \Rightarrow$ **vanishish**
- $\sigma > 1 \Rightarrow$ **exploding**

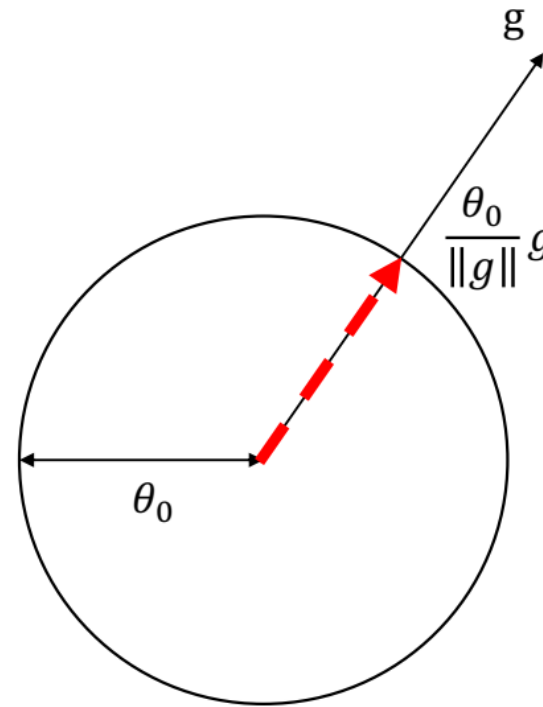


UNIVERSITÀ DI PISA

Gradient Clipping for Exploding Gradients

- Take $g = \frac{\delta L_t}{\delta W}$
- If $\|g\| > \theta_0$ then $g = \frac{\theta_0}{\|g\|} g$

Rescaling does not work for gradient vanish as total gradient is a sum of time dependent gradients (preserving relative contribution from each time makes it exponentially decay)



$$\frac{\delta L_t}{\delta W} = \sum_{k=1}^t \frac{\delta L_t}{\delta h_t} \frac{\delta h_t}{\delta h_k} \dots$$

Recap - Simple RNN

1. Expressive and general model.
2. Hard to train due to gradient propagation issues.
3. Relatively fast at inference time. However, the recurrence limits the parallelization opportunities.

The next models we will see will solve (2) and (3). We will also see how to improve (1).



UNIVERSITÀ DI PISA

Constant Error Propagation

- Solution seems to be having the **Jacobian** where $\sigma = 1$
 - Change the activation function
 - Constrain the recurrent weights matrix

$$\frac{\delta \mathbf{h}_{l+1}}{\delta \mathbf{h}_l} = \mathbf{D}_{l+1} \mathbf{W}_h^T$$



UNIVERSITÀ DI PISA

Activation Function

- Popular choices (sigmoid, tanh) are always contractive ($\sigma < 1$)
- Alternatives: modReLU
- Much simpler alternative: no activation function (identity)

Tanh activation

$$\mathbf{h}_{l+1} = \tanh(\mathbf{W}_h^T \mathbf{h}_l + \mathbf{W} \mathbf{x}_{l+1})$$

$$\frac{\delta \mathbf{h}_{l+1}}{\delta \mathbf{h}_l} = \mathbf{D}_{l+1} \mathbf{W}_h^T$$

Linear activation

$$\mathbf{h}_{l+1} = \mathbf{W}_h^T \mathbf{h}_l + \mathbf{W} \mathbf{x}_{l+1}$$

$$\frac{\delta \mathbf{h}_{l+1}}{\delta \mathbf{h}_l} = \mathbf{I} \mathbf{W}_h^T = \mathbf{W}_h^T$$



UNIVERSITÀ DI PISA

Recurrent Weights

- It is possible to achieve $\sigma = 1$
 - Orthogonal matrices: $\mathbf{W}^T \mathbf{W} = \mathbf{I}$
 - Unitary matrices (complex domain): $\mathbf{W}^H \mathbf{W} = \mathbf{I}$
 - Identity matrix: $\mathbf{W} = \mathbf{I}$

Orthogonal Matrix + linear activation:

$$\mathbf{h}_{l+1} = \mathbf{W}_h^T \mathbf{h}_l + \mathbf{W} \mathbf{x}_{l+1}$$

$$\frac{\delta \mathbf{h}_{l+1}}{\delta \mathbf{h}_l} = \mathbf{I} \mathbf{W}_h^T = \mathbf{W}_h^T$$

$$\left\| \frac{\delta \mathbf{h}_{l+1}}{\delta \mathbf{h}_l} \right\| = \left\| \mathbf{I} \mathbf{W}_h^T \right\| = 1$$

Constant Error Proagation

- Identity activation function
- Identity weight matrix

$$\mathbf{h}_t = \mathbf{h}_{t-1} + \hat{c}(\mathbf{x}_t)$$

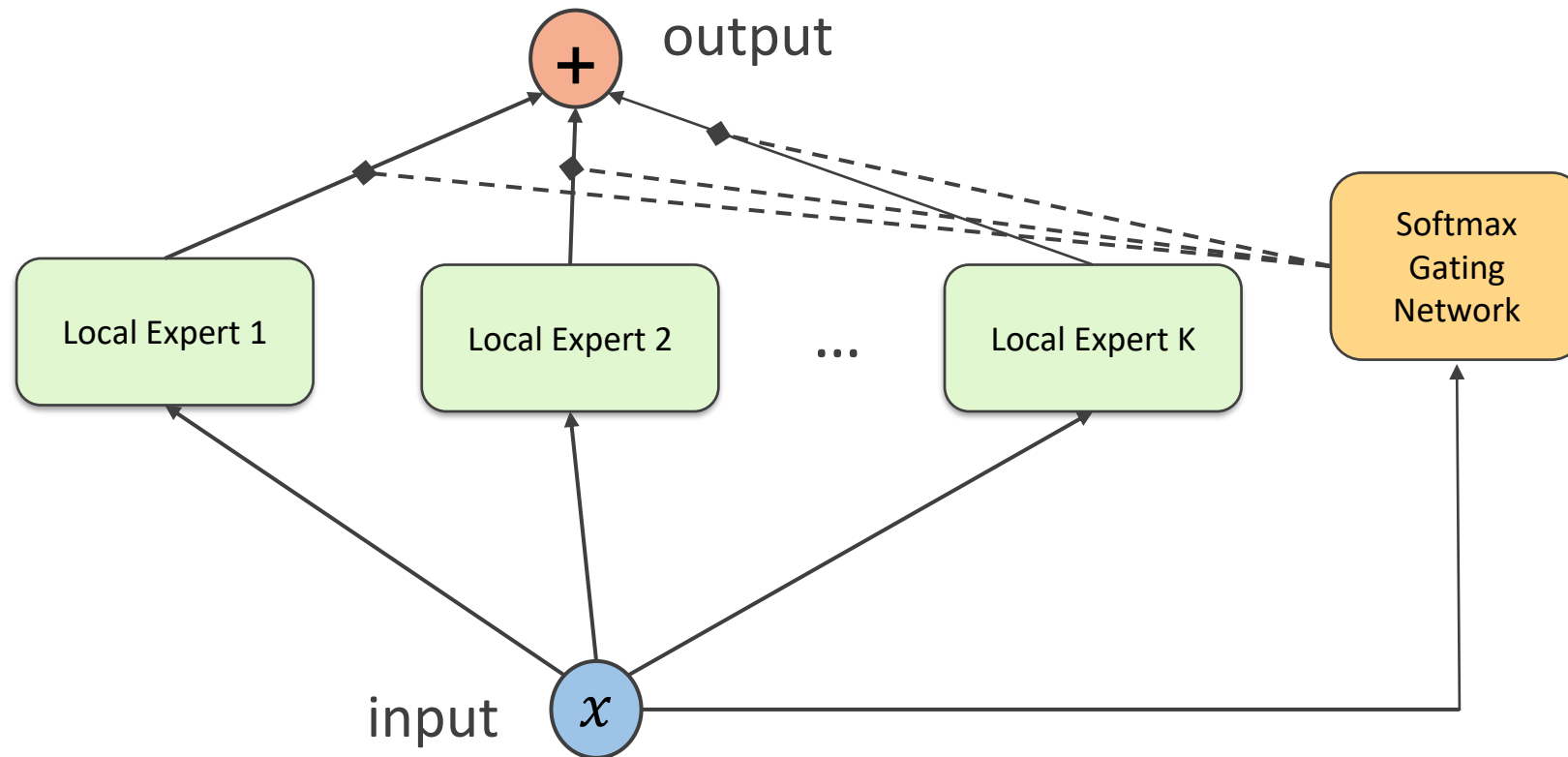
*Has the desired spectral properties but does not work in practice as it quickly **saturates memory** (e.g. with replicated/non-useful inputs and states). We want to be able to “control the forgetting”.*



UNIVERSITÀ DI PISA

Gating Units

Mixture of experts \Rightarrow the origin of gating



Jacobs et al (1991), Adaptive Mixtures of Local Experts, ...

Forget gate

Constant Error Carousel (CEC)

- Identity activation function
- Identity weight matrix

$$\mathbf{h}_t = \mathbf{h}_{t-1} + \hat{c}(\mathbf{x}_t)$$

- No forgetting
- Hidden state saturation

CEC + forget gate

- CEC
- Forget gate to “soft reset” units

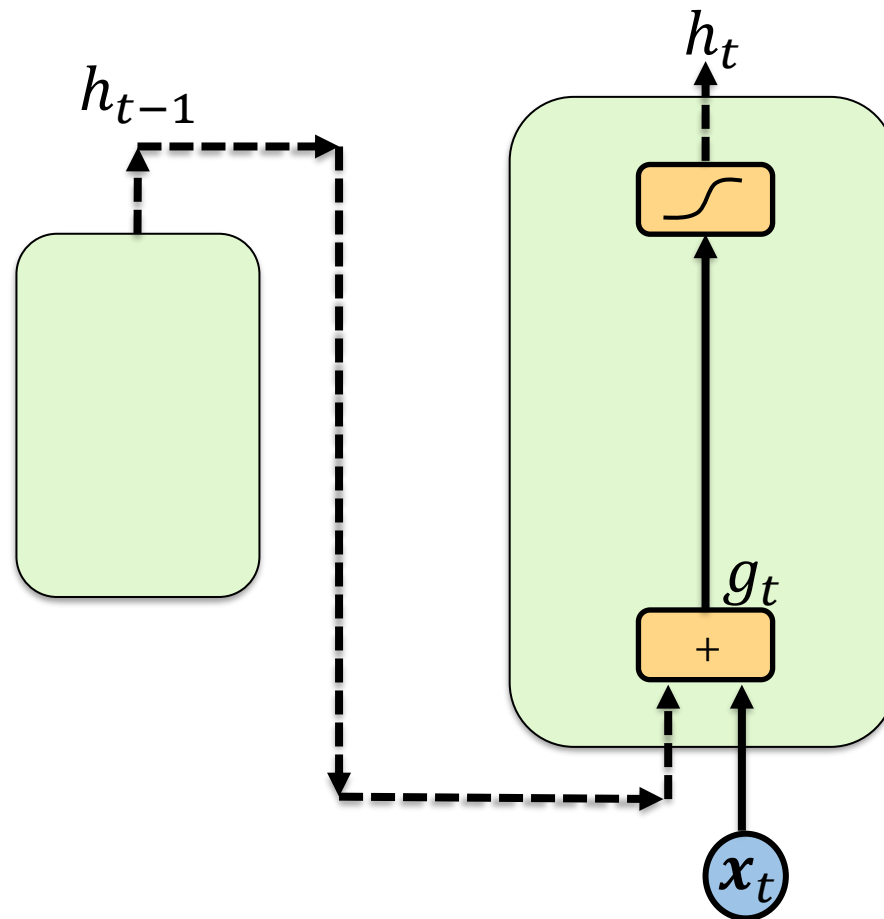
$$\mathbf{f}_t = \sigma(\mathbf{W}_{fh}\mathbf{h}_{t-1} + \mathbf{W}_{fx}\mathbf{x}_t + \mathbf{b}_f)$$

$$\mathbf{h}_t = \mathbf{f}_t \odot \mathbf{h}_{t-1} + \hat{c}(\mathbf{x}_t)$$

- Adaptively forgets the past
- Avoid saturation
- No guarantees about constant propagation



Long-Short Term Memory (LSTM) Cell



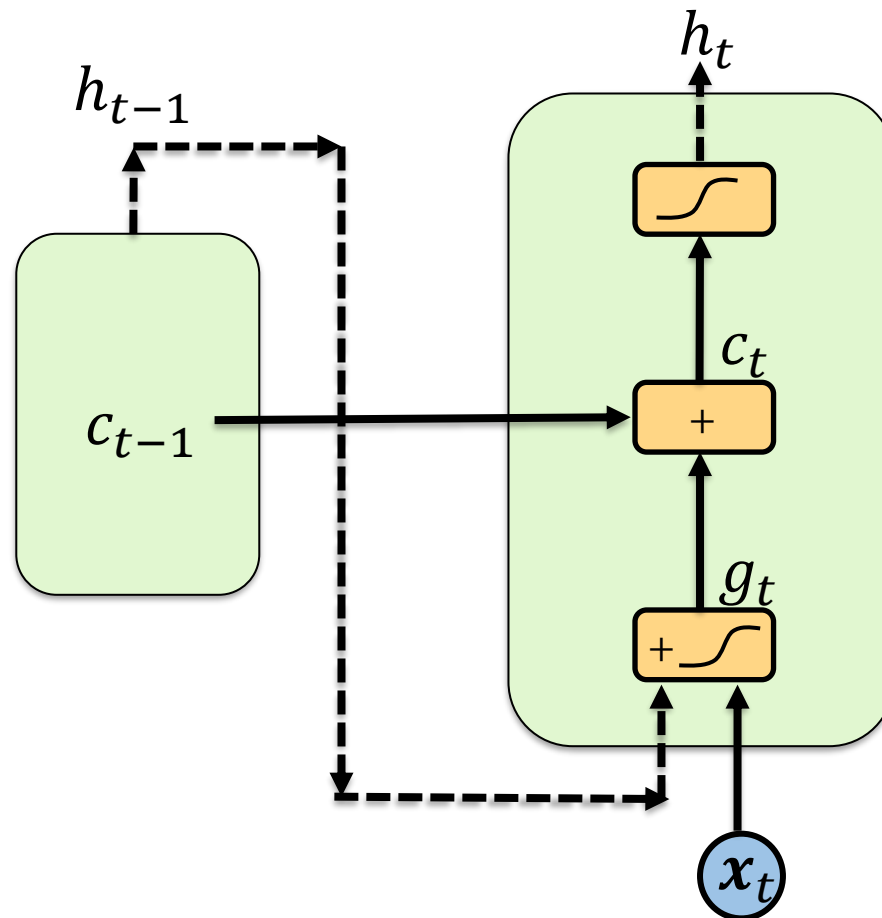
Let's start from the vanilla RNN unit

S. Hochreiter, J. Schmidhuber, Long short-term memory". Neural Computation, Neural Comp. 1997



UNIVERSITÀ DI PISA

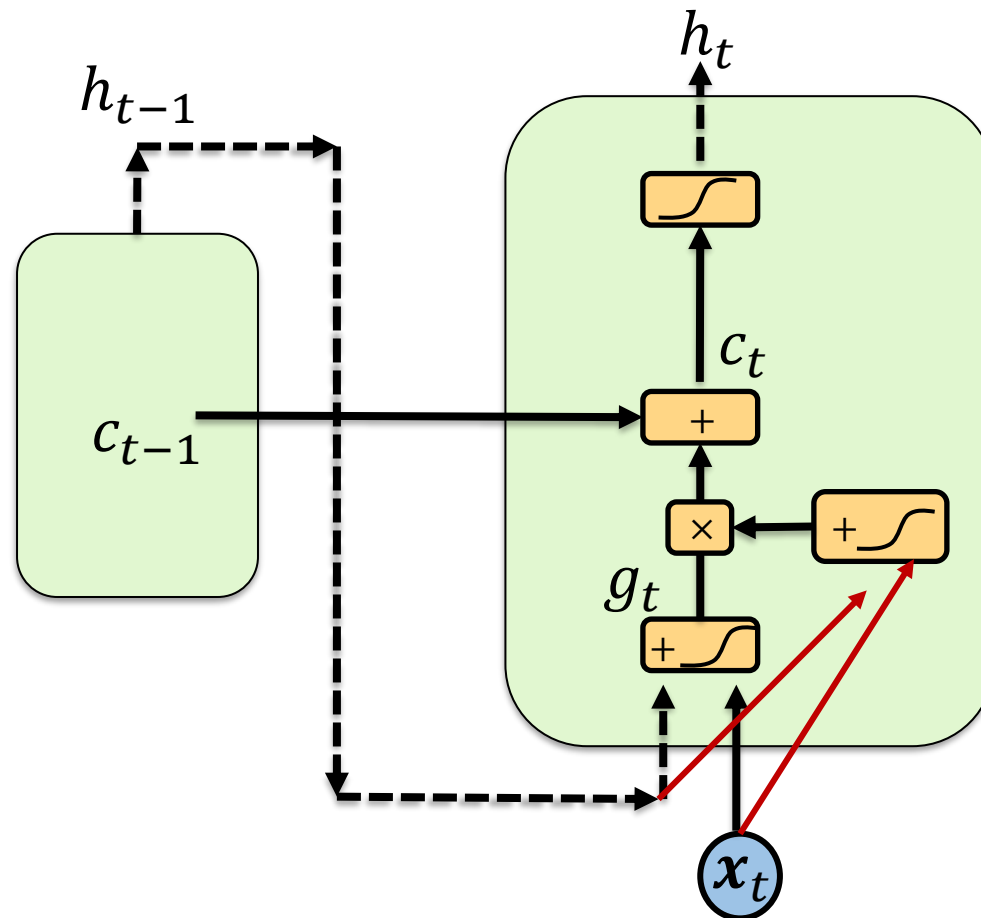
LSTM Design – Step 1



Introduce a linear/
identity memory c_t

Combines past **internal state** c_{t-1} with current input x_t

LSTM Design – Step 2 (Gates)



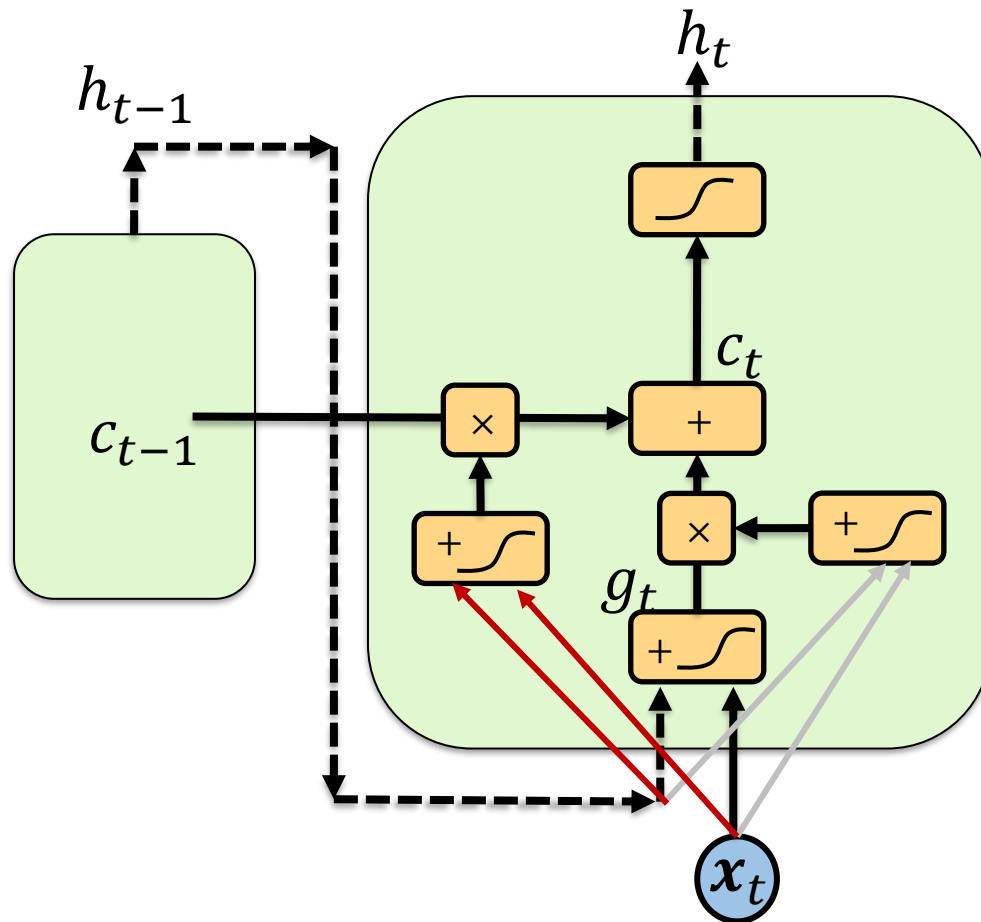
Input gate

Controls how inputs contribute to the internal state

$$I_t(x_t, h_{t-1})$$

Logistic sigmoid

LSTM Design – Step 2 (Gates)



Forget gate

Controls how past internal state c_{t-1} contributes to c_t

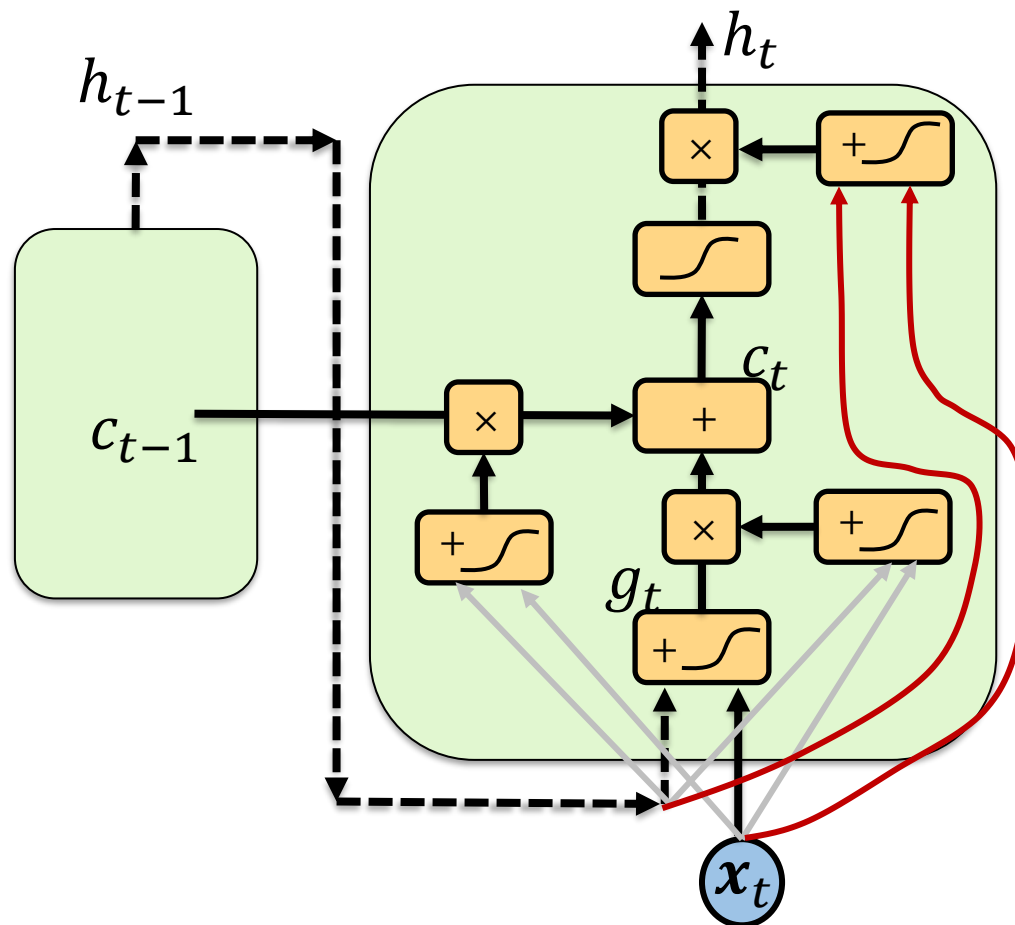
$$F_t(x_t, h_{t-1})$$

Logistic sigmoid



UNIVERSITÀ DI PISA

LSTM Design – Step 2 (Gates)



Output gate

Controls what part of the internal state is propagated out of the cell

$$O_t(x_t, h_{t-1})$$

Logistic sigmoid



UNIVERSITÀ DI PISA

LSTM in Equations

1) Compute activation of input and forget gates

$$I_t = \sigma(W_{Ih}h_{t-1} + W_{Iin}x_t + b_I)$$
$$F_t = \sigma(W_{Fh}h_{t-1} + W_{Fin}x_t + b_F)$$

2) Compute input potential and internal state

$$g_t = \tanh(W_h h_{t-1} + W_{in}x_t + b_h)$$
$$c_t = F_t \odot c_{t-1} + I_t \odot g_t$$

\odot element-wise multiplication

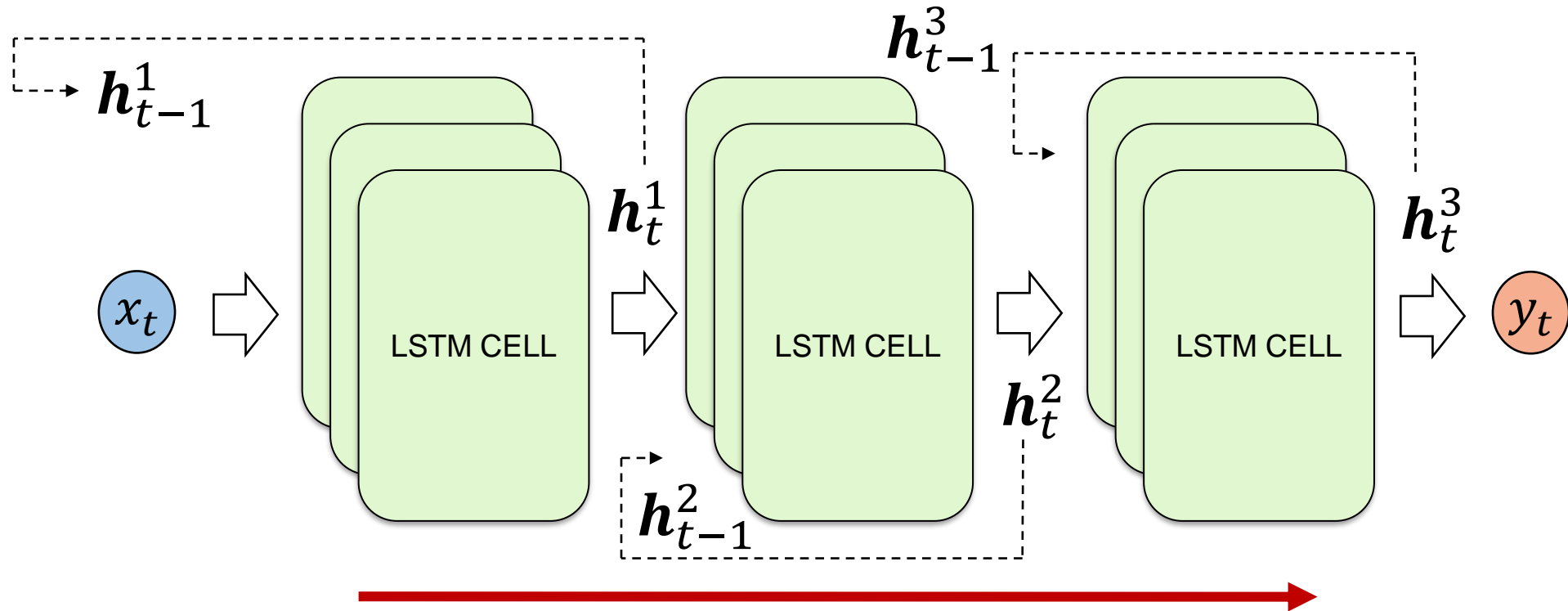
3) Compute output gate and output state

$$O_t = \sigma(W_{Oh}h_{t-1} + W_{Oin}x_t + b_O)$$
$$h_t = O_t \odot \tanh(c_t)$$



UNIVERSITÀ DI PISA

Deep LSTM



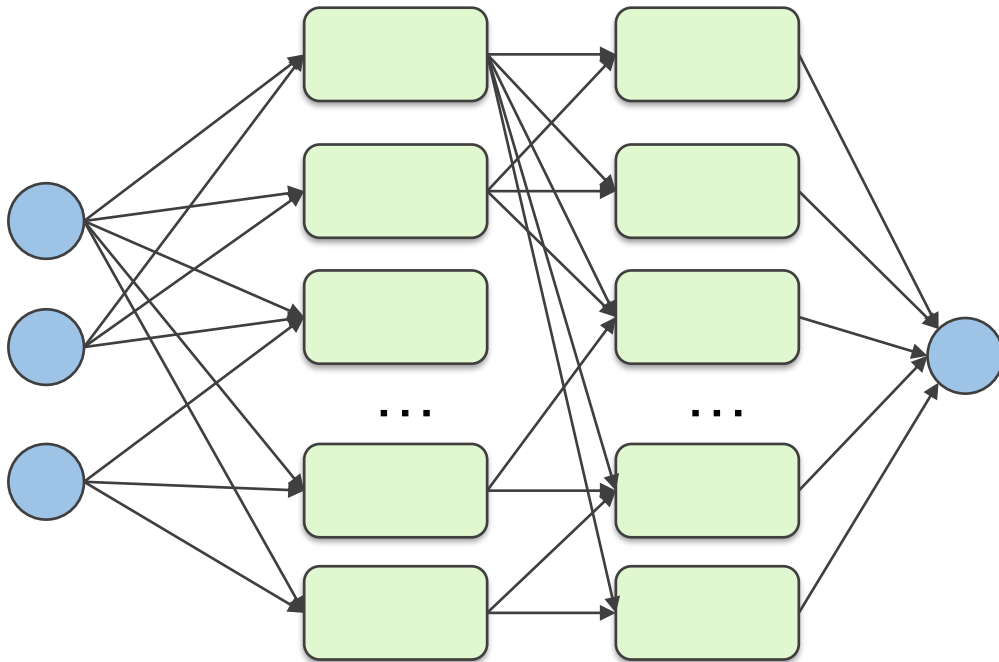
LSTM layers extract information at **increasing levels of abstraction** (enlarging context)

Training LSTM

- Original LSTM training algorithm was a mixture of RTRL and BPTT
 - BPTT on internal state gradient
 - RTRL-like truncation on other recurrent connections
 - No exact gradient calculation!
- All current LSTM implementation use full BPTT training
 - Introduced by Graves and Schmidhuber in 2005
 - Typically use Adam or RMSProp optimizer

Regularizing LSTM - Dropout

Randomly disconnect units from the network during training



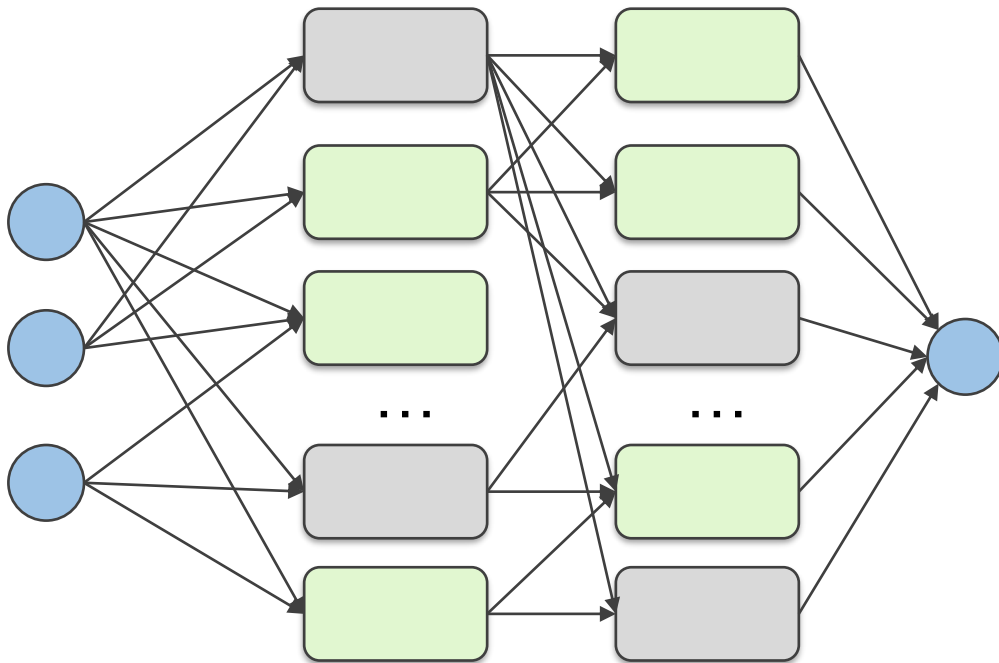
N. Srivastava et al, Dropout: A Simple Way to Prevent Neural Networks from Overfitting, JLMR 2014



UNIVERSITÀ DI PISA

Regularizing LSTM - Dropout

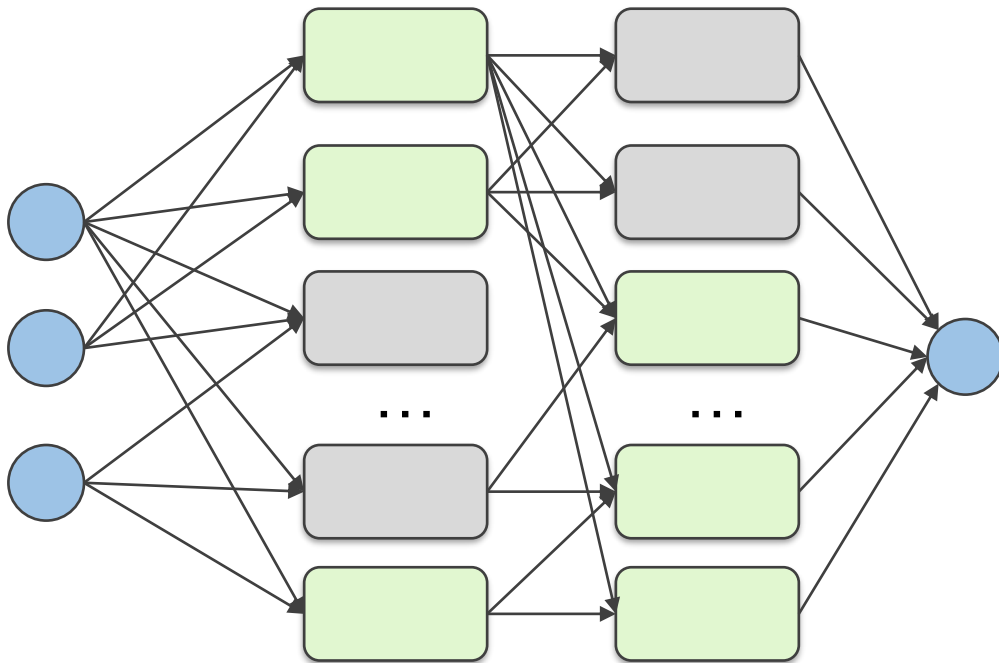
Randomly disconnect units from the network during training



N. Srivastava et al, Dropout: A Simple Way to Prevent Neural Networks from Overfitting, JLMR 2014

Regularizing LSTM - Dropout

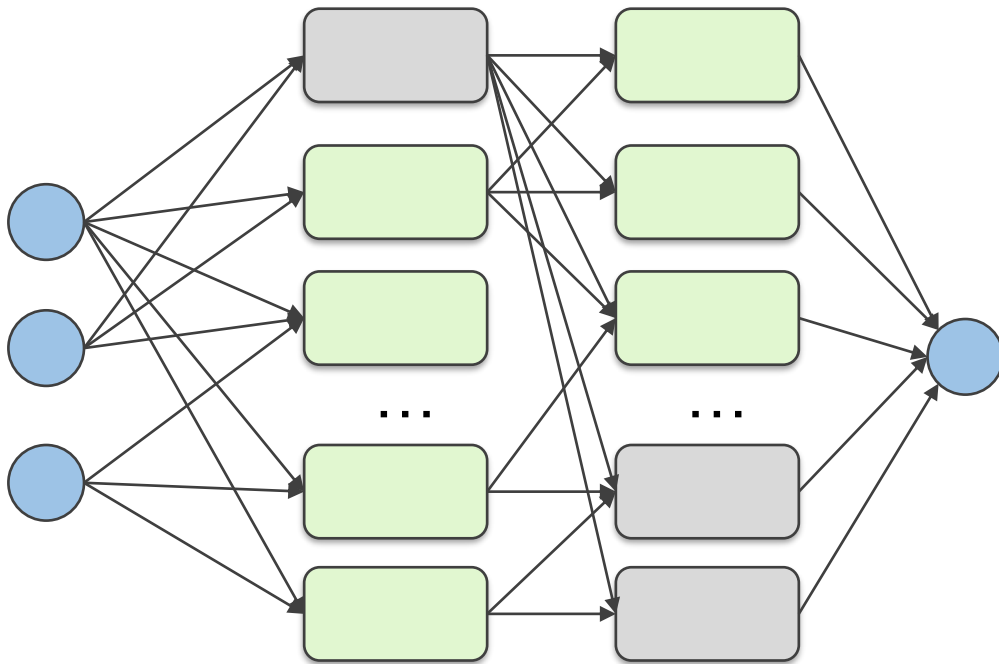
Randomly disconnect units from the network during training



N. Srivastava et al, Dropout: A Simple Way to Prevent Neural Networks from Overfitting, JLMR 2014

Regularizing LSTM - Dropout

Randomly disconnect units from the network during training

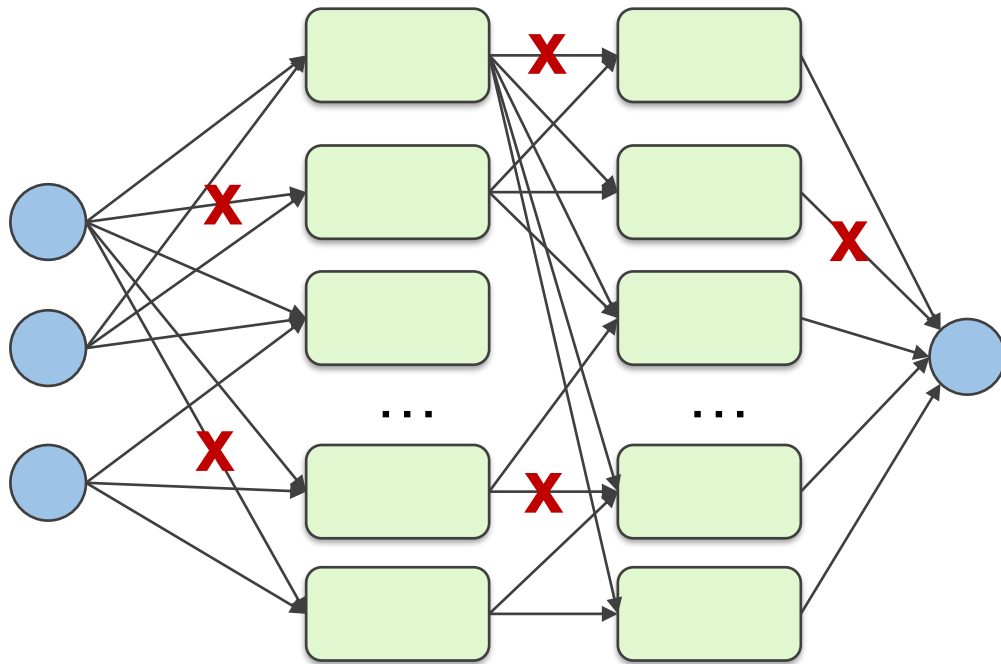


- Regulated by unit **dropping hyperparameter**
- Prevents unit **coadaptation**
- Committee machine effect
- Need to adapt **prediction phase**
- Drop units for the whole sequence!

N. Srivastava et al, Dropout: A Simple Way to Prevent Neural Networks from Overfitting, JLMR 2014

Regularizing LSTM - Dropout

Randomly disconnect units from the network during training



- Regulated by unit **dropping hyperparameter**
- Prevents unit **coadaptation**
- Committee machine effect
- Need to adapt **prediction phase**
- Drop units for the whole sequence!

You can also **drop single connections** (dropconnect)

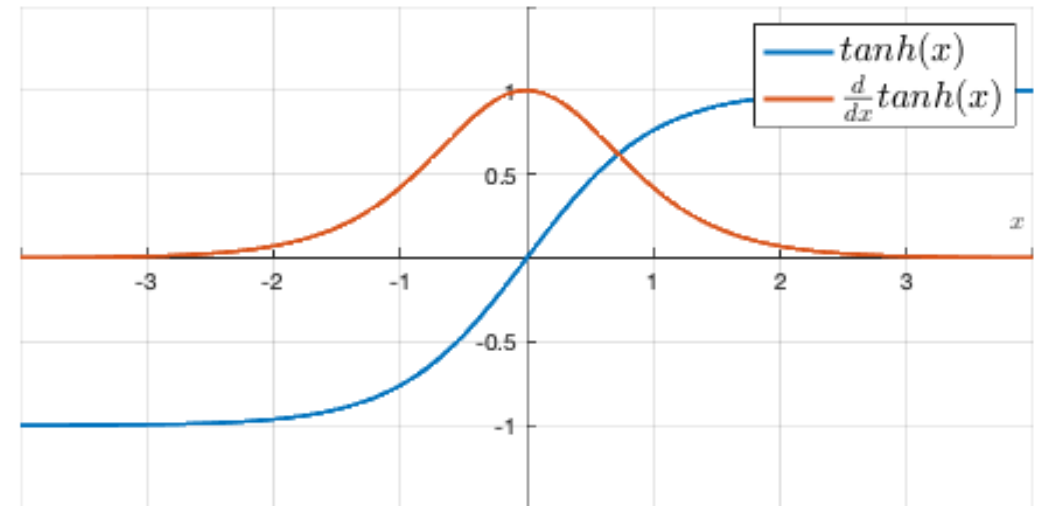
Activity Regularization

- Penalize the model's activations
- High or low activations saturate the tanh function (zero gradient)
- L2 activity regularization:

$$\alpha ||m \odot h^t||_2^2$$

- Temporal activity regularization:

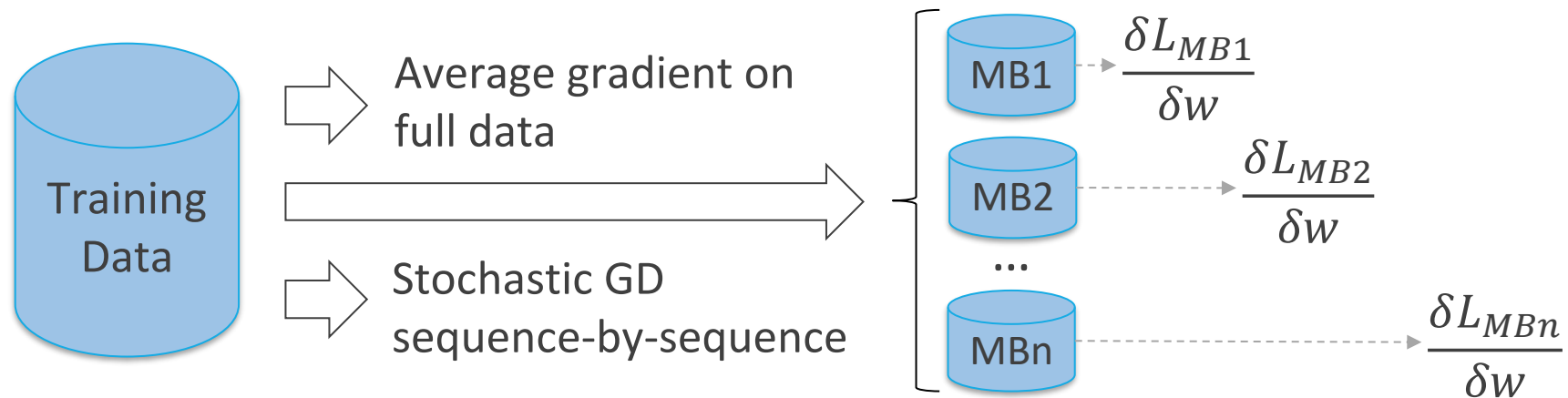
$$\beta ||h^t - h^{t+1}||_2^2$$



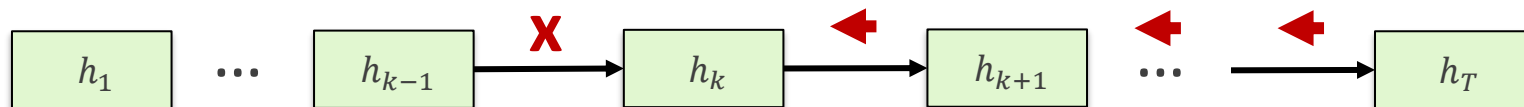
Revisiting Activation Regularization for Language RNNs, <https://arxiv.org/pdf/1708.01009.pdf>

Practicalities – Minibatch and Truncated BP

Minibatch (MB)



Truncated gradient propagation



Gated Recurrent Unit (GRU)

Reset acts directly on output state (no internal state and no output gate)

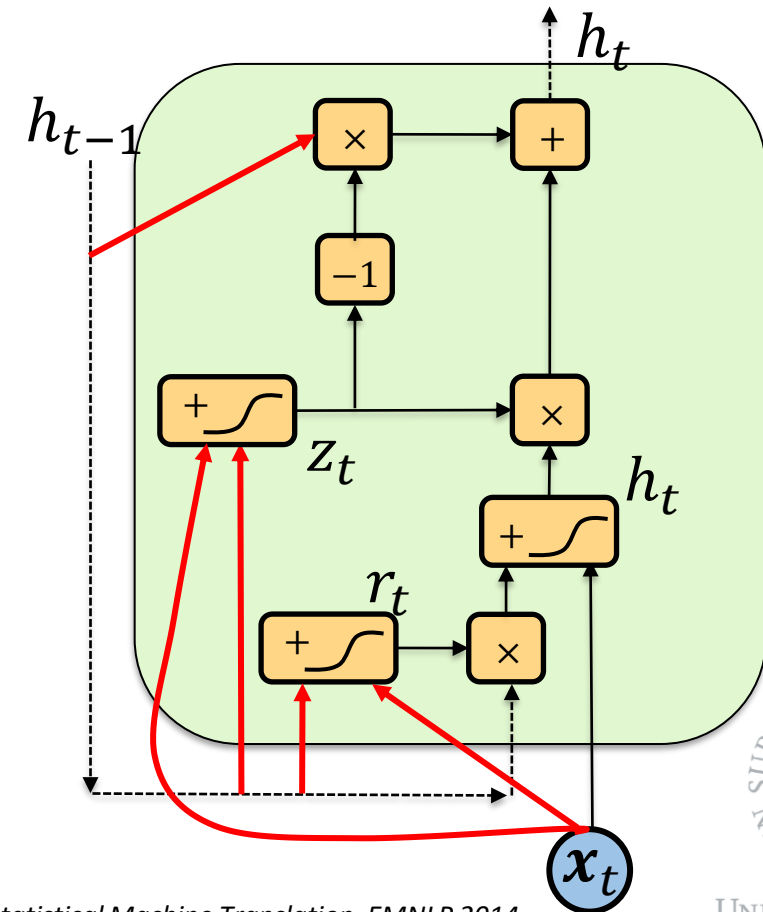
$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \mathbf{h}_t$$

$$\mathbf{h}_t = \tanh(\mathbf{W}_{hh}(\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{W}_{hin}\mathbf{x}_t + \mathbf{b}_h)$$

Reset and **update** gates when coupled act as input and forget gates

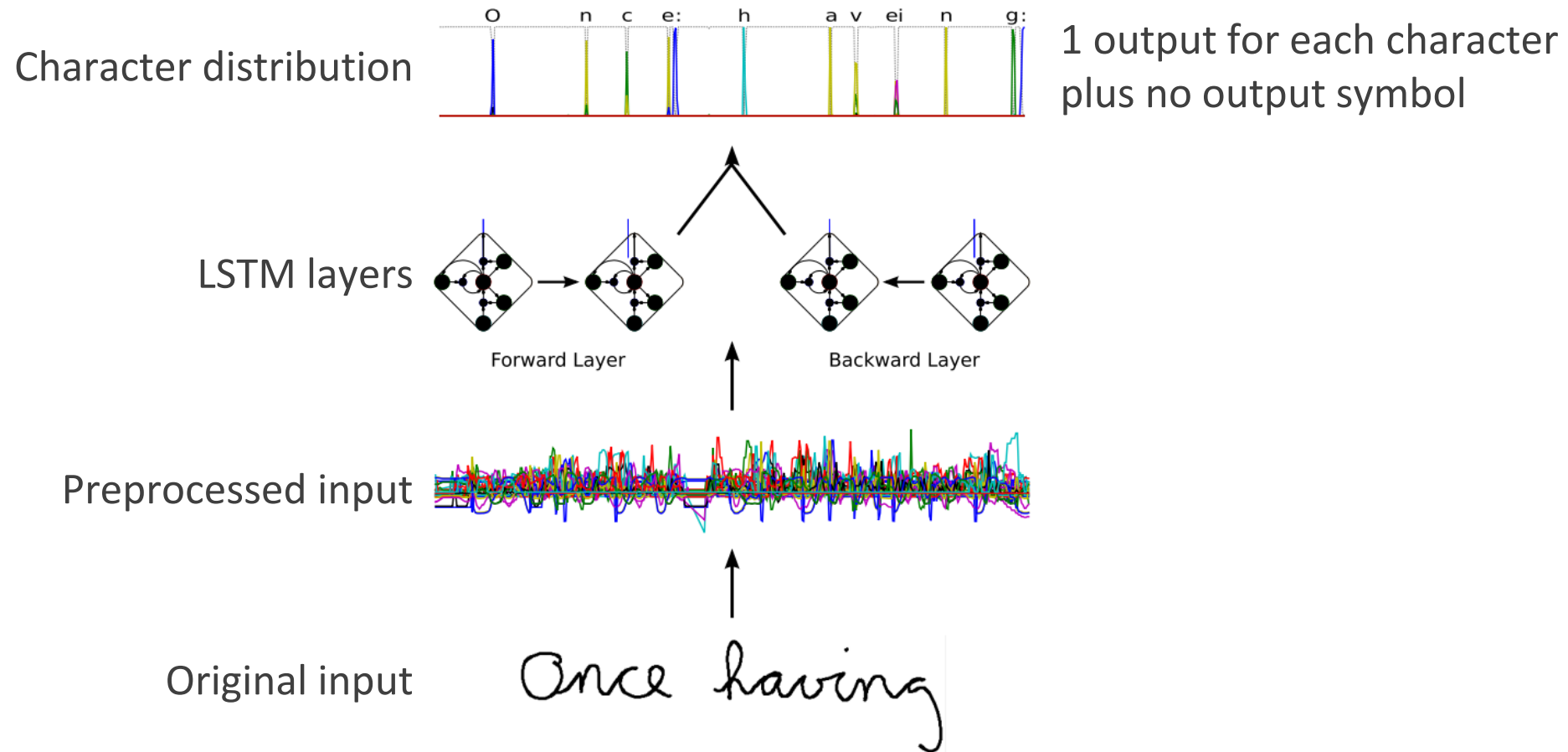
$$\mathbf{z}_t = \sigma(\mathbf{W}_{zh}\mathbf{h}_{t-1} + \mathbf{W}_{zin}\mathbf{x}_t + \mathbf{b}_z)$$

$$\mathbf{r}_t = \sigma(\mathbf{W}_{rh}\mathbf{h}_{t-1} + \mathbf{W}_{rin}\mathbf{x}_t + \mathbf{b}_r)$$



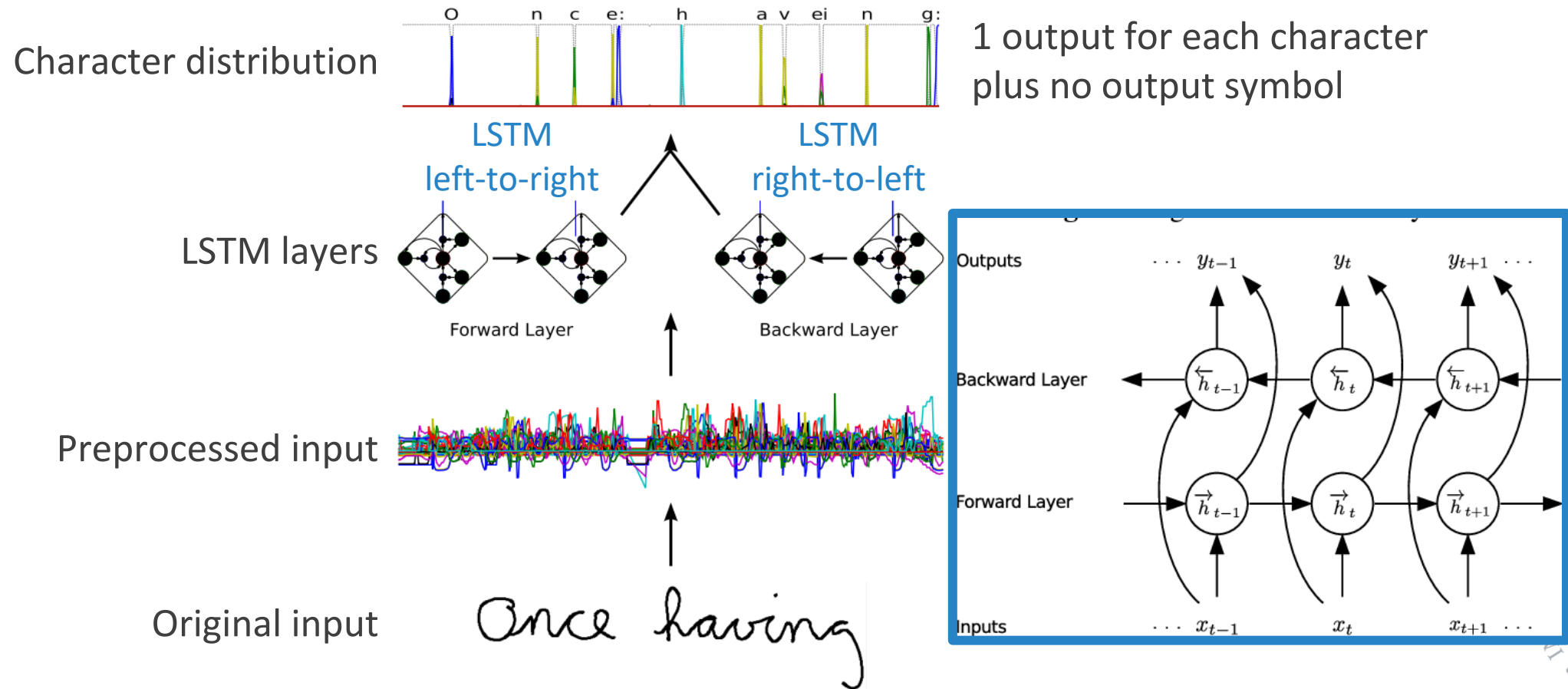
C. Kyunghyun et al, Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation, EMNLP 2014

Bidirectional LSTM – Character Recognition



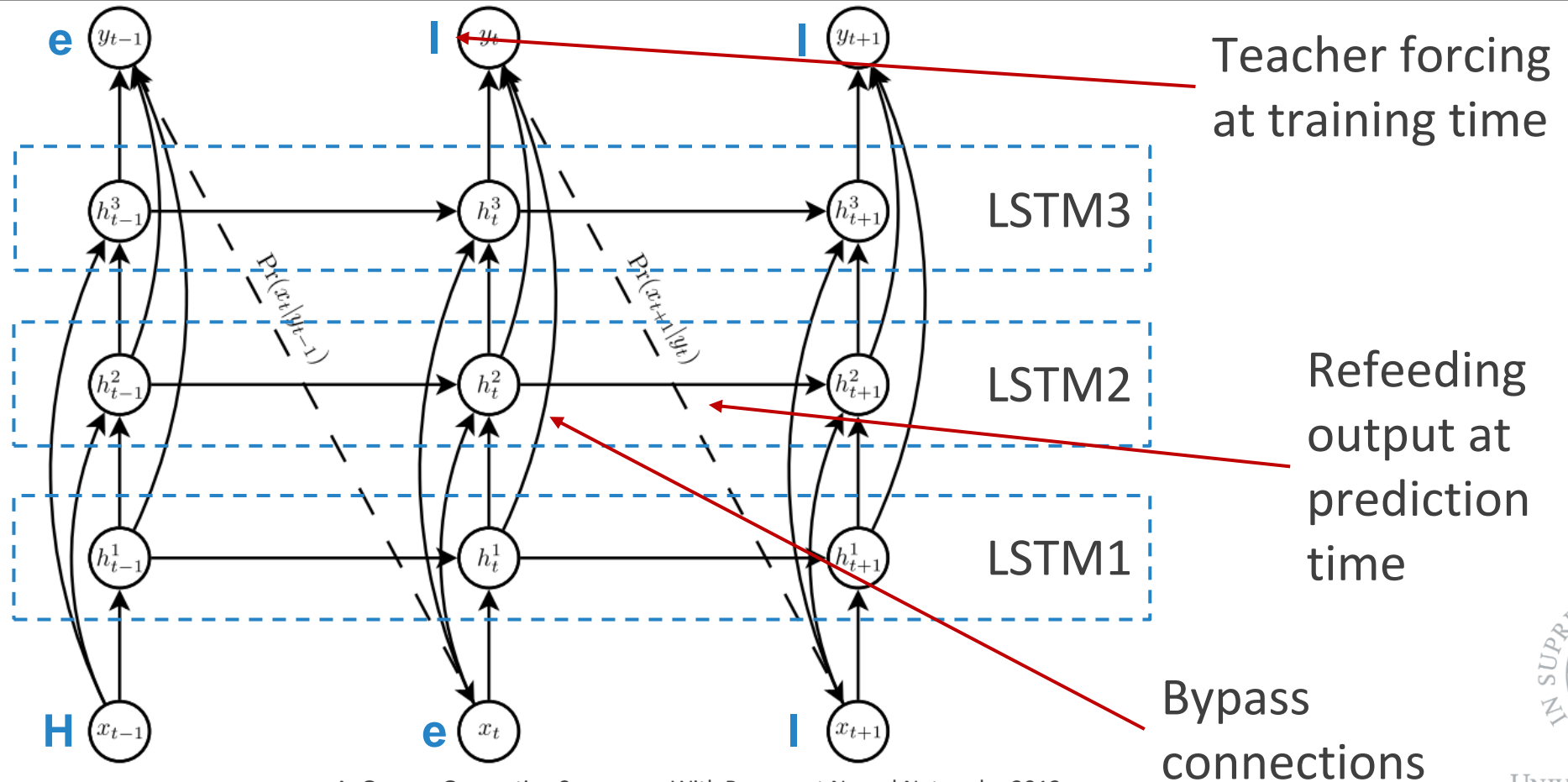
A. Graves, A novel connectionist system for unconstrained handwriting recognition, TPAMI 2009

Bidirectional LSTM – Character Recognition



A. Graves, A novel connectionist system for unconstrained handwriting recognition, TPAMI 2009

LSTM – (Primitive) Language Modeling



A. Graves, Generating Sequences With Recurrent Neural Networks, 2013

Character Generation Fun

Linux Kernel Code

```
/*  
 * If this error is set, we will need anything right after that BSD.  
 */  
static void action_new_function(struct s_stat_info *wb)  
{  
    unsigned long flags;  
    int lel_idx_bit = e->edd, *sys & ~((unsigned long) *FIRST_COMPAT);  
    buf[0] = 0xFFFFFFFF & (bit << 4);  
    min(inc, slist->bytes);  
    printk(KERN_WARNING "Memory allocated %02x/%02x, "  
        "original MLL instead\n"),  
        min(min(multi_run - s->len, max) * num_data_in),  
        frame_pos, sz + first_seg);  
    div_u64_w(val, inb_p);  
    spin_unlock(&disk->queue_lock);  
    mutex_unlock(&s->sock->mutex);  
    mutex_unlock(&func->mutex);  
    return disassemble(info->pending_bh);  
}
```

<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>



UNIVERSITÀ DI PISA

Generate Sad Jokes

A 3-LSTM layers neural network to generate English jokes character by character

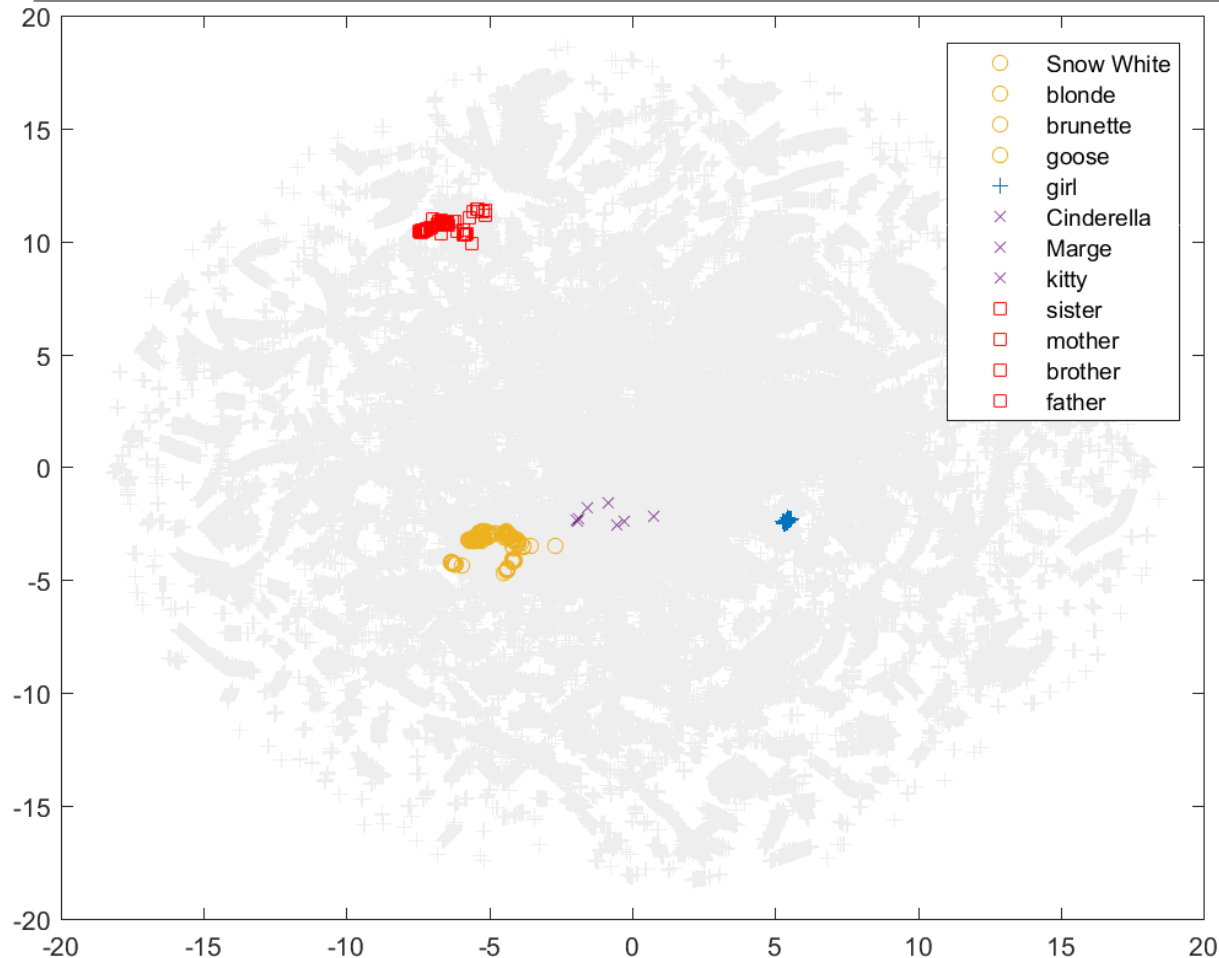
Why did the boy stop his homework?
Because they're bunny boo!



What do you get if you cross a famous California little boy with an elephant for players?
Market holes.

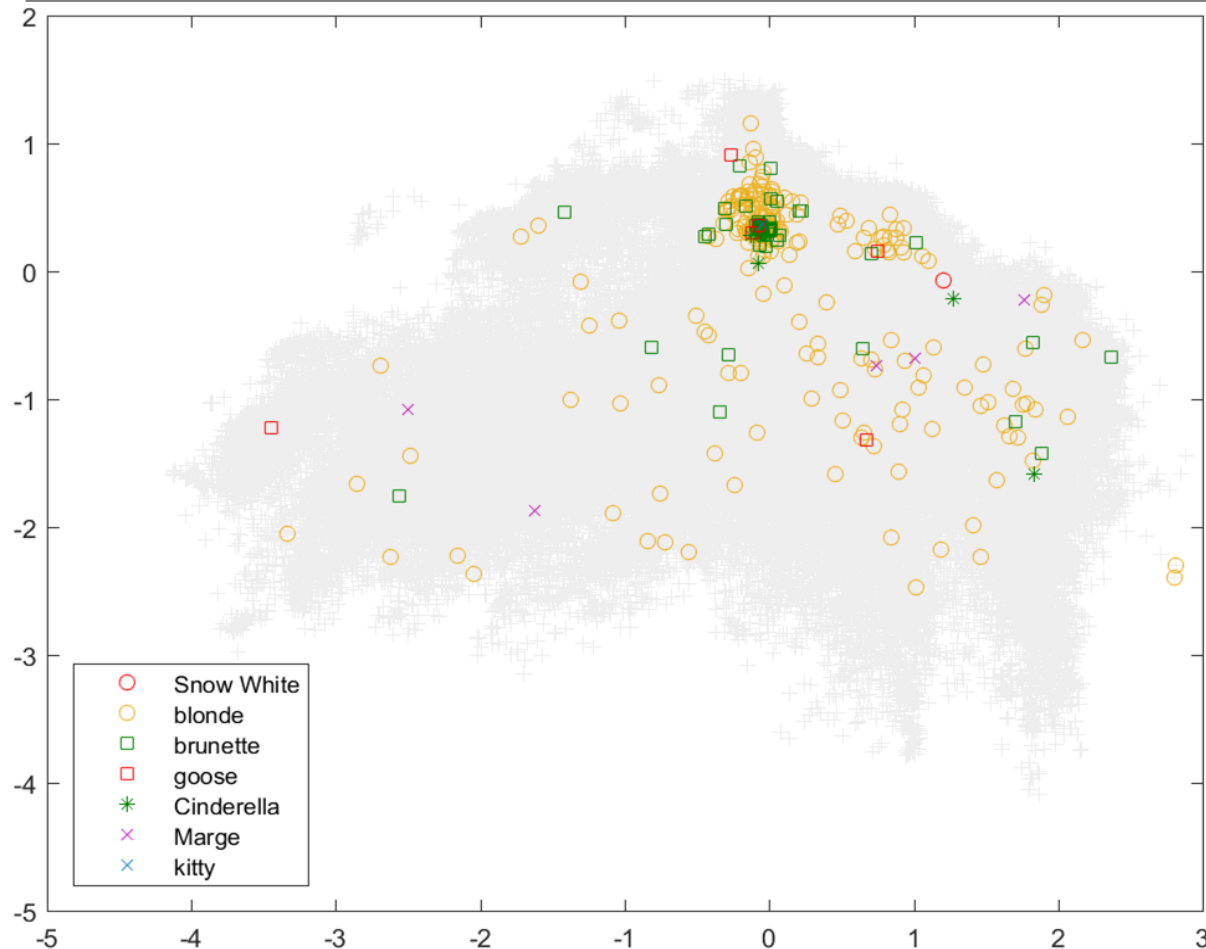
Q: Why did the death penis learn string?
A: Because he wanted to have some roasts case!

Understanding Memory Representation



At **Layer-3** neuron show
some form of **context**
induced representation of
subsequences (words)

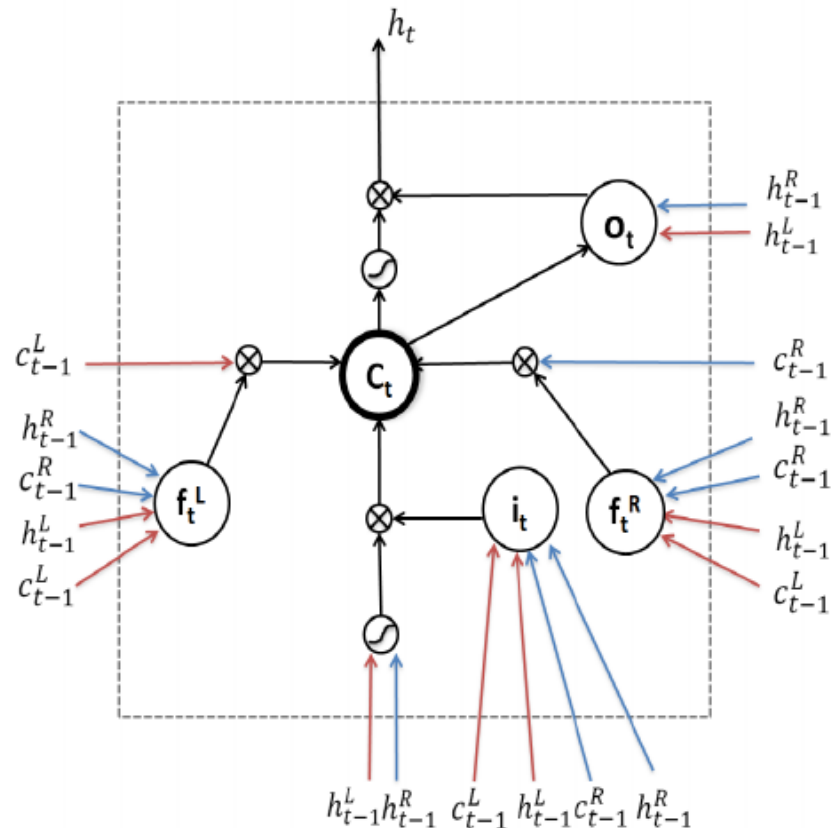
Understanding Memory Representation



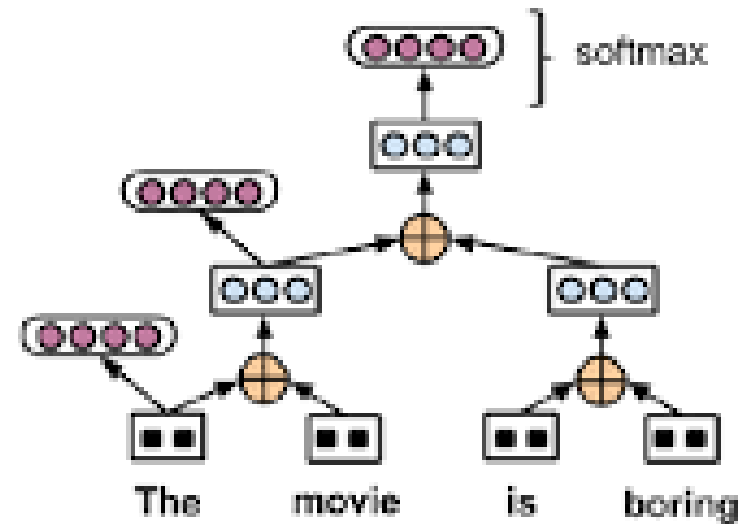
Neurons in **early recurrent layers** tend to organize according to **sequence suffix**

Recursive Gated Networks

Recursive LSTM cell for binary trees



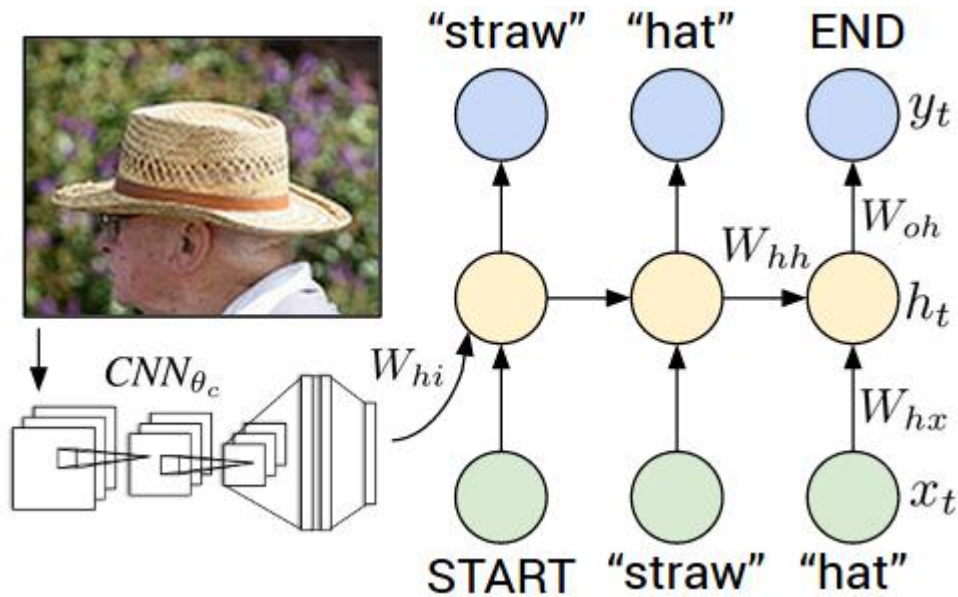
Unfolding on parse trees for sentiment analysis



UNIVERSITÀ DI PISA

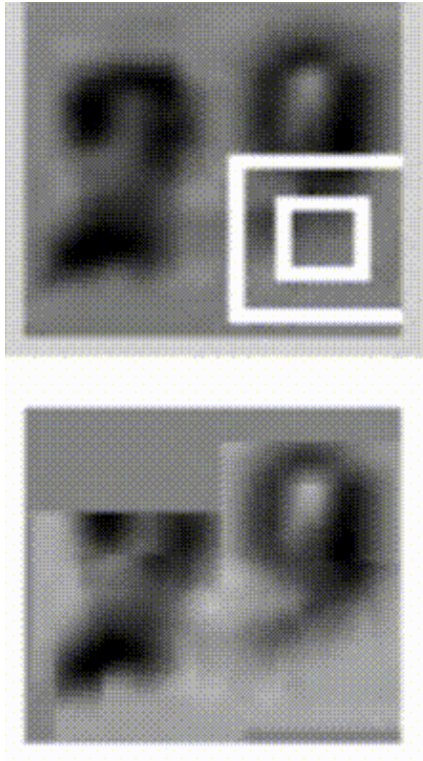
Differentiable Compositions

CNN-LSTM Composition for **image-to-sequence** (NeuralTalk)

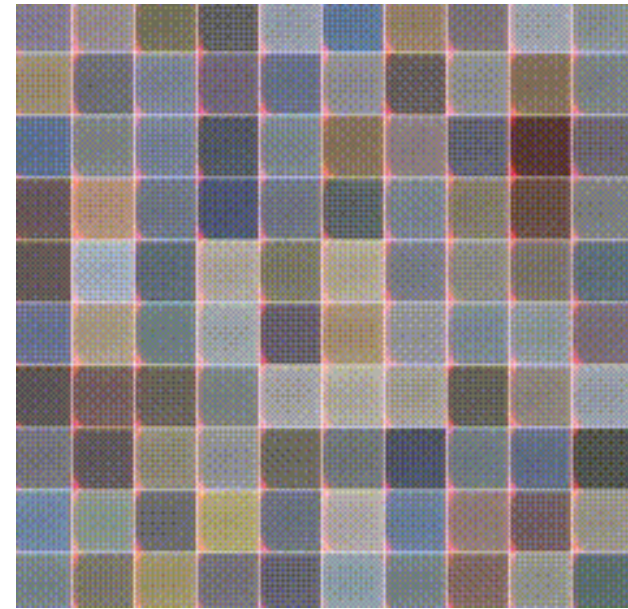


A. Karpathy and L. Fei-Fei, Deep Visual-Semantic Alignments for Generating Image Descriptions, CVPR 2015
<https://github.com/karpathy/neuraltalk2>

RNN – A Broader View



RNN are only for sequential/structured data?



a recurrent network *generates* images of digits by learning to sequentially add color to a canvas ([Gregor et al.](#))

an algorithm learns a recurrent network policy that steers its attention around an image; In particular, it learns to read out house numbers from left to right ([Ba et al.](#)).

RNN as **stateful** systems

Software

- Standard LSTM and GRU are available in all deep learning frameworks (Python et al) as well as in Matlab's Neural Network Toolbox
- If you want to play with one-element ahead sequence generation try out char-RNN implementations
 - <https://github.com/karpathy/char-rnn> (ORIGINAL)
 - <https://github.com/sherjilozair/char-rnn-tensorflow>
 - <https://github.com/crazydonkey200/tensorflow-char-rnn>
 - http://pytorch.org/tutorials/intermediate/char_rnn_generation_tutorial.html

Take Home Messages

- Learning **long-term dependencies** can be difficult due to gradient vanish/explosion
- Gated RNN solution
 - Gates are neurons whose **output is used to scale** another neuron's output
 - Use gates to determine what information can enter (or exit) the **internal state**
 - Training gated RNN non always straightforward
- Deep RNN can be used in **generative** mode
 - Can seed the network with **neural embeddings**
- Deep RNN as **stateful** and **differentiable** machines

Lecture Plan – Next Week

- Tue 15/04 - Seq2seq, Attention & Transformers
- Wed 16/04 - Coding I
- Thu 17/04 – Coding II
- Apr. 18 – Apr 28 – Spring Break (no lectures)
- Bonus track
 - AI Meets Psychiatry: fMRI-Based Multi-Disorder Diagnosis
 - Lecture by Elisa Ferrari at the AI for Health course
 - Tue 15/04/2025 h. 16.15-17.30 – Room L1