

●○ 상황별음성 과제

회의 음성 데이터



●○ 개요: 회의음성 데이터란?

- 한국인의 음성을 문자로 바꾸어 주고, 문맥을 이해하는 한국어 음성언어처리 기술 개발을 위한 AI 학습용 한국어 음성 DB를 구축.
- 한국어로 된 회의영상/음성을 인식하여 자동으로 자막을 생성해주고, 내용을 이해하는 서비스를 위한 한국어 회의 음성DB를 구축.
- AI 학습용 데이터 구축량 : 한국교육방송공사(EBS)로부터 제공받은 3,000 방송시간 이상 분량의 토론/토크 콘텐츠로 20~40분 내외로 구성되며 3인 이상의 화자로 구성.
- 화상회의 대화록 작성 기술은 아래 그림과 같은 데이터의 특징을 작고 있으며, 각 채널에 대해 회의록을 작성하는 것을 목표로 함.
- 화자별로 채널이 분리된 상황에서 서비스는 크게 파일 형식으로 회의 데이터 전체 파일을 입력으로 받아 회의록을 생성하는 서비스와 실시간 화상회의 상황에서 회의록을 작성하는 서비스로 분리할 수 있음.
- 작성된 회의록은 음성의 시작 시간, 끝 시간, 입력 채널 및 텍스트로 구성되어 있으며, 이를 자막으로 변환하여 각 화자 채널 화면에 나타내어 시각 효과를 강화하고자 함.

| 시작시간 | 끝시간 | 채널 | 텍스트 |
|-----------|-----------|-----|-------------|
| 0.309622 | 2.456449 | 1 | 안녕하세요 |
| 2.403249 | 4.582167 | 2 | 반갑습니다 |
| 5.638745 | 9.213648 | 4 | 오늘 주제가 뭔가요 |
| 10.151987 | 14.121577 | 3 | 화상회의록 작성입니다 |
| ... | ... | ... | ... |

그림 | 화상회의 회의록 작성 결과

●○ 데이터셋의 구성

| 구분 | 회의 음성 데이터 |
|-----------|---|
| 시간(hr) | • 3,000 |
| 데이터 선정 기준 | <ul style="list-style-type: none"> • 3인 이상 화자로 구성된 토론/토크 콘텐츠에서 수집한 데이터 • 전체 회의 시간(20분, 30분, 40분)별 3종 • 교육, 문화예술, 가족, 교양, 금융, 시사, IT, 토크 총 8개 분야 선정 • 400시간 X 6개 분야 = 2,400시간 • 300시간 X 2개 분야 = 600시간 |
| 수집항목 | • 회의(토론/토크) 음성, 주제, 연령, 참여인원, 회의시간 등 |

- 문장별로 전사된 시간 기준으로 데이터량을 산정했으며, 모든 데이터셋은 음성, 스크립트로 구성되어 있음.
- 전사된 문장의 시간에는 앞뒤의 묵음 시간이 포함되지 않으므로 실제 음성 신호만으로 데이터 구축량 산정 가능
- 음성 신호만으로 3,00 시간 이상의 데이터 구축
- 방송 참여자 메타데이터를 바탕으로 연령, 지역 등을 고려하여 균형 있는 데이터 구축

●○ 데이터셋의 설계 기준과 분포

- 최근 AI를 활용한 음성인식 기술이 인공지능 비서를 포함한 다양한 서비스에 적용되고 있음
- 하지만 AI 음성인식 모델이 회의 음성 데이터로 학습이 이루어지지 않을 경우, 회의에서 발생하는 다양한 노이즈 및 화자 간 목소리가 겹치는 등의 문제로 음성인식 성능이 떨어지는 경향을 보임
- 회의 데이터는 회의 참여자 간의 민감한 정보를 다루는 경우가 많아 외부 유출이 가능한 경우가 드물어 데이터셋 구축에 어려움이 있음
- 회의 테이블의 크기, 회의 참여자의 위치와 마이크의 개수와 배치에 따라 수집하는 음성의 품질이 다양해 질 수 있음
- 본 과제에서는 회의 음성 도메인에 대한 AI 음성인식 성능 향상을 위한 한국교육방송공사(EBS)로부터 제공받은 토론/토크 콘텐츠를 활용하여 3,000 시간 이상의 회의 음성 데이터셋을 구축함
- 회의 음성 도메인은 넓은 범위를 다룰 필요가 있으며, 이를 위해 교육, 문화예술, 가족, 교양, 금융, 시사, IT 토크 등 총 8개 도메인을 선정하였으며, 도메인 별 최소 300시간 이상의 데이터로 구성함

| 구분 | 교육 | 문화예술 | 가족 | 교양 | 금융 | 시사 | IT | 토크 |
|-------|-----|------|-----|-----|-----|-----|-----|-----|
| 계획 시간 | 400 | 400 | 400 | 400 | 400 | 400 | 300 | 300 |
| 합계 | 400 | 400 | 400 | 400 | 400 | 400 | 300 | 300 |

- 구축한 강의 음성 데이터셋은 1,000명 이상의 발화자로 구성하며, 음성 별 발화자에 대한 메타 데이터를 제공함
- 녹음은 약 128m², 400m², 600m², 800m² 공간에서 슈어 제조사의 UR4D+ 모델 또는 제나이저 SK5212/SKM5200 모델을 사용하여 이루어졌으며, 발화자와 마이크 사이의 거리는 약 20cm에서 50cm 사이임
- 데이터 전사 과정에서 앞뒤의 묵음구간을 포함하도록 싱크를 조절하였으나, 일부 데이터의 경우 발화 사이의 묵음구간이 짧을 수 있음
- 1,4절 음성 전사 규칙에 발음겹침, 잡음, 말더듬, 텍스트정규화/비정규화, 발화자 처리 등에 대한 처리 방안을 제시함
- 발음겹침 음성의 경우 메인 발화자 이외의 발화자에 대한 텍스트 전사에 상당한 어려움이 있어, 메인 발화자만 텍스트 전사하고 해당 어절에 전사 규칙에 따라 ‘+’를 표시함
- 잡음이 포함된 음성의 경우 전사자와 검수자의 주관에 따라 충분히 인식이 가능한 경우만 포함하였음
- 데이터셋 구축 과정에서 성명, 전화번호, 주소 등의 개인정보가 포함된 음성은 배제함
- 성차별, 정치적 성향, 종교 등 사회적 민감 정보를 포함한 음성은 배제함
- 공개된 한국어 강의 데이터셋은 적합한 콘텐츠 구매과정을 통하였으며, 저작권자의 허가를 받아 데이터 사용 및 활용이 자유로움
- 데이터셋의 효용성을 측정하기 위해 구축 데이터 기반의 상용화 가능한 수준의 AI 응용 서비스 개발을 진행함

●○ 데이터 구조

| No | 항목 | | 타입 |
|-------|-----------------------------|-------------|--------|
| | 한글명 | 영문명 | |
| | 데이터셋 | dataSet | |
| 1 | 데이터셋 버전 | version | String |
| 2 | 녹취된 음원의 URL | mediaUrl | String |
| 3 | 녹취된 날짜 | date | String |
| 4 | 음원 데이터 상세 정보 | typeInfo | |
| 4-1 | 음원 카테고리 정보 : 강의, 회의, 고객응대 등 | category | String |
| 4-2 | 음원 서브카테고리 | subcategory | String |
| 4-3 | 음원 녹취 장소 | place | String |
| 4-4 | 화자 목록 | speakers | List |
| 4-3-1 | 화자 유형 : 강사, 상담사, 고객, 기타 | type | String |
| 4-3-2 | 인입 유형 : 유선, 모바일, 인터넷 등 | gender | String |
| 4-5 | 화자 성별 : 남성, 여성 | inputType | String |
| 5 | 전사 데이터 목록 : 화자가 변경될 때마다 생성 | dialogs | List |
| 5-1 | 화자 아이디 : speakers 에 등록된 순번 | speaker | String |
| 5-2 | 전사된 텍스트 | text | String |
| 5-3 | 전사된 텍스트의 음원 재생 시작 위치 | startTime | String |
| 5-4 | 전사된 텍스트의 음원 재생 끝 위치 | endTime | String |
| 5-5 | 전사된 텍스트 문장과 관련된 태그 리스트 | tags | String |

●○ 데이터 예시

```

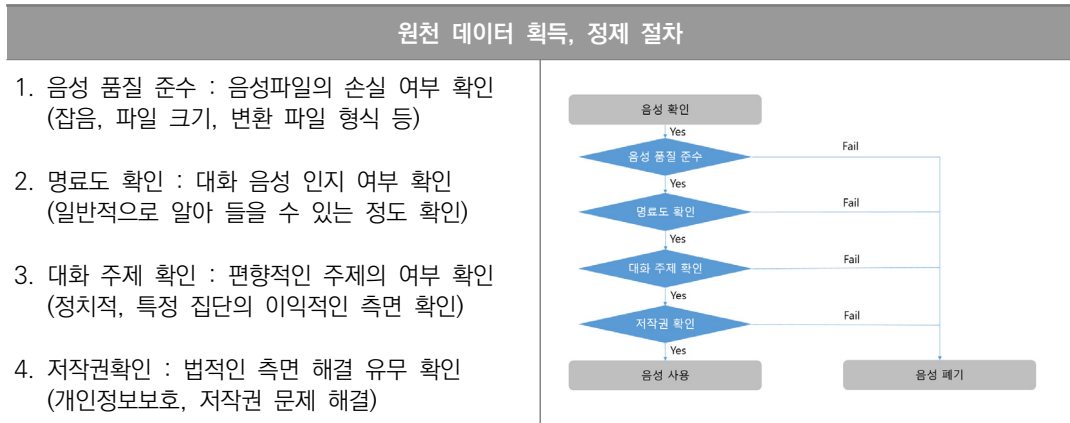
{
  "metadata": {
    "title": "CG2D1450041",
    "creator": "윤현식 (yh7335@naver.com)",
    "distributor": "",
    "year": "2020년",
    "category": "",
    "annotation_level": [],
    "sampling": ""
  },
  "speaker": [
    {
      "no": 1,
      "name": "김일중",
      "age": null,
      "sex": "남",
      "shortcut": 1,
      "occupation": "지식가이드 MC"
    },
    {
      "no": 2,
      "name": "한수연",
      "age": null,
      "sex": "여",
      "shortcut": 2,
      "occupation": "배우"
    },
    {
      "no": 3,
      "name": "남인숙",
      "age": null,
      "sex": "여",
      "shortcut": 3,
      "occupation": "작가"
    },
    {
      "no": 4,
      "name": "M1",
      "age": null,
      "sex": "남",
      "shortcut": null,
      "occupation": "지식가이드"
    },
    {
      "no": 5,
      "name": "김유진",
      "age": null,
      "sex": "남",
      "shortcut": null,
      "occupation": "지식가이드 2"
    },
    {
      "no": 6,
      "name": "etc",
      "age": null,
      "sex": null,
      "shortcut": null
    },
    {
      "no": 7,
      "name": "M2",
      "age": "",
      "sex": "남",
      "shortcut": null,
      "occupation": "오징어빵 장수"
    },
    {
      "no": 8,
      "name": "M3",
      "age": null,
      "sex": "남",
      "shortcut": null,
      "occupation": "용가리 장수"
    }
  ],
  "setting": {
    "relation": ""
  },
  "utterance": [
    {
      "id": "a02e0e9a-a499-4f52-9d77-3c9eef013965",
      "start": 25.75,
      "end": 29.8,
      "speaker_id": 1,
      "form": "n/ 여행자를 위한 지식 가이드 숨은 한국 찾기 김일중입니다. 반갑습니다.",
      "standard form": "n/ 여행자를 위한 지식 가이드 숨은 한국 찾기 김일중입니다. 반갑습니다.",
      "dialect form": "n/ 여행자를 위한 지식 가이드 숨은 한국 찾기 김일중입니다. 반갑습니다.",
      "note": ""
    },
    {
      "id": "f02fc60d-aa1d-499c-875f-5c06e8c70eae",
      "eojjeol": "n/",
      "standard": "n/",
      "begin": 0,
      "end": 2,
      "isDialect": false
    },
    {
      "id": "297495ee-959a-4715-a6f1-88601e7d79be",
      "eojjeol": "여행자들",
      "standard": "여행자들",
      "begin": 3,
      "end": 7,
      "isDialect": false
    },
    {
      "id": "64b9e323-5543-4ee6-8dcc-ba990fb6cae3",
      "eojjeol": "위한",
      "standard": "위한",
      "begin": 8,
      "end": 10,
      "isDialect": false
    },
    {
      "id": "a892fa06-a41a-41cc-83cb-727e350d66bf",
      "eojjeol": "지식",
      "standard": "지식",
      "begin": 11,
      "end": 13,
      "isDialect": false
    },
    {
      "id": "0abfd631-bb99-4918-a1a6-cd43561b4e91",
      "eojjeol": "가이드",
      "standard": "가이드",
      "begin": 14,
      "end": 17,
      "isDialect": false
    },
    {
      "id": "6f0193d8-4eca-4e12-bb20-67a63fd5f217",
      "eojjeol": "숨은",
      "standard": "숨은",
      "begin": 18,
      "end": 20,
      "isDialect": false
    },
    {
      "id": "98762c12-b3f7-439e-93e2-14659a1dfe3d",
      "eojjeol": "한국",
      "standard": "한국",
      "begin": 21,
      "end": 23,
      "isDialect": false
    },
    {
      "id": "46203d50-20d5-4992-acd2-a36ded20962f",
      "eojjeol": "찾기",
      "standard": "찾기",
      "begin": 24,
      "end": 26,
    }
  ]
}

```

- 상기데이터는 전체 전사한 JSON 파일의 일부임

●○ 데이터 구축 과정

• 획득 정제 절차



• 획득 정제 기준

| 구분 | 절차 | 기준 |
|----|-------------|--|
| 획득 | 1. 음성 품질 준수 | <ul style="list-style-type: none"> • Extension: PCM • Precision: 16-bit • Sample rate: 16kHz • Channel: mono • Sample Encoding: 16-bit Signed Integer PCM |
| 정제 | 2. 명료도 확인 | <ul style="list-style-type: none"> • 클리핑, Frame drop 등 손실된 음성데이터 제외 • 심한 잡음으로 사람도 인식하기 어려운 음성데이터 제외 • 비음성구간 제외 |
| | 3. 대화 주제 확인 | <ul style="list-style-type: none"> • 민감한 이슈 (정치적 견해, 개인정보, 특정인물 비하, 성적인 표현) 발언이 포함된 경우 제외 |
| | 4. 저작권 확인 | <ul style="list-style-type: none"> • 지적재산권 및 개인정보보호 관련 사항 해결 유무 |
| | 5. 최종 결정 | <ul style="list-style-type: none"> • 4가지 절차에 모두 이상이 없는 경우 |

●○ 검수화 품질 확보

- 데이터 검수는 양질의 데이터를 얻는 데 필요한 작업으로 정성적/정량적 평가를 통해 데이터의 유효함을 판별함
- 웹 저작도구를 사용하여 전사한 어노테이션 결과에 대해 1차, 2차, 3차 검수자를 걸쳐 음성품질 / 어노테이션 정확도 / 대화 주제 및 저작권확인 / 목표데이터 수집량 달성 여부를 확인함
- 1, 2차 검수자는 음성품질 및 어노테이션 정확도를 위주로 검수를 진행하며, 교차 검증을 통해 데이터의 오류를 최소화함
- 3차 검수자는 대화 주제 및 저작권확인, 목표데이터의 도메인별 수집량 달성 여부를 위주로 검수를 진행하며, 데이터의 다양성을 보장할 수 있도록 함
- 각 차수의 검수자는 피드백을 통해 데이터가 한쪽에 치우치지 않도록 하며, 전사자가 숙달할 수 있도록 지원함
- 3차 검수 완료 이후 TTA를 통한 데이터 검증 절차를 걸치며, 동시에 AI 모델 학습을 진행하여 데이터의 적합성과 유효성을 최종 판별함
- 어노테이션 결과 피드백 기준

| 대항목 | 피드백 목표 | 비고 |
|---------------|--|-------------------------------------|
| 음성품질 | <ul style="list-style-type: none"> • 음성 파일의 샘플레이트, 채널, 인코딩의 적절성 여부 • 음성의 명료함 정도 (SNR) | 전사자와 검수자 주관에 맡기며, 과반수 동의 시 적절함으로 판정 |
| 어노테이션 정확도 | <ul style="list-style-type: none"> • 음성에 대응하는 텍스트 어노테이션 정확성 여부 • 음성에 대응하는 발화자 정보의 정확성 여부 • 음성과 텍스트의 싱크 정확성 여부 | 상기와 동일 |
| 대화 주제 및 저작권확인 | <ul style="list-style-type: none"> • 대화 주제의 편향성 여부 • 저작권 및 개인정보 침해 여부 | 상기와 동일 |
| 목표데이터 수집량 달성 | <ul style="list-style-type: none"> • 데이터 도메인, 화자 등 수집량의 적절성 여부 | 상기와 동일 |

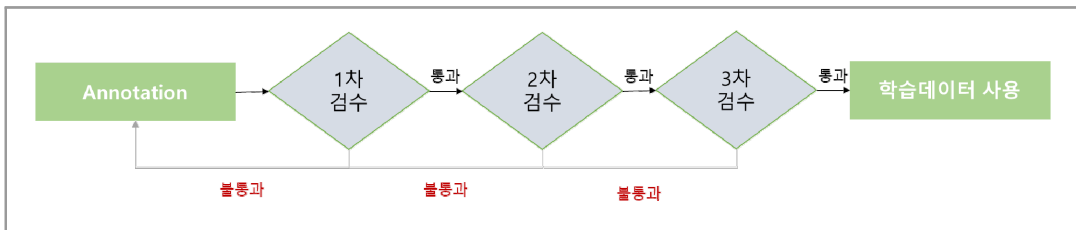


그림 | 어노테이션 결과 검증 절차

• 어노테이션 데이터 항목별 세부 판정 기준

| 대항목 | 데이터 항목 | 어노테이션 | 판정 기준 |
|----------------|------------|----------------------|---------------------------------------|
| 음성품질 | 음성 파일 형식 | • N/A | • 전사도구 내 음성 파일 체크를 통한 자동 검수 |
| | 명료함 | • N/A | • 3인 이상의 검수자 판단하에 과반수 통과 |
| 어노테이션 정확도 | 텍스트 | • 전사규칙에 따라 전사된 발화 내용 | • 3인 이상의 검수자 판단하에 과반수 통과 |
| | 발화자 | • 화자 고유번호, 화자 성별, 나이 | • 3인 이상의 검수자 판단하에 과반수 통과 |
| | 싱크 | • 음성과 텍스트의 시작과 끝 | • 3인 이상의 검수자 판단하에 과반수 통과 |
| 대화 주제 및 저작권 확인 | 대화주제 | • 회의 도메인에 8종에 대한 태깅 | • 3인 이상의 검수자 판단하에 과반수 통과 |
| | 저작권 및 개인정보 | • N/A | • 3인 이상의 검수자에 의한 저작권 및 개인정보 침해여부 판단 |
| 목표데이터 수집량 달성 | 도메인 | • 회의 도메인 별 음성 길이 | • 음성 발화자 및 싱크 정보를 바탕으로 알고리즘을 통한 자동 계산 |
| | 화자 | • 화자 별 음성 길이 | • 음성 발화자 및 싱크 정보를 바탕으로 알고리즘을 통한 자동 계산 |

| 승인 이력 | | | |
|-------|------|---|---------------------|
| 번호 | 상태 | 내용 | 날짜 |
| 1 | 승인거절 | 10:37 ~ 10:48 영어로 -> (0)/(영)으로 11:10 ~ 11:16 마이너스에 (2) 분의 -> 마이너스 (A 분의)/(에이 분의) 15:27 ~ 15:29 최고 왕의 -> 최고 차왕의 18:44 ~ 18:54 숫자 자들의 -> 숫자 왕의 19:28 ~ 19:42 세 계급만 접근 -> 격분 33:40 ~ 33:48 (25)/(이 십 오) -> (25)/(이십 오) 열에 4도 똑같이 34:05 ~ 34:11 위에 것과 똑같이 바꿔주세요 34:11 ~ 34:20 위에 것 똑같이 35:02 ~ 35:08 위에 것 똑같이 36:08 ~ 36:11 37:26 ~ 37:36 (66)/(육십육), (61)/(육 십 일) -> (66)/(육십 육), (61)/(육십 일) 37:39 ~ 37:50 40:28 ~ 40:34 42:15 ~ 42:28 | 2020.10.29 17:10:04 |

그림 | 검수자 피드백 예시