



AGH UNIVERSITY OF SCIENCE
AND TECHNOLOGY

Deep Learning in Speech Technology

Lectures

Jakub Galka, Department of Electronics ©
jgalka@agh.edu.pl

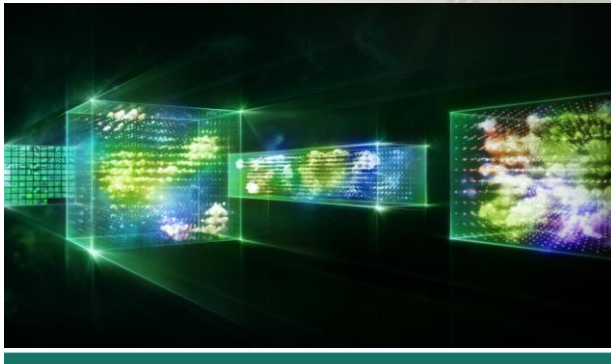


Outline

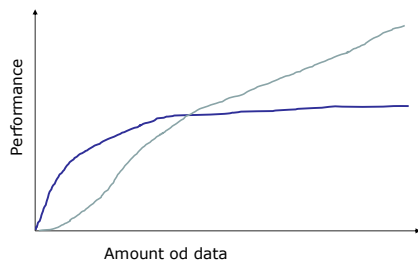
1. What is Deep Learning
2. Deep Learning tools
3. Deep Network Architectures and Applications
 - Feed-Forward Networks (FF)
 - Recurrent Neural Networks (RNN)
 - Long Short-Term Memory Networks (LSTM)
 - Convolutional Neural Networks (CNN)
 - Residual Neural Networks
 - Autoencoders
4. Learning architectures and approaches
 - Classification, discrimination, regression
 - Transfer Learning
 - Generative-Adversarial Networks (GAN)



Why going deep ?



Why Deep Learning ?



Why Now ? (2010+)

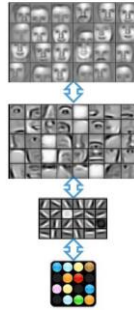
- More data
- More computational power
- More interest, more people
- Better algorithms
- Better results
- More applications



What is Deep Learning

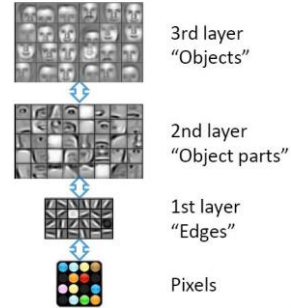
Wiki:

- Deep learning is a class of [machine learning algorithms](#)
- use a cascade of multiple layers of [nonlinear processing](#) units for [feature extraction](#) and transformation. Each successive layer uses the output from the previous layer as input.
- learn in [supervised](#) (e.g., classification) and/or [unsupervised](#) (e.g., pattern analysis) manners.
- learn multiple levels of representations that correspond to different levels of abstraction; tl levels form a hierarchy of concepts.



Layered abstraction

Feature representation



3rd layer
"Objects"

2nd layer
"Object parts"

1st layer
"Edges"

Pixels

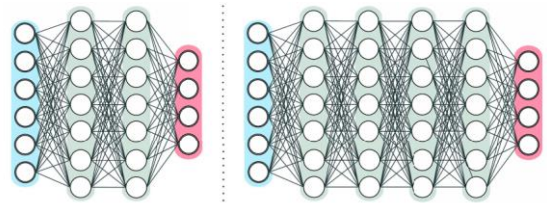


How To „Deep Learning“ ?



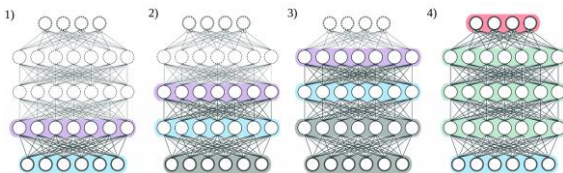
Normal vs Deep Feed-Forward Network (fully connected)

• Input hidden ... hidden output



Deep network training

- Network pretraining (1-3) and Fine-tuning (4)
- Unsupervised and supervised training (loss function)



- Regularization (L2, L1, Dropout)
- Data augmentation
- Meta-parameters optimization (Batch size, net size)
- Evaluation (Metric)

```

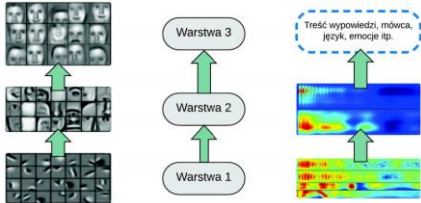
1 from keras.models import Sequential
2 from keras.layers import Dense
3 import numpy
4 # fix random seed for reproducibility
5 numpy.random.seed(7)
6
7 dataset = numpy.loadtxt("my_dataset.csv", delimiter=",")
8 # input (X) and output (Y)
9 X = dataset[:,0:39]
10 Y = dataset[:,39]
11
12 # create model
13 # 12, 8, 4, and 1 neuron in consecutive layers
14 model = Sequential()
15 # input layer
16 model.add(Dense(12, input_dim=39, activation='relu'))
17 model.add(layers.Dropout(0.2, noise_shape=None, seed=None))
18 model.add(Dense(8, activation='relu'))
19 model.add(layers.Dropout(0.1, noise_shape=None, seed=None))
20 model.add(Dense(4, activation='relu'))
21 model.add(Dense(1, activation='sigmoid'))
22 # Output layer have one neuron - '1' for speech, '0' for non-speech frame
23
24 # Compile model
25 model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
26
27 # Fit the model
28 model.fit(X, Y, epochs=100, batch_size=20)
29
30 # evaluate the model
31 scores = model.evaluate(X, Y)
32 print("\n%s: %.2f%%" % (model.metrics_names[1], scores[1]*100))
33
34 # predict output
35 predictions = model.predict(X)

```



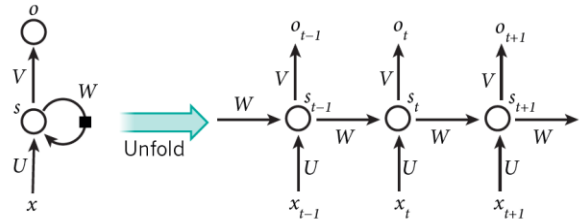
Typical Deep FFN architectures for speech processing

- Input: MFCC, Mel-filterbank, FFT
- Frame stacking at input (broader context)
- Outputs: softmax (prob), classification decision
- Feature extractor
- Problem: No direct support for time-wise/process learning



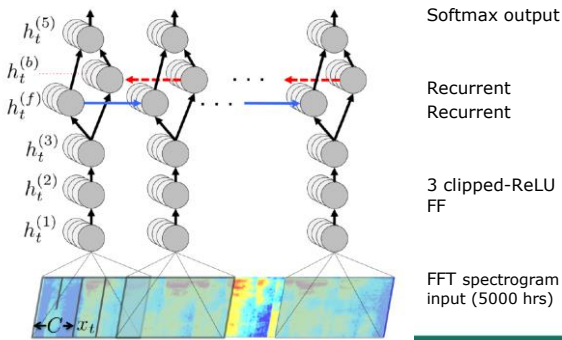
Recurrent Neural Networks

- Ability to remember information



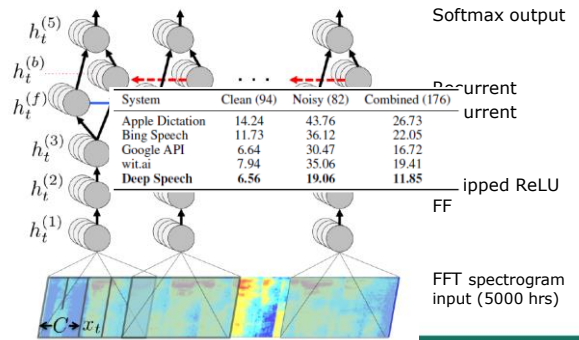
Recurrent Neural Networks

BAIDU (12.2014) Deep Speech: Scaling up end-to-end ASR



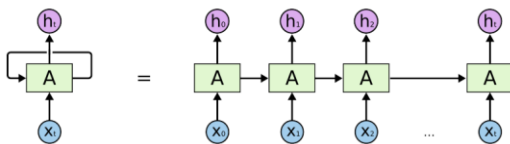
Recurrent Neural Networks

BAIDU (12.2014) Deep Speech: Scaling up end-to-end ASR

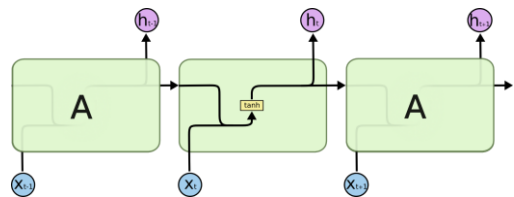


RNN vs LSTM

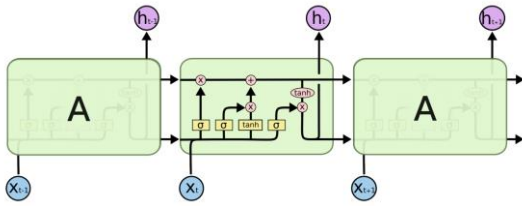
- Vanishing gradient problem
- RNN acts as a very deep FF network



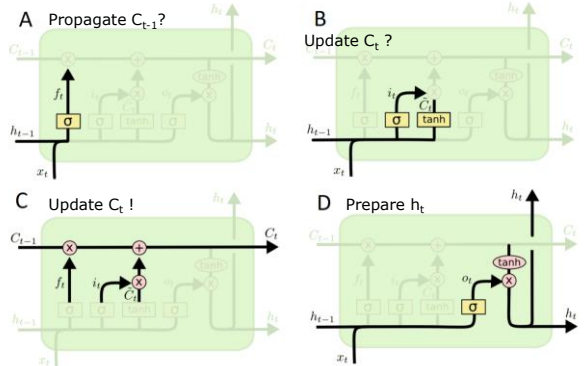
Standard RNN unwrapped



Long Short-Term Memory Networks (1997)

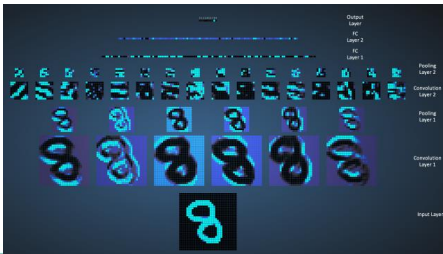


Steps of LSTM operation

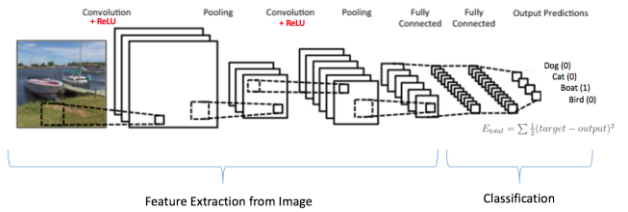


Convolutional Neural Networks

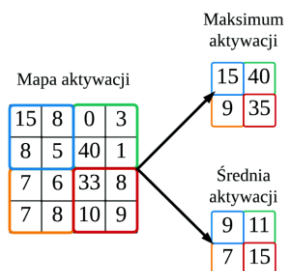
- In fully connected FF networks number of weights is too big for huge inputs.
- Same features can be observed in different places of the same signal (as in images)
- We can **group neurons** to analyse different parts of data -> lower number of weights



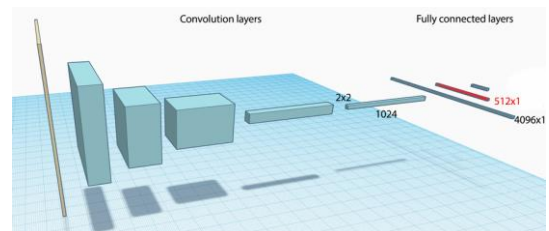
CNN workflow



Max-pooling, Average-pooling

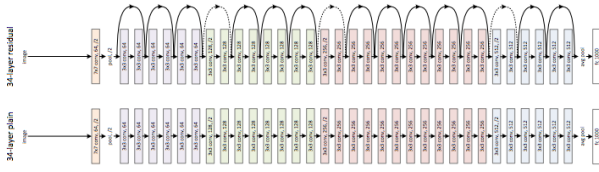


Exemple of typical CNN-FC model

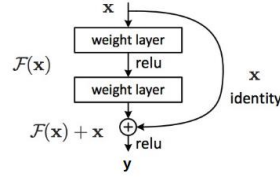




ResNet - Residual Networks



ResNet - Residual Networks



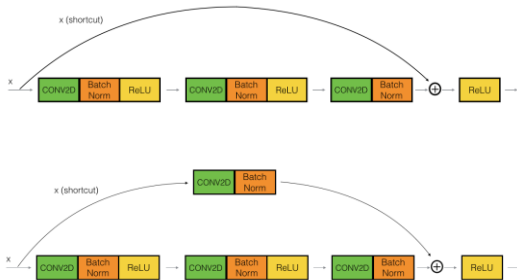
$$y = x + F(x)$$

$$\frac{\delta E}{\delta x} = \frac{\delta E}{\delta y} * \frac{\delta y}{\delta x} = \frac{\delta E}{\delta y} * (1 + F'(x))$$

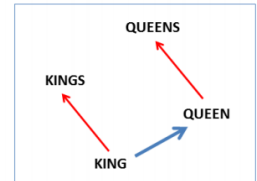
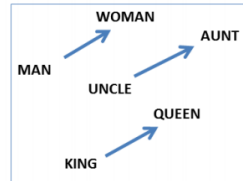
$$= \frac{\delta E}{\delta y} + \frac{\delta E}{\delta y} * F'(x)$$



X_shortcut = X # Store the initial value of X in a variable
 ## Here perform convolution + batch norm operations on X
 X = Add([X, X_shortcut]) # SKIP Connection

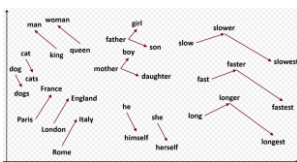


Word embeddings - latent space



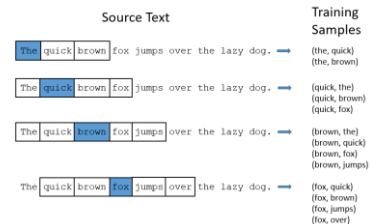
Word2Vec - word embeddings (Google, Mikolov)

- Converting text into high-dimensional vector of numbers (latent space embedding)
- Preserves linguistic and pragmatic information
- Embeddings are easy to manipulate and use in ML algorithms



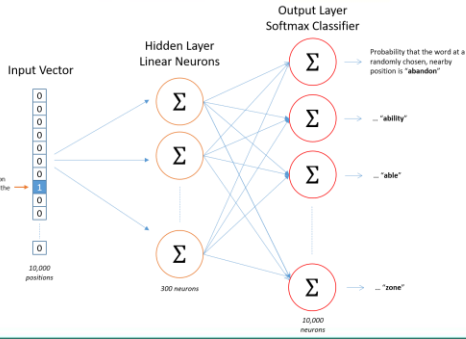
Word2vec algorithm

- Train your network using hot-word representation
- Use skip-gram method, or continuous bag of words



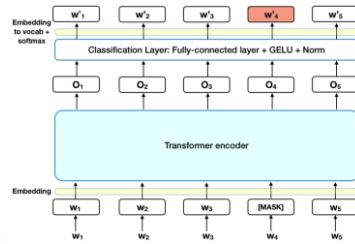
- Use bottleneck feature as embeddings

One-hot as NN input

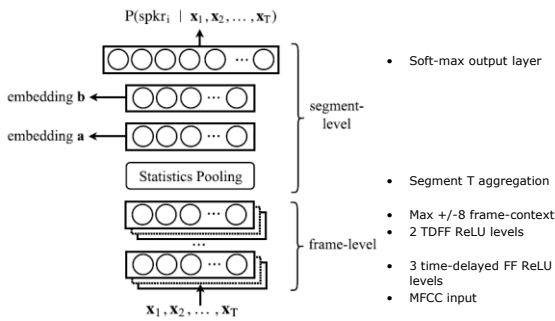


BERT Language Model (Google) Bidirectional Encoder Representations from Transformers

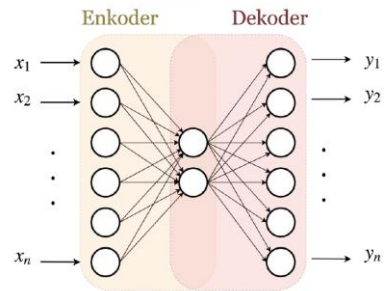
- Transformer is an attention mechanism that learns contextual relations between words



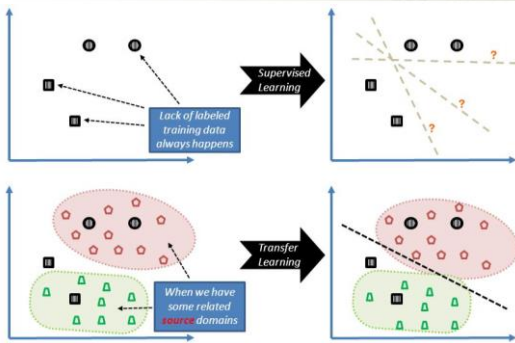
Deep Neural Network Embeddings for Text-Independent Speaker Verification



Autoencoders

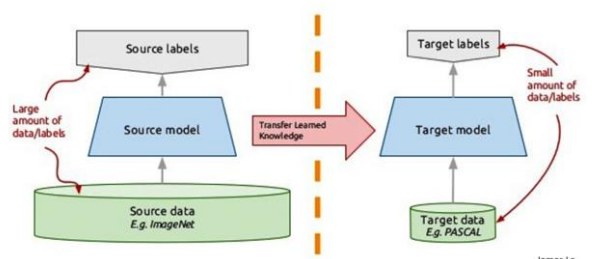


Transfer Learning

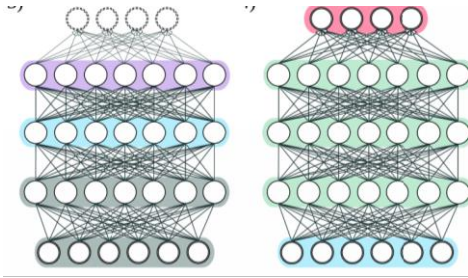


Transfer Learning

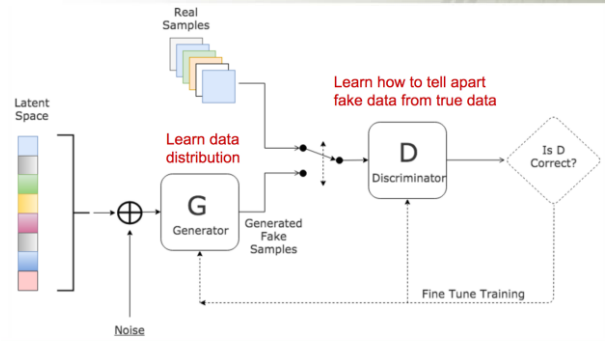
Transfer learning: idea



Transfer Learning possible training approaches

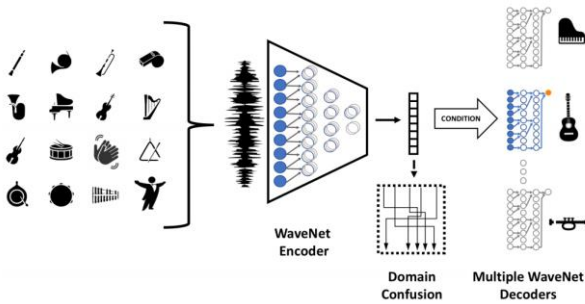


Generative Adversarial Networks



GAN Example from FB A Universal Music Translation Network

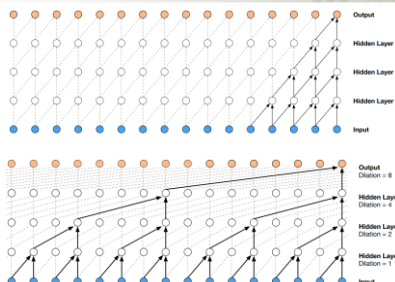
<https://www.youtube.com/watch?v=vdxCGNWTpUs&feature=youtu.be>



Speech Synthesis/Conversion using Deep Learning

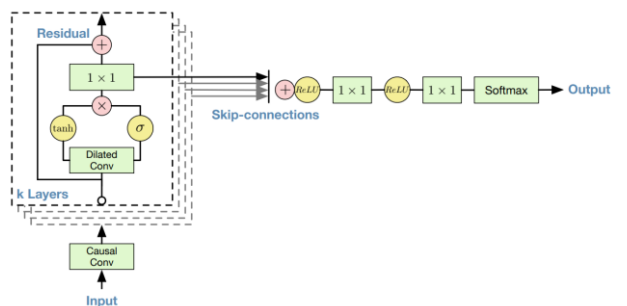
- WaveNet baseline
- Distilled Wavenet, ClariNet (Baidu)
- Tacotron2 (NATURAL TTS SYNTHESIS BY CONDITIONING WAVENET ON MEL SPECTROGRAM PREDICTIONS)
- Lot of implementations (NVIDIA), and downloadable models (eg. Wavenet generator)

Wavenet – sample by sample convolution



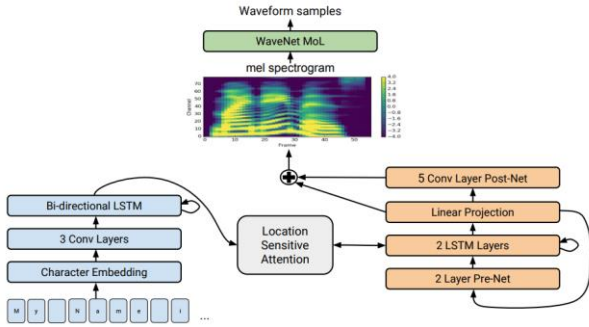
- „At training time, the conditional predictions for all timesteps can be made in parallel because all timesteps of ground truth x are known. When generating with the model, the predictions are sequential: after each sample is predicted, it is fed back into the network to predict the next sample.“

Wavenet training architecture



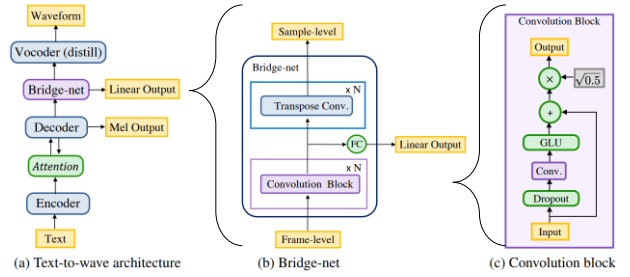
Tacotron 2 – Google TTS

<https://google.github.io/tacotron/publications/tacotron2/index.html>



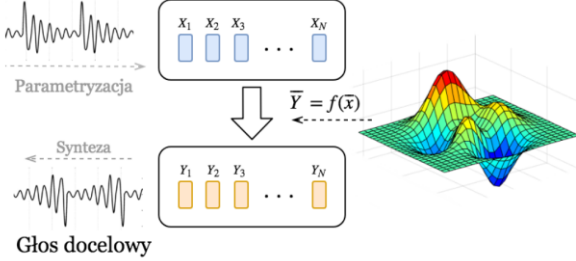
ClariNet: Parallel Wave Generation in End-to-End Text-to-Speech, Baidu

<https://clarinet-demo.github.io/>



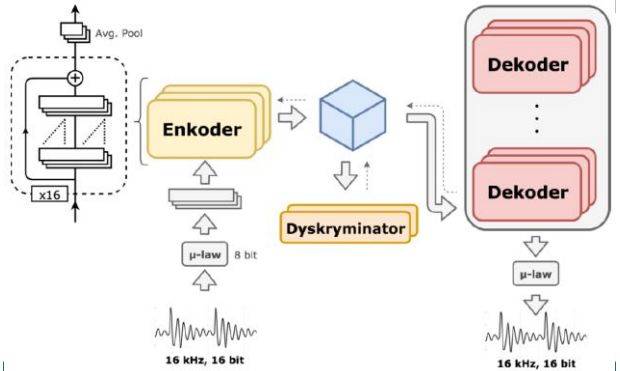
Voice Conversion

Głos źródłowy



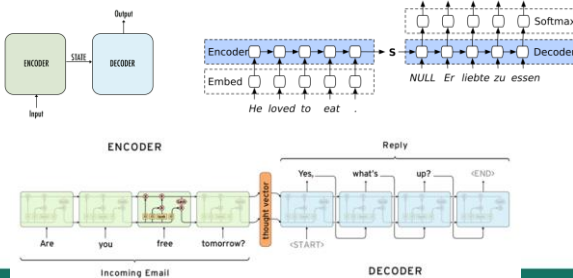
Voice Conversion with GAN

Hubert Siuzdak (IA, 2019)



Sequence-to-sequence modeling

- Used in machine translation
- Dialog modeling



LSTM Sequence-to-sequence translation example

