

artificial

using

Machine

neurons

advanced

extract

sentences

structure

range

pre-defined

indispensable

places

paragraph

shallow

high

mappings

informal

trained

volunteering

things



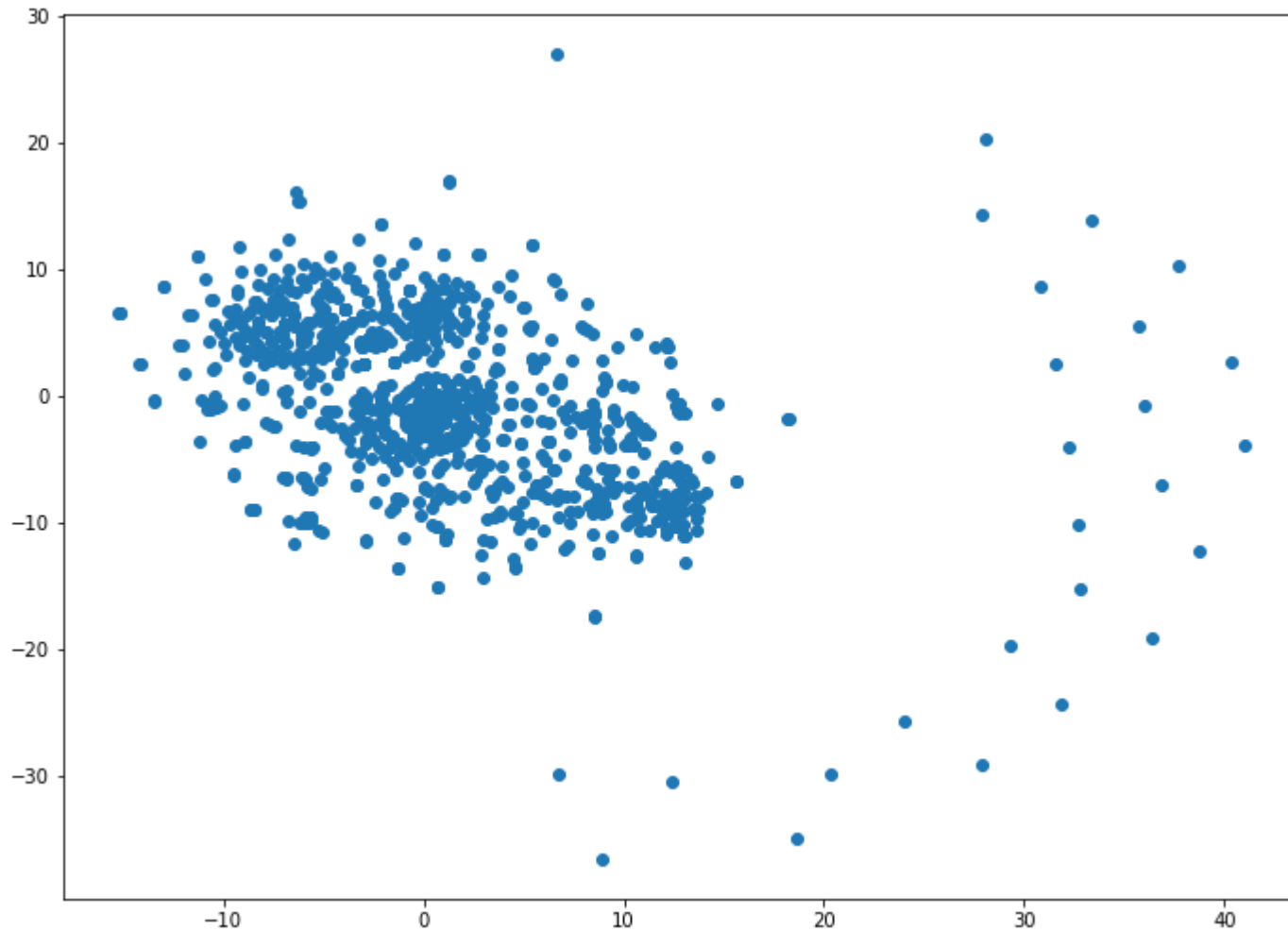
# Word2vec in Python

- Visualization like t-SNE can be used to inspect the pretrained word2vec resulting vectors

```
# Collect pretrained word2vec vectors of all words in the
# 20newsgroups vocabulary (use zeros for words we do not find)
word2vec_allvectors=numpy.zeros((len(remainingvocabulary),300))
for i in range(len(remainingvocabulary)):
    try:
        tempvector=word2vec_pretrainedvectors.wv[remainingvocabulary[i]]
        word2vec_allvectors[i,:]=tempvector
    except:
        continue
# Take a random subset of 1000 words
wordsubsetindices=numpy.random.permutation(len(remainingvocabulary))
# Create and plot a 2D t-SNE visualization
import sklearn.manifold
tsnemodel = sklearn.manifold.TSNE(n_components=2, verbose=1,
perplexity=20, n_iter=400)
tsneplot3 = tsnemodel.fit_transform(\
    word2vec_allvectors[wordsubsetindices[0:1000],:])
myfigure, myaxes = matplotlib.pyplot.subplots();
myaxes.scatter(tsneplot3[:,0],tsneplot3[:,1]);
```

# Word2vec in Python

- Result seems to show several subgroups of words with close-by meanings, plus outliers:



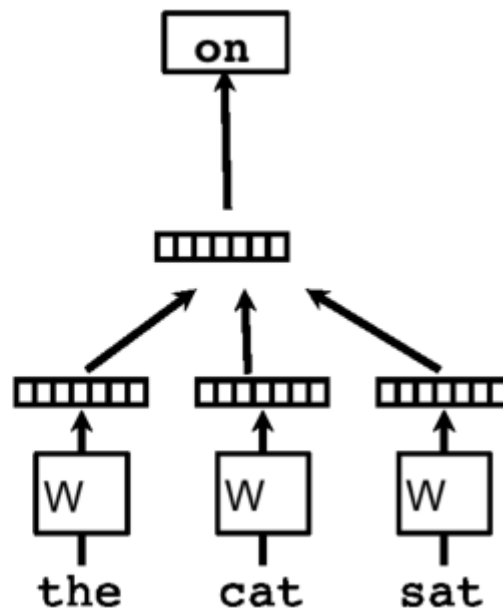
# Paragraph vector models

- Quoc Le and Tomas Mikolov. Distributed Representations of Sentences and Documents. International Conference on Machine Learning, 2014. ArXiv: <https://arxiv.org/pdf/1405.4053.pdf>
- Reminder: the continuous bag of words (CBOW) model predicted a central word based on context words near it.

Classifier

Average/Concatenate

Word Matrix



CBOW model (picture from Le & Mikolov 2014). Here the context is three previous words "the cat sat" and the central word is the next word (here "on").

$$p(w|C(w); \theta) = \frac{e^{\text{affinity}(w, C(w))}}{\sum_{v=1}^V e^{\text{affinity}(v, C(w))}}$$

$$\text{affinity}(w, C(w)) = \phi_w^T \left( \frac{1}{2R} \sum_{c \in C(w)} \psi_c \right) + \phi_{0,w}$$

Usually includes also a bias term

# Paragraph vector models

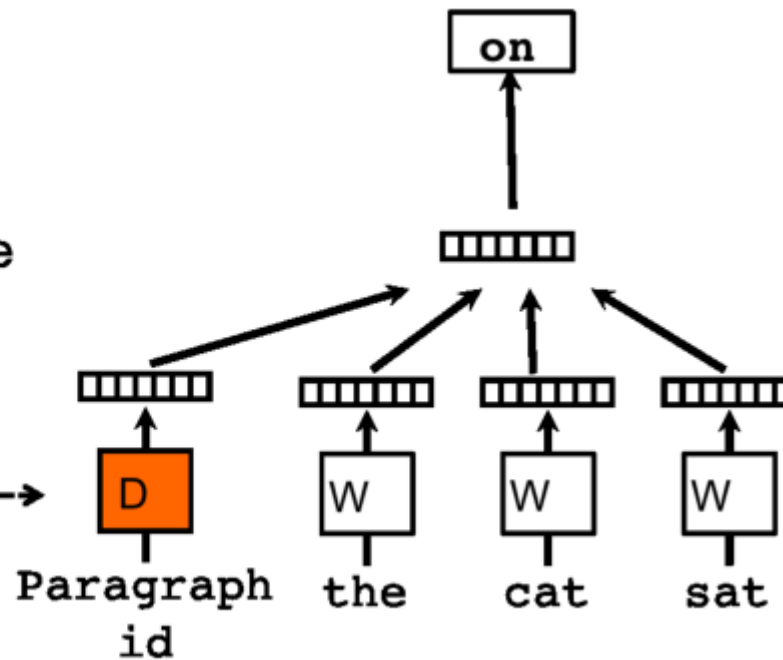
- **Paragraph vector model** (also called "Distributed Memory Model of Paragraph Vectors, PV-DM"):
  - in addition to the context words, the central word is predicted also based on a paragraph id number.
  - Context words are mapped to vectors like before
  - Every paragraph is mapped to a unique **paragraph vector**.
  - The paragraph vector is shared for all contexts in the same paragraph, but is different for different paragraphs
- The context-word vectors and the paragraph vector are either **averaged** or **concatenated** to form the context representation.
- In the concatenation version, the paragraph vectors can have different dimensionality from the context-word vectors

# Paragraph vector models

Classifier

Average/Concatenate

Paragraph Matrix----->



Paragraph vector model (picture from Le & Mikolov 2014). Here the context is "the cat sat" + the paragraph id and the word to be predicted is the next word (here "on")

## • Equations:

$$p(w|C(w); \theta) = \frac{e^{\text{affinity}(w, C(w))}}{\sum_{v=1}^V e^{\text{affinity}(v, C(w))}}$$

$$\text{affinity}(w, C(w)) = \phi_w^T \left( \frac{1}{2R+1} \left( \sum_{c \in C(w)} \psi_c + \gamma_{\text{Paragraph-id}(w)} \right) \right) + \phi_{0,w} \quad \text{Averaging version}$$

$$\text{affinity}(w, C(w)) = \phi_w^T [\psi_{c_1}, \dots, \psi_{c_{|C(w)|}}, \gamma_{\text{Paragraph-id}(w)}] + \phi_{0,w}$$

Concatenation version. Here  $\phi_w$  has a much higher dimensionality than in the averaging version!

**Parameters are found by maximizing log-likelihood**

# Paragraph vector models

- Parameters are found by maximizing the likelihood (probability of the central words given the contexts):

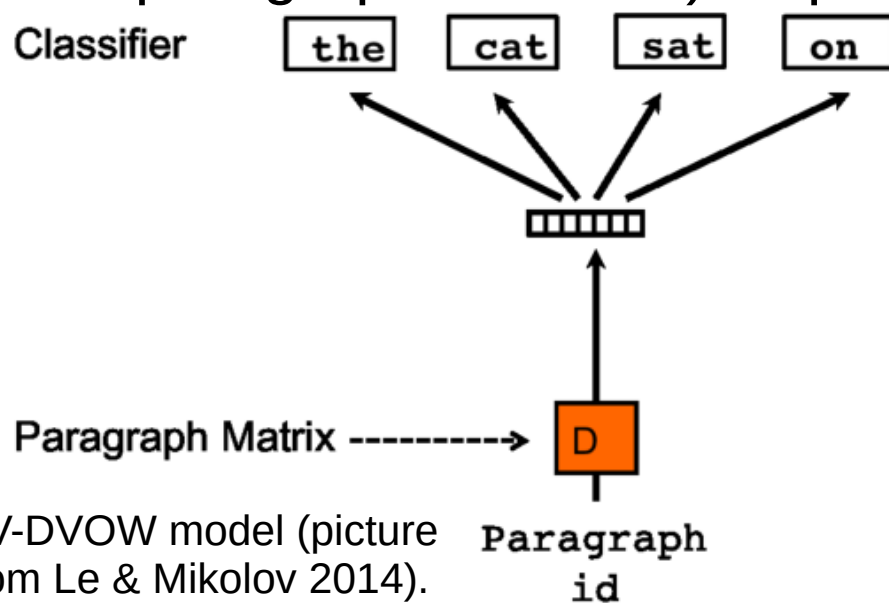
$$\prod_{s=1}^M \prod_{i=1}^{N^{(s)}} p(w_i^{(s)} | \mathbf{C}(w_i^{(s)}); \boldsymbol{\theta})$$

with respect to the word vectors, context-word vectors, and paragraph vectors.

- In a new prediction situation (predicting the next word in a new paragraph), the paragraph vector would not be available yet. It must be optimized first:
  - Word vectors and context-word vectors are kept at their previously optimized values.
  - Maximize the probability of the known words in the new paragraph (given their contexts) with respect to the paragraph vector.
  - Then use the paragraph vector to predict the new word
- In this model the context-word vectors  $\psi_c$  are used as the embedding of words, since the word vectors  $\phi_w$  are high-dimensional

# Paragraph vector models

- In the concatenation version, the order of concatenating the context-words affects the result. Taking word order into account can be good.
- Le and Mikolov (2014) suggest a "Distributed Bag of Words version of Paragraph Vector (PV-DBOW)" that does not consider word order
- It uses only the paragraph vector to predict randomly picked words from the paragraph. Similar to the skip-gram model, which used a central word to predict context words.
- For each paragraph, 1) sample a random text window, and a random context-word word from the window. 2) Predict the context-word given the paragraph vector. 3) Repeat several times.



$$\prod_{i \in \text{paragraphs}} \prod_{c \in \text{samples}_i} p(c|i; \theta) \quad \text{Likelihood}$$

$$p(c|i; \theta) = \frac{e^{\text{affinity}(c,i)}}{\sum_{u=1}^U e^{\text{affinity}(u,i)}} \quad \text{probability}$$

$$\text{affinity}(c,i) = \psi_c^T \mathbf{y}_i + \psi_{0,c}$$

affi-  
nity



# Paragraph vector models

- Even though they are called "paragraph vectors", the idea applies to any blocks of words: for example, sentences.
- Example application: sentiment classification of sentences (positive or negative). Data: Stanford Sentiment Treebank Dataset, 11855 sentences from Rotten Tomatoes movie reviews.
- In training data, use sentences as "paragraphs", and learn paragraph vectors  $\mathbf{y}_i$  for them.
- Learn a logistic regression classifier to predict classes of training sentences:

$$p(\text{positive}|i) = \frac{1}{1 + e^{-(\mathbf{w}_i^T \mathbf{y}_i + w_0)}}$$

Class probability

$$\sum_{i \in \text{positive}} \log p(\text{positive}|i) + \sum_{i \in \text{negative}} \log (1 - p(\text{positive}|i))$$

Log-likelihood:  
optimize with  
respect to weights  $\mathbf{w}$

- For test sentences, first optimize a paragraph vector for them, then feed it to the logistic regression classifier to predict class probability
- Le and Mikolov use a concatenation of paragraph vectors from the two models: PV-DBOW and PV-DM
- Result: 12.2% classification error rate.

# Paragraph vector models

- Example application 2: sentiment classification of documents (positive or negative). Data: IMDB movie reviews data set, 25000 train + 25000 labeled test reviews + 50000 unlabeled reviews. Probably the same as the one at <https://ai.stanford.edu/~amaas/data/sentiment/>
- In training data, use entire reviews as "paragraphs", and learn paragraph vectors  $\mathbf{y}_i$  for them.
- Learn a logistic regression classifier to predict classes of training reviews, like before.
- For test reviews, first optimize a paragraph vector for them, then feed it to the logistic regression classifier to predict class probability
- Result: 7.42% classification error rate.

# Paragraph vector models

- Example application 3: information retrieval.
- Data: triplets of search result snippets:  $(s1, s2, s3)$ . In each,  $s1$  and  $s2$  are snippets of webpages from the first page of search results for the same query;  $s3$  is a snippet from another, randomly chosen query.
- Learn paragraph vectors for each snippet, then compute Euclidean distances between  $s1$ ,  $s2$ , and  $s3$ .
- Count a retrieval success if  $s1$  was closer to  $s2$  than  $s3$ , error otherwise
- Result: 3.82% error rate
- Le and Mikolov (2014) compared this to:
  - bag-of-words (unclear what data representation and distance metric was used - e.g. Euclidean distance of TF-IDF vectors?). 8.10% error rate
  - bag-of-bigrams (unclear what data representation and distance metric was used - e.g. Euclidean distance of TF-IDF bigram vectors?). 7.28% error rate
  - bag-of-bigrams version that tried to learn a weighting matrix for keeping  $s1$  closer to  $s2$  than  $s3$ . 5.67% error rate
  - learning word vectors for words, and averaging the word vectors for each snippet. 10.25% error rate

# Paragraph vector models

- Python's gensim library includes learning of paragraph vector models (there called "doc2vec")
- They use a "hierarchical softmax" to compute log-likelihood over the output vocabulary, or a negative sampling model (unclear what exactly that means for the paragraph vector case)

```
import gensim
# We need to create a tagged version of each document
gensim_tagged_docs=[]
for k in range(len(mycrawled_prunedtexts)):
    doctag='doc'+str(k)
    tagged_document= \
        gensim.models.doc2vec.TaggedDocument( \
            mycrawled_prunedtexts[k], [doctag])
    gensim_tagged_docs.append(tagged_document)
# Create a dictionary from the documents
gensim_dictionary = gensim.corpora.Dictionary(gensim_docs)
```



# Paragraph vector models

- Python's gensim library includes learning of paragraph vector models (there called "doc2vec")
- They use a "hierarchical softmax" to compute log-likelihood over the output vocabulary, or a negative sampling model (unclear what exactly that means for the paragraph vector case)

```
# Train the word2vec model
# The dm_concat parameter controls whether to concatenate
# or average word vectors when learning the paragraph
# vector (see slides 5 and 6).
doc2vecmodel =
gensim.models.doc2vec.Doc2Vec(gensim_tagged_docs, \
    vector_size=10, window=5, min_count=1, \
    workers=4, dm_concat=0)
doc2vecmodel['doc986']
Out[296]:
array([-0.33598384, -0.29383156, -0.3040801 , -0.11860801, -0.15689434,
        0.12553768,  0.02579052, -0.05977593,  0.17080739, -0.02232333],
      dtype=float32)
```

# Neural language models

- Models like word2vec and paragraph vector models use continuous-valued vector representations internally, but combine them only linearly, except a final softmax nonlinearity to get output probabilities.
- Neural language models extend this to use nonlinearities in intermediate computations.
- Generally, neural language models:
  - use as input a string of words, each encoded as 1-of-V vectors (where the only 1 is at the vocabulary index of the word, other elements are zero).
  - use a softmax equation at the output to compute probabilities that sum to 1 over several possibilities.
  - optimize parameters to maximize the likelihood (probability of observations)

$$p(w|\text{inputs}; \theta) = \frac{e^{\phi_w^T f(\text{inputs}) + \phi_{w,0}}}{\sum_{v=1}^V e^{\phi_v^T f(\text{inputs}) + \phi_{v,0}}}$$

inputs, where  $\mathbf{f}$  is computed in different ways by different neural models

weights  $\phi_w, \phi_{w,0}$  are network parameters

# Neural language models

- In principle,  $\mathbf{f}$  could be computed in any way from the inputs
- The basic idea in neural networks is to compute  $\mathbf{f}$  as a composition of many functions:  $\mathbf{f} = \mathbf{f}_5(\mathbf{f}_4(\mathbf{f}_3(\mathbf{f}_2(\mathbf{f}_1(\text{inputs})\dots)))$ , so that each of the individual functions has a simple form and tunable parameters
- In a feedforward network, each  $\mathbf{f}_i$  is a layer of the network, and its  $K$ -dimensional output is computed by  $K$  neurons:

$$\mathbf{f}_i(\mathbf{x}) = [\text{neuron}_{i1}(\mathbf{x}), \dots, \text{neuron}_{iK}(\mathbf{x})]$$

- A typical neuron computes a weighted linear sum of inputs, and one of several typical nonlinear transformations of the sum:

$$\text{neuron}_{ij}(\mathbf{x}) = \text{nonlinearity}(\mathbf{w}_{ij}^T \mathbf{x} + w_{ij,0})$$

the weights  $w$  are parameters of the neuron

$$\text{nonlinearity}(y) = \frac{1}{1 + \exp(-y)}$$

logistic nonlinearity

$$\text{nonlinearity}(y) = \tanh(y)$$

hyperbolic tangent nonlinearity

$$\text{nonlinearity}(y) = \max(0, y)$$

"rectified linear unit" nonlinearity

# Recurrent networks

- In feedforward neural networks the input of each neuron are the set of outputs of the previous layers: the input for f5 is f4, the input for f4 is f3, the input for f3 is f2, the input for f2 is f1, and the input for f1 is the input text.
- However, this means that if e.g. the input is a window of 10 words, whatever happened earlier than those 10 words does not affect the network output predictions

- Recurrent networks use neurons that have a "memory": at time t, their own previous outputs at time t-1 are used as another input

$$neuron_{ij}(\mathbf{x}_t) = nonlinearity(\mathbf{w}_{ij}^T[\mathbf{x}, neuron_{ij}(\mathbf{x}_{t-1})] + w_{ij,0})$$

- More generally, the neuron can depend on previous outputs of the entire layer

$$neuron_{ij}(\mathbf{x}_t) = nonlinearity(\mathbf{w}_{ij}^T[\mathbf{x}, neuron_{i1}(\mathbf{x}_{t-1}), \dots, neuron_{iK}(\mathbf{x}_{t-1})] + w_{ij,0})$$

- Various other architectures allow connections to previous outputs of other layers too.



# LSTM (long short-term memory)

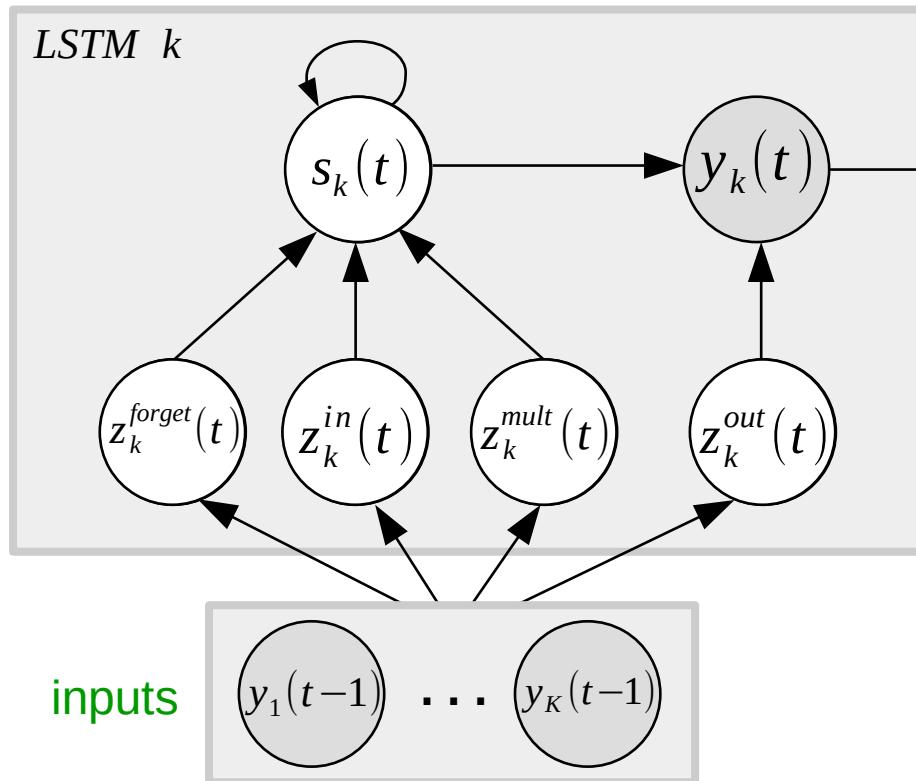
- References:

- S. Hochreiter and J. Schmidhuber. Long Short-Term Memory. *Neural Computation* 9 (8), 1735–1780, 1997.
- F. A. Gers. Learning to Forget: Continual Prediction with LSTM. Proc. ICANN 1999, pages 850–855, 1999.
- F. A. Gers, N. N. Schraudolph, and J. Schmidhuber. Learning Precise Timing with LSTM Recurrent Networks”, *Journal of Machine Learning Research* 3, 115–143, 2002.
- Martin Sundermeyer, Ralf Schlüter, and Hermann Ney. LSTM Neural Networks for Language Modeling. In Proc. Interspeech 2012, 2012.

- Idea: replace the usual recurrent neuron by a more complicated unit that involves several nonlinearities  $g$  and  $h$

# LSTM (long short-term memory)

- Mathematics:



$$z_k^{forget}(t) = g^{forget} \left( \sum_u w_k^{forget} y_u(t-1) \right)$$

$$z_k^{in}(t) = g^{in} \left( \sum_u w_k^{in} y_u(t-1) \right)$$

$$z_k^{mult}(t) = g^{mult} \left( \sum_u w_k^{mult} y_u(t-1) \right)$$

$$z_k^{out}(t) = g^{out} \left( \sum_u w_k^{out} y_u(t-1) \right)$$

$$s_k(t) = s_k(t-1) z_k^{forget}(t) + z_k^{in}(t) z_k^{mult}(t)$$

$$y_k(t) = s_k(t) h(z_k^{out}(t))$$