

Motivation

Why document-level NMT?

- Not considering the document context and discourse connections affects coherence and cohesion of a text.

Why hierarchical attention networks?

- Different abstraction levels: word-level and sentence-level.
- Allows dynamic access to the context for each predicted word.

Other advantages in our approach

- Joint optimization of multiple sentences.
- Shared hidden representations across sentence translations.
- Exploiting source and target context.
- Multi-head attention to capture different discourse phenomena.

Document-level NMT

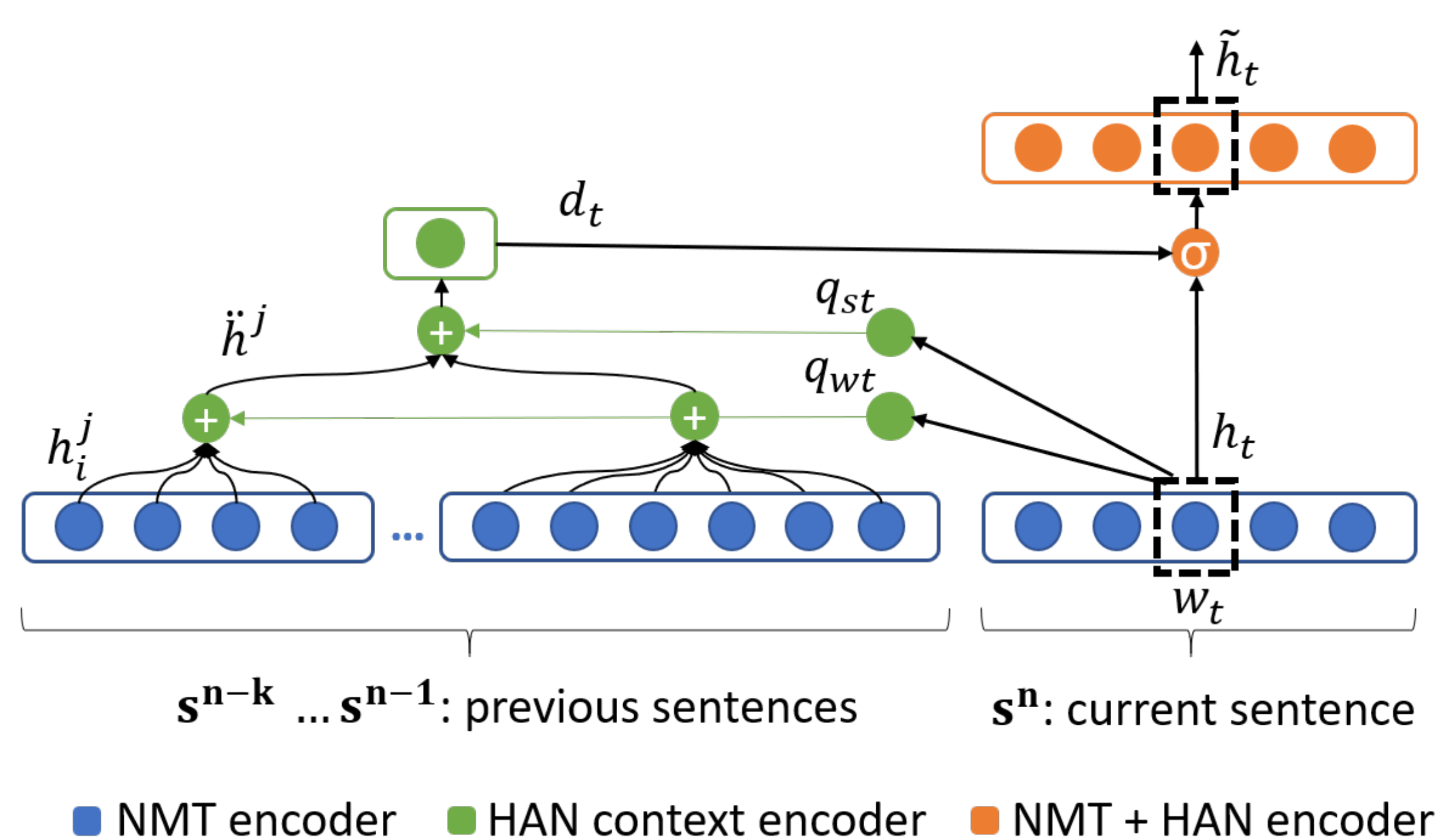
$$\max_{\Theta} \sum_{n=1}^N \log(P_{\Theta}(t^n | s^n)) \rightarrow \max_{\Theta} \sum_{n=1}^N \log(P_{\Theta}(t^n | s^n, D_{s^n}, D_{t^n}))$$

Baseline NMT

Document-level NMT

- s^n : source sentence, and $D_{s^n} = (s^{n-k}, \dots, s^{n-1})$: source context
- t^n : target sentence, and $D_{t^n} = (t^{n-k}, \dots, t^{n-1})$: target context.
- Context (k previous sentences) is modeled by HANs:

Hierarchical Attention Network (HAN)



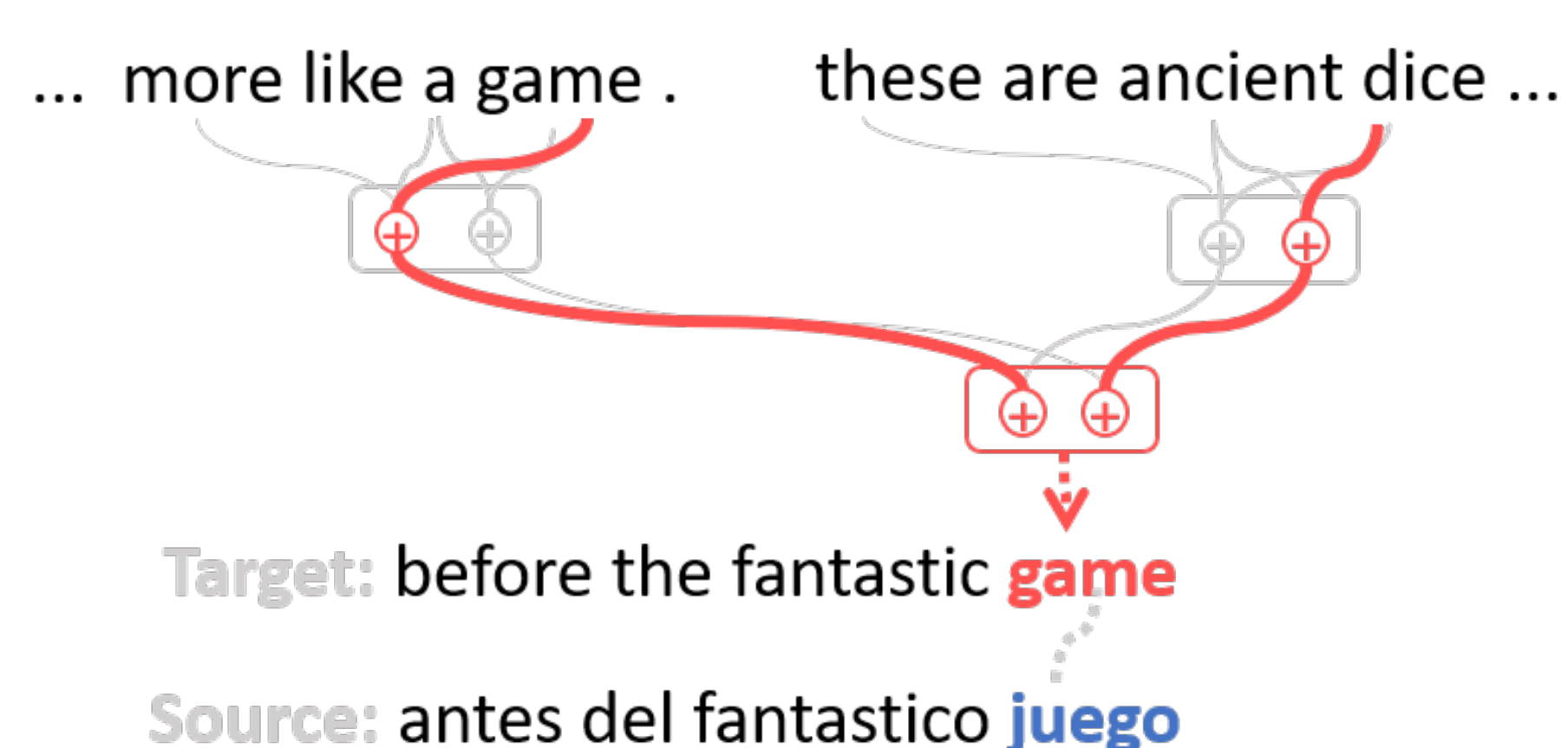
Word-level attention:

$$\tilde{h}^j = \text{MultiHead}(q_{wt}, h_i^j) \quad q_{wt} = f_w(h_t)$$

Sentence-level attention:

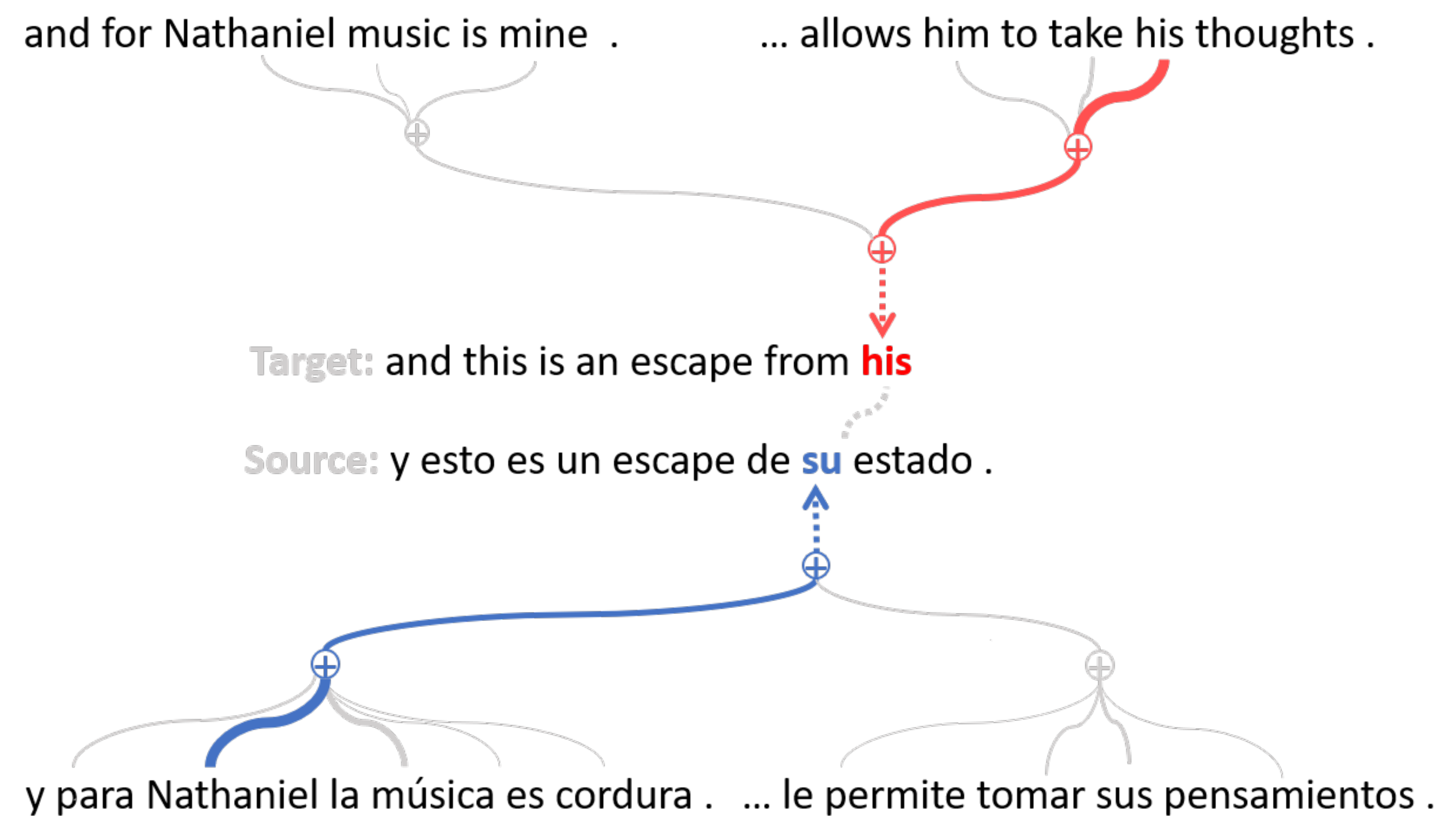
$$d_t = \text{FFN}(\text{MultiHead}(q_{st}, \tilde{h}^j)) \quad q_{st} = f_s(h_t)$$

Multi-head attention



- “juego” can be translated as “game” or “set”

Source and Target Sides Context



- “su” can be translated as “his”, “her”, or “its”

Experimental Results

	TED Talks		Subtitles		News
	Zh-En	Es-En	Zh-En	Es-En	Es-En
NMT transformer	16.87	35.44	28.60	35.20	21.36
+ cache	17.32 ***	36.46 ***	28.86	35.49	22.36 ***
+ HAN encoder	17.61 ***	36.91 ***	29.35 †	35.96 †	22.36 ***
+ HAN decoder	17.39 ***	37.01 ***	29.21 *	35.50	22.62 ***
+ HAN joint	17.79 ***	37.24 ***	29.67 **	36.23 **	22.76 ***

BLEU scores. Significance with respect to NMT *, and to cache model †.
P-values: † < .05, ** < .01, *** < .001.

- Significant improvement over strong baselines on multiple data sets.
- Context from source and target sides are complementary.

Discourse Evaluation

	Coherence	Lexical Cohesion	Pronouns	Nouns
NMT transformer	28.42	47.98	62.84	52.50
+ HAN encoder	28.60	48.35	64.48	53.61
+ HAN decoder	28.78	48.51	64.04	53.55
+ HAN joint	28.82	48.61	64.32	54.19
Human reference	29.79	52.94	100.0	100.0

- HAN decoder helps in lexical cohesion and coherence.
- HAN encoder helps in pronoun and noun disambiguation.

Conclusion

- We proposed a hierarchical multi-head attention model for document-level context.
- It directly connects representations from previous sentence translations into the current sentence translation.
- It significantly outperforms two competitive baselines.
- It improves cohesion and coherence, and noun/pronoun translation.
- We show that target and source context is complementary.

Our multi-head HAN could be used to model context in other NLP tasks. Code available at https://github.com/idiap/HAN_NMT