

# CSE 576: Topics in Natural Language Processing

## Dataset - 3

### Phase 2 Report

**Task Statement:** Understand the task and think about alternative ways to represent the same task. To create a variant task, you can write alternate definitions by changing writing style, vocabulary, length of definition etc. and change the instances generated by the dataset

**Solution:** We created 1-5 variants of the task definitions by paraphrasing the original definitions.

**Created by:** Anshita Singh Bais, Prutha Gaherwar

#### List of tasks:

task072\_abductivenli\_answer\_generation.json  
task073\_commonsenseqa\_answer\_generation.json  
task074\_squad1.1\_question\_generation.json  
task075\_squad1.1\_answer\_generation.json  
task076\_splash\_correcting\_sql\_mistake.json  
task077\_splash\_explanation\_to\_sql.json  
task078\_splash\_sql\_to\_explanation.json  
task079\_conala\_concat\_strings.json  
task080\_piqa\_answer\_generation.json  
task081\_piqa\_wrong\_answer\_generation.json  
task082\_babi\_t1\_single\_supporting\_fact\_question\_generation.json  
task083\_babi\_t1\_single\_supporting\_fact\_answer\_generation.json  
task084\_babi\_t1\_single\_supporting\_fact\_identify\_relevant\_fact.json  
task085\_unnatural\_addsub\_arithmetic.json  
task086\_translated\_symbol\_arithmetic.json  
task087\_new\_operator\_addsub\_arithmetic.json  
task088\_identify\_typo\_verification.json  
task089\_swap\_words\_verification.json  
task090\_equation\_learner\_algebra.json  
task092\_check\_prime\_classification.json  
task093\_conala\_normalize\_lists.json  
task094\_conala\_calculate\_mean.json  
task095\_conala\_max\_absolute\_value.json  
task096\_conala\_list\_index\_subtraction.json  
task097\_conala\_remove\_duplicates.json

task098\_conala\_list\_intersection.json  
task102\_commongen\_sentence\_generation.json  
task103\_facts2story\_long\_text\_generation.json  
task104\_ semeval\_2019\_task10\_closed\_vocabulary\_mathematical\_answer\_generation.json  
task105\_story\_cloze-rocstories\_sentence\_generation.json  
task106\_scruples\_ethical\_judgment.json  
task107\_splash\_question\_to\_sql.json  
task108\_contextualabusedetection\_classification.json  
task109\_smsspamcollection\_spamsmsdetection.json  
task110\_logic2text\_sentence\_generation.json  
task111\_asset\_sentence\_simplification.json  
task112\_asset\_simple\_sentence\_identification.json  
task113\_count\_frequency\_of\_letter.json  
task114\_is\_the\_given\_word\_longest.json  
task115\_help\_advice\_classification.json  
task116\_com2sense\_commonsense\_reasoning.json  
task117\_spl\_translation\_en\_de.json  
task118\_ semeval\_2019\_task10\_open\_vocabulary\_mathematical\_answer\_generation.json  
task119\_ semeval\_2019\_task10\_geometric\_mathematical\_answer\_generation.json  
task119\_zest\_text\_modification.json  
task120\_zest\_text\_modification.json

### **For Tasks 72-95:**

Number of variants created per task -

5

Tools used -

DiverseParaphraser  
Casual2Former  
Former2Casual  
Quillbot.com  
Styleformer  
MySecondTransformation

Number of Instances Generated -

In the range 200 - 6500

New Instances generated for -

Task 72, 73, 74, 75, 80, 81

Tasks for which the instances were shuffled as all the instances of the dataset were used -

Tasks 76, 77, 78, 79, 82, 83, 84, 85, 86, 87, 88, 89, 90, 92, 93, 94, 95

Datasets Used -

Abductive NLI  
CommonsenseQA  
SQuad1.1

Splash  
Conala  
PIQA  
BABI  
Synthetic

Types of Definitions Generated-

Paraphrase, Long, Short, Formal, Casual, Summary

**For Tasks 96-120:**

Number of variants created per task - 1

Tools used -

sentence\_reordering  
style\_paraphraser  
concat\_monolingual  
styleformer

New instances were generated for the tasks -

102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 115, 116, 117, 118, 119, 119, 120

New instances couldn't be generated for the tasks due to lesser size of the training dataset -

96, 97, 98, 113, 114

Number of instances generated -

In the range 150 to 2500

Datasets used -

CoNala  
commongen  
Fact2Story  
semeval 2019 task10  
story cloze and ROCStories  
scruples  
splash  
context abuse detection  
SMS spam collection v.1  
Logic2Text  
Synthetic dataset  
Help dataset from EMNLP 2020 paper  
SPL  
zest  
asset  
Com2Sense

Types of Definitions Generated-

Paraphrase, Long, Short, Formal, Casual, Summary