

귀가노니 - 출퇴근길에 듣는 인공지능 뉴스 팟캐스트

#뉴스요약 #팟캐스트 #TTS #갓생

NLP-05

(ㅇㄱㅇ팀)

CONTENTS

01 팀 소개

02 프로젝트 소개

03 모듈별 소개

04 시연 영상

05 향후 개선사항

06 Q&A

01

팀 소개

귀가노니 - 출퇴근길에 듣는 인공지능 뉴스 팟캐스트

01 팀 소개



권혜빈

캠퍼 ID: T2014
#Clustering #TTS



김민진

캠퍼 ID: T2029
#Summarization



김상현

캠퍼 ID: T2034
#Summarization



 **김성한**

캠퍼 ID: T2039
#PM #Crawling #Serving



김제우

캠퍼 ID: T2052
#Crawling #Serving #TTS



이노아

캠퍼 ID: T2152
#Clustering #TTS



이다솔

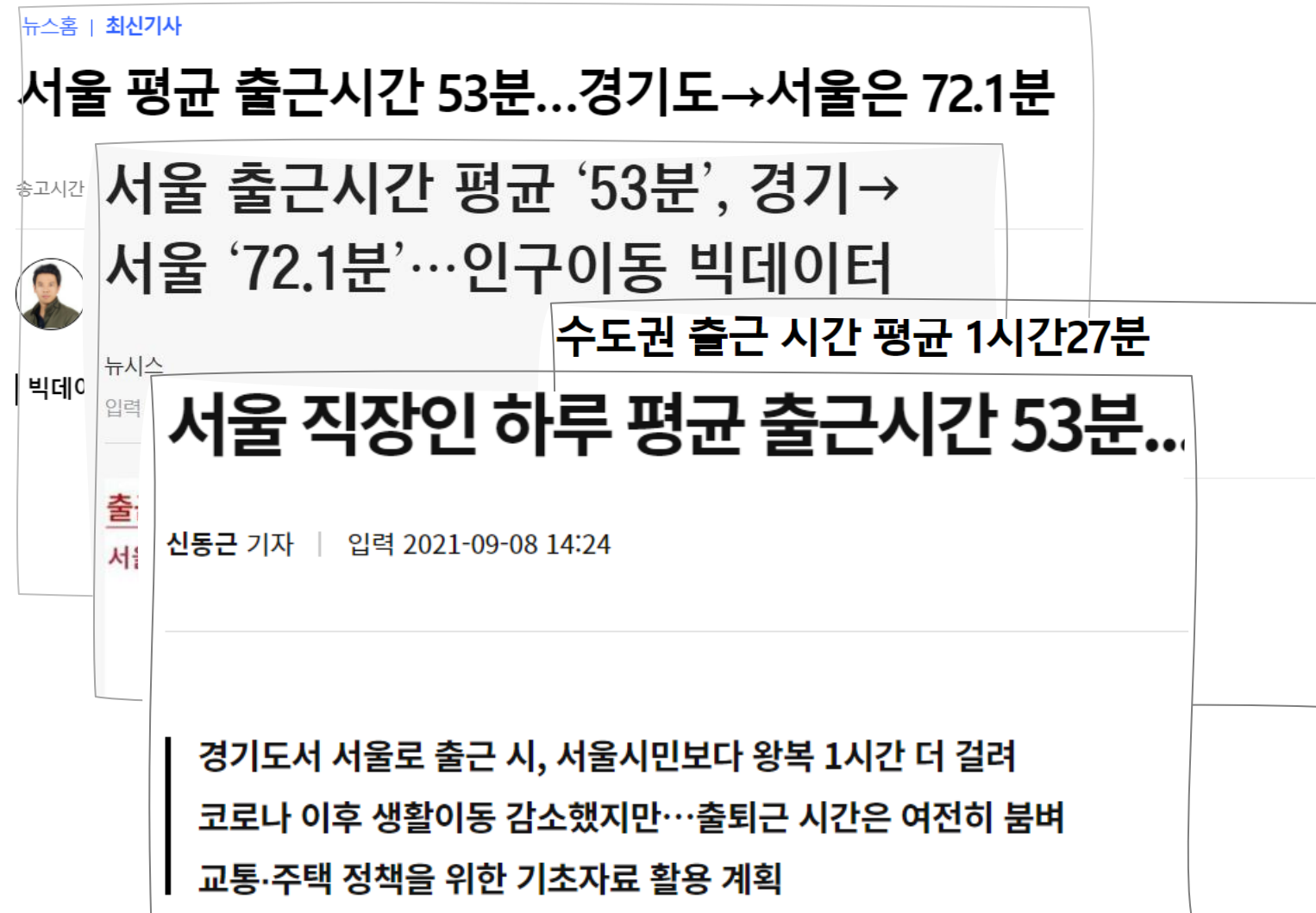
캠퍼 ID: T2154
#PPT #Clustering #TTS

02

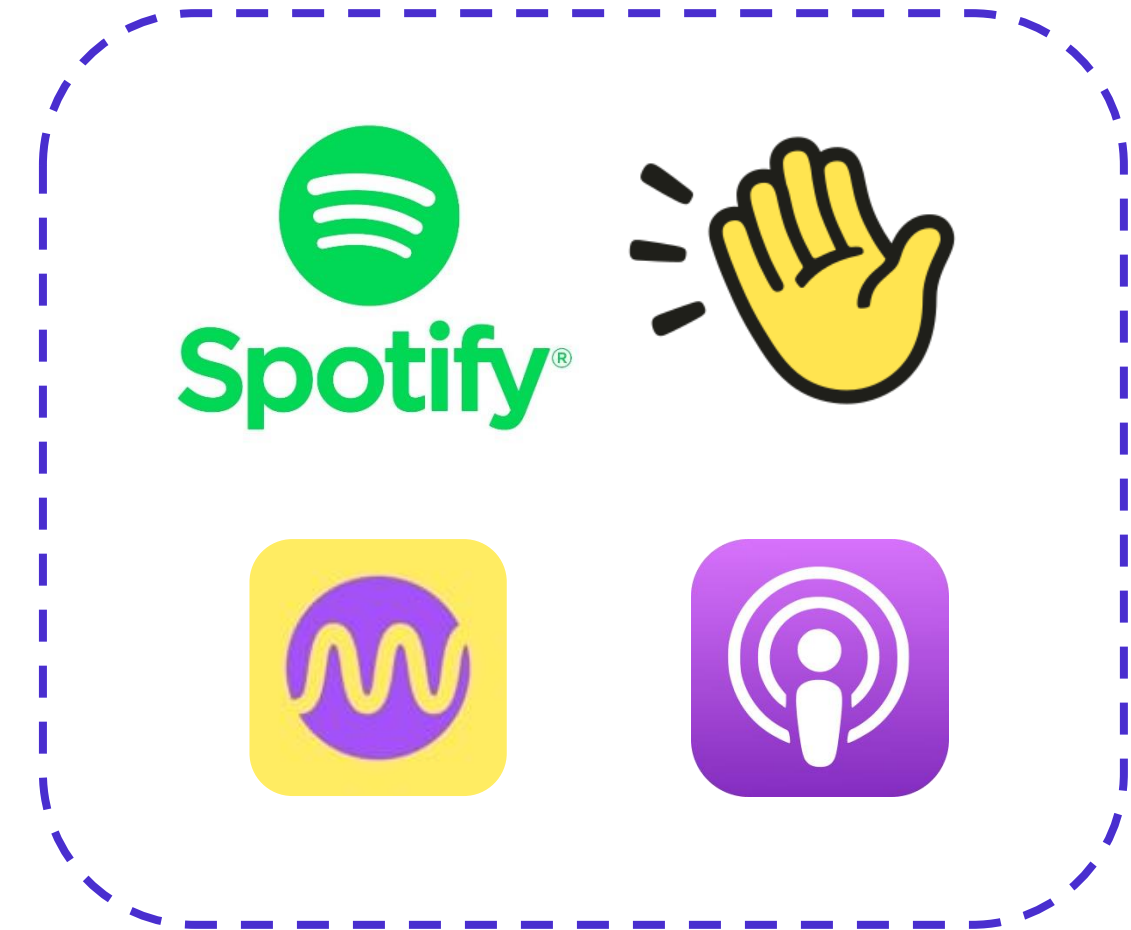
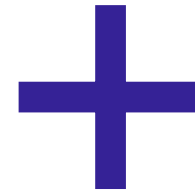
프로젝트 소개

귀가노니 - 출퇴근길에 듣는 인공지능 뉴스 팟캐스트

02 프로젝트 소개



서울 평균 출근시간 53분, 경기도→서울은 72분
이 시간을 유익하게 쓸 수는 없을까?



밀리의 서재, 클럽하우스 등
오디오 서비스 플랫폼의 등장



어제 기사들에 대한 주제별 클러스터링을 이용한
개인맞춤형 인공지능 TTS 팟캐스트, **귀가노니**

03

모듈별
소개

귀가노니 - 출퇴근길에 듣는 인공지능 뉴스 팟캐스트

Clustering

FastAPI

EasyBART

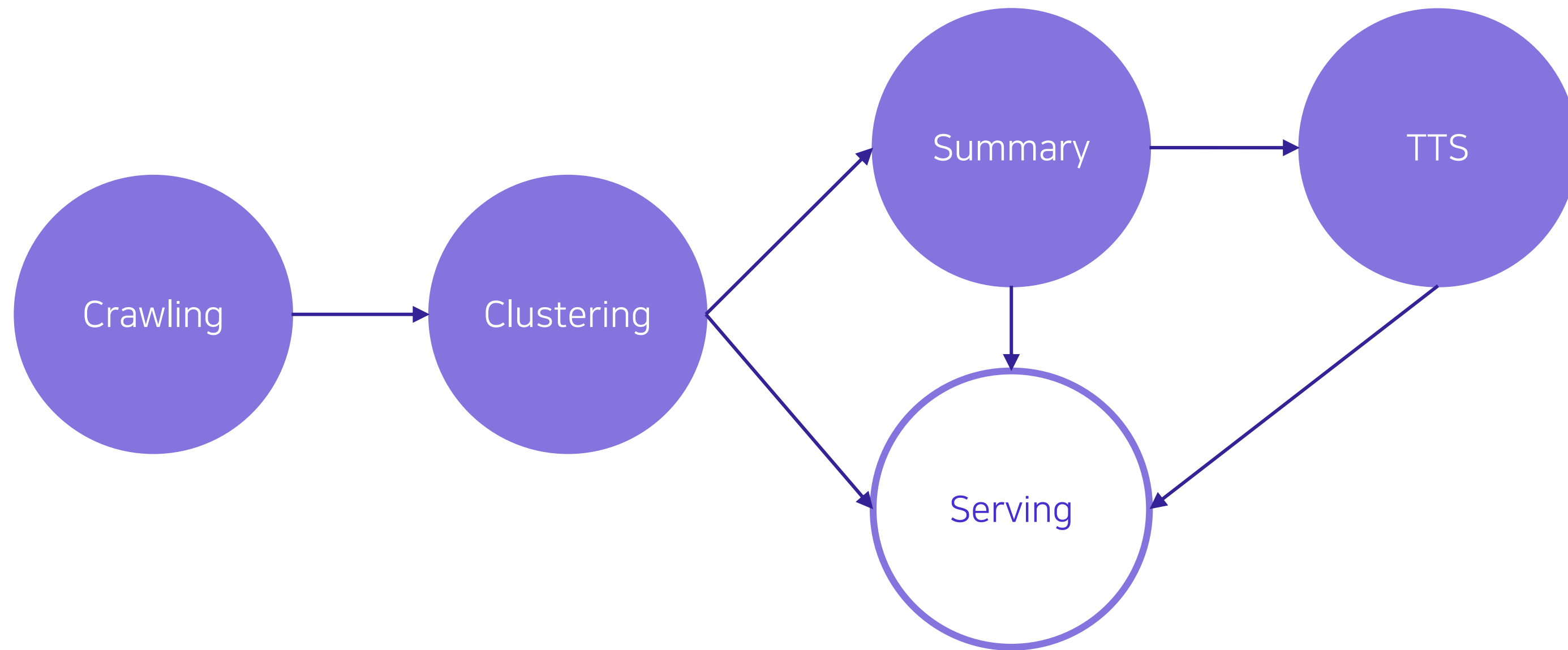
Crawling

Model Serving

TTS

03 모듈별 소개

셸 스크립트에 각각의 모듈을 순차적으로 실행하는 파이썬 명령어 정의
Crontab을 통해 매일 밤 12시에 어제 날짜에 해당하는 데이터가 생성되도록 지정



<데이터 생성 배치 프로세싱>

01 Crawling

02 Clustering

03 Summarization

04 Text to Speech

05 Serving

전체기사

< 2021. 12. 17 > 오늘

8개 카테고리

- 최신
- 사회
- 정치
- 경제
- 국제
- 문화
- 연예
- 스포츠
- IT
- 칼럼
- 보도자료
- 자동생성기사
- 전체기사
- 금융
- 기업산업
- 취업직장인
- 경제일반
- 자동차
- 주식
- 시황분석
- 공시
- 해외증시
- 채권선물
- 외환
- 주식일반
- 부동산
- 생활경제
- 국제경제

영업제한이라도 매출 감소해야..추가 지원은 '미지수'

[앵커] 이렇게 방역 강화와 함께 내놓은 소상공인 지원책. 실제로 받게 될 소상공인들 반응은 어떨고 또 짊어볼 내용은 없는지 산업과학부 김지숙 기자와 이야기 나눠...

매출감소 확인되면 지원금 100만 원..손실보상 확대

[앵커] 정부가 방역조치 강화로 피해를 보게되는 소상공인을 위한 대책을 내놨습니다. 방역지원금 100만원을 지원하고, 손실보상 대상 업종도 확대한다는 내용입니다. ...

[2021 대한민국 신성장경영대상] 에스넷시스템, 산업통상자원부장관 표창

에스넷시스템이 제18회 대한민국 신성장 경영대상에서 산업통상자원부 장관상을 수상했습니다. 에스넷시스템은 통상 신장 사업을 넘어 시스템통합과, 솔루션, 사물인터넷 ...

11월 울산 통관 기준 수출 67억 달러..34%↑

[KBS 울산]울산세관은 통관 기준 지난달 울산지역 수출액은 67억 7천만 달러로 지난해 11월에 비해 34.4% 증가했다고 밝혔습니다. 지난달 수입액도 전년 동월 대비 76.4...

울산 주력산업 내년 업황 전망도 '맑음'

[KBS 울산] [앵커] 조선과 자동차 정유 등 울산의 주력산업이 내년에도 높은 성장세를 이어갈 것이라는 전망이 나왔습니다. 다만 코로나19 장기화 가능성과 원자재가격 ...

[2021 대한민국 신성장경영대상] 엘엔벤처그룹, 매경미디어그룹 회장상

엘엔벤처그룹이 제18회 대한민국 신성장 경영대상에서 우수상인 매경미디어그룹 회장상을 수상했습니다. 대한민국 신성장 경영대상은 MBN과 매일경제신문이 산업통상자...

카테고리당 약 3,000 ~ 4,000개 기사

```
{
  "category": "경제",
  "id": "20211217235422809",
  "source": "한국경제TV",
  "publish_date": "2021-12-17 23:54",
  "extractive": [0],
  "abstractive": [],
  "title": "초반 기술주 매도세 지속..리비안 급락 [뉴욕증시 나우]",
  "text": [
    [
      {
        "index": 0,
        "sentence": "[한국경제TV 신인규 기자]여기는 미국 동부시간 17일 오전 9시 30분입니다."
      },
      {
        "index": 1,
        "sentence": "어제에 이어 오늘도 미국 내에서 기술주에 대한 매도세가 나타나면서 개장 전 3대 지수 선물에 ..."
      },
      {
        "index": 2,
        "sentence": "우선 기술 섹터 부문에서 종목들의 개장 전 거래 움직임을 살펴보면 애플이 1.2%대, 테슬라가 ..."
      },
      {
        "index": 3,
        "sentence": "프리마켓에서 나스닥 100 지수에 편입된 거래량 상위 10개 종목 가운데 9개 종목이 하락중입니 ..."
      },
      {
        "index": 4,
        "sentence": "그래도 프리장에서 개장을 앞두고 낙폭이 조금씩 줄어들고 있는 모습도 보였습니다."
      }
    ],
    [
      {
        "index": 5,
        "sentence": "FOMC 발표 직후에 장이 반짝 상승했다 그 다음날부터 하락세가 지속되고 있지만 단순히 금리 인 ..."
      },
      {
        "index": 6,
        "sentence": "보통은 채권시장에서 금리 인상 신호가 나오면 2년물 국채 금리가 오르는데 2년물 미 국채 수익 ..."
      }
    ]
  ],
}
```

AI Hub 포맷의 JSON 파일

01 Crawling

02 Clustering

03 Summarization

04 Text to Speech

05 Serving

전체기사

8개 카테고리

최신 사회 정치 경제 국제 문화 연예 스포츠 IT 칼럼 보도자료 자동생성기사

전체기사 금융 기업산업 취업직장인 경제일반 자동차 주식 시황분석 공시 해외증시 채권선물 외환 주식일반 부동산 생활경제 국제경제

영업제한이라도 매출 감소해야.추가 지원은 '미지수'

[앵커] 이렇게 방역 강화와 함께 내놓은 소상공인 지원책. 실제로 받게 될 소상공인들 반응은 어떨고 또 짊어볼 내용은 없는지 산업과학부 김지숙 기자와 이야기 나눠...

매출감소 확인되면 지원금 100만 원..손실보상 확대

[앵커] 정부가 방역조치 강화로 피해를 보게되는 소상공인을 위한 대책을 내놓고, 손실보상 대상 업종도 확대한다는 내용입니다...

[2021 대한민국 신성장경영대상] 에스넷시스템, 산업통상자원부장관 표창

에스넷시스템이 제18회 대한민국 신성장 경영대상에서 산업통상자원부 장관상을 수상했습니다. 에스넷시스템은 통신판매 사업을 넘어 시스템통합과, 솔루션, 사물인터넷...

11월 울산 통관 기준 수출 67억 달러..34%↑

[KBS 울산]울산세관은 통관 기준 지난달 울산지역 수출액은 67억 7천만 달러로 지난해 11월에 비해 34% 증가했다고 밝혔습니다. 지난달 수입액도 전년 동월 대비 76.4%...

울산 주력산업 내년 업황 전망도 '맑음'

[KBS 울산] [앵커] 조선과 자동차 정유 등 울산의 주력산업이 내년에도 높은 성장세를 이어갈 것이라는 전망이 나왔습니다. 다만 코로나19 장기화 가능성과 원자재가격...

[2021 대한민국 신성장경영대상] 엘엔벤처그룹, 매경미디어그룹 회장상

엘엔벤처그룹이 제18회 대한민국 신성장 경영대상에서 우수상인 매경미디어그룹 회장상을 수상했습니다. 대한민국 신성장 경영대상은 MBN과 매일경제신문이 산업통상자...

{

"category": "경제",

"id": "20211217235422809",

"source": "한국경제TV",

"publish_date": "2021-12-17 23:54",

"extractive": [0],

"abstractive": [],

"title": "초반 기술주 매도세 지속..리비안 급락 [뉴욕증시 나우]",

"text": [

{

"index": 0,

"sentence": "[한국경제TV 신인규 기자]여기는 미국 동부시간 17일 오전 9시 30분입니다."

},

{

"index": 1,

"sentence": "[아제에 이어 오늘도 미국 내에서 기술주에 대한 매도세가 나타나면서 개장 전 3대 지수 선물에"

},

"sentence": "우선 기술 섹터의 주요 종목들의 개장 전 거래 움직임을 살펴보면 애플이 1.2%대, 테슬라가

},

{

"index": 3,

"sentence": "프리마켓에서 나스닥 100 지수에 편입된 거래량 상위 10개 종목 가운데 9개 종목이 하락중입니

},

"index": 4,

"sentence": "장에서 개장을 앞두고 낙폭이 조금씩 줄어들고 있는 모습도 보였습니다."

},

{

"index": 5,

"sentence": "FOMC 발표 직후에 장이 반짝 상승했다 그 다음날부터 하락세가 지속되고 있지만 단순히 금리 인

},

{

"index": 6,

"sentence": "보통은 채권시장에서 금리 인상 신호가 나오면 2년물 국채 금리가 오르는데 2년물 미 국채 수익

},

],

}

카테고리당 약 3,000 ~ 4,000개 기사

AI Hub 형태의 JSON 파일

boostcamp AI Tech - Networking Day

10

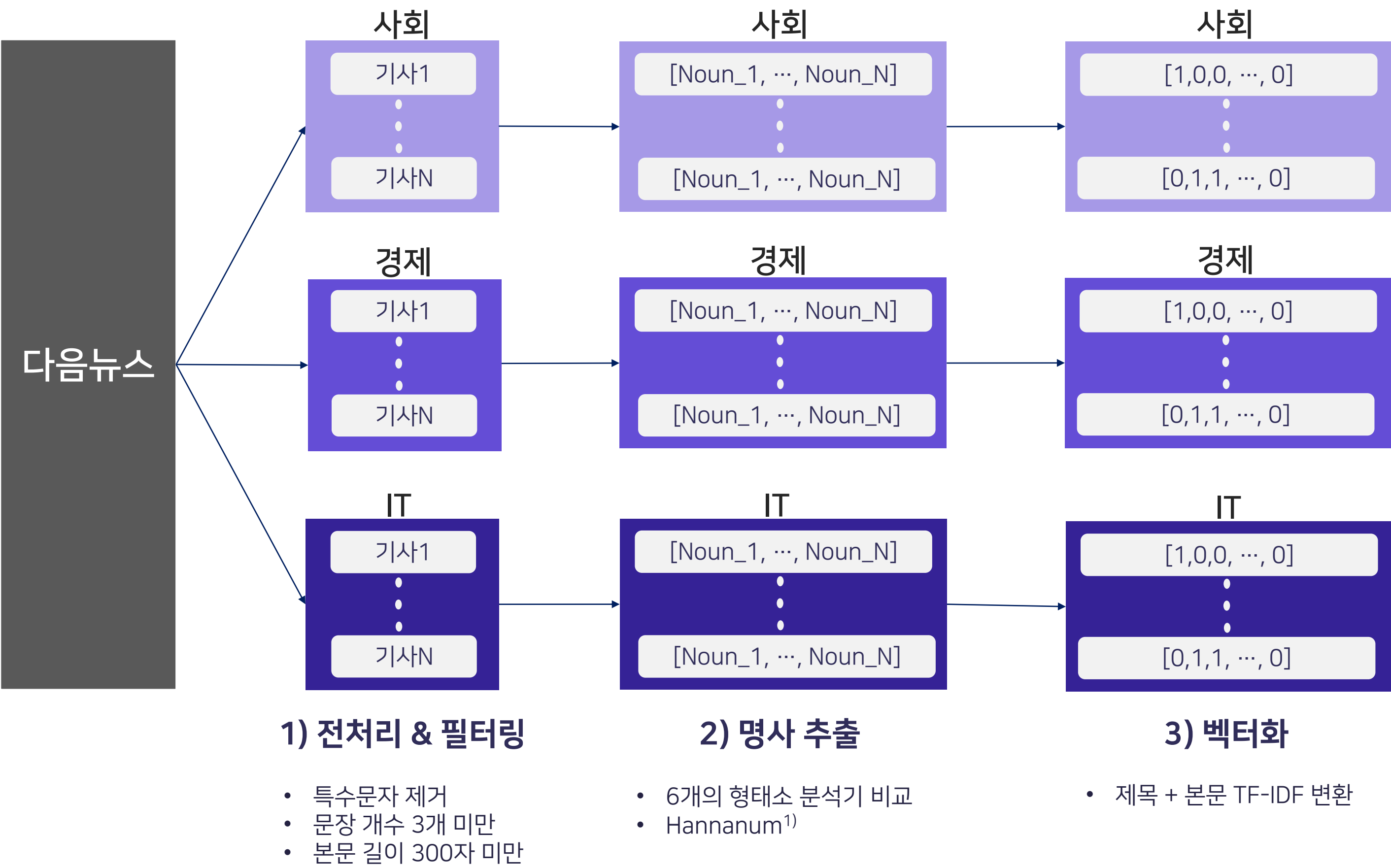
01 Crawling

02 Clustering

03 Summarization

04 Text to Speech

05 Serving



1) <http://semanticweb.kaist.ac.kr/hannanum/>

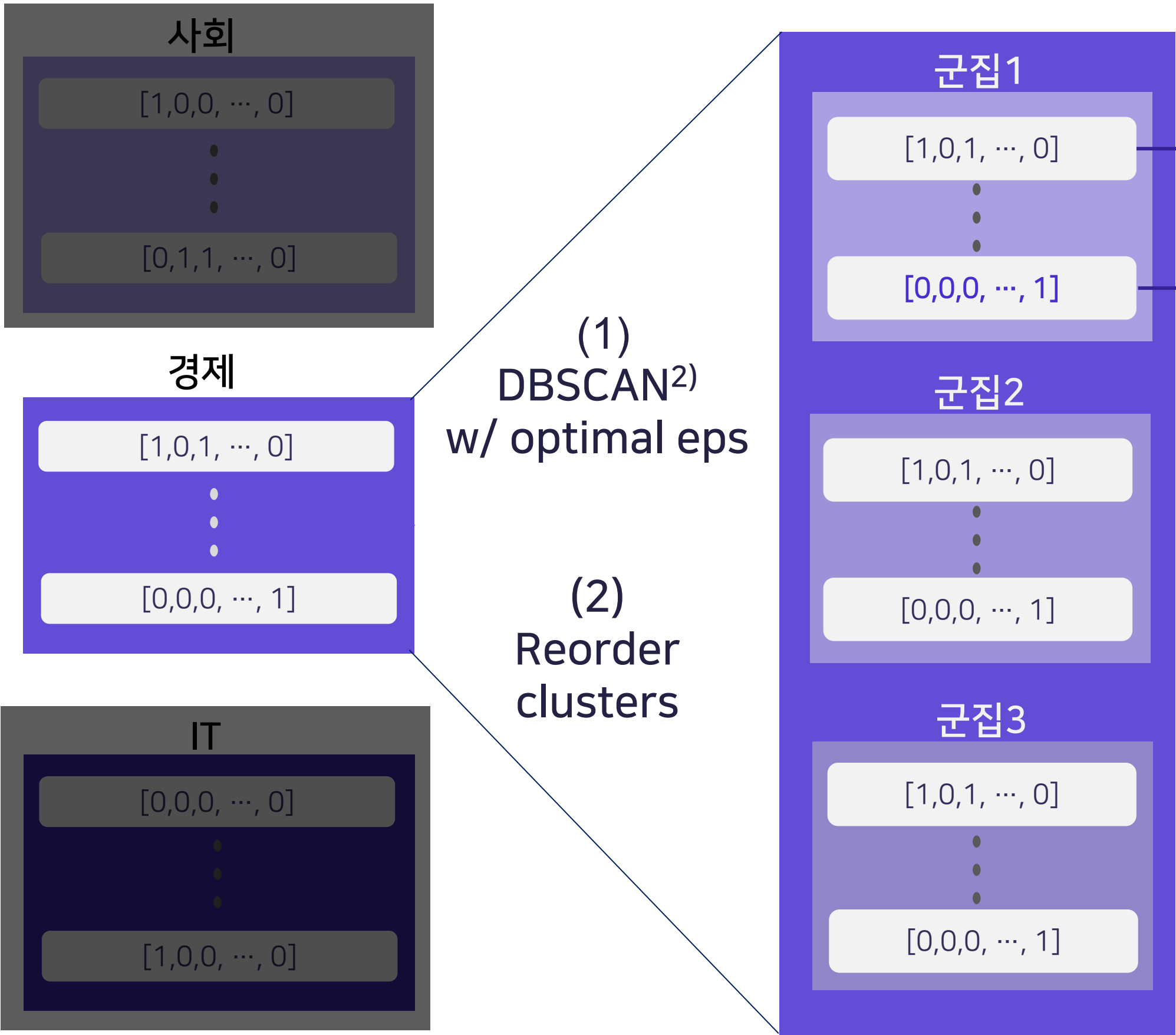
01 Crawling

02 Clustering

03 Summarization

04 Text to Speech

05 Serving



center \bar{V}_j

$$center_j = \bar{V}_j$$
$$= \frac{1}{N_j} \sum_{i=1}^{N_j} V_{ij}$$

$J = [1, \dots, k]$

$N = [N_1, \dots, N_k]$

featured article $_j$

$$\underset{x \in J}{argmin} dist(V_{xj}, \bar{V}_j)$$
$$\Rightarrow V_{xj}$$

top k keywords $_j$

$$\bar{V}_j = [x_1, \dots, x_M]$$
$$\Rightarrow [x_{top_1}, x_{top_2}, \dots]$$

2) Ester, Kriegel, et al. "A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise" (1996)

01 Crawling

02 Clustering

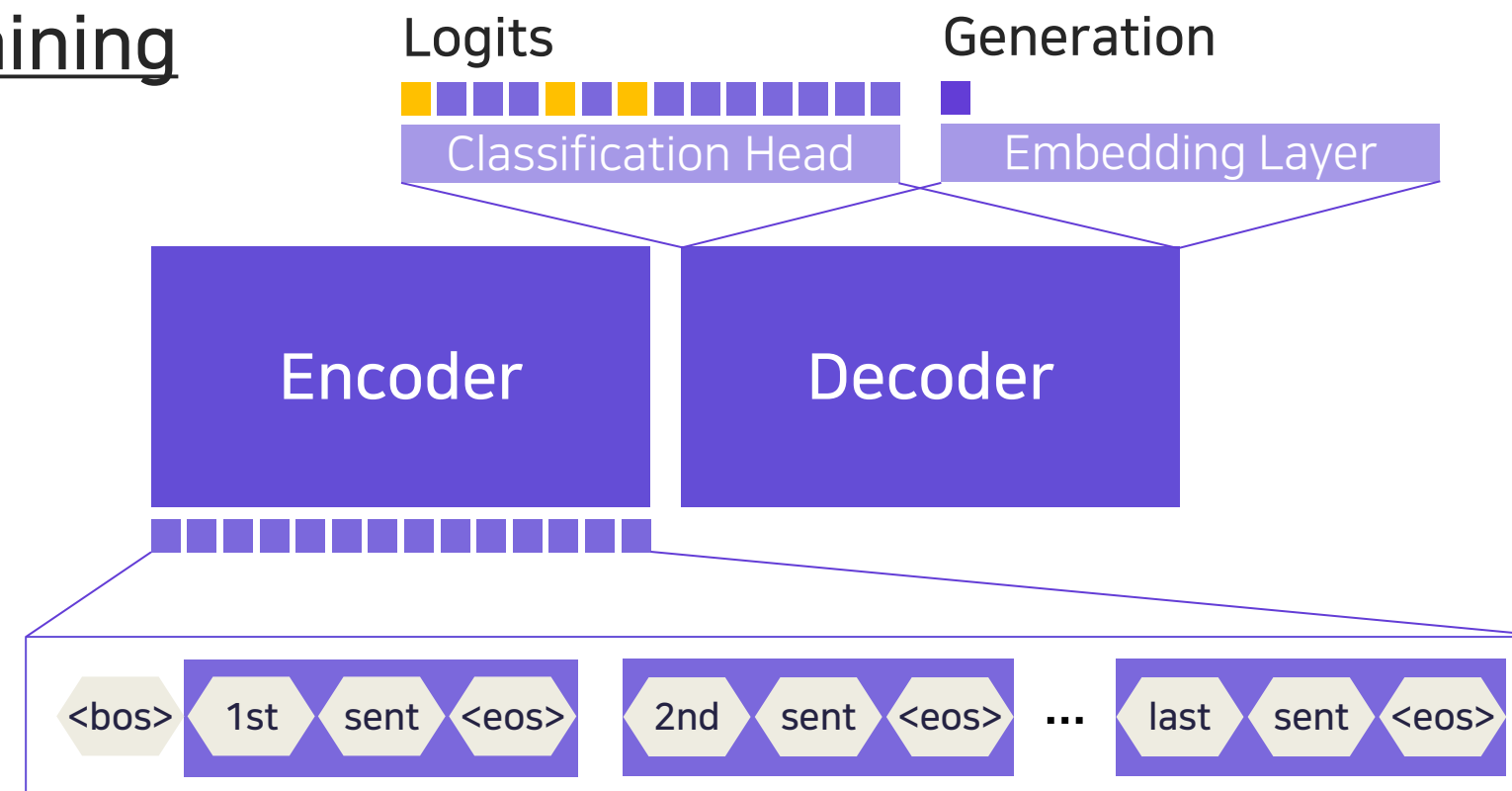
03 Summarization

04 Text to Speech

05 Serving

EasyBART – Extractive/Abstractive SummarY

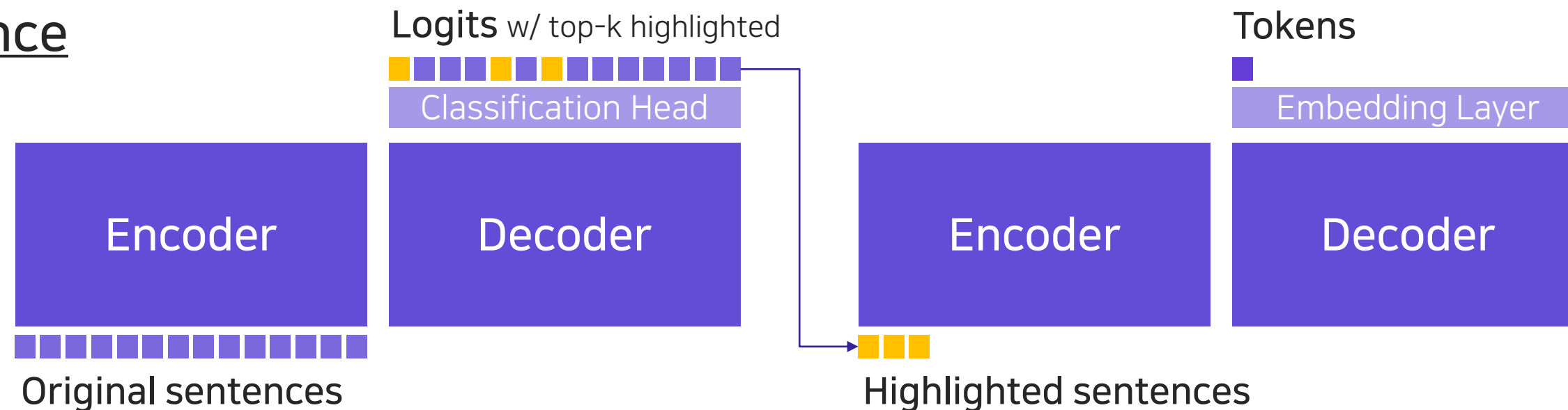
Training



“ One-Step Training & Two-Step Inference ”

- BART³⁾에 **custom classification head**를 추가하고, inference logic을 수정한 모델
- 추출/생성 요약 동시 진행 + 모델 구조와 훈련이 쉽다는 점에서 “**EasyBART**”로 명명
- Input 형태가 BERTSum⁴⁾과 유사하지만, `<sep>` 대신 `<eos>` 토큰을 활용

Inference



3) Lewis, Mike, et al. "Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension." *arXiv preprint arXiv:1910.13461* (2019).

4) Liu, Yang. "Fine-tune BERT for extractive summarization." *arXiv preprint arXiv:1903.10318* (2019).

01 Crawling

02 Clustering

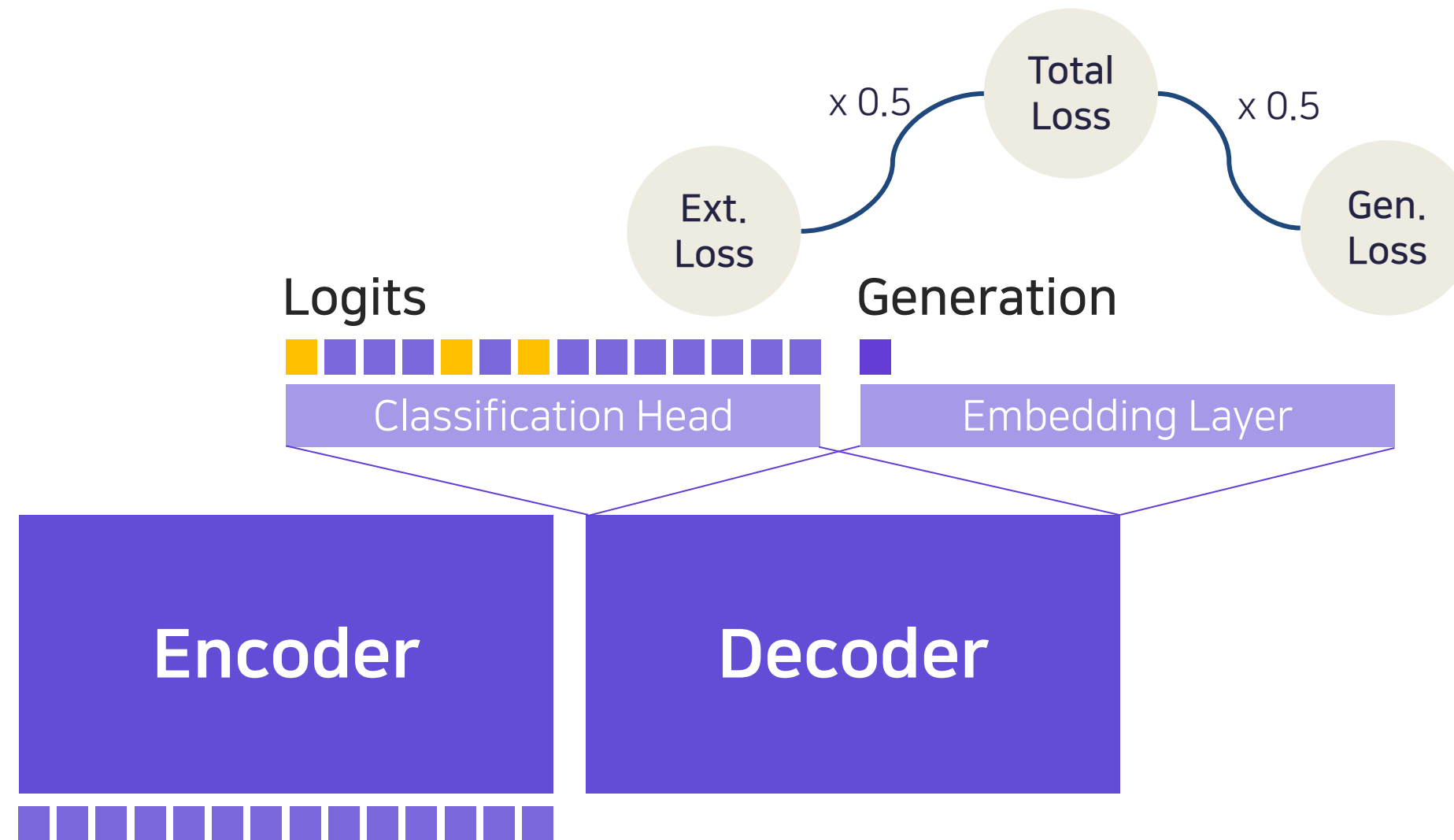
03 Summarization

04 Text to Speech

05 Serving

EasyBART – Extractive/Abstractive SummarY

Training



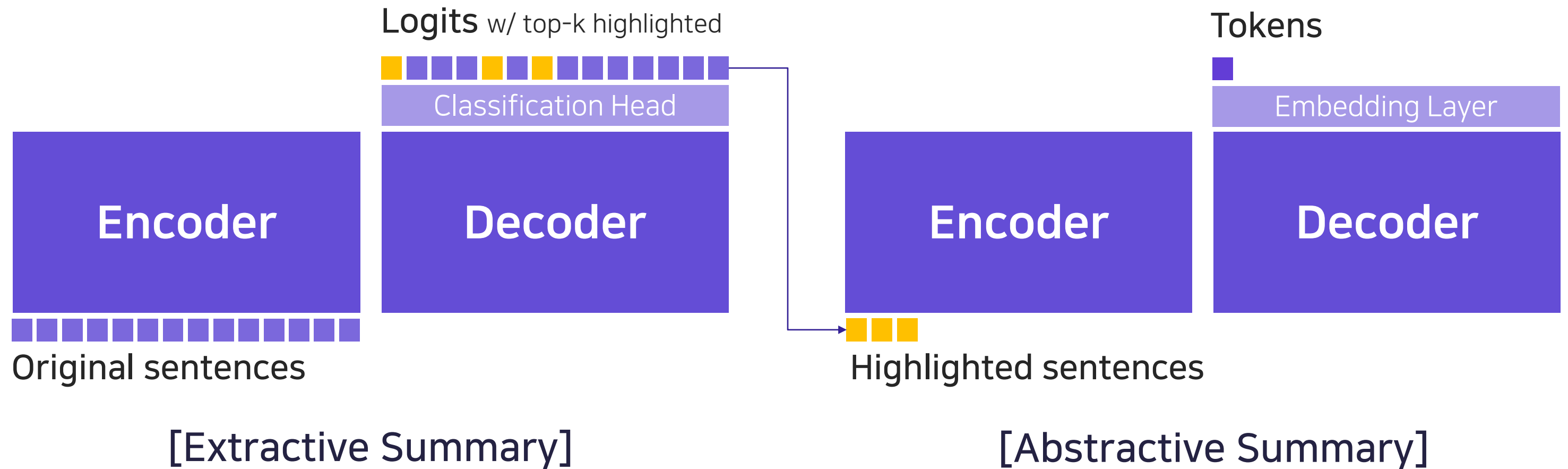
- 학습 데이터: AI Hub의 “문서 요약 텍스트”⁵⁾ 말뭉치 - 추출/생성요약 라벨 모두 존재
- 추출 요약문 기반의 생성 요약문 라벨은 존재하지 않으므로, 생성 요약 시에도 원문 그대로를 input으로 활용

5) 본 AI 데이터는 한국지능정보사회진흥원의 사업 결과입니다. <https://aihub.or.kr/aidata/8054>

- 01 Crawling
- 02 Clustering
- 03 Summarization**
- 04 Text to Speech
- 05 Serving

EasyBART – Extractive/Abstractive SummarY

Inference



- 선정된 기사 본문에서 주요 문장 top-k개 추출 후, 추출된 문장으로 요약문 생성
- 원문이 긴 경우에도 recursive하게 top-k를 반복적으로 추출하도록 개선 (서비스에는 미적용)

Example #1

01 Crawling

02 Clustering

03 Summarization

04 Text to Speech

05 Serving

[KoBART-summarization]

"김경남 소속사 제이알이엔티는 배우 김경남 측이 측간 소음으로 이웃에게 피해를 준 데 대해 17일 공식입장을 통해 먼저 좋지 않은 일로 심려를 끼쳐 죄송하다고 사과했다."

[EasyBART]

"김경남 소속사 제이알이엔티는 배우 김경남 측이 측간 소음으로 이웃에게 피해를 준 데 대해 17일 공식입장을 통해 먼저 좋지 않은 일로 심려를 끼쳐 죄송하다고 사과했고, 이날 저녁 김경남 배우가 당사자 분을 찾아가 이야기를 나눴다며 진심으로 사과드리고 앞으로는 더 주의하겠다고 말씀드렸다고 전했다."

01 Crawling

02 Clustering

03 Summarization

04 Text to Speech

05 Serving

EasyBART – Extractive/Abstractive SummarY

A/B Test 결과 분석

4가지 요약문에 대한 Friedman Test

	EasyBART			BART
	Top-3	Top-5	Top-7	baseline ⁶⁾
접수합*	1755.5	1629.5	1841.5	1673.5
순위합*	296.0	255.5	333.0	265.5

- $\chi^2 = 21.452$ ($df = 3$)
- $p < 0.005$

→ 4가지 요약문의 품질 차이가 존재

Top-3 vs. Baseline Sign Test

	EasyBART	BART	Total
	Top-3	baseline	
응답수*	610	310	920
비율	66.3%	33.7%	

- $\hat{p} = 0.663$
- $p < 0.005$

→ Top-3를 압도적으로 선호

“k에 따른 **readability**와 **information** 사이의 trade-off가 존재”

일종의 hyperparameter로, 도메인 및 목적에 따른 tuning 가능성
Top-5의 성능이 낮은 이유도 두 측면 모두 열등하기 때문

* 높을수록 긍정적인 평가를 의미함. 자세한 방법론은 Appendix 참조.
6) <https://huggingface.co/gogamza/kobart-summarization>

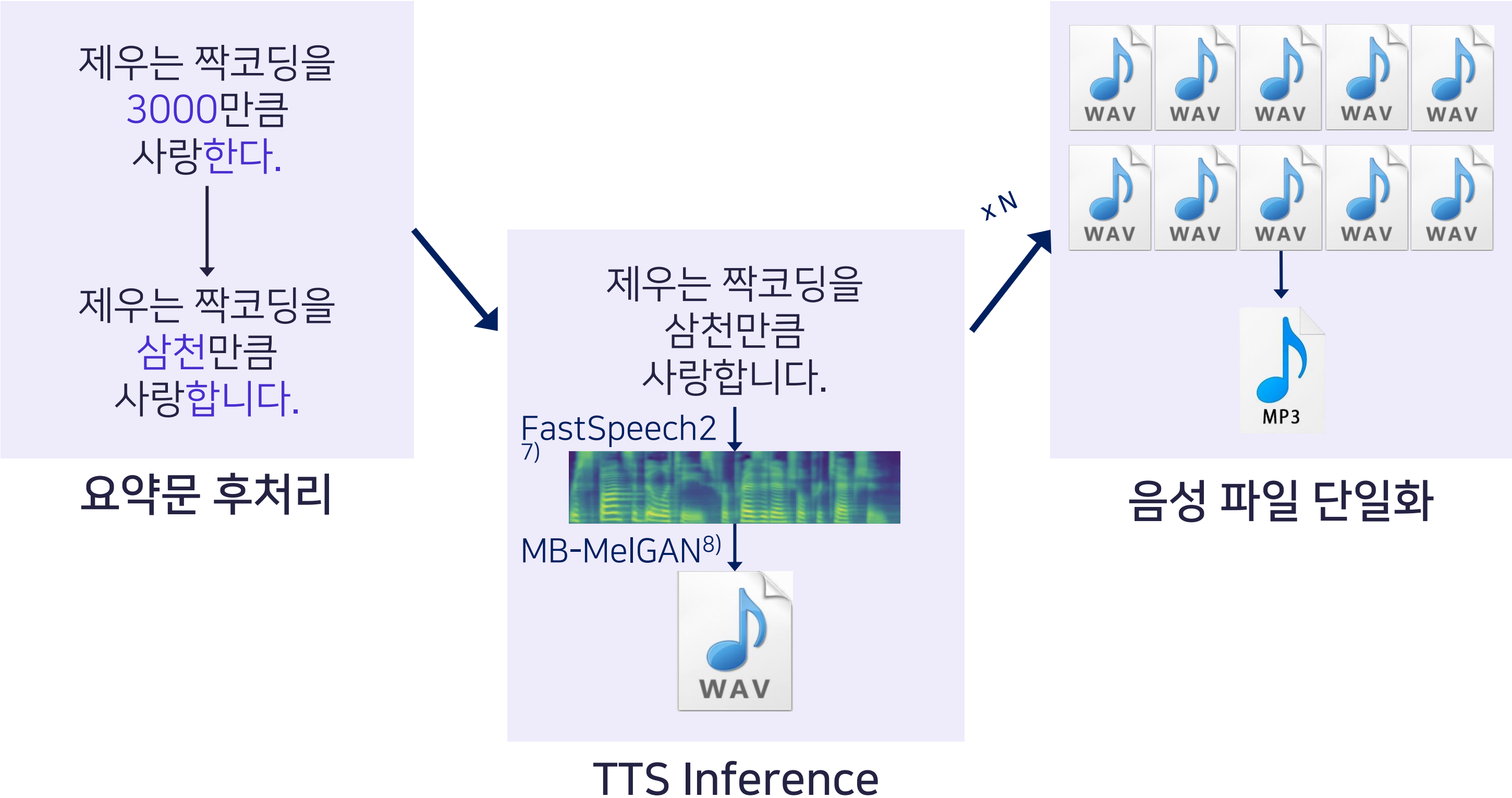
01 Crawling

02 Clustering

03 Summarization

04
Text to Speech

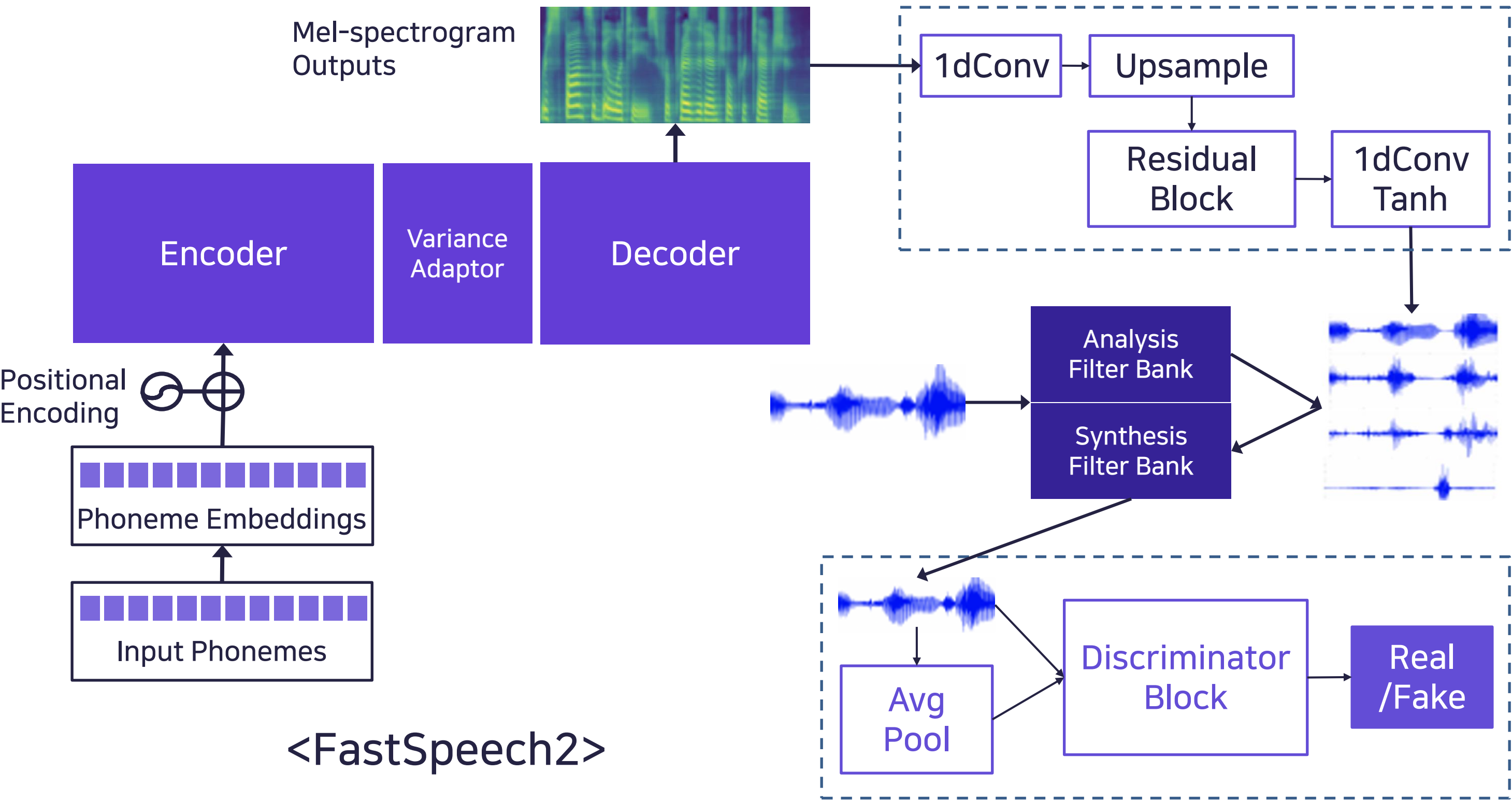
05 Serving



7) <https://huggingface.co/tensorspeech/tts-fastspeech2-kss-ko>
8) https://huggingface.co/tensorspeech/tts-mb_melgan-kss-ko

TTS 모델 구조

- 01 Crawling
- 02 Clustering
- 03 Summarization
- 04 Text to Speech**
- 05 Serving



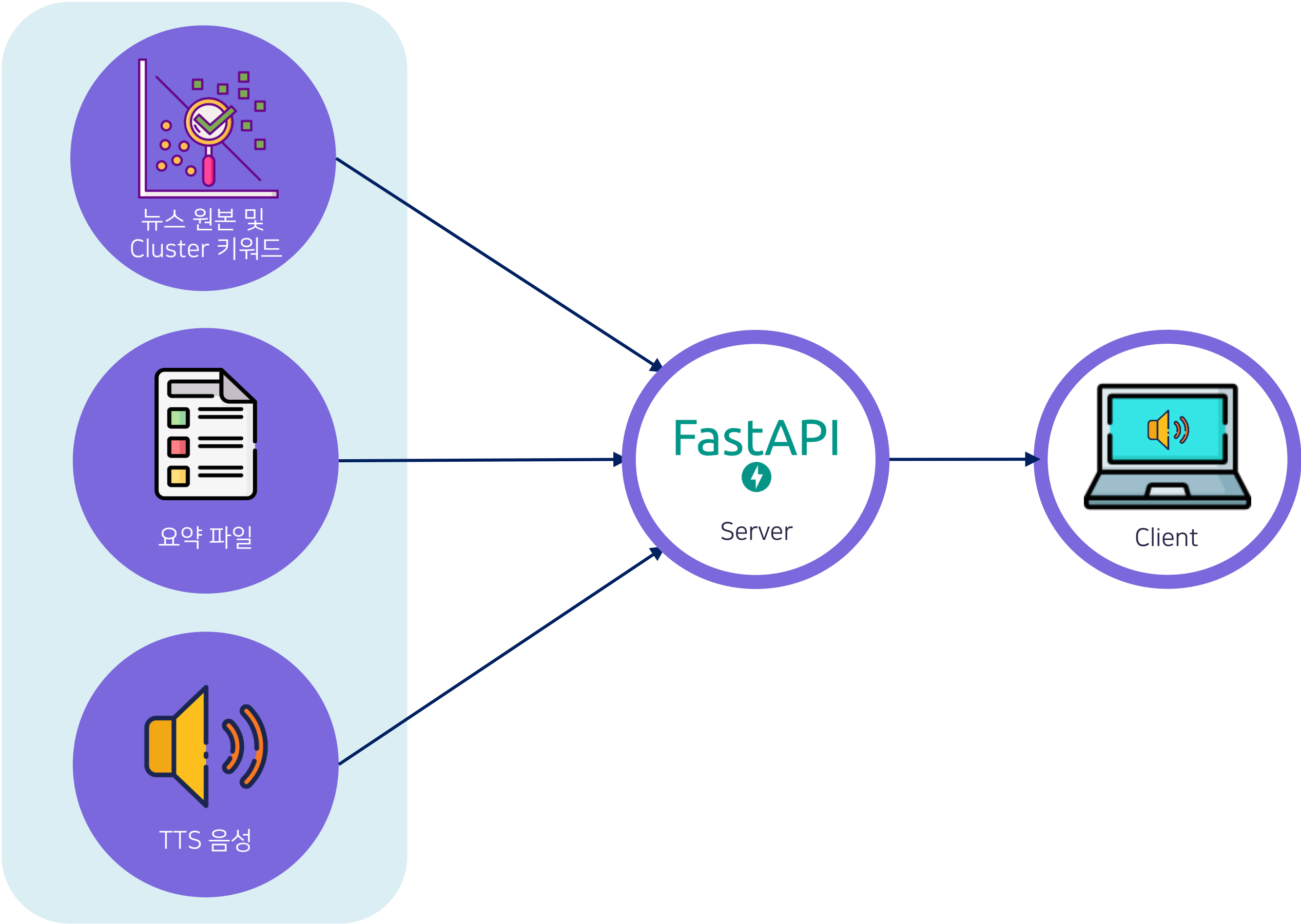
01 Crawling

02 Clustering

03 Summarization

04 Text to Speech

05
Serving



04

시연
영상

귀가노니 - 출퇴근길에 듣는 인공지능 뉴스 팟캐스트

귀가노니

출퇴근길에 듣는 인공지능 뉴스 팟캐스트

HomeCategory

[IT] 2021년 12월 23일

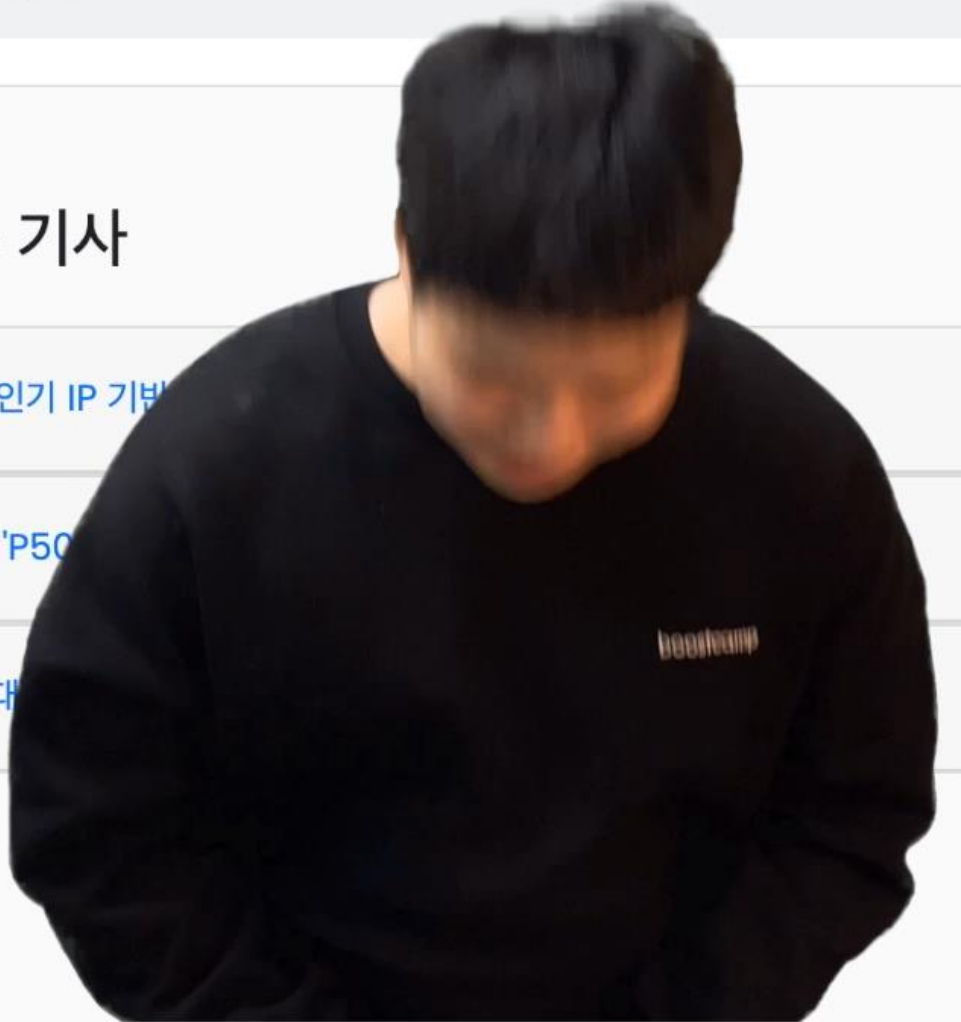
날짜 선택

▶ 0:00 / 1:04

🔊 ⋮

요약문 및 원본 기사

'1조 매출' 콘솔시장, 인기 IP 기반	아이뉴스24
화웨이, 새 폴더블폰 'P50	뉴스1
마비노기, '대교역시대	게임동아

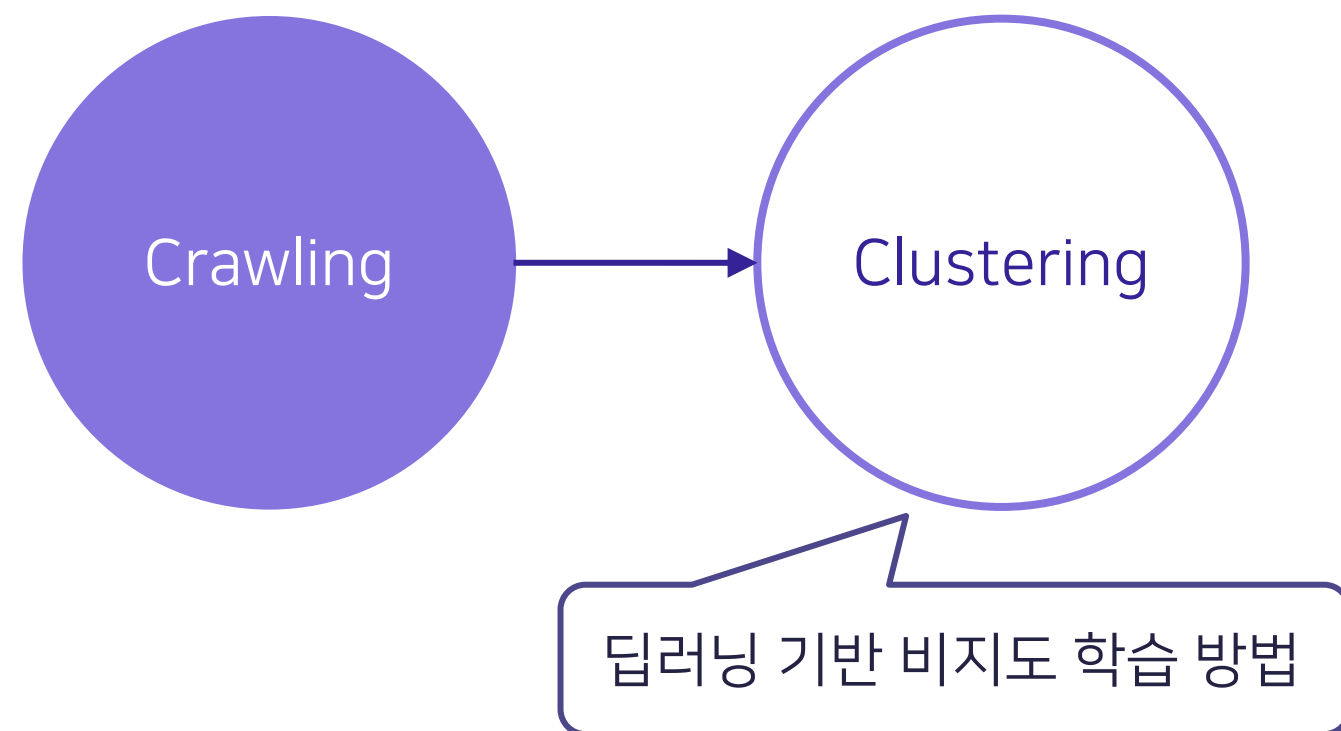


05

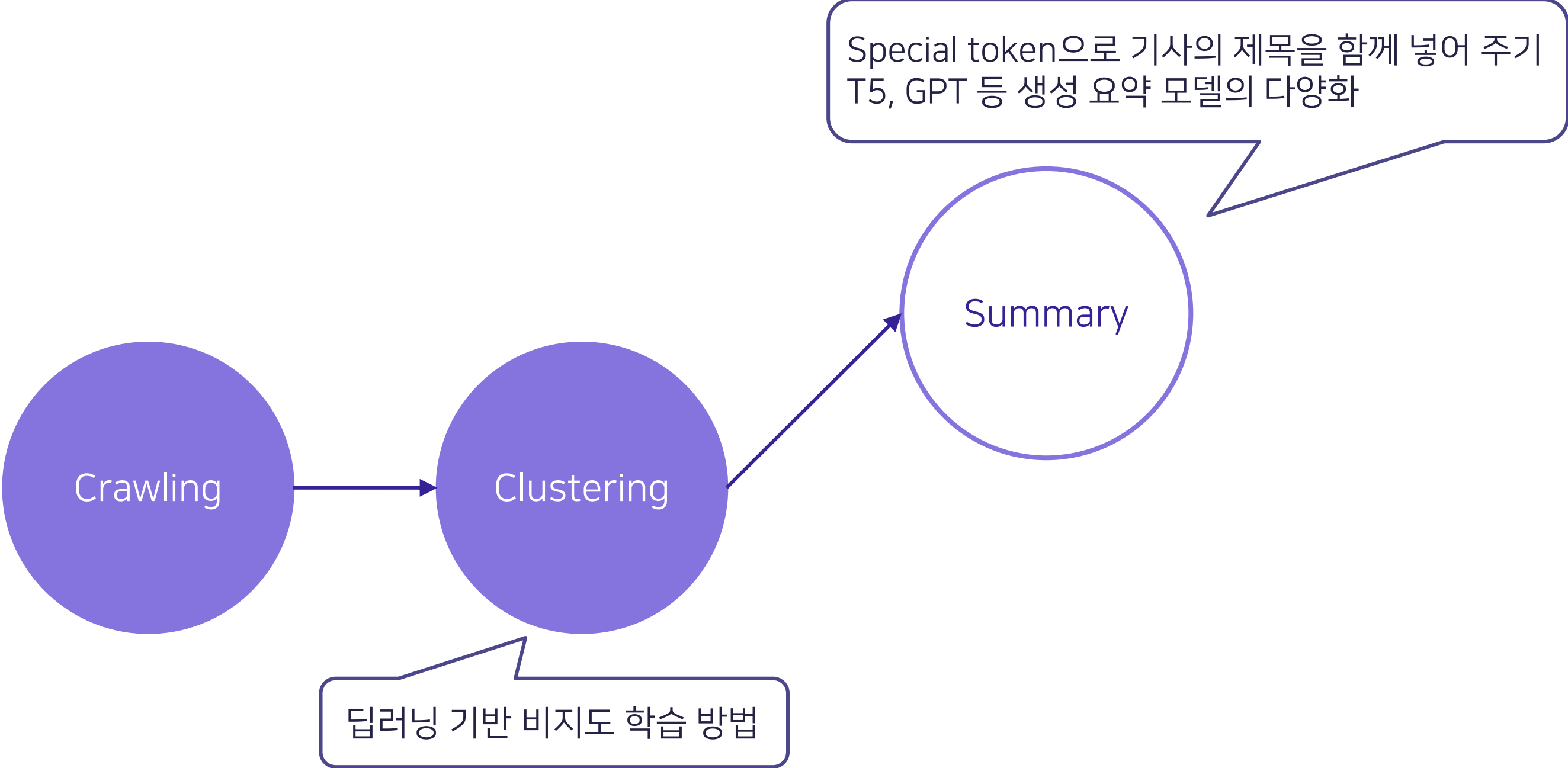
향후 개선사항

귀가노니 - 출퇴근길에 듣는 인공지능 뉴스 팟캐스트

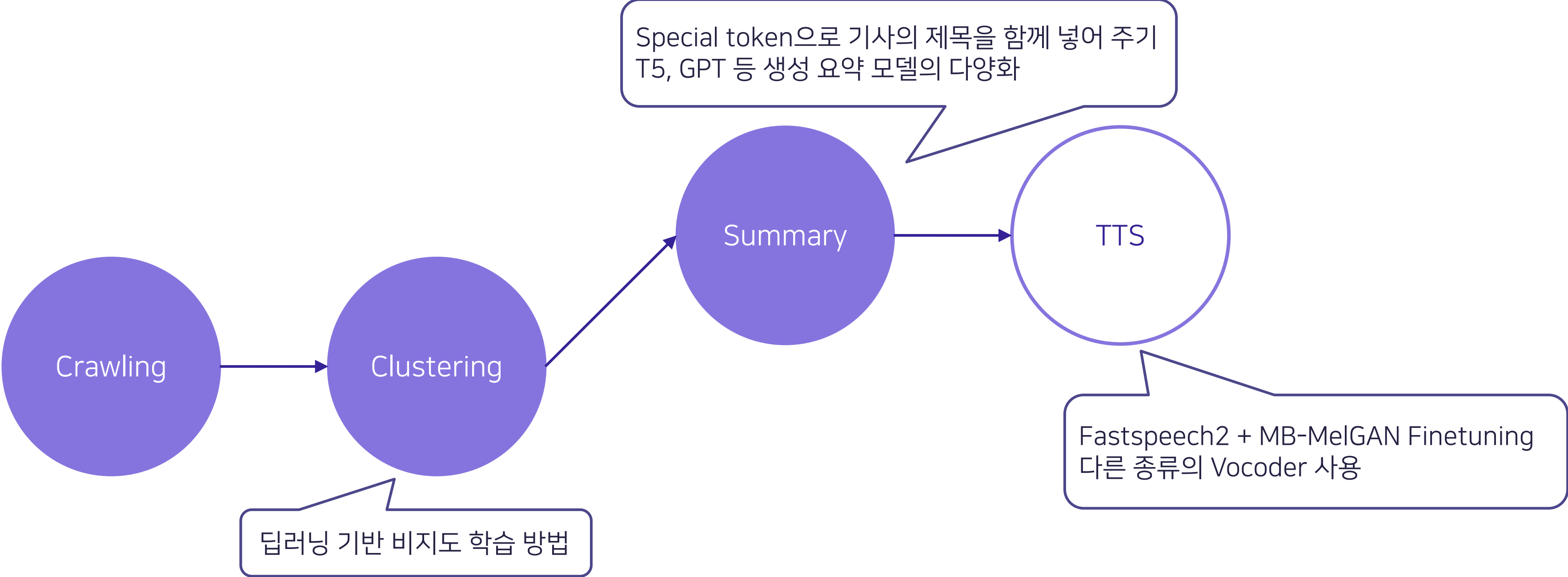
04 향후 개선 사항



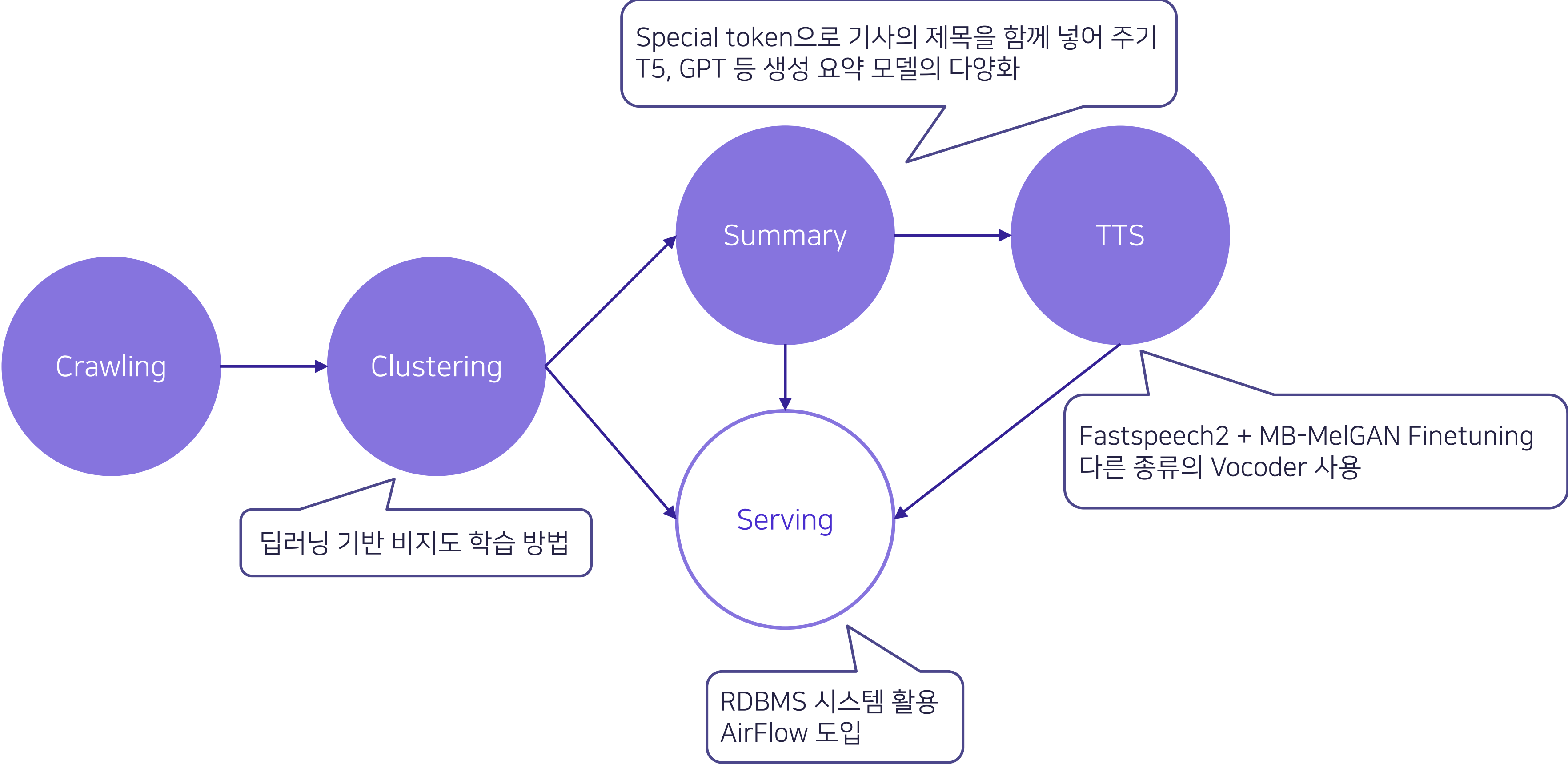
04 향후 개선 사항



04 향후 개선 사항



04 향후 개선 사항



06

Q & A

귀가노니 - 출퇴근길에 듣는 인공지능 뉴스 팟캐스트

감사합니다.

boostcamp^{ai}tech

Example #2: 유럽 에너지 대란에 산업계도 초비상.. “금속, 비료 생산 감축” (국제)

[KoBART-summarization]

유럽 천연가스 공급난으로 가격이 급등하면서 금속 제련, 비료 제조 관련 기업들은 수익성이 낮아져 생존까지 위협받고 있으며 블룸버그통신은 "유럽 천연가스 공급난으로 가격이 급등하면서 산업계 전체가 막대한 재정적 피해를 보고 있다"며 "금속 제련소와 비료 제조업체는 이미 생산량 감축에 들어갔다"고 보도했다.

[EasyBART]

22일(현지시간) 블룸버그통신은 "유럽 천연가스 공급난으로 가격이 급등하면서 산업계 전체가 막대한 재정적 피해를 보고 있다"며 "금속 제련소와 비료 제조업체는 이미 생산량 감축에 들어갔다"고 보도했다.

- ✓ 주로 학습된 도메인에서 벗어나 새로운 단어들이 등장하는 경우 반복되는 표현 억제됨

Example #3: 한은 “집값 폭등해 금융불균형 심화시 성장률 3퍼센트 추락” (사회)

[KoBART-summarization]

"한국은행은 국내외 금융불균형 상황에서 실물경제 하방리스크를 점검하기 위해 국내외 금융취약성지수를 활용한 GaR(최대성장감소율) 분석을 실시한 결과 국내 금융불균형 상황은 실물경제 하방리스크를 확대시키고, 특히 주요국 금융불균형을 감안할 경우 국내 실물경제 하방 리스크가 더 커지는 것으로 나타났다."

[EasyBART]

"한국은행이 국내외 금융불균형 상황에서 실물경제 하방리스크를 점검하기 위해 국내외 금융취약성지수를 활용한 GaR(최대성장감소율) 분석을 실시한 결과 부동산 가격 등 자산가격 상승에 따른 '금융불균형'이 심화될 경우 우리나라 경제성장률(GDP)이 최악의 경우 3퍼센트를 기록할 수 있다는 전망을 내놨다."

✓ 문법적(주술 관계)으로 더 나은 문장, 제목을 주지 않았음에도 원문의 핵심이 담김

Example #3: 한은 “집값 폭등해 금융불균형 심화시 성장률 3퍼센트 추락” (사회)

[원문]

한국은행이 부동산 가격 등 자산가격 상승에 따른 '금융불균형'이 심화될 경우 우리나라 경제성장률(GDP)이 최악의 경우 3퍼센트를 기록할 수 있다는 전망을 내놨다.

한국은행이 23일 발표한 '2021년 하반기 금융안정 보고서'에 따르면 한은은 국내외 금융불균형 상황에서 실물경제 하방리스크를 점검하기 위해 국내외 금융취약성지수를 활용한 GaR(최대성장감소율) 분석을 실시했다.

분석 결과 국내 금융불균형 상황은 실물경제 하방리스크를 확대시키고, 특히 주요국 금융불균형을 감안할 경우 국내 실물경제 하방 리스크가 더 커지는 것으로 나타났다.

한은 시산 결과 올 3분기 GDP는 국내 금융취약성지수만을 고려하면 10퍼센트의 확률로 1.4퍼센트를 기록하고, 여기에 주요국의 금융취약성지수를 추가하면 3.0퍼센트로 크게 낮아지는 것으로 나타났다.

국내 금융불균형은 최근 들어 일부 개선되고 있으나 과거에 비해서는 부동산 시장을 중심으로 크게 누중된 상황이다.

(...)

Example #4: KT, 고속 양자암호통신 독자 기술 개발 (IT)

[KoBART-summarization]

"22일 **KT는** 빛의 가장 작은 단위인 광자를 이용해 정보를 전달하는 차세대 통신 기술인 고속 양자암호통신용 **독자** 기술을 **독자적**으로 개발했다고 밝혔다."

[EasyBART]

"22일 **KT는** 20kbps(초당 킬로비트) 속도의 **고속 양자암호통신 기술**을 **독자적**으로 개발했다고 밝혔으며 정부의 디지털뉴딜 양자암호통신 속도 **기준(5bps)에 비해 4000배 빠르다.**"

- ✓ 주로 학습된 도메인에서 벗어나 새로운 단어들이 등장하는 경우 반복되는 표현 억제와 핵심을 잘 요약함

A/B Testing 시행 방법 소개

대조군: Top-3, Top-5, Top-7 (추출 + 생성 요약문)

비교군: KoBART-Summarization (생성 요약문)

- 12월 17일 실제 기사 12건의 요약문을 이용
- 위의 4가지 case 요약문이 동일한 경우를 제외하고 A/B Test 10문항 구성
- 카테고리(정치, 경제, 사회, 연예)와 요약문 case를 고르게 반영
- n-gram (n=1, 2, 3)을 이용하여 요약문 간 중첩도를 계산 → 가장 멀리 떨어진 두 요약문을 A, B로 선정
- 이후 n-gram 중첩도와 사람이 판단했을 때 유사한 수준을 고려 → 나머지 2개의 문장도 A, B 그룹에 속하도록 함

Coding Scheme & 검정 방법

Coding Scheme

- 요약문들이 비슷하기 때문에 최소한의 문항 수를 통해 최대한의 정보를 얻도록 설계
- 요약문1과 요약문2가 동일한 집단에 속한다면 각각 0.5, 0.5점을 부여 (우열이 없으므로)
- 요약문1과 요약문2가 다른 집단에 속한다면 응답에 따라 1, 0점 혹은 0, 1점을 부여
- 따라서, 랜덤하게 응답했다면 각 요약문의 기댓값이 동일함

검정 방법: Friedman Test, Quade Test

- Dependent case에 사용되는 nonparametric statistics 방법론
- Friedman Test: 개인이 block으로 사용되어 순위를 매기므로 개인별 차이가 0이 됨
- Quade Test: 개인 내 평가 차이에 대해서도 순위를 매겨, 차이가 크고 적음을 반영하여 더 검정력이 큼
- 정규분포 가정이 필요 없으므로 위의 같은 데이터 형태에 적합

Coding Scheme 도식화



응답자가 응답 A를 선택했을 때,

- “응답 A”에 속한 요약문
- “응답 B”에 속한 요약문

고려사항 & 한계점

문항 작성 시 고려사항

- 설문조사가 온라인 익명으로 진행되고 긴 요약문을 읽어야하는 특성을 고려
- 응답률을 높이기 위해 문항 수를 최대한 적지만 많은 정보를 추출할 수 있도록 설계함

한계점 및 논의

- 기사 원문을 제공하지 않아 의미적/맥락적 요소를 살피기 어려움
- 많은 응답이 문법적 요소 등은 고려되지 않음
- “요약”문보다는 “생성된” 문장 자체의 품질을 평가하는 문항

Friedman Test & Quade Test

4가지 요약문에 대한 Friedman test 결과

- Friedman Chi-Squared = 21.452
- Degree of Freedom = 3
- p-value < 0.005

4가지 요약문에 대한 Quade test 결과

- Quade F = 7.9312
- df1 = 3, df2 = 342
- p-value < 0.005

	EasyBART			BART
	Top-3	Top-5	Top-7	baseline
접수합*	1755.5	1629.5	1841.5	1673.5
순위합*	296.0	255.5	333.0	265.5

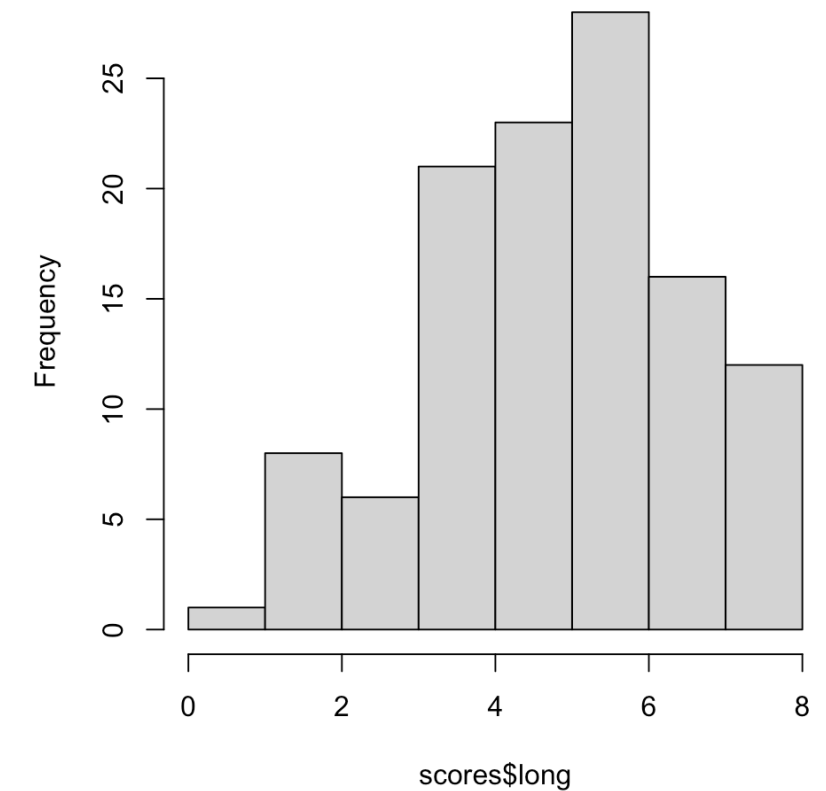
- $\chi^2 = 21.452 (df = 3)$
- $p < 0.005$

따라서, 4가지 요약문에 대한 품질 차이가 존재함

Binomial Test

Top-3 vs. Baseline

- 8개 응답은 Top-3와 Baseline을 비교할 수 있는 문항
- 920개(115 * 8) 응답에 대해 총 610개 응답이 Top-3를 선택
 - $\hat{p} = 0.663$
- Binomial Test 결과 p-value < 0.005으로 유의한 차이를 보임 -> Top-3 선호



긴 요약문 vs. 짧은 요약문

- A/B로 제시된 요약문 중 긴 것을 택한 경우를 평가
- 전체 1150개 중 526 응답
 - $\hat{p} = 0.457$
- Binomial test 결과 p-value < 0.005 로 유의한 차이를 보임 -> 짧은 요약문 선호

