# Solstice: Multi-Agent System for Medical Document Fact-Checking

An AI-Native Approach to Evidence Extraction and Verification

## 1 Introduction

Solstice is a sophisticated multi-agent system designed to automatically fact-check medical claims against scientific literature and clinical documents. By combining state-of-the-art layout detection, multimodal language models, and orchestrated agent pipelines, the system extracts and verifies evidence from complex medical PDFs containing text, tables, and figures.

The system is designed to handle real-world medical documentation at scale, processing complex clinical documents containing text, tables, and figures.

## 2 System Architecture

### 2.1 Document Ingestion Pipeline

The ingestion pipeline transforms unstructured PDFs into queryable structured documents through multiple stages:

1. **Layout Detection**: Uses Detectron2-based models (Faster R-CNN with ResNet-50 backbone) to identify document elements with 95%+ mAP on medical documents

2. **Box Consolidation**: Resolves overlapping detections through IoU-based merging (threshold: 0.7) and hierarchical nesting

3. **Text Extraction**: Employs PyMuPDF for precise text extraction within detected regions, with OCR artifact correction using SymSpell

4. **Figure/Table Extraction**: Saves visual elements as PNG images at 300 DPI for multimodal analysis

5. **Reading Order**: Determines logical flow using column detection and vertical positioning algorithms

### 2.2 Multi-Agent Fact-Checking System

The fact-checking pipeline employs specialized agents orchestrated to work together:

- **Evidence Extraction Agent**: Searches document text for claim-relevant quotes using GPT-4 with temperature=0. Preserves exact quotes while correcting OCR errors and returns structured evidence with relevance explanations.

- **Evidence Verification Agent (V2)**: Validates that extracted quotes exist in the source document. Uses semantic matching to handle OCR variations and filters out tangentially related content, achieving high verification rates.

- **Completeness Checker**: Reviews verified evidence for gaps and searches for additional supporting content. Increases evidence coverage by identifying missing perspectives and feeding new findings back to verification.

- **Image Evidence Analyzer**: Analyzes figures and tables using vision models to identify supporting visual evidence. Processes images in parallel with semaphore control (max 5 concurrent) and provides detailed explanations.

- **Evidence Presenter**: Consolidates all verified text and image evidence into structured JSON and HTML reports. Assesses overall evidence coverage (complete/partial/none) and produces human-readable summaries.

## 3 Technical Implementation

### 3.1 Orchestration Layer

The system uses asynchronous Python with careful orchestration:

- **Claim Orchestrator**: Manages pipeline for single claim across documents

- **Study Orchestrator**: Coordinates multiple claims with configurable parallelism (default: 2 concurrent)

- **Resource Management**: Memory-aware semaphores prevent exhaustion

- **Error Recovery**: Automatic retry with exponential backoff

### 3.2 Model Integration

Different models are selected based on task requirements:

- **GPT-4**: Primary model for text analysis (evidence extraction/verification)

- **O1-mini**: Complex reasoning tasks requiring deeper analysis

- **Claude-3**: Multimodal analysis of tables/figures

- **Dynamic Selection**: Agent-specific model configuration via centralized registry

## 3.3 Robust Error Handling

The system implements multiple layers of reliability:
- **JSON Parsing**: Handles markdown-wrapped and truncated responses
- **Retry Logic**: Structured prompts with error feedback (max 2 retries)
- **Validation**: Pydantic models with custom validators ensure consistency
- **Token Management**: Removed output limits to prevent truncation

# 4 Key Innovations

## 4.1 Multimodal Evidence Integration

Unlike traditional text-only fact-checkers, Solstice analyzes tables and figures—critical for medical evidence where key data often appears in visual formats. The support rate for images demonstrates both the challenge and importance of visual analysis.

## 4.2 Hierarchical Verification

The two-stage verification process (extraction then verification) with an additional completeness check ensures both precision and recall, achieving high verification rates while maintaining comprehensive evidence coverage.

## 4.3 OCR-Aware Processing

The system explicitly handles OCR artifacts common in medical PDFs, using semantic matching rather than exact string comparison, with custom rules for common substitutions (e.g., "0" vs "O").

## 4.4 Production-Ready Architecture

Smart orchestration with configurable parallelism, resource monitoring, and comprehensive caching enables processing at scale while respecting system resources and API limits.

# 5 Real-World Performance

Analysis of the processed corpus reveals:
- **Document Complexity**: Handles multi-page documents with mixed content types
- **Evidence Distribution**: Multiple verified quotes per claim on average
- **Visual Analysis**: Comprehensive image analysis across all documents
- **Processing Efficiency**: 2-3x speedup with parallel claim processing
- **Cache Efficiency**: All intermediate results cached for reproducibility

Current applications include:
- Vaccine efficacy claims against clinical trial reports
- Drug safety statements cross-referenced with FDA documents
- Medical device specifications validated against regulatory filings

# 6 Future Directions

- **Enhanced Table Understanding**: Structured extraction of tabular data with cell-level analysis
- **Cross-Document Reasoning**: Evidence synthesis across multiple sources with conflict resolution
- **Confidence Scoring**: Probabilistic assessment incorporating source reliability
- **Interactive Verification**: Human-in-the-loop validation with active learning

# 7 Conclusion

Solstice demonstrates the power of combining modern AI capabilities—layout understanding, multimodal analysis, and orchestrated agents—to tackle the complex challenge of medical fact-checking. Processing multiple documents with numerous claims has validated the architecture's robustness and scalability for real-world medical documentation.