# DARING DUCKS
## Varsha Alangar, Harshitha Venkata

# Question Answering System

- Uses NLP to process a question and extract it's answer from a story or text corpus.
- Significant because the ability to answer questions about a story is the hallmark to text understanding.

**TECHNIQUES WE EXPLORED:**

1) Rule-based methods
2) Machine Learning

**IMPLEMENTED RULE-BASED APPROACH**

**OUR CHALLENGES:**

1) Extracting syntactic and semantic understanding of the text
2) Understanding the question-answer patterns

Set of Questions

WORD MATCH:
1) Remove stop words
2) Stem verbs in question and story sentence – give score of 12
3) For each match word match in story and question – give score of 4

Story

Why Rules | Where Rules | When Rules | What Rules | Who Rules | How Rules

Sentences Scored

Get maximum scored sentence

ANSWER PROCESSING:
1) Prune "where" rules
2) Prune "when" rules
3) Prune "what time" rules

Display the answers

| CORPUS | F-SCORE |
|---|---|
| Develop set | 0.3069 |
| Test Set 1 | 0.3968 |
| Test Set 2 | 0.4299 |

**STEP 1:** Word match heuristic scores each sentence in the story when the verb stems match or when there is an exact word to word match.
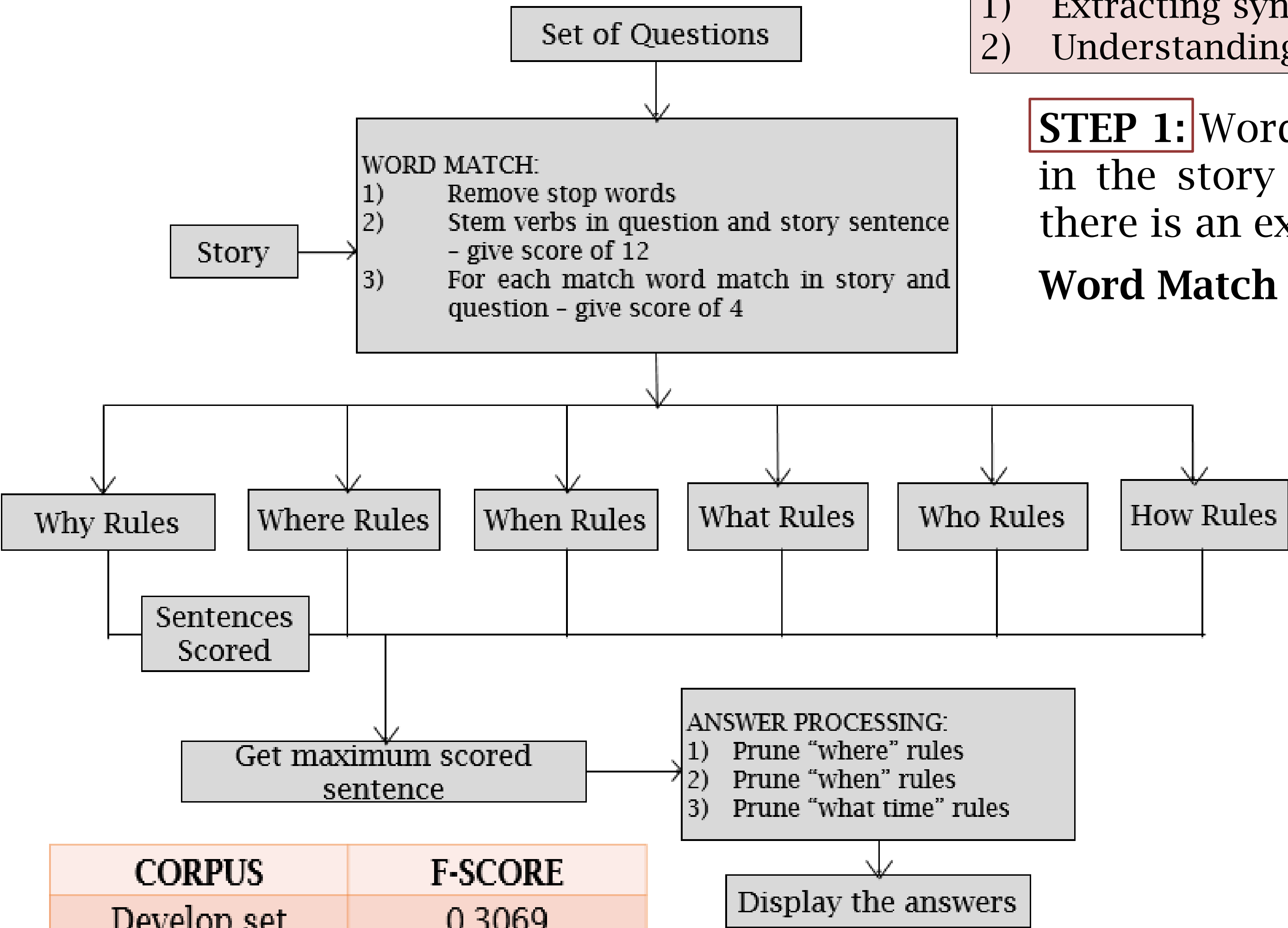
**Word Match : stemming -> score +=12**

STEMMING APPROACHES WE USED:

**Lancaster stemmer (NLTK version):** (most accurate stemmer) - re-write rule interpreter which may be configured as a FSA by using appropriate rules.

**Look up stemming:**
Collection of 5000 verbs in English language with the root verbs and it's derivatives acting as a look-up.
**Works better for our system than the Lancaster and gave a significant improvement in accuracy.**

Other stemmers explored (ineffective):
1) Porter stemmer
2) WordNet lemmatization

**PRONOUN RESOLUTION:**

**Rule:** If sentence has a pronoun and previous sentence has a noun → name or location. Score += 6

**STEP 2:** APPLYING RULES:

| RULE | CONDITION | SCORE |
|---|---|---|
| WHO,WHOM, WHOSE: #1 | Sentence has a name or has a person reference | 6 |
| #2 | Question does not contain a name or occupation Sentence does. | 4 |
| #3 | Word "name" in sentence and a name in sentence | 4 |
| #4 | Sentence has both a name and an occupation | 20 |
| HOW MUCH: #1 | Question has enquiring cost words, answer has words describing currency | 10 |
| HOW MANY #2 | question has the words "how many" and sentence has a number | 6 |
| HOW TALL HOW BIG #3 | Question has words "how tall" or "how big" and sentence has words that describe size. | 6 |

| RULE | CONDITION | SCORE |
|---|---|---|
| HOW OLD: #4 | Question has the words "how old" and sentence has a number and words that describe age. | 10 |
| WHAT: #1 | Question contains "what time" and sentence has words describing time. | 15 |
| #2 | Sentence has "known", called" and question has the words " name", "named". | 10 |
| #3 | Question has name followed by a preposition and sentence has a name or occupation. | 10 |
| WHEN | Sentence has words describing time. | 6 |
| WHERE | Sentence has words that indicate a location and question has location prepositions. | 6 |
| WHY | Sentence has words like "because" ," due to" and sentence and word have maximum word match. | 10 |

**STEP 3:** PRUNING ANSWERS:

1) We choose the answer with BEST SCORE from the rules and prune it, increasing the precision greatly.
2) Common words between the sentence and question, punctuations are removed.
3) If question is a "when" or a "what time" question, then displaying just the time, month, etc. will suffice.
4) If the question is a "where" question, we display only the location words in the sentence.

**THINGS THAT DIDN'T WORK:**

Removing stop words from answers reduced accuracy. Stop words seem significant in the answers.
Anaphora Resolution Algorithms in NLTK like Discourse Resolution Theory (DRT), gave faulty results.
Lack of relevant corpus for Machine Learning Techniques to resolve answers from the question patterns.

Riloff, E. and Thelen, M., (2000) "A Rule-based Question Answering System for Reading Comprehension Tests", *ANLP/NAACL-2000 Workshop on Reading Comprehension Tests as Evaluation for Computer-Based Language Understanding Systems.*