

# 역사 텍스트 디지털화 입문

OCR에서 데이터베이스까지

---

이미지 → OCR 추출 → 검수 → 구조화 → 데이터베이스

2026년 3월 3일 | 줌 비대면 특강

도구: Google AI Studio (Gemini) | 준비물: Google 계정

# 오늘의 목표



강의가 끝나면: **내 손으로 만든 검색 가능한 문장 DB 1개**

도구: Google AI Studio ([aistudio.google.com](https://aistudio.google.com))

준비물: Google 계정 (AI Studio + Google Sheets 공용)

방식: 강사 화면을 보며 동시에 따라하기

# 라이브 시연: AI가 읽습니다

## 1931년 국한문혼용 세로쓰기 철학서

이돈화, 《新人哲學》(천도교중앙종리원, 1931)

→ AI가 이것을 읽어냅니다

AI는 "자연어로 지시하는 연구 조수"  
코딩을 배우는 것이 아닙니다  
→ 무엇을 원하는지 말할 줄 알면 됩니다

⚠ 단, AI는 그럴듯하게 틀릴 수 있음  
→ 최종 판단은 반드시 연구자 본인

### 新人哲學

#### 第一編 宇宙觀

##### 第一章 宇宙

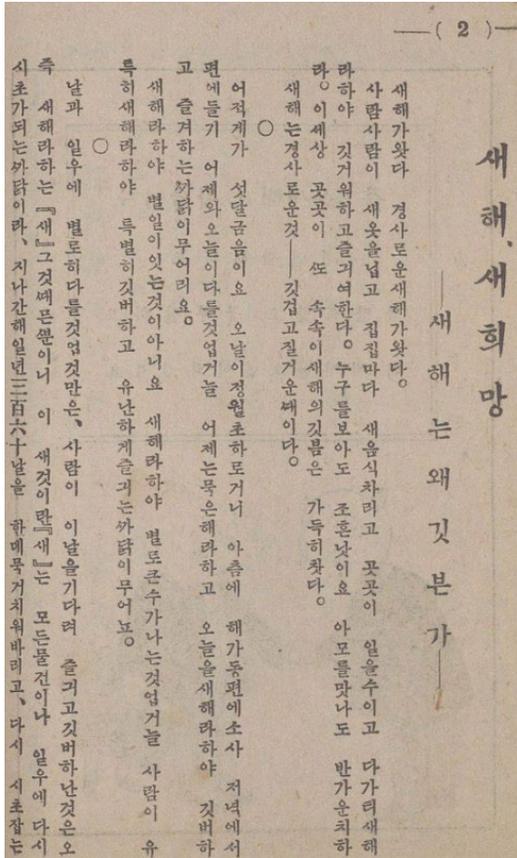
###### 甲、宇宙의 名義

『無窮한 이理致를 無窮히 살피면 無窮한 이 우주에 無窮한 내아닌가』 이것은 水雲의 노태이다  
水雲은 無窮의 意義를 自我의 中에서 發見하고 그를 吟味하여 스스로 無窮을 讚美한 것이다  
『無窮한 이물』이라 함은 곳 『한을』을 가트치하느말이니 無窮은 『한』을 意味한 말이며 한은  
『크』다느 뜻이다 朝鮮말의 尊 尊을 한길이라 하며 한아 버지들 큰아 버지라 함은 『한』과  
『은』 것이 同 한 것을 意味한 것이다 그러하여 『은』이라느 뜻은 量的 意味에 서는 範圍를 表  
象한 것으로 解釋할 수 있는데 空間上으로는 無窮의 範圍와 時間上으로는 通三界의 範圍를 總合한

第一編 宇宙觀

1

# 실습 ① OCR: 어린이 「새해 새희망」 (1924)



이 이미지는 1924년에 발행된 어린이 잡지입니다.  
세로쓰기 텍스트를 오른쪽→왼쪽 순서로,  
원문 그대로 추출해 주세요.

## AI Studio 세팅 (첫 실습 전 한 번만)

- ① 모델: Gemini 3.1 Pro (화면 상단)
- ② Temperature: 0.0 (오른쪽 Run settings)
- ③ Safety Settings: Off (오른쪽 하단)

※ 이미지 업로드: 채팅창 + 버튼 또는 드래그 앤 드롭

# OCR 결과 검수: 원문 대조는 필수

## 옛 표기 ≠ OCR 오류

정확히 읽은 것 (옛 표기 그대로)  
"깃븐가" → 기븐가 | "왓다" → 왔다  
"닙고" → 입고 | "조흔" → 좋은  
1920년대 표기법이 지금과 다를 뿐입니다

### ✗ AI가 틀릴 수 있는 부분

- 비슷한 한자 혼동: 大↔太, 己↔己, 日↔日
- 글자 누락 또는 없는 글자 생성

## 🔍 검수 원칙

1. 결과를 메모장에 복사
2. 원문 이미지와 나란히 놓고 대조
3. 글자 단위로 확인

**"반드시 원문과 대조 검수가  
필요합니다"**





# 전반부 정리

① 어린이 (1924) 한글 세로쓰기 [쉬움]

② 신인철학 (1931) 국한문혼용 세로쓰기 [중간]

③ 문교의조선 (1925) 일본어 구자체 [시연]

**OCR은 완벽하지 않습니다 — 반드시 원문과 대조 검수가 필요합니다**

▶ 후반부: 신인철학 텍스트를 검색 가능한 문장 단위 DB로 만듭니다



# 10분 휴식

Google AI Studio에 로그인이 아직 안 되신 분은 이 시간에 시도해 주세요  
[aistudio.google.com](https://aistudio.google.com)

Google Sheets도 열어두시면 좋습니다  
[sheets.google.com](https://sheets.google.com)

전반부에서 신인철학 C001+C002 OCR 결과를  
메모장에 저장해 두셨는지 확인해 주세요

# 구조화란 무엇인가

## 新人哲學 目次

### 第一編 宇宙觀

#### 第一章 한울

甲、한울의名義..... 一

乙、한울本體의屬性..... 三

甲、科學的考察로부터한울의進化..... 六

乙、科學的進化說과水雲主義의進化說과의差異..... 一

第三章 質的한울과氣一元實在體..... 一

甲、認識方法의一片..... 一

텍스트 뭉치 → 문장 DB = 검색·분석 가능  
서류 더미 → 엑셀 표로 정리하는 것과 같습니다

### 《新人哲學》 계층 구조 (A01 목차 기준)

第一編 宇宙觀 → E01-C00-S00-P00-S00

第一章 한울 → E01-C01-S00-P00-S00

甲, 한울의名義 → E01-C01-S01-P00-S00

乙, 한울本體의屬性 → E01-C01-S02-P00-S00

문장 한 줄 한 줄이 표의 한 행이 됩니다

# 문장 DB: 4열 구성

열 이름	설명	예시
local_id	계층적 위치 (E=편, C=장, S=절, P=문단, S=문장)	E01-C01-S01-P01-S01
page_info	원본 페이지 태그	C001
line_class	행 유형	STRUCT / TEXT
kr_text	텍스트 내용	『無窮한이理致를...』

**이 4열이면 검색·필터·정렬이 모두 가능합니다**

(완성 DB는 10열이지만, 핵심은 이 4열 — 나머지는 필요할 때 추가)

# 완성 DB 시연: 인내천요의 2,254행

같은 저자 이돈화의 다른 책 《人乃天要義》 (1924)

10월 2,254행의 완성된 문장 DB

Google Sheets에서 시연할 내용:

- ① 필터: line\_class로 STRUCT만 보기 → 목차 구조가 한눈에
- ② 검색: Ctrl+F → "한울" → 해당 문장들이 바로 찾아짐
- ③ 정렬: local\_id로 정렬하면 원문 순서 복원

▶ "오늘 신인철학 첫 두 페이지로 이것을 시작합니다"

# 실습 ③ 구조화: 신인철학 → 문장 단위 TSV

다음은 《新人哲學》(이돈화, 1931) 본문의 OCR 텍스트입니다.  
이 텍스트를 문장 단위로 분절하여 탭으로 구분된 TSV 형식으로 출력해 주세요.

열 구성:

- local\_id: 계층적 위치 (예: E01-C01-S01-P01-S01)  
E##=편, C##=장, S##=절, P##=문단, S##=문장
- page\_info: 원본 페이지 태그 (예: C001)
- line\_class: STRUCT(제목/소제목) 또는 TEXT(본문)
- kr\_text: 텍스트 내용

규칙:

- 편명, 장명, 절명은 line\_class를 STRUCT로 분류
- 본문은 의미가 완결되는 문장 단위로 분절하여 각각 한 행으로
- 첫 줄은 헤더(열 이름)로 출력 / 구분자는 탭(Tab)

C001 → C002 순서로 처리 (헤더 행은 한 번만)

⚠ "문장 경계가 여러분과 다를 수 있습니다"  
→ 이것이 연구자 판단이 필요한 지점

# DB 완성: Google Sheets

## ① 붙여넣기

- Google Sheets 새 스프레드시트 열기
- AI가 출력한 TSV 전체 복사 → A1 셀에 붙여넣기
- 탭 구분자로 열이 자동 분리되는 것 확인

## ② 필터 + 검색

- 첫 행 선택 → 데이터 → 필터 만들기
- line\_class 필터: STRUCT만 → 목차 구조 추출
- Ctrl+F → "한울" 검색 → 문장 바로 찾기

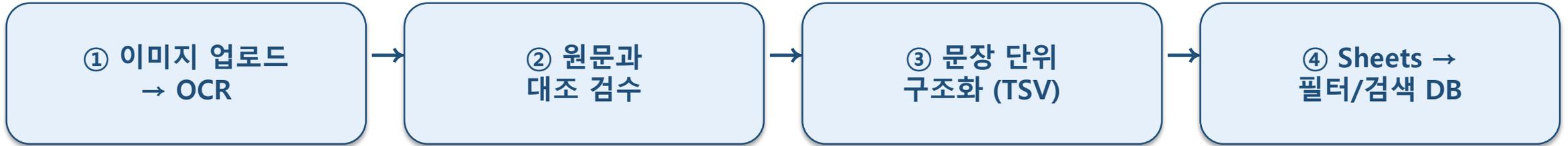
## ③ 성과 확인

- 이 스프레드시트 = 내가 만든 신인 철학 문장 DB
- 같은 방법 → 책 전체 → 2,000행 이상



"내 자료가 검색 가능한 데이터베이스가 되었다"

# 정리 + 다음 단계



## 다음 단계 안내

단계	내용	배포 자료
다음 단계 1	대량 정밀 OCR + 번역	일본 사상서 판독 및 번역 가이드
다음 단계 2	텍스트 정제·일괄 편집	VS Code 멀티커서 가이드
다음 단계 3	논문/자료 코퍼스 관리	VS Code 논문 검색 가이드

모든 자료(핸드아웃, 이미지, 결과물, 다음 단계 가이드)는 사후 배포됩니다