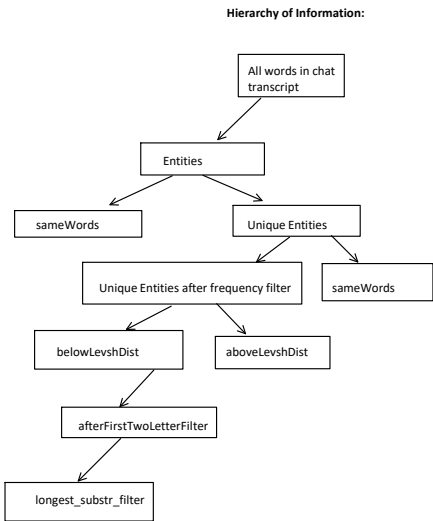


Structure of data from Levana's dataset (testSpellError.csv)



	Definition
Entities	Can be one of the following <ul style="list-style-type: none"><li>- Action:<ul style="list-style-type: none"><li>Handled by testSpellingErrorAct.py file</li><li>Fourth column from the dataset Levana sent</li><li>Compared against the action entity library</li></ul></li><li>- General Entity<ul style="list-style-type: none"><li>Handled by testSpellingErrorProd.py file</li><li>Fifth column from the dataset Levana sent</li><li>All entities in fifth column with confidence "middle"</li><li>Compared against the general entity library</li></ul></li><li>- Product Entity<ul style="list-style-type: none"><li>Handled by testSpellingErrorProd.py file</li><li>Fifth column from the dataset Levana sent</li><li>All entities in fifth column with confidence "high"</li><li>Compared against the product entity library; compares against both entities and synonyms to see if there is a match</li></ul></li></ul>
aboveLevshDist	All <b>Entities</b> with a Levenshtein distance higher than the Levenshtein threshold compared to ANY of the entity library words <ul style="list-style-type: none"><li>- Informative to look through to see which ones to add to the library because these are the entities that were different from any of the current library words that they could be completely different terms</li></ul>
Unique Entities	Subset of <b>belowLevshDist</b> that all have Levenshtein distance > 0 compared to library word (i.e, filters out all identified library words) <ul style="list-style-type: none"><li>- These are instances, so can be repetitions of the same unique entity word/phrase</li><li>- Metric for how representative the library words are; if number unique entities is significantly less than total entities, that means library has strong coverage of the identified entities</li></ul>
Unique entities after frequency filter	All unique entities after filtering out the most frequent occurrences of the same entity words <ul style="list-style-type: none"><li>- Frequency analysis handled in FreqDist.py file; removes words that occur greater than or equal to 2 (or 3 depending on entity) times in the entire dataset. <b>Rationale:</b> The more frequent that the same word (same ordering of same set of letters) occurs, the more unlikely it is to be a spelling error because multiple users have applied that same spelling</li></ul>
belowLevshDist	All <b>unique entities after frequency filter</b> with a Levenshtein distance <= the Levenshtein threshold compared to AT LEAST ONE entity library word <ul style="list-style-type: none"><li>- List that matches each instance of an entity with an entity library word that meets the above requirement</li></ul>
sameWords	Subset of <b>Unique Entities</b> that have Levenshtein distance of 0 with an action entity library word (in other words, an action entity library word has been identified in the chat transcript)
afterFirstTwoLetterFilter	Subset of <b>belowLevshDist</b> ; must contain the same first letter as the action entity library word matched to
longest_substr_filter	Subset of <b>afterFirstTwoLetterFilter</b> . Final filter in the process that requires matches to have substring of at least x characters in common, where x depends on the number of letters in the string of the identified entity <ul style="list-style-type: none"><li>- Output: <i>[entity with potential spelling, a potential match from the entity library based on similarity requirements above]</i></li></ul>

SpellCheck Output Key Highlights:

	Number of identified entities in dataset	Number of Unique entities (instances)	Number of unique entities after frequency filter	After first two letter filters	After final filter (longest substring filter)	Entities above Levenshtein distance
Product Entities	30,402	1,134	7	['window hello', 'windows hello'], ['window hello', 'windows hello'], ['Windows IoT Core', 'windows 10 iot core'], ['Windows IoT Core', 'windows 10 iot core'], ['Windows IoT Core', 'windows 10 iot core'], ['Windows IoT Core', 'windows 10 iot core'], ['window hello', 'windows hello'], ['window hello', 'windows hello'], ['window PC', 'windows pc'], ['window PC', 'windows pc'], ['window PC', 'windows pc']	['window hello', 'window hello'], ['Windows IoT Core', 'windows 10 iot core'], ['Windows IoT Core', 'windows 10 iot core'], ['window hello', 'window hello'], ['window PC', 'window PC']	['window hello', 'windows hello'], ['Windows IoT Core', 'windows 10 iot core'], ['Windows IoT Core', 'windows 10 iot core'], ['window hello', 'windows hello'], ['window PC', 'windows pc']
General Entities	29,120	10,079	3,781	['account id', 'accounting'], ['account id', 'accounting'], ['subscription id', 'subscription'], ['subscription id', 'subscriptions'], ['transcript id', 'transcripts']	['account id', 'accounting'], ['account id', 'accounting'], ['subscription id', 'subscription'], ['subscription id', 'subscriptions'], ['transcript id', 'transcripts']	See separate document
**Actions	59,523	(after frequency elimination) 63	(after same words filter) 23	['eliminate', 'eliminated'], ['activate', 'achieve'], ['activate', 'activate'], ['activate', 'activate'], ['associate', 'associate'], ['associate', 'associate'], ['associate', 'associate'], ['associate', 'associate'], ['generate', 'generate'], ['separate', 'separate'], ['separate', 'separate'], ['support', 'support'], ['acced', 'accept'], ['acced', 'access'], ['create', 'create'], ['create', 'created'], ['support', 'support'], ['clash', 'clashes'], ['lapse', 'lapses']	['eliminate', 'eliminated'], ['activate', 'activate'], ['activate', 'activate'], ['associate', 'associate'], ['associate', 'associate'], ['associate', 'associate'], ['associate', 'associate'], ['generate', 'generate'], ['separate', 'separate'], ['separate', 'separate'], ['support', 'support'], ['acced', 'accept'], ['acced', 'access'], ['create', 'create'], ['create', 'created'], ['support', 'support'], ['clash', 'clashes'], ['lapse', 'lapses'], ['earnrenew', 'earn'], ['earnrenew', 'renew']	['action', 'disassociation', 'tear'] <b>**These are spelling errors that are not being caught by the algorithm, perhaps because it still differs from the vocab list</b>

\*\*filtered out most common actions first, then library word occurrences mainly for performance and efficiency reasons

- Given proportion of action entity instances and actual unique action entities, would have been inefficient to compare every instance against library word first
- Also want to compare the lemmatized version of the action to the library words, would take way too long to lemmatize 59,523 words instead of 63

Frequency Distribution of different entities:

General Entity Frequency Analysis:

- FreqDist with 3754 samples and 29120 outcomes (so 3654 unique general entities out of 29120 total instances of general entities identified)

Product Entity Frequency Analysis:

- FreqDist with 271 samples and 30402 outcomes (so 271 unique product entities out of 30402 total instances of product entities identified)

Action Entity Frequency Analysis:

- <FreqDist with 319 samples and 59523 outcomes> (so 391 unique action entities out of 59523 total instances of action entities identified)

However, don't know the frequency distribution

Analysis of above results:

Product Entities:

- The fact that the majority of the identified entities in the dataset were eliminated after searching for identical library words indicates the current state of the product entity library is highly representative of the product entities identified from a given user chat transcript.
- Only 7 unique entities after eliminating the product entities with the most frequent occurrences, showing there is less likely to be spelling errors among the product entities. 203 out of the 271 samples (unique product entities) occurred more than once (i.e, multiple occurrences of the same ordering of the same set of letters); only 68 samples had a singular frequency in the dataset
- Final filter after the longest substring evaluation shows that filtering process is effective in eliminating noisy data and outputting results that would actually be potential spell errors

General Entities:

- 34% of the total entities identified in dataset are instances of unique entities (not in the library). Around 37% of the unique entities remain after eliminating the most frequent ones. Shows decent coverage of the general entity library in its current state but could be improved with more manual review
- Most are eliminated after the spell check filtering process as there are only 2 out of 3,781 total entities being evaluated for spell check that were identified to be similar enough to any of the library words. Indicates either that library is not representative enough of actual general entities, or that there are rarely ever spelling errors

Action Entities:

- Likely to hold highest accuracy rate in detecting actual spelling errors, given that majority of final output are actual misspelled words (whereas the output in product or general entities were correctly spelled words that were just similar *enough* in distance to some entity library word to be flagged).
- Failed to catch the spelling errors in aboveLevshDist however