

## └─ Nội dung

Nội dung mà em trình bày gồm n phần. Đầu tiên, em xin được giới thiệu vấn đề.

## Báo cáo Đồ án Môn học

- └ Đặt vấn đề và giới thiệu bài toán
  - └ Đặt vấn đề
    - └ Đặt vấn đề

- Là một bài toán có nhiều ứng dụng thực tế
  - phân loại và đánh giá trạng thái
  - lọc thư rác, tin nhắn rác
  - phân loại sách báo
  - ...
- Chưa có nhiều công cụ, thư viện hỗ trợ phân loại văn bản Tiếng Việt

Trong thời đại bùng nổ thông tin hiện nay, những hướng nghiên cứu về nhận dạng và xử lý văn bản là rất cần thiết, là bước tiền xử lý cho nhiều quá trình để trích xuất thông tin có ích, đặc biệt là các văn bản Tiếng Việt. Tuy nhiên, Tiếng Việt lại chưa nhận được sự hỗ trợ nhiều như Tiếng Anh. Do vậy, nhóm em quyết định thực hiện đề tài này.

# Báo cáo Đồ án Môn học

└ Đặt vấn đề và giới thiệu bài toán

└ Phát biểu bài toán

└ Đặc điểm của đề tài

## Đặc điểm của đề tài

- Vấn đề tách từ cho Tiếng Việt
  - Xâu đầu vào: "Đây là một ví dụ đơn giản minh họa cho việc sử dụng công cụ Đồng Du để tách từ."
  - Xâu đầu ra: "Đây là một ví dụ đơn giản minh họa cho việc sử dụng công cụ Đồng Du để tách từ."
- Vấn đề trích ra tập các từ khóa
  - Xâu đầu vào: "Việt Nam tài khố, định chủ quyền, với hai quần, đảo Hoàng Sa và Trường Sa, sau khi Chủ tịch Trung Quốc nói rằng các quần đảo này là của Trung Quốc."
  - Các từ ý nghĩa thập: tài, đồ, với, và, sau, khi, là, ...
  - Xâu đầu ra: "Việt Nam chủ quyền Hoàng Sa Trường Sa, Chủ tịch Trung Quốc Trung Quốc"

Trong Tiếng Việt, dấu cách không được sử dụng như một kí hiệu phân tách từ, nó chỉ có tác dụng phân tách các tiếng với nhau. Vấn đề từ đơn, từ ghép.

## Báo cáo Đồ án Môn học

└ Định hướng giải pháp

└ Tách từ Tiếng Việt

└ Tách từ Tiếng Việt

- phương pháp ghép cực đại: sử dụng từ điển, đặt các từ vào câu sao cho phù hợp câu đó
- các phương pháp học máy như đồ thị hóa, hidden Markov, conditional random field, maximum entropy
  - phương pháp-Support Vector Machine (SVM) đi kèm pointwise => thư viện Dong Du (C/C++)
- So sánh dung lượng và vnTokenizer

| Tên cơ           | vnTokenizer | Dung lượng |
|------------------|-------------|------------|
| Dữ chính xác     | 97.5%       | 98.3%      |
| Thời gian (giây) | 194.672     | 26.2       |
| RAM (MB)         | 93.8        | 15.1       |

Bảng: So sánh dung lượng và vnTokenizer

ưu điểm: nhanh; nhược điểm: độ chính xác thấp, không xử lý được từ có trong từ điển Sử dụng bản sửa đổi của thư viện Đông Du để tiến hành tách từ trên Python