

From OCR to Multimodal LLMs

Building a Textbook Corpus with ALICE

Steven Denney
Leiden University

ALICE-SHARK User Meeting 2025
June 3, 2025

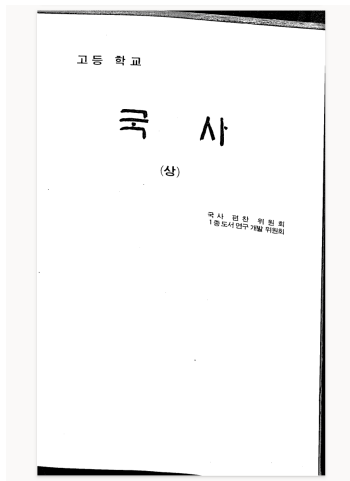
Why This Project?

- Textbooks are key to understanding civic education, national identity, and state discourse, especially over long periods of time.
- OCR quality on historical text in non-Latin scripts (e.g., Korean) is often poor, especially when using English-centric tools.
- By working with Korean text and a Korean-trained LLM (EXAONE), this project demonstrates the viability of non-English, script-specific LLM pipelines.
- ALICE enables local, high-performance language model inference – ideal for multilingual document workflows.

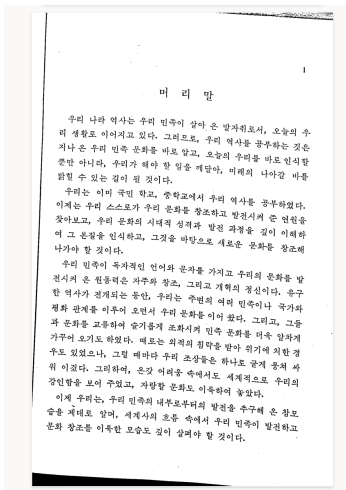
Corpus Construction Goals

- Build a clean, searchable corpus from scanned Korean textbooks.
- Compare and evaluate two workflows:
 1. Pipeline A: Traditional OCR + LLM cleaning
 2. Pipeline B: Vision-language model + LLM post-editing
- Optimize for ALICE: GPU-accelerated, fully local, efficient.

What We're Working With Here



Cover page



Introduction page

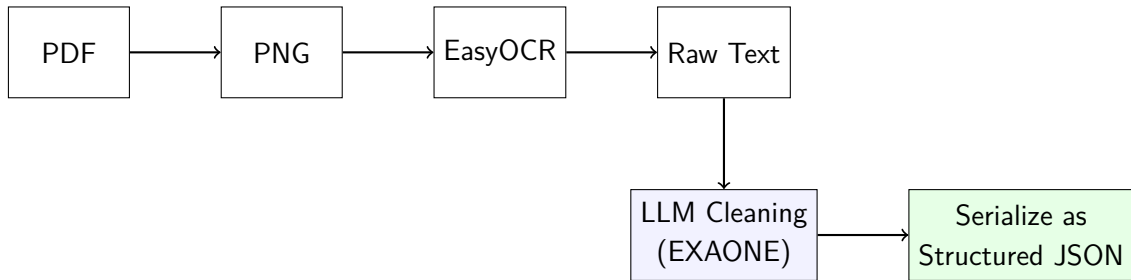
Language, Modularity, and Scale

- This pipeline is prototyped on Korean text — ideal for testing non-Latin scripts.
- EXAONE is used because it is trained on Korean — but other LLMs could be substituted.
- Pipeline is modular: Vision model, LLM, and post-processing components are interchangeable.
- Goal: develop a scalable, multilingual corpus-building framework for diverse scripts.

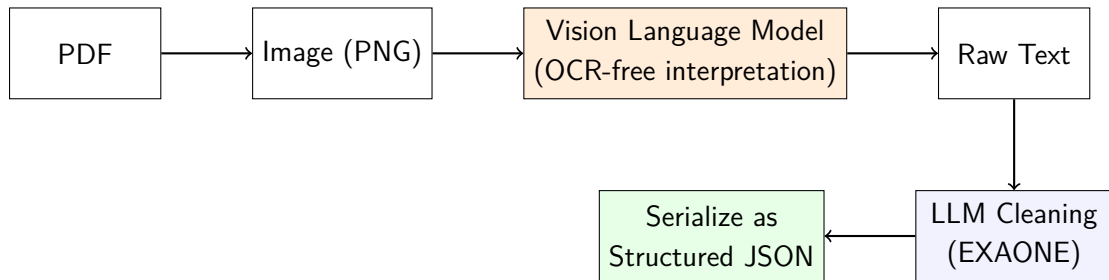
Pipeline Typical (Legacy Approach)

- The traditional pipeline relies entirely on OCR:
 1. PDF → Image conversion
 2. Image → OCR engine (e.g., Tesseract or EasyOCR)
 3. Output stored as raw or segmented plain text
 4. Optionally post-processed and structured manually or semi-manually into JSON or XML
- Works reasonably well for clean Latin-script sources, but performs poorly on layout-rich or non-Latin texts.

Pipeline A: OCR-Based Workflow



Pipeline B: Prototype Workflow (Vision Language Model)



Model Infrastructure: Ideal and Tested

Ideal setup for scalable, high-quality inference:

- EXAONE-Deep 70B (Q5_K_M) – Korean-specific LLM used for post-processing and normalization. High performance; fits on ALICE GPU nodes via quantization.
- LLaVA-1.5 + Mistral 7B – Open-source multimodal model for direct image-to-text interpretation. Supports OCR-free pipeline design, and runs efficiently with quantized weights on ALICE GPUs.

Tested alternatives via local prototyping:

- EXAONE 7B – Used for early-stage post-cleaning; faster but with slightly reduced performance.
- MiniGPT-4, LLaVA-1.5 – Lightweight VLMs tested for proof-of-concept inference. Run on consumer GPU, but slow(er) and less reliable than ideal models.

Comparing the Pipelines

Feature	Pipeline A (OCR)	Pipeline B (VLM)
Accuracy	Moderate, layout lost	Promising, preserves layout
Speed	Stable but multi-step	Faster per step, but newer
Robustness	OCR is (semi-)reliable	VLMs still evolving
Deployment	Simple batch processing	Needs better fallback/error handling

Pipeline B: Built for ALICE?

- Entire pipeline executes on ALICE — no cloud dependencies, full data control.
- Uses quantized models (e.g., EXAONE 70B Q5_K_M) for efficient GPU inference on single nodes.
- Vision-Language inference and LLM post-processing both benefit from ALICE's parallel architecture.
- Enables testing, tuning, and benchmarking of large multilingual models in-house.
- Future: add automated fallback from VLM to OCR — and batch-parallel scaling across textbook collections.

Next Steps and Broader Considerations

- Fine-tune VLM interpretation on Korean textbook samples.
- Release annotated corpus + pipeline documentation for other researchers/socializing.
- Assess tradeoffs: Is multimodal extraction truly necessary? Does it justify the extra compute and complexity?

Thank You

Project materials and code (work in progress):

<https://github.com/scdenney/textbook-pipeline/tree/main>

Contact: s.c.denney@hum.leidenuniv.nl