

## 4. WordNet – Antonyms (APairs)

### 1. Antonyms form WordNet

- From WordNet:
  - Directory: \${WORDNET}
  - File: \${WORDNET}/data/output/**WnAPairs.data.3.0**
  - **WnAPairs.data.3.0**: 12,248 aPairs
  - format: base1|base2|pos
  - These are antonym pairs from WordNet directly.
  - Antonyms must have the same POS
  - Antonyms include single words and multiwords
- **Task 4.1: Generate aPair candidates from WordNet aPairs.**
  - Directory: \${LEXICON\_ANTONYM}/../WordNetAPairs
  - Input file: WnAPairs.data.3.0
  - Output File: \${LEXICON\_ANTONYM}/data/6.WordNet/2021/outData/WordNet/Cand
  - Program: GenAPairCand.java
  - Used Lexicon.2021 for the initial baseline data (Log.WnApCand.2021)

Steps	Description	Notes
1	Unify sPair and remove duplicates	<ul style="list-style-type: none"> <li>• WordNet aPairs are duplicated and are not in a standard format. They are unified and standardized in this step.</li> <li>• WnAPairs.unique.data (8,242)</li> </ul>
2	Remove Lexicon SpVar	<ul style="list-style-type: none"> <li>• WordNet aPairs include spVars. SpVars are not aPairs and need to be removed.</li> <li>• Wn.aPairCand.data.trap.spVar (2)</li> </ul>
3	Apply combined filters	<ul style="list-style-type: none"> <li>• Words in aPairs must be valid words. Use combined filter to remove invalid words.</li> <li>• wn.aPairCand.data.trap.cf (107)</li> </ul>
4	Filter out Lexicon sPairs (2022)	<ul style="list-style-type: none"> <li>• Remove Lexicon sPairs.</li> <li>• wn.aPairCand.data.trap.sp (1)</li> </ul>
5	Baseline output (sorted)	<ul style="list-style-type: none"> <li>• wn.aPairCand.data.baseline (8,132)</li> </ul>
6	Filter out Lexicon tagged aPairs	<ul style="list-style-type: none"> <li>• wn.aPairCand.data.yes (964)</li> <li>• wn.aPairCand.data.no (380)</li> <li>• wn.aPairCand.data.tbd (6,788)</li> </ul>

- The baseline candidate file (used for paper) use Lexicon.2021. This is used to calculate precision for paper.
- The final candidate file sent to linguist should use the latest Lexicon, by re-running the program to get the latest results.
- The procedures are:
  - Send all word candidates from Task below (Task 4.2) to linguist for LexBuilding.

- Update inflVars (inflVars.data.current) and notBaseWord.data once above step is done
    - Go through annual aPairs process to update aPairs.
    - Re-run this step with latest Lexicon and notBaseWord.data, The tbd should become zero because invalid aPairs are excluded from candidates or in the tagged file and valid aPairs are in the tagged file.
  - For the baseline, there are 1344 (= 964+380) among 8,132 WordNet aPair candidates (16.53% = 1344/8132) are overlapped). The precision is estimate as 71.13% (=964/(964+380)). This implies there are much more aPairs could be retrieved from the WordNet, it is a good resource for Lexicon aPair retrieval.
  - WordNet includes 2 pairs of spVars as aPair:
    - kern|kern|verb
    - ravel|ravel|verb
  - WordNet includes 1 sPair as aPair:
    - mushroom|toadstool|noun
- **Task 4.2: Generate word candidates from aPairs in WordNet (To be updated)**
    - Directory: \${LEXICON\_ANTONYM}/../6.WordNet
    - Input file:
      - \${LEXICON\_ANTONYM}/data/6.WordNet/2021/outData/Cand/wn.aPairCand.data.2021
    - Output File:
      - \${LEXICON\_SYNONYM}/data/6.WordNet/2021/outData/Cand/wn.Ap.wordCand
    - Program: GenWordCandFromAPairs.java
    - Used Lexicon.2021 for the initial baseline data (Log.WnApWordCand.2021)

Steps	Description	Notes
0	Retrieve words from sPairs	Get unique words from aPairs (9,904)
1	Combined filters <ul style="list-style-type: none"> <li>• Standard filter</li> <li>• Single word filters</li> <li>• Multiword filters</li> </ul>	<ul style="list-style-type: none"> <li>• Apply combined filter to filter out invalid word</li> <li>• wn.Ap.wordCand.1.cfTrap (7,595)</li> </ul>
2	WordNet Previous cand filter <ul style="list-style-type: none"> <li>• verb complement</li> <li>• zeroD</li> <li>• suffixD</li> </ul>	<ul style="list-style-type: none"> <li>• Exclude words from previous WordNet word candidates</li> <li>• wn.Ap.wordCand.2.pcTrap (373)</li> </ul>
3	MEDLINE n-gram filter	<ul style="list-style-type: none"> <li>• Use both MEDLINE 1-gram and n-grams (WC &gt;=30)</li> <li>• wn.Ap.wordCand.3.mlTrap (695)</li> </ul>
4	UMLS CUI filter	<ul style="list-style-type: none"> <li>• Use UMLS CUI</li> <li>• wn.Ap.wordCand.4.cuiTrap (1,355)</li> </ul>
	Output	<ul style="list-style-type: none"> <li>• <b>wn.Ap.wordCand.cand.2021 (1,936)</b></li> <li>• wn.Ap.wordCand.cuiPass (581, have CUIs)</li> </ul>

- Use Lexicon.2021 for the baseline in the paper.
- Use Lexicon.current for practice and send to linguists.
- Combined filters: Lexicon, general, pattern, single-word, multiword-lead-end-term filters.
- Previous Candidate filters: not in the previous candidate list in the WordNet model (verb-complement, zeroD and suffixD)
- MEDLINE filter: use both 1-gram for single words and n-gram (wc >= 30) for multiwords
- CUI filter: must have UMLS CUIs
- Output: **wn.Ap.wordCand.cand**
  - Use Lexicon current (recompile) and WordNet.aPair.2021 to generate this file.
  - Send to linguists for lexBuilding (**TBD**)
    - No need to tag
    - Analyze the precision on WC and ML and CUI
  - Format: word|POS|WC|Tags
  - Includes words in/not in MEDLINE, with/without CUIs, maybe use a cutoff WC for initial test (**wn.Ap.wordCand.cand.wc10000, wc > 10000**)
  - Tags:
    - TAG\_ML\_1G\_CUI\_NO, TAG\_ML\_1G\_CUI\_YES
    - TAG\_ML\_5G\_CUI\_NO, TAG\_ML\_5G\_CUI\_YES
    - TAG\_NOT\_ML\_CUI\_NO, TAG\_NOT\_ML\_CUI\_YES

===== What we used before this new Model (To BE Deleted) =====

- The antonym criteria (STI, CUI, Single word, EUI, etc., 693 aPairs) are not used here.
- If not use STI, CUI criteria, it extended to 8,252 aPairs (can be used in the future)
- Check if in aPair tagged file (216 aPairs from previous tag, 693-216=477): 477 aPairs.
- Convert to aPair candidate format (include ants without EUIs): 477 aPairs.
- Shell> GetAntonyms 2021

71

- Output: \${ANTONYM}/data/6.WordNet/outpt/Cand/APairCand.data
- Tagged: [N]:155, [Y]: 215, [TBD]: 107, total 370.