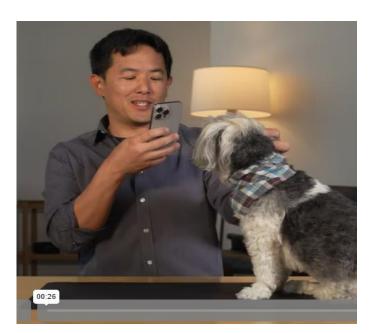
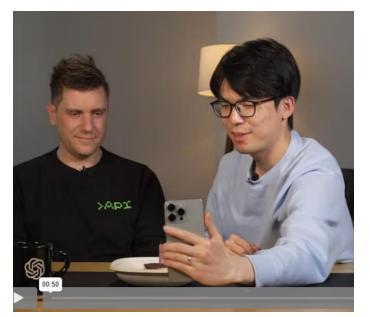
신뢰할 수 있는 데이터, 가치있는 데이터

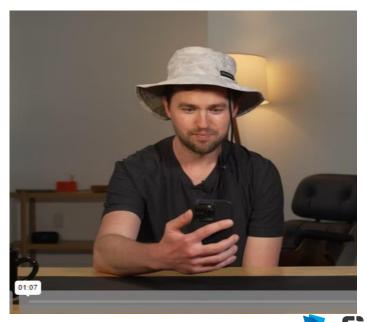
May 13, 2024

Hello GPT-40

We're announcing GPT-4o, our new flagship model that can reason across audio, vision, and text in real time.



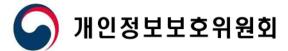




비정형 데이터 가명처리 기준

비정형데이터 가명처리 기준 주요 내용

2024. 2.



AX를 위한 대응

■ 데이터 3법 (2020.1)

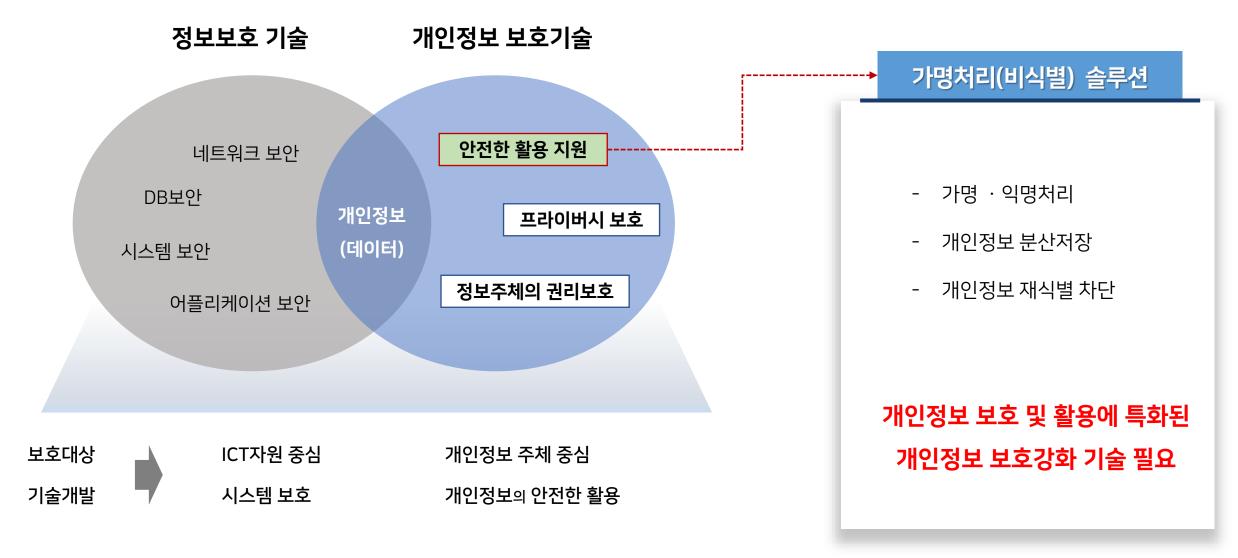
- 데이터 활용을 위한 가명정보 개념 도입
- 가명처리(비식별) 조치 후 동의없이 사용가능
- 통계작성 및 과학적 연구, 공익적 기록보존 목적

■ 패러다임의 변화

- AI시대 대응을 위해
- 비정형 데이터를 안전하게 활용



개인정보 가명처리 기술의 필요성





비정형 데이터 가명처리 예시

1. 이미지 + 이미지 內 포함된 텍스트

항목 가명처리 前 가명처리 後 흉부CT 이미지 1. 흉부촬영 부분 1. 그대로 유지 2. 환자관련 정보 2. 환자관련 정보 2.1 환자이름 2.1 가명처리로 대체 2.2 생년월일 2.2 연도정보로 변경 2.3 환자성별 2.3 성별값 "S"로 대체 2.4 환자번호 2.4 일련번호로 대체

2. 영상 + 이미지

항목	가명처리 前	가명처리 後	
사람 얼굴			
차량 번호판			
처리 기술	이미지 필터링(Blurring) 기술		

3. 음성 : STT 처리 된 텍스트(json)

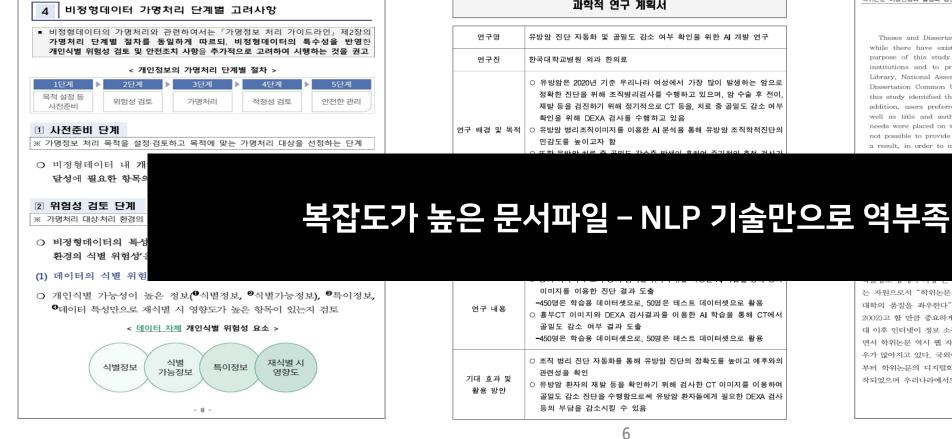
	항목	가명처리 前	가명처리 後
	이름	<u>홍길동</u> 고객님 본인 맞으실까요? <고객이름 '김행복 ' 으로 치 환>	'[NAME_CUSTOMER]' 고객님 본인 맞으실까요? → '김행복' 고객님 본인 맞으실까 요?
		네 감사합니다. 이상 상담원 <u>이</u> 선정이었습니다. <상담원 이름 '김신뢰'로 치환>	네 감사합니다. 이상 상담원 '[NAME_CUSTOMER]' 이었습니다. → 네 감사합니다. 이상 상담원 '김신뢰' 였습니다.
	생년 월일	가입해주신 고객님 혹시 1987 년 7월 22일생 맞으실까요?	가입해주신 고객님 혹시 '[BIRTH]' 생 맞으실까요?
		<메타데이터의 범주화 된 연령 (30~40대) 내 랜덤값 치환>	→ 가입해주신 고객님 혹시 '1991년 6월 8일 ' 생 맞으실 까요?
	주소 정보	고객님 서울시 신뢰동 신뢰아파 트 백일동 구백삼호로 배송해 드리겠습니다. <메타데이터의 주거지역(서울) 내 가상주소로 치환>	고객님 '[ADDR]' 로 배송해 드리겠습니다. → 고객님 '서울시 신뢰동 신뢰아 파트 1동 1호 '로 배송해 드리겠습니다.

문서 파일에 대한 가명처리 방안은?

4. 문서를 포함함 텍스트



- 자연어처리(NLP) 기술을 이용하여 텍스트 가명처리
- 정규표현식 및 주석달기(Annotation) 기술 등을 이용하여 마스킹 및 대체 처리



파악식 연구 계획서					
연구명	유방암 진단 자동화 및 골밀도 감소 여부 확인을 위한 AI개발 연구				
연구진 한국대학교병원 외과 한의료					
연구 배경 및 목적	○ 유방암은 2020년 기준 우리나라 여성에서 가장 많이 발생하는 암으로 정확한 진단을 위해 조직범리검사를 수행하고 있으며, 암 수술 후 전이, 재발 등을 검진하기 위해 정기적으로 CT 등을, 치료 중 골밀도 감소 여부 확인을 위해 DEXA 검사를 수행하고 있음 ○ 유방암 병리조직이미지를 이용한 AI 분석을 통해 유방암 조직학적진단의 민감도를 높이고자 함				

ABSTRACT

Theses and Dissertation(TD) have been considered one of valuable scholarly resources. purpose of this study is to investigate users' perception on six TD services from five Library, National Assembly Library, RISS, dCollections, NDSL, and Council of Theses and Dissertation Common Use, Based on the survey results from 151 users, the findings of addition, users preferred keyword, full text, department, abstract, table of contents, as well as title and author over other bibliographic information. More importantly, users needs were placed on whether specific TD services provide full text or not, In case that is not possible to provide full text, users have a preference for full text link information, As

이 인쇄본과 더불어 전자 형))을 생산하고 있다. 하지만

관리 주체, 제출 시스템, 파일 포맷(File Format), 저작권, 아카이빙(Archiving), 온라인 접근 문제 등 아직 해결해야 할 문제가 많다. 최근 10여 년간 이러한 문제들 에 관하여 다양한 관점에서 다룬 연구들이 진 행되어 왔으며, 학위논문의 전자형태 및 제출, 제공시스템의 표준을 통일하고 하나의 시스템 으로 통합하여 총괄하는 국가적인 네트워크 형성이 필요하다는 것으로 초점이 모아지고

는 자원으로서 "학위논문의 디지털화가 미래 대학의 품질을 좌우한다"(Edminster et al. 2002)고 할 만큼 중요하게 인식된다. 1990년 대 이후 인터넷이 정보 소통의 주요 수단이 되 면서 학위논문 역시 웹 자원으로 활용되는 경 우가 많아지고 있다. 국외에서는 1990년대 초

부터 학위논문의 디지털화에 관한 논의가 시

작되었으며 우리나라에서도 1990년대 중반부

-30 -

-450명은 학습용 데이터셋으로, 50명은 테스트 데이터셋으로 활용

-450명은 학습용 데이터셋으로, 50명은 테스트 데이터셋으로 활용

○ 조직 병리 진단 자동화를 통해 유방암 진단의 정확도를 높이고 예후와의

○ 유방암 환자의 재발 등을 확인하기 위해 검사한 CT 이미지를 이용하여

골밀도 감소 진단을 수행함으로써 유방암 환자들에게 필요한 DEXA 검사

○ 흉부CT 이미지와 DEXA 검사결과를 이용한 AI 학습을 통해 CT에서

이미지를 이용한 진단 결과 도출

골밀도 감소 여부 결과 도출

등의 부담을 감소시킬 수 있음

연구 내용

기대 효과 및

활용 방안

문서에 있는 '개인정보 가명처리' 어떻게 해야 하는가?

문서에서의 가명처리 방안

1 문서의 '구조' 이해

- 다양한 문서포맷
- 한글, 오피스, PDF, 기타
- 다양한 버전
- 수많은 문서 유형
- 문서 작성 Style : Lay-Out
- 문서 유형 : 보고서, 증명서 등
- Object
- Font, Line, Table, Image
- 위치(좌표), Size, Color 등

2 문장의 '내용' 이해

- 구문 분석
- 문장을 이루고 있는 구성 성분 으로 분해 → 관계 분석
 - → 문장의 구조 결정
- 문장화 처리
- 장 / 절 구분
- 데이터 처리
- 형태소 분석
- 사전(Dictionary) 분석
- 기타

3 '결과물' 활용 방안

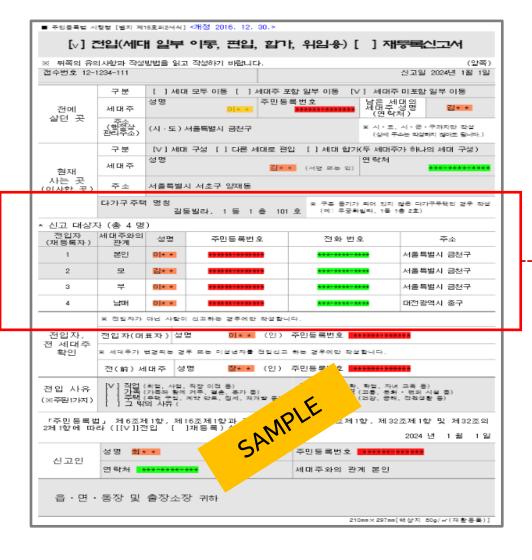
가명처리 된 결과를 어떻게 제공(output)할 것인가?



'반출할 수 있는 파일'



■ 대상 : 한글, MS오피스,PDF + 이미지



1. META DATA 추출 및 구조화

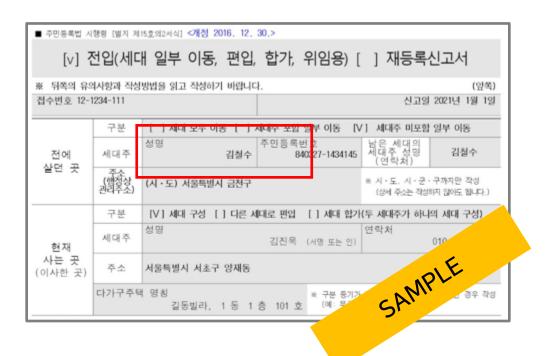
page":{"pagenum":1,"sz":{"cx":59528,"cy":84188},"body":{"l":5669,"t":5669,"r":53859,"b":81354},"ct":[ct":[]},{"ot":"shape","pos":{"x":4252,"y":4252},"sz":{"cx":51024,"cy":78519},"type":"pageborder" sz":{"cx":48198,"cy":75681},"tr":[{"sz":{"cx":48198,"cy":1182},"tc":[{"cols":11,"pos":{"x":5952,"y":595" "pos":{"x":6235,"y":6858},"rect":{"l":6235,"t":6171,"r":6979,"b":6971},"tid":0,"face":"唇음제","sz":800, ("1":6979,"t":6171,"r":7351,"b":6971},"tid":1,"adv":"372","v":" "},{"ot":"t","pos":{"x":7351,"y":6858} v":"주민등록법"},{"ot":"t","pos":{"x":11071,"y":6858},"rect":{"l":11071,"t":6171,"r":11443,"b":6971},"t "t":6171,"r":13675,"b":6971},"tid":4,"adv":"744,744,744","v":"시행령"},{"ot":"t","pos":{"x":13675,"y":68 ["},{"ot":"t","pos":{"x":14419,"y":6858},"rect":{"l":14419,"t":6171,"r":15907,"b":6971},"tid":6,"adv" "t":6171,"r":16279,"b":6971},"tid":7,"adv":"372","v":" "},{"ot":"t","pos":{"x":16279,"y":6858},"rect": "pos":{"x":17023,"y":6858},"rect":{"l":17023,"t":6171,"r":17767,"b":6971},"tid":9,"adv":"372,372","v" "b":6971},"tid":10,"adv":"744,744","v":"호의"},{"ot":"t","pos":{"x":19255,"y":6858},"rect":{"1":19255, {"x":19627,"y":6858},"rect":{"1":19627,"t":6171,"r":21115,"b":6971},"tid":12,"adv":"744,744","v":"서식 b":6971},"tid":13, adv":"372,372","v":"] "},{"ot":"t","pos":{"x":21859,"y":6858},"rect":{"1":21859,"t {"ot":"t","pos":{"x":22309,"y":6858},"rect":{"1":22309,"t":6085,"r":24109,"b":6985},"tid":15,"adv":"900 cy":1182},"ct":[{"ot":"p","align":"right","ct":[]}]}]},{"sz":{"cx":48198,"cy":3966},"tc":{"cols":15,"p" 7134 L 53952 7134 L 53952 11063 L 5952 11063 Z","fillPr":{"type":"clr","clrs":[{"clr":"ffe5e5e5", "pos" "rect":{"1":9076,"t":8395,"r":9546,"b":9995},"tid":17,"face":"휴먼옛체","sz":1600,"clr":"ff000000","adv "r":10363,"b":9825},"tid":18,"face":"HY신명조","sz":1048,"adv":"817","v":"V"},{"ot":"t","pos":{"x":10363 "face":"휴면옛체","sz":1600,"adv":"470,699","v":"] "},{"ot":"t","pos":{"x":11532,"y":9677},"rect":{"1":1 "pos":{"x":14330,"y":9677},"rect":{"l":14330,"t":8395,"r":14844,"b":9995},"tid":21,"adv":"514","v":"("} "b":9995},"tid":22,"adv":"1399,1399","v":"세대"},{"ot":"t","pos":{"x":17642,"y":9677},"rect":{"1":17642, {"x":18341,"y":9677},"rect":{"1":18341,"t":8395,"r":21139,"b":9995},"tid":24,"adv":"1399,1399","v":"일부 b":9995},"tid":25,"adv":"699","v":" "},{"ot":"t","pos":{"x":21838,"y":9677},"rect":{"1":21838,"t":8395,

Meta Data 추출 : 문서의 내용을 구성하는 주요정보

Text	Object	Lay-Out
- txt & 좌표 - font & size &	- Line & 좌표 - Table & 속성	- 문서 유형 - 문서 형식
color 등	- Image 등	



■ 대상 : 한글, MS오피스,PDF + 이미지



2. 개인정보 검출

```
"align" : "right",
"ct" : [
        "ot": "t",
         "pos": {
             "x": 25796,
            "y": 19265
        "rect": {
            "1": 25796,
            "t": 18492,
            "r": 28496,
            "b": 19392
              : "900,900,900",
```

3. 개인정보 가명처리

```
"ot" : "p",
"align" : "right",
        "ot": "t",
        "pos": { ···
        "rect": { ...
        "tid": 99,
        "sz": 900,
        "adv": "900,900,900",
        "v": "김철수",
        "pim": {
            "use": "y",
            "config_use": "y",
            "bgcolor": "#FE7F50"
             "deidtype": ["휴리스틱1"],
            "findlist": ["김철수"],
             "word idxs": [99],
            "indexOf"
            "convertedword":
```



■ 논문 : 다단 구조 형식

제33회 한글 및 한국어 정보처리 학술대회 논문집 (2021년)

의존 구문 분석을 활용한 자연어 추론

Natural Language Inference using Dependency Parsing

Seul-gi Kim^{0.1}, Hong-Jin Kim², Hark-Soo Kim^{1,2}
Konkuk University Department of Computer and Communications Engineering¹,
Konkuk University Department of Artificial Inteligence²

9 9

자연이 추론은 두 문장 사이의 외미 관계를 문휴하는 작업이다. 본 논문에서 제안하는 외미 추운 방법은 의존 구용 분석을 사용하여 동일한 구문 정보나 기능 정보를 가진 두 개의 (고기비소, 지배스) 여행 방에서 하나의 어떻이 경찰 때 두 피지배소를 하나의 정크로 만들어라고 청고 기문으로 만들어진 의존 구문 논석을 사용하여 자연여 주론 작업을 수명하는 방법을 의미한다. 이리한 의미 주론 방법을 통해 만들어진 정크와 구문 구조 정보를 Staffine Attention을 사용하여 한 문장에 대한 정고 한위의 구문 구조 정보를 반영하고 구문 구조 정보가 반영된 후 문장을 Billinear을 통해 관계를 예곡하는 시스템을 제안한다. 실험 결과 정확도 90,75%로 가장 높은 성논을 보였다.

주제어: 의미 주론 방법, 의준 구문 분석, 자연어 주론, 정크(Chunk)

1. 서론

자연이 추론(Natural Language Inference)은 자연이 이해(Natural Language Understanding)를 기반으로 모델 의 추본 능력을 평가하는 작업으로 두 문장 사이의 의미 관계를 함의(Entailment), 모순(Contradiction), 중립 (Neutral)으로 분류하는 문장 쌍 분류(Sentence-Pair Classification)의 일종이다. 두 문장은 전제(Premise) 와 가설(Hypothesis)로 나누어지는데 전제를 잠이라고 가정할 때 가설의 내용이 참(함의)인지, 거짓(모순)인 지,혹은 알 수 없는지(중립)에 따라 두 문장의 관계가

에를 들어, 한국어 NLI 데이터 세트 KLUE-NLI**에는 전제 문장이 "부산은 블록제인 기술을 활용한 실증사업이 금융분야까지 확장하고, 대전은 바이오 스타트의 범원제 공용연구시설을 공유하는 등 일부 기조 동사업이 추가됐다."이고 가설 문장이 "부산과 기준 복구에는 사업들이 모두 폐지됐다."라 있다. 전제 문장에서 ("부산은", "대전", "사업이 추가됐다.")를 통해 구어도", "사업이 추가됐다.")를 통해

한국어 데이터 세트 KLUE-NLI,

https://aistages-prod-server-public.s3.amazonaws.com/a pp/Competitions/000068/data/klue-nli-v1.1.tar.gz

그리 1 이조 그로 그조 예시

제 문장이 모순됨을 알 수 있다. 본 논문에서는 KLUE-NA 의 주론 방법을 표 1과 같이 나는다. 표 법 중 의미 주론은 전제와 가설 문장 파악해야 해결할 수 있는 경우를 뜻는 문장의 의미를 더 정확하게 파악하 구문 분석을 활용한다. 먼저, 의존 구문 소설 중 각 어절(피지배소)이 지배소와 가질 수의존 관계 태그는 구문 정보와 기능 정보로 나눠진지, 이 중 구문 정보는 각 어절이 제연인지, 용연인지, 부사인지 등을 나타낸다. 기능 정보는 각 어절이 지배소와 가지는 관계가 주격인지, 목적격인지, 관정적인지 등을 나타낸다. 의존 관계 태그는 그림 1에서와같이 "-"를 기준으로 앞은 구문 정보, 뒤는 기능 정보가 표기되

■ 개인정보 검출

검출정보

- 김슬기
- 김홍진
- 김학수
- cludyju@konkuk.ac.kr
- jin34@gmail.com
- mlpdrkm@konkuk.ac.kr

■ 개인정보 가명처리

가명처리 결과

- 김*
- 김*
- 김**
- ***********
- ******@******
- ************



■ 이력서 : 표형식 + 이미지 객체



■ 개인정보 검출

검출정보

- 김일중
- **9**10101-1222222
- 02-111-2222
- **•** 010-1111-2222
- 김일중

■ 개인정보 가명처리

가명처리 결과

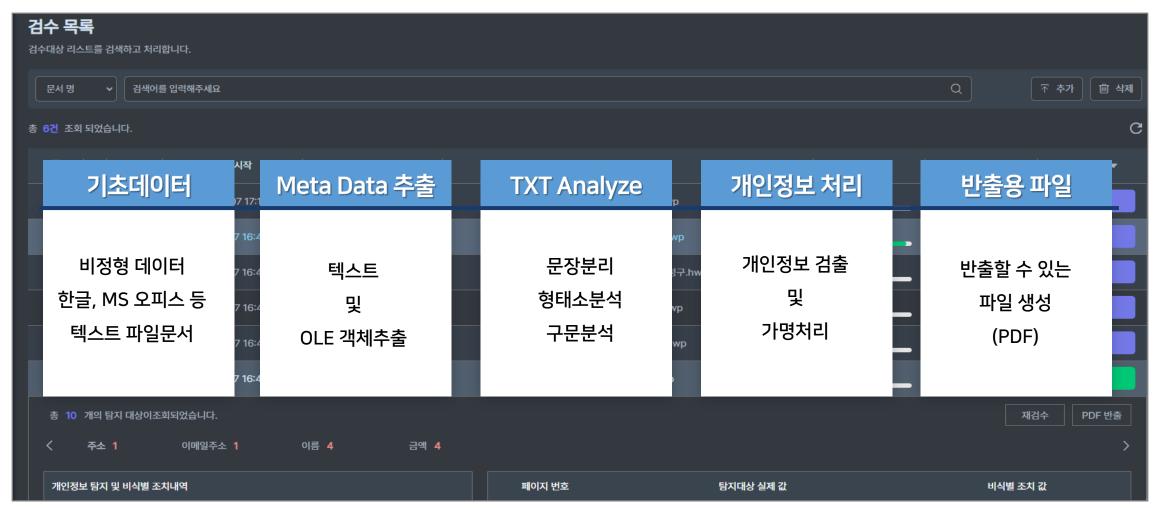
- 김**
- **_***
- ***_****
- 김**



반출할 수 있는 개인정보 가명처리 솔루션 Docu-Guard

Docu-Guard Process

☑ 반출(재사용)할 수 있는 구조

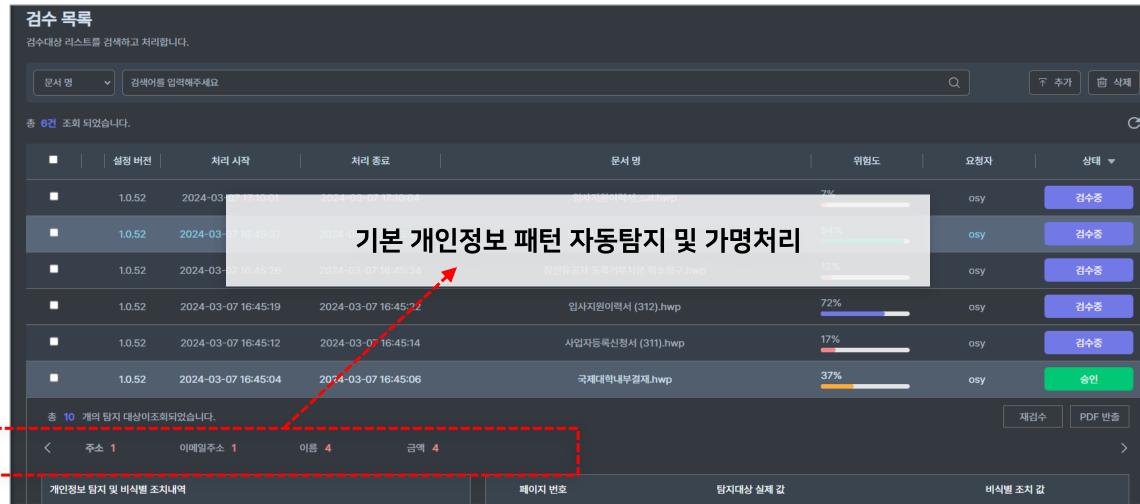




개인정보 자동 가명처리

\square

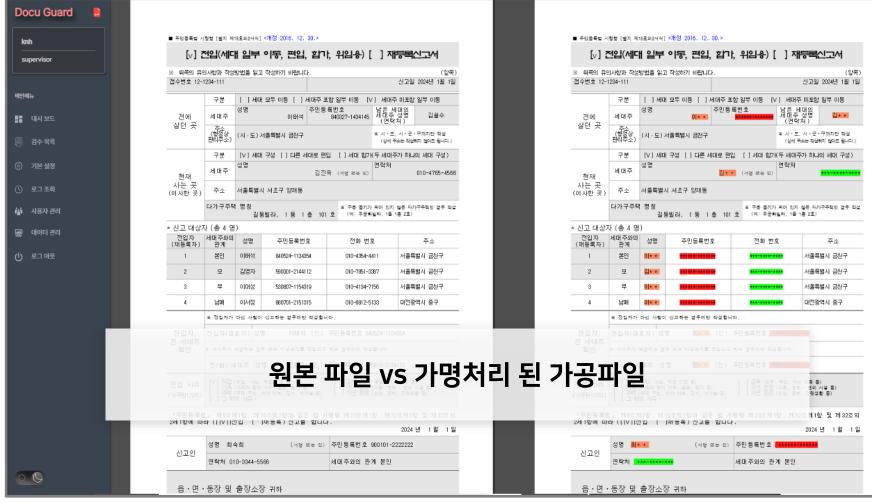
개인정보 자동 가명처리



개인정보 자동 가명처리

\square

쉽게 확인 할 수 있는 가명처리 정보



■ Dual Mode 적용

- HTML5 지원
- 한 화면에서 원본 vs 가공파일 비교
- 기본 Base Pattern 적용



개인정보 사용자 처리방식

\square

사용자 정의에 의한 수동 가명처리 방식



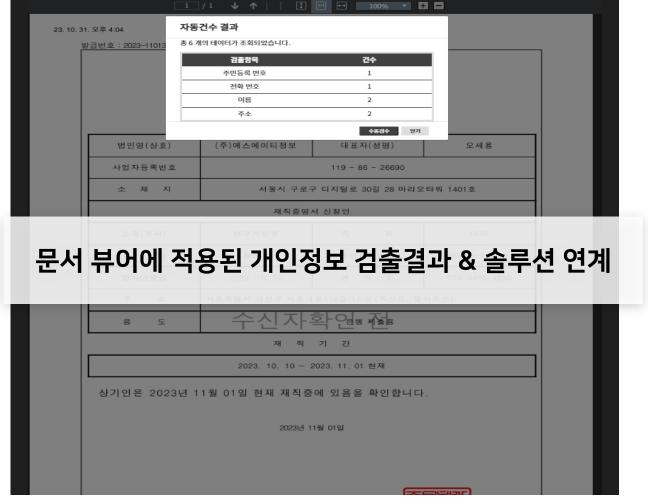
■ Drag 방식

- Base Pattern 이외의 텍스트 정보 등을 drag하여 손쉽게 가 명처리
- 신규 등록 패턴에 대한 일괄적 용 가능
- 문서 內 이미지 객체에 대해 인 식 & 가명처리 가능



사용성 강화 - 他 솔루션 연계

🗹 문서 뷰어 및 타 개인정보보호 솔루션 연계를 통한 사용성 강화



■ API 제공

1. Viewer 연계

- 기관에서 이용중인 문서 뷰어 연계 가능
- 파일 열람 시, 뷰어에서 개인정보 자동 검출 기능이용
- 자동 검출 된 개인정보 내역 확인 후,
- Docu-Guard 연계 API를 이용하여 사용 가능

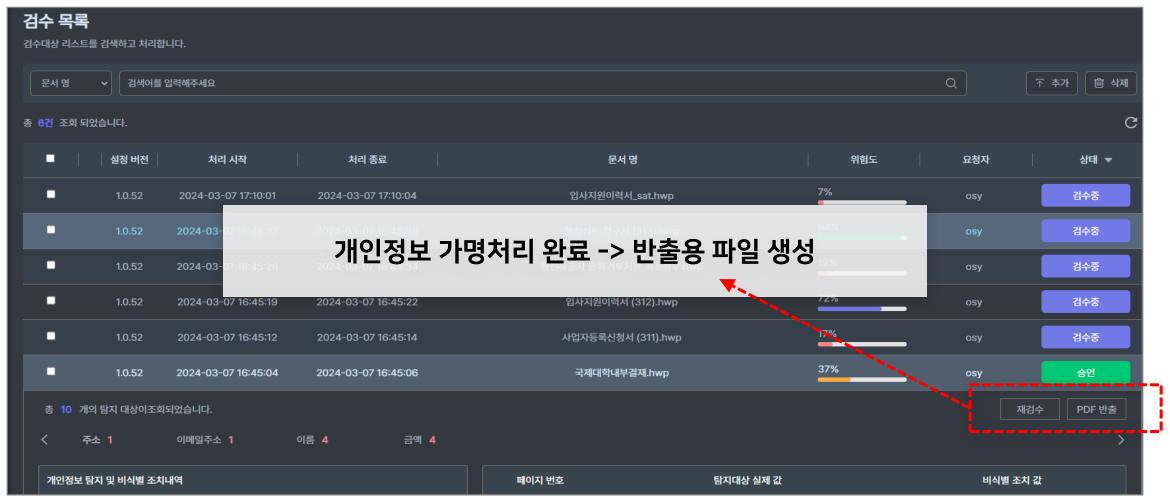
2. 개인정보 탐지 솔루션 연계

- 기관에서 이용중인 개인정보 탐지 솔루션 연계 가능
- 기 사용중인 솔루션을 통해 개인정보 탐지 결과를 제공받고,
- 이를 메타데이터와 연계하여 원본 파일에 가명처리



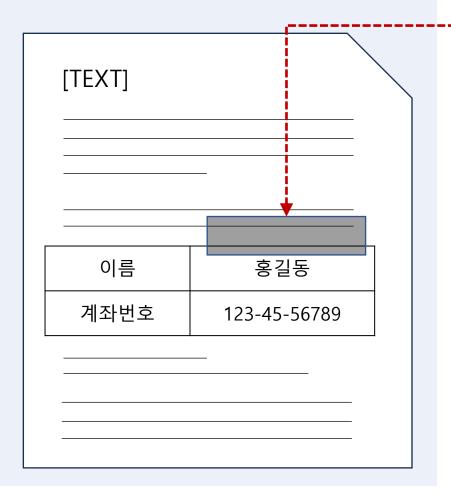
반출용 파일 생성

☑ 반출용 파일(PDF) 생성



가명처리 기술비교

■ 한글, OFFICE, PDF 등 텍스트 파일



1. Black Masking 방식

- 원본 → 이미지 or PDF 파일로 변환

을 좌표 위에 덧씌우는 방식

- 좌표가 어긋나는 현상 & 기본 패턴 외 개인정보 처리 난항

2. Text 추출 방식

- 원본 → 텍스트 파일을 추출
- 텍스트 파일 → 개인정보 가명처리 → 가명처리 된 별도의 파일생성
- 가명처리 정보를 원본에 재결합 하는 방법 난항

3. Meta Data 추출 방식

- 원본 → Meta Data 추출
- Source Data → 직접 가명처리 → 가명처리 된 사본 생성
- 원본형식을 유지한 사본 생성

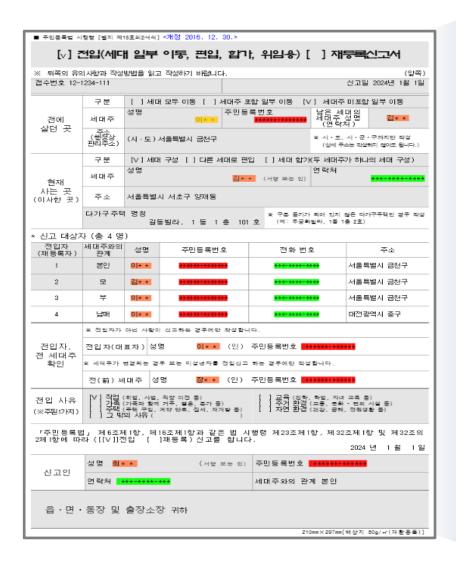


레퍼런스 예시



가명처리된 데이터 어떻게 사용할 것인가?

가명처리 된 반출용 파일 - 활용방안



- **7** 정보를 '공유' 할 수 있는 데이터로 사용가능
 - 가명처리 된 파일 → 외부와의 '정보 공유'
- **2** 정보를 '공개' 할 수 있는 데이터로 사용가능
 - 소극적 민원대응 → '적극적 민원대응' 으로
- **3** 신뢰할 수 있는 '데이터' 로 사용가능
 - 가명처리 된 파일 → 'AI를 활용할 수 있는 기초데이터' 로 활용



가명처리 된 반출용 파일 - 활용방안

가명처리 데이터

AI를 위한 기초데이터

Machine Readable(정형) 데이터

Fine Tunning

SLM



AI 활용을 위한 '가치있는 데이터' 로 사용가능

- 주요 정보자산 활용으로 '데이터 부족' 문제 해결
 - 비정형 데이터는 조직의 주요자산
 - AI를 Business Domain으로 사용하기 위해 필수 데이터
- 보다 '안전한 AI 생태계' 조성 기반
 - 출처가 명확한 데이터 사용으로 Hallucination 최소화
 - 개인정보 오남용 문제 최소화



Thank You

