





	Encoder-Decoder MT				
	<ul> <li>Sutskever seq2seq paper: first major application of LSTMs to NLP</li> <li>Basic encoder-decoder with beam search</li> </ul>				
	Method test BLEU score (ntst14)				
Noural MT	Single forward LSTM, beam size 12 26.17				
	Single reversed LSTM, beam size 12 30.59				
	Ensemble of 5 reversed LSTMs, beam size 1 33.00				
	Ensemble of 2 reversed LSTMs, beam size 12 33.27				
	Ensemble of 5 reversed LSTMs, beam size 2 34.50				
	Ensemble of 5 reversed LSTMs, beam size 12 34.81				
	SOTA = 37.0 — not all that competitive Sutskever el				



## Results: WMT English-French

## 12M sentence pairs

Classic phrase-based system: ~33 BLEU, uses additional target-language data

Rerank with LSTMs: **36.5** BLEU (long line of work here; Devlin+ 2014)

Sutskever+ (2014) seq2seq single: 30.6 BLEU

Sutskever+ (2014) seq2seq ensemble: 34.8 BLEU

Luong+ (2015) seq2seq ensemble with attention and rare word handling: **37.5** BLEU

But English-French is a really easy language pair and there's tons of data for it



	MT Examples	Bac	<pre>ktranslation</pre>
src	Wegen der von Berlin und der Europäischen Zentralbank verhängten strengen Sparpolitik in Verbindung mit der Zwangsjacke , in die die jeweilige nationale Wirtschaft durch das Festhal- ten an der gemeinsamen Währung genötigt wird , sind viele Menschen der Ansicht , das Projekt Europa sei zu weit gegangen	<ul> <li>Classical MT methods used a tage monolingual corpus T' do the same?</li> </ul>	ilingual corpus of sentences B = (S, T) and to train a language model. Can neural MT
ref 	The austerity imposed by Berlin and the European Central Bank, coupled with the straitjacket imposed on national economies through adherence to the common currency, has led many people to think Project Europe has gone too far. Because of the strict austerity measures imposed by Berlin and the European Central Bank in	<ul> <li>Approach 1: force the system generate T' as targets from nu inputs</li> </ul>	<ul> <li>Approach 2: generate synthetic</li> <li>sources with a T-&gt;S machine</li> <li>translation system (backtranslation)</li> </ul>
base	<ul> <li>connection with the straitjacket in which the respective national economy is forced to adhere to the common currency, many people believe that the European project has gone too far.</li> <li>Because of the pressure imposed by the European Central Bank and the Federal Central Bank with the strict austerity imposed on the national economy in the face of the single currency ,</li> </ul>	s <sub>1</sub> , t <sub>1</sub> s <sub>2</sub> , t <sub>2</sub>	$s_{1}, t_{1}$ $s_{2}, t_{2}$
▶ be	many people believe that the European project has gone too far . st = with attention, base = no attention	 [null], t' <sub>1</sub> [null], t' <sub>2</sub>	MT(t'1), t'1 MT(t'2), t'2
	Luong et al. (2015)		Sennrich et al. (2015)

Backtranslation								
name training I data instances tst2011 tst201						tst2014		
baseline (Gülçe	18.4	18.8	19.9	18.7				
deep fusion (Gü	ilçehre et al., 2015)	7.2m	20.2	20.2	21.3	20.6		
parallel <sub>synth</sub>	parallel/parallel <sub>synth</sub>	6m/6m	18.6	20.4	20.1	20.0		
Gigaword <sub>mono</sub>	parallel/Gigaword <sub>mono</sub>	7.6m/7.6m	18.8	19.6	19.4	18.2		
Gigaword <sub>synth</sub>	parallel/Gigaword <sub>synth</sub>	8.4m/8.4m	21.2	21.1	21.8	20.4		
<ul> <li>Gigaword: large monolingual English corpus</li> </ul>								
<ul> <li>parallel<sub>synth</sub>: backtranslate training data; makes additional noisy source sentences which could be useful</li> </ul>								
Sennrich et al. (2015)								









		BLEU		
Π	D system	100k	3.2M	
1	phrase-based SMT	$15.87\pm0.19$	$26.60\pm0.00$	
- 2	NMT baseline	$0.00\pm0.00$	$25.70\pm0.33$	
3	2 + "mainstream improvements" (dropout, tied embeddings, layer normalization, bideep RNN, label smoothing)	$7.20\pm0.62$	31.93 ± 0.05	
4	$3 + \text{reduce BPE vocabulary (}14k \rightarrow 2k \text{ symbols)}$	$12.10\pm0.16$	-	
4	4 + reduce batch size (4k $\rightarrow$ 1k tokens)	$12.40\pm0.08$	$31.97\pm0.26$	
(	5 + lexical model	$13.03\pm0.49$	$31.80\pm0.22$	
5	5 + aggressive (word) dropout	$15.87\pm0.09$	<b>33.60</b> ± 0.14	
8	<ul> <li>7 + other hyperparameter tuning (learning rate, model depth, label smoothing rate)</li> </ul>	$\textbf{16.57} \pm 0.26$	$32.80\pm0.08$	
ç	8 + lexical model	$16.10\pm0.29$	$33.30\pm0.08$	

Sennrich and Zhang (2019)

## Frontiers in MT: Low-Resource > Particular interest in deploying MT systems for languages with little or no parallel data Burmese, Indonesian, Turkish BLEU BPE allows us to transfer Transfer $My \rightarrow En Id \rightarrow En Tr \rightarrow En$ models even without 4.0 20.6 19.0 baseline (no transfer) 17.8 27.4 20.3 transfer, train training on a specific 13.3 25.0 20.0 transfer, train, reset emb, train language 3.6 18.0 19.1 transfer, train, reset inner, train Table 3: Investigating the model's capability to restore Pre-trained models can its quality if we reset the parameters. We use $En \rightarrow De$ help further as the parent.

Aji et al. (2020)







## Properties of Self-Attention

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	O(1)	<i>O</i> (1)
Recurrent	$O(n \cdot d^2)$	O(n)	O(n)
Convolutional	$O(k \cdot n \cdot d^2)$	O(1)	$O(log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	O(1)	O(n/r)

- n = sentence length, d = hidden dim, k = kernel size, r = restricted neighborhood size
- Quadratic complexity, but O(1) sequential operations (not linear like in RNNs) and O(1) "path" for words to inform each other

```
Vaswani et al. (2017)
```









