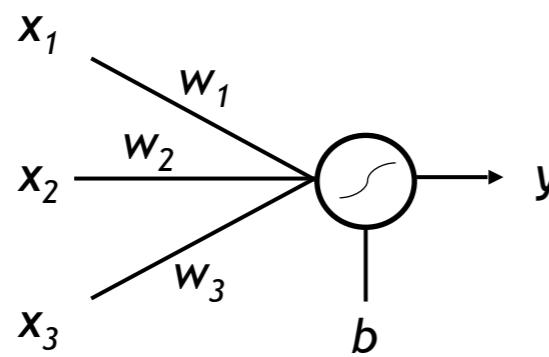


Recurrent Neural Networks

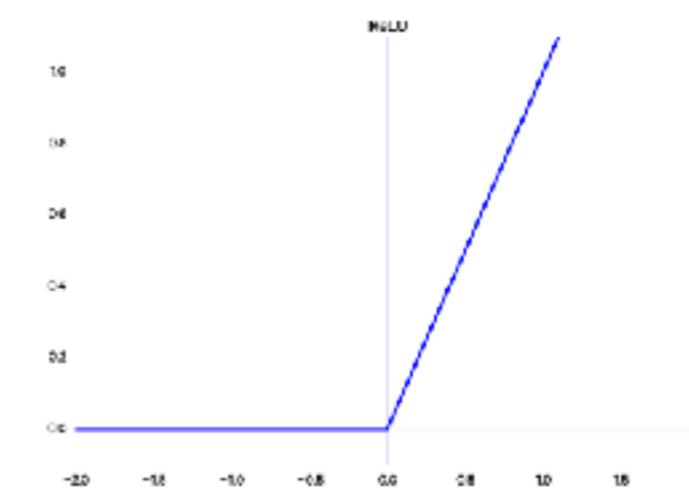
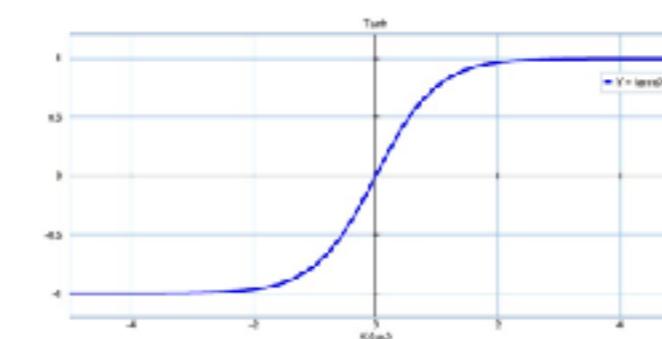
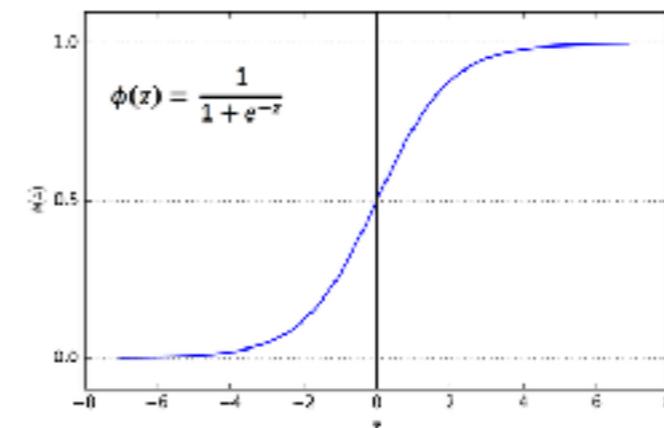
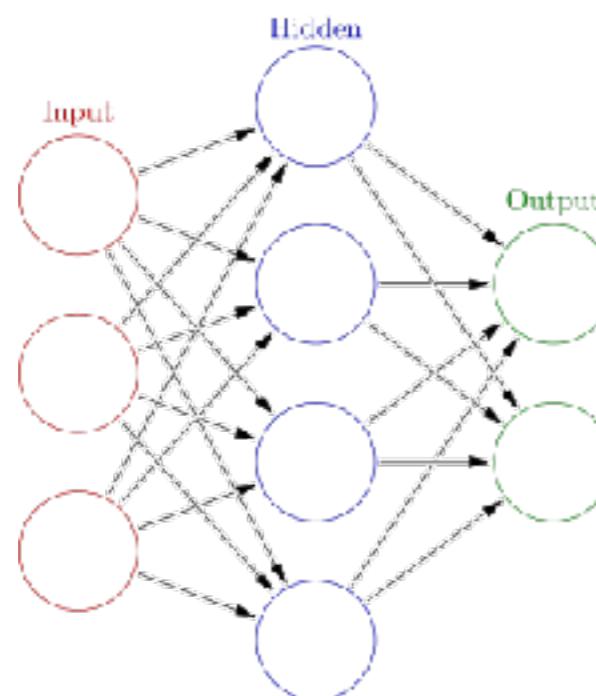
Pramod Kompalli
OLA

Recap

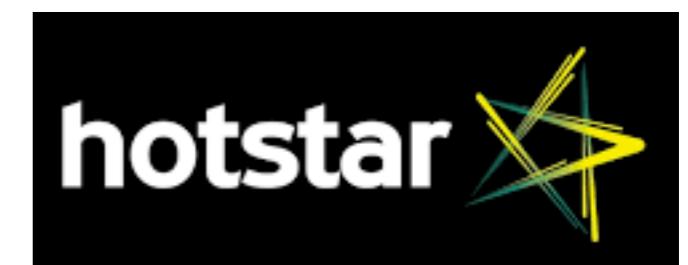


$$y = \begin{cases} 1, & g(w \cdot x + b) > 0 \\ 0, & \text{otherwise} \end{cases}$$

$g(\cdot)$ is a non-linear function



Videos



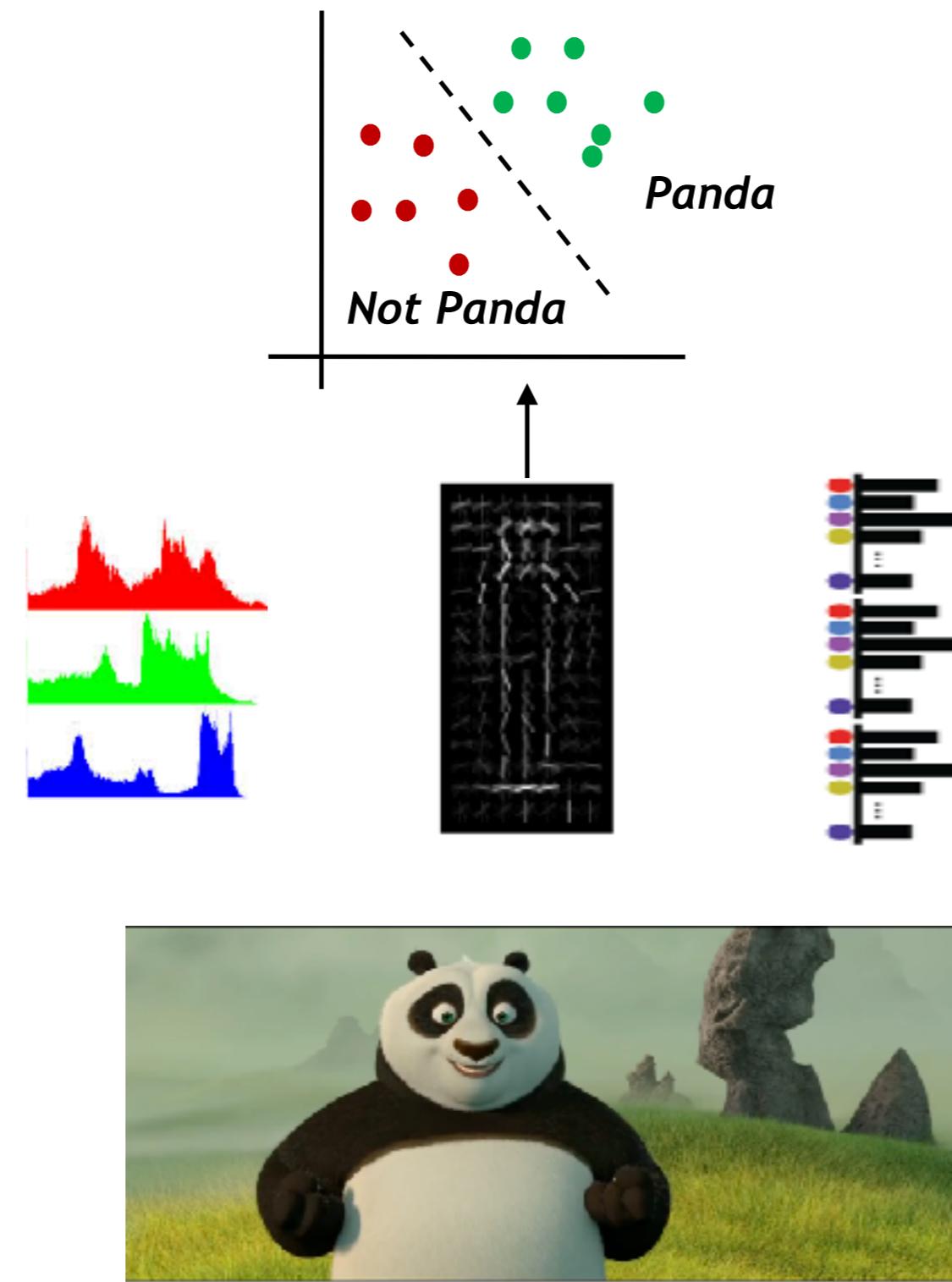
Videos



Image vs. Video

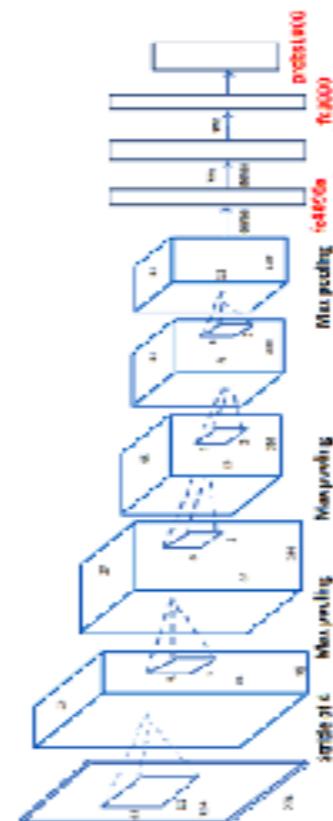
- Arbitrary length in T
- Scale
- Deeper Semantics

Representing Images

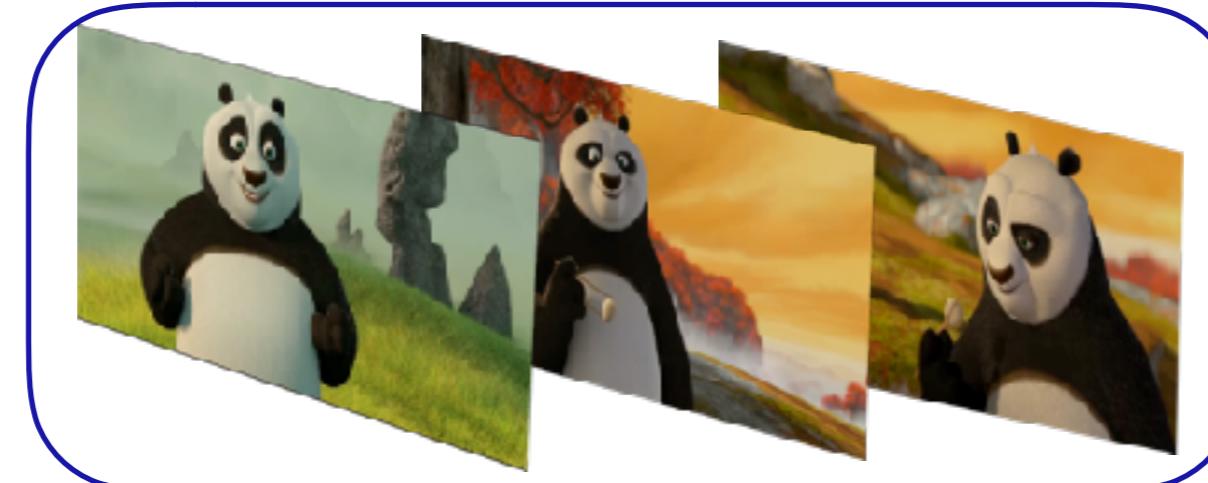
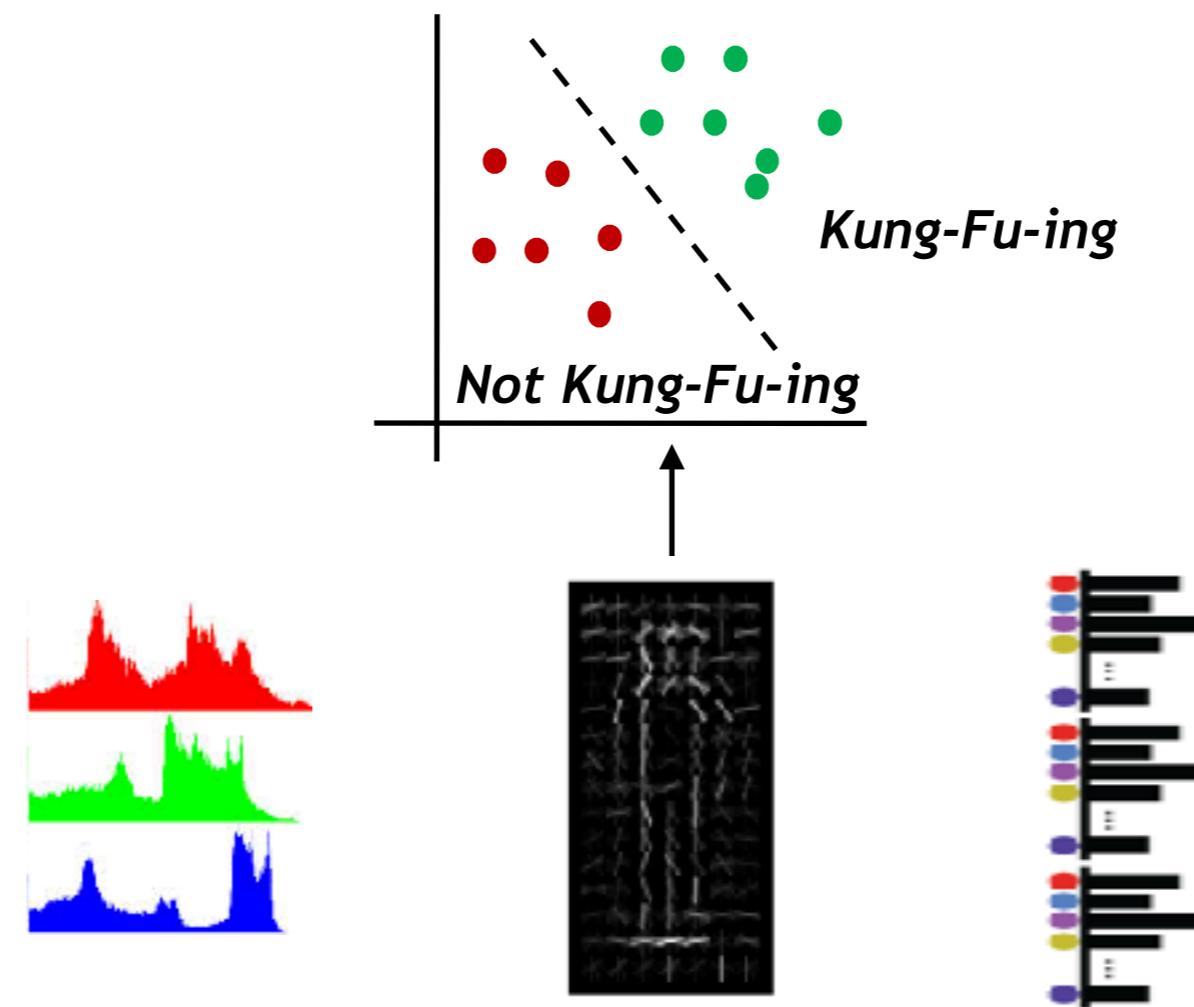


Modern Image Rep.

Panda

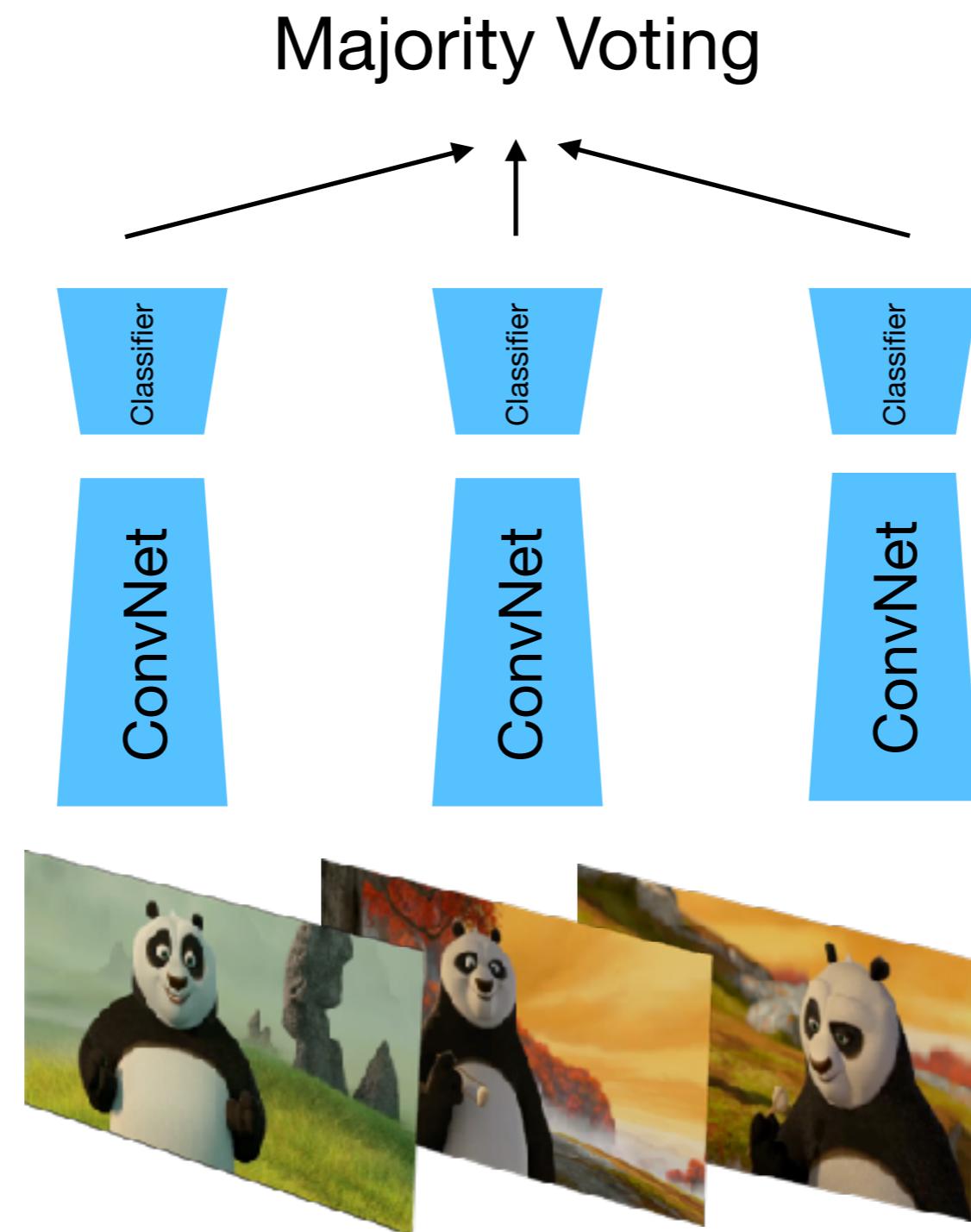


Simple Video Rep.

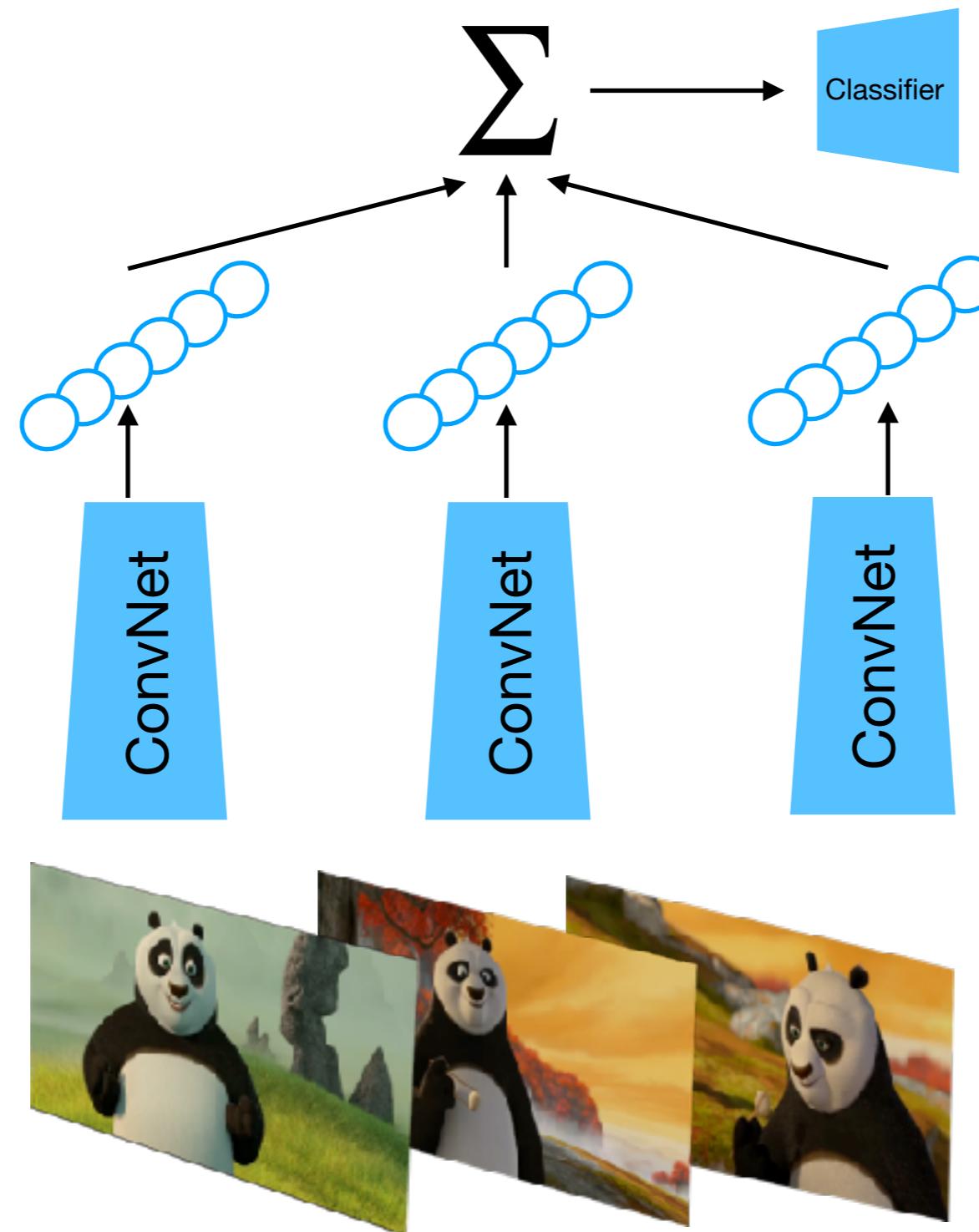


**What could be
the issues?**

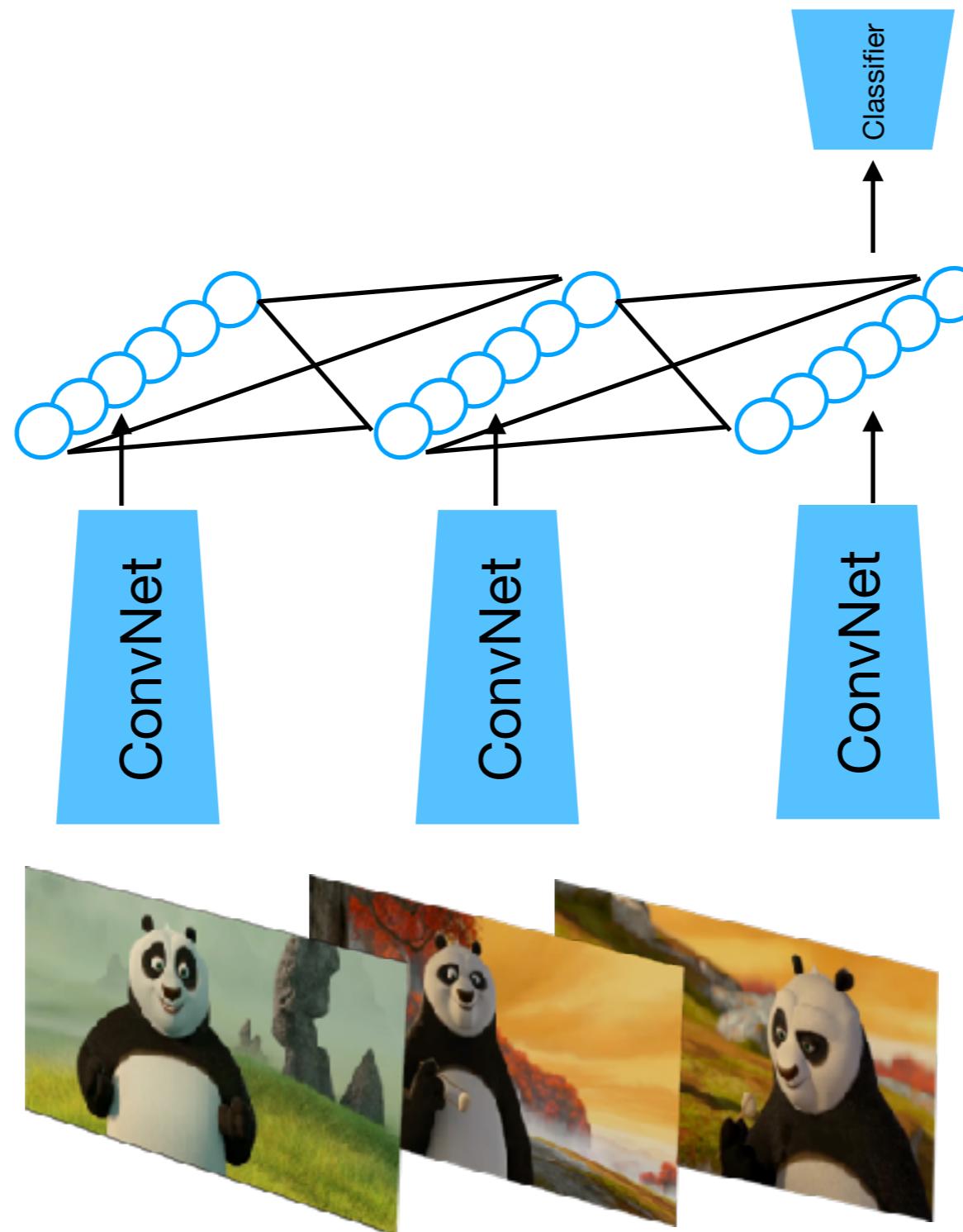
Late Fusion



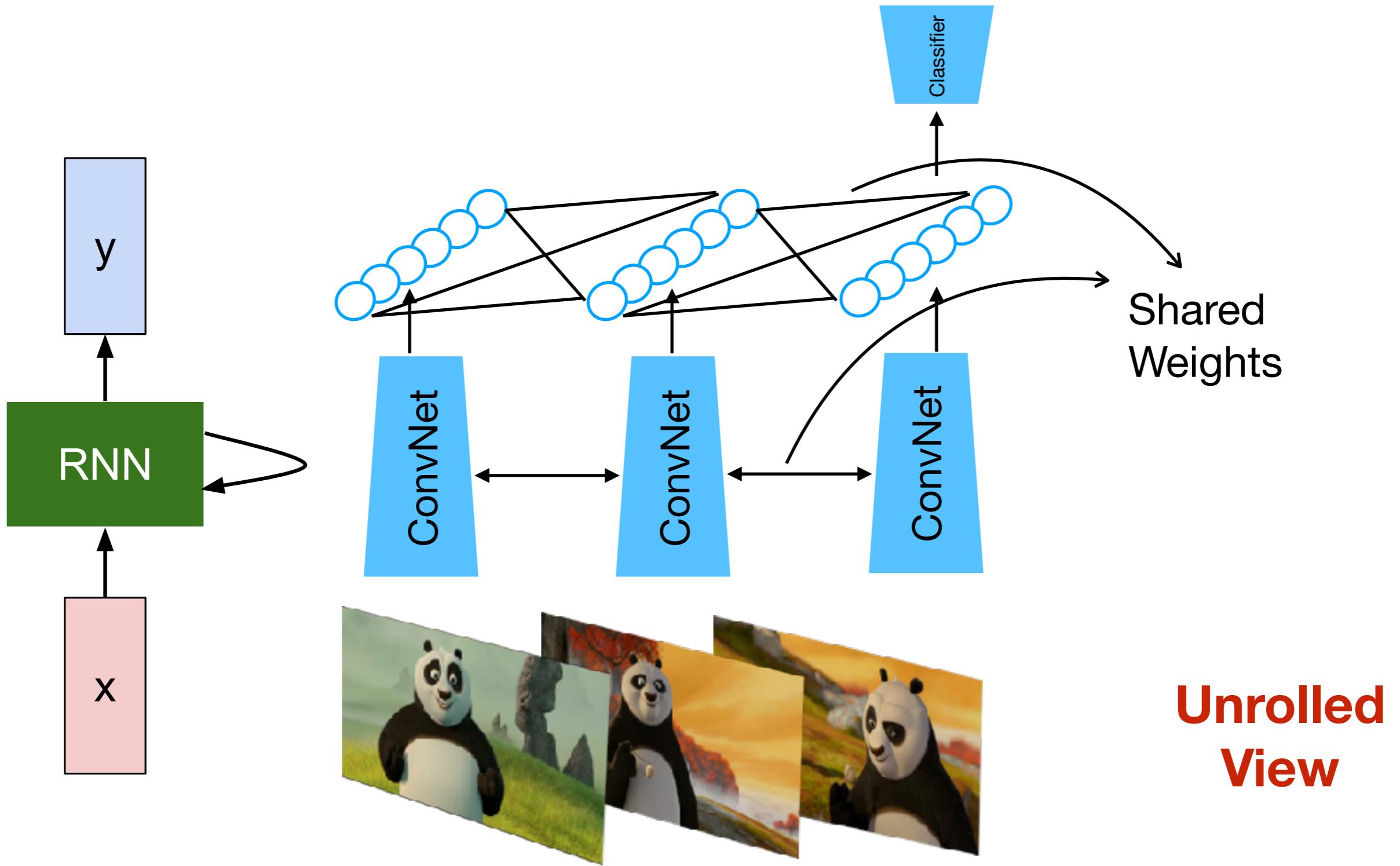
Slightly-Early Fusion



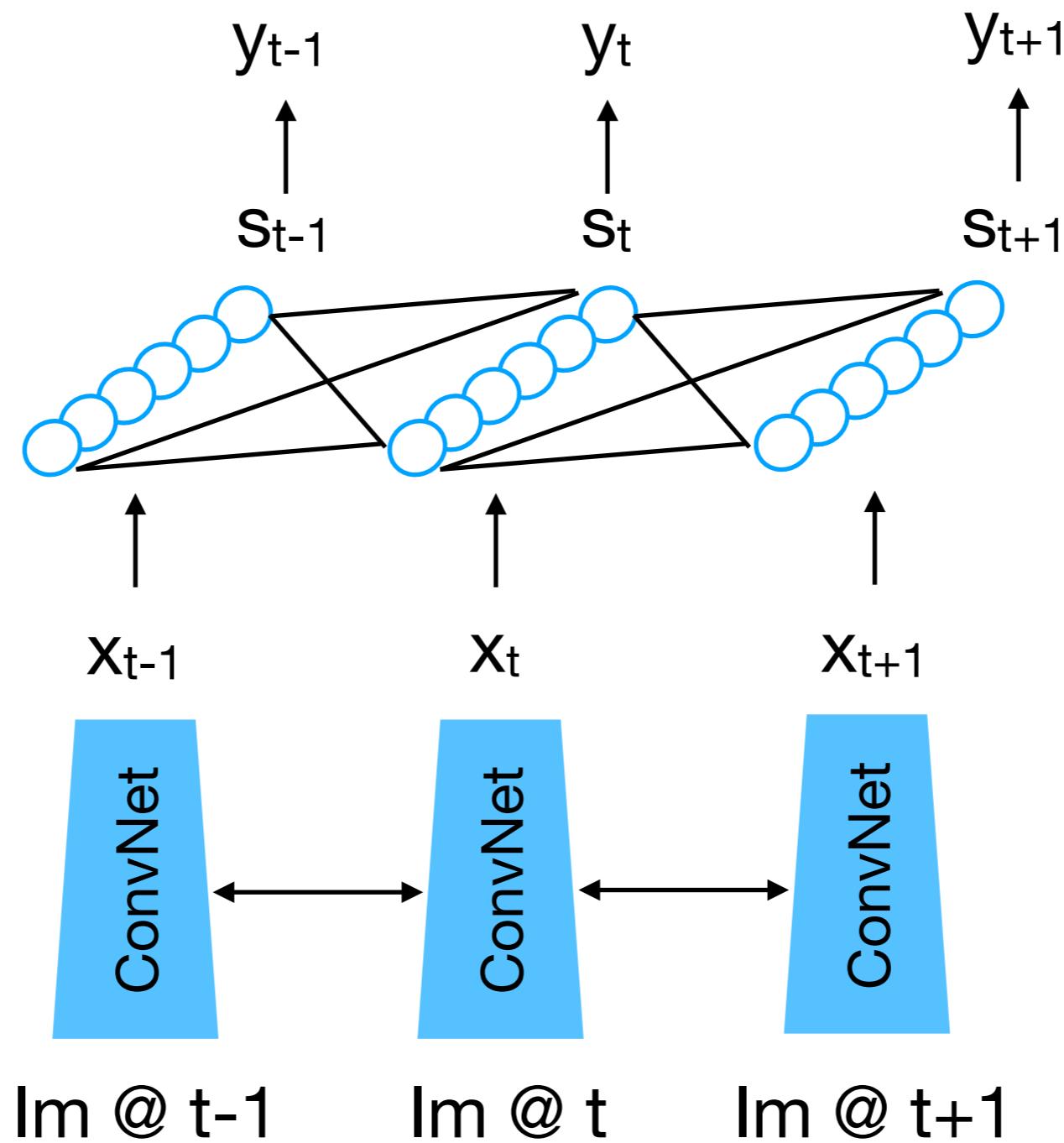
Fusing Representations



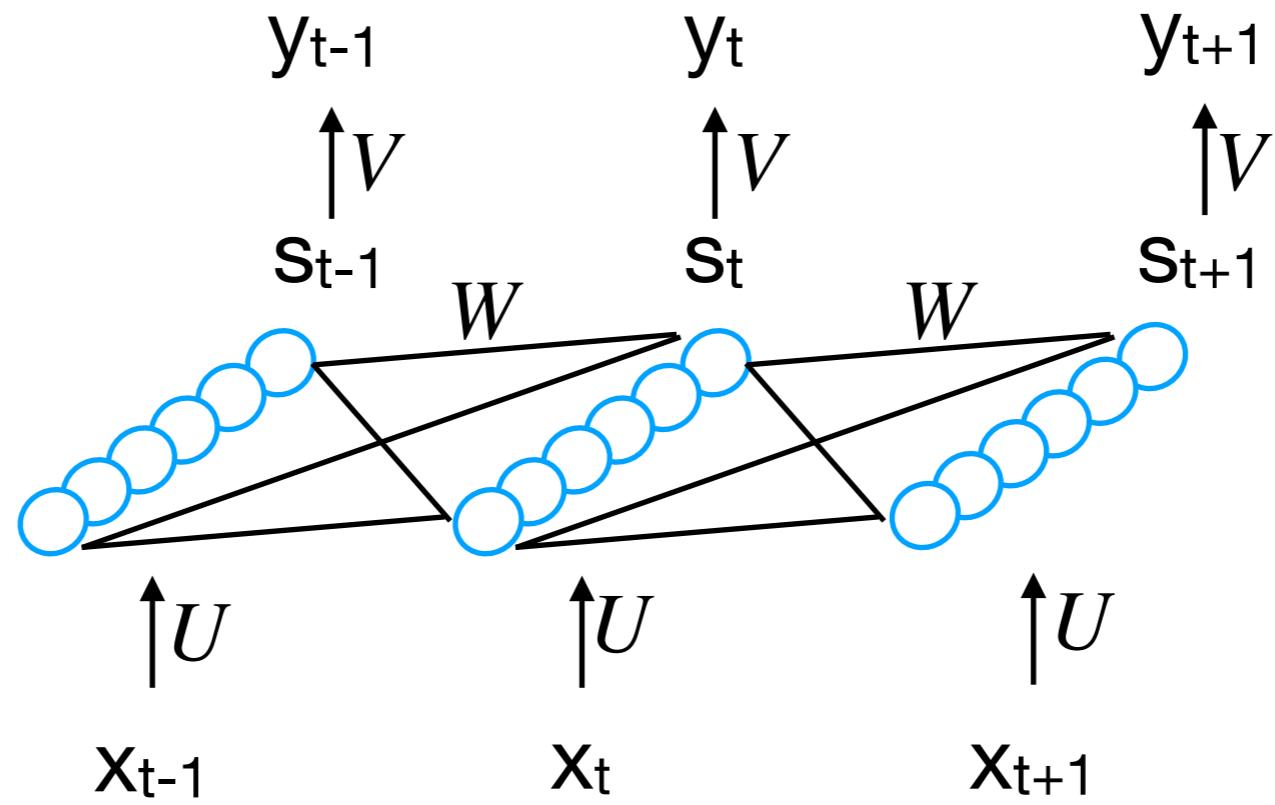
Fusing Representations



RNN Formulation



RNN Formulation

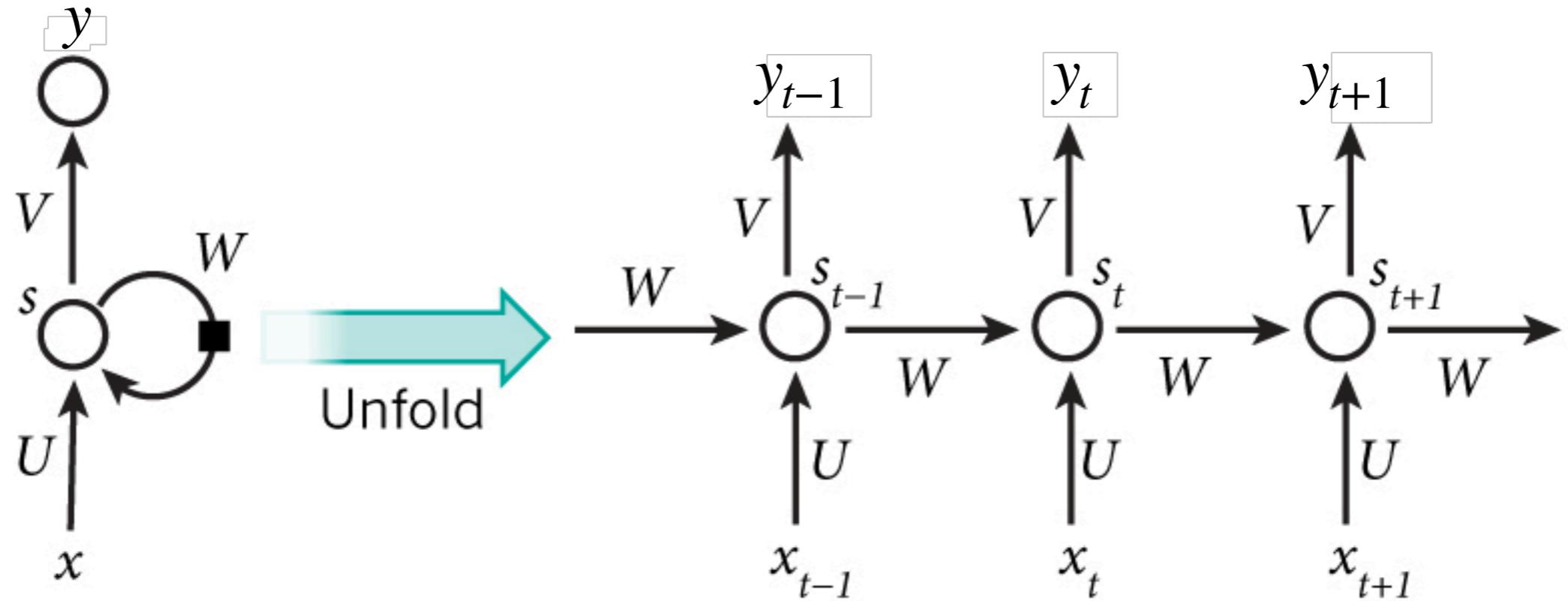


$$s_t = f_\theta(x_t, s_{t-1})$$

$$= \tanh(Ux_t + Ws_{t-1})$$

$$y_t = Vs_t$$

RNN Formulation



$$s_t = f_W(x_t, s_{t-1})$$

$$= \tanh(Ux_t + Ws_{t-1})$$

$$y_t = Vs_t$$

Depth in Layers
vs. Depth in Time

Training RNNs

- Typical loss: Cross Entropy (Predicted Label, True Label)

$$E_t = -\hat{y}_t \cdot \log(y_t)$$

- Back-propagation

$$\begin{aligned}\frac{\partial E_3}{\partial V} &= \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial V} \\ &= (y_3 - \hat{y}_3) \otimes s_3\end{aligned}$$

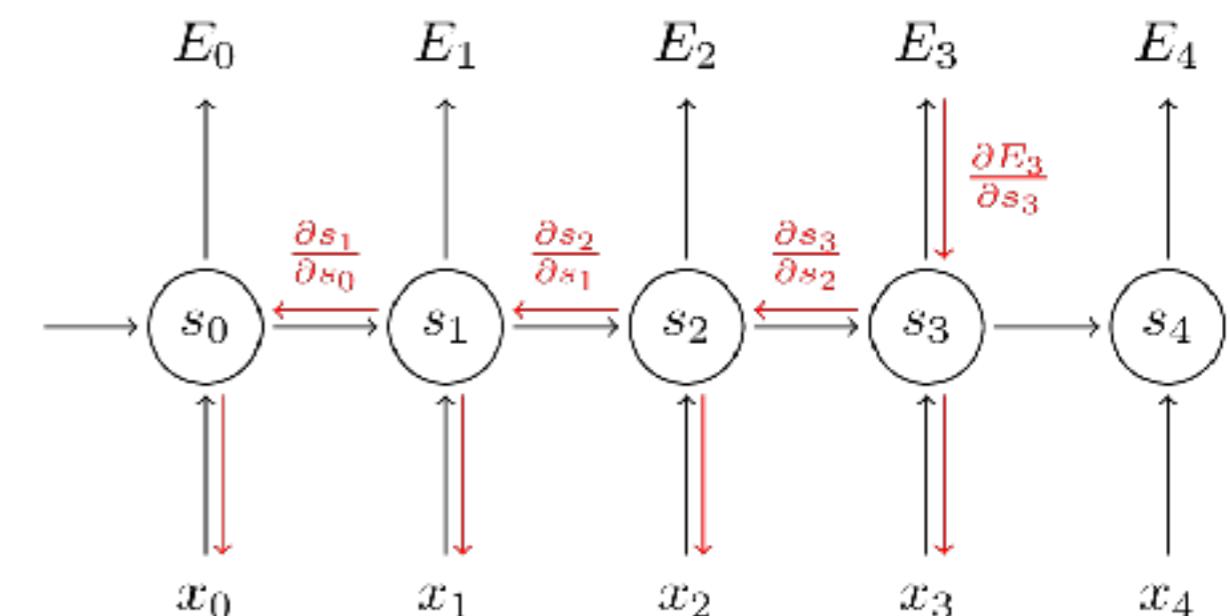
Training RNNs

- Back-propagation through Time (BPTT)

$$\frac{\partial E_3}{\partial W} = \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial W}$$

$$= \sum_{k=0}^3 \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_k} \frac{\partial s_k}{\partial W}$$

$$= \sum_{k=0}^3 \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \boxed{\prod_{j=k+1}^3 \frac{\partial s_j}{\partial s_{j-1}}} \frac{\partial s_k}{\partial W}$$



$$\frac{\partial s_3}{\partial s_1} = \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1}$$

Vanishing/Exploding Gradients

$$\prod_{j=k+1}^3 \frac{\partial s_j}{\partial s_{j-1}}$$

If Scalar?

If Gradient $\rightarrow 0$, no information flows back
long term information is lost

If Gradient $\rightarrow \infty$, information overflows
historical information is lost
(could be handled with “Clipping”)

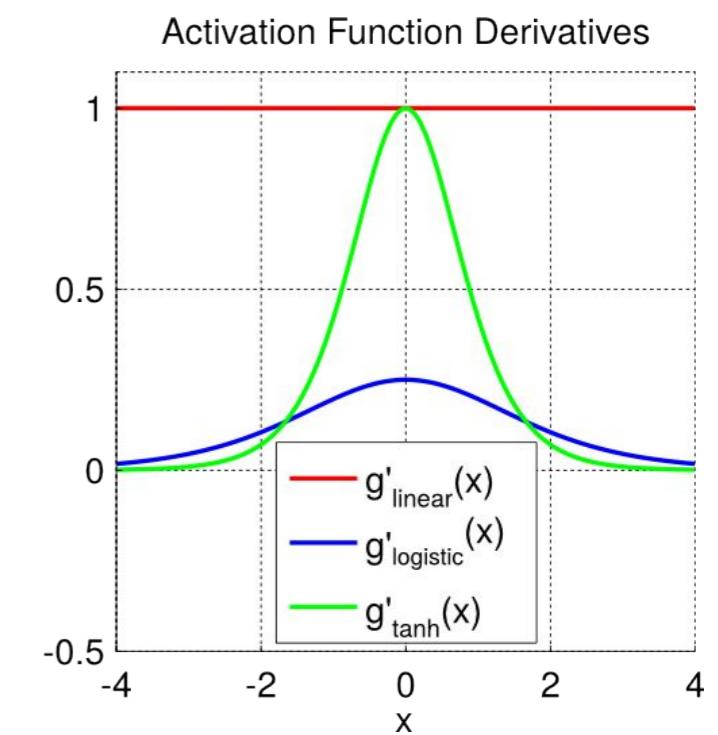
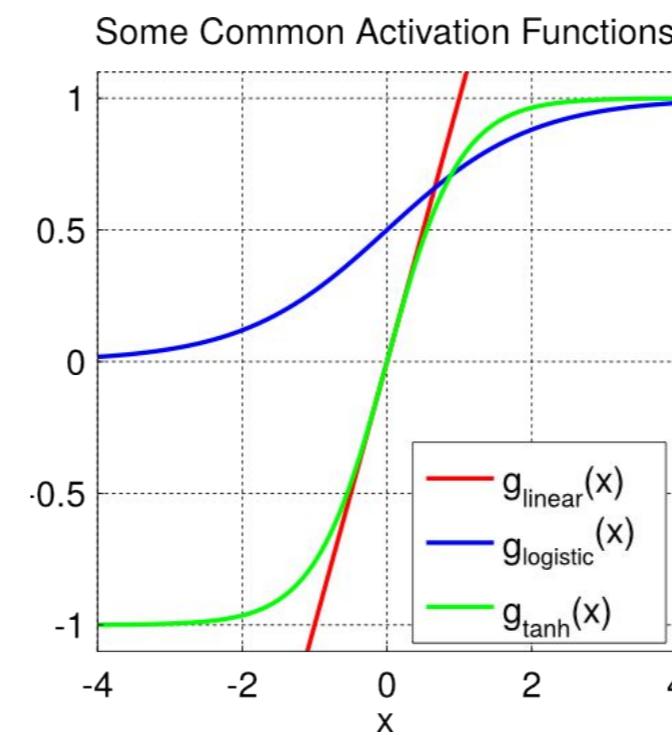
Vanishing/Exploding Gradients

$$\frac{d\mathbf{f}}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \dots & \frac{\partial \mathbf{f}}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$\left\| \frac{\partial h_t}{\partial h_k} \right\| = \left\| \prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}} \right\| \leq (\beta_W \beta_h)^{t-k}$$

Vanishing/Exploding Gradients

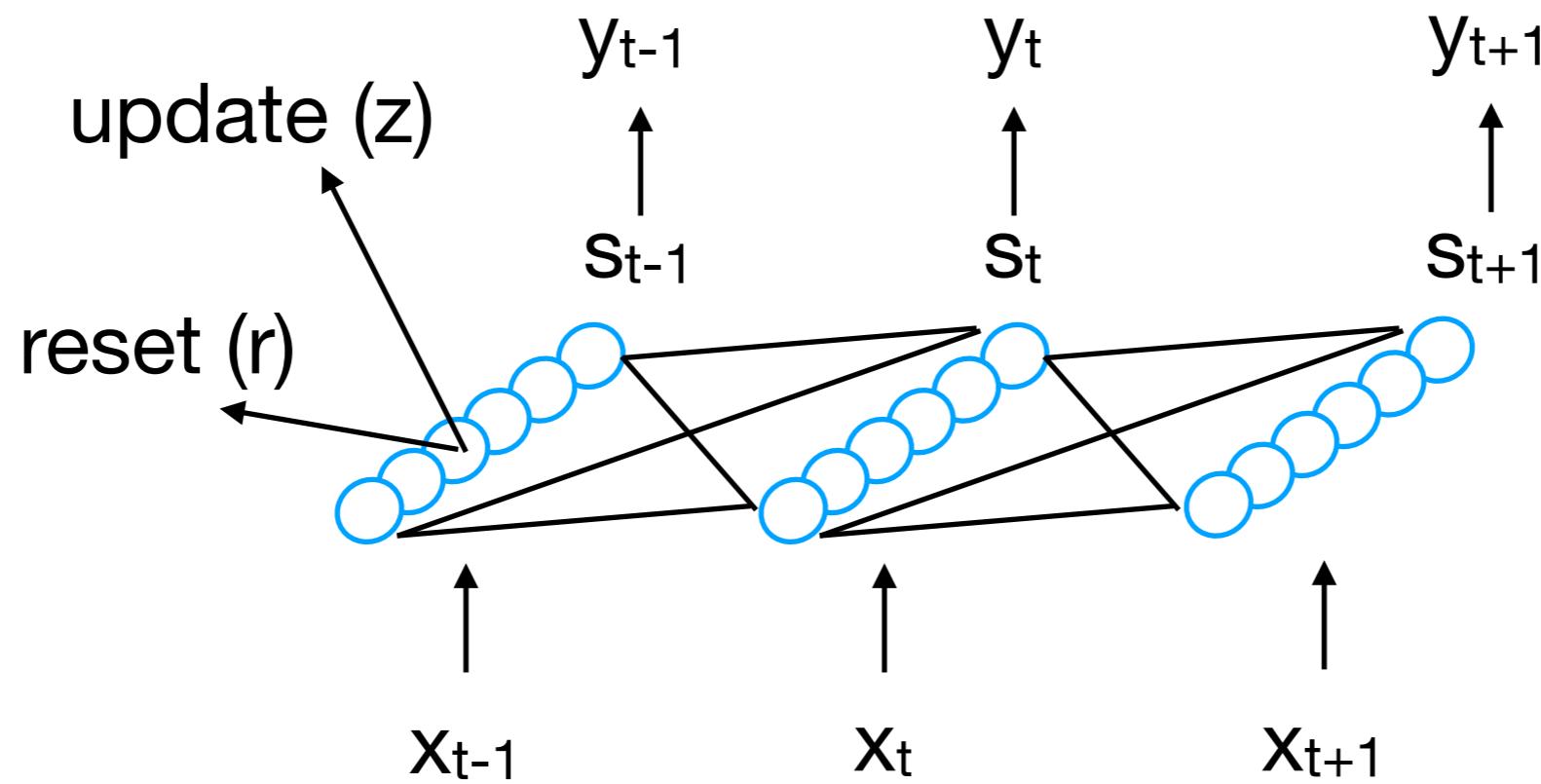
- Tanh is slightly better, derivative bounded by 1
- Sigmoid is bounded by 1/4
- ReLU?
 - Need good init
 - Just use Better Units



Picture courtesy: Girish Varma

Gated Recurrent Units

- Core Idea: Each “Unit” or neuron keeps track of a different amount of “History”



Gated Recurrent Units

- Core Idea: Each “Unit” or neuron keeps track of a different amount of “History”
- Update Gate: How much of previous state vector we want to keep

$$z = \sigma(x_t U^z + s_{t-1} W^z)$$

- Reset Gate: How to combine current input with previous state vector

$$r = \sigma(x_t U^r + s_{t-1} W^r)$$

Gated Recurrent Units

$$z = \sigma(x_t U^z + s_{t-1} W^z)$$

$$r = \sigma(x_t U^r + s_{t-1} W^r)$$

$$g = \tanh(x_t U^g + (s_{t-1} \circ r) W^g)$$

$$s_t = (1 - z) \circ g + z \circ s_{t-1}$$

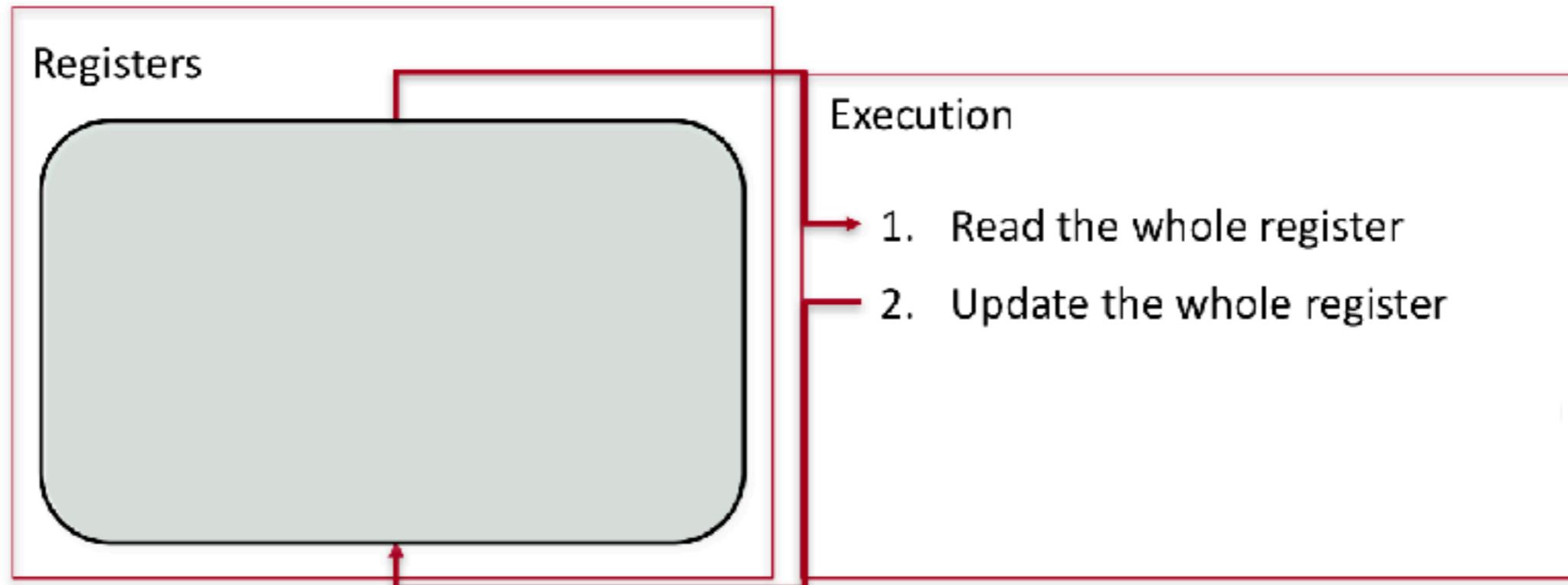
$$s_t = \tanh(x_t U + s_{t-1} W) \quad \text{Vanilla RNN Unit}$$

Gated Recurrent Units

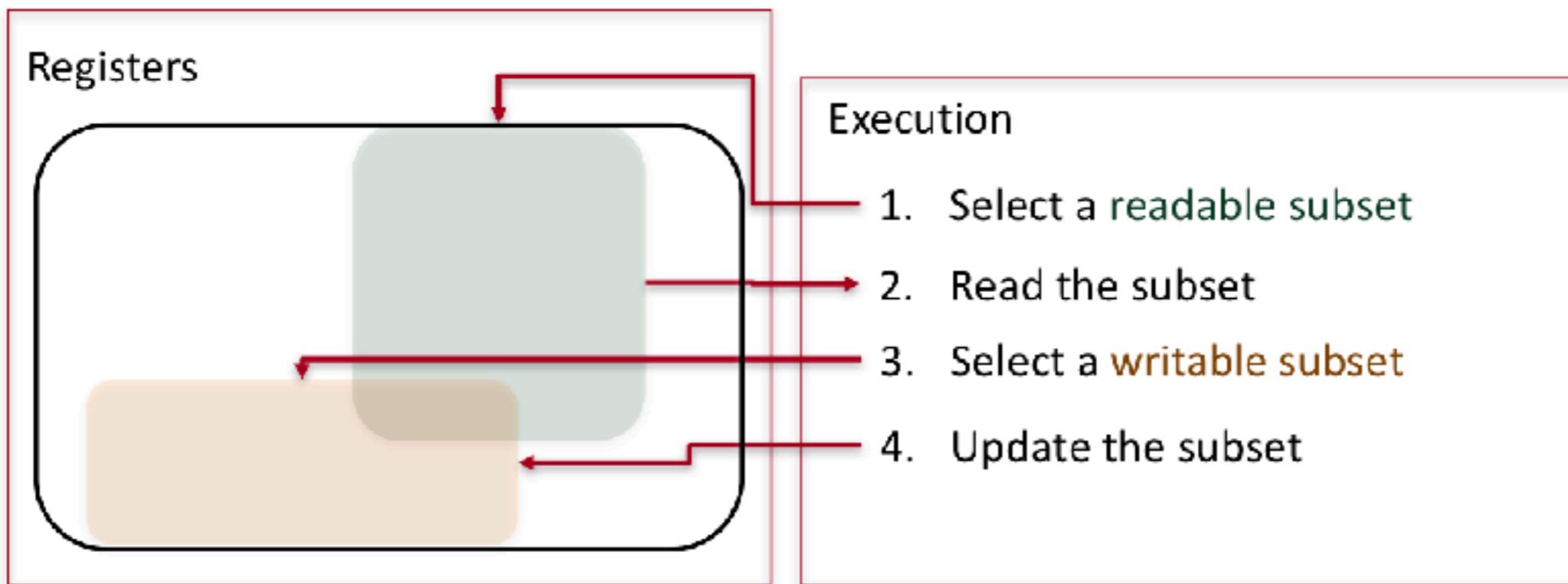
If Reset Gate $\rightarrow 0$, ignore previous state,
use only current input
History is forgotten

If Update Gate $\rightarrow 1$, ignore current input,
simply copy previous state
Gradient doesn't vanish

Vanilla RNN

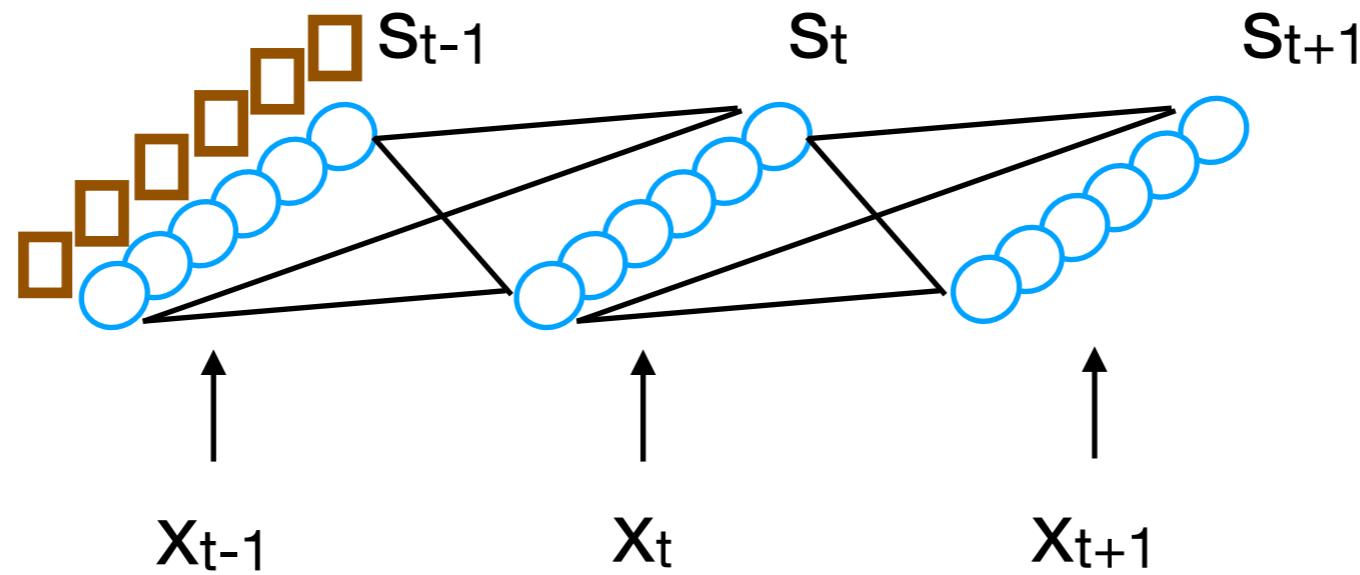


GRU

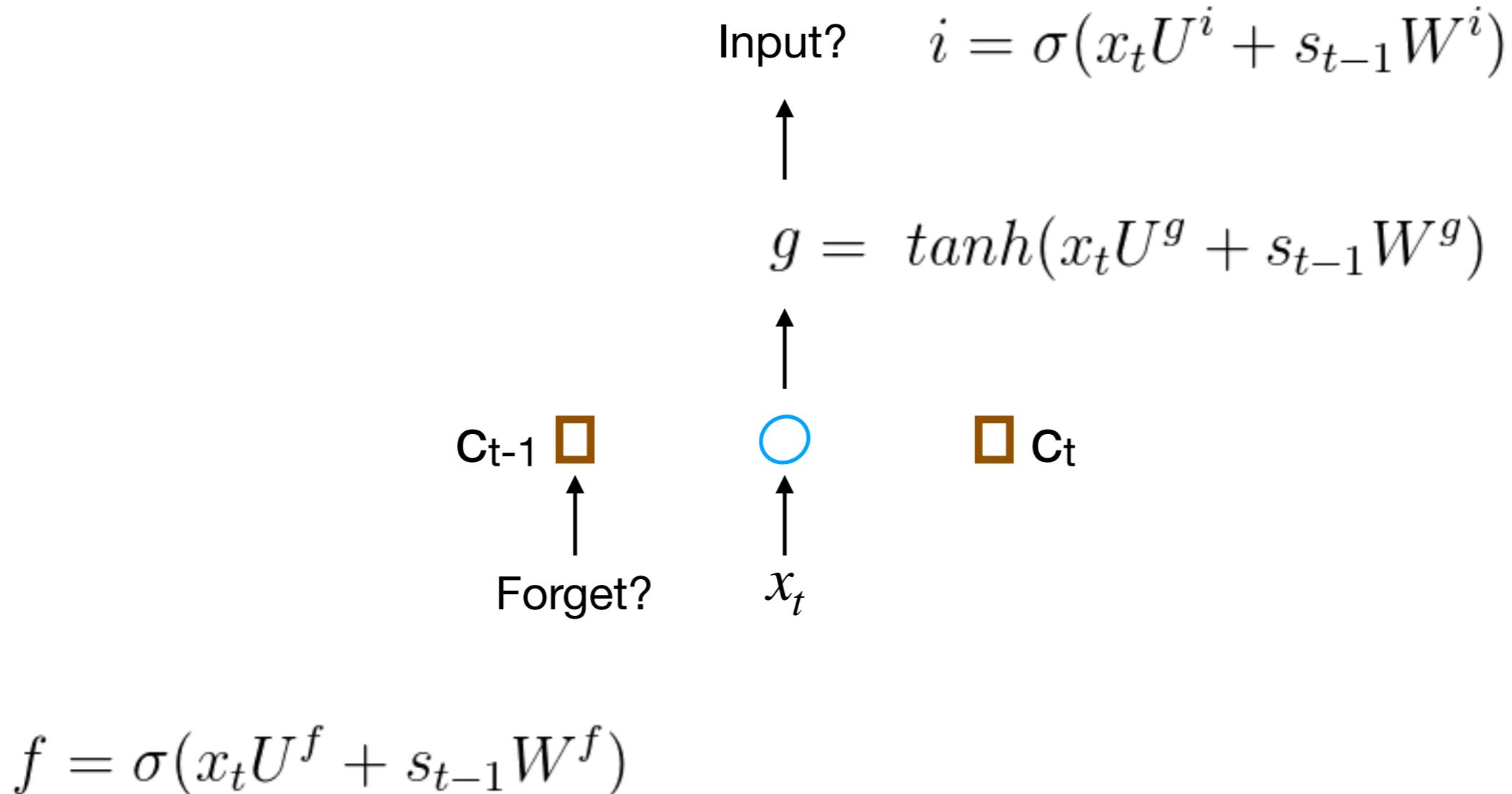


Long-Short Term Memory

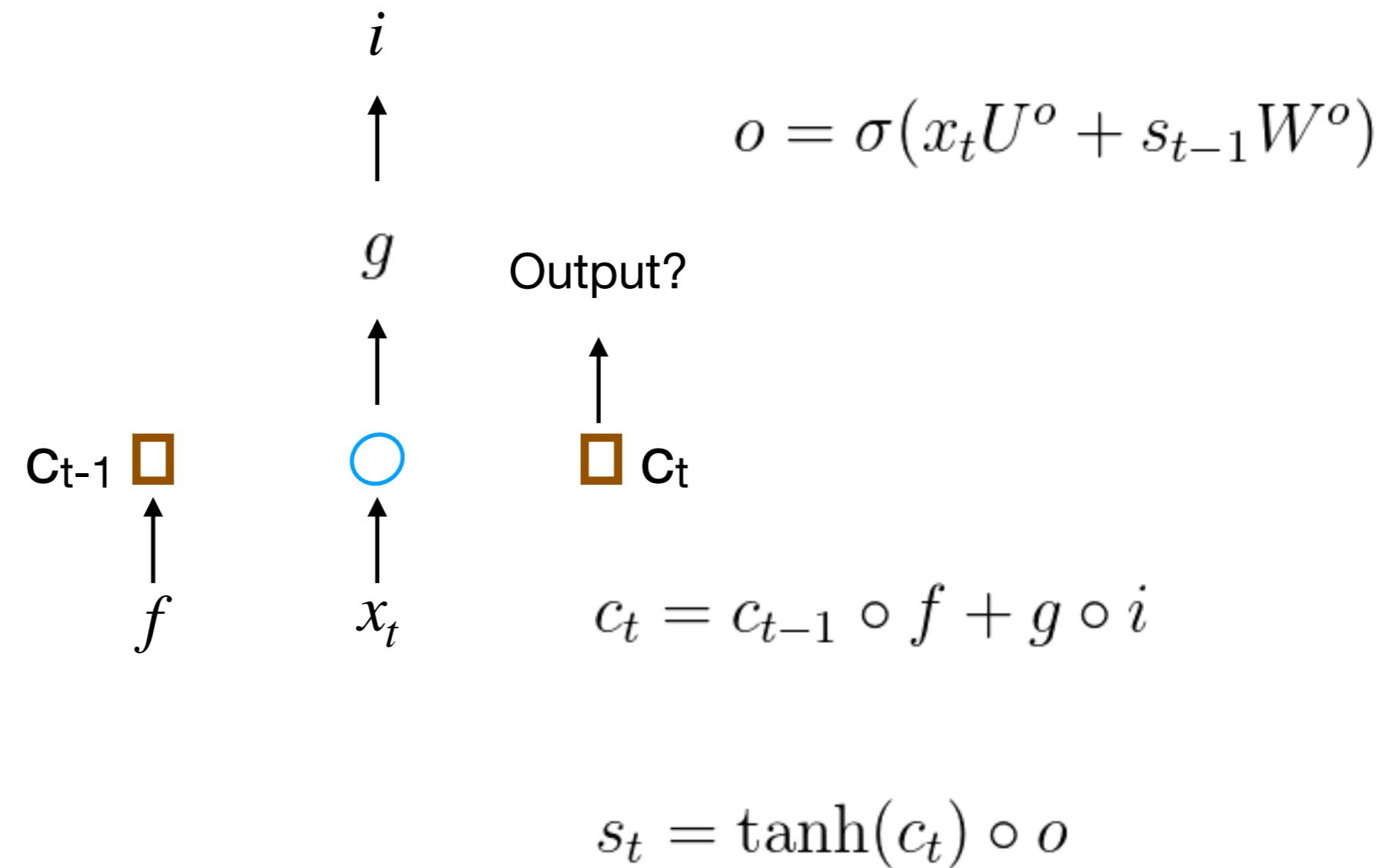
- Core Idea: Each “Unit” or neuron has a “Memory” cell in it



LSTM Unit



LSTM Unit



LSTM Unit

- Input Gate: How much of current input to use
- Forget Gate: How much of previous memory to use
- Output Gate: How much of current memory to use
- These gates explicitly model long-term dependencies
- Additive interaction

$$c_t = c_{t-1} \circ f + g \circ i$$

LSTM Unit

$$s_t = \tanh(x_t U + s_{t-1} W)$$

Vanilla RNN Unit

$$g = \tanh(x_t U^g + s_{t-1} W^g)$$

$$i = \sigma(x_t U^i + s_{t-1} W^i)$$

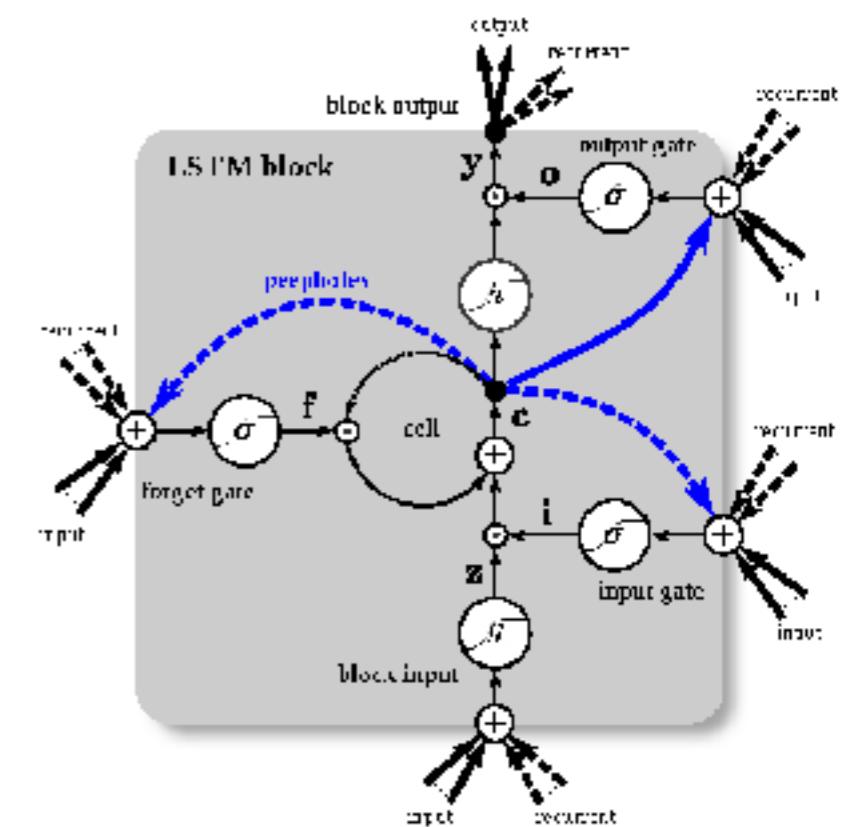
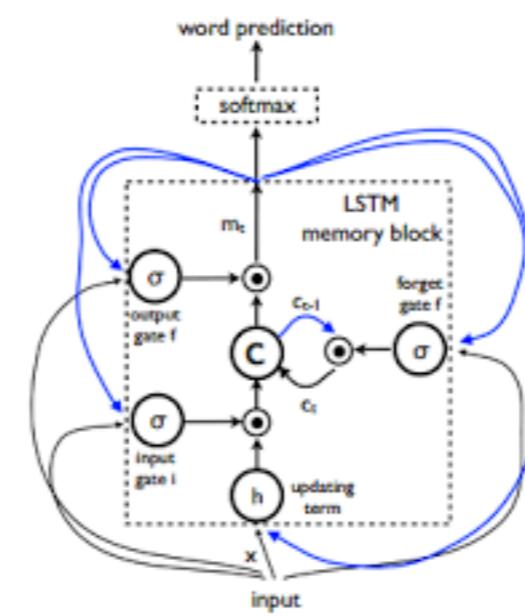
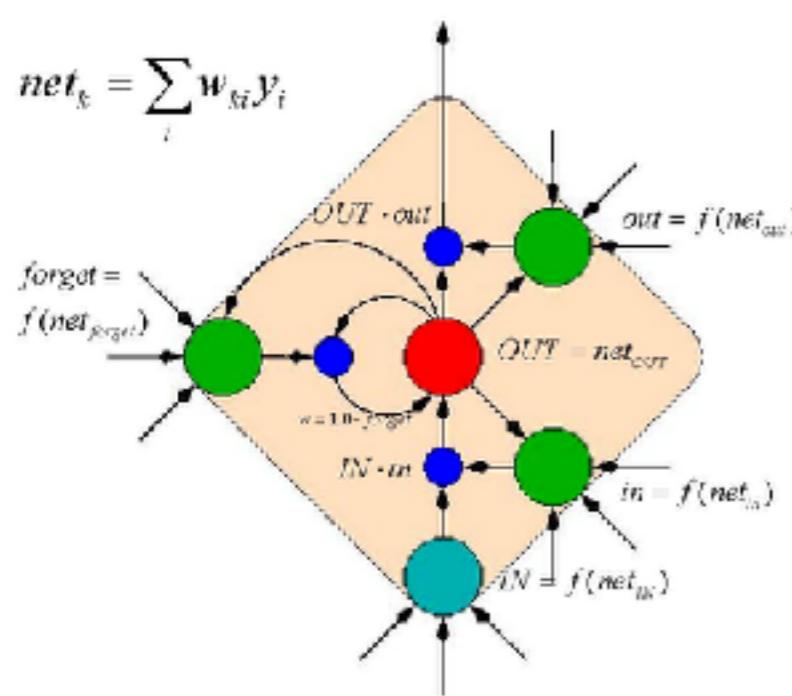
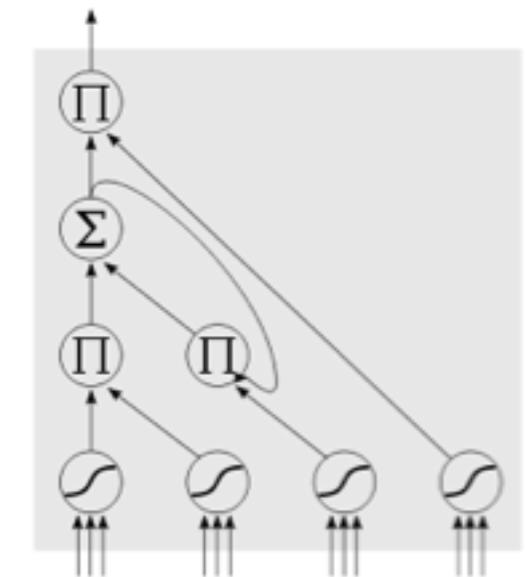
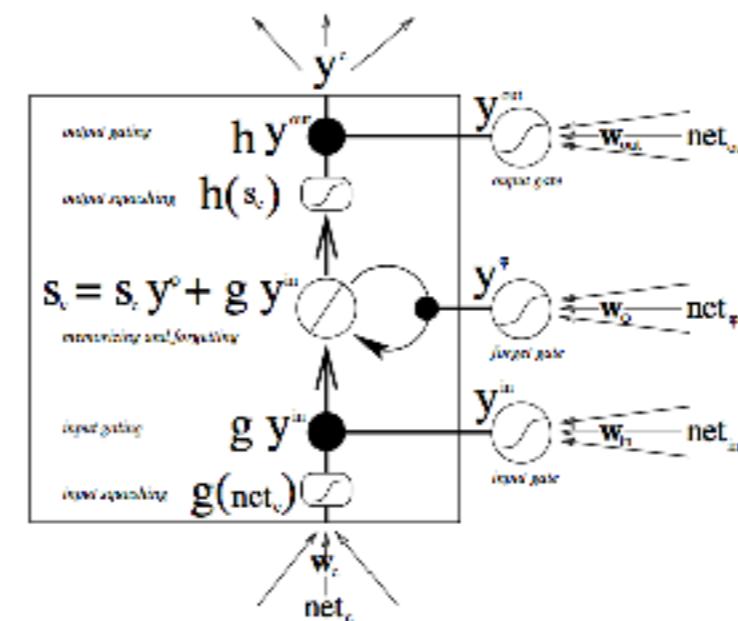
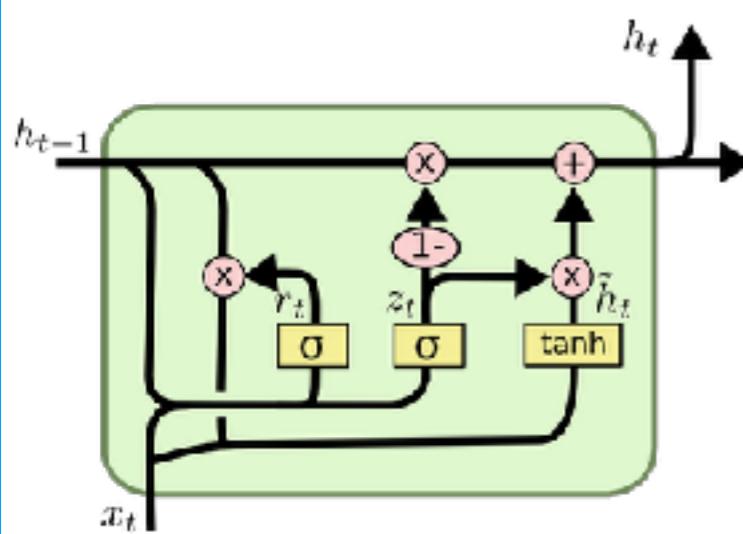
$$f = \sigma(x_t U^f + s_{t-1} W^f)$$

$$o = \sigma(x_t U^o + s_{t-1} W^o)$$

$$c_t = c_{t-1} \circ f + g \circ i$$

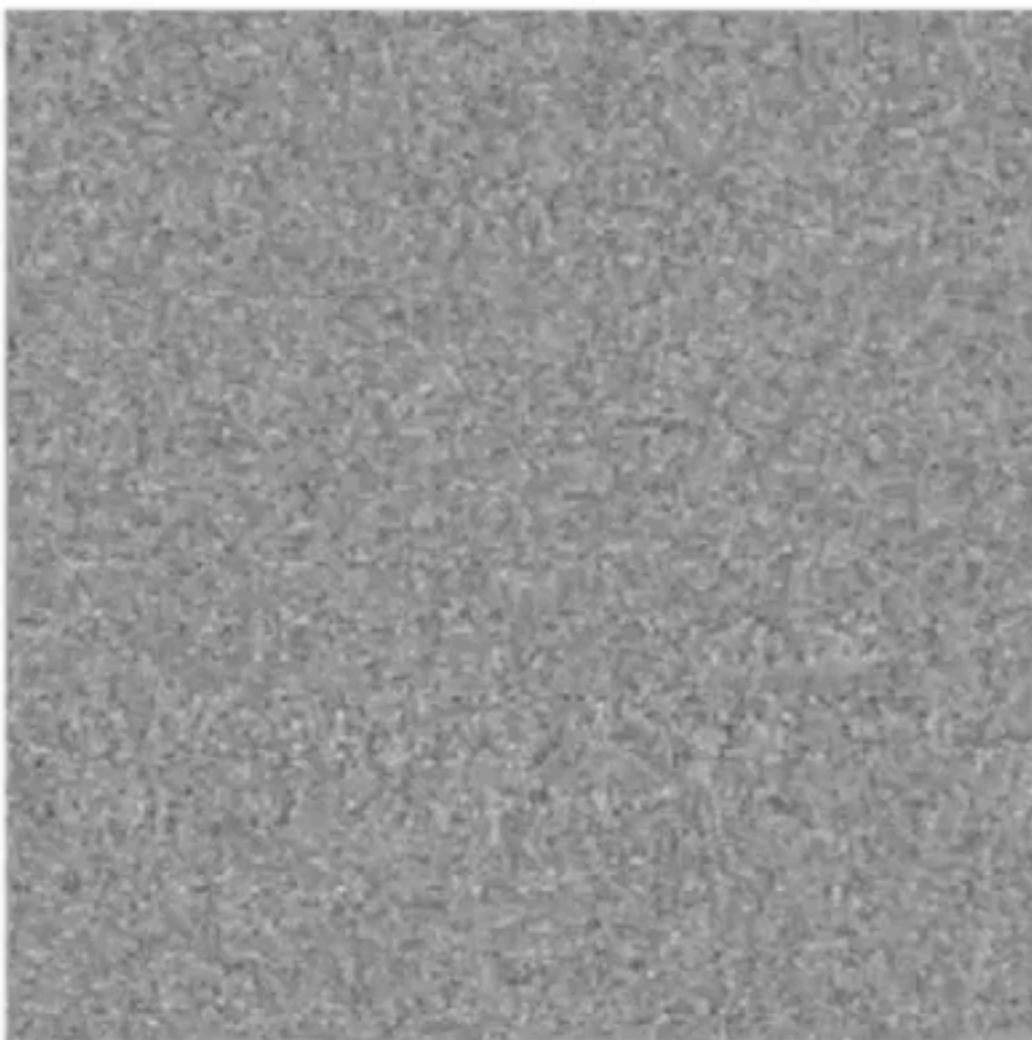
$$s_t = \tanh(c_t) \circ o$$

LSTM Unit

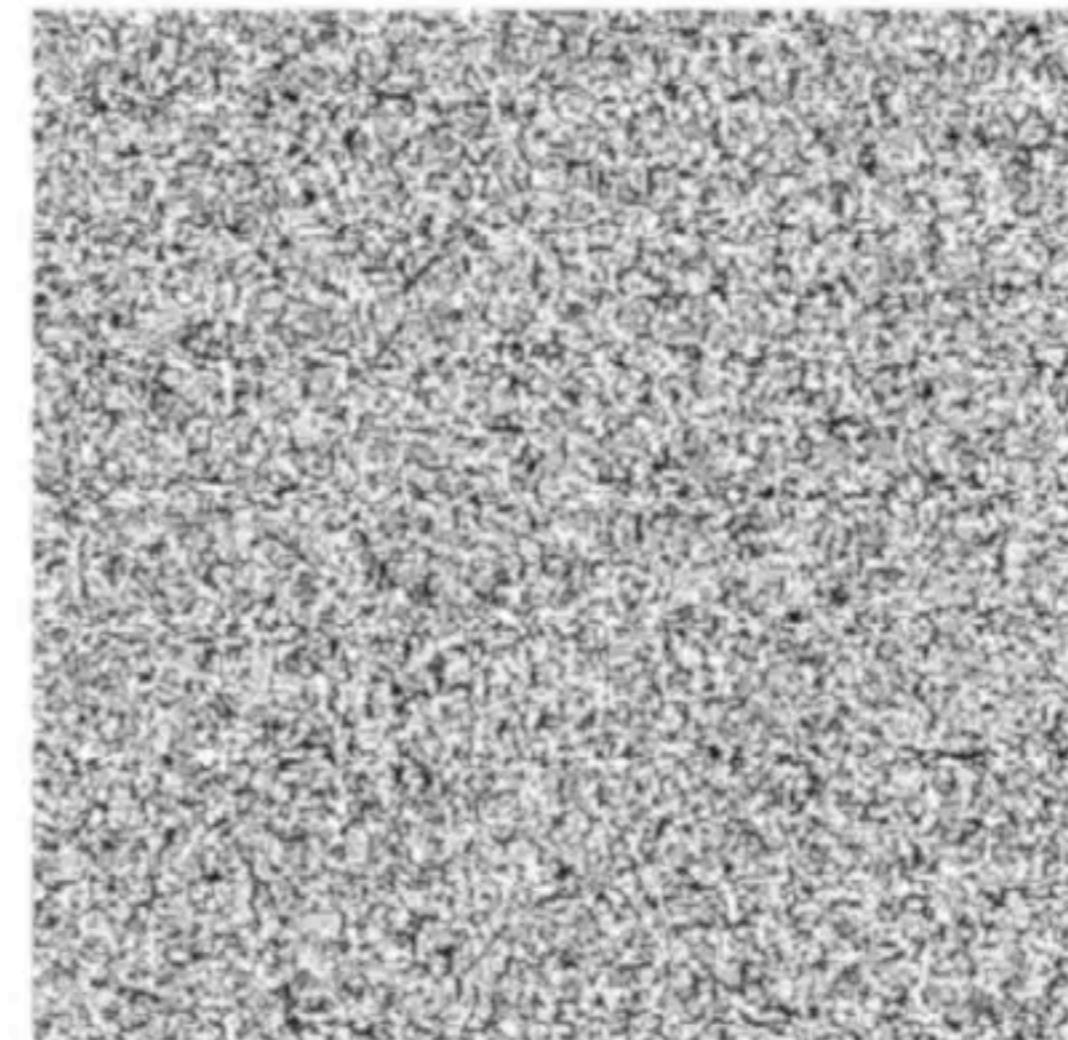


LSTM vs. RNN

127



127



<https://imgur.com/gallery/vaNahKE>

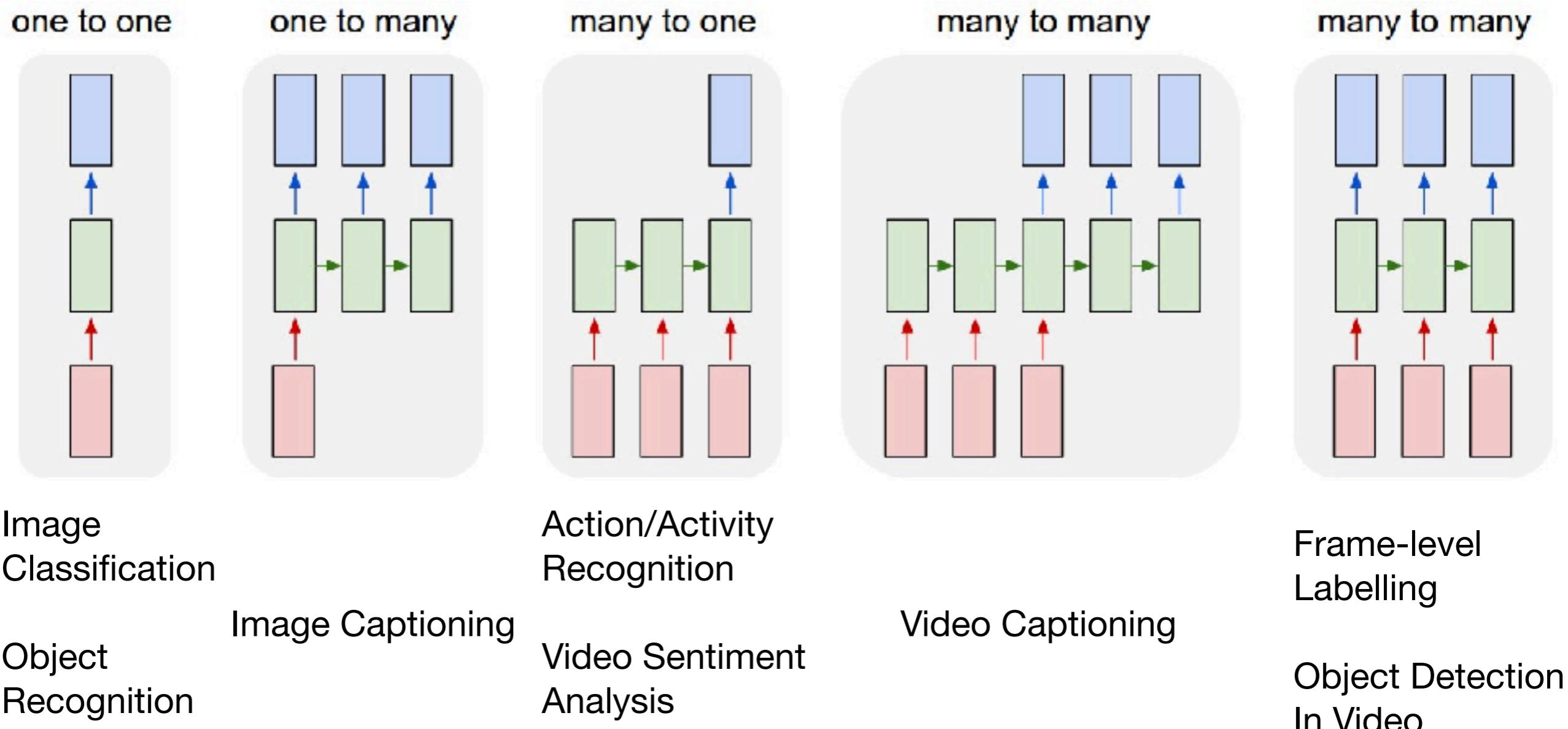
LSTM vs GRU

- 3 Gates
- Contains additional memory for each cell
- Two Non-Linearities in computing output
- More expressive power
 - Better results with more data
- 2 Gates
- No memory beyond cell state
- No second Non-Linearity
- Slightly lesser parameters
 - Quicker to Train
 - Smaller datasets

LSTM Training Advice

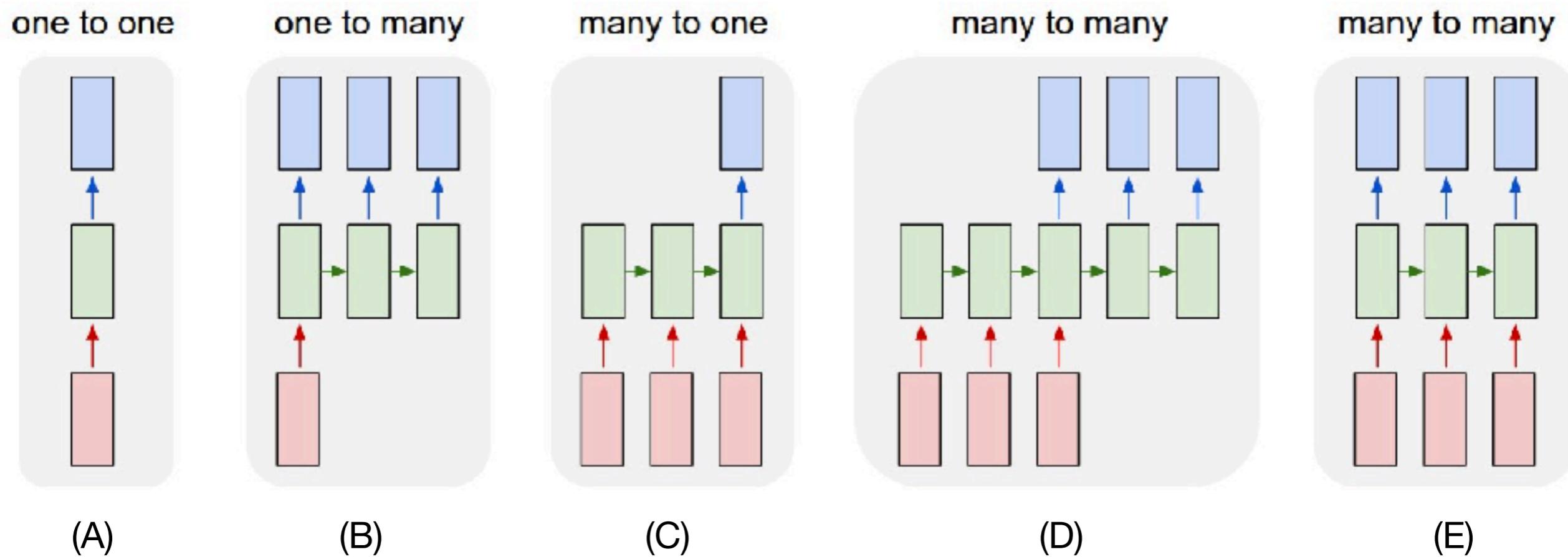
- Use an LSTM or GRU: it makes your life so much simpler!
- Initialize recurrent matrices to be orthogonal
- Initialize other matrices with a sensible (small!) scale
- Initialize forget gate bias to 1: default to remembering
- Use adaptive learning rate algorithms: Adam, AdaDelta, ...
- Clip the norm of the gradient: 1–5 seems to be a reasonable threshold when used together with Adam or AdaDelta.
- Either only dropout vertically or learn how to do it right
- Be patient!

RNN Architectures



Picture courtesy: Andrej Karpathy

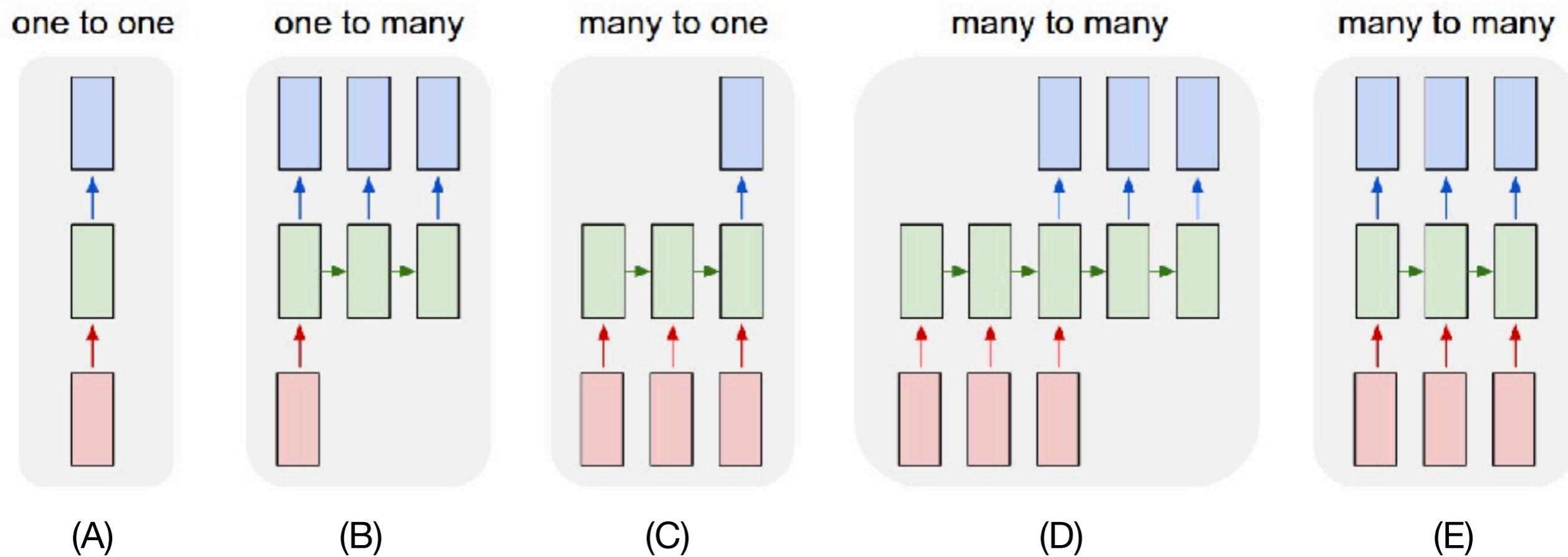
Beyond Vision



Speech Recognition
Speaker Recognition

Sentiment from a Tweet
Translate sentence from English to Telugu

Beyond Vision



Text to Speech (Alexa)
Answering a Question

Controlling a Self-Driving Car
Generate Trailer for a Movie

Modeling Text

- Given a partial sentence, predict the next word

My name is _____

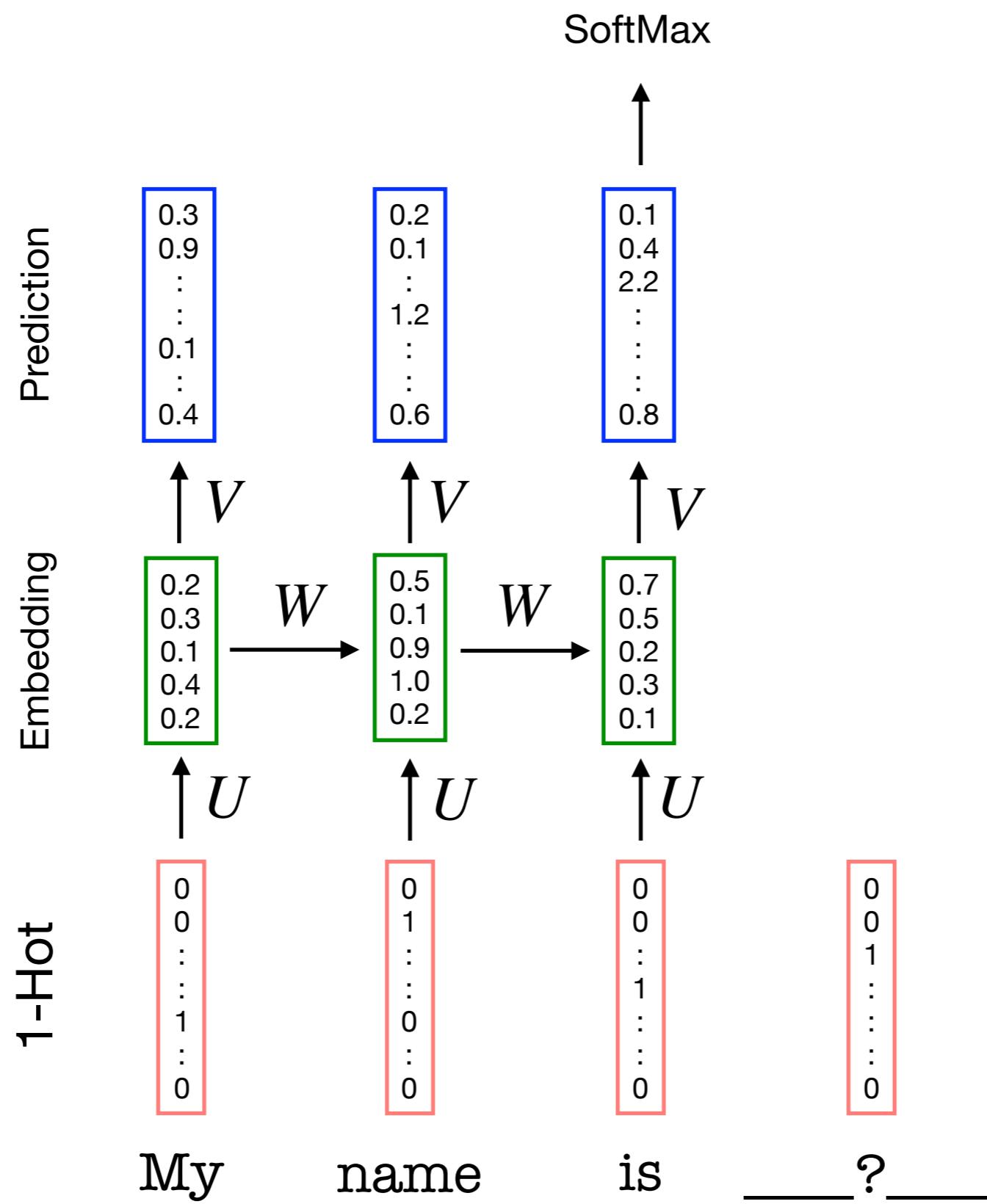
My name is Anthony _____

My name is Anthony Gonzalves _____

My name is Anthony Gonzalves main _____

My name is Anthony Gonzalves main duniya _____

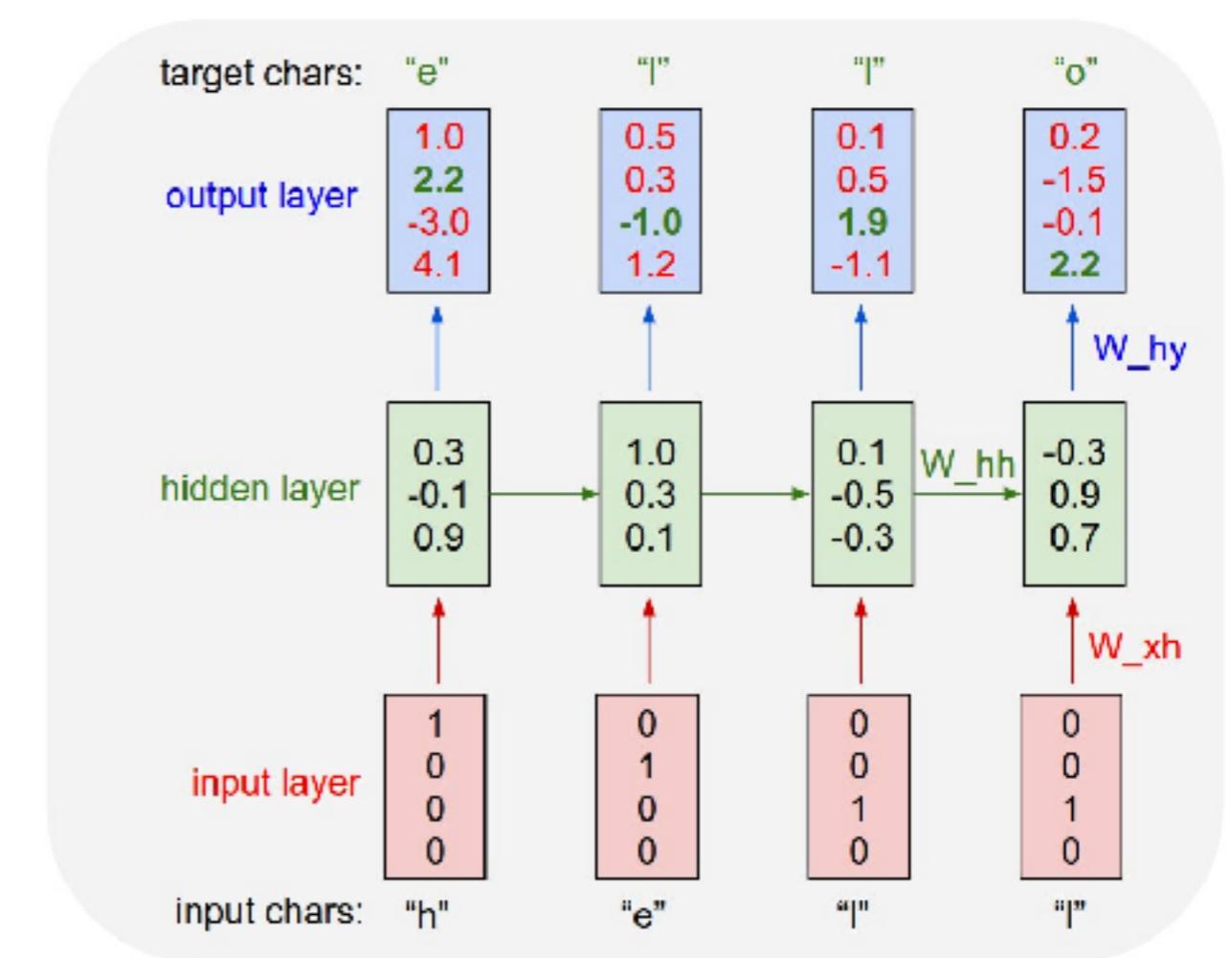
My name is Anthony Gonzalves main duniya mein _____



Char-RNN

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

Easier at Character level?



all hail andrej karpathy

The Shakespeare Hallucination

at first:

tyntd-iafhatawiaoihrdemot 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 shuwy 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.

The Shakespeare Hallucination

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.

Generated Wikipedia Text

Naturalism and decision for the majority of Arab countries' capitalide was grounded by the Irish language by [[John Clair]], [[An Imperial Japanese Revolt]], associated with Guangzham's sovereignty. His generals were the powerful ruler of the Portugal in the [[Protestant Immineners]], which could be said to be directly in Cantonese Communication, which followed a ceremony and set inspired prison, training. The emperor travelled back to [[Antioch, Perth, October 25|21]] to note, the Kingdom of Costa Rica, unsuccessful fashioned the [[Thrales]], [[Cynth's Dajoard]], known in western [[Scotland]], near Italy to the conquest of India with the conflict. Copyright was the succession of independence in the slop of Syrian influence that was a famous German movement based on a more popular servicious, non-doctrinal and sexual power post. Many governments recognize the military housing of the [[Civil Liberalization and Infantry Resolution 265 National Party in Hungary]], that is sympathetic to be to the [[Punjab Resolution]] (PJS) [<http://www.humah.yahoo.com/guardian.cfm/7754800786d17551963s89.htm>]

Official economics Adjoint for the Nazism, Montgomery was swear to advance to the resources for those Socialism's rule, was starting to signing a major tripad of aid exile.]]

Generated C Code

```
/*
 * Increment the size file of the new incorrect UI_FILTER group information
 * of the size generatively.
 */
static int indicate_policy(void)
{
    int error;
    if (fd == MARN_EPT) {
        /*
         * The kernel blank will coeld it to userspace.
         */
        if (sa->segment < mem_total)
            unblock_graph_and_set_blocked();
        else
            ret = 1;
        goto bail;
    }
    segaddr = in_SB(in.addr);
    selector = seg / 16;
    setup_works = true;
    for (i = 0; i < blocks; i++) {
        seq = buf[i++];
        bpf = bd->bd.next + i * search;
        if (fd) {
            current = blocked;
        }
    }
    rw->name = "Getjbbregs";
    bprm_self_clear1(&iv->version);
    regs->new = blocks[(BPF_STATS << info->historidac)] | PFMR_CLOBATHINC_SECONDS << 12;
    return segtable;
}
```

Generated Algebraic Geometry

For $\bigoplus_{n=1, \dots, m} \mathcal{L}_{m,n} = 0$, 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 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 $x, x' \in S'$ 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 X' , and \mathcal{T}_i is an object of $\mathcal{F}_{X/S}$ for $i > 0$ and \mathcal{F}_0 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

$$\widehat{M}^* = \mathcal{I}^* \otimes_{\text{Spec}(S)} \mathcal{O}_{S,s} - i_X^* \mathcal{F}$$

is a unique morphism of algebraic stacks. Note that

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

and

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

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_{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 $X = \lim X_i$ (by the formal open covering X and a single map $\text{Proj}_X(\mathcal{A}) = \text{Spec}(B)$ 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 $\mathcal{J}_{n,0} \circ \mathcal{A}_2$ 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 δ_{n+1} is a scheme over S . \square

Interpreting LSTM Cells

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.

Length Tracker

"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."

Inside/Outside Quotes

Interpreting LSTM 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;
}
```

Inside/Outside IF

```
#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;
}
```

Depth of Expression

Interpreting LSTM Cells

```
/* 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.
     */
}
```

Uninterpretable

“Deep” RNNs

RNN:

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

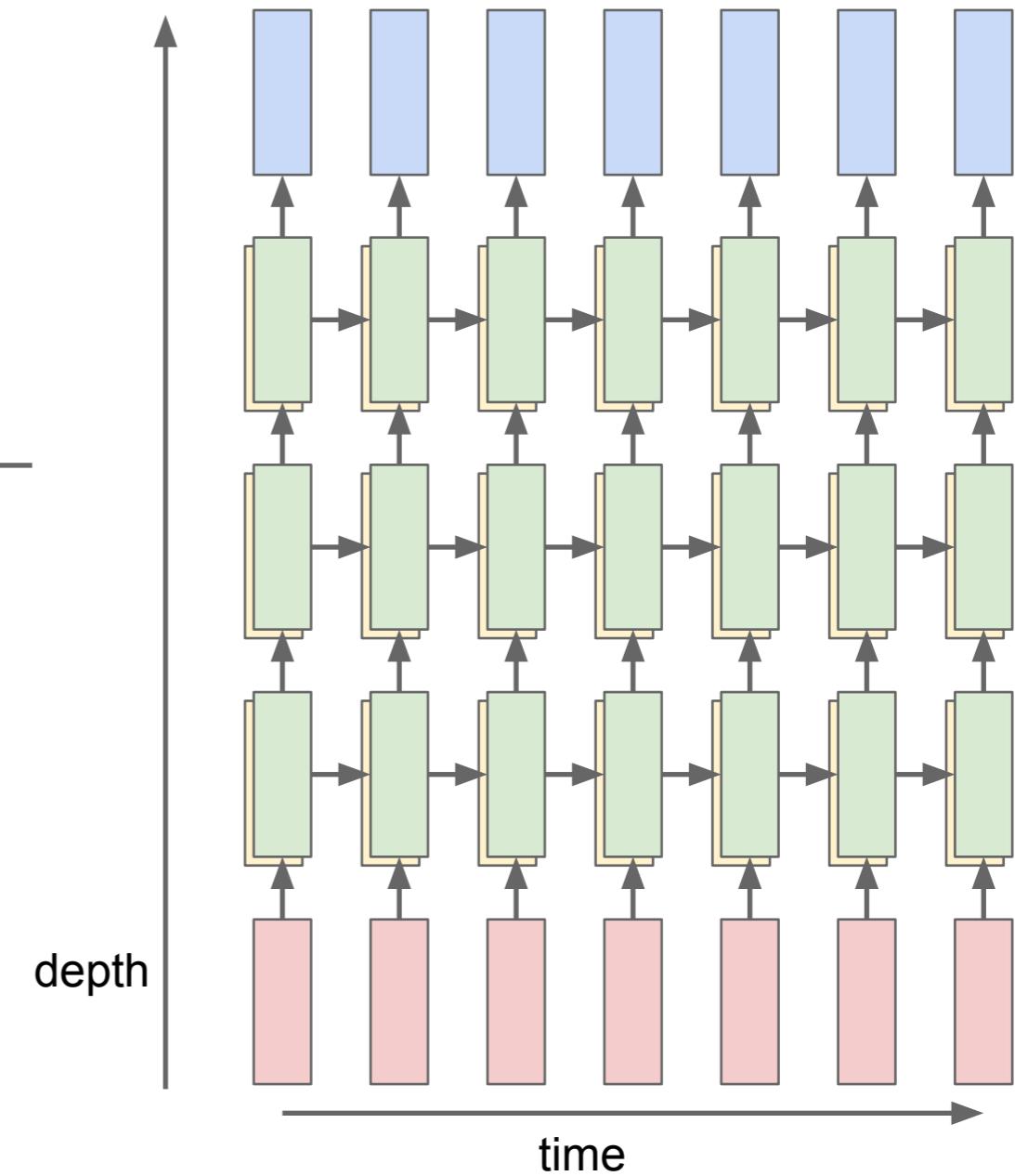
$h \in \mathbb{R}^n$ $W^l [n \times 2n]$

LSTM:

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

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

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

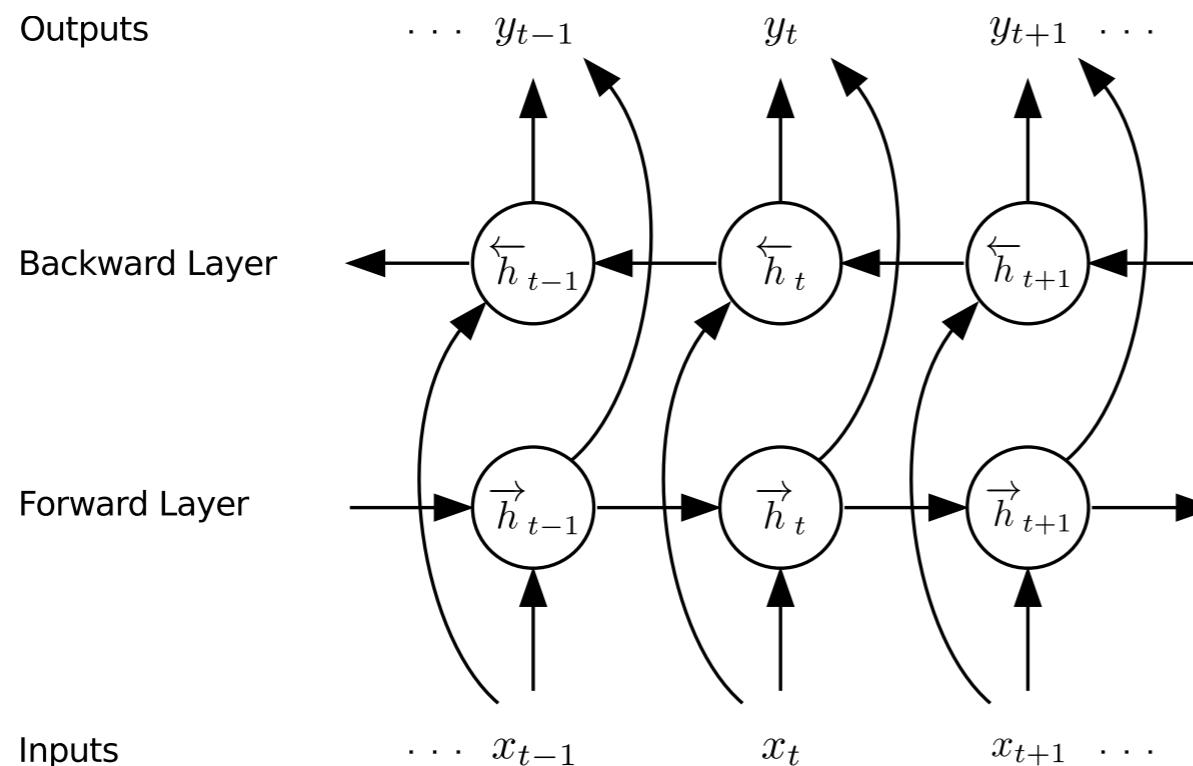


Note: notations are different

all hail andrej karpathy

Other Fancy Architectures

- BRNNs



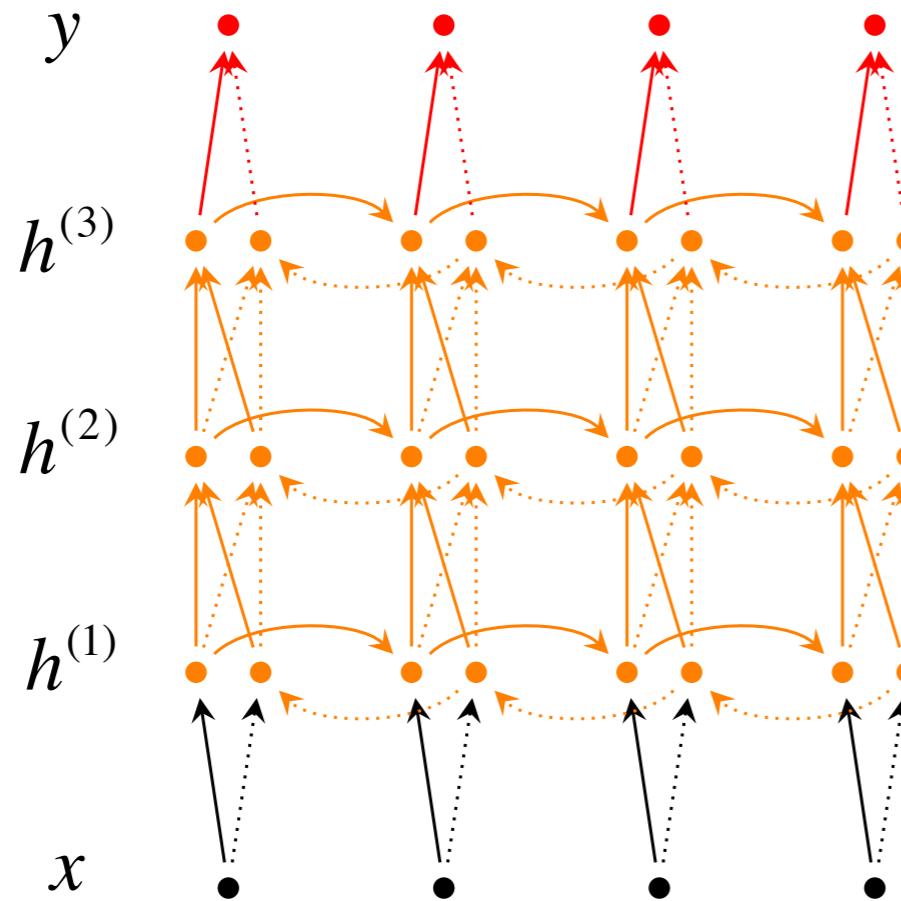
$$\begin{aligned}\vec{h}_t &= f(\vec{W}x_t + \vec{V}\vec{h}_{t-1} + \vec{b}) \\ \overleftarrow{h}_t &= f(\overleftarrow{W}x_t + \overleftarrow{V}\overleftarrow{h}_{t+1} + \overleftarrow{b}) \\ y_t &= g(U[\vec{h}_t; \overleftarrow{h}_t] + c)\end{aligned}$$

Note: notations are different

A. Graves and J. Schmidhuber, “Framewise Phoneme Classification with Bidirectional LSTM and Other Neural Network Architectures,” *Neural Networks*, June/July 2005.

Other Fancy Architectures

- Deep BRNNs



$$\vec{h}_t^{(i)} = f(\vec{W}^{(i)} \vec{h}_t^{(i-1)} + \vec{V}^{(i)} \vec{h}_{t-1} + \vec{b}^{(i)})$$

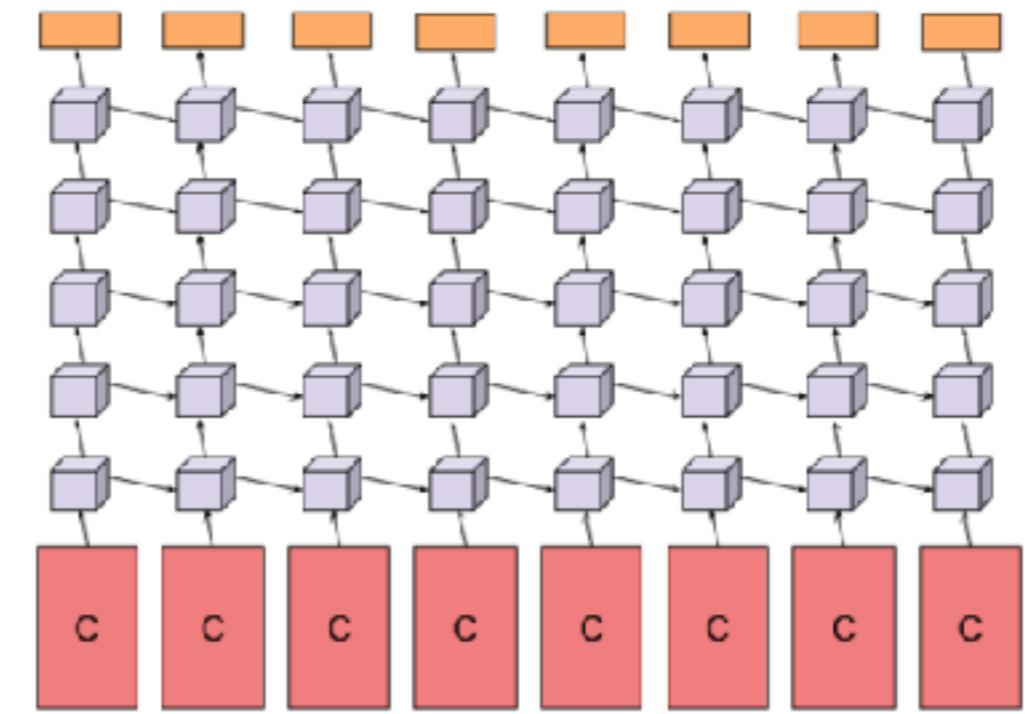
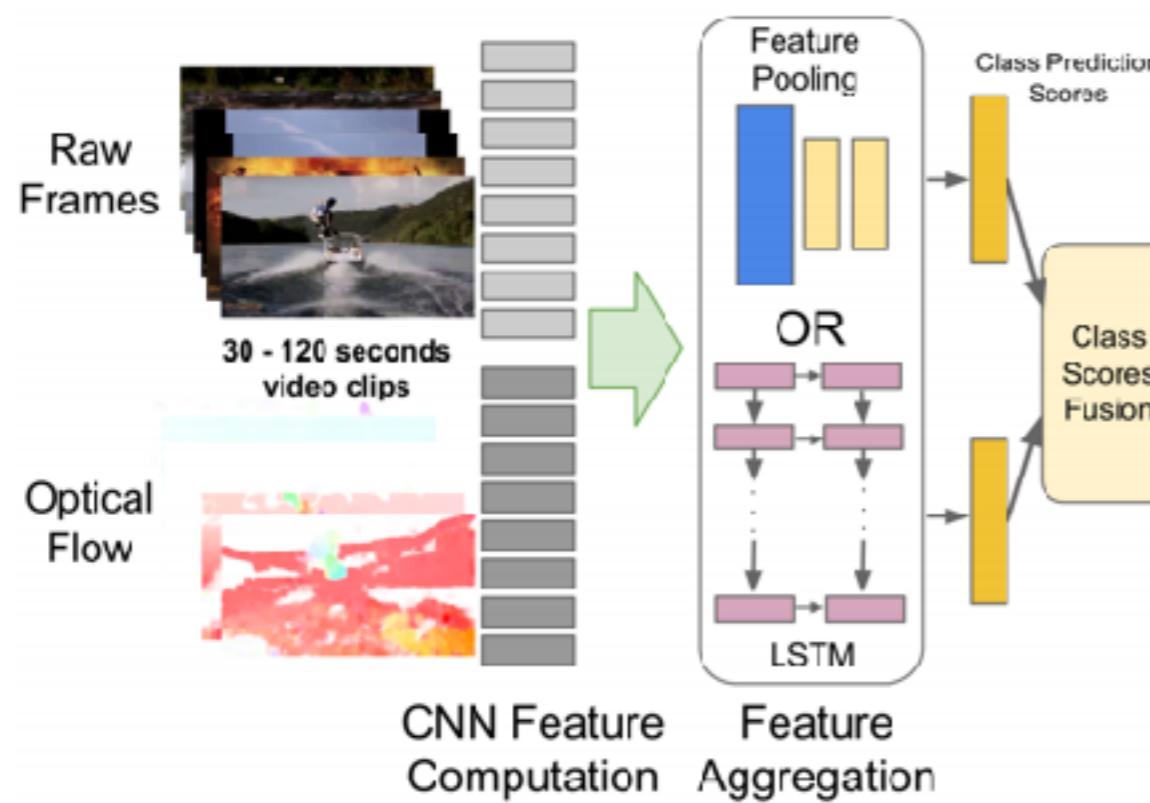
$$\overset{\leftarrow}{h}_t^{(i)} = f(\overset{\leftarrow}{W}^{(i)} \overset{\leftarrow}{h}_t^{(i-1)} + \overset{\leftarrow}{V}^{(i)} \overset{\leftarrow}{h}_{t+1} + \overset{\leftarrow}{b}^{(i)})$$

$$y_t = g(U[\vec{h}_t^{(L)}; \overset{\leftarrow}{h}_t^{(L)}] + c)$$

Note: notations are different

Picture courtesy: Richard Socher

Other Fancy Architectures



Ng et al., Beyond Short Snippets: Deep Networks for Video Classification, CVPR 2015

Next Session



"man in black shirt is playing guitar."

Image Captioning



What is the mustache made of?

Visual Question Answering

Summary

- RNNs & LSTMs are Awesome
 - Took us 20 years to realise it!
 - Remember those LSTM equations
- Can you design your own architectures?
 - And Cells?

Online References

- <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
- <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>
- <http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/>

References

LSTMs

- Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, November 1997
- Gref et al., LSTM: A Search Space Odyssey, *IEEE Trans. Neural Networks and Learning Systems*, Oct. 2017.
- Pascanu et al., On the difficulty of training recurrent neural networks, *ICML* 2013

References

GRUs

- Cho et al., Learning Phrase Representations using RNN Encoder Decoder for Statistical Machine Translation. *arXiv:1406.1078*, 2014

RNNs

- Bengio et al., Advances in optimizing recurrent networks, IEEE ICASSP, 2013

CNNs for Video

- Simonyan & Zisserman, Two-Stream Convolutional Networks for Action Recognition in Videos, NIPS 2014

Thank You

Questions?