

Language Correction



Group 6

Alex Melnick
Michael Harkess
Reza Sajjadinasab

Background

Big Picture

User Input

- The user input can be in form of a file text or a sentence.
- It can be any natural language or programming language.

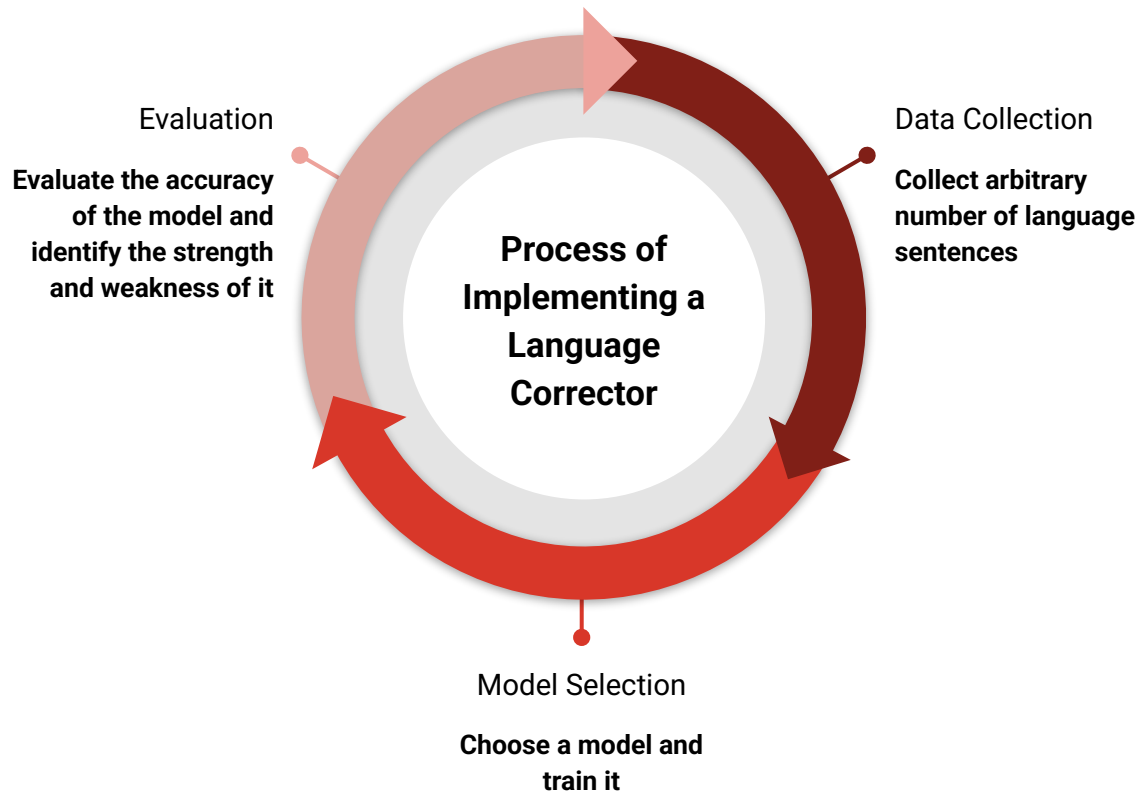
Black Box



Corrected Output

- identifying suspicious phrases
- Assign a confidence score
- Provide suggestions on how to correct them
- Correct them

What We Need?



Minimum Requirements

1

Crawler

- Crawls through English-language web pages
- A Command-line User Interface

2

Checker

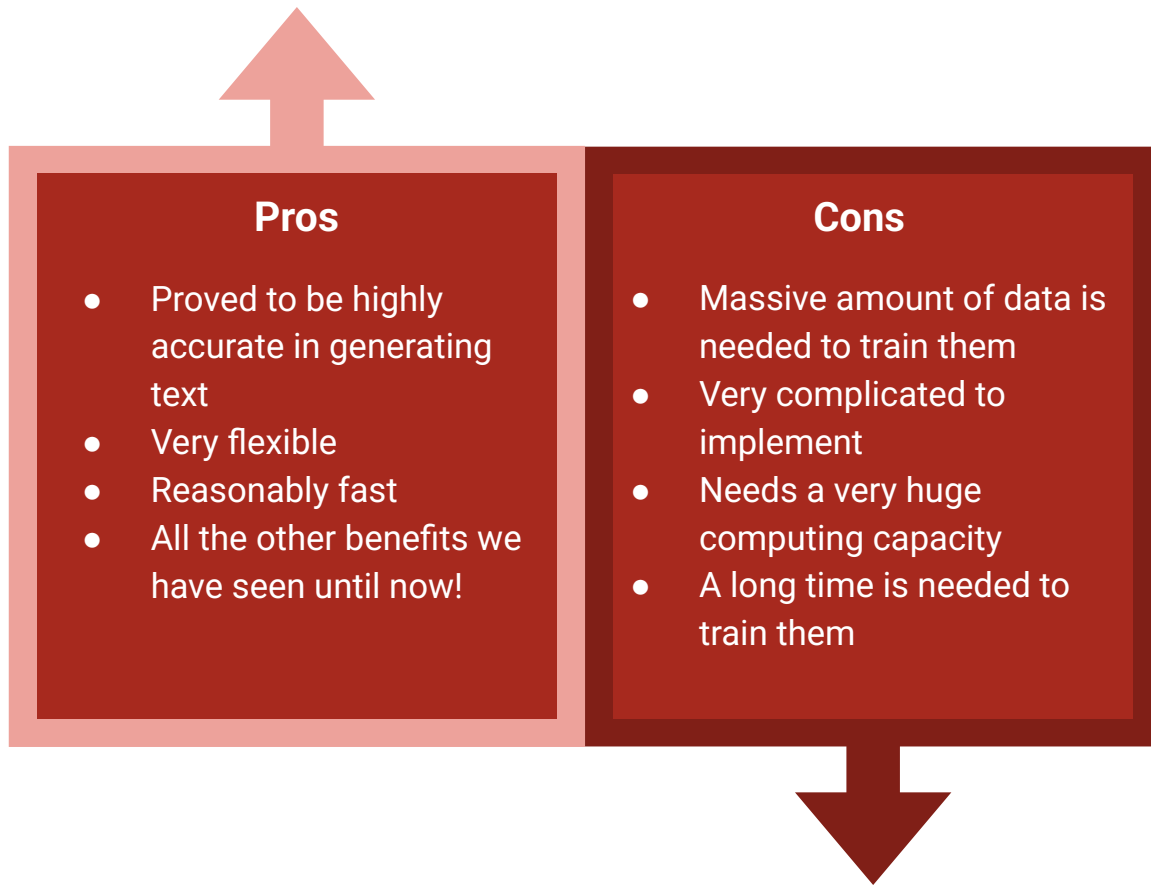
- Assign confidence to the various language structures
- Provide convincing experimental evidence
- A Command-line User Interface

3

Corrector

- Provides a reasonable correction to a provided sentence

Black Box? LLMs?



Features in FosWiki Description

Implemented

- **Real-time status** and statistics feedback for the crawler
- List of **reasonable corrections** to a suspicious text
- Graphical User Interface that **highlights suspicious** and non-suspicious textual elements
- **Crawling social media** posts of some large network
- **Graphical human feedback system** for deciding among possible phrase corrections
- **Android client** for your checker.
- Extension the system to two other language (**Turkish and Dutch**)
- **Translation** from English to another language

Not Implemented

- Extend your system to **computer source code** in some well-used language (e.g., C++, Java).

Design

Our Challenges

- Two people from our group dropped the course
- Computing capacity and time limitations
- The goal of the project is a fundamental problem
- The requirement of using Java



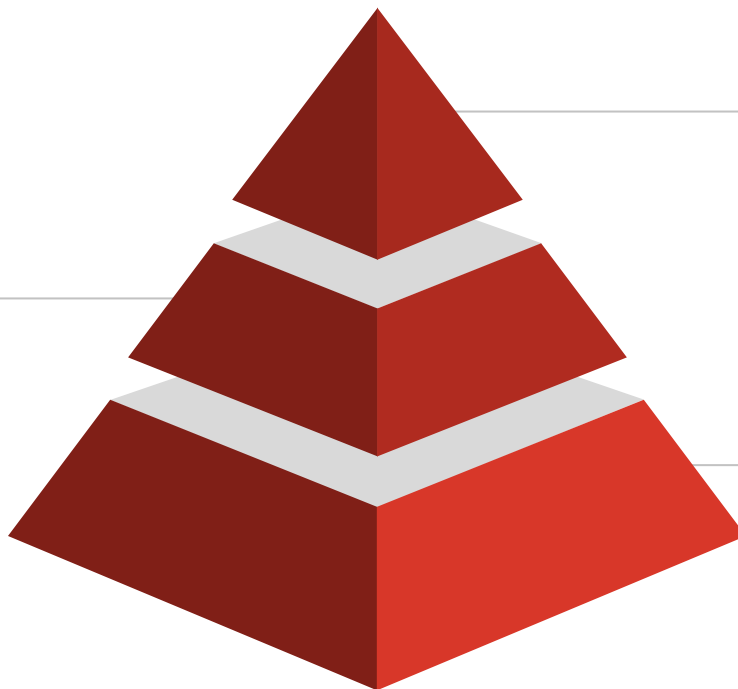
Crawler

- Major focus on **obtaining natural language** and future links-
 - Excluded from crawling all direct links to media (images, videos, music, etc.)
 - Heavily restricted by default size requirements
 - The default 1KB maximum data crawled per page is approximately 200 words
- Aggressive parsing to separate natural language and links to external web-pages from crawled data

Checker/Corrector

Sentence Structure Checker and Correction

Using a **state machine** to see if a sentence follow a valid structure



Typo Correction

1 Finding **distance** of a given word with a words in database.

N-Grams Probability and Similarity Correction

3 Compared phrase of a sentence with seen phrases stored in **hash table** in a SQL database

Requirements of Checker/Corrector

Typo Corrector

Needs to at least see a word once

State Machine

Needs to know the role of the majority the words in a sentence

N-Grams & Similarity

Needs to see combination of different usage of a word

Pros & Cons of Checker/Corrector



Pros

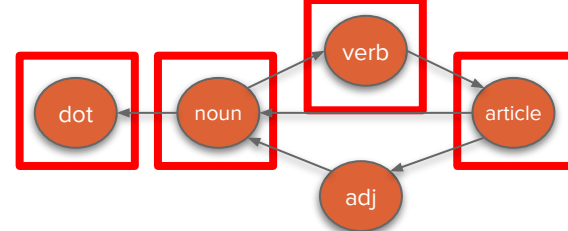
Typo corrector	Relatively easy to implement and fast in computation
State Machine	Effective even without having knowing all of the roles
N-grams & Similarity	Only need several phrases to get trained

Cons

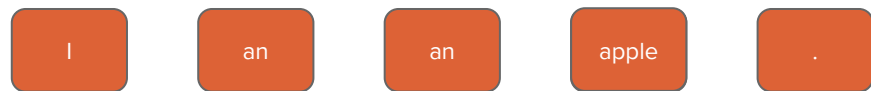
Typo corrector	May forced a wrong decision
State Machine	May significantly change the structure of the sentence
N-grams & Similarity	Requires a massive amount of data to be effective



Example of Correction Flow



Typo Corrector



State Machine



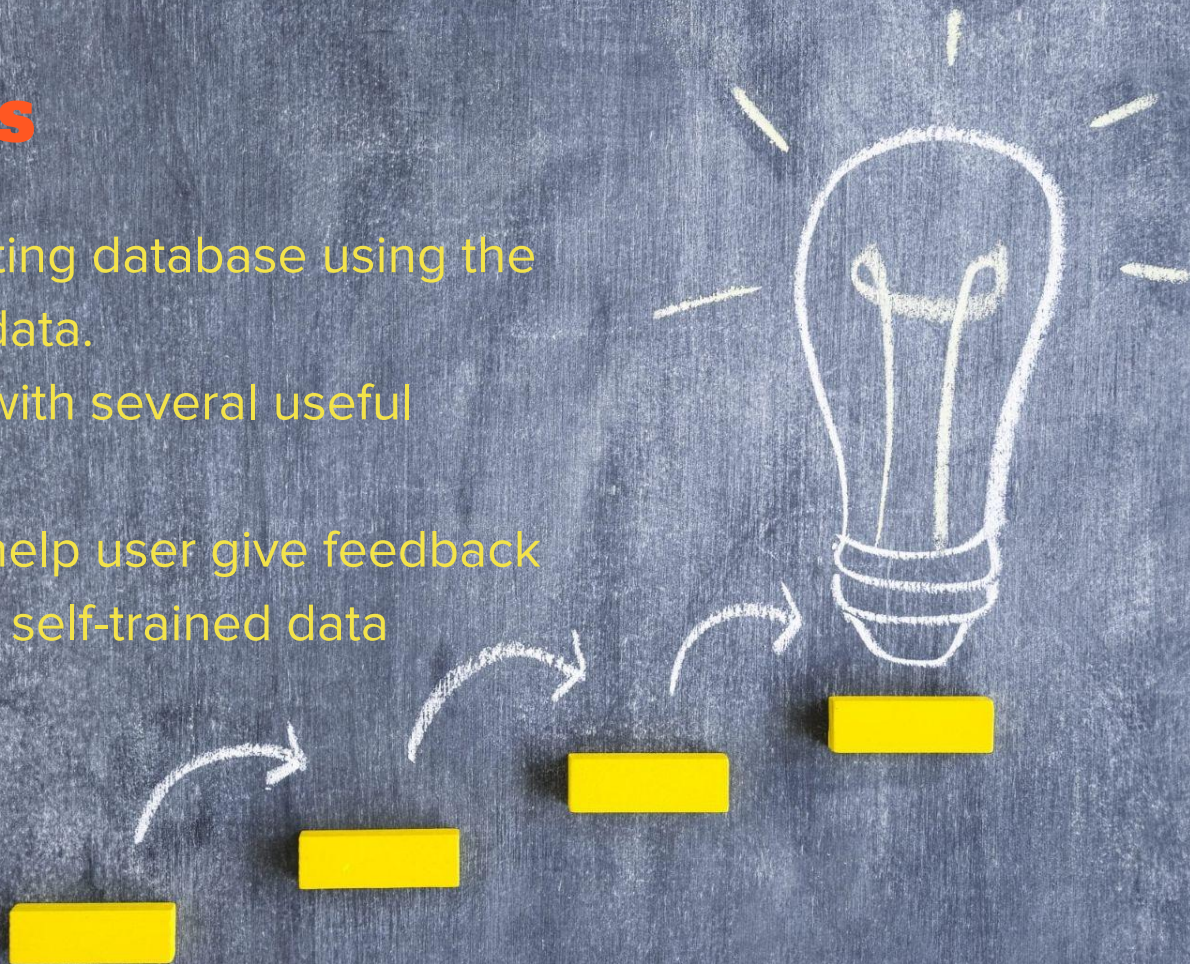
N-Grams &
Similarity

Look up in the database for a better
suggestion that seen more



Other Features

- Four ways of updating database using the provided crawled data.
- Sophisticated CLI with several useful options
- Interactive GUI to help user give feedback to model based on self-trained data



Analysis

Crawler - Time & Space Complexity

- Data structure choice was **driven by speed and uniqueness**
 - Hash tables were the obvious choice
 - $\Theta(1)$ search, insert, and delete; $\Theta(n)$ space
 - $O(n)$ search, insert, and delete; $O(n)$ space
 - ArrayLists were used to parse the crawled data
 - $\Theta(n)$ search and delete, $\Theta(1)$ insert; $\Theta(n)$ space
 - $O(n)$ search and delete, $O(1)$ insert; $O(n)$ space

Checker/Corrector - Time & Space Complexity

Typo Corrector

- $O(s)$
 - s is the size of database

* Space for all of the data structure is constant without considering the database,

Structure Corrector

- Checker $\rightarrow O(k \log(s))$
 - k : Number of words a
 - s : Size of database
- Corrector $\rightarrow O(N^k)$! \rightarrow linear behavior!
 - N : Total number of states/tokens
 - k : Number of words

Similarity Corrector

- $O(k \log(s))$
 - K : Number of words
 - s : Size of database

Comparison

Other groups

None of our group member where a reviewer for them, so as of now we don't know what are they doing

Grammarly

The tool is not open source
It can detect most of the grammar errors and typos and suggest some better language structure in the premium version **only for English**

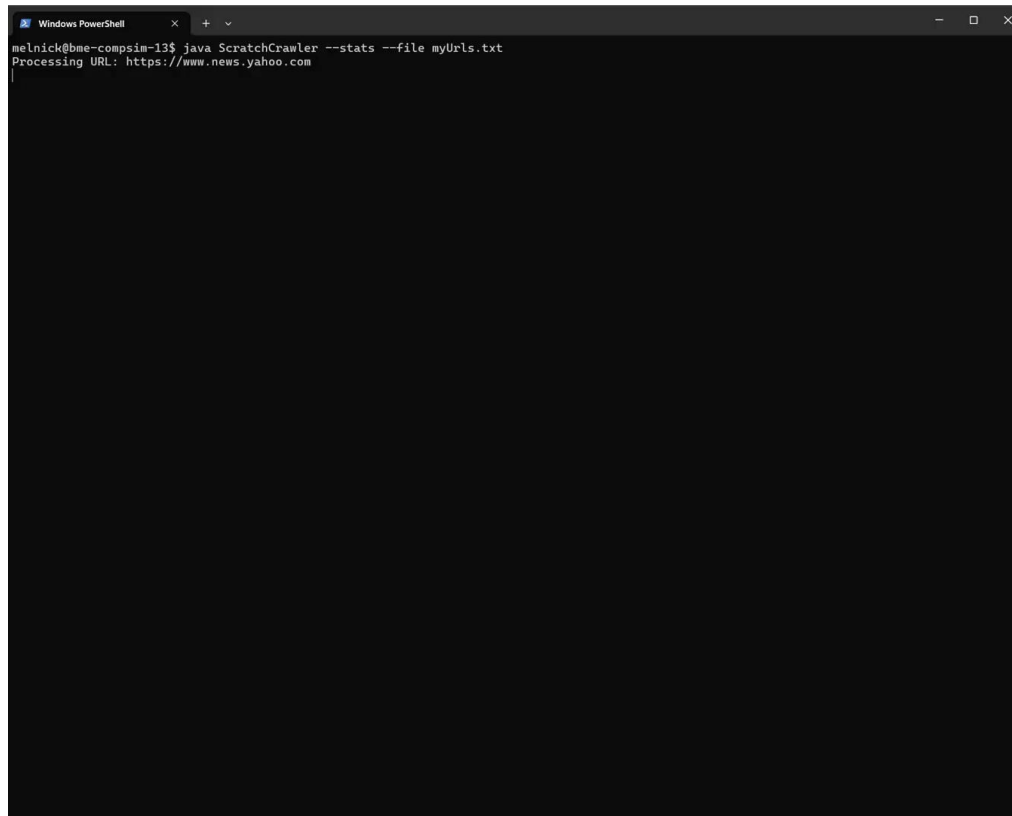
LLMs

Almost all of the existing LLMs are able to **detect language errors** and correct them with a **acceptable accuracy** for **different languages**.

This tools are using **transformers** and complicated language model which **requires massive amount of data** and training.

Demonstration

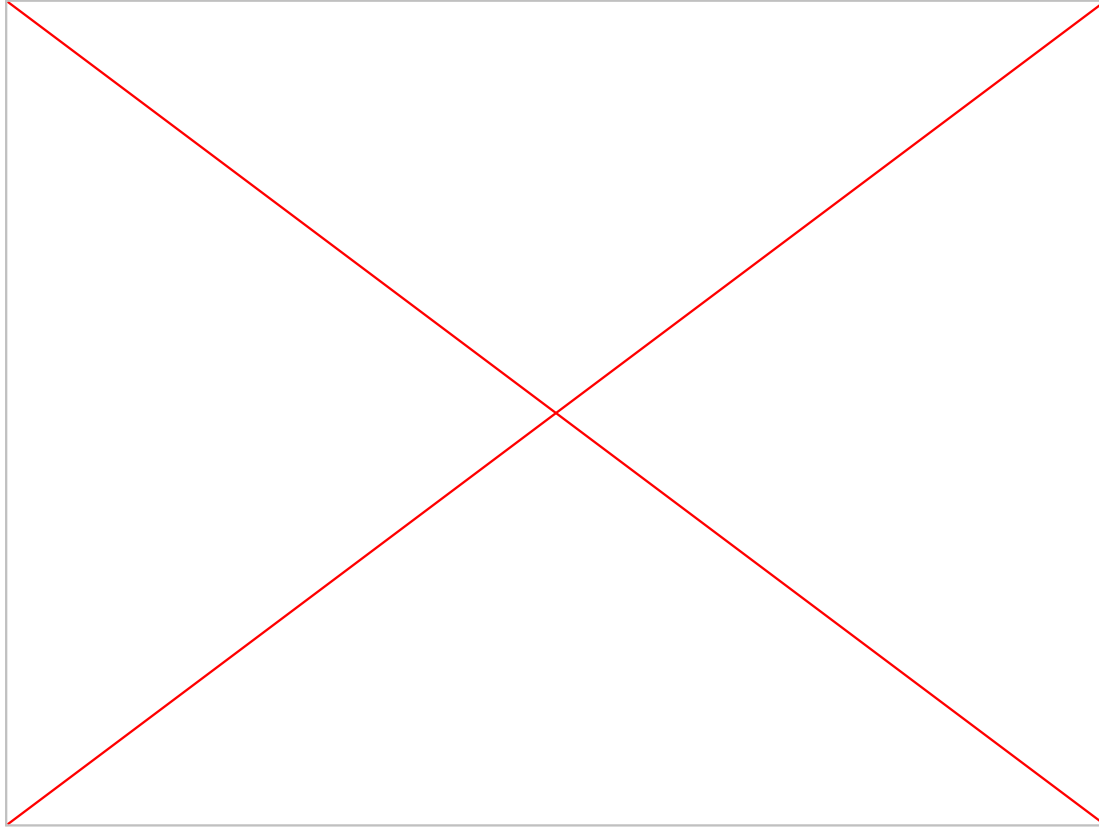
Crawler



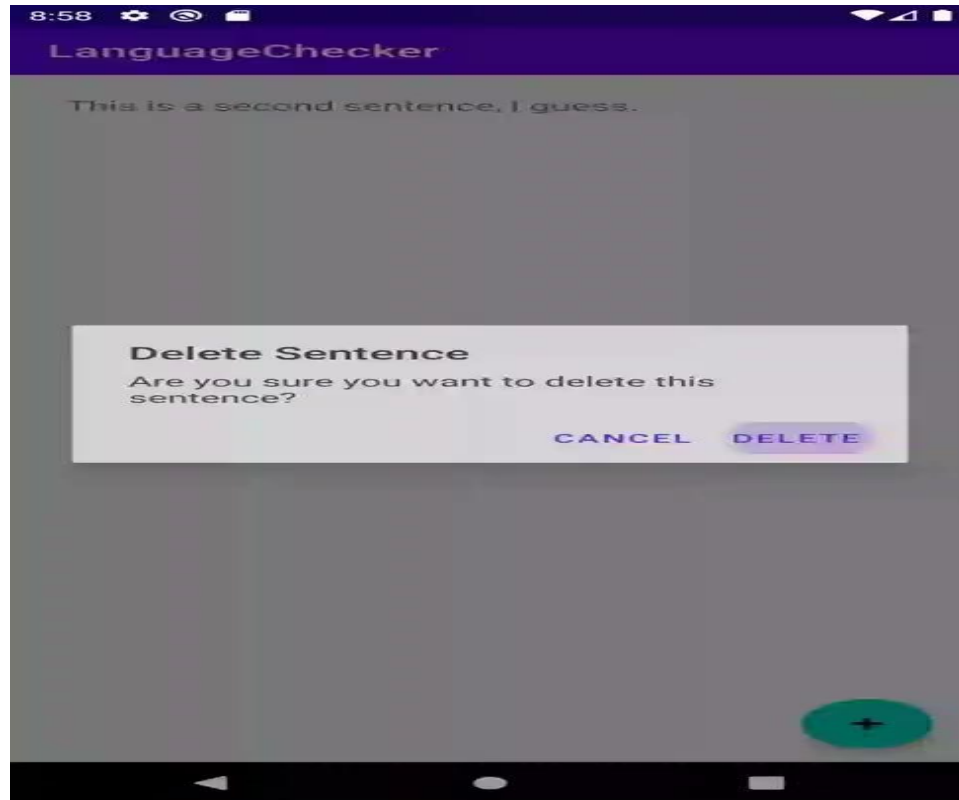
A screenshot of a Windows PowerShell terminal window. The window has a title bar that says "Windows PowerShell" and standard window controls (minimize, maximize, close). The terminal content shows a user prompt "melnick@bme-compsim-13\$" followed by the command "java ScratchCrawler --stats --file myUrls.txt". Below the command, the output "Processing URL: https://www.news.yahoo.com" is displayed. The rest of the terminal is black with a white cursor line at the end of the first line.

```
melnick@bme-compsim-13$ java ScratchCrawler --stats --file myUrls.txt
Processing URL: https://www.news.yahoo.com
```

Checker/Corrector



Android App



Questions

**Thanks for your
time!**