Attention



- Encoder hidden states capture contextual source word identity
- Decoder hidden states are now mostly responsible for selecting what to attend to
- Doesn't take a complex hidden state to walk monotonically through a sentence and spit out word-by-word translations

Attention



Neural MT



- 12M sentence pairs
- Classic PBMT system: ~33 BLEU, uses additional target-language data PBMT + rerank w/LSTMs: **36.5** BLEU (long line of work here; Devlin+ 2014) Sutskever+ (2014) seq2seq single: **30.6** BLEU (input reversed) 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! Does this approach work for anything harder?





- 4.5M sentence pairs
- Classic phrase-based system: **20.7** BLEU Luong+ (2014) seq2seq: **14** BLEU
- Not nearly as good in absolute BLEU, but BLEU scores aren't really comparable across languages
- French, Spanish = easiest German, Czech = harder

Results: WMT English-German

Luong+ (2015) seq2seq ensemble with rare word handling: 23.0 BLEU

Japanese, Russian = hard (grammatically different, lots of morphology...)



MT Examples

src	In einem Interview sagte Bloom jedoch
ref	However, in an interview, Bloom has s
best	In an interview, however, Bloom said
base	However, in an interview, Bloom said

- best = with attention, base = no attention
- phrase-based doesn't do this

, dass er und Kerr sich noch immer lieben .

said that he and *Kerr* still love each other.

that he and *Kerr* still love.

that he and **Tina** were still $\langle unk \rangle$.

NMT systems can hallucinate words, especially when not using attention

Luong et al. (2015)





MT Examples

Wegen der von Berlin und der Europäis
Verbindung mit der Zwangsjacke, in die
ten an der gemeinsamen Währung genötig
Europa sei zu weit gegangen
The austerity imposed by Berlin and the
imposed on national economies through a
to think Project Europe has gone too far .
Because of the strict austerity measures
connection with the straitjacket in which
the common currency, many people belie
Because of the pressure imposed by the E
with the strict austerity imposed on the
many people believe that the European pro-

best = with attention, base = no attention

schen Zentralbank verhängten strengen Sparpolitik in e die jeweilige nationale Wirtschaft durch das Festhalgt wird, sind viele Menschen der Ansicht, das Projekt

European Central Bank, coupled with the straitjacket dherence to the common currency, has led many people

imposed by Berlin and the European Central Bank in the respective national economy is forced to adhere to eve that the European project has gone too far. uropean Central Bank and the Federal Central Bank e national economy in the face of the single currency, oject has gone too far.

Luong et al. (2015)







- Words are a difficult unit to work with: copying can be cumbersome, word vocabularies get very large
- Character-level models don't work well
- Compromise solution: use thousands of "word pieces" (which may be full words but may also be parts of words)

Input: _the _eco tax _port i co _in _Po nt - de - Bu is ...

Output: _le _port ique _éco taxe _de _Pont - de - Bui s

Can achieve transliteration with this, subword structure makes some translations easier to achieve Sennrich et al. (2016)

Handling Rare Words





- for i in range(num_merges): pairs = get_stats(vocab) cooccurrences best = max(pairs, key=pairs.get) vocab = merge_vocab(best, vocab)

- many whole words
- Most SOTA NMT systems use this on both source + target

Byte Pair Encoding (BPE)

Start with every individual byte (basically character) as its own symbol

- Count bigram character
- Merge the most frequent pair of adjacent characters

Doing 8k merges => vocabulary of around 8000 word pieces. Includes

Sennrich et al. (2016)







	Original:	furiously			
(a)	BPE:	_fur	_fur iously ((
	Unigram LM:	_fur	ious	ly	
	Original:	Comp	letely p	reposter	r
(c)	BPE:	_Com	ple t	ely	
	Unigram LM:		mplete	l ly	

BPE produces less linguistically plausible units than another technique based on a unigram language model

Byte Pair Encoding (BPE)

Original: tricycles BPE:_tricycUnigram LM:_tricycles (b)cles rous suggestions

_prep | ost | erous | _suggest ions nplete | ly | _pre | post | er | ous | _suggestion | s

Bostrom and Durrett (2020)







8-layer LSTM encoder-decoder with attention, word piece vocabulary of 8k-32k

Google's NMT System





English-French:

- 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

Google's NMT System





Human Evaluation (En-Es)

200

100

0

Similar to human-level 400 performance on English-Spanish 300 Count (total 500)



PBMT - GNMT - Human





Source	She was spotted three days later by a
PBMT	Elle a été repéré trois jours plus tard j
GNMT	Elle a été repérée trois jours plus tard
Human	Elle a été repérée trois jours plus tard coincée dans la carrière

Gender is correct in GNMT but not in PBMT

Google's NMT System



