











![](_page_3_Figure_0.jpeg)

- One loss term for each target-sentence word, feed the correct word regardless of model's prediction
- Length of gold sequence is known, can run the whole encoder-decoder in one computation graph and compute losses

![](_page_3_Figure_3.jpeg)

Bengio et al. (2015)

![](_page_3_Figure_5.jpeg)

![](_page_4_Figure_0.jpeg)

Problems with Neural MT Models	Problems with Neural MT Models
Encoder-decoder models like to repeat themselves:	Unknown words:
Un garçon joue dans la neige → A boy plays in the snow <b>boy plays boy plays</b>	<i>en</i> : The <u>ecotax</u> portico in <u>Pont-de-Buis</u> , [truncated], was taken down on Thursday morning <i>fr</i> : Le <u>portique écotaxe</u> de <u>Pont-de-Buis</u> , [truncated], a été <u>démonté</u> jeudi matin
<ul> <li>Often a byproduct of training these models poorly</li> </ul>	nn: Le <u>unk</u> de <u>unk</u> à <u>unk</u> , [truncated], a été pris le jeudi matin
<ul> <li>Solution: include coverage in the model so we don't repeat stuff: Haitao Mi et al. (2016) for MT, See and Manning (2017) for summarization</li> </ul>	<ul> <li>We restricted the target vocabulary to 80,000 — that throws out a lot!</li> <li>Fixed vocabulary is too restrictive, especially around named entities</li> </ul>

![](_page_5_Figure_0.jpeg)

![](_page_5_Figure_1.jpeg)

![](_page_6_Figure_0.jpeg)

## Machine Translation Results

WMT English-French: 12M sentence pairs, 80,000 word target vocab

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

Bahdanau+ (2014) seq2seq with attention: 28.5 BLEU

• But English-French is a really easy language pair!

Results from Luong et al. (ACL 2015)

![](_page_6_Picture_10.jpeg)

![](_page_7_Figure_0.jpeg)

![](_page_7_Figure_1.jpeg)

![](_page_8_Figure_0.jpeg)

## Google's NMT System

English-French:

۲

Google's phrase-based system: 37.0 BLEU Luong+ (2015) seq2seq ensemble with rare word handling: 37.5 BLEU Google's 32k word pieces: 38.95 BLEU

English-German:

Google's phrase-based system: 20.7 BLEU Luong+ (2015) seq2seq ensemble with rare word handling: 23.0 BLEU Google's 32k word pieces: 24.2 BLEU

Wu et al. (2016)

![](_page_8_Figure_7.jpeg)

![](_page_9_Figure_0.jpeg)

Takeaways	
<ul> <li>RNNs are effective at machine translation, but lots of tricks to to work right</li> </ul>	get them
Attention is a critical way to get a better representation of the	input
Handling rare words is important, lots of techniques here	
<ul> <li>Encoder-decoder models can be successfully applied to most t where you generate language as output</li> </ul>	asks