

# CS182/282A: Designing, Visualizing and Understanding Deep Neural Networks

---

**John Canny**

Spring 2020

Lecture 9: Recurrent Networks, LSTMs and  
Applications

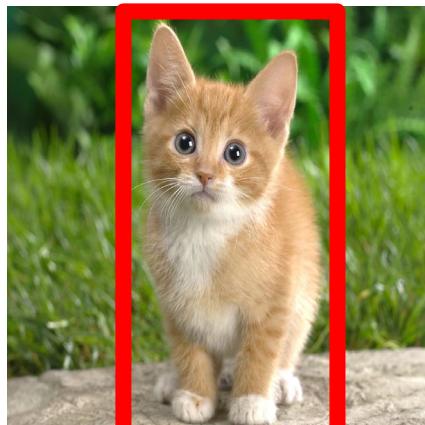
# Last time: Localization and Detection

**Classification**



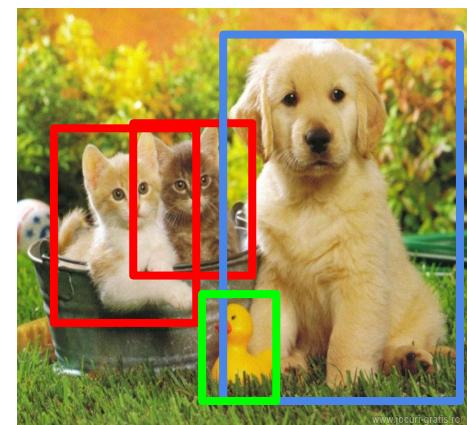
CAT

**Classification + Localization**



Single  
object

**Object Detection**



CAT, DOG, DUCK

**Instance Segmentation**



CAT, DOG, DUCK

Multiple  
objects

# Updates

Midterm 1 is coming up on 2/26, from 7-8:30pm. Its closed-book with one 2-sided sheet of notes.

Practice MT solutions coming soon. There should be a review session next week. Still working on scheduling.

Assignment 1 is due tomorrow.

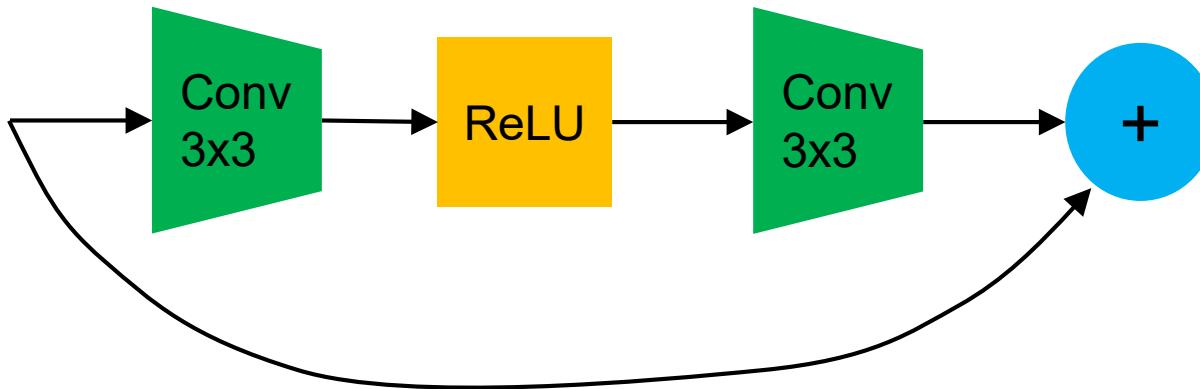
CS282A proposals are due today.

CS182 proposals (really team formation) are due next Monday.

# Neural Network structure

Standard Neural Networks are DAGs (Directed Acyclic Graphs). That means they have a topological ordering.

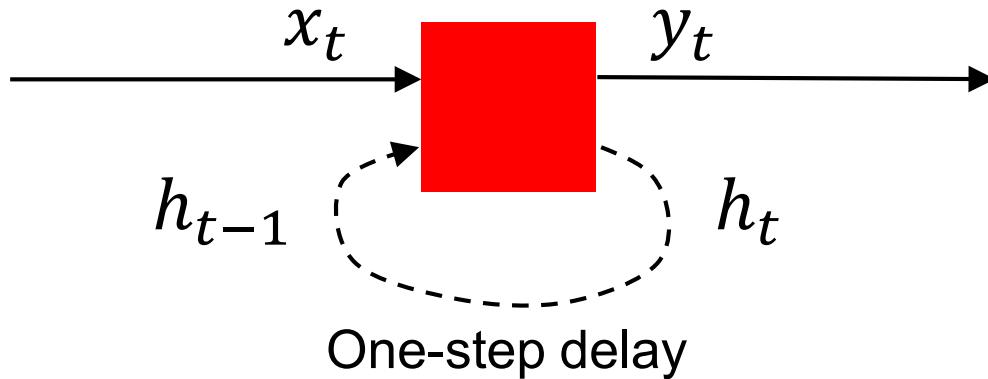
- The topological ordering is used for activation propagation, and for gradient back-propagation.



- These networks process one input minibatch at a time.

# Recurrent Neural Networks (RNNs)

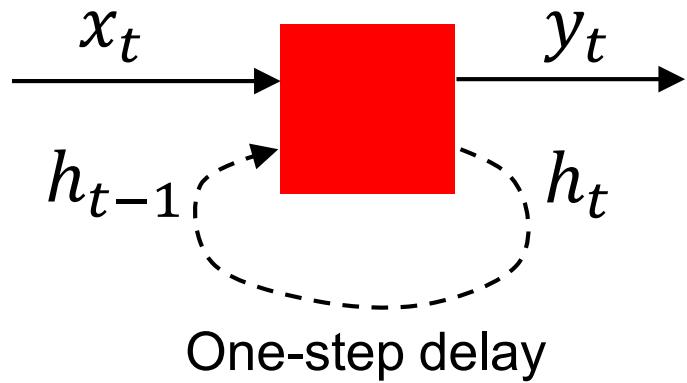
Recurrent networks introduce cycles and a notion of time.



- They are designed to process sequences of data  $x_1, \dots, x_n$  and can produce sequences of outputs  $y_1, \dots, y_m$ .

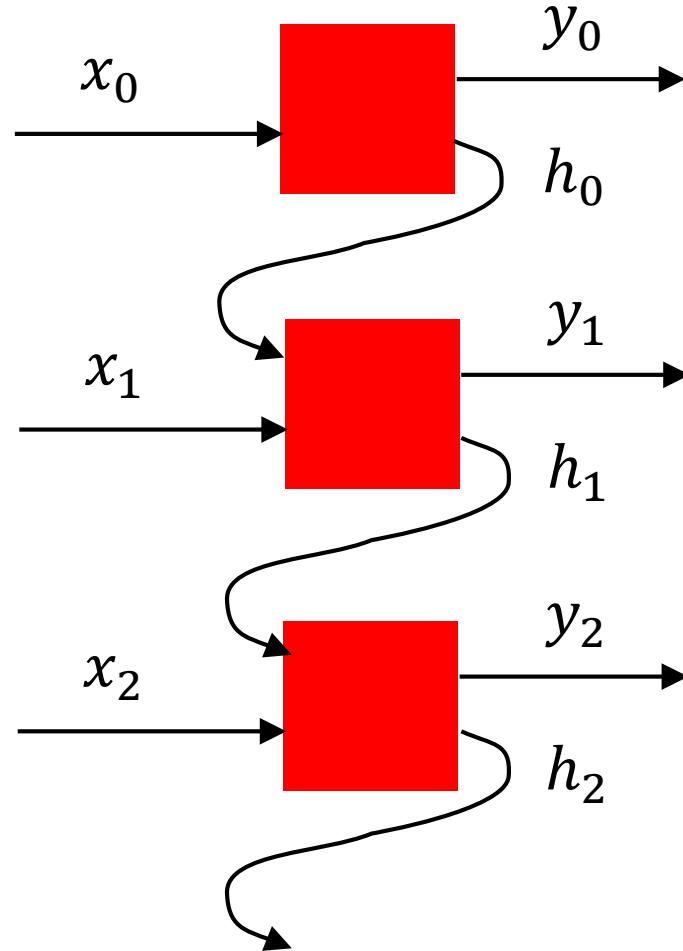
# Unrolling RNNs

RNNs can be unrolled across multiple time steps.



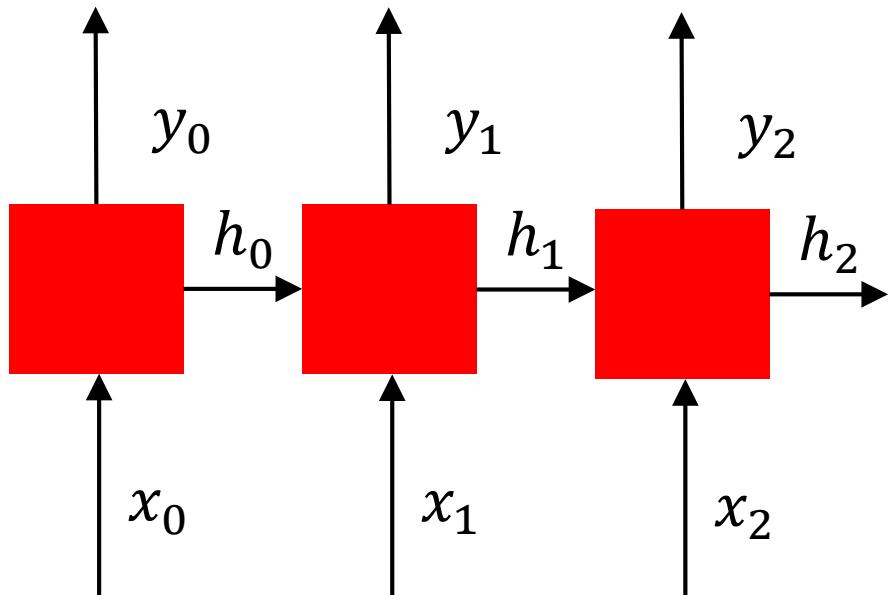
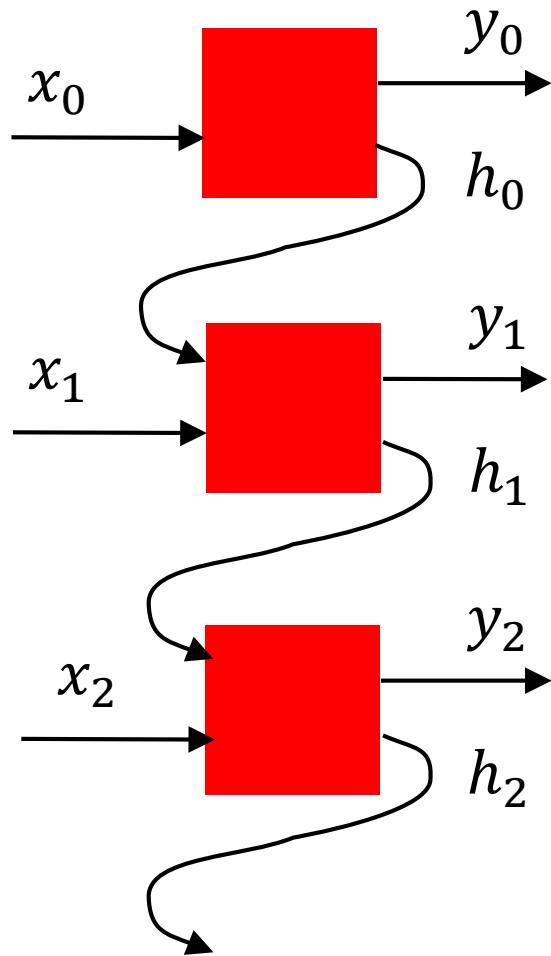
This produces a DAG which supports backpropagation.

But its size depends on the input sequence length.



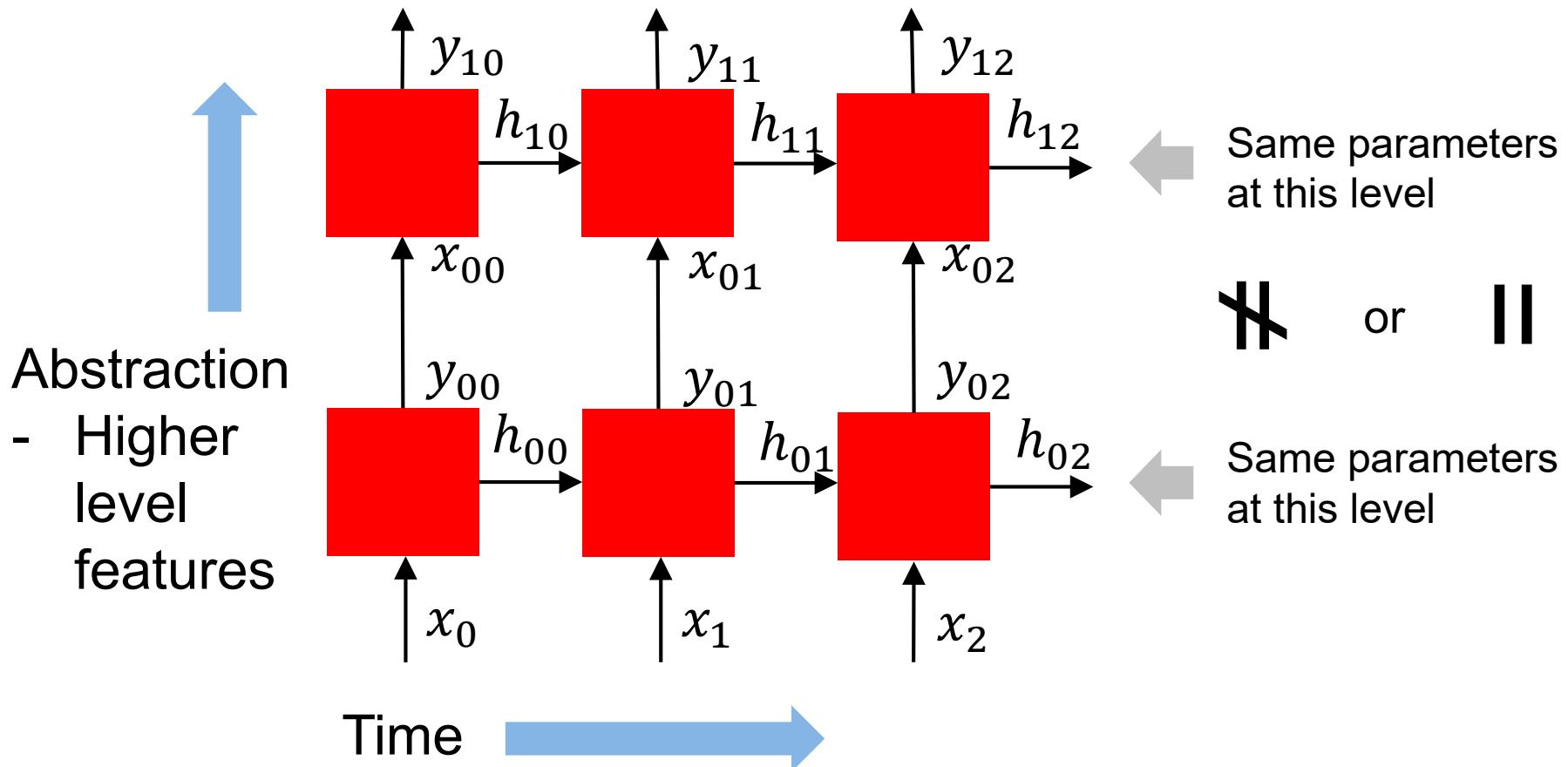
# Unrolling RNNs

Usually drawn as:



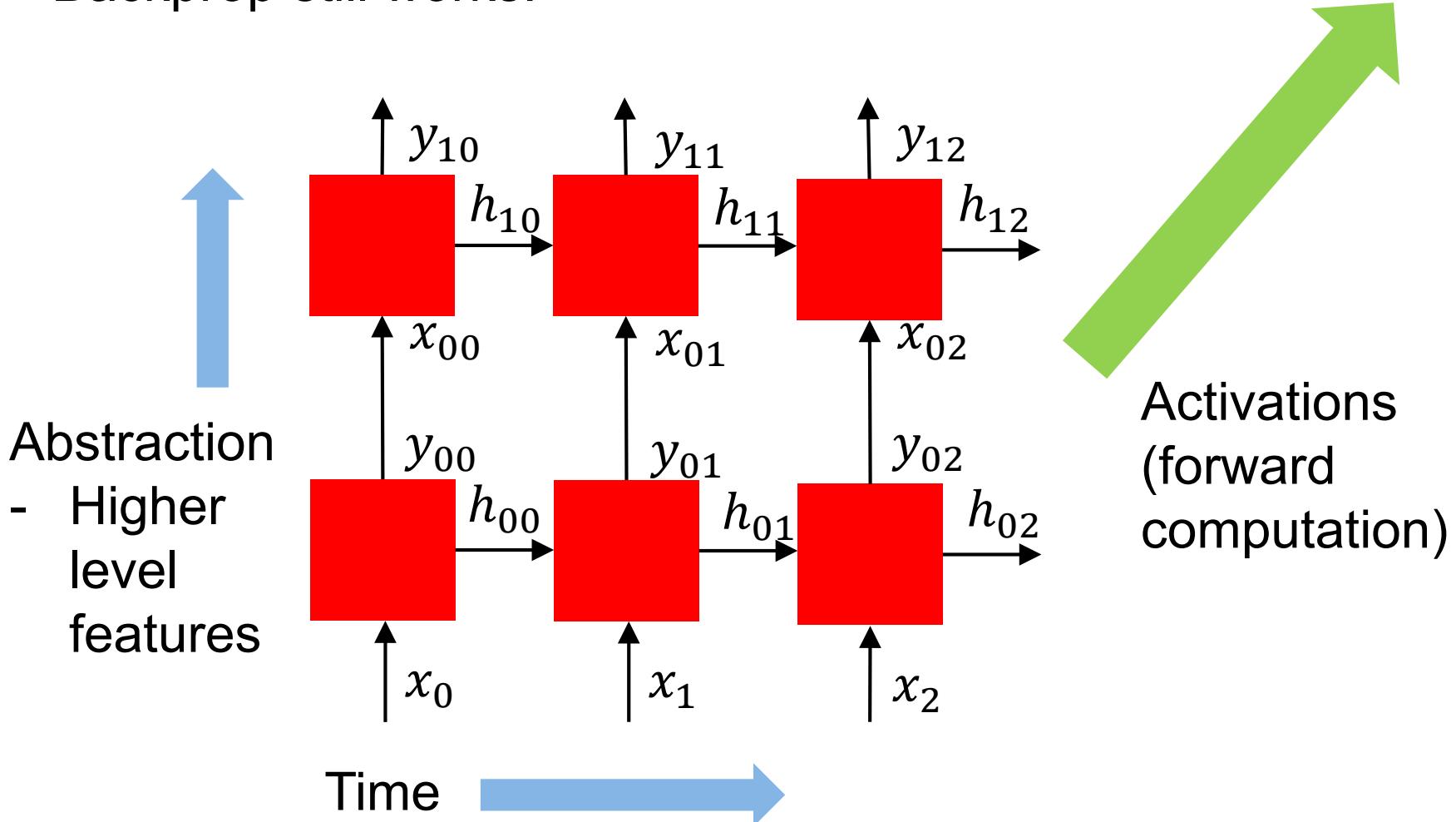
# RNN structure

Often layers are stacked vertically (deep RNNs):



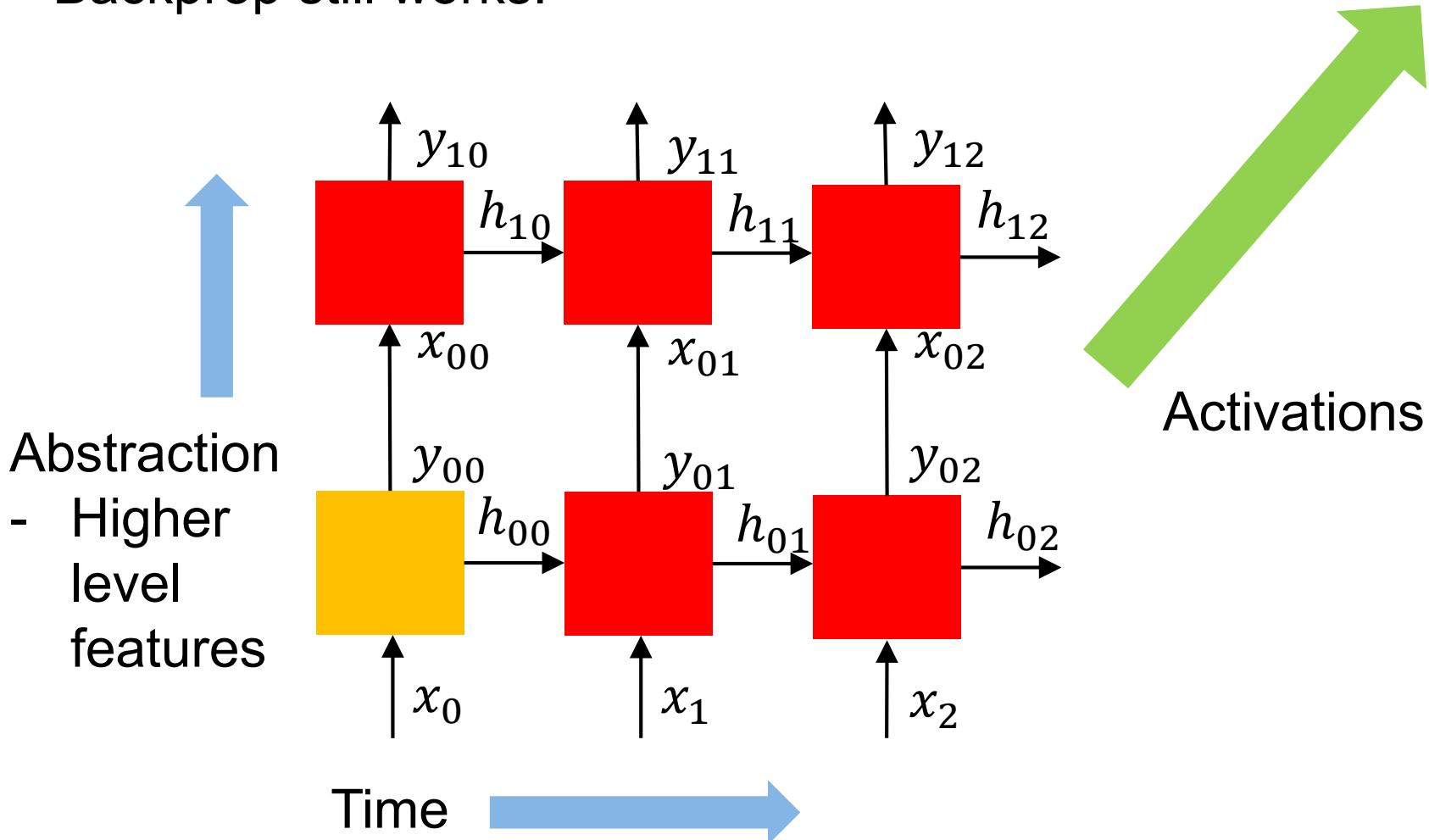
# RNN structure

Backprop still works:



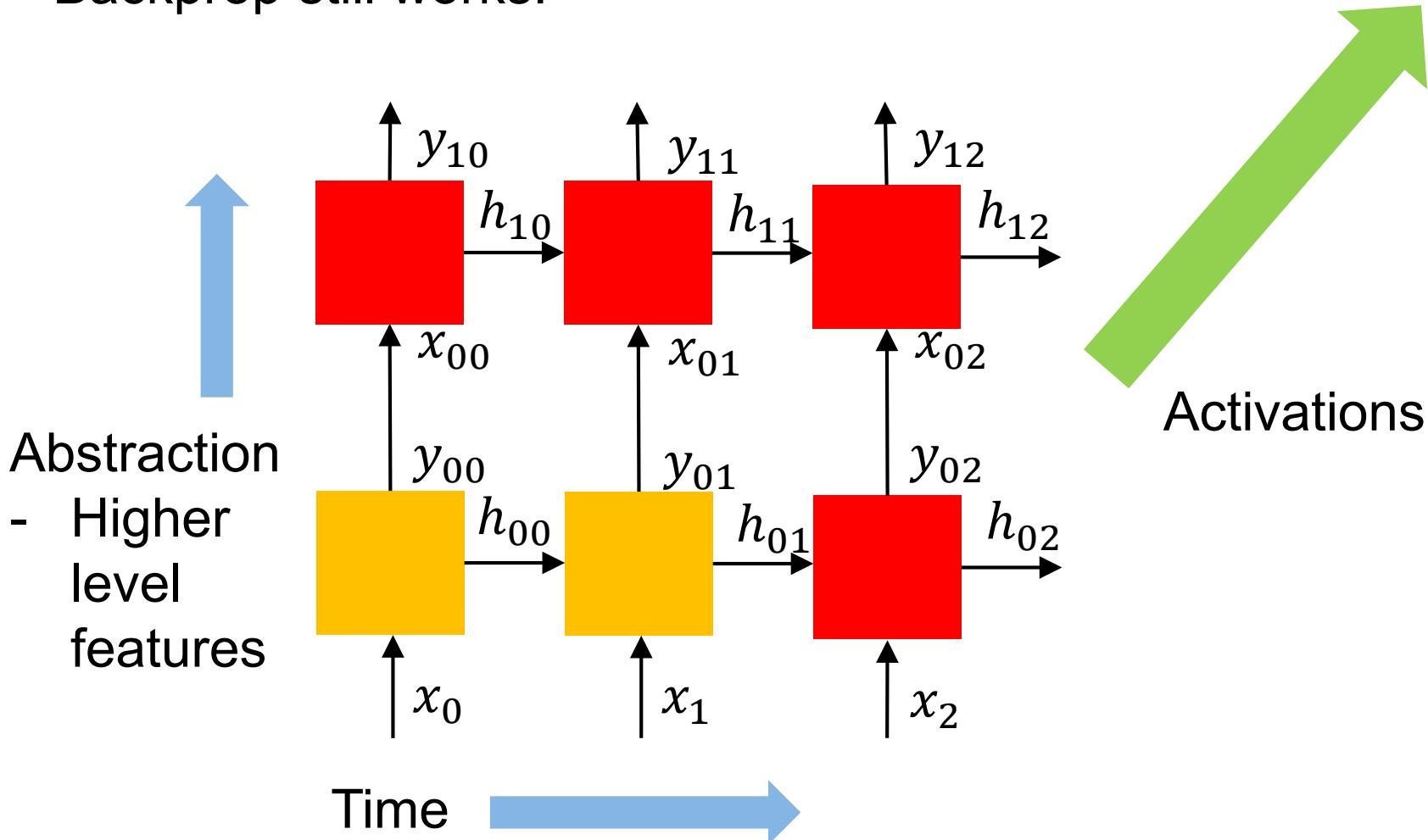
# RNN structure

Backprop still works:



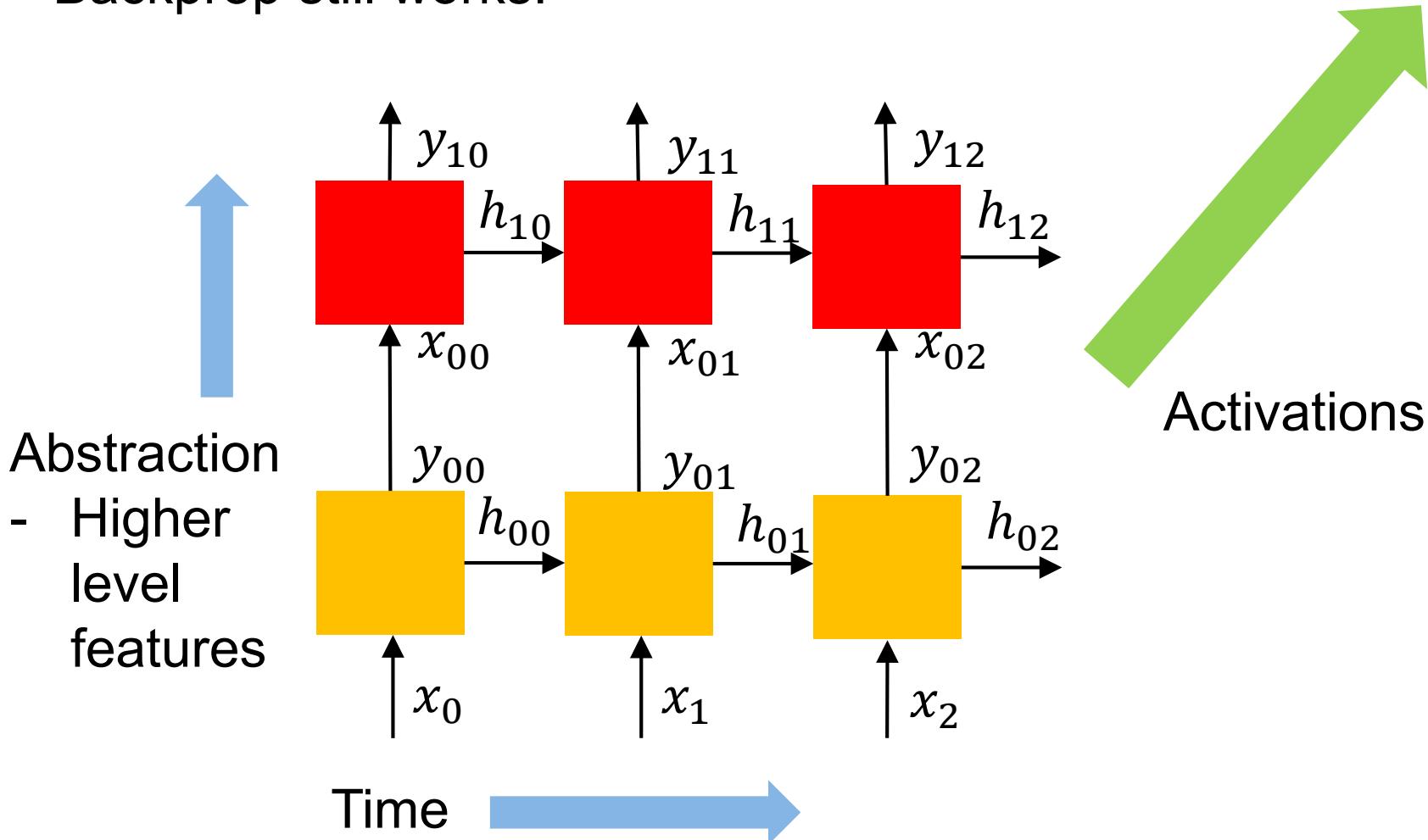
# RNN structure

Backprop still works:



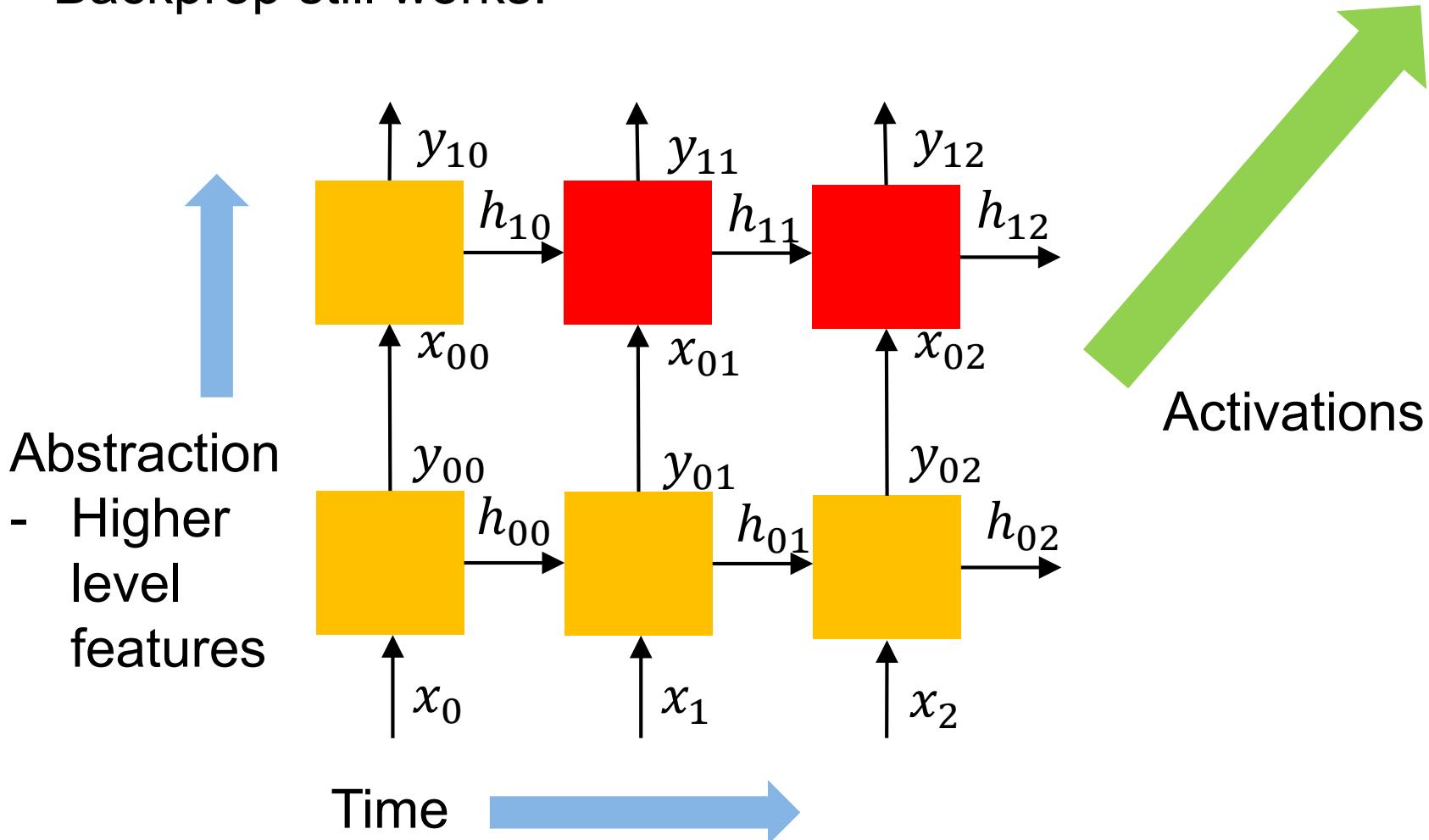
# RNN structure

Backprop still works:



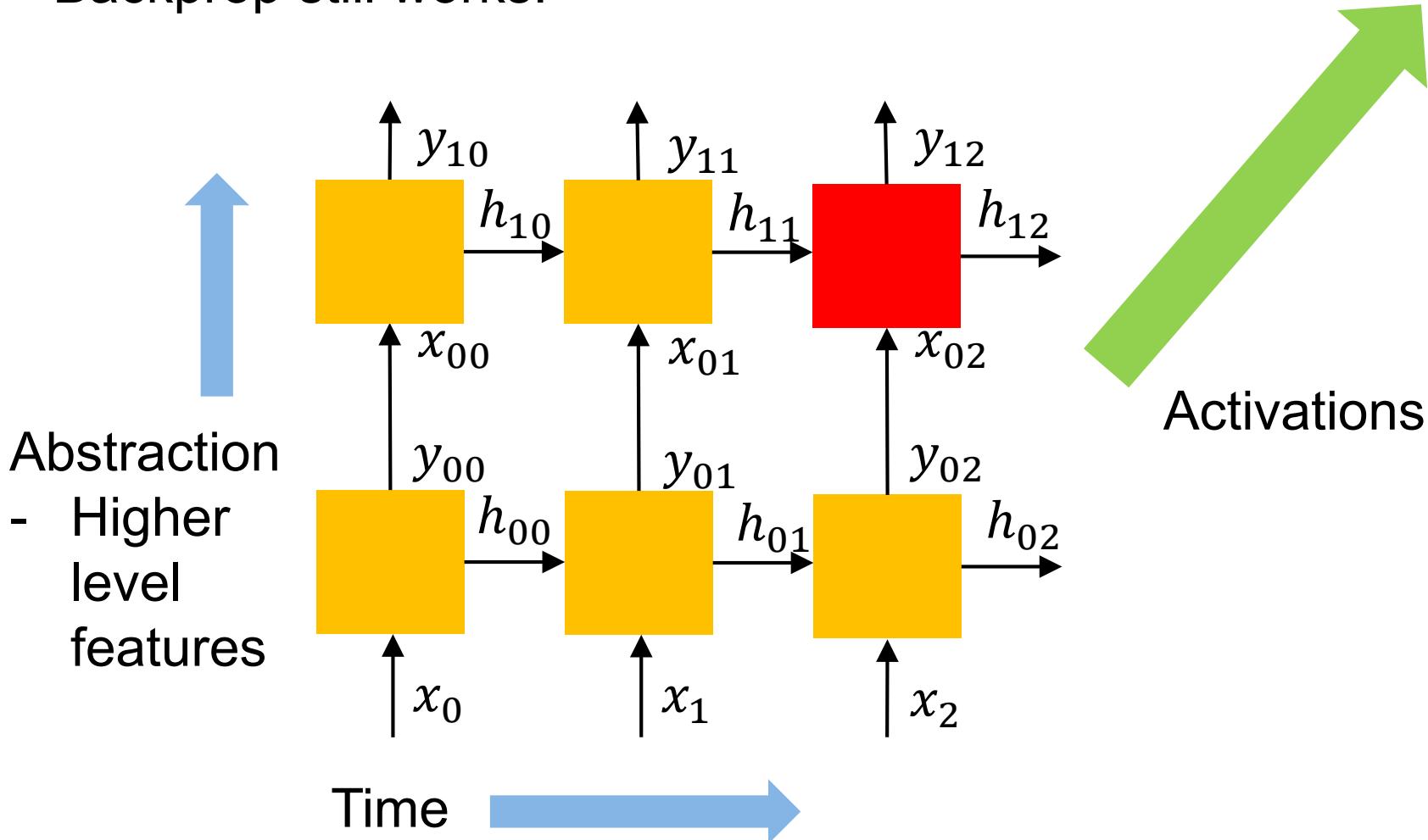
# RNN structure

Backprop still works:



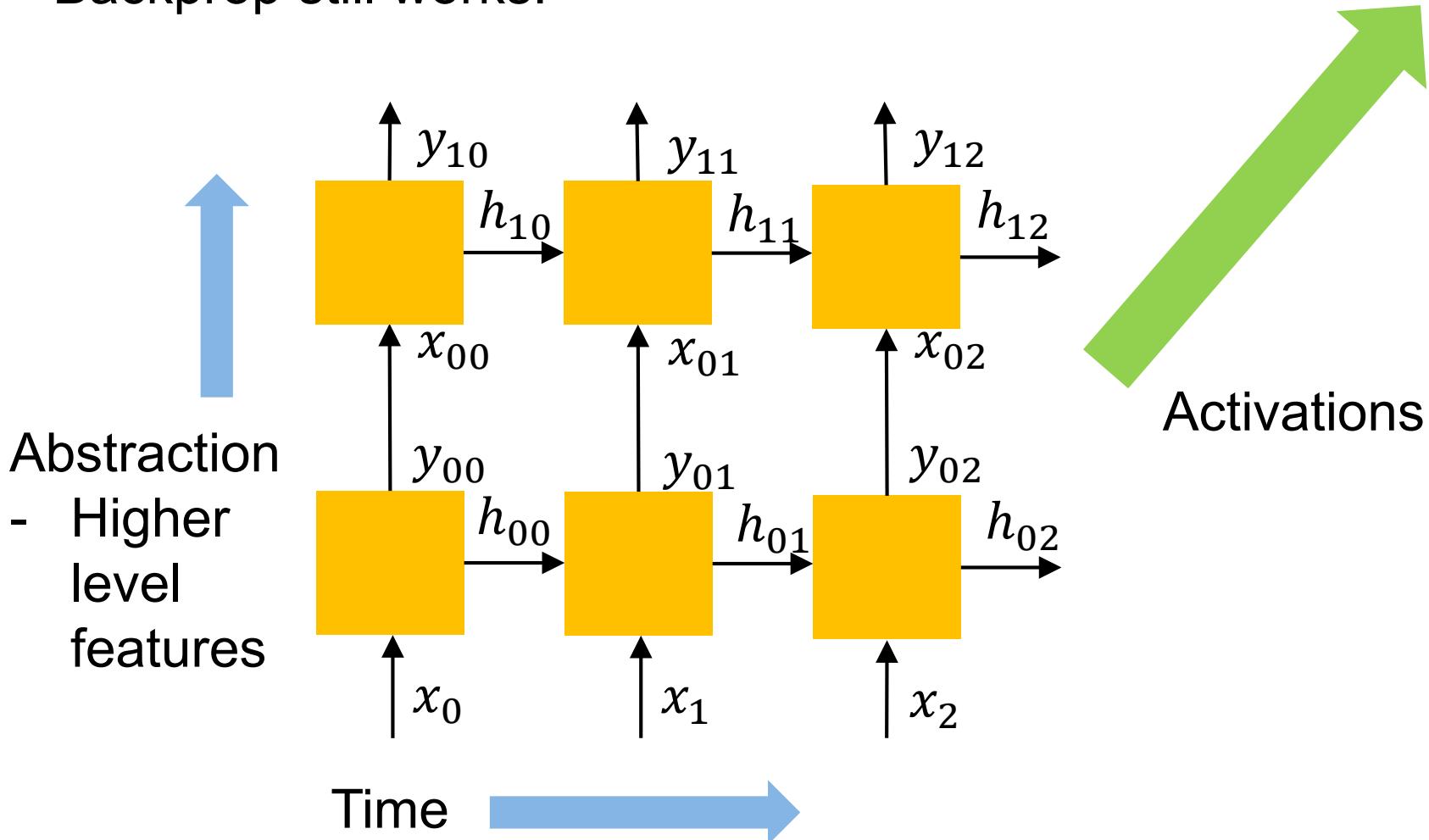
# RNN structure

Backprop still works:



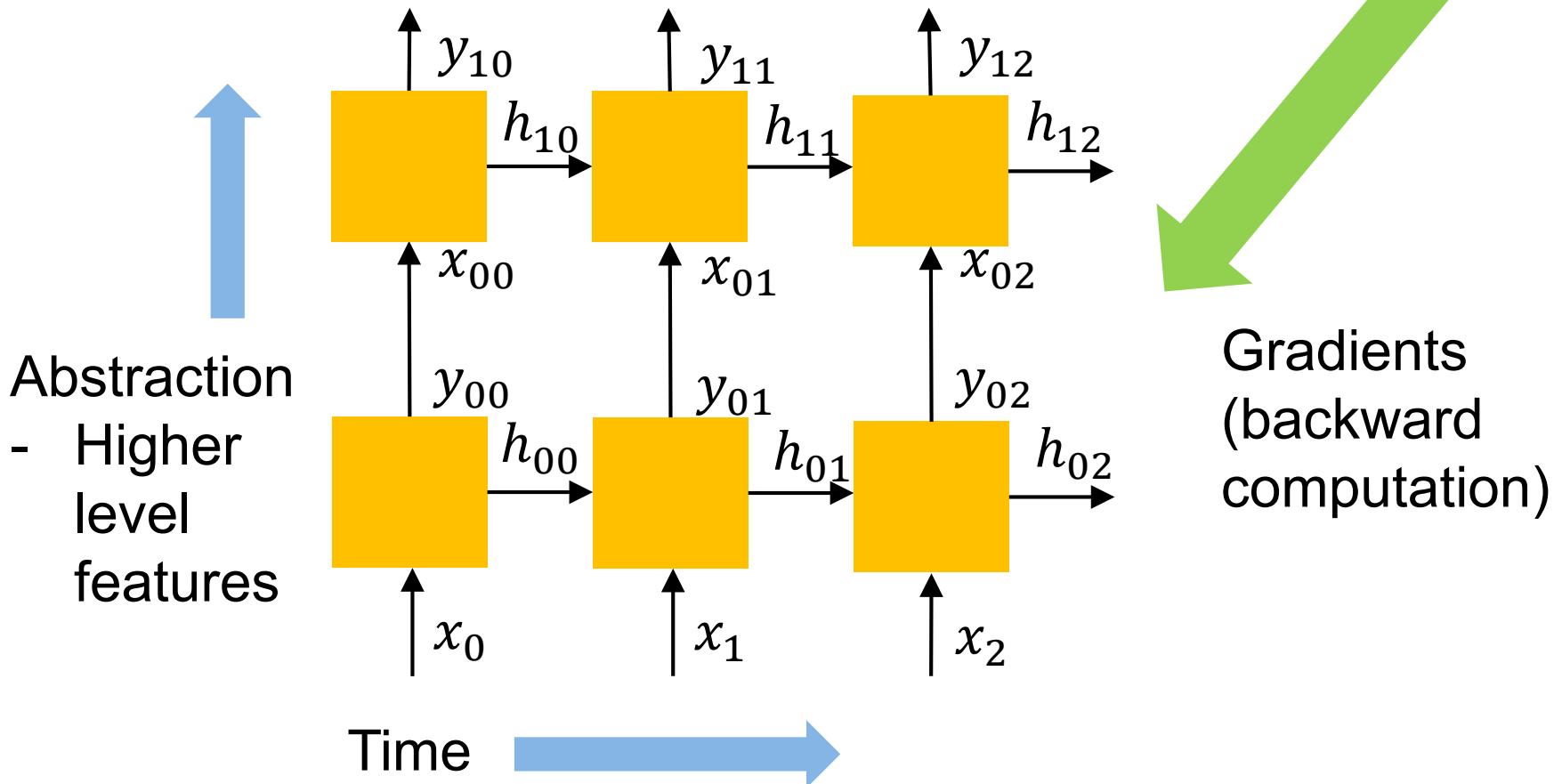
# RNN structure

Backprop still works:



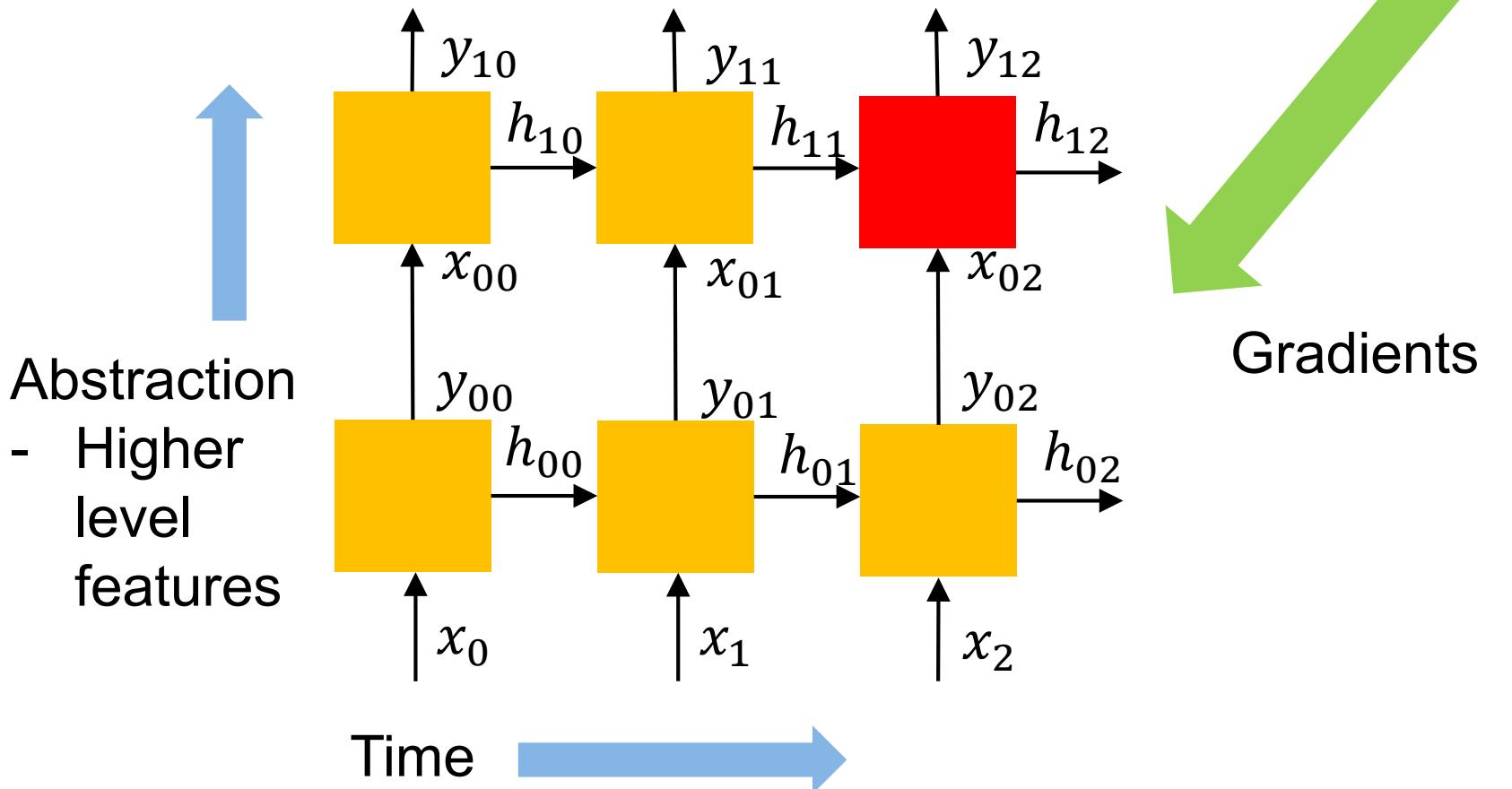
# RNN structure

Backprop still works:



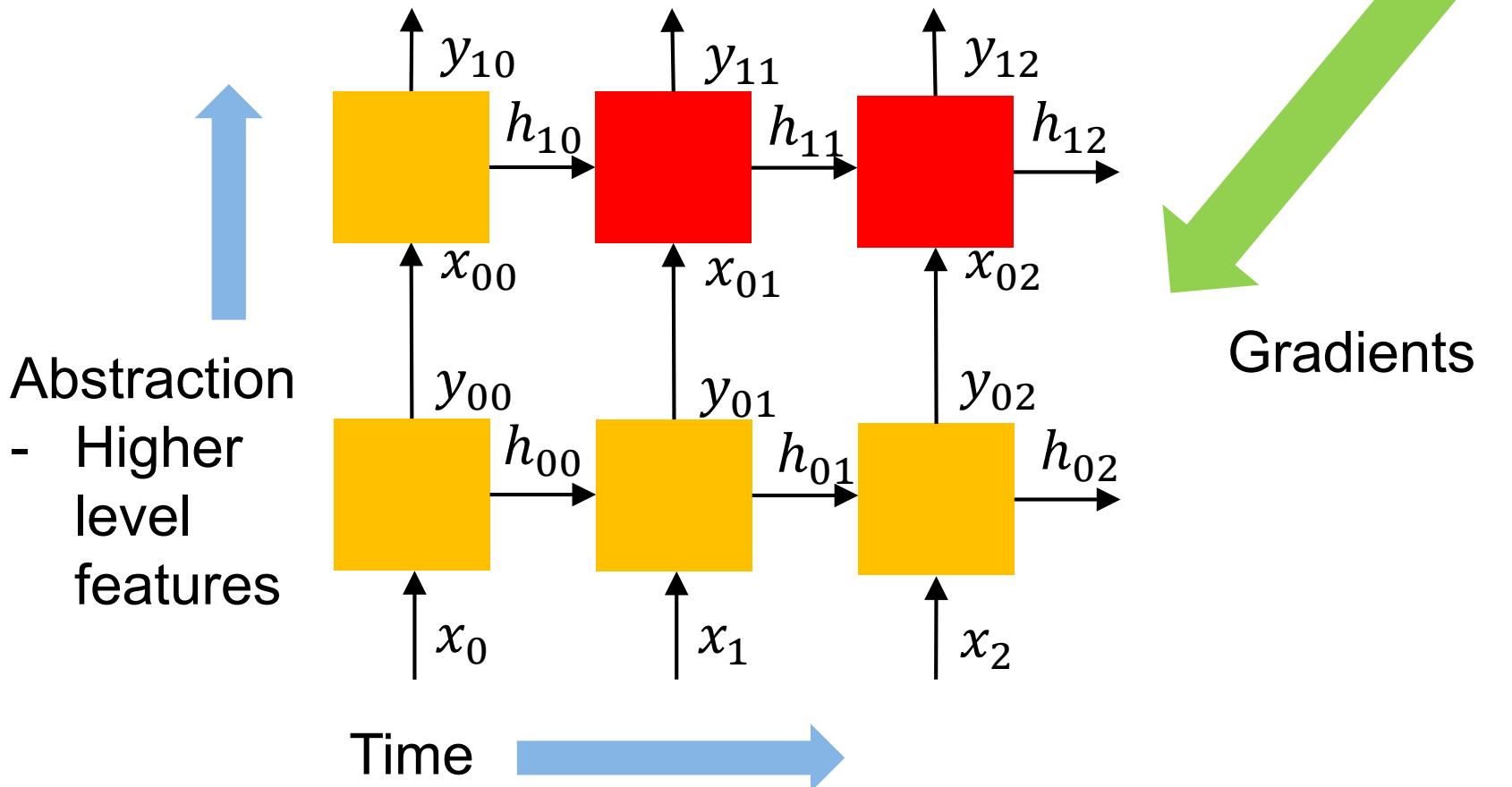
# RNN structure

Backprop still works:



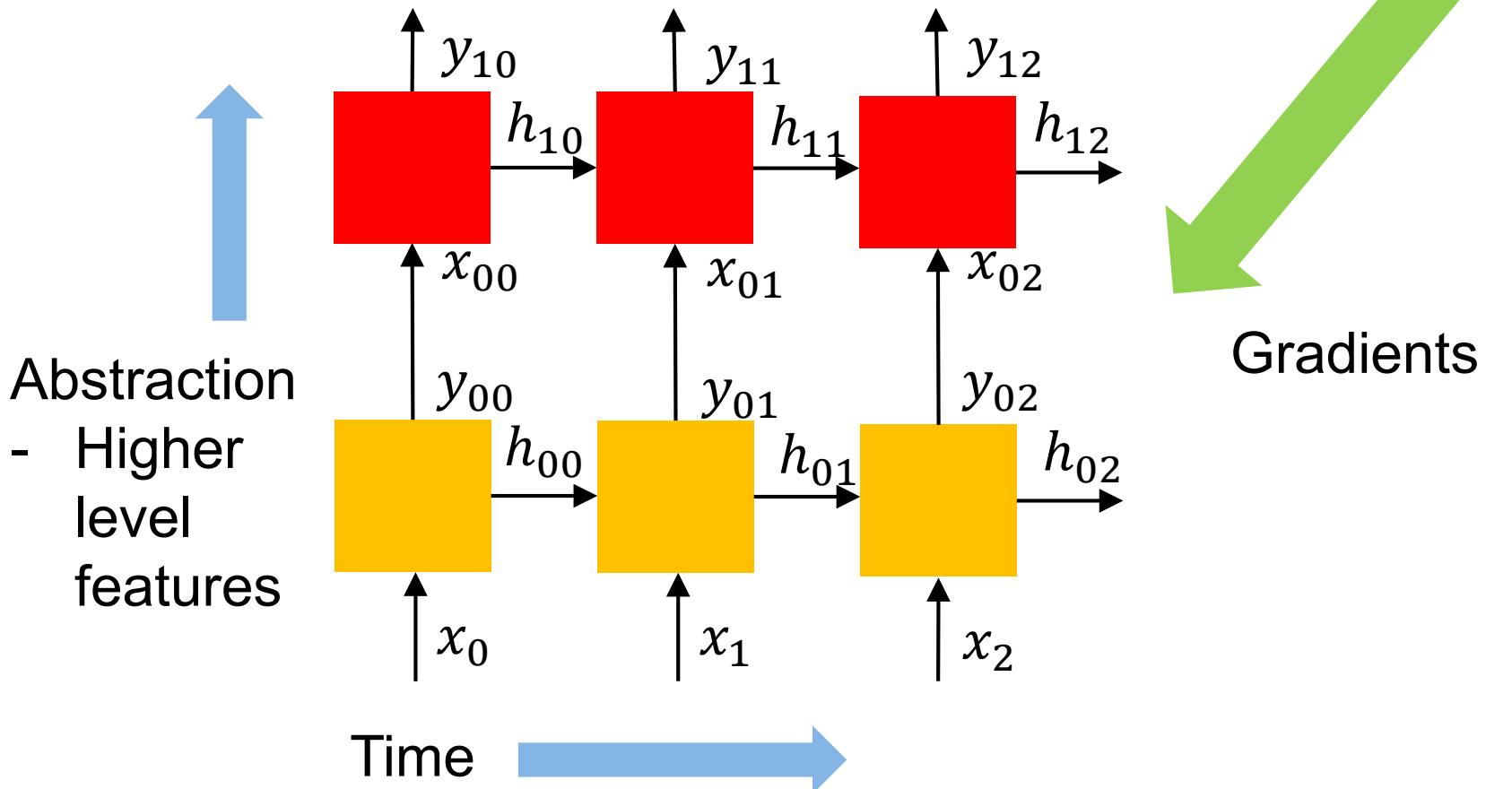
# RNN structure

Backprop still works:



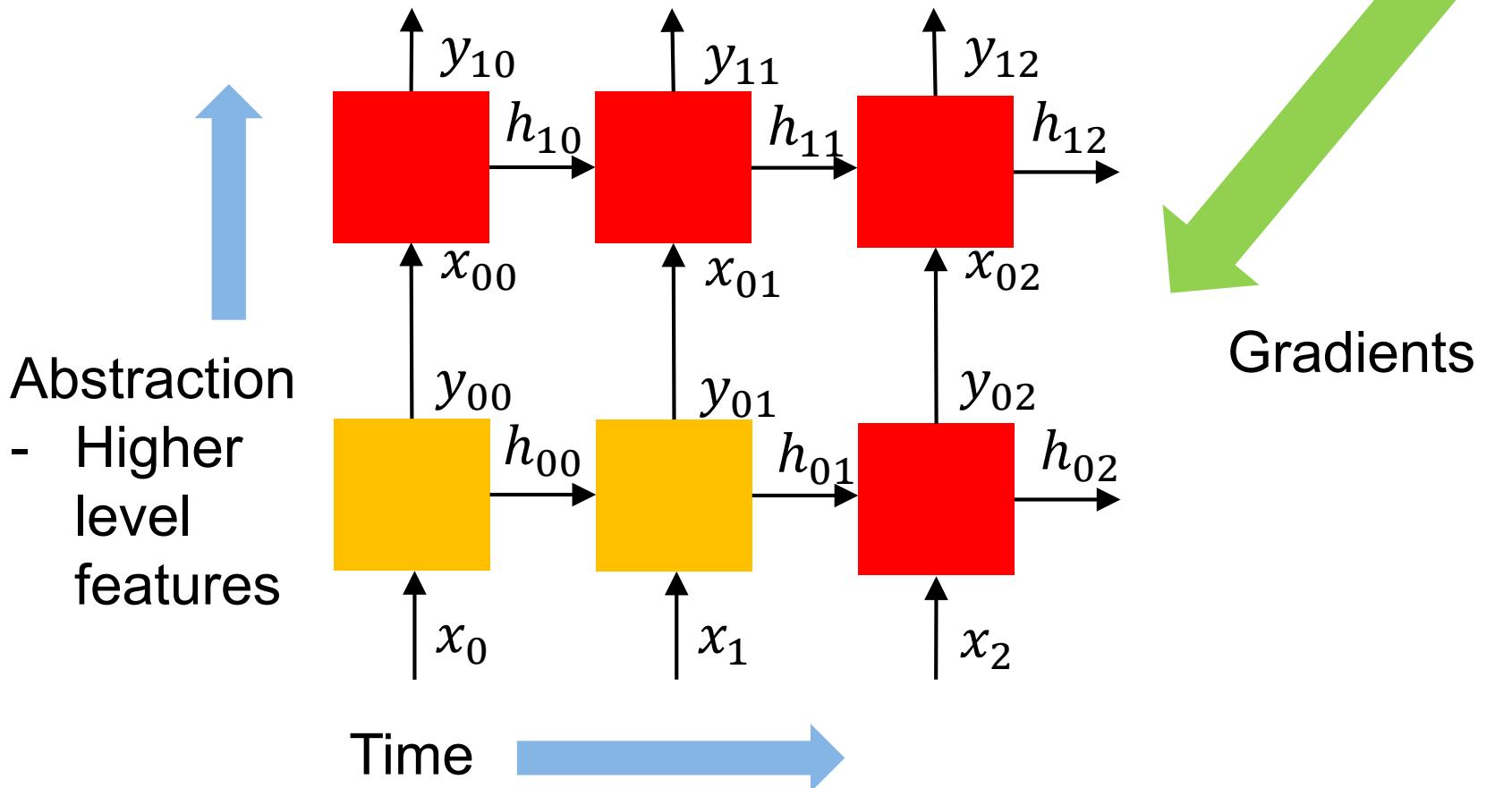
# RNN structure

Backprop still works:



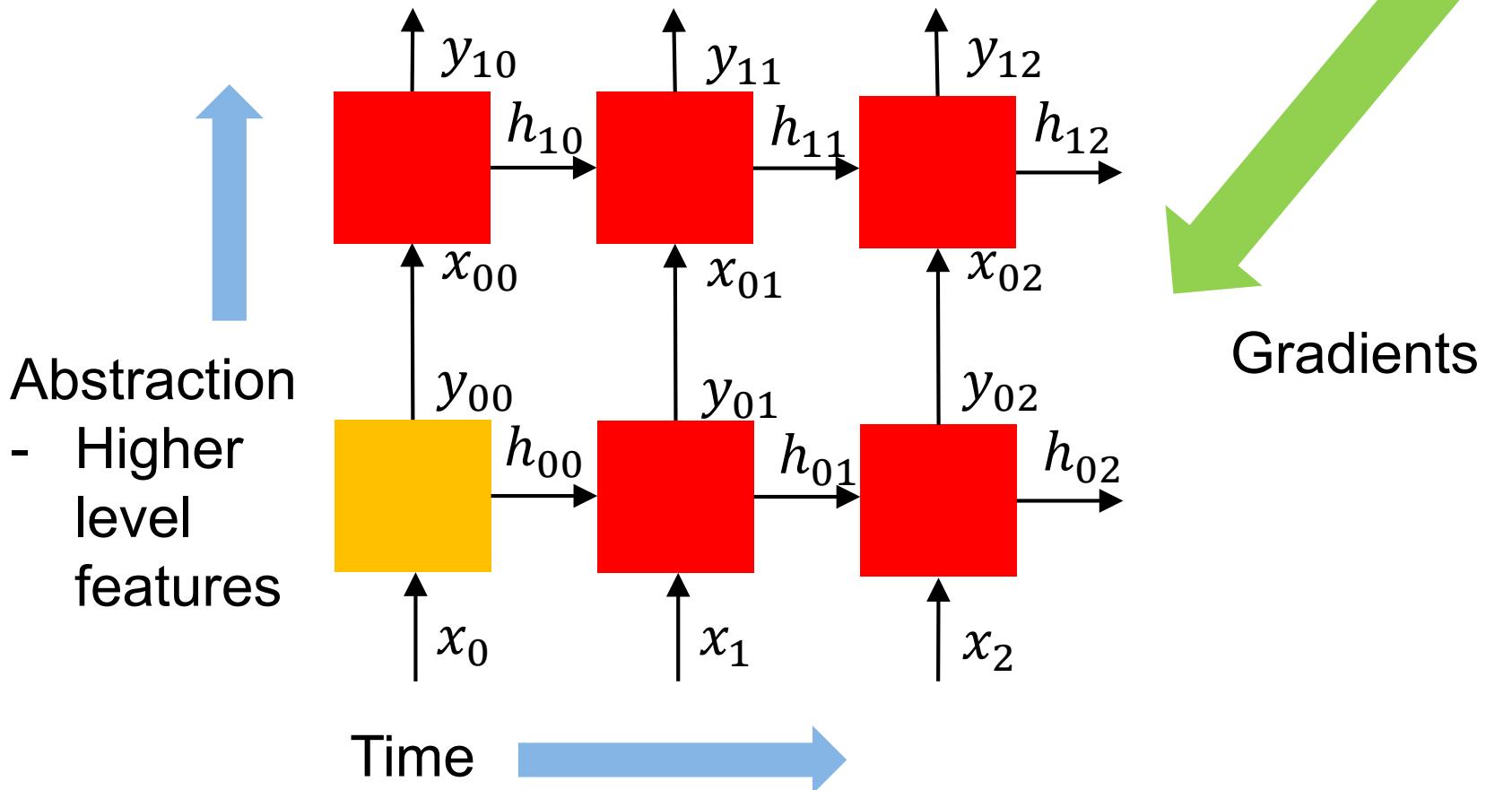
# RNN structure

Backprop still works:



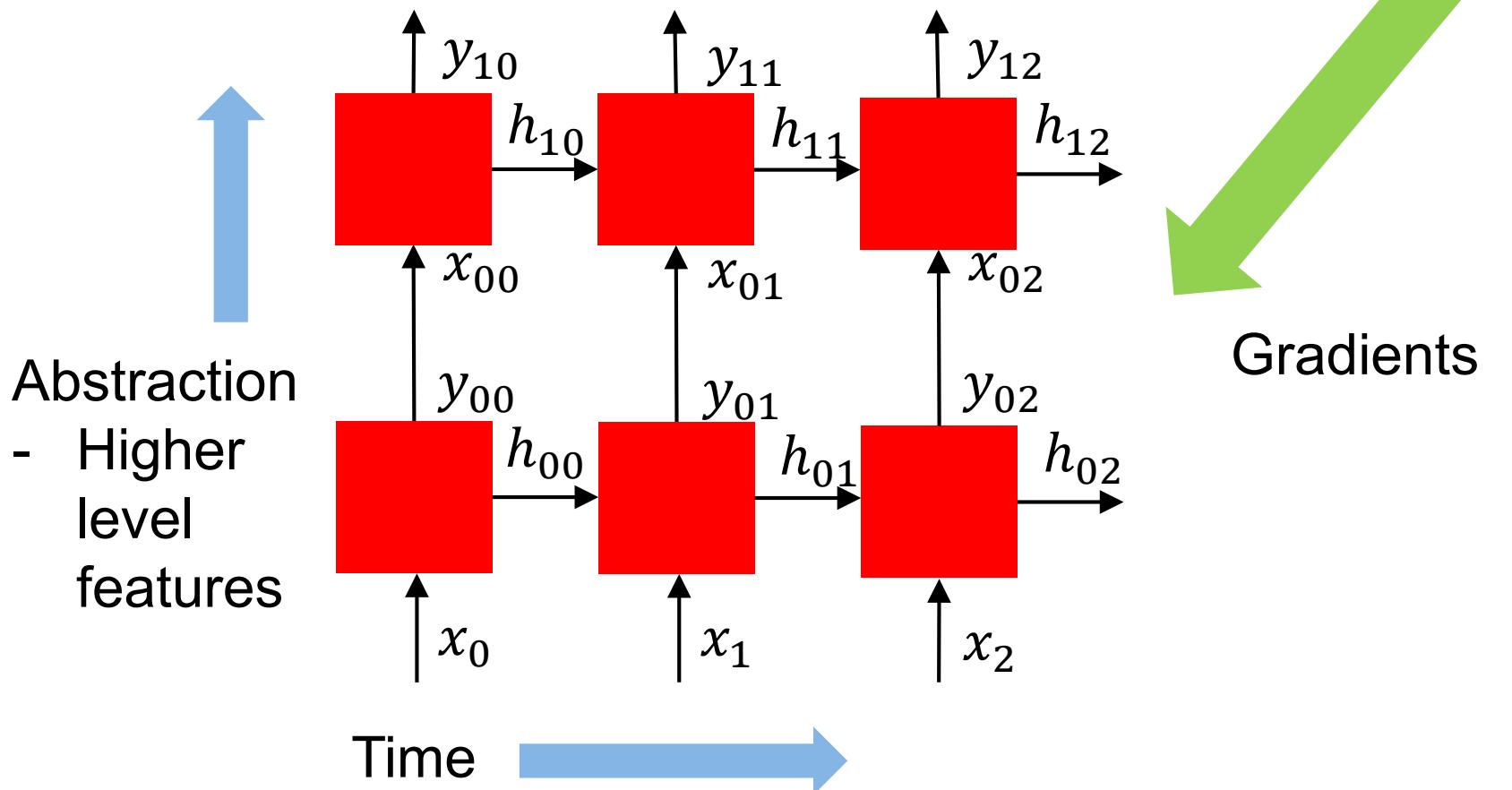
# RNN structure

Backprop still works:



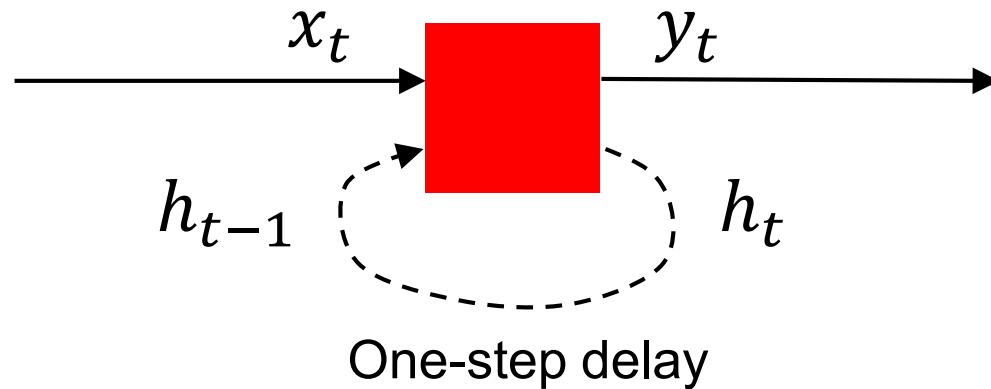
# RNN structure

Backprop still works:



# RNN unrolling

Question: Can you run forward/backward inference on an unrolled RNN, i.e. keeping only one copy of its state?

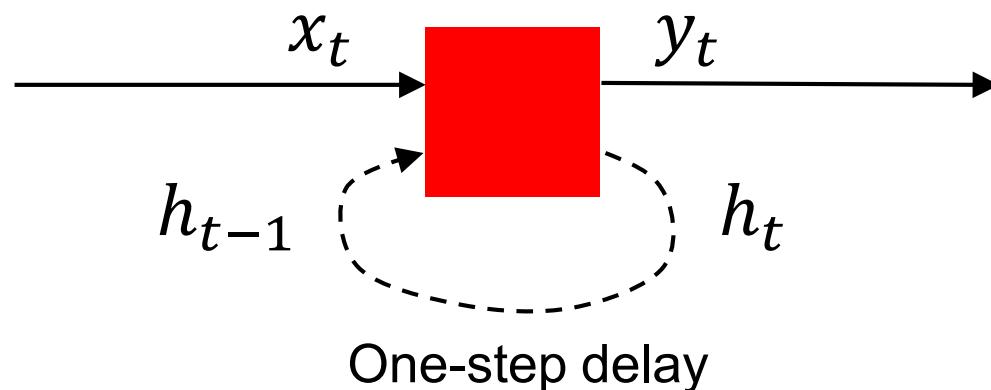


Run Forward Y/N, Run Backward Y/N

- A. No, No
- B. No, Yes
- C. Yes, No
- D. Yes, Yes

# Oops!

Question: Can you run forward/backward inference on an unrolled RNN, i.e. keeping only one copy of its state?



Run Forward Y/N, Run Backward Y/N

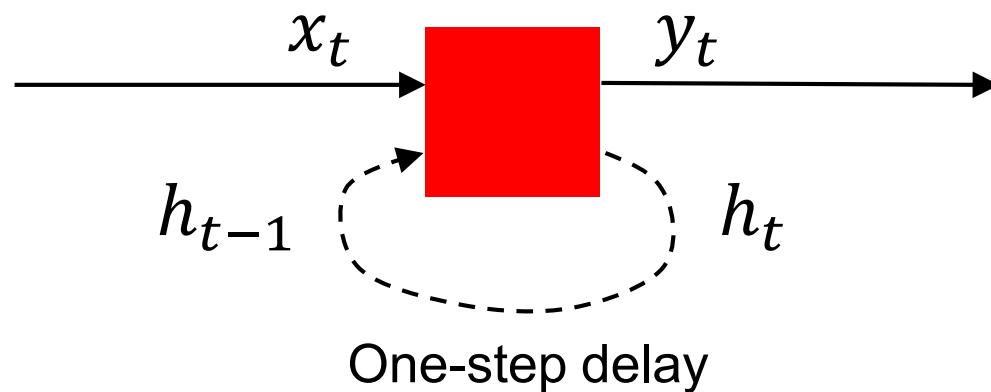
A. No, No

Try Again

Continue

# Oops!

Question: Can you run forward/backward inference on an unrolled RNN, i.e. keeping only one copy of its state?



Run Forward Y/N, Run Backward Y/N

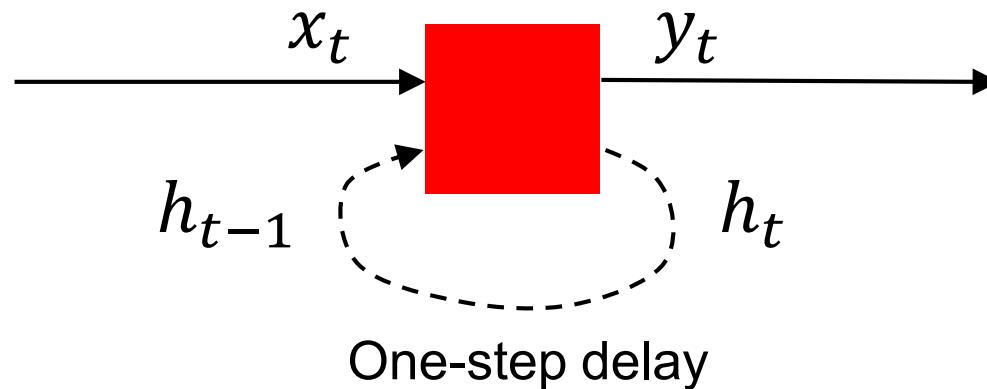
B. No, Yes

Try Again

Continue

# Correct!

Question: Can you run forward/backward inference on an unrolled RNN, i.e. keeping only one copy of its state?



Run Forward Y/N, Run Backward Y/N

C. Yes, No

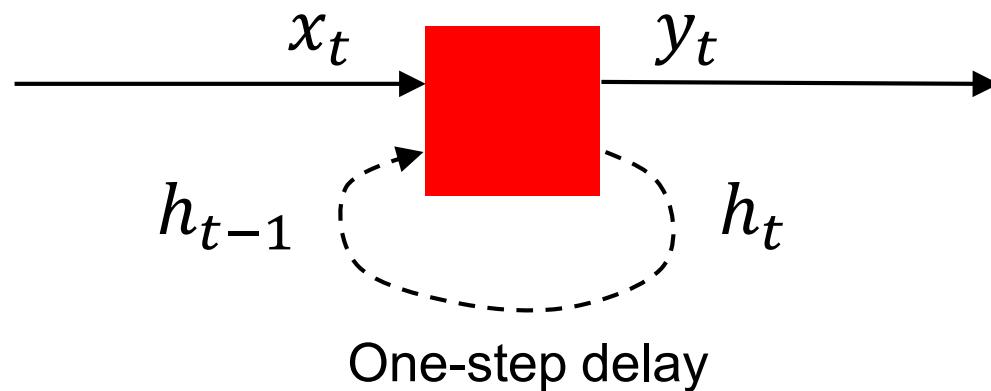
Forward computation only requires activations from the last state. Backprop requires gradients from later states.

Try Again

Continue

# Oops!

Question: Can you run forward/backward inference on an unrolled RNN, i.e. keeping only one copy of its state?



Run Forward Y/N, Run Backward Y/N

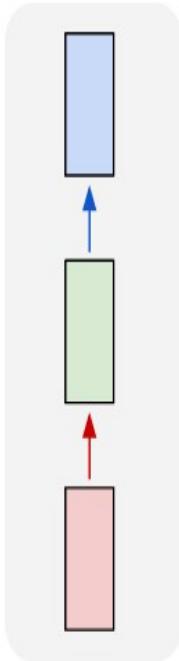
- D. Yes, Yes

Try Again

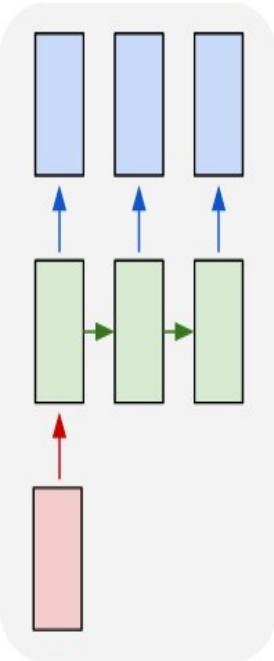
Continue

# Recurrent Networks offer a lot of flexibility:

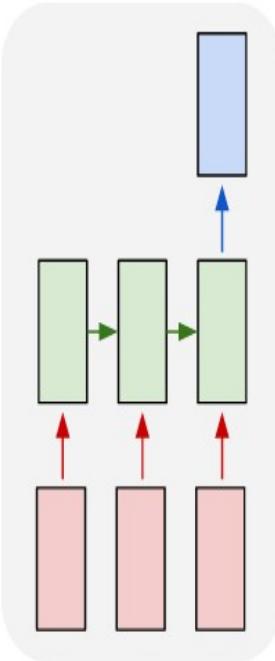
one to one



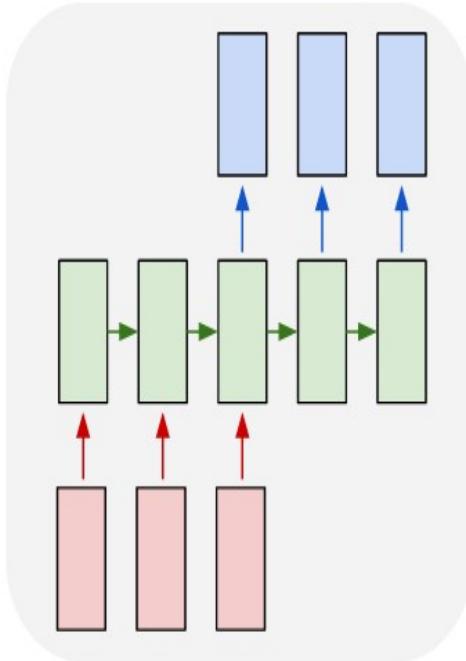
one to many



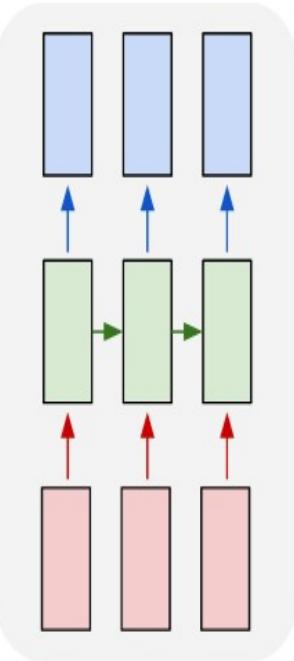
many to one



many to many



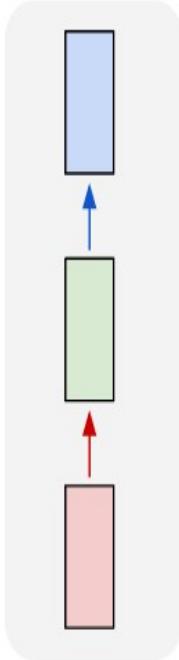
many to many



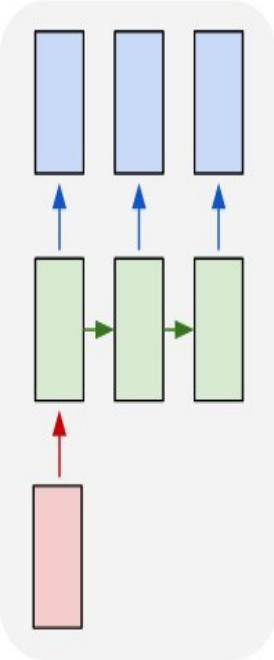
**Vanilla Neural Networks**

# Recurrent Networks offer a lot of flexibility:

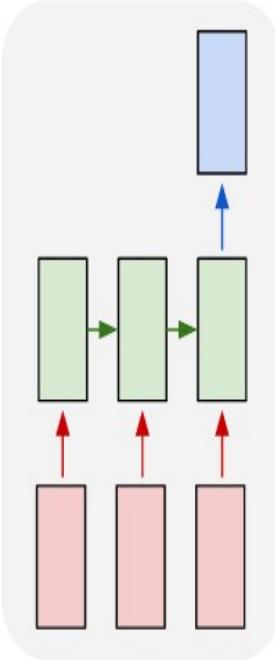
one to one



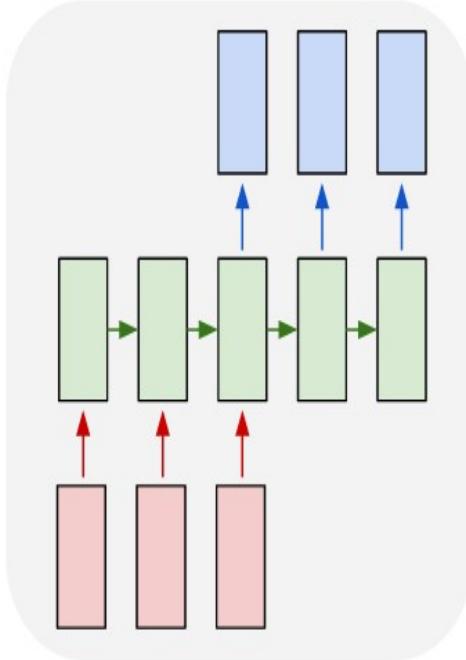
one to many



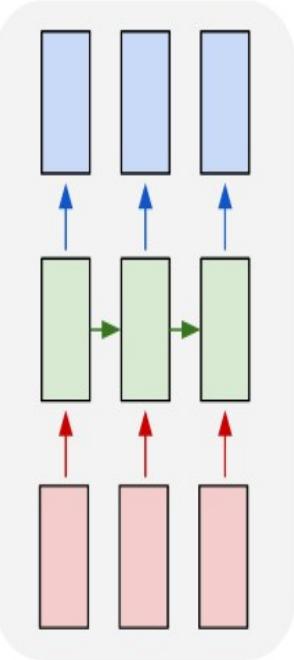
many to one



many to many



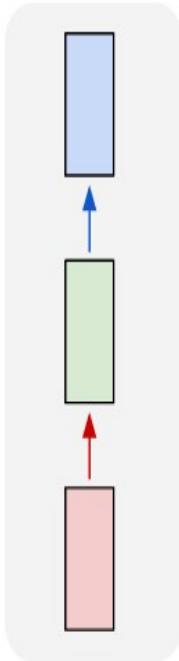
many to many



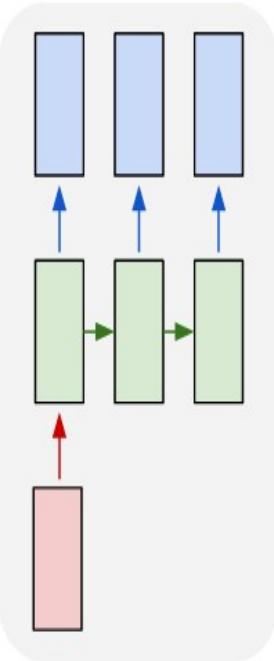
e.g. **Image Captioning**  
image -> sequence of words

# Recurrent Networks offer a lot of flexibility:

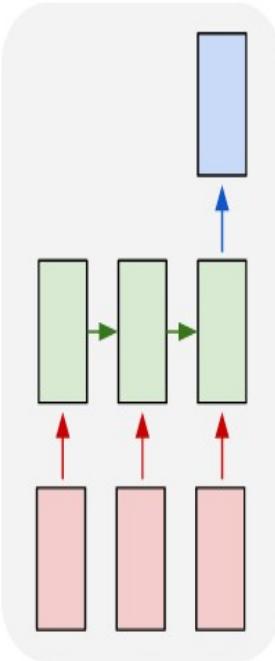
one to one



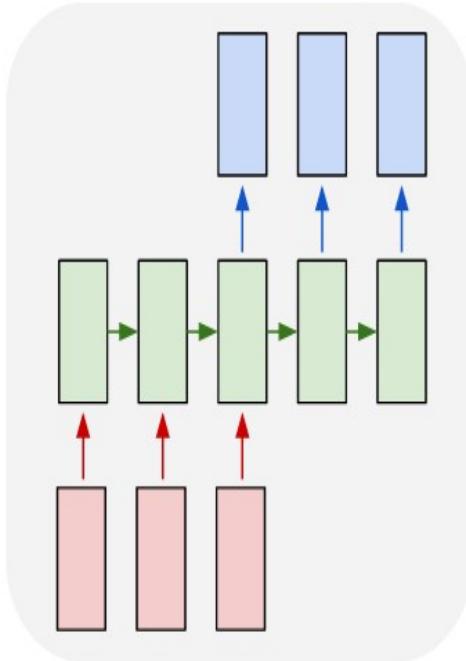
one to many



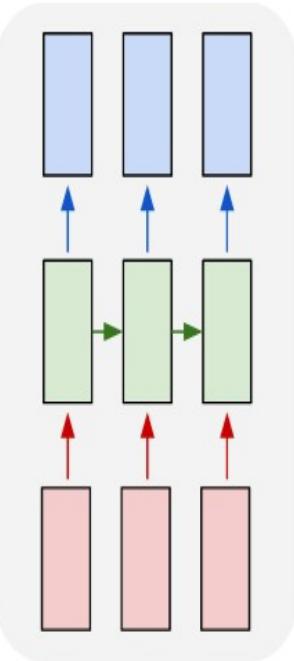
many to one



many to many



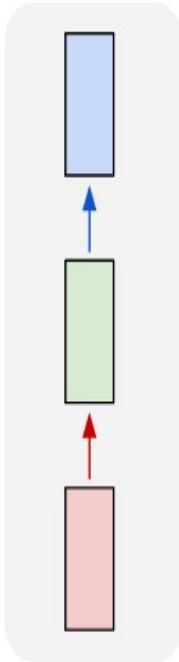
many to many



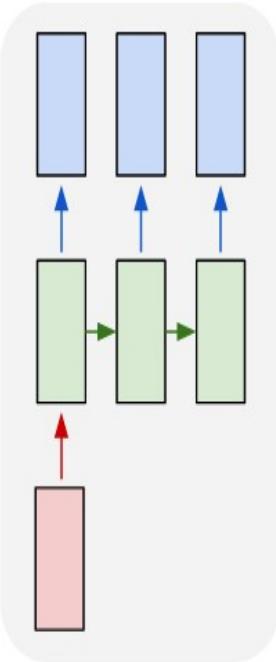
e.g. **Sentiment Classification**  
sequence of words -> sentiment

# Recurrent Networks offer a lot of flexibility:

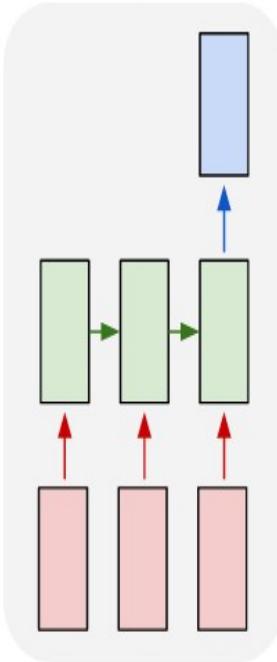
one to one



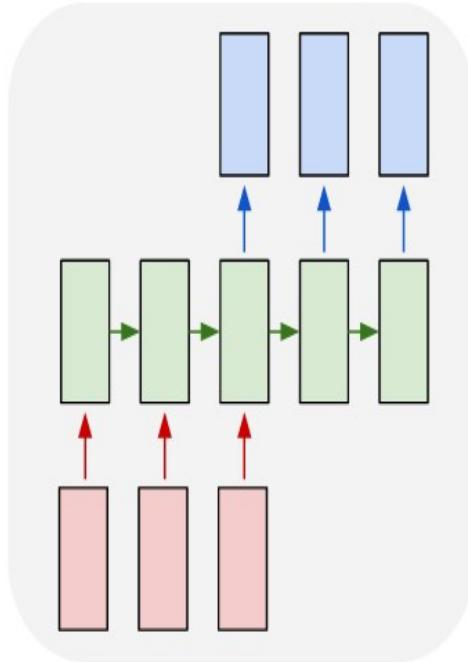
one to many



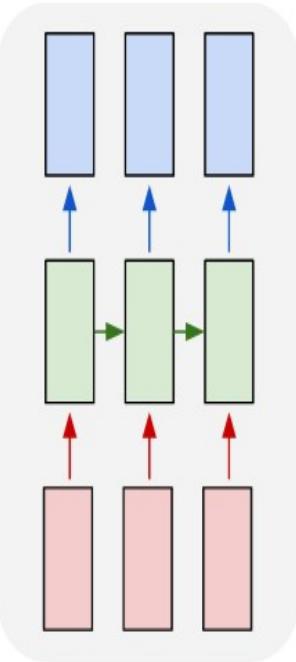
many to one



many to many



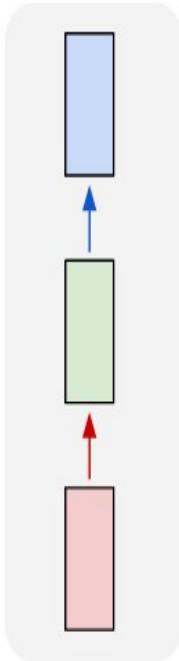
many to many



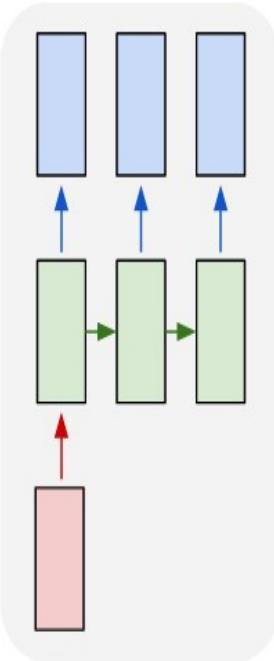
e.g. **Machine Translation**  
seq of words -> seq of words

# Recurrent Networks offer a lot of flexibility:

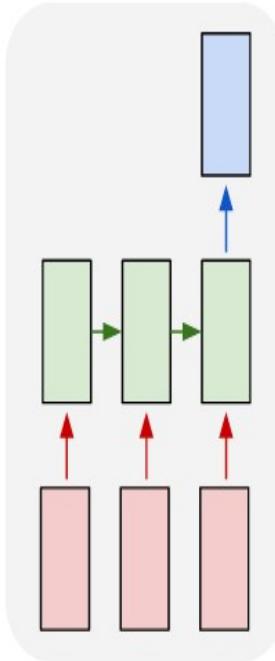
one to one



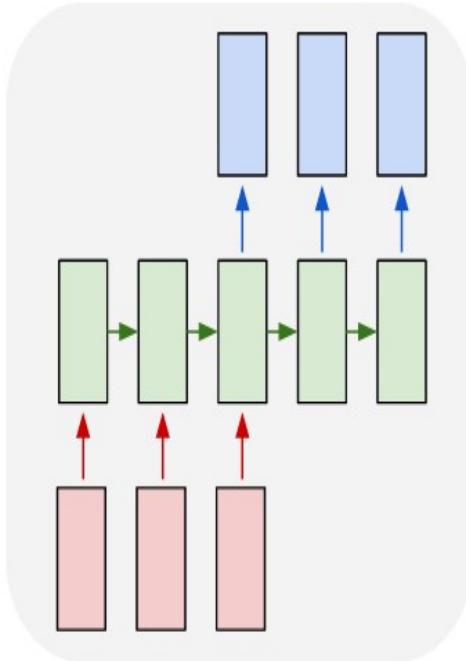
one to many



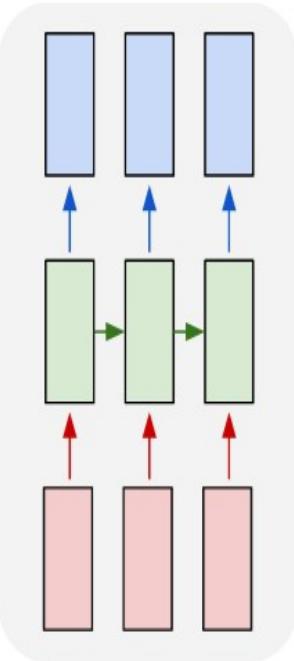
many to one



many to many



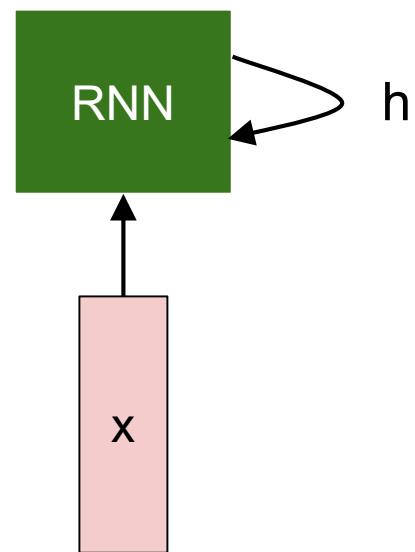
many to many



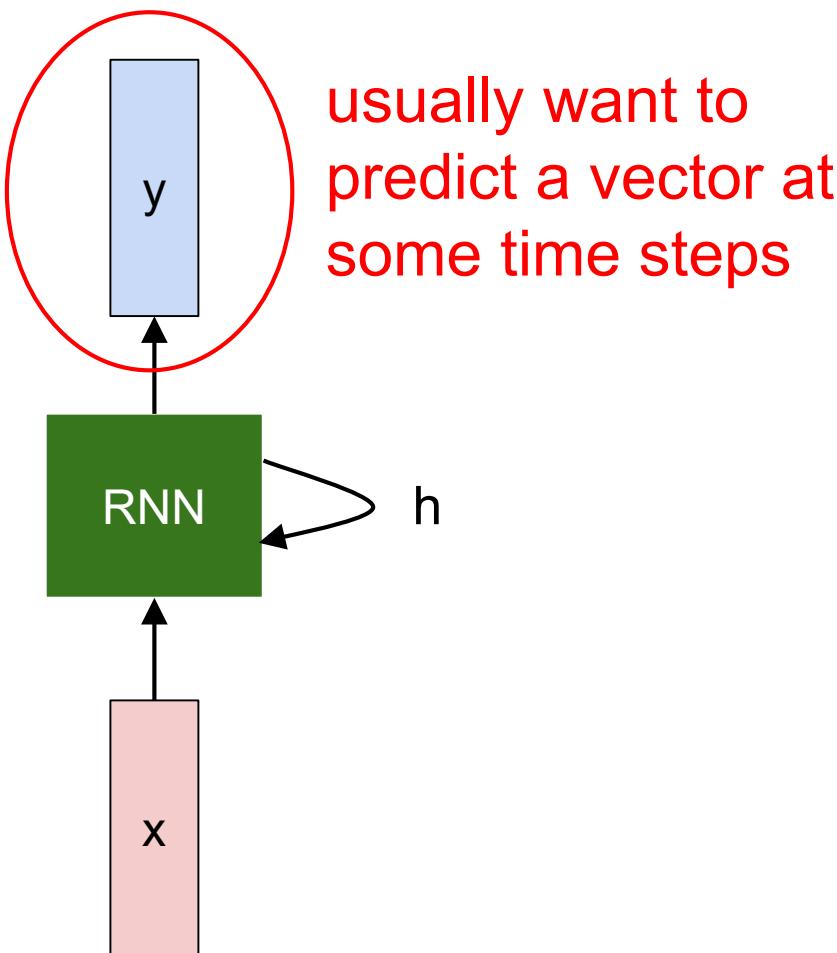
e.g. **Video classification on frame level**



# Recurrent Neural Network



# Recurrent Neural Network

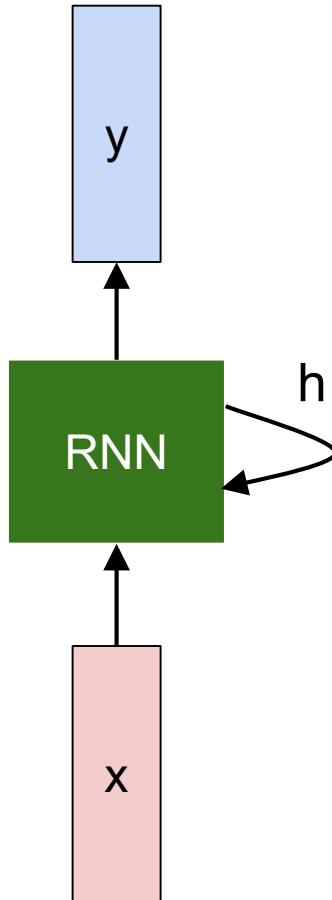


# Recurrent Neural Network

We can process a sequence of vectors  $\mathbf{x}$  by applying a recurrence formula at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

new state      old state      input vector at  
some function      some time step  
with parameters W

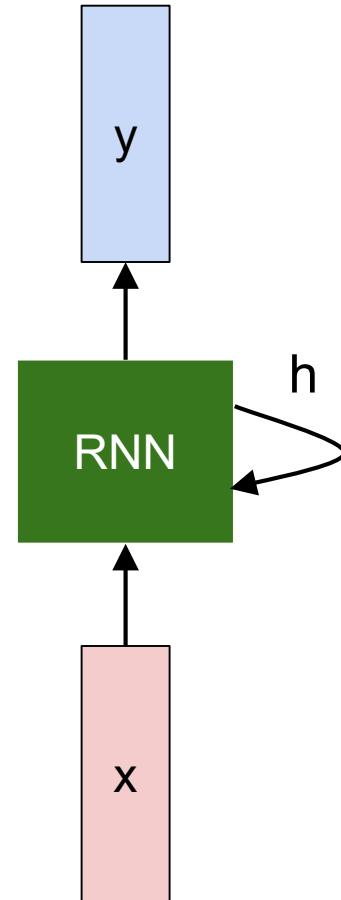


# Recurrent Neural Network

We can process a sequence of vectors  $\mathbf{x}$  by applying a recurrence formula at every time step:

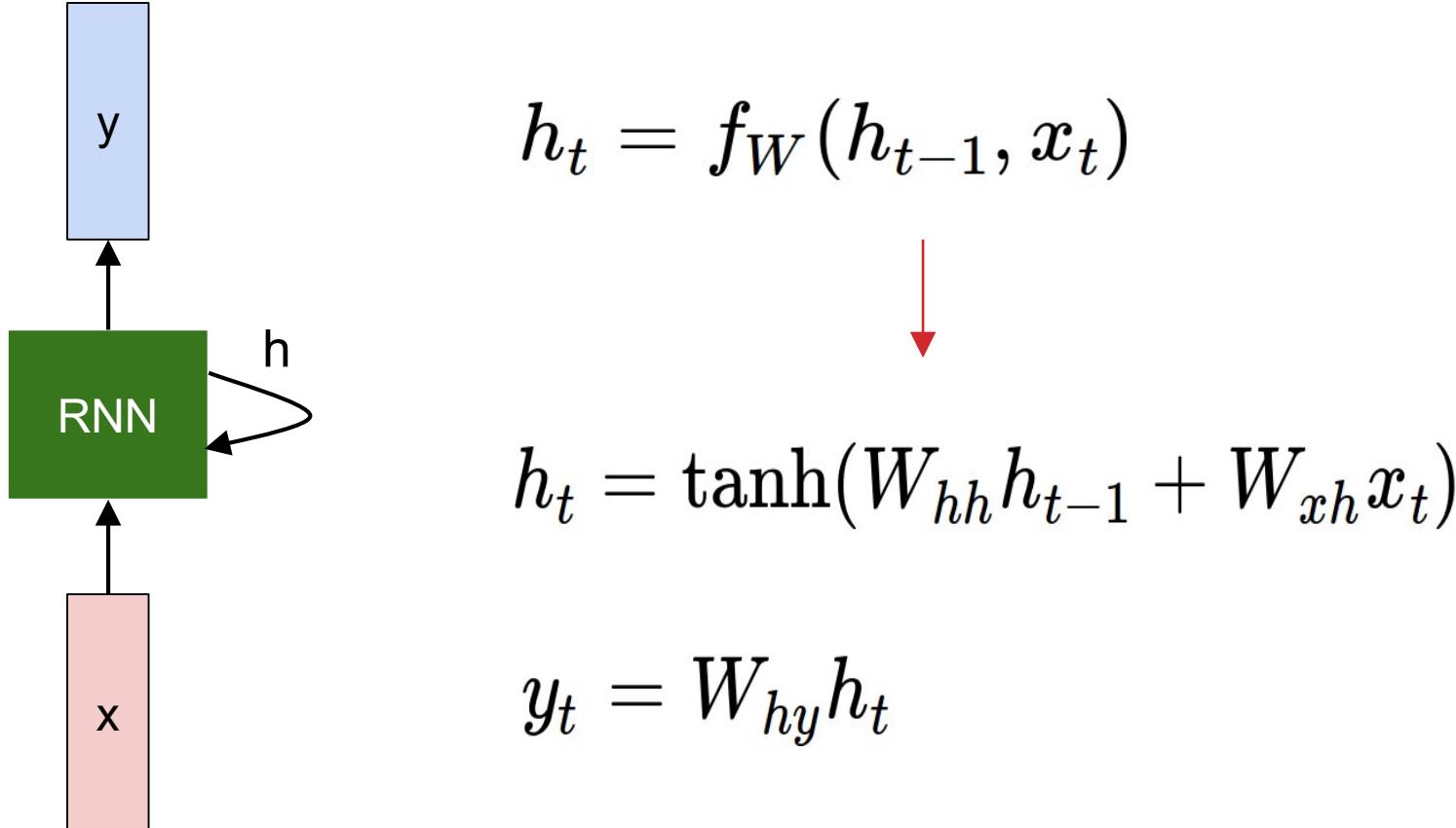
$$h_t = f_W(h_{t-1}, x_t)$$

Notice: the same function and the same set of parameters are used at every time step.



# (Vanilla) Recurrent Neural Network

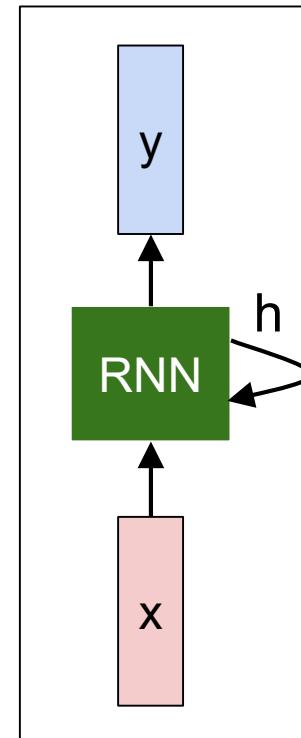
The state consists of a single “*hidden*” vector  $\mathbf{h}$ :



# Character-level language model example

Vocabulary:  
[h,e,l,o]

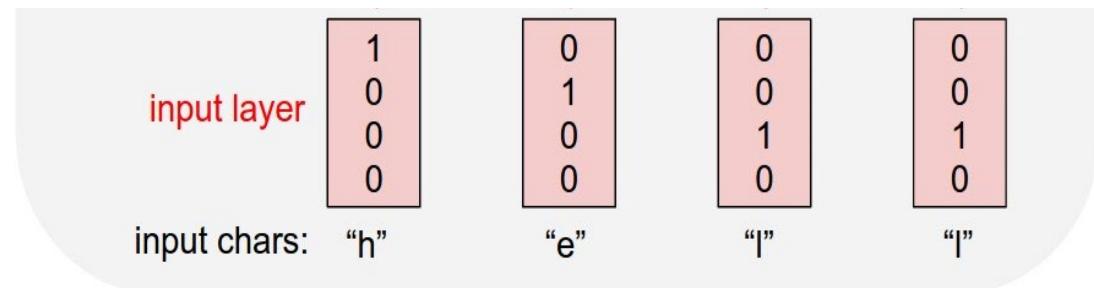
Example training  
sequence:  
**“hello”**



# Character-level language model example

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
**“hello”**

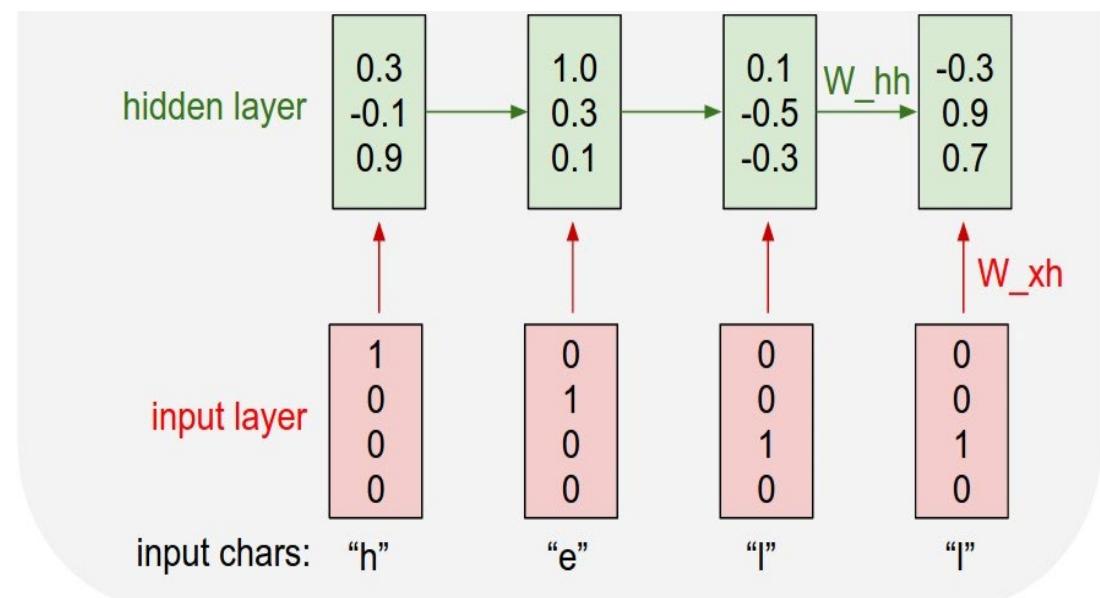


# Character-level language model example

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
**“hello”**

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

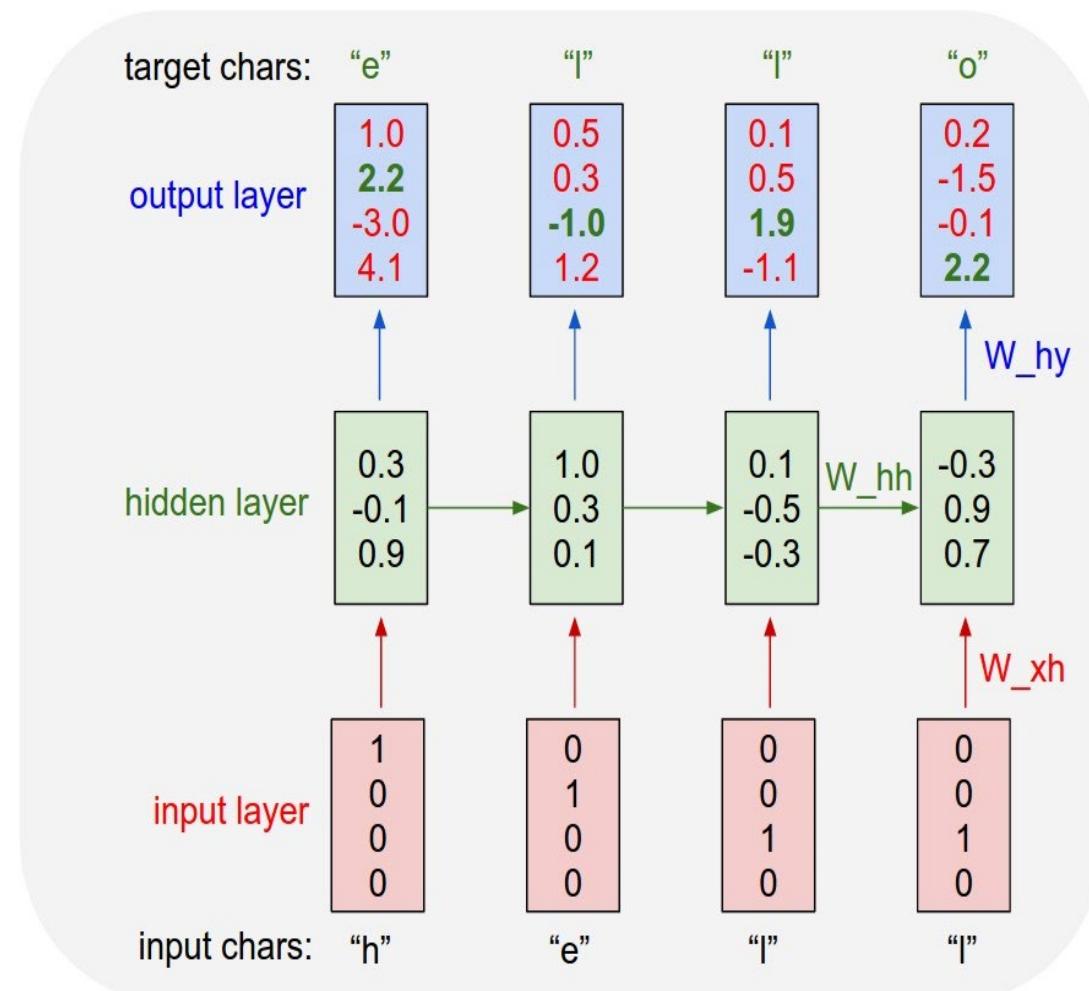


# Character-level language model example

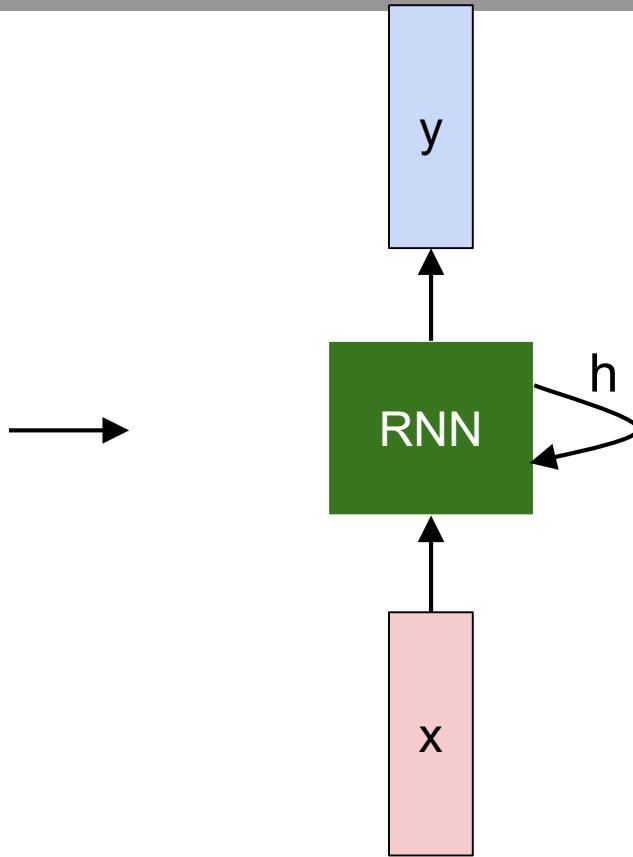
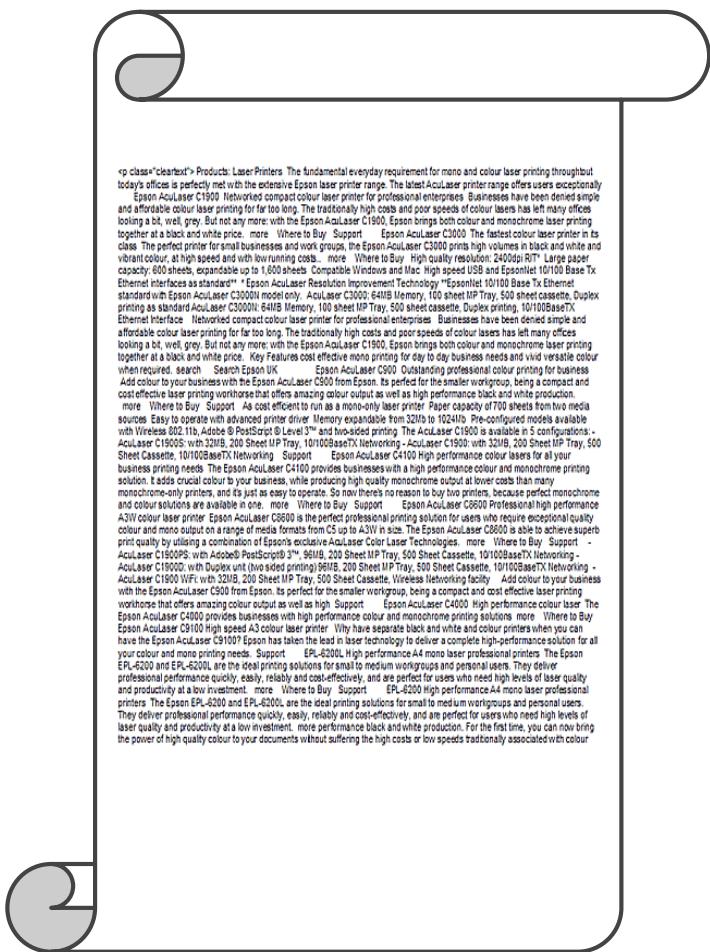
Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
**“hello”**

$$y_t = W_{hy} h_t$$



# Applications - Poetry



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

Based on cs231n by Fei-Fei Li & Andrej Karpathy & Justin Johnson

# Character-Level Recurrent Neural Network

## Sonnet 116 – Let me not ...

*by William Shakespeare*

Let me not to the marriage of true minds  
    Admit impediments. Love is not love  
Which alters when it alteration finds,  
    Or bends with the remover to remove:  
O no! it is an ever-fixed mark  
    That looks on tempests and is never shaken;  
It is the star to every wandering bark,  
    Whose worth's unknown, although his height be taken.  
Love's not Time's fool, though rosy lips and cheeks  
    Within his bending sickle's compass come:  
Love alters not with his brief hours and weeks,  
    But bears it out even to the edge of doom.  
If this be error and upon me proved,  
    I never writ, nor no man ever loved.

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

Based on cs231n by Fei-Fei Li & Andrej Karpathy & Justin Johnson

At first:

tyntd-iafhatawiaoihrdemot lytdws e ,ftti, astai f ogoh eoase rrranbyne 'nhthnee e  
plia tkldg 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 shuwyl 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.

# And later:

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.

Clown:

Come, sir, I will make did behold your worship.

VIOLA:

I'll drink it.

VIOLA:

Why, Salisbury must find his flesh and thought  
That which I am not aps, not a man and in fire,  
To show the reining of the raven and the wars  
To grace my hand reproach within, and not a fair are hand,  
That Caesar and my goodly father's world;  
When I was heaven of presence and our fleets,  
We spare with hours, but cut thy council I am great,  
Murdered and by thy master's ready there  
My power to give thee but so much as hell:  
Some service in the noble bondman here,  
Would show him to her wine.

KING LEAR:

O, if you were a feeble sight, the courtesy of your law,  
Your sight and several breath, will wear the gods  
With his heads, and my hands are wonder'd at the deeds,  
So drop upon your lordship's head, and your opinion  
Shall be against your honour.

# Open source textbook on algebraic geometry

The Screenshot shows the homepage of The Stacks Project. At the top, there is a navigation bar with links: home, about, tags explained, tag lookup, browse, search, bibliography, recent comments, blog, and add slogans. Below the navigation bar, there is a section titled "Browse chapters". This section contains a table with two columns: "Part" and "Chapter". The "Part" column lists categories like Preliminaries, Topics in Scheme Theory, etc., and the "Chapter" column lists specific chapters such as Introduction, Conventions, Set Theory, Categories, Topology, Sheaves on Spaces, Sites and Sheaves, Stacks, Fields, and Commutative Algebra. Each chapter entry includes three links: online, TeX source, and view pdf. To the right of the main content area, there is a sidebar with a "Parts" section containing a numbered list of topics from 1 to 8, and a "Statistics" section providing project metrics.

Part	Chapter	online	TeX source	view pdf
Preliminaries	1. Introduction	<a href="#">online</a>	<a href="#">TeX</a>	<a href="#">pdf</a>
	2. Conventions	<a href="#">online</a>	<a href="#">TeX</a>	<a href="#">pdf</a>
	3. Set Theory	<a href="#">online</a>	<a href="#">TeX</a>	<a href="#">pdf</a>
	4. Categories	<a href="#">online</a>	<a href="#">TeX</a>	<a href="#">pdf</a>
	5. Topology	<a href="#">online</a>	<a href="#">TeX</a>	<a href="#">pdf</a>
	6. Sheaves on Spaces	<a href="#">online</a>	<a href="#">TeX</a>	<a href="#">pdf</a>
	7. Sites and Sheaves	<a href="#">online</a>	<a href="#">TeX</a>	<a href="#">pdf</a>
	8. Stacks	<a href="#">online</a>	<a href="#">TeX</a>	<a href="#">pdf</a>
	9. Fields	<a href="#">online</a>	<a href="#">TeX</a>	<a href="#">pdf</a>
	10. Commutative Algebra	<a href="#">online</a>	<a href="#">TeX</a>	<a href="#">pdf</a>

Parts

1. [Preliminaries](#)
2. [Schemes](#)
3. [Topics in Scheme Theory](#)
4. [Algebraic Spaces](#)
5. [Topics in Geometry](#)
6. [Deformation Theory](#)
7. [Algebraic Stacks](#)
8. [Miscellany](#)

Statistics

The Stacks project now consists of

- o 455910 lines of code
- o 14221 tags (56 inactive tags)
- o 2366 sections

Latex source

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

Based on cs231n by Fei-Fei Li & Andrej Karpathy & Justin Johnson

# Generated math

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

For  $\bigoplus_{n=1,\dots,m}$  where  $\mathcal{L}_{m,\bullet} = 0$ , hence we can find a closed subset  $\mathcal{H}$  in  $\mathcal{H}$  and any sets  $\mathcal{F}$  on  $X$ ,  $U$  is a closed immersion of  $S$ , then  $U \rightarrow T$  is a separated algebraic space.

*Proof.* Proof of (1). It also start we get

$$S = \text{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by  $\coprod Z \times_U U \rightarrow V$ . Consider the maps  $M$  along the set of points  $\text{Sch}_{fppf}$  and  $U \rightarrow U$  is the fibre category of  $S$  in  $U$  in Section, ?? and the fact that any  $U$  affine, see Morphisms, Lemma ???. Hence we obtain a scheme  $S$  and any open subset  $W \subset U$  in  $\text{Sh}(G)$  such that  $\text{Spec}(R') \rightarrow S$  is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that  $f_i$  is of finite presentation over  $S$ . We claim that  $\mathcal{O}_{X,x}$  is a scheme where  $x, x', s'' \in S'$  such that  $\mathcal{O}_{X,x'} \rightarrow \mathcal{O}'_{X',x'}$  is separated. By Algebra, Lemma ?? we can define a map of complexes  $\text{GL}_{S'}(x'/S'')$  and we win.  $\square$

To prove study we see that  $\mathcal{F}|_U$  is a covering of  $\mathcal{X}'$ , and  $\mathcal{T}_i$  is an object of  $\mathcal{F}_{X/S}$  for  $i > 0$  and  $\mathcal{F}_p$  exists and let  $\mathcal{F}_i$  be a presheaf of  $\mathcal{O}_X$ -modules on  $\mathcal{C}$  as a  $\mathcal{F}$ -module. In particular  $\mathcal{F} = U/\mathcal{F}$  we have to show that

$$\widetilde{\mathcal{M}}^\bullet = \mathcal{I}^\bullet \otimes_{\text{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1} \mathcal{F}$$

is a unique morphism of algebraic stacks. Note that

$$\text{Arrows} = (\text{Sch}/S)^{opp}_{fppf}, (\text{Sch}/S)_{fppf}$$

and

$$V = \Gamma(S, \mathcal{O}) \longmapsto (U, \text{Spec}(A))$$

is an open subset of  $X$ . Thus  $U$  is affine. This is a continuous map of  $X$  is the inverse, the groupoid scheme  $S$ .

*Proof.* See discussion of sheaves of sets.  $\square$

The result for prove any open covering follows from the less of Example ???. It may replace  $S$  by  $X_{\text{spaces},\text{étale}}$  which gives an open subspace of  $X$  and  $T$  equal to  $S_{\text{Zar}}$ , see Descent, Lemma ???. Namely, by Lemma ?? we see that  $R$  is geometrically regular over  $S$ .

**Lemma 0.1.** Assume (3) and (3) by the construction in the description.

Suppose  $X = \lim |X|$  (by the formal open covering  $X$  and a single map  $\underline{\text{Proj}}_X(\mathcal{A}) = \text{Spec}(B)$  over  $U$  compatible with the complex

$$\text{Set}(\mathcal{A}) = \Gamma(X, \mathcal{O}_{X,\mathcal{O}_X}).$$

When in this case of to show that  $\mathcal{Q} \rightarrow \mathcal{C}_{Z/X}$  is stable under the following result in the second conditions of (1), and (3). This finishes the proof. By Definition ?? (without element is when the closed subschemes are catenary. If  $T$  is surjective we may assume that  $T$  is connected with residue fields of  $S$ . Moreover there exists a closed subspace  $Z \subset X$  of  $X$  where  $U$  in  $X'$  is proper (some defining as a closed subset of the uniqueness it suffices to check the fact that the following theorem

(1)  $f$  is locally of finite type. Since  $S = \text{Spec}(R)$  and  $Y = \text{Spec}(R)$ .

*Proof.* This is form all sheaves of sheaves on  $X$ . But given a scheme  $U$  and a surjective étale morphism  $U \rightarrow X$ . Let  $U \cap U = \coprod_{i=1,\dots,n} U_i$  be the scheme  $X$  over  $S$  at the schemes  $X_i \rightarrow X$  and  $U = \lim_i X_i$ .  $\square$

The following lemma surjective restrocomposes of this implies that  $\mathcal{F}_{x_0} = \mathcal{F}_{x_0} = \mathcal{F}_{x,\dots,0}$ .

**Lemma 0.2.** Let  $X$  be a locally Noetherian scheme over  $S$ ,  $E = \mathcal{F}_{X/S}$ . Set  $\mathcal{I} = \mathcal{J}_1 \subset \mathcal{I}'_n$ . Since  $\mathcal{I}^n \subset \mathcal{I}^n$  are nonzero over  $i_0 \leq p$  is a subset of  $\mathcal{J}_{n,0} \circ \overline{A}_2$  works.

**Lemma 0.3.** In Situation ???. Hence we may assume  $q' = 0$ .

*Proof.* We will use the property we see that  $p$  is the next functor (??). On the other hand, by Lemma ?? we see that

$$D(\mathcal{O}_{X'}) = \mathcal{O}_X(D)$$

where  $K$  is an  $F$ -algebra where  $\delta_{n+1}$  is a scheme over  $S$ .  $\square$

*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  $X_{\text{étale}}$  we have

$$\mathcal{O}_X(\mathcal{F}) = \{\text{morph}_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

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 ??.

**Lemma 0.3.** Let  $S$  be a scheme. Let  $X$  be a scheme and  $X$  is an affine open covering. Let  $\mathcal{U} \subset \mathcal{X}$  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

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

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.

*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  $\mathcal{C}$ . The functor  $\mathcal{F}$  is a “field”

$$\mathcal{O}_{X,x} \longrightarrow \mathcal{F}_{\bar{x}} \dashrightarrow (\mathcal{O}_{X_{\text{étale}}}) \longrightarrow \mathcal{O}_{X'_\ell}^{-1} \mathcal{O}_{X_\lambda}(\mathcal{O}_{X_n}^v)$$

is an isomorphism of covering of  $\mathcal{O}_{X_i}$ . 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  $\mathcal{O}_{X_\lambda}$  is a closed immersion, see Lemma ???. This is a sequence of  $\mathcal{F}$  is a similar morphism.

# Train on Linux code

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

The screenshot shows the GitHub repository page for `torvalds/linux`. The top navigation bar includes links for 'Explore', 'Gist', 'Blog', and 'Help'. On the right, there are user profile icons for `karpathy` and various repository management buttons like '+', 'Watch', 'Star', 'Fork', and 'Code'.

The repository name `torvalds / linux` is displayed prominently. Below it, metrics are shown: 520,037 commits, 1 branch, 420 releases, and 5,039 contributors. A dropdown menu indicates the current branch is `master`.

The main content area displays a list of recent commits, each with a author, commit message, date, and a link to the commit details. Some commits are from specific branches like `Documentation`, `arch`, `block`, etc.

On the right sidebar, there are links for 'Code', 'Pull requests' (74), 'Pulse', and 'Graphs'. Below these are sections for cloning the repository via 'HTTPS clone URL' (with a copy icon) and options to 'Clone in Desktop' or 'Download ZIP'.

Author	Commit Message	Date
<code>torvalds</code>	authored 9 hours ago	latest commit 4b1706927d
<code>Documentation</code>	Merge git://git.kernel.org/pub/scm/linux/kernel/git/nab/target-pending	6 days ago
<code>arch</code>	Merge branch 'x86-urgent-for-linus' of git://git.kernel.org/pub/scm/l...	a day ago
<code>block</code>	block: discard bdi_unregister() in favour of bdi_destroy()	9 days ago
<code>crypto</code>	Merge git://git.kernel.org/pub/scm/linux/kernel/git/herbert/crypto-2.6	10 days ago
<code>drivers</code>	Merge branch 'drm-fixes' of git://people.freedesktop.org/~airlied/linux	9 hours ago
<code>firmware</code>	firmware/ihex2fw.c: restore missing default in switch statement	2 months ago
<code>fs</code>	vfs: read file_handle only once in handle_to_path	4 days ago
<code>include</code>	Merge branch 'perf-urgent-for-linus' of git://git.kernel.org/pub/scm/...	a day ago
<code>init</code>	init: fix regression by supporting devices with major:minor:offset fo...	a month ago
<code>iommu</code>	IOMMU: New interface for virtio IOMMU driver to handle memory references	a month ago

Based on cs231n by Fei-Fei Li & Andrej Karpathy & Justin Johnson

# Generated C code

```
static void do_command(struct seq_file *m, void *v)
{
    int column = 32 << (cmd[2] & 0x80);
    if (state)
        cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
    else
        seq = 1;
    for (i = 0; i < 16; i++) {
        if (k & (1 << 1))
            pipe = (in_use & UMXTHREAD_UNCCA) +
                ((count & 0x00000000fffff8) & 0x0000000f) << 8;
        if (count == 0)
            sub(pid, ppc_md.kexec_handle, 0x20000000);
        pipe_set_bytes(i, 0);
    }
    /* Free our user pages pointer to place camera if all dash */
    subsystem_info = &of_changes[PAGE_SIZE];
    rek_controls(offset, idx, &soffset);
    /* Now we want to deliberately put it to device */
    control_check_polarity(&context, val, 0);
    for (i = 0; i < COUNTER; i++)
        seq_puts(s, "policy ");
}
```

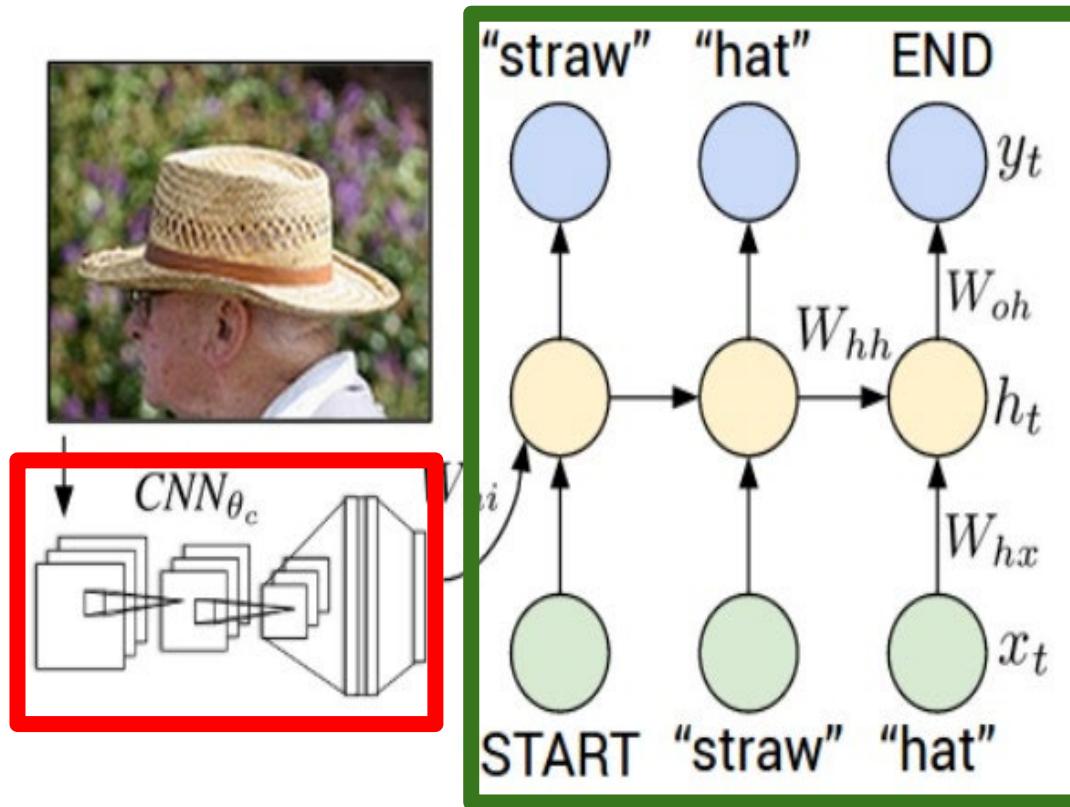
# Generated C code

```
/*
 * 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>
#include <linux/platform_device.h>
#include <linux/multi.h>
#include <linux/ckevent.h>

#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/seteew.h>
#include <asm/pgproto.h>
```

# Image Captioning



## Convolutional Neural Network

test image



image



test image



conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

FC-1000

softmax

image



test image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

FC-1000

softmax

image



test image



conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

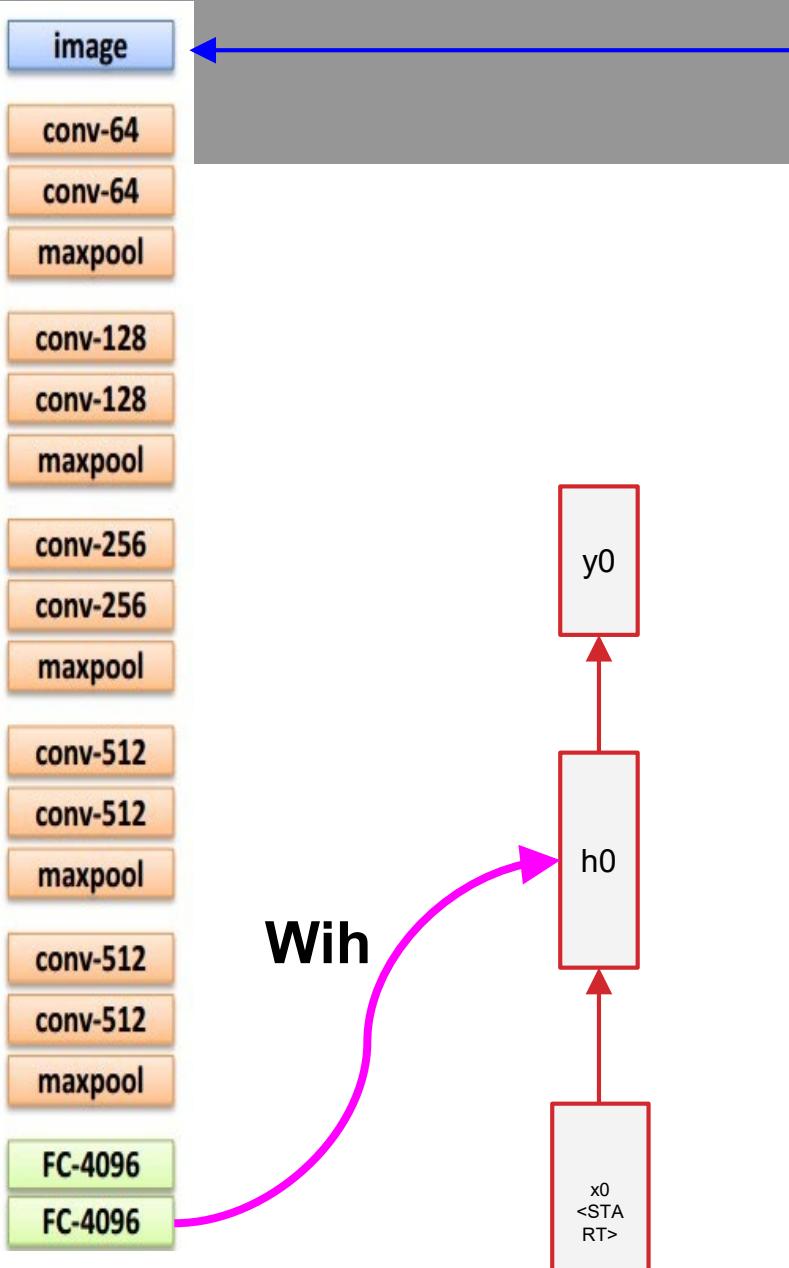
maxpool

FC-4096

FC-4096

x0  
<STA  
RT>

<START>



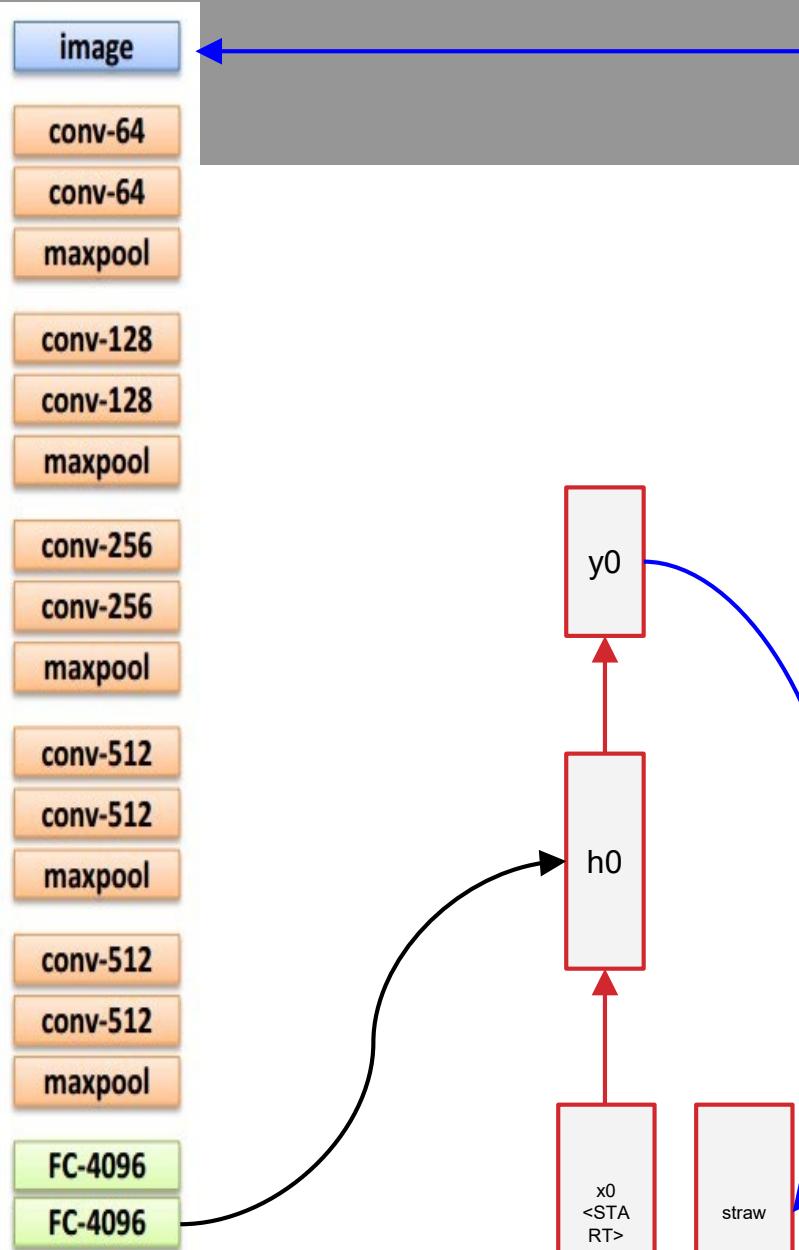
test image

**before:**

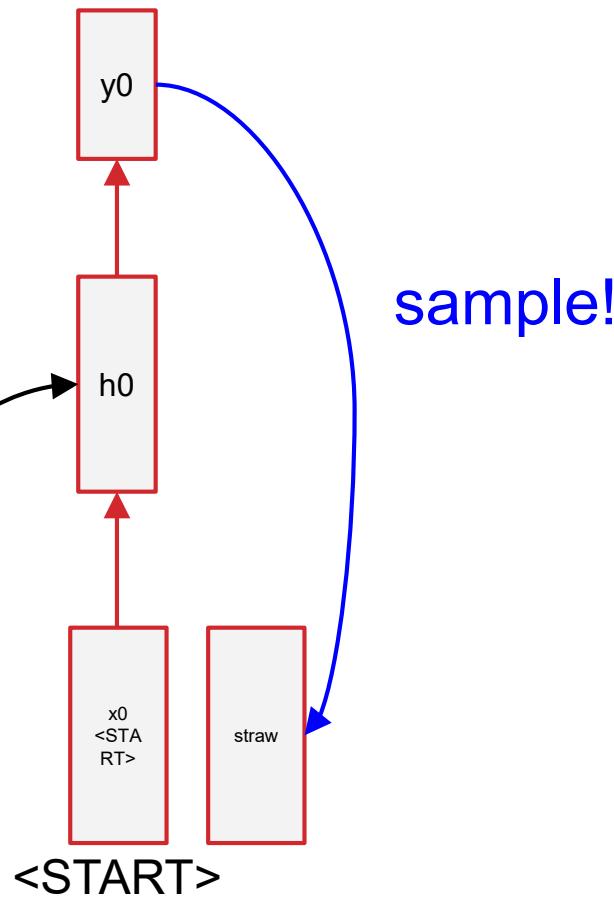
$$h = \tanh(Wxh * x + Whh * h)$$

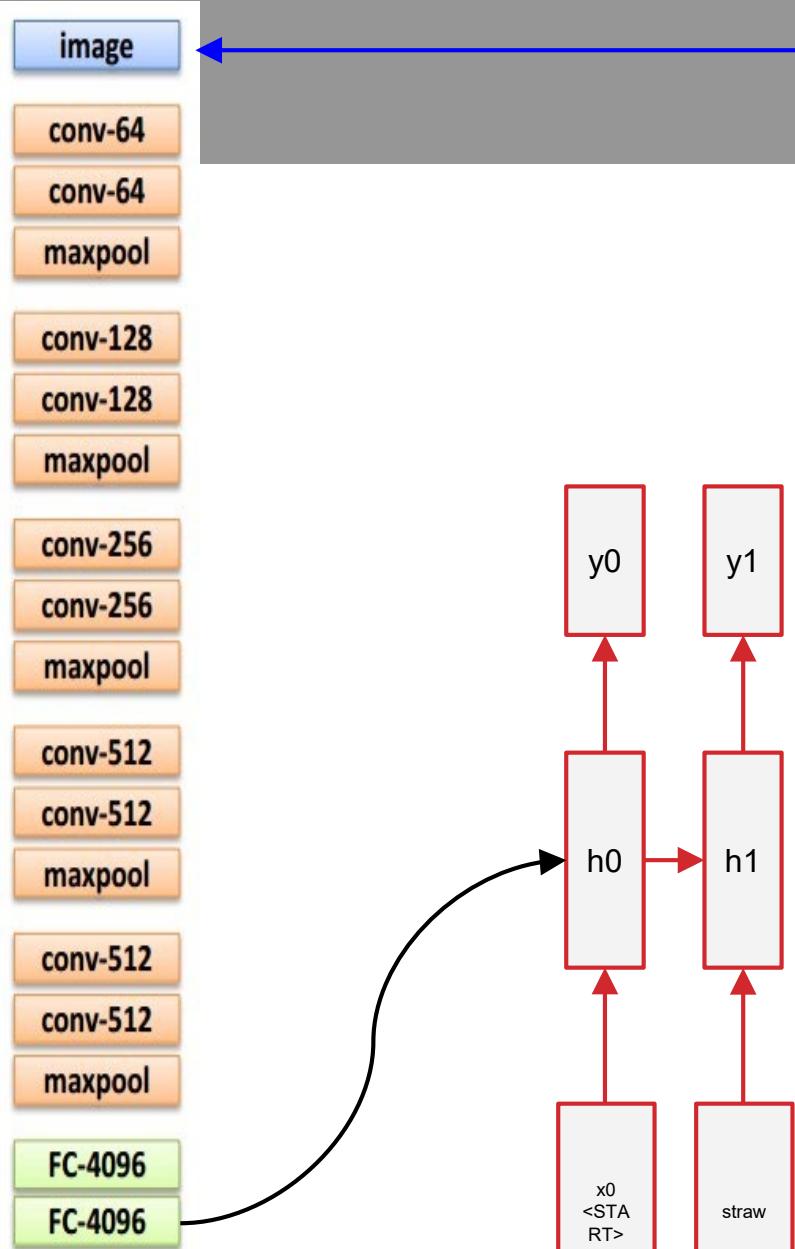
**now:**

$$h = \tanh(Wxh * x + \textcolor{magenta}{Wi}_h * v)$$

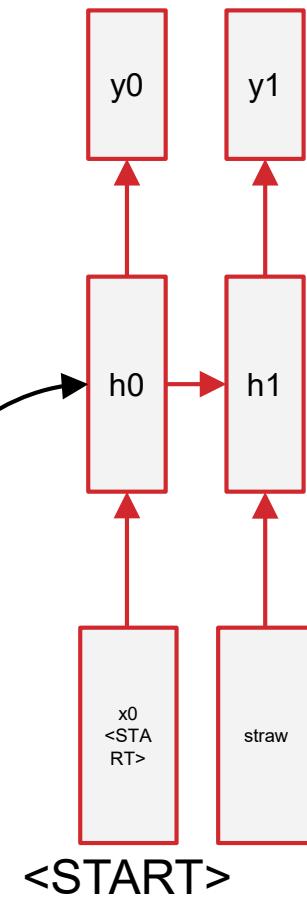


test image



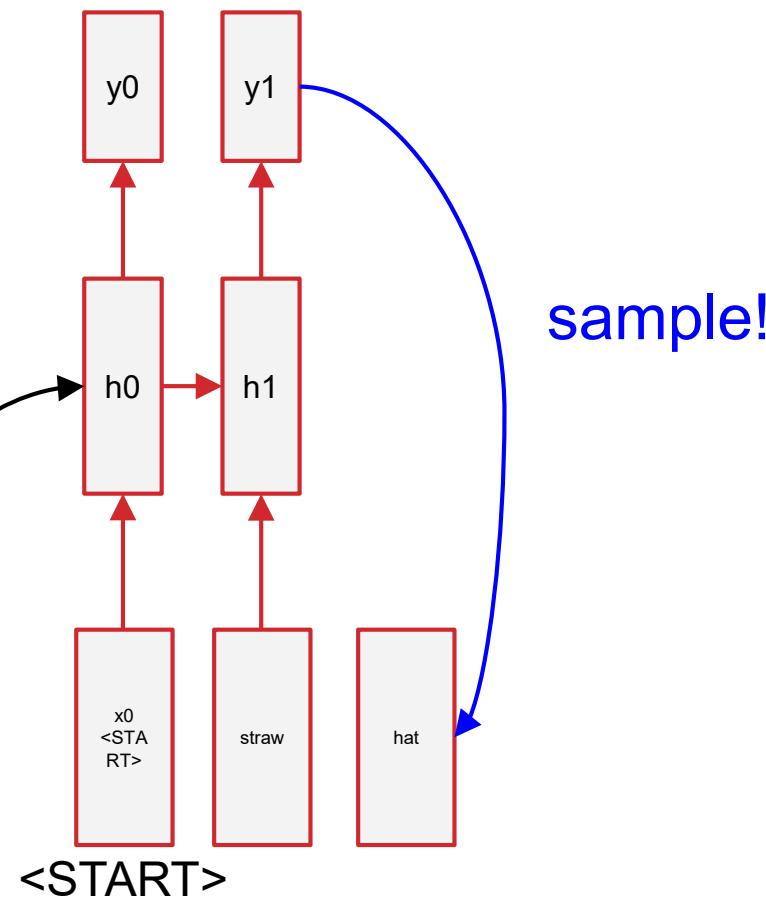


test image



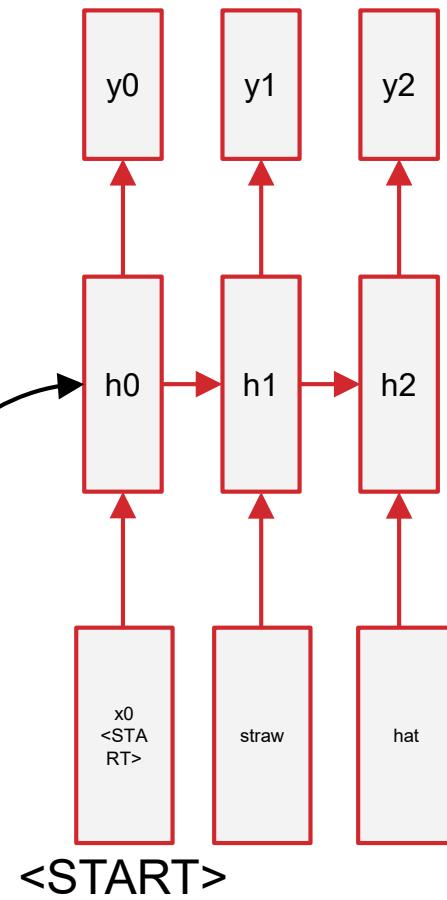


test image



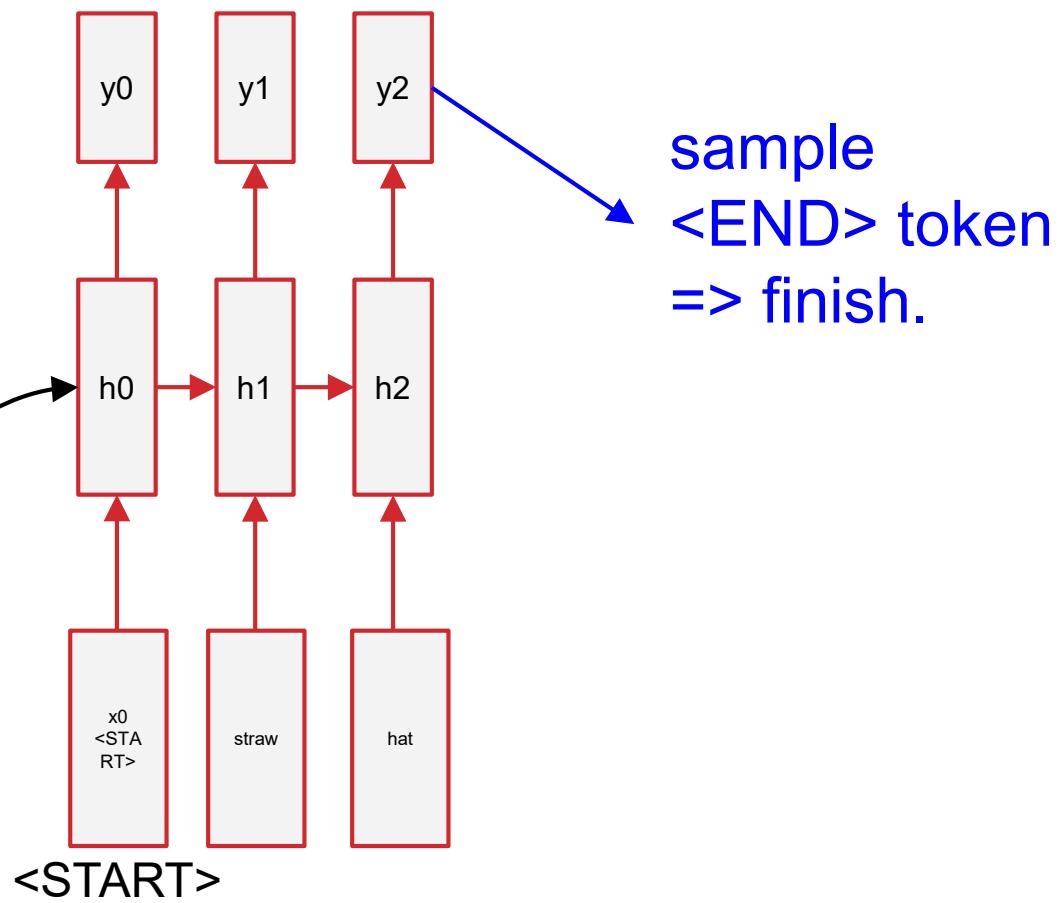


test image





test image



# RNN sequence generation

Greedy (most likely) symbol generation is not very effective.

Typically, the top-k sequences generated so far are remembered and the top-k of their one-symbol continuations are kept for the next step (beam search)

However, k does not have to be very large (7 was used in Karpathy and Fei-Fei 2015).

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

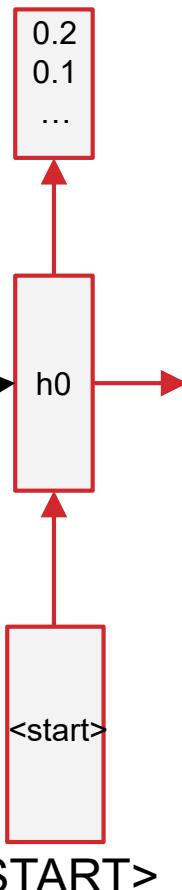
FC-4096

# Beam search with k = 2



Top-k most likely words

straw 0.2  
man 0.1  
...



image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

# Beam search with k = 2



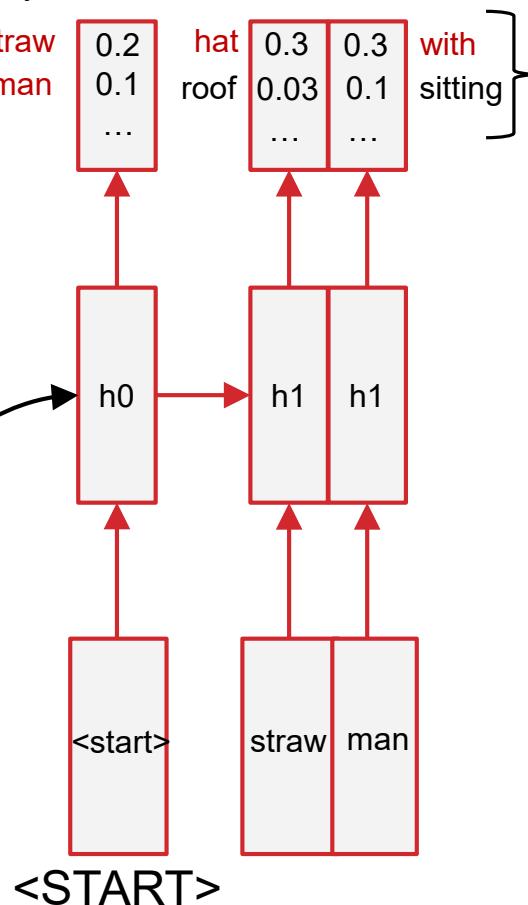
Top-k most likely words

straw	0.2
man	0.1
	...
hat	0.3
roof	0.03
	...
	0.1
	...

Need Top-k most likely words  
Because they might be part of  
the top-K most likely subsequences

Top sequences so far are:

straw hat ( $p = 0.06$ )  
man with ( $p = 0.03$ )



image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

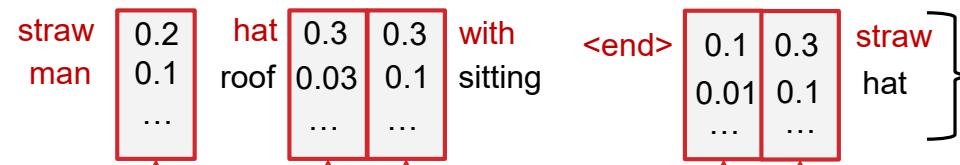
FC-4096

FC-4096

# Beam search with k = 2



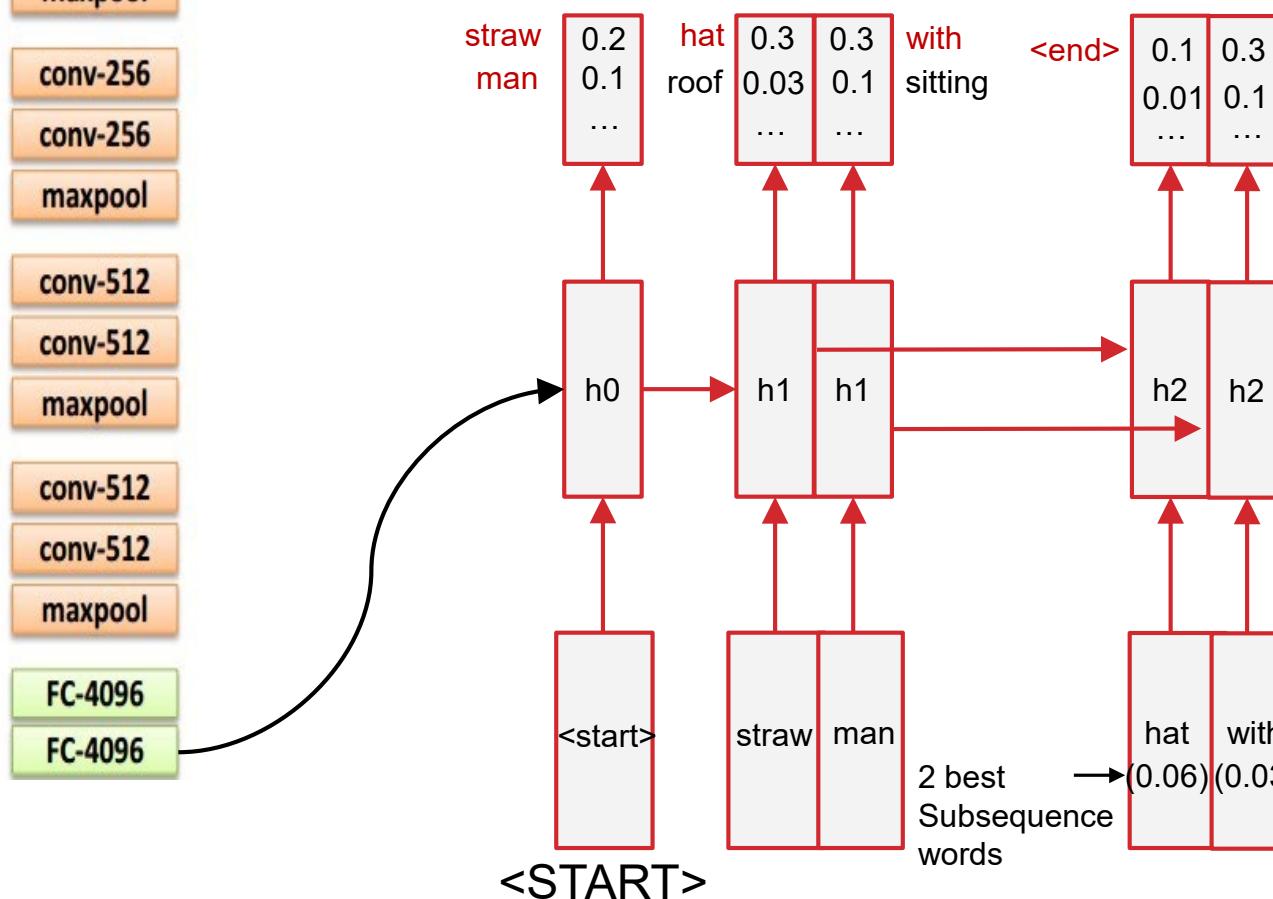
Top-k most likely words



Best 2 words this pos'n

Top sequences so far:

straw hat <end> (0.006)  
man with straw (0.009)



image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

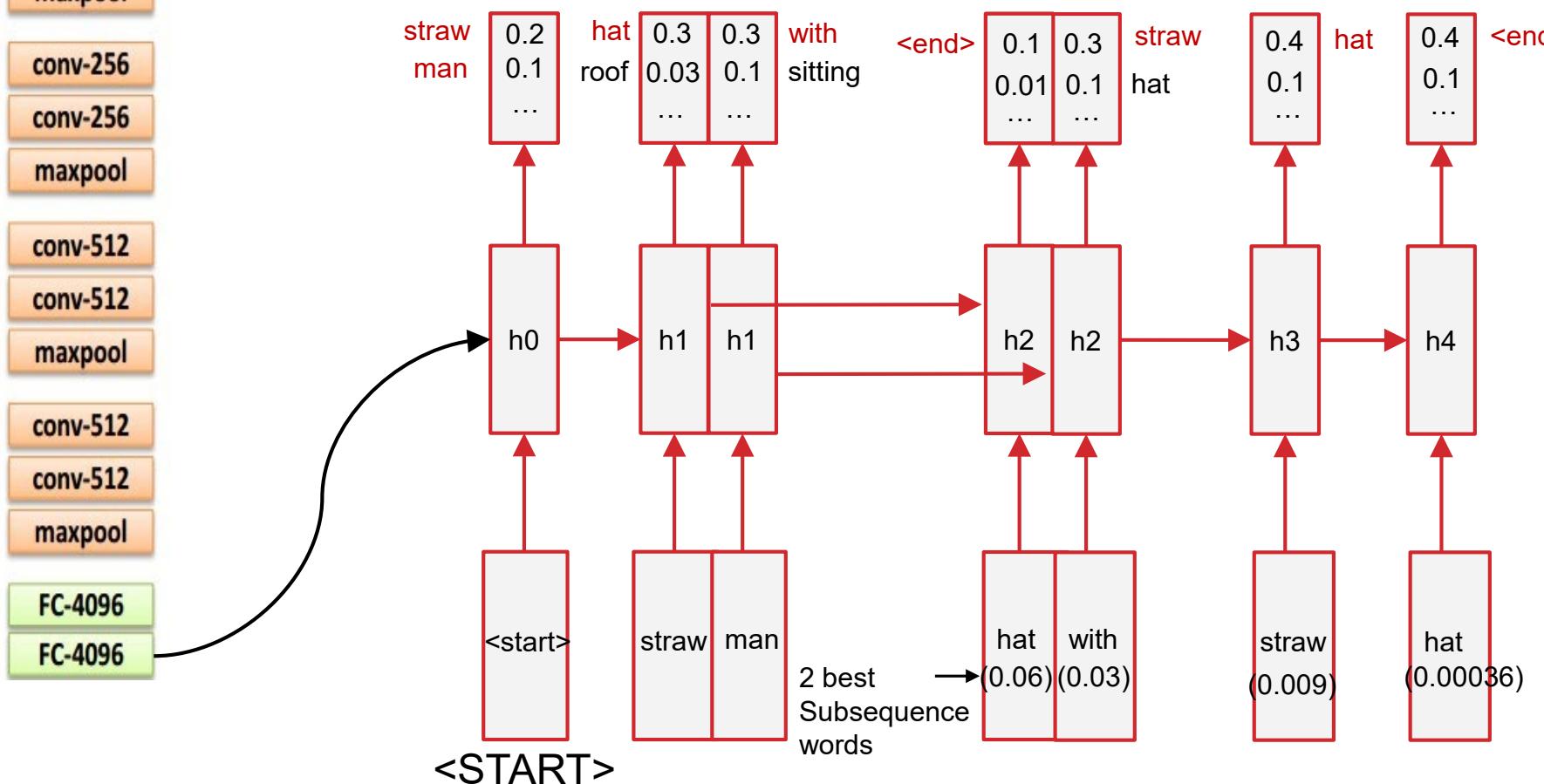
FC-4096

# Beam search with k = 2



Top-k most likely words

straw	0.2	hat	0.3	0.3	with	sitting	<end>	0.1	0.3	straw	0.4	hat	0.4	<end>
man	0.1	roof	0.03	0.1				0.01	0.1	hat	0.1	0.1	0.1	
	...		...	...				...	...		...	...	...	



# Image Sentence Datasets

a man riding a bike on a dirt path through a forest.  
bicyclist raises his fist as he rides on desert dirt trail.  
this dirt bike rider is smiling and raising his fist in triumph.  
a man riding a bicycle while pumping his fist in the air.  
a mountain biker pumps his fist in celebration.



**Microsoft COCO**  
*[Tsung-Yi Lin et al. 2014]*  
[mscoco.org](http://mscoco.org)

currently:  
~120K images  
~5 sentences each



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"a young boy is holding a baseball bat."



"a cat is sitting on a couch with a remote control."

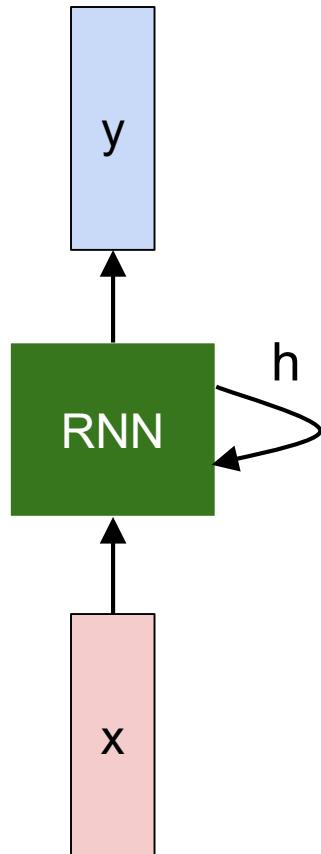


"a woman holding a teddy bear in front of a mirror."



"a horse is standing in the middle of a road."

# Vanilla RNN: Exploding/Vanishing Gradients



$$h_t = \tanh(\underbrace{W_{hh}h_{t-1} + W_{xh}x_t}_{\hat{h}})$$
$$y_t = W_{hy}h_t$$

Using Jacobians:

$$J_L(h_{t-1}) = J_L(h_t) J_{h_t}(h_{t-1})$$

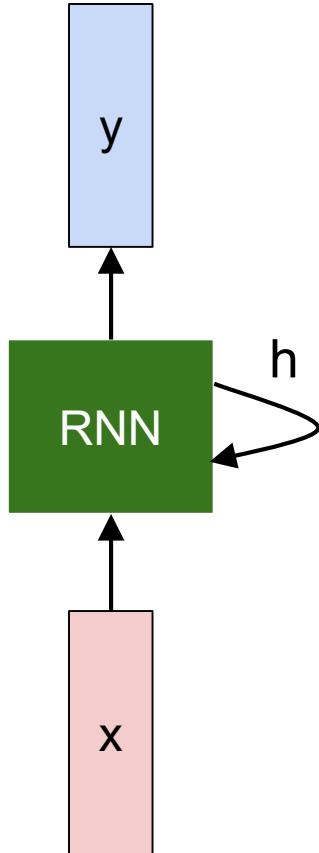
$$J_{h_t}(h_{t-1}) = \frac{1}{\cosh^2 \hat{h}} W_{hh}$$

# Vanilla RNN: Exploding/Vanishing Gradients

Using Jacobians:

$$J_L(h_{t-1}) = J_L(h_t) J_{h_t}(h_{t-1})$$

$$J_{h_t}(h_{t-1}) = \frac{1}{\cosh^2 \hat{h}} W_{hh}$$



Now  $\frac{1}{\cosh^2 \hat{h}}$  is  $\leq 1$ , and  $W_{hh}$  can be arbitrarily large/small.

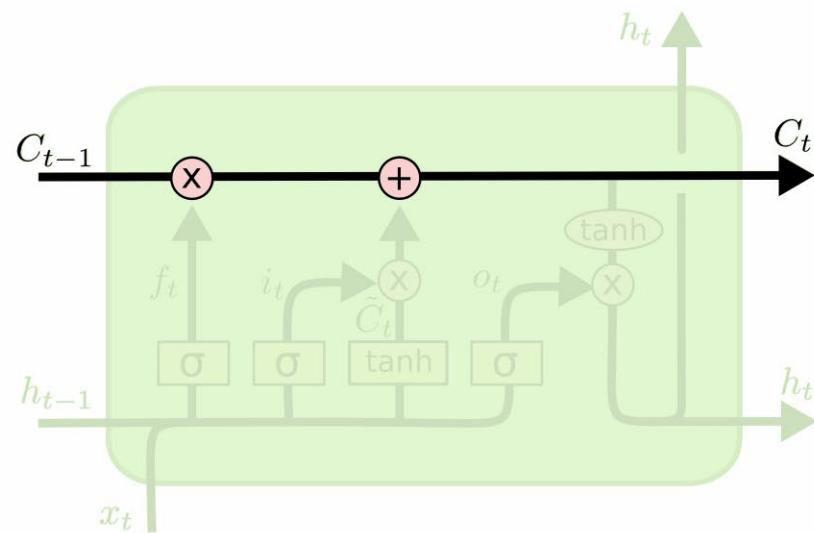
$W_{hh}$  is multiplied by the gradient at the next step, so the gradient across time steps is roughly a power of  $W_{hh}$ .

If the largest eigenvalue of  $W_{hh} > 1$ , the gradients will grow exponentially with time (exploding).

If the largest eigenvalue of  $W_{hh} < 1$ , the gradients will shrink exponentially with time (vanishing).

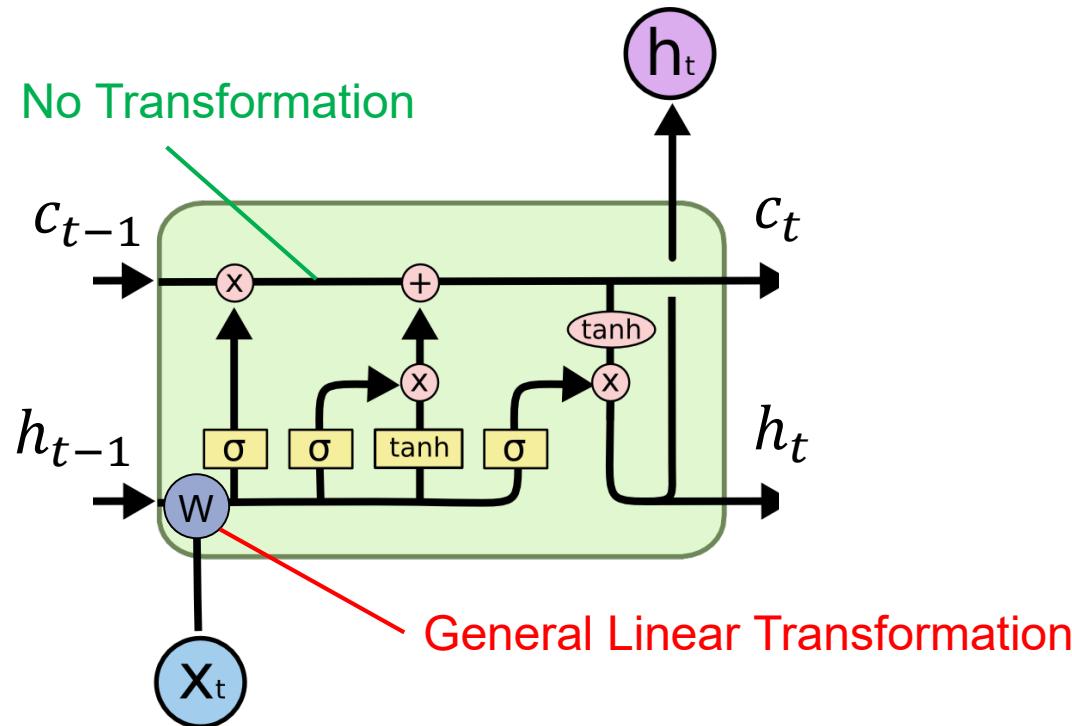
# Better RNN Memory: LSTMs

LSTMs (Long Short-Term Memory) units have a memory cell  $c_i$  which is gated element-wise (no linear transform) from one time step to the next.



# Better RNN Memory: LSTMs

They also have linearly transformed hidden states like a standard RNN.



# LSTM: Long Short-Term Memory

Vanilla RNN:

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

Input from below  
Input from left

$$h \in \mathbb{R}^n \quad W^l [n \times 2n]$$

---

LSTM: (much more widely used)

$$W^l [4n \times 2n]$$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

Input from below  
Input from left

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

# LSTM: Anything you can do...

## Vanilla RNN:

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$h \in \mathbb{R}^n$        $W^l [n \times 2n]$

Input from below  
Input from left

## LSTM: (much more widely used)

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$
$$c_t^l = f \odot c_{t-1}^l + i \odot g$$
$$h_t^l = o \odot \tanh(c_t^l)$$

Input from below  
Input from left

An LSTM can emulate a simple RNN by remembering with  $i$  and  $o$

# LSTM: Anything you can do...

## Vanilla RNN:

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$h \in \mathbb{R}^n$        $W^l [n \times 2n]$

Input from below  
Input from left

## LSTM: (much more widely used)

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$
$$c_t^l = f \odot c_{t-1}^l + i \odot g$$
$$h_t^l = o \odot \tanh(c_t^l)$$

Input from below  
Input from left

An LSTM can emulate a simple RNN by remembering with  $i$  and  $o$

# LSTM: Anything you can do...

## Vanilla RNN:

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$h \in \mathbb{R}^n$        $W^l [n \times 2n]$

Input from below  
Input from left

## LSTM: (much more widely used)

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$
$$c_t^l = f \odot c_{t-1}^l + i \odot g$$
$$h_t^l = o \odot \tanh(c_t^l)$$

Input from below  
Input from left

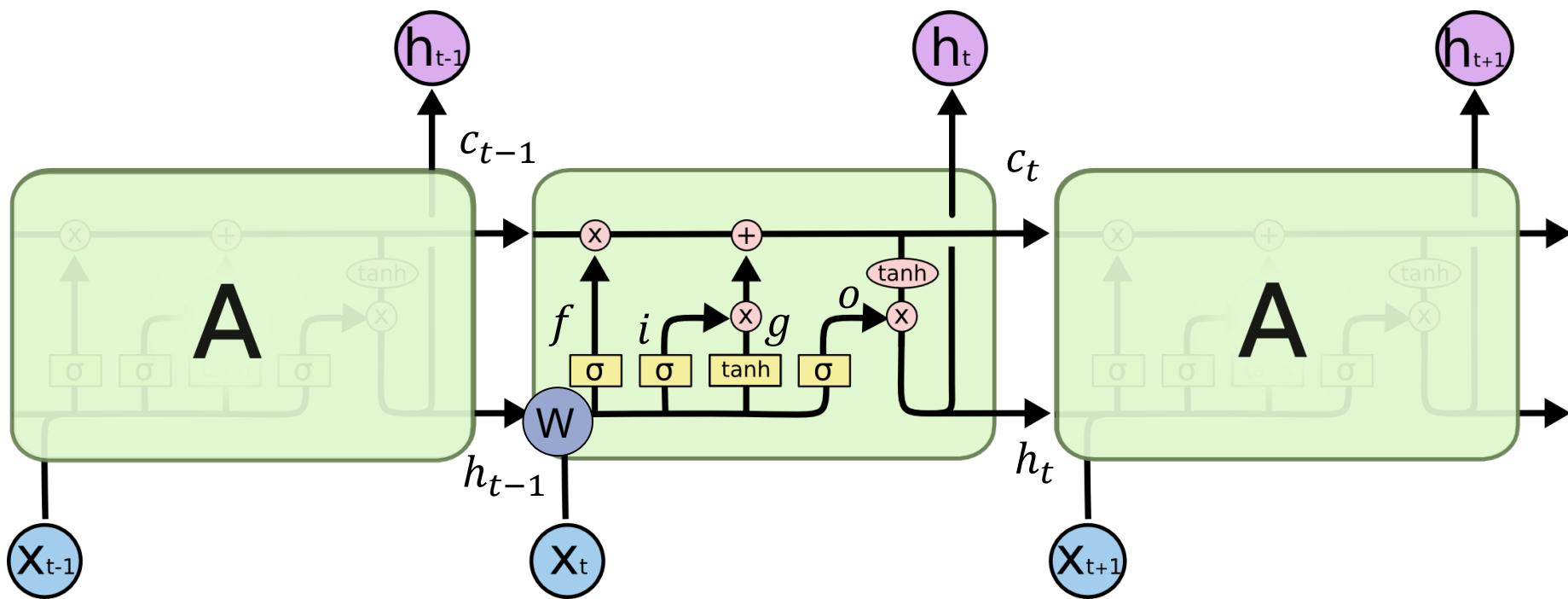
An LSTM can emulate a simple RNN by remembering with  $i$  and  $o$ , almost

# LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} x_t \\ h_{t-1}^l \end{pmatrix}$$
$$c_t = f \circ c_{t-1} + i \circ g$$
$$h_t = o \circ \tanh(c_t)$$

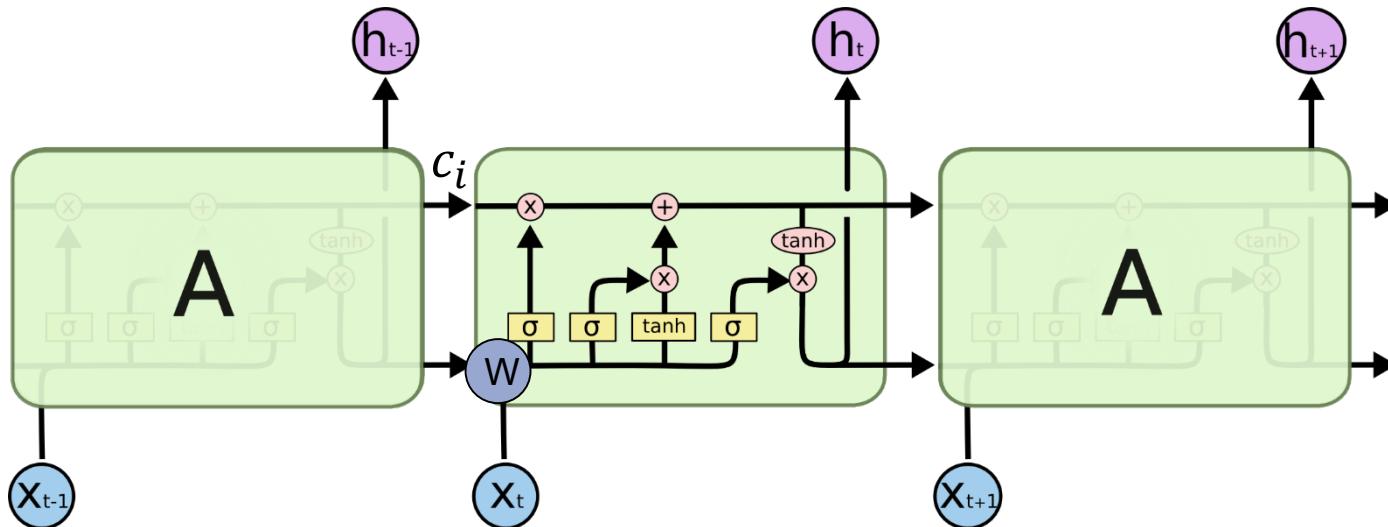
Input from below  
Input from left

cell Input (left)



# LSTM Arrays

- We now have two recurrent nodes,  $c_i$  and  $h_i$ .
- $h_i$  plays the role of the output in the simple RNN, and is recurrent.
- The cell state  $c_i$  is the cell's *memory*, it undergoes no transform.
- When we compose LSTMs into arrays, they look like this:

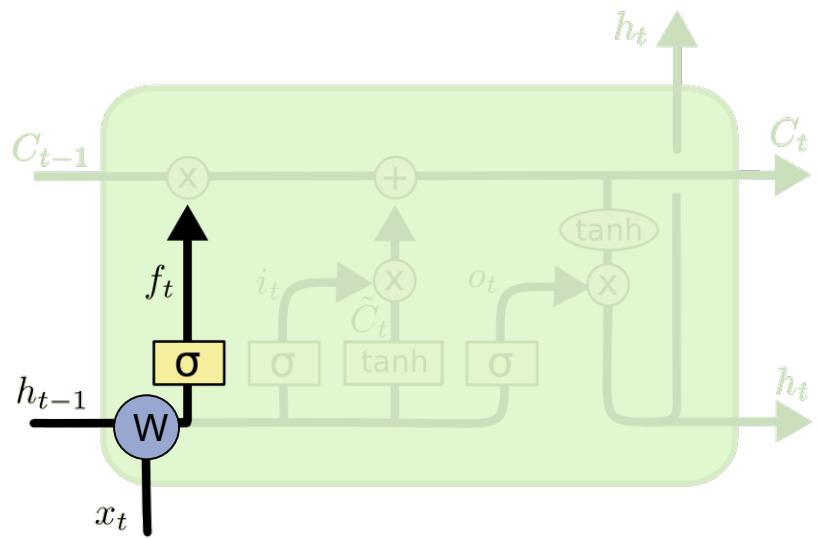


- For stacked arrays, the hidden layers ( $h_i$ 's) become the inputs ( $x_i$ 's) for the layer above.

Figure courtesy Chris Olah <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

# LSTM Walkthrough

First compute a  $[0, 1]$ -value  $f_t$  to “remember” the memory state.

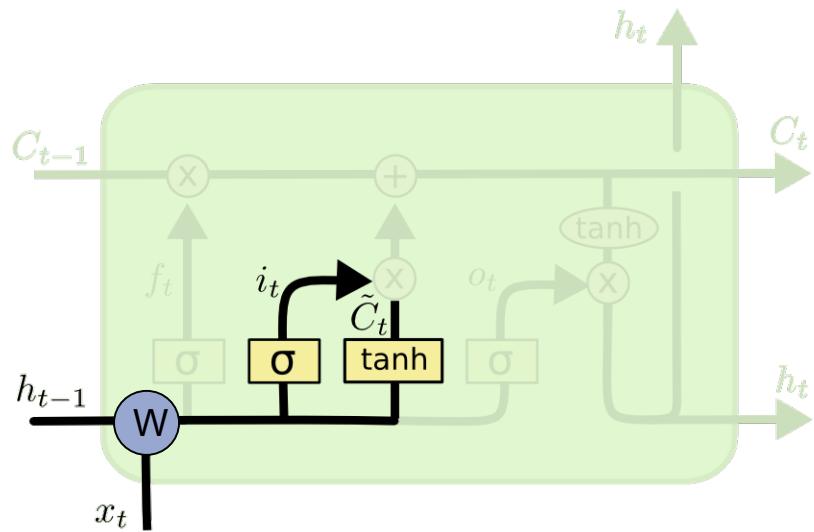


$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

Figure courtesy Chris Olah <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

# LSTM Walkthrough

Next compute an update to the memory cell value  $C_t$ , with its own gate value  $i_t$ .



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Figure courtesy Chris Olah <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

# LSTM Walkthrough

Now apply the memory gating and add the new update for the memory cell.

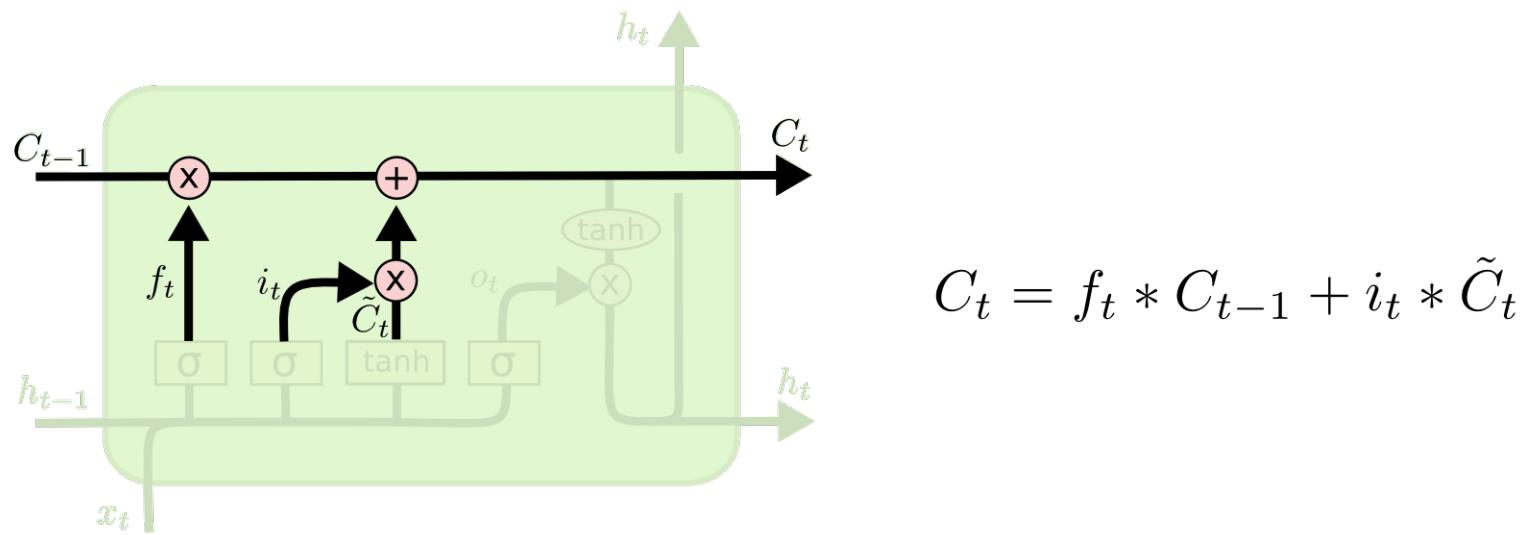
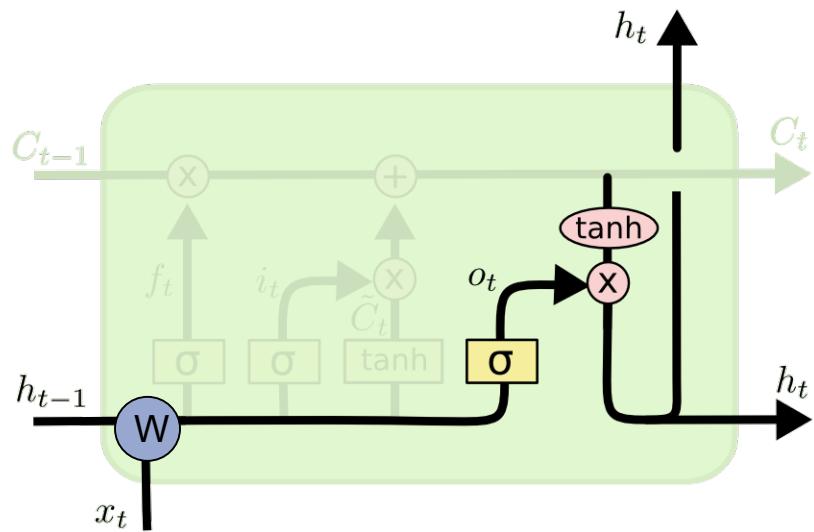


Figure courtesy Chris Olah <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

# LSTM Walkthrough

Finally compute a new hidden state gating value  $o_t$  and a new “output” value in  $[-1, 1]$  as  $\tanh(C_t)$



$$o_t = \sigma (W_o [ h_{t-1}, x_t ] + b_o)$$
$$h_t = o_t * \tanh (C_t)$$

Figure courtesy Chris Olah <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

# LSTM Summary

1. Decide what to forget
2. Decide what new things to remember
3. Decide what to output

# Searching for interpretable cells

"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."

quote detection cell

(LSTM, tanh(c), red = -1, blue = +1)

[*Visualizing and Understanding Recurrent Networks, Andrej Karpathy\*, Justin Johnson\*, Li Fei-Fei*]

# Searching for interpretable cells

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.

line length tracking cell

[Visualizing and Understanding Recurrent Networks, Andrej Karpathy\*, Justin Johnson\*, Li Fei-Fei]

# Searching for interpretable cells

```
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;
}
```

if statement cell

[Visualizing and Understanding Recurrent Networks, Andrej Karpathy\*, Justin Johnson\*, Li Fei-Fei]

# Searching for interpretable cells

```
/* Duplicate LSM field information. The lsm_rule is opaque, so
 * re-initialized. */
static inline int audit_dupe_lsm_field(struct audit_field *df,
                                       struct audit_field *sf)
{
    int ret = 0;
    char *lsm_str;
    /* our own copy of lsm_str */
    lsm_str = kstrdup(sf->lsm_str, GFP_KERNEL);
    if (unlikely(!lsm_str))
        return -ENOMEM;
    df->lsm_str = lsm_str;
    /* our own (refreshed) copy of lsm_rule */
    ret = security_audit_rule_init(df->type, df->op, df->lsm_str,
                                   (void **) &df->lsm_rule);
    /* Keep currently invalid fields around in case they
     * become valid after a policy reload. */
    if (ret == -EINVAL) {
        pr_warn("audit rule for LSM \\'%s\\' is invalid\n",
               df->lsm_str);
        ret = 0;
    }
    return ret;
}
```

quote/comment cell

[Visualizing and Understanding Recurrent Networks, Andrey Karpathy\*, Justin Johnson\*, Li Fei-Fei]

# Searching for interpretable cells

```
#ifdef CONFIG_AUDITSYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```

code depth cell

# LSTM stability

How fast can the cell value  $c$  grow with time?

How fast can the backpropagated gradient of  $c$  grow with time?

Cell Value Growth?, Backprop gradient growth?

- A. Linear, Linear
- B. Linear, Exponential
- C. Exponential, Linear
- D. Exponential, Exponential

# Correct!

How fast can the cell value  $c$  grow with time?

How fast can the backpropagated gradient of  $c$  grow with time?

Cell Value Growth?, Backprop gradient growth?

A. Linear, Linear

Yes, the  $c$ -path includes only multiplication by a factor in  $[0, 1]$  and the addition of a factor in  $[-1, 1]$

Try Again

Continue

# Oops!

How fast can the cell value  $c$  grow with time?

How fast can the backpropagated gradient of  $c$  grow with time?

Cell Value Growth?, Backprop gradient growth?

B. Linear, Exponential

Try Again

Continue

# Oops!

How fast can the cell value  $c$  grow with time?

How fast can the backpropagated gradient of  $c$  grow with time?

Cell Value Growth?, Backprop gradient growth?

C. Exponential, Linear

Try Again

Continue

# Oops!

How fast can the cell value  $c$  grow with time?

How fast can the backpropagated gradient of  $c$  grow with time?

Cell Value Growth?, Backprop gradient growth?

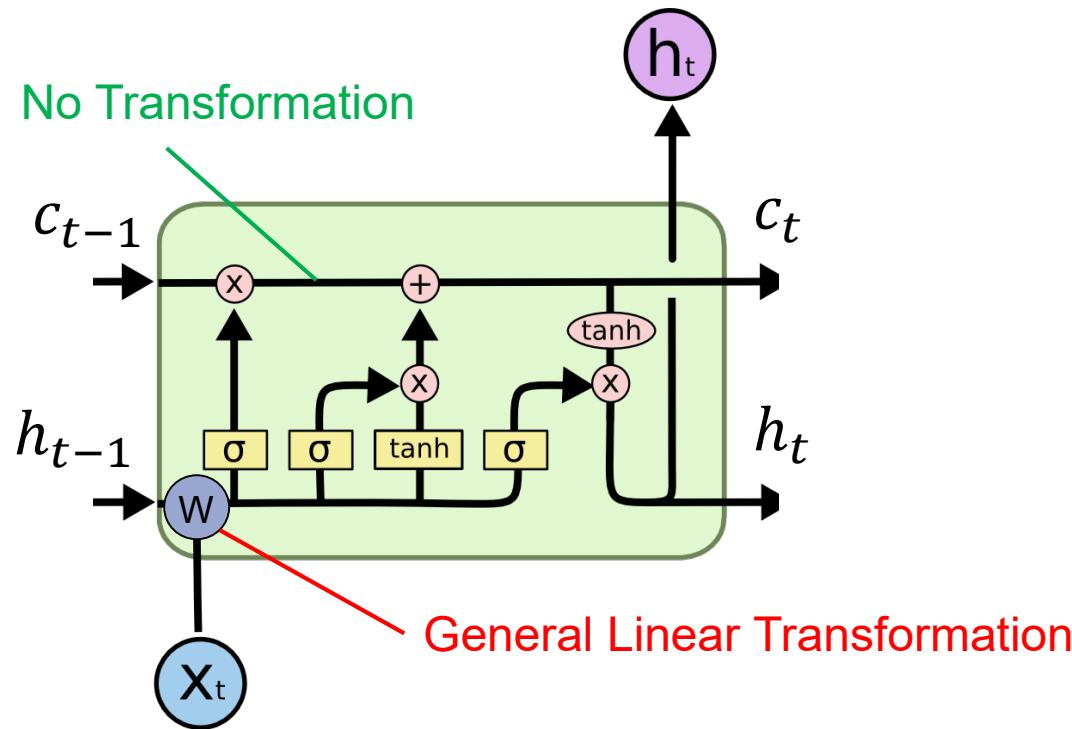
D. Exponential, Exponential

Try Again

Continue

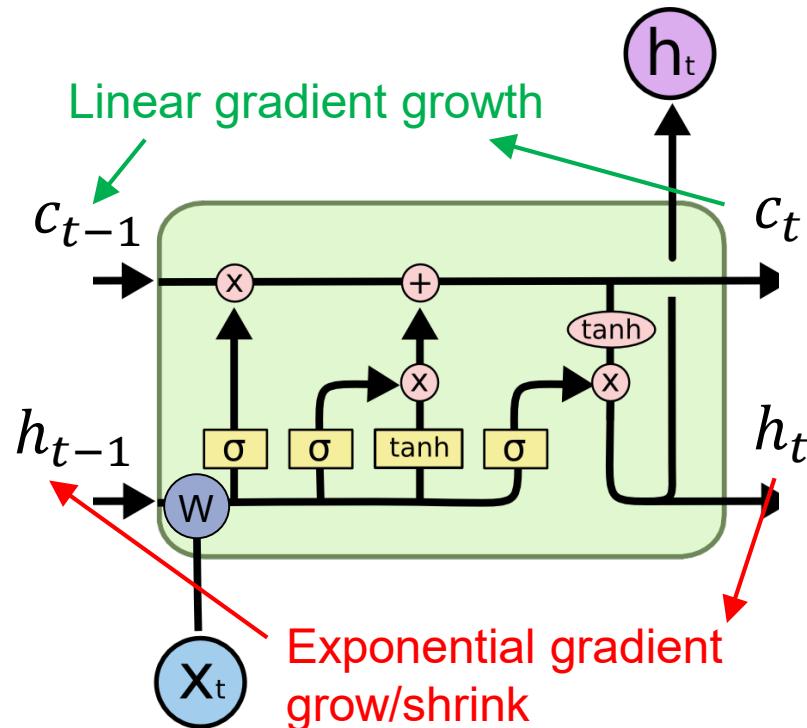
# LSTM Stability

Remember that the  $h$  path has a linear transform  $W_{hh}$  at each node.



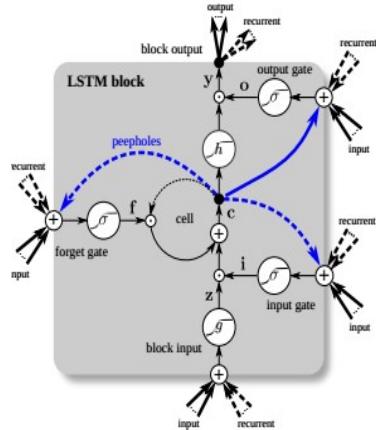
# LSTM Stability

Gradients are well-behaved along the  $c$  path but not the  $h$  path. Luckily LSTMs learn to rely mostly on  $c$  for long-term memory.



# LSTM variants and friends

[*An Empirical Exploration of Recurrent Network Architectures*, Jozefowicz et al., 2015]



[*LSTM: A Search Space Odyssey*,  
Greff et al., 2015]

**GRU** [*Learning phrase representations using rnn encoder-decoder for statistical machine translation*, Cho et al. 2014]

$$\begin{aligned} r_t &= \text{sigm}(W_{xr}x_t + W_{hr}h_{t-1} + b_r) \\ z_t &= \text{sigm}(W_{xz}x_t + W_{hz}h_{t-1} + b_z) \\ \tilde{h}_t &= \tanh(W_{xh}x_t + W_{hh}(r_t \odot h_{t-1}) + b_h) \\ h_t &= z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t \end{aligned}$$

MUT1:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + b_z) \\ r &= \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + \tanh(x_t) + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

MUT2:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + W_{hz}h_t + b_z) \\ r &= \text{sigm}(x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

MUT3:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + W_{hz}\tanh(h_t) + b_z) \\ r &= \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

# LSTM variants and friends

These designs combine the function of  $c$  and  $h$ ,  
and (similar to ResNet) always include an identity  
path to the previous  $h$  (memory path).

## GRU:

These designs combine the function of  $c$  and  $h$ ,  
and (similar to ResNet) always include an identity  
path to the previous  $h$  (memory path).

MUT1:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + b_z) \\ r &= \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + \tanh(x_t) + b_h) \odot z \\ &\quad + h_t \odot (1 - z) \end{aligned}$$

MUT2:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + W_{hz}h_t + b_z) \\ r &= \text{sigm}(x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\ &\quad + h_t \odot (1 - z) \end{aligned}$$

MUT3:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + W_{hz}\tanh(h_t) + b_z) \\ r &= \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\ &\quad + h_t \odot (1 - z) \end{aligned}$$

# LSTM variants and friends

They also combine the functions of  $f$  (forgetting gate) and  $i$  (input) in a single  $z$  gate.

The more you add from the current state (large  $z$ ) the less  $(1 - z)$  you remember from previous  $h$ , and vice versa.

## GRU:

$$r_t = \text{sigm}(W_{xr}x_t + W_{hr}h_{t-1} + b_r)$$

$$z_t = \text{sigm}(W_{xz}x_t + W_{hz}h_{t-1} + b_z)$$

$$\tilde{h}_t = \tanh(W_{xh}x_t + W_{hh}(r_t \odot h_{t-1}) + b_h)$$

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t$$

## MUT1:

$$z = \text{sigm}(W_{xz}x_t + b_z)$$

$$r = \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r)$$

$$h_{t+1} = \tanh(W_{hh}(r \odot h_t) + \tanh(x_t) + b_h) \odot z \\ + h_t \odot (1 - z)$$

## MUT2:

$$z = \text{sigm}(W_{xz}x_t + W_{hz}h_t + b_z)$$

$$r = \text{sigm}(x_t + W_{hr}h_t + b_r)$$

$$h_{t+1} = \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\ + h_t \odot (1 - z)$$

## MUT3:

$$z = \text{sigm}(W_{xz}x_t + W_{hz}\tanh(h_t) + b_z)$$

$$r = \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r)$$

$$h_{t+1} = \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\ + h_t \odot (1 - z)$$

# Summary

- RNNs are a widely used model for sequential data, including text
- RNNs are trainable with backprop when unrolled over time
- RNNs learn complex and varied patterns in sequential data
- Vanilla RNNs are simple but don't work very well
  - Backward flow of gradients in RNN can explode or vanish
- Common to use LSTM or GRU. Memory path makes them stably learn long-distance interactions.