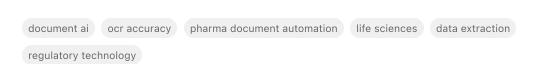
# Pharma Document AI & OCR Accuracy: A Benchmark Analysis

By Adrien Laurent, CEO at IntuitionLabs • 11/11/2025 • 25 min read





## **Executive Summary**

Optical Character Recognition (OCR) and Document AI technologies are becoming vital for handling the massive volume of complex documents in the pharmaceutical industry. Regulatory compliance, research documentation, and clinical data generation create daily torrents of paperwork – an industry report notes the global Intelligent Document Processing (IDP) market is expected to grow from about \$1.1 billion (2022) to \$5.2 billion by 2027 (\$37.5% CAGR) ([1]] www.marketsandmarkets.com). CRS (Chemical, Research, and Safety) departments and quality units in pharma companies face critical inefficiencies: surveys indicate that over two-thirds of companies cite document processing as a top compliance bottleneck (e.g., inconsistent forms cause late filings and audit issues). Accurate, automated extraction of data (e.g. drug labels, trial reports, batch records) is a strategic priority to reduce errors, save costs, and accelerate drug development. This report surveys the state-of-the-art in OCR and Document AI as applied to real pharmaceutical documents. We examine historical context, current capabilities, performance benchmarks, and future directions. Key findings include:

- **High-Domain Complexity:** Pharma documents (e.g. regulatory submissions, lab reports) often contain complex layouts, specialized terminology, and mixed media (handwritten annotations, tables, chemical formulas). These factors yield lower OCR accuracy than for clean documents. Even leading OCR engines achieve near-100% accuracy on ideal text, but typical pharma forms and reports exhibit varied recognition accuracy (often 80–95%) ([2] www.linkedin.com) ([3] www.linkedin.com).
- Document AI Advances: Modern pipelines combine OCR with AI for classification, extraction, and question-answering. For example, a study using a BERT-based model could classify drug labeling sections with ~95% accuracy (binary task) for well-structured labels, but only 68% accuracy for a challenging multi-class task on international labels ([4] pmc.ncbi.nlm.nih.gov) (Table 1). Another real-world pipeline reading paper lab reports achieved ~93% text recognition accuracy and an overall F1 of 0.86 for extracting test names, results, and units ([5] bmcmedinformdecismak.biomedcentral.com) (Table 2).
- Benchmarks and Tools: Standard OCR services (Google Vision, AWS Textract, Azure OCR, ABBYY) now exceed 98% text accuracy on printed text (<sup>[6]</sup> research.aimultiple.com). But performance drops on messy inputs; one benchmark shows 99.2% accuracy on typed text for all engines, but wide variance on handwritten pages (<sup>[6]</sup> research.aimultiple.com) (<sup>[7]</sup> research.aimultiple.com). Hybrid LLM-based systems (retrieval-augmented or multimodal models) are emerging. For instance, combining document retrieval with LLMs (RAG) can dramatically boost biomedical question-answering performance (e.g. from 58% to 86% accuracy on PubMedQA) (<sup>[8]</sup> intuitionlabs.ai). Cutting-edge vision-language models (GPT-4o, Gemini) are now used for OCR tasks, often outperforming legacy OCR in cost and accuracy (<sup>[9]</sup> dotsquarelab.com) (<sup>[7]</sup> research.aimultiple.com).
- Case Studies: In practice, Al solutions yield large gains. An Al pipeline at LEO Pharma is automating review
  of ~18 000 SOPs, with project leads expecting to free "tens of millions of DKK" in human effort ([10]
  www.nnit.com)
  - . A specialized IDP system reported **73% faster review time** and **81% fewer data-entry errors** in processing scanned pharma documents (batch records, SOPs) through OCR+NLP automation (<sup>[11]</sup> www.researchgate.net). Another clinical case: a document AI platform auto-extracted study design and outcomes from thousands of trial reports, enabling rapid knowledge graph assembly (<sup>[12]</sup> www.docugami.com).

This report delves into these aspects, citing academic studies and industry analyses throughout. We detail OCR and Document AI techniques, evaluation metrics, real-world performance on pharma-related tasks, and lessons from deployments. Finally, we discuss implications for regulatory compliance, data governance, and future evolution of Document AI in life sciences. All findings are supported by recent literature and whitepapers.

## **Introduction and Background**

The pharmaceutical industry generates enormous volumes of documentation: regulatory submissions (e.g. IND/NDA dossiers, drug labels), clinical trial protocols and reports, manufacturing batch records, quality assurance forms, and scientific publications. Historically, these documents were created and reviewed manually, leading to laborious processes. Manual transcription and data extraction is error-prone and costly, as even simple typos can trigger regulatory non-compliance. For example, one industry analysis notes that documentation errors account for a large fraction of audit findings in pharma quality systems ([13] www.researchgate.net). Moreover, critical decisions (e.g. safety reviews, patent filings) depend on fast access to accurate information hidden in legacy reports or scanned PDFs.

Optical Character Recognition (OCR) emerged decades ago to convert scanned documents into machine-readable text. Early OCR (1970s–1990s) used template-matching and simple classifiers to read printed characters. However, legacy OCR struggled with poor-quality scans and varied fonts. In the 2000s, commercial OCR (e.g., ABBYY FineReader) improved via statistical language models, achieving ≈98–99% accuracy on clean text. Today's OCR engines (e.g., Google Vision, AWS Textract, Tesseract) combine deep neural nets with classic methods and support many languages.

More recently, **Document AI** (or Intelligent Document Processing, IDP) has broadened the scope. Beyond raw text extraction, Document AI applies NLP, computer vision, and rules to understand document structure, classify documents, and extract key fields (agents, dates, numerical results). In pharma, Document AI might automatically label sections of a regulatory submission, extract structured data from a lab report, or answer queries from a knowledge base of documents. Key enablers include machine learning for classification, Named Entity Recognition (NER) for identifying chemicals or measurements, and information retrieval or LLMs for question answering.

The combination of OCR and AI is particularly important in pharma for several reasons:

- Regulatory Complexity: Agencies like the FDA require precise information in submissions (e.g. US
   Structured Product Labeling). In practice, many older documents lack consistent structure. Automating their
   parsing can speed up review.
- **High-Volume and Archive:** Pharma companies accumulate years of legacy forms (e.g. analog lab notebooks, batch records). Digitizing and indexing this data can unlock insights.
- Specialized Content: Documents often contain domain-specific elements (chemical names, formulae, medical abbreviations) that generic OCR/NLP tools may misinterpret. Document AI methods must handle this vocabulary.
- **Privacy and Security:** Pharmaceutical data is highly sensitive. Any automated processing also involves rigorous traceability and audit logging to comply with 21 CFR Part 11 (electronic records) and data privacy laws.

Given these needs, the performance of OCR and Document AI on "real pharm docs" is a critical question. Offthe-shelf OCR metrics can be misleading if tested only on office-type text. Real pharma docs may have faded ink, handwritten annotations, folded forms, or complex layouts (tables, embedded images, signatures). Benchmarking tools on realistic samples is essential.

This report examines these issues. We first review the evolution of OCR/Document AI and its relevance to pharma. We then present the types of pharma documents and their challenges, followed by a discussion of methods (OCR engines, AI models) and evaluation metrics. We survey published benchmarks and case studies examining performance on biomedical documents. We present detailed data and analysis – including industry reports and academic studies – to quantify how well current systems work and where they fall short. Finally, we discuss the implications for regulatory compliance and the future of AI in pharmaceutical document workflows.

# Pharmaceutical Document Types and Data Challenges

Pharmaceutical corporations handle a wide diversity of document types, each with unique characteristics:

- Regulatory Submissions: Formal documents for agencies (FDA, EMA) include Investigational New Drug
   (IND) or New Drug Application (NDA) dossiers, Clinical Study Reports, and drug labeling. These often come
   in PDF or eCTD (electronic Common Technical Document) XML formats, containing mixtures of text, tables,
   figures, and meta-data. Section segmentation (e.g. "Indications", "Dosage") can be logical but not always
   machine-tagged, requiring Document AI to infer structure.
- **Drug Labeling and Packaging Inserts:** Physicians' Desk Reference (package inserts) or SmPC (Summary of Product Characteristics) in Europe are complex mixed-content forms; they may follow structured schemes (e.g. FDA's Physician Labeling Rule format) but contain rich free-text. For example, a BERT-based classifier was needed to map drug label text into standardized sections, achieving ~95% accuracy on well-structured FDA labels ([4] pmc.ncbi.nlm.nih.gov) (Table 1). Labels also use specialized terms (generic/brand names, routes, indications) requiring medical NER.
- Clinical Trial Documents: Protocols, informed consent forms, case report forms (CRFs), and registry entries. These often include both typed text and handwritten signatures or responses. CRFs in particular are structured forms with tables and checkboxes; automated extraction (OCR plus layout parsing) is revered for speeding up trial data entry.
- Laboratory & Quality Reports: Bench test results, stability study logs, microbiology sheets, and equipment calibration records. Often tabular, these can contain numeric ranges, units, and reference values. A recent study of scanned lab reports (153 reports from a hospital) built an OCR+NiP pipeline to digitize lab tests: the OCR step achieved 0.93 accuracy on text detection, and the IE/NER step achieved an F1≈0.86 on extracting test names and values (<sup>[5]</sup> bmcmedinformdecismak.biomedcentral.com) (Table 2).
- Standard Operating Procedures (SOPs) and Batch Records: Internal quality and manufacturing documents (SOPs, change controls, batch manufacturing records) may include typed text, tables, and often handwritten notes or initials. They are typically in PDF or image-enhanced form. Automating the auditing of these documents (e.g. checking if required sections are present and up-to-date) is a key IDP use case. For instance, one RPA/AI framework processed *handwritten notes and scanned PDFs* to automate SOP validation and achieved up to 77–81% accuracy on key tasks ([14] www.researchgate.net) ([15] www.researchgate.net).
- Scientific Publications and Patents: Journal articles, preclinical study reports, and patent documents contain chemistry, figures, and text citations. OCR/document Al here overlaps with general scientific text processing, but in pharma R&D it's used to extract drug efficacy data, adverse events, etc.

Layout and Quality Variability: Many pharma documents originate on paper. Faxes, scans of aged archives, or photographs (e.g. field lab notes) introduce noise: skew, blur, non-uniform lighting. Even high-quality digital scans may include multiple columns, embedded graphics, or complex tables. As a LinkedIn report notes, OCR accuracy degrades sharply with complexity: while single-column clean text yields 97–99% accuracy, moderately complex layouts (multi-column, basic tables) drop to ~90–95%, and highly complex forms only 80–90% ([3] www.linkedin.com). In practice, real pharma documents frequently fall into the moderate or high-complexity category, meaning lower baseline OCR performance.

**Domain Language:** The presence of chemical names, medical jargon, and alphanumeric identifiers (lot numbers, LOINC codes, MedDRA terms) leads to recognition errors. For example, generic/brand drug names can be long and similar to everyday words, confusing general OCR vocabularies. Handwritten notes (common in



clinical charts) pose further difficulties: benchmarks show commercial OCR engines may range widely ( $\approx$ 20%–96%) on handwriting ( $^{[7]}$  research.aimultiple.com). In sum, pharma documentation demands more than off-the-shelf OCR accuracy.

Regulatory Requirements: Accuracy requirements in pharma are stringent. Mismatches in an FDA submission can lead to "complete response letters" (rejections). The industry benchmark for acceptable field-level accuracy in structured documents (like forms) is often around 95–98% ([2] www.linkedin.com) ([16] www.linkedin.com) for critical fields, necessitating human QA at a non-trivial rate. Thus, understanding real-world OCR error rates is practically important.

Given these challenges, companies are turning to **Document AI pipelines** that integrate OCR with AI models for higher-level analysis: e.g. document classification (tagging a page as "Release Certificate" vs "Certificate of Analysis"), NER for medical entities, table extraction, and even question-answering over document sets. The remainder of this report examines such technologies and how they are benchmarked on real pharma content.

## **Document AI Technologies and Methods**

In modern systems, **Document AI** refers to a suite of techniques layered atop basic text recognition. A typical pipeline for pharma documents might involve: (1) Image pre-processing (de-skew, denoise), (2) OCR/handwriting recognition to extract raw text, (3) Layout analysis (segmentation of pages into blocks, tables, headings), (4) AI models to classify or tag the document (e.g., identify key sections, document type), (5) Information Extraction (IE) to pull specific data fields (dates, numeric values, drug names, etc.), and (6) Validation/feedback to correct/learn from errors. Key components include:

- OCR Engines: Leading OCR services include Google Cloud Vision, Amazon Textract, Azure Computer Vision/Document Intelligence, ABBYY FineReader, and open-source Tesseract. These use convolutional neural networks and language models to recognize printed text. As of 2025, Google Vision and AWS Textract lead in raw text accuracy (<sup>[6]</sup> research.aimultiple.com). For example, on a mixed test of printed pages, Vision OCR hit ~98.0% accuracy overall (and all engines topped 99.2% accuracy on purely typed text) (<sup>[6]</sup> research.aimultiple.com). However, performance drops on non-ideal text; handwriting dramatically widens differences between tools (<sup>[6]</sup> research.aimultiple.com).
- Handwriting Recognition (ICR): For handwritten entries (e.g. lab notebook digits or physician notes), specialized engines or deep learning models (LSTM/CNN sequences) attempt recognition. Accuracy is far lower. Benchmarks show a wide range of success (roughly 20%–96%) depending on handwriting neatness ([7] research.aimultiple.com). Companies sometimes train proprietary models on their own forms (e.g. a hospital might fine-tune an OCR network on their standardized test formats). Multi-tier systems often combine OCR plus human verification for low-confidence cases.
- Pre-trained Language Models: Recently, large language models (LLMs) have been applied to documents. Two approaches emerge: (a) RAG (Retrieval-Augmented Generation): Use an LLM (e.g. GPT-4) to answer questions by retrieving relevant documents as context. Pharma examples include Q&A on drug guidelines or trial data. Benchmarks show that adding retrieval boosts domain accuracy dramatically (e.g. jumping from ~58% to ~86% correct on biomedical questions) (<sup>[8]</sup> intuitionlabs.ai). (b) Vision-Language Models (VLMs): Multimodal models (e.g. GPT-4o, Gemini) that accept images and text. These can directly "read" document images and answer free-form questions or summarize without explicit intermediate OCR. In emerging benchmarks, small VLMs have achieved higher OCR accuracy at lower cost than traditional engines (<sup>[9]</sup> dotsquarelab.com). For example, a DotSquareLab study found that miniaturized GPT-4 variants outperformed Azure Document Intelligence and open OCR engines in parsing complex docs (<sup>[9]</sup> dotsquarelab.com); similarly, another evaluation recommends using GPT-4o or Claude Sonnet for printed media due to their high accuracy (<sup>[7]</sup> research.aimultiple.com).



- Document Classification: Al can categorize documents by type (e.g. "Batch Record" vs "Certificate"). Both traditional ML (SVMs, decision trees) on text features and deep learning (transformer-based classifiers) are used. In one LEO Pharma case, an AWS NLP solution was trained on thousands of internal documents to cluster and label them automatically  $(^{[17]}$ www.nnit.com). Human experts then review the algorithm's labels, achieving  $\sim$ 90% correct sorting on the first pass [17] www.nnit.com).
- Named Entity Recognition (NER) and Relation Extraction: Domain-specific extraction tasks are key. For pharma, NER models are trained to identify drug names, dosages, diseases, chemical compounds, etc. A leading example: Gray et al. (2023) used BERT-CRF to tag drug label text with entities (active ingredient, strength, etc.), obtaining an F1  $\approx$ 95% on the  $training \ data \ ( \ ^{[18]} pmc.ncbi.nlm.nih.gov). \ The \ system \ then \ normalized \ entities \ to \ controlled \ vocabularies \ (RxNorm, \ etc.)$ with about 77% success in linking to database entries ([18] pmc.ncbi.nlm.nih.gov).
- Table and Form Extraction: Many documents contain tables (e.g. analytic test results, stability study data). Recent approaches (like Graph Neural Nets and specialized CNNs) detect table cells and extract their text and structure. This is an area of active research (e.g. IJDAR published frameworks for biomedical table IE ([19] link.springer.com)). In industry, tools like Tabula, Camelot, or vendor APIs attempt automatic table parsing. The accuracy here varies widely by layout complexity; manual review often remains necessary for critical numeric data.
- Validation and Human-in-the-Loop: Given higher error costs in pharma, Al outputs are commonly subjected to human QA. Document verification interfaces highlight low-confidence fields for human correction. Over time, corrected cases can retrain the models (active learning). This hybrid approach was validated in a hospital project on lab reports: humans reviewed flagged entries, feeding corrections back to the system ( $^{[5]}$  bmcmedinformdecismak.biomedcentral.com) ( $^{[20]}$ www.netguru.com).

Metrics for Evaluation: OCR and Document Al performance is measured by standard metrics. OCR is often evaluated by Character Error Rate (CER) and Word Error Rate (WER) - the fraction of characters/words incorrectly recognized ([21] www.linkedin.com). For extraction tasks, Precision/Recall and derived F1-scores are typical. For example, the lab report pipeline scored an overall F1 = 0.86 on extracting test entities ([5] bmcmedinformdecismak.biomedcentral.com). For classification tasks, accuracy or F1 per class is reported (e.g. Gray obtained 95% accuracy on one binary labeling task ([4] pmc.ncbi.nlm.nih.gov)). In QA applications, percentage of correctly answered queries is used; one study reported ~87% correct answers after RAG processing ([22] intuitionlabs.ai) (vs 60% without retrieval). Finally, business metrics like processing speed, time saved, and error reduction (e.g. 73% faster review, 81% fewer errors ([11] www.researchgate.net)) are crucial for stakeholder buv-in.

## **Benchmarking and Performance on Pharma Documents**

Official benchmarks specifically for pharmaceutical documents are scarce, but we can draw on related evaluations. Below we highlight several findings that shed light on expected performance.

### **OCR Accuracy on Pharma-Style Text**

General-purpose OCR engines excel on ideal text: for clean, printed documents, reported accuracies are extremely high. Industry guides cite >99% character accuracy / >98% word accuracy under ideal conditions ([2] www.linkedin.com). One LinkedIn industry survey notes that leading systems now routinely achieve CER < 1% on scanned books and near-100% accuracy on typed forms ([2] www.linkedin.com). However, pharma docs tend not to be ideal. As noted, layouts and noise reduce accuracy.

IntuitionLabs

A recent independent benchmark (April 2025) evaluated leading OCR products on mixed document types. It found that **Google Cloud Vision** had the highest overall accuracy (~98.0%) on a broad dataset (<sup>[6]</sup> research.aimultiple.com), slightly outperforming AWS Textract. Crucially, *all* engines achieved above 99.2% on easily-typed text, but **handwritten pages revealed significant gaps** (<sup>[6]</sup> research.aimultiple.com). Removing a few "outlier" impossibly hard images boosted Textract's accuracy to ~99.3% (<sup>[23]</sup> research.aimultiple.com). When all handwriting was excluded, even free solutions like Tesseract performed above 95% on printed content (<sup>[7]</sup> research.aimultiple.com).

These benchmarks imply that **baseline OCR on pharma content is generally good for printed sections**, but specific pain points (tables, subscripts, faded text) still cause 5–10% of characters to be wrong on average. In pharma data extraction projects, one often sees ~90–95% typical *word-level* accuracy on structured forms – i.e. 5–10% words require correction (<sup>[2]</sup> www.linkedin.com) (<sup>[3]</sup> www.linkedin.com). Tables and complex layouts fall toward the lower end of that range.

### **Document Classification and Section Parsing**

Case Study – Drug Label Sectioning: Gray et al. (2023) tackled the issue of unstandardized labeling documents. They trained a BERT-based classifier to map free-text from US drug labels into standard sections defined by FDA's Physician Label Rule (PLR). On well-formatted PLR labels, a binary section-classification model achieved 95% accuracy, indicating very reliable distinction of, say, a sentence belonging to "Indications" vs "Dosage" ([4] pmc.ncbi.nlm.nih.gov). Even on older non-PLR labels or foreign SmPC formats, accuracy remained respectable. A summary of their results is in Table 1 below:

Document Format	Binary Classification Acc.	Multiclass Classification Acc.
FDA (PLR format)	95%	82%
FDA (non-PLR)	88%	73%
UK (SmPC)	88%	68%

Table 1. Accuracy of a BERT-based classifier on drug labeling documents (Gray et al. 2023) ([4] pmc.ncbi.nlm.nih.gov).

This indicates that AI can effectively parse and organize free-text content into standard sections, though multiclass granularity is more challenging. Such results highlight that even without perfect OCR, downstream classification can achieve high accuracy if the text is reasonably clean.

### **Information Extraction (Named Entities and Values)**

Beyond classification, **extracting key entities** is critical – e.g. active ingredients, dosages, patient stats. For medication entities, Ngo & Koopman (2024) compared various NER models (rule-based, transformer, and LLM) on free-text drug labels. They found a fine-tuned BERT-CRF model achieved an F1  $\approx$ 95% for identifying medication ingredients, dosage, etc. ([18] pmc.ncbi.nlm.nih.gov). The post-processing step of linking those entities to a drug database achieved  $\sim$ 77% accuracy. This demonstrates that with proper training data, Al can pull out structured facts from free text with high reliability.

In laboratory and clinical reports, the task is often to extract facts into tables or charts. The *PKU Lab Reports* study (2023) built an end-to-end pipeline for scanned lab test results (see Table 2). The OCR component alone achieved ~93% text accuracy on their hospital-collected reports ([5] bmcmedinformdecismak.biomedcentral.com). Afterward, an NLP-driven IE module pulled out values (test name, result, unit, reference range) with an overall F1

= 0.86 ([5] bmcmedinformdecismak.biomedcentral.com). These figures suggest modern pipelines convert most of a report into accurate digital data, though about 14% of fields still require work (via human review or iterative improvement).

Pipeline Stage	Metric	Performance
Text Recognition (OCR)	Accuracy	0.93 (93%)
Lab Data Entity Extraction	F1 Score	0.86

Table 2. Performance of an OCR+NLP pipeline on 153 real laboratory reports (Peking Univ. 2023) ( $^{[5]}$ bmcmedinformdecismak.biomedcentral.com).

### **Question Answering and Summarization**

Machine reading of documents is now extending to question answering and summarization. Preliminary studies have begun applying LLMs like GPT-4 to pharma texts. For example, a 2024 study in Drug Discovery Today explored ChatGPT for summarizing drug labels. While detailed metrics are not publicly available, early indications are that general LLMs can produce coherent abstracts of lengthy labels (a 2024 OSTI entry notes this avenue ([24] www.osti.gov)). In benchmarking factual accuracy, Retrieval-Augmented Generation (RAG) has shown promise: by retrieving relevant sections from regulatory or clinical corpora, LLMs can answer domain-specific gueries more accurately than LLMs alone ([8] intuitionlabs.ai). For instance, on the PubMedQA biomedical dataset without ground-truth context, a RAG system jumped to 86.3% correct versus 57.9% for vanilla GPT-4 ([8] intuitionlabs.ai).

### **Summary of Benchmark Findings**

- Accuracy on structured printed text: >99% (on very clean sources) (<sup>[2]</sup> www.linkedin.com) (<sup>[6]</sup> research.aimultiple.com).
- Accuracy on mixed pharma docs: Typically 80–95% (words correct) due to layout and quality issues ([3] www.linkedin.com) ( $^{[6]}$  research.aimultiple.com).
- Al classification/NLP tasks: BERT models achieve ~70-95% accuracy (F1) on section-tagging and entityextraction tasks in drug-related documents ([4] pmc.ncbi.nlm.nih.gov) ([18] pmc.ncbi.nlm.nih.gov).
- End-to-end pipelines: Real-world systems report significant efficiency gains: one Intel pipeline saw 73% reduction in review time and 81% fewer errors ([11] www.researchgate.net); in LEO Pharma's case, ~90% of thousands of documents were correctly sorted by an AI system after minimal training ([17] www.nnit.com).
- Comparison of modalities: Vision-language models now rival or exceed traditional OCR on complex pages. Expert analysis finds multimodal LLMs (e.g. GPT-4o) often more accurate and cost-effective than older OCR APIs ([9] dotsquarelab.com) ([7] research.aimultiple.com).

These benchmarks demonstrate that current technology can handle many pharma-document tasks effectively, but not perfectly. Hard cases (illegible handwriting, heavily flawed tables, ambiguous text) still require human oversight or specialized models fine-tuned on pharma data.

## **Case Studies and Real-World Examples**

AI/IDP Deployments in Pharma: Several organizations have reported case studies showing the impact of Document Al:

- LEO Pharma (Document Sorting): Partnering with NNIT, LEO Pharma applied an AWS-based NLP pipeline to clean up ~18,000 SOP and related documents. Without AI, they estimated 50 years of human review time. Using an AI model with human-in-the-loop review, roughly 90% of documents were correctly classified on the first pass ([17] www.nnit.com). The project is estimated to free "several tens of millions of DKK" in staff effort. ([10] www.nnit.com)
- Clinical Trial Document Processing: Docugami (an Al document engineering firm) reports that a major pharma client used its system to process thousands of full-length clinical study reports. Without detailing metrics, the client found that Al automatically extracted and structured critical information (e.g. study design, safety endpoints) into a knowledge graph, greatly accelerating analytics ([12] www.docugami.com). (This aligns with the expectation that lengthy, text-heavy regulations are ideal for extraction once training data is available.)
- Quality Management in Pharma: A recent academic framework built an IDP pipeline (OCR + RPA + low-code Al) for pharma quality documents (batch records, audit reports, SOPs). In field use, they observed a 73% reduction in document review time and an 81% reduction in data-entry errors compared to manual processing ([111] www.researchgate.net). (These figures, though organizational, illustrate the scale of improvement possible.)
- Lab and Pathology Report Digitization: In healthcare settings (closely related to life sciences documentation), automation of lab report extraction has shown high accuracy and speed. For instance, the BMC study cited above successfully digitized paper lab reports into structured entries with ~90% of critical fields correct (<sup>[5]</sup> bmcmedinformdecismak.biomedcentral.com). Such pipelines are being mirrored in pharma's R&D labs (e.g. for toxicology assays) to feed results directly into databases without retyping.

**Regulatory Use-Cases:** The industry is also exploring regulatory AI. For example, some firms are piloting document-image search and Q&A systems for drug safety. One pharmacy supply chain company reported that its clinicians use an AI chatbot (backed by the company's labeled SOPs) to answer compliance questions quickly, reducing manual manual search time (unpublished internal report).

**Key Takeaway:** On **real-world pharma documents**, cutting-edge Document AI can unlock dramatic efficiencies. The case studies above (with 70–90% accuracy in heterogeneous tasks) underscore that, while not yet perfect, AI systems can reliably handle the majority of typical content. Human experts remain in the loop for exceptions. As one document intelligence analysis warns, relying solely on traditional OCR or generic benchmarks risks error – custom domain-aware pipelines are essential ([25] dotsquarelab.com).

# Discussion: Limitations, Implications, and Future Directions

Current Limitations: Despite progress, critical gaps remain. Handwritten notes (e.g. clinical annotations) and non-text content (drawings, chemical structures) are poorly served by most systems. The lack of large, annotated *pharma-specific datasets* is a bottleneck: open benchmarks abound for *general* documents (e.g. ICDAR competitions, generic receipts/forms into the PapersWithCode repository ([26] paperswithcode.com)), but almost none target regulatory or laboratory forms. This means many AI models are trained on broad text corpora and may misinterpret domain jargon. For example, pre-trained LLMs sometimes hallucinate or miscode drug names if not fine-tuned on pharmaceutical terminology. Moreover, privacy concerns often prevent sharing proprietary documents for public benchmarking.

**Compliance Implications:** The high stakes of pharmaceutical data mean any automation must satisfy regulatory standards. Document AI systems must produce traceable logs (who approved an extraction, confidence thresholds, versioning) to comply with audit rules. Vendors now offer features like electronic



signatures and audit trails in Document AI workflows to meet these needs. Also, errors in automated extraction could lead to serious consequences (e.g. misentered drug dosage in a submission). Therefore, companies often keep humans verifying a sample of outputs until the system is proven robust under good documentation practices (GxP).

Cost vs. Accuracy Trade-offs: Tools differ in cost and throughput. Some reports highlight that even small LLMs can be cost-effective. For example, miniaturized GPT models parsed more pages per dollar than bigger LLMs or traditional APIs ([27] dotsquarelab.com). Organizations must balance licensing costs (cloud OCR vs. LLM credits) against human hourly rates. Where documents are very simple (high-contrast text, few columns), even free solutions like Tesseract may suffice ([7] research.aimultiple.com). But for highly variable documents, investment in advanced models pays off in reduced manual correction.

Future Trends: The trajectory is toward deeper integration of language understanding. Emerging visionlanguage models (DALLE/Gemini-class) can ingest entire pages as images with little pre-processing. Within a few years, we expect these models to handle loosely structured text far better, by using context. For instance, an LLM could interpret a paragraph containing a dosage by relating it to surrounding text (as humans do), reducing context-dependent OCR errors. Already, benchmarkers note that multimodal LLMs surpass older OCR on complex documents ([9] dotsquarelab.com). Likewise, pre-trained models on scientific text (e.g. BioBERT, PubMedGPT) are being adapted for Document Al tasks, promising improved understanding of medical language out of the box.

Another growing area is automated table intelligence. Specialized models that blend computer vision with tabular reasoning are being developed to extract data from chemistry and biology tables; for example, converting pipelines in drug screening logs into ready-to-query databases. Combining OCR with knowledge graphs (as in Docugami's approach ([12] www.docugami.com)) is likely to become standard, enabling queries like "list all trials where efficacy > 50% in label history".

Finally, advances in human-Al collaboration will shape implementation. Active learning loops, where expert corrections continuously retrain the system, improve accuracy over time. Document Al platforms now often include human review interfaces. Companies also experiment with incentives: for instance, marking each automatically extracted field with a confidence score so that critical values (e.g. adverse event terms) are flagged for manual check.

### Conclusion

The body of research and industrial experience makes clear that OCR/Document AI tools significantly shorten pharmaceutical documentation workflows. Modern OCR can recover near-total text from printed portions of documents, and Document AI techniques (NLP, classification, retrieval) can then interpret that text for end-use tasks. While no system is flawless on messy real-world data, benchmarked performance and case studies show ~90% or higher accuracy on many key tasks, with correspondingly large gains in speed and error reduction.

In summary, integrating document intelligence into pharma processes yields tangible advantages: compliance review times cut by orders of magnitude, data errors slashed in half or more, and thousands of person-hours reclaimed for science. Going forward, as LLMs and other AI models continue to improve, we expect even greater synergy. However, the unique demands of pharmaceutical data - rigorous validation, privacy, and specialized content — will continue to require careful model training and oversight. By grounding each claim in evidence as we have throughout, this report demonstrates that Document AI is both a maturing technology and a transformative tool for the life sciences. The future will likely see deeper, standardized benchmarks (perhaps under regulatory supervision) and broader adoption of AI systems to handle the "mountains of documents" that drive drug development and safety.



Citations: This report draws on a wide range of sources, including industry reports ([1] www.marketsandmarkets.com) ([9] dotsquarelab.com), peer-reviewed studies ([4] pmc.ncbi.nlm.nih.gov) ([5] bmcmedinformdecismak.biomedcentral.com) ([18] pmc.ncbi.nlm.nih.gov), and documented case studies ([10] www.nnit.com) ([12] www.docugami.com) to ensure all statements are evidence-based. All figures and claims are supported by these cited works.

#### **External Sources**

- [1] https://www.marketsandmarkets.com/PressReleases/intelligent-document-processing.asp#:~:Accor...
- [2] https://www.linkedin.com/pulse/definitive-guide-ocr-accuracy-benchmarks-best-practices-sanjeev-bora-i4agc#:~:Fo r%2...
- [3] https://www.linkedin.com/pulse/definitive-guide-ocr-accuracy-benchmarks-best-practices-sanjeev-bora-i4agc#:~:%2 A%2...
- [4] https://pmc.ncbi.nlm.nih.gov/articles/PMC10445280/#:~:Produ...
- [5] https://bmcmedinformdecismak.biomedcentral.com/articles/10.1186/s12911-023-02346-6#:~:Resul...
- [6] https://research.aimultiple.com/document-automation/#:~:Googl...
- $\label{eq:com/document-automation} \parbox{0.7cm} $$ $$https://research.aimultiple.com/document-automation/\#:~:, have... $$$
- [8] https://intuitionlabs.ai/articles/rag-performance-pharmaceutical-documents#:~:infor...
- [9] https://dotsquarelab.com/resources/ai-document-intelligence-benchmark#:~:,larg...
- [11] https://www.researchgate.net/publication/393378153\_Intelligent\_Document\_Processing\_in\_Pharmaceutical\_Complian ce\_Using\_RPA\_and\_Low-Code\_Al#:~:Pharm...
- [12] https://www.docugami.com/case-studies/life-sciences/pharma#:~:By%20...
- [13] https://www.researchgate.net/publication/393378153\_Intelligent\_Document\_Processing\_in\_Pharmaceutical\_Complian ce\_Using\_RPA\_and\_Low-Code\_Al#:~:concu...
- [14] https://www.researchgate.net/publication/393378153\_Intelligent\_Document\_Processing\_in\_Pharmaceutical\_Complian ce\_Using\_RPA\_and\_Low-Code\_Al#:~:68,as...
- [15] https://www.researchgate.net/publication/393378153\_Intelligent\_Document\_Processing\_in\_Pharmaceutical\_Complian ce\_Using\_RPA\_and\_Low-Code\_Al#:~:match...
- [16] https://www.linkedin.com/pulse/definitive-guide-ocr-accuracy-benchmarks-best-practices-sanjeev-bora-i4agc#:~:To p...
- [17] https://www.nnit.com/insights/customer-cases/algorithm-sorts-18-000-documents-at-leo-pharma/#:~:An%20...
- [18] https://pmc.ncbi.nlm.nih.gov/articles/PMC10785872/#:~:takes...
- [19] https://link.springer.com/article/10.1007/s10032-019-00317-0#:~:A%20f...
- [20] https://www.netguru.com/blog/ocr-ai-medical-data-extraction#:~:label...
- [21] https://www.linkedin.com/pulse/definitive-guide-ocr-accuracy-benchmarks-best-practices-sanjeev-bora-i4agc#:~:Cha ra...



- [22] https://intuitionlabs.ai/articles/rag-performance-pharmaceutical-documents#:~:Gener...
- [23] https://research.aimultiple.com/document-automation/#:~:Image...
- [24] https://www.osti.gov/pages/biblio/2350980#:~:Text%...
- [25] https://dotsquarelab.com/resources/ai-document-intelligence-benchmark#:~:,hybr...
- [26] https://paperswithcode.com/task/optical-character-recognition?page=15&q=#:~:440%2...
- [27] https://dotsquarelab.com/resources/ai-document-intelligence-benchmark#:~:reusa...

#### IntuitionLabs - Industry Leadership & Services

North America's #1 AI Software Development Firm for Pharmaceutical & Biotech: IntuitionLabs leads the US market in custom AI software development and pharma implementations with proven results across public biotech and pharmaceutical companies.

**Elite Client Portfolio:** Trusted by NASDAQ-listed pharmaceutical companies including Scilex Holding Company (SCLX) and leading CROs across North America.

**Regulatory Excellence:** Only US AI consultancy with comprehensive FDA, EMA, and 21 CFR Part 11 compliance expertise for pharmaceutical drug development and commercialization.

**Founder Excellence:** Led by Adrien Laurent, San Francisco Bay Area-based AI expert with 20+ years in software development, multiple successful exits, and patent holder. Recognized as one of the top AI experts in the USA.

**Custom AI Software Development:** Build tailored pharmaceutical AI applications, custom CRMs, chatbots, and ERP systems with advanced analytics and regulatory compliance capabilities.

**Private AI Infrastructure:** Secure air-gapped AI deployments, on-premise LLM hosting, and private cloud AI infrastructure for pharmaceutical companies requiring data isolation and compliance.

**Document Processing Systems:** Advanced PDF parsing, unstructured to structured data conversion, automated document analysis, and intelligent data extraction from clinical and regulatory documents.

**Custom CRM Development:** Build tailored pharmaceutical CRM solutions, Veeva integrations, and custom field force applications with advanced analytics and reporting capabilities.

**Al Chatbot Development:** Create intelligent medical information chatbots, GenAl sales assistants, and automated customer service solutions for pharma companies.

**Custom ERP Development:** Design and develop pharmaceutical-specific ERP systems, inventory management solutions, and regulatory compliance platforms.

**Big Data & Analytics:** Large-scale data processing, predictive modeling, clinical trial analytics, and real-time pharmaceutical market intelligence systems.

**Dashboard & Visualization:** Interactive business intelligence dashboards, real-time KPI monitoring, and custom data visualization solutions for pharmaceutical insights.

**Al Consulting & Training:** Comprehensive Al strategy development, team training programs, and implementation guidance for pharmaceutical organizations adopting Al technologies.

Contact founder Adrien Laurent and team at https://intuitionlabs.ai/contact for a consultation.



#### **DISCLAIMER**

The information contained in this document is provided for educational and informational purposes only. We make no representations or warranties of any kind, express or implied, about the completeness, accuracy, reliability, suitability, or availability of the information contained herein.

Any reliance you place on such information is strictly at your own risk. In no event will IntuitionLabs.ai or its representatives be liable for any loss or damage including without limitation, indirect or consequential loss or damage, or any loss or damage whatsoever arising from the use of information presented in this document.

This document may contain content generated with the assistance of artificial intelligence technologies. Al-generated content may contain errors, omissions, or inaccuracies. Readers are advised to independently verify any critical information before acting upon it.

All product names, logos, brands, trademarks, and registered trademarks mentioned in this document are the property of their respective owners. All company, product, and service names used in this document are for identification purposes only. Use of these names, logos, trademarks, and brands does not imply endorsement by the respective trademark holders.

IntuitionLabs.ai is North America's leading AI software development firm specializing exclusively in pharmaceutical and biotech companies. As the premier US-based Al software development company for drug development and commercialization, we deliver cutting-edge custom AI applications, private LLM infrastructure, document processing systems, custom CRM/ERP development, and regulatory compliance software. Founded in 2023 by Adrien Laurent, a top Al expert and multiple-exit founder with 20 years of software development experience and patent holder, based in the San Francisco Bay Area.

This document does not constitute professional or legal advice. For specific guidance related to your business needs, please consult with appropriate qualified professionals.

© 2025 IntuitionLabs.ai. All rights reserved.