

딥러닝을 활용한 실시간 보이스피싱 탐지 서비스

2018. 11. 21



과학기술정보통신부

NIA 한국정보화진흥원



IBK 기업은행



Contents

I 사업개요

II 서비스 소개

III 기술 소개

IV 향후 발전 방향

V Q & A



Contents

I



사업개요

II

서비스 소개

III

기술 소개

IV

향후 발전 방향

V

Q & A

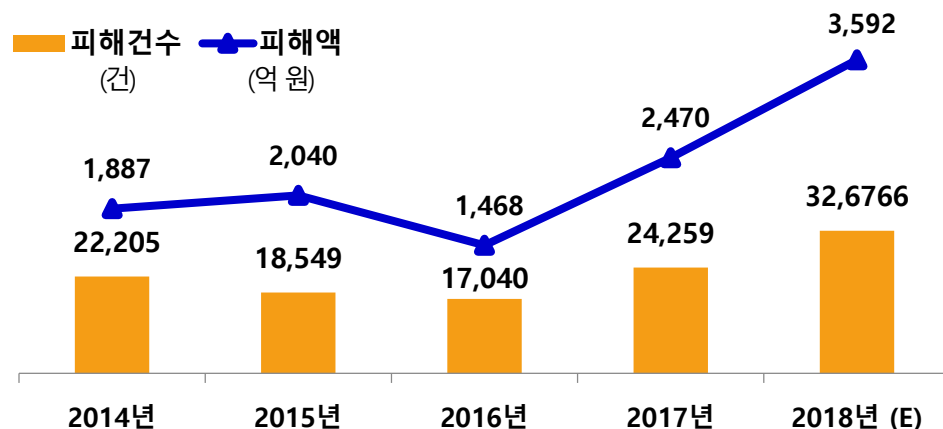


정우·일면 스님, 조계종 총무원장 후보 사퇴

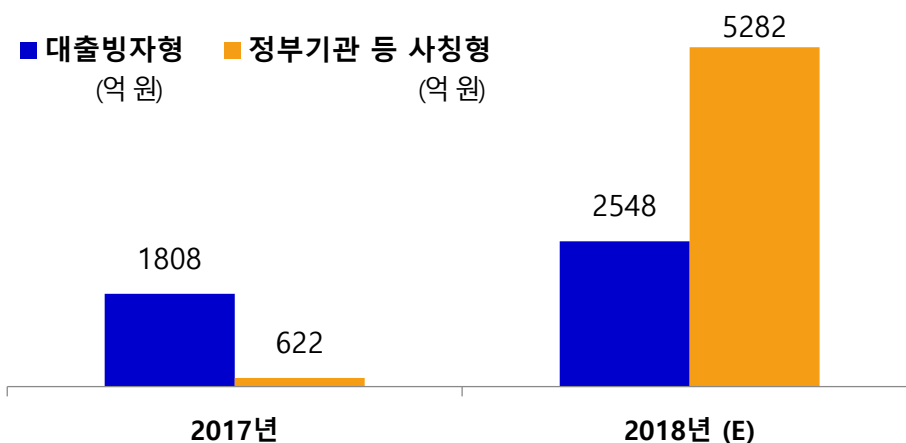
▶ 속리산 에밀레 박물관 S&P500 2905.97 ▼ 9.59

5년간 보이스피싱 피해액은 1,116억 1,860만 원이며, 발생 중
 성별로는 여성 피해액이 남성의 2.4배에 달하며, 전 연령대에서 고루 발생한 것으로 조사됨

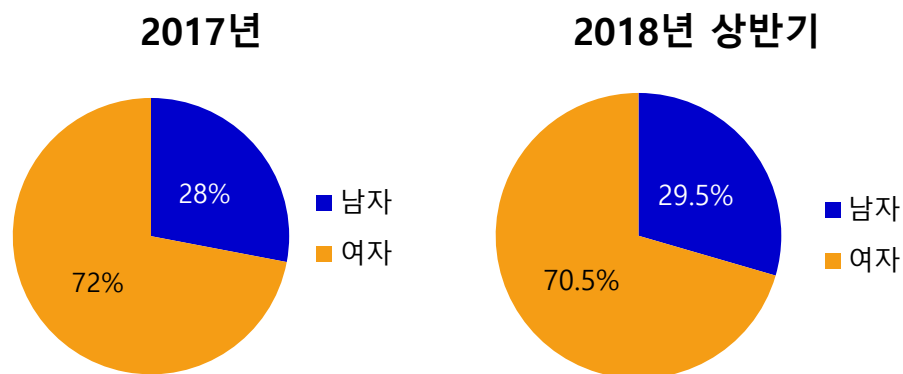
년간 보이스피싱 피해 상황 (출처 : 경찰청)



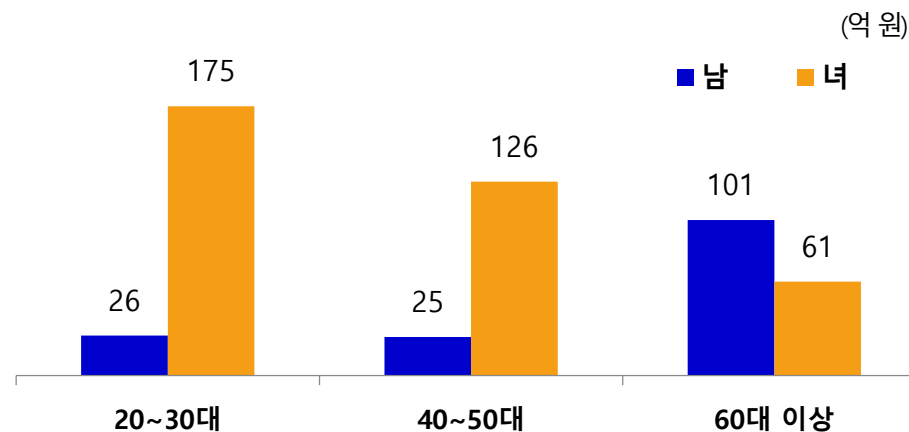
년간 보이스피싱 유형별 피해액 현황 (출처 : 경찰청)



보이스피싱 남/녀 피해 비율 (출처 : 경찰청)



'18년 성별/연령별 보이스피싱 피해액 (출처 : 경찰청)



현재 통화 App산업은 T전화, 후후, 후스콜 3사가 가장

하고 있음

이는 통화 App 시장에

됨. 또한, 기업 이미지 제고 및 광고효과를 가지고 있음



후스콜

제작사 : gogolook (Naver)

출시일 : 2010. 08

다운로드 수 : 5,500만 건 (2015)

비고 : 전세계 31개국 진출

2016 구글 최고 App 선정



후후

제작사 : KT CS (KT)

출시일 : 2013. 08

다운로드 수 : 3,500만 건 (2017)

비고 : 국내 사용률 압도적 1위

711만명, 66.5% 사용



T전화

제작사 : SK Telecom (SKT)

출시일 : 2014. 02

다운로드 수 : 2,500만 건 (2017)

비고 : 어플 선택제로 사용자 多

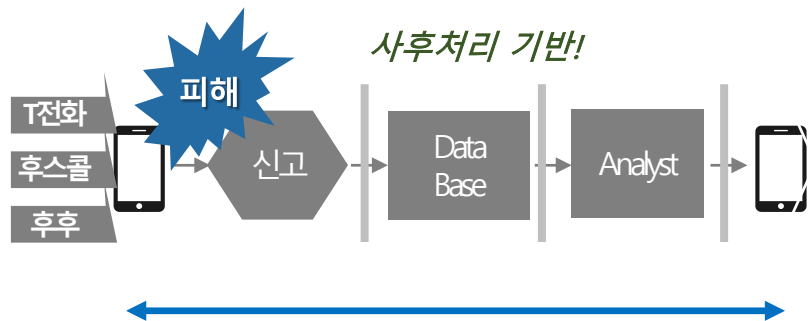
현행 통화 App은

이 맞춰져 있으며, 실시간으로

「IBK 실시간 보이스피싱 AI탐지 서비스」는 음성인식 및 딥러닝 기술을 활용하여 실시간

AS-IS

현재는 등록된 Data 기반으로만 동작하므로
실시간 보이스피싱 탐지 불가



- 사용자 : 전화 수신
- [사후처리] 통화 후 보이스피싱 전화 신고
- Database에 신고내용 기록
- 분석 및 검증 작업 및 배포
- 신고된 번호 탐지 및 알림 (보이스피싱)

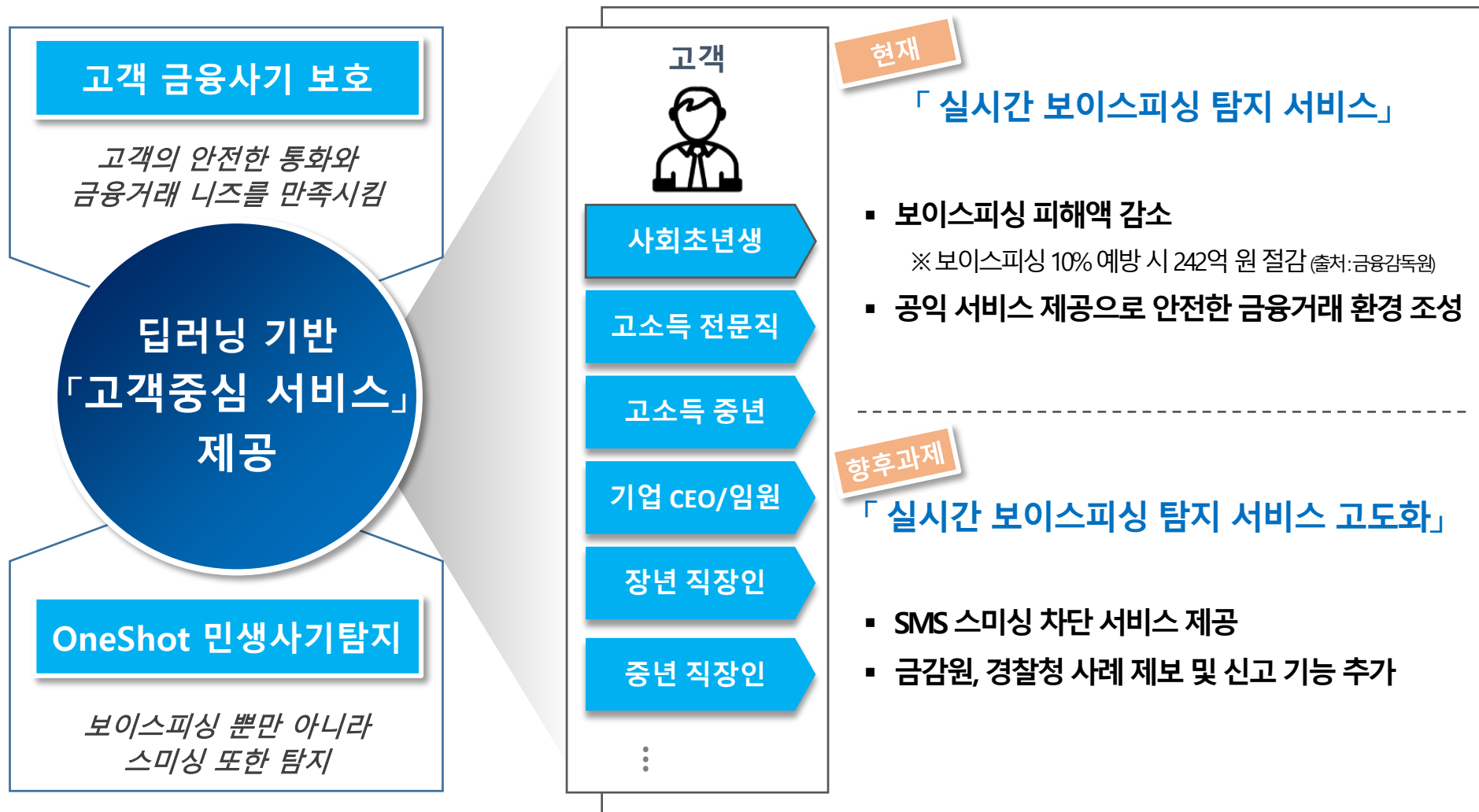
TO-BE

음성인식 및 딥러닝 기반 기술을 통한
실시간 보이스피싱 탐지 가능



- 사용자 : 전화 수신
- [실시간] 딥러닝을 통해 보이스피싱 확률 산출
- 딥러닝을 통해 실시간 보이스피싱 동시 400건 산출
- 통화 중, 보이스피싱 의심 경고음 및 진동 알림
- 기계학습을 활용한 보이스피싱 탐지 정확도 향상

「딥러닝을 활용한 보이스피싱 탐지 서비스」를 통해 **매일 변화하는 보이스피싱**으로부터 고객을 보호하고,
기술 고도화를 통해 보이스피싱을 넘어 SMS 스미싱 방지 등 딥러닝 서비스로



IBK, 금융감독원, 한국정보화진흥원 등 협업을 통해 「딥러닝을 활용한 보이스피싱 탐지 서비스」 개발 사업의 성공 완수를 목표로 함. 이를 통해 **해외 판로개척** 등 추진 목표를 가짐



금융감독원

IBK

한국정보화
진흥원

「딥러닝을 활용한 보이스피싱 탐지 서비스」
성공적인 사업 완료

국제 특허출원, 해외 판로 개척
민생사기탐지 Platform



Contents

I 사업개요

II ✖ 서비스 소개

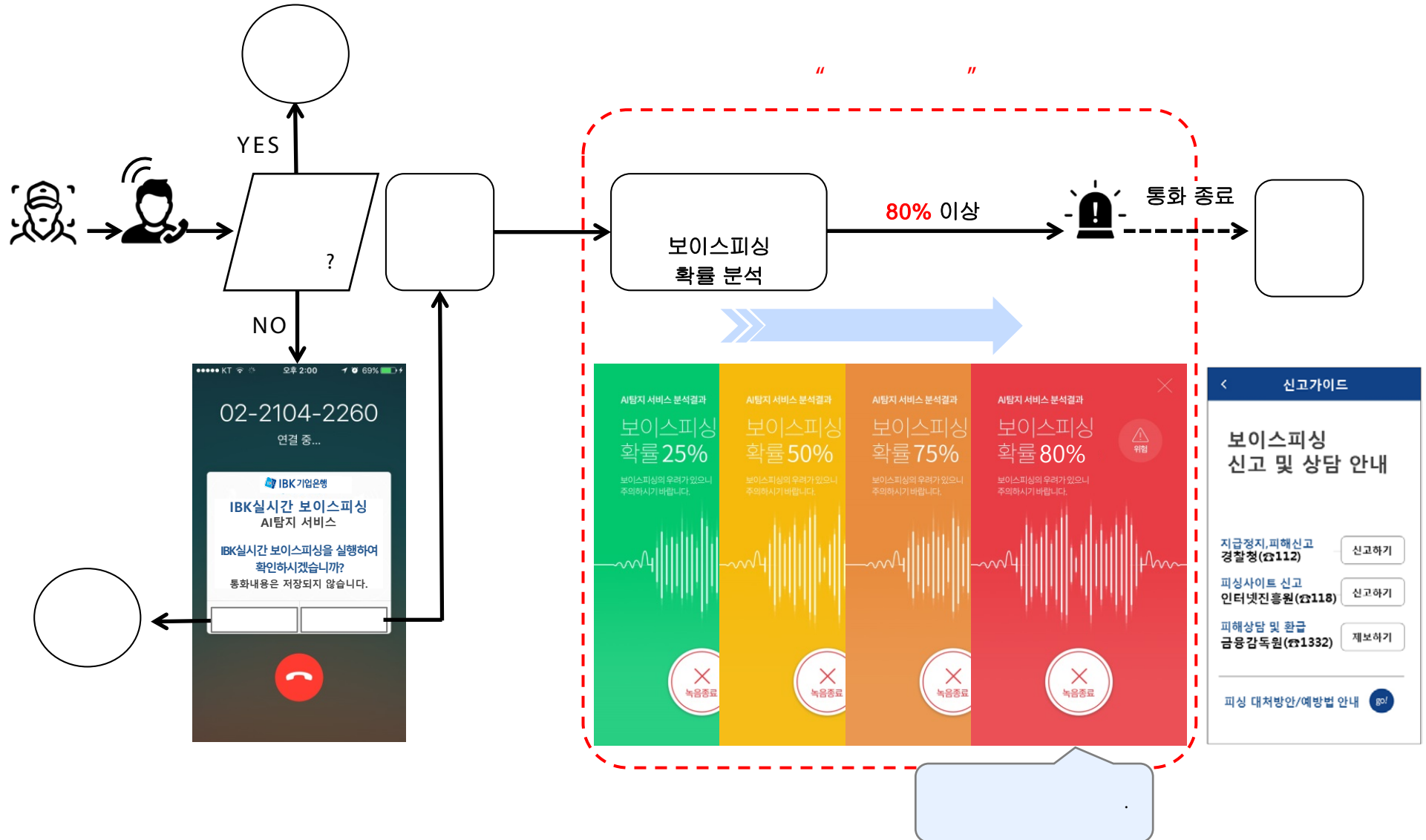
III 기술 소개

IV 향후 발전 방향

V Q & A



2016년 10월 24일 13시 36분
금융감독원에 제보된 **실제 사례**로
재연 해보았습니다.





Contents

I 사업개요

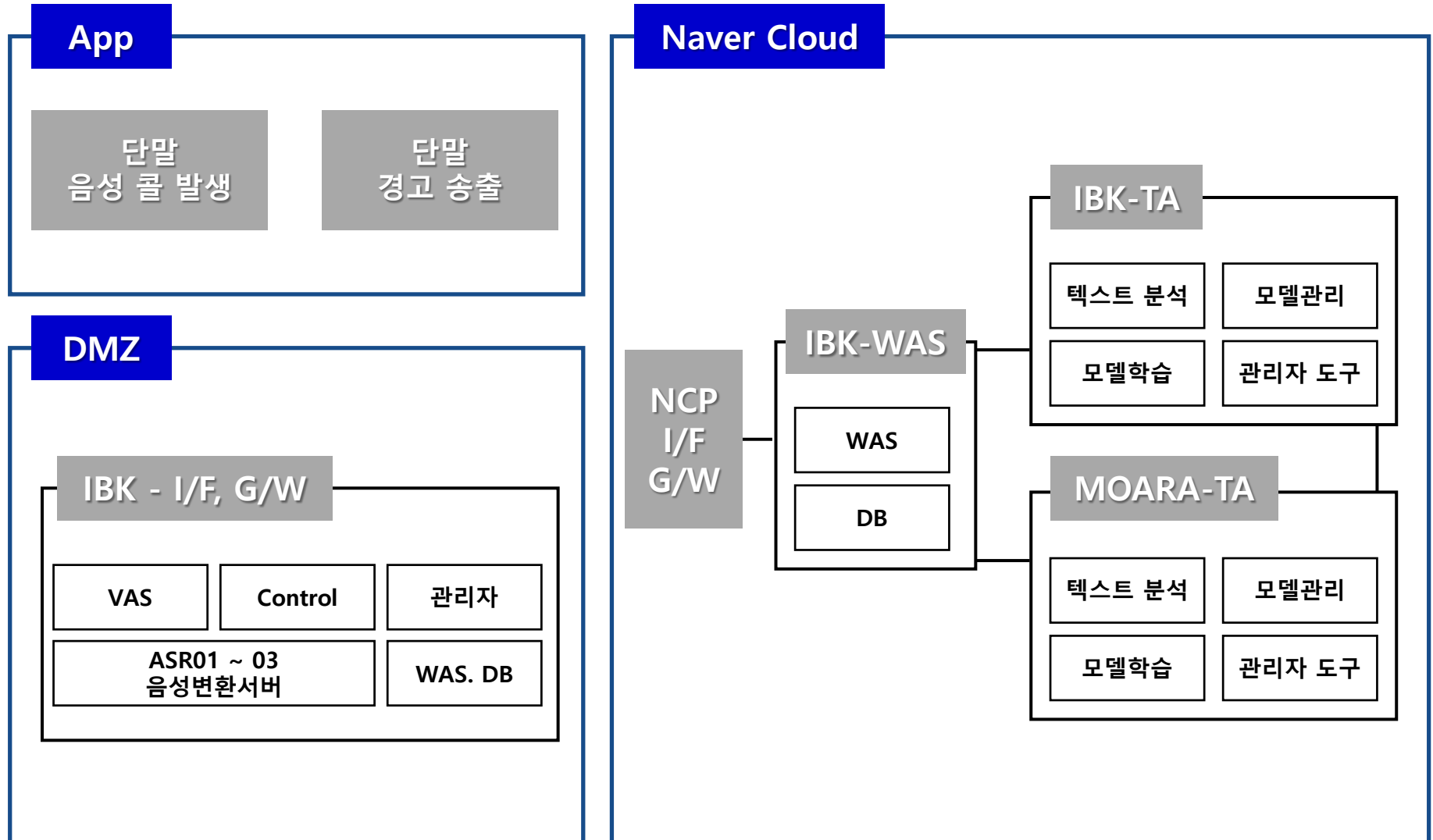
II 서비스 소개

III  기술 소개

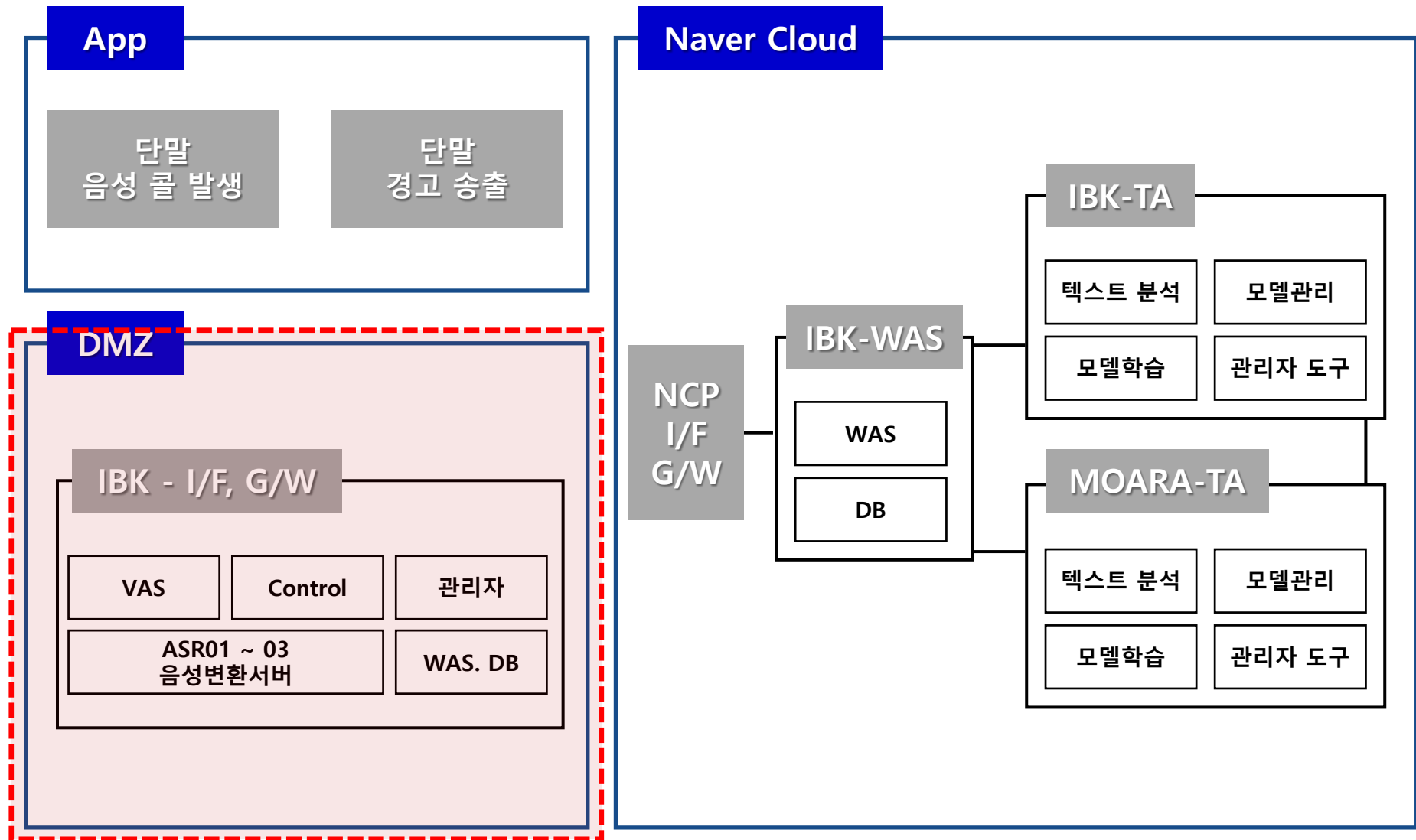
IV 향후 발전 방향

V Q & A

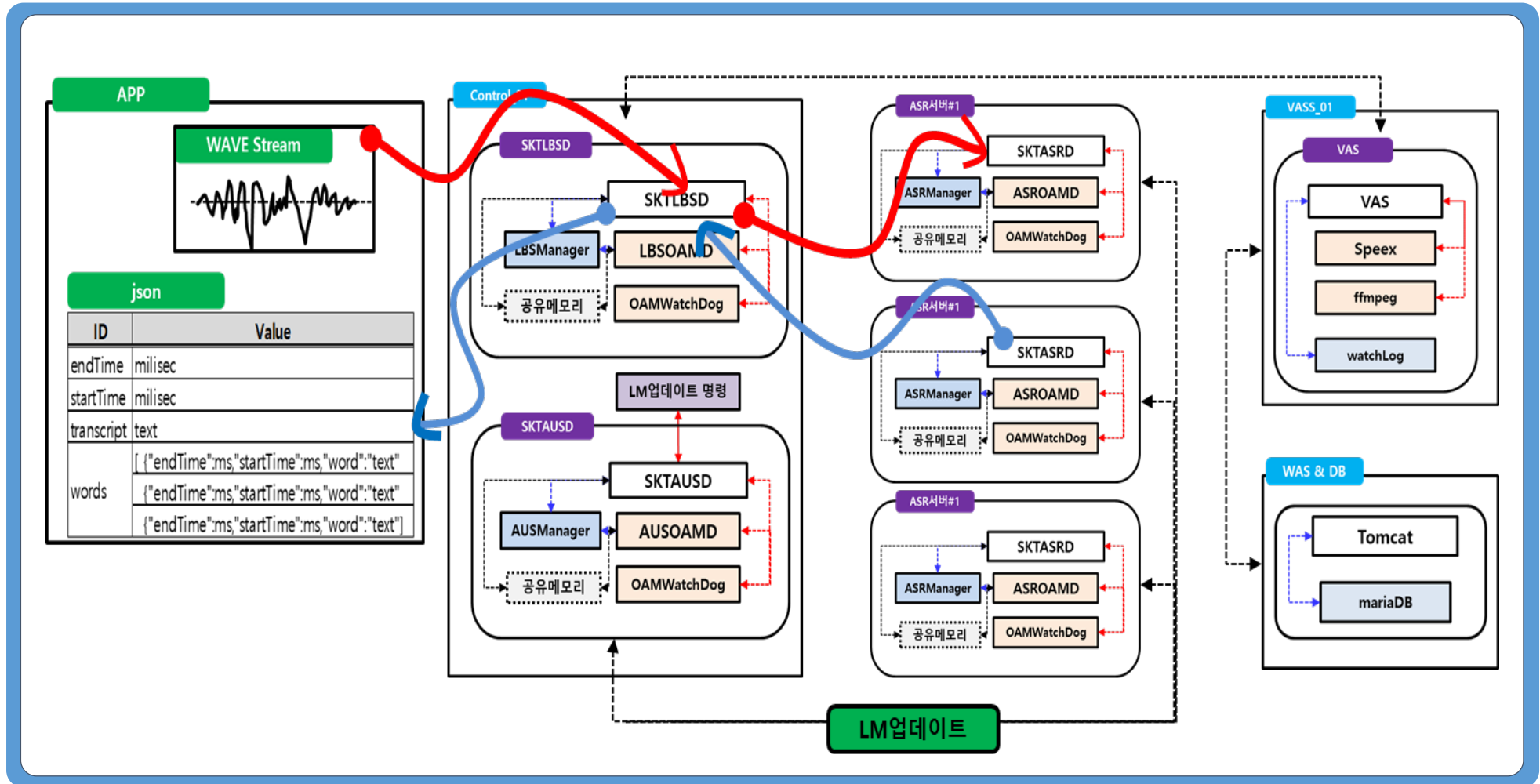
III 전체 시스템 구성



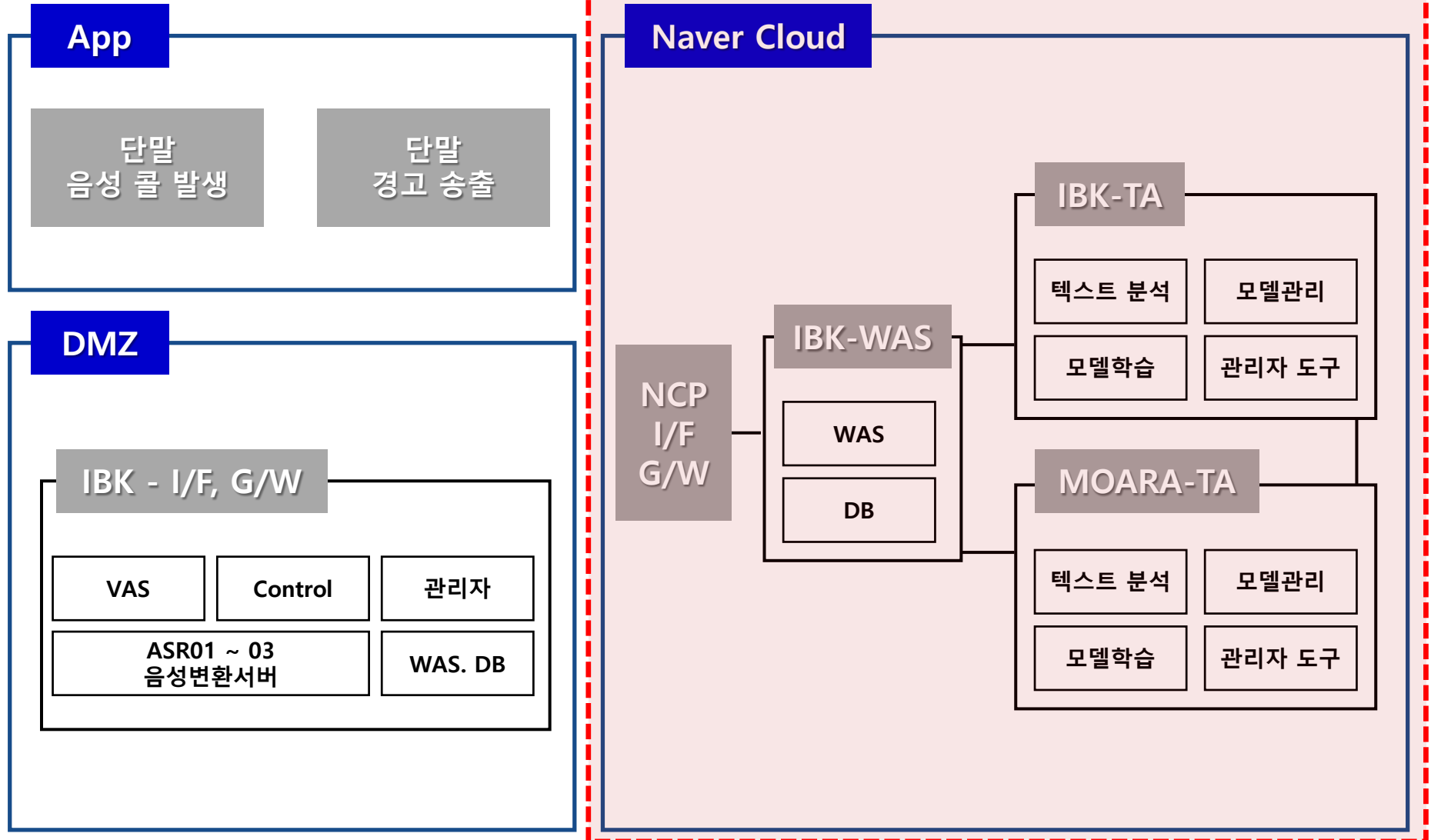
III 음성 인식



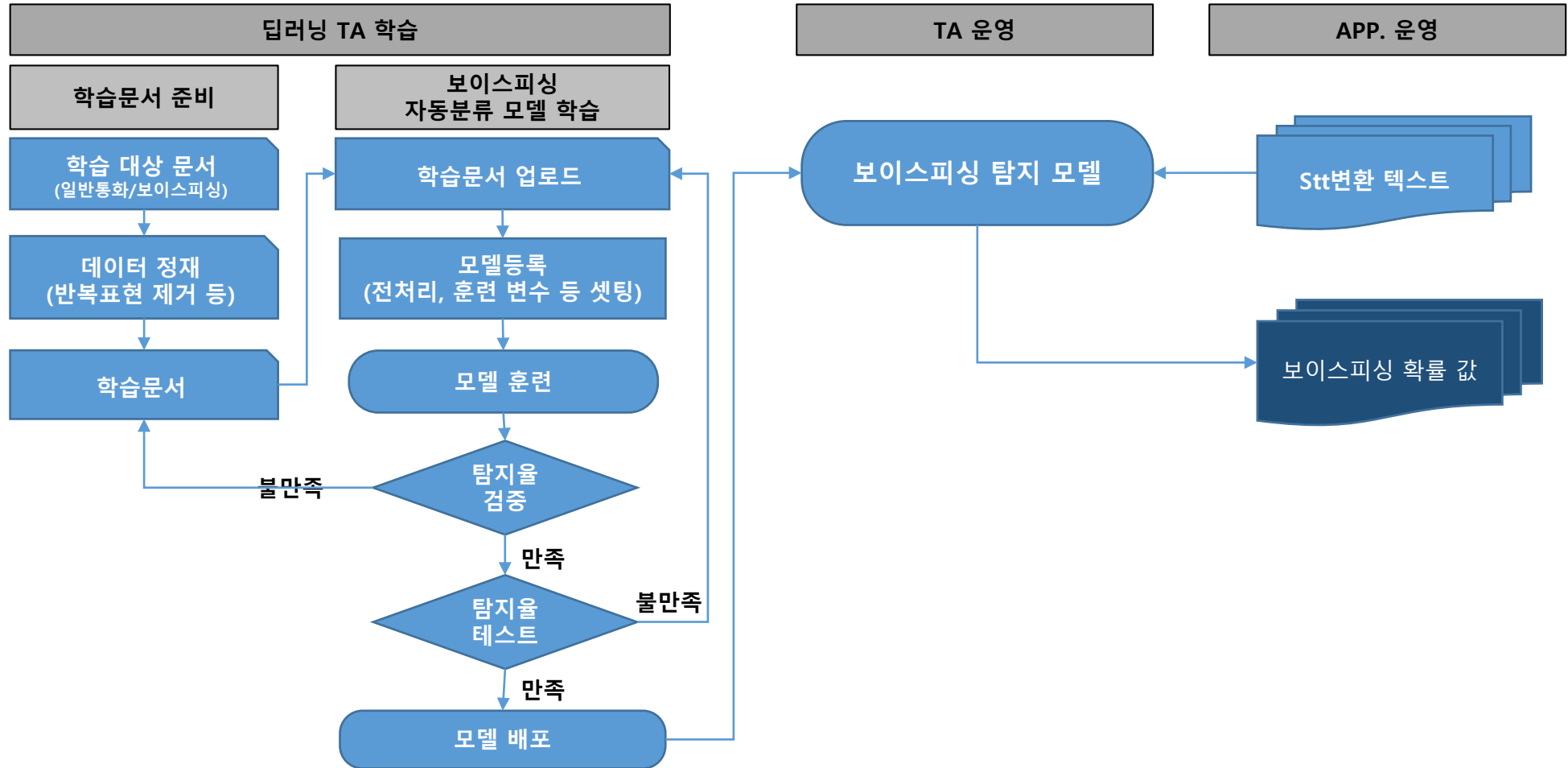
실시간 STT 시스템 흐름도



- 머신러닝을 기반으로 생성된 언어모델을 통해 음성을 문자로 실시간 변환
- SKT Nugu 음성인식 엔진 활용

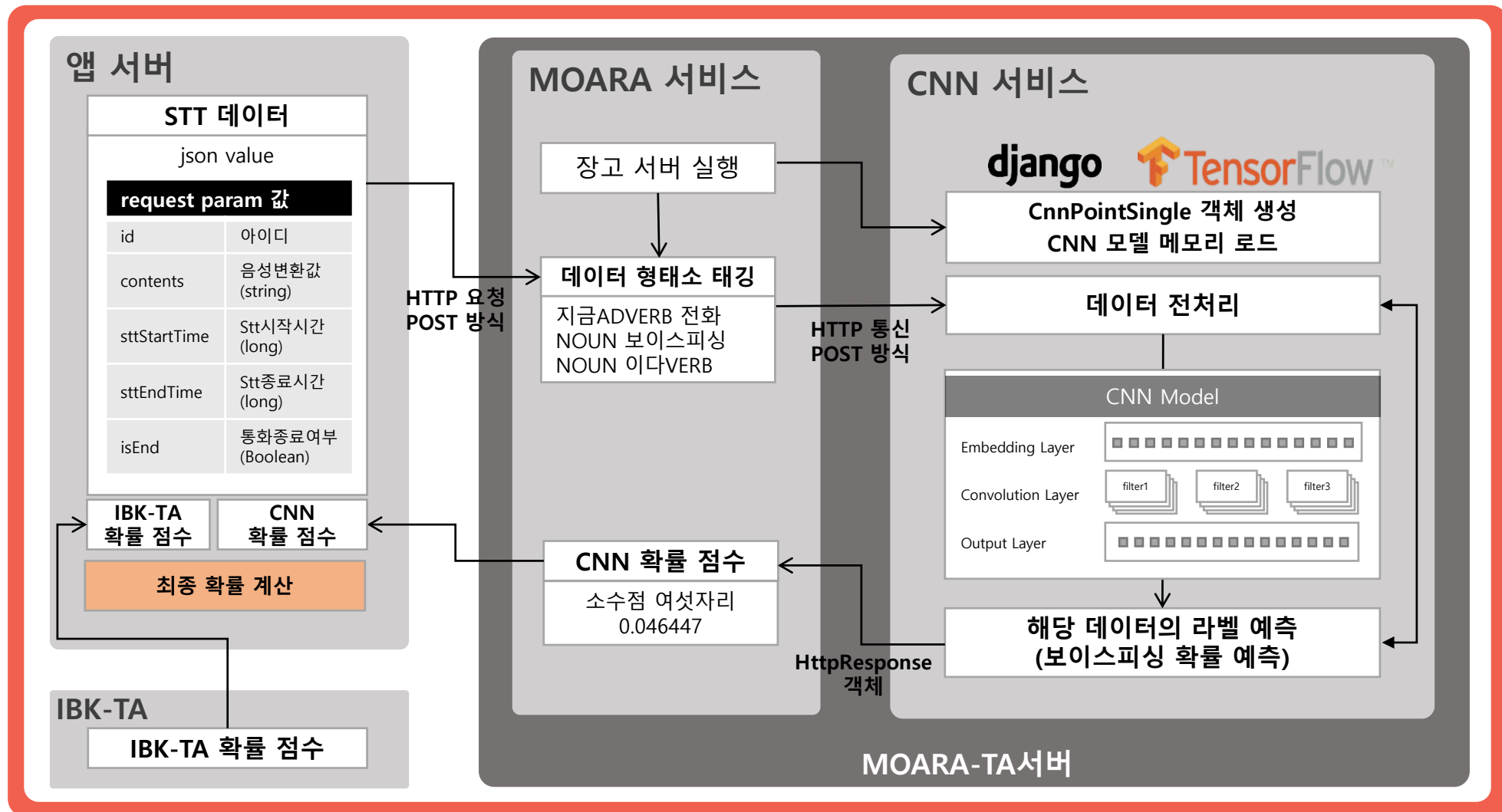


MOARA-TA 딥러닝을 활용한 보이스피싱 탐지 학습 및 운영 흐름도



- 딥러닝을 활용한 보이스 피싱 탐지 모델 개발 (CNN 알고리즘 적용)
- 딥러닝을 활용하여 보이스피싱 모델을 지속적으로 학습하고 관리 할 수 있는 관리 시스템 개발

MOARA-TA 딥러닝을 활용한 보이스피싱 탐지 시스템 흐름도



MOARA-TA 딥러닝(CNN)을 활용한 보이스피싱 탐지 모델 학습 결과 (4차 ~ 6차 학습)

Test Data 글자수	2000자	3000자	4000자
200	78.79%	85.02%	85.86%
400	83.84%	89.56%	90.07%
600	86.53%	91.08%	90.24%
800	85.69%	91.58%	91.58%
1000	87.37%	92.42%	91.75%
1200	87.21%	93.10%	92.42%
1400	87.71%	93.10%	92.59%
1600	88.22%	93.43%	92.59%
1800	88.89%	93.77%	92.76%
2000	88.55%	93.60%	93.10%
2200	88.55%	93.94%	93.43%
2400	88.55%	94.11%	93.43%
2600	88.55%	94.11%	93.43%
2800	88.55%	94.11%	93.60%
3000	88.55%	94.28%	93.60%
3200	88.55%	94.11%	93.43%
3400	88.55%	94.11%	93.43%
3600	88.55%	93.94%	93.43%
3800	88.55%	93.94%	93.60%
4000	88.55%	93.94%	93.77%

<글자수 별 탐지 모델 정확도>

훈련 데이터 글자수	전체 정확도
2000	85.35%
3000	94.50%
4000	90.21%

<탐지 모델별 전체(모든글자) 정확도>

탐지단계	적용 보이스피싱 확률
안전	~ 34.99 %
보통	35 ~ 49.99 %
위험	50 ~ 69.99 %
경고	70 % ~

<탐지 단계별 확률 적용 범위>

- 2,000 ~ 4,000 자 모델로 학습하여 탐지율을 비교한 결과 3,000자 모델이 가장 정확 하였음
- 3,000자 모델을 기준으로 탐지 단계별 보이스피싱 확률 적용

| Doc2Vec, Jaro-Wrinkler 문서간 유사도 검색

- 전처리, 형태소분석, 오타변환, 치환, 벡터화

- 기계학습된 기존 보이스피싱 스크립트와 유사도 판단

1. Cos 유사도 산출

→ 기존 학습된 스크립트의 벡터와 보이스피싱 통화의 벡터와의 내적을 구하여 얼마나 가까운지(유사한지) 산출

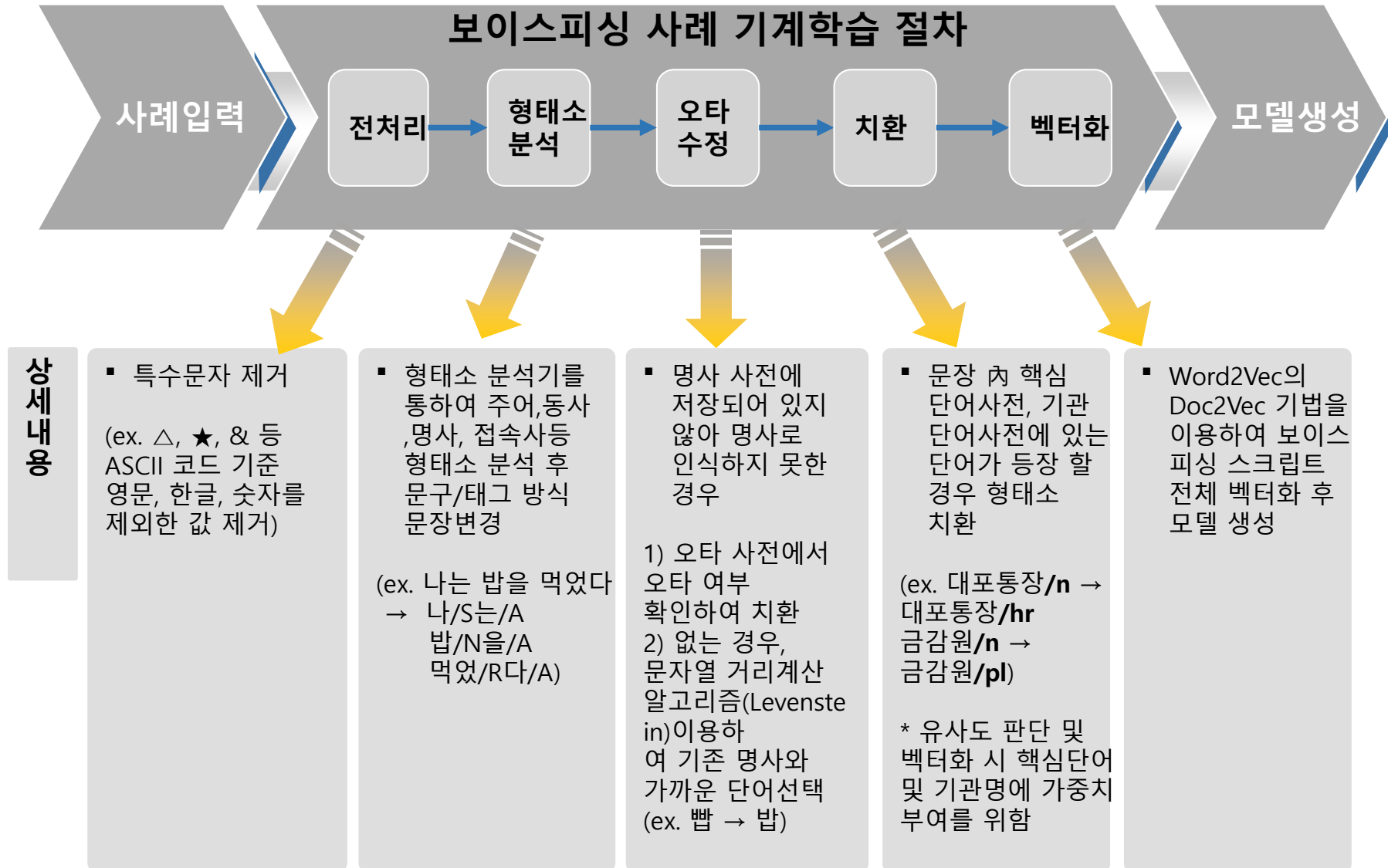
2. JaroWinkler 알고리즘 이용 유사도 산출

→ 문장 內 단어들의 위치를 전이시켜 일치하는 정도를 파악하는 문자열 유사도 판단 알고리즘

3. 최종유사도 판단

→ $\text{Cos 유사도} * 0.3 + \text{JaroWinkler} * 0.7$ 으로 유사도 판단

(학습데이터의 양이 적어, Cos 유사도에 적은 가중치 부여)



| Jaro-Wrinkler 알고리즘 고도화

The Jaro distance d_j of two given strings S_1 and S_2 is

$$d_j = \begin{cases} 0 & \text{if } m = 0 \\ \frac{1}{3} \left(\frac{m}{|s_1|} + \frac{m}{|s_2|} + \frac{m-t}{m} \right) & \text{otherwise} \end{cases}$$

where:

- m is the number of *matching characters* (see below);
- t is half the number of *transpositions* (see below).

변경사항

1. s_1, s_2 (sentence) $\rightarrow d_1, d_2$ (document) 로 변형
2. m, d_1, d_2 변수에 각 단어 별 가중치 부여 (가중치가 높을수록 높은 비율 차지)
3. 단어사전에 등록된 단어일 경우에만 t (transpositions) 에 포함

→ 변경전 대비 정확도 15% 향상

Jaro-Wrinkler 모델 실험보고서

VPDS_builder x +

← → ↺ 주의 요약 | 182.173.185.28:8090/frame?user=vpds&project=vpds&emno=anonymous&model_type=2&screen_type=2#

텍스트등록 유형관리 단어등록 단어분석 실험보고서 로그관리 환경설정

실험보고서 x 유형관리

보이스피싱 4그룹 조회 Q 엑셀다운로드 정확도계산 %

1 2 3 4 5 6 7 8 9 10 >

결과확률 75%

유사문서 75% 74% 61% 61% 59%

여보세요 예 예 맞습니까 예 수고하십니까 서울지방경찰청 경제범죄수사팀 조병준 경장입니다 예 다름이 아니고 이제 본인 앞으로 연루된 사건이 있어서 연락드렸구요 뭐 사건이요 사기 및 전자 금융거래 위반사건으로 고소가 들어왔어요 그게 뭐예요 제가 이제 다 설명드릴텐데 예 예 혹시 어떤 사이시죠 그 누군지 모르는데요 모르시구요 예 아 저희가 왜 여쭙 봤냐면 저희 특별수사팀에서 지난주 금요일에 금융사기단 팔 명을 검거했어요 예 금융사기단 저희가 팔 명을

더보기

결과확률 57%

유사문서 57% 56% 55% 55% 53%

잠깐 저기 소리 다시 한 번만 얘기해 주겠어요 여기 소리가 잘 안 들리는데요 네 본인께서 혹시 아십니까 몰라요 출신의 남성입니다 네 출신의 남성입니다 모르는데요 아 저희가 금융범죄사기단을 검거했습니다 네 검거한 현장에서 대량의 신용카드와 대포통장을 압수했구요 아 그래요 네 본인 명의로 된 제일은행 통장과 외환은행 통장이 발견돼서 조사자 연락을 드린 겁니다 저는 제일은행하고 외환은행은 관리를 안 아 저기 거래 안 하는데요 개설하신 적도 전혀 없

더보기

결과확률 79%

유사문서 79% 75% 57% 55% 53%

연락드렸는데 통화 가능하십니까 네 당신의 상황 때문에 연락드린 건 아니구요 개인정보 유출 건 때문에 몇 가지만 여쭙보겠습니다 혹시 최근 사이에 신분증 면허증 여권 등을 분실하거나 타인에게 양도하신 적 있으십니까 아니요 없는데요 그럼 각 은행이나 카드사에서 개인정보 유출이 많이 됐는데요 그걸로 인해 피해보신 사례는 있으십니까 아니요 전혀 없으시다 말씀이세요 네 없는데요 그런 적 여보세요 네 그런 적 없는데요 이거 여쭙보는 이유는요 저희가 이전 실험

더보기

결과확률 75%

유사문서 75% 75% 55% 54% 54%

문서단어추출

데이터업로드

실험보고서생성

- 최근 업로드된 데이터

추출어휘: 명사

데이터그룹: 1

보이스피싱 데이터 입력 시

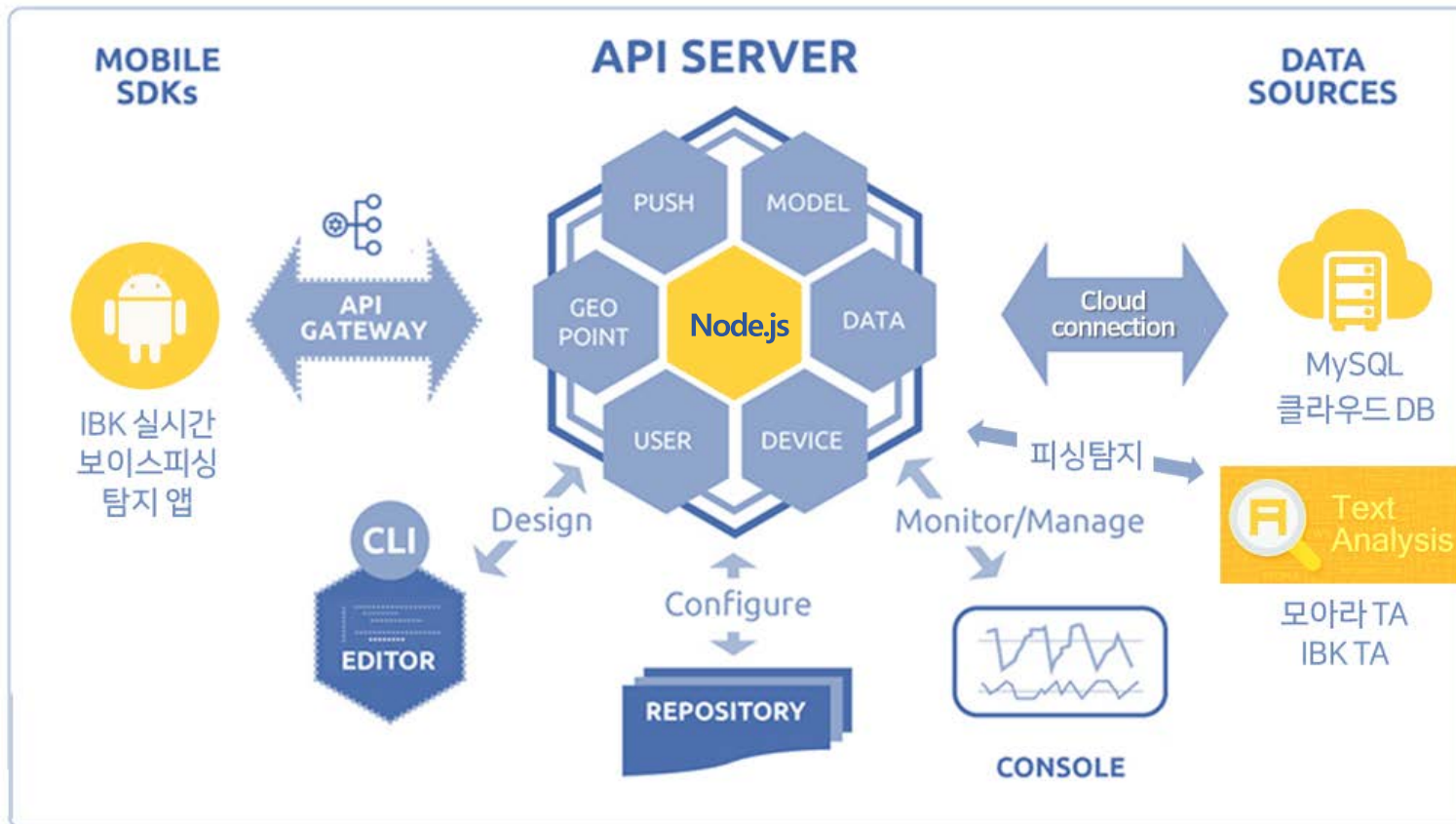
확률(%)	문서 수	비율
0~35	3	4%
35~50	2	2%
50~70	49	66%
70~	20	27%

일반통화 데이터 입력 시

확률(%)	문서 수	비율
0~35	117	61%
35~50	50	26%
50~70	22	11%
70~	0	0%

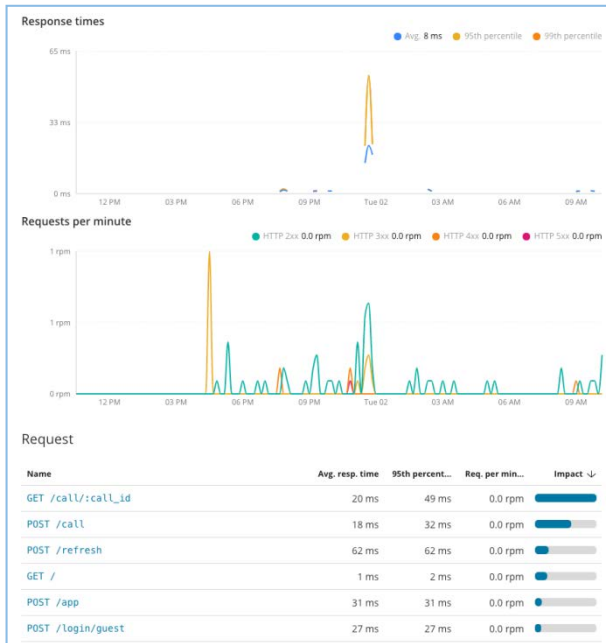
- 종합적으로 88~90% 정확도 달성

수행 내용 : 소프트웨어 아키텍처



“성능 극대화를 위한 비동기(Async) 아키텍처”

수행 내용 : 운영관리 소프트웨어



실시간 거래성능
모니터링

GitLab Projects Groups More

dev was +

History Find file Web IDE

내부 이미지 연동
김성신 authored 3 days ago

Name	Last commit	Last update
api	가중치 설정	4 days ago
bin	appkill.sh 추가	4 days ago
config	내부 이미지 연동	3 days ago
database	분석서버 인터페이스 연동	2 weeks ago
lib	분석서버 인터페이스 연동	2 weeks ago
util	SNS 로그인 기능 추가	1 month ago
views	index 파일 수정	4 days ago
.gitignore	provider 검증 예외처리 추가	3 weeks ago
app.ts	내부 이미지 연동	3 days ago
package.json	노드 APM 구동	3 days ago
route.ts	통화분석 API 추가	3 weeks ago
tsconfig.json	types-express 변경	2 months ago

앱, AP, TA, ML
통합 형상관리

Jenkins

새로운 Item

사람

빌드 기록

프로젝트 연관 관계

파일 링거프린트 확인

Jenkins 관리

My Views

Credentials

New View

빌드 대기 목록

빌드 대기 항목이 없습니다.

빌드 실행 상태

1 대기 중

2 대기 중

S	W	Name
		JOB-P-WAS-01
		pip
		sah-test

아이콘: S M L

빌드 및 배포 관리
자동화

“클라우드 모니터링, 형상관리 및 빌드/배포 자동화 구현”



Contents

I 사업개요

II 서비스 소개

III 기술 소개

IV  향후 발전 방향

V Q & A

IV

향후 발전 방향

서비스 고도화로 실시간 보이스피싱 탐지를 넘어, SMS 스미싱까지 개발하여

Platform으로 발전

민생사기탐지 Platform

보이스피싱

딤러닝

스미싱

**>SMS 스미싱 방지 App**

이상 URL 탐지를 DB 및 딤러닝을 활용하여 SMS 스미싱 방지 Process를 구축

>보이스피싱 방지 App

딤러닝을 활용한 실시간 보이스피싱 탐지



등 통계를 활용하여

서비스 개발 및 제공

82%

- 20대 ~ 30대
- 40대 ~ 50대
- 60대 이상 ~

01



사용자 연령대별 보이스피싱
방지 대안 마련

02

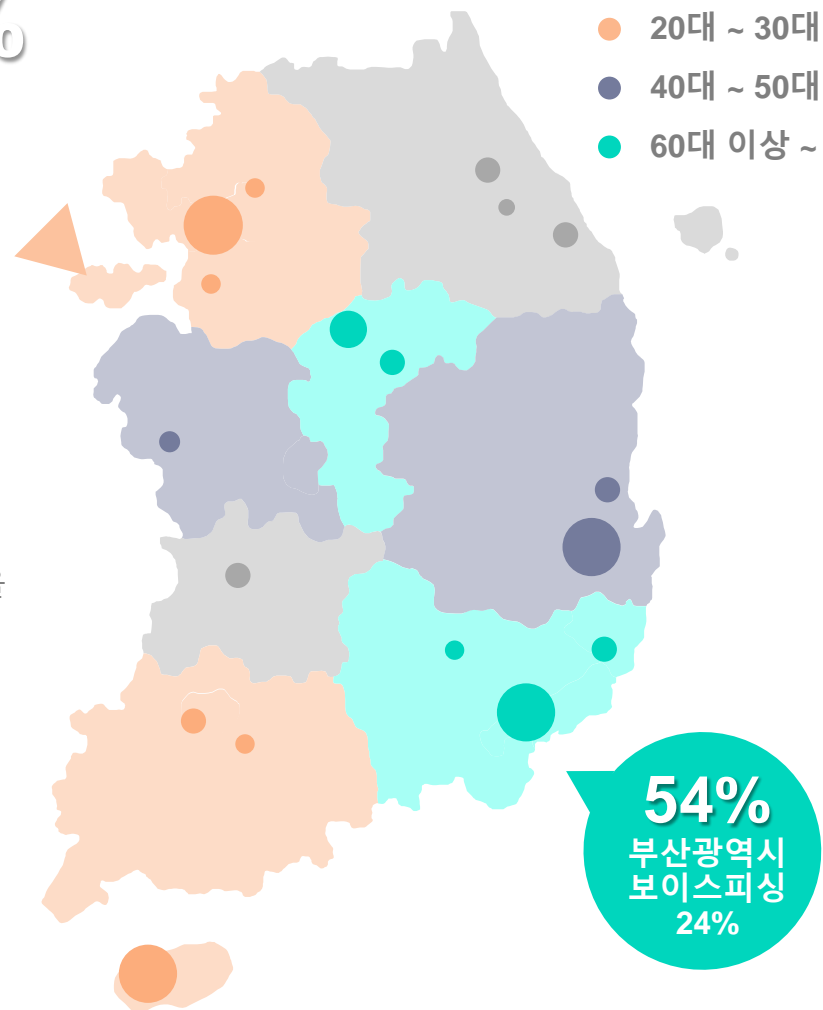


지역별 보이스피싱 발생 확률을
통해 방지 대안 마련

03



직업군별 분석하여 그에 맞는
보이스피싱 방지 대안 마련





Contents

I 사업개요

II 서비스 소개

III 기술 소개

IV 향후 발전 방향

V Q & A





Q&A

감사합니다

2018. 11. 21

