

# Medical Device Document Retrieval Benchmark Report

Performance Evaluation on IFUs, Recalls, FSNs, and MAUDE Reports

<div>92.4%</div> <div>Table Extraction</div>	<div>5.4×</div> <div>Diagram Retrieval</div>	<div>91%</div> <div>Precision@3</div>	<div>14ms</div> <div>Per Page</div>
----------------------------------------------	----------------------------------------------	---------------------------------------	-------------------------------------

Version	2.0
Date	December 2025
Test Corpus	500+ IFUs, FSNs, Recalls; 5,000 MAUDE Reports
Benchmark Status	Production Ready

# Executive Summary

This report presents benchmark results for Mixpeek's medical device document retrieval system, evaluated on a corpus of 500+ regulatory documents (IFUs, FSNs, recalls) and 5,000 MAUDE adverse event reports. The system demonstrates state-of-the-art performance across three critical capabilities: table extraction, diagram retrieval, and cross-document search.

## Key Findings

- **Table Extraction:** 92.4% cell-level accuracy on FMEA matrices, specification tables, and risk matrices—outperforming GPT-4 Vision (63.1%) and Google Document AI (70.4%).
- **Diagram Retrieval:** 5.4× improvement over text-only search when retrieving assembly drawings, flowcharts, and schematics using CLIP-BGE hybrid approach.
- **Retrieval Precision:** 91% Precision@3 on regulatory queries, compared to 78% for vector-only and 69% for keyword-only baselines.
- **Processing Speed:** 14ms per page average, enabling real-time search across document collections of 100,000+ pages.

## Why This Matters

Medical device regulatory documents present unique challenges that general-purpose AI systems fail to address: complex table structures (FMEA with merged cells), technical diagrams that contain critical specifications, and the need for cross-document correlation (linking recalls to IFUs to MAUDE reports). This benchmark validates that specialized extraction and retrieval significantly outperforms general-purpose alternatives.

## Performance vs. Alternatives

Capability	Mixpeek	GPT-4 Vision	Google Doc AI	Text-Only RAG
Table cell accuracy	92.4%	63.1%	70.4%	N/A
Diagram retrieval	5.4× baseline	1.2×	1.0×	1.0×
Precision@3	91%	—	—	78%
Cross-doc linking	✓	✗	✗	✗
Source attribution	✓	Partial	Partial	✓
Regulatory tuning	✓	✗	✗	✗

Table 1: Performance comparison across key capabilities

# Methodology

## Test Corpus

The benchmark corpus was assembled from publicly available FDA sources and manufacturer documentation to represent the full range of documents encountered in regulatory workflows.

Document Type	Count	Source	Characteristics
Instructions for Use (IFUs)	312	Manufacturer portals	Multi-page, tables, diagrams
Field Safety Notices (FSNs)	89	FDA, manufacturer archives	Structured, cross-references
Recall Notices	156	FDA recall database	HTML/PDF, variable format
MAUDE Reports	5,000	FDA MAUDE bulk download	Narrative text, device codes
510(k) Summaries	47	FDA 510(k) database	Structured sections
Total	5,604		

Table 2: Benchmark corpus composition

## Evaluation Metrics

### Table Extraction Accuracy

Cell-level accuracy measured as the percentage of correctly extracted cells, including headers, data values, and handling of merged cells. Ground truth established by manual annotation of 200 tables across FMEA matrices, specification tables, and risk assessment tables.

### Diagram Retrieval Improvement

Measured as the ratio of Mean Reciprocal Rank (MRR) for our system vs. text-only BM25 baseline. Test set of 150 diagram-seeking queries (e.g., 'assembly diagram for pump cassette', 'flowchart showing sterilization process').

### Retrieval Precision@K

Percentage of queries where at least one relevant document appears in the top K results. Relevance judged by regulatory affairs professionals (2 annotators, Cohen's kappa = 0.81).

## Baseline Systems

System	Configuration	Purpose
BM25	Elasticsearch, default parameters	Keyword baseline
Vector-only	BGE-large-en-v1.5, cosine similarity	Dense retrieval baseline
GPT-4 Vision	gpt-4-vision-preview, structured output	Commercial LLM baseline
Google Document AI	Document OCR + Form Parser	Commercial extraction baseline
Mixpeek	CLIP-BGE hybrid, cross-encoder rerank	Our system

Table 3: Baseline systems used for comparison

# Results: Table Extraction

Table extraction was evaluated on 200 manually annotated tables from IFUs and regulatory documents, including FMEA risk matrices, electrical specification tables, and material composition tables.

## Overall Accuracy

System	Cell Accuracy	95% CI	Header Detection	Merged Cell Handling
Mixpeek	92.4%	[90.1%, 94.7%]	96.8%	88.3%
Google Document AI	70.4%	[66.2%, 74.6%]	82.1%	54.2%
GPT-4 Vision	63.1%	[58.7%, 67.5%]	78.4%	41.7%
Tesseract + heuristics	52.3%	[47.8%, 56.8%]	61.2%	28.9%

Table 4: Table extraction accuracy by system (n=200 tables)

## Performance by Table Type

Table Type	Count	Mixpeek	GPT-4V	Google
FMEA / Risk Matrix	45	89.2%	51.4%	62.8%
Specification Table	62	94.7%	68.3%	74.1%
Material Composition	28	93.1%	72.6%	71.9%
Comparison Table	35	91.8%	65.2%	69.4%
Procedure Checklist	30	95.2%	71.8%	76.2%

Table 5: Extraction accuracy by table type

**Key insight:** The largest performance gap occurs on FMEA tables, where merged cells and multi-level headers cause general-purpose systems to fail. Mixpeek's table-specific extraction achieves 89.2% on FMEA tables vs. 51.4% for GPT-4 Vision—a 74% relative improvement.

# Results: Diagram Retrieval

Diagram retrieval was evaluated on 150 queries seeking visual content: assembly drawings, flowcharts, schematics, and labeled component diagrams. Ground truth established by regulatory professionals.

## Retrieval Performance

System	MRR	Improvement	P@1	P@3	P@5
Mixpeek (CLIP-BGE Hybrid)	0.847	5.4×	78.0%	89.3%	94.0%
Vector-only (BGE)	0.328	2.1×	24.7%	41.3%	52.7%
BM25 (text only)	0.156	1.0×	10.0%	21.3%	29.3%
CLIP-only (no text)	0.412	2.6×	32.0%	49.3%	60.0%

Table 6: Diagram retrieval performance (n=150 queries)

## Hybrid Approach Analysis

The CLIP-BGE hybrid approach combines visual similarity (CLIP embeddings of diagram images) with textual context (BGE embeddings of captions, OCR text, and surrounding content). Fusion weights were learned on a held-out validation set.

Component	Weight	Contribution to MRR
CLIP visual embedding	0.45	+0.381
BGE text embedding (caption + OCR)	0.35	+0.296
BM25 keyword match	0.20	+0.170

Table 7: Learned fusion weights and MRR contribution

**Key insight:** Neither visual nor textual features alone achieve optimal performance. The hybrid approach captures both 'what the diagram looks like' (visual) and 'what the diagram is about' (text), enabling queries like 'assembly diagram for pump cassette' to match diagrams by both visual similarity to other assembly drawings and textual relevance to 'pump cassette'.

# Results: Retrieval Precision

End-to-end retrieval was evaluated on 300 regulatory queries across four categories: keyword lookup, semantic search, cross-modal (diagram/table seeking), and multi-hop (cross-document).

## Overall Precision

System	P@1	P@3	P@5	P@10	MRR
Mixpeek (Hybrid + Rerank)	84%	91%	94%	97%	0.891
Vector-only (BGE)	68%	78%	83%	89%	0.742
BM25 (keyword)	52%	69%	76%	84%	0.614
Naive RAG (chunk + embed)	61%	74%	80%	87%	0.683

Table 8: Retrieval precision across all query types (n=300)

## Performance by Query Type

Query Type	Example	Count	Mixpeek P@3	Vector P@3
Keyword lookup	"K123456 510k clearance"	75	96%	89%
Semantic search	"battery thermal warnings"	100	93%	81%
Cross-modal	"assembly diagram cassette"	75	89%	64%
Multi-hop	"recalls for devices with similar MAUDE events" 50		82%	71%

Table 9: Precision by query type

**Key insight:** The largest improvement is on cross-modal queries (89% vs. 64%), where users seek visual content using natural language. These queries are common in regulatory workflows ('find the diagram showing...') but impossible for text-only systems to handle effectively.

## Processing Speed

Stage	p50 Latency	p95 Latency
Query encoding	8ms	12ms
Sparse retrieval (BM25)	15ms	28ms
Dense retrieval (BGE)	22ms	38ms
Visual retrieval (CLIP)	18ms	32ms
Fusion + Reranking	45ms	78ms
Total	108ms	188ms

Table 10: Query latency breakdown (corpus: 50,000 documents)

# Conclusions

---

## Summary of Results

This benchmark demonstrates that specialized extraction and retrieval significantly outperforms general-purpose alternatives for medical device regulatory documents:

- **92.4% table extraction accuracy**—47% relative improvement over GPT-4 Vision, 31% over Google Document AI
- **5.4x improvement in diagram retrieval**—enabling natural language search for visual content
- **91% Precision@3**—17% absolute improvement over vector-only, 32% over keyword-only baselines
- **Sub-200ms query latency**—real-time search at scale

## Implications for Regulatory Teams

These results translate directly to time savings for RA, QA, and PMS workflows:

- **Audit preparation:** Find all relevant documentation across IFUs, recalls, and risk files with a single query
- **CAPA investigations:** Surface similar incidents and related design controls instantly
- **Adverse event correlation:** Identify patterns across MAUDE reports using semantic similarity
- **FSN development:** Retrieve precedents and related warnings from historical documentation

## Reproducibility

Benchmark methodology, test corpus composition, and evaluation metrics are documented in this report. The corpus includes publicly available FDA documents and can be reconstructed from source links provided. Statistical confidence intervals are reported for all accuracy metrics ( $\alpha = 0.05$ ).

**Ready to see these results on your documents?**

Schedule a demo or request a pilot with your own regulatory corpus.