

OCR 개념 & 원리

1 page

< OCR 단계 >

이미지 전처리 (노이즈제거, 기울기, 대비조정) → 텍스트 영역 감지 (box)

텍스트 추출

OCR이란? → 텍스트 인식

이미지나 스캔본 문서에서 텍스트를 인식해 디지털 텍스트로 변환

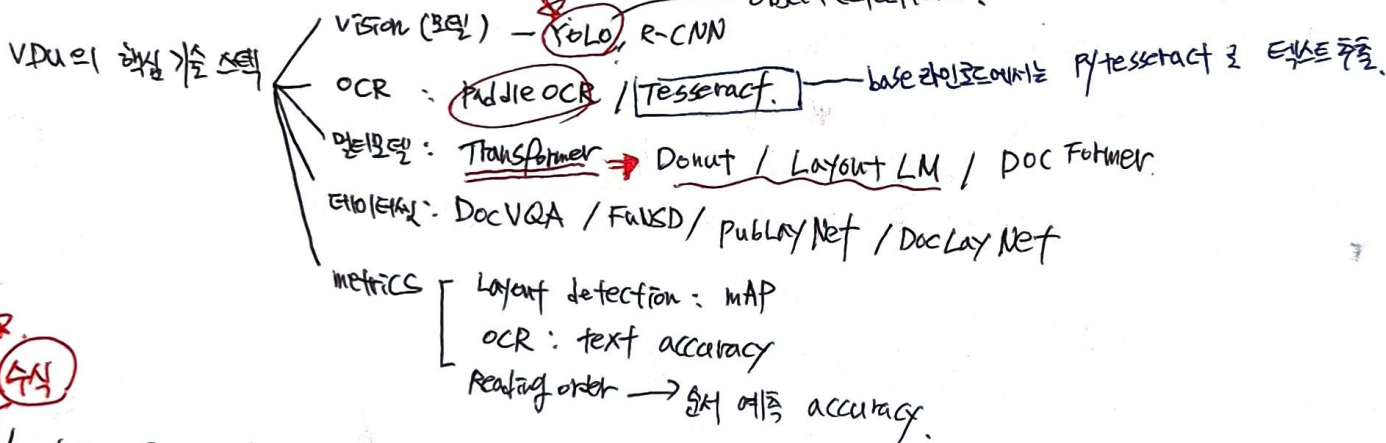
VDU (visually-rich document understanding) 모델에서 OCR은 필수이다.

→ 단순히 텍스트만 읽는 것이 아닌 문서의 시각요소 (레이아웃, 구조, 텍스트, 이미지)까지 포함해서 문서의 의미와 의도를 파악함!

NLP : 순수 텍스트 데이터 (txt, csv) 만 취급.

VDU : 문서의 위치정보, 시각적 배치, 각요소 간 관계를 함께 이해. (시각적 문맥을 포함하는 AI)

- VDU
 - 문서 레이아웃
 - OCR (텍스트 추출) ⇒ ~~Tesseract~~ / Paddle OCR ⇒ 답례 가능한 높은 정확도, 한글 포함 다국어 지원, GPU 가속기 / 텍스트 방향 / 수식 OCR / 테이블 OCR 결의
 - 사람이 읽는 순서대로 문서요소 정렬 (reading order prediction)
 - 멀티모달 정보 통합 (시각정보 + OCR) → 문서의 의미와 문맥 구조까지 이해 가능. object detection.



수식

Latex 추출 방법 (일반 OCR로는 안됨) : lm2Latex 모델 (CNN + attention 기반 딥러닝 모델)

Paddle OCR (+ 수식 전용 모델 fine tuning)

OCR 성능 개선

- 이미지 전처리
 - 이진화 (텍스트와 background 대비 극대화)
 - denoising (노이즈제거)
 - Thresholding (임계값 설정)
 - Deskewing (기울어진 문서를 회전 보정)

모델 성능 극대화 part.

OCR & 이미지 전처리

- 모델 최적화
 - Paddle OCR의 rec-model은 언어별로 최적화된 모델로 변경

- 성능 튜닝
 - 구조적 요소와 위치에 맞춘 det-model 선택.
 - 수식, 테이블 등 특수데이터에는 별도 특화모델 사용.

base OCR 모델

- 성능 튜닝
 - 이미지 전처리
 - 이미지 해상도를 최소 300dpi로 보정.
 - 밝기, 대비 자동 보정 / 다국어 문서는 OCR 언어설정 (lang = "kor+eng")

Multimodal: 서로 다른 type의 데이터를 동시에 처리하는 AI.

→ text / image / tabular / equation / speech (실제 음성)

<Transformer>

각각은 인코딩한다음, Transformer로 통합해서 복잡하게 처리한다.
(Attention 메커니즘)

시작
이미지 → 작은 패치로 분할 → 임베딩 (벡터화)
텍스트 토큰화 → 임베딩
동일한 Transformer Encoder에 이미지 & 텍스트 벡터 입력
이미지와 텍스트가 서로 attention 해서 의미를 갖음. (관계성)

Attention: 텍스트 나 이미지 같이 서로 다른 type의 데이터를 동일한 Transformer 임베딩 공간에 임베딩해서 서로의 위치관계에 집중할 수 있도록 함. (self attention & multi attention)

문서 레이아웃 분석 (object detection)

문서 (보판, PDF, 논문)에서 제목, 본문, 표, 이미지, 수식, 캡션 등 시각적 요소를 인식하고 구분함.

VDU에 경우 단순 OCR만 하게 되면 문서의 구조 읽기 때문에 레이아웃 정보 반드시 필요.

~~Bounding Box~~ (bbox) Confidence Score
⇒ 여러 bbox가 실제 해당 클래스인 확률 (임계값 이상이면 해당 클래스로 판단)
이미지에서 특정 개체를 감지하는 시작점

~~MAP~~ (object detection 평가치)
MAP 값이 높을수록 precision 높음. (성능 good)

<Yolo 기반 문서 분석> YOLOv11m - doclaynet (문서 전용 데이터셋으로 pretrained 된 모델)

Yolo: 단일 신경망으로 이미지 전체에서 객체를 한번에 감지.
문서 분석에서 텍스트, 표, 이미지, 수식 등 (따라가 감지) 가능. → 다양한 객체 동시 탐색 가능.

Yolo 구조: 이미지를 입력으로 변환 / 각 영역이 객체 존재 여부 (bbox) / 클래스 예측.

<Transformer 기반 접근>

• Layout LM (Microsoft)

OCR로 추출한 텍스트 토큰

각 토큰의 BBox 좌표

이미지 임베딩

→ Attention으로 문서복합 이해

• Donut

네이버가 개발한 OCR-free 모델

CNN + Transformer

OCR단위 생략 → 속도 / 정확도 개선

영역 예측 (텍스트의 논리 순서)

성능 평가 지표 35% 차지

가장 기본 방식 box 좌표값 정렬

1차 정렬: y좌표(세로)

2차 정렬: x좌표(가로)

한치: 왼쪽면 다 읽기 전에 오른쪽면이 먼저 정렬될.
테이블 내 순서 오류

⇒ 언타일링 문서에 적용 / 간단하고 빠름

• 그래프 기반 순서 추론 (Text flow Graph) ⇒ 구조적, 논리 순서에 강함. / 느리고 튜닝 필요 (hyper parameter)

Text block은 노드, 블록 간 읽기 관계를 엣지로 표현한 graph 구성

위상 정렬 (Topological Sort) 사용해 순서 예측.

각 텍스트 블록의 위치, 크기, 정렬 상태 비교

연접 블록 간 비방향성 연결 (아래에 있는 블록 후순위)

모든 노드를 그래프 구성

위상 정렬로 순서 결정

상용 캘리그래피 자동화: YOLO로 BBox 추출 → PaddleOCR로 Text 추출 → Text flow graph order prediction.
object detection ⊕ 텍스트 순서 예측

TTA (Test time Augmentation)

TTA: 주변화에 입력 이미지를 여러 방향으로 변형 (증강) 해서 여러 결과를 추출 → 평균 값 향상시키기

TTA 효과: 있잖아 있는 텍스트 블록 → 정렬 매개변수 → TTA로 정확도 개선
OCR 에러로 Text 유실 → 여러 변형에서 텍스트 비교해 가장 신뢰도 높은 순서 추출.
더 정확한 예측을 하는 것.

해리진 / flip / 해상도 변형 등 Augmentation → Voting / 평균 순서로 최종 결정

Text flow graph ⊕ TTA 고도화 전략

• 관계 무작 → edge weight 채서 soft graph 구성.

• 텍스트 정보 기반 보정.

★ 임베딩 기반 의미 연결: LayoutLM처럼 텍스트 ⊕ BBox 임베딩
이 블록이 다음 블록과 논리적으로 이어지는 것을 판단하는
classification 모델로 붙이는 것도 가능.

• 2단 레이어아웃, 테이블이 섞인 문서 → 클러스터링 + 지역 순서 정렬.

그래프 기반 Topo Sort: 위치 좌표 위주 정렬 (bbox)

LayoutLM 블록: 텍스트 의미 + 위치까지 고려한 정밀한 순서 판단 가능.

→ 서로 다른 기준으로 정렬 ⇒ Score voting / hybrid / 앙상블 구조로 결합

- 멀티모달 AI**
- Vision + Language 모델: 이미지와 텍스트 각각 임베딩 → 하나의 임베딩 공간에서 처리 (이미지 & 텍스트)
 - Layout LM: 텍스트 토큰 + BBox (문서 구조) + 이미지 특징은 하나의 Transformer로 처리.
 - Donut: 이미지를 보면 바로 구조화된 텍스트를 디코딩 (OCR free 모델)
 - ↳ 메인 대신 보조 라인 (엔드투엔드 푸른 앙상블 & 산패 케이스 백업용)
 - BIP2 (이미지 & 텍스트 pre-trained model): 캡셔닝 / QA 등 범용 멀티모달

Reading order prediction: Graph + Layout LM KNN (후보제한)

상-리소스 제복 → LORA 3 파라미터수 / VRAM 시간 대폭 감소 (메모리 가중치 제복)
↓
정중 fine-tuning

Graph prior 후보셋이 같은 가중치를 그래프 제한
↓
Layout LM pair score (A, B) 쌍은 Layout LM 입력으로 만들어 이진화 확률 추정 (hard / random negative 포함)

↓
Score fusion (late ensemble) ⇒ geometric ⊕ LM score
Topo sort 와 Layout LM score ensemble
↓
우로 방향 가중치 만들고 사이클 제거 (가중치 4는 edgedrop)
위상정렬 & local tree breaker로 안정화

Score Fusion:

- Topo sort (가중치 score)
- 가중치로 점수 정렬
- KNN의 후보제한 (리소스 절약)
- 0°/90°/180°/270° TTA 적용
- Layout LM (이미지 score)

↓
Pair score (Layout LM 출력 점수)
Lora (정중 fine-tuning)
hard / random negative

가중합 ensemble ⇒ 최종 그래프 edge 후보 ⇒ Topo sort로 최종 순서 결정.

사이클 제거: 가중치 큰 edge부터 끊고, 사이클 생기면 그 엣지는 재스
로인 타이머: in degree = 0 후보가 많을 때 자전스런 임기순서에 맞춰 우선순위를 정해서 블록 결정

PPTX → 이미지 변환

Lo / OCR / Layout LM 모두 페이지 단위 이미지를 만든다.

해상도 / 색공간 / 페이지 순서가 뒤섞이면 뒤에 전색영향 ⇒ "이미지 전체가 중요!"

<Alignments 데이터 증강>

- 레이아웃 보존
- YOLO BBox와 동일 변환을 적용하여 GT가 혼동되지 X
- OCR 학습에 광학 (Photometric) 변화가 특히 중요

YOLO 학습증강 (레이아웃 인식증강, 보수적 변화가 강함)

⊕

OCR 크롭증강 (노이즈가 적은 텍스트 boxes 간 백색으로 광학노이즈 변형)

OCR 크롭 주의사항 {

- 삽입 회전 / 퍼스펙티브 X
- train/validation split 전 증강함. ⇒ 검증/평가에는 미적용.



OCR 후 텍스트 증강 {

- 8자 정제화
- 화이트 스페이스/행바꿈 / 하이픈 처리 ⇒ 자연스러운 문장복합 (후처리)
- LaTeX/수학 수식 정제화 (수식의 경우)

- BBox/order/text 라벨링관리 ⇒ 도판간 호환 단위 0-1000 정제화 (점수) 권장
- Ranking order V_1, V_2 모두 자함 ⇒ V_2 최종제출 (V_1 은 디버깅용)
- ID = doc-id + 페이지 index + block index → 재식별을 위한 ID 라벨

YOLO 증서처리 (레이아웃 보정) → OCR 크롭 (텍스트 광학보정) → Paddle OCR (텍스트 정제) → 텍스트 후처리 (정제)

YOLO fine-tuning (문서 레이아웃 학습)

모델 = yolov11m - doclayout.pt

imagesz = 1024 / batch = 8~16

클래스 불균형이면 class weights 적용 / 과적합시 imagesz 960/896으로 낮춤.

모델 최적화

- YOLO 조정: 80~120epoch, best AP 기준 선택
- OCR 가이드 retry threshold 튜닝
- clean-text 를 유효성 점검

Layout LM pair classifier 학습 방법

Score function 가중치 하이퍼 파라미터 튜닝

- 11점 TTA (geo) 하이퍼 임계값 조정
- 시간 예산 simulation → fail-safe 구조 설계

precision : 맞힌 예측한 것 중 실제로 맞은 확률

$$\frac{TP}{TP+FP}$$

recall : 실제로 맞은 것을 모델이 맞히 예측한 확률

$$\frac{TP}{TP+FN}$$