

Machine Intelligence: Deep Learning

Modelling sequence data and using ideas from CNN

&

Recurrent Neural Networks

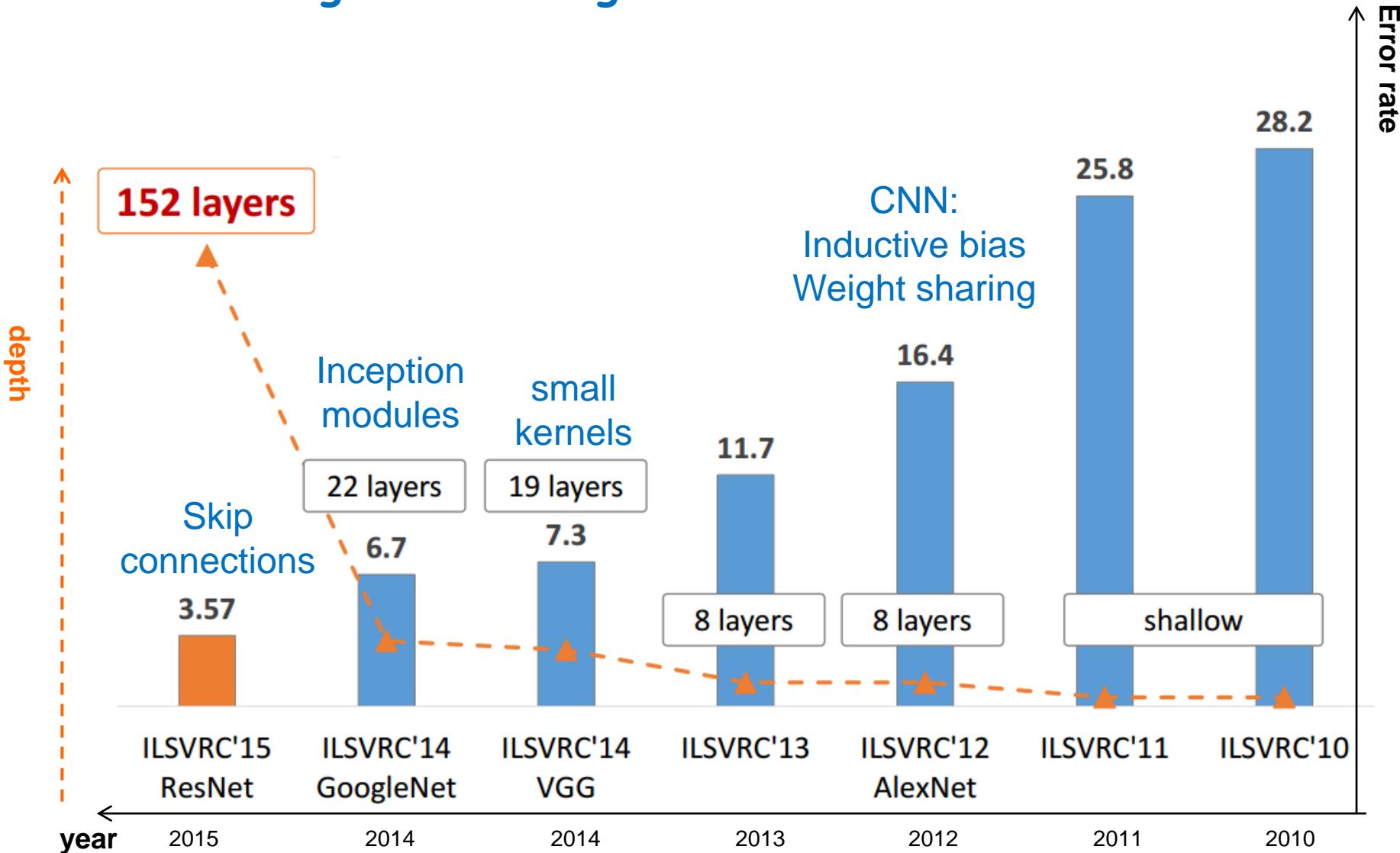
Beate Sick
sick@zhaw.ch

Remark: Much of the material has been developed together with Elvis Murina and Oliver Dürr

Topics

- **Famous tricks in challenge winning CNN architectures**
 - Inductive Bias, weight sharing, Inception moduls, gradient highway via skip-connections
- **Introduction of recurrent neural networks (RNN)**
 - possible use cases
 - memory in recurrent NN
- **Working with an RNN**
 - stepping through an RNN
 - loss construction in an RNN
 - RNNs in Keras
- **Tricks of the trade**
 - common and different tricks to train deep CNNs and RNNs

Review of ImageNet winning CNN architectures



Going deeper is easy – or not?



The challenge is to design a network in which the gradient can reach all the layers of a network which might be dozens, or even hundreds of layers deep.

This was achieved by some recent improvements, such as ReLU and batch normalization, and by designing the architecture in a way which allows the gradient to reach deep layer, e.g. by additional skip connections.

“Oxford Net” or “VGG Net” 2nd place

- 2nd place in the imageNet challenge
- More traditional, easier to train
- Small pooling
- Stacked 3x3 convolutions before maxpooling
-> large receptive field
- ReLU after conv. and FC (batchnorm was not used)
- More weights than GoogLeNet

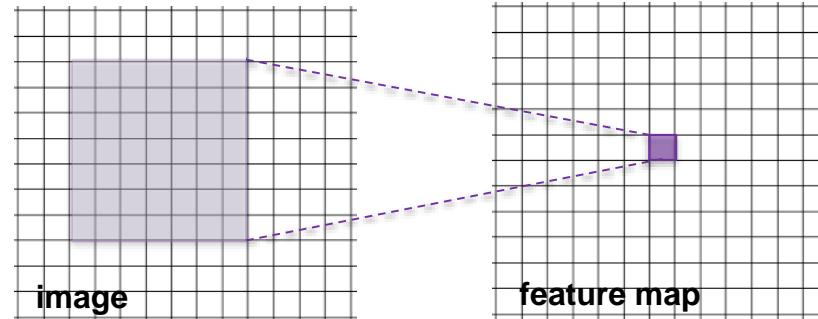


<http://arxiv.org/abs/1409.1556>

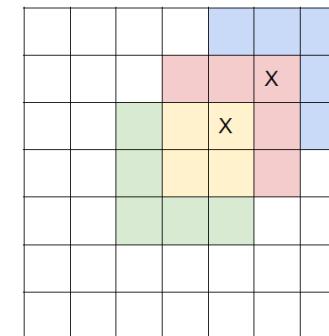
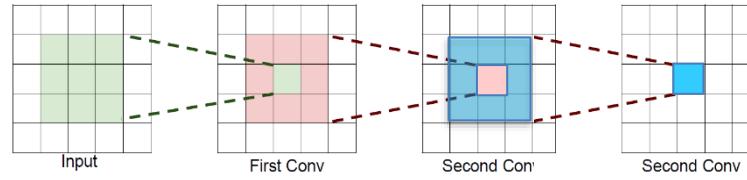
The trend in modern CNN architectures goes to small filters

How best to achieve a receptive field of 7x7 pixels?

- 1) Working with **one**
7x7 conv layers (stride 1)
49 weights

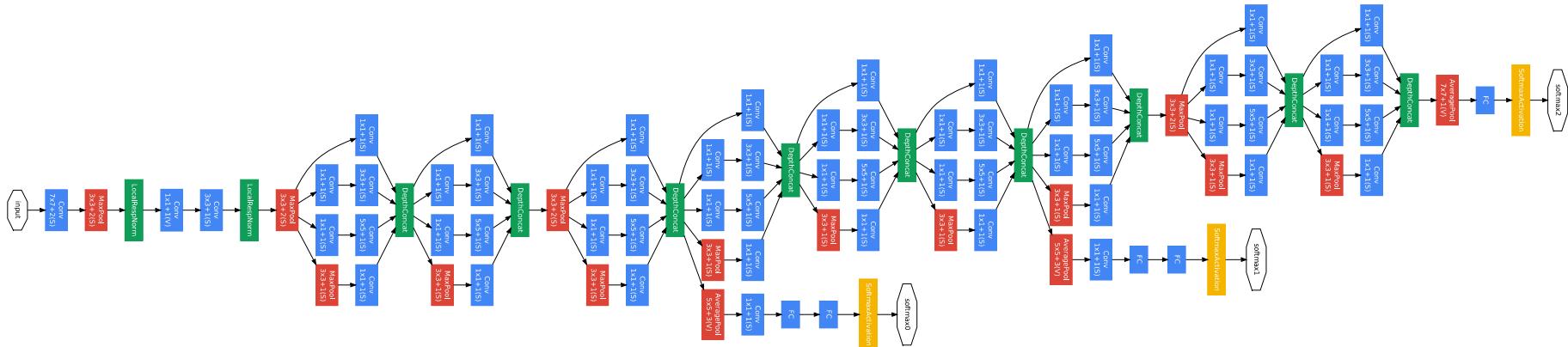


- 2) Working with **three** **3x3 conv layers** (stride 1)
 $3 \times 9 = 27$ weights

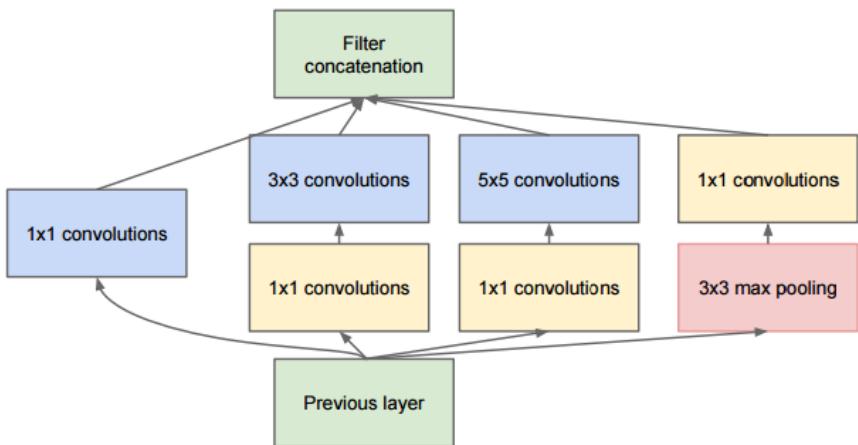


When working stacking conv-layers with small kernels we can achieve with less weights and more non-linear combinations the same receptive field

Winning architecture (GoogLeNet, 2014)

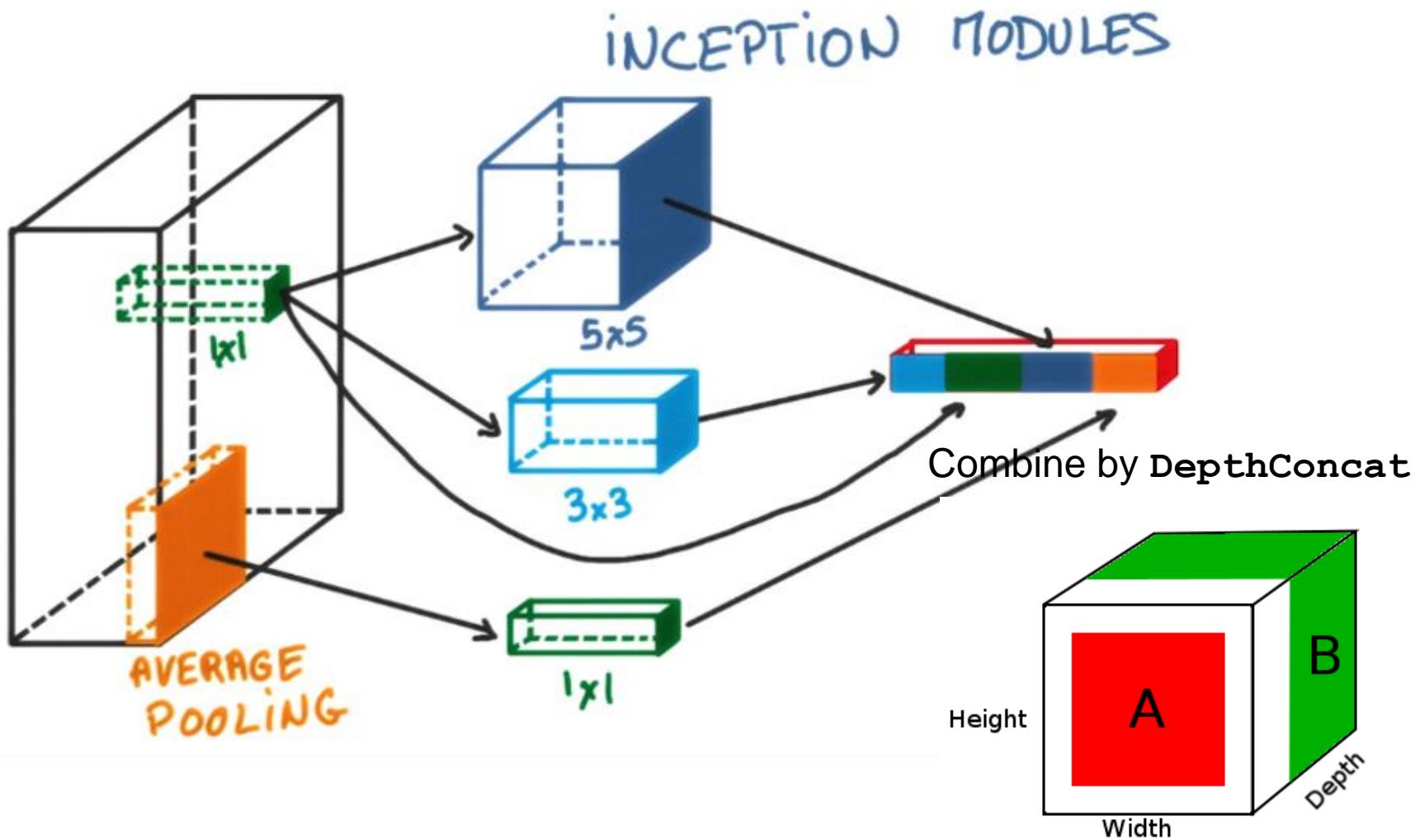


The inception module:
use in 1 layer in parallel different kernels and combine their results



Few parameters, hard to train.
Comments see [here](#)

The idea of inception modules

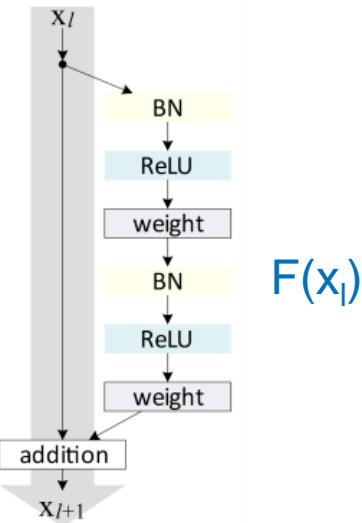


"ResNet" from Microsoft 2015 winner of imageNet

152
layers

ResNet basic design (VGG-style)

- add shortcut connections every two
- all 3x3 conv (almost)



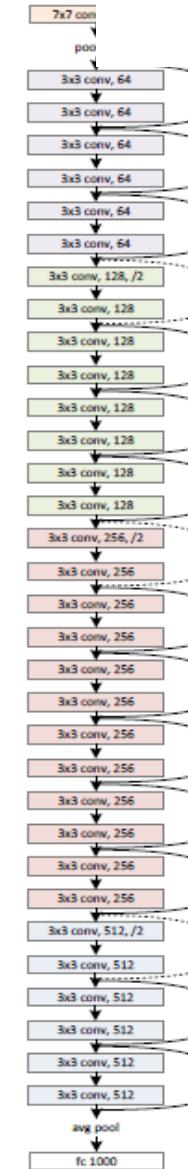
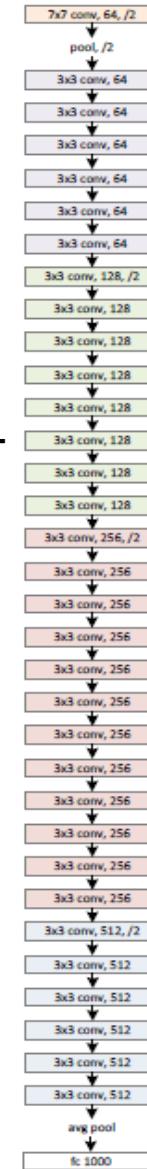
$$H(x_i) = x_{i+1} = x_i + F(x_i)$$

$F(x)$ is called "residual" since it only learns the "delta" which is needed to add to x to get $H(x)$

152 layers:
Why does this train at all?

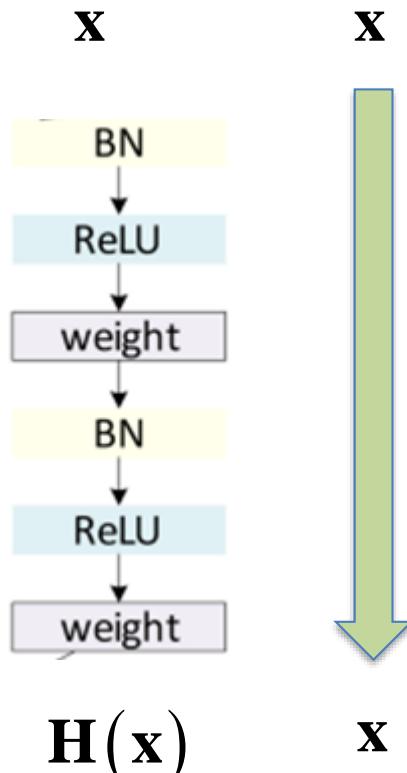
This deep architecture
could still be trained, since
the gradients can skip
layers which diminish the
gradient!

plain VGG



ResNet

Highway Networks with skip connections: providing a highway for the gradient



Idea: Use nonlinear transform T to determine how much of the output y is produced by H or the identity mapping. Technically we do that by:

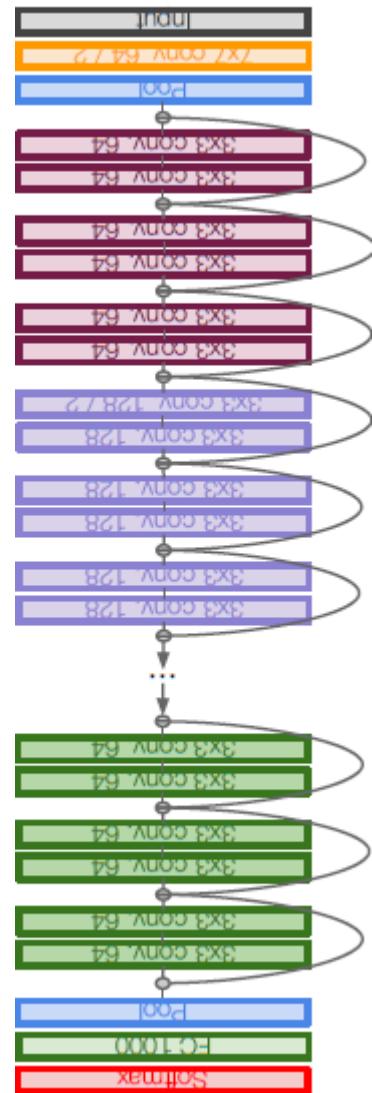
$$y = H(x, \mathbf{W}_H) \cdot T(x, \mathbf{W}_T) + x \cdot (1 - T(x, \mathbf{W}_T)).$$

Special case:

$$y = \begin{cases} x, & \text{if } T(x, \mathbf{W}_T) = 0 \\ H(x, \mathbf{W}_H), & \text{if } T(x, \mathbf{W}_T) = 1 \end{cases}$$

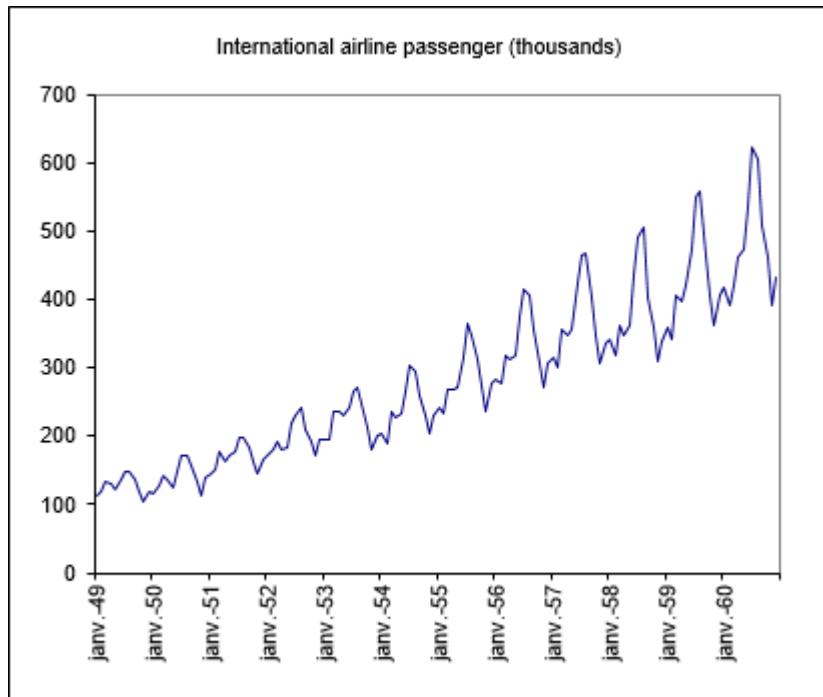
This opens a highway for the gradient:

$$\frac{dy}{dx} = \begin{cases} \mathbf{I}, & \text{if } T(x, \mathbf{W}_T) = 0, \\ H'(x, \mathbf{W}_H), & \text{if } T(x, \mathbf{W}_T) = 1. \end{cases}$$

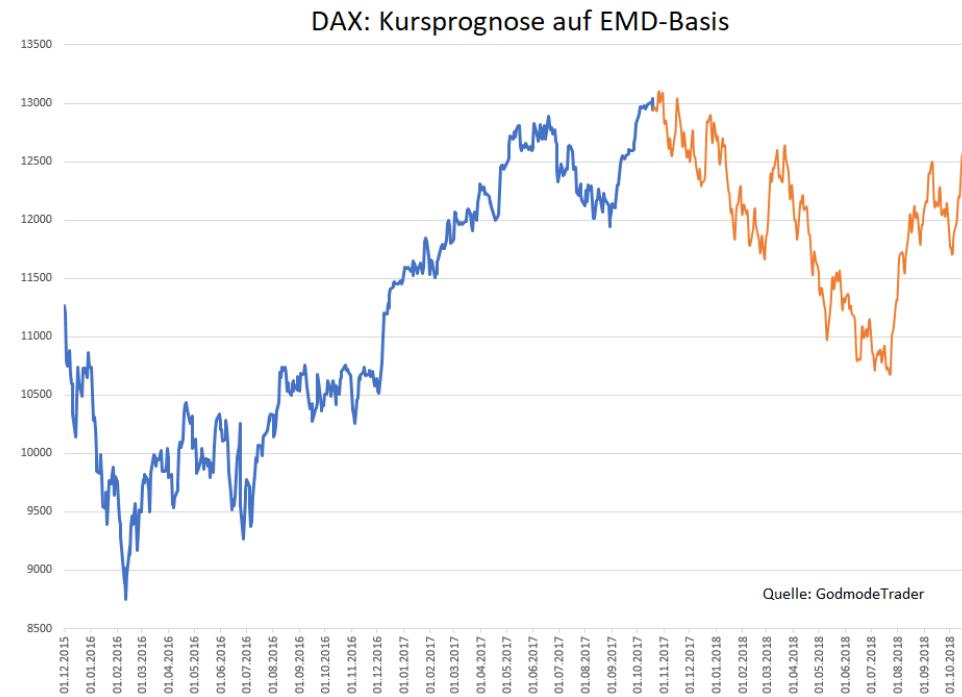


Sequence Data

Example Sequence Data: time-series

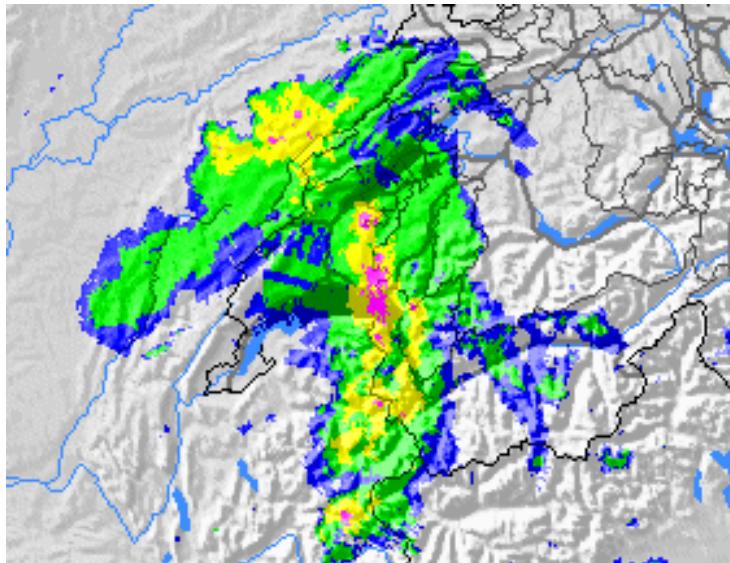


How many passenger will we have next month?



What will be the DAX value tomorrow?

Example Sequence Data: videos



https://www.metradar.ch/2009_exp/pc/service.php

How much does it rain in
Winterthur within the next hour?



<https://www.pinterest.ch/pin/460704236854535539/>

Who will be the next star?

Example sequence data: speech translation



Speech Recognition Breakthrough for the Spoken, Translated Word

2012: [Microsoft Chief Research Officer Rick Rashid demonstrates](#) breakthrough in DL based translation that converts his spoken English words into computer-generated Chinese language.

2017: Language translator is available as [mobile app](#)

Example Sequence Data: text

Guten Tag, Sick Beate (sick)

siehe Anhang

gescanntes Dokument.:

[http://dildosatisfaction.com/Rechnungs-Details-98773504333/Sick Beate \(sick\)](http://dildosatisfaction.com/Rechnungs-Details-98773504333/Sick Beate (sick))

Mit freundlichen GrÃ¼ÃŸe

I know the sender very well!

Should I open the attachment?

Example: Produce sequence of words as caption



"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."



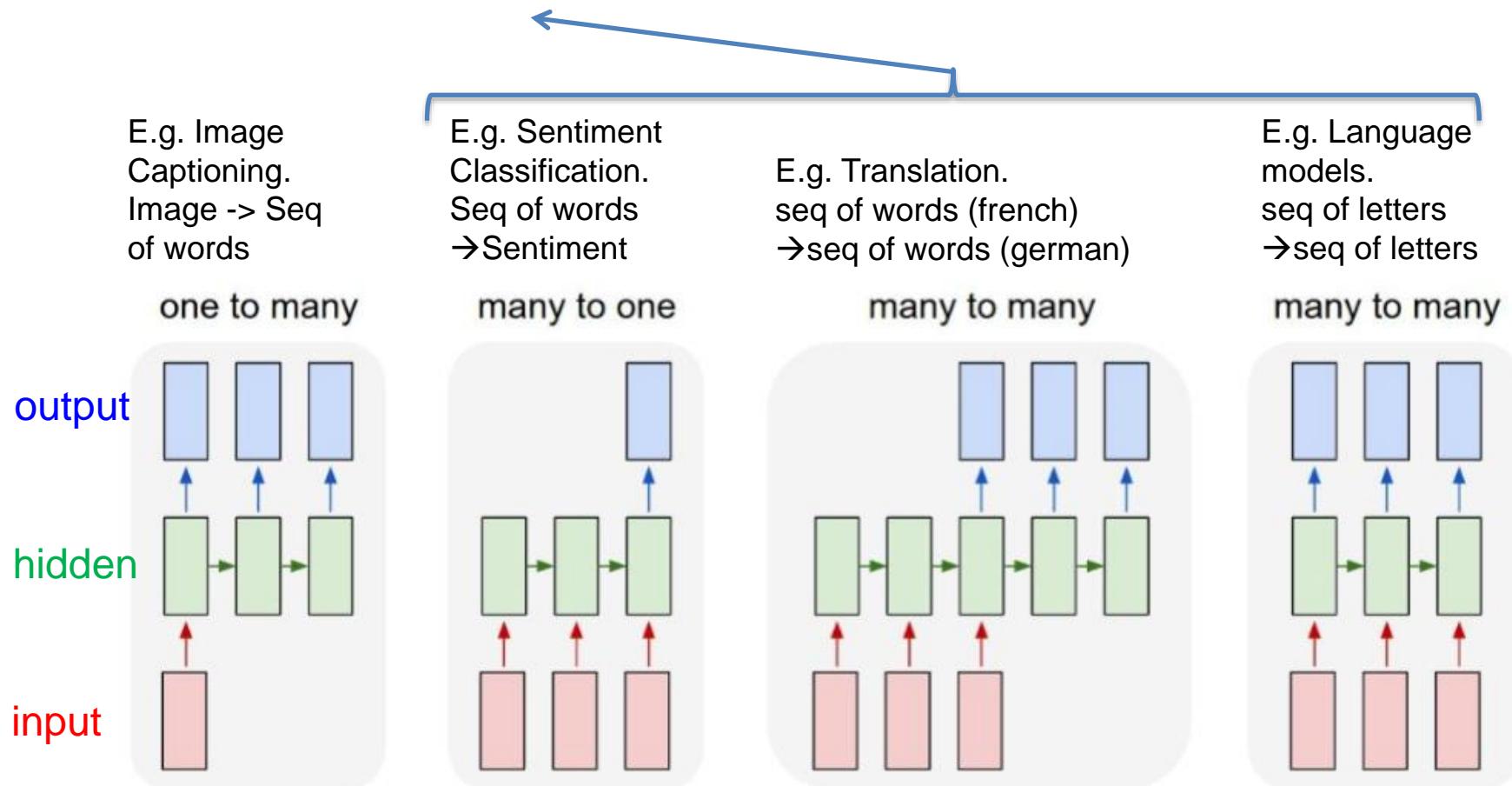
"a woman holding a teddy bear in front of a mirror."



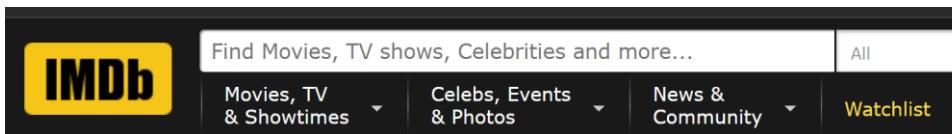
"a horse is standing in the middle of a road."

Modeling sequence data

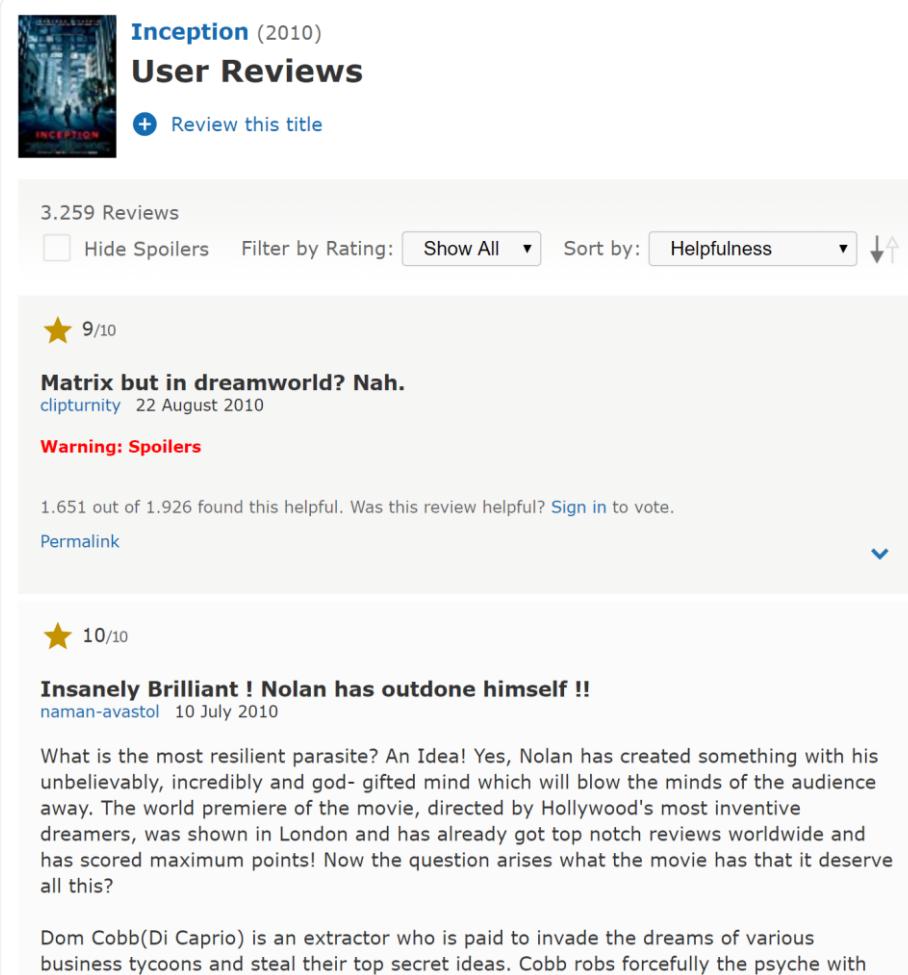
To learn with sequence data we need a **memory** about the seen former parts.



Possible task: Sentiment analysis with movie reviews



The image shows the IMDb website's header. It features the yellow 'IMDb' logo on the left. To its right is a search bar with the placeholder text 'Find Movies, TV shows, Celebrities and more...'. Below the search bar are four main navigation categories: 'Movies, TV & Showtimes', 'Celebs, Events & Photos', 'News & Community', and 'Watchlist'. Each category has a dropdown arrow next to it.



The image shows the 'User Reviews' section for the movie 'Inception' (2010) on IMDb. At the top, there is a thumbnail image of the movie poster. Below it, the title 'Inception (2010)' and 'User Reviews' are displayed. A blue link '+ Review this title' is visible. Underneath, there are 3,259 reviews. A 'Hide Spoilers' checkbox is checked. The sorting options are 'Show All' and 'Helpfulness' (sorted by descending order). The average rating is 9/10. A review by 'Matrix but in dreamworld? Nah.' from 'clipturny' on 22 August 2010 is shown, with a warning about spoilers. Another review by 'Insanely Brilliant ! Nolan has outdone himself !!' from 'naman-avastol' on 10 July 2010 is also shown, along with a detailed summary of the plot.

	review	sentiment
0	I went and saw this movie last night after bei...	1
1	Actor turned director Bill Paxton follows up h...	1
2	As a recreational golfer with some knowledge o...	1
3	I saw this film in a sneak preview, and it is ...	1
4	Bill Paxton has taken the true story of the 19...	1

Challenges:

- 1) We need to find a numeric representation of words (e.g. bag of words, or embedding)
- 2) We need to be able to handle inputs of different length.

Bag of words: Ignoring word order

- Count vectors or “bag of words”
 - Determine vocabulary (or alphabet, or word, or token)

Example:

Document 1: “The cat sat on the hat”

Document 2: “The dog ate the cat and the hat”

Bag of words (=word count vector):

document	the	cat	sat	on	hat	dog	ate	and
1	2	1	1	1	1	0	0	0
2	3	1	0	0	1	1	1	1

Represent document as word count vector ignoring order of words (token).

This allows to represent each sentence as a numeric vector of the same length!

This can be seen as feature-vector and can be used for traditional classifiers as RF.

Getting Bag of words in sklearn: ignoring word order

```
from sklearn.feature_extraction.text import CountVectorizer  
vectorizer = CountVectorizer(max_features=6000, min_df=5, max_df=0.7)  
  
X = vectorizer.fit_transform(documents).toarray()
```

The diagram illustrates the parameters of the CountVectorizer class with three annotations:

- An annotation for `max_features=6000` points to the first parameter in the constructor. It is described as "consider only the 6000 most frequent words".
- An annotation for `min_df=5` points to the second parameter. It is described as "take only words that appear in at least 5 different reviews".
- An annotation for `max_df=0.7` points to the third parameter. It is described as "ignore all words that appear in more than 70% of all reviews".

Get from text to ordered numeric vector

Step 1) Tokenize text, 1-hot-encoding

- Determine the size N of the relevant vocabulary
- Each token (word) is represented by a numeric value between 0 and N-1
- Corresponding 1-hot-encoded representations have length N

Example: Look at vocabulary with N=9 to tokenize the 2 text-samples below

Vocabulary: {'the': 1, 'cat': 2, 'sat': 3, 'on': 4, 'mat': 5, 'dog': 6, 'ate': 7, 'my': 8, 'homework': 9}

Text-sample: “The cat sat on the hat”

“The dog ate my homework”

Tokenized: [1, 2, 3, 4, 1, 5].

[1, 6, 7, 8, 9]

1-hot:

1	0	0	0	0	0	0	0	0
0	1	0	0	0	0	0	0	0
0	0	1	0	0	0	0	0	0
0	0	0	1	0	0	0	0	0
1	0	0	0	0	0	0	0	0
0	0	0	0	1	0	0	0	0

the
cat
sat
on
the
hat

1	0	0	0	0	0	0	0	0
0	0	0	0	0	0	1	0	0
0	0	0	0	0	0	0	1	0
0	0	0	0	0	0	0	0	1
0	0	0	0	0	0	0	0	1

the
dog
ate
my
homework

Tokenizing in keras

```
from keras.preprocessing.text import Tokenizer  
  
samples = ['The cat sat on the mat.', 'The dog ate my homework.]
```

Creates a tokenizer, configured to only take into account the 1,000 most common words

```
tokenizer = Tokenizer(num_words=1000)
```



```
tokenizer.fit_on_texts(samples)
```

Turns strings into lists of integer indices

```
sequences = tokenizer.texts_to_sequences(samples)
```



```
one_hot_results = tokenizer.texts_to_matrix(samples, mode='binary')
```

```
word_index = tokenizer.word_index
```



```
print('Found %s unique tokens.' % len(word_index))
```

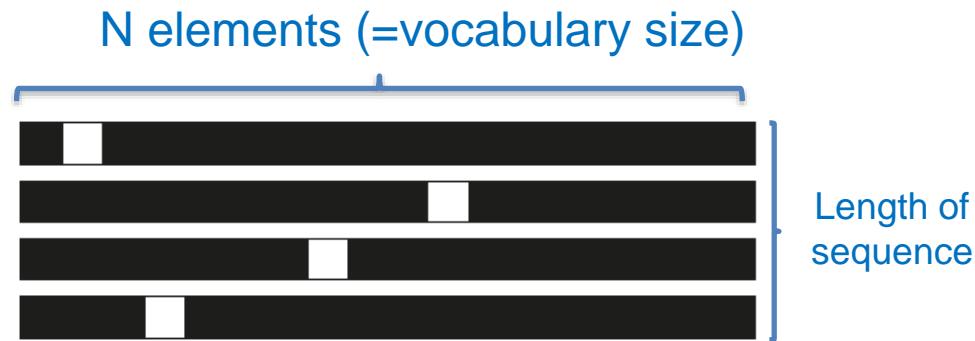
How you can recover the word index that was computed

You could also directly get the one-hot binary representations. Vectorization modes other than one-hot encoding are supported by this tokenizer.

Go from 1-hot-encodings to word embeddings

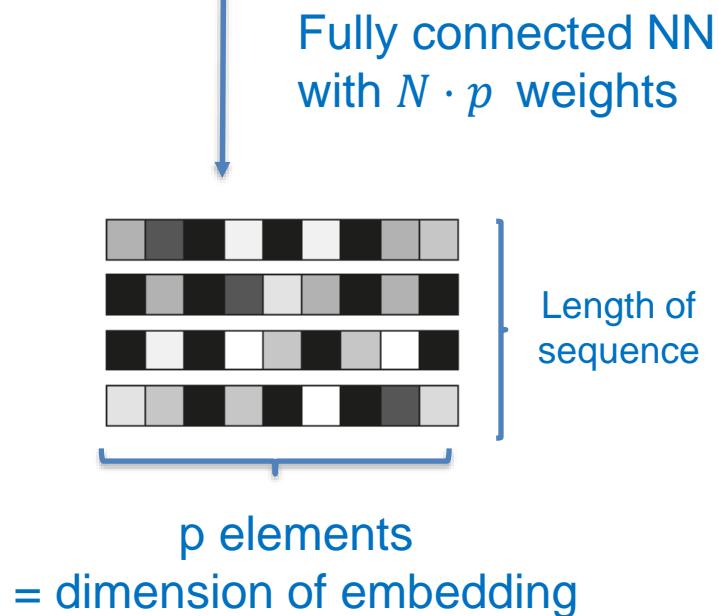
1-hot encodings

- Based on vocabulary of size N
- sparse: one 1 and $N-1$ zeros
- High-dim: vector-length = N



Word embedding's are

- Dense
- Low-dimensional: vector-length = p
- Learned from data via fcNN $N \rightarrow p$



Wordembedding layer in keras

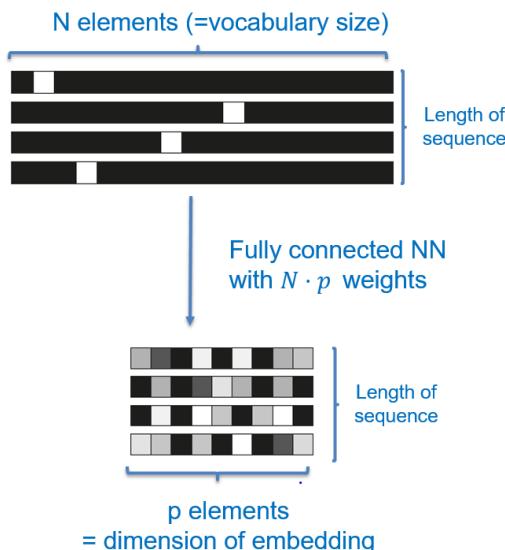
```
from keras.layers import Embedding  
embedding_layer = Embedding(1000, 64)
```

N p

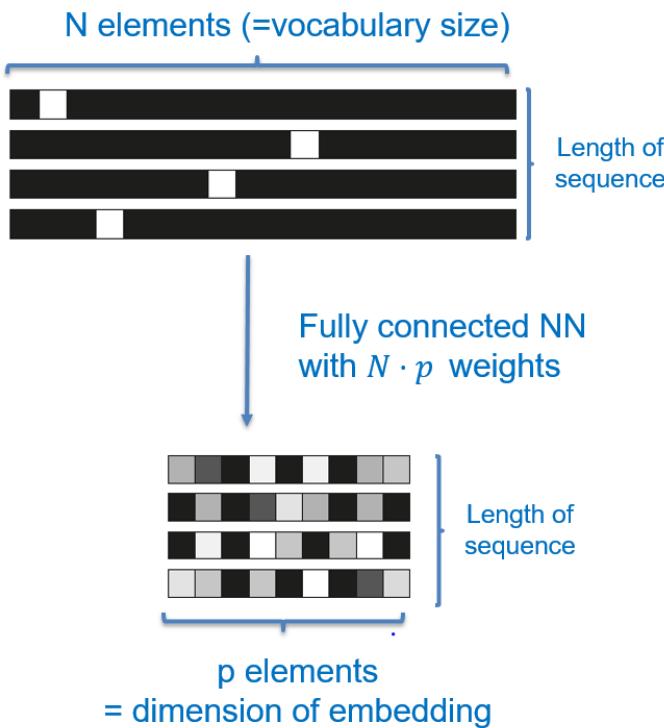
The Embedding layer takes at least two arguments: the number of possible tokens (here, 1,000: 1 + maximum word index) and the dimensionality of the embeddings (here, 64).

The embedding layer returns a 3D floating-point tensor of shape (samples, sequence_length, embedding_dimensionality).

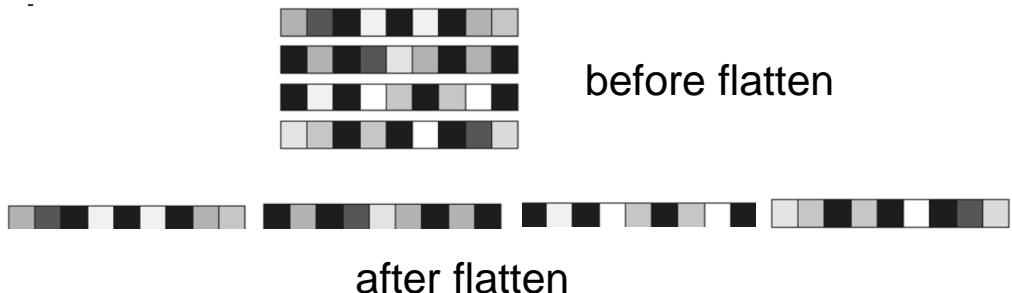
Look at 1 sample:



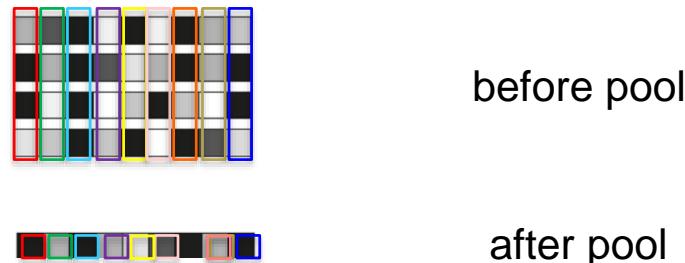
How to proceed with embeddings



- **Flatten** results in vector of length: $\text{seq-length} * p$

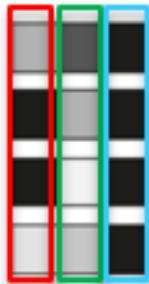


- **Pool** over sequence length results in vector of length: p

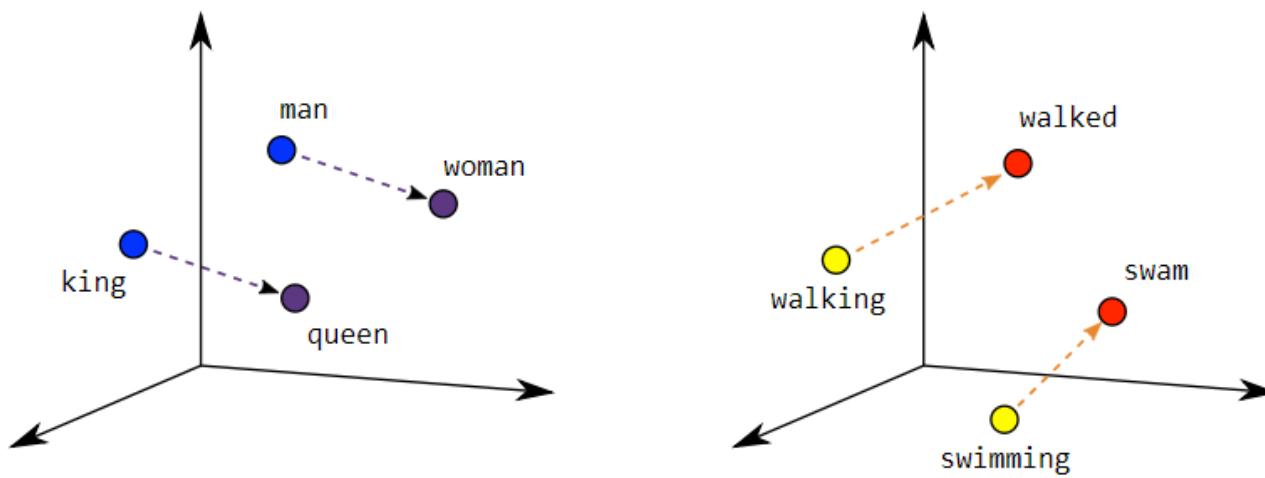


- Input to RNN layer or a 1D convolution layer

Word embedding space



Sequence with 4 words, represented by **3-dim embedding space**



A simple text classification NN using flattened embeddings

Specifies the maximum input length to the Embedding layer so you can later flatten the embedded inputs. After the Embedding layer, the activations have shape (samples, maxlen, 8).

```
from keras.models import Sequential
from keras.layers import Flatten, Dense

model = Sequential()
model.add(Embedding(10000, 8, input_length=maxlen))

model.add(Flatten())
model.add(Dense(1, activation='sigmoid'))
model.compile(optimizer='rmsprop', loss='binary_crossentropy', metrics=['acc'])
model.summary()

history = model.fit(x_train, y_train,
                     epochs=10,
                     batch_size=32,
                     validation_split=0.2)
```

Flattens the 3D tensor of embeddings into a 2D tensor of shape (samples, maxlen * 8)

Adds the classifier on top

A text classification NN using pooled embeddings

```
from keras.models import Sequential
from keras.layers import Dense, Embedding, GlobalAveragePooling1D, Dropout

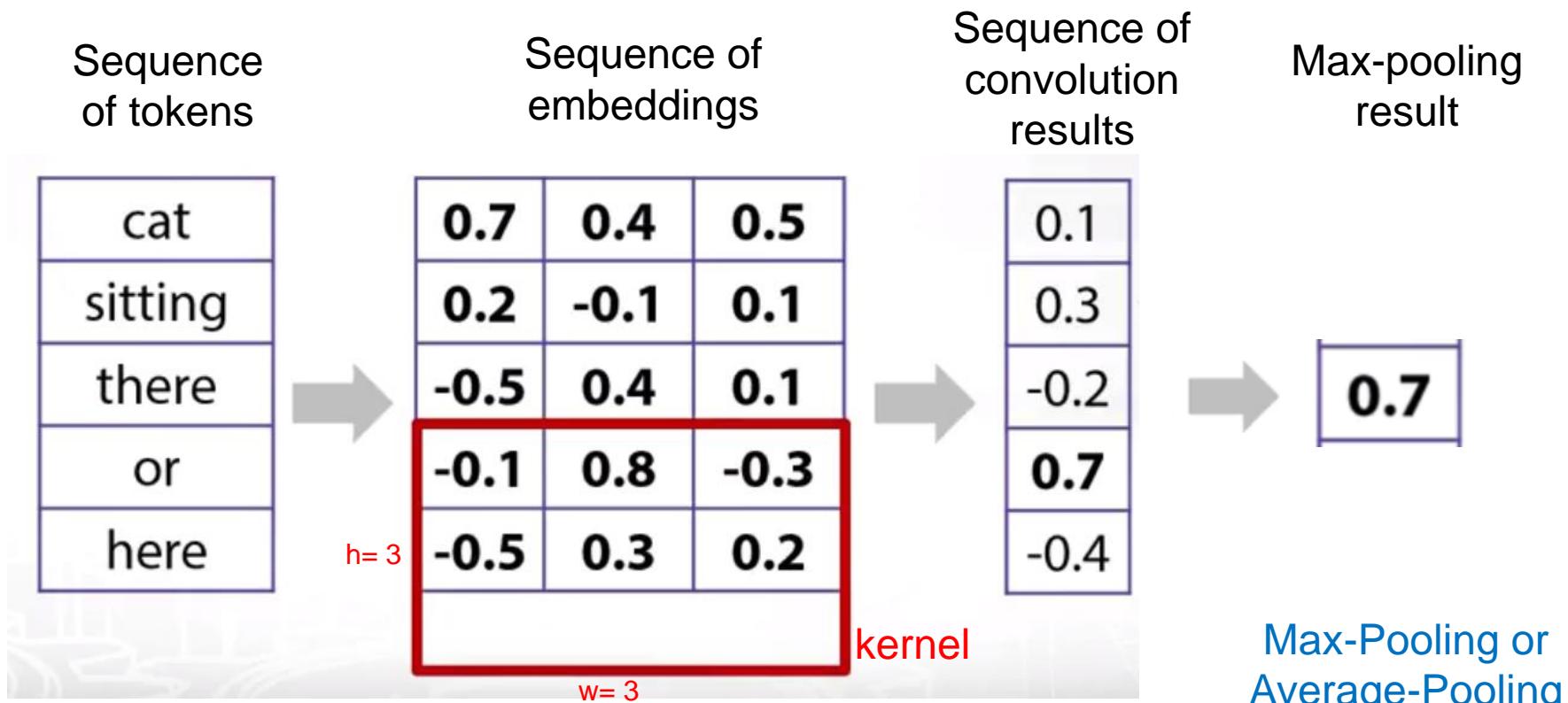
EMBEDDING_DIM = 30

model = Sequential()
model.add(Embedding(vocab_size, EMBEDDING_DIM, input_length=(None)))
model.add(GlobalAveragePooling1D())
model.add(Dropout(0.5))
model.add(Dense(20, activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(1, activation='sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])

model.summary()
```

Layer (type)	Output Shape	Param #
embedding_1 (Embedding)	(None, None, 30)	1868910
global_average_pooling1d_1 (GlobalAveragePooling1D)	(None, 30)	0
dropout_1 (Dropout)	(None, 30)	0
dense_1 (Dense)	(None, 20)	620
dropout_2 (Dropout)	(None, 20)	0
dense_2 (Dense)	(None, 1)	21

Applying 1D convolution embedding sequences



Credits: <https://www.youtube.com/watch?v=wNBaNhvL4pg>

w=width of kernel = embedding dimension
h=height of kernel = #neighboring tokes

Max-Pooling or Average-Pooling makes result independent of sequence length

We slide $h \times w$ kernel only in 1 direction \rightarrow 1D convolution

A kernel can intakes h embeddings (=h-grams, h neighboring words)

A text classification NN feeding embeddings to 1D conv

```
from keras.models import Model
from keras.layers import Input, Dense, Concatenate, Dropout, Embedding, Conv1D, GlobalMaxPooling1D, GlobalAveragePooling1D
EMBEDDING_DIM = 30

a = Input(shape=(max_length,))
x = Embedding(vocab_size, EMBEDDING_DIM)(a)
x1 = Conv1D(filters=50, kernel_size=(3), activation="relu", padding="same")(x)
x2 = Conv1D(filters=50, kernel_size=(5), activation="relu", padding="same")(x)
x3 = Conv1D(filters=50, kernel_size=(7), activation="relu", padding="same")(x)

g1 = GlobalAveragePooling1D()(x1)
g2 = GlobalAveragePooling1D()(x2)
g3 = GlobalAveragePooling1D()(x3)
conc = Concatenate()([g1, g2, g3])
conc = Dropout(0.3)(conc)
conc = Dense(50, activation='relu')(conc)
conc = Dropout(0.3)(conc)
out = Dense(1, activation='sigmoid')(conc)
model = Model(inputs=a, outputs=out)
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
model.summary()
```

Inception idea:
Use several filter-heights in parallel.

wiederkehrend

Recurrent Neural Networks

Resources

- Many figures are taken from the following resources:
 - Deep Learning Book chap10
 - <http://www.deeplearningbook.org/contents/rnn.html>
 - Other online DL courses
 - Lecture on RNN: http://cs231n.stanford.edu/slides/winter1516_lecture10.pdf
 - Video to CS231n <https://www.youtube.com/watch?v=iX5V1WpxxkY>
 - [CS 598 LAZ](#) Lecture 2 and 3 on RNN
 - Blog Posts
 - Karpathy, May 2015: The unreasonable effectiveness of Recurrent Neural Networks <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
 - Colah, August 2015: Understanding LSTM Networks
<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Recurrent NN have memory

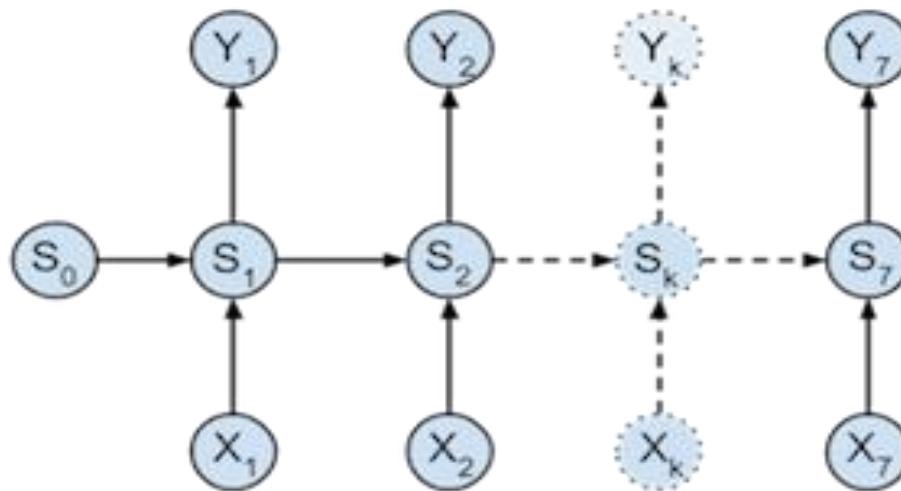
Challenge 2: How to handle input of different lengths?

Each text (e.g. e-mail) can have a different number of words!

By reading from word to word our belief in different categories (e.g. spam or not-spam) can change and we want to update our classification.

We need a model which can memorize the information from former inputs.

Output: $y_1 = \text{prob}_{\text{spam}}$, $y_2 = \text{prob}_{\text{spam}}$,



Input: $x_1 = \text{vector}_{\text{word}1}$, $x_2 = \text{vector}_{\text{word}2}$,

Hidden memory State:

$S_1 = \text{state after first word}$

$S_2 = \text{state after second word}$

...

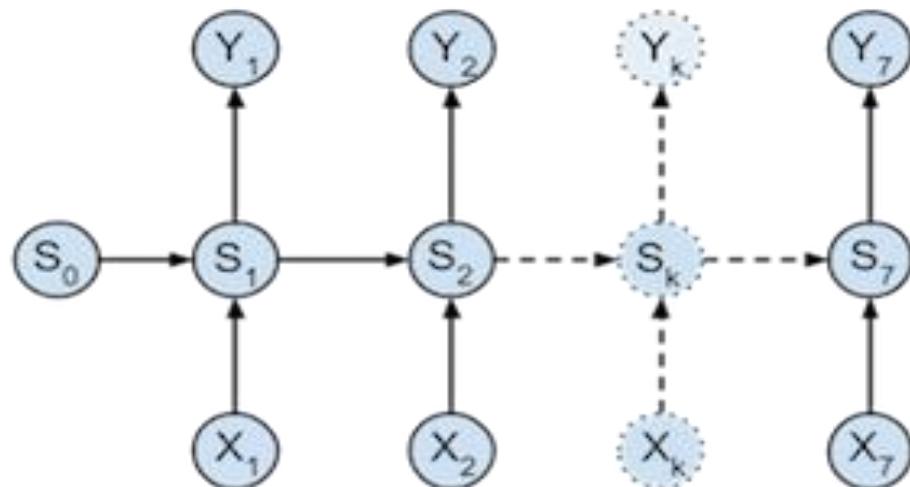
How to choose the dimensions of the hidden state?

The input \mathbf{x} is a vector of the size of our vocabulary.

The output \mathbf{y} is a vector of size 2 (spam/no-spam or passed/failed).

The **hidden state \mathbf{h} (or \mathbf{S})** should have enough capacity to capture different concepts of spams (e.g. fraud, sex, conference) – we could choose a vector of length 3.

Output: $y_1 = (p_1, p_2)_{t1}$ $y_2 = (p_1, p_2)_{t2}$...



Input: $x_1 = \text{vector}_{\text{word}1}$, $x_2 = \text{vector}_{\text{word}2}$, ...

Hidden memory State:

$S_1 = \text{state after first word}$

$S_2 = \text{state after second word}$

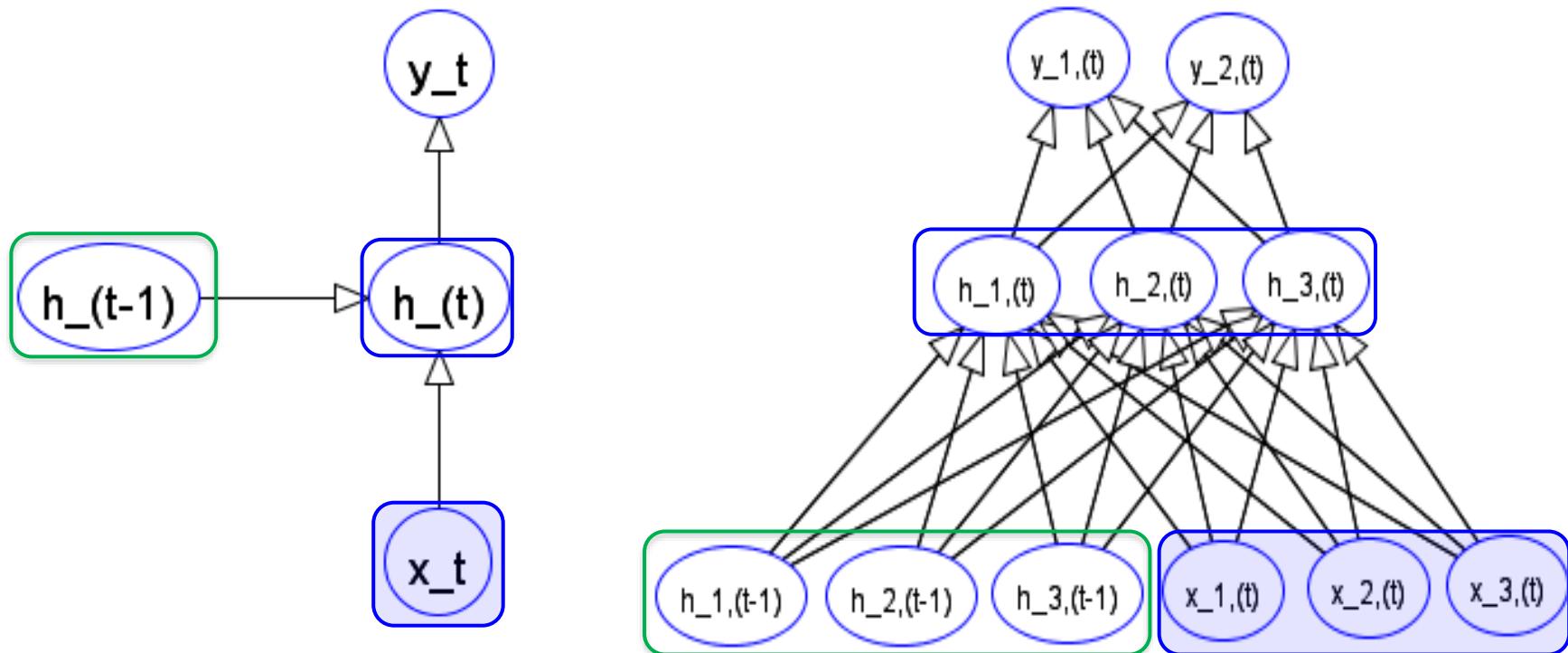
...

Two representations of a RNN at time t

x has 3 components – a very small 1-hot-encoded vocabulary or embedding vector ;-)

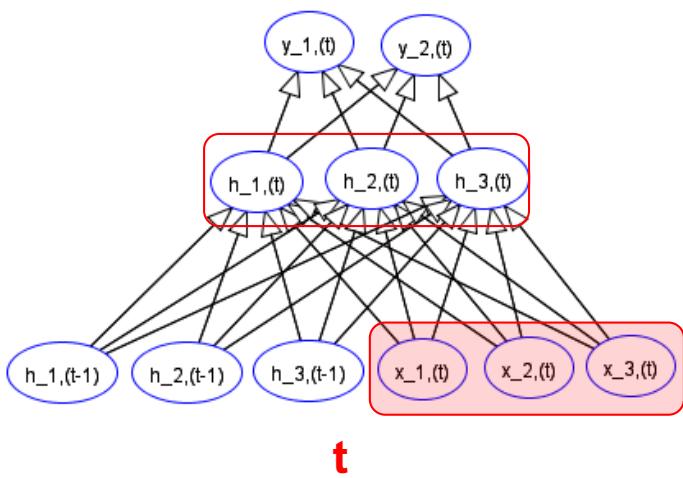
y has 2 components (spam/no-spam for mails, passed/failed for scoring essays).

h 3 components to capture abstract concepts (e.g. fraud/sex/conference for emails) and is initialized at $t=0$ with $(0,0,0)$.

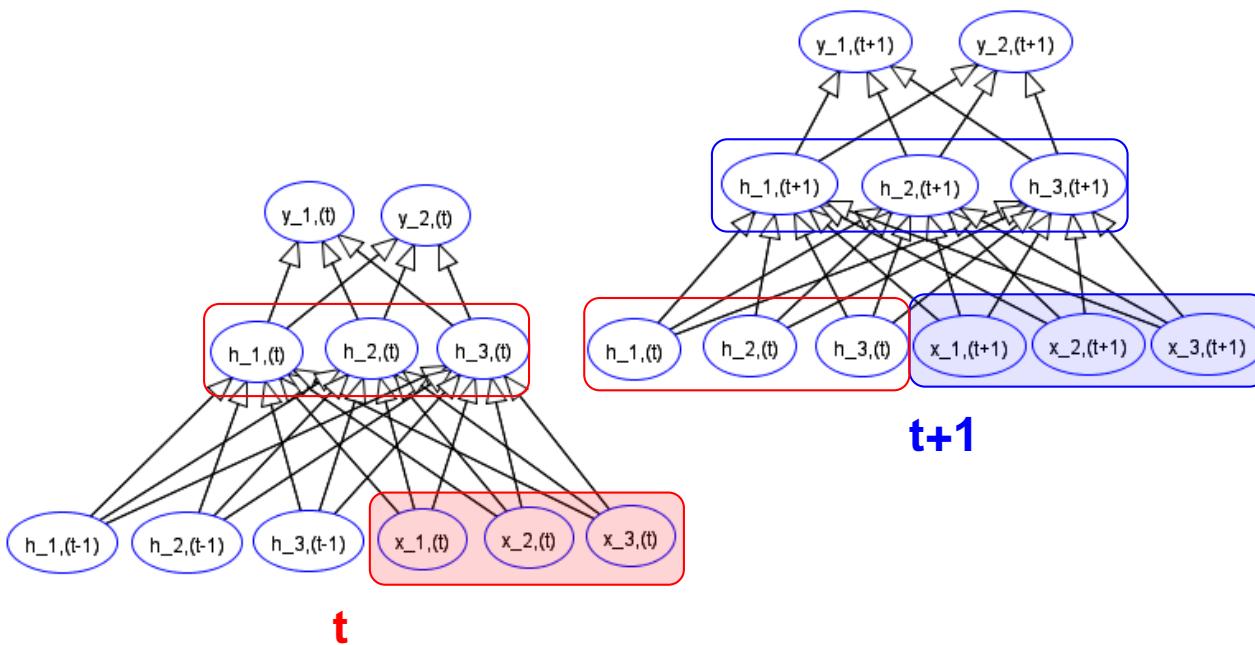


Stepping through an RNN

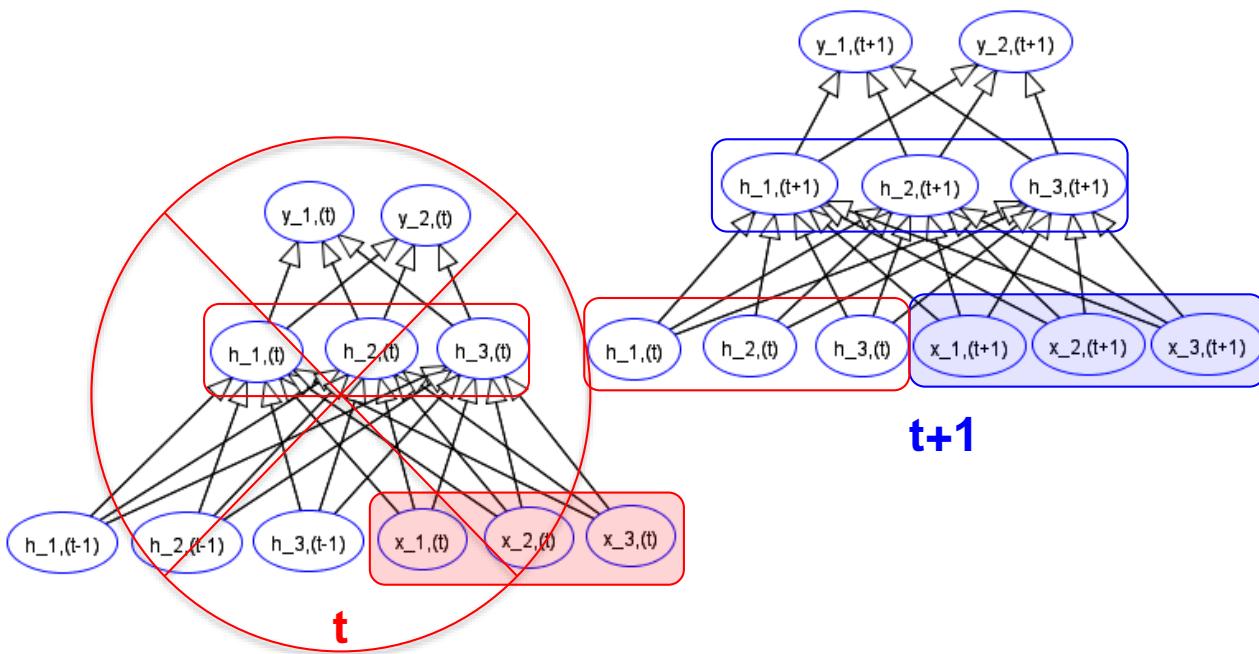
A simple RNN at 3 successive time steps



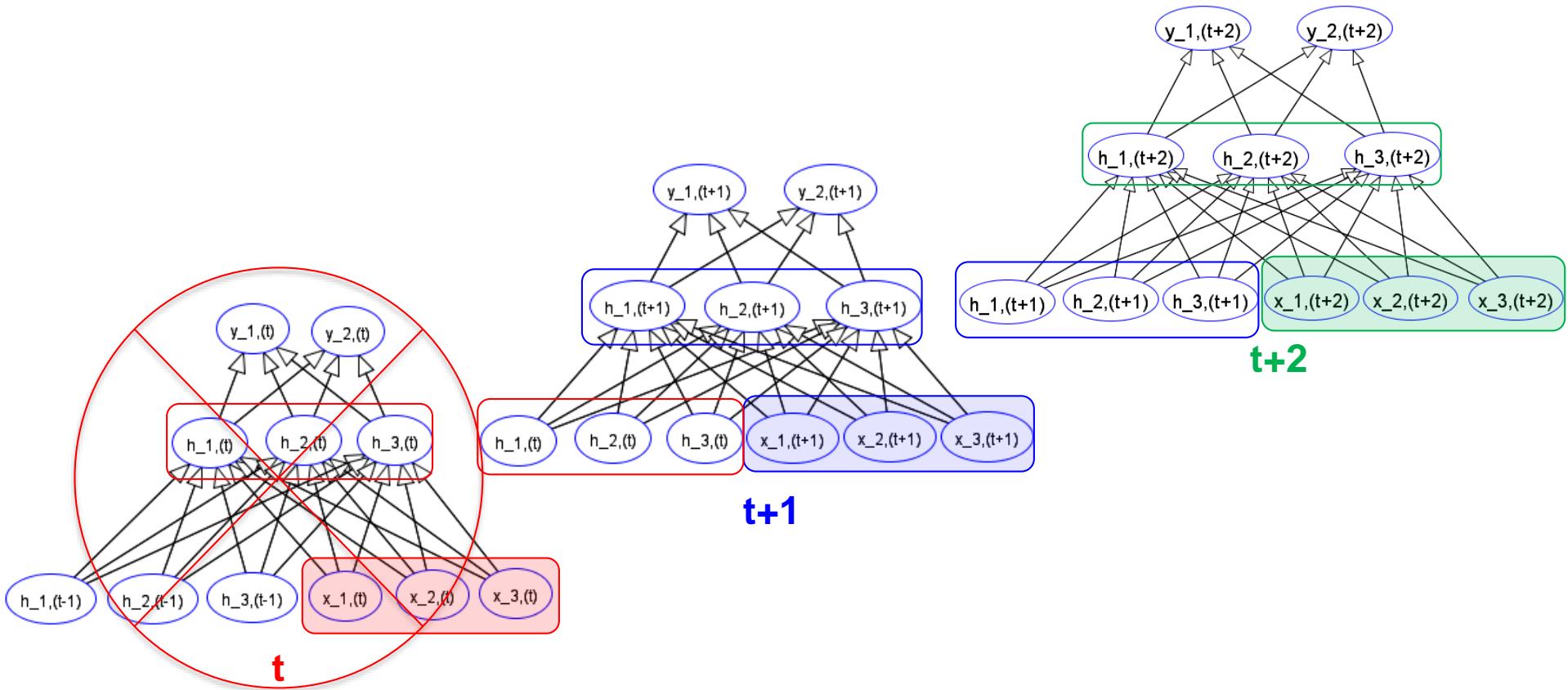
A simple RNN at 3 successive time steps



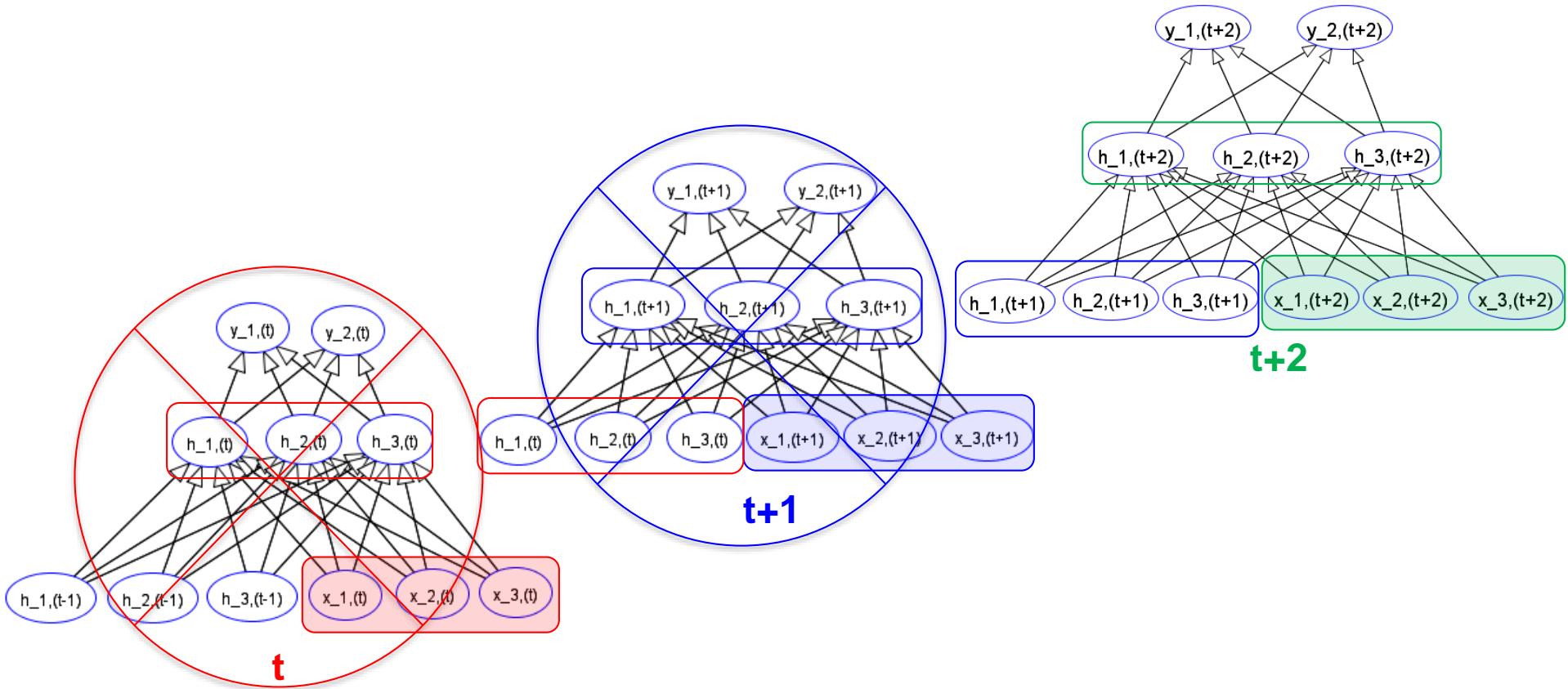
A simple RNN at 3 successive time steps



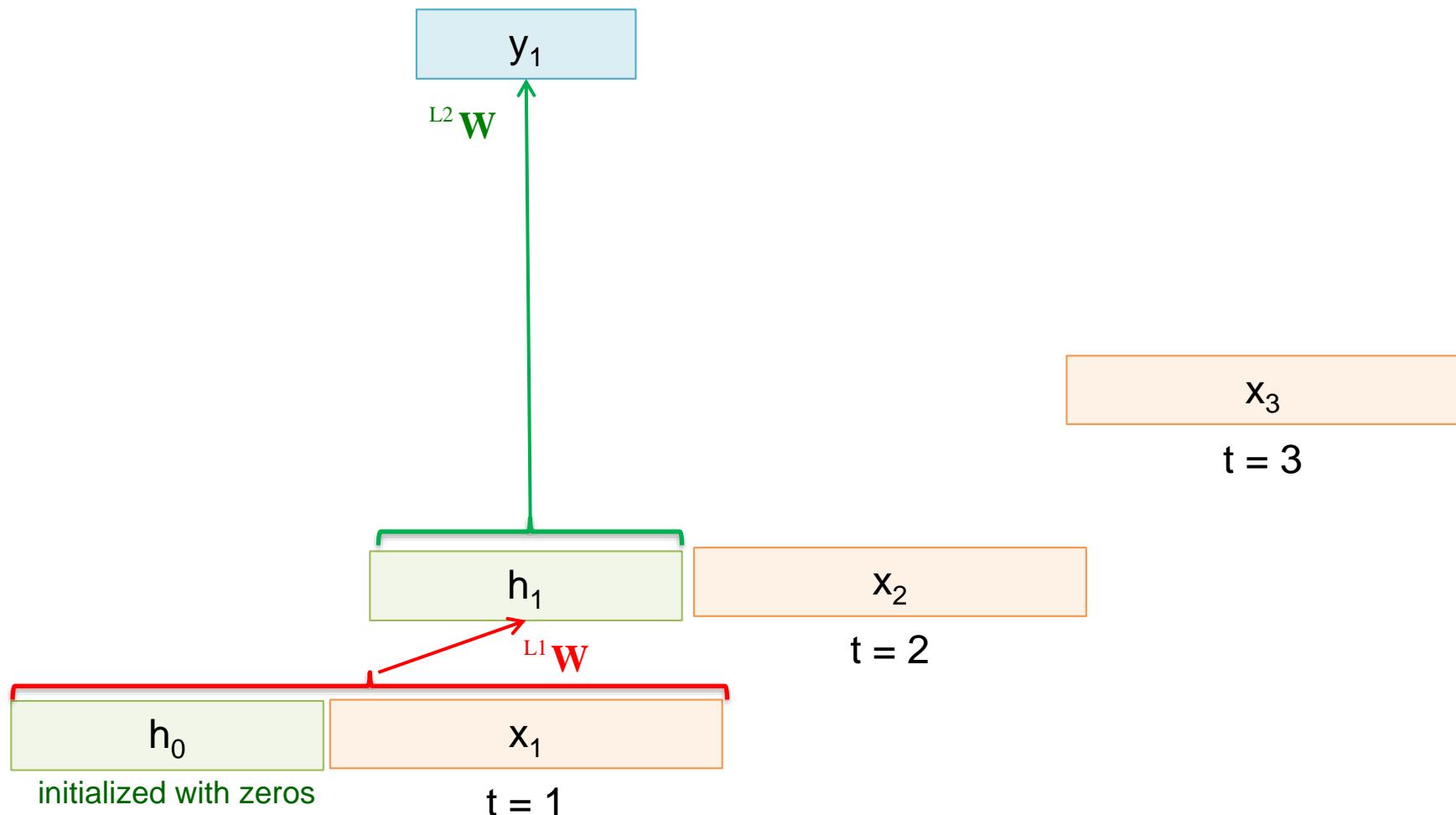
A simple RNN at 3 successive time steps



A simple RNN at 3 successive time steps

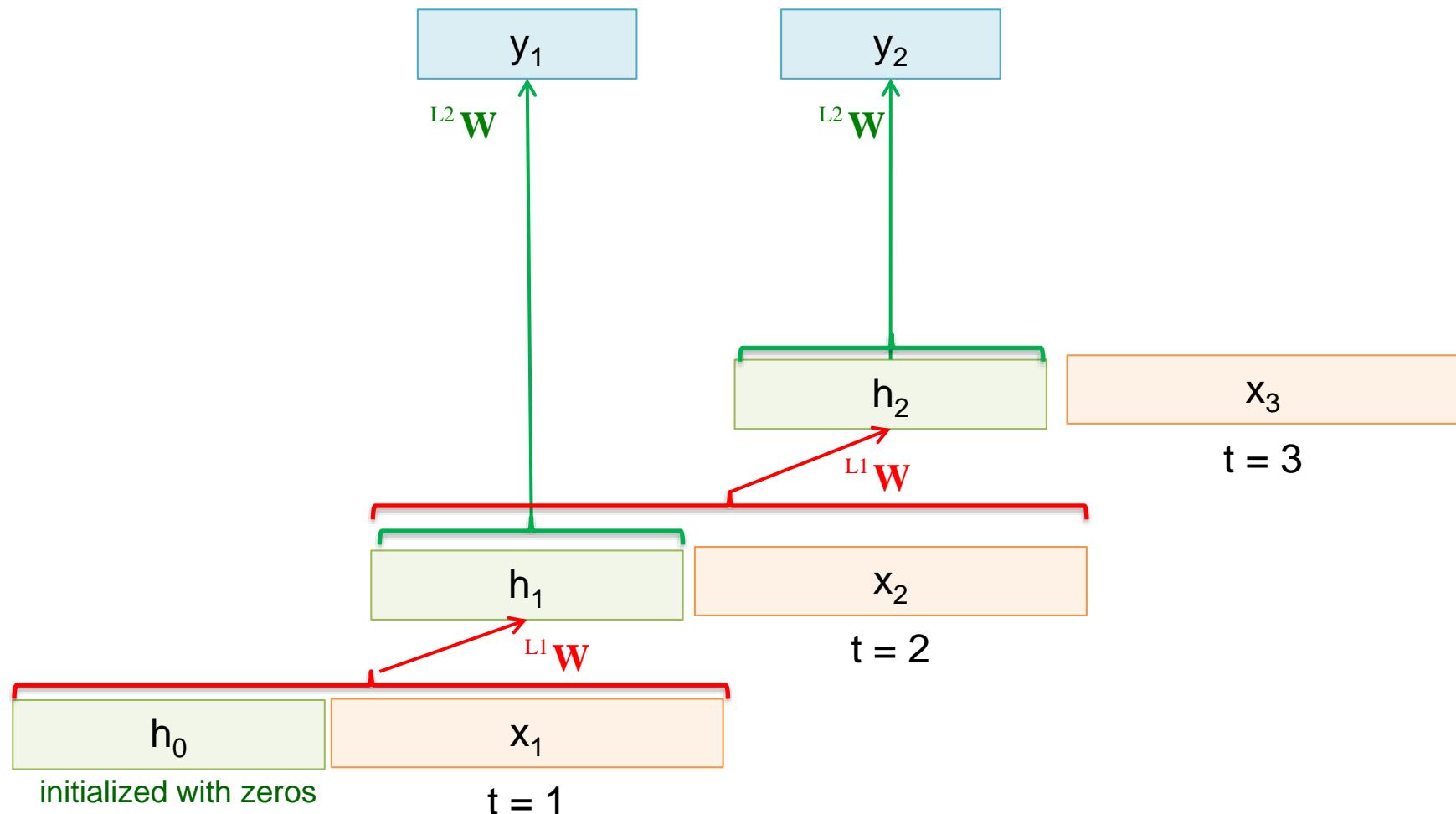


An RNN shares weights across all time steps



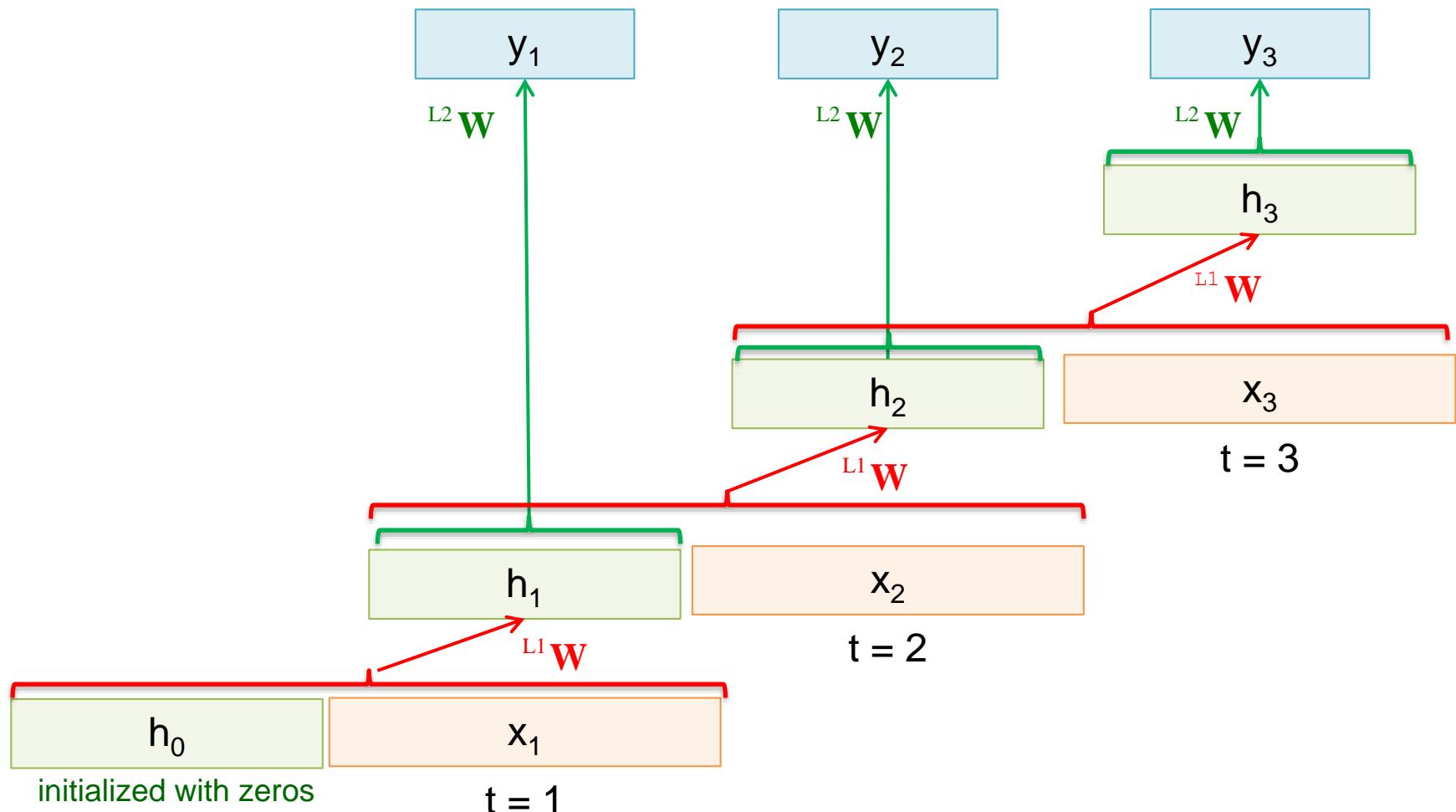
Imagine a trained RNN with fixed weight matrices in layer 1 and layer 2.

An RNN shares weights across all time steps



We use at each time step the same weight matrices between the layers!

An RNN shares weights across all time steps



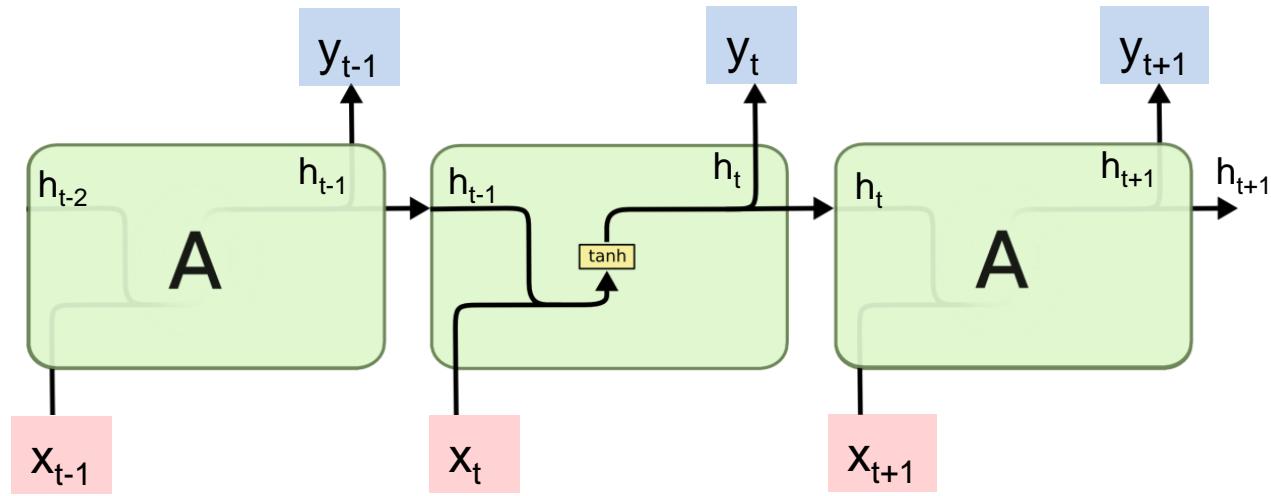
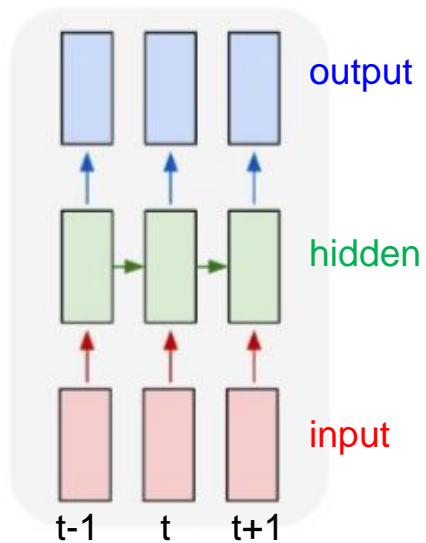
The length of the input sequence can have arbitrary length.

We just reuse (keras: distribute) the same NN for each instance in the sequence!

Therefor it is called recurrent network!!

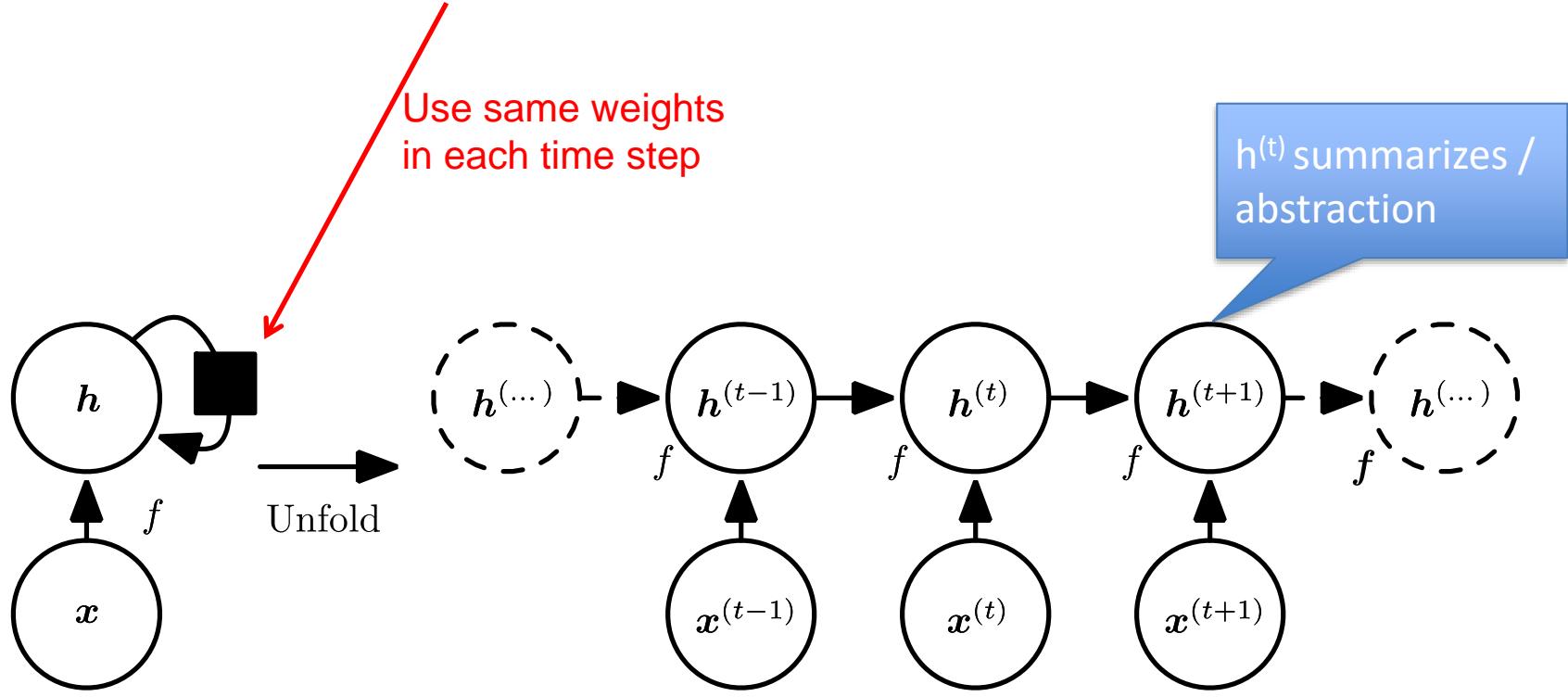
Using diagrams to represent an RNN

many to many



$$\mathbf{A} = f_{\mathbf{W}}(\mathbf{h}_{t-1}, \mathbf{x}_t) = \tanh([\mathbf{h}_{t-1}, \mathbf{x}_t] \cdot \mathbf{W} + \mathbf{b})$$

Using a circuit diagram to represent an RNN

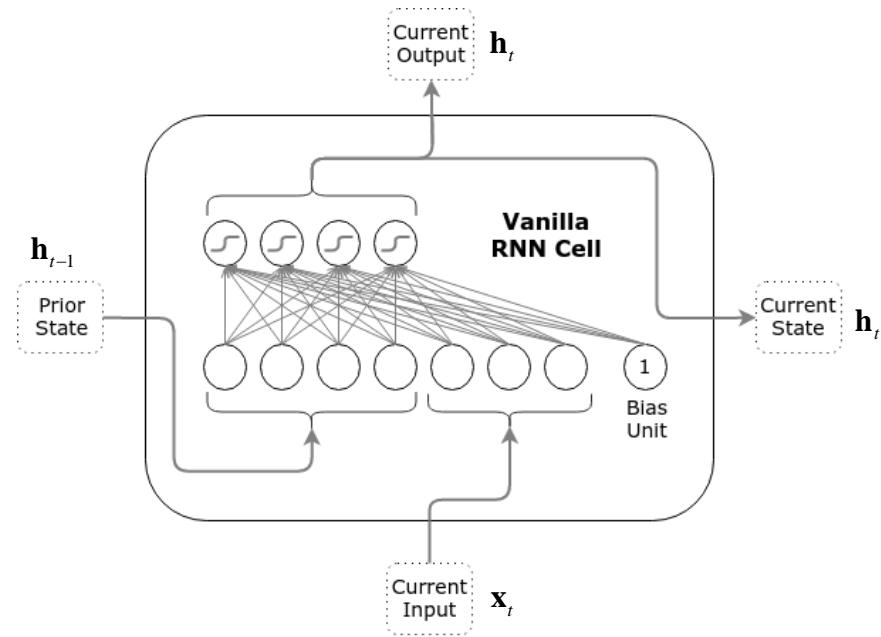
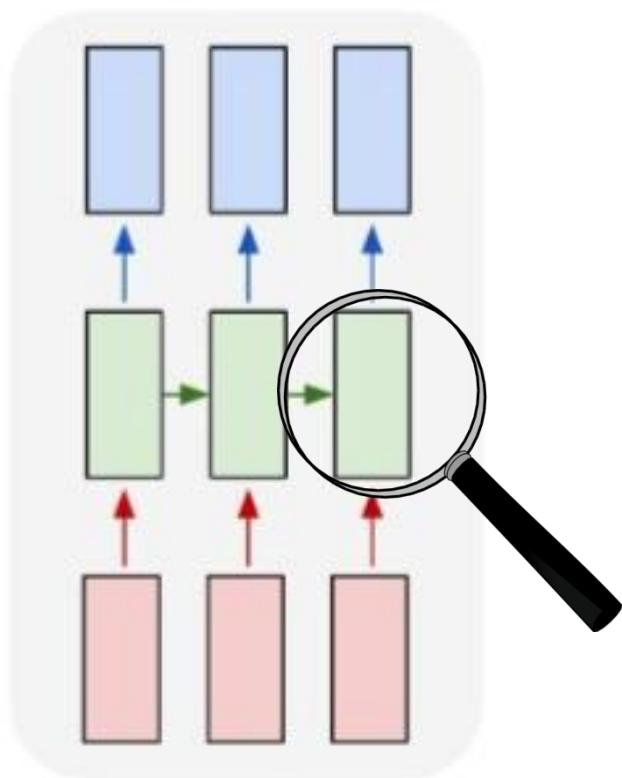


Left: Circuit Diagram (black square delay of one time step)

Right: **Unrolled / unfolded**

Looking into a RNN “cell”

many to many



$$\text{output} = \mathbf{h}_t = \tanh([\mathbf{h}_{t-1}, \mathbf{x}_t] \cdot \mathbf{W} + \mathbf{b})$$

A simple forward pass



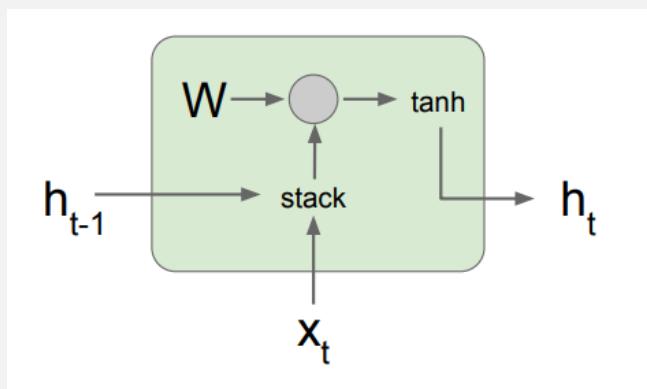
- Given the hidden state at t-1, the input x at t and the weight matrix as:

$$h_{t-1} = \begin{pmatrix} 0, & 1 \end{pmatrix}$$

$$x_t = \begin{pmatrix} 1, & 0 \end{pmatrix}$$

$$W = \begin{pmatrix} 1.5 & -2 \\ 0 & 1 \\ -1 & 0.5 \\ 2 & 0 \end{pmatrix}$$

$$b = \begin{pmatrix} 0 & 0 \end{pmatrix}$$



$$A = f_W(h_{t-1}, x_t) = \tanh([h_{t-1}, x_t] \cdot W + b) \dots$$

- Calculate the activation A of the hidden state h_t at time t.

Solution

$$h \in \mathbb{R}^2 \quad (0, 1)$$

$$x \in \mathbb{R}^2 \quad (1, 0)$$

$$W = \begin{pmatrix} 1.5 & -2 \\ 0 & 1 \\ -1 & \frac{1}{2} \\ 2 & 0 \end{pmatrix}$$

$$(0, 1, 1, 0) \begin{pmatrix} 1.5 & -2 \\ 0 & 1 \\ -1 & \frac{1}{2} \\ 2 & 0 \end{pmatrix} = (-1, 1.5)$$

$$\Rightarrow h = Wh((-1, 1.5)) \approx (-0.76, 0.91)$$

Loss construction in an RNN

Determine the loss contribution of instance 1

mini-batch of size M=8

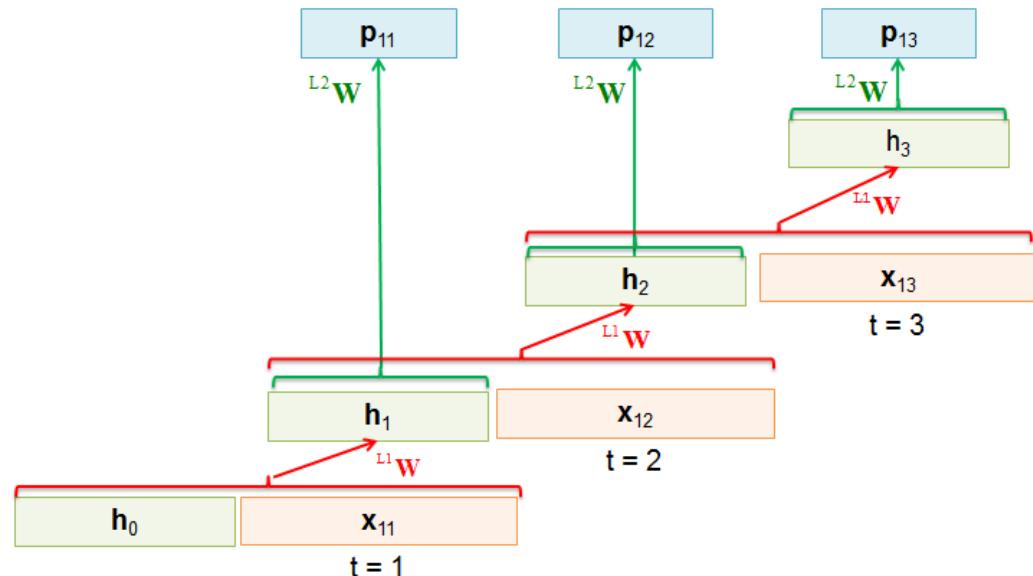
train data input (S=len(seq)=3):

instance_id	seq_t1	seq_t2	seq_t3
1	x_{11}	x_{12}	x_{13}
2	x_{21}	x_{22}	x_{23}
3	x_{31}	x_{32}	x_{33}
⋮	⋮	⋮	⋮
8	x_{81}	x_{82}	x_{83}

train data target (2 classes, K=2):

instance id	y t1	y t2	y t3
1	(1,0)	(1,0)	(0,1)
2	(0,1)	(1,0)	(0,1)
3	(0,1)	(0,1)	-1
⋮	⋮	⋮	⋮
8	(1,0)	(1,0)	(1,0)

instance 1:



x-entropy is used to determine distance between 2-dim p-vector and 2-dim y-vector at each of the 3 positions in the sequence:

$$\text{Loss_1} = \sum_{s=1}^3 \left(-\sum_{k=1}^2 y_{1sk} \cdot \log(p_{1sk}) \right)$$

Determine the loss contribution of instance 2

mini-batch of size M=8

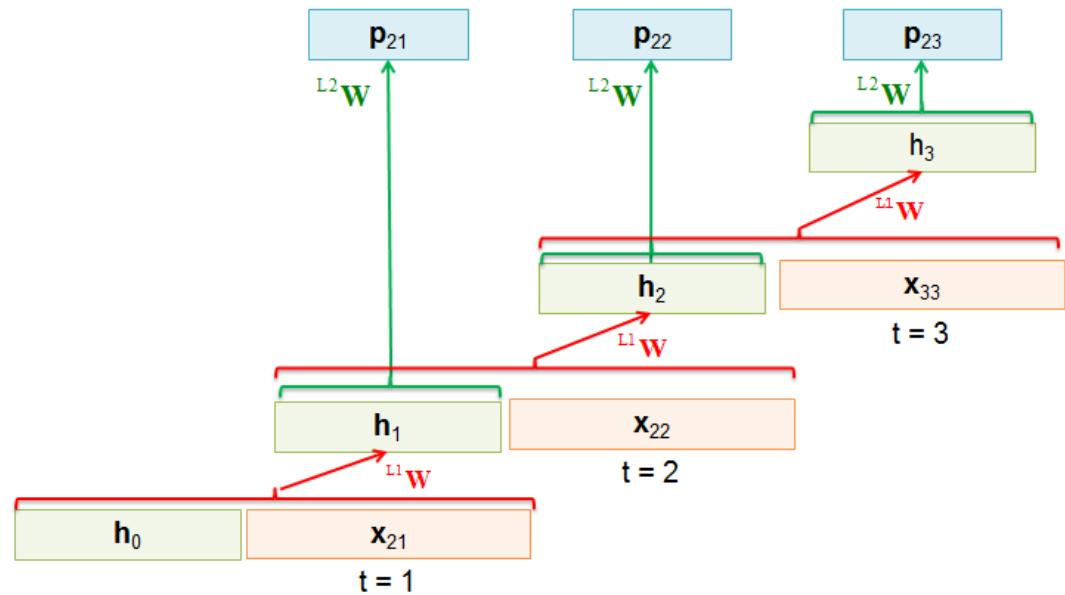
train data input (S=len(seq)=3):

instance_id	seq_t1	seq_t2	seq_t3
1	x ₁₁	x ₁₂	x ₁₃
2	x ₂₁	x ₂₂	x ₂₃
3	x ₃₁	x ₃₂	x ₃₃
⋮	⋮	⋮	⋮
8	x ₈₁	x ₈₂	x ₈₃

train data target (2 classes, K=2):

instance_id	y_t1	y_t2	y_t3
1	(1,0)	(1,0)	(0,1)
2	(0,1)	(1,0)	(0,1)
3	(0,1)	(0,1)	-1
⋮	⋮	⋮	⋮
8	(1,0)	(1,0)	(1,0)

instance 2:



x-entropy is used to determine distance between 2-dim p-vector and 2-dim y-vector at each of the 3 positions in the sequence:

$$\text{Loss_2} = \sum_{s=1}^3 \left(-\sum_{k=1}^2 y_{2sk} \cdot \log(p_{2sk}) \right)$$

Determine the loss of the whole mini-batch

mini-batch of size M=8

train data input (S=len(seq)=3):

instance_id	seq_t1	seq_t2	seq_t3
1	x_{11}	x_{12}	x_{13}
2	x_{21}	x_{22}	x_{23}
3	x_{31}	x_{32}	x_{33}
⋮	⋮	⋮	⋮
8	x_{81}	x_{82}	x_{83}

train data target (2 classes, K=2):

instance_id	y_t1	y_t2	y_t3
1	(1,0)	(1,0)	(0,1)
2	(0,1)	(1,0)	(0,1)
3	(0,1)	(0,1)	-1
⋮	⋮	⋮	⋮
8	(1,0)	(1,0)	(1,0)

Cost C or Loss is given by the cross-entropy averaged over all instances in mini-batch:

$$\text{Loss} = \frac{1}{8} \sum_{m=1}^8 \text{Loss}_m$$

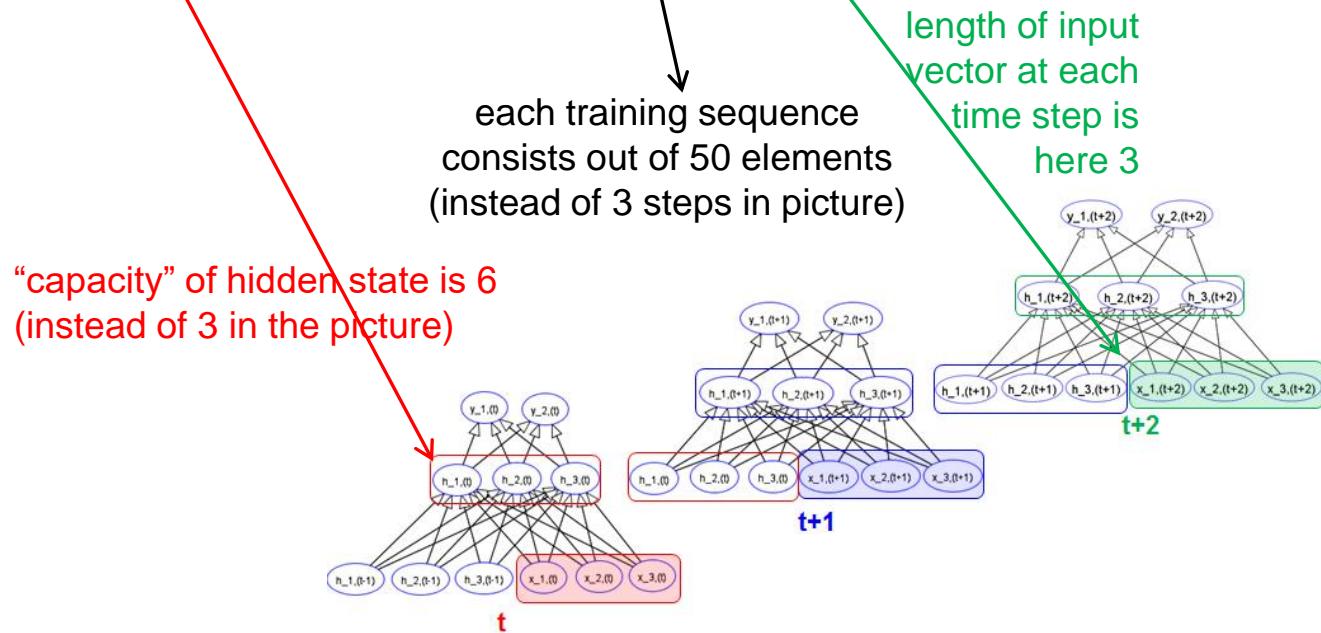
$$\text{Loss} = \frac{1}{8} \sum_{m=1}^8 \left[\sum_{s=1}^3 \left(-\sum_{k=1}^2 y_{msk} \cdot \log(p_{msk}) \right) \right]$$

Based on the mini-batch loss the weights in the two weight matrices of layer 1 and layer 2 are updated.

RNN in Keras

```
from keras.layers import Activation, Dense, SimpleRNN, TimeDistributed
```

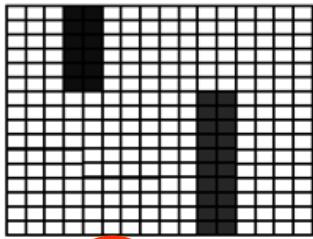
```
model = keras.models.Sequential()  
  
model.add(SimpleRNN(6, batch_input_shape=(None, 50, 3), return_sequences=True))  
model.add(TimeDistributed(Dense(2)))  
model.add(Activation('softmax'))  
  
model.compile(loss='categorical_crossentropy', optimizer='rmsprop', metrics=['accuracy'])
```



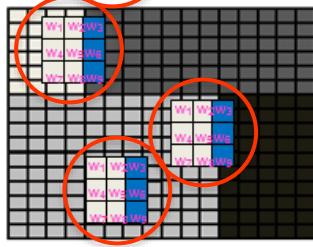
Common tricks in RNN & CNN and some differences

CNN and Recurrent Network share weights

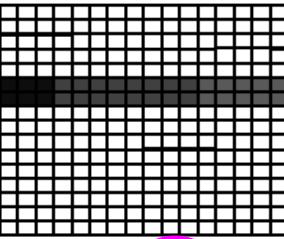
activation map



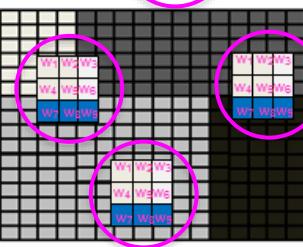
filter 1



activation map

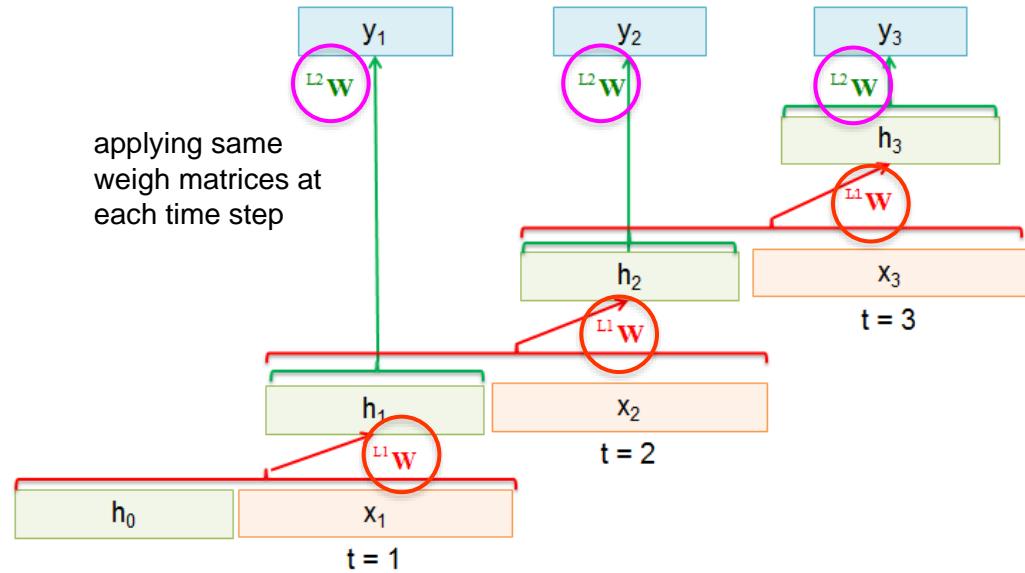


filter 2



apply filter with same
weights at each position

applying same
weigh matrices at
each time step

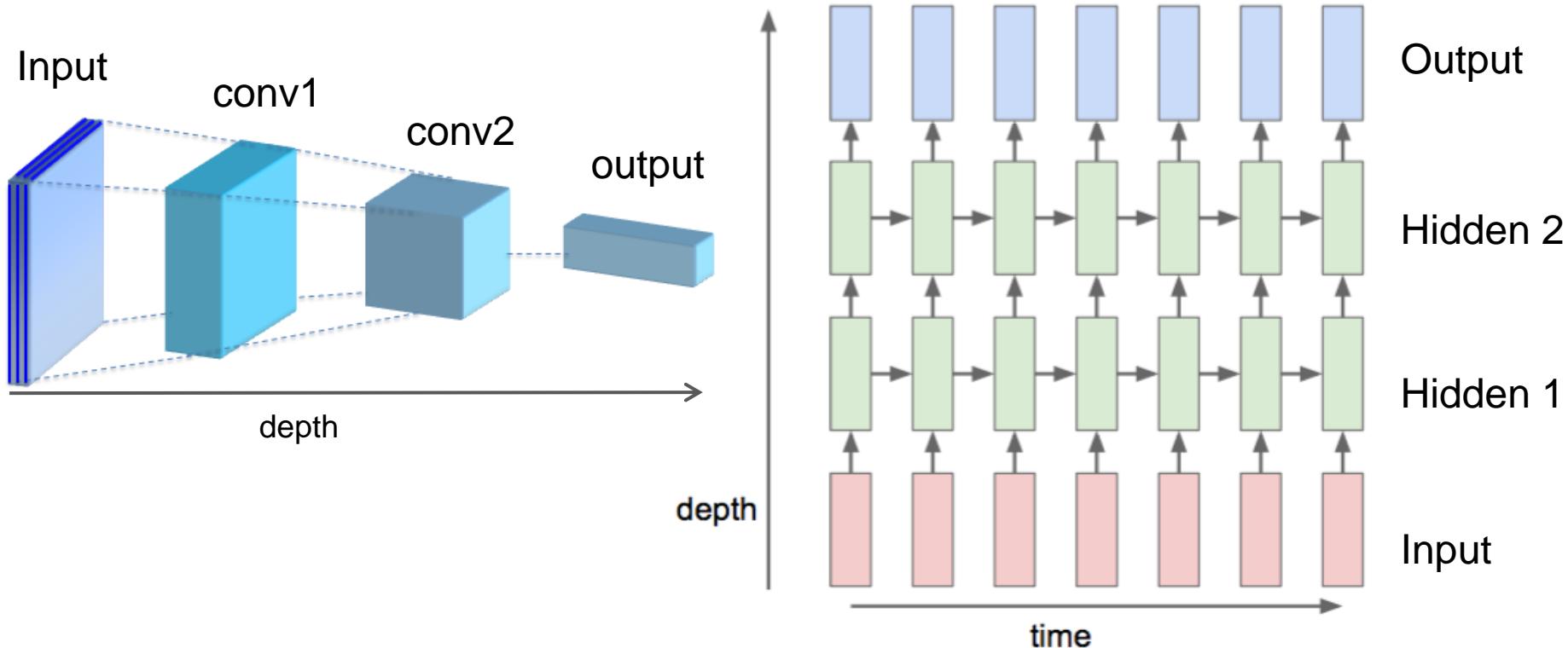


CNN share weights between
different local regions of the image

RNN share weights between time steps

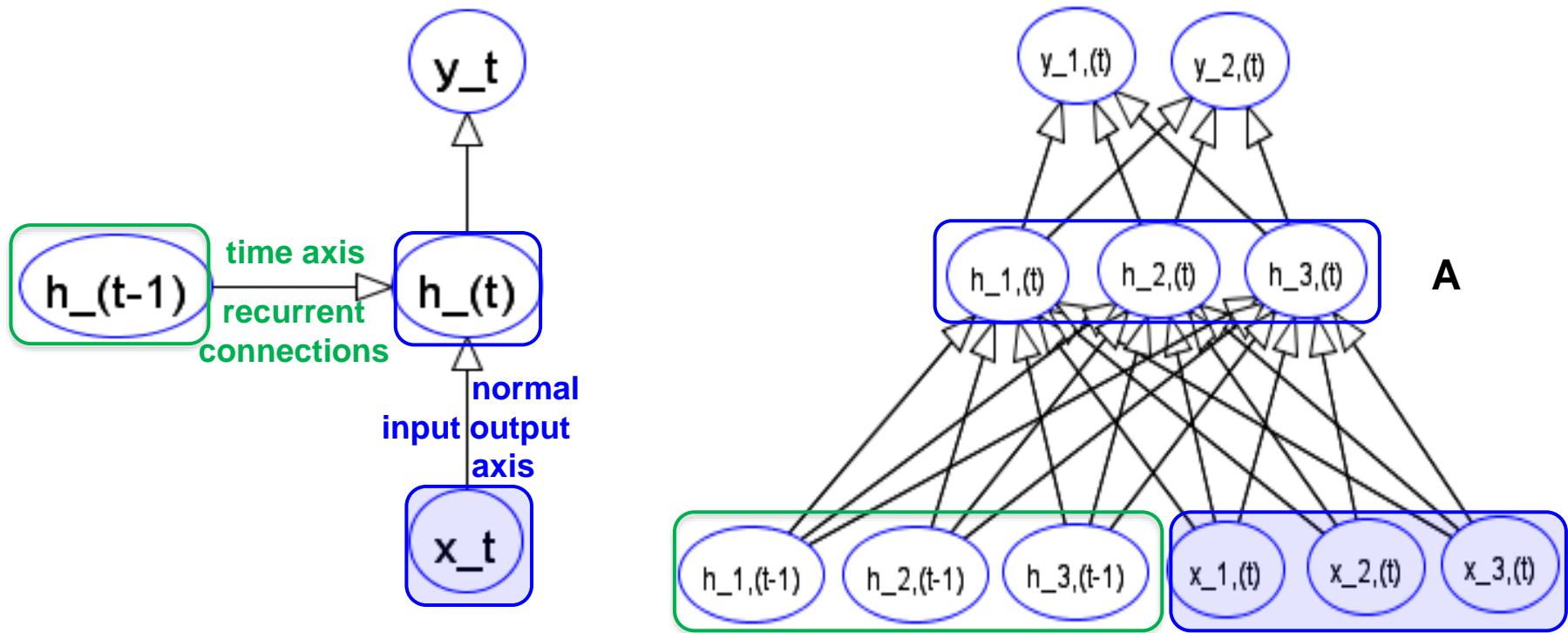
Remark: no weight sharing in fully connected NN

Also in RNN we can go deep for hierarchical features



Usually we see only 1-4 hidden layers in an RNN compared to usually 4-100 stacked hidden convolutional blocks in CNNs.

Dropout in recurrent architectures allow to choose different dropout rates for recurrent and normal connections



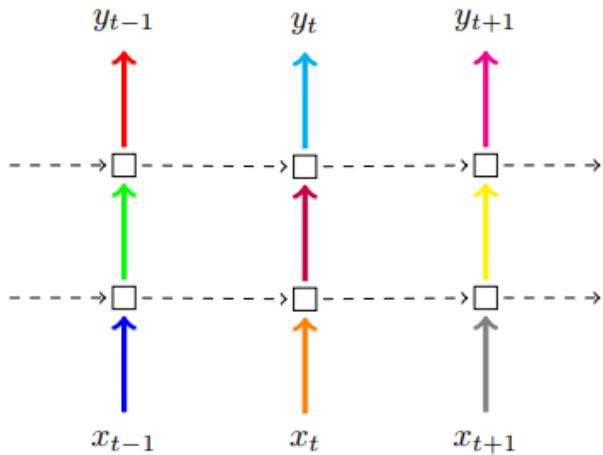
$$\mathbf{A} = f_{\mathbf{W}}(\mathbf{h}_{t-1}, \mathbf{x}_t) = \tanh([\mathbf{h}_{t-1}, \mathbf{x}_t] \cdot \mathbf{W} + \mathbf{b}) = \tanh(\mathbf{h}_{t-1} \cdot \mathbf{W}_h + \mathbf{x}_t \cdot \mathbf{W}_x + \mathbf{b})$$

$$\mathbf{W} = \begin{pmatrix} \mathbf{W}_h \\ \mathbf{W}_x \end{pmatrix}$$

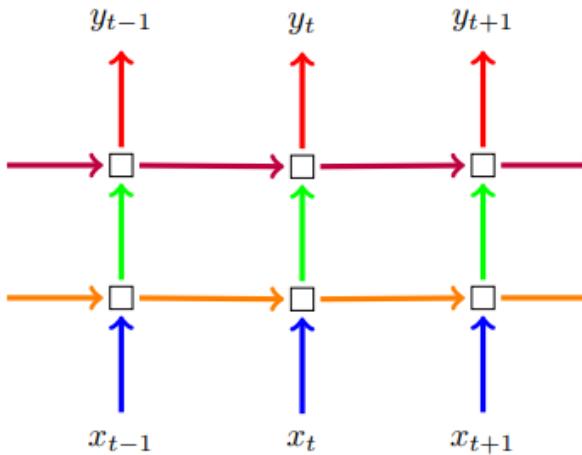
Dimensions in example: $\mathbf{W}:6 \times 3$, $\mathbf{W}_h:3 \times 3$, $\mathbf{W}_x:3 \times 3$

Dropout in recurrent architectures

It is important to **use identical dropout masks** (marked by arrows with same color) **at different time steps** in recurrent architectures like GRU or LSTM.



(a) Naive dropout RNN



(b) Variational RNN

same arrow
color indicates
identical dropout
mask

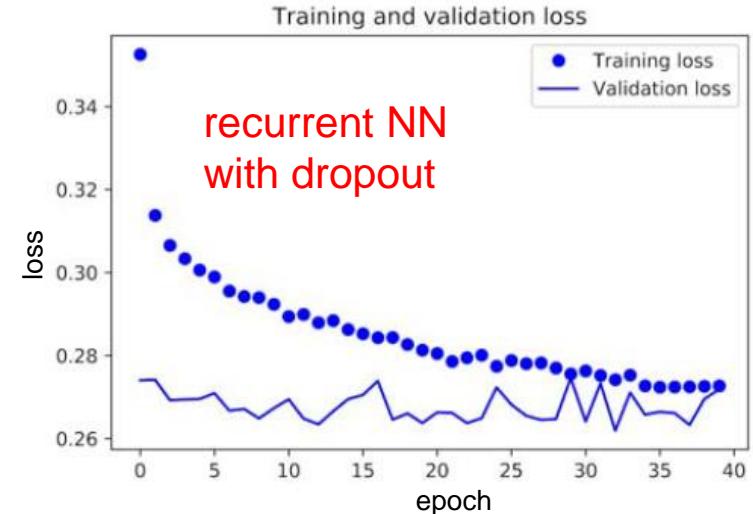
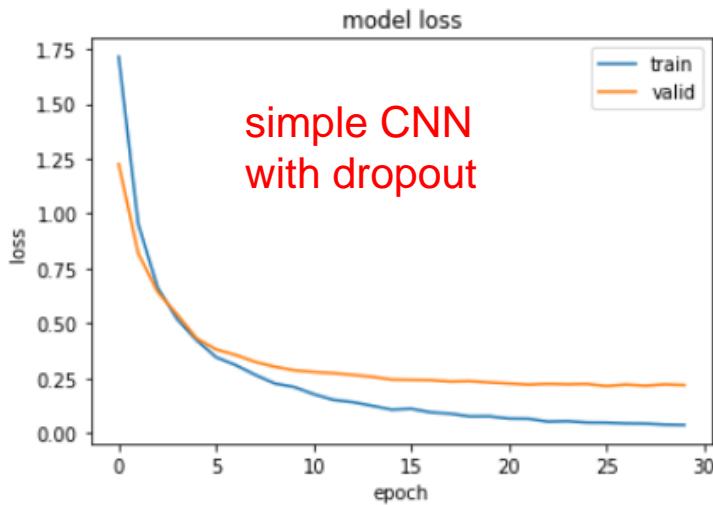
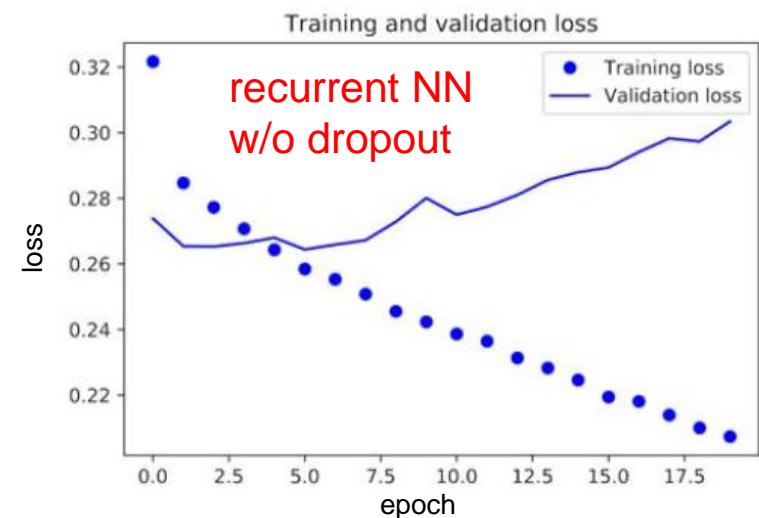
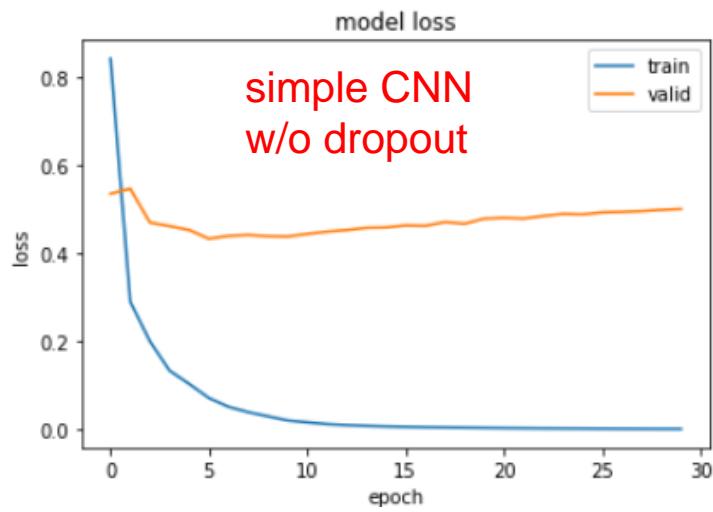
Figure 1: **Depiction of the dropout technique following our Bayesian interpretation (right) compared to the standard technique in the field (left).** Each square represents an RNN unit, with horizontal arrows representing time dependence (recurrent connections). Vertical arrows represent the input and output to each RNN unit. Coloured connections represent dropped-out inputs, with different colours corresponding to different dropout masks. Dashed lines correspond to standard connections with no dropout. Current techniques (naive dropout, left) use different masks at different time steps, with no dropout on the recurrent layers. The proposed technique (Variational RNN, right) uses the same dropout mask at each time step, including the recurrent layers.

[Gal2016](#)

In keras:

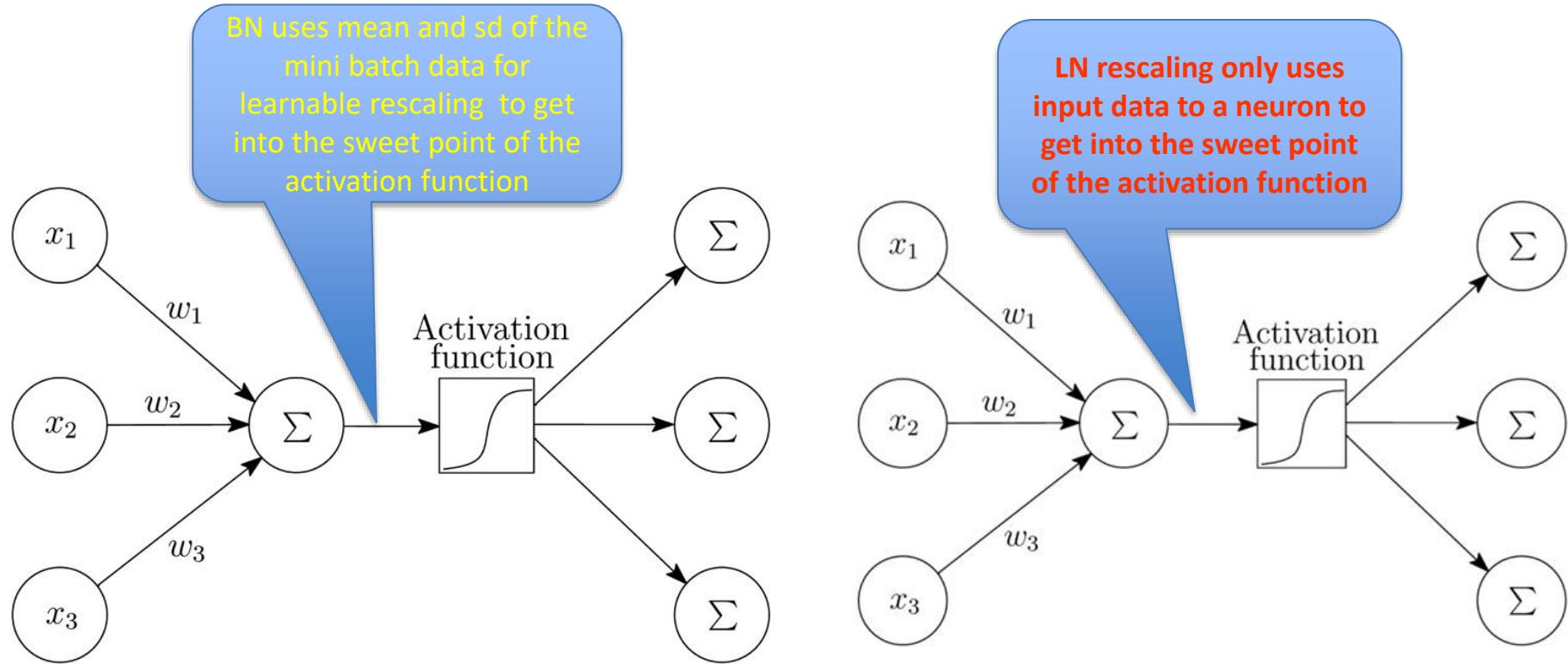
```
model.add(layers.GRU(32, dropout=0.2, recurrent_dropout=0.2, input_shape=(None, ...)))
```

Dropout can fight overfitting in CNN and recurrent NN



Batchnormalization is crucial to train deep CNNs

Layernormalization is beneficial in RNN: LN \neq BN



Applying BN to RNN would not take into account the recurrent architecture of the NN over which statistics of the input to a neuron might change considerable within the same mini batch. In LN the mean and variance from all of the summed inputs to the neurons in a layer on a single training case are used for normalization .

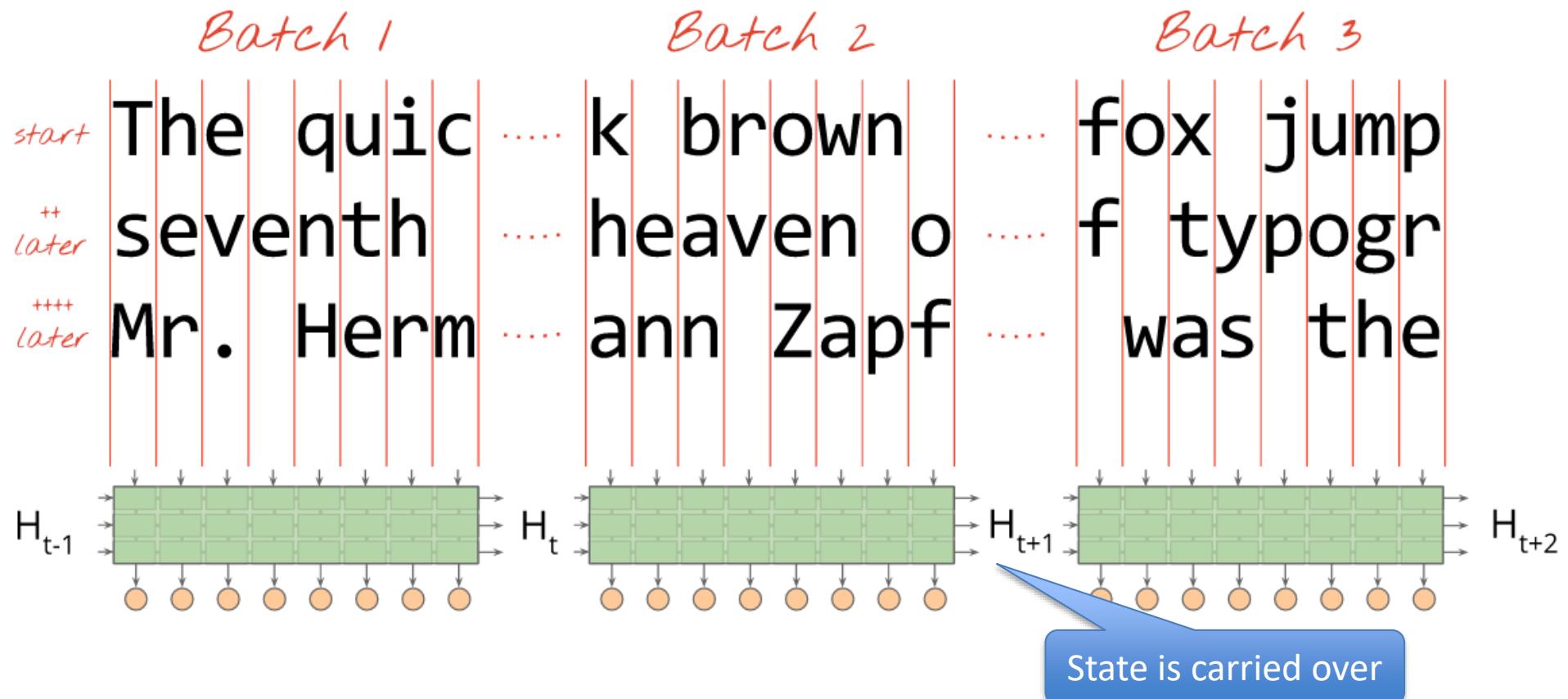
Stateful RNN model

Training a stateful RNNs

- RNN are often trained on sequence data with inherent order
- Sequences are often very long and need to be cut between mini-batches
- By default the hidden state is initialized with zeros in each mini-batch
- In stateful RNN we connect sequences in the right order between mini-batches allowing to make use of the hidden state learned so far
- This requires a careful construction of the mini-batches and an appropriate transfer of the hidden state between mini-batches

Mini-batches in statefull RNN

The gradient is propagated back a fixed amount of steps defined by the size of a mini-batch. In stateful RNNs the hidden state is carried over between mini-batches and hence between connecting sequences given appropriate batches.



Vanishing/Exploding Gradient problem during training a RNN

Recall: Loss of a mini-batch is used to determine update

mini-batch of size M=8

train data input (S=len(seq)=3):

instance_id	seq_t1	seq_t2	seq_t3
1	x_{11}	x_{12}	x_{13}
2	x_{21}	x_{22}	x_{23}
3	x_{31}	x_{32}	x_{33}
⋮	⋮	⋮	⋮
8	x_{81}	x_{82}	x_{83}

train data target (2 classes, K=2):

instance_id	y_t1	y_t2	y_t3
1	(1,0)	(1,0)	(0,1)
2	(0,1)	(1,0)	(0,1)
3	(0,1)	(0,1)	-1
⋮	⋮	⋮	⋮
8	(1,0)	(1,0)	(1,0)

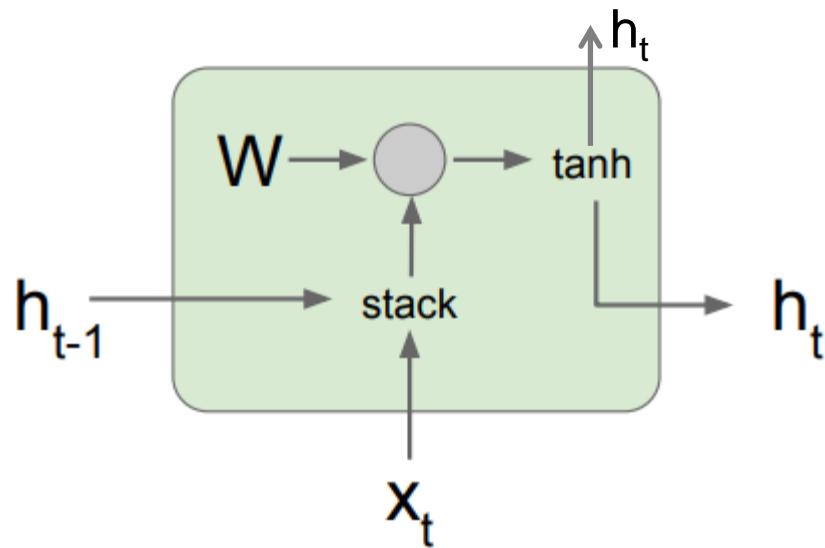
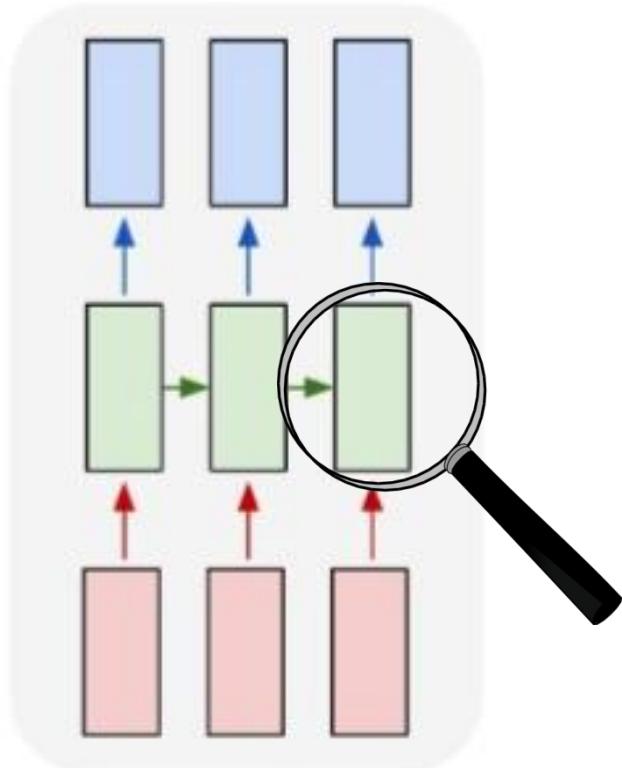
Cost C or Loss is given by the cross-entropy averaged over all instances in mini-batch:

$$\text{Loss} = \frac{1}{8} \sum_{m=1}^8 \left[\sum_{s=1}^3 \left(- \sum_{k=1}^2 y_{\text{msk}} \cdot \log(p_{\text{msk}}) \right) \right]$$

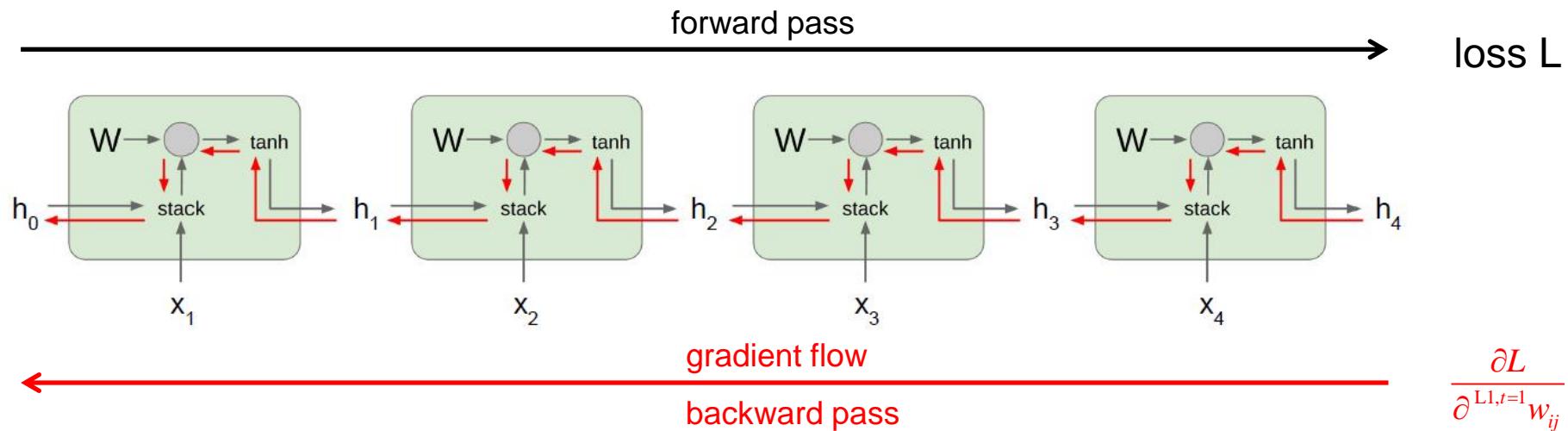
Based on the mini-batch loss the weights in the tow weight matrices of layer 1 and layer 2 are updated.

Recall: Design of a RNN “cell”

many to many



Backpropagation in RNNs: Gradient is multiplied at each time step with same factor: Gradient explosion/vanishing



Propagating the gradient of the cost function via chain rule to the first time point involves multiplying at each time step with W^T (and the derivation of tanh).

⇒ Vanishing gradient if we multiply at each time step with a number < 1
(more precisely we have only a number if W is a scalar, otherwise we need to look on the first singular value of W^T)

⇒ Exploding gradient if we multiply at each time step with a number > 1
(more precisely we have only a number if W is a scalar, otherwise we need to look on the first singular value of W^T)

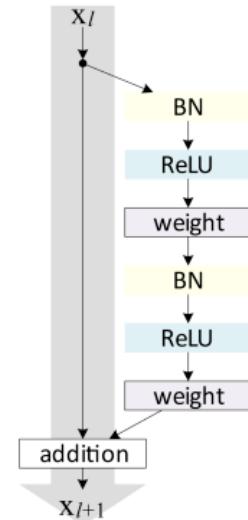
Solution: gradient clipping (hack), or use better architecture like LSTM or GRU!

GRU and LSTM cells to avoid vanishing/exploding gradients

Recall: ResNet

- use ResNet like architectures allowing for a gradient highway

(in CNN also batch-normalization and ReLU helped to train deep NN, but cannot naively transferred to recurrent NN)



$$x_{l+1} = x_l + F(x_l)$$

ResNet basic design (VGG-style)

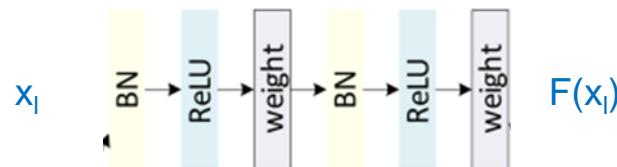
- add shortcut connections every two
- all 3x3 conv (almost)

152 layers:
Why does this train at all?

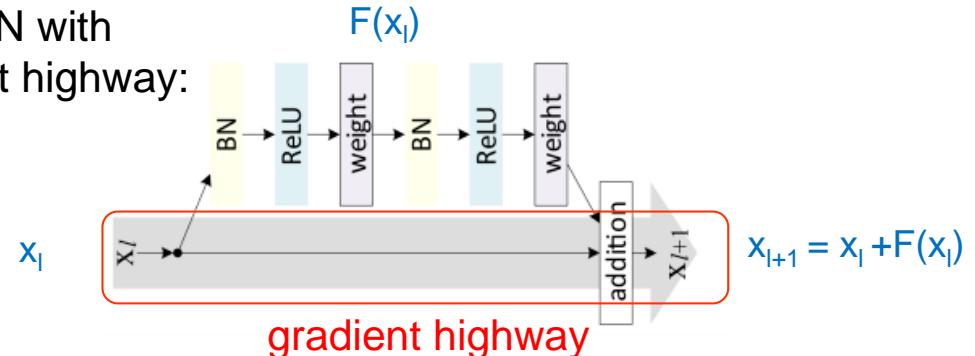
This deep architecture
could still be trained, since
the gradients can skip
layers which diminish the
gradient!

Provide gradient highway also in recurrent NN: GRU, LSTM

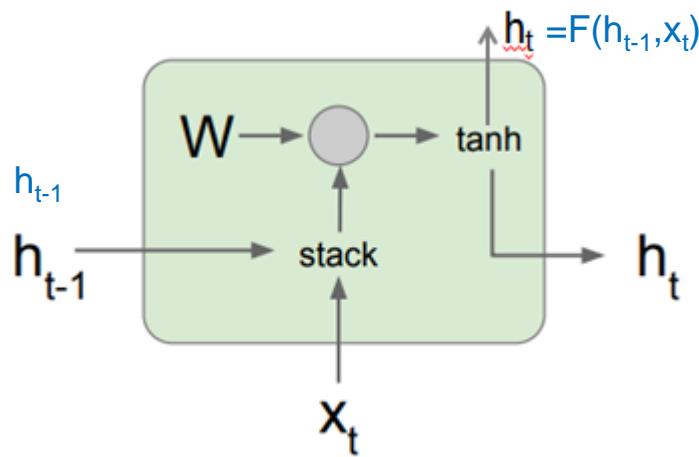
CNN classic:



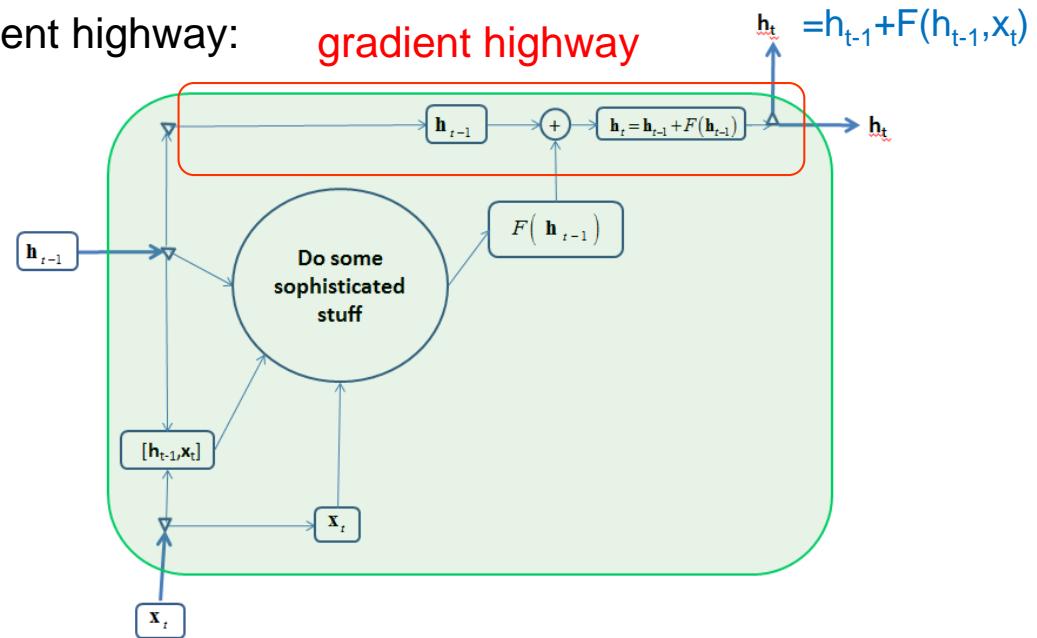
CNN with gradient highway:



RNN classic:

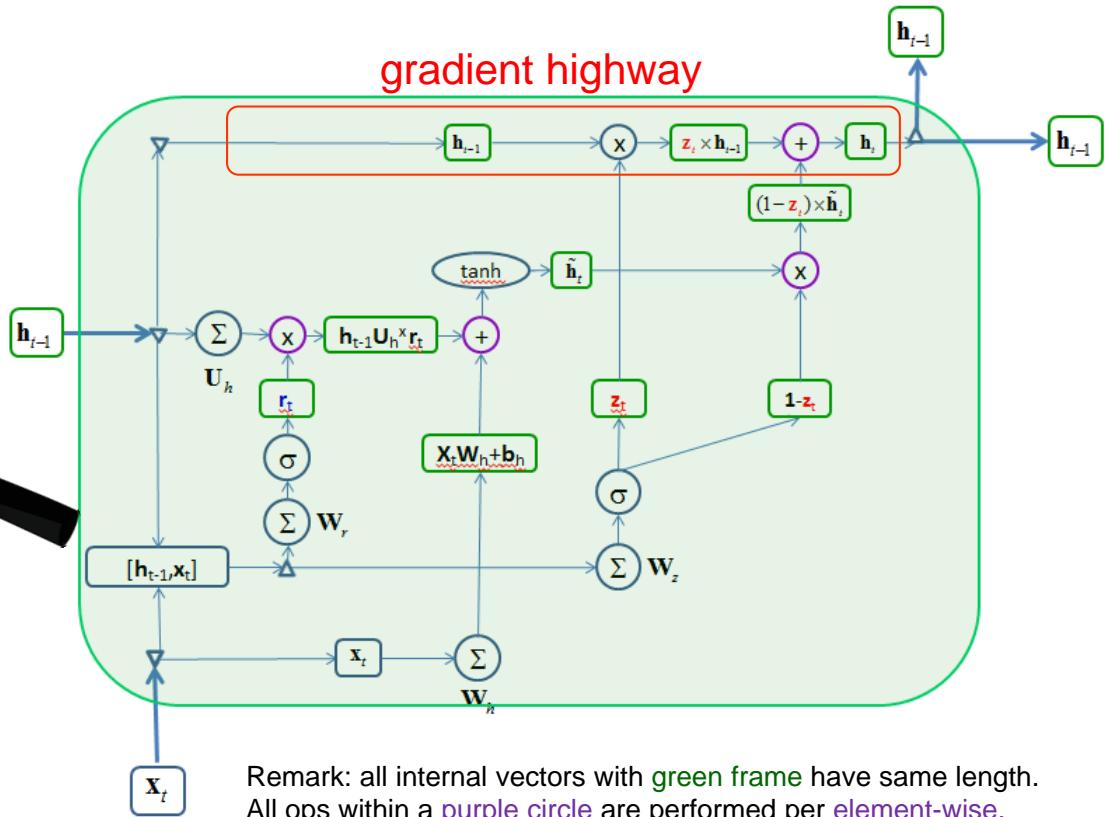
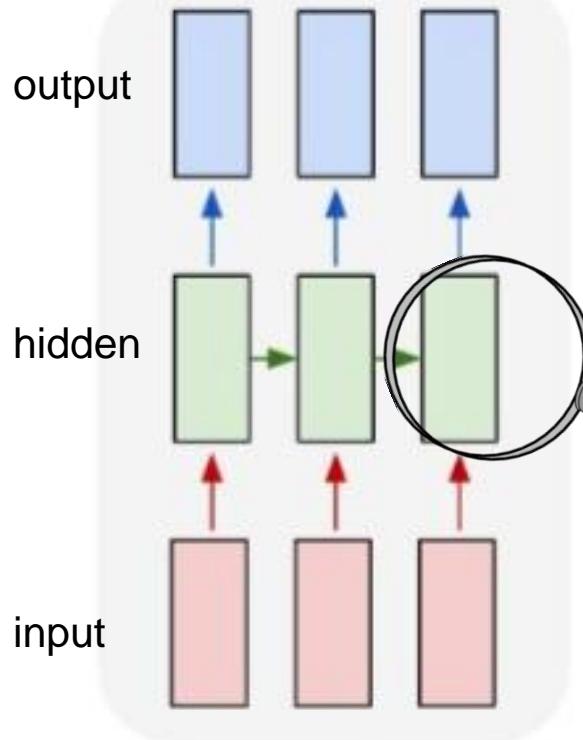


RNN with gradient highway:



Solution via “highway allowing” architecture: GRU

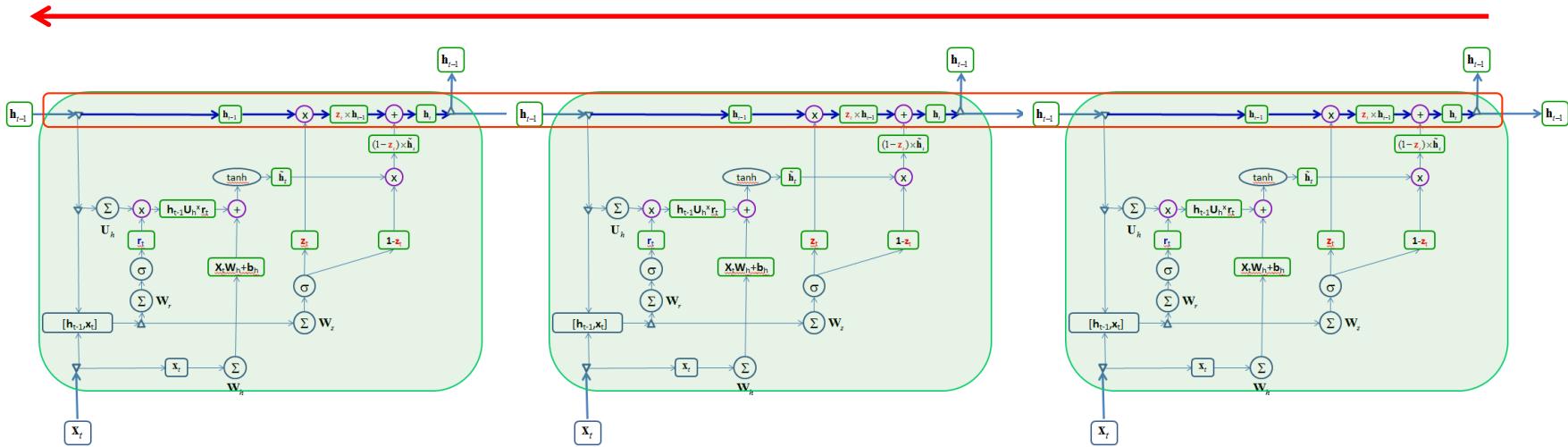
many to many



The gradient highway avoids gradient vanishing. The GRU also avoids gradient explosion since the element-wise operations on vector-elements that change over the time steps, avoids multiplying the gradients with the same number in each step.

The Gated Recurrent Unit (GRU): Gradient Flow

Uninterrupted gradient flow!



Relevant gate:

$$\mathbf{r}_t = \sigma([\mathbf{h}_{t-1}, \mathbf{x}_t] \cdot \mathbf{W}_r + \mathbf{b}_r)$$

Update gate:

$$\mathbf{z}_t = \sigma([\mathbf{h}_{t-1}, \mathbf{x}_t] \cdot \mathbf{W}_z + \mathbf{b}_z)$$

Proposed hidden state: $\tilde{\mathbf{h}}_t = \tanh(\mathbf{x}_t \cdot \mathbf{W}_h + \mathbf{b}_h + \mathbf{h}_{t-1} \mathbf{U}_h \otimes \mathbf{r}_t)$

New hidden state is:

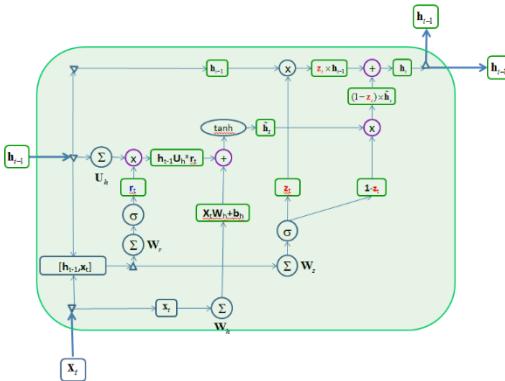
$$\mathbf{h}_t = (1 - \mathbf{z}_t) \otimes \tilde{\mathbf{h}}_t + \mathbf{z}_t \otimes \mathbf{h}_{t-1}$$

GRU in keras

```
from keras.models import Sequential
from keras import layers
from keras.optimizers import RMSprop

model = Sequential()
model.add(layers.GRU(32, input_shape=(None, float_data.shape[-1])))
model.add(layers.Dense(1))

model.compile(optimizer=RMSprop(), loss='mae')
history = model.fit_generator(train_gen,
                               steps_per_epoch=500,
                               epochs=20,
                               validation_data=val_gen,
                               validation_steps=val_steps)
```



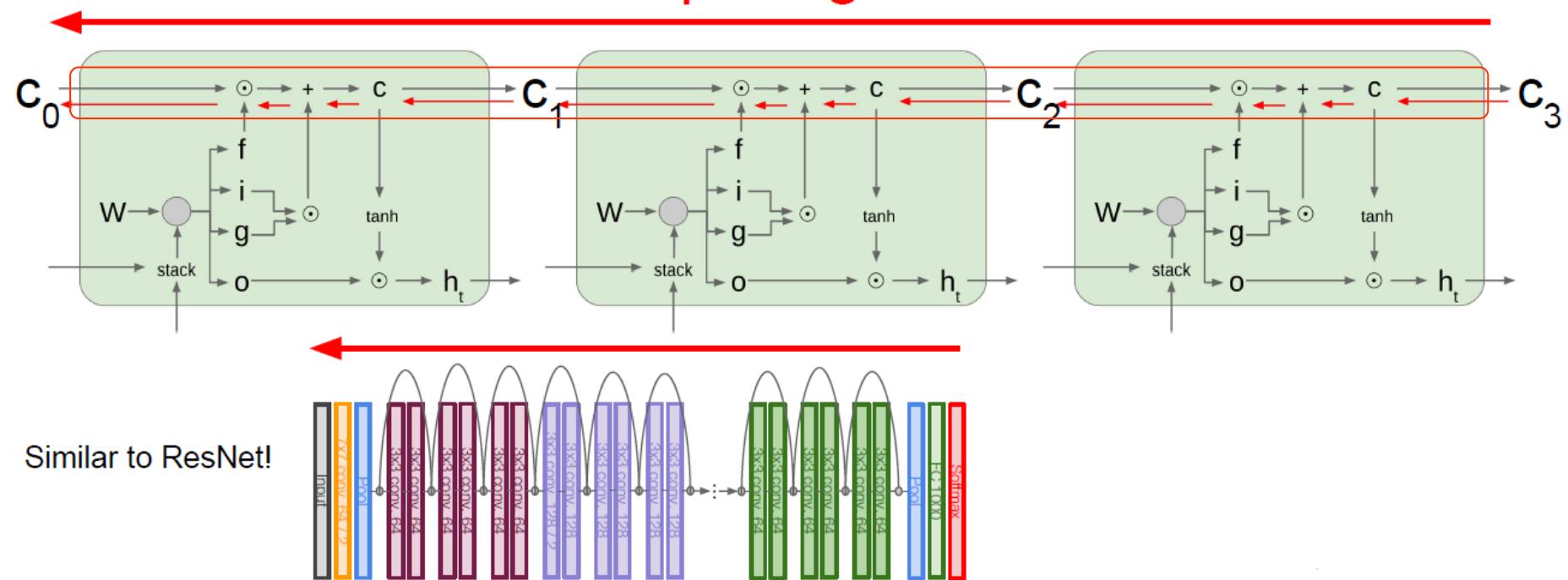
length of internal state

Remark: all internal vectors within GRU green frame have same length.
All ops within a purple circle are performed per element-wise on the ingoing vectors

Long Short Term Memory (LSTM): Gradient Flow

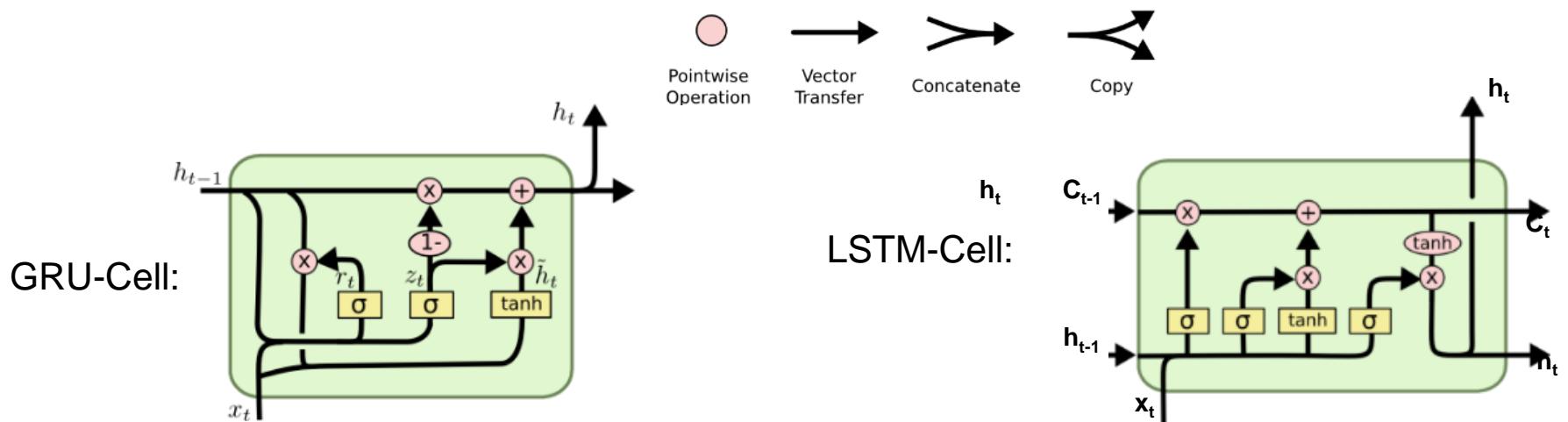
LSTM has an additional cell state C for a “long term memory”.

Uninterrupted gradient flow!



LSTM: [Hochreiter et al., 1997](#)

Long Short Term Memory cell (LSTM) as GRU-extension



2 gates, 1 cell states (h)

$$\text{Relevant gate: } \mathbf{r}_t = \sigma([\mathbf{h}_{t-1}, \mathbf{x}_t] \cdot \mathbf{W}_r + \mathbf{b}_r)$$

$$\text{Update gate: } \mathbf{z}_t = \sigma([\mathbf{h}_{t-1}, \mathbf{x}_t] \cdot \mathbf{W}_z + \mathbf{b}_z)$$

$$\text{Proposed hidden state: } \tilde{\mathbf{h}}_t = \tanh(\mathbf{x}_t \cdot \mathbf{W}_h + \mathbf{b}_h + \mathbf{h}_{t-1} \mathbf{U}_h \otimes \mathbf{r}_t)$$

$$\text{New hidden state is: } \mathbf{h}_t = (\mathbf{1} - \mathbf{z}_t) \otimes \tilde{\mathbf{h}}_t + \mathbf{z}_t \otimes \mathbf{h}_{t-1}$$

3 gates, 2 cell states (S:h, L:C)

$$\text{Forget gate: } \mathbf{f}_t = \sigma([\mathbf{h}_{t-1}, \mathbf{x}_t] \cdot \mathbf{W}_f + \mathbf{b}_f)$$

$$\text{Input gate: } \mathbf{i}_t = \sigma([\mathbf{h}_{t-1}, \mathbf{x}_t] \cdot \mathbf{W}_i + \mathbf{b}_i)$$

$$\text{Output gate: } \mathbf{o}_t = \sigma([\mathbf{h}_{t-1}, \mathbf{x}_t] \cdot \mathbf{W}_o + \mathbf{b}_o)$$

$$\text{Proposed cell state: } \tilde{\mathbf{C}}_t = \tanh([\mathbf{h}_{t-1}, \mathbf{x}_t] \cdot \mathbf{W}_C + \mathbf{b}_C)$$

$$\text{New L cell state: } \mathbf{C}_t = \mathbf{f}_t \otimes \mathbf{C}_{t-1} + \mathbf{i}_t \otimes \tilde{\mathbf{C}}_t$$

$$\text{New S hidden state: } \mathbf{h}_t = \mathbf{o}_t \otimes \tanh(\mathbf{C}_t)$$

Credits: Colah blog <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

Long Short Term Memory networks (LSTM) were introduced by [Hochreiter & Schmidhuber \(1997\)](#).

A simplified variation, the Gated Recurrent Unit, or GRU, introduced by [Cho, et al. \(2014\)](#).

Long Short Term Memory (LSTM) in keras

```
from keras.layers import LSTM  
  
model = Sequential()  
model.add(Embedding(max_features, 32))  
model.add(LSTM(32))  
model.add(Dense(1, activation='sigmoid'))  
  
model.compile(optimizer='rmsprop',  
              loss='binary_crossentropy',  
              metrics=['acc'])  
  
history = model.fit(input_train, y_train,  
                     epochs=10,  
                     batch_size=128,  
                     validation_split=0.2)
```

dimension of vocabulary

dimension of embedding

Model zoo: many pretrained NN are out there

<https://modelzoo.co/>

The screenshot shows a grid of 45 project cards, each representing a different neural network implementation. The projects are categorized by framework (PyTorch, Keras, Chainer, TensorFlow, etc.) and domain (CV, NLP, Generative, etc.). The first two rows of projects are circled in red, and the third row is also circled in red.

Base pretrained models and datasets in pytorch (MNIST, SVHN, CIFAR10, CIFAR100, STL10, AlexNet, VGG16, VGG19, ResNet, Inception, SqueezeNet)

Project	Framework	Description
Graphics code generating model using Processing	PyTorch	PyTorch - processing code generator
Imagenet-vgg	PyTorch	PyTorch Image Classification with Kaggle Dogs vs Cats Dataset
chainer-Variational-AutoEncoder	Chainer	Variational autoencoder (VAE)
Hierarchical Attention Network for Document Classification	PyTorch	A faster and up to date implementation is in my other repo
Simple Generative Adversarial Networks	PyTorch	Run the sample code by typing:
SqueezeNet	PyTorch	Top-1 Acc@1 on ImageNet, without any finetuning compared with SqueezeNet v1.1.
chainer-gan-denosing-feature-matching	Chainer	Generative Adversarial Networks with Denoising Feature Matching
mxnet-audio	MXNet	Implementation of music genre classification, audio-to-text, song recommender, and music search in mxnet
MXSeq2Seq(Gluon)	MXNet	Python implementation of mxnet's Seq2Seq
PyTorch Image Classification with Kaggle Dogs vs Cats Dataset	PyTorch	Classifies an image as containing either a dog or a cat, based on a public dataset, but could easily be extended to other image classification problems.
Bidirectional LSTM on the IMDB dataset	Keras	LSTM-based network on the bAbI dataset
1D CNN on the IMDB dataset	Keras	Memory network on the bAbI dataset (reading comprehension question answering)
1D CNN-LSTM on the IMDB dataset	Keras	Sequence to sequence learning for performing additions of strings of digits
LSTM text generation	Keras	Using pre-trained word embeddings
FastText on the IMDB dataset	Keras	Structurally constrained recurrent nets text generation
Simple CNN on MNIST	Keras	Inception v3
Simple CNN on CIFAR10 with data augmentation	Keras	Neural Style Transfer
Visualizing the filters learned by a CNN	Keras	Deep dreams
Stateful LSTM	Keras	Stateful LSTM
Siamese network	Keras	Siamese network
DeconvNet	PyTorch	Learning Deconvolution Network for Semantic Segmentation
Pixel-wise Segmentation on VOC2012 Dataset using PyTorch	PyTorch	Pixel-wise segmentation on the VOC2012 dataset using
generative-models	PyTorch	Amortized, unrolled and visually interpretable PyTorch implementations of VAE, BRNN, RNN, GAN, DCGAN, LSGAN, BEGAN, REGAN, Im2GAN, ISGAN, FisherGAN
V-Net	PyTorch	Fully Convolutional Neural Networks for Volumetric Medical Image Segmentation
Evolution Strategies	PyTorch	Evolution Strategies as a Scalable Alternative to Reinforcement Learning
CycleGAN and Semi-Supervised GAN	PyTorch	PyTorch implementation of CycleGAN and Semi-Supervised GAN for Domain Transferring
Adversarial Generator-Encoder Network	PyTorch	This repository contains code for the paper
Recurrent Variational Autoencoder	PyTorch	Recurrent Variational Autoencoder that generates variational data implemented in pytorch, 1511.05848, 1609.0615
AttGAN	TensorFlow	Tensorflow implementation of AttGAN - Arbitrary Facial Attribute Editing. Only Change What You Want
img-classification_pk_pytorch	PyTorch	Quickly compare your image classification model with the state-of-the-art models (such as DenseNet, ResNet, ...)
PNASNet.pytorch	PyTorch	PyTorch implementation of PNASNet-S on Imagenet
U-Net	TensorFlow	For Brain Tumor Segmentation
CNN-LSTM-CTC	PyTorch	I realize that there are many models for text recognition, and I think this is one of CTC loss layer to realize no segmentation for text images