

## Modelling sequence data and using ideas from CNN & Recurrent Neural Networks

*Beate Sick*

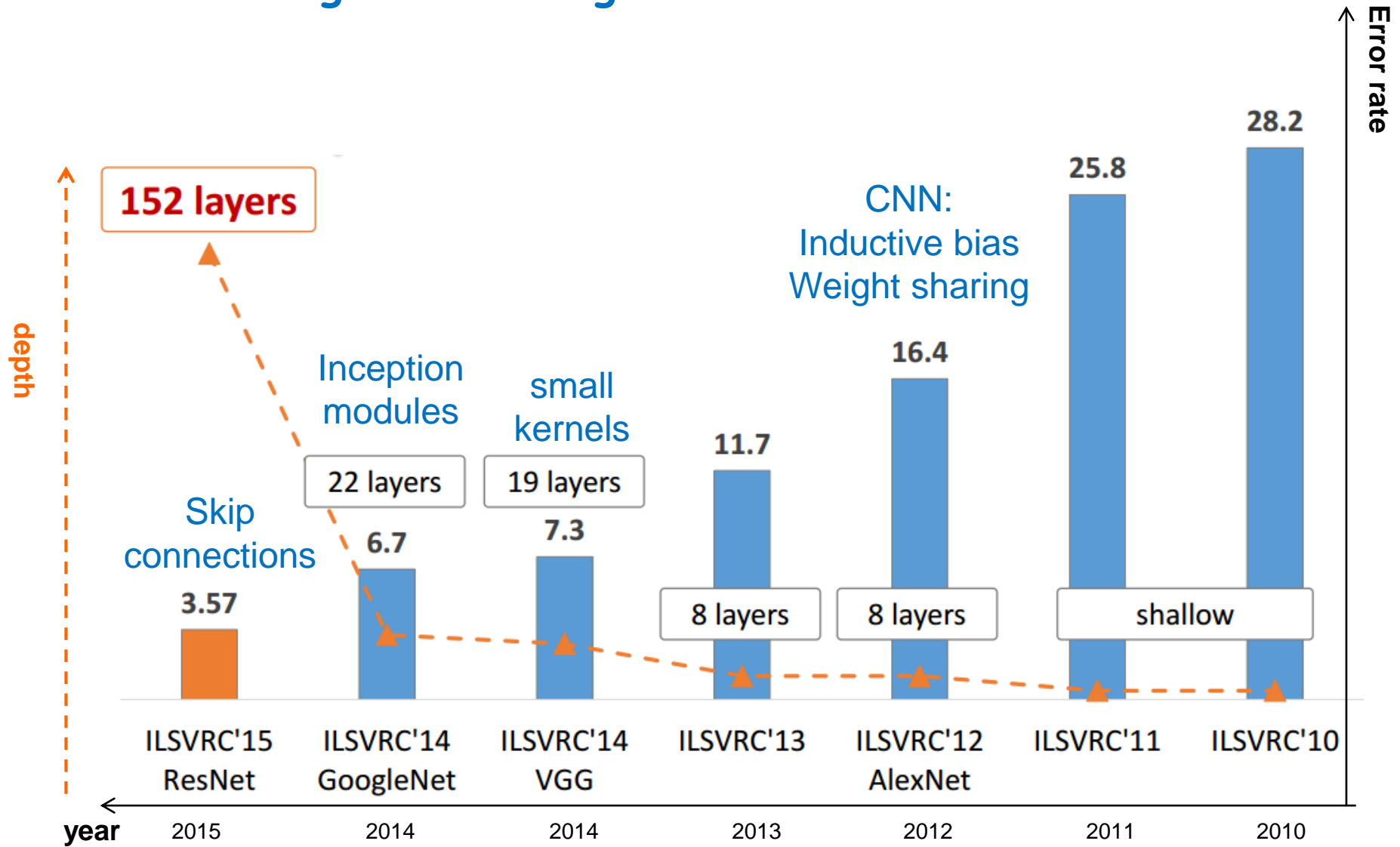
[sick@zhaw.ch](mailto:sick@zhaw.ch)

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

# Topics

- **Famous tricks in challenge winning CNN architectures**
  - Inductive Bias, weight sharing, Inception moduls, gradient highway via skip-connections
- **Working with text data**
  - Bag of words
  - Word-embedding
- **Recurrent neural networks (RNN)**
  - possible use cases
  - memory in recurrent NN
  - loss construction in an RNN
  - RNNs in Keras

# 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” 2014 2<sup>nd</sup> place

- 2<sup>nd</sup> 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

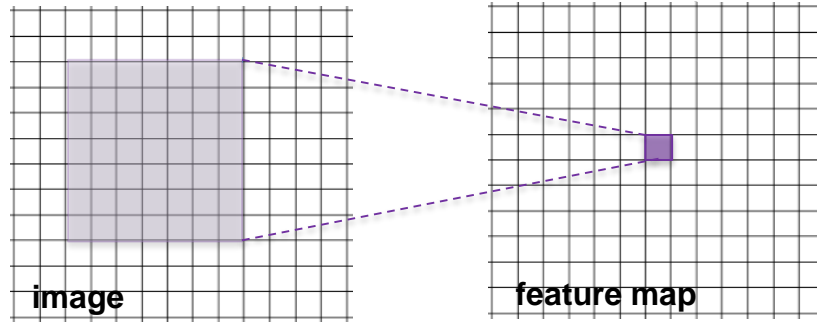
<http://arxiv.org/abs/1409.1556>



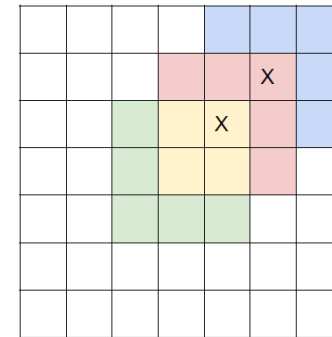
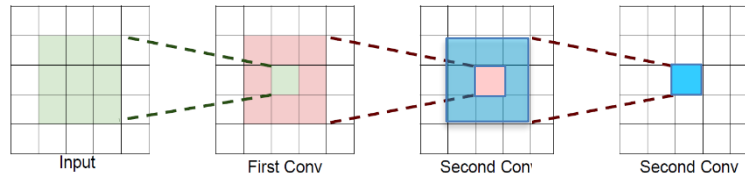
# The trend in modern CNN architectures goes to small filters

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

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

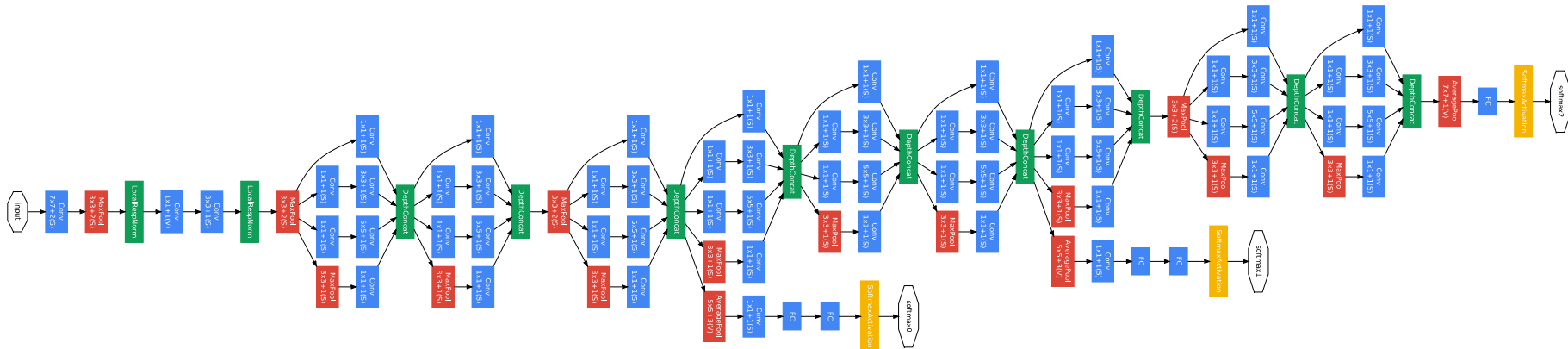


- 2) Working with **three 3x3**  
**conv layers** (stride 1)  
**3\*9=27 weights**

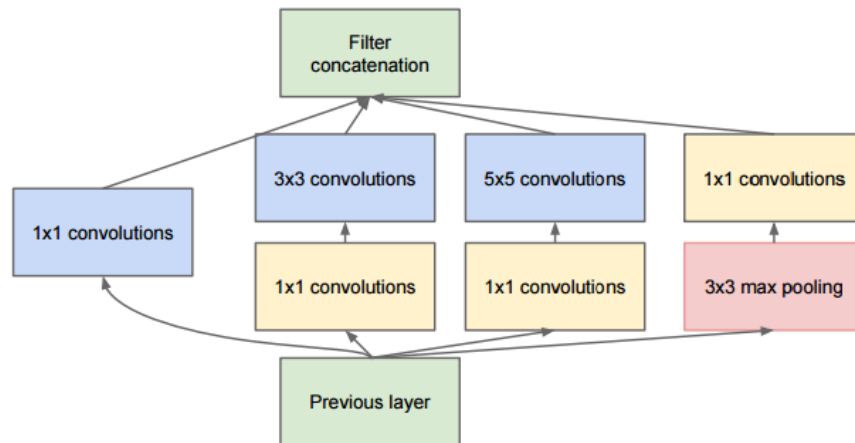


When working stacking conv-layers with small kernels we can achieve with less weights and more non-linear combinations the same receptive field

# Winning architecture (GoogLeNet, 2014)



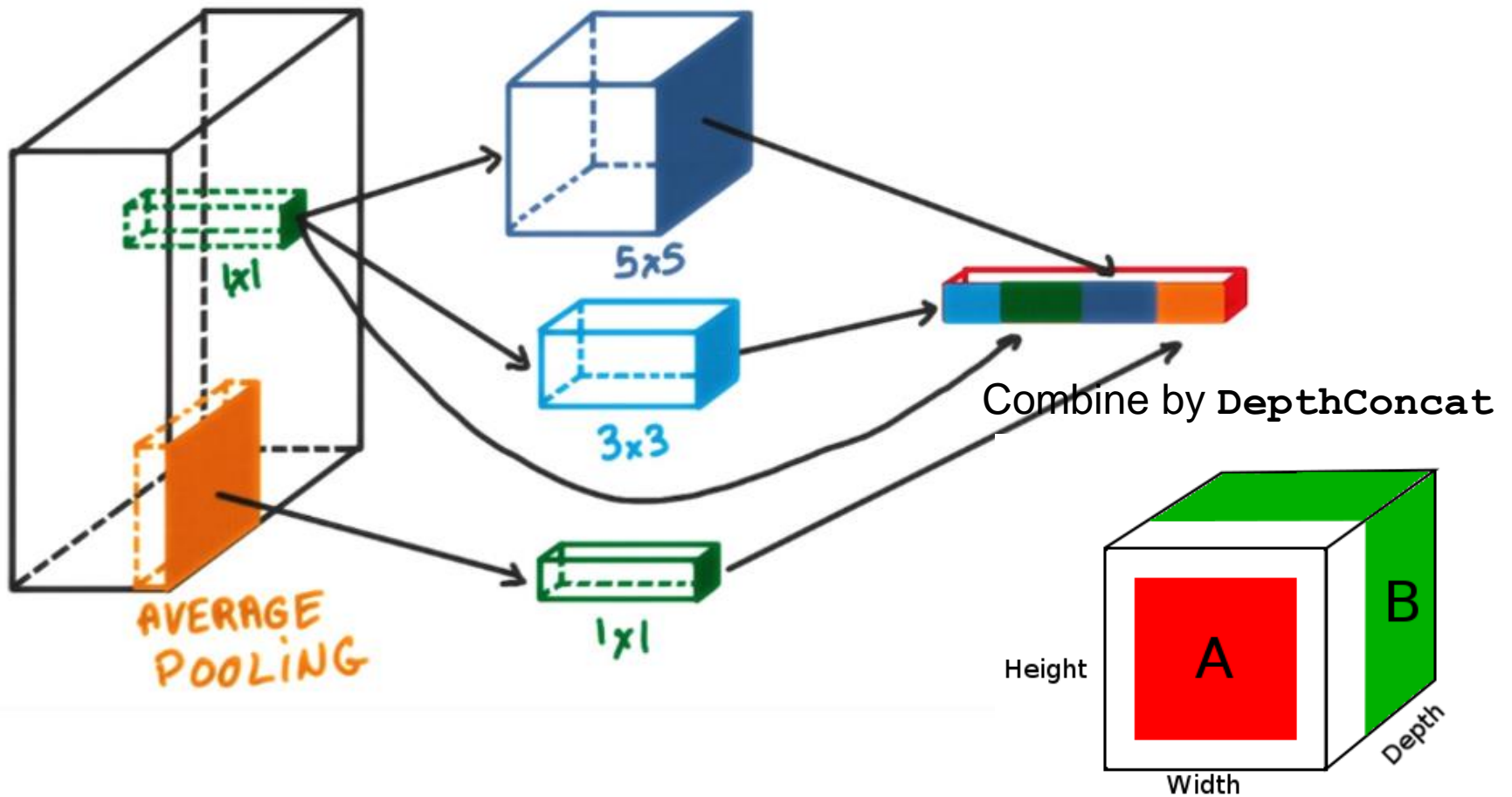
The inception module:  
use in 1 layer in parallel different kernels and combine their results



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

# The idea of inception modules

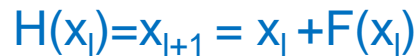
## INCEPTION MODULES





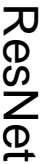
152  
layers

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

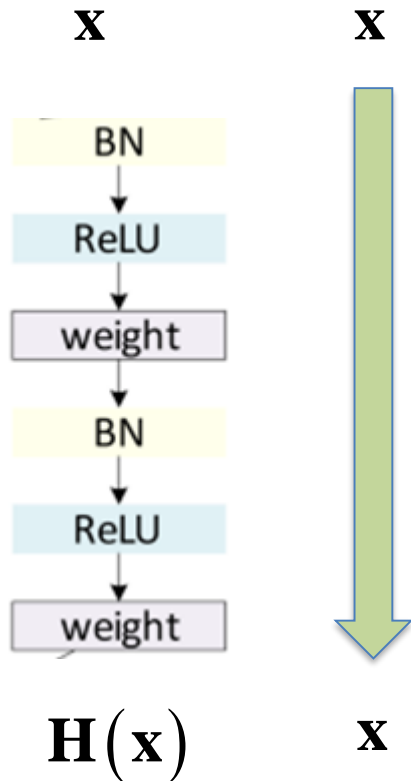


152 layers:  
Why does this train at all?

# plain VGG



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



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

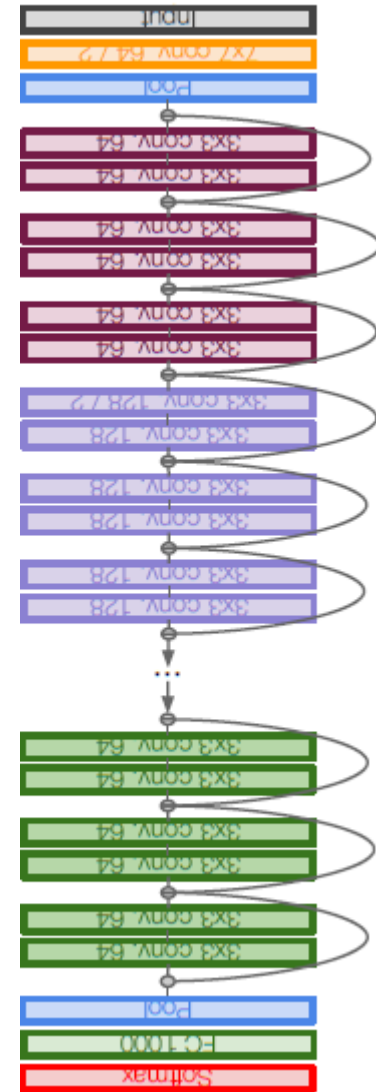
$$\mathbf{y} = \mathbf{H}(\mathbf{x}, \mathbf{W}_H) \cdot T(\mathbf{x}, \mathbf{W}_T) + \mathbf{x} \cdot (1 - T(\mathbf{x}, \mathbf{W}_T)).$$

Special case:

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

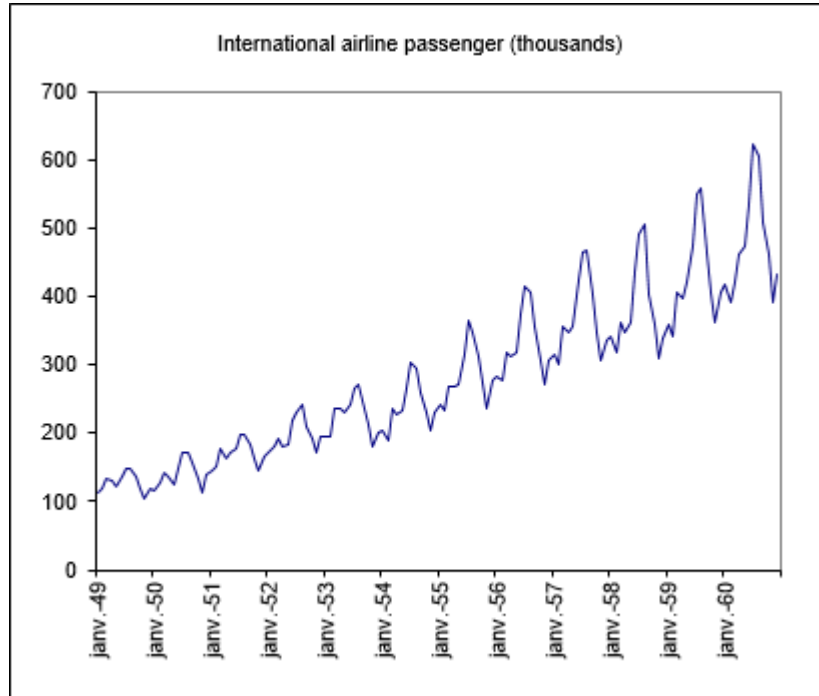
This opens a highway for the gradient:

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

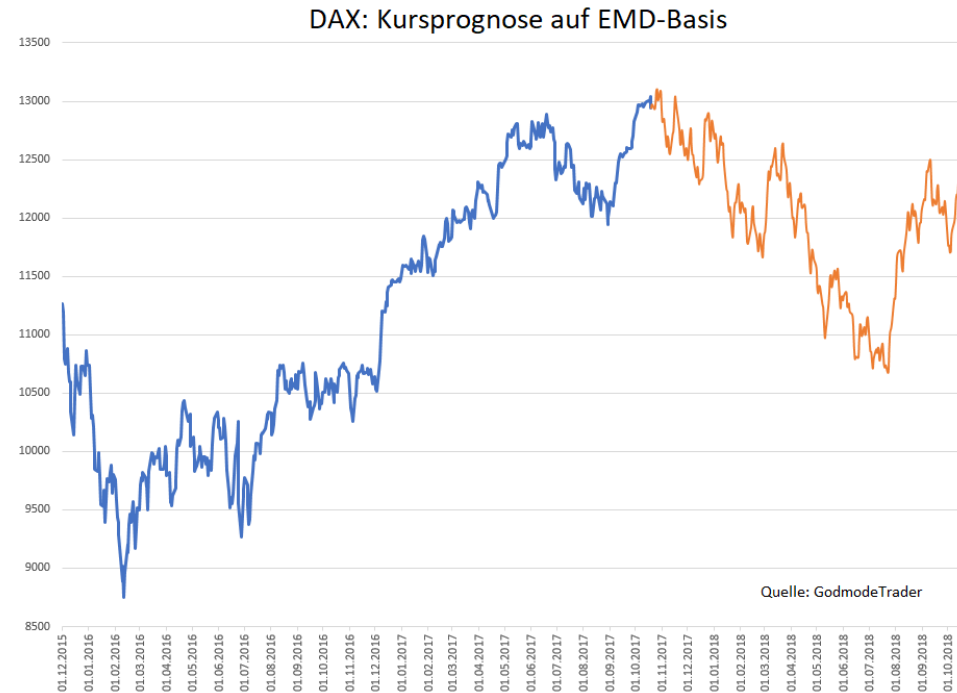


# Sequence Data

# Example Sequence Data: time-series

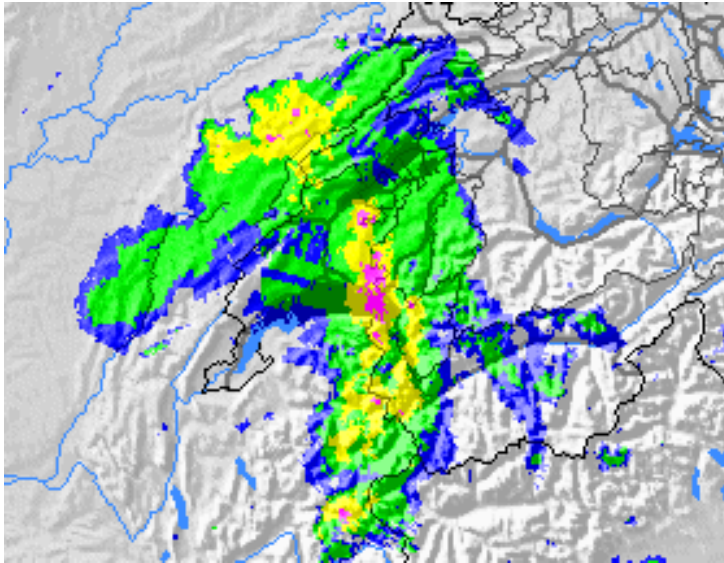


How many passenger will we have next month?



What will be the DAX value tomorrow?

# Example Sequence Data: videos



[https://www.metradar.ch/2009\\_exp/pc/service.php](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 GrÃ¼Ãe

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."

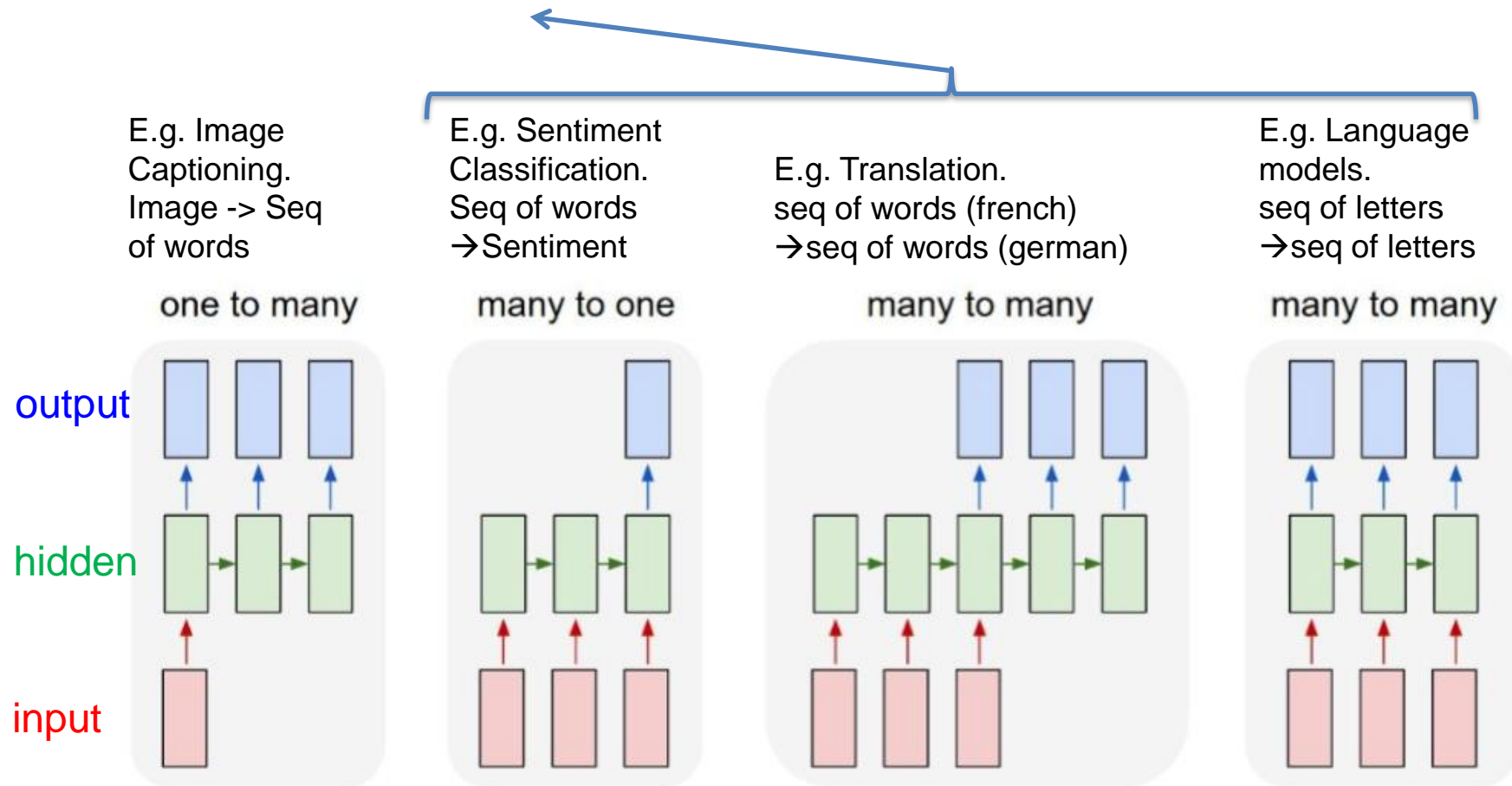


"a horse is standing in the middle of a road."

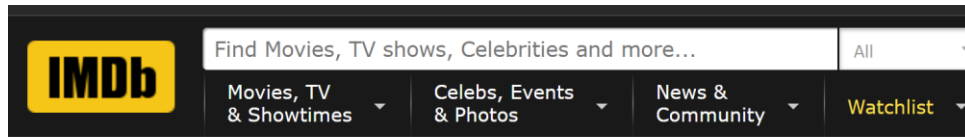


# Modeling sequence data

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



# Possible task: Sentiment analysis with movie reviews



## Inception (2010) User Reviews

[+ Review this title](#)

3,259 Reviews



Hide Spoilers

Filter by Rating:

Show All

Sort by:

Helpfulness



9/10

### Matrix but in dreamworld? Nah.

[clipturnity](#) 22 August 2010

**Warning: Spoilers**

1,651 out of 1,926 found this helpful. Was this review helpful? [Sign in](#) to vote.

[Permalink](#)



10/10

### Insanely Brilliant ! Nolan has outdone himself !!

[naman-avastol](#) 10 July 2010

What is the most resilient parasite? An Idea! Yes, Nolan has created something with his unbelievably, incredibly and god- gifted mind which will blow the minds of the audience away. The world premiere of the movie, directed by Hollywood's most inventive dreamers, was shown in London and has already got top notch reviews worldwide and has scored maximum points! Now the question arises what the movie has that it deserve all this?

Dom Cobb(Di Caprio) is an extractor who is paid to invade the dreams of various business tycoons and steal their top secret ideas. Cobb robs forcefully the psyche with

	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).

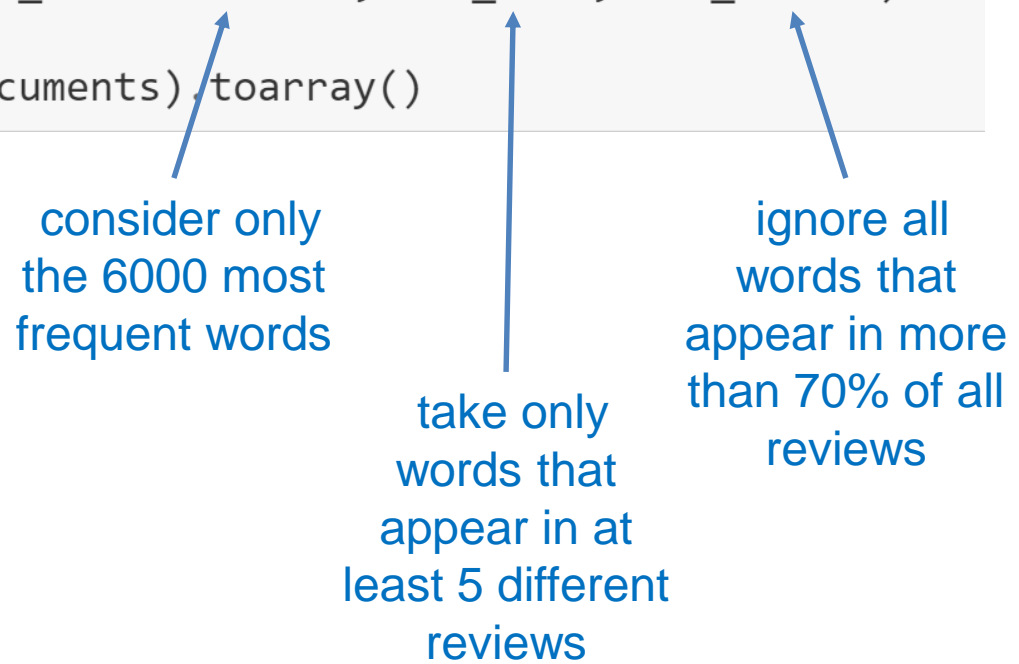
This allows to **represent each sentence as a numeric vector of the same length!**

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()
```

consider only  
the 6000 most  
frequent words



take only  
words that  
appear in at  
least 5 different  
reviews

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”

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

1-hot:

1	0	0	0	0	0	0	0	0	the
0	1	0	0	0	0	0	0	0	cat
0	0	1	0	0	0	0	0	0	sat
0	0	0	1	0	0	0	0	0	on
1	0	0	0	0	0	0	0	0	the
0	0	0	0	1	0	0	0	0	hat




$N$  elements (=vocabulary size)

“The dog ate my homework”

[1, 6, 7, 8, 9]

1	0	0	0	0	0	0	0	0	the
0	0	0	0	0	1	0	0	0	dog
0	0	0	0	0	0	1	0	0	ate
0	0	0	0	0	0	0	1	0	my
0	0	0	0	0	0	0	0	1	homework



$N$  elements (=vocabulary size)

# Tokenizing in keras

```
from keras.preprocessing.text import Tokenizer
```

```
samples = ['The cat sat on the mat.', 'The dog ate my homework.']
```

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

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

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

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

**Turns strings into lists of integer indices**

```
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.**

**Builds the word index**

# Learning an embedding in keras

```
# fit tokenizer on all reviews
total_reviews = documents
tokenizer = Tokenizer()
tokenizer.fit_on_texts(total_reviews)
```

```
# transform tokens to a sequence of integers
X_train_tokens = tokenizer.texts_to_sequences(X_train)
X_val_tokens = tokenizer.texts_to_sequences(X_val)
X_test_tokens = tokenizer.texts_to_sequences(X_test)
```

```
# zeropad the sequences to have the "same" length
X_train_pad = pad_sequences(X_train_tokens, maxlen=max_length, padding='post')
X_val_pad = pad_sequences(X_val_tokens, maxlen=max_length, padding='post')
X_test_pad = pad_sequences(X_test_tokens, maxlen=max_length, padding='post')
```

1) Generate token-representation and use zero-padding to get same length for all sequences

```
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()
```

2) Define embedding-dimension: Within the embedding layer the token-representation is transformed into 1-hot-encoding

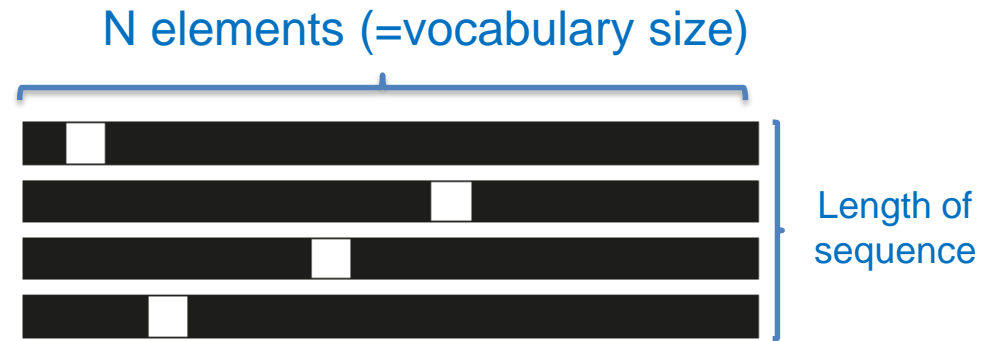
3) Train NN including weight of embedding-layer from 1-hot-encoding to embedding-representation

```
history=model.fit(X_train_pad, y_train, batch_size=64, epochs=40, validation_data=(X_val_pad, y_val), verbose=1)
```

# Go from 1-hot-encodings to word embeddings

## 1-hot encodings

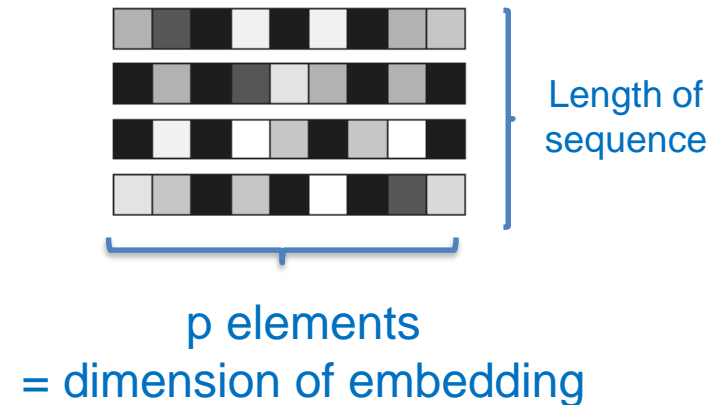
- Based on vocabulary of size  $N$
- sparse: one 1 and  $N-1$  zeros
- High-dim: vector-length =  $N$



Fully connected NN  
with  $N \cdot p$  weights  
(for details see below)

## Word embedding's are

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





# Wordembedding layer in keras

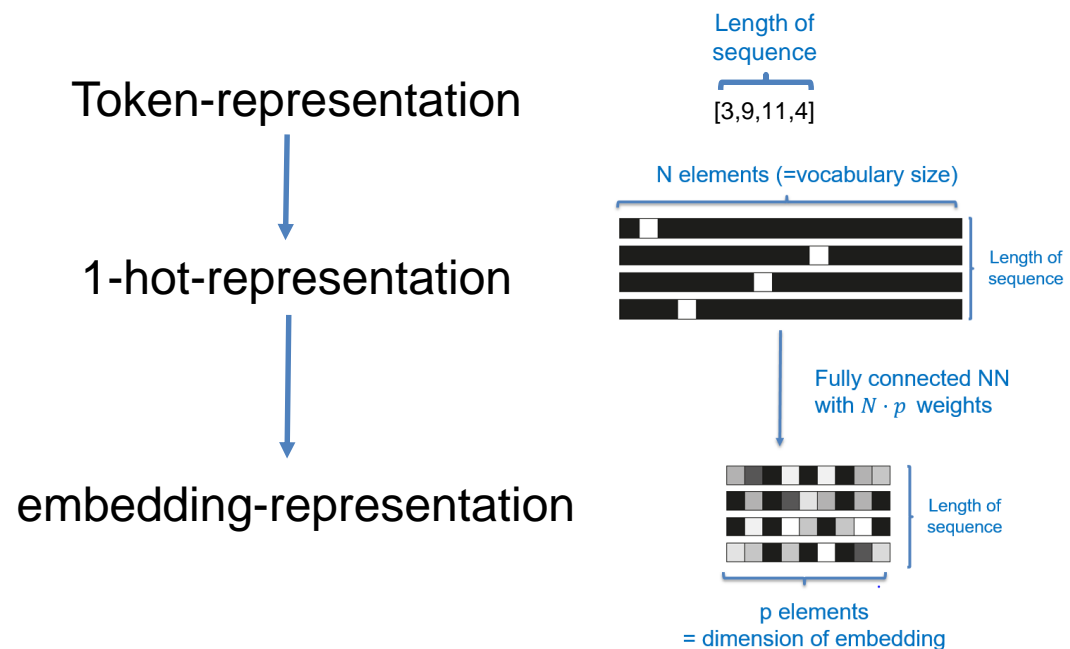
```
from keras.layers import Embedding  
embedding_layer = Embedding(1000, 64)
```

$N$        $p$

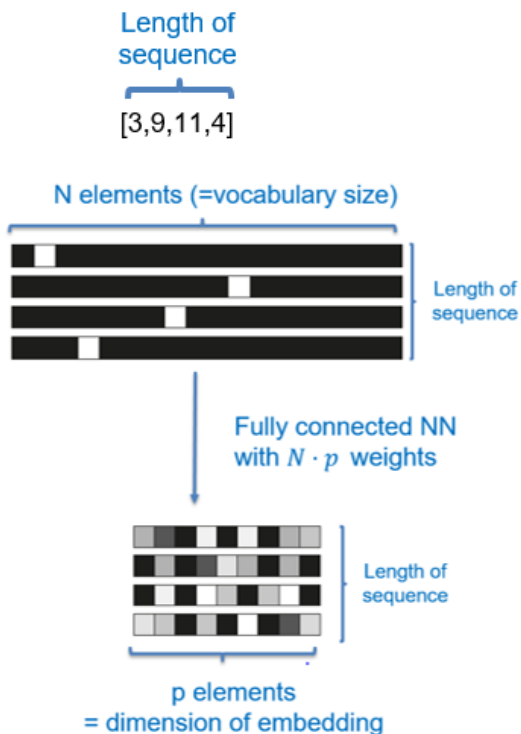
The Embedding layer takes at least two arguments: the number of possible tokens (here, 1,000: 1 + maximum word index) and the dimensionality of the embeddings (here, 64).

The embedding layer returns a 3D floating-point tensor of shape (samples, sequence\_length, embedding\_dimensionality).

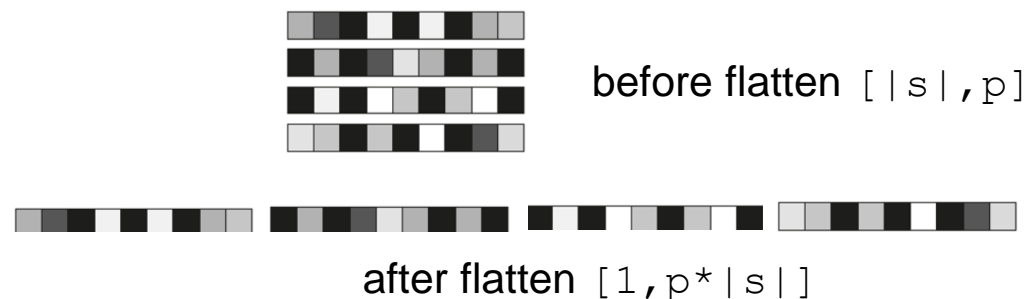
Look at 1 sample:



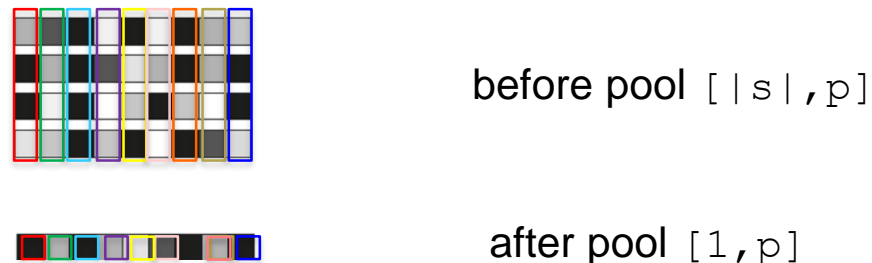
# How to proceed with embeddings



- **Flatten** results in vector of length:  $|s| \cdot p$



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



- Use as **input to 1D-conv**, RNN, GRU, LSTM

# A text classification NN using pooled embeddings

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

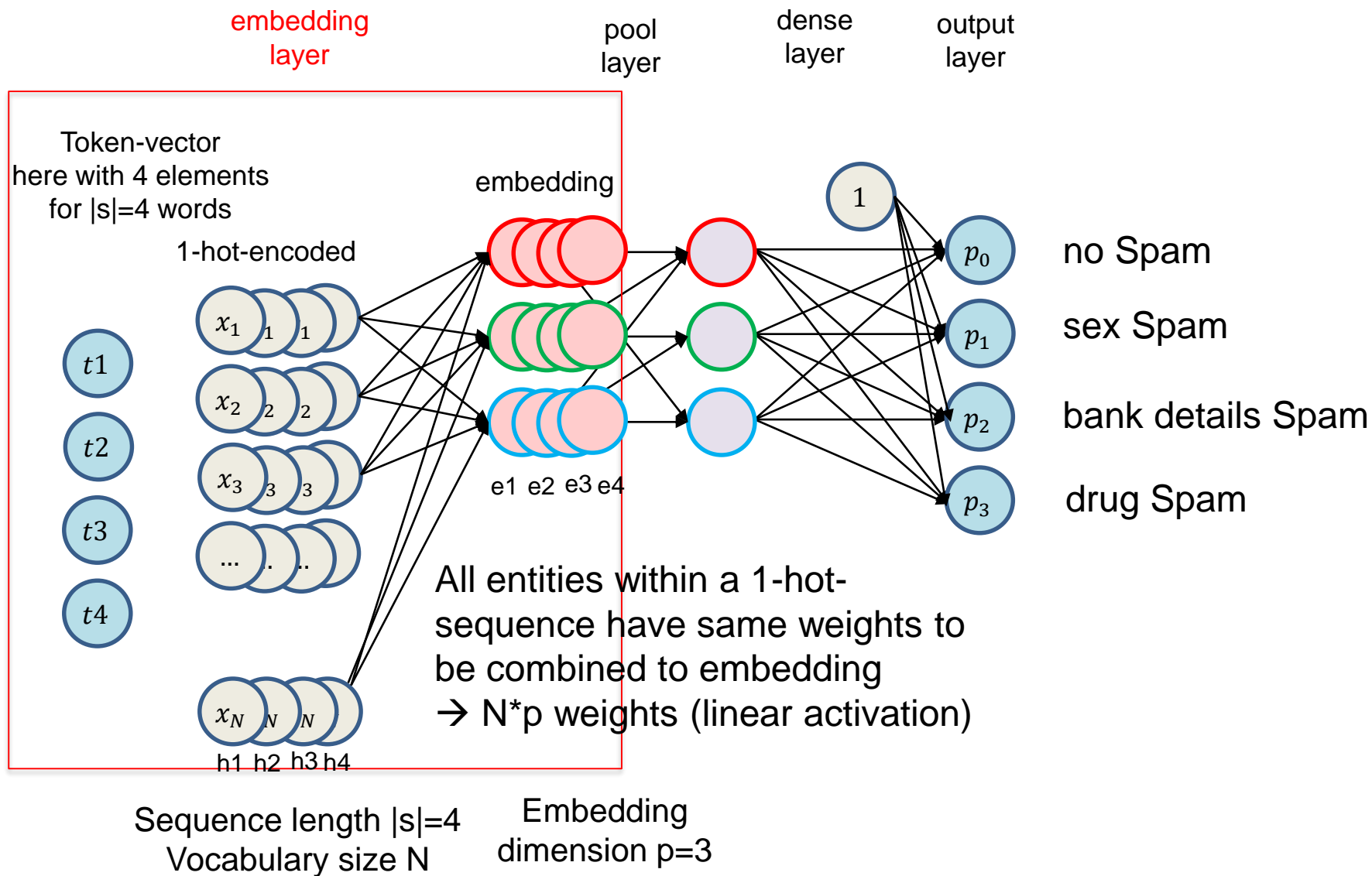
EMBEDDING_DIM = 30

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
<hr/>		
global_average_pooling1d_1 ( (None, 30)		0
<hr/>		
dropout_1 (Dropout)	(None, 30)	0
<hr/>		
dense_1 (Dense)	(None, 20)	620
<hr/>		
dropout_2 (Dropout)	(None, 20)	0
<hr/>		
dense_2 (Dense)	(None, 1)	21
=====		

# Learning embeddings within a NN based on pooling



→ During pooling we lose information on the order in the sequence

# 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).**

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

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

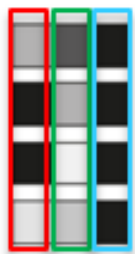
model = Sequential()
model.add(Embedding(10000, 8, input_length=maxlen))
model.add(Flatten())
model.add(Dense(1, activation='sigmoid'))
model.compile(optimizer='rmsprop', loss='binary_crossentropy', metrics=['acc'])
model.summary()

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

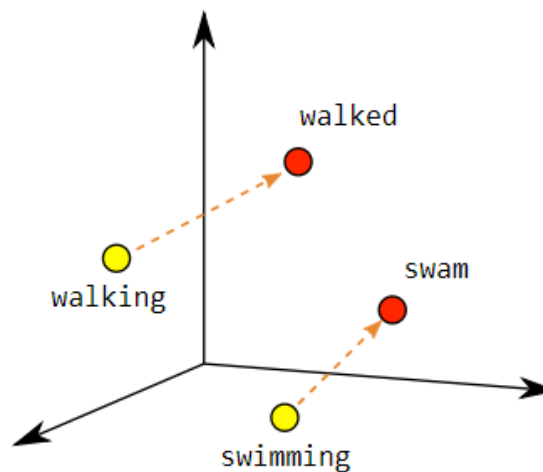
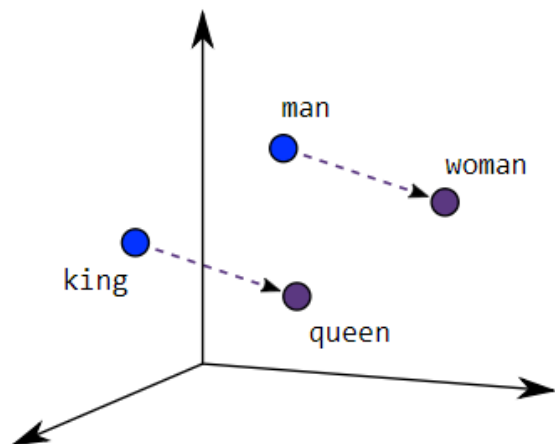
**Adds the classifier on top**

→ Flattening preserves information on the order in the sequence

# Word embedding space

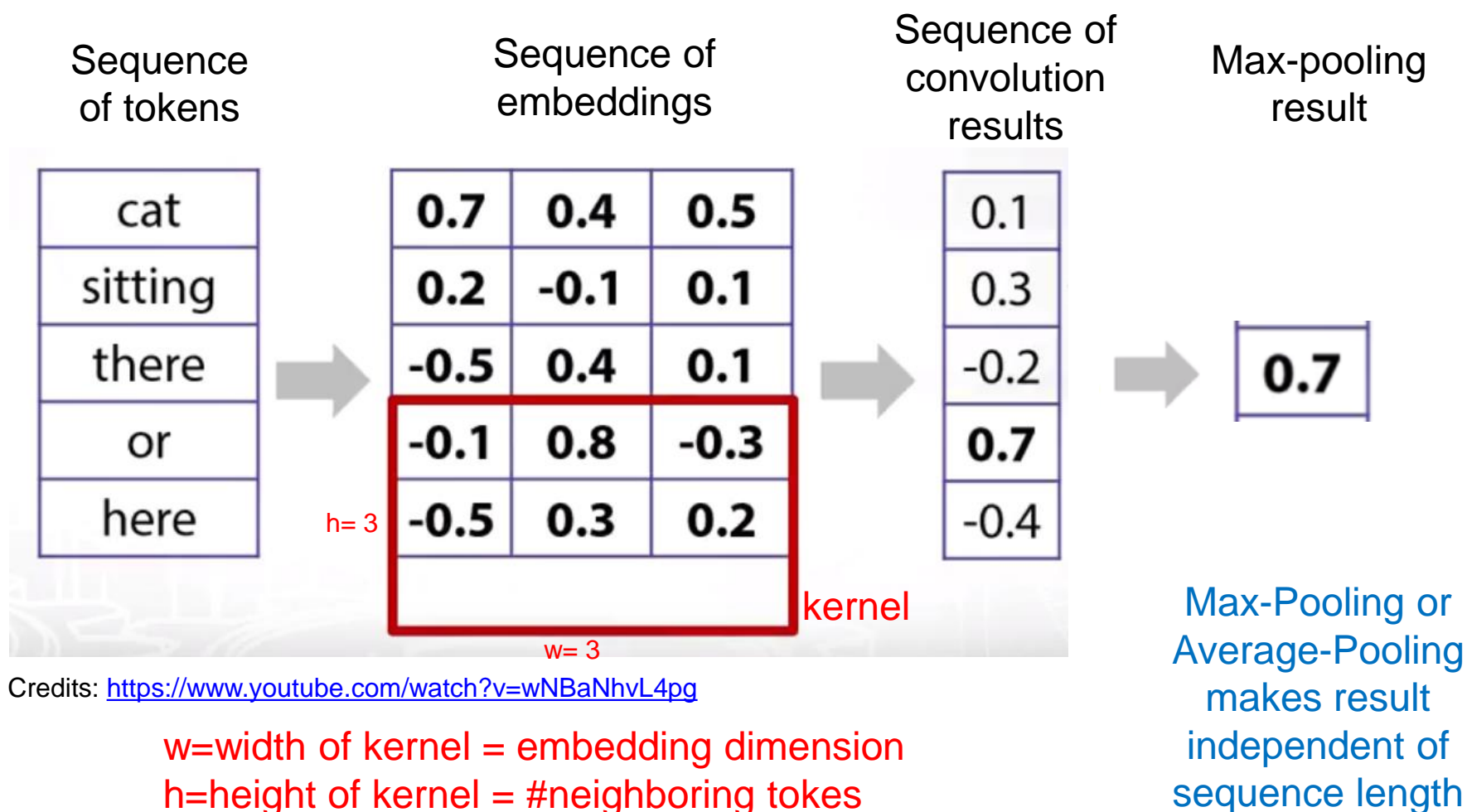


Sequence with 4 words, represented by **3-dim embedding space**



Remark: The result of the embedding depends on the task where it was trained-  
Task in pre-trained embedding is usually predicting context word.

# Applying 1D convolution on n-grams of embeddings



We slide  $h \times w$  kernel only in 1 direction  $\rightarrow$  1D convolution  
A kernel can intake  $h$  embeddings (=h-grams,  $h$  neighboring words)

# Use embeddings as input 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.



# Do your own sentiment analysis

- Work through the instructions in sentiment exercise in day 6 using [13 sentiment analysis with imdb reviews.ipynb](#)



wiederkehrend

# Recurrent Neural Networks

Recurrent NN have memory

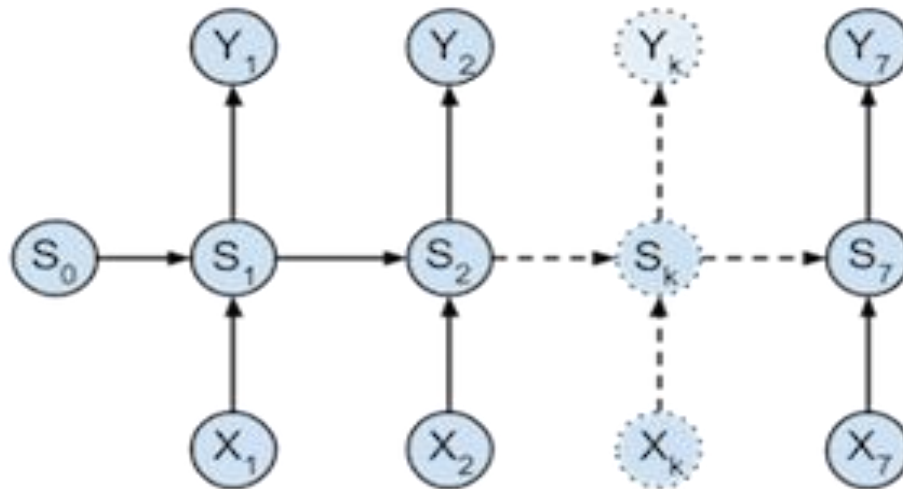
## Challenge 2: How to handle input of different lengths?

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

By reading from word to word our believe 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:  $y_1 = \text{prob}_{\text{spam}}$ ,  $y_2 = \text{prob}_{\text{spam}}$ , ....



**Hidden memory State:**

**$S_1$ =state after first word**

**$S_2$ =state after first word**

...

Input:  $x_1 = \text{vector}_{\text{word1}}$ ,  $x_2 = \text{vector}_{\text{word2}}$ , ....

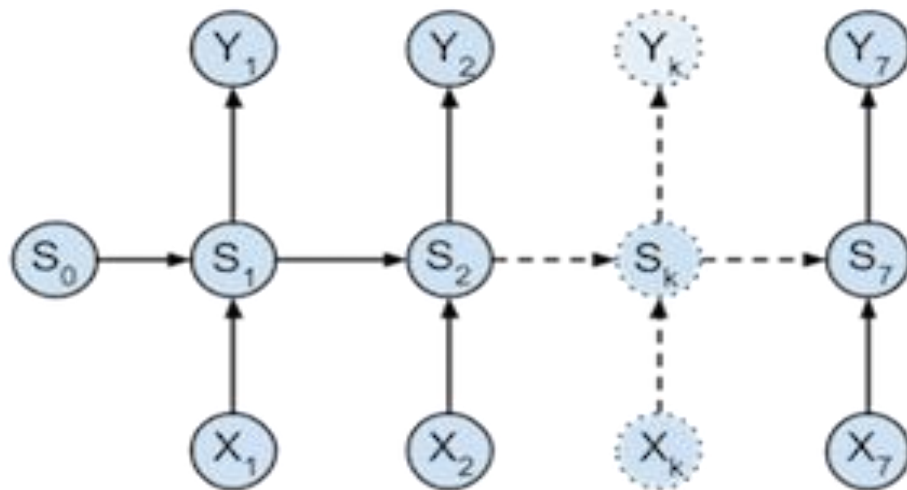
# How to choose the dimensions of the hidden state?

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



Hidden memory State:

$S_1$ =state after first word

$S_2$ =state after first word

...

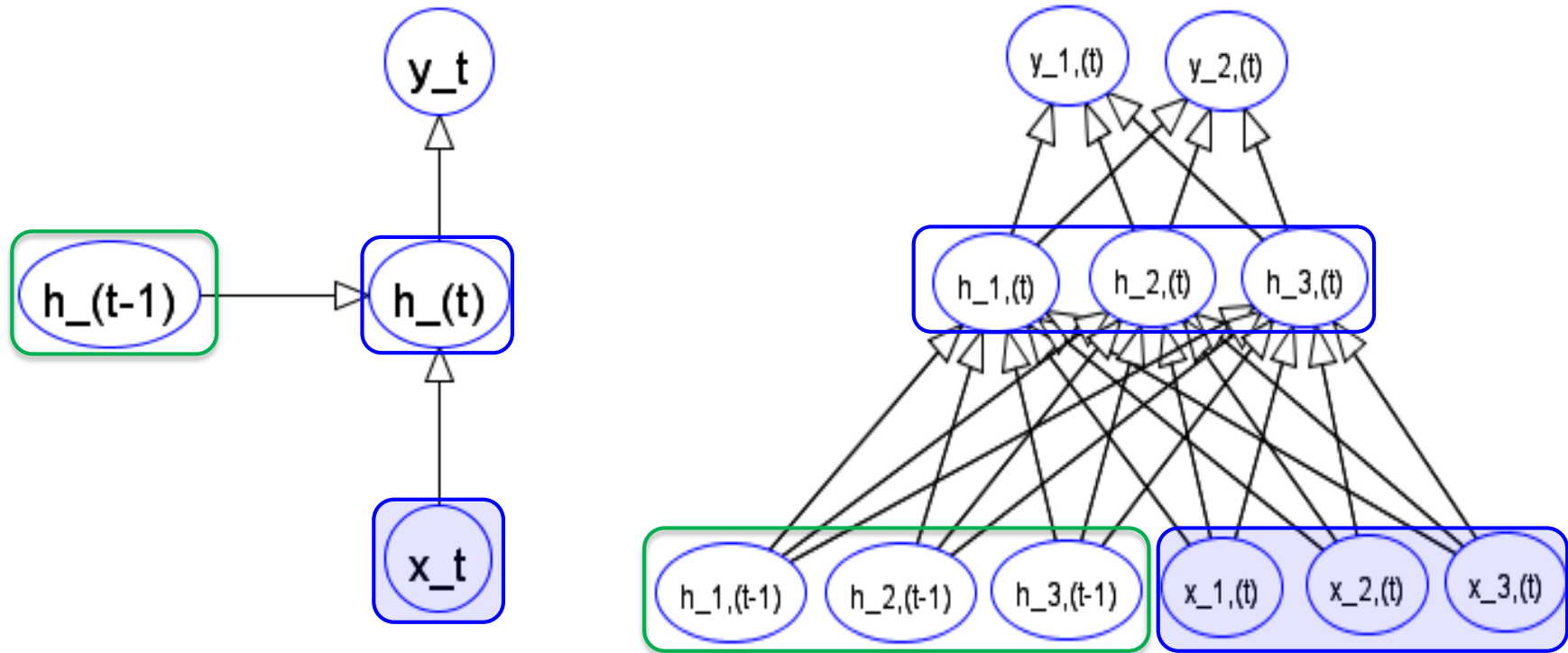
Input:  $x_1=\text{vector}_{\text{word1}}$ ,  $x_2=\text{vector}_{\text{word2}}$ , ....

# Two representations of a RNN at time $t$

$\mathbf{x}$  has 3 components – a very small 1-hot-encoded vocabulary or embedding vector ;-)

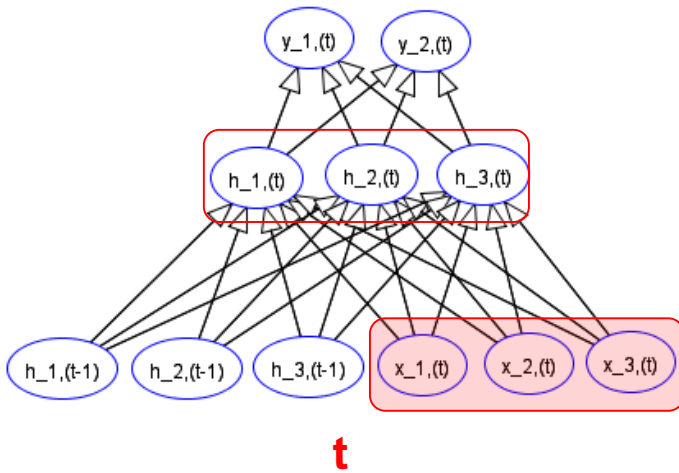
$\mathbf{y}$  has 2 components (spam/no-spam for mails, passed/failed for scoring essays).

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

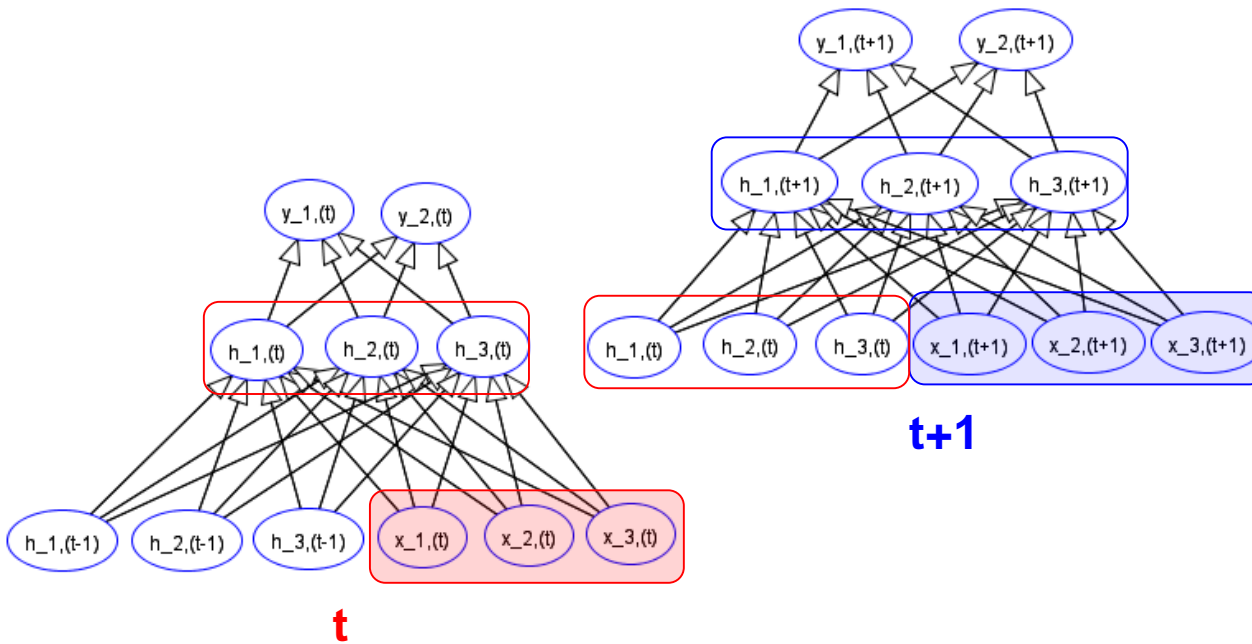


# Stepping through an RNN

# A simple RNN at 3 successive time steps

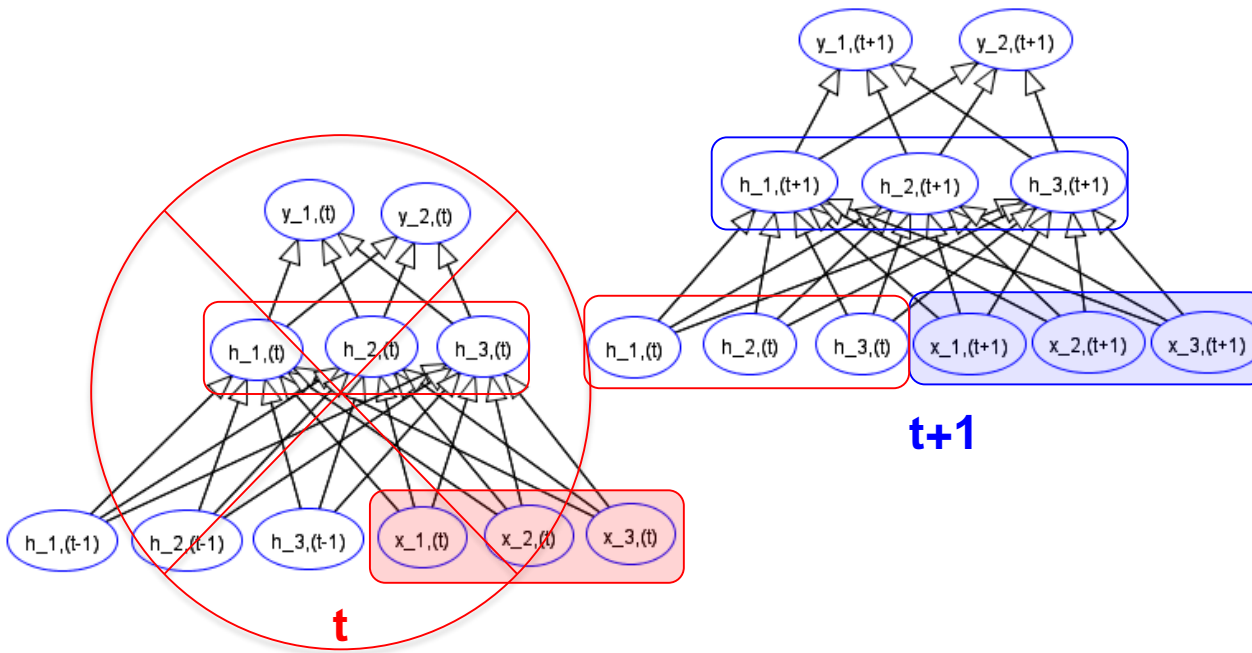


# A simple RNN at 3 successive time steps

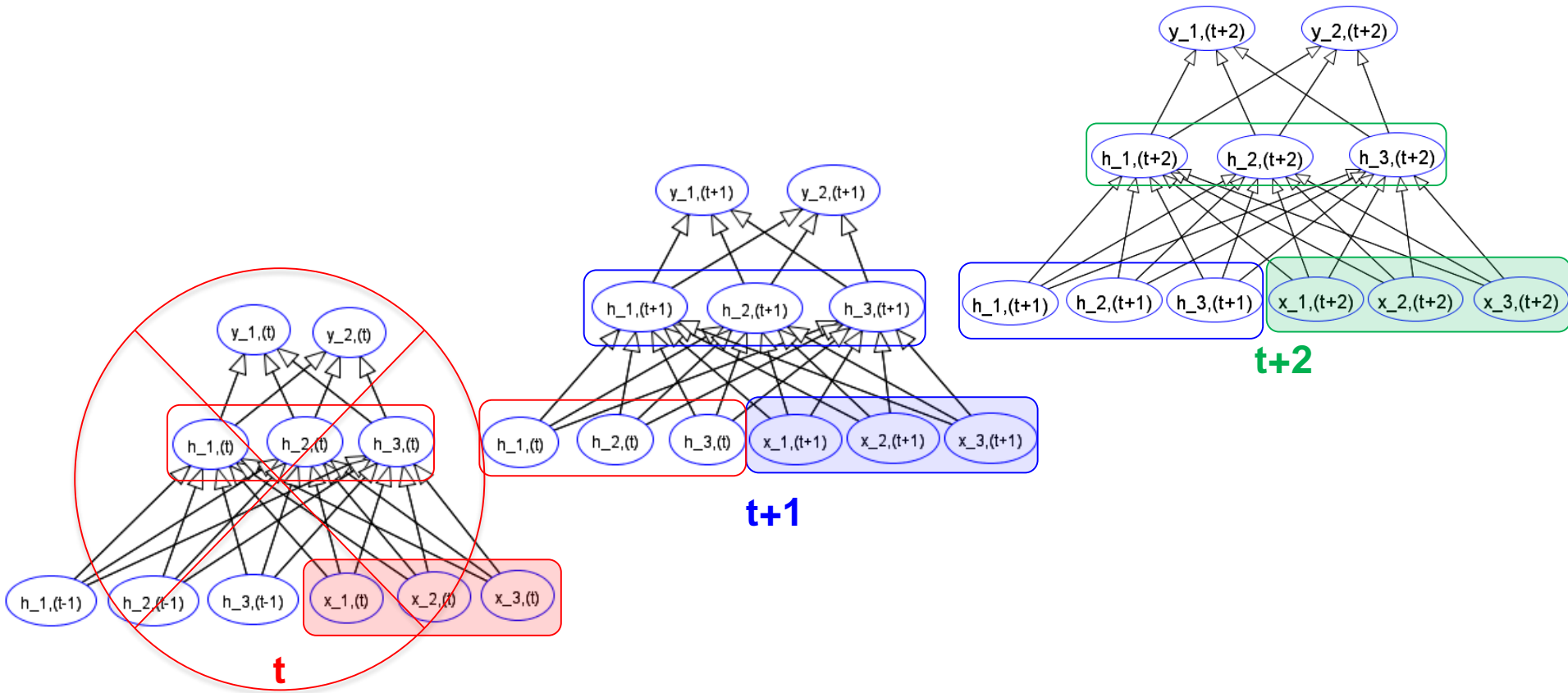




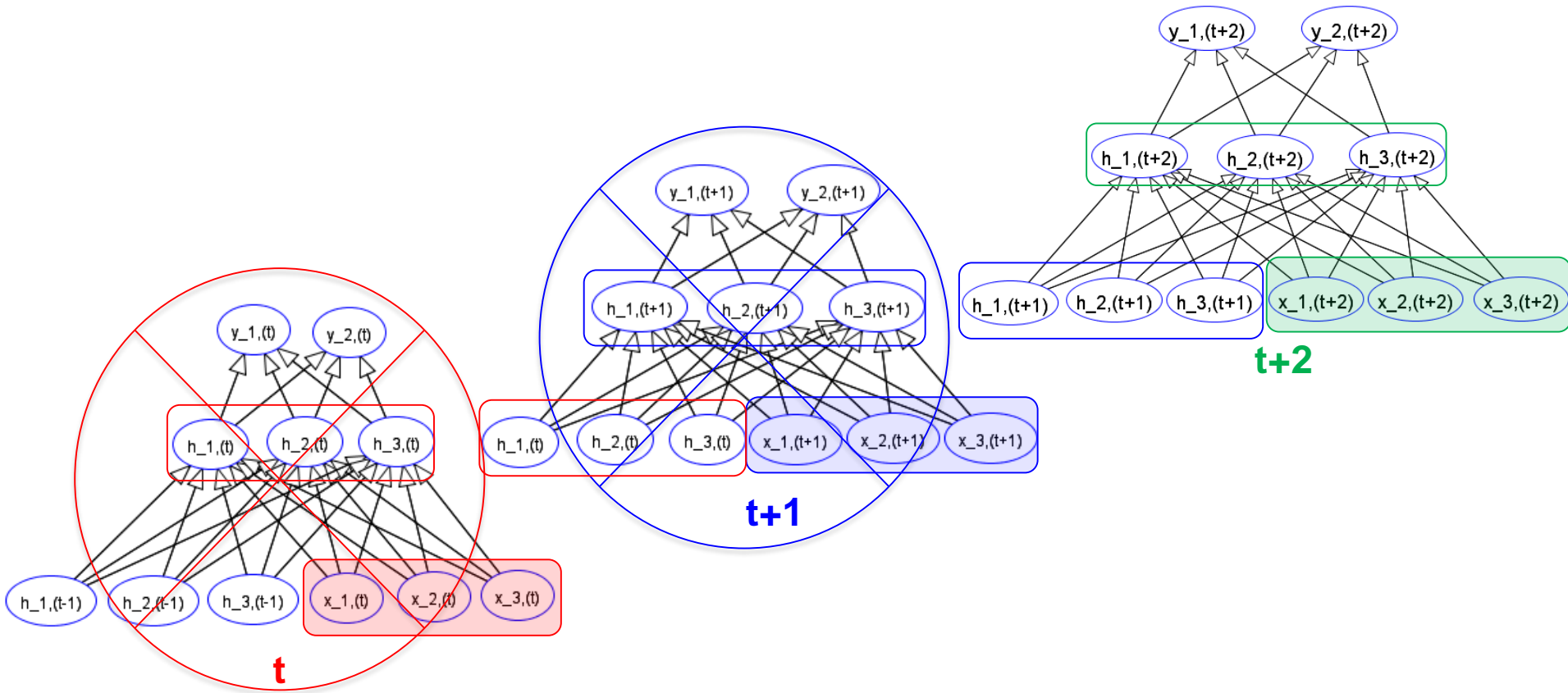
# A simple RNN at 3 successive time steps



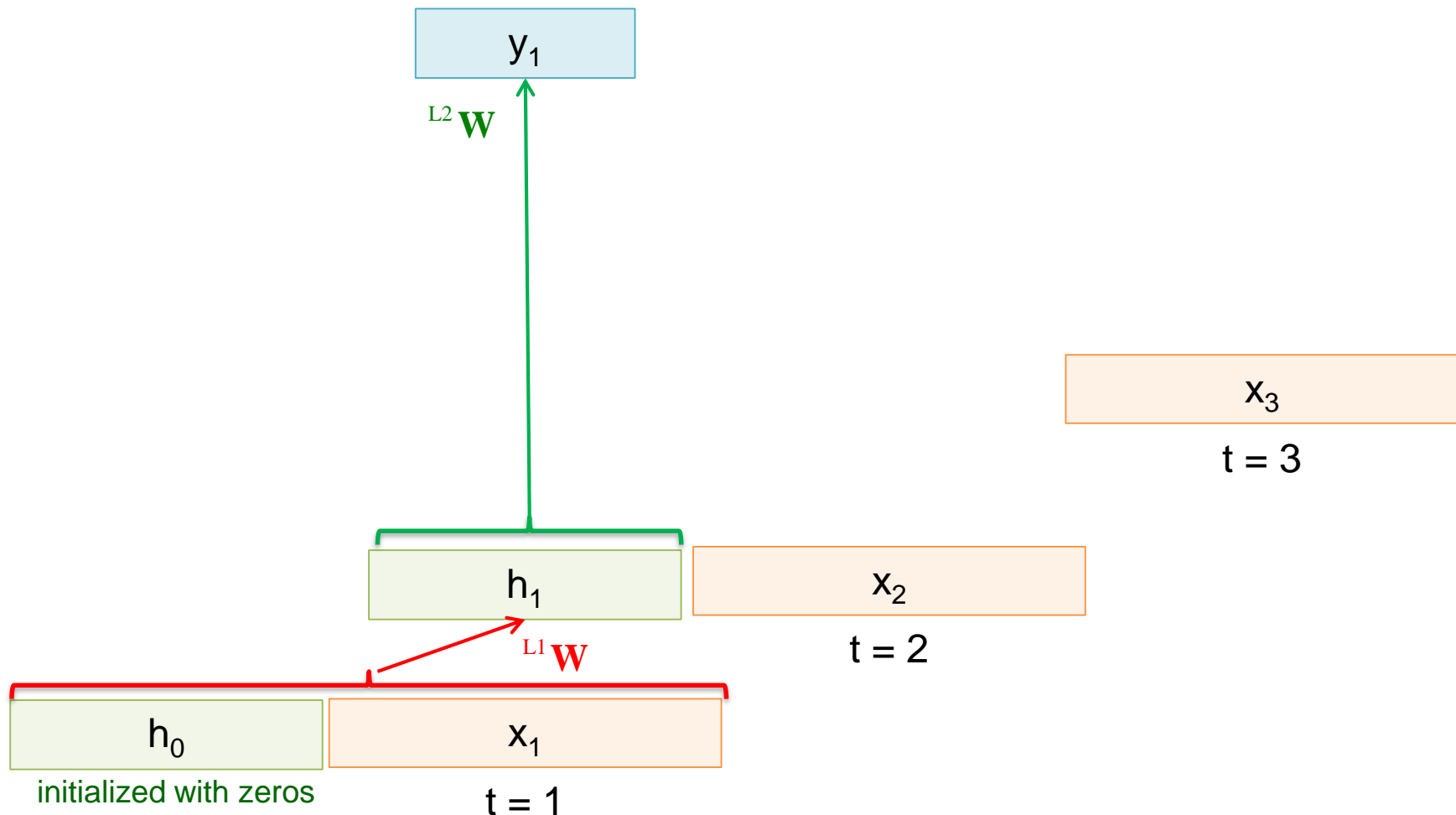
# A simple RNN at 3 successive time steps



# A simple RNN at 3 successive time steps

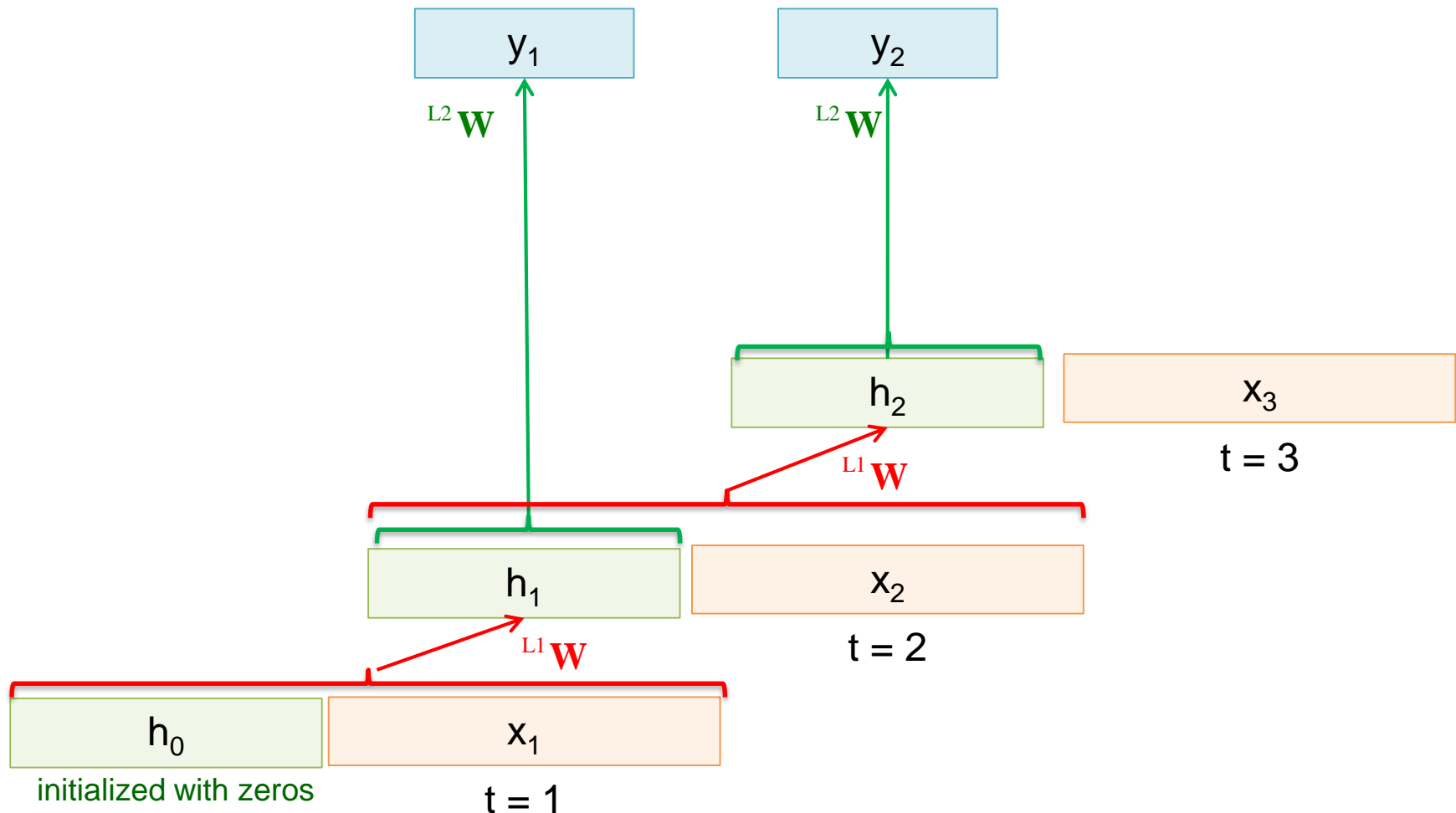


# An RNN shares weights across all time steps



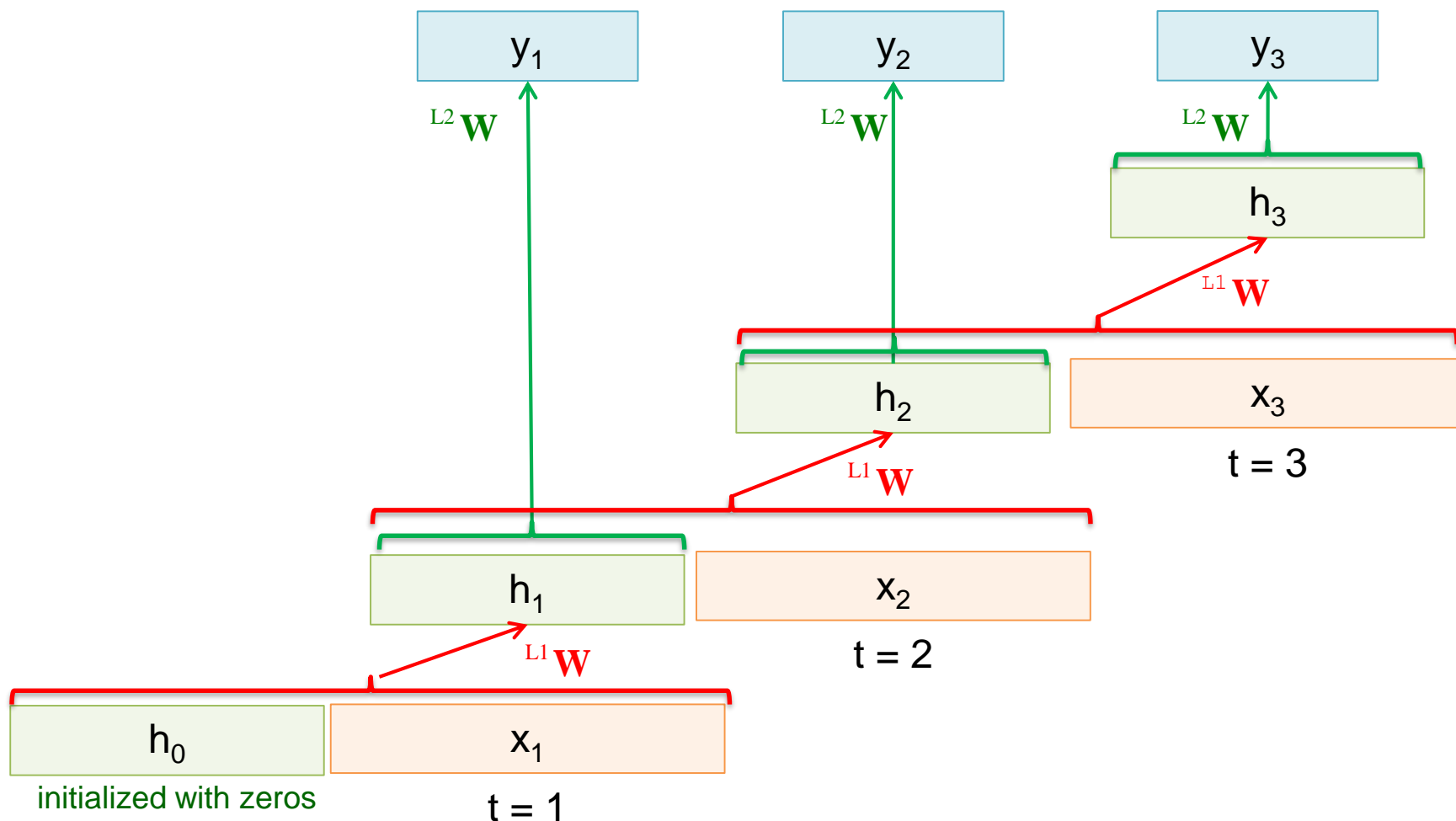
Imagine a trained RNN with fixed weight matrices in layer 1 and layer 2.

# An RNN shares weights across all time steps



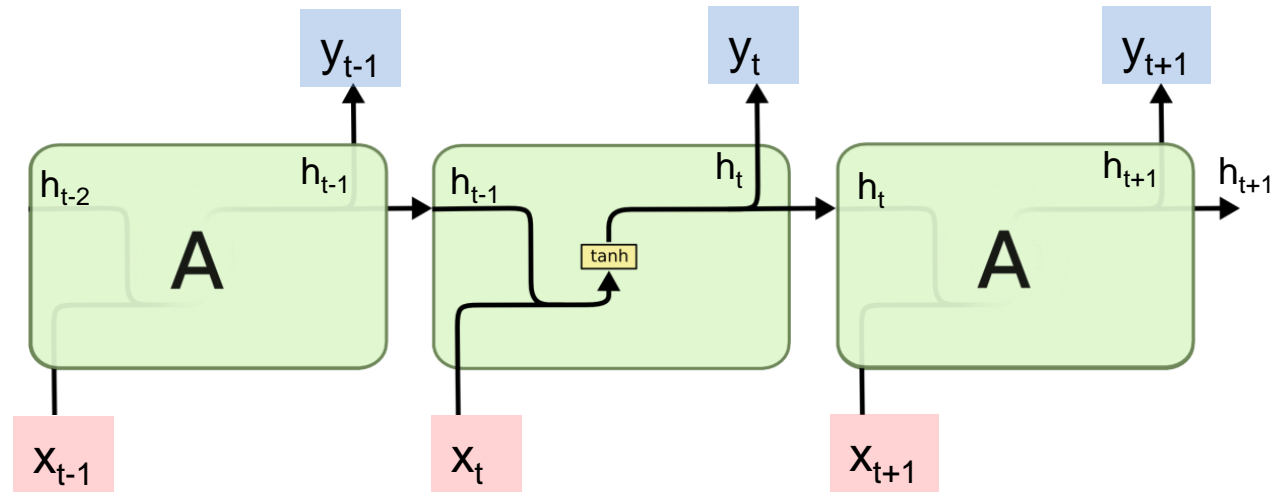
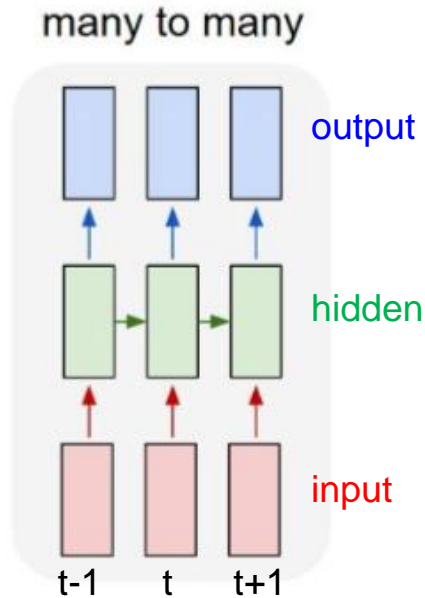
We use at each time step the same weight matrices between the layers!

# 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!  
Therefore it is called recurrent network!!

# Using diagrams to represent an RNN

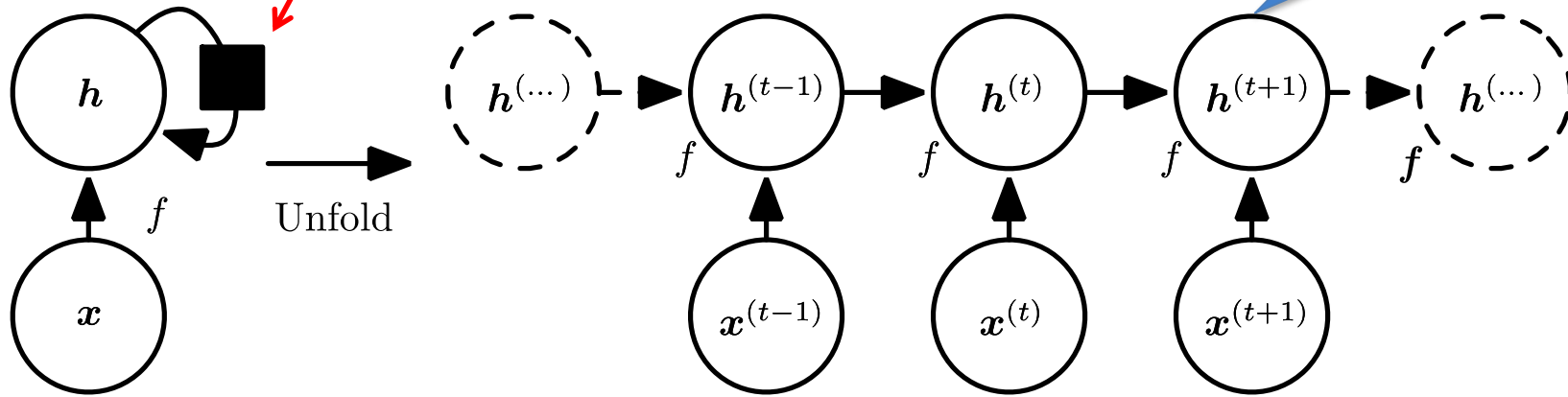


$$\mathbf{A} = f_{\mathbf{W}}(\mathbf{h}_{t-1}, \mathbf{x}_t) = \tanh([\mathbf{h}_{t-1}, \mathbf{x}_t] \cdot \mathbf{W} + \mathbf{b})$$

# Using a circuit diagram to represent an RNN

Use same weights  
in each time step

$h^{(t)}$  summarizes /  
abstraction

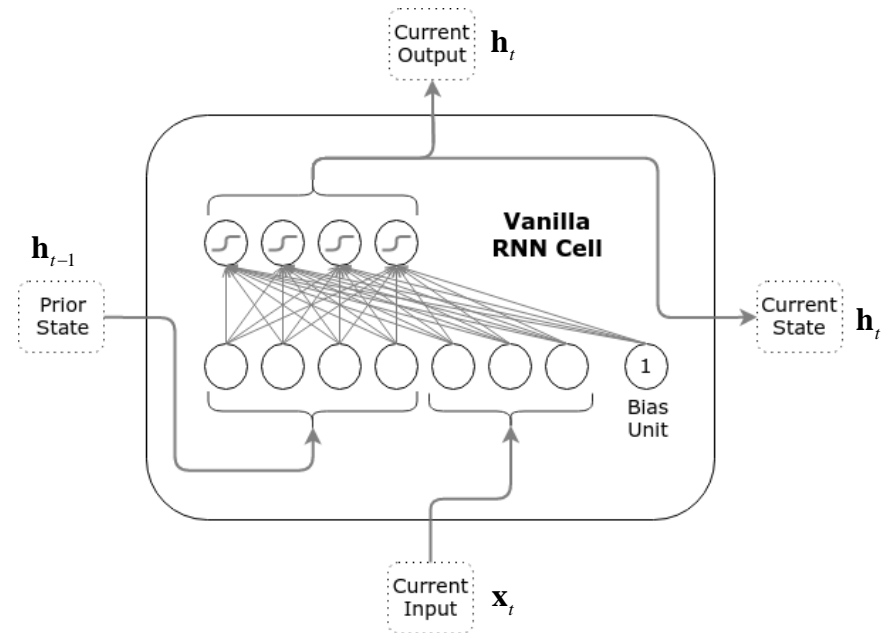
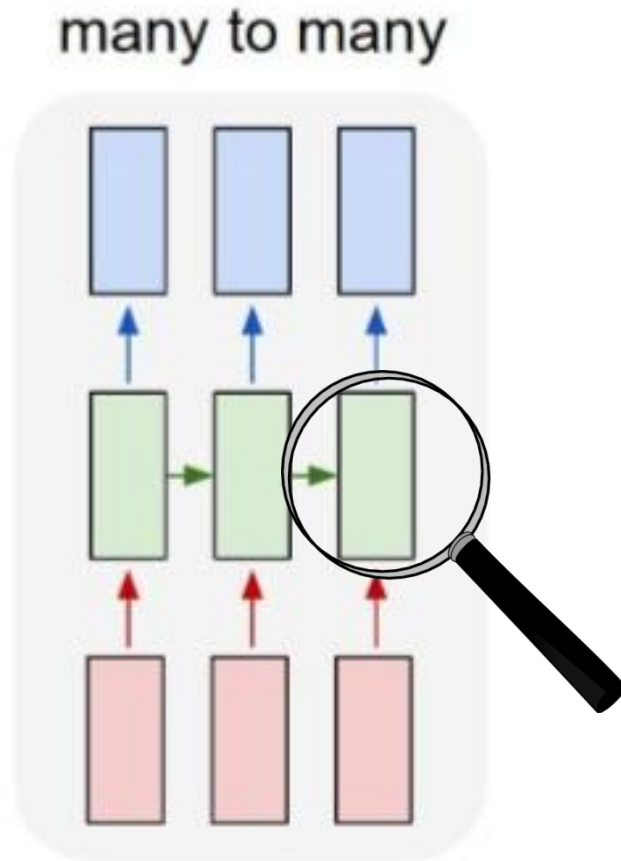


Left: Circuit Diagram (black square delay of one time step)

Right: Unrolled / unfolded



# Looking into a RNN “cell”



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

# A simple forward pass



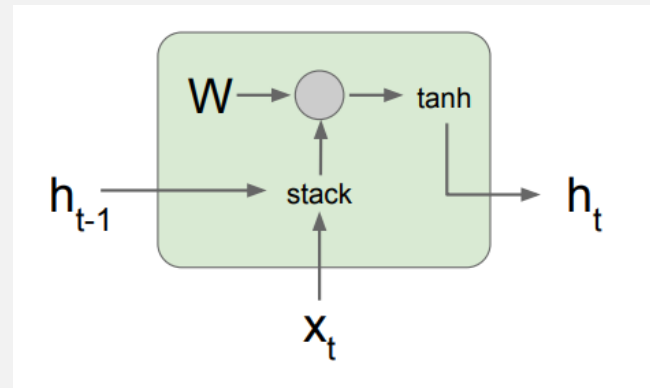
- Given the hidden state at t-1, the input x at t and the weight matrix as:

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

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

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

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



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

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

## Solution

$$h_{t-1} = (0, 1)$$

$$x_t = (1, 0)$$

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

$$(0, 1, 1, 0) \begin{pmatrix} 1.5 & -2 \\ 0 & 1 \\ -1 & 1/2 \\ 2 & 0 \end{pmatrix} = (-1, 1.8)$$

$$\Rightarrow h_t = \text{tanh}((-1, 1.8)) \approx (-0.76, 0.91)$$

# Loss construction in an RNN

# Determine the loss contribution of instance 1

mini-batch of size M=8

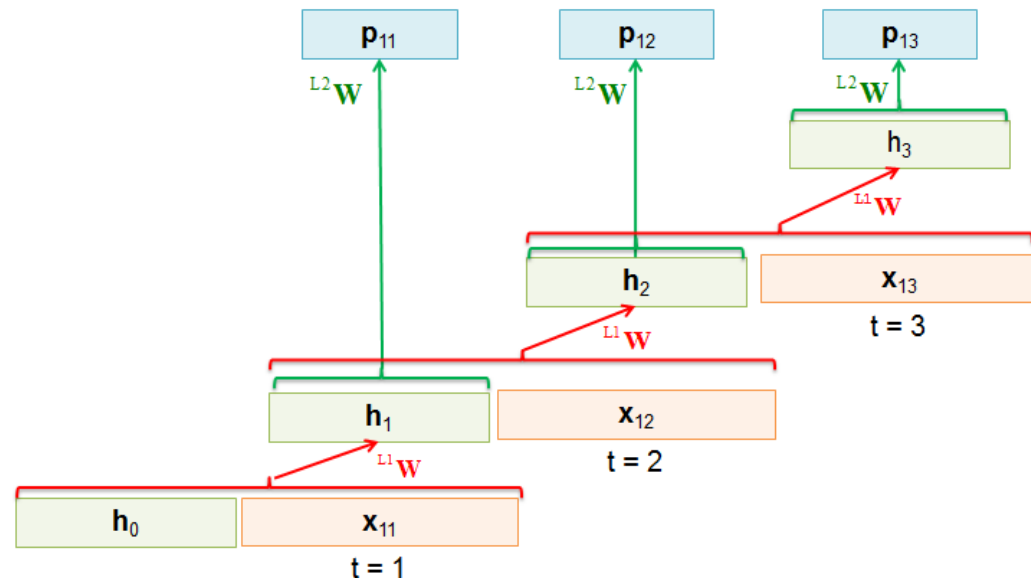
train data input (S=len(seq)=3):

instance_id	seq_t1	seq_t2	seq_t3
1	$\mathbf{x}_{11}$	$\mathbf{x}_{12}$	$\mathbf{x}_{13}$
2	$\mathbf{x}_{21}$	$\mathbf{x}_{22}$	$\mathbf{x}_{23}$
3	$\mathbf{x}_{31}$	$\mathbf{x}_{32}$	$\mathbf{x}_{33}$
⋮	⋮	⋮	⋮
8	$\mathbf{x}_{81}$	$\mathbf{x}_{82}$	$\mathbf{x}_{83}$

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

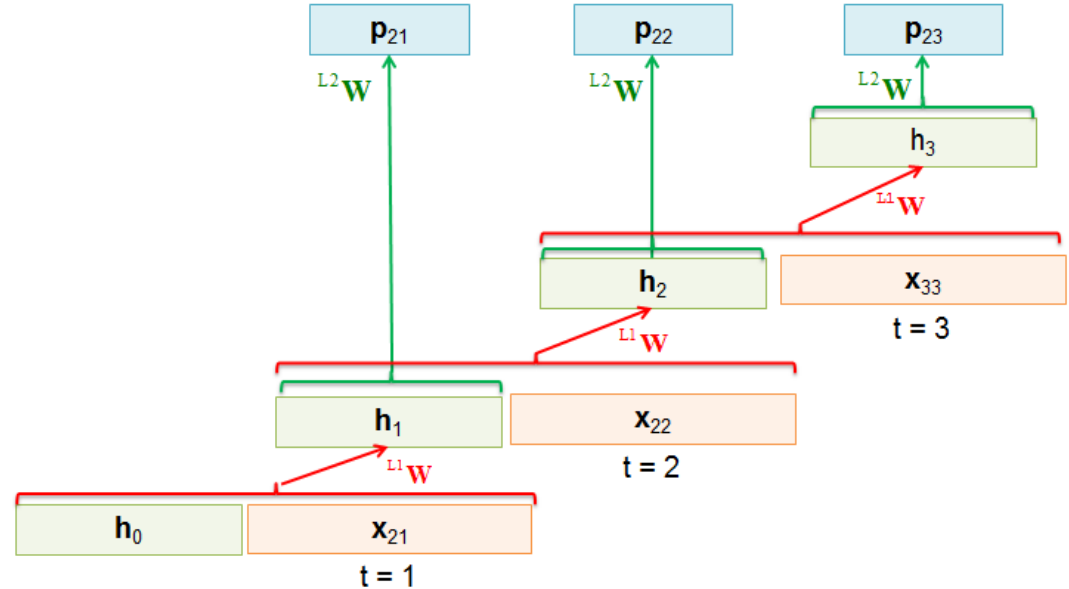
train data input (S=len(seq)=3):

instance_id	seq_t1	seq_t2	seq_t3
1	$\mathbf{x}_{11}$	$\mathbf{x}_{12}$	$\mathbf{x}_{13}$
2	$\mathbf{x}_{21}$	$\mathbf{x}_{22}$	$\mathbf{x}_{23}$
3	$\mathbf{x}_{31}$	$\mathbf{x}_{32}$	$\mathbf{x}_{33}$
⋮	⋮	⋮	⋮
8	$\mathbf{x}_{81}$	$\mathbf{x}_{82}$	$\mathbf{x}_{83}$

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	$\mathbf{x}_{11}$	$\mathbf{x}_{12}$	$\mathbf{x}_{13}$
2	$\mathbf{x}_{21}$	$\mathbf{x}_{22}$	$\mathbf{x}_{23}$
3	$\mathbf{x}_{31}$	$\mathbf{x}_{32}$	$\mathbf{x}_{33}$
⋮	⋮	⋮	⋮
8	$\mathbf{x}_{81}$	$\mathbf{x}_{82}$	$\mathbf{x}_{83}$

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:

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

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

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.models import Sequential
from keras.layers import Dense, SimpleRNN, Activation

model = Sequential()

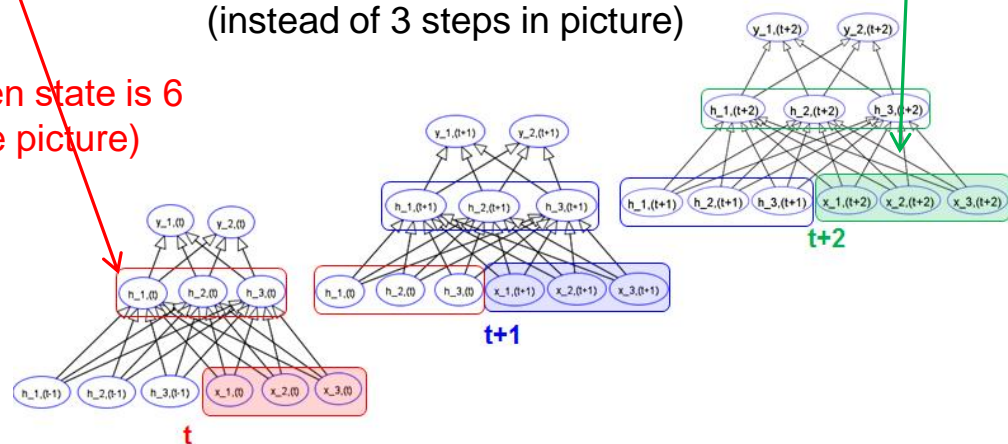
model.add(SimpleRNN(6, return_sequences=True, batch_input_shape=(None, 50, 3)))
model.add(Dense(2))
model.add(Activation("softmax"))

model.summary()
model.compile(loss='categorical_crossentropy', optimizer='adam')
```

“capacity” of hidden state is 6  
(instead of 3 in the picture)

each training sequence  
consists out of 50 elements  
(instead of 3 steps in picture)

length of input  
vector at each  
time step is  
here 3





# Do your own time series predictions

- Work through the instructions in the second exercise in day 6 using [https://github.com/tensorchiefs/dl\\_course\\_2018/blob/master/notebooks/12\\_LSTM\\_vs\\_1DConv.ipynb](https://github.com/tensorchiefs/dl_course_2018/blob/master/notebooks/12_LSTM_vs_1DConv.ipynb)

