

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

Beate Sick

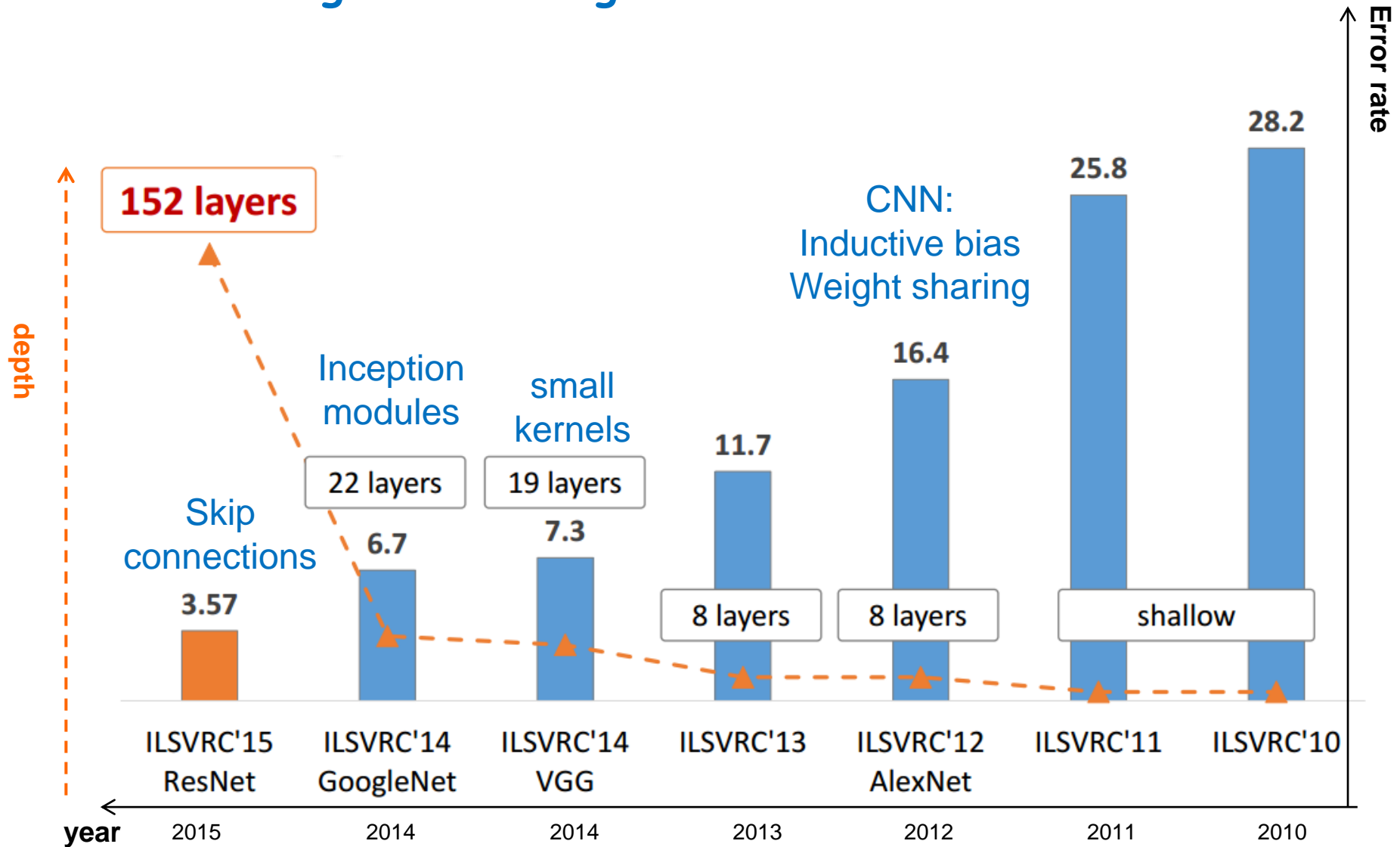
sick@zhaw.ch

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

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” 2014 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

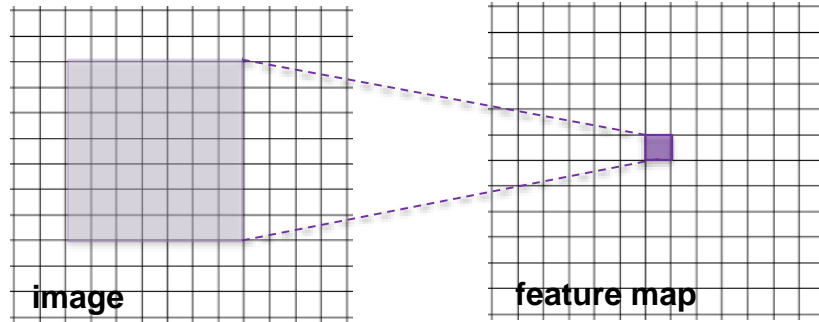
<http://arxiv.org/abs/1409.1556>



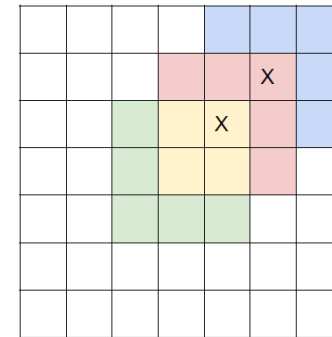
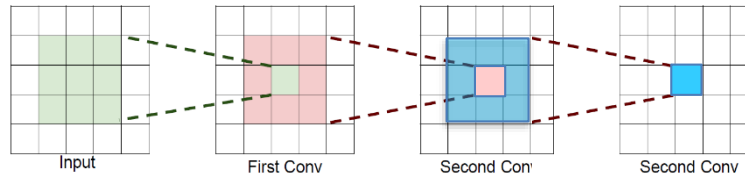
The trend in modern CNN architectures goes to small filters

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

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

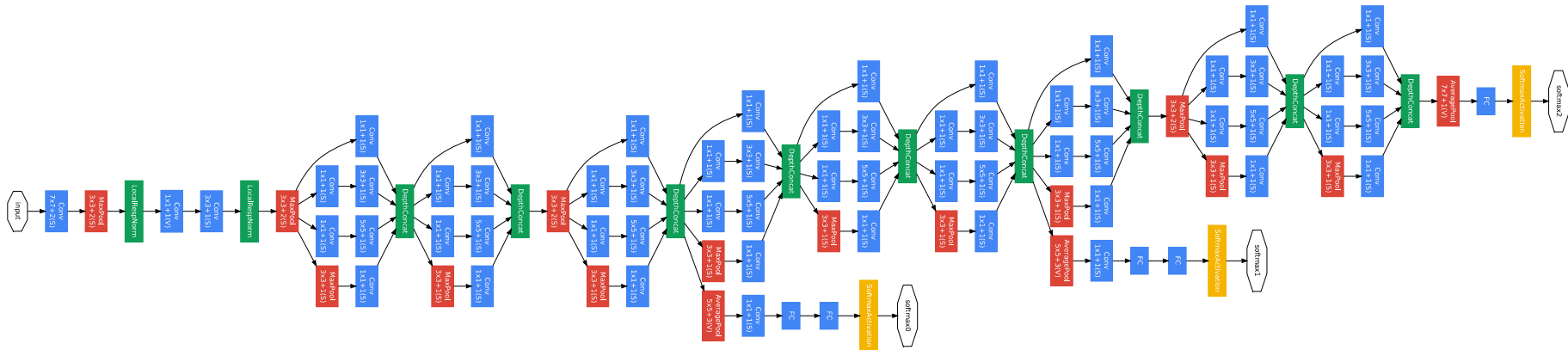


- 2) Working with **three 3x3**
conv layers (stride 1)
3*9=27 weights

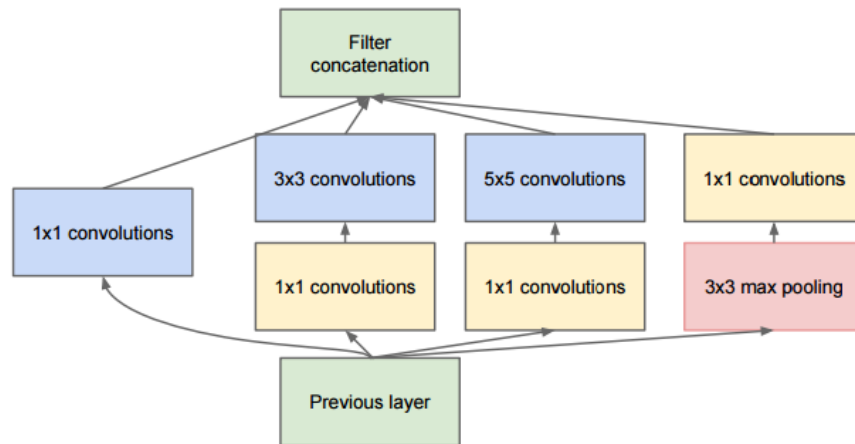


When working stacking conv-layers with small kernels we can achieve with less weights and more non-linear combinations the same receptive field

Winning architecture (GoogLeNet, 2014)



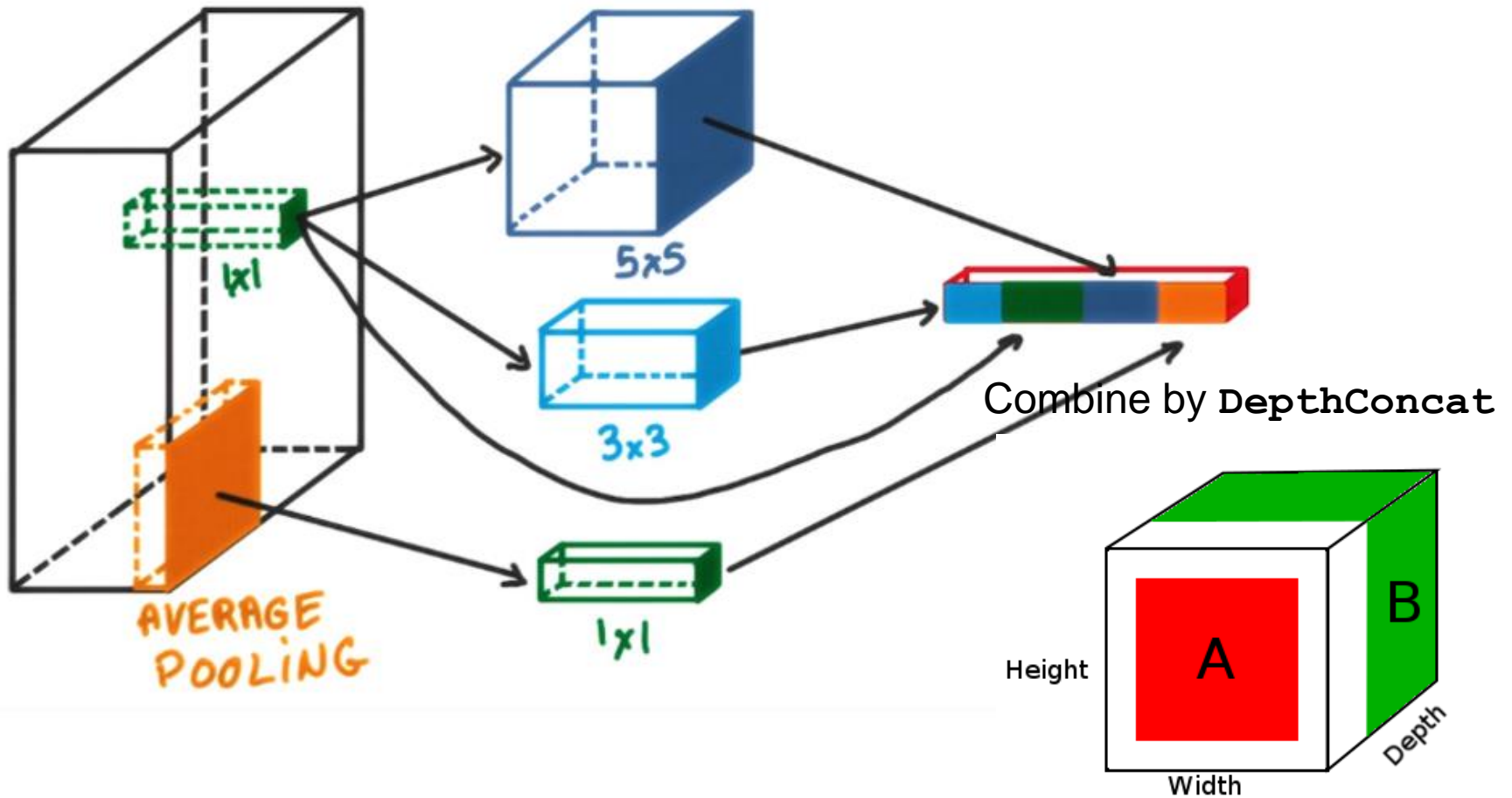
The inception module:
use in 1 layer in parallel different kernels and combine their results



Few parameters, hard to train.
Comments see [here](http://arxiv.org/abs/1409.4842)

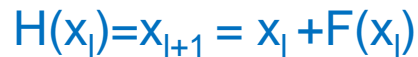
The idea of inception modules

INCEPTION MODULES



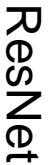
152
layers

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

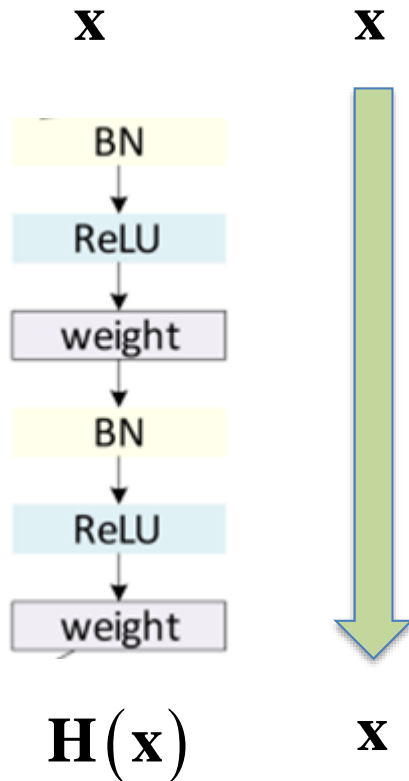


152 layers:
Why does this train at all?

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:

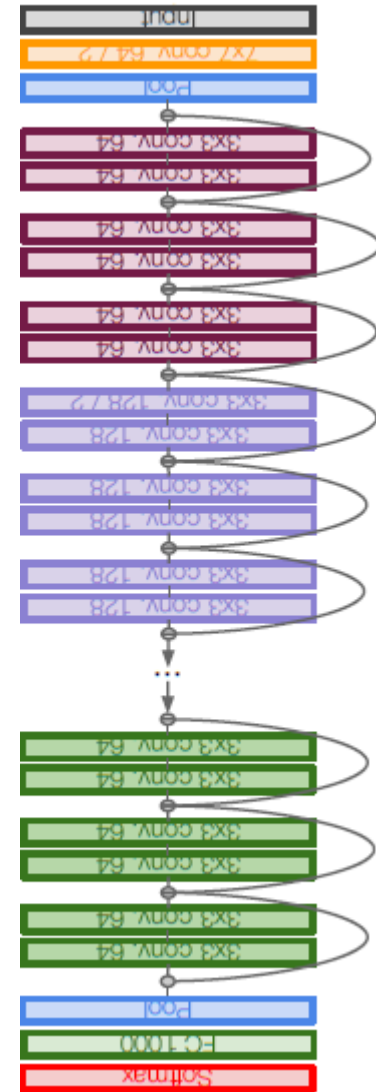
$$\mathbf{y} = \mathbf{H}(\mathbf{x}, \mathbf{W}_H) \cdot T(\mathbf{x}, \mathbf{W}_T) + \mathbf{x} \cdot (1 - T(\mathbf{x}, \mathbf{W}_T)).$$

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:

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



Model zoo: many pretrained NN are out there

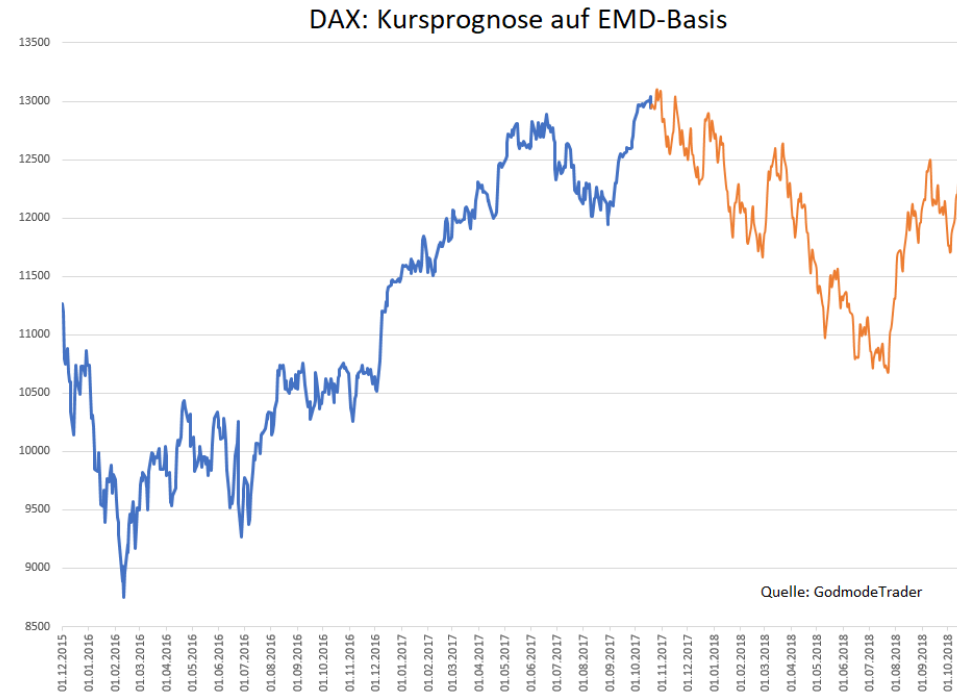
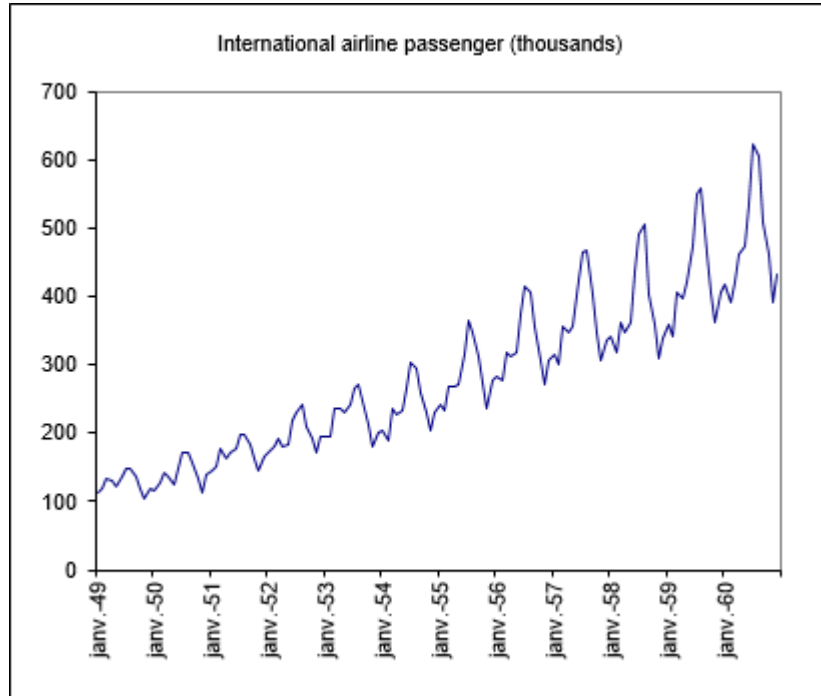
<https://modelzoo.co/>

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

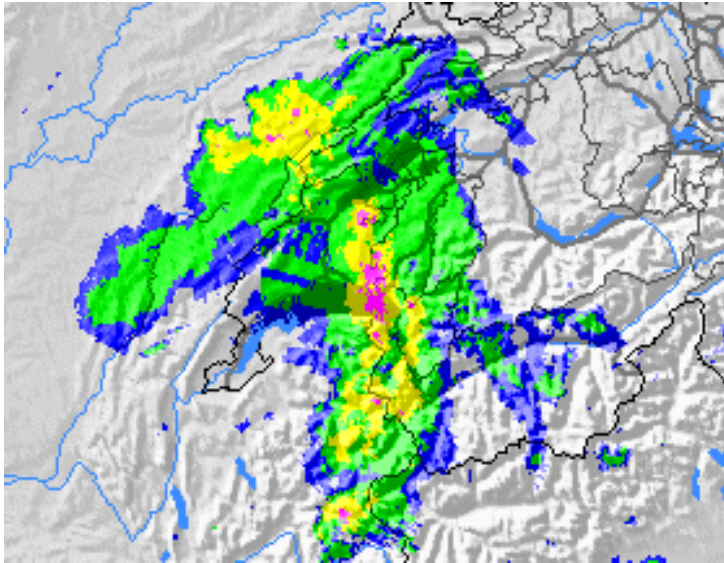
Model Name	Star Rating	Description	Frameworks
Graphics code generating model using Processing	★ 31	swatchnet - processing code generator	PyTorch, CV
imagenet-vgg	★ 31		PyTorch, CV
PyTorch Image Classification with Kaggle Dogs vs Cats Dataset	★ 30	Classifies an image as containing either a dog or a cat using a pre-trained dataset. Out could easily be extended to other image classification problems.	PyTorch, CV
Bidirectional LSTM on the IMDB dataset			Keras
1D CNN on the IMDB dataset			Keras
1D CNN-LSTM on the IMDB dataset			Keras
chainervariational-VAE	★ 31	Variational autoencoder (VAE)	PyTorch, CV, Generative
Hierarchical Attention Network for Document Classification	★ 31	A faster and up to date implementation is in my other repo	PyTorch, NLP
improved-gan	★ 29	code for the paper	PyTorch, Generative
Simple Generative Adversarial Networks	★ 29	Run the sample code by typing:	PyTorch, Generative
multilabel	★ 24	This is the implement of the multilabel image classification in MNIST.	PyTorch, CV
VAE	★ 22	Variational AutoEncoder	PyTorch, CV, Generative
SqueezeNet	★ 20	Top-1 Acc@1 0% on ImageNet without any sacrificing compared with SqueezeNet v1.1.	PyTorch, CV
chainervariational-encoder	★ 19	Neural Encoder-Decoder Machine Translation	PyTorch, CV, NLP
ADDA	★ 18	Adversarial Discriminative Domain Adaptation	PyTorch, CV, Generative
chainervariational-gan-denoising-feature-matching	★ 17	Generative Adversarial Networks with Denoising Feature Matching	PyTorch, CV, Generative
mxnet-audio	★ 17	Implementation of music genre classification, audio-to-text, song recommender and music search in mxnet	PyTorch, CV, Audio
MXSeq2Seq(Gluon star)	★ 16	python seq2seq.py	PyTorch, CV, NLP
Inception v3			Keras
Neural Style Transfer			Keras
Visualizing the filters learned by a CNN			Keras
Deep dreams			Keras
Stateful LSTM			Keras
Siamese network			Keras
DeconvNet	★ 10	Learning DeepConvolution Networks for Semantic Segmentation	PyTorch, CV
Pixel-wise Segmentation on VOC2012 Dataset using PyTorch	★ 10	Pixel-wise segmentation on the VOC2012 dataset using	PyTorch, CV
generative-models	★ 10	Amplified, understandable, and visually interpretable PyTorch implementations of VAE, BBVAE, VGGAN, UNetGAN, VGGAN, VGGAN, LSGAN, SRGAN, BEGAN, ResGAN, InceptionGAN, FGAN, FGGAN	PyTorch, Generative
V-Net	★ 10	Poly Convolutional Neural Networks for Volumetric Medical Image Segmentation	PyTorch, CV
Evolution Strategies	★ 10	Evolution Strategies as a Scalable Alternative to Reinforcement Learning	PyTorch, CV
CycleGAN and Semi-Supervised GAN	★ 10	PyTorch implementation of CycleGAN and Semi-Supervised GAN for Domain Transfer	PyTorch, Generative
Adversarial Generator-Encoder Network	★ 10	This repository contains code for the paper:	PyTorch, CV, Generative
Recurrent Variational Autoencoder	★ 10	Recurrent Variational Autoencoder that generates sequential data implemented in pytorch. 1911.06348, 1908.06619	PyTorch, CV, Generative
AttGAN	★ 10	Tensorflow implementation of AttGAN - arbitrary face attribute editing. Only Change What You Want	TensorFlow, Generative
BEGAN in PyTorch	★ 10	BEGAN: Boundary Equilibrium Generative Adversarial Networks 1902.10717	PyTorch, CV, Generative
Neural machine translation between the writings of Shakespeare and modern English using TensorFlow	★ 10	This performs a morphological translation, going from modern English to Shakespeare and vice versa.	TensorFlow, NLP
img_classification_pk_pytorch	★ 10	Quickly comparing your image classification models with the state-of-the-art models (such as ResNet, VGG, ...)	PyTorch, CV
PNASNet.pytorch	★ 10	PyTorch implementation of PNASNet-S on ImageNet	PyTorch, CV
U-Net	★ 10	For Brain Tumor Segmentation	TensorFlow, CV
CNN-LSTM-CTC	★ 10	I realize three different models for text recognition, and all of them consist of CTC loss layer to realize the segmentation for text images.	PyTorch, CV

Sequence Data

Example Sequence Data: time-series



Example Sequence Data: videos



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



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

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

Important Notice

Please note that starting from October 30, 2015 we will be introducing new online banking authentication procedures in order to protect the information of our online banking users.

You are required to confirm your personal details with us as you will not be able access our online service until this has been done. As you're already registered for online banking all you need to do is to confirm your online banking details.

Get Started

Once you've completed this process you'll be able to have full access to our online banking service.

Best regards,
Natwest Online Banking Team

Guten Tag, Sick Beate (sick)

siehe Anhang

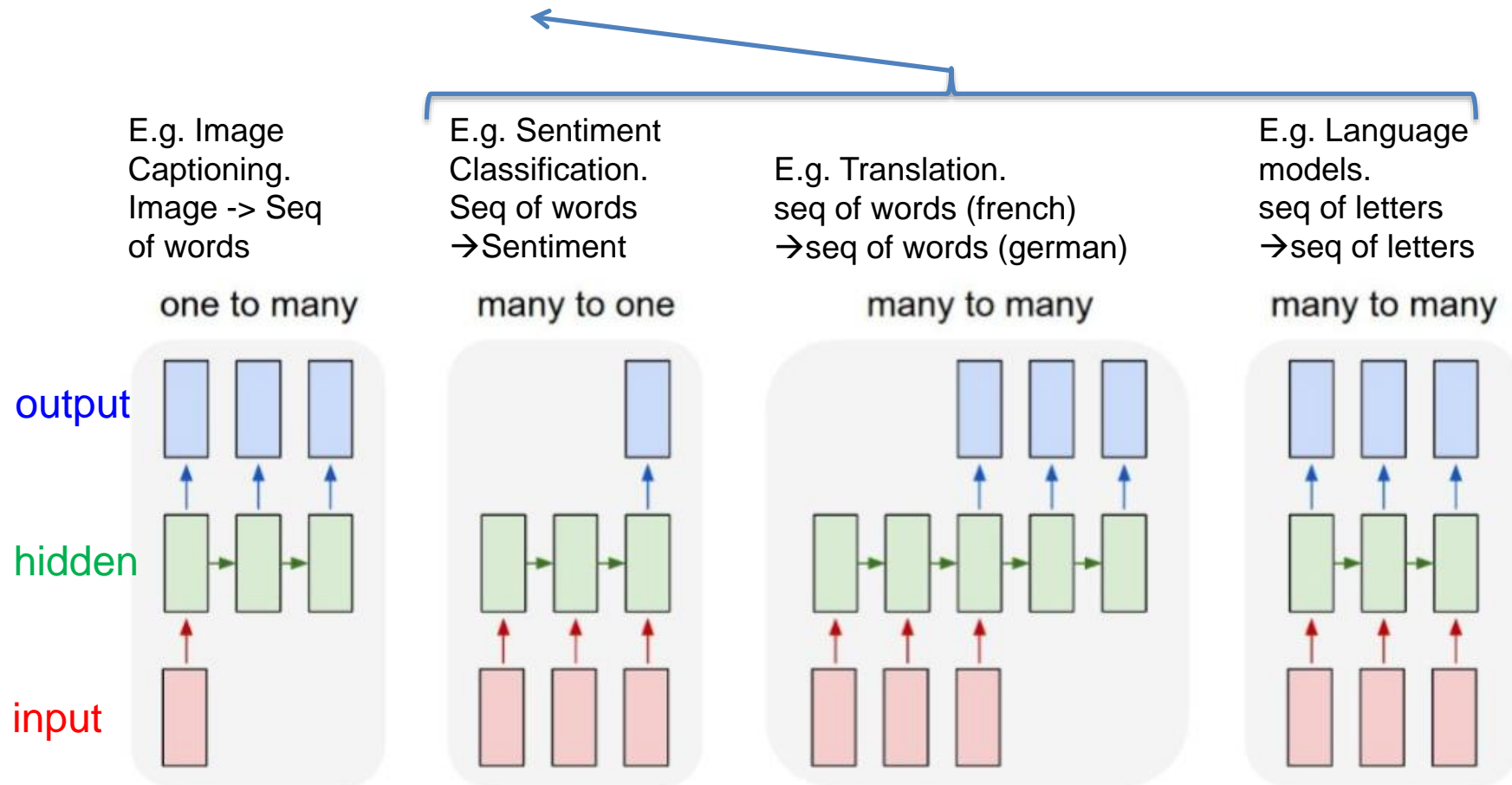
gescanntes Dokument.:

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

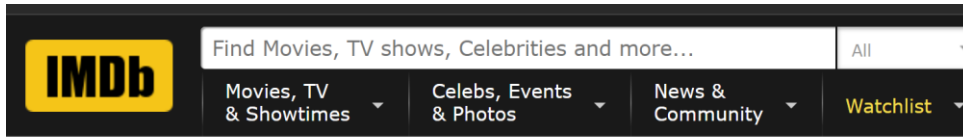
Mit freundlichen GrÃ¼Ãen

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



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](#)



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).

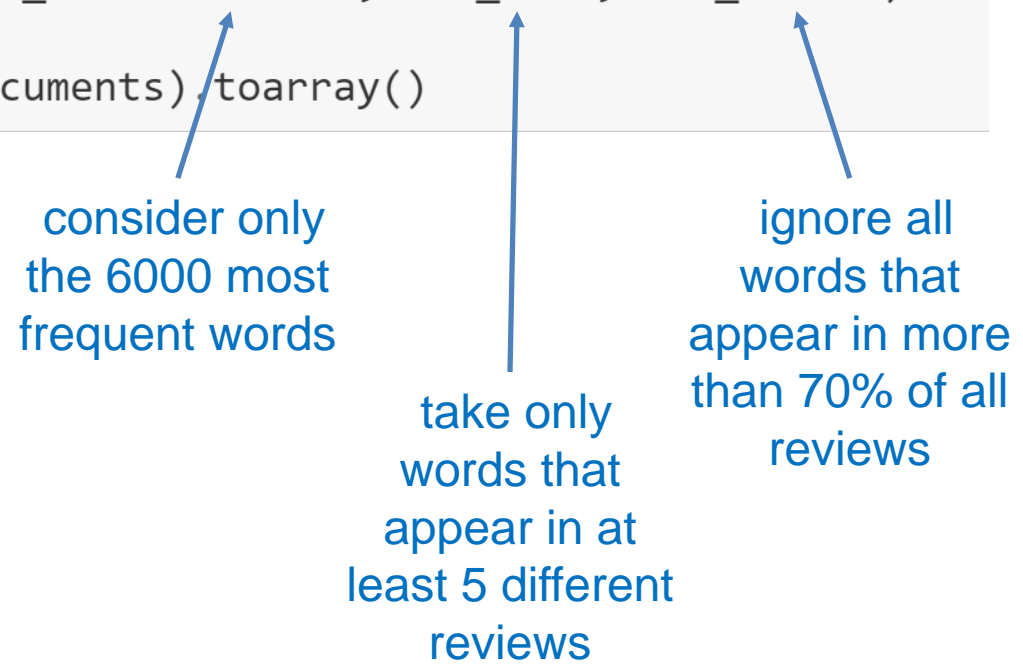
It allows to **represent each sentence as a numeric vector of the same length!**

It 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

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)
```

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

```
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))
```

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

Turns strings into lists of integer indices

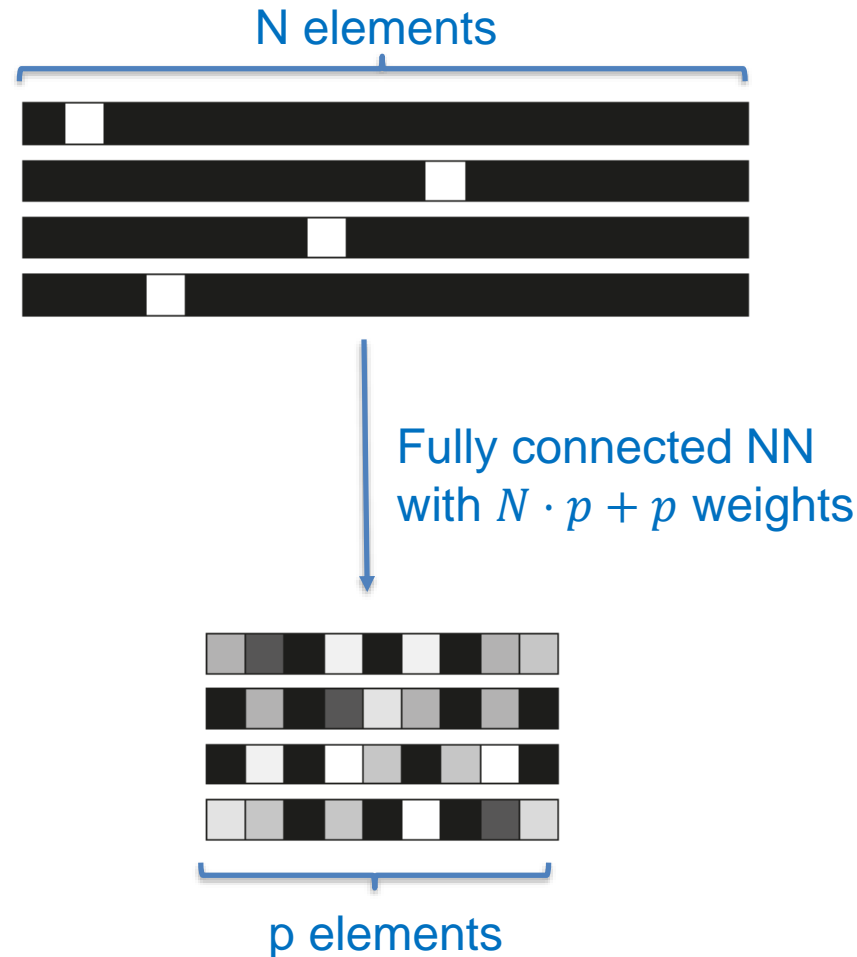
How you can recover the word index that was computed

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 $N \rightarrow p$

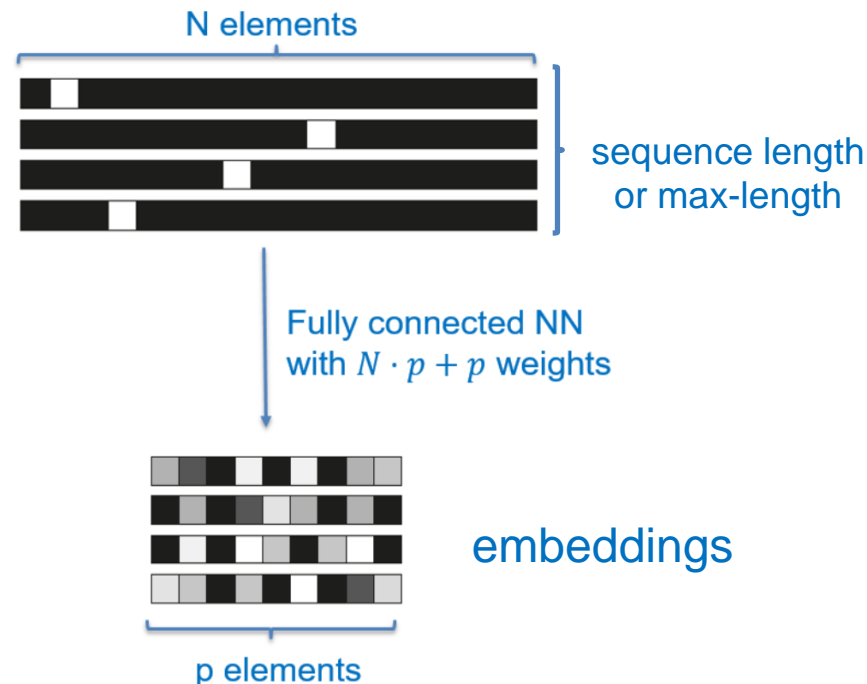
Wordembedding layer in keras

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

← 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).

- Flatten results in vector of length: $\text{seq-length} \cdot p$
- Pool over sequence length results in vector of length: p
- Input to RNN layer or a 1D convolution layer



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

A text classification NN using pooled embeddings

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

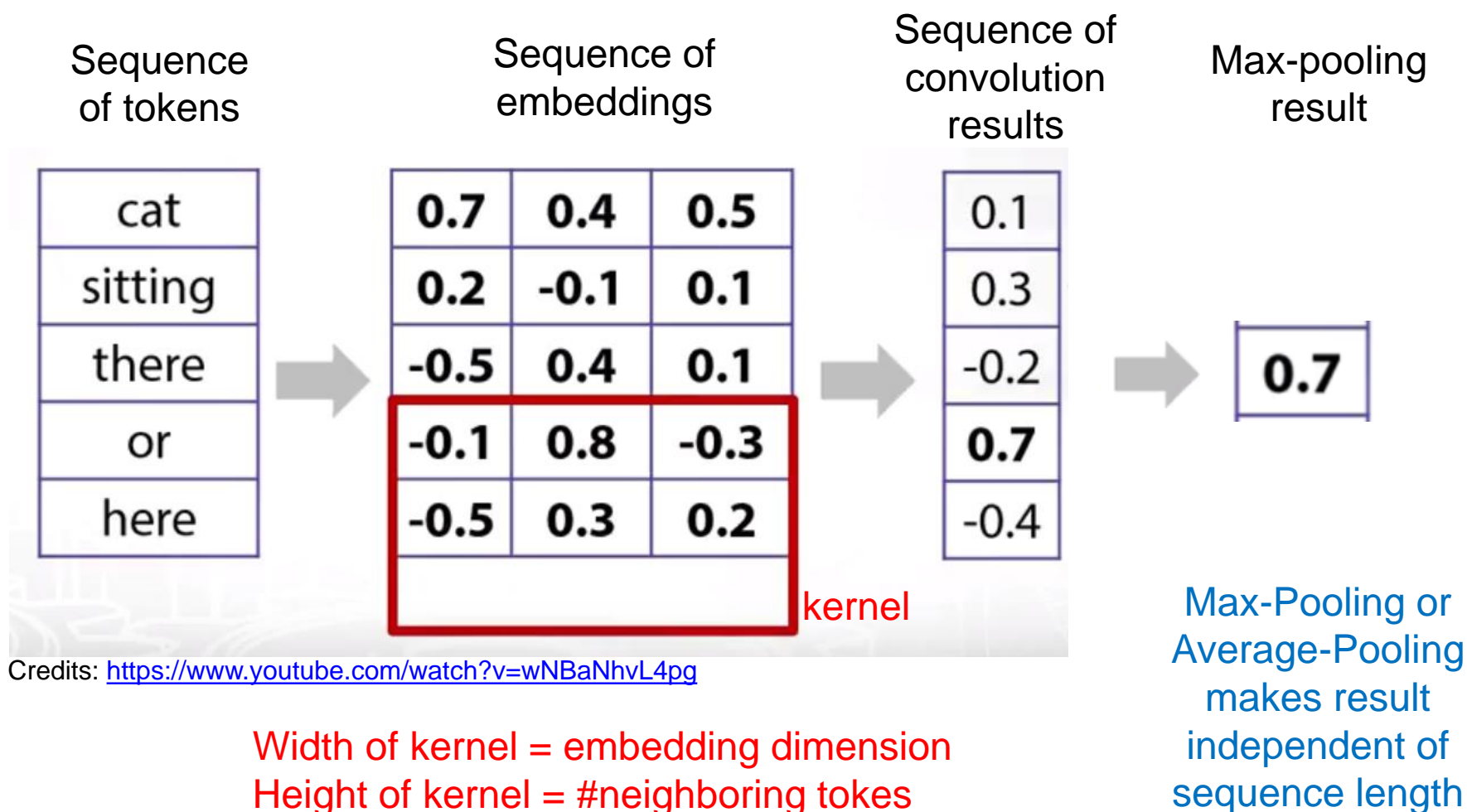
EMBEDDING_DIM = 30

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
=====		

Applying 1D convolution embedding sequences



We slide kernel only in direction → 1D convolution

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

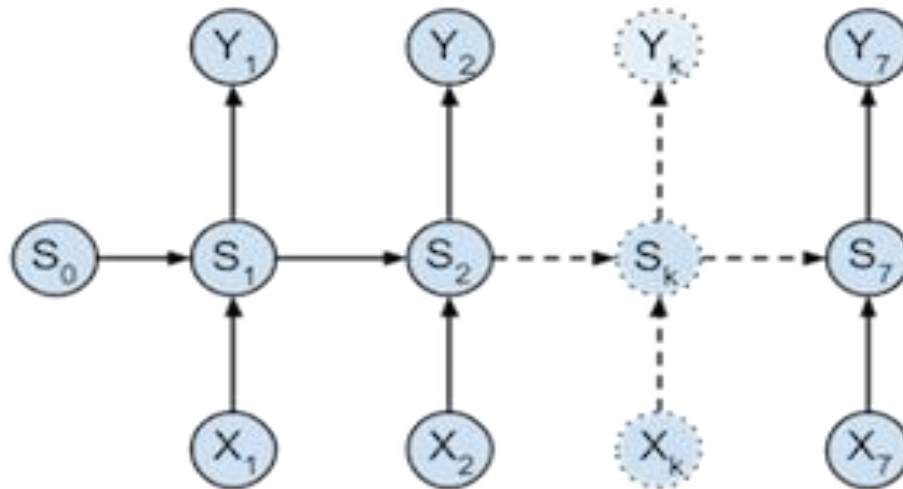
Challenge 2: How to handle input of different lengths?

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

By reading from sentence to sentence 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 sentence

S_2 =state after first sentence

...

Input: $x_1 = \text{vector}_{\text{sentence1}}$, $x_2 = \text{vector}_{\text{sentence2}}$,

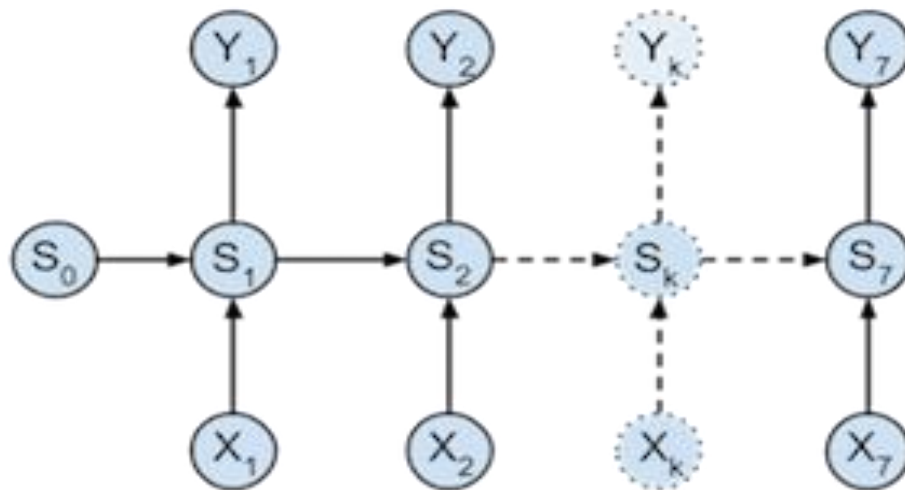
How to choose the dimensions of the hidden state?

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 sentence
 S_2 =state after first sentence
...

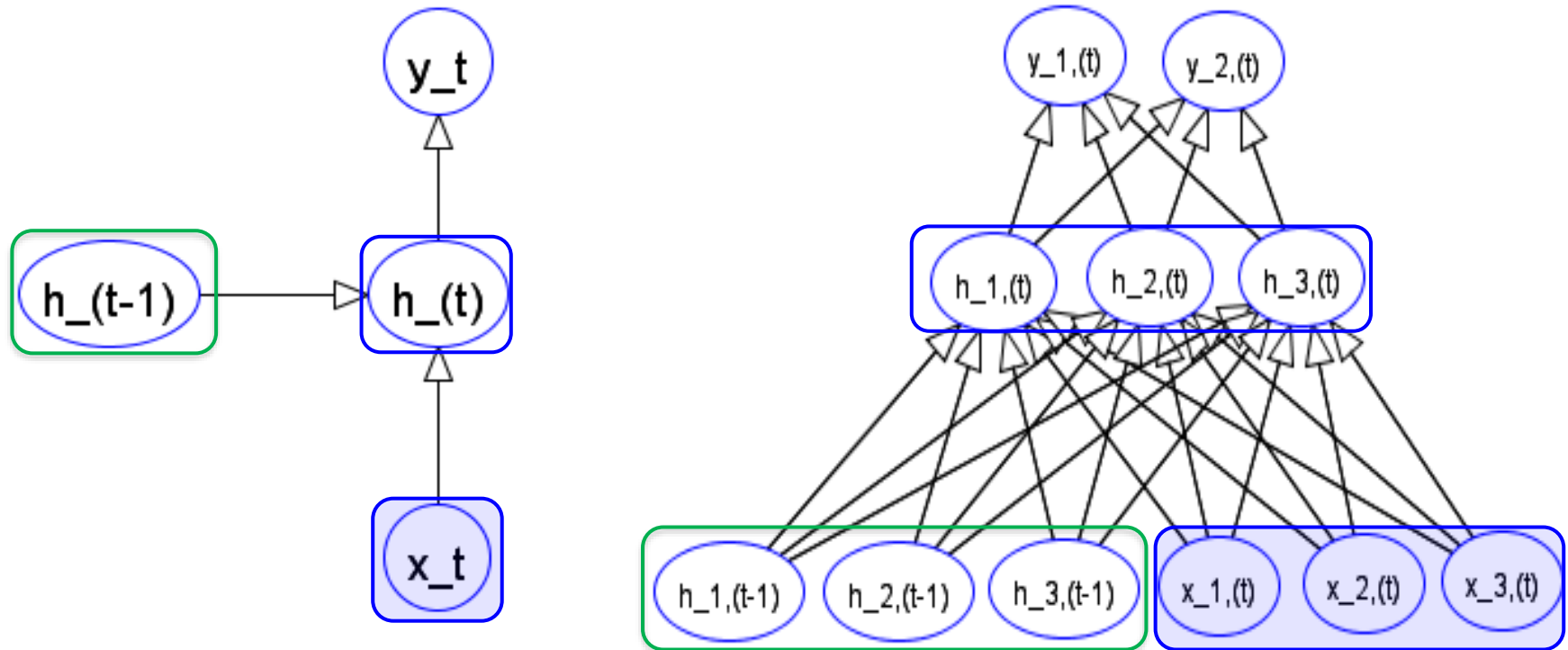
Input: x_1 =vector_{sentence1}, x_2 =vector_{sentence2},

Two representations of a RNN at time t

\mathbf{x} has 3 components – a very small vocabulary ;-)

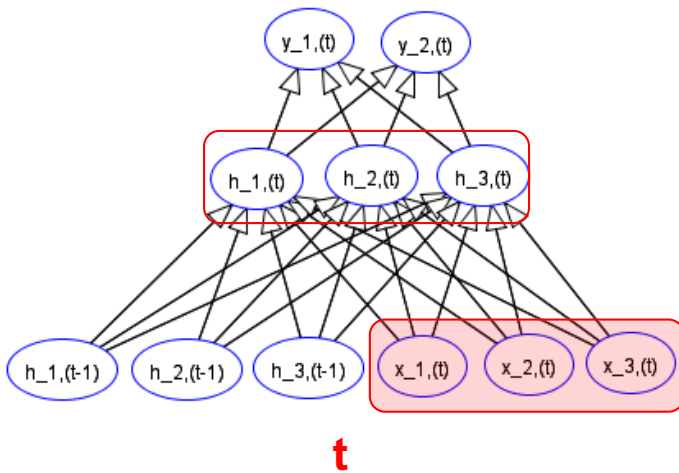
\mathbf{y} has 2 components (spam/no-spam for mails, passed/failed for scoring essays).

\mathbf{h} 3 components to capture abstract concepts (e.g. fraud/sex/conference for emails, or copied/boring/original for essays) and is initialized at $t=0$ with $(0,0,0)$.

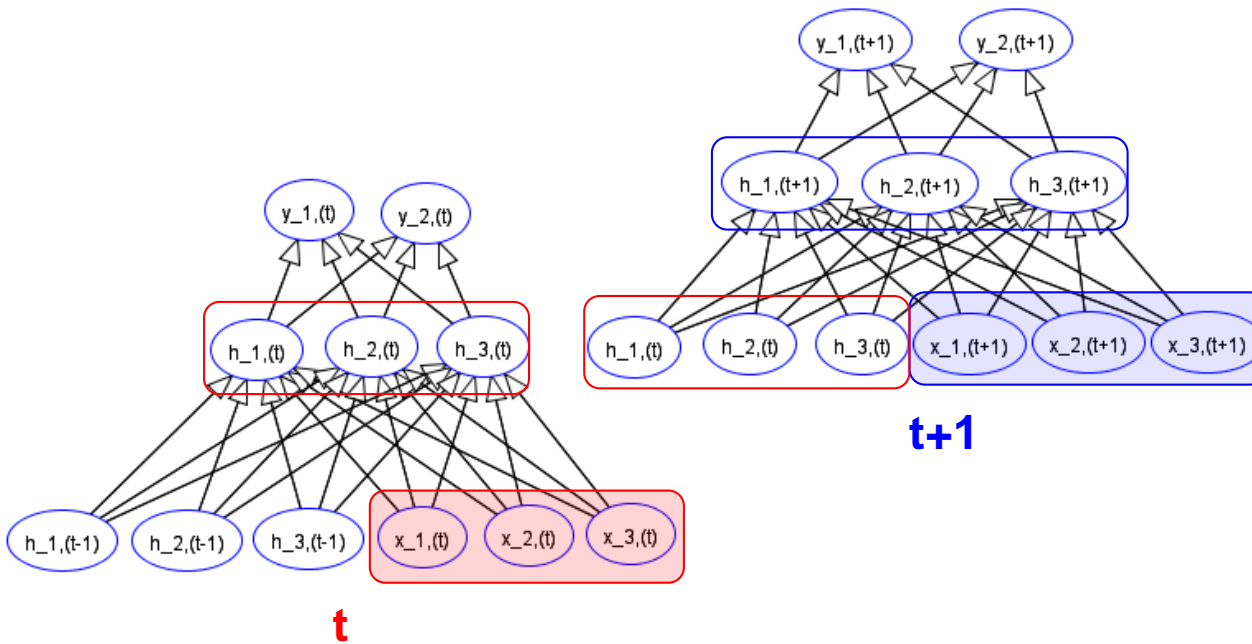


Stepping through an RNN

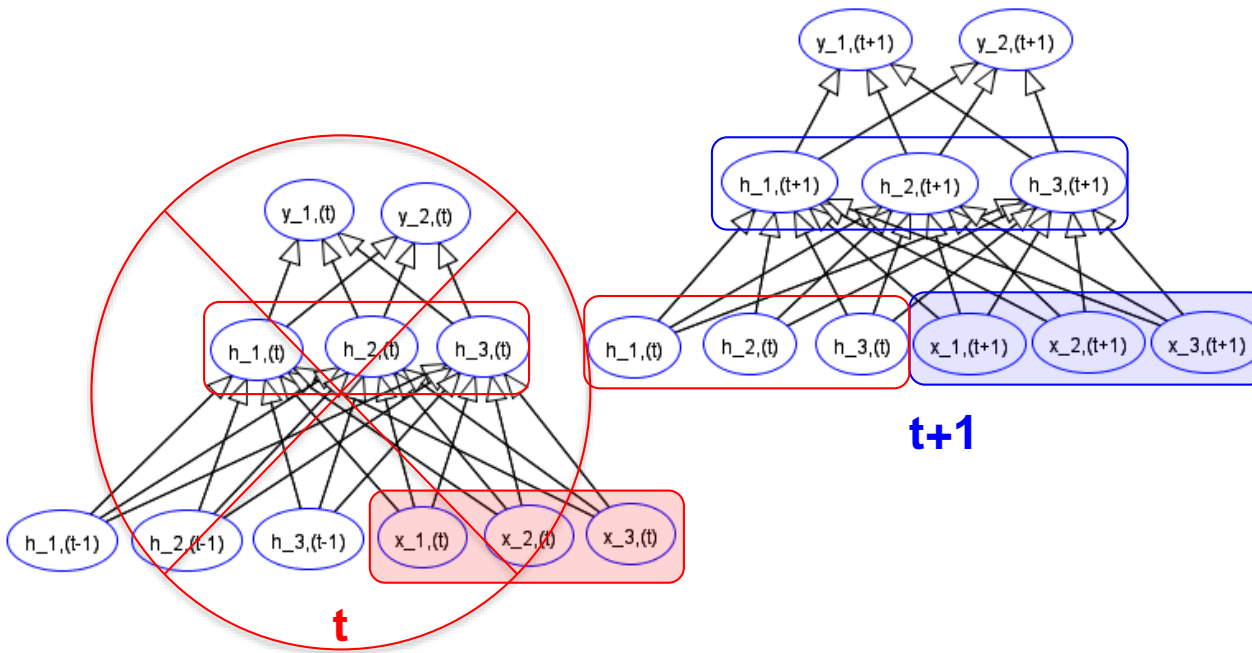
A simple RNN at 3 successive time steps



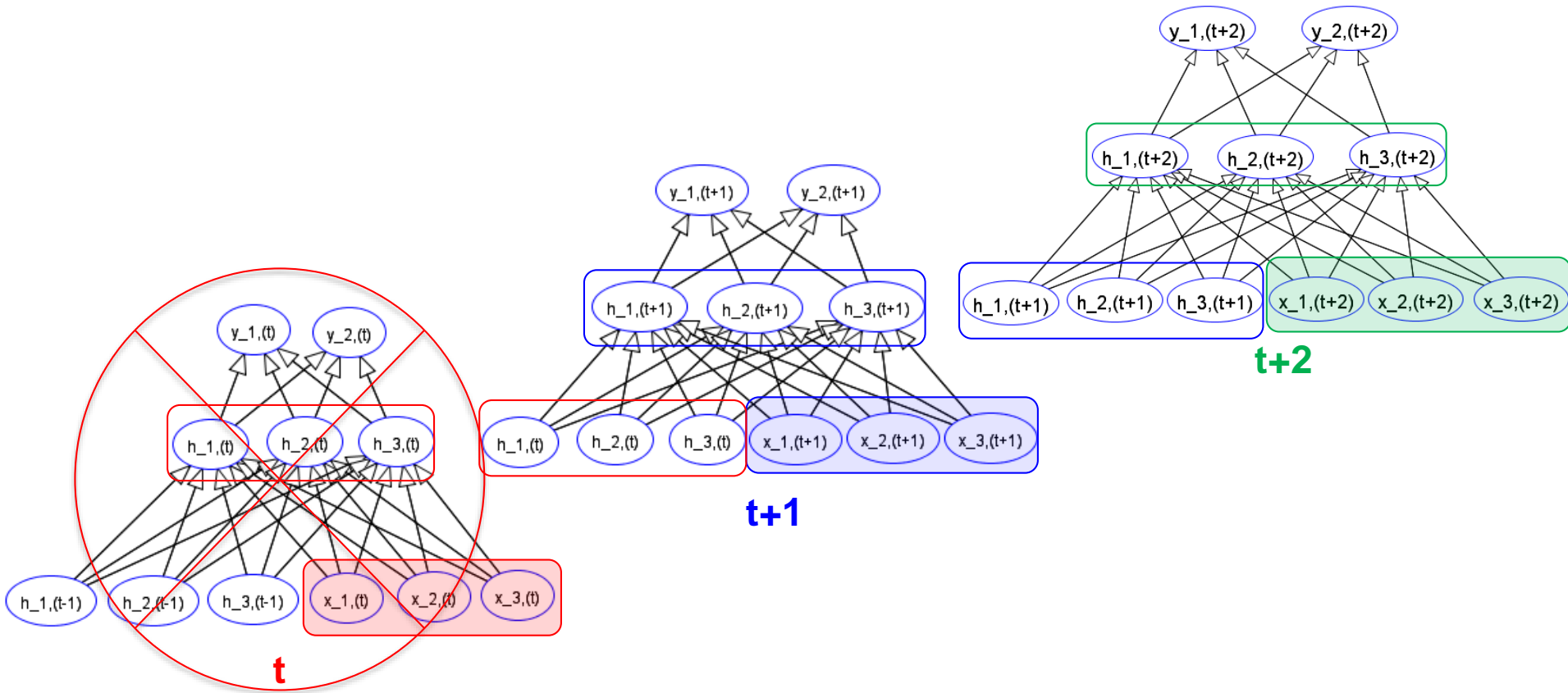
A simple RNN at 3 successive time steps



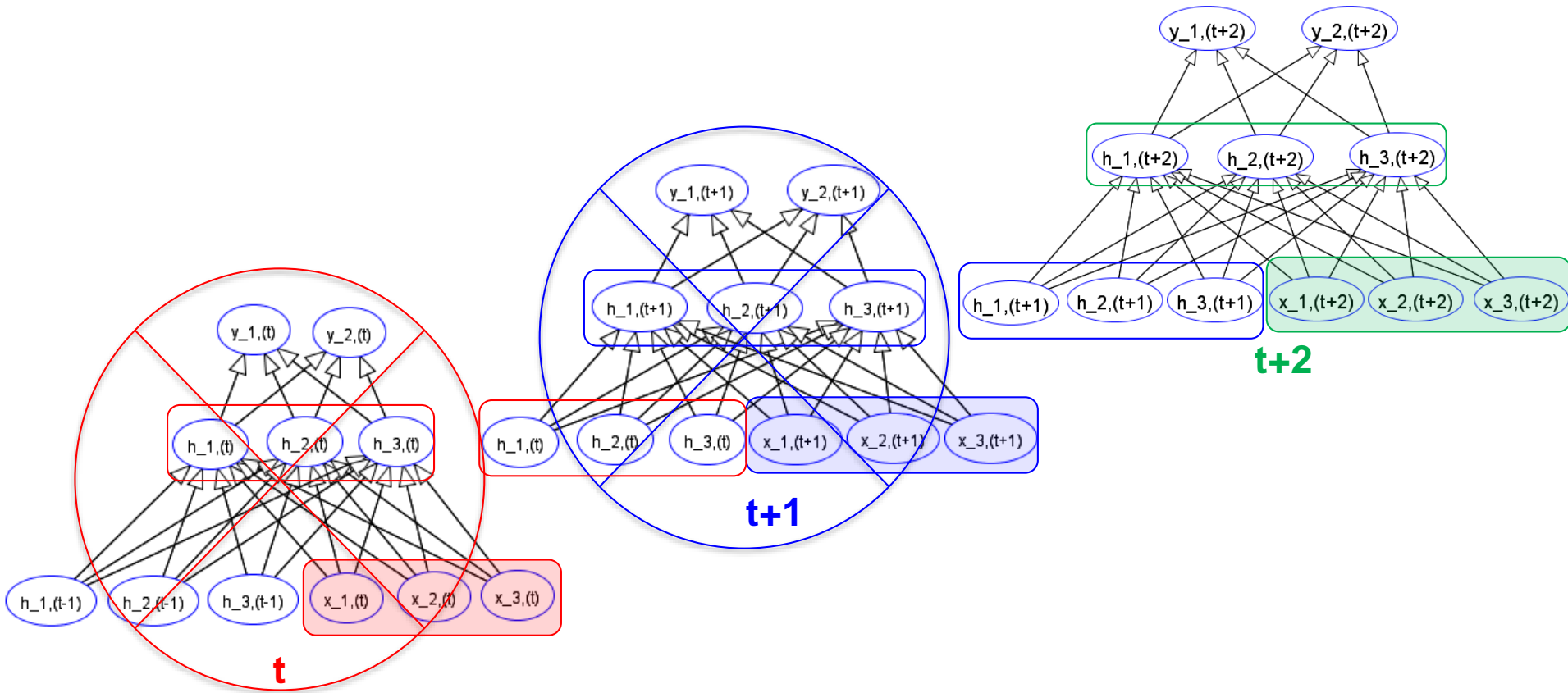
A simple RNN at 3 successive time steps



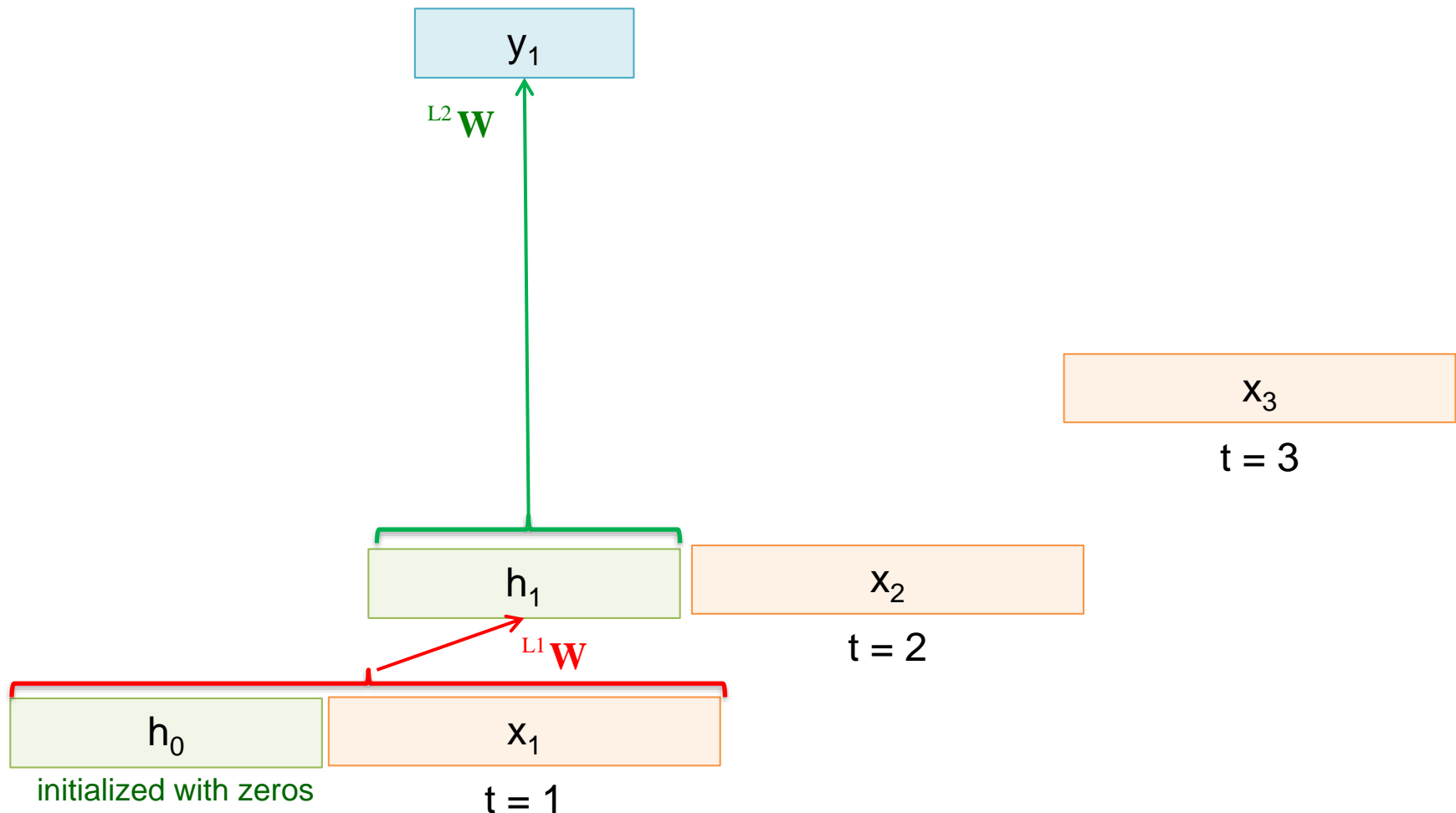
A simple RNN at 3 successive time steps



A simple RNN at 3 successive time steps

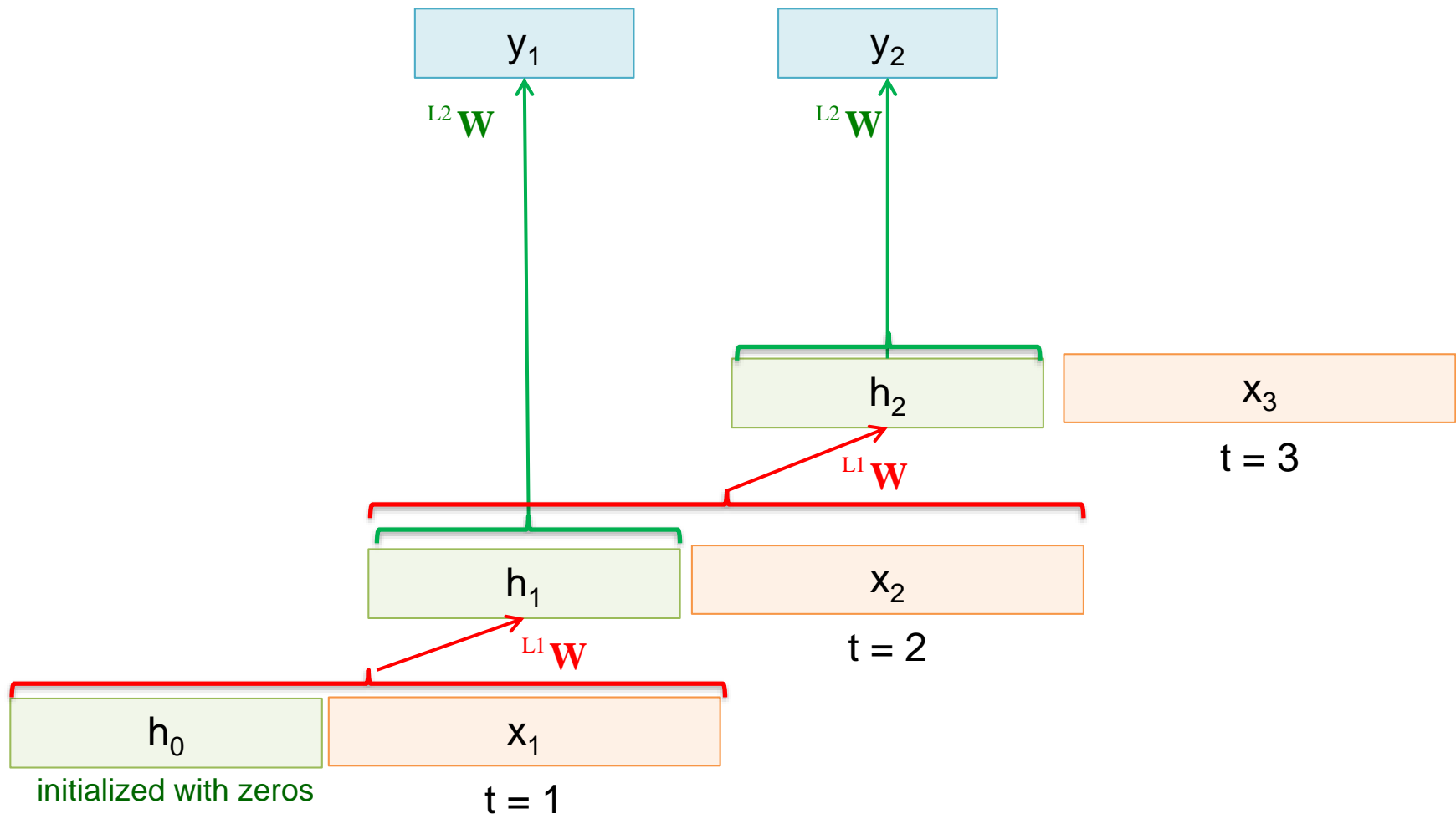


An RNN shares weights across all time steps



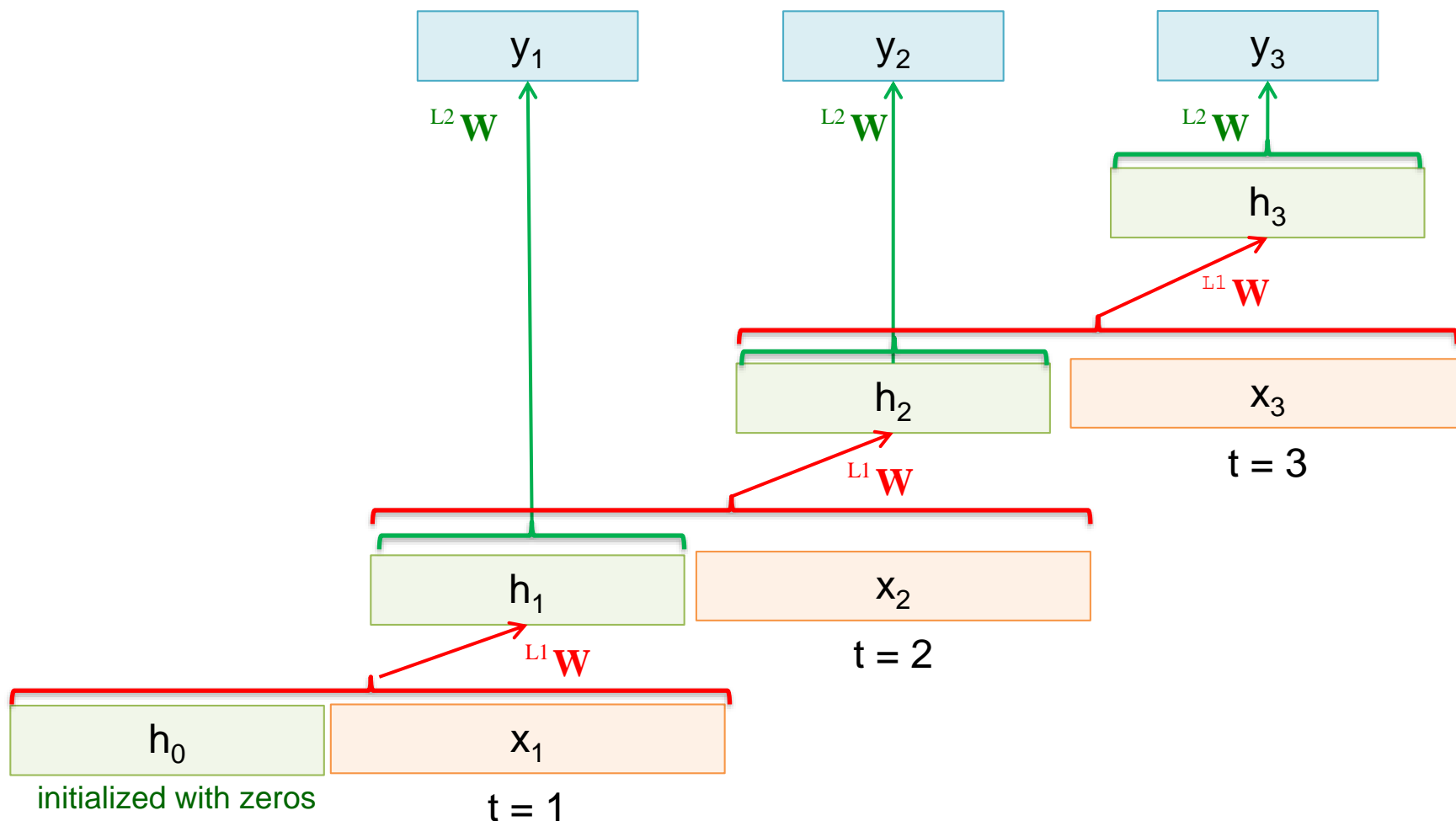
Imagine a trained RNN with fixed weight matrices in layer 1 and layer 2.

An RNN shares weights across all time steps



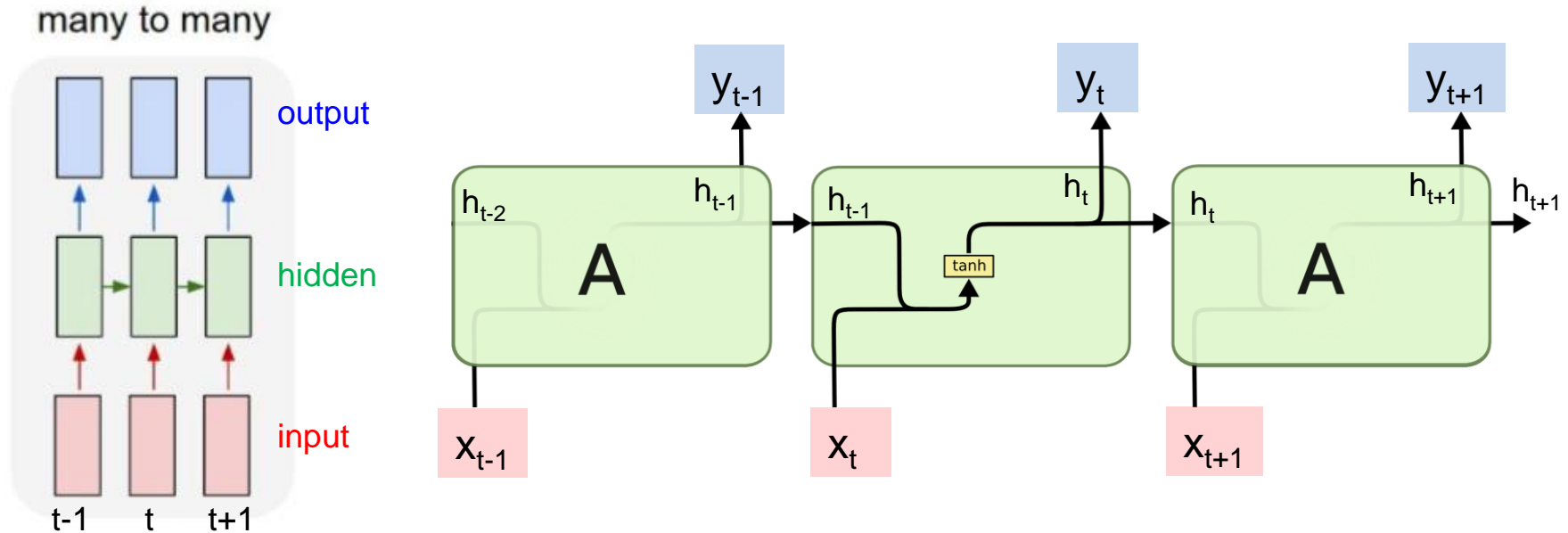
We use at each time step the same weight matrices between the layers!

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

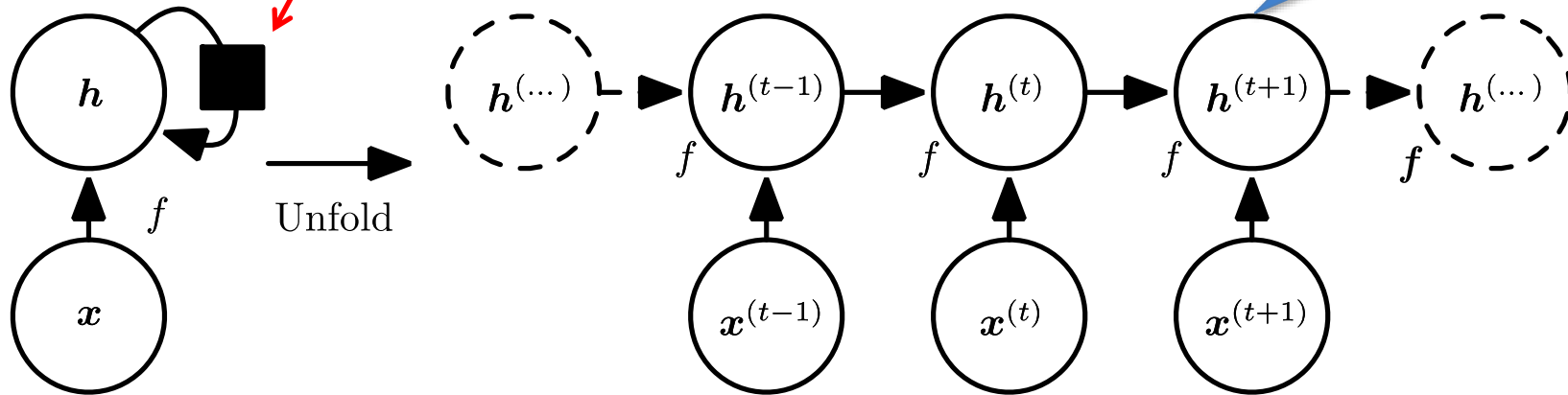


$$\mathbf{A} = f_{\mathbf{W}}(\mathbf{h}_{t-1}, \mathbf{x}_t) = \tanh([\mathbf{h}_{t-1}, \mathbf{x}_t] \cdot \mathbf{W} + \mathbf{b})$$

Using a circuit diagram to represent an RNN

Use same weights
in each time step

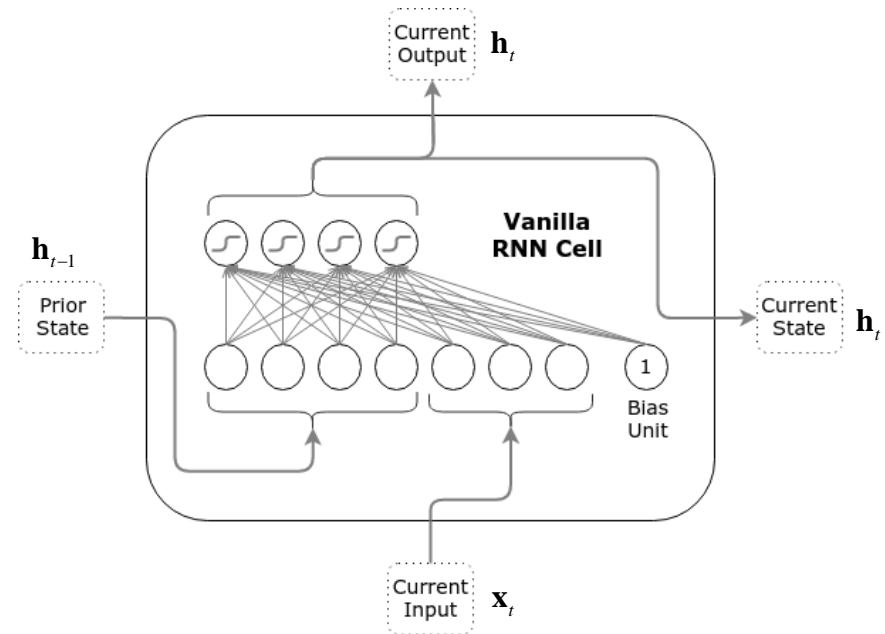
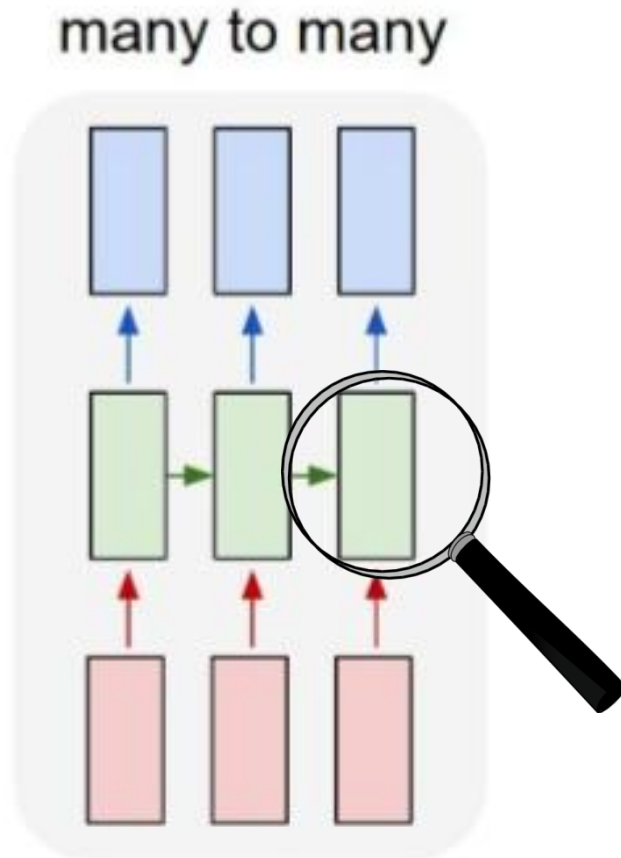
$h^{(t)}$ summarizes /
abstraction



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

Right: Unrolled / unfolded

Looking into a RNN “cell”



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

Loss construction in an RNN

Determine the loss contribution of instance 1

mini-batch of size M=8

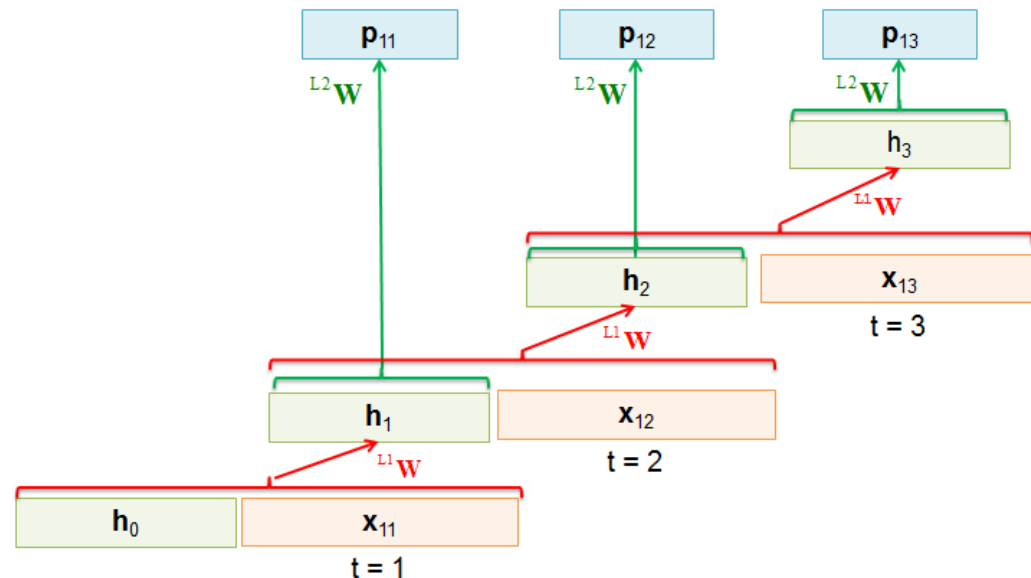
train data input (S=len(seq)=3):

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

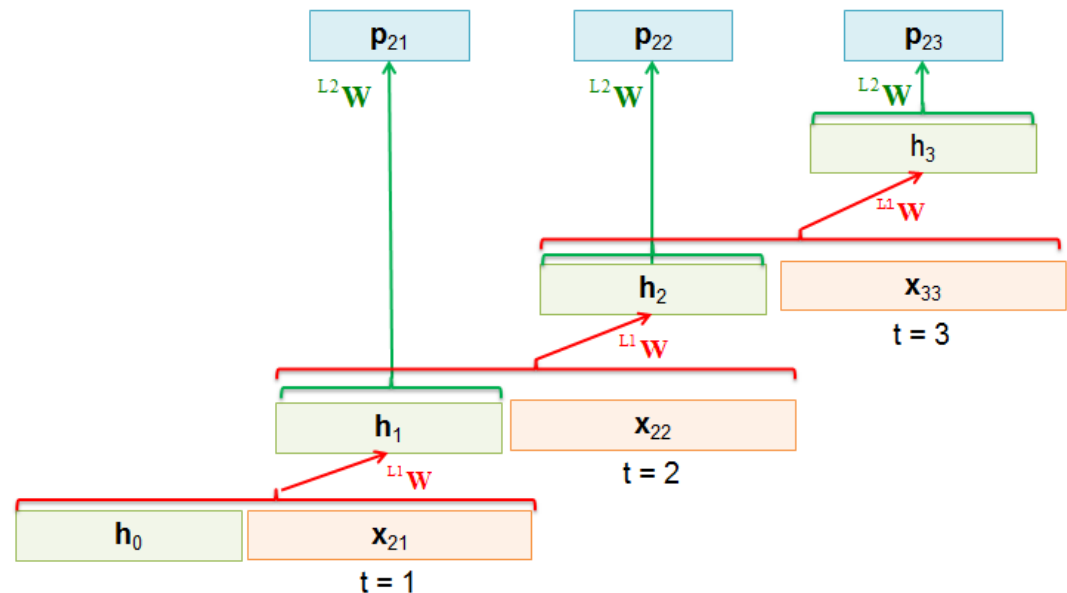
train data input (S=len(seq)=3):

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.

A simple forward pass



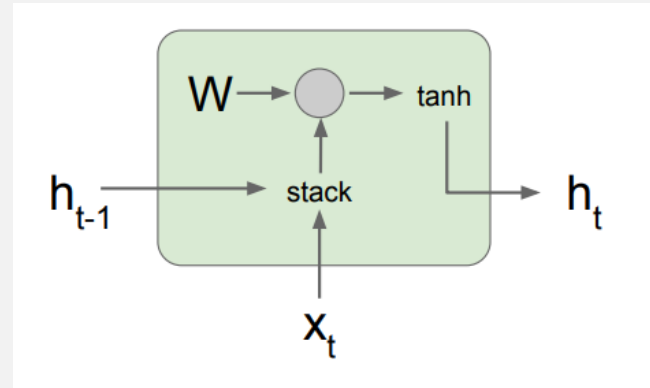
- Given the hidden state at t-1, the input x at t and the weight matrix as:

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

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

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

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



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

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

Solution

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

$$x_t = (1, 0)$$

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

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

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

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'))
```

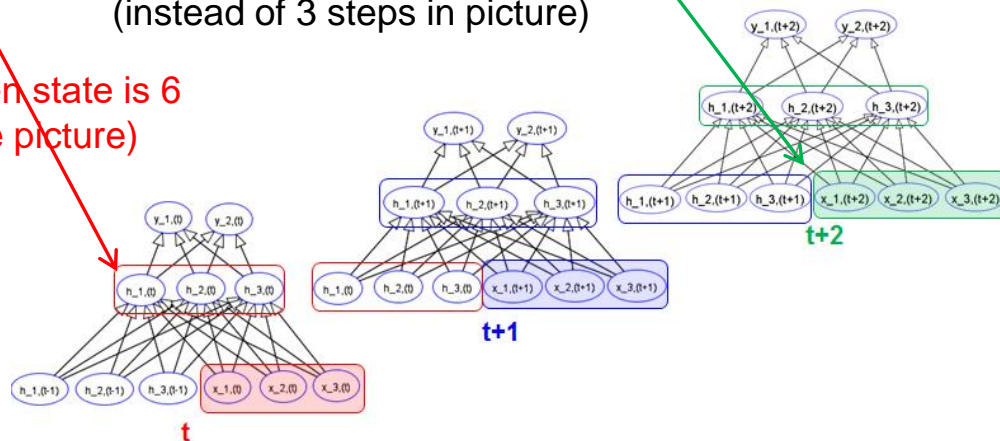
```
model.compile(loss='categorical_crossentropy', optimizer='rmsprop', metrics=['accuracy'])
```

at each step we use hidden state as input to a fcNN with 2 output nodes

length of input vector at each time step is here 3

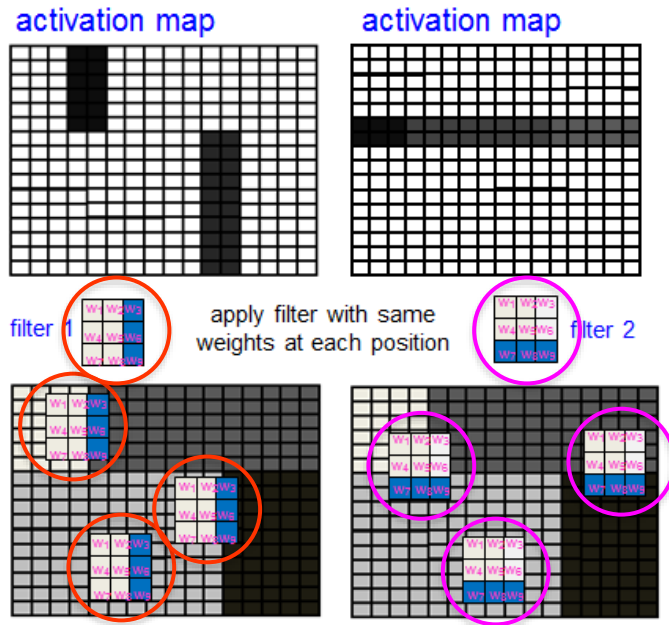
each training sequence consists out of 50 elements (instead of 3 steps in picture)

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

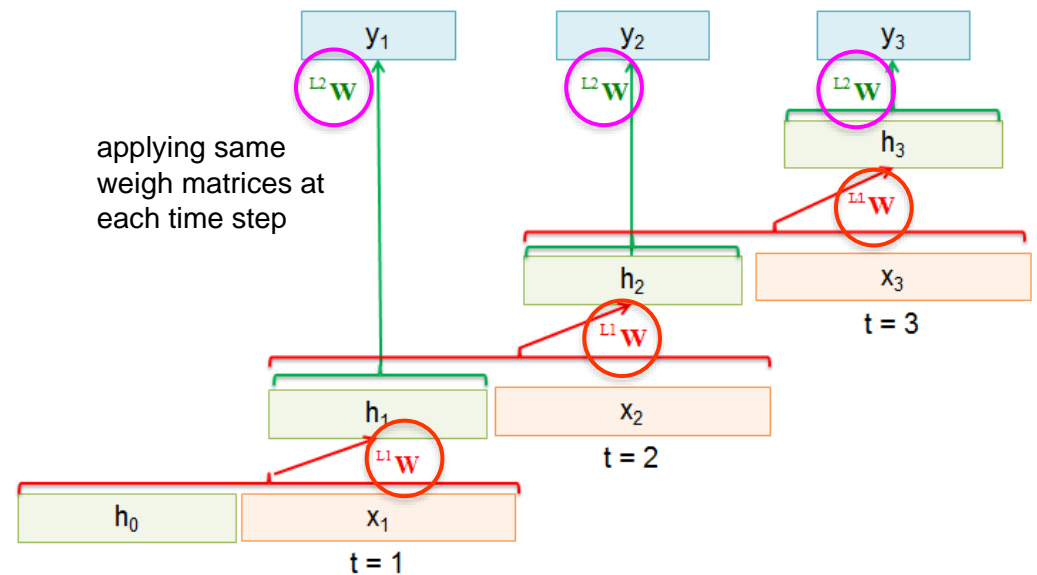


Common tricks in RNN & CNN and some differences

CNN and Recurrent Network share weights



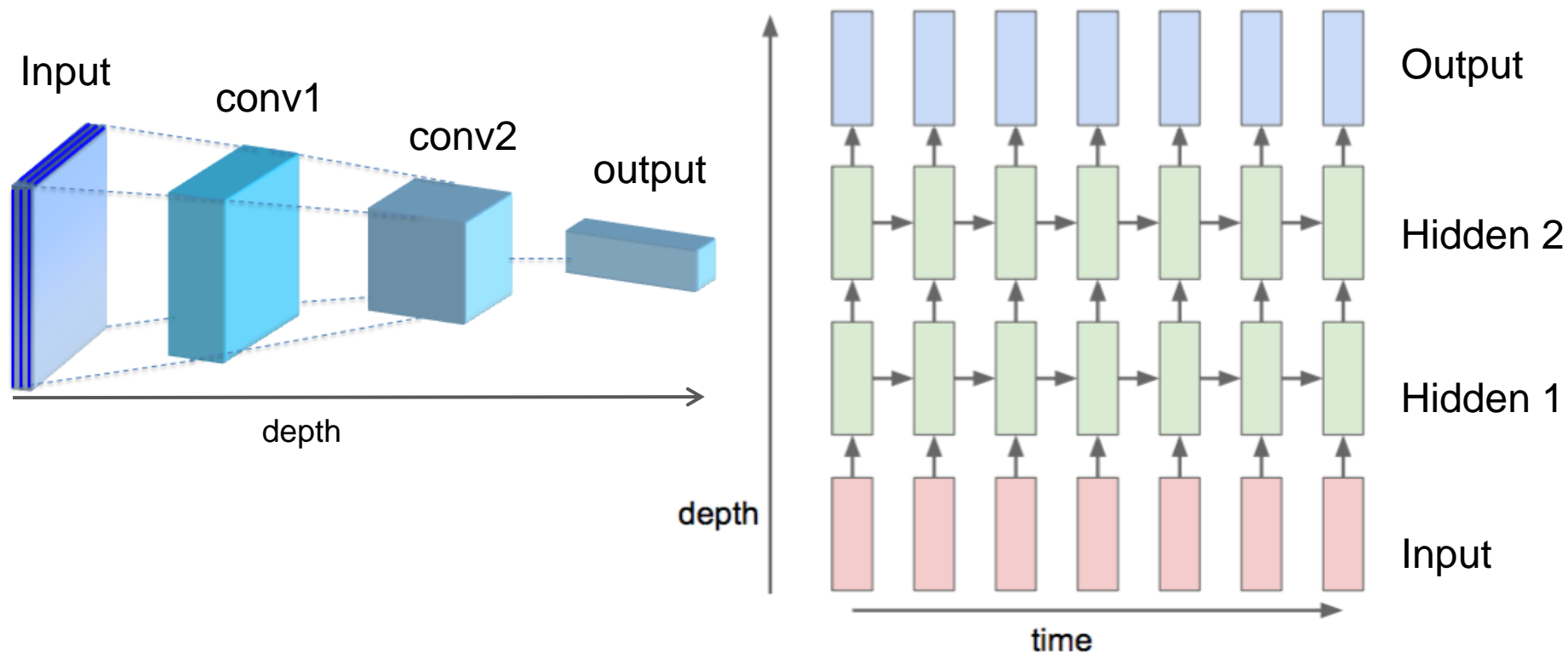
CNN share weights between different local regions of the image



RNN share weights between time steps

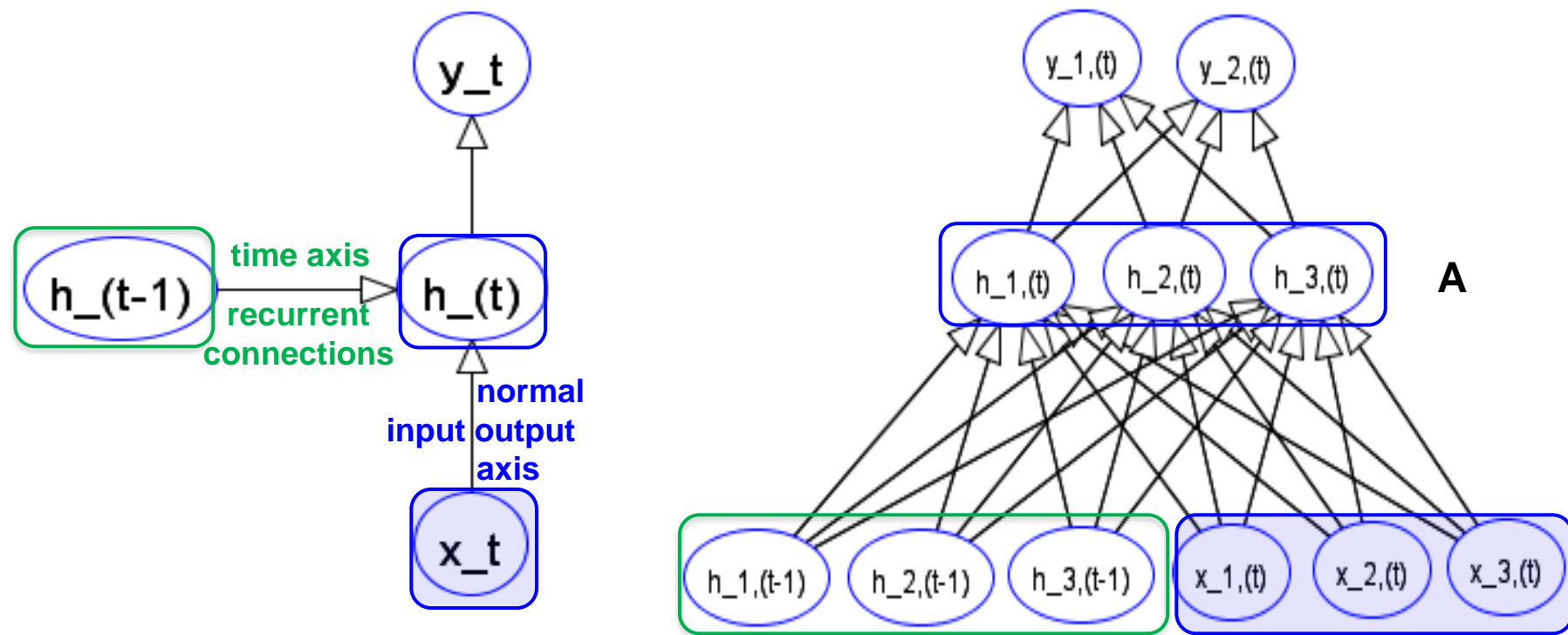
Remark: no weight sharing in fully connected NN

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 different dropout rates for recurrent and normal connections



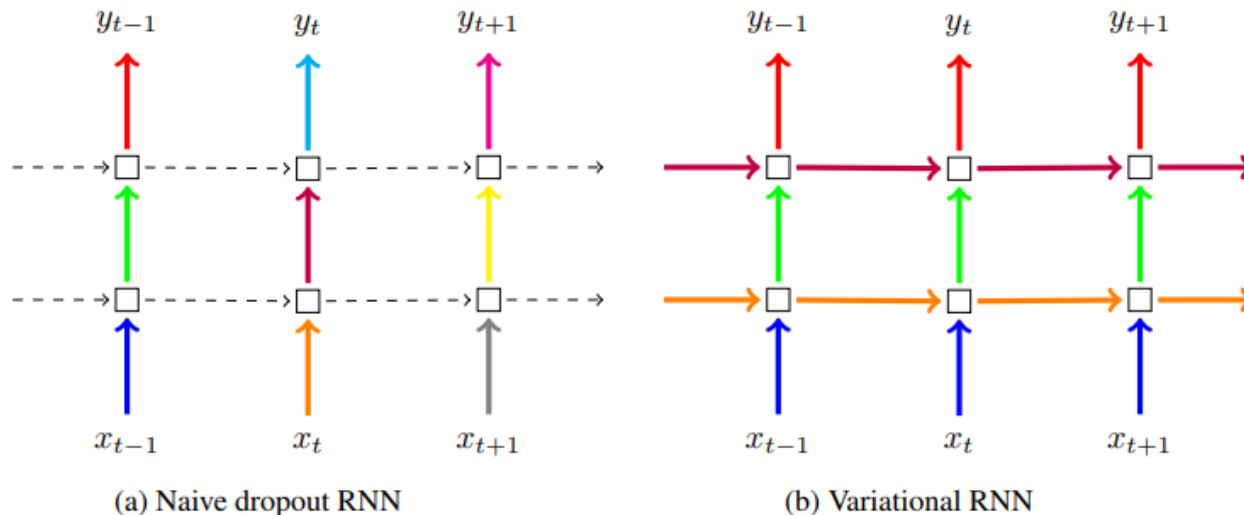
$$\mathbf{A} = f_{\mathbf{W}}(\mathbf{h}_{t-1}, \mathbf{x}_t) = \tanh([\mathbf{h}_{t-1}, \mathbf{x}_t] \cdot \mathbf{W} + \mathbf{b}) = \tanh(\mathbf{h}_{t-1} \cdot \mathbf{W}_h + \mathbf{x}_t \cdot \mathbf{W}_x + \mathbf{b})$$

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

Dimensions in example: \mathbf{W} :6x3, \mathbf{W}_h :3x3, \mathbf{W}_x :3x3

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.



same arrow
color indicates
identical dropout
mask

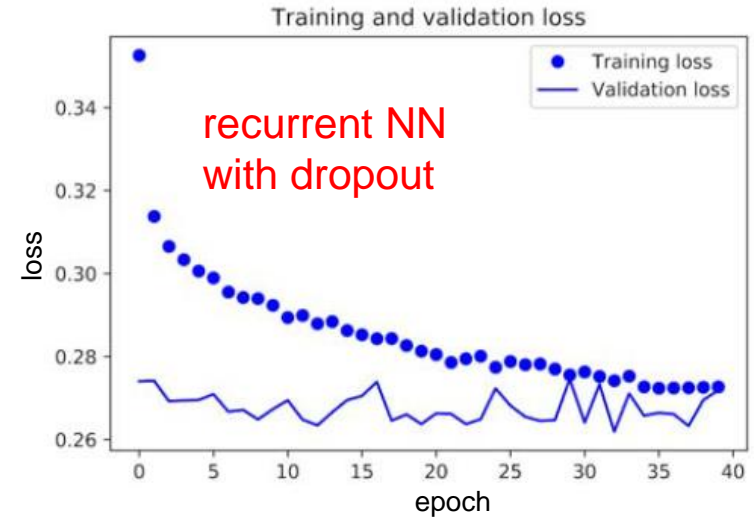
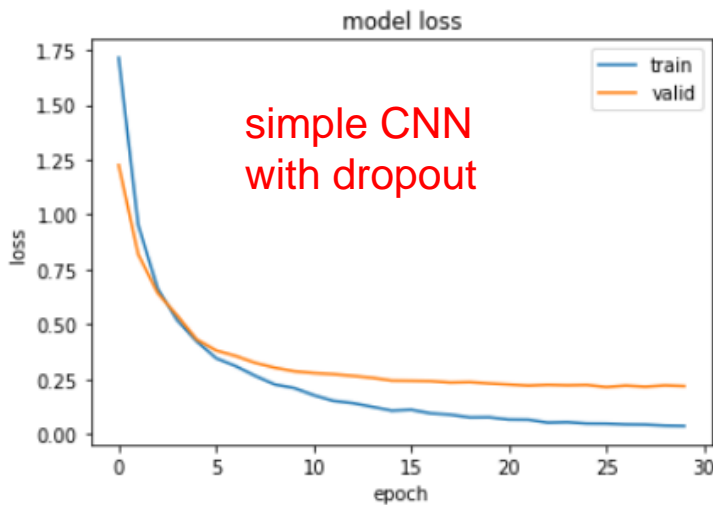
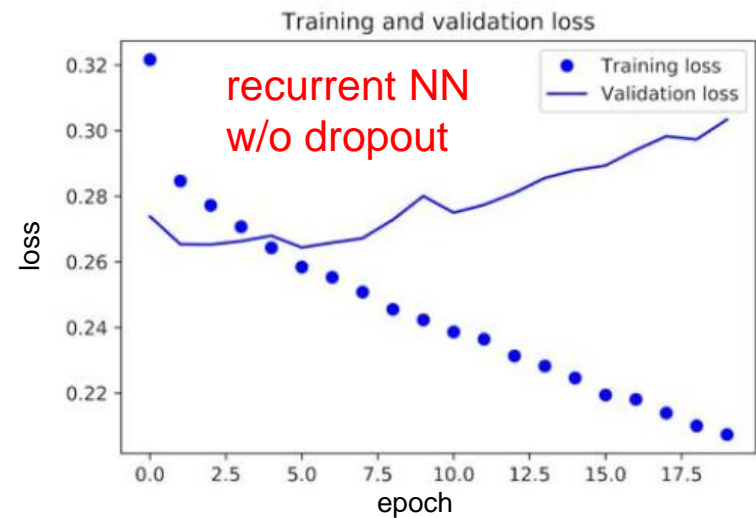
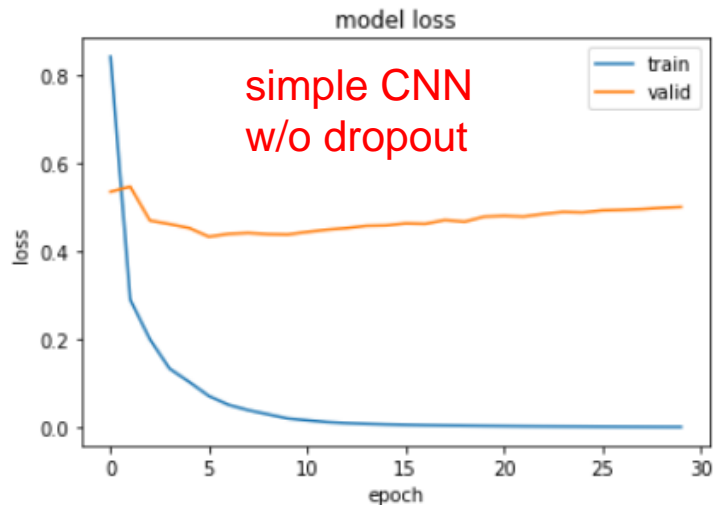
Figure 1: **Depiction of the dropout technique following our Bayesian interpretation (right) compared to the standard technique in the field (left).** Each square represents an RNN unit, with horizontal arrows representing time dependence (recurrent connections). Vertical arrows represent the input and 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:

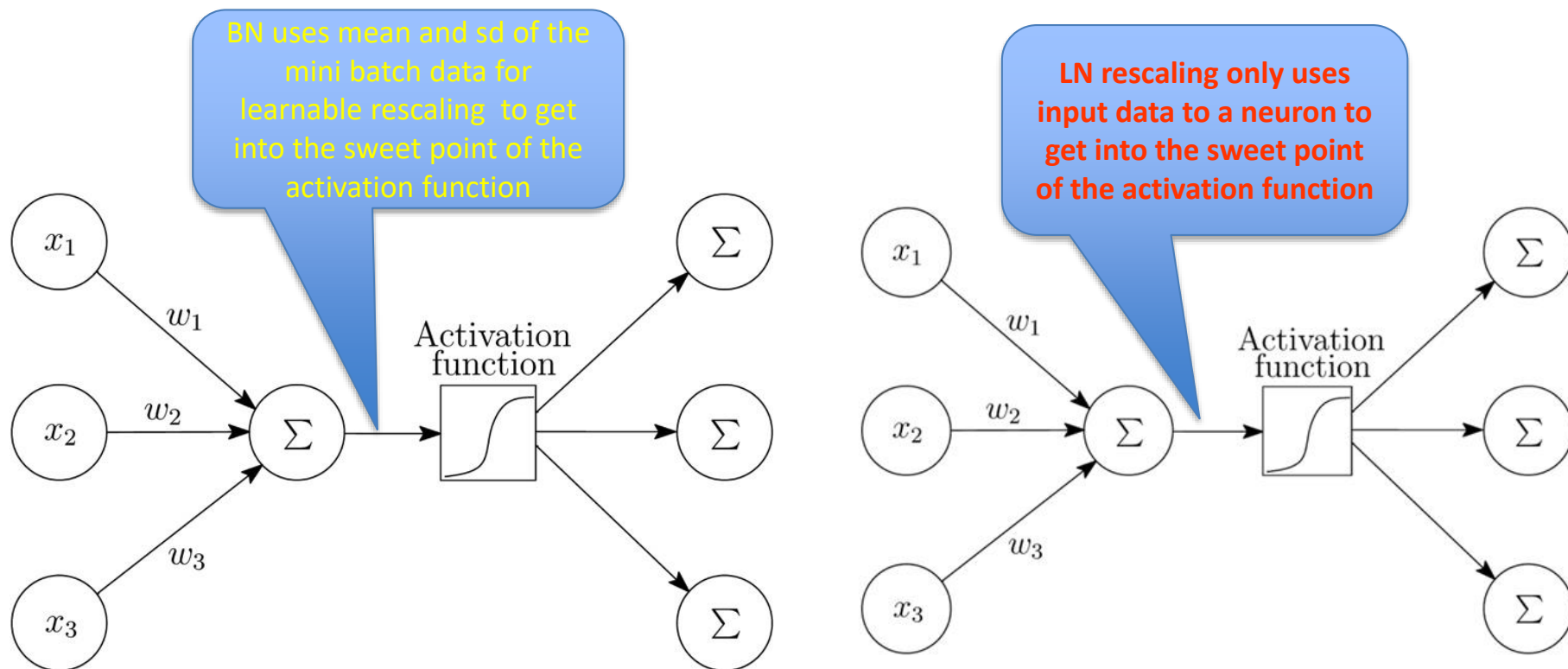
```
model.add(layers.GRU(32, dropout=0.2, recurrent_dropout=0.2, input_shape=(None, ...)))
```

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 .