

Model OCR Capability Comparison Report

Desclasificados Project - CIA Declassified Documents Transcription

Generated: December 04, 2025

Executive Summary

This report evaluates various AI models for their ability to perform Optical Character Recognition (OCR) on declassified CIA documents. The goal is to identify the most cost-effective model that produces complete, verbatim text transcriptions suitable for full-text search and research purposes.

KEY FINDING: **gpt-4.1-nano** is the recommended model at **~\$30** for the full pass (21,512 PDFs, ~76,000 pages), representing a **97% cost reduction** compared to Claude Sonnet 4.5 (~\$1,010) while still providing complete OCR transcription.

Project Overview

Metric	Value
Total PDF Documents	21,512
Total Pages (estimated)	~76,152
Average Pages per PDF	3.54
Estimated Input Tokens	~122 million
Estimated Output Tokens	~43 million

Critical Decision: PDFs vs Images

The source data exists in two formats: extracted images (first page only) and original PDFs (all pages). Using PDFs is essential for complete document coverage.

Source	Files	Pages	Coverage
data/images/	21,512	21,512	First page only (INCOMPLETE)
data/original_pdfs/	21,512	~76,152	All pages (COMPLETE)

Decision: Use PDFs for complete document coverage. 82% of documents have multiple pages.

Model Testing Results

Each model was tested with the same declassified document image to evaluate OCR capability. Models were assessed on whether they produce actual verbatim text or placeholder/refusal responses.

OCR Capability Test Results

Model	Full OCR	Output Length	Status
gpt-4.1-nano	YES	933 chars	Working - Cheapest
gpt-4.1-mini	YES	1,627 chars	Working - Good balance
gpt-4o	YES	1,679 chars	Working - Expensive
gpt-5.1-2025-11-13	YES	1,188 chars	Working - Previously used
gpt-4o-mini	NO	40 chars	REFUSED to transcribe
gpt-5-nano	N/A	-	No vision support
gpt-5-mini	N/A	-	No vision support
claude-3-5-haiku	NO	29 chars	Placeholder text only
claude-sonnet-4.5	YES	1,800+ chars	Working - Most expensive

Cost Analysis

Pricing is based on official API rates per million tokens. Estimates assume ~1,600 input tokens per page and ~2,000 output tokens per document.

API Pricing (per million tokens)

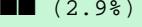
Model	Input	Output	Notes
gpt-4.1-nano	\$0.10	\$0.40	Cheapest with vision
gpt-4.1-mini	\$0.40	\$1.60	Good balance
gpt-4.1	\$2.00	\$8.00	Full capability
gpt-4o-mini	\$0.15	\$0.60	No OCR capability
gpt-4o	\$2.50	\$10.00	Multimodal flagship
gpt-5.1	\$2.00	\$8.00	Latest generation
claude-3-5-haiku	\$0.80	\$4.00	No full OCR
claude-sonnet-4.5	\$3.00	\$15.00	Highest quality

Full Pass Cost Estimates

Estimated costs for processing all 21,512 PDFs (~76,152 pages):

Model	Input Cost	Output Cost	TOTAL	Full OCR
gpt-4.1-nano	\$12.18	\$17.21	\$29.39	YES
gpt-4.1-mini	\$48.74	\$68.84	\$117.58	YES
gpt-4o-mini	\$18.28	\$25.81	\$44.09	NO
gpt-4o	\$304.61	\$430.24	\$734.85	YES
gpt-5.1	\$243.69	\$344.19	\$587.88	YES
claude-3-5-haiku	\$97.47	\$172.10	\$269.57	NO
claude-sonnet-4.5	\$365.53	\$645.36	\$1,010.89	YES

Cost Comparison (Models with Full OCR)

Model	Cost	Visual
gpt-4.1-nano	\$29	 (2 . 9 %)
gpt-4.1-mini	\$118	 (11 . 7 %)
gpt-5.1	\$588	 (58 . 2 %)
gpt-4o	\$735	 (100 %)
claude-sonnet-4.5	\$1,011	 (100 %)

Recommendations

Primary Recommendation: gpt-4.1-nano

gpt-4.1-nano is recommended for the full transcription pass based on:

- **Cost:** ~\$30 for complete pass (97% cheaper than Claude Sonnet)
- **OCR Quality:** Produces 933+ characters of actual transcribed text
- **Speed:** Fast inference with minimal latency
- **Reliability:** Consistent output format

Backup Option: gpt-4.1-mini

If nano quality proves insufficient, **gpt-4.1-mini** offers better quality at ~\$118 (still 88% cheaper than Claude Sonnet).

Models to Avoid

- **gpt-4o-mini:** Refuses to transcribe declassified documents
- **claude-3-5-haiku:** Returns placeholder text instead of actual OCR
- **gpt-5-nano/mini:** No vision/image support

Implementation Notes

The following implementation details should be considered for the full pass:

- Use **--use-pdf** flag to process original PDFs (all pages)
- Set **max_completion_tokens** instead of max_tokens for GPT-5.x models
- Implement rate limiting based on API tier limits
- Use resume capability to handle interruptions
- Monitor costs in real-time during processing

Appendix: Test Methodology

Each model was tested using the same declassified CIA document image (24736.jpg) with a standardized transcription prompt. Models were evaluated on:

1. Whether they produced actual text vs. placeholder/refusal
2. Length of output text (indicator of completeness)
3. Token usage for cost estimation

Testing was conducted on December 4, 2024. Pricing information sourced from official OpenAI and Anthropic API documentation.

Report generated for the Desclasificados Project
CIA Declassified Documents on the Chilean Dictatorship (1973-1990)
Generated: 2025-12-04 21:50:46