NPFL114, Lecture 9



Recurrent Neural Networks III

Milan Straka

🖬 April 29, 2019





EUROPEAN UNION European Structural and Investment Fund Operational Programme Research, Development and Education Charles University in Prague Faculty of Mathematics and Physics Institute of Formal and Applied Linguistics



unless otherwise stated

Recurrent Neural Networks



Single RNN cell



Unrolled RNN cells







Given an input $m{x}^{(t)}$ and previous state $m{s}^{(t-1)}$, the new state is computed as

$$oldsymbol{s}^{(t)} = f(oldsymbol{s}^{(t-1)},oldsymbol{x}^{(t)};oldsymbol{ heta}).$$

One of the simplest possibilities is

$$oldsymbol{s}^{(t)} = anh(oldsymbol{U}oldsymbol{s}^{(t-1)} + oldsymbol{V}oldsymbol{x}^{(t)} + oldsymbol{b}).$$

Attention

NMT

3/49

Basic RNN Cell

Basic RNN cells suffer a lot from vanishing/exploding gradients (*the challenge of long-term dependencies*).

If we simplify the recurrence of states to

$$oldsymbol{s}^{(t)} = oldsymbol{U}oldsymbol{s}^{(t-1)},$$

we get

$$oldsymbol{s}^{(t)} = oldsymbol{U}^t oldsymbol{s}^{(0)}.$$

If U has eigenvalue decomposition of $oldsymbol{U}=oldsymbol{Q}oldsymbol{\Lambda}oldsymbol{Q}^{-1}$, we get

$$oldsymbol{s}^{(t)} = oldsymbol{Q}oldsymbol{\Lambda}^toldsymbol{Q}^{-1}oldsymbol{s}^{(0)}.$$

The main problem is that the *same* function is iteratively applied many times.

Several more complex RNN cell variants have been proposed, which alleviate this issue to some degree, namely **LSTM** and **GRU**.

NPFL114, Lecture 9

Refresh CTC Word2vec

Subword Embeddings

Seq2seq Attention

Long Short-Term Memory

Later in Gers, Schmidhuber & Cummins (1999) a possibility to *forget* information from memory cell c_t was added.

$$\begin{split} \boldsymbol{i}_t &\leftarrow \sigma(\boldsymbol{W}^i \boldsymbol{x}_t + \boldsymbol{V}^i \boldsymbol{h}_{t-1} + \boldsymbol{b}^i) \\ \boldsymbol{f}_t &\leftarrow \sigma(\boldsymbol{W}^f \boldsymbol{x}_t + \boldsymbol{V}^f \boldsymbol{h}_{t-1} + \boldsymbol{b}^f) \\ \boldsymbol{o}_t &\leftarrow \sigma(\boldsymbol{W}^o \boldsymbol{x}_t + \boldsymbol{V}^o \boldsymbol{h}_{t-1} + \boldsymbol{b}^o) \\ \boldsymbol{c}_t &\leftarrow \boldsymbol{f}_t \cdot \boldsymbol{c}_{t-1} + \boldsymbol{i}_t \cdot \tanh(\boldsymbol{W}^y \boldsymbol{x}_t + \boldsymbol{V}^y \boldsymbol{h}_{t-1} + \boldsymbol{b}^y) \\ \boldsymbol{h}_t &\leftarrow \boldsymbol{o}_t \cdot \tanh(\boldsymbol{c}_t) \end{split}$$

NPFL114, Lecture 9

Refresh CTC

Word2vec Subword Embeddings

Seq2seq Attention



Long Short-Term Memory





NPFL114, Lecture 9

Refresh CTC

Word2vec Subword Embeddings

Seq2seq Attention

ntion NMT

Gated Recurrent Unit





$$oldsymbol{r}_t \leftarrow \sigma(oldsymbol{W}^roldsymbol{x}_t + oldsymbol{V}^roldsymbol{h}_{t-1} + oldsymbol{b}^r)$$

 $oldsymbol{u}_t \leftarrow \sigma(oldsymbol{W}^uoldsymbol{x}_t + oldsymbol{V}^uoldsymbol{h}_{t-1} + oldsymbol{b}^u)$
 $oldsymbol{\hat{h}}_t \leftarrow ext{tanh}(oldsymbol{W}^holdsymbol{x}_t + oldsymbol{V}^h(oldsymbol{r}_t \cdot oldsymbol{h}_{t-1}) + oldsymbol{b}^h)$
 $oldsymbol{h}_t \leftarrow oldsymbol{u}_t \cdot oldsymbol{h}_{t-1} + (1 - oldsymbol{u}_t) \cdot oldsymbol{\hat{h}}_t$
TC Word2vec Subword Embeddings Seq2seq Attention NMT

NPFL114, Lecture 9

Refresh СТС Subword Embeddings

Seq2seq Attention 7/49

Gated Recurrent Unit





$$z_t = \sigma \left(W_z \cdot [h_{t-1}, x_t] \right)$$
$$r_t = \sigma \left(W_r \cdot [h_{t-1}, x_t] \right)$$
$$\tilde{h}_t = \tanh \left(W \cdot [r_t * h_{t-1}, x_t] \right)$$
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-var-GRU.png

NMT

NPFL114, Lecture 9

Refresh CTC

Word2vec Subword Embeddings

Attention

Seq2seq

Word Embeddings



One-hot encoding considers all words to be independent of each other.

However, words are not independent – some are more similar than others.

Ideally, we would like some kind of similarity in the space of the word representations.

Distributed Representation

The idea behind distributed representation is that objects can be represented using a set of common underlying factors.

We therefore represent words as fixed-size *embeddings* into \mathbb{R}^d space, with the vector elements playing role of the common underlying factors.

Word2vec

Word Embeddings

The word embedding layer is in fact just a fully connected layer on top of one-hot encoding. However, it is important that this layer is *shared* across the whole network.



Refresh

Subword Embeddings

Seg2seg Attention



Word Embeddings for Unknown Words



Recurrent Character-level WEs



Figure 1 of paper "Finding Function in Form: Compositional Character Models for Open Vocabulary Word Representation", https://arxiv.org/abs/1508.02096.

NPFL114, Lecture 9 Refresh CTC Word2vec Subword Embeddings Seq2seq Attention NMT

Word Embeddings for Unknown Words



Convolutional Character-level WEs



NPFL114, Lecture 9

Refresh

CTC

Word2vec Subword Embeddings S

s Seq2seq Attention

NMT

12/49

Basic RNN Applications



Sequence Element Classification

Use outputs for individual elements.



Sequence Representation

Use state after processing the whole sequence (alternatively, take output of the last element).

Structured Prediction



Consider generating a sequence of $y_1,\ldots,y_N\in Y^N$ given input $oldsymbol{x}_1,\ldots,oldsymbol{x}_N$.

Predicting each sequence element independently models the distribution $P(y_i | \boldsymbol{X})$.

However, there may be dependencies among the y_i themselves, which is difficult to capture by independent element classification.

Word2vec

NMT

Seq2seq

Linear-Chain Conditional Random Fields (CRF)



Linear-chain Conditional Random Fields, usually abbreviated only to CRF, acts as an output layer. It can be considered an extension of a softmax – instead of a sequence of independent softmaxes, CRF is a sentence-level softmax, with additional weights for neighboring sequence elements.

$$egin{aligned} &s(oldsymbol{X},oldsymbol{y};oldsymbol{ heta},oldsymbol{A}) = \sum_{i=1}^N ig(oldsymbol{A}_{y_{i-1},y_i} + f_{oldsymbol{ heta}}(y_i|oldsymbol{X})ig) \ &p(oldsymbol{y}|oldsymbol{X}) = ext{softmax}_{oldsymbol{z}\in Y^N}ig(s(oldsymbol{X},oldsymbol{z})ig)_{oldsymbol{z}} \ & ext{og}\,p(oldsymbol{y}|oldsymbol{X}) = s(oldsymbol{X},oldsymbol{y}) - ext{logadd}_{oldsymbol{z}\in Y^N}(s(oldsymbol{X},oldsymbol{z})) \end{aligned}$$

NPFL114, Lecture 9

Refresh CTC

Subword Embeddings

Word2vec

Seq2seq

Computation

We can compute p(y|X) efficiently using dynamic programming. If we denote $\alpha_t(k)$ as probability of all sentences with t elements with the last y being k.

The core idea is the following:

$$lpha_t(k) = f_{oldsymbol{ heta}}(y_t = k | oldsymbol{X}) + ext{logadd}_{j \in Y}(lpha_{t-1}(j) + oldsymbol{A}_{j,k}).$$

For efficient implementation, we use the fact that

$$\ln(a+b) = \ln a + \ln(1+e^{\ln b - \ln a}).$$

NPFL114, Lecture 9

Word2vec



Decoding

We can perform optimal decoding, by using the same algorithm, only replacing logadd with max and tracking where the maximum was attained.

Applications

CRF output layers are useful for span labeling tasks, like

- named entity recognition
- dialog slot filling

Word2vec

NMT

Seq2seq



Connectionist Temporal Classification

Ú F_ÁL

Let us again consider generating a sequence of y_1, \ldots, y_M given input x_1, \ldots, x_N , but this time $M \leq N$ and there is no explicit alignment of x and y in the gold data.



Word2vec

Attention NMT

Sea2sea

Connectionist Temporal Classification

Ú F_AL

We enlarge the set of output labels by a - (blank) and perform a classification for every input element to produce an *extended labeling*. We then post-process it by the following rules (denoted \mathcal{B}):

- 1. We remove neighboring symbols.
- 2. We remove the -.

Because the explicit alignment of inputs and labels is not known, we consider *all possible* alignments.

Denoting the probability of label l at time t as p_l^t , we define

Word2vec

$$lpha^t(s) \stackrel{\scriptscriptstyle{ ext{def}}}{=} \sum_{ ext{labeling } oldsymbol{\pi}: \mathcal{B}(oldsymbol{\pi}_{1:t}) = oldsymbol{y}_{1:s}} \prod_{t'=1}^t p_{oldsymbol{\pi}_{t'}}^{t'}.$$

NPFL114, Lecture 9

CRF and CTC Comparison



In CRF, we normalize the whole sentences, therefore we need to compute unnormalized probabilities for all the (exponentially many) sentences. Decoding can be performed optimally.

In CTC, we normalize per each label. However, because we do not have explicit alignment, we compute probability of a labeling by summing probabilities of (generally exponentially many) extended labelings.

Word2vec Subword Embeddings

Seq2seq A

Computation

When aligning an extended labeling to a regular one, we need to consider whether the extended labeling ends by a *blank* or not. We therefore define

$$egin{aligned} lpha_{-}^t(s) & \stackrel{ ext{def}}{=} \sum_{ ext{labeling } m{\pi}: \mathcal{B}(m{\pi}_{1:t}) = m{y}_{1:s}, \pi_t = - \prod_{t'=1}^t p_{m{\pi}_{t'}}^{t'} \ lpha_*^t(s) & \stackrel{ ext{def}}{=} \sum_{ ext{labeling } m{\pi}: \mathcal{B}(m{\pi}_{1:t}) = m{y}_{1:s}, \pi_t
eq - \prod_{t'=1}^t p_{m{\pi}_{t'}}^{t'} \end{aligned}$$

and compute $lpha^t(s)$ as $lpha^t_-(s)+lpha^t_*(s).$

NPFL114, Lecture 9

Word2vec

Connectionist Temporal Classification

Computation

We initialize α s as follows:

- $lpha_{-}^{1}(0) \leftarrow p_{-}^{1}$
- $lpha^1_*(1) \leftarrow p^1_{y_1}$

We then proceed recurrently according to:

 $\begin{array}{l} \bullet \ \ \alpha^t_-(s) \leftarrow p^t_-(\alpha^{t-1}_-(s) + \alpha^{t-1}_*(s)) \\ \bullet \ \ \alpha^t_*(s) \leftarrow \begin{cases} p^t_{y_s}(\alpha^{t-1}_*(s) + \alpha^{t-1}_*(s-1) + a^{t-1}_-(s-1)), \text{ if } y_s \neq y_{s-1} \\ p^t_{y_s}(\alpha^{t-1}_*(s) + a^{t-1}_-(s-1)), \text{ if } y_s = y_{s-1} \end{cases} \end{array}$



Figure 7.3 of the dissertation "Supervised Sequence Labelling with Recurrent Neural Networks" by Alex Graves.



CTC Decoding

Unlike CRF, we cannot perform the decoding optimally. The key observation is that while an optimal extended labeling can be extended into an optimal labeling of a larger length, the same does not apply to regular (non-extended) labeling. The problem is that regular labeling coresponds to many extended labelings, which are modified each in a different way during an extension of the regular labeling.



CTC Decoding



Beam Search

To perform beam search, we keep k best regular labelings for each prefix of the extended labelings. For each regular labeling we keep both α_{-} and a_{*} and by *best* we mean such regular labelings with maximum $\alpha_{-} + \alpha_{*}$.

To compute best regular labelings for longer prefix of extended labelings, for each regular labeling in the beam we consider the following cases:

• adding a *blank* symbol, i.e., updating both α_{-} and α_{*} ;

Word2vec

• adding any non-blank symbol, i.e., updating α_* .

Finally, we merge the resulting candidates according to their regular labeling and keep only the k best.

Unsupervised Word Embeddings

Ú F_ÁL

The embeddings can be trained for each task separately.

However, a method of precomputing word embeddings have been proposed, based on *distributional hypothesis*:

Words that are used in the same contexts tend to have similar meanings.

The distributional hypothesis is usually attributed to Firth (1957).

Refresh

Seq2seq

Word2Vec



Word2vec

Subword Embeddings

eddings Seq2seq

Word2Vec



Table 8: Examples of the word pair relationships, using the best word vectors from Table 4 (Skipgram model trained on 783M words with 300 dimensionality).

Relationship Example 1		Example 2	Example 3	
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee	
big - bigger	small: larger	cold: colder	quick: quicker	
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii	
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter	
Sarkozy - France	arkozy - France Berlusconi: Italy		Koizumi: Japan	
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium	
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack	
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone	
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs	
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza	

Table 8 of paper "Efficient Estimation of Word Representations in Vector Space", https://arxiv.org/abs/1301.3781.

NMT

Word2vec

Seq2seq

Word2Vec – SkipGram Model





Considering input word w_i and output w_o , the Skip-gram model defines

$$p(w_o|w_i) \stackrel{ ext{\tiny def}}{=} rac{e^{oldsymbol{W}_{w_o}^ op oldsymbol{V}_{w_i}}}{\sum_w e^{oldsymbol{W}_w^ op oldsymbol{V}_{w_i}}}.$$

NPFL114, Lecture 9

Refresh CTC Word2vec

Subword Embeddings

s Seq2seq

Word2Vec – Hierarchical Softmax



Instead of a large softmax, we construct a binary tree over the words, with a sigmoid classifier for each node.

If word w corresponds to a path n_1, n_2, \ldots, n_L , we define

$$p_{ ext{HS}}(w|w_i) \stackrel{\scriptscriptstyle ext{def}}{=} \prod_{j=1}^{L-1} \sigma([+1 ext{ if } n_{j+1} ext{ is right child else -1}] \cdot oldsymbol{W}_{n_j}^ op oldsymbol{V}_{w_i}).$$

Seq2seq

Word2Vec – Negative Sampling

Instead of a large softmax, we could train individual sigmoids for all words. We could also only sample the *negative examples* instead of training all of them. This gives rise to the following *negative sampling* objective:

$$l_{ ext{NEG}}(w_o, w_i) \stackrel{ ext{\tiny def}}{=} \log \sigma(oldsymbol{W}_{w_o}^ op oldsymbol{V}_{w_i}) + \sum_{j=1}^k \mathbb{E}_{w_j \sim P(w)} \logig(1 - \sigma(oldsymbol{W}_{w_j}^ op oldsymbol{V}_{w_i})ig).$$

For P(w), both uniform and unigram distribution U(w) work, but

Word2vec

 $U(w)^{3/4}$

outperforms them significantly (this fact has been reported in several papers by different authors).

NPFL114, Lecture 9

Refresh CTC

Subword Embeddings

Seq2seq Attention



Recurrent Character-level WEs



increased	John	Noahshire	phding
reduced	Richard	Nottinghamshire	mixing
improved	George	Bucharest	modelling
expected	James	Saxony	styling
decreased	Robert	Johannesburg	blaming
targeted	Edward	Gloucestershire	christening

Table 2: Most-similar in-vocabular words under the C2W model; the two query words on the left are in the training vocabulary, those on the right are nonce (invented) words.

Figure 1 of paper "Finding Function in Form: Compositional Character Models for Open Vocabulary Word Representation", https://arxiv.org/abs/1508.02096.

NPFL114, Lecture 9

Refresh

Convolutional Character-level WEs



	In Vocabulary				Out-of-Vocabulary			
	while	his	уои	richard	trading	computer-aided	misinformed	loooook
	although	your	conservatives	jonathan	advertised		_	_
LSTM-Word	letting	her	we	robert	advertising	_	_	_
	though	my	guys	neil	turnover	_	_	_
	minute	their	i	nancy	turnover	_	—	_
	chile	this	your	hard	heading	computer-guided	informed	look
LSTM-Char	whole	hhs	young	rich	training	computerized	performed	cook
(before highway)	meanwhile	is	four	richer	reading	disk-drive	transformed	looks
	white	has	youth	richter	leading	computer	inform	shook
	meanwhile	hhs	we	eduard	trade	computer-guided	informed	look
LSTM-Char	whole	this	your	gerard	training	computer-driven	performed	looks
(after highway)	though	their	doug	edward	traded	computerized	outperformed	looked
	nevertheless	your	i	carl	trader	computer	transformed	looking

Table 6: Nearest neighbor words (based on cosine similarity) of word representations from the large word-level and character-level (before and after highway layers) models trained on the PTB. Last three words are OOV words, and therefore they do not have representations in the word-level model.

Table 6 of paper "Character-Aware Neural Language Models", https://arxiv.org/abs/1508.06615.

NPFL114, Lecture 9

Refresh CTC

Word2vec

Subword Embeddings

Seq2seq Atte

Character N-grams

Another simple idea appeared simultaneously in three nearly simultaneous publications as <u>Charagram</u>, <u>Subword Information</u> or <u>SubGram</u>.

A word embedding is a sum of the word embedding plus embeddings of its character *n*-grams. Such embedding can be pretrained using same algorithms as word2vec.

The implementation can be

- dictionary based: only some number of frequent character *n*-grams is kept;
- hash-based: character *n*-grams are hashed into K buckets (usually $K \sim 10^6$ is used).

Word2vec

Seq2seq



query	tiling	tech-rich	english-born	micromanaging	eateries	dendritic
sisg	tile	tech-dominated	british-born	micromanage	restaurants	dendrite
	flooring	tech-heavy	polish-born	micromanaged	eaterie	dendrites
sg	bookcases	technology-heavy	most-capped	defang	restaurants	epithelial
	built-ins	.ixic	ex-scotland	internalise	delis	p53

Table 7: Nearest neighbors of rare words using our representations and skipgram. These hand picked examples are for illustration.

Table 7 of paper "Enriching Word Vectors with Subword Information", https://arxiv.org/abs/1607.04606.

NPFL114, Lecture 9

Refresh CTC

Word2vec

Subword Embeddings

Seq2seq Attention

Charagram WEs





Figure 2: Illustration of the similarity between character *n*-grams in out-of-vocabulary words. For each pair, only one word is OOV, and is shown on the *x* axis. Red indicates positive cosine, while blue negative.
 Figure 2 of paper "Enriching Word Vectors with Subword Information", https://arxiv.org/abs/1607.04606.



re 9 Refresh CTC Word2vec Subword Embeddings Seq2seq Attention NMT

35/49



NPFL114, Lecture 9

Refresh CTC

Word2vec Subword Embeddings

Seq2seq

Attention NMT





NPFL114, Lecture 9

Refresh CTC

Word2vec Subword Embeddings

Seq2seq Atte

Attention NMT



Decoder



Figure 1 of paper "Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation", https://arxiv.org/abs/1406.1078.



Refresh

Subword Embeddings

Seq2seq Attention



Training

The so-called *teacher forcing* is used during training – the gold outputs are used as inputs during training.

Inference

During inference, the network processes its own predictions.

Usually, the generated logits are processed by an $\arg \max$, the chosen word embedded and used as next input.



Word2vec

Seq2seq Attention

Tying Word Embeddings





NPFL114, Lecture 9

Refresh CTC

Word2vec

Subword Embeddings

Seq2seq Attention

on NMT

Attention



As another input during decoding, we add *context vector* c_i :

$$oldsymbol{s}_i = f(oldsymbol{s}_{i-1},oldsymbol{y}_{i-1},oldsymbol{c}_i).$$

We compute the context vector as a weighted combination of source sentence encoded outputs:

$$oldsymbol{c}_i = \sum_j lpha_{ij}oldsymbol{h}_j$$

The weights $lpha_{ij}$ are softmax of e_{ij} over j,

 $oldsymbol{lpha}_i = ext{softmax}(oldsymbol{e}_i),$

with e_{ij} being

$$e_{ij} = oldsymbol{v}^ op anh(oldsymbol{V}oldsymbol{h}_j + oldsymbol{W}oldsymbol{s}_{i-1} + oldsymbol{b}).$$

Word2vec



NMT



to Align and Translate", https://arxiv.org/abs/1409.0473.

Attention



Figure 3 of paper "Neural Machine Translation by Jointly Learning to Align and Translate", https://arxiv.org/abs/1409.0473.

NPFL114, Lecture 9

Refresh

CTC

Word2vec

Subword Embeddings

dings Seq2seq

eq Attention

NMT

42/49

Subword Units



Translate *subword units* instead of words. The subword units can be generated in several ways, the most commonly used are

BPE – Using the *byte pair encoding* algorithm. Start with characters plus a special end-of-word symbol ·. Then, merge the most occurring symbol pair A, B by a new symbol AB, with the symbol pair never crossing word boundary.

Considering a dictionary with words *low, lowest, newer, wider*:

Word2vec

 $egin{array}{ccc} r & \cdot
ightarrow r \cdot \ l & o
ightarrow lo \ lo & w
ightarrow low \ e & r \cdot
ightarrow er \cdot \end{array}$

• Wordpieces – Joining neighboring symbols to maximize unigram language model likelihood.

Usually quite little subword units are used (32k-64k), often generated on the union of the two vocabularies (the so-called *joint BPE* or *shared wordpieces*).

NPFL114, Lecture 9

Refresh CTC

Subword Embeddings

Seq2seq Attention

Google NMT





NPFL114, Lecture 9

Refresh

CTC

Word2vec

Seq2seq Attention

Google NMT



Figure 5 of paper "Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation", https://arxiv.org/abs/1609.08144.

NPFL114, Lecture 9

CTC Word2vec

Refresh

Subword Embeddings

Seq2seq Attention

Google NMT





Beyond one Language Pair

A person riding a motorcycle on a dirt road.



A group of young people playing a game of frisbee.



A herd of elephants walking across a dry grass field.







Two hockey players are fighting over the puck.



A close up of a cat laying on a couch.



A skateboarder does a trick



A little girl in a pink hat is



A red motorcycle parked on the



A dog is jumping to catch a



A refrigerator filled with lots of food and drinks.



A yellow school bus parked





Describes with minor errors

Somewhat related to the image

Unrelated to the image

Fig. 5. A selection of evaluation results, grouped by human rating.

Word2vec

Figure 5 of "Show and Tell: Lessons learned from the 2015 MSCOCO...", https://arxiv.org/abs/1609.06647.



Refresh CTC

Subword Embeddings

Seq2seq



Beyond one Language Pair







What vegetable is the dog chewing on? MCB: carrot GT: carrot



What kind of dog is this? MCB: husky GT: husky



What kind of flooring does the room have? MCB: carpet GT: carpet





What color is the traffic light? MCB: green GT: green

Word2vec



Is this an urban area? MCB: yes GT: yes



Figure 6 of "Multimodal Compact Bilinear Pooling for VQA and Visual Grounding", https://arxiv.org/abs/1606.01847.



Subword Embeddings

ldings Seq2seq

Multilingual Translation

Many attempts at multilingual translation.

- Individual encoders and decoders, shared attention.
- Shared encoders and decoders.



Word2vec

