

TRƯỜNG ĐẠI HỌC BÁCH KHOA HÀ NỘI

ĐỒ ÁN TỐT NGHIỆP

Chuyển đổi cuộc gọi sang đoạn hội thoại tiếng Việt

Khuất Ngọc Sơn

son.kn204602@sis.hust.edu.vn

Ngành: Khoa học máy tính

Giảng viên hướng dẫn: PGS.TS. Lê Thanh Hương

Chữ kí GVHD

Khoa: Khoa học máy tính

Trường: Công nghệ Thông tin và Truyền thông

HÀ NỘI, 01/2025

LỜI CAM KẾT

Họ và tên sinh viên: Khuất Ngọc Sơn
Điện thoại liên lạc: 0379135639
Email: son.kn204602@sis.hust.edu.vn
Lớp: Khoa học máy tính 03-K65
Hệ đào tạo: Hệ Kỹ sư chính quy

Tôi – *Khuất Ngọc Sơn* – cam kết Đồ án Tốt nghiệp (ĐATN) là công trình nghiên cứu của bản thân tôi dưới sự hướng dẫn của *PGS.TS. Lê Thanh Hương*. Các kết quả nêu trong ĐATN là trung thực, là thành quả của riêng tôi, không sao chép theo bất kỳ công trình nào khác. Tất cả những tham khảo trong ĐATN – bao gồm hình ảnh, bảng biểu, số liệu, và các câu từ trích dẫn – đều được ghi rõ ràng và đầy đủ nguồn gốc trong danh mục tài liệu tham khảo. Tôi xin hoàn toàn chịu trách nhiệm với dù chỉ một sao chép vi phạm quy chế của nhà trường.

Hà Nội, ngày 06 tháng 01 năm 2025

Tác giả ĐATN

Sơn

Khuất Ngọc Sơn

LỜI CẢM ƠN

Đồ án này được em hoàn thành tại Trường Công nghệ Thông tin và Truyền thông, Đại học Bách Khoa Hà Nội dưới sự hướng dẫn của giáo viên hướng dẫn của PGS.TS. Lê Thanh Hương.

Lời cảm ơn đầu tiên em xin gửi tới PGS.TS. Lê Thanh Hương, người đã giúp đỡ và hướng dẫn em tận tình trong khoảng thời gian vừa qua để em có thể hoàn thành môn học cuối cùng này.

Tiếp theo, em xin được gửi lời cảm ơn chân thành đến đại học Bách Khoa Hà Nội, nơi chấp cánh ước mơ của bao thế hệ sinh viên và giúp chúng em nên người và trở thành những công dân tốt và có ích cho nước nhà.

Em xin được cảm ơn quý giảng viên đã và đang làm việc tại đây, đặc biệt là những thầy cô giảng viên trong hơn 4 năm qua đã bảo ban và truyền đạt cho em bao nhiêu kiến thức hay và bổ ích để em có thể phát triển bản thân mình.

Và cuối cùng không thể thiếu, em xin cảm ơn gia đình nơi điểm tựa tinh thần lớn lao để em bước tiếp trên con đường của chính mình. Cảm ơn bạn bè đã đồng hành cùng em trong những năm học Đại học và tương lai xa nữa.

Trong quá trình hoàn thành báo cáo và đồ án khó có thể tránh khỏi sai sót, em mong thầy cô và bạn đọc có thể góp ý để em hoàn thiện hơn nữa đồ án tốt nghiệp này.

TÓM TẮT NỘI DUNG ĐỒ ÁN

Trong thời đại công nghệ số phát triển mạnh mẽ, việc tự động hóa các quy trình trong doanh nghiệp ngày càng trở nên quan trọng. Một trong những lĩnh vực nổi bật là chuyển đổi dữ liệu từ các cuộc gọi telesale sang văn bản và trích xuất thông tin quan trọng nhằm tối ưu hóa hoạt động kinh doanh. Hiện nay, mặc dù có nhiều giải pháp xử lý âm thanh và nhận diện giọng nói, việc áp dụng chúng trong ngữ cảnh cụ thể như telesale vẫn còn nhiều thách thức. Các phương pháp truyền thống như ghi âm và phân tích thủ công không chỉ tốn thời gian mà còn dễ dẫn đến sai sót. Các giải pháp tự động, dù đã đạt được một số tiến bộ, vẫn gặp hạn chế về độ chính xác, đặc biệt trong việc nhận diện giọng nói vùng miền hoặc trong các đoạn hội thoại phức tạp.

Để giải quyết vấn đề này, nghiên cứu lựa chọn hướng tiếp cận bằng sử dụng các mô hình xử lý âm thanh tiên tiến như Whisper, Wav2Vec hay Whisperx. Lý do lựa chọn hướng này là khả năng mở rộng và tiềm năng cải thiện độ chính xác khi áp dụng vào các bài toán thực tiễn.

Giải pháp đề xuất bao gồm ba giai đoạn chính: (i) chuyển đổi âm thanh cuộc gọi sang văn bản bản thông qua các mô hình nhận diện giọng nói, (ii) hậu xử lý đoạn văn bản hội thoại và rút được văn bản tóm tắt, và (iii) sử dụng mô hình GPT để trích xuất thông tin quan trọng như tên khách hàng, sản phẩm quan tâm, và kết quả cuộc gọi. Hệ thống được thiết kế với khả năng tùy chỉnh theo nhu cầu của từng doanh nghiệp.

Đóng góp chính của đề tài là phát triển một quy trình tích hợp từ nhận diện giọng nói đến trích xuất thông tin, giúp tăng hiệu suất và giảm chi phí xử lý thủ công. Kết quả thử nghiệm cho thấy hệ thống đạt độ chính xác cao trong cả hai bước nhận diện và trích xuất, hứa hẹn mang lại giá trị thực tiễn cao trong lĩnh vực telesale và chăm sóc khách hàng.

Sinh viên thực hiện

ABSTRACT

In the era of rapidly advancing digital technology, automating business processes has become increasingly important. One prominent area is the transformation of telesales call data into text and the extraction of key information to optimize business operations. Although numerous solutions for audio processing and speech recognition are currently available, their application in specific contexts such as telesales still presents many challenges. Traditional methods, such as recording and manual analysis, are not only time-consuming but also prone to errors. While automated solutions have made some progress, they still face limitations in accuracy, especially when dealing with regional accents or complex conversations.

To address this issue, this study adopts an approach based on advanced audio processing models like Whisper, Wav2Vec, and Whisperx. These models are chosen for their scalability and potential to improve accuracy when applied to real-world problems.

The proposed solution comprises three main phases: (i) converting call audio into text using speech recognition models, (ii) post-processing the conversational text to generate summaries, and (iii) employing GPT models to extract critical information such as customer names, products of interest, and call outcomes. The system is designed to be customizable to meet the specific needs of different businesses.

The primary contribution of this study lies in developing an integrated process that spans from speech recognition to information extraction, enhancing efficiency while reducing manual processing costs. Experimental results demonstrate high accuracy in both recognition and extraction steps, promising significant practical value in telesales and customer service domains.

MỤC LỤC

CHƯƠNG 1. GIỚI THIỆU ĐỀ TÀI.....	1
1.1 Đặt vấn đề.....	1
1.2 Các giải pháp hiện tại và hạn chế	2
1.3 Mục tiêu và định hướng giải pháp	3
1.3.1 Mục tiêu.....	3
1.3.2 Định hướng giải pháp	3
1.4 Đóng góp của đề án	4
1.5 Bố cục đề án	4
CHƯƠNG 2. NỀN TẢNG LÝ THUYẾT	6
2.1 Ngữ cảnh của bài toán.....	6
2.1.1 Âm vị	6
2.1.2 Tín hiệu âm thanh	8
2.2 Các kết quả nghiên cứu tương tự	10
2.2.1 Mô hình chuyển từ âm thanh sang văn bản	10
2.2.2 Mô hình phiên âm giọng nói chính xác theo thời gian.....	10
2.3 Mô hình Whisper.....	11
2.4 Mô hình Whisperx.....	13
2.5 Mô hình NVIDIA NeMo	15
CHƯƠNG 3. PHƯƠNG PHÁP ĐỀ XUẤT.....	17
3.1 Tổng quan giải pháp.....	17
3.2 Mô đun xử lý âm thanh thành văn bản sử dụng Whisper Model kết hợp LoRA Adapter	20
3.3 Mô đun chuyển đổi thành đoạn hội thoại tiếng Việt.....	21

CHƯƠNG 4. ĐÁNH GIÁ THỰC NGHIỆM.....	24
4.1 Các tham số đánh giá	24
4.2 Phương pháp thí nghiệm.....	24
4.2.1 Chuẩn bị dữ liệu.....	24
4.2.2 Chuẩn bị môi trường.....	25
4.2.3 Tinh chỉnh mô hình Whisper	26
4.2.4 Tinh chỉnh mô hình Whisper-large-v3	26
4.2.5 Đánh giá kết quả tinh chỉnh.....	27
CHƯƠNG 5. PHÁT TRIỂN THÀNH HỆ THỐNG XỬ LÝ CUỘC GỌI BÁN HÀNG	29
5.1 Khảo sát hiện trạng	29
5.2 Triển khai hệ thống xử lý cuộc gọi bán hàng	31
5.2.1 Phân tích yêu cầu	31
5.2.2 Tổng quan chức năng.....	31
5.2.3 Đặc tả chức năng.....	32
5.3 Công nghệ sử dụng	36
5.3.1 Frontend.....	36
5.3.2 Backend	37
CHƯƠNG 6. KẾT LUẬN	38
6.1 Kết luận	38
6.2 Hướng phát triển trong tương lai	38
CHƯƠNG 7. TÀI LIỆU THAM KHẢO	40

DANH MỤC HÌNH VẼ

Hình 2.1	Bảng âm vị của 22 phụ âm	7
Hình 2.2	Bảng âm vị của 2 bán nguyên âm	7
Hình 2.3	Bảng âm vị của 13 nguyên âm đơn và 3 nguyên âm đôi làm âm chính	8
Hình 2.4	Bảng âm vị của 9 âm cuối	8
Hình 2.5	Biểu đồ dạng sóng	9
Hình 2.6	Mô hình Whisper	11
Hình 2.7	Định dạng đào tạo đa nhiệm	12
Hình 2.8	Kiến trúc mô hình Whisperx	13
Hình 2.9	Multi-Scale Diarization Decoder	15
Hình 3.1	Quy trình xử lý âm thanh thành đoạn hội thoại gốc	17
Hình 3.2	Quy trình xử lý âm thanh	19
Hình 3.3	Whisper Model kết hợp LoRA	20
Hình 3.4	Nhật ký thời gian mẫu	22
Hình 4.1	Công thức tính toán CER	24
Hình 4.2	Biểu đồ thống kê số lượng câu theo thời gian	25
Hình 4.3	Biểu đồ thống kê số lượng câu theo độ dài	25
Hình 4.4	Mô hình LoRA finetuning	26
Hình 5.1	Biểu đồ use case tổng quan của hệ thống	32
Hình 5.2	Danh sách use case	32
Hình 5.3	Đặc tả UC01 - Đăng nhập	33
Hình 5.4	Đặc tả UC02 - Chuyển cuộc gọi thành đoạn hội thoại	33
Hình 5.5	Đặc tả UC03 - Tóm tắt đoạn hội thoại	34
Hình 5.6	Đặc tả UC04 - Trích xuất thông tin bản tóm tắt	35
Hình 5.7	Đặc tả UC05 - Lưu thông tin cuộc gọi	35
Hình 5.8	Đặc tả UC06 - Thống kê	36
Hình 5.9	Giao diện hệ thống	36
Hình 5.10	Tương tác với cơ sở dữ liệu	37

DANH MỤC BẢNG BIỂU

Bảng 4.1	Tỷ lệ CER của model Whiser-large-v3.	27
Bảng 4.2	Tỷ lệ CER của model Whiser-large-v3.	28
Bảng 4.3	Tỷ lệ CER của model Whisper-large-v3 đã finetune bằng LoRA	28

DANH MỤC THUẬT NGỮ VÀ TỪ VIẾT TẮT

Thuật ngữ	Ý nghĩa
ĐATN	Đồ án tốt nghiệp
AI	Trí tuệ nhân tạo (Artificial Intelligence)
API	Giao diện lập trình ứng dụng (Application Programming Interface)
ASR	Nhận dạng tiếng nói tự động (Automatic Speech Recognition)
CER	Tỷ lệ lỗi ký tự (Character Error Rate))
DTW	Dynamic Time Warping)
NLP	Xử lý ngôn ngữ tự nhiên (Natural Language Processing)
OCR	Nhận dạng ký tự quang học (Optical Character Recognition))
VAD	Phát hiện hoạt động giọng nói (Voice Activity Detection)