# Improved Machine Translation for CN-EN Novels through Targeted Dictionary Substitutions of Idioms and a Multi-Word-Phrase Synonym-Focused Automated Metric

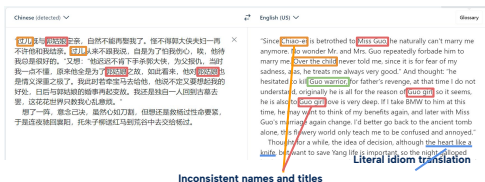Lisa Liu, Ryan Liu, Angela Tsai     Advisor: Jingbo Shang

**UC San Diego**
**HALICIOĞLU DATA SCIENCE INSTITUTE**

## BACKGROUND

Machine translation is frequently used in the CN-EN webnovel translation scene, as human translators cannot keep up with an ever-growing body of texts. MT struggles with idioms, slang, and other culturally nuanced phrases found in webnovels. These phrases have meanings beyond the actual words used, but machines often incorrectly translate them literally. As such, our project aims to build a MT model that more accurately identifies and translates Chinese idioms could help reduce the amount of human post-editing needed to produce higher-quality translations, and increase the rate at which CN literature can be accessed by EN audiences.

## SAMPLE TRANSLATION

**Direct DeepL Translation:**



Inconsistent names and titles
Literal idiom translation

**Pipelined DeepL Translation:**

"**Guo'er** is betrothed to **Miss Guo**, naturally he can't marry me anymore. No wonder Mr. and Mrs. **Hero Guo** have repeatedly forbidden him to marry me. **Guo'er** never told me, because he was afraid that I would be sad, alas, he always treats me very well." And thought: "He hesitated to kill **Hero Guo**, for father's revenge, at that time I did not understand, the original he was all for **Miss Guo's** sake, so it seems, he is also **Miss Guo** is also a very deep love. If I took the BMW to him, he might want to remember my favor again, and the marriage with Miss Guo should be changed again in the future. I'd better go back to the ancient tomb alone, this flowery world only teaches me to be disturbed."
After thinking for a while, I decided to **put my heartache aside** and focus on saving Yang's life...
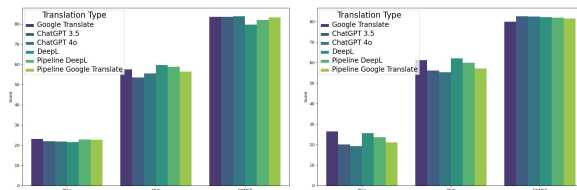
## METHODOLOGY

1. Use Jieba to identify names with its built-in dictionary
2. Select sample subset of sentences that do not contain any names identified by Jieba, and use LLM to find and translate any missed names
    a. Weight sentences earlier in the novel higher for selection, as characters are more likely to be introduced then
3. Combine newly identified names with Jieba's name dictionary; replace common suffixes and titles with standardized translation
4. Using names dictionary, replace all instances of names in original text with their translations
5. Using CC-CEDICT, replace all instances of idioms with their translations
6. Use low-cost translation model such as Google Translate to translate remaining text

## CONCLUSIONS

1. Pipeline introduces much more consistency in resulting translations
    a. Consistent suffixes and names
    b. Prevent literal idiom translations
2. Compromise of cost/efficiency of NMTs and performance of LLMs
3. Currently only feeds 20% of the text into GPT for finding names
    a. 15% from first ¾, and 5% from last ¼
    b. Rest of text handled by NMTs
4. Can perform better than GPT in word-level evaluations (BLEU, ChrF)
5. Can perform better than NMTs in predicted human reception (COMET)

## RESULTS



## FUTURE WORK

In addition to improving the MT process with targeted idiom translations, one of our original goals was to also create our own automatic evaluation metric to more accurately judge the translation. Going forward, we aim to build an automated metric focused on multi-word synonyms, using another metric. METEOR, as a base METEOR is an existing metric that performs synonym matching using WordNet's synsets, which only contain unigrams, or single words. We will use a dictionary to form multi-word synsets for another stage of matching.

## ACKNOWLEDGEMENTS