

# Dataprep

---

Standardize data format. Retrieve source sentences.

**all\_annotations\_v1.json** contains littrans data is from <https://github.com/marzenakrp/LiteraryTranslation>

**wmt23/** contains en-de and de-en datasets from the WMT2023 testsets

## **json2csv\_littrans.py**

Parses a json file with annotations and creates csv files for each book (language pair): para (gpt3, human) and sent formatted as para (gpt3, nmt). Extracts human preferences into a csv file  
output/littrans\_annotators\_choices.csv

## **txt2csv\_wmt23.py**

Converts the WMT23 txt files to csv files, merging the source and target languages into one file. Para (human, gpt4), sent (nmt).

## **run\_csv2json4Llama.sh** -> csv2json4Llama.py

Iterates through all\_csv/{lang}.para.human.csv and extracts source paragraphs into json files formatted for Llama.

## **split\_sents.py** needs GPU

Iterates through all\_csv/{lang}.para.human.csv files, preprocesses source texts and standardizes punctuation based on lang prior to segmentation. Splits source texts into sentences. Writes json files formatted for Llama, writes txt files.

# Create translations

## **translate\_gpt.sh** -> translate\_with\_openAI.py

Iterates through /dataprep/source\_\${level}\_json/\*.json. Uses OpenAI API to produce translations with GPT-3 and GPT-4. Saves files to translated/\${level}-level. Script needs to be manually adjusted depending on level and model. Read annotation.

## Translating with Llama

1. Translate (needs GPU)
2. **run\_json2csv4Llama.sh** -> json2csv.py

Converts json files from llama\_translations/llama\_{level}\_json to  
llama\_translations/llama\_{level}\_csv

3. **clean\_Llama\_with\_gpt4.py**

Iterates through llama\_translations/llama\_{level}\_csv and flags missing translations with NO TRANSLATION FOUND. Flagged lines are sent back to the model for re-evaluation, which

produces flags: <<WRONG STATEMENT, TRANSLATION FOUND>>, <<INACCURATE TRANSLATION>>, and <<CORRECT STATEMENT, NO TRANSLATION FOUND, because>>  
Writes files to llama\_translations/llama\_{level}\_gpt4\_cleaned/ Make sure to indicate the "id" number of the line where to start processing file.

4. Feed flagged src-tgt pairs back to Llama for re-translation.

5. **remove\_gpt4\_flags.py**

6.