# Socio-Technical Solution to Large-Scale Formal Specification Mining

## May1620

Alex Dana, Cody Hanika

Client/Adviser: Professor Hridesh Rajan

**Problem:**
- ❏ Formal specifications are not widely available to programmers
- ❏ Many programmers do not know how or why to use formal specifications
- ❏ It can be difficult to write new formal specifications without a lot of experience using them (even then it is quite time consuming)
- ❏ It is difficult to automatically generate formal specifications for the code one is writing

```
/*
*    This function returns a divided by b. An exception will
*    be thrown if b is 0.
*/
public static int divide(int a, int b) {
    return a/b;
}

/*@ public normal_behavior
@      requires b != 0;
@      ensures \result == a / b;
@ also
@   public exceptional_behavior
@      requires b == 0;
@      signals_only ArithmeticException
*/
public static int divide2(int a, int b) {
    return a/b;
}
```
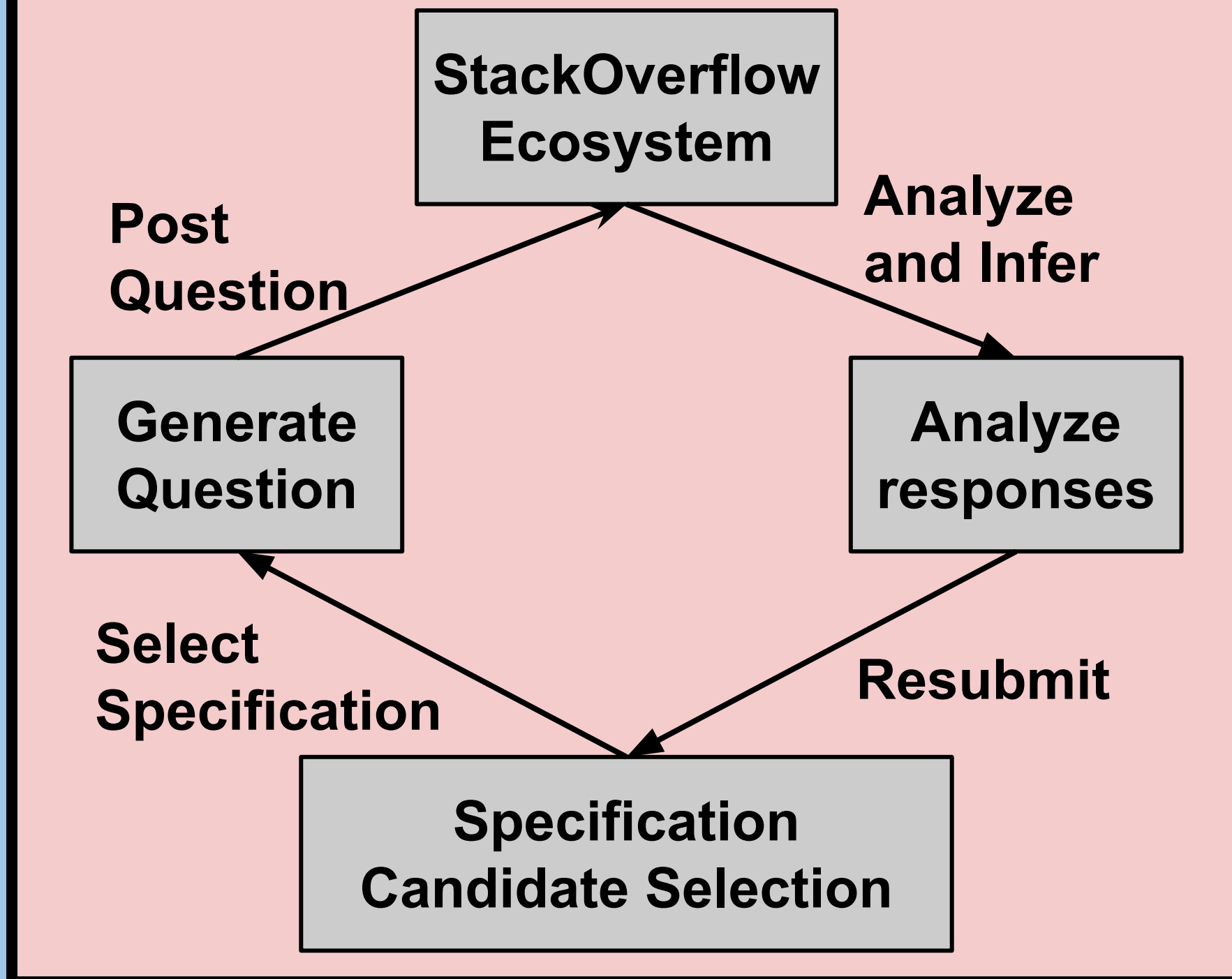*Informal vs. Formal JML Specifications*

**Operating Environment:**
Our main Java systems will be platform independent with our MySQL database running on Linux.
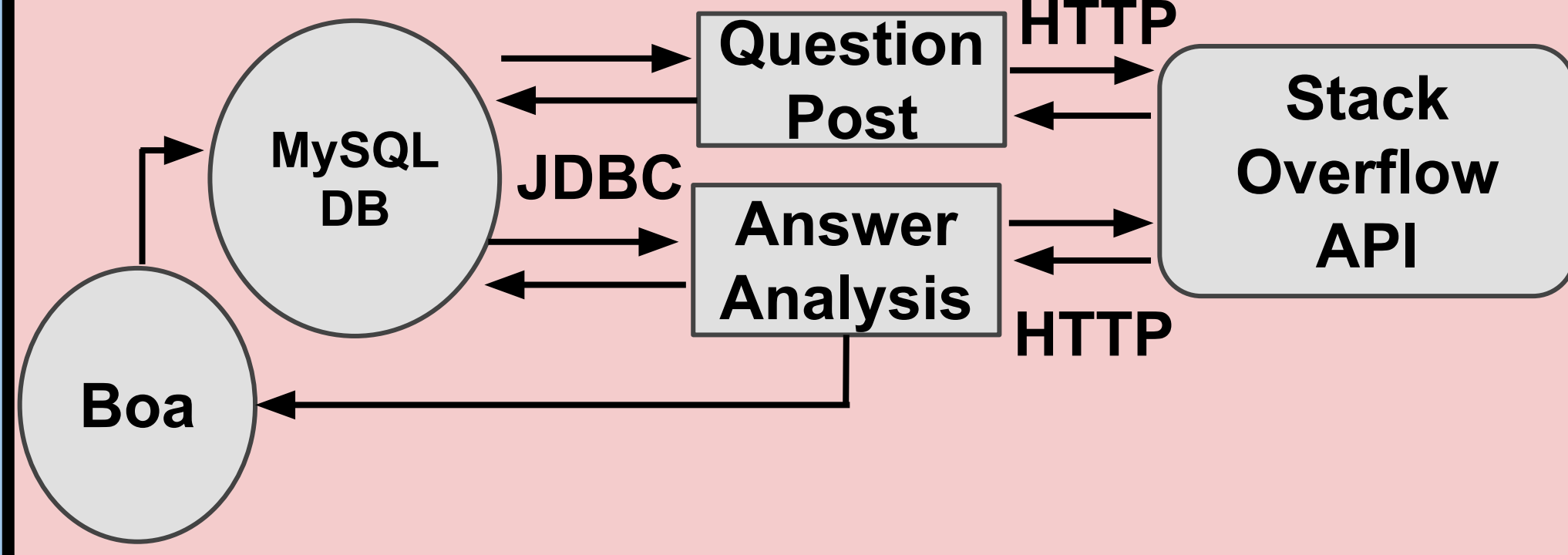
**Approach:**
- ❏ We tackled this problem using 2 unique Research Stages
  - ❏ **Stage 1** - Attempt to mine JML (Java Specification Language) specifications from StackOverflow questions
    - ❏ We generate these questions via http requests using the SO API
  - ❏ **Stage 2** - Use current javadoc to generate JML utilizing a natural language processor
    - ❏ The Stanford Natural Language Processor (SNLP) was chosen

---

## Stage 1

**Initial Concept**

StackOverflow Ecosystem — Generate Question — Specification Candidate Selection — Analyze responses
Post Question — Analyze and Infer — Select Specification — Resubmit

**System Overview**

MySQL DB — JDBC — Question Post — HTTP — Stack Overflow API
Boa — Answer Analysis — HTTP

**Question Post:** Retrieve unfinished spec from database (JDBC). Java http request to generate and post question via SO API.
**Answer Analysis:** Retrieve in progress spec. Java http request to retrieve question via SO API. Analyze using Java.
**Database:** MySQL database to store specs in need of improvement as well as specs currently posted.
**Boa System:** Large software repository knowledge base
**StackOverflow API:** StackOverflow provides a robust API for posting, answering, and retrieving questions.

**Functional Requirements:**
- ❏ Generate valid JML specifications
- ❏ Automatically generate SO questions

**Non Functional Requirements:**
- ❏ Run continuously
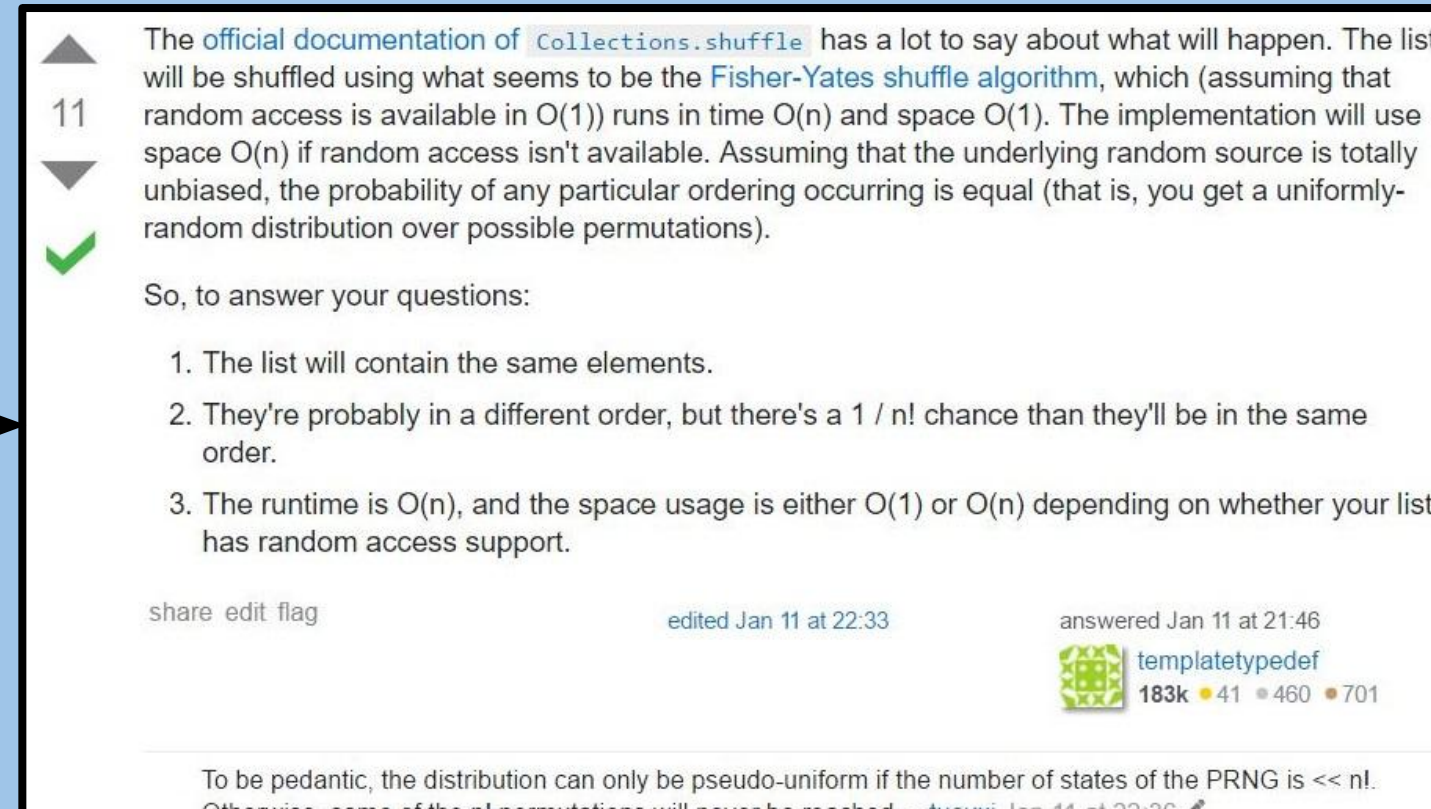- ❏ Little/No user input

**Test Plan:**
JUnit used to ensure the components of our system interfaced properly. Manual tests to validate http request success.

**Results:**
Research concluded that stage 1 was not a viable solution to the problem. A general lack of quality answers on StackOverflow led us to abandon this approach in favor of stage 2.
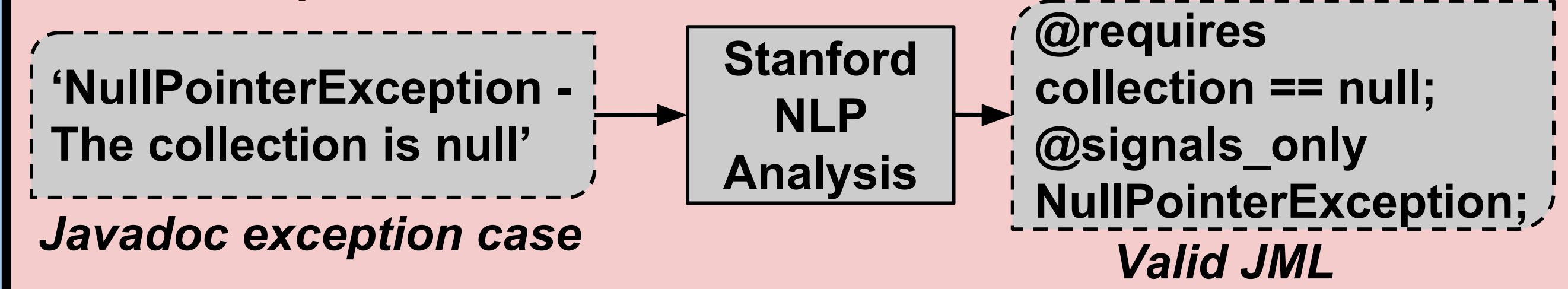
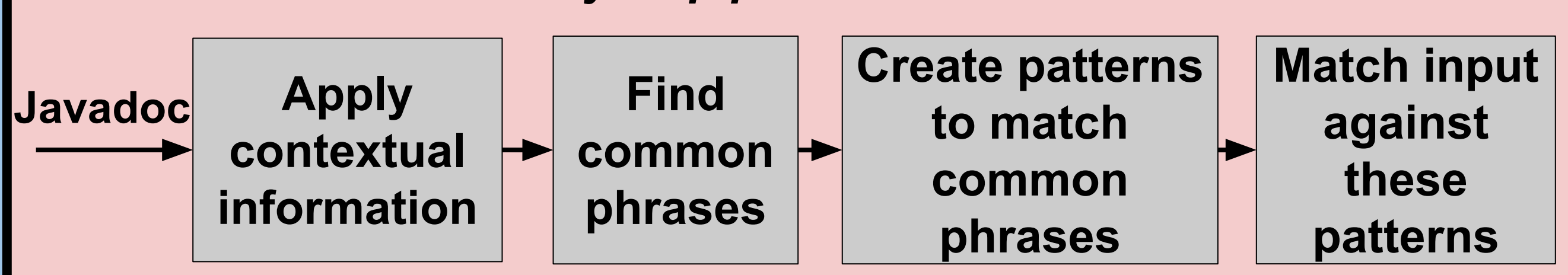*Question asked regarding the 'shuffle' method*

*Valid but non-structured answer*

---

## Stage 2

**Initial Concept**

'NullPointerException - The collection is null' → Stanford NLP Analysis → @requires collection == null; @signals_only NullPointerException;
*Javadoc exception case* — *Valid JML*

**Java documentation analysis pipeline**

Javadoc → Apply contextual information → Find common phrases → Create patterns to match common phrases → Match input against these patterns

**Functional Requirements:**
- ❏ Generate valid JML specifications
- ❏ Obtain javadoc by parsing source code

**Non Functional Requirements:**
- ❏ Analysis runs in a 'reasonable' amount of time (less than 2 minutes)

**Our Process:**
- ❏ Parse javadoc of jdk
- ❏ Create a parse tree via the analysis pipeline
- ❏ Extract frequent parse trees to discover common informal specifications
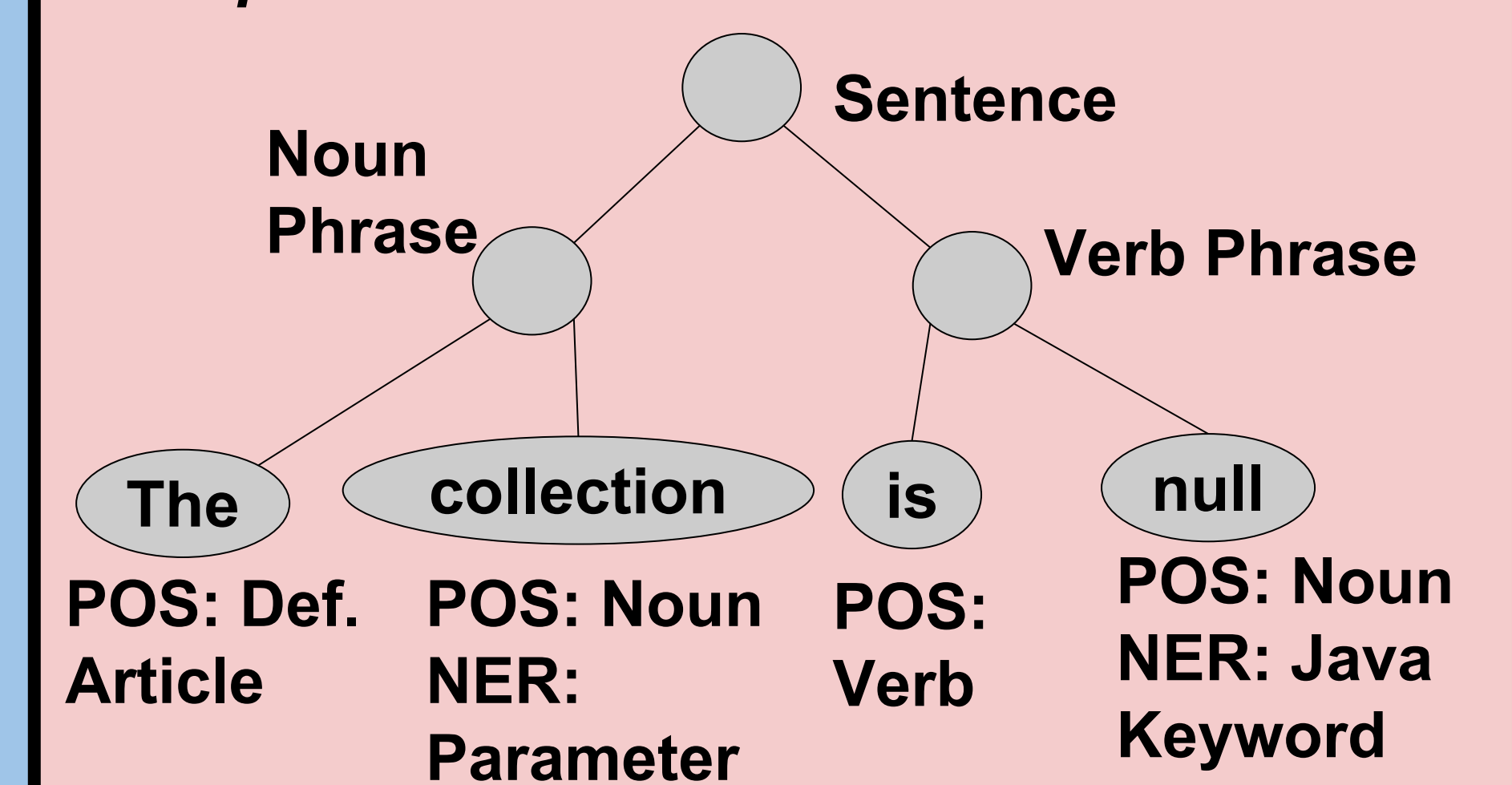- ❏ Create 'tokensregex' rule files to generate JML from these informal specifications

**Parse Tree:**
We create a parse tree of the javadoc using SNLP. All subtrees are hashed and then we look for common tree structures.

**Tokensregex Rule Files:**
We create 'rules' using regular expressions to tell SNLP how to handle certain groups of words, POS, NERs, etc. which we use to generate JML.

**SNLP parse tree of "The collection is null."**

Sentence — Noun Phrase — Verb Phrase
The (POS: Def. Article) — collection (POS: Noun NER: Parameter) — is (POS: Verb) — null (POS: Noun NER: Java Keyword)

**Test Plan:**
Manually ran analysis over Java's standard library to identify percentage of methods we found matches for. Then verify how many of them were valid.

**Results:**
Using SNLP we were able to successfully generate JML for certain types of Exceptions, such as NullPointerException. We had less success with documentation that didn't describe a boolean expression.