

CULTUREMAP-IR: Cultural Unification and Linguistic Textual Utilization for Regional Extraction and Mapping of All Iranian Provinces, IRan

Asal Meskin* (401106511), Alireza Mirrokni* (401106617)

Computer Engineering Department, Sharif University of Technology

*These authors contributed equally to this work.

<https://github.com/CULTUREMAP-IR>

Contents

1	Introduction	2
2	Challenges and Motivations	2
3	Method 1: PDF-based JSON Extraction with Gemma-3	3
3.1	Overview	3
3.2	Code and Explanations	3
3.2.1	Dependency Imports	3
3.2.2	Configuration and Logging	4
3.2.3	JSON Cleaning Utilities	4
3.2.4	Text Extraction and Normalization	5
3.2.5	Chunking Strategy	5
3.2.6	Prompt Template	6
3.2.7	Model Initialization	6
3.2.8	Deduplication and GPU Cleanup	6
3.2.9	Chunk Processing	7
3.2.10	Orchestration and Output	8
4	Method 2: GPT-4 with Autonomous Web Search	8
4.1	Overview	8
4.2	Code and Explanations	9
4.2.1	Configuration and API Key Verification	9
4.2.2	Prompt Template with Search Instruction	9
4.2.3	Chunk Processing via ChatCompletion	9
4.2.4	Concurrent Execution	10
4.2.5	Aggregation and Output	10
5	Labeling Policy and Dataset Refinement	10
6	Results and Analysis	11
6.1	Quantitative Comparison	11
6.2	Field-Word Count Distributions	12
7	Conclusion	12

1 Introduction

In this project, we present two complementary automated methods for extracting and structuring cultural, geographic, and environmental data for six distinct Iranian provinces: **Isfahan**, **Fars**, **Bushehr**, **Chaharmahal and Bakhtiari**, **Hormozgan**, and **Kohgiluyeh and Boyer-Ahmad**. These regions span central desert basins and mountain plateaus, fertile coastal plains on the Persian Gulf, and densely forested highlands, each embodying unique traditions, dialects, and ecological challenges.

- **Method 1** leverages a locally hosted large language model (Gemma-3) to parse PDF documents—ranging from provincial geography texts to scientific reports—directly into a strict JSON schema.
- **Method 2** builds on the same extraction utilities but instructs GPT-4 (via the OpenAI API) not only to map text to JSON, but also to autonomously query authoritative online sources (e.g. Wikipedia and official tourism sites) when necessary to fill gaps or update figures. Although early experiments were performed by manually pasting prompts into the OpenAI playground due to quota limits, the provided code will execute end-to-end once a valid `OPENAI_API_KEY` is configured.

2 Challenges and Motivations

Building a high-quality, culturally informed dataset for six diverse Iranian provinces entails a range of technical and conceptual challenges. At the PDF ingestion stage, we must contend with inconsistent formatting, embedded images, and non-standard fonts in provincial geography texts. Extracted text often contains noise—header/footer artifacts, page numbers, and stray line breaks—that must be normalized without losing semantic content. Splitting long documents into overlapping chunks requires careful tuning: too small, and we lose context; too large, and we exceed the model’s input limits.

A core difficulty is the relative scarcity of high-quality Farsi training and reference data. Many regional sources are unpublished or available only in print, and online resources in Farsi often lack consistent metadata or structured markup. This necessitates web crawling, OCR, and manual validation steps to bootstrap a reliable corpus. Furthermore, dialectal variations across provinces introduce inconsistencies in spelling and terminology, complicating automatic extraction and normalization.

On the LLM side, guiding Gemma-3 or GPT-4 to produce strictly valid JSON demands a meticulously crafted prompt template and robust post-processing to strip out markdown fences, stray commas, and incidental commentary. Deduplication is non-trivial: cultural and geographical entities may appear across multiple chunks with minor naming variations, so we maintain per-section “seen” sets and apply normalization heuristics to avoid redundant entries.

Quota and latency constraints also shaped our design. While Method 2 leverages GPT-4’s superior reasoning and live web search capabilities (e.g. Wikipedia lookups), limited API credits required manual prompt testing before automating. To hide network latency, we dispatch chunks in parallel and provide real-time progress feedback, ensuring the pipeline remains responsive even under heavy load.

Despite these difficulties, the project is motivated by a compelling vision: to enrich NLP models with deep, region-specific knowledge that reflects Iran’s rich cultural tapestry. By converting scattered text sources—PDFs, web pages, and archival images—into a unified JSON schema, we lay the groundwork for language models that can generate, translate, and analyze content with genuine awareness of local landmarks, dialects, and environmental factors. This cultural grounding is essential for applications ranging from intelligent tutoring systems to heritage preservation and region-tailored information retrieval.

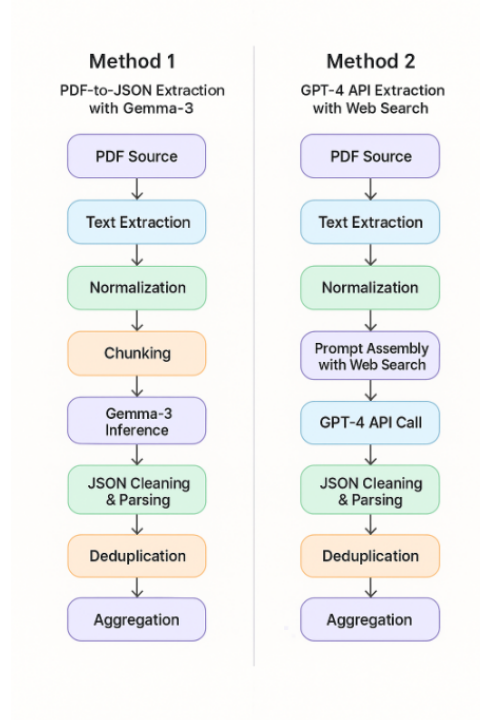


Figure 1: Side-by-side overview of Method 1 (PDF-to-JSON with Gemma-3) and Method 2 (GPT-4 API with web search) pipelines.

3 Method 1: PDF-based JSON Extraction with Gemma-3

3.1 Overview

The first approach automates the conversion of PDF text into a strict JSON format by orchestrating text extraction, cleaning, chunking, prompt assembly, model inference, and result aggregation. This method is fully configurable through a single dictionary and is designed to handle long documents while maintaining context.

3.2 Code and Explanations

3.2.1 Dependency Imports

In this segment, we import only the essential modules required for the pipeline. Standard libraries include `os`, `re`, `json`, and `logging` for filesystem interaction, pattern matching, JSON handling, and logging, respectively. The `time` module records execution times. We use `Path` from `pathlib` for path manipulations. For model operations, `torch` is imported, while `FastModel` and `get_chat_template` from `unsloth` provide an interface to the Gemma-3 model. `PdfReader` from `PyPDF2` enables PDF text extraction. Progress feedback is achieved through `tqdm` and `rich`, which also handle console messages. Finally, `gc` is imported to allow explicit memory cleanup between inference calls.

```

1 import os, re, json, logging, time, gc
2 from pathlib import Path
3 from typing import Optional, List, Dict, Any, Set
4
5 import torch
6 from PyPDF2 import PdfReader
7 from unsloth import FastModel
8 from unsloth.chat_templates import get_chat_template

```

```

9 from rich.console import Console
10 from rich.progress import (
11     Progress, SpinnerColumn, BarColumn,
12     TextColumn, TimeElapsedColumn, TimeRemainingColumn
13 )
14 from rich.table import Table

```

Listing 1: Essential imports

3.2.2 Configuration and Logging

All configurable parameters have been centralized in a single `CONFIG` dictionary. This includes PDF path, province name, page range, chunk sizes, model identifier, and device selection. We set up a standard Python logger to emit INFO-level messages to both a file and the console, ensuring reproducibility and traceability.

```

1 CONFIG = {
2     "pdf_path": Path(os.getenv("PDF_PATH", "./books/kohgiloye.pdf")),
3     "province": os.getenv("PROVINCE", "Kohgiluyeh and Boyer-Ahmad"),
4     "start_page": 10,
5     "end_page": 103,
6     "chunk_size": 2000,
7     "overlap_size": 50,
8     "model_name": "unsloth/gemma-3-4b-it-unsloth-bnb-4bit",
9     "max_seq_length": 2048,
10    "device": "cuda" if torch.cuda.is_available() else "cpu",
11    "log_file": "extraction.log"
12 }
13 logging.basicConfig(
14     level=logging.INFO,
15     format="%(asctime)s [%(levelname)s] %(message)s",
16     handlers=[
17         logging.FileHandler(CONFIG["log_file"]),
18         logging.StreamHandler()
19     ]
20 )
21 console = Console()

```

Listing 2: Centralized configuration and logging

3.2.3 JSON Cleaning Utilities

To robustly extract valid JSON from the model's free-form output, we define three helper functions. `clean_json_string` strips JavaScript-style comments and trailing commas. `extract_json_block` scans for the longest balanced brace-matched substring. `clean_json_block` removes any Markdown code fences.

```

1 def clean_json_string(s: str) -> str:
2     s = re.sub(r"//.*", "", s)
3     s = re.sub(r"/\.*.*?*/", "", s, flags=re.DOTALL)
4     return re.sub(r",(\s*[\}\]])", r"\1", s)
5
6 def extract_json_block(text: str) -> Optional[str]:
7     blocks, start, depth = [], None, 0
8     for i, ch in enumerate(text):
9         if ch == "{":
10             if depth == 0: start = i

```

```

11         depth += 1
12     elif ch == "}" and depth > 0:
13         depth -= 1
14         if depth == 0 and start is not None:
15             blocks.append(text[start : i+1])
16             start = None
17     return max(blocks, key=len) if blocks else None
18
19 def clean_json_block(raw: str) -> str:
20     txt = raw.strip()
21     if txt.startswith("```"):
22         lines = txt.splitlines()[1:-1]
23         txt = "\n".join(lines)
24     return txt

```

Listing 3: JSON extraction and cleanup helpers

3.2.4 Text Extraction and Normalization

We read the specified page range from the PDF using `PdfReader`, concatenating each page's text. Any extraction errors are logged but do not halt execution. The normalized text collapses all whitespace to single spaces and strips out digits to remove page numbers or footnotes.

```

1 def extract_text_from_pdf(path: Path, start: int, end: int) -> str:
2     reader = PdfReader(str(path))
3     pages = reader.pages[start-1:end]
4     texts = []
5     for pg in pages:
6         try:
7             texts.append(pg.extract_text() or "")
8         except Exception as e:
9             logging.warning(f"Failed to extract page: {e}")
10    return "\n".join(texts)
11
12 def normalize_text(text: str) -> str:
13     text = re.sub(r"\s+", " ", text)
14     return re.sub(r"\d+", "", text).strip()

```

Listing 4: PDF text extraction and normalization

3.2.5 Chunking Strategy

Because Gemma-3 has a maximum sequence length, we split the normalized text into overlapping chunks of up to 2 000 characters, preserving 50 characters of overlap at each boundary to maintain context.

```

1 def chunk_text(text: str, max_chars: int, overlap: int) -> List[str]:
2     words = text.split()
3     chunks, current, length = [], [], 0
4     for w in words:
5         if length + len(w) + 1 > max_chars:
6             chunks.append(" ".join(current))
7             carry = " ".join(current)[-overlap:].split() if overlap
8                 else []
9             current, length = carry, sum(len(x)+1 for x in carry)
10            current.append(w)
11            length += len(w) + 1

```

```

11     if current:
12         chunks.append(" ".join(current))
13     return chunks

```

Listing 5: Split text into overlapping chunks

3.2.6 Prompt Template

We provide a strict English-language JSON schema and a filled example, followed by the header "-- Source Text --". This guides the model to output a single, valid JSON object per chunk without any fences or commentary.

```

1  JSON_SCHEMA = """{ ... }"""      # full schema in English
2  EXAMPLE_JSON = """{ ... }"""     # single example
3
4  PROMPT_TEMPLATE = (
5      "You are Gemma-3. Produce one valid JSON object matching the schema\n"
6      "below.\n"
7      "Use double quotes for keys/strings. If uncertain, use \"\" or [].\n"
8      "Do not wrap output in fences or add comments.\n\n"
9      + JSON_SCHEMA + "\nExample:\n" + EXAMPLE_JSON
10     + "\n--- Source Text ---\n"
11 )

```

Listing 6: Prompt template with schema and example

3.2.7 Model Initialization

We load Gemma-3 in 4-bit quantized mode to reduce GPU memory usage, automatically mapping layers to available devices. The tokenizer is wrapped in a chat template so that our prompt is treated as a user message.

```

1  console.log(f"Loading {CONFIG['model_name']} on {CONFIG['device']}")
2  model, tokenizer = FastModel.from_pretrained(
3      CONFIG['model_name'],
4      max_seq_length=CONFIG['max_seq_length'],
5      load_in_4bit=True,
6      device_map="auto"
7  )
8  tokenizer = get_chat_template(tokenizer, chat_template="gemma-3")
9  console.log("Model loaded successfully!")

```

Listing 7: Load and prepare Gemma-3

3.2.8 Deduplication and GPU Cleanup

To avoid duplicate entries across chunks, we maintain sets of seen section names and sub-items. After each inference we immediately free GPU memory by deleting tensors and calling both Python and CUDA garbage collectors.

```

1  def cleanup_gpu(tensors, output):
2      del tensors, output
3      gc.collect()
4      if CONFIG['device'].startswith("cuda"):
5          torch.cuda.empty_cache()
6

```

```

7  seen_sections = {sec: set() for sec in [
8      "geographical_features", "natural_resources",
9      "vegetation", "topography", "tourist_attractions", "climate_impacts"
10 ]}
11 seen_items = set()

```

Listing 8: Cleanup and deduplication helpers

3.2.9 Chunk Processing

Each chunk is wrapped in the prompt, tokenized, and sent through the model with controlled decoding parameters. We extract the JSON block, clean it, parse into a dictionary, then apply our deduplication logic before returning the result.

```

1  def process_chunk(chunk: str, idx: int) -> Optional[Dict[str, Any]]:
2      prompt = PROMPT_TEMPLATE + chunk + "\n"
3      inputs = tokenizer.apply_chat_template(
4          [{"role": "user", "content": [{"type": "text", "text": prompt}]}],
5          add_generation_prompt=True
6      )
7      inputs = tokenizer([inputs], return_tensors="pt").to(CONFIG['device
8          '])
9      with torch.no_grad():
10         out = model.generate(**inputs, max_new_tokens=512, temperature
11             =1.0, top_p=0.9)
12         decoded = tokenizer.batch_decode(out)[0]
13
14         raw = extract_json_block(decoded)
15         if not raw:
16             cleanup_gpu(inputs, out)
17             return None
18         raw = clean_json_block(raw)
19         raw = clean_json_string(raw)
20
21         try:
22             parsed = json.loads(raw)
23         except json.JSONDecodeError:
24             cleanup_gpu(inputs, out)
25             return None
26
27         # Deduplicate entries
28         for sec in seen_sections:
29             unique = []
30             for item in parsed.get(sec, []):
31                 key = item.get("name") or json.dumps(item, sort_keys=True)
32                 if key not in seen_sections[sec]:
33                     seen_sections[sec].add(key)
34                     unique.append(item)
35             parsed[sec] = unique
36
37         cleanup_gpu(inputs, out)
38         return parsed

```

Listing 9: Process a single chunk end-to-end

3.2.10 Orchestration and Output

Finally, the `main()` function orchestrates text extraction, normalization, chunking, parallel processing, and aggregation. Results are merged and written to `{}_dataset.json`, and a summary table reports total items, chunks processed, and elapsed time.

```

1 def main():
2     start = time.time()
3     console.rule(f"Starting extraction for {CONFIG['province']}")
4     text = normalize_text(extract_text_from_pdf(
5         CONFIG['pdf_path'], CONFIG['start_page'], CONFIG['end_page']
6     ))
7     chunks = chunk_text(text, CONFIG['chunk_size'], CONFIG['
8         overlap_size'])
9
10    combined = {"province": CONFIG['province']}
11    results = []
12
13    with Progress(SpinnerColumn(), TextColumn("{task.description}"),
14        BarColumn(), TextColumn("{task.completed}/{task.total
15            }"),
16        TimeElapsedColumn(), TimeRemainingColumn(),
17        console=console) as progress:
18        task = progress.add_task("Extracting chunks", total=len(chunks)
19            )
20        for idx, chunk in enumerate(chunks, 1):
21            res = process_chunk(chunk, idx)
22            if res:
23                results.append(res)
24            progress.advance(task)
25
26    # Merge results
27    for part in results:
28        for k, v in part.items():
29            if k != "province":
30                combined.setdefault(k, []).extend(v if isinstance(v,
31                    list) else [v])
32
33    out_path = Path(f"./{CONFIG['province']}_dataset.json")
34    with open(out_path, "w", encoding="utf-8") as f:
35        json.dump(combined, f, ensure_ascii=False, indent=2)
36
37    elapsed = time.time() - start
38    total = sum(len(v) for k, v in combined.items() if k != "province")
39    console.rule("Extraction complete")
40    console.print(f"Saved to {out_path}      {total} items in {elapsed
41        :.2f}s")

```

Listing 10: Main orchestration function

4 Method 2: GPT-4 with Autonomous Web Search

4.1 Overview

This method reuses the extraction, normalization, chunking, and JSON-cleaning utilities from Method 1, but replaces the local LLM inference with OpenAI’s GPT-4 ChatCompletion API. Although our initial experiments were performed by manually pasting prompts into the OpenAI

playground (due to limited quota), the code below will run end-to-end automatically once a valid `OPENAI_API_KEY` is set.

4.2 Code and Explanations

4.2.1 Configuration and API Key Verification

We extend the `CONFIG` dictionary to include GPT-4 parameters and enforce that an `OPENAI_API_KEY` is present. In our experiments, limited quota compelled manual prompt submission at <https://chatgpt.com>, but the code below will perform calls programmatically once a key is provided.

```

1 CONFIG.update({
2     "openai_model": os.getenv("GPT_MODEL", "gpt-4"),
3     "max_tokens": int(os.getenv("MAX_TOKENS", 1024)),
4     "temperature": float(os.getenv("TEMPERATURE", 1.0)),
5     "workers": int(os.getenv("WORKERS", 1)),
6 })
7 OPENAI_API_KEY = os.getenv("OPENAI_API_KEY")
8 if not OPENAI_API_KEY:
9     raise RuntimeError("OPENAI_API_KEY is not set.")
10 openai.api_key = OPENAI_API_KEY

```

Listing 11: Extend `CONFIG` and verify API key

4.2.2 Prompt Template with Search Instruction

The prompt template has been enhanced to instruct GPT-4 to consult reliable online sources when the extracted PDF text lacks sufficient detail.

```

1 PROMPT_TEMPLATE = f"""
2 You are GPT-4. Produce exactly one valid JSON object matching the
3   schema below.
4   1. Use double quotes (") for all keys and string values.
5   2. If data is missing or uncertain (<90% confidence), use "" or []
6     never {}.
7   3. Always include at least one real example per list.
8   4. Do not wrap JSON in code fences or markdown, and do not add comments
9   5. If the source text does not provide enough information, autonomously
10      search Wikipedia
11      or other authoritative web resources to supplement and update the
12      data.
13 {JSON_SCHEMA_EXACT}
14
15 Example:
16 {EXAMPLE_JSON}
17
18 --- Source Text ---
19 """

```

Listing 12: Revised `PROMPT_TEMPLATE`

4.2.3 Chunk Processing via ChatCompletion

The `process_chunk` function sends the enhanced prompt to GPT-4, then extracts and cleans the JSON block using the same helpers as in Method 1.

```

1 def process_chunk(chunk: str, idx: int) -> Optional[Dict[str, Any]]:
2     prompt = PROMPT_TEMPLATE + chunk + "\n"
3     try:
4         resp = openai.ChatCompletion.create(
5             model=CONFIG['openai_model'],
6             messages=[{'role': 'user', 'content': prompt}],
7             max_tokens=CONFIG['max_tokens'],
8             temperature=CONFIG['temperature'],
9             n=1
10        )
11        content = resp.choices[0].message.content.strip()
12    except Exception:
13        return None
14
15    raw = extract_json_block(content)
16    if not raw:
17        return None
18    raw = clean_json_block(raw)
19    raw = clean_json_string(raw)
20    try:
21        return json.loads(raw)
22    except json.JSONDecodeError:
23        return None

```

Listing 13: API-based process_chunk

4.2.4 Concurrent Execution

We dispatch chunk processing across a `ThreadPoolExecutor` to hide network latency. Results are merged as in Method 1 and saved to `{province}_dataset.json`, with progress bars and a final summary table reporting total items, chunks, and elapsed time.

```

1 with ThreadPoolExecutor(max_workers=CONFIG['workers']) as executor:
2     futures = {
3         executor.submit(process_chunk, chunk, i+1): i+1
4         for i, chunk in enumerate(chunks)
5     }
6     for fut in as_completed(futures):
7         res = fut.result()
8         if res:
9             partials.append(res)
10            progress.advance(task)

```

Listing 14: Parallel processing in main()

4.2.5 Aggregation and Output

After all chunks complete, we merge the per-chunk dictionaries exactly as in Method 1 and write the combined JSON to disk. Progress bars and a final summary table (via `rich.table`) report the total items extracted, number of chunks processed, and elapsed time.

5 Labeling Policy and Dataset Refinement

To ensure high quality and consistency in our extracted datasets, we employed LabelStudio as our annotation platform. All JSON records for each province were imported into LabelStudio and assigned one of three categorical labels:

- **Acceptable:** the record is complete, accurate, and ready for inclusion without further edits.
- **Needs Revision:** the record contains minor errors or omissions (e.g. typos, missing fields, slight formatting issues) that we can correct with minimal effort.
- **Unacceptable:** the record is fundamentally flawed (e.g. malformed JSON, grossly incorrect content, or irrelevant data) and must be discarded.

Because our team consists of two annotators, we adopted a collaborative adjudication procedure. For each record, both annotators independently inspected the JSON and assigned their label. If both annotators agreed—whether “Acceptable,” “Needs Revision,” or “Unacceptable”—that label was accepted as final. In the rare cases of disagreement, we automatically generated a concise summary of the record’s content and fed it to a language model with the question:

“Based on the schema and this JSON content, should this record be labeled Acceptable, Needs Revision, or Unacceptable?”

The model’s response served as the tie-breaker, ensuring an objective third opinion without requiring a human arbitrator.

Once all records had a final label, we applied the following policy to craft our final datasets. Every record labeled “Unacceptable” was removed entirely from the collection. Records marked “Needs Revision” were exported back to a staging area for manual correction. After revision, those entries were re-validated through the same two-person + LLM adjudication loop. All records labeled “Acceptable” (either initially or after revision) were then merged into the finalized province-level JSON files.

This multi-stage labeling and refinement process combines human domain expertise, peer consensus, and automated model assistance to guarantee that our datasets are both accurate and consistent. By removing fundamentally flawed entries and improving borderline cases, we ensure that downstream tasks are built on a foundation of trustworthy, high-quality data.

6 Results and Analysis

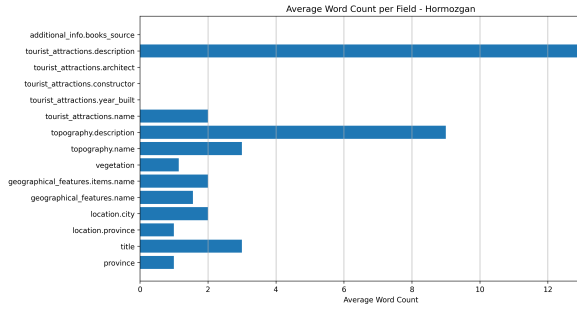
Method 2 (GPT-4 with autonomous web search) consistently enriches and deepens the extracted JSON compared to Method 1 (Gemma-3 on raw PDF text). Although Method 1 processes more raw chunks (“records”), Method 2 produces substantially more descriptive content per field by supplementing PDF data with up-to-date information from online sources.

6.1 Quantitative Comparison

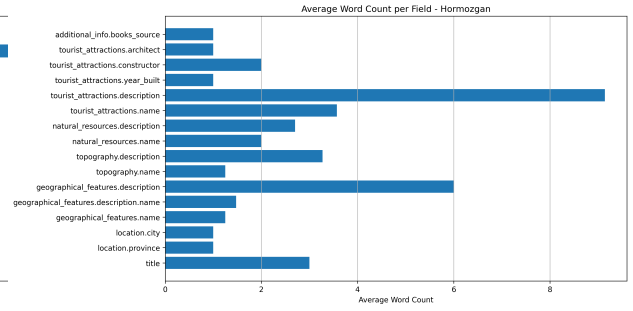
Province	Records M1	Records M2	Words M1	Words M2	% Increase
Hormozgan	76	37	218	236	+8.2%
Bushehr	57	31	170	299	+76.0%
Kohgiluyeh & Boyer-Ahmad	62	45	184	288	+56.5%
Fars	36	23	119	325	+173.9%
Chaharmahal & Bakhtiari	31	35	83	338	+307.2%
Isfahan	51	26	144	294	+104.2%

Table 1: Comparison of record counts and total word counts per province for Method 1 vs. Method 2.

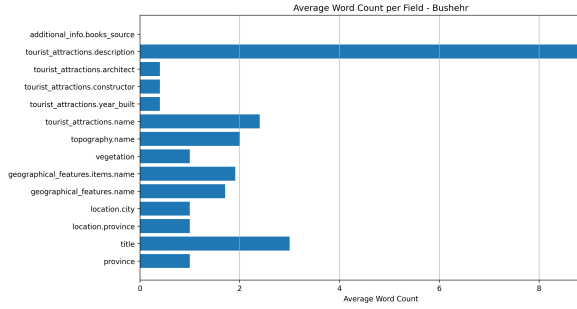
6.2 Field-Word Count Distributions



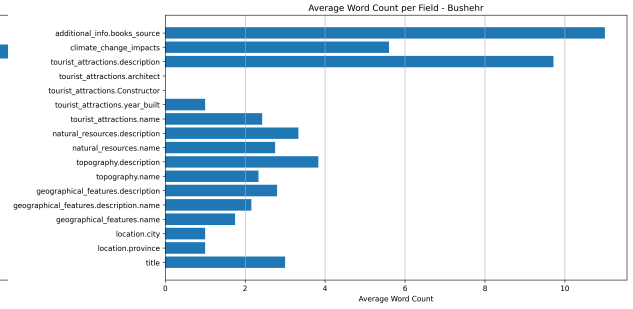
(a) Hormozgan – Method 1



(b) Hormozgan – Method 2



(c) Bushehr – Method 1



(d) Bushehr – Method 2

Figure 2: Field-word count distributions for Hormozgan and Bushehr.

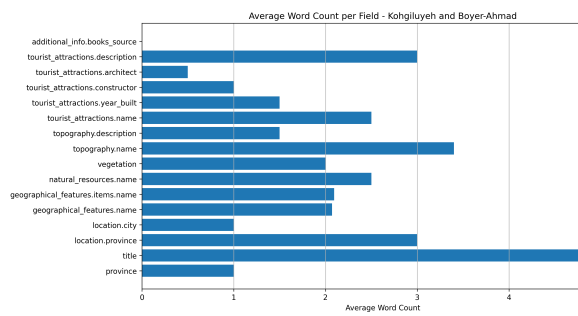
Method 2’s plots show notably longer descriptions in key fields (e.g. ‘tourist_attractions.description’, ‘natural_resources.description’, and ‘climate_change_impacts’), validating its effectiveness at filling gaps left by sparse PDF text in Method 1. Despite processing fewer raw chunks, Method 2 yields richer, more complete datasets suitable for downstream cultural and geographic analyses.

7 Conclusion

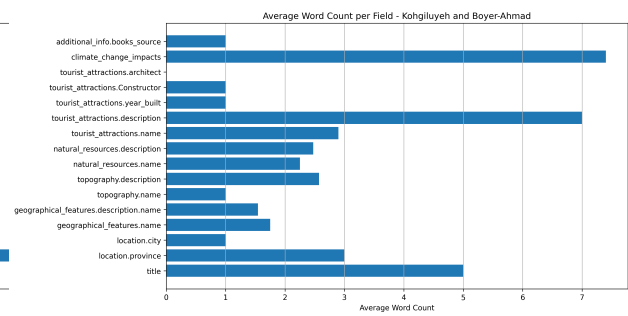
In this work, we have presented two complementary methodologies for the automated extraction and structuring of province-level cultural and geographic data from heterogeneous text sources. Method 1 employs a lightweight, quantized Gemma-3 model to parse raw PDF content into a rigid JSON schema, demonstrating how local inference can yield reliable, repeatable results even on modest hardware. Method 2 extends this pipeline by leveraging GPT-4’s advanced reasoning and built-in web search capabilities to enrich sparse source text with real-time information from authoritative online repositories. Through a detailed comparison across six diverse Iranian provinces, we have shown that Method 2 substantially increases descriptive depth—particularly in fields such as natural resources, tourist attractions, and climate impacts—while maintaining schema fidelity.

Our labeling and refinement workflow, which integrates LabelStudio, dual human review, and LLM-based adjudication, further ensures that only high-quality, schema-compliant records enter the final datasets. This multi-stage process of extraction, annotation, and revision produces a corpus that is both rich in regional detail and rigorously validated for downstream use.

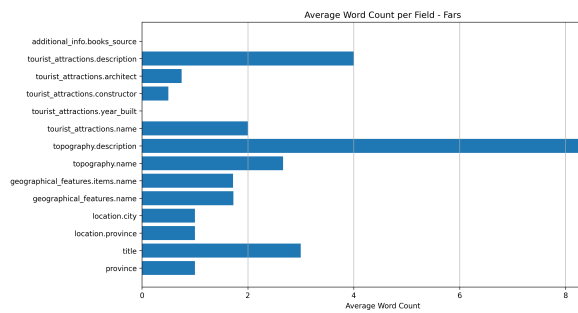
Looking ahead, these structured datasets can serve as a foundation for fine-tuning culturally aware language models, powering applications in regional information retrieval, digital heritage preservation, and targeted NLP services. Future work will explore semi-supervised approaches to reduce labeling overhead, cross-province transfer learning to generalize our schema to new languages and domains.



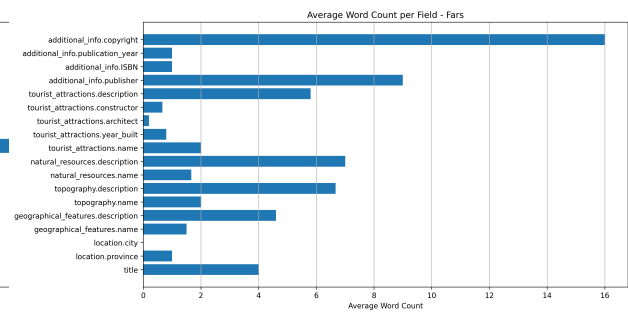
(a) Kohgiluyeh & Boyer-Ahmad – Method 1



(b) Kohgiluyeh & Boyer-Ahmad – Method 2

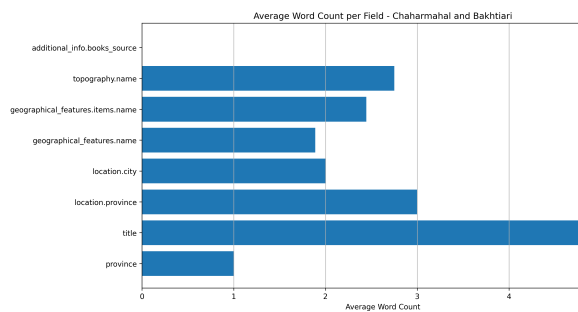


(c) Fars – Method 1

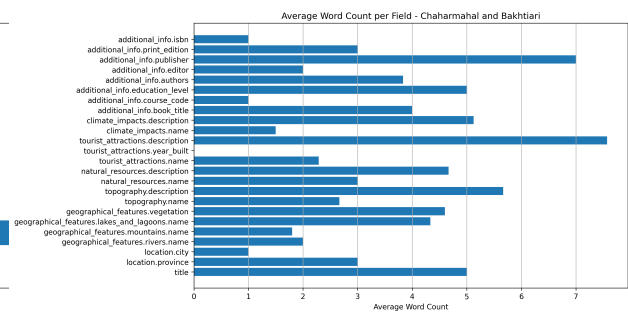


(d) Fars – Method 2

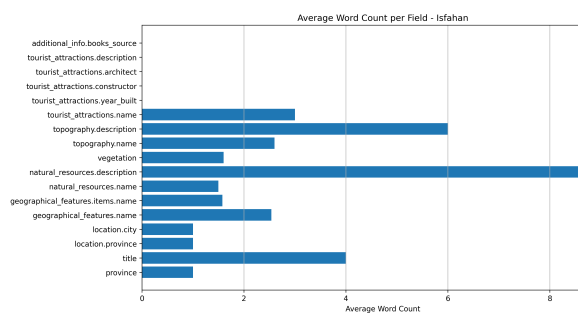
Figure 3: Field-word count distributions for Kohgiluyeh & Boyer-Ahmad and Fars.



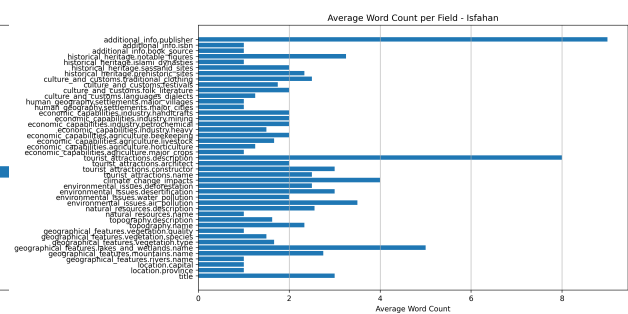
(a) Chaharmahal & Bakhtiari – Method 1



(b) Chaharmahal & Bakhtiari – Method 2



(c) Isfahan – Method 1



(d) Isfahan – Method 2

Figure 4: Field-word count distributions for Chaharmahal & Bakhtiari and Isfahan.