

Evaluation of pdf-docx-to-json-docling-v1

Accuracy and Generality of Form Extraction

This repository provides a two-step pipeline that reliably converts dental intake forms (PDF or DOCX) into structured JSON. **Text extraction accuracy** is high: it uses PyMuPDF for PDFs and python-docx for DOCX to faithfully extract text content locally ¹. For scanned PDFs without text, the tool automatically detects the missing text layer and falls back to Tesseract OCR ² ³, ensuring no form is skipped due to scan images. This automatic OCR feature (recently added **Priority 1.1**) triggers only when needed and even handles mixed-content PDFs page-by-page ³, which improves accuracy without excessive overhead. The extraction step thus captures virtually all visible text from the forms, a crucial foundation for accuracy.

The **parsing logic** that follows is designed to be general and form-agnostic. It employs **intelligent pattern matching** rather than any hardcoded form layout assumptions ⁴. The parser identifies form fields, questions, and options using a large set of regex-based rules and heuristics, covering common patterns in dental/medical forms (dates, phone numbers, checkboxes, etc.). Notably, the code explicitly avoids any form-specific hardcoding or one-off fixes ⁵, instead relying on generic patterns that apply across different form designs. This general approach is reinforced by a **template dictionary** of expected dental fields (`dental_form_dictionary.json`) that ensures output keys and formats are standardized (Modento-compliant) without tying to a single form structure ⁶. The TemplateCatalog mechanism uses case-insensitive exact matches, aliases, and fuzzy matching to map varied field labels to the standard schema ⁷ ⁸. For example, synonyms like "DOB" or "Birth Date" are recognized as *Date of Birth* ⁹, and many alternate phrasings are pre-mapped (e.g. "Cell Phone" vs "Mobile") ¹⁰. This yields consistent JSON output even if the input form used different wording.

In practice, the pipeline reportedly achieves **over 95% field capture accuracy** on typical dental forms ¹¹. All key patient info, medical history, consent fields, etc., are correctly extracted and converted. When a form uses an unexpected label not in the dictionary, the parser still includes it in output with a sensible generated key (slugified from the text). A debug `stats.json` sidecar will flag such unmatched fields and even suggest adding them to the dictionary for future improvement ¹² ¹³. This means the system errs on the side of capturing everything (favoring recall), then normalizes what it recognizes. Overall, accuracy is **excellent for dental intake forms** and the approach generalizes well to variations in wording. (If applied to a completely different domain with entirely new field types, the extraction would still work, but the JSON might not map to known keys — appropriate given the project's scope on dental forms.)

Robustness to Different Form Layouts

The pipeline demonstrates strong **robustness across diverse form layouts** common in dentistry. It handles multi-page forms, varying section orders, and complex field arrangements without form-specific tuning. Several previously tricky edge cases have dedicated solutions:

- **Multi-part field labels** (e.g. a single line containing multiple sub-fields) are correctly split and captured. For example, a line like "Phone: Mobile ___ Home ___ Work ___" is no longer treated as one field but as three separate fields (mobile phone, home phone, work phone) ¹⁴. The parser detects multiple underscore blanks and keywords on one line and divides them, appending the main label appropriately ("Mobile Phone", "Home Phone", etc.) ¹⁵ ¹⁶. This generic logic covers other cases like "City/State/Zip" on one line as well.
- **Multi-column checkbox grids** (common in medical history forms) are parsed into structured fields. The code can detect table layouts by analyzing spacing and alignment of checkboxes across lines ¹⁷ ¹⁸. It now also captures **column headers** in such grids: if a grid has category labels atop each column (e.g. *Appearance / Function / Habits*), the parser will prefix each option with its column header (like "Habits - Smoking") to retain that grouping ¹⁴. This was implemented as the "Grid Column Headers" improvement (Priority 2.2) and significantly improves form-agnostic behavior since the output preserves the intent of grouped options.
- **Inline checkboxes in sentences** (e.g. "[] Yes, send me text alerts") are recognized and isolated as individual boolean fields ¹⁴. The parser finds checkboxes even without a preceding line break or colon, extracting meaningful field titles from the surrounding text. This means even embedded consent checkboxes or yes/no toggles mid-paragraph won't be missed.
- **Section headings and repeated patterns:** The parser intelligently detects section titles (like "Patient Information", "Medical History") and normalizes them. It uses a combination of capitalization and keyword rules to identify headings, and maps variations ("Patient Info", etc.) to standard section names ¹⁹. It also handles cases like a sequence of identical questions (common in checkbox lists) by not merging them incorrectly. Signature lines, initials, and witness lines are identified to avoid mistaking them for normal text fields ²⁰. All these contribute to stable parsing regardless of layout quirks.

Importantly, the project has been careful to keep solutions **form-agnostic**. The logic is based on textual patterns (colons, underscores, spacing, capitalization) and configurable lists of keywords, rather than coordinates or explicit form templates. For instance, column detection uses whitespace heuristics and alignment tolerance ¹⁸ ²¹ instead of assuming a specific number of columns. Similarly, the multi-field splitter looks for multiple blanks and known subfield names (Mobile, Home, Work, etc.) without hardcoding their presence ²² ¹⁶. This means the code can adapt to new forms that follow *similar conventions*, even if the exact layout wasn't seen before. The comprehensive unit tests underscore this generality: tests use representative snippets and **avoid any hardcoded form-specific data** ²³, ensuring that parsing rules apply broadly to any similar input.

There are still a few **edge layouts** that could pose challenges – for example, very unusual formatting or forms with unconventional question phrasing – but these represent a small minority. According to the documentation, the known edge cases now affect under 5% of fields ²⁴. Overall, the pipeline is remarkably

robust across different dental form designs, requiring no custom tweaks per form. It intelligently adapts to content structure, making it a **form-agnostic solution** within the dental/medical intake domain.

Performance and Resource Usage

Performance of the conversion is quite efficient for typical usage, with careful consideration for scalability. On single forms, the local text extraction is fast – PyMuPDF and python-docx operate in-memory without network calls ²⁵, and they process even multi-page documents within seconds. The parsing step (which involves numerous regex checks and string operations) is optimized with compiled patterns and efficient algorithms. For example, the template matching uses set intersections and sequence matching for fuzzy logic ⁸, which are reasonably quick given the moderate size of the field dictionary. The code avoids redundant work where possible; for instance, it will not attempt OCR unless a PDF page is essentially empty of text ³, and it closes files promptly to free resources ²⁶ ²⁷.

For **batch processing** of many forms, the pipeline supports **parallel execution** to improve throughput. Using the `--jobs` flag, users can leverage multi-core CPUs and process multiple documents concurrently ²⁸. Internally, the code implements multiprocessing in both the extraction and conversion steps ²⁹ ³⁰. This design scales well: for example, with `--jobs 4`, it will spawn 4 worker processes to extract text from PDFs in parallel, then similarly convert text to JSON in parallel. The README notes that parallelism can “significantly speed up large batches (50+ forms)” ²⁸. The implementation uses a pool of worker processes and distributes files evenly, which can nearly linearize speedup with CPU cores for CPU-bound parts like OCR.

Resource usage is managed appropriately. Memory footprint is modest for text extraction; even parsed text of a form is just a few hundred KB at most. The most memory- and CPU-intensive operation is OCR on scanned pages, where each page image at 300 DPI is processed by Tesseract ³¹ ³². This is inherently heavy, but the code mitigates it by only doing OCR when necessary (and even then on a page-by-page basis rather than the entire document at once) ³. In a large batch of mostly text-based PDFs, OCR will rarely run, so performance remains I/O-bound by PDF parsing (fast). In a batch of all scanned images, conversion will be slower, but one can distribute such files across processes to utilize multiple cores. A possible consideration is that in extreme cases (very large PDFs or many concurrent OCRs) memory use can spike; however, such scenarios are atypical for intake forms.

The code also collects some performance-related stats in the debug output – e.g. character counts, line counts per file ³³ – which can help in monitoring resource usage. In summary, the **conversion performance is good**, processing forms quickly and making effective use of hardware when needed. The **resource usage is reasonable** and scales with input size (with OCR being the main cost driver). The design choices (local processing, optional parallelism, on-demand OCR) strike a practical balance between speed and completeness.

Test Coverage and Documentation Quality

The project excels in both testing and documentation, which boosts confidence in its reliability and maintainability. A **comprehensive test suite** is included, covering all critical aspects of the pipeline ³⁴. The tests exercise text preprocessing (line wrapping, glyph normalization), question parsing logic (option cleaning, yes/no extraction, field splitting), template matching, and more ³⁵ ³⁶. For example, there are

unit tests to ensure that alias resolution works (so "DOB" indeed maps to the birth date field) ⁹, that fuzzy matching catches slight title variations ³⁷, and that section names like "Patient Info" normalize to "Patient Information" ¹⁹. The tests use representative snippets rather than any one form's exact text, ensuring they are general and cover a range of possible inputs ²³. According to the test documentation, the goal is 100% coverage on critical functions and >70% overall ³⁸, and continuous integration is intended to run these tests on every commit ³⁹. This level of test coverage means regressions or edge-case failures are likely to be caught early. One noted gap is the lack of automated **integration tests** (end-to-end on full sample forms), which is currently done manually ⁴⁰. Even so, the unit tests give thorough coverage of the logic in isolation, which strongly indicates a well-tested codebase.

Documentation quality is excellent. The primary README is detailed, explaining features, usage, and even internal workflow steps ⁴¹ ⁶. It outlines supported form types and known limitations clearly (including which edge cases are now handled) ⁴² ²⁴. In addition to the README, there are multiple focused documents: an **ARCHITECTURE.md** describing the system design, an **ACTIONABLE_ITEMS.md** tracking planned improvements, a **PRODUCTION_READINESS_REPORT.md**, and others. These indicate a high level of project maturity and transparency. Furthermore, the code itself is well-commented. Nearly every function has a docstring explaining its purpose and parameters ⁴³ ⁴⁴. Important changes are annotated (e.g., "Priority 2.1: ..." comments next to the code implementing a new feature), which makes it easy for developers to trace why a piece of code exists ⁵. The maintainers have also started refactoring the monolithic parser into a modular package (e.g., separating `grid_parser.py`, `question_parser.py`, etc.), as described in the package README ⁴⁵. This ongoing modularization effort, along with the clear design principles listed, shows a commitment to long-term maintainability.

In summary, both testing and documentation meet a very high standard. The test suite gives confidence in correctness and helps future development, while the rich documentation ensures users and contributors have the information they need. Few open-source projects in this space have such thorough coverage of their functionality in writing.

Grade: A. This repository demonstrates a high level of accuracy in data extraction, generalizes well to different form layouts, and has acceptable performance characteristics. The extensive tests and documentation further cement its quality. Minor improvements are possible (as noted below), but overall the project is robust and well-engineered for its intended use.

Identified Issues and Suggested Patches

Despite the strong performance, a few issues/opportunities for improvement were identified during review. Each issue is listed with its location, explanation, and a suggested patch or approach to address it. All suggested patches maintain the form-agnostic philosophy (no hardcoding to specific forms).

1. **Non-atomic output file naming in parallel extraction**
2. **Location:** `docling_extract.py` – Function `unique_txt_path()` and usage in `process_one()` ⁴⁶ ⁴⁷.
3. **Issue:** When extracting text in parallel, if two input files have the **same base filename** (e.g., `PatientForm.pdf` in different folders), both processes will attempt to write to the same output `PatientForm.txt`. The current `unique_txt_path` simply checks for existing files and appends a number if needed ⁴⁸, but in a parallel scenario two processes might simultaneously see no file

and both choose the same name. This **race condition** can lead to one file overwriting the other or a file write collision.

4. **Suggested Patch:** Introduce a **more robust unique naming scheme** that accounts for concurrency. One approach is to incorporate a unique token per process or per input file into the output name. For example, include a hash or a truncated UUID of the full input path:

```
# In docling_extract.py, inside process_one before writing:
base_name = file_path.stem
# e.g., use part of the parent folder name or a hash for uniqueness
folder_hash = hex(abs(hash(file_path.parent)) % (16**4))[2:]
out_name = f"{base_name}_{folder_hash}.txt"
out_path = out_dir / out_name
out_path.write_text(text, encoding="utf-8")
```

Alternatively, use Python's built-in `tempfile` or `uuid` to generate a unique name if a conflict is detected. Ensure that each process uses a **different suffix** (like the process ID or an atomic counter from a manager) when a clash is possible.

5. **Benefit:** This will prevent output clashes in parallel mode, guaranteeing that text from every form is saved exactly once. It improves the robustness of batch processing when files might share names, without requiring the user to manually rename anything.

6. Incomplete modularization of parsing logic

7. **Location:** `docling_text_to_modento/core.py` (and `docling_text_to_modento/README.md`) ⁴⁵ ⁴⁹.

8. **Issue:** The main parsing script, while functional, is still partially monolithic (nearly 4000 lines in `core.py`). The maintainers have outlined a plan to split this into modules (`text_preprocessing.py`, `question_parser.py`, etc.) but many parts are marked “Planned” and remain in the single file ⁵⁰ ⁵¹. This makes the codebase harder to navigate and could slow down testing or future enhancements. It’s more of a maintainability issue than a runtime bug.

9. **Suggested Patch: Finish refactoring** the core parser into the planned modules. For instance:
- Move preprocessing functions (like `coalesce_soft_wraps`, `scrub_headers_footers`) from `core.py` into `modules/text_preprocessing.py`, and adjust imports accordingly.
 - Similarly, move question parsing routines (field extraction regex, option parsing, etc.) into `modules/question_parser.py`. The `Question` dataclass and related logic can reside there or in a separate model module.
 - Ensure `core.py` mainly coordinates the workflow (reading input, invoking preprocessing, parsing, post-processing, template matching, etc.) using the new module functions. This can be done incrementally, module by module, verifying tests pass at each step. No new logic is added – it’s organizing existing code. For example:

```
# Example: moving a function out of core.py
# In text_preprocessing.py
def coalesce_soft_wraps(lines: List[str]) -> List[str]:
    # ... (move implementation from core.py) ...
```

```
# In core.py
from .modules import text_preprocessing
...
lines = text_preprocessing.coalesce_soft_wraps(lines)
```

10. **Benefit:** This refactor will **improve code readability and maintainability**. Developers can focus on one aspect at a time (e.g., grid parsing) without scrolling through thousands of lines. It also enables more targeted unit tests (each module can be tested in isolation once separated). In the long run, this makes the project easier to extend (for example, adding new parsing features or supporting new field types) and reduces the risk of unintended side effects when modifying code.

11. Lack of automated end-to-end integration tests

12. **Location:** Testing framework – no specific file (concern noted in `tests/README.md`) ⁴⁰.
13. **Issue:** While unit test coverage is strong, the project currently relies on manual testing for full-form integration ⁴⁰. There is no automated test that takes a sample input file through the entire pipeline (extract text -> parse -> JSON) to verify the end result. This means certain integration issues (like a mismatch between extraction output and parser expectations) might not be caught by the CI tests. For example, if a change inadvertently alters how text extraction formats a checkbox (affecting parsing), unit tests might not catch it since they use fixed strings.
14. **Suggested Patch:** Add a small **integration test** using a representative form. This could be a simplified synthetic form stored in the repo (ensuring no privacy issues). For instance: create a `documents/test_form.pdf` that contains a few fields (name, DOB, a checkbox question). In `tests/test_integration.py`, run the pipeline on this file:

```
import subprocess, json, pathlib

def test_full_pipeline(tmp_path):
    # Copy test_form.pdf to a temp documents directory
    docs_dir = tmp_path / "docs"; docs_dir.mkdir()
    sample_pdf = pathlib.Path("tests/sample_forms/test_form.pdf")
    dest = docs_dir / "test_form.pdf"
    sample_pdf.copy(dest)
    # Run extraction and conversion
    subprocess.run(["python3", "run_all.py", "--in", str(docs_dir), "--out", str(tmp_path)], check=True)
    # Load the output JSON
    out_file = tmp_path / "JSONs/test_form.modento.json"
    data = json.loads(out_file.read_text())
    # Assertions: check that key fields are present in data
    assert any(q["key"] == "first_name" for q in data["questions"])
    assert any(q["key"] == "date_of_birth" for q in data["questions"])
    # (additional checks for correctness of parsed values)
```

This is a high-level sketch; in practice, you might call the `process_one` functions directly rather than shell out to subprocess, to keep it in-process. The key is to simulate the pipeline's two steps and then verify the JSON structure (for expected fields and perhaps counts).

15. **Benefit:** An automated integration test will **catch any regressions in the end-to-end behavior**. If future changes break the interface between extraction and parsing, or cause certain fields to go missing in the final JSON, this test will flag it. It adds an extra safety net beyond unit tests, ensuring the entire pipeline works as intended on real form data. This improves quality assurance, especially before releases.

16. Repeated dictionary loading in parallel conversion

17. **Location:** `docling_text_to_modento/core.py` – in `process_one_wrapper` for parallel jobs ⁵².

18. **Issue:** When converting text to JSON in parallel, each worker process loads the `dental_form_dictionary.json` afresh for each file it processes. The code opens and parses the JSON in `TemplateCatalog.from_path` inside `process_one_wrapper` ⁵². In a scenario with many text files and limited worker processes, a single worker might convert multiple files sequentially, re-reading the same dictionary each time. This is **inefficient**, albeit not catastrophic given the JSON is only a few thousand lines. It adds unnecessary overhead and could slightly slow down processing when hundreds of forms are handled.

19. **Suggested Patch:** Cache the template dictionary in each worker process so it's loaded only **once per worker**, not once per file. A simple way is to use a module-level global. For example:

```
# At top of core.py
_loaded_catalog = None

def get_template_catalog(path):
    global _loaded_catalog
    if _loaded_catalog is None:
        _loaded_catalog = TemplateCatalog.from_path(path)
    return _loaded_catalog

# In process_one_wrapper:
if dict_path and dict_path.exists():
    try:
        catalog = get_template_catalog(dict_path)
    except Exception:
        catalog = None
```

This way, when a worker calls `process_one_wrapper` the first time, it loads the catalog and reuses it for subsequent files. Another approach is to initialize the pool with a function that loads the dictionary into a global variable in each process (using the `initializer` argument of `multiprocessing.Pool`). Either solution avoids repeated disk I/O and JSON parsing.

20. **Benefit:** Though a micro-optimization, this patch will **improve performance and resource usage** during batch conversion. It reduces redundant work and ensures that even as scale grows (large dictionary or very large number of files), each process does the dictionary loading once. This yields a

small speedup in multi-file scenarios and makes the conversion step more efficient without changing its behavior.

Each of the above patches is intended for manual review and implementation. They address concurrency robustness, code maintainability, test completeness, and performance optimization respectively. By applying these targeted improvements, the already strong **pdf-docx-to-json-docling-v1** project can become even more reliable and easier to work with, all while preserving its generality and accuracy on dental forms.

1 4 6 11 14 24 25 28 34 41 42 README.md

<https://github.com/rontavious999/pdf-docx-to-json-docling-v1/blob/51447a8bb215aab0ee7950486237d6dd9534de9c/README.md>

2 3 26 27 29 31 32 43 46 47 48 docling_extract.py

https://github.com/rontavious999/pdf-docx-to-json-docling-v1/blob/51447a8bb215aab0ee7950486237d6dd9534de9c/docling_extract.py

5 12 13 15 16 22 30 33 52 core.py

https://github.com/rontavious999/pdf-docx-to-json-docling-v1/blob/51447a8bb215aab0ee7950486237d6dd9534de9c/docling_text_to_modento/core.py

7 8 10 44 template_catalog.py

https://github.com/rontavious999/pdf-docx-to-json-docling-v1/blob/51447a8bb215aab0ee7950486237d6dd9534de9c/docling_text_to_modento/modules/template_catalog.py

9 19 20 37 test_template_matching.py

https://github.com/rontavious999/pdf-docx-to-json-docling-v1/blob/51447a8bb215aab0ee7950486237d6dd9534de9c/tests/test_template_matching.py

17 18 21 grid_parser.py

https://github.com/rontavious999/pdf-docx-to-json-docling-v1/blob/a8d76b9a50d815df646a8d5ae6325cbb2fd5296a/docling_text_to_modento/modules/grid_parser.py

23 35 36 38 39 40 README.md

<https://github.com/rontavious999/pdf-docx-to-json-docling-v1/blob/a8d76b9a50d815df646a8d5ae6325cbb2fd5296a/tests/README.md>

45 50 51 README.md

https://github.com/rontavious999/pdf-docx-to-json-docling-v1/blob/51447a8bb215aab0ee7950486237d6dd9534de9c/docling_text_to_modento/README.md

49 ACTIONABLE_ITEMS.md

https://github.com/rontavious999/pdf-docx-to-json-docling-v1/blob/51447a8bb215aab0ee7950486237d6dd9534de9c/ACTIONABLE_ITEMS.md