

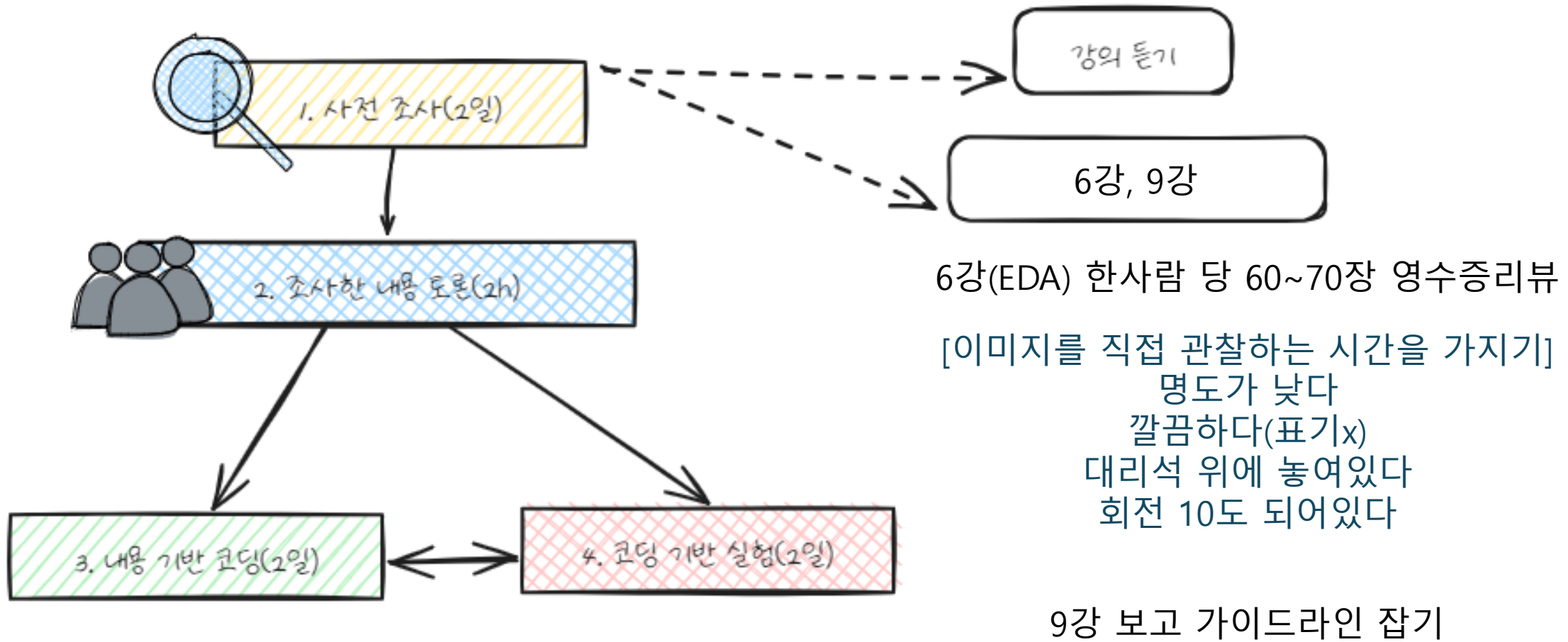
Datacentric week11 주간 회고록

# **NAVER Boostcamp AI Tech**

## **CV-21**

김한얼 김보현 김성주 윤남규 정수현 허민석

# Mon



# Tue

EDA_train	0	extractor.zh.in_house.appen2_001047_page0001.j	QR코드나 바코드가 존재함	보현
	1	extractor.zh.in_house.appen2_001050_page0001.j	확대해도 흐릿함 박스가 겹침	보현
EDA_test	68	extractor.zh.in_house.appen_000700_page0001.	그림자가 존재함 QR코드나 바	성주
	69	extractor.zh.in_house.appen_000701_page0001.	확대해도 흐릿함	성주
	135	extractor.ja.in_house.appen_000309_page0001.j	그림자가 존재함	한얼
	136	extractor.ja.in_house.appen_000316_page0001.j	그림자가 존재함 텍스트가 로고나 Q	한얼
	236	extractor.th.in_house.appen_000368_page0001.j	텍스트가 로고나 QR코드 안...	남규
	237	extractor.th.in_house.appen_000377_page0001.j	그림자가 존재함 박스가 잘못됨	남규
	267	extractor.th.in_house.appen_000660_page0001.j	스캔본임	수현
	268	extractor.th.in_house.appen_000666_page0001.j	그림자가 존재함 QR코드나 바	수현
	398	extractor.th.in_house.appen_000184_page0001.j	스캔본임	민석
	399	extractor.th.in_house.appen_000110_page0001.j		민석

데일리 스크럼에서 Annotation rule을 세웠다.  
 전원 10강까지 다듣고 train,test 이미지를 인당 80장씩 나누어서 EDA를 했다.

# Wed

## 2.데이터셋을 어떻게 쓸지

(public-ufo로 가져오기, synthetic)

public: 언어 비중에 대한 EDA만 해서 몇개 있는지 파악해서 가져오

synthetic: 문서를 이미지파일로 저장해보고싶다..

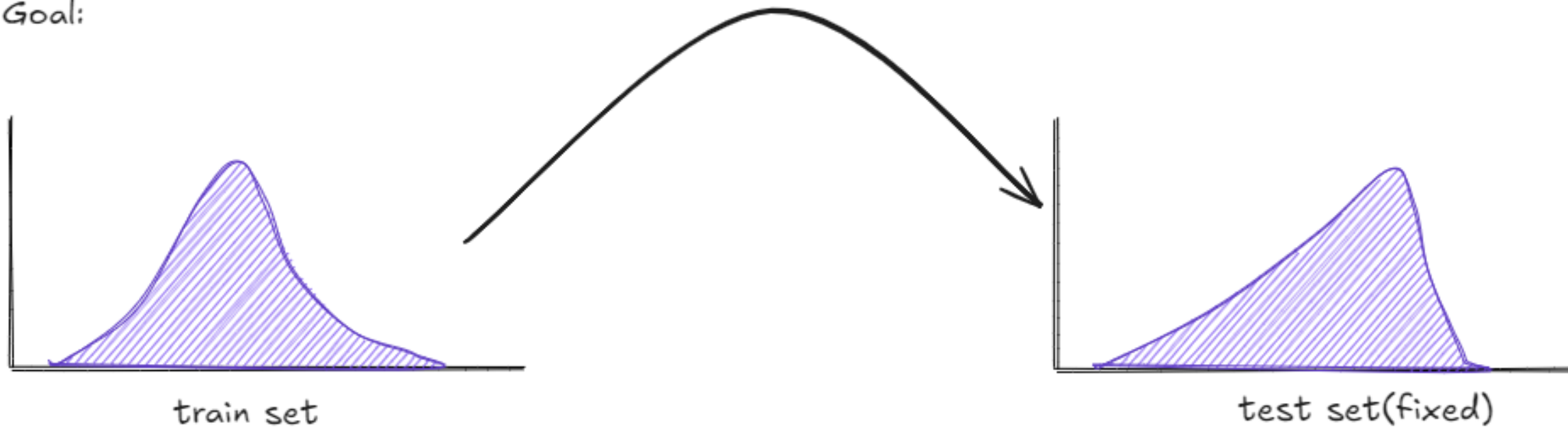
 영수증에 들어갈 단어를 어떻게 생성할지

## 3.증강을 추가적으로 어떻게 할지 고민(Test Data 보기)

- 깨끗하게
- 더럽게

# Thu

Goal:



Methods

- 10~15도 회전
- 그림자 부여 또는 밝기 조절
- Normalization: 테스트 데이터의 평균과 분산에 맞춰야 함

# Thu

## Methods

6명이 공통으로 써봄직한 방법

- 10~15도 회전
- 그림자 부여 또는 밝기 조절
- *Normalization*: 테스트 데이터의 평균과 분산에 맞춰야 함
- 가우시안 노이즈, 솔트앤페퍼vs디노이즈?
- 배경을 추가 vs 배경을 제거?
- *Superresolution* vs 일부러 *Blur*?
- 잘린 글씨
- 작은 이미지 그냥 확대

# Thu

6명이 각자 써봄직한 방법

- 학습을 언어 종류 별로 진행하기
  - 중국어:
  - 일본어: 가로 선을 제거??  
세로쓰기가 없으면 세로데이터 제거
  - 태국어:
  - 베트남어: 스캔본이 많고 공백이 있으므로 padding
- 공개 데이터 가져오기
  - 퍼블릭 데이터
  - 데이터 크롤링: 슈도 라벨링 사용
- 합성 데이터 만들기: 휘어진 글자 추가, 회전된 글자 등

결과별로 CSV 통합

Ensemble

그 외

- 특수한 데이터 처리 고민: 흰색 줄 처리, 보안코드 처리, 분홍 선, 투명 배경, 붉은 영수증  
접힌 자국. ---oOo---, 손가락, 영수증 외부 글씨, 인식되지 않은 경계선? 등
- 앙상블: 앙상블도 데이터를 보면서 잘 해보려는 시도

# Fri

因为手机或网络交易比重增加所以线下工  
作岗位减少了:了解此次事件的  
一位相关人士透露:前李副社长曾公开表示自  
己在青瓦台和金融监督院有人脉  
所以不用担心并表示Lime资产的问题也会得到  
解决还展示了该调查书拒绝监查意见是  
像不能相信公司财务报表一样不能相信根  
据材料的情况下提出的意见:全国公共工  
会联盟釜山地区本部决定介绍釜山银  
行提供的多种金融商品和服务:据悉除此之  
外检察机关还掌握了李前副总经理等收购多  
家投资企业债券时收取非法手续费的状况  
正在进行调查:实际上上月日伊朗为反对  
美军空袭而实施报复空袭比特币价格暴涨了  
达到万韩元以下自上月日中国发生的新冠  
状病毒在全球蔓延的忧虑正式爆发以来  
已经稳定在万韩元水平:同期的本期  
净损失为67,270亿韩元亏损急剧增加:  
朴部长法官解释了发放理由说:有可能逃  
走和毁灭证据:以初中毕业学历开始创业并  
创造了2,000亿韩元销售额企业的他以自己  
的经验为基础说明了世界考虑未来时我们  
创造未来的经营蓝图:韩亚金  
融控股公司副会长咸英株于日进入在首  
尔松坡区首尔牙山医院殡仪馆里已故  
乐天集团名誉会长申格浩灵堂:Telechips  
集团董事长李秀仁(音)解释道

중국어 말뭉치 전처리 txt파일 10개입니다

Zip ▼



chinese.zip

Zip

일본어 말뭉치 전처리 10개 txt 압축파일, 합성 이미지 생성 위한 중국어 폰  
트, 태국어 말뭉치 남기고 가겠습니다,, 좋은 하루 되세요 😊

2개 파일 ▼



japanese.zip

Zip

01

ZhouFangRiMingTi-2.otf

이진

수현님의 공공데이터 해체show



# 감사합니다.

