

CS 466/566

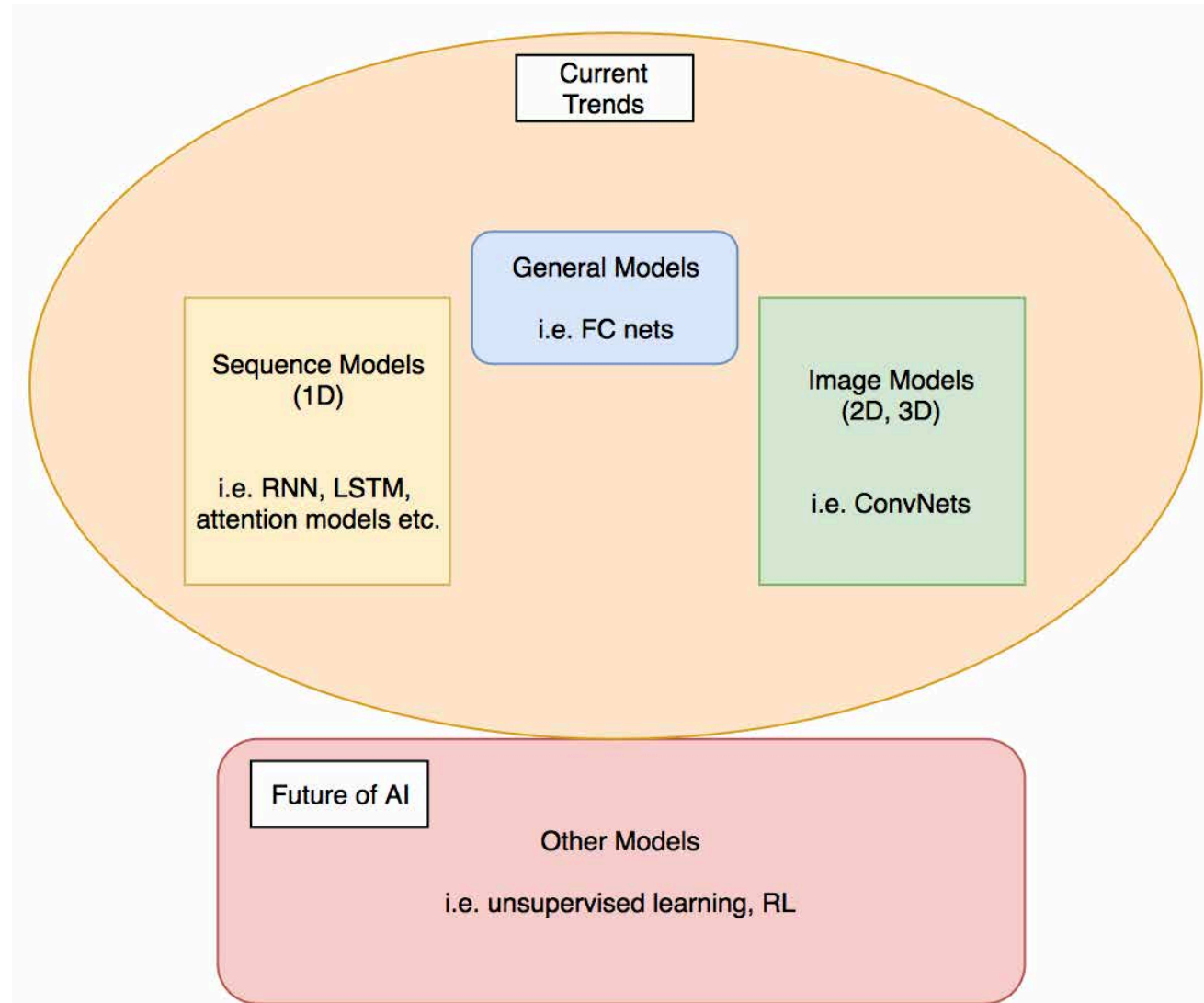
Introduction to Deep Learning

Lecture 11 – Recurrent Neural Networks – Part 1

Recall: ConvNet Demo

<http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

What were we doing? Where are we?



Sequence Modelling

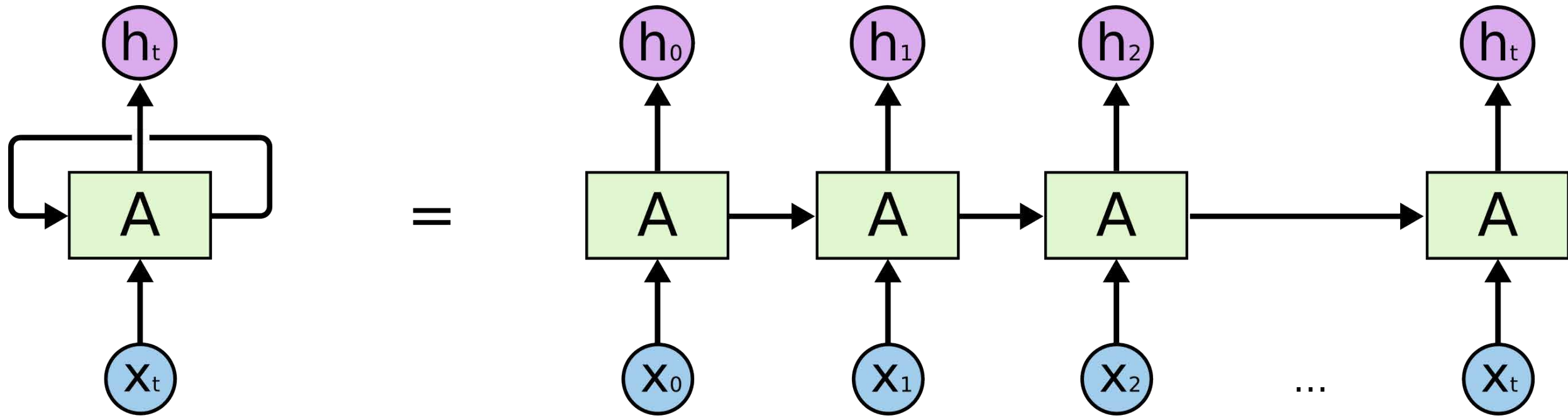
- Humans don't start their thinking from scratch every now and then.
 - You understand each word based on your understanding of previous words.
 - We don't throw everything away and restart thinking from scratch. Thoughts have persistence.
- Traditional neural networks approach can't do this, and it is a major shortcoming.
 - Imagine classifying what kind of event is happening at every point in a movie.
 - It's unclear how a traditional neural network could use its reasoning about previous events in the film to inform later ones.

Sequence Modelling

- Humans don't start their thinking from scratch every now and then.
 - You understand each word based on your understanding of previous words.
 - We don't throw everything away and restart thinking from scratch. Thoughts have persistence.
- Traditional neural networks approach can't do this, and it is a major shortcoming.
 - Imagine classifying what kind of event is happening at every point in a movie.
 - It's unclear how a traditional neural network could use its reasoning about previous events in the film to inform later ones.

Recurrent Neural Networks (RNNs)

- RNNs solve this issue. They are networks with loops in them, allowing information to persist.



Recurrent Neural Networks (RNNs)

- Limitation of Traditional NNs (including CNNs) is that they are too constrained:
 - accept a fixed-sized vector as input (e.g. an image)
 - produce a fixed-sized vector as output (e.g. probabilities of different classes).
 - perform this mapping by using a fixed amount of computational steps (e.g. the number of layers in the model).
- The core reason that recurrent nets are more exciting is that they allow us to operate over *sequences* of vectors: Sequences in the input, the output, or in the most general case both.

Recurrent Neural Networks (RNNs)

one to one

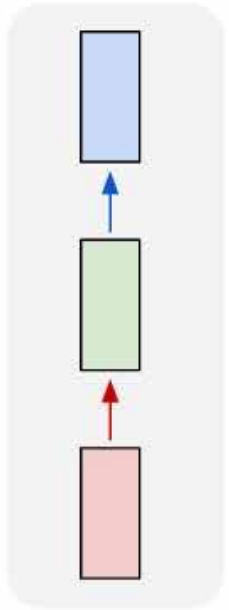


Image
Classification

one to many

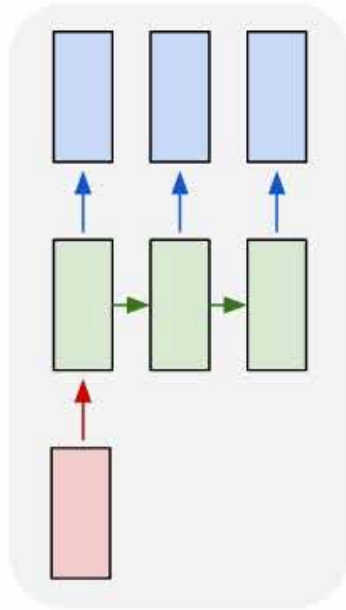
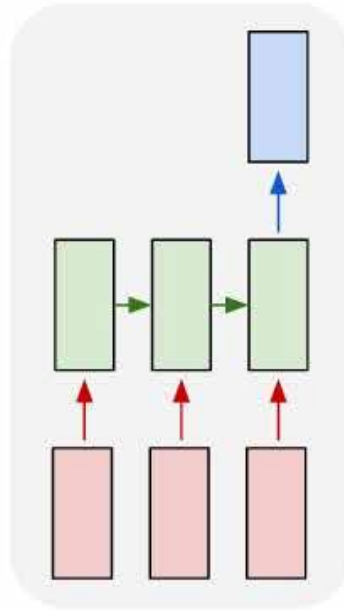


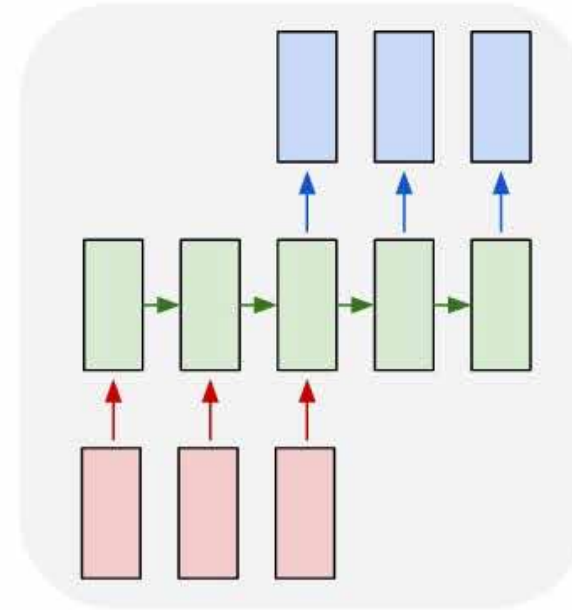
Image
Captioning

many to one



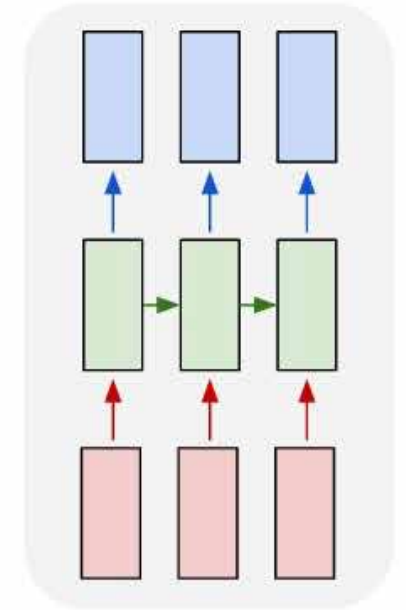
Sentiment
Analysis

many to many



Google
Translate

many to many



Video Frame
Classification

red: input, blue: output, green: RNN state

Recurrent Neural Networks (RNNs)

- The sequential operation is much more powerful compared to fixed networks.
- Moreover, RNNs combine the input vector with their state vector with a fixed (but learned) function to produce a new state vector.
- This can in programming terms be interpreted as running a fixed program with certain inputs and some internal variables.
- Viewed this way, RNNs are essentially similar to programs. They can theoretically simulate programs.

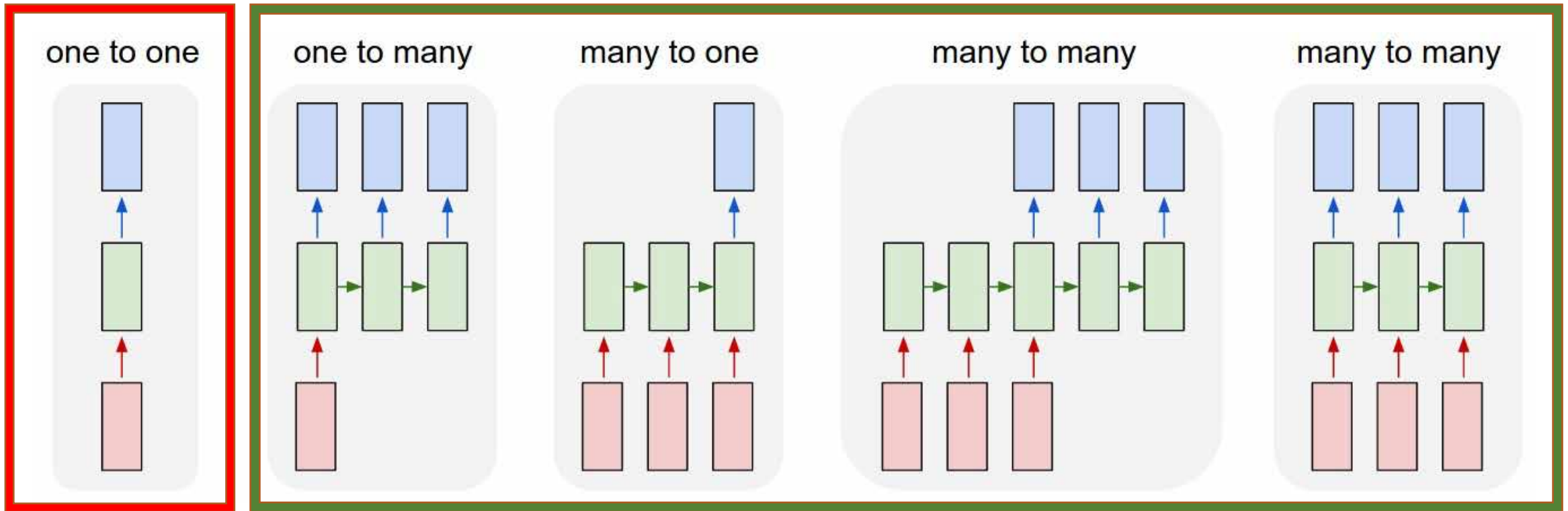
Recurrent Neural Networks (RNNs)

If training vanilla neural nets is optimization over functions, training recurrent nets is optimization over programs.

Recurrent Neural Networks (RNNs)

**Optimization
over functions**

Optimization over programs



Recurrent Neural Network Handwriting Generation Demo

- <http://www.cs.toronto.edu/~graves/handwriting.html>

Four effective ways to learn an RNN

- **Long Short Term Memory (LSTM)**

Make the RNN out of little modules that are designed to remember values for a long time.

- **Hessian Free Optimization:**

Deal with the vanishing gradients problem by using a fancy optimizer that can detect directions with a tiny gradient but even smaller curvature.

- The HF optimizer (Martens & Sutskever, 2011) is good at this.

- **Echo State Networks:**

Initialize the input→hidden and hidden→hidden and output→hidden connections very carefully so that the hidden state has a huge reservoir of weakly coupled oscillators which can be selectively driven by the input.

- ESNs only need to learn the hidden→output connections.

- **Good initialization with momentum**

Initialize like in Echo State Networks, but then learn all of the connections using momentum.

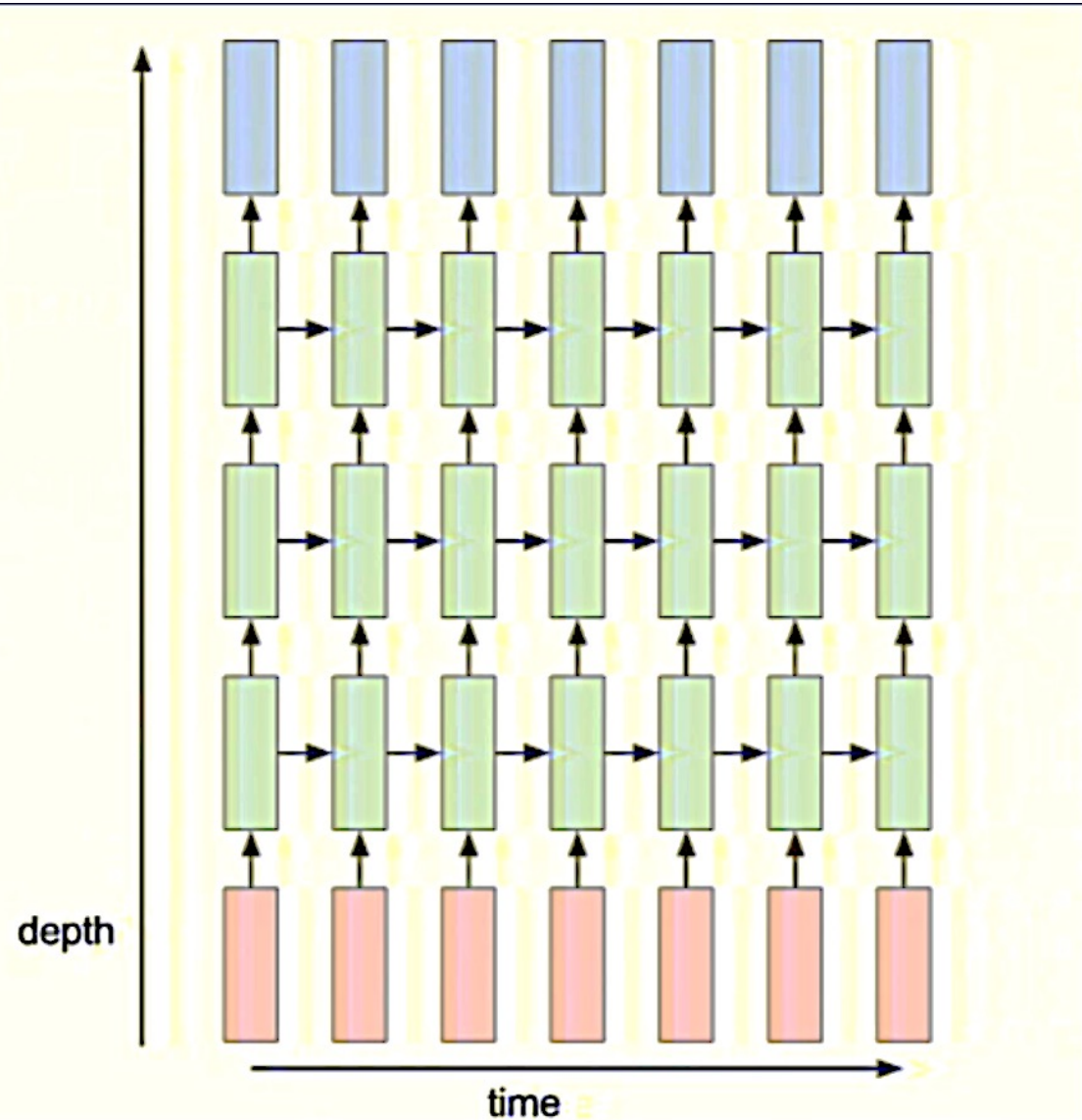
Deep Recurrent Neural Networks (RNNs)

RNN:

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

$h \in \mathbb{R}^n$

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



Deep Recurrent Neural Networks (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:

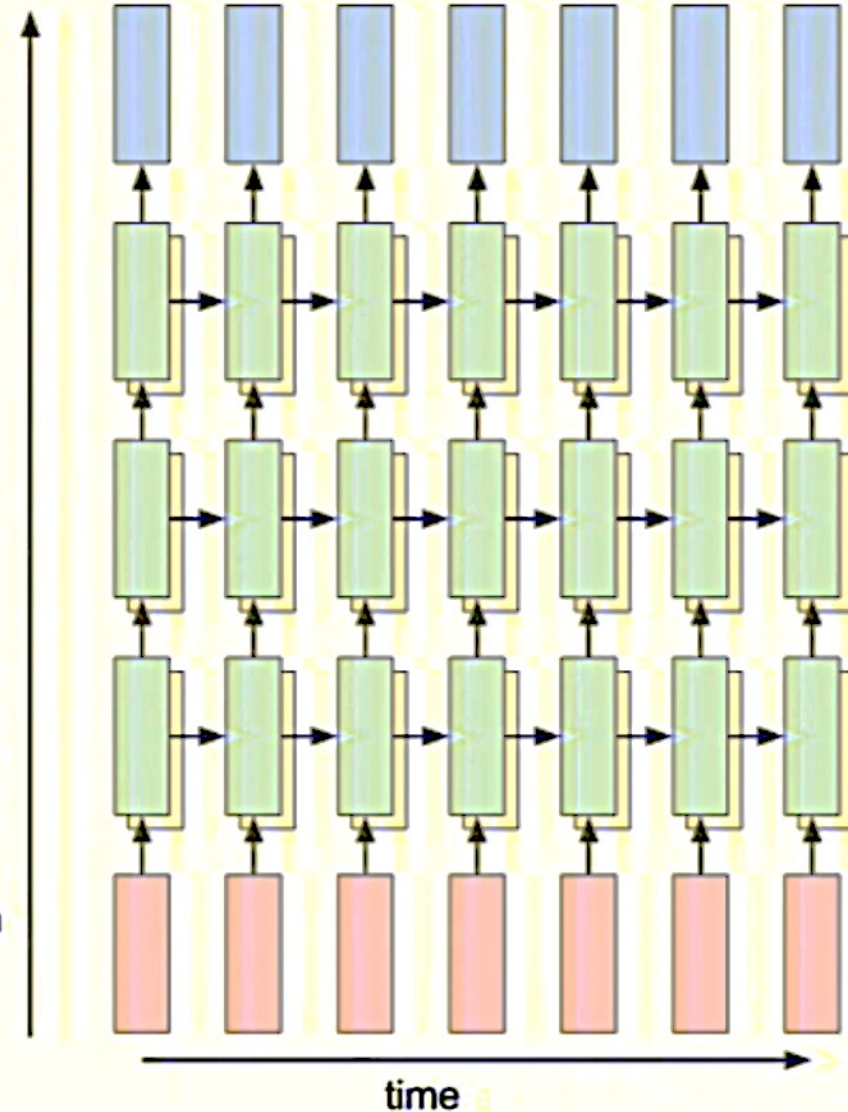
$W^l [4n \times 2n]$

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

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

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

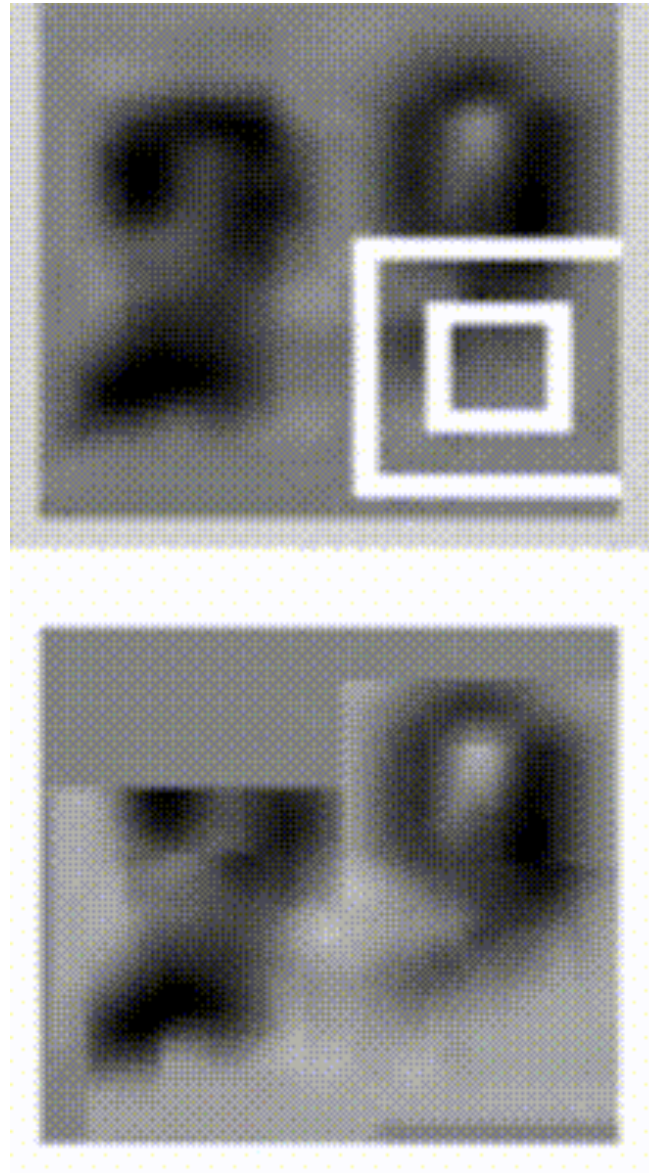
depth



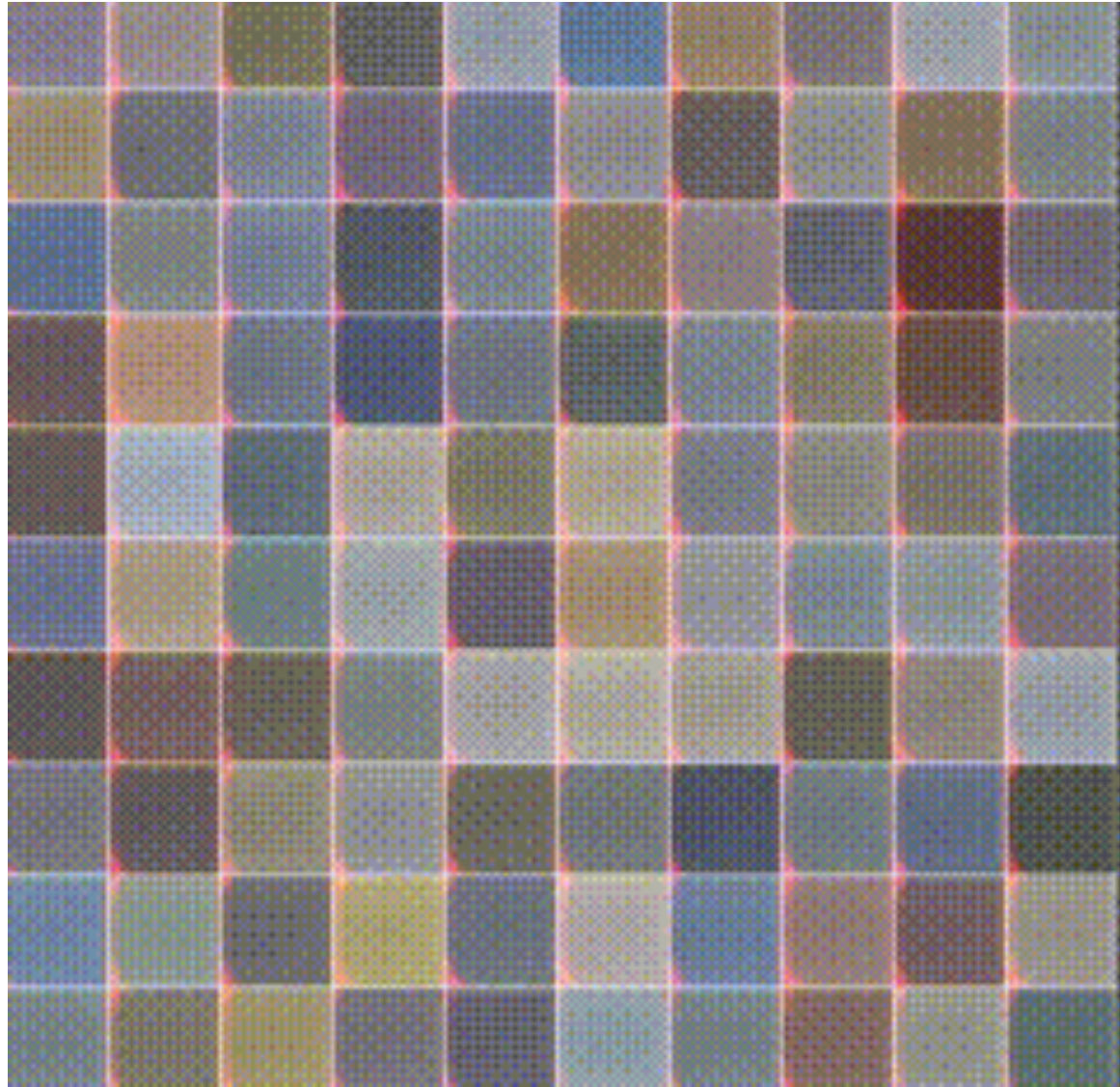
Finding sequence in absence of sequence

- Even if your data is not in form of sequences, you can still formulate and train powerful models that learn to process it sequentially.
- You're actually learning stateful programs that process your fixed-sized data.

Finding sequence in absence of sequence

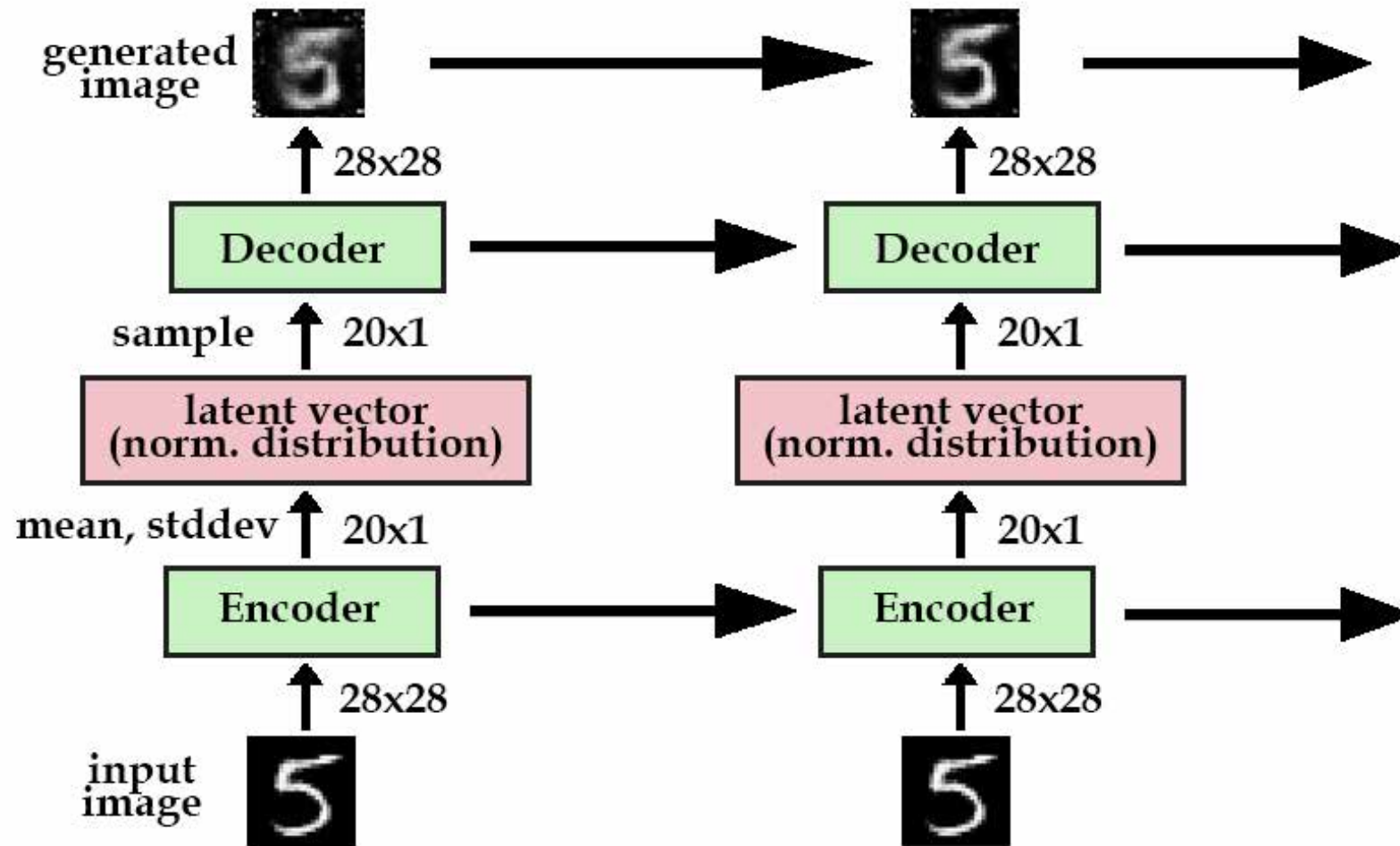


Finding sequence in absence of sequence

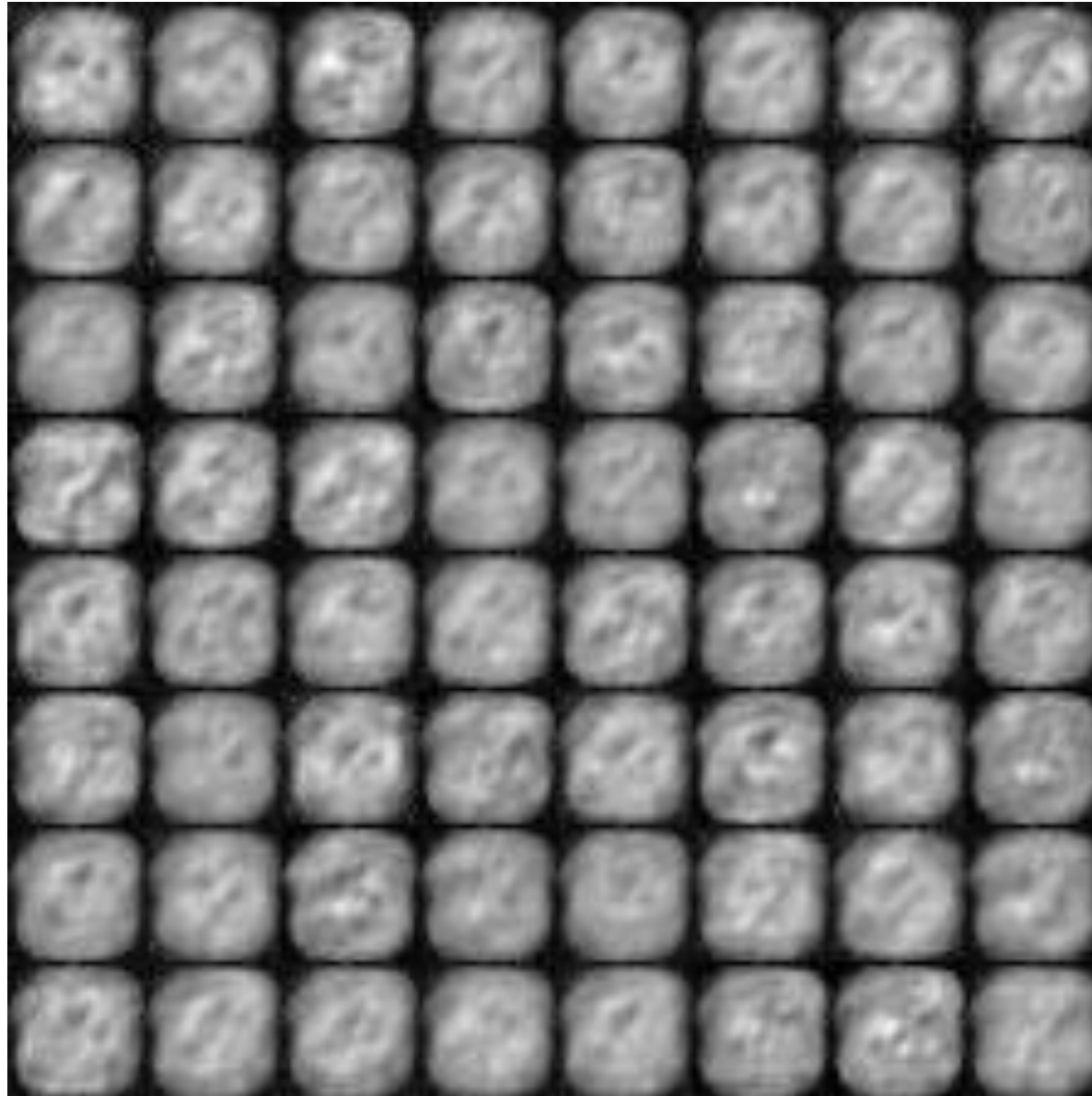


Teaser #1: DRAW: Deep Recurrent Attentive Writer

- It is a recurrent version of Variational Autoencoder!

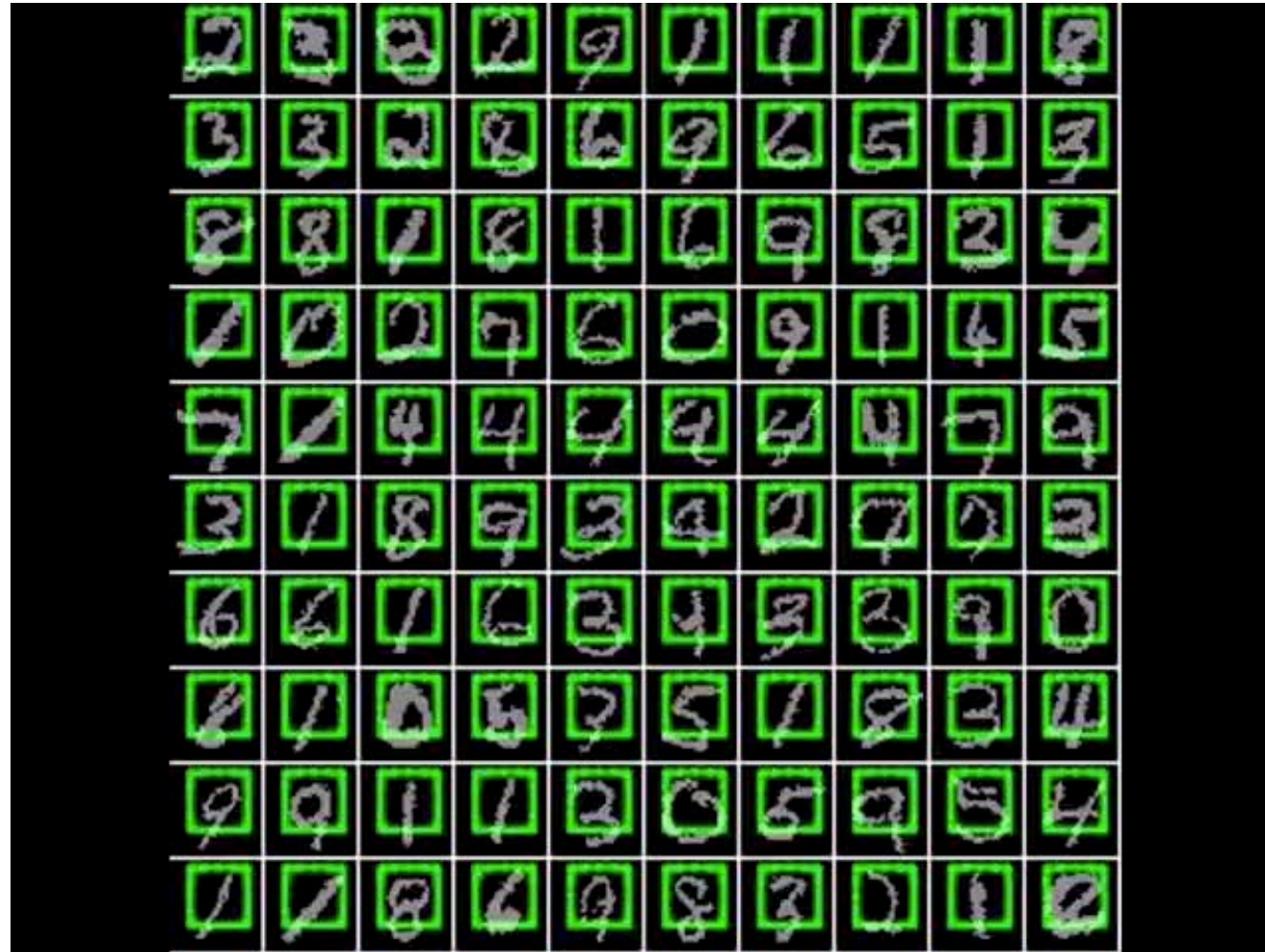


DRAW: Deep Recurrent Attentive Writer

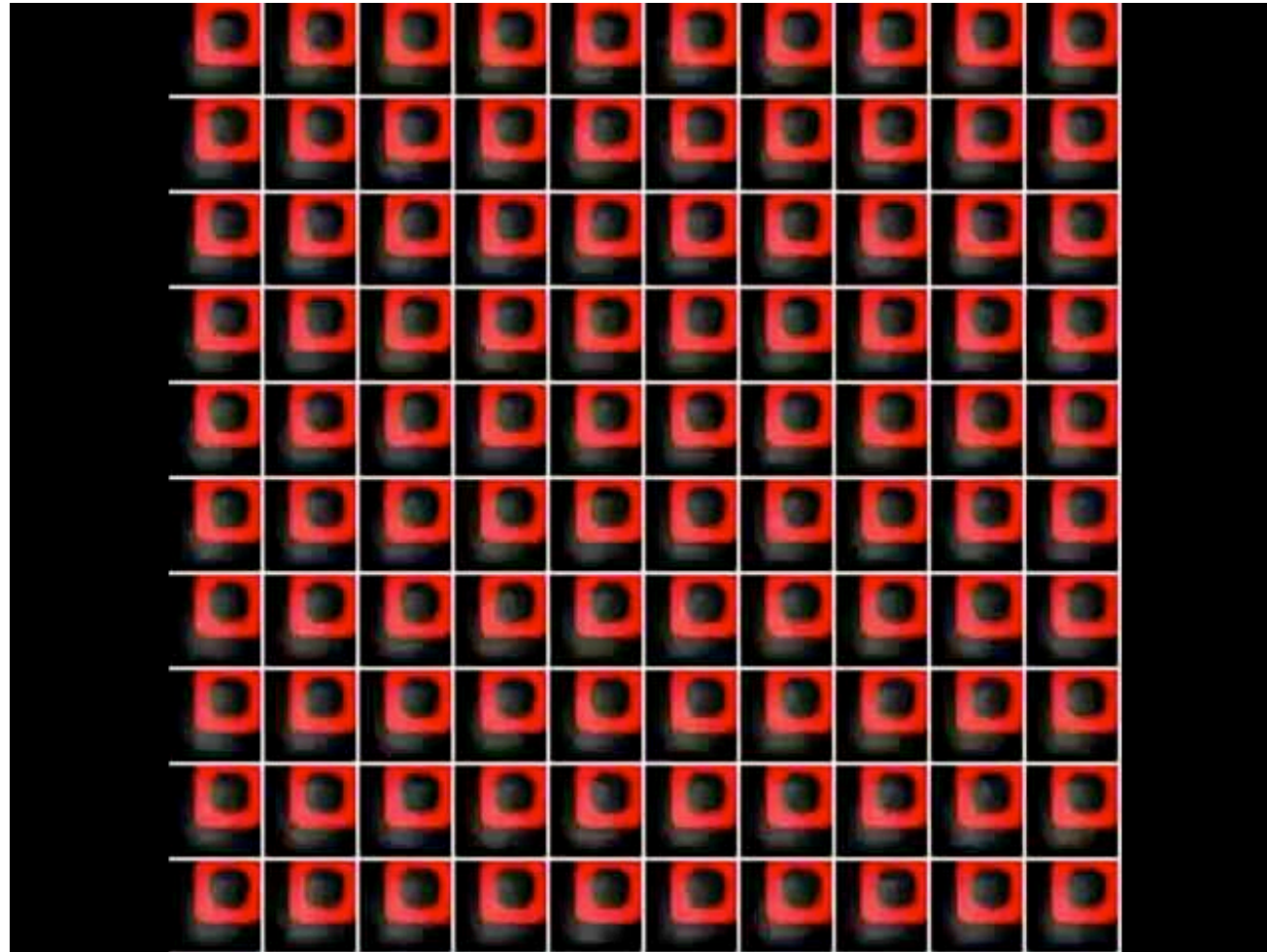


DRAW: Deep Recurrent Attentive Writer

Attention Gate and Encoder



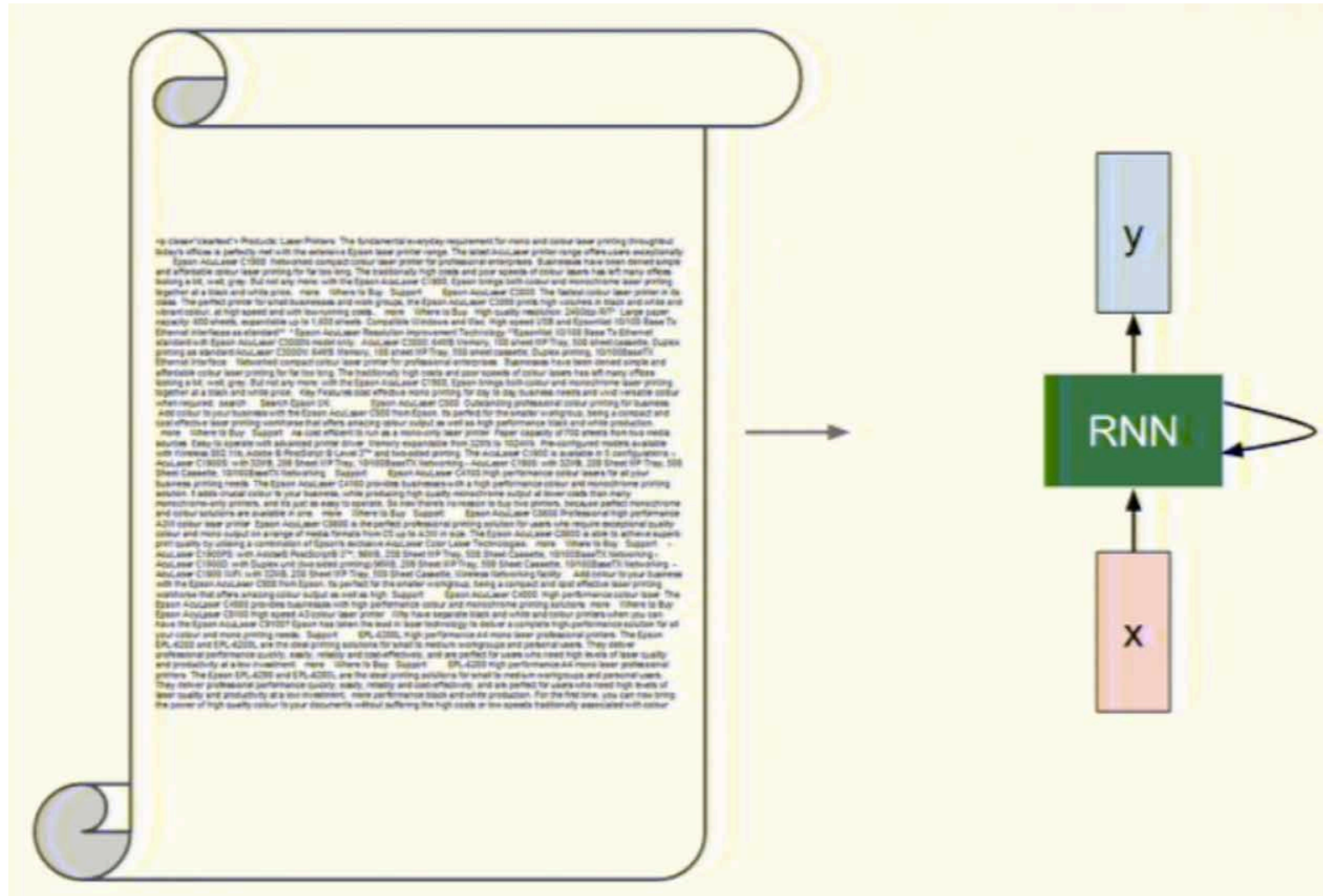
DRAW: Deep Recurrent Attentive Writer Attention Gate and Decoder



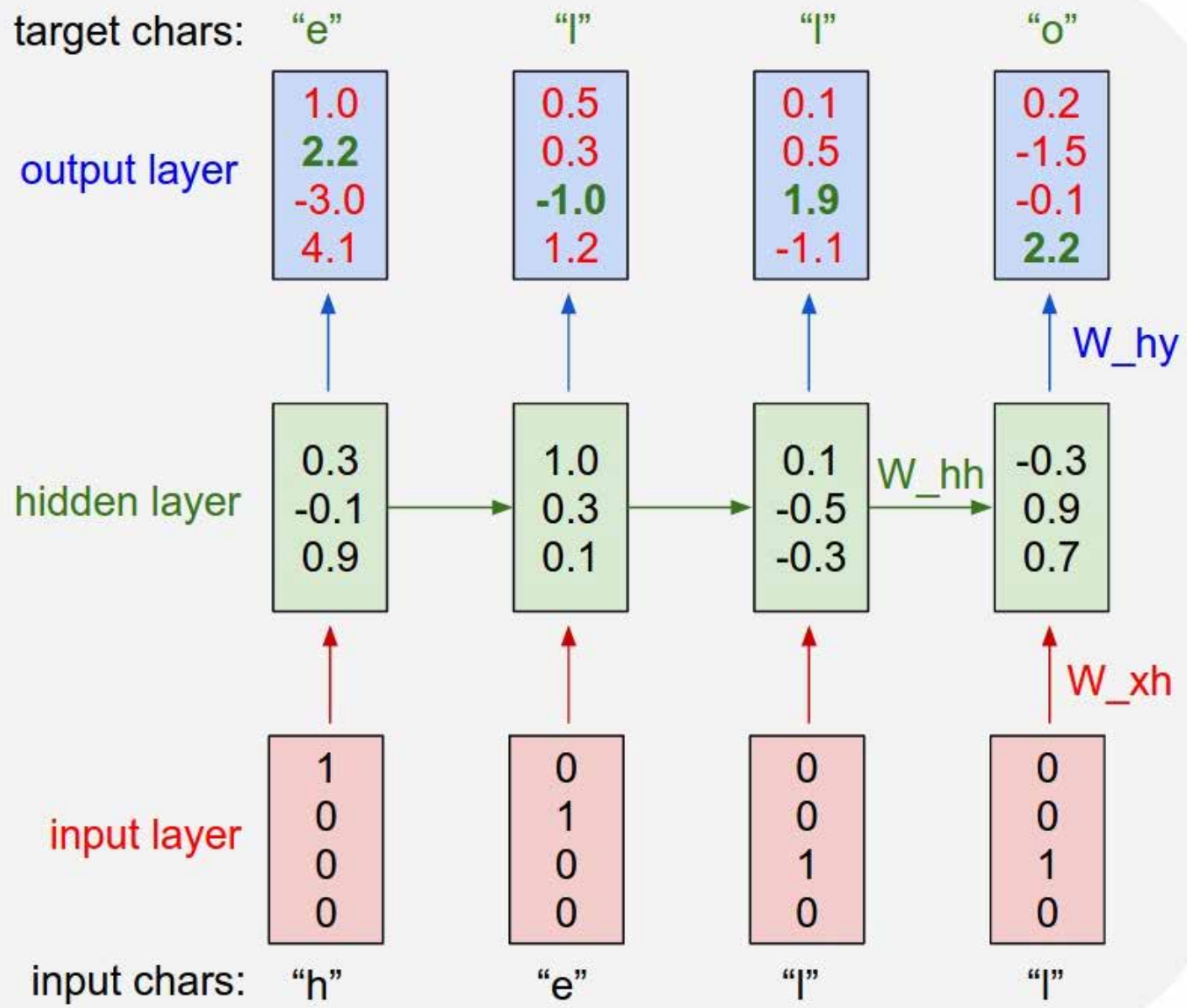
Teaser #2: Fun example: Character Level Networks

- give the RNN a huge chunk of text and ask it to model the probability distribution of the next character in the sequence given a sequence of previous characters.
- generate new text one character at a time.
- assume we have a vocabulary of 4 possible letters: “helo”
- Now, train an RNN on a sequence like: “hello”
- For each letter, the context will be the letters before it.

Fun example: Character Level Networks



Borrowed from Andrej Karpathy



We want:

Green numbers to be high
Red numbers to be low

We use the standard Softmax classifier (cross-entropy loss) on every output vector simultaneously.

Code available:
<https://gist.github.com/karpathy/d4dee566867f8291f086>

An example RNN with 4-dimensional input and output layers, and a hidden layer of 3 units (neurons)

Evolution of Samples during Training

- At iteration 100 the model samples random jumbles:

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

- At 300 iterations the model starts to get an idea about quotes and periods:

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

- At 500 iterations the words are now also separated with spaces and the model starts to get the idea about periods at the end of a sentence.

```
we counter. He stutn co des. His stanted out one ofler that concossions and was  
to gearang reay Jotrets and with fre colt of f paitt thin wall. Which das stimn
```

Evolution of Samples during Training

- At iteration 700 we're starting to see more and more English-like text emerge:

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 offer.

- At iteration 1200 we're now seeing use of quotations and question/exclamation marks. Longer words have now been learned as well:

"Kite vouch!" he repeated by her
door. "But I would be done and quarts, feeling, then, son is people...."

- we start to get properly spelled words, quotations, names, and so on by about iteration 2000:

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

Visualizing predictions and “neuron” firings in RNN

t t p : / / w w w . y n e t n e w s . c o m /] E n g l i s h - l a n g u a g e w e b s i t e o f I s r a e l ' s l a r
t p : / / w w w . b a c a h e t s . c o m / - x g l i s h l i n g u a g e s a i r s i t e o f t s l a e l i s s i n g
d : x n e . w a e a . . a w a t o a . s & n t i a c a - s a r d e e l h o a n t b i s a n f a n r e i f ' a a t d
m w - 2 p i i i s o e s s i s . / e r n . c] (d c e e n e p e s a a i k i i e e l e d h , i r t h r a o n s e , c o s e
d r . < : a h b - n p t w t . x i g h / m a) T v d r y z i c o u e d l s u : t h a - o o t u , s t u i f l v e p e r y
s t p , t c o a 2 d r u l w o c l e n s r] p . l l v a o d , , e y t c - n d m - o i b u v s] b b i m s u l t a t l y b n

g e s t n e w s p a p e r ' ' [[Y e d i o t h A h r o n o t h]] ' ' ' ' H e b r e w - l a n g u a g e p e r i o d
e t a a w s p a p e r s o ' [[T e l t i (f e a n e m t i) ' ' * ' ' [e r r e w s l e n g u a g e : a r o s o d i
i r s c o e e n a i T T h A o a i n n h S r m u w] e y s [' i n e i a ' s i w d d e ' h s o l r i f r :
u s . . s e t l g o r s . a s a t C a r e e g ' a C l r i s z] i e ' : : , # : T A a a a a t B a s e e i l o ' i a n f v l
- t u a e v r t i d , t B A m S u s y u t]] A s a o i g s]] , . : s M B o l o u s : T o u a - n : d w o a p n u
a , d , i i u i t i c p .] (l S v H v t u s u i e D n o e g a n o . ,] : { C C u i b o h e C y b k s l s : r - e p c n t s

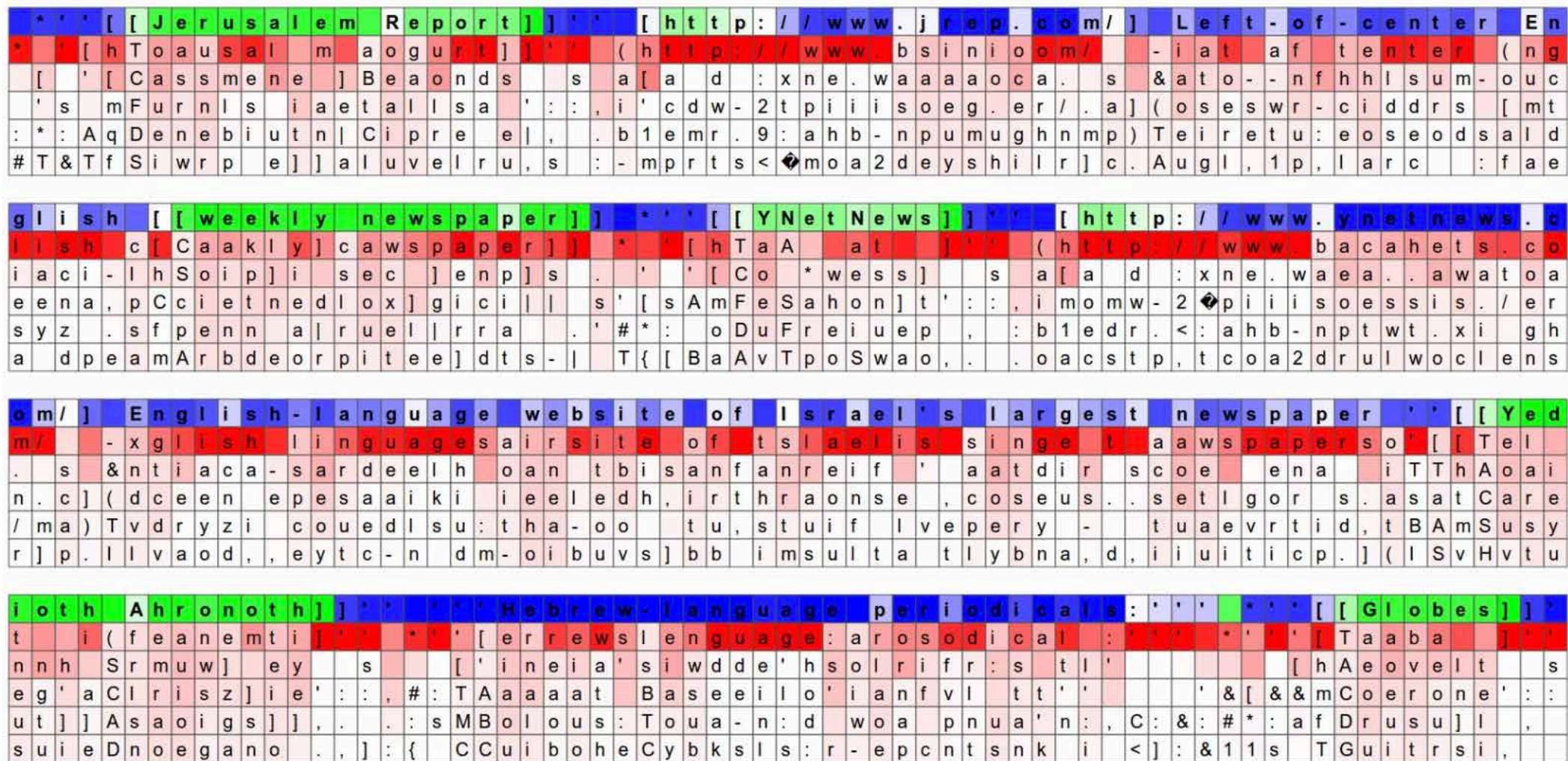
i c a l s : ' ' ' * ' ' [[G l o b e s]] ' ' [h t t p : / / w w w . g l o b e s . c o . i l /] b u s i n e s s d a
c a l : ' ' ' * ' ' [T a a b a] ' ' ([t t p : / / w w w . b u o b a l . c o m u n / s A - y t i n e s s a e t
s t l ' [h A e o v e l t s a h a d : x g e . w a o i r . r t o a . e l . i T & a i e g e o o y
t t ' ' ' & [& & m C o e r o n e ' : : , i ' o d w . , : n i i i s a a u e . e n i / o m l c C . (e f t g i r i i u
a ' n : , C : & : # * : a f D r u s u] l , . o m e l p < , d h a ; d e u o o t / i h n c s i f S , u r h o s t , t u n
n k i <] : & 1 1 s T G u i t r s i , : b a c m r - x t p o b - g r e s i s l e r l n a f a D] l o s p t a d , i f r m

i l y * ' ' [[H a a r e t z | H a ' A r e t z]] ' ' [h t t p : / / w w w . h a a r e t z . c o . i l /] R e l a t i v
l y * ' ' [[T e r r d n F e r a n t a h]] ' ' ([t t p : / / w w w . b o n m d s t . c o m u n / s - e s a t e o i
r e ' ' ' h A i l n n t t e H a l s r c n o l ' s a h a d : x n e . w a a m r t d h e o h . o l . c & o p i n i v e
k i . : * s C O S a n l t h i T i m ' l i] e : , i m c d w - 2 p h i i s e r d i t . i n a / c m f i . (a f l c a n a
d s - ! [t B T C o m m g d]] W o n a a e , : . b a e r r . < t a i b - d u l c n n c / a r n e s i] l i c e y s t o
n d s # & : G l D u v c c s a o S u c l t e l] z l , : o ' o m t] , : e o a 2 n i v f s r o o e i u n a l a) u v v r o

The neuron highlighted in this image seems to get very excited about URLs and turns off outside of the URLs. The LSTM is likely using this neuron to remember if it is inside a URL or not.

Borrowed from Andrej Karpathy

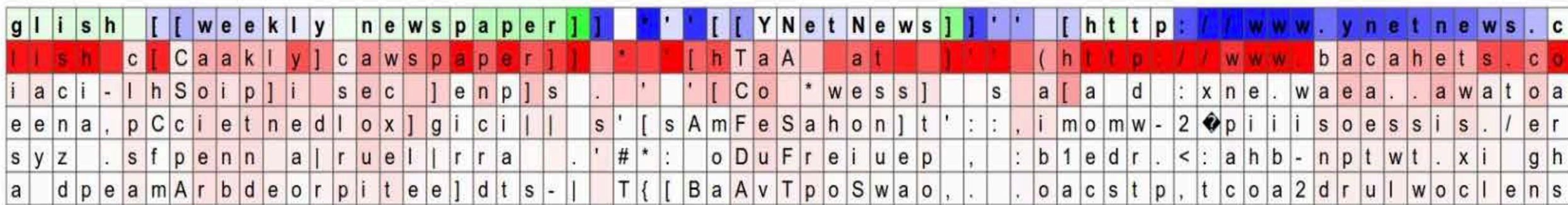
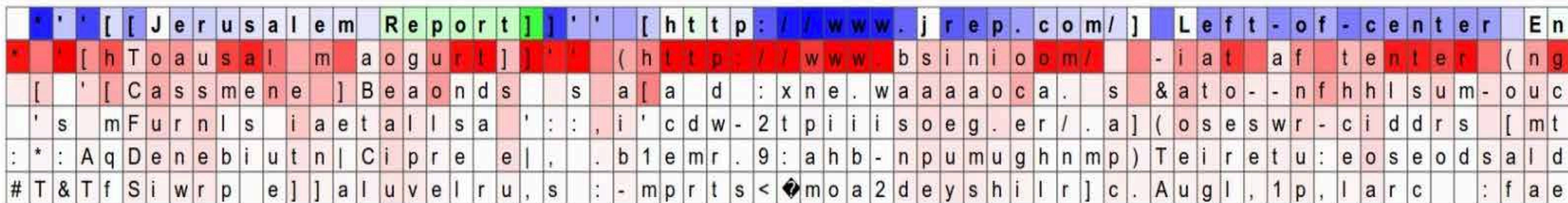
Visualizing predictions and “neuron” firings in RNN



The highlighted neuron here gets very excited when the RNN is inside the `[[]]` markdown environment and turns off outside of it.

Interestingly, the neuron can't turn on right after it sees the character `[`, it must wait for the second `[` and then activate. This task of counting whether the model has seen one or two `[` is likely done with a different neuron.

Visualizing predictions and “neuron” firings in RNN



Here we see a neuron that varies seemingly linearly across the `[[]]` environment. In other words its activation is giving the RNN a time-aligned coordinate system across the `[[]]` scope. The RNN can use this information to make different characters more or less likely depending on how early/late it is in the `[[]]` scope (perhaps?).

Visualizing predictions and “neuron” firings in RNN

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."

Visualizing predictions and “neuron” firings in RNN

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;
}
```


Visualizing predictions and “neuron” firings in RNN

Cell that turns on inside comments and quotes:

```
/* 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;
}
```

Visualizing predictions and “neuron” firings in RNN

Cell that is sensitive to the depth of an expression:

```
#ifdef CONFIG_AUDIT_SYSCALL
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;
}
```

Visualizing predictions and “neuron” firings in RNN

Cell that might be helpful in predicting a new line. Note that it only turns on for some “)”:

```
char *audit_unpack_string(void **bufp, size_t *remain, si
{
    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.
     */
    if (len > PATH_MAX)
        return ERR_PTR(-ENAMETOOLONG);
    str = kmalloc(len + 1, GFP_KERNEL);
    if (unlikely(!str))
        return ERR_PTR(-ENOMEM);
    memcpy(str, *bufp, len);
    str[len] = 0;
    *bufp += len;
    *remain -= len;
    return str;
}
```

How to generate character strings from the model

- Start the model with its default hidden state.
- Give it a “burn-in” sequence of characters and let it update its hidden state after each character.
- Then look at the probability distribution it predicts for the next character.
- Pick a character randomly from that distribution and tell the net that this was the character that actually occurred.
 - i.e. tell it that its guess was correct, whatever it guessed.
- Continue to let it pick characters until bored.
- Look at the character strings it produces to see what it “knows”.

He was elected President during the Revolutionary War and forgave Opus Paul at Rome. The regime of his crew of England, is now Arab women's icons in and the demons that use something between the characters' sisters in lower coil trains were always operated on the line of the **ephemeral** street, respectively, the graphic or other facility for deformation of a given proportion of large segments at RTUS). The B every chord was a "strongly cold internal palette pour even the white blade."

Some completions produced by the model

- Sheila thrungs **s** (most frequent)
- People thrunge (most frequent next character is space)
- Shiela, Thrunge **lini del Rey** (first try)
- The meaning of life is **literary recognition.** (6th try)
- The meaning of life is **the tradition of the ancient human reproduction: it is less favorable to the good boy for when to remove her bigger.**
(one of the first 10 tries for a model trained for longer).



RNN Bible @RNN_Bible · 18 Jun 2016



24:3 And there went up the captivity of Jacob his son, and slew the trumpet of the God of Israel.



RNN Bible @RNN_Bible · 18 Jun 2016



107:33 Therefore they that were slain by the sword, and did eat and drink; and they stood before them.



RNN Bible @RNN_Bible · 18 Jun 2016



23:1 And Joash said unto them, Thus shall it be down in the land of Egypt, and they shall come to the most High;



RNN Bible @RNN_Bible · 18 Jun 2016



30:22 For I will break them away captive to the mountains, and there shall be no more a prince among the nations.



RNN Bible @RNN_Bible · 18 Jun 2016

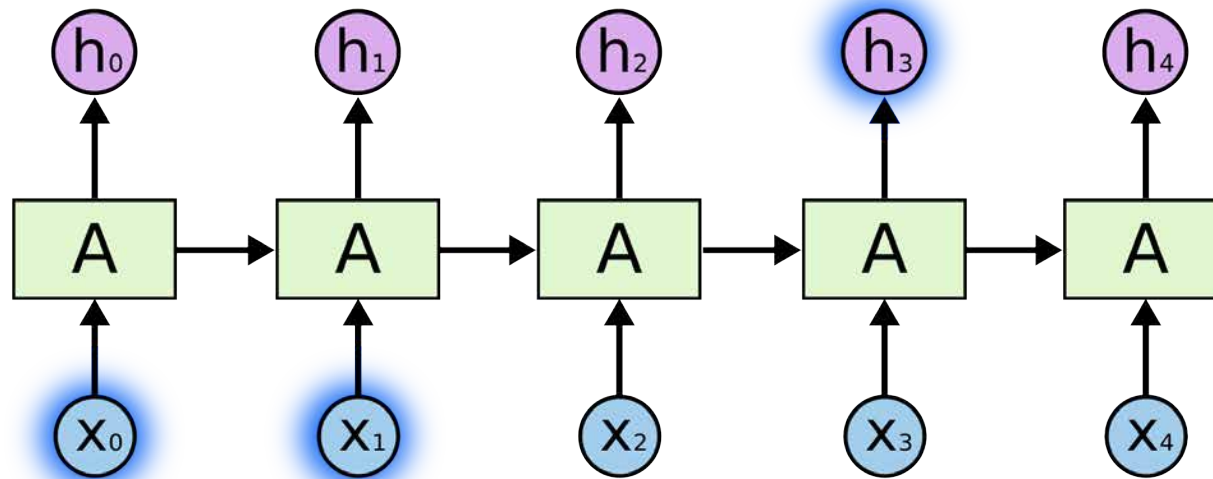


2:3 And it shall come to pass, that when thou goest to possess it, that they may not eat thereof.



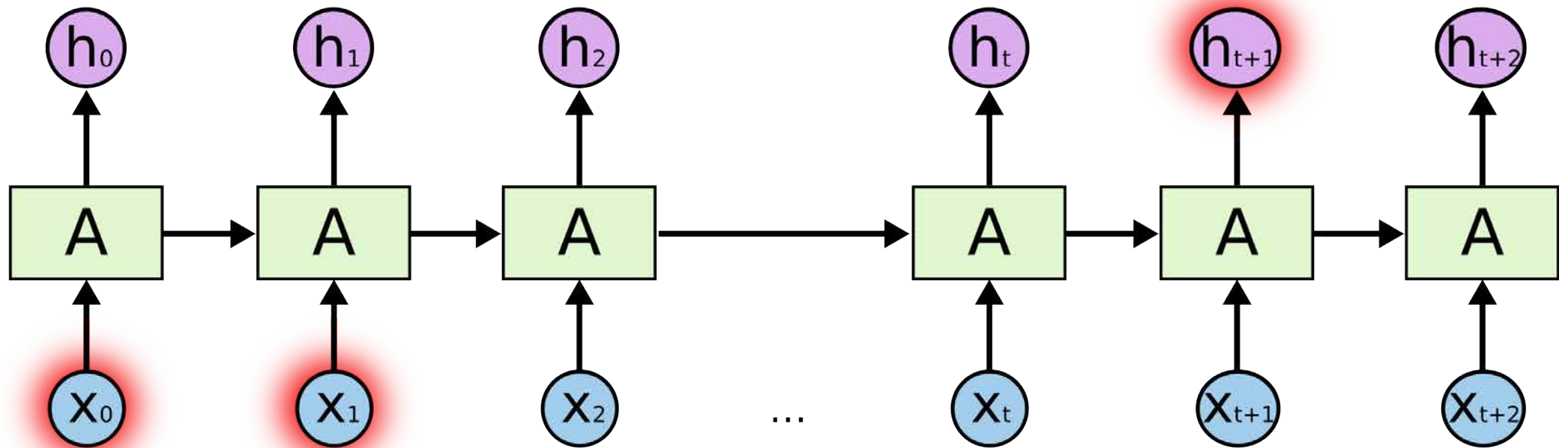
The Problem of Long-Term Dependencies

- Sometimes, we only need to look at recent information to perform the present task.
- Predict last word in the sentence:
 - “the clouds are in the *sky*” (relevant information is very close)



The Problem of Long-Term Dependencies

- There are also cases where we need more context.
- Predict last word in the sentence:
 - “I grew up in France... I speak fluent *French*.” (relevant information is far away)



Long Short-Term Memory Networks (LSTM)

Long Short-Term Memory, 1997

Sepp Hochreiter, J. Schmidhuber

<http://www.mitpressjournals.org/doi/abs/10.1162/neco.1997.9.8.1735>

Learning to Forget: Continual Prediction with LSTM, 2000

F. A. Gers, J. Schmidhuber, F. Cummins

<http://www.mitpressjournals.org/doi/abs/10.1162/089976600300015015>

LSTM recurrent networks learn simple context-free and context-sensitive languages, 2001

F. A. Gers, J. Schmidhuber

<http://ieeexplore.ieee.org/abstract/document/963769>