

2020.7.23

# 음성인식을 이용한 자동자막 프로그램

## (Design and Implementation of Real-time subtitle System using Voice Recognition)

참여인원 권태헌  
이윤성  
자줄리

현재 COVID-19로 인해 국제 재난 수준의 문제가 야기됨에 따라, 이에 많은 문제점들이 확인
---

<p><b>프로젝트 간략 내용</b></p>	<p>되고 있다. 그 중 온라인 비대면을 통한 회의, 수업 등에 쓰이는 일부 영상들이 청각 장애인을 고려하지 않아 참여가 어려운 상태이며, 이러한 상황을 개선할 만한 지원이 이루어지지 않아 불편을 겪는 사람들에 대한 문제가 대두된다. 이에 따른 상황을 개선하기 위해 수화, 수어번역 담당 도우미, 영상 자막 대입 등 많은 노력이 있지만 아직까지 온라인상의 불편함은 남아있다. 문제를 좀 더 확인해 보면 실시간 환경이 대한 대응이 부족하고 수어번역 담당 도우미를 구하거나 매 영상 자막을 대입하는 일은 과도한 인건비를 생산한다. 이런 문제점을 모집한 결과에 따라서 본 시스템은 RT-STT(Real Time Speech to Text)를 적용하여 실시간 온라인 회의, 강의 부분에 대한 실질적 대안을 마련하고자 한다.</p>
<p><b>아키텍처</b></p>	<pre> graph LR     WaveOutput[Wave Output] --&gt; DAC[DAC]     DAC --&gt; EndpointList[Endpoint List select]     EndpointList --&gt; GetAudioSession[get AudioSession]     EndpointList --&gt; ChangeEndpoint[Change Endpoint to VB-cable]     GetAudioSession --&gt; DataSampling[Data Sampling]     ChangeEndpoint --&gt; DataSampling     DataSampling --&gt; Send[Send]     Send --&gt; WSService[WSService]     WSService --&gt; Receive[Receive]     Receive --&gt; View[View]     View --&gt; Export[Export]          subgraph ASR_Model [ASR Model]         FeatureExtraction[Feature Extraction] --&gt; Decoding[Decoding]         AcousticModel[Acoustic Model] --&gt; Decoding         WordLexicon[Word Lexicon] --&gt; Decoding         LanguageModel[Language Model] --&gt; Decoding     end          subgraph WSService [WSService]         Init[Init] --&gt; WSConnect[WSConnect]         WSConnect --&gt; Send         Send --&gt; Receive     end          subgraph Subtitle [Subtitle]         View         Export     end </pre> <p>The diagram illustrates the system architecture. It starts with 'Wave Output' leading to a 'DAC' (Digital-to-Analog Converter). The signal then goes to 'Endpoint List select', which branches into 'get AudioSession' and 'Change Endpoint to VB-cable'. Both lead to 'Data Sampling'. 'Data Sampling' feeds into 'Send', which connects to the 'WSService' block. Inside 'WSService', the flow is 'Init' → 'WSConnect' → 'Send' → 'Receive'. 'Receive' feeds into the 'Subtitle' block, which contains 'View' and 'Export'. Additionally, 'Data Sampling' feeds into the 'ASR Model' block. The 'ASR Model' contains 'Feature Extraction', 'Acoustic Model', 'Decoding', 'Word Lexicon', and 'Language Model'. 'Feature Extraction' and 'Acoustic Model' feed into 'Decoding', which also receives input from 'Word Lexicon' and 'Language Model'. 'Decoding' feeds back into 'Data Sampling'.</p>
<p><b>개발내용</b></p>	<p>Zeroth를 이용한 음성인식 모델 개발 및 학습 ELC 개발</p>
<p><b>이슈</b></p>	<p>음성인식 모델에 대한 개념과 용어가 생소하여 학습에 오랜 시간이 걸림.</p>
<p><b>시연연상</b></p>	<p><a href="https://youtu.be/7hMhH9JVMJc">https://youtu.be/7hMhH9JVMJc</a></p>