

Project Progress Report

프로젝트 진행 상황 보고서

Carbot팀 안성규, 박상범, 김세희, 문벼리

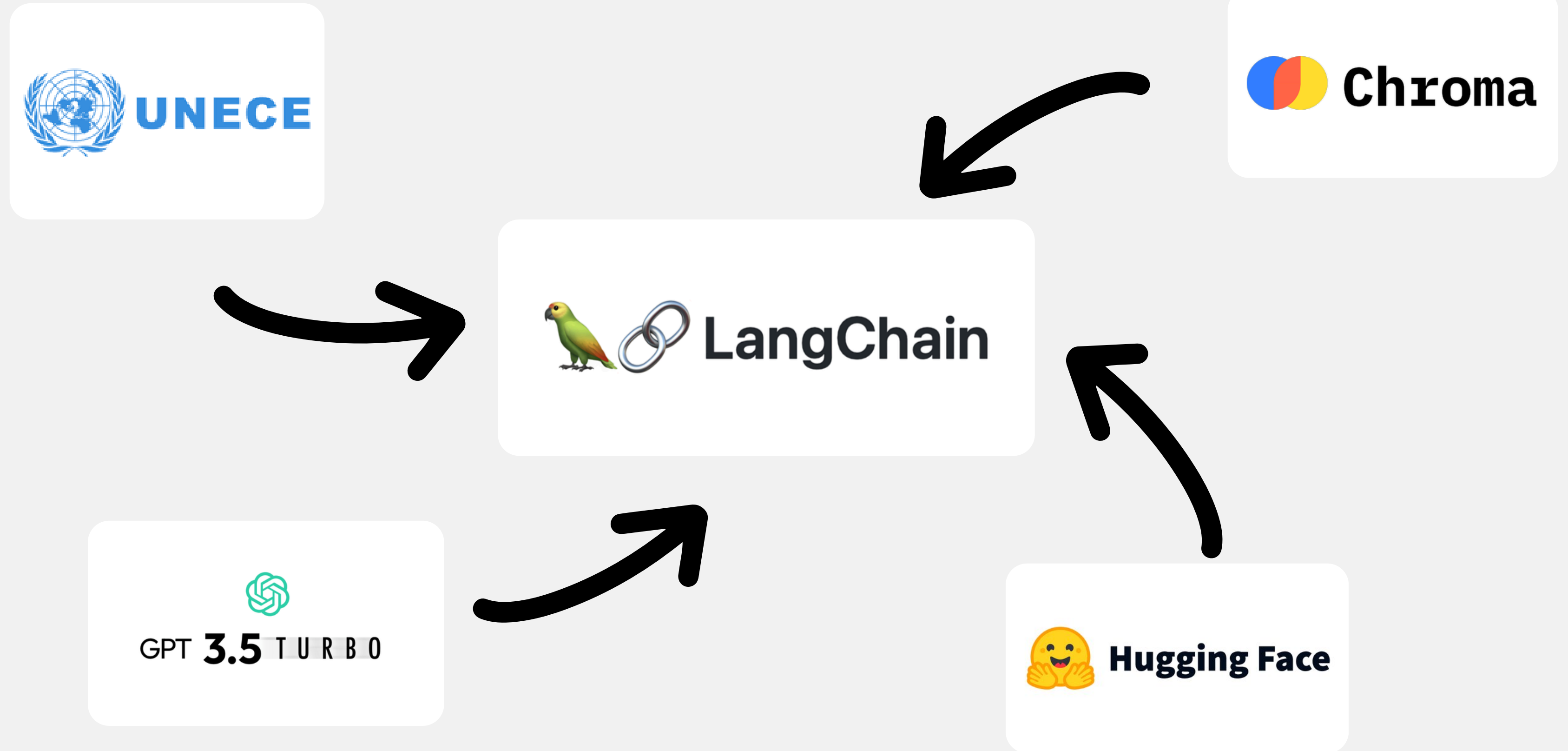


Contents

- 01 RAG 시스템 현황
- 02 데이터 전처리
- 03 고려해야 할 점
- 04 앞으로의 계획 및 발전 방향



RAG 모델 구현



RAG 성능 평가

정확성, 적합성, 완전성, 일관성, 이해도 순으로 고려

개념

Q1. How does the testing process for braking systems work?

Q2. What is an Antilock Brake System (ABS)?

Q3. How do I configure the parking brake system?

Q4. Tell me all the things I should consider if I collide with a sedan (Hyundai Sonata) in front of me in the same lane.

4개의 질문 -> 대체로 잘 대답

- 확인 사항
 - 부록의 내용을 정확하게 얘기할 수 있는지 확인 필요.
 - 차 이름만 넣어도 분류를 할 수 있는지 확인 필요. (아마 안 될 것으로 추정)
- 더 만들 질문
 - 다른 구성, 개념, 과정에 대한 질문
 - 4번 질문처럼 자세하고 구체적인 질문

표

2개의 질문 -> 1개는 질문에 맞는 답변
1개는 질문의 포인트가 답이 없는 부분이었기에 애매함.

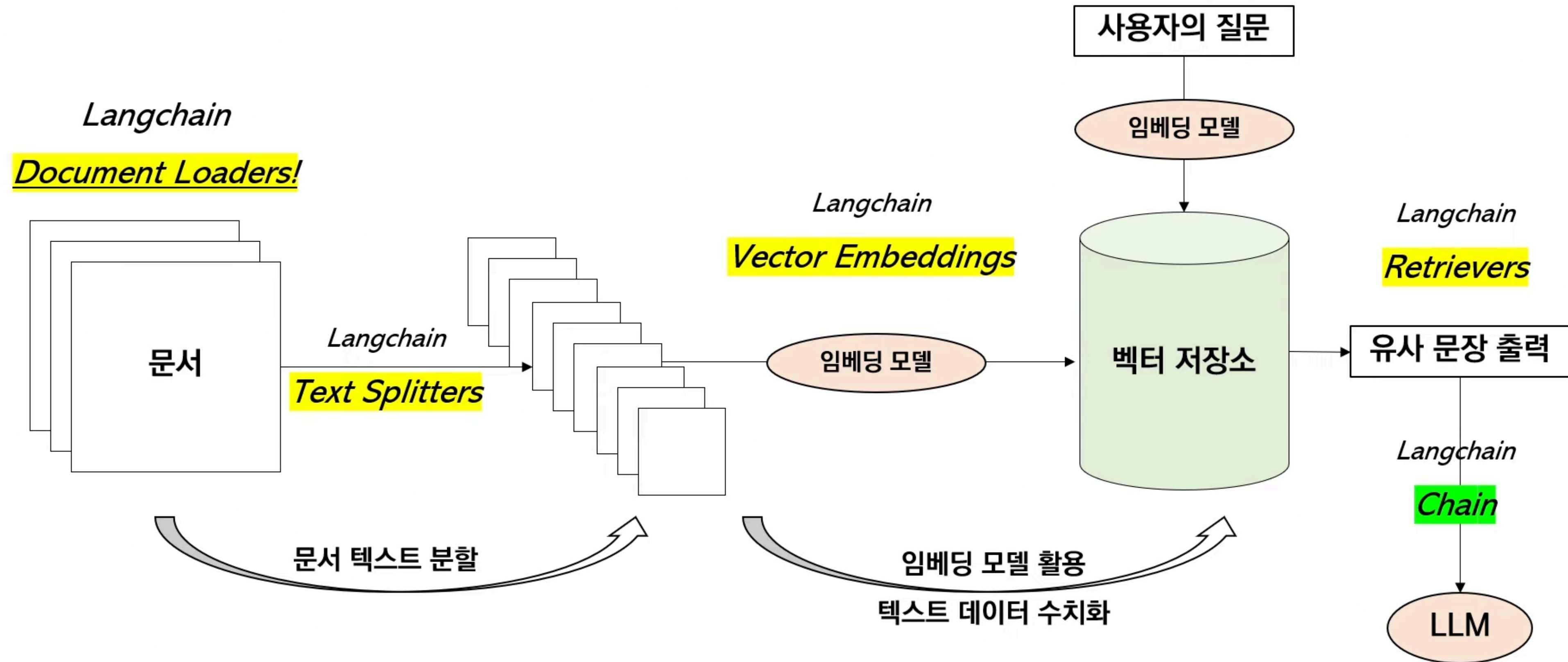
- 확인 사항
 - 좀 더 검증이 필요함.

그림&수식

그림 관련 2개, 간단한 수식 관련 1개의 질문
-> 답변 어려움.

- 확인 사항
 - 수식-텍스트 변환을 거쳐도 동일한지 확인

RAG 프로세스 모형



RAG 프로세스 최적화 기법

텍스트 포매팅

Json 혹은 Markdown 형식으로 규격화

텍스트 분할

의미론적 유사도에 따른 텍스트 분리

텍스트 임베딩

임베딩 모델 변환
임베딩 벡터 캐싱

벡터스토어

리트리버에 더욱 용이한 FAISS 사용

리트리버

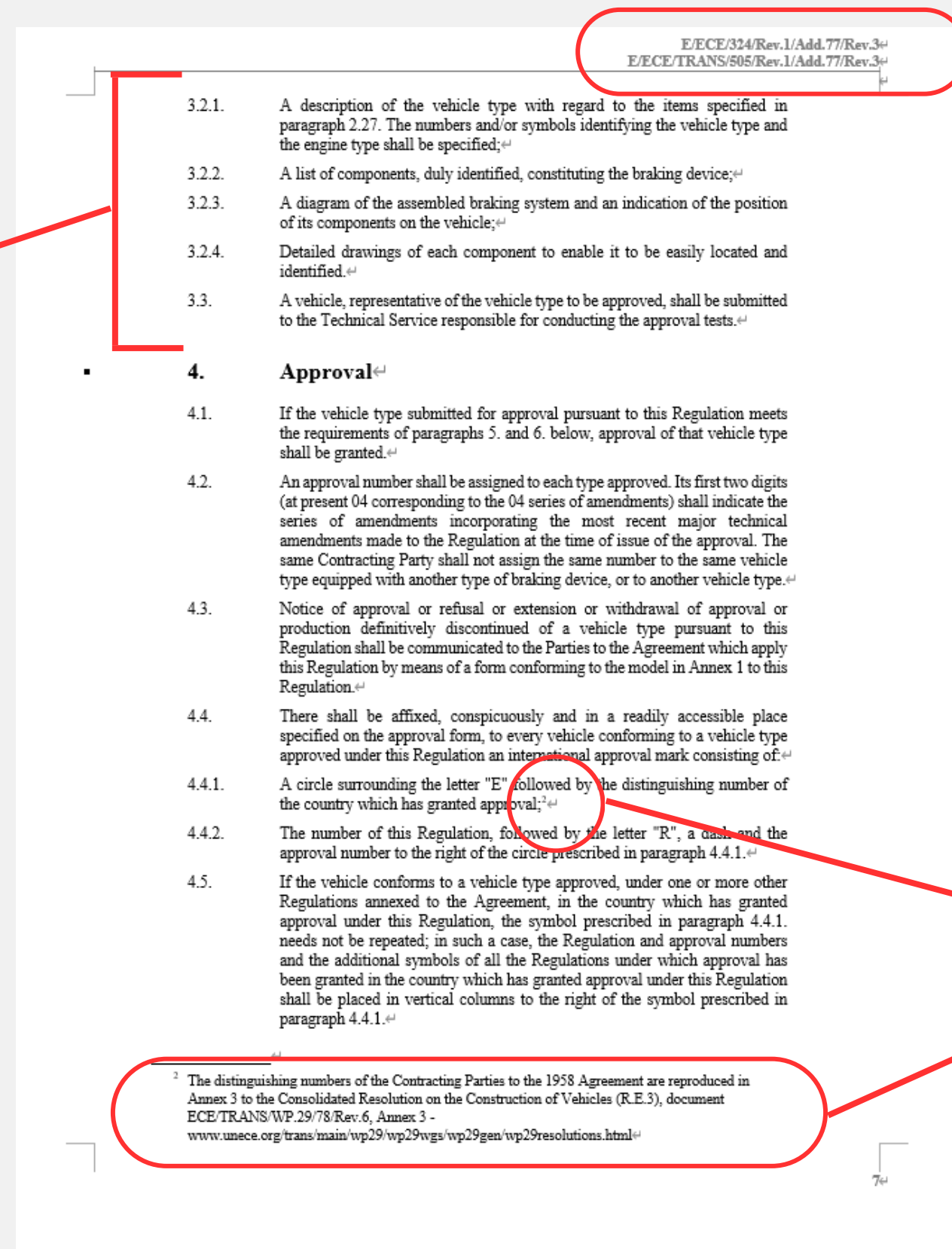
기반 모델 고도화
체인 기법에 따른 텍스트 품질 고도화
프롬프트 엔지니어링을 통한 답변 생성 고도화

데이터 전처리 - 현재 방식의 문제점

정보가 부분적으로 입력되어
온전한 맥락을 참고하지 못함

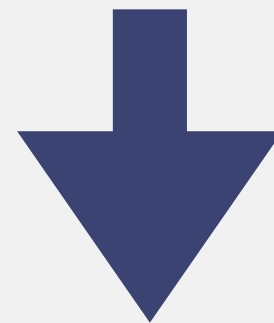
필요 없는 정보가 반복적으로 입력

보충 정보(주석)과 거리가 있어
모델에게 혼란을 줄 가능성



데이터 전처리 - 현재 방식의 문제점

<i>Vehicle decelerations</i>	<i>Signal generation</i>
$\leq 0.7 \text{ m/s}^2$	The signal shall not be generated
$> 0.7 \text{ m/s}^2$ and $\leq 1.3 \text{ m/s}^2$	The signal may be generated
$> 1.3 \text{ m/s}^2$	The signal shall be generated



테이블의 경우, 텍스트로 추출이 가능하나
해당 텍스트가 “테이블”임을 인지하지 못할 가능성 존재

Vehicle decelerations Signal generation
 0.7 m/s^2 The signal shall not be generated
 $> 0.7 \text{ m/s}^2$ and $\leq 1.3 \text{ m/s}^2$ The signal may be generated
 $> 1.3 \text{ m/s}^2$ The signal shall be generated

데이터 전처리 - 현재 방식의 문제점

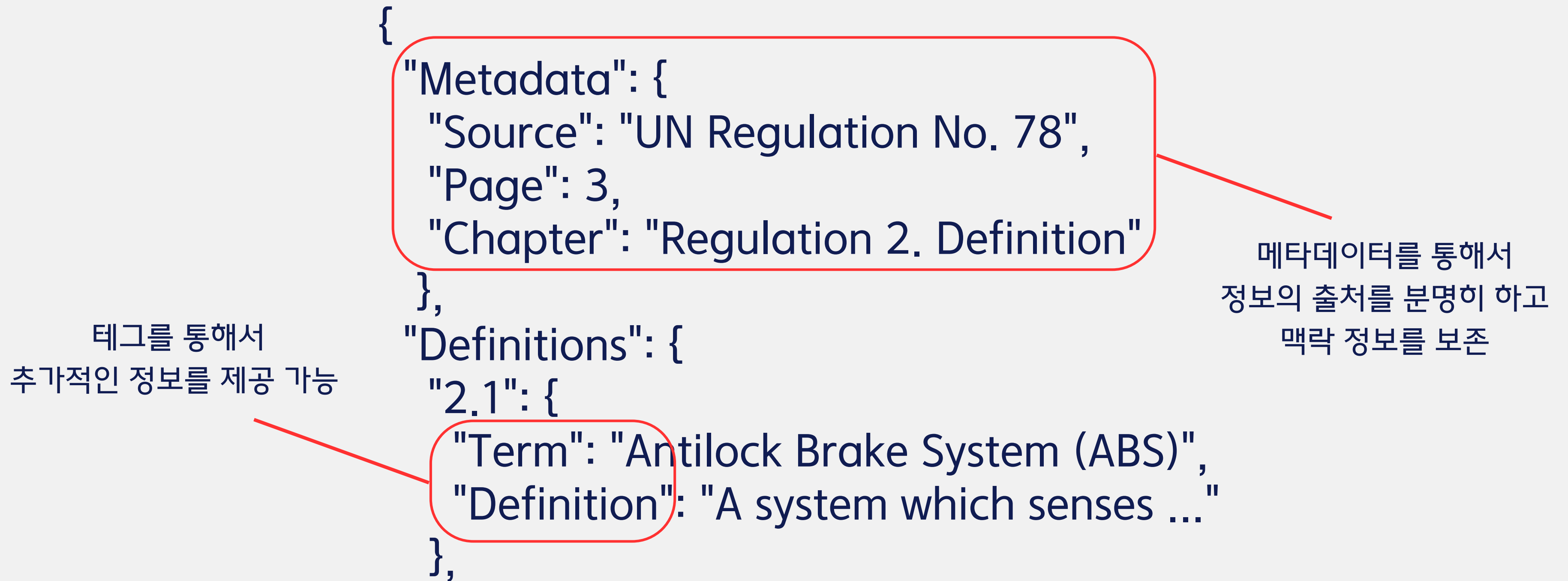
$$d_m = \frac{V_b^2 - V_e^2}{25.92 \cdot (S_e - S_b)} \quad \text{in m/s}^2$$



수식의 구조가 완전히 와해되고
대부분의 정보를 상실함

92.252 2
b ee b
mS SV Vd-□-=
in m/s2

데이터 전처리 - 포매팅으로 얻을 수 있는 이점들



예시: jsonl 형식

데이터 전처리 - 포매팅으로 얻을 수 있는 이점들

- 필요한 정보만 선별하여 모델의 혼란을 방지
- 분산 되어 있는 정보(주석, 부록 참조)를 하나로 취합
- 특수한 데이터(테이블, 이미지, 수식)의 의미를 온전히 보존

데이터 전처리 - 후보 포맷 형식들

Jsonl

```
{
  "Table 2-1": [
    {
      "Deceleration": "≤ 0.7 m/s2",
      "SignalGeneration": "The signal shall not be generated"
    },
    {
      "Deceleration": "> 0.7 m/s2 and ≤ 1.3 m/s2",
      "SignalGeneration": "The signal may be generated"
    },
    {
      "Deceleration": "> 1.3 m/s2",
      "SignalGeneration": "The signal shall be generated"
    }
  ]
}
```

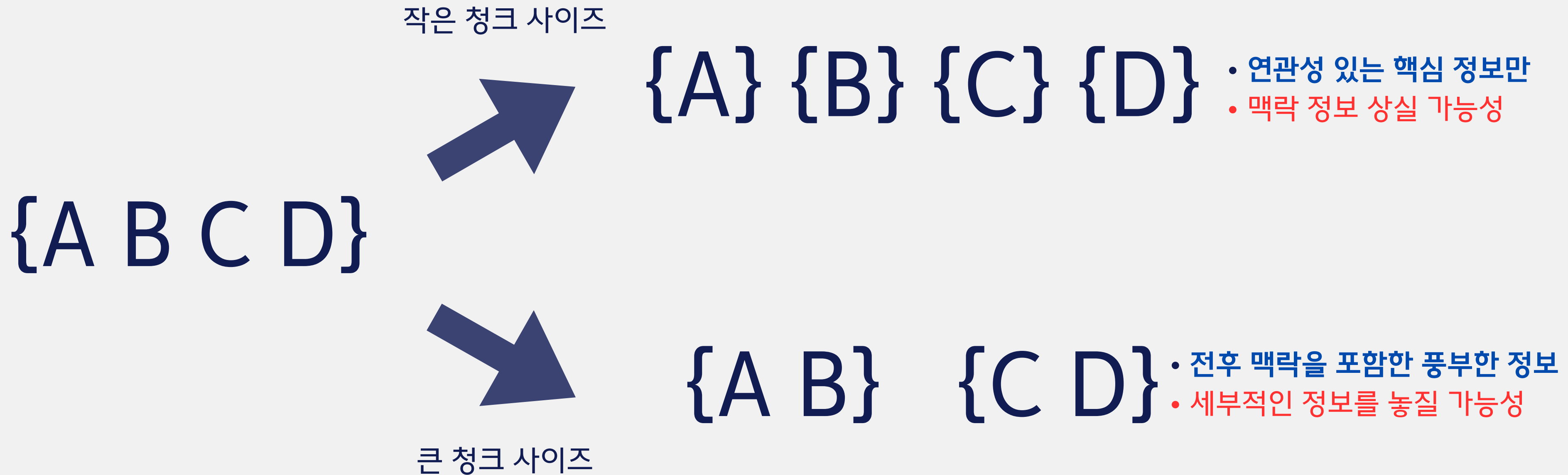
- 높은 유연성
- 메타데이터 기입 용이

Markdown

+-----+	+-----+	+-----+
> *Vehicle decelerations*	> *Signal generation*	
+=====+	+=====+	+=====+
> ≤ 0.7 m/s ²	> The signal shall not be generated	
+-----+	+-----+	+-----+
> \> 0.7 m/s ² and ≤ 1.3 m/s ²	> The signal may be generated	
+-----+	+-----+	+-----+
> \> 1.3 m/s ²	> The signal shall be generated	
+-----+	+-----+	+-----+

- 높은 정형성
- 문서 학습량이 높음

텍스트 분할 (칭킹) 사이즈 조절



분야별 발전 방향

데이터 전처리

- 기본적인 데이터 포매팅 코드 구현을 통해 json, markdown 성능 비교
- 의미적 손실을 줄일 수 있는 최적화된 chunk size 찾기
- 문서 내에서 상호 참조하는 부분 구현 방법 찾기

RAG

- RAG 성능 평가를 위해 자동차 도메인 특화 질문 만들기
- Gpt -4 와 기반 모델 및 임베딩 모델 간의 비교

파인 튜닝

- 파인튜닝 오픈 소스 모델 선정
- 파인 튜닝 코드 구현

앞으로의 계획

항목	과업	7월			8월			
		2주차	3주차	4주차	1주차	2주차	3주차	4주차
데이터 전처리	데이터 chunk size 결정	<div><div></div></div> <div>7/9 ~ 16</div>						
	json, markdown 포맷 변경	<div><div></div></div> <div>7/9 ~ 19</div>						
RAG	코드 구현	<div><div></div></div> <div>7/10 ~ 23</div>						
	모델 평가 질문 선정	<div><div></div></div> <div>7/10 ~ 29</div>						
	최적화		<div><div></div></div> <div>7/15 ~ 8/4</div>					
파인튜닝	모델 선정 및 코드 구현			<div><div></div></div> <div>7/24 ~ 8/12</div>				
	최적화				<div><div></div></div> <div>8/2 ~ 8/26</div>			



감사합니다