

# Deep Learning

## Lecture 8 – Sequence Models

Kumar Bipin

BE, MS, PhD (MMMTU, IISc, IIIT-Hyderabad)

Robotics, Computer Vision, Deep Learning, Machine Learning, System Software



# Agenda

**8.1** Recurrent Networks

**8.2** Recurrent Network Applications

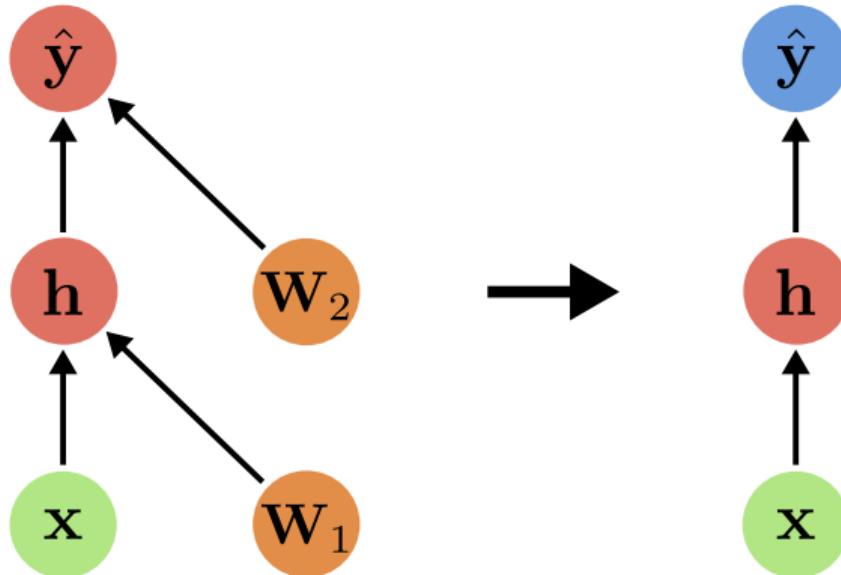
**8.3** Gated Recurrent Networks

**8.4** Autoregressive Models

# 8.1

## Recurrent Networks

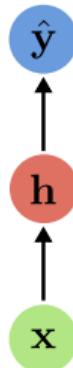
## Recap: Computation Graphs



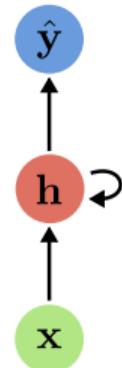
- ▶ Example of a **feedforward neural network** (here: MLP with 1 hidden layer)
- ▶ In the following, we will drop all parameter nodes and highlight outputs in blue

# Feedforward vs. Recurrent Neural Networks

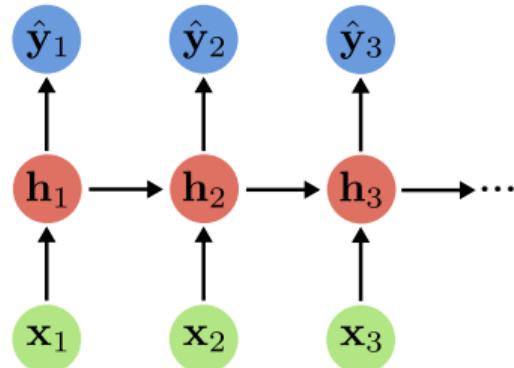
Feedforward  
Neural Network



Recurrent Neural Network (RNN)  
with feedback connection



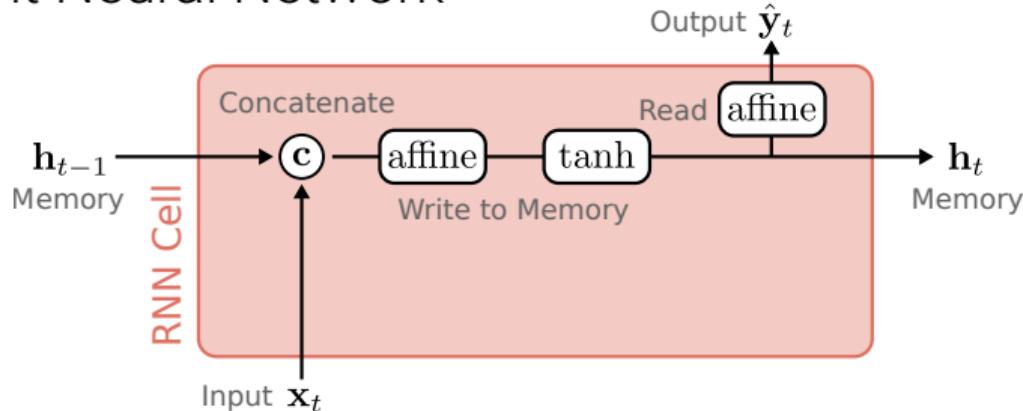
Recurrent Neural Network (RNN)  
unrolled over time (index = time t)



## Recurrent Neural Networks (RNNs):

- ▶ Core idea: update **hidden state  $h$**  based on input and previous hidden state using same update rule (same/shared parameters) at each **time step**
- ▶ Allows for processing **sequences of variable length**, not only fixed-sized vectors
- ▶ **Infinite memory:**  $h$  is function of all previous inputs (long-term dependencies)

# Basic Recurrent Neural Network



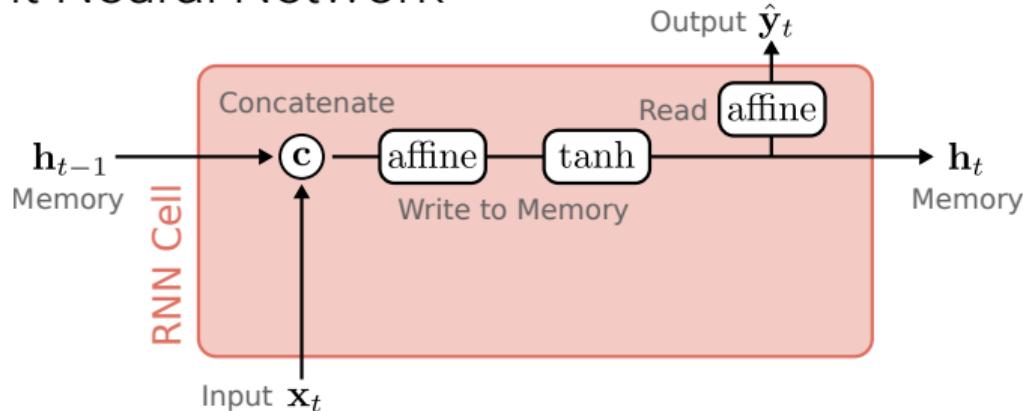
**General Form:**

$$\mathbf{h}_t = f_h(\mathbf{h}_{t-1}, \mathbf{x}_t)$$

$$\hat{\mathbf{y}}_t = f_y(\mathbf{h}_t)$$

- We use  $t$  as the **time index** (in feedforward networks we used  $i$  as layer index)
- General form does not specify the form of the hidden and output mappings
- Important:  $f_h$  and  $f_y$  **do not change over time**, unlike in layers of feedforward net

# Basic Recurrent Neural Network



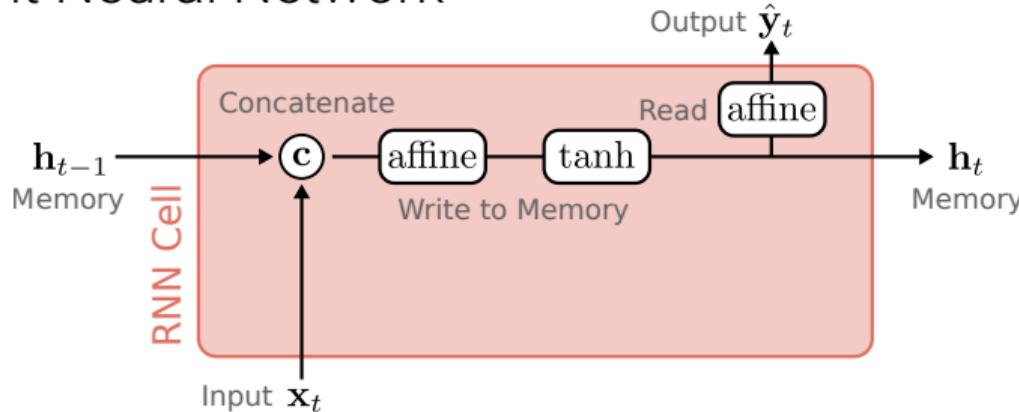
## Single-Layer RNN:

$$\mathbf{h}_t = \tanh(\mathbf{A}_h \mathbf{h}_{t-1} + \mathbf{A}_x \mathbf{x}_t + \mathbf{b}_h)$$

$$\hat{\mathbf{y}}_t = \mathbf{A}_y \mathbf{h}_t + \mathbf{b}_y$$

- ▶ Hidden state  $\mathbf{h}_t$  = linear combination of input  $\mathbf{x}_t$  and previous hidden state  $\mathbf{h}_{t-1}$
- ▶ Output  $\hat{\mathbf{y}}_t$  = linear prediction based on current hidden state  $\mathbf{h}_t$
- ▶  $\tanh(\cdot)$  is commonly used as activation function (data is in the range  $[-1, 1]$ )
- ▶ Parameters  $\mathbf{A}_h, \mathbf{A}_x, \mathbf{A}_y, \mathbf{b}_h, \mathbf{b}_y$  are constant over time

# Basic Recurrent Neural Network



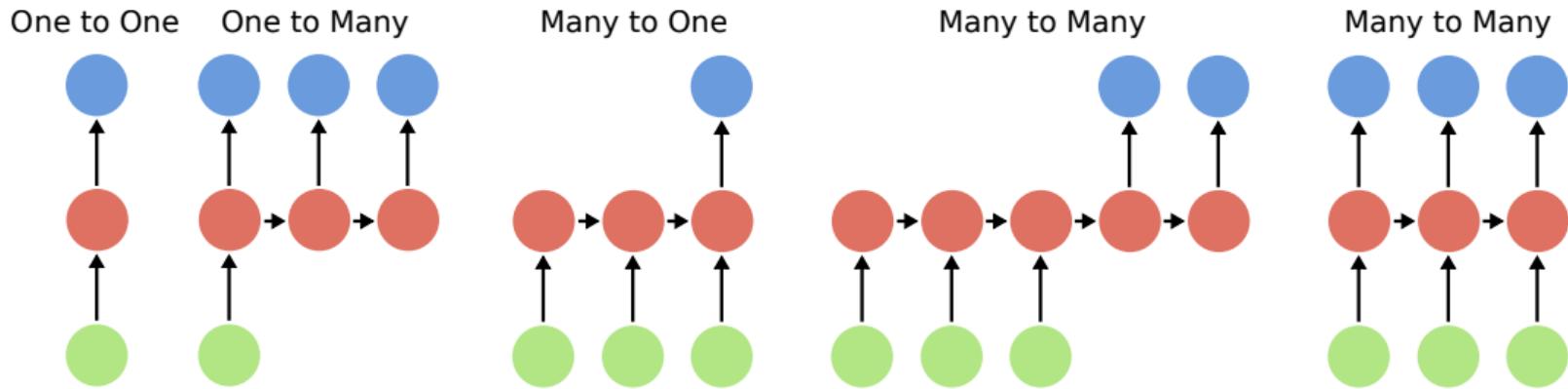
## Single-Layer RNN in Einstein Notation:

$$\begin{aligned} H_t[b, c_{out}] &= \tanh(A_h[c_{out}, C_{in}] H_{t-1}[b, C_{in}] + A_x[c_{out}, C_{in}] X_t[b, C_{in}] + b[c_{out}]) \\ &= \tanh\left(A[c_{out}, C_{in}] \begin{pmatrix} H_{t-1}[b, C_{in}] \\ X_t[b, C_{in}] \end{pmatrix} + b[c_{out}]\right) \end{aligned}$$

$$\hat{Y}_t[b, c_{out}] = A_y[c_{out}, C_{in}] H_t[b, C_{in}] + b[c_{out}]$$

Note:  $c_{in}$  and  $c_{out}$  are (per-product) index names (not dimensionality).  $\mathbf{A}_h$  is square.

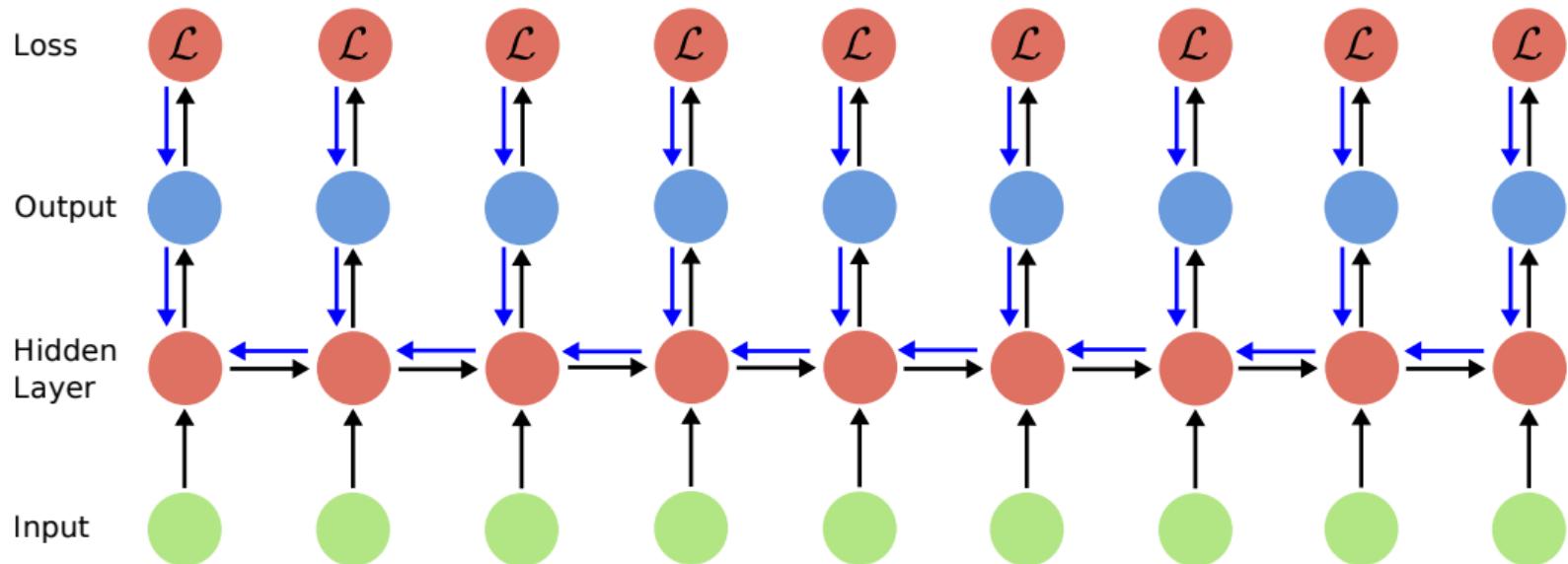
# Mapping Types



RNNs allow for processing **variable length** inputs and outputs:

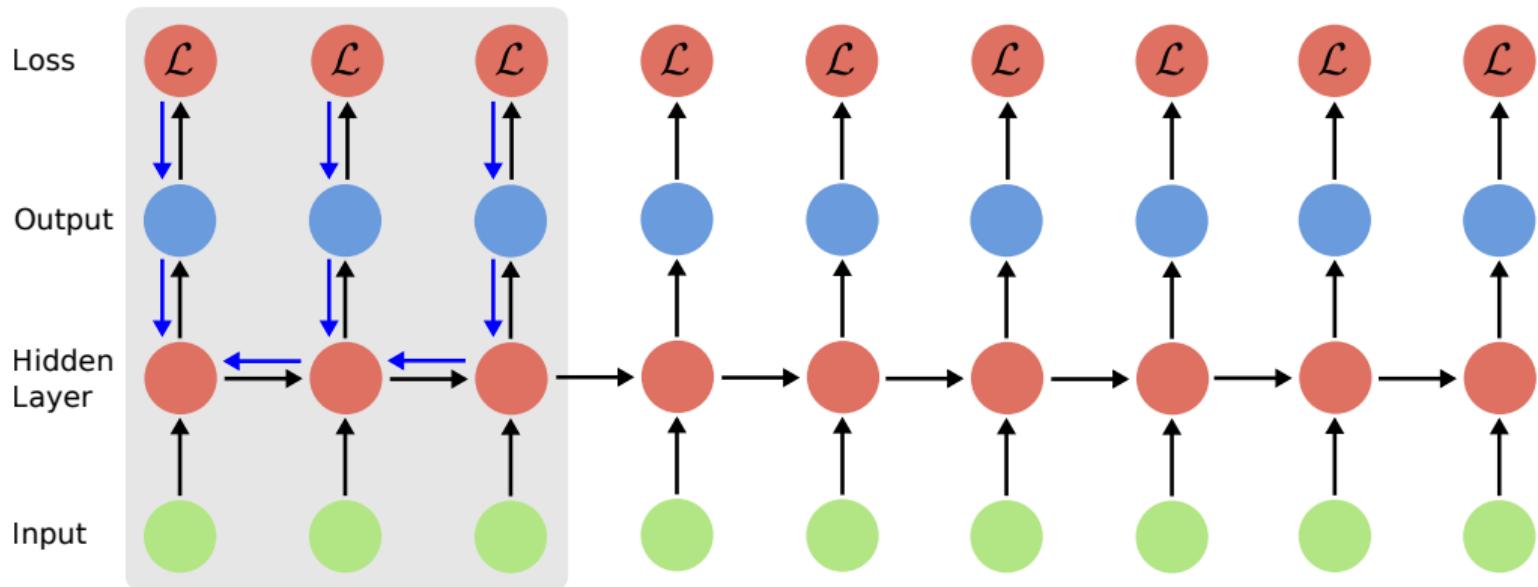
- ▶ **One to Many:** Image captioning (image to sentence)
- ▶ **Many to One:** Action recognition (video to action)
- ▶ **Many to Many:** Machine translation (sentence to sentence)
- ▶ **Many to Many:** Object tracking (video to object location per frame)
- ▶ To determine the length of the output sequence, a **stop symbol** can be predicted

# Backpropagation through Time



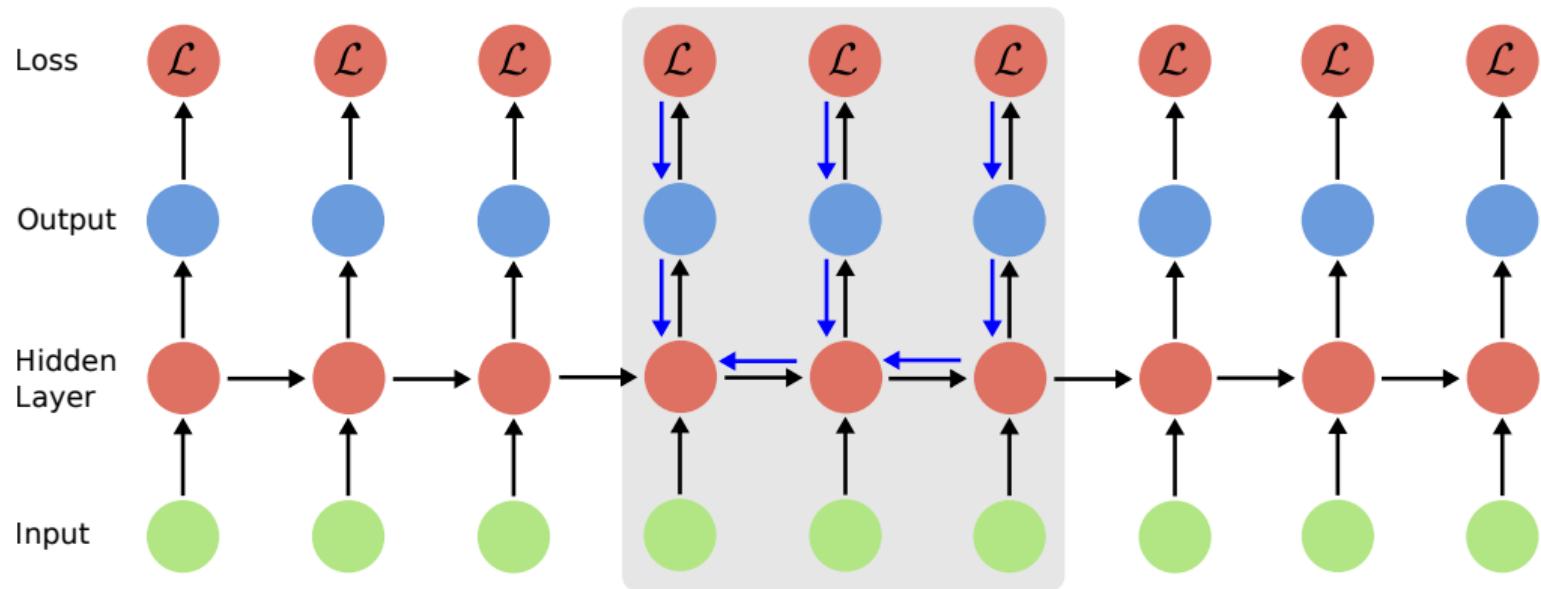
- To train RNNs, we **backpropagate gradients through time**
- As all hidden RNN cells share their parameters, gradients get accumulated
- However, very quickly intractable (memory) for larger sequences (e.g., Wikipedia)

# Truncated Backpropagation through Time



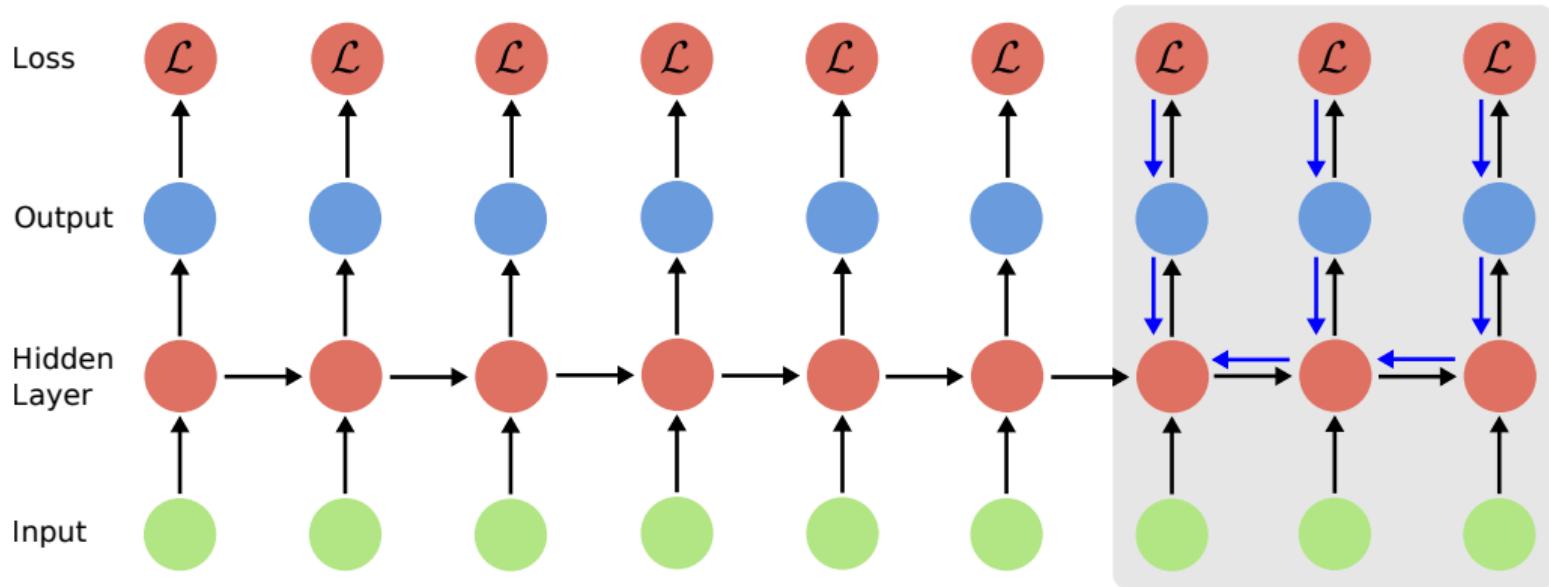
- ▶ Thus, one typically uses **truncated backpropagation through time** in practice
- ▶ Carry hidden states forward in time forever, but stop backpropagation earlier
- ▶ Total loss is sum of individual loss functions (= negative log likelihood)

# Truncated Backpropagation through Time



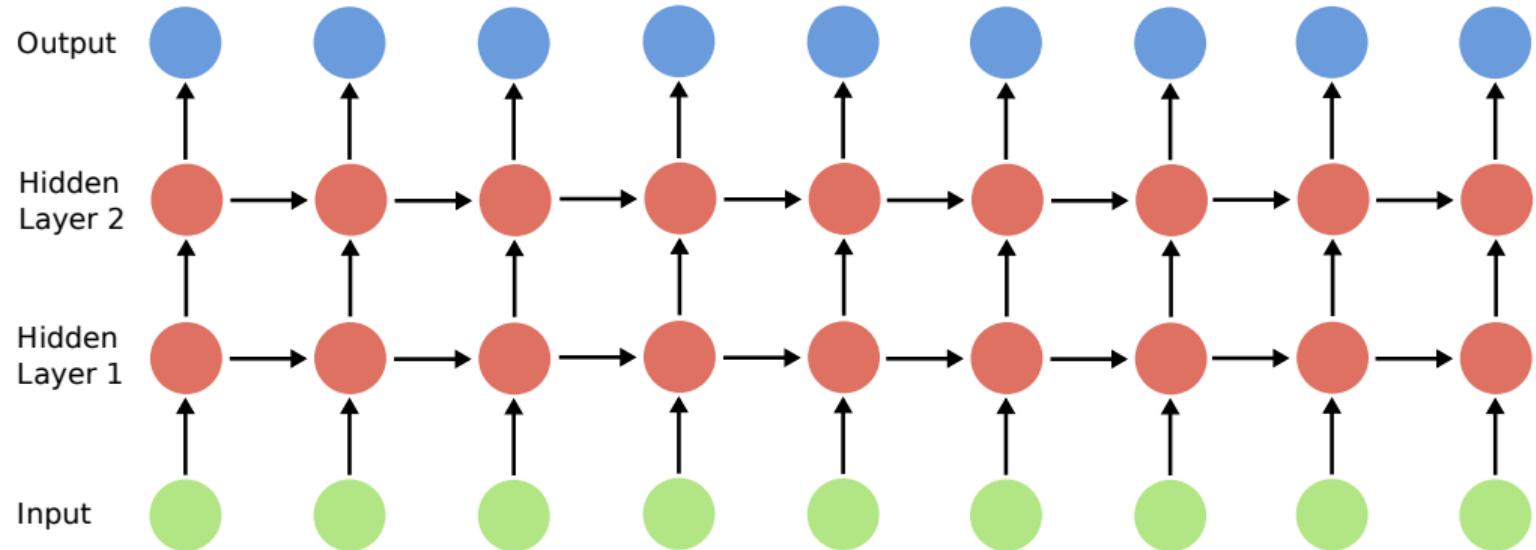
- ▶ Thus, one typically uses **truncated backpropagation through time** in practice
- ▶ Carry hidden states forward in time forever, but stop backpropagation earlier
- ▶ Total loss is sum of individual loss functions (= negative log likelihood)

# Truncated Backpropagation through Time



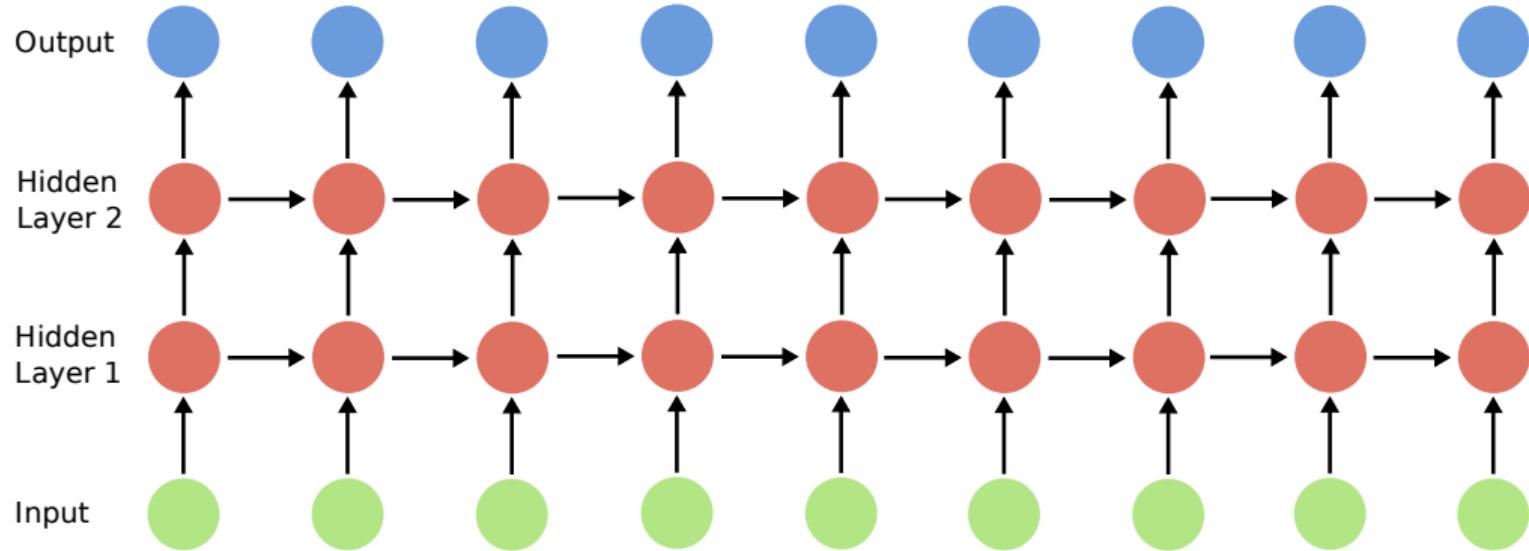
- ▶ Thus, one typically uses **truncated backpropagation through time** in practice
- ▶ Carry hidden states forward in time forever, but stop backpropagation earlier
- ▶ Total loss is sum of individual loss functions (= negative log likelihood)

# Multi-Layer RNNs



- Deeper **multi-layer RNNs** can be constructed by stacking RNN layers
- An alternative is to make each individual computation (=RNN cell) deeper
- Today, often combined with residual connections in vertical direction

# Multi-Layer RNNs



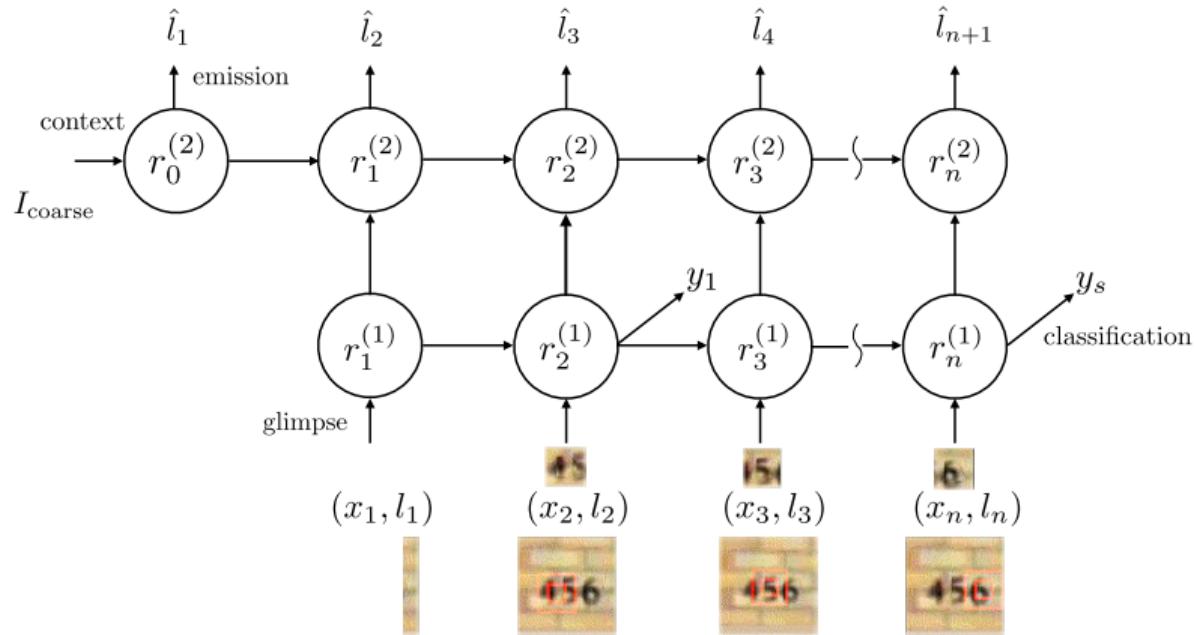
$$H_t^1[b, c_{out}] = \tanh (A_h^1[c_{out}, C_{in}] H_{t-1}^1[b, C_{in}] + A_x^1[c_{out}, C_{in}] X_t[b, C_{in}] + b[c_{out}])$$

$$H_t^2[b, c_{out}] = \tanh (A_h^2[c_{out}, C_{in}] H_{t-1}^2[b, C_{in}] + A_x^2[c_{out}, C_{in}] H_t^1[b, C_{in}] + b[c_{out}])$$

$$\hat{Y}_t[b, c_{out}] = A_y[c_{out}, C_{in}] H_t^2[b, C_{in}] + b[c_{out}]$$

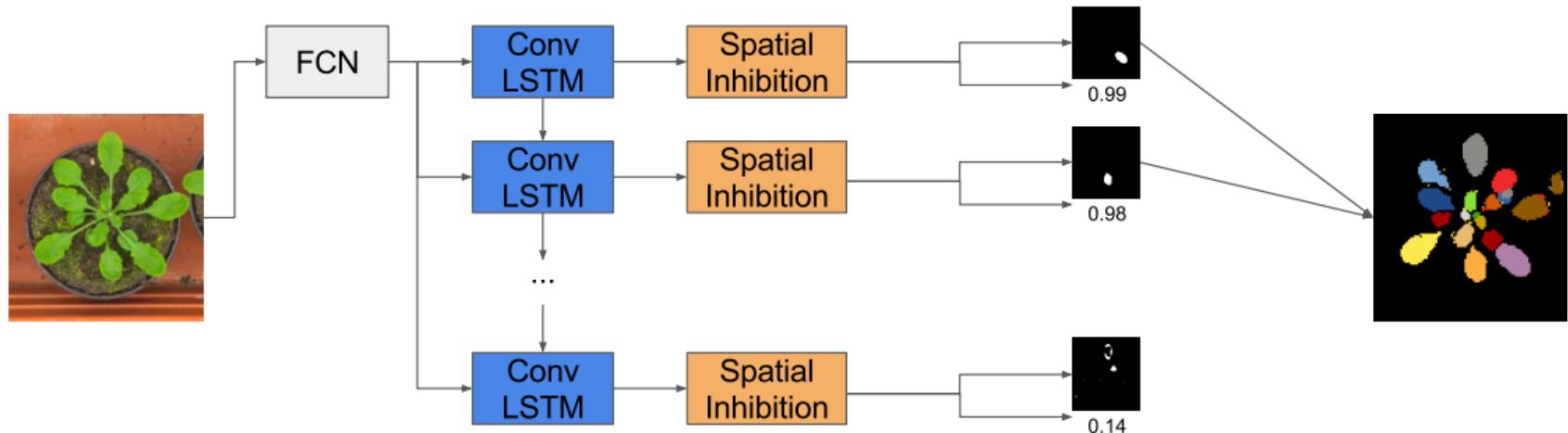


# Multiple Object Recognition



- At each time step, perceive a **glimpse** (= image region) and predict a **saccade**

# Recurrent Instance Segmentation



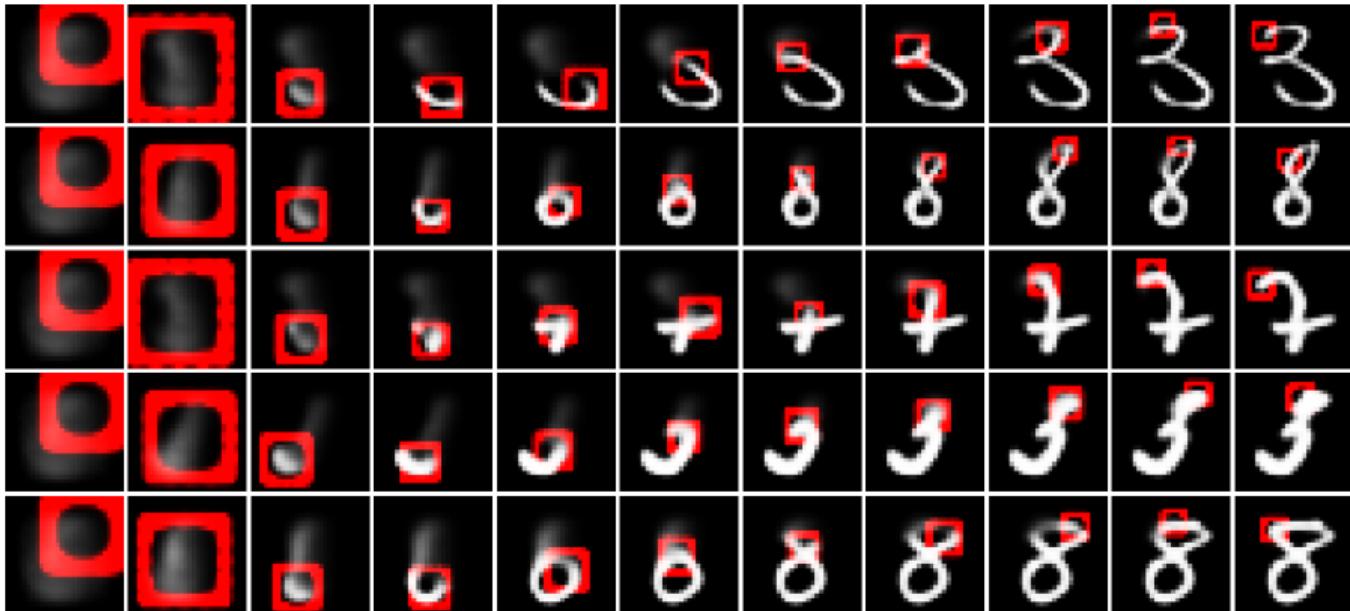
- At each time step, **segment** the next (not yet segmented) **part** of an object

# Object Tracking



- **Tracking of multiple objects** by updating each object's hidden state using an RNN

# Image Generation



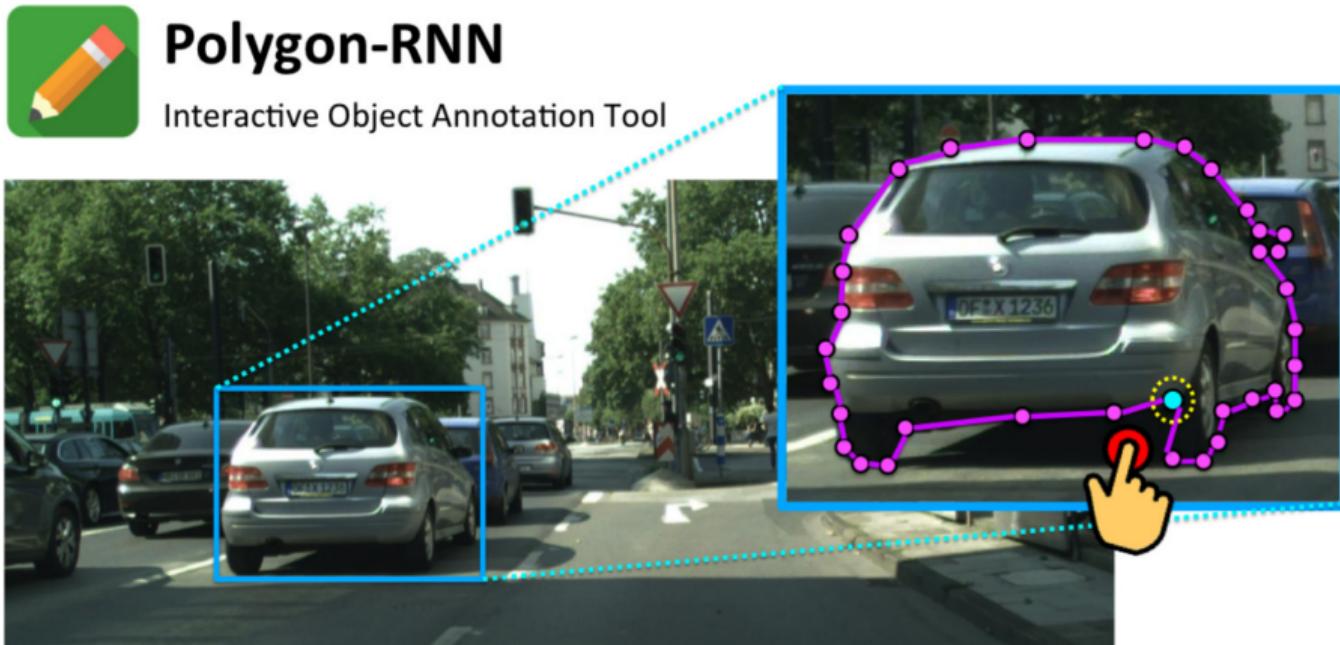
- ▶ Model for **sequential image generation** (red rectangle = attended region)
- ▶ Demonstrates that RNNs can also process/generate non-sequential data

# Image Generation



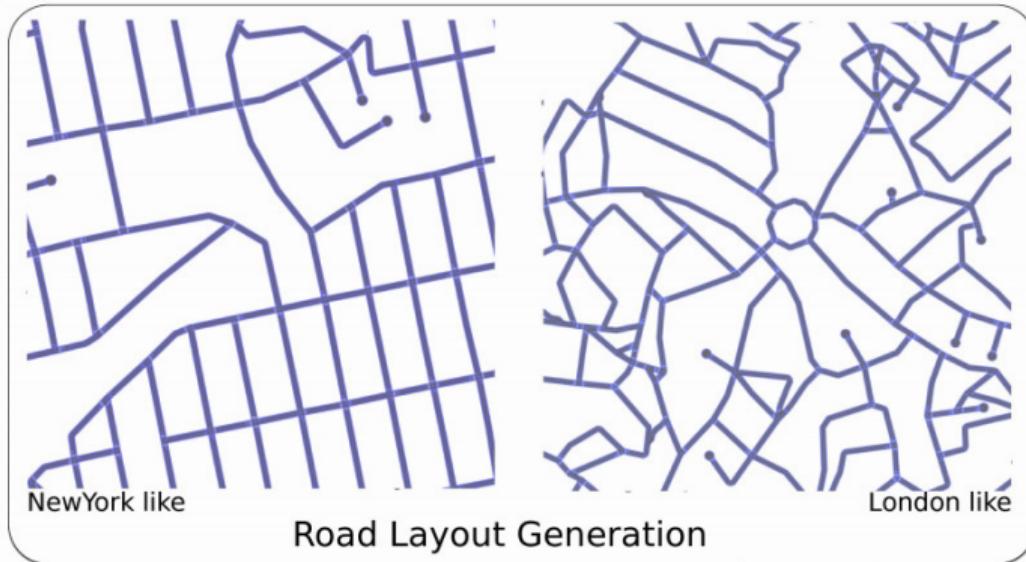
- ▶ Generated images based on partially occluded inputs
- ▶ Demonstrates that RNNs can also process/generate non-sequential data

# Image Annotation



- ▶ Iteratively **annotate** the outline of 2D object instances in an image with **polygons**

# Modeling Road Layouts



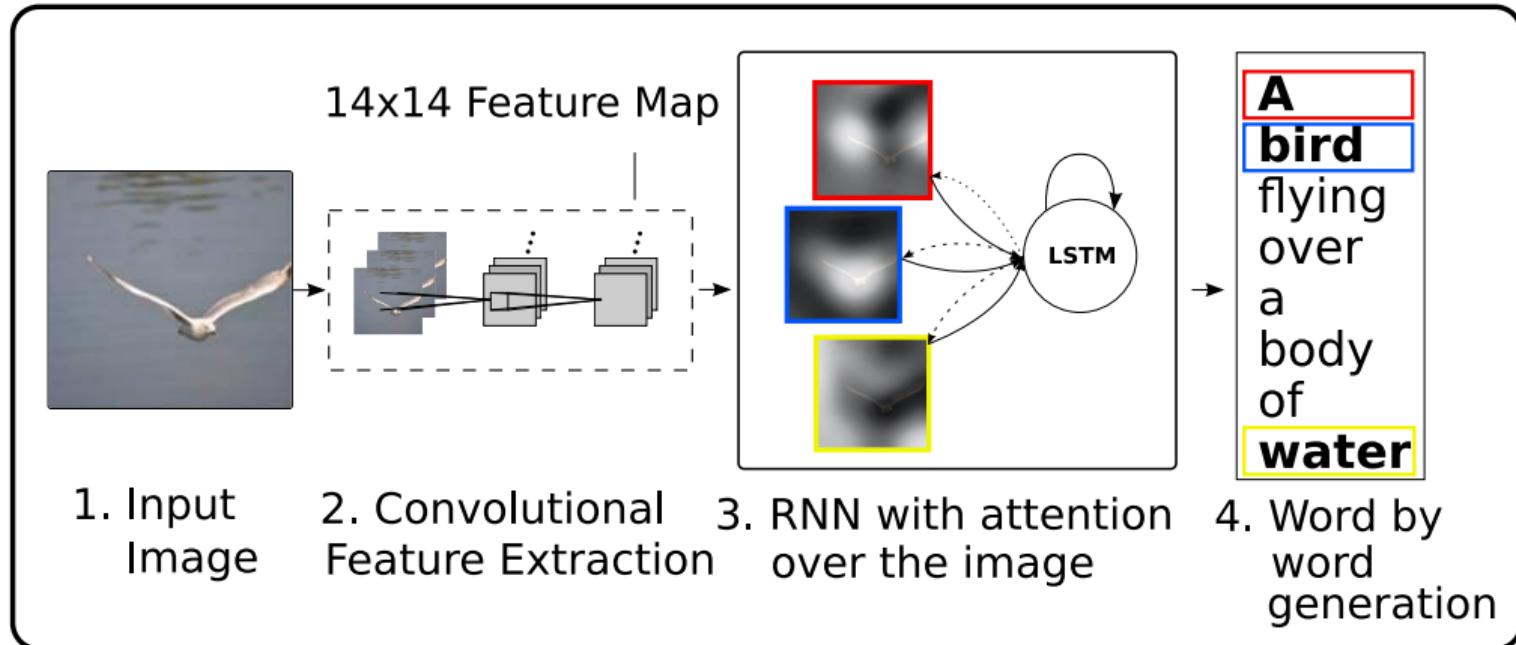
Road Layout Generation



Aerial Road Parsing

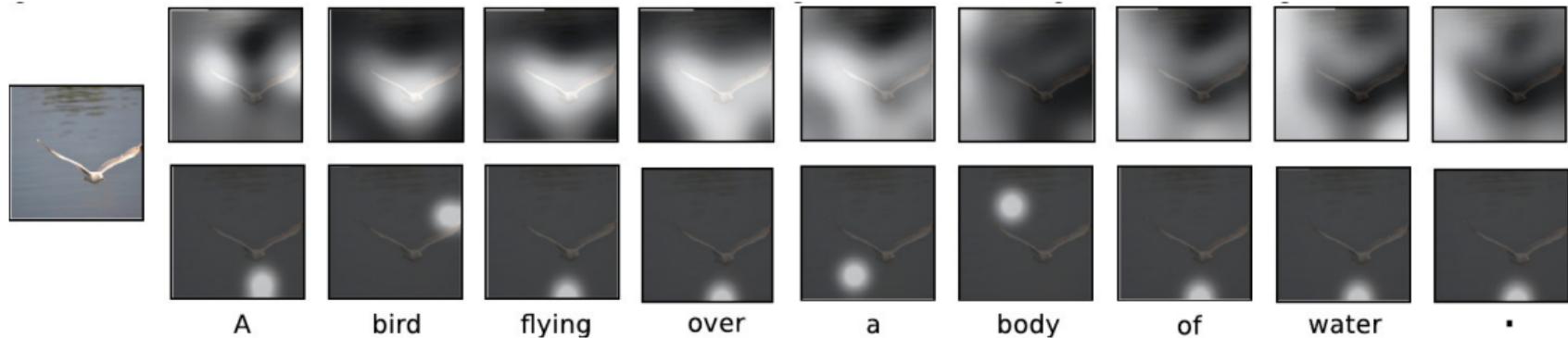
- ▶ Iteratively generate or infer **road layouts** as spatial graphs

# Image Captioning



- Generate **image description** by sequentially looking at an image

# Image Captioning



- ▶ **Attention over time**
- ▶ Top: soft attention. Bottom: hard attention

# Image Captioning



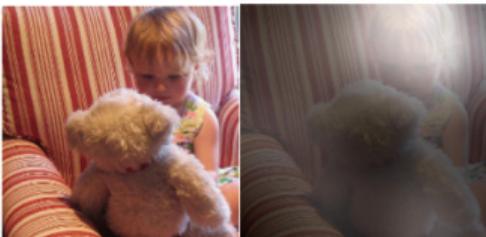
A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



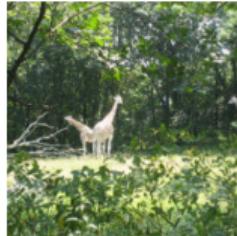
A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

## ► Successful caption generations

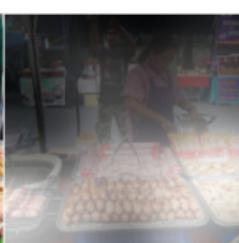
# Image Captioning



A large white bird standing in a forest.

A woman holding a clock in her hand.

A man wearing a hat and a hat on a skateboard.



A person is standing on a beach with a surfboard.

A woman is sitting at a table with a large pizza.

A man is talking on his cell phone while another man watches.

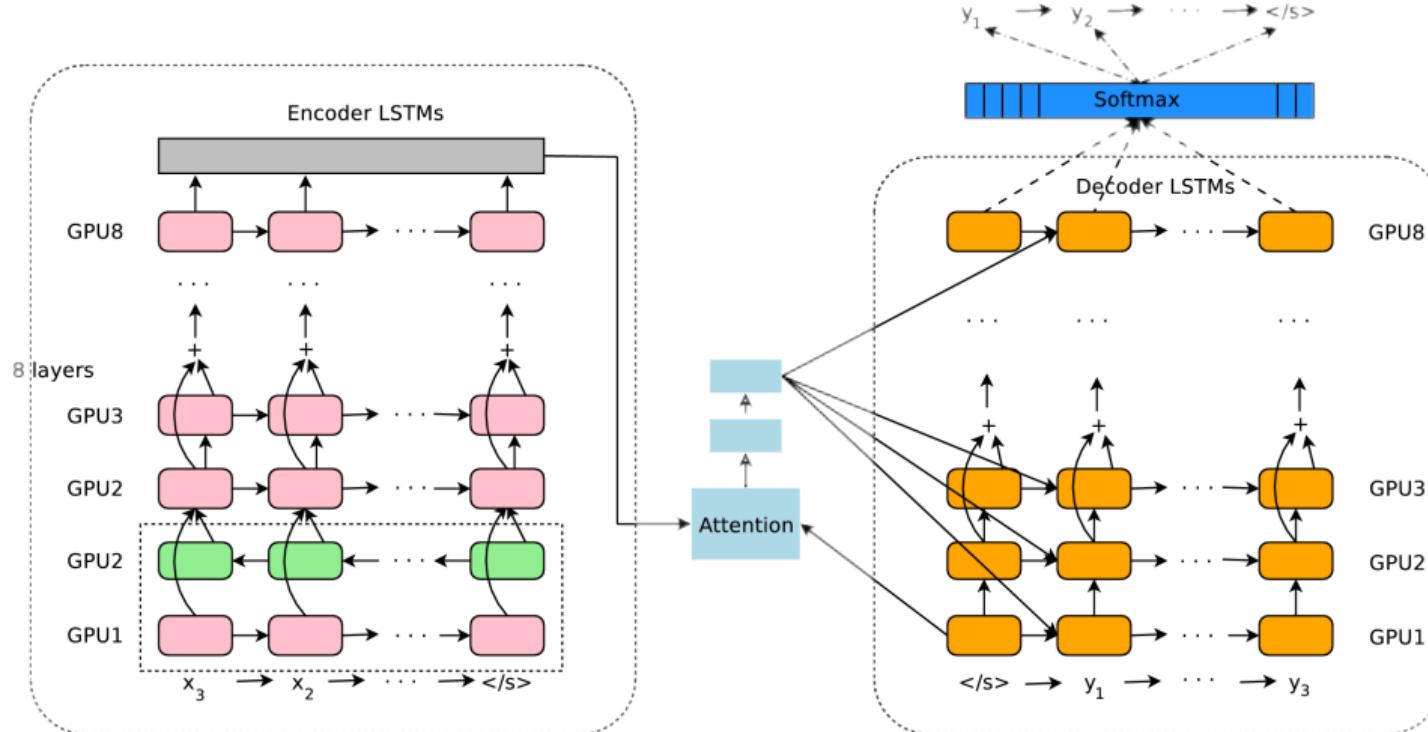
- ▶ Wrongly generated captions. Attention can reveal insights into what went wrong.

# Visual Question Answering

Image						
Multiple Choices	<p>Q: Who is behind the batter?</p> <p>A: Catcher. A: Umpire. A: Fans. A: Ball girl.</p>	<p>Q: What adorns the tops of the post?</p> <p>A: Gulls. A: An eagle. A: A crown. A: A pretty sign.</p>	<p>Q: How many cameras are in the photo?</p> <p>A: One. A: Two. A: Three. A: Four.</p>	<p>Q: Why is there rope?</p> <p>A: To tie up the boats. A: To tie up horses. A: To hang people. A: To hit tether balls.</p>	<p>Q: What kind of stuffed animal is shown?</p> <p>A: Teddy Bear. A: Monkey. A: Tiger. A: Bunny rabbit.</p>	<p>Q: What animal is being petted?</p> <p>A: A sheep. A: Goat. A: Alpaca. A: Pig.</p>
w/ Image w/o Image	<p>H: Catcher. ✓ M: Umpire. ✗</p>	<p>H: Gulls. ✓ M: Gulls. ✓</p>	<p>H: Three. ✗ M: One. ✓</p>	<p>H: To hit tether balls. ✗ M: To hang people. ✗</p>	<p>H: Monkey. ✗ M: Teddy Bear. ✓</p>	<p>H: A sheep. ✓ M: A sheep. ✓</p>
w/ Image	<p>H: Catcher. ✓ M: Catcher. ✓</p>	<p>H: Gulls. ✓ M: A crown. ✗</p>	<p>H: One. ✓ M: One. ✓</p>	<p>H: To tie up the boats. ✓ M: To hang people. ✗</p>	<p>H: Teddy Bear. ✓ M: Teddy Bear. ✓</p>	<p>H: Goat. ✗ M: A sheep. ✓</p>

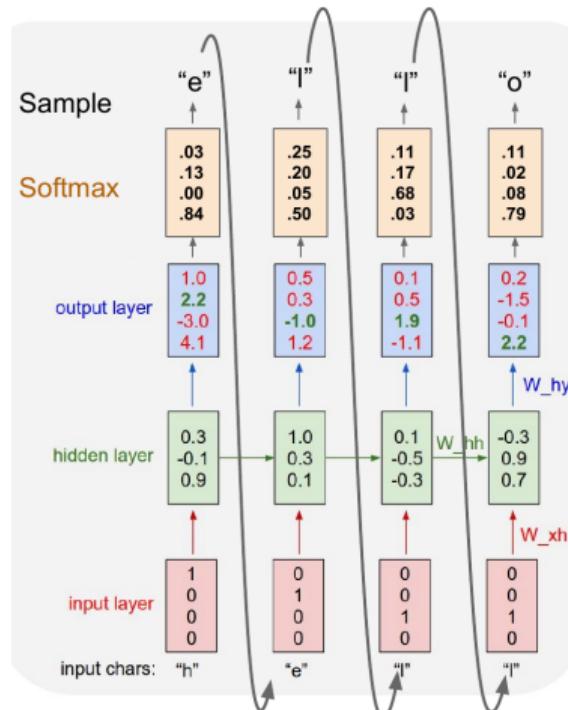
► Input: Image and question sentence. Output: answer sentence.

# Neural Machine Translation



# Character-level Language Models

- ▶ Generating natural language text **character by character** with an RNN
- ▶ In the example on the right: Alphabet with 4 characters ('h', 'e', 'l', 'o')
- ▶ Each character is represented by a **1 hot vector**, e.g.:  $(1, 0, 0, 0)^\top$
- ▶ Model predicts **distribution over next character** via softmax function
- ▶ Character drawn from distribution is **fed as input** to RNN at next time step



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

# Character-level Language Models

- ▶ 3-layer RNN
- ▶ 512 hidden nodes
- ▶ Trained on all works  
of William Shakespeare
- ▶ 4.4 Mio characters

PANDARUS:

Alas, I think he shall be come approached and the day  
When little strain would be attain'd into being never fed,  
And who is but a chain and subjects of his death,  
I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,  
Breaking and strongly should be buried, when I perish  
The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

Second Lord:

They would be ruled after this chamber, and  
my fair nues begun out of the fact, to be conveyed,  
Whose noble souls I'll have the heart of the wars.

# Character-level Language Models

tyntd-iafhatawiaoahrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e  
plia tkldrgd t o idoe ns,smtt h ne etie h,hregtrs nigtike,aoaenns lng

↓ train more

"Tmont thithey" fomesscerliund  
Keushey. Thom here  
sheulke, anmerenith ol sivh I lalterthend Bleipile shuw fil on aseterlome  
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

↓ train more

Aftair fall unsuch that the hall for Prince Velzonski's that me of  
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort  
how, and Gogition is so overelical and ofter.

↓ train more

"Why do what that day," replied Natasha, and wishing to himself the fact the  
princess, Princess Mary was easier, fed in had oftened him.  
Pierre aking his soul came to the packs and drove up his father-in-law women.

# Character-level Language Models

*Proof.* Omitted.  $\square$

**Lemma 0.1.** Let  $\mathcal{C}$  be a set of the construction.

Let  $\mathcal{C}$  be a gerber covering. Let  $\mathcal{F}$  be a quasi-coherent sheaves of  $\mathcal{O}$ -modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

*Proof.* This is an algebraic space with the composition of sheaves  $\mathcal{F}$  on  $X_{etale}$  we have

$$\mathcal{O}_X(\mathcal{F}) = \{morph_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where  $\mathcal{G}$  defines an isomorphism  $\mathcal{F} \rightarrow \mathcal{F}$  of  $\mathcal{O}$ -modules.  $\square$

**Lemma 0.2.** This is an integer  $\mathcal{Z}$  is injective.

*Proof.* See Spaces, Lemma ???.  $\square$

**Lemma 0.3.** Let  $S$  be a scheme. Let  $X$  be a scheme and  $X$  is an affine open covering. Let  $\mathcal{U} \subset X$  be a canonical and locally of finite type. Let  $X$  be a scheme. Let  $X$  be a scheme which is equal to the formal complex.

The following to the construction of the lemma follows.

Let  $X$  be a scheme. Let  $X$  be a scheme covering. Let

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

be a morphism of algebraic spaces over  $S$  and  $Y$ .

*Proof.* Let  $X$  be a nonzero scheme of  $X$ . Let  $X$  be an algebraic space. Let  $\mathcal{F}$  be a quasi-coherent sheaf of  $\mathcal{O}_X$ -modules. The following are equivalent

- (1)  $\mathcal{F}$  is an algebraic space over  $S$ .
- (2) If  $X$  is an affine open covering.

Consider a common structure on  $X$  and  $X$  the functor  $\mathcal{O}_X(U)$  which is locally of finite type.  $\square$

This since  $\mathcal{F} \in \mathcal{F}$  and  $x \in \mathcal{G}$  the diagram

$$\begin{array}{ccccc}
 S & \xrightarrow{\quad} & & & \\
 \downarrow & & & & \\
 \xi & \longrightarrow & \mathcal{O}_{X'} & & \\
 gor_s & & \uparrow & \searrow & \\
 & & = a' \longrightarrow & & X \\
 & & \downarrow & & \downarrow \\
 & & = a' \longrightarrow a & & \text{Spec}(K_\psi) \qquad \text{Mor}_{Sets} \qquad d(\mathcal{O}_{X_{f/k}}, \mathcal{G})
 \end{array}$$

is a limit. Then  $\mathcal{G}$  is a finite type and assume  $S$  is a flat and  $\mathcal{F}$  and  $\mathcal{G}$  is a finite type  $f_*$ . This is of finite type diagrams, and

- the composition of  $\mathcal{G}$  is a regular sequence,
- $\mathcal{O}_{X'}$  is a sheaf of rings.

$\square$

*Proof.* We have see that  $X = \text{Spec}(R)$  and  $\mathcal{F}$  is a finite type representable by algebraic space. The property  $\mathcal{F}$  is a finite morphism of algebraic stacks. Then the cohomology of  $X$  is an open neighbourhood of  $U$ .  $\square$

*Proof.* This is clear that  $\mathcal{G}$  is a finite presentation, see Lemmas ???. A reduced above we conclude that  $U$  is an open covering of  $C$ . The functor  $\mathcal{F}$  is a  $\text{field}$

$$\mathcal{O}_{X,x} \rightarrow \mathcal{F}_{\bar{x}} \dashv \mathcal{I}(\mathcal{O}_{X_{etale}}) \rightarrow \mathcal{O}_{X_{\bar{x}}}^{-1} \mathcal{O}_{X_{\bar{x}}}(\mathcal{O}_{X_{\bar{x}}})$$

is an isomorphism of covering of  $\mathcal{O}_{X_{\bar{x}}}$ . If  $\mathcal{F}$  is the unique element of  $\mathcal{F}$  such that  $X$  is an isomorphism.

The property  $\mathcal{F}$  is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme  $\mathcal{O}_X$ -algebra with  $\mathcal{F}$  are opens of finite type over  $S$ .

If  $\mathcal{F}$  is a scheme theoretic image points.  $\square$

If  $\mathcal{F}$  is a finite direct sum  $\mathcal{O}_{X_{\bar{x}}}$  is a closed immersion, see Lemma ???. This is a sequence of  $\mathcal{F}$  is a similar morphism.

# Character-level Language Models

```
/*
 * Copyright (c) 2006-2010, Intel Mobile Communications. All rights reserved.
 *
 * This program is free software; you can redistribute it and/or modify it
 * under the terms of the GNU General Public License version 2 as published by
 * the Free Software Foundation.
 *
 * This program is distributed in the hope that it will be useful,
 * but WITHOUT ANY WARRANTY; without even the implied warranty of
 * MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
 *
 * GNU General Public License for more details.
 *
 * You should have received a copy of the GNU General Public License
 * along with this program; if not, write to the Free Software Foundation,
 * Inc., 675 Mass Ave, Cambridge, MA 02139, USA.
 */
#include <linux/kexec.h>
#include <linux/errno.h>
#include <linux/io.h>
```

```
#define REG_PG    vesa_slot_addr_pack
#define PFM_NOCOMP AFSR(0, load)
#define STACK_DDR(type)      (func)

#define SWAP_ALLOCATE(nr)      (e)
#define emulate_sigs() arch_get_unaligned_child()
#define access_rw(TST)  asm volatile("movd %%esp, %0, %3" : : "r" (0)); \
    if (_type & DO_READ)

static void stat_PC_SEC __read_mostly offsetof(struct seq_argsqueue, \
                                              pC>[1]);

static void
os_prefix(unsigned long sys)
{
#endif CONFIG_PREEMPT
    PUT_PARAM_RAID(2, sel) = get_state_state();
    set_pid_sum((unsigned long)state, current_state_str(),
                (unsigned long)-1->lr_full; low;
}
```

- ▶ 3-layer RNN trained for several days on **Linux source code** (474 MB)
- ▶ Sampled code snippets do not compile but look reasonable overall
- ▶ Learned that code starts with license, uses correct syntax, adds comments

# Character-level Language Models

## **Samples from RNN trained on 8000 baby names:**

Rudi Levette Berice Lussa Hany Mareanne Chrestina Carissy Marylen Hammine Janye  
Marlise Jacacie Hendred Romand Charienna Nenotto Ette Dorane Wallen Marly  
Darine Salina Elvyn Ersia Maralena Minoria Ellia Charmin Antley Nerille Chelon Walmor  
Evena Jeryly Stachon Charisa Allisa Anatha Cathanie Geetra Alexie Jerin Cassen  
Herbett Cossie Velen Daurenge Robester Shermond Terisa Licia Roselen Ferine Jayn  
Lusine Charyanne Sales Sanny Resa Wallon Martine Merus Jelen Candica Wallin Tel  
Rachene Tarine Ozila Ketia Shanne Arnande Karella Roselina Alessia Chasty Deland  
Berther Geamar Jackein Mellisand Sagdy Nenc Lessie Rasemy Guen Gavi Milea  
Annedo Margoris Janin Rodelin Zeanna Elyne Janah Ferzina Susta Pey Castina

# Character-level Language Models

Cell sensitive to position in line:

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.

Cell that turns on inside quotes:

"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."

- The behavior of some hidden neurons ( $\sim 5\%$ ) is logical and human interpretable

# Character-level Language Models

Cell that robustly activates inside if statements:

```
static int __dequeue_signal(struct sigpending *pending, sigset_t *mask,
    siginfo_t *info)
{
    int sig = next_signal(pending, mask);
    if (sig) {
        if (current->notifier) {
            if (sigismember(current->notifier_mask, sig)) {
                if (!(*current->notifier)(current->notifier_data)) {
                    clear_thread_flag(TIF_SIGPENDING);
                    return 0;
                }
            }
        }
        collect_signal(sig, pending, info);
    }
    return sig;
}
```

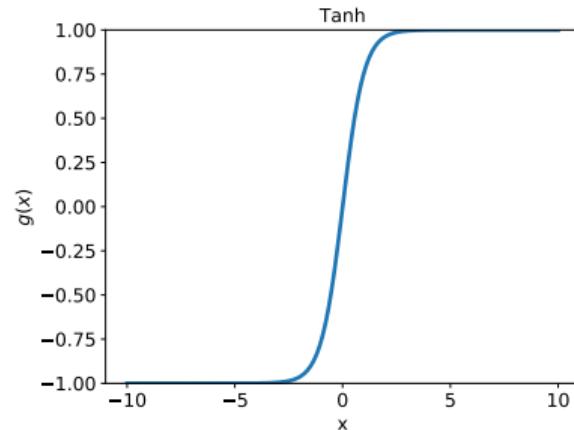
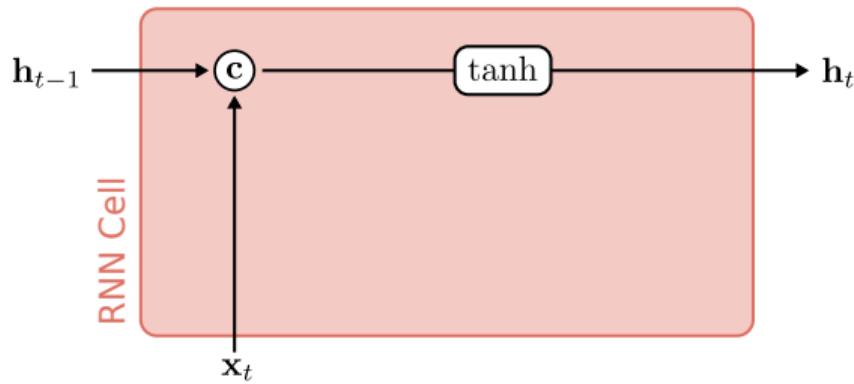
A large portion of cells are not easily interpretable. Here is a typical example:

```
/* Unpack a filter field's string representation from user-space
 * buffer. */
char *audit_unpack_string(void **bufp, size_t *remain, size_t len)
{
    char *str;
    if (!*bufp || (len == 0) || (len > *remain))
        return ERR_PTR(-EINVAL);
    /* Of the currently implemented string fields, PATH_MAX
     * defines the longest valid length.
    */
```

## 8.3

# Gated Recurrent Networks

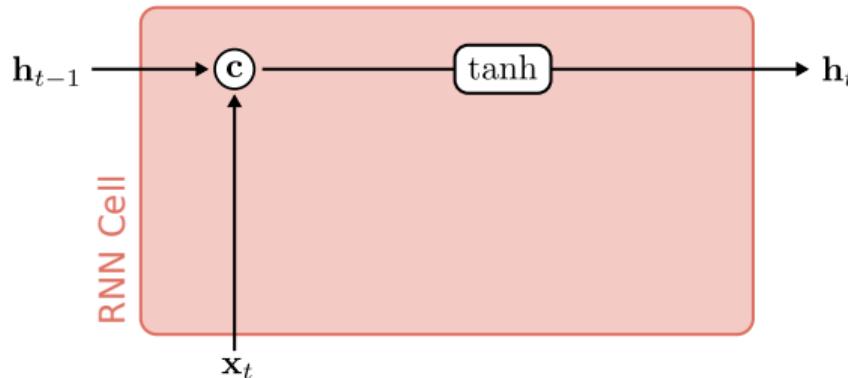
# Recurrent Neural Network



$$H_t[b, c_{out}] = \tanh(A_h[c_{out}, C_{in}] H_{t-1}[b, C_{in}] + A_x[c_{out}, C_{in}] X_t[b, C_{in}] + b[c_{out}])$$

- ▶ The state update  $H_t$  is modeled using a zero-centered  $\tanh(\cdot)$
- ▶  $\tanh(\cdot)$  assumes that the processed data is in the range  $[-1, 1]$
- ▶ Remark: we omit the affine transformations and the output layer for clarity

# Vanishing / Exploding Gradients



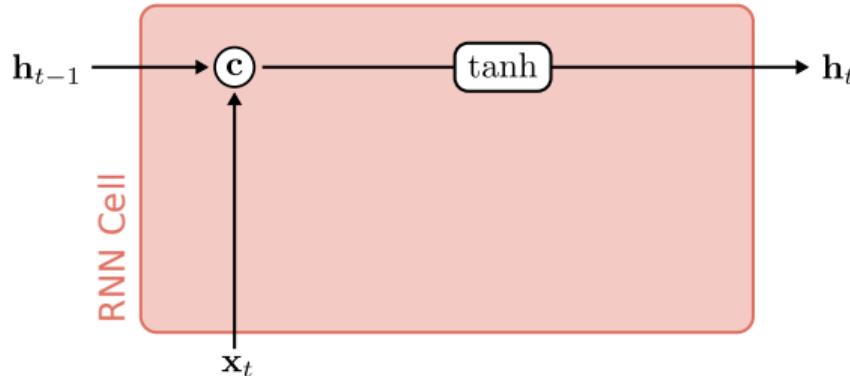
$$h_t = \tanh(a_h h_{t-1} + a_x x_t + b)$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \tanh' a_h \text{ with } \tanh' = \frac{\partial \tanh(x)}{\partial x}$$

## What is the problem with vanilla RNNs?

- ▶ Let us consider an RNN with one dimensional hidden state  $h_t \in \mathbb{R}$
- ▶ We have: 
$$\frac{\partial h_t}{\partial h_{t-k}} = \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} \dots \frac{\partial h_{t-k+1}}{\partial h_{t-k}} = \left( \prod_{i=t-k+1}^t \tanh'_i \right) a_h^k$$
- ▶ Thus, the gradient vanishes if  $\tanh(\cdot)$  saturates as in feedforward networks
- ▶ RNNs require careful initialization to avoid saturating activation functions

# Vanishing / Exploding Gradients



$$h_t = \tanh(a_h h_{t-1} + a_x x_t + b)$$

$$\approx a_h h_{t-1} + a_x x_t + b$$

$$\frac{\partial h_t}{\partial h_{t-1}} \approx a_h$$

**However, gradients might still misbehave:**

- ▶ Let us now assume that the RNN has been initialized well such that the activation functions are not saturated  $\Rightarrow a_h h_{t-1} + a_x x_t + b \in [-1, 1] \Rightarrow \tanh(x) \approx x$
- ▶ We now have:

$$\frac{\partial h_t}{\partial h_{t-k}} = \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} \dots \frac{\partial h_{t-k+1}}{\partial h_{t-k}} = a_h^k$$

# Vanishing / Exploding Gradients

$$\frac{\partial h_t}{\partial h_{t-k}} = \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} \dots \frac{\partial h_{t-k+1}}{\partial h_{t-k}} = a_h^k$$

For  $a_h > 1$  gradients will **explode** (become very large, cause divergence):

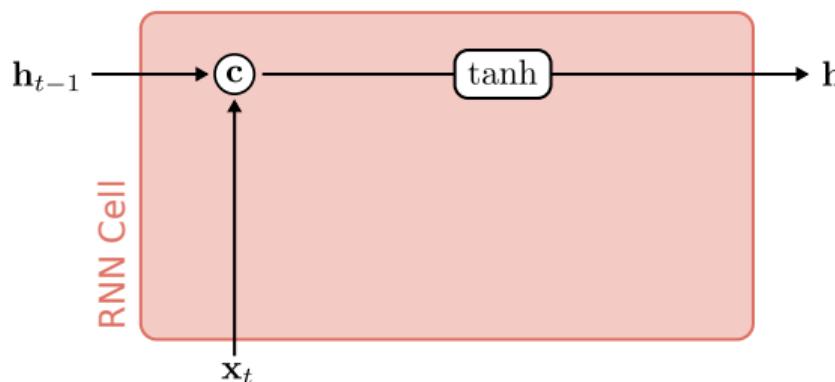
- ▶ Example: For  $a_h = 1.1$  and  $k = 100$  we have  $\partial h_t / \partial h_{t-k} = a_h^k = 13781$
- ▶ This problem is often addressed in practice using **gradient clipping**
- ▶ Forward values do not explode due to bounded  $\tanh(\cdot)$  activation function

For  $a_h < 1$  gradients will **vanish** (no learning in earlier time steps):

- ▶ Example: For  $a_h = 0.9$  and  $k = 100$  we have  $\partial h_t / \partial h_{t-k} = a_h^k = 0.0000266$
- ▶ Avoiding this problem requires an **architectural change**
- ▶ But residual connections do not work here as the parameters are shared across time and the input and desired output at each time step are different

# Vanishing / Exploding Gradients

**Similar situation for vector-valued hidden states:**



$$\begin{aligned} \mathbf{h}_t &= \tanh(\mathbf{A}_h \mathbf{h}_{t-1} + \mathbf{A}_x \mathbf{x}_t + \mathbf{b}) \\ &\approx \mathbf{A}_h \mathbf{h}_{t-1} + \mathbf{A}_x \mathbf{x}_t + \mathbf{b} \\ \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} &\approx \mathbf{A}_h \quad \Rightarrow \quad \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-k}} \approx \mathbf{A}_h^k \end{aligned}$$

- ▶ Let  $\mathbf{A}_h = \mathbf{Q}\Lambda\mathbf{Q}^{-1}$  be the eigendecomposition of the square matrix  $\mathbf{A}_h$
- ▶ We have  $\mathbf{A}_h^k = (\mathbf{Q}\Lambda\mathbf{Q}^{-1})^k = \mathbf{Q}\Lambda^k\mathbf{Q}^{-1}$  with diagonal eigenvalue matrix  $\Lambda$
- ▶ Components with eigenvalue  $< 1 \Rightarrow$  vanishing gradient
- ▶ Components with eigenvalue  $> 1 \Rightarrow$  exploding gradient
- ▶ Note that  $\mathbf{A}_h$  is shared across time (unlike in layers of a feedforward network)

# Gradient Clipping

The effect of exploding gradients can be dampened by a simple heuristic which clips the gradients to  $\|A \cdot \text{grad}\|_2 \leq \tau$  before applying the gradient update during SGD:

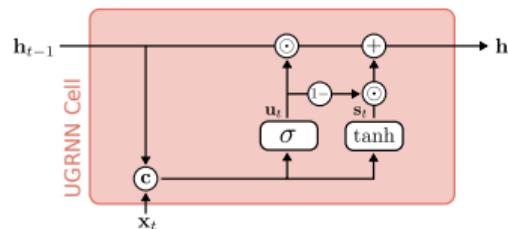
$$A \cdot \text{grad} = \begin{cases} A \cdot \text{grad} & \text{if } \|A \cdot \text{grad}\|_2 \leq \tau \\ \tau \frac{A \cdot \text{grad}}{\|A \cdot \text{grad}\|_2} & \text{otherwise} \end{cases}$$

The maximal gradient magnitude  $\tau$  is a hyperparameter, often  $\tau \in [1, 10]$ .

What can we do to avoid vanishing gradients?

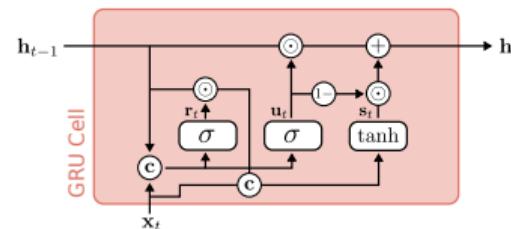
# Gated Recurrent Networks

**UGRNN**



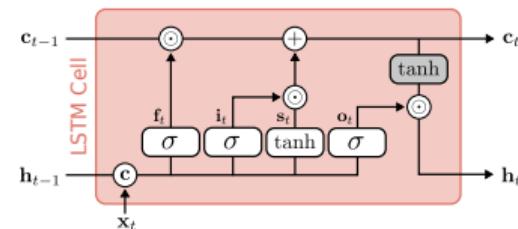
Collins, 2017

**GRU**



Cho, 2014

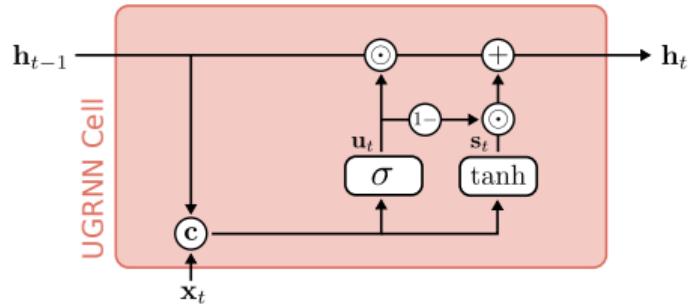
**LSTM**



Hochreiter, 1997

- ▶ **UGRNN:** Update Gate Recurrent Neural Network
- ▶ **GRU:** Gated Recurrent Unit
- ▶ **LSTM:** Long Short-Term Memory
- ▶ LSTM was the first and most transformative (revolutionized NLP in 2015, e.g. at Google), but also most complex model. UGRNN and GRU work similarly well.
- ▶ Common to all architectures: **gates** for filtering information

# Update Gate RNN



$$U_t[b, c_{out}] = \sigma(A_{uh}[c_{out}, C_{in}] H_{t-1}[b, C_{in}] + A_{ux}[c_{out}, C_{in}] X_t[b, C_{in}] + b_u[c_{out}])$$

$$S_t[b, c_{out}] = \tanh(A_{sh}[c_{out}, C_{in}] H_{t-1}[b, C_{in}] + A_{sx}[c_{out}, C_{in}] X_t[b, C_{in}] + b_s[c_{out}])$$

$$H_t[b, c_{out}] = U_t[b, c_{out}] H_{t-1}[b, c_{out}] + (1 - U_t[b, c_{out}]) S_t[b, c_{out}]$$

- ▶  $U_t$  is called **update gate** as it determines if the hidden state  $H$  is updated or not
- ▶  $S_t$  is the next **target state** that is added to  $H_{t-1}$  with element-wise weights  $U_t$
- ▶ Remark: Gates use sigmoid ( $\in [0, 1]$ ), state computation uses tanh ( $\in [-1, 1]$ )

# Hadamard Product

This expression

$$H_t[b, c_{out}] = U_t[b, c_{out}] H_{t-1}[b, c_{out}] + (1 - U_t[b, c_{out}]) S_t[b, c_{out}]$$

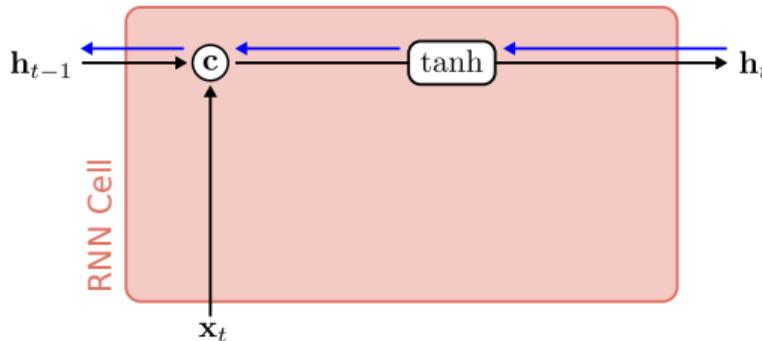
can be written as

$$H_t[b, C_{out}] = U_t[b, C_{out}] \odot H_{t-1}[b, C_{out}] + (1 - U_t[b, C_{out}]) \odot S_t[b, C_{out}]$$

where  $\odot$  denotes the **Hadamard product** (elementwise product).

We use this symbol to denote elementwise products in this lecture.

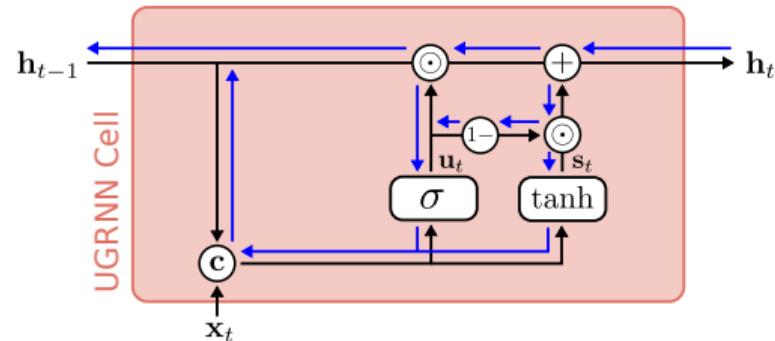
# Gradient Flow



$$\mathbf{h}_t = \tanh(\mathbf{A}_h \mathbf{h}_{t-1} + \mathbf{A}_x \mathbf{x}_t + \mathbf{b})$$

$$\mathbf{h}'_t = \tanh' \mathbf{A}_h \approx \mathbf{A}_h$$

$$\text{with } \mathbf{h}'_t = \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}}$$



$$\mathbf{u}_t = \sigma(\mathbf{A}_{uh} \mathbf{h}_{t-1} + \mathbf{A}_{ux} \mathbf{x}_t + \mathbf{b}_u)$$

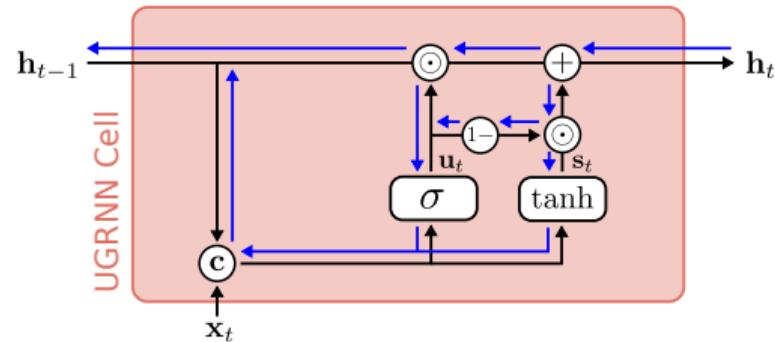
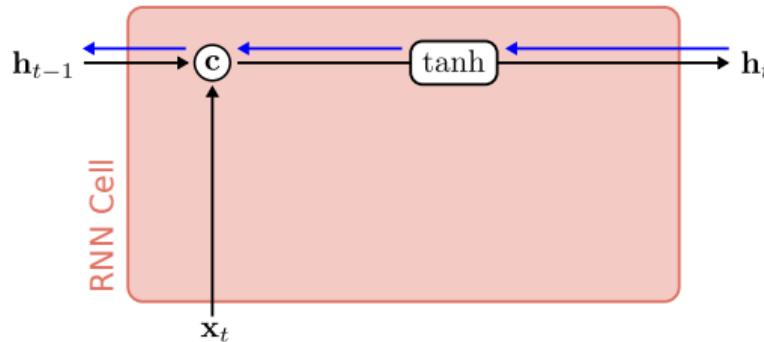
$$\mathbf{s}_t = \tanh(\mathbf{A}_{sh} \mathbf{h}_{t-1} + \mathbf{A}_{sx} \mathbf{x}_t + \mathbf{b}_s)$$

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

$$\mathbf{h}'_t = \mathbf{u}'_t \odot \mathbf{h}_{t-1} + \mathbf{u}_t + (1 - \mathbf{u}'_t) \odot \mathbf{s}_t + ..$$

- UGRNN can maintain gradient flow despite small  $\mathbf{A}_h$  by setting its gate to  $u \approx 1$

# Gradient Flow



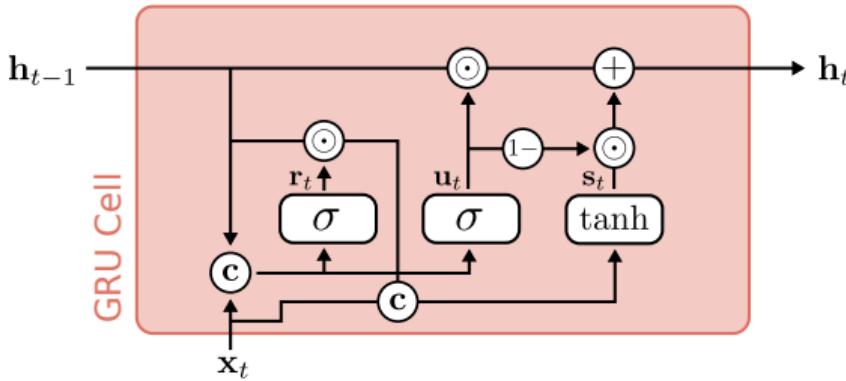
Cell that turns on inside quotes:

"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."

- ▶ UGRNN is able to **keep the state** of a variable **over a long time horizon** ( $u \approx 1$ )
- ▶ For example, it can implement a logic to keep track of being inside a quote or not

# Gated Recurrent Unit



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

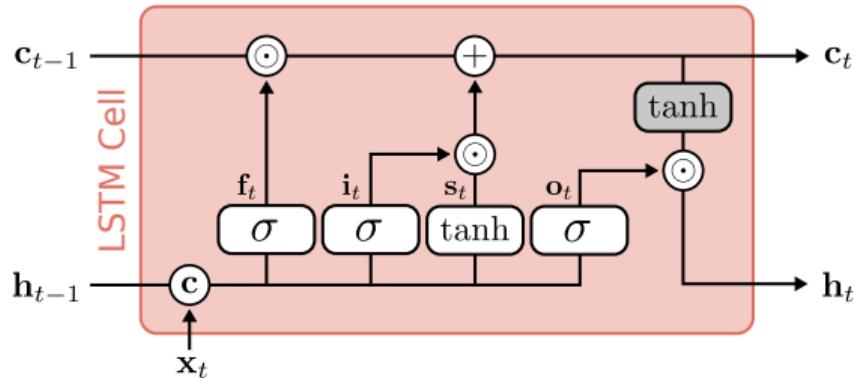
$$\mathbf{u}_t = \sigma(\mathbf{W}_{uh} \mathbf{h}_{t-1} + \mathbf{W}_{ux} \mathbf{x}_t + \mathbf{b}_u)$$

$$\mathbf{s}_t = \tanh(\mathbf{W}_{sh} (\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{W}_{sx} \mathbf{x}_t + \mathbf{b}_s)$$

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

- ▶ **Reset gate** controls which parts of the state are used to compute next target state
- ▶ **Update gate** controls how much information to pass from previous time step

# Long Short-Term Memory



$$f_t = \sigma(\mathbf{W}_{fh} h_{t-1} + \mathbf{W}_{fx} x_t + \mathbf{b}_f)$$

$$i_t = \sigma(\mathbf{W}_{ih} h_{t-1} + \mathbf{W}_{ix} x_t + \mathbf{b}_i)$$

$$o_t = \sigma(\mathbf{W}_{oh} h_{t-1} + \mathbf{W}_{ox} x_t + \mathbf{b}_o)$$

$$s_t = \tanh(\mathbf{W}_{sh} h_{t-1} + \mathbf{W}_{sx} x_t + \mathbf{b}_s)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot s_t$$

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

Passes along an additional cell state  $\mathbf{c}$  in addition to the hidden state  $\mathbf{h}$ . Has 3 gates:

- ▶ **Forget gate** determines information to erase from cell state
- ▶ **Input gate** determines which values of cell state to update
- ▶ **Output gate** determines which elements of cell state to reveal at time  $t$

Remark: Cell update  $\tanh(\cdot)$  creates new target values  $s_t$  for cell state

# UGRNN vs. GRU vs. LSTM

## UGRNN:

- ▶ One gate
- ▶ Expose entire state
- ▶ Single update gate
- ▶ Few parameters

## GRU:

- ▶ Two gates
- ▶ Expose entire state
- ▶ Single update gate
- ▶ Medium parameters

## LSTM:

- ▶ Three gates
- ▶ Control exposure
- ▶ Input/forget gates
- ▶ Many parameters

## A systematic study [Collins et al., 2017] states:

"Our results point to the GRU as being the most learnable of gated RNNs for shallow architectures, followed by the UGRNN."

# 8.4

## Autoregressive Models

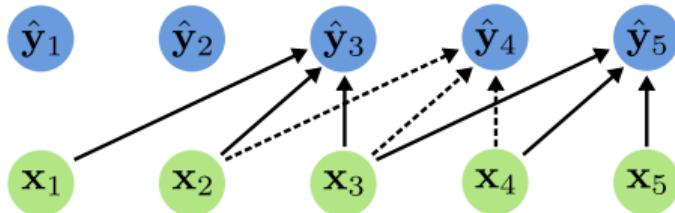
# Autoregressive Models



$$\mathbf{x}_t = f(\mathbf{x}_{t-1}, \mathbf{x}_{t-2}, \dots, \mathbf{x}_{t-k})$$

- ▶ A  $k$ 'th order **autoregressive model** is a **feedforward model** which predicts the next variable  $\mathbf{x}_t$  in a time series based on the  $k$  previous variables  $\mathbf{x}_{t-1}, \mathbf{x}_{t-2}, \dots$
- ▶ As in RNNs, **parameters are shared** across time (same function  $f(\cdot)$  at each  $t$ )
- ▶ Autoregressive models make a strong **conditional independence** assumption

# Autoregressive Models

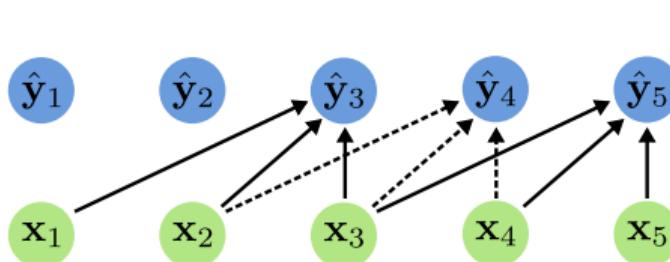


$$\hat{y}_t = f(x_t, x_{t-1}, \dots, x_{t-k})$$

- ▶ This concept can be extended to situations with **different inputs and outputs**
- ▶  $\hat{y}_t$  is **dependent** on  $\{x_i \mid t - k \leq i \leq t\}$  and **independent** of  $\{x_i \mid i < t - k\}$
- ▶ Remark: dashed arrows are only used to visually highlight different time steps

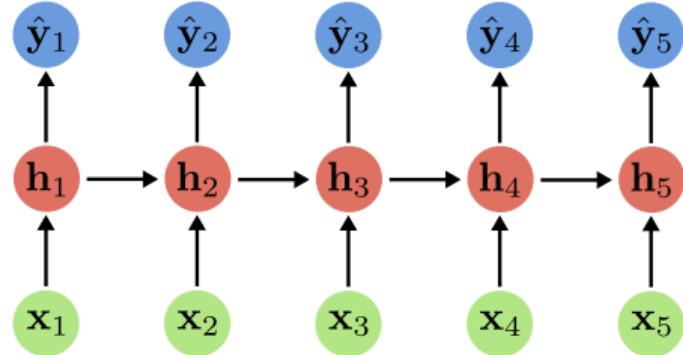
# Autoregressive Models vs. RNNs

Autoregressive Network



$$\hat{y}_t = f(\mathbf{x}_t, \mathbf{x}_{t-1}, \dots, \mathbf{x}_{t-k})$$

Recurrent Neural Network

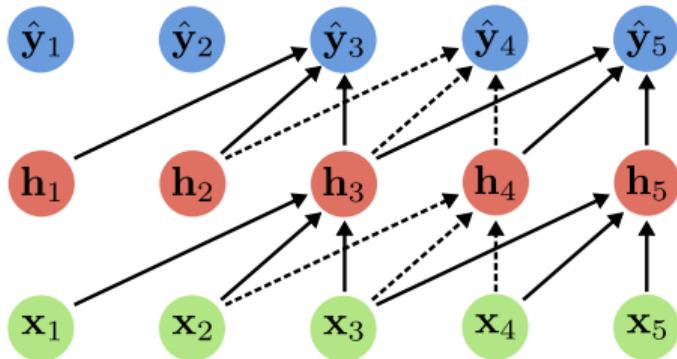


$$\mathbf{h}_t = f_h(\mathbf{h}_{t-1}, \mathbf{x}_t)$$

$$\hat{y}_t = f_y(\mathbf{h}_t)$$

- ▶ Recurrent models summarize past information through their **hidden state h**
- ▶ In contrast to auto-regressive models, RNNs have **infinite memory**
- ▶ Autoregressive models are **easier to train** (no backprop through time)

# Multi-Layer Autoregressive Models

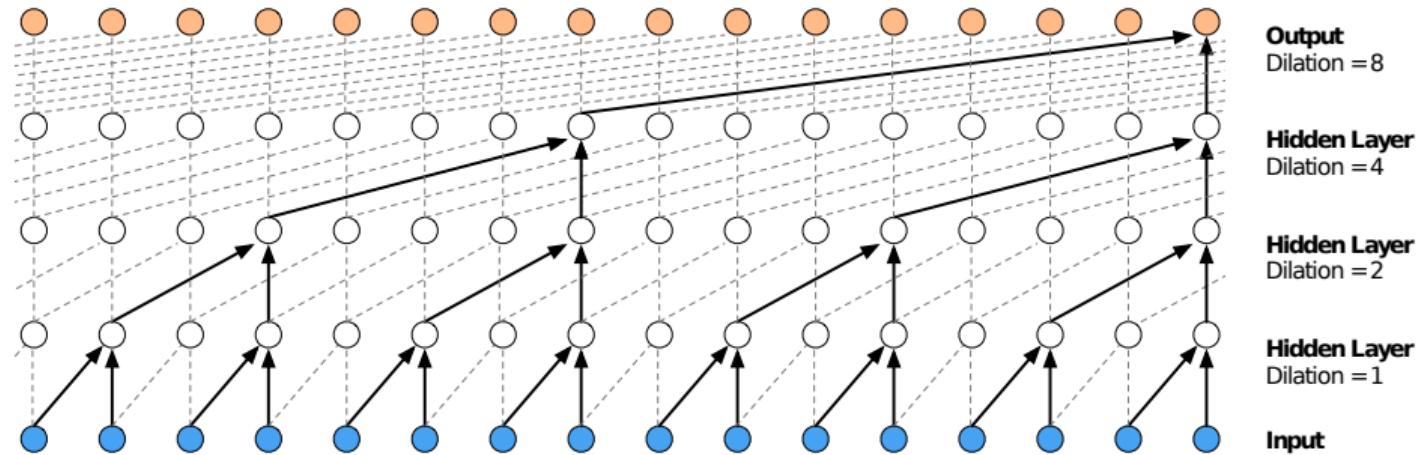


$$\mathbf{h}_t = f_1(\mathbf{x}_t, \mathbf{x}_{t-1}, \dots, \mathbf{x}_{t-k})$$

$$\hat{\mathbf{y}}_t = f_2(\mathbf{h}_t, \mathbf{h}_{t-1}, \dots, \mathbf{h}_{t-k})$$

- ▶ Autoregressive models can be extended to deep models with **multiple layers**
- ▶ They effectively perform multiple causal **temporal convolutions**
- ▶ They can be combined with residual connections and dilated convolutions

# WaveNet



## WaveNet:

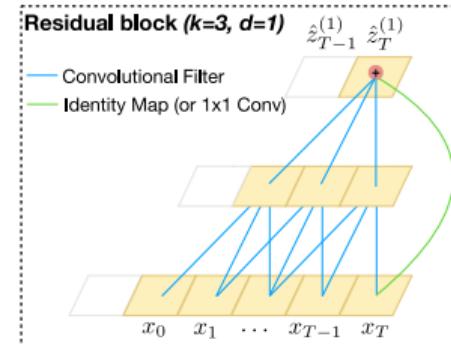
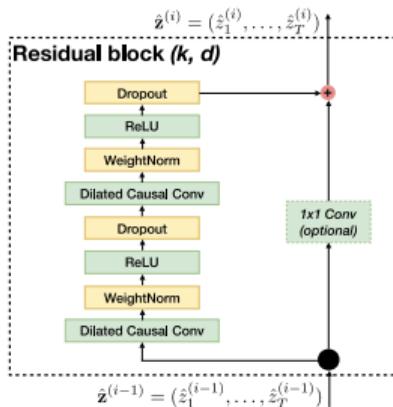
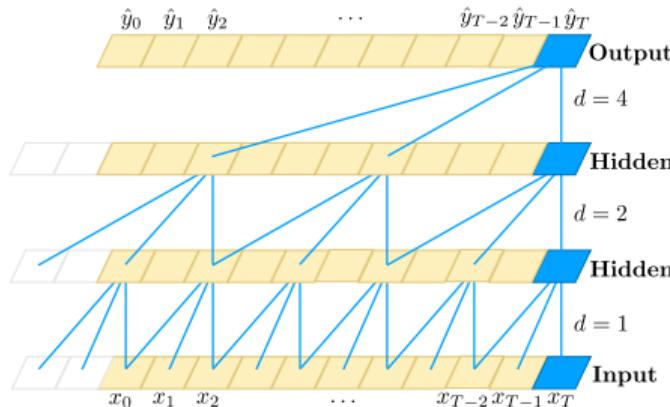
- Generative model for **raw audio waveforms** which generates realistic speech
- Combines **dilated convolutions** with **residual and skip connections**
- Audio is hard: 16k samples per second, structure at **multiple time scales**

# WaveNet

Speech samples	Subjective 5-scale MOS in naturalness	
	North American English	Mandarin Chinese
LSTM-RNN parametric	3.67 ± 0.098	3.79 ± 0.084
HMM-driven concatenative	3.86 ± 0.137	3.47 ± 0.108
<b>WaveNet (L+F)</b>	<b>4.21</b> ± 0.081	<b>4.08</b> ± 0.085
Natural (8-bit $\mu$ -law)	4.46 ± 0.067	4.25 ± 0.082
Natural (16-bit linear PCM)	4.55 ± 0.075	4.21 ± 0.071

- WaveNet outperforms traditional LSTMs on speech synthesis by large margin

# Temporal Convolution Networks



## Temporal Convolution Networks:

- Autoregressive model with **zero-padding** to handle sequences of arbitrary length
- **Deep multi-layer** network with **residual layers**, but otherwise simple architecture
- **Dilated convolutions** to increase the receptive field size (often called “context”)

# Temporal Convolution Networks

Sequence Modeling Task	Model Size ( $\approx$ )	Models			
		LSTM	GRU	RNN	TCN
Seq. MNIST (accuracy <sup>h</sup> )	70K	87.2	96.2	21.5	<b>99.0</b>
Permuted MNIST (accuracy)	70K	85.7	87.3	25.3	<b>97.2</b>
Adding problem $T=600$ (loss <sup>ℓ</sup> )	70K	0.164	<b>5.3e-5</b>	0.177	<b>5.8e-5</b>
Copy memory $T=1000$ (loss)	16K	0.0204	0.0197	0.0202	<b>3.5e-5</b>
Music JSB Chorales (loss)	300K	8.45	8.43	8.91	<b>8.10</b>
Music Nottingham (loss)	1M	3.29	3.46	4.05	<b>3.07</b>
Word-level PTB (perplexity <sup>ℓ</sup> )	13M	<b>78.93</b>	92.48	114.50	88.68
Word-level Wiki-103 (perplexity)	-	48.4	-	-	<b>45.19</b>
Word-level LAMBADA (perplexity)	-	4186	-	14725	<b>1279</b>
Char-level PTB (bpcl)	3M	1.36	1.37	1.48	<b>1.31</b>
Char-level text8 (bpcl)	5M	1.50	1.53	1.69	<b>1.45</b>

- TCNs outperform LSTMs, GRUs and RNNs on **several different sequence tasks**

## Discussion

- ▶ The unlimited context offered by recurrent models is not strictly necessary for language modeling [Dauphin, 2016]
- ▶ Models that truncate at sequence lengths of 13 [Chelba, 2017] or 25 [Dauphin, 2016] are often competitive with models that have infinite memory
- ▶ The “infinite memory” advantage of RNNs is largely absent in practice, even when considering a variety of different models and tasks [Bai, 2018]
- ▶ Theoretical result: If the recurrent model is stable (meaning the gradients can not explode), then the model can be well-approximated by a feed-forward network for the purposes of both inference and training [Miller, 2018]
- ▶ <http://www.offconvex.org/2018/07/27/approximating-recurrent/>