

Vũ Quang Minh

18110150

MÔN HỌC: PHÁT TRIỂN PHẦN MỀM HƯỚNG ĐỐI TƯỢNG

Đề 3. Thiết kế và cài đặt một ứng dụng cho phép phân tích một tập tin văn bản, xuất ra tần số xuất hiện của các từ có trong tập tin văn bản. Các ký tự ngăn cách bao gồm khoảng trắng chuẩn và các ký tự ngăn cách câu (dấu ‘,’ ‘.’ ‘;’ ‘:’ ‘?’ ‘!’ ‘...’). Từ là từ có nghĩa trong từ điển từ tiếng Anh, có xét đến số ít, số nhiều, quá khứ (ví dụ: mouse, mice được tính là cùng một từ, study, studies, studied, studying được tính là cùng một từ).

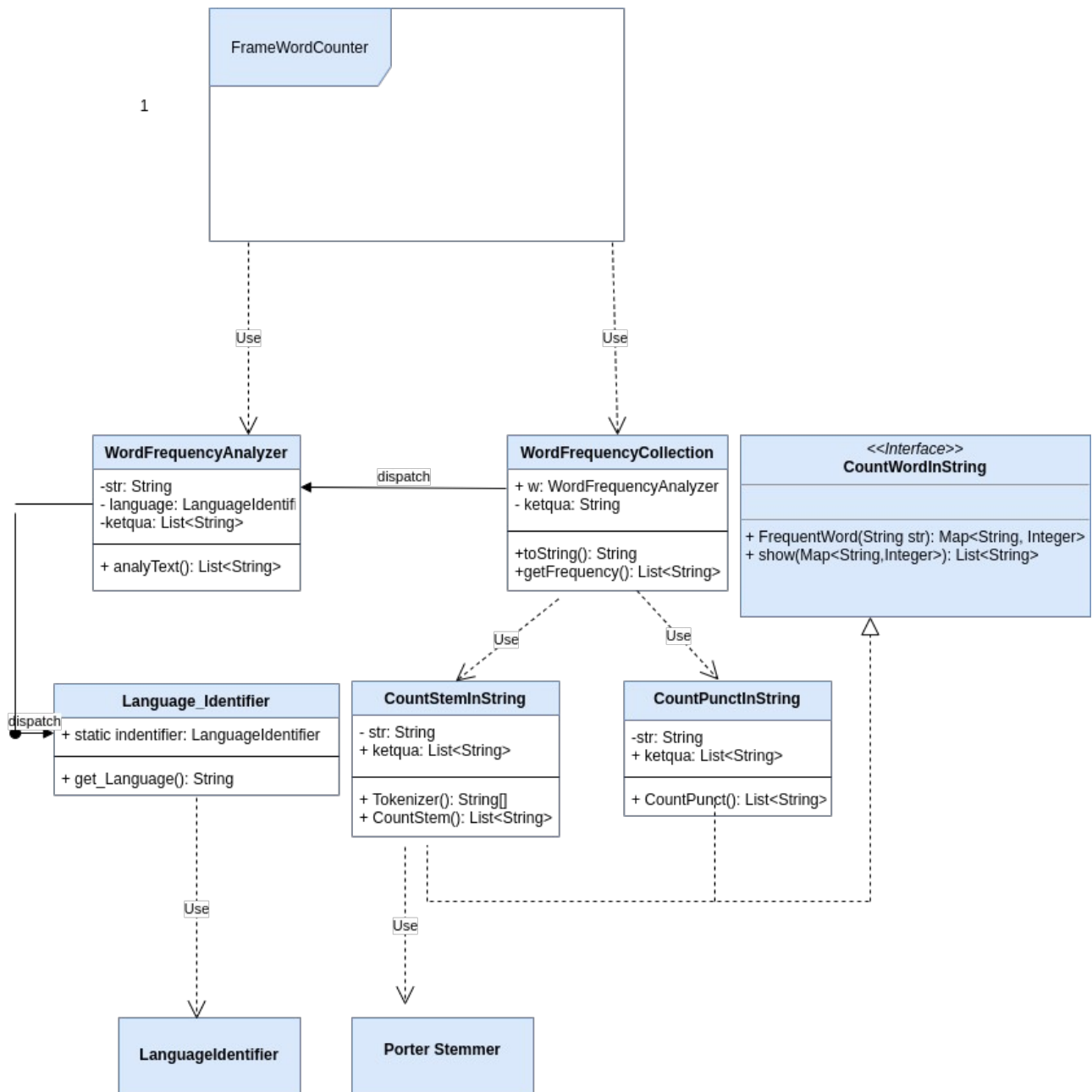
A) Specification

- Người dùng sẽ chọn **FrameWordCounter** và nhấn nút Run As Application để chạy chương trình.
- Chương trình sẽ cho mở file bất kì ở vị trí nào trên bộ nhớ bằng nút(Browse) nhưng chỉ đọc được File txt nếu gặp file định dạng khác thì không đọc được.
- Chương trình cũng có phần xuất file cũng được lưu ở vị trí bất kì trên bộ nhớ chỉ lưu được file dưới dạng txt. Nếu lưu dưới dạng khác thì sẽ không chạy được file.
- Sau khi mở file thì đường link sẽ được lưu trên ô Open File, và nội dung sẽ xuất hiện trên ô xuất nội dung.
- Nhấn nút **Xuất kết quả** để xử lý chương trình.
- Chương trình định dạng được ngôn ngữ tiếng Anh.
- Chương trình có thể đếm được các ký tự(.,?;!...)
- Chương trình sẽ liệt kê tất cả các từ và các loại từ của chúng(ở đây là số ít, số nhiều, quá khứ, hiện tại, tương lai).
- Chương trình sẽ đếm các loại từ.

B) Phạm vi đề tài

- Chương trình định dạng được ngôn ngữ tiếng Anh.
- Chương trình có thể đếm được các ký tự(.,?;!...).
- Chương trình sẽ liệt kê tất cả các từ và các loại từ của chúng(ở đây là số ít, số nhiều, quá khứ, hiện tại, tương lai).
- Chương trình sẽ đếm các loại từ.
- Chương trình có phần giao diện từ được hỗ trợ bởi java1.8 và WindowBuilder được tải từ Marketplace.
- Chương trình có 3 thư viện được thêm vào *org.apache.tika* để định dạng ngôn ngữ, *opennlp-uima-1.9.4.jar* và *opennlp-tools-1.9.4.jar* để xử lý các từ cùng loại.

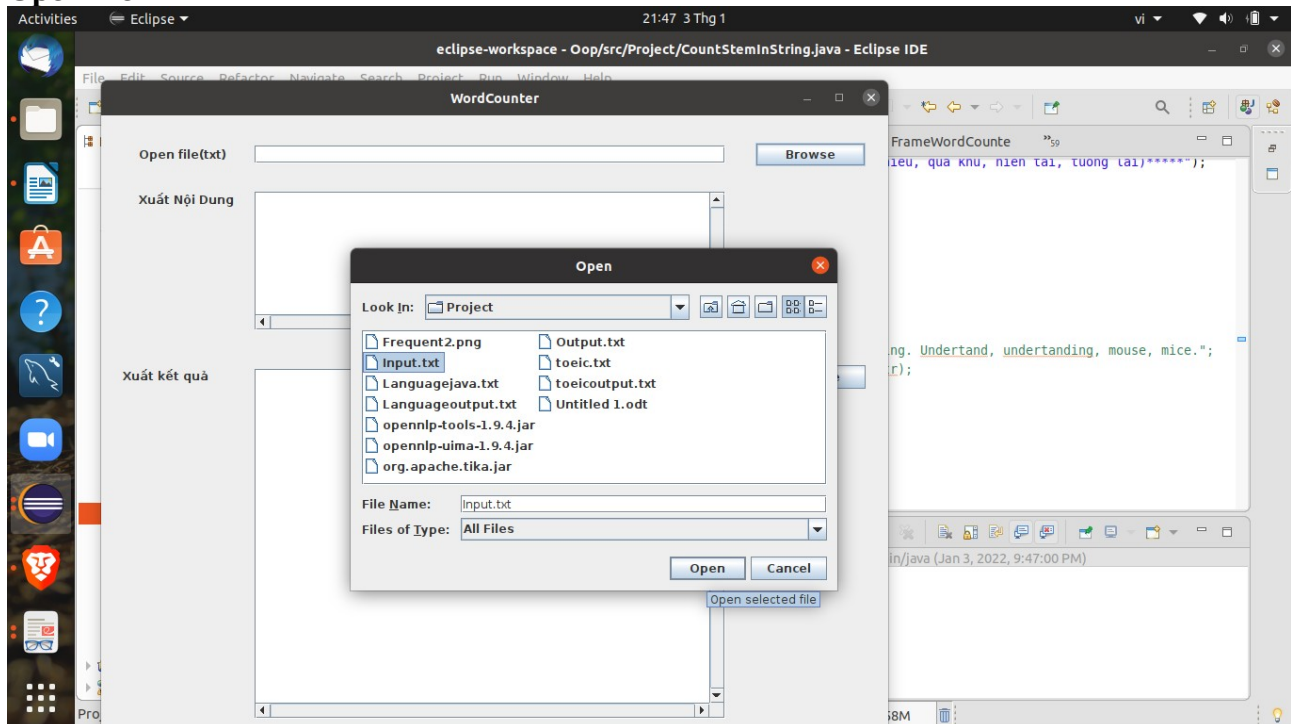
C) Sơ đồ và Code



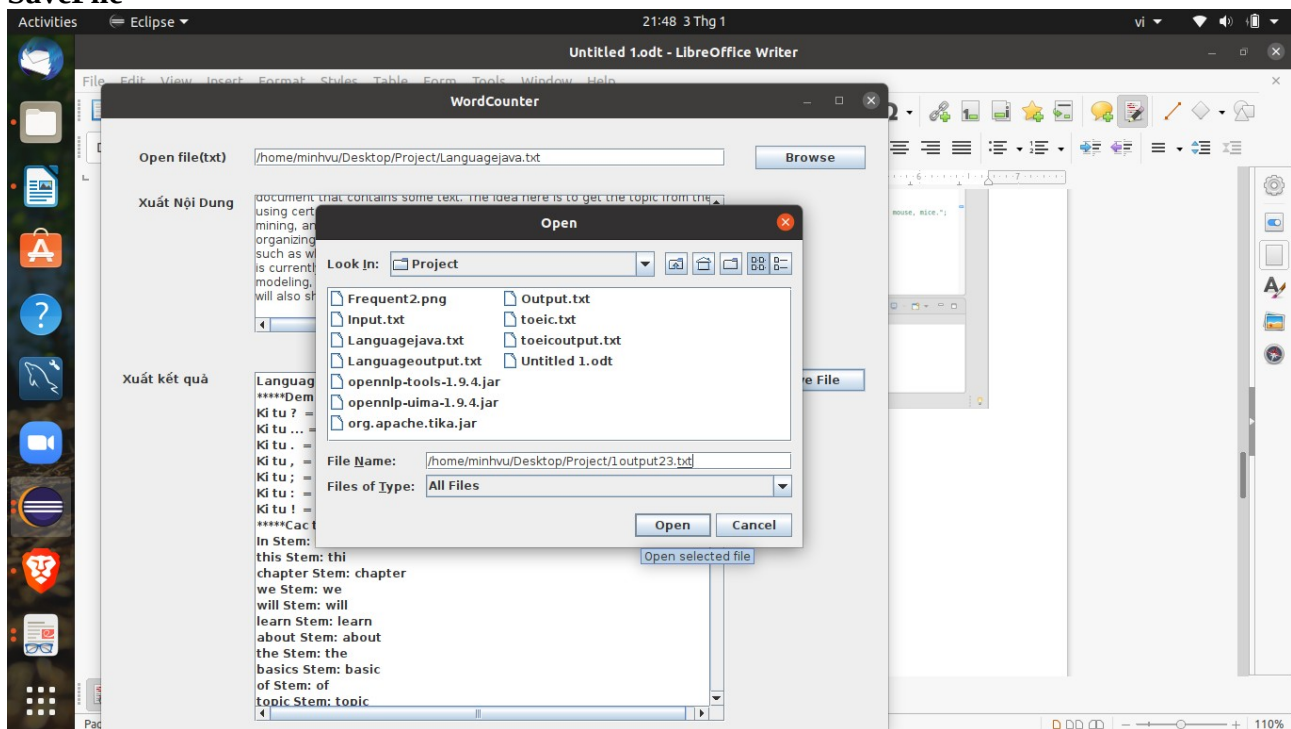
- Người dùng sẽ chọn *FrameWordCounter* để chạy chương trình.
- *FrameWordCounter* sẽ sử dụng *WordFrequencyCollection* để đếm kí tự và các từ cùng loại và *WordFrequencyAnalyzer* để định dạng ngôn ngữ.
- *WordFrequencyCollection* sẽ sử dụng *CountStemInString* để đếm các từ cùng loại và *CountPunctInString* để đếm các kí tự. *CountStemInString* sẽ sử dụng *PorterStemmer* để xử lý các từ cùng loại và được hỗ trợ bởi thư viện *opennlp*.
- Hai hàm *FrequentWord(String str)* và *show(Map<String, Integer>)* của *CountStemInString* và *CountPunctInString* được implement từ interface *CountWordInString*.
- *WordFrequencyCollection* sẽ được tham chiếu từ *WordFrequencyAnalyze*.
- *WordFrequencyAnalyze* khi khởi tạo sẽ tham chiếu đến *Language_Identifier* để định dạng ngôn ngữ.

D) Demo và hướng dẫn sử dụng

OpenFile



SaveFile



Demo vài đoạn chương trình

