

Digitizing the Last Mile: OCR in Research Administration

Luke Sheneman

1. The Silent Crisis of Data Entry in Sponsored Programs

The fundamental currency of the scientific enterprise is not merely ideas, but the funding that enables them. In the United States alone, federal aid (primarily via grants) to state, local, and territorial governments totaled approximately \$1.2 trillion in fiscal year 2022, representing roughly [19% of total federal spending](#). Behind this immense flow of capital stands the profession of Research Administration (RA), a discipline tasked with the compliant stewardship of these funds. Yet, despite the scale of investment and the sophistication of the science supported, the administrative infrastructure of grant management remains tethered to a surprisingly archaic workflow: the manual transcription of data from static documents.

For the typical Sponsored Programs Office (SPO) at a major university or research institute, the ["Notice of Award" \(NoA\) is one definitive source of truth](#). Whether issued by the National Institutes of Health (NIH), the National Science Foundation (NSF), or the Department of Energy (DOE), this document transforms a proposal into a [legally binding agreement](#). It delineates the authorized budget, the period of performance, reporting requirements, and [specific terms and conditions](#) such as restrictions on human subjects' research or foreign involvement.

The crisis lies in the format of this truth. Despite legislative efforts like the [Federal Grant and Cooperative Agreement Act of 1977](#) and the [DATA Act of 2014](#)¹, which mandated standardized, machine-readable data, the practical reality for RAs is a deluge of Portable Document Format (PDF) files. A Research Administrator often finds themselves in a *"swivel chair"* scenario: a PDF of the award letter open on one monitor, and the institution's Enterprise Resource Planning (ERP) or Electronic Research Administration (eRA) system open on the other. They must read, interpret, and manually retype dates, dollar amounts, and grant numbers. This manual swivel-chair process is not only inefficient, it is also especially prone to accidental human error, consistent with research showing that manual data transcription tasks commonly produce measurable error rates that can [materially compromise data integrity in knowledge work](#).

This manual bridge is the single greatest bottleneck in the award setup process. It introduces **latency** by preventing investigators from hiring staff or purchasing equipment immediately and creates significant compliance risk. A **typographical error** in a grant number or dollar amount can lead to rejected financial drawdowns; a missed clause buried in a 40-page attachment can result in audit findings or the repayment of disallowed costs. As the complexity of federal regulations grows, with new requirements regarding [data sharing](#) and [research security](#), the

¹ Also: <https://www.gao.gov/assets/gao-15-241t.pdf>

reliance on human eyes to parse thousands of pages of legalese has simply become unrealistic and untenable.

2. The "PDF Trap" and the Failure of Data Standardization

To understand the necessity of advanced Optical Character Recognition (OCR) in this domain, one must first appreciate the failure of purely digital interoperability. The federal government has long recognized the inefficiency of document-based data exchange. The [**Federal Grant and Cooperative Agreement Act of 1977**](#) attempted to distinguish and standardize the legal instruments of assistance. Decades later, the [**Digital Accountability and Transparency Act \(DATA Act\) of 2014**](#) sought to establish government-wide data standards to make federal spending information accessible and machine-readable.

However, the implementation of these mandates has been uneven. While agencies use sophisticated portals like [Grants.gov](#) for *intake* (application submission), the *output* provided to recipients often reverts to **unstructured document formats**. A [GAO report](#) from 2014 highlighted that roughly \$619 billion in assistance awards were not properly reported in machine-readable systems, forcing reliance on the documents themselves.

The [DATA Act](#) matters to university research administrators because it shows that federal grant data remain incomplete, inconsistent, and poorly standardized despite the DATA Act's mandate for common, machine-readable data standards. Although the Act focused on government reporting, ***its core contribution was establishing shared data elements and formats across agencies, which is the same foundation needed for reliable, automated data exchange between funders and recipients***. Without this standards-based, machine-readable exchange, campuses bear added reconciliation burden, higher compliance risk, and greater exposure to data quality errors in ERP and eRA workflows.

Emails and PDFs from different funding agencies fail to achieve the goals of the DATA Act because they are designed for human reading, not standardized, machine-readable data exchange. Each agency's award letter format, layout, and terminology varies, and critical fields (award number, dates, amounts, compliance terms) are embedded in unstructured text or tables that software cannot reliably parse without error-prone OCR and custom logic. This defeats the DATA Act's intent to enable consistent data standards, interoperability, and automated reuse of award information across systems, forcing institutions to rely on brittle extraction workflows and manual verification instead of trusted, structured data flows.

The "PDF Trap" manifests in three primary forms within grant documentation:

- **Born-Digital PDFs with "Spaghetti" Layouts:** These documents are generated from word processors. They contain a text layer, but the logical structure is often lost. A multi-column budget table in an NSF award might copy-paste as a jumbled stream of numbers if the underlying encoding does not strictly define the reading order.
- **Rasterized Scans:** Older awards, or sub-award agreements from less technologically mature partners, are frequently flat images wrapped in a PDF container. To the computer, these are merely [grids of colored pixels](#); they contain no searchable text. This is common in "wet-signed" contracts.

- **Hybrid Documents:** A modern digital cover sheet combined with scanned attachments, such as a negotiated indirect cost rate agreement (NICRA) or a specific budget modification. These can include both vector and raster components, complicating our ability to parse them.

The persistence of this internal PDF variability necessitates a technological intervention. The goal is to move from simple "document management" to "document intelligence," where the system extracts, validates, and integrates the data within the document automatically.

Understanding this transition requires a deep dive into the technology that makes it possible: ***Optical Character Recognition or OCR.***

3. The Evolution of Optical Character Recognition

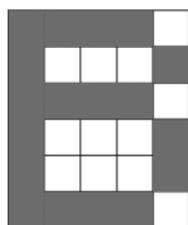
The technology capable of unlocking the data trapped in federal grant PDFs has a rich history, evolving from rudimentary pattern matching to the vision-language models of the artificial intelligence era. Understanding this trajectory is helpful for Research Administrators to evaluate current tools; knowing *why* a legacy system fails on a complex budget table requires understanding the architectural limitations of its era.

3.1. Matrix Matching (1920s–1950s)

The origins of OCR predate modern computing. In the early 20th century, inventors like Emanuel Goldberg and Gustav Tauschek developed "reading machines" utilizing photo-electric cells. Goldberg's "[Statistical Machine](#)" (1931) searched microfilm archives using an optical code recognition system. These early devices relied on ***Matrix Matching*** (also known as template matching).

In this paradigm, the machine compares the image of a character pixel-by-pixel against a stored library of templates (glyphs). If the input image of the letter "B" overlaps sufficiently with the stored template for "B," a match is registered as shown in ***Figure 1***.

Input sample character



Template characters in template database

		1		
	1	1	1	
	1		1	
1	1	1	1	1
1				1
1				1

A

1	1	1	1	
1				1
1	1	1	1	
1				1
1				1
1	1	1	1	

B

	1	1	1	
1				1
1				
1				
1				1
	1	1	1	

C

Digitized character

1	1	1	1	0
1	0	0	0	1
1	1	1	1	0
1	0	0	0	1
1	0	0	0	1
1	1	1	1	0

Digitized sample compared with stored templates

X	X	1	X	
X	X	X	X	X
X	1	X	1	
1	X	X	X	1
1				1
1	X	X	X	X

17 mismatched points

1	1	1	1	
1				1
1	1	1	1	
1				1
1				1
1	1	1	1	

0 mismatched points

X	1	1	1	
1				1
1	X	X	X	
1				X
1				1
X	1	1	1	

6 mismatched points

Figure 1 Template matching of digitized characters. Image source: semanticscholar.org — An Implementation of OCR System Based on Skeleton Matching

Relevance to RA: While obsolete, the concept of matrix matching explains why early digitization efforts failed so spectacularly with varying fonts. If a grant letter was printed in *Courier* but the system was calibrated for *Helvetica*, the pixel overlap would fail. This brittleness made automation impossible for the diverse typography of federal correspondence.

3.2. The Age of Feature Extraction (1960s–1990s)

The introduction of digital computers allowed for a more abstract approach: **Feature Extraction**. Instead of matching pixels, software began to analyze the geometric properties or "features" of a character (**Figure 2**). Ray Kurzweil, a pioneer in this field, developed "[omni-font](#)" OCR in the 1970s.

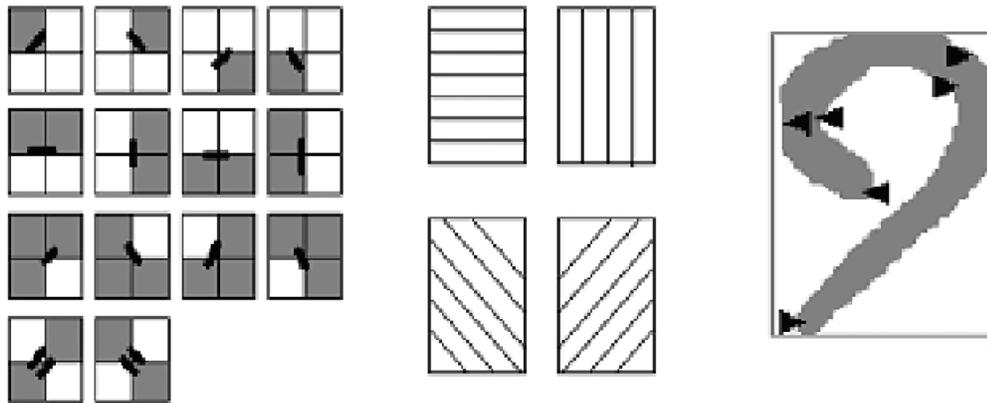


Figure 2: Contour direction and bending features Image source: semanticscholar.org — An Overview of Feature Extraction Techniques in OCR for Indian Scripts Focused on Offline Handwriting

- **Mechanism:** The system identifies a character not by its exact shape, but by its components: lines, loops, and intersections. An "H" is defined as two vertical lines connected by a horizontal mid-stroke.
- **Impact:** This allowed systems to read virtually any printed font, opening the door for commercial applications in law and finance.
- **Limitation:** While these systems could recognize characters, they struggled profoundly with *layout*. They treated a page as a sequence of characters. A multi-column NSF program solicitation would be read line-by-line across the entire width of the page, merging the text of column A with column B into incoherent gibberish. This "reading order" problem remains a plague in legacy OCR tools used in many university archives today.

3.3 The Machine Learning Era and Tesseract (1990s–2015)

The release of [Tesseract](#) as an open-source project by Google in 2005 marked the democratization of OCR. Initially developed by Hewlett-Packard, Tesseract evolved to incorporate machine learning techniques. By version 4, it implemented [Long Short-Term Memory \(LSTM\) networks](#), a type of Deep Learning Network optimized for sequence data.

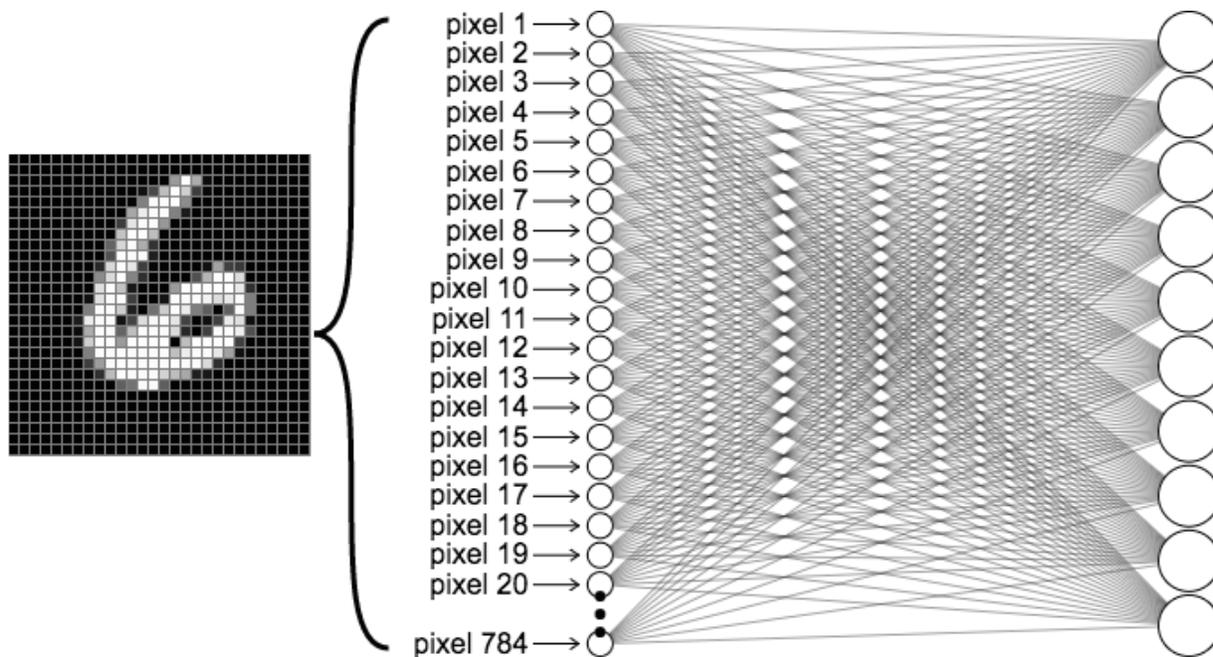


Figure 3 - 1-layer **neural network** for character recognition of numeric digits.

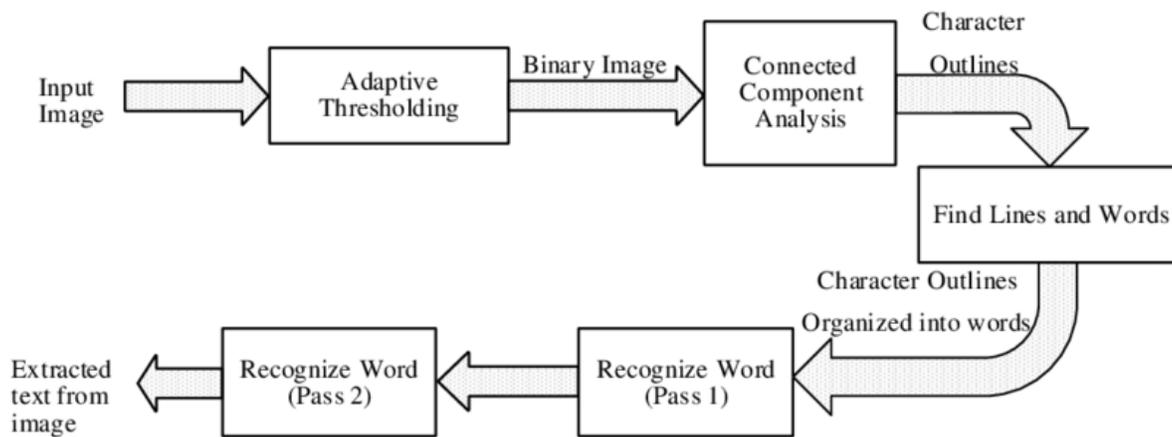


Figure 4 - The multi-stage OCR pipeline used by Tesseract. Image source: <https://arXiv:1811.06193> (2018)

- Mechanism:** LSTMs allowed the OCR engine to consider **context**. It didn't just look at a single character; it looked at the sequence. It learned that "Grant" is a likely word, while "Gr8nt" is not, using linguistic probability to correct character-level ambiguities.

The "Pipeline" Problem: Despite improved character accuracy, these systems operated in a rigid pipeline:

- **Preprocessing:** Binarization (black/white conversion) and de-skewing.
- **Layout Analysis:** Heuristic algorithms (like XY-cuts) to find text blocks.
- **Recognition:** Converting blocks to text.

For Research Administrators, the pipeline is a point of failure. If the layout analysis step fails to detect the faint gridlines of a budget table in a scanned PDF, the subsequent text recognition step will mash the "Salaries" and "Equipment" columns together. The system has *no semantic understanding* that "Year 1" is a header that should govern the data below it.

3.4 Deep Learning Revolution: Vision-Language Models (2018–Present)

The current state-of-the-art has shifted from "*reading text*" to "*seeing documents*". The advent of the Transformer architecture has given rise to **Vision-Language Models (VLMs)**.

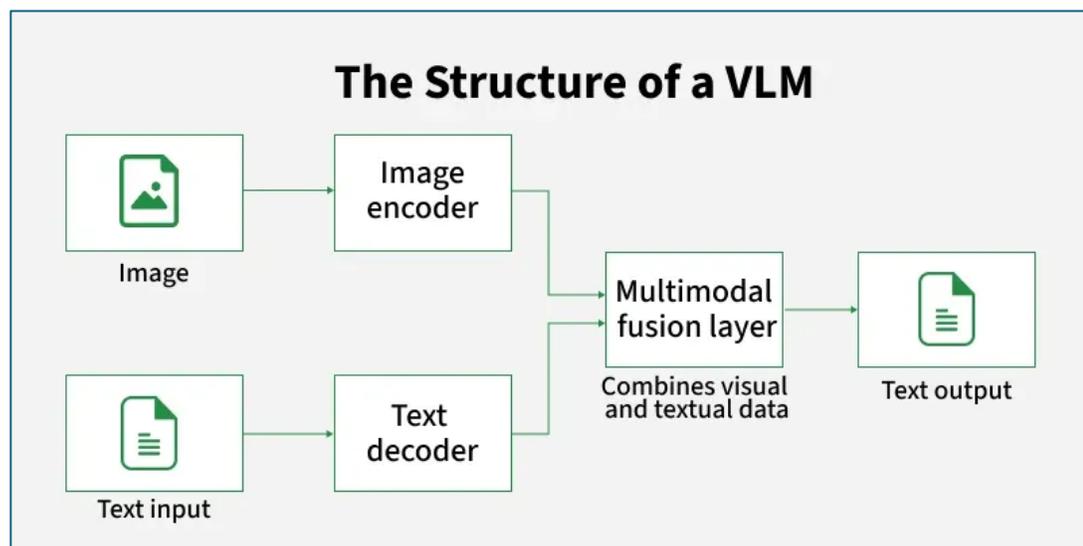


Figure 5 - The architecture of a Vision Language Model. Image source: [GeeksforGeeks](https://www.geeksforgeeks.org/vision-language-models/)

- **Mechanism:** These models process the document image holistically. They do not separate layout detection from text recognition. A model like [dots.OCR](#) or [GPT-5.2](#) takes the image as input and can generate a structured text representation (like Markdown or JSON) as output.
- **Cognitive Shift: *This mimics human perception.*** When a Research Administrator looks at a Notice of Award, they don't draw bounding boxes around words; they instantly recognize "That is the header," "That is the budget table," and "That is the signature." *VLMs effectively replicate this cognitive process, learning to associate the visual structure of a table with the semantic meaning of the data it contains.*

4 The Technical Anatomy of Grant Documentation

To select the right OCR solution, Research Administrators must understand the specific adversaries they face: the technical complexity of the documents generated by federal agencies.

4.1 The Raster vs. Vector Dichotomy: *Not all PDFs are created equal.*

Not all PDFs are equally **machine-readable** and treating them as interchangeable inputs in an OCR or data extraction pipeline is a common and costly mistake. The underlying structure of a PDF determines whether text can be directly extracted, must be inferred from pixels, or falls somewhere in between. Understanding these distinctions is essential for building reliable, automated document processing workflows.

- **Vector PDFs:** These are "born digital," created directly from software like Microsoft Word. ***The text exists as mathematical instructions (draw a character here and there). In theory, data extraction should be easy.*** However, the *encoding* can be messy. *Custom Identity (CID) fonts used by some government systems can result in copy-paste errors where the text appears normal to the eye but copies as random symbols.*
- **Raster PDFs:** These are images (scans) wrapped in a PDF. They rely entirely on OCR. The quality of the scan (Dots Per Inch or DPI) is critical. A faxed sub-award agreement at 96 DPI will suffer from "bleed," where characters merge, causing OCR engines to misinterpret "rn" as "m" or "cl" as "d".
- **Hybrid PDFs:** Many PDFs are hybrid and contain vector and raster components. Extracting text from a hybrid PDF without OCR would lead to an incomplete extraction.



Figure 6 - Vector vs. Raster PDFs. Vector PDFs usually contain easily extractable text, but raster PDFs are just pictures and require OCR. Many PDFs are hybrid and contain vector and raster components.

4.2 The Table Recognition Bottleneck

Tables are the nemesis of traditional OCR. Extracting data from tables using Optical Character Recognition (OCR) is considered one of the most challenging tasks in document processing because it requires not only character recognition but also **structural understanding**: mapping text to specific rows and columns.



Figure 7 - Optical Character Recognition on tables is notoriously difficult due to issues of layout out, sparsity, spanning headers, and more. These things defeat conventional OCR, requiring modern AI-driven OCR tools.

In a **Notice of Award**, the financial data (the most critical data for account setup) is almost always tabular.

- **Borderless Tables:** Modern design often omits gridlines, using whitespace to align columns. Heuristic OCR engines (like Tesseract) often fail to detect these "ghost" columns, collapsing the data into a single string.
- **Spanning Headers:** A header row might say "Approved Budget" centered over three sub-columns: "Direct," "Indirect," and "Total." A pipeline system often fails to associate the spanning header with the sub-columns, leading to data extraction that lacks context.¹⁵

- **Nested Information:** Federal forms often include instructions or footnotes *inside* table cells. Distinguishing between the data (e.g., "\$50,000") and the instruction ("*Subject to 10% cap") requires semantic understanding, not just geometric analysis.
- **Sparse Tables:** Some tables include many empty or conditionally populated cells (e.g., multi-year budgets where only Year 1 is filled). Traditional OCR and rule-based table extractors often misinterpret missing values as layout breaks, causing column shifts, row merges, or spurious data insertion that corrupts the financial structure.

4.3 Multilingual and Special Character Complexity

While US federal grants are English-dominated, global research collaborations introduce multilingual complexity. Sub-awards to foreign institutions may involve documents in French, Spanish, or Chinese. Additionally, scientific grants (NSF, DOE) often contain mathematical formulas. **Traditional OCR engines strip out mathematical symbols or garble them (e.g., turning a summation symbol Σ into an 'E')**. Specialized models are required to preserve the integrity of scientific notation in grant abstracts and technical reports.

5 Generative AI Approaches to Document Understanding

The integration of Generative AI (GenAI) has bifurcated the OCR landscape into "Extractive" vs. "Generative" approaches. This distinction is vital for RAs deciding between a secure, local solution and a cloud-based AI service.

5.1 RAG (Retrieval-Augmented Generation) vs. Native Vision

Before diving into modern "*Chat with your PDF*" tools, it is important to distinguish between systems that treat documents as plain text and those that reason directly over visual structure. Most production systems today rely on Retrieval-Augmented Generation (RAG), where OCR output becomes the sole interface between the document and the LLM. This architectural choice creates a hard dependency on OCR quality and layout fidelity, especially for tables and forms.

- A traditional OCR tool (like Tesseract) scrapes the text from the PDF.
- The text is chopped into "chunks" and stored in a database.
- When a user asks, "What is the total award amount?", the system finds the relevant text chunk and sends it to an LLM (like GPT-5.2 or Llama 3.1) to answer.
- **The Fatal Flaw:** If the initial OCR step fails to capture the table structure (reading across rows instead of down columns), the text chunk becomes nonsensical ("**Salary \$50,000 Equipment \$10,000**" becomes "**Salary Equipment \$50,000 \$10,000**"). The LLM, no matter how smart, cannot repair this fundamental loss of structure. "Garbage in, garbage out" applies strictly here (**Figure 8**).



Figure 8 - Garbage In, Garbage Out: You must have impeccable OCR of your source documents with RAG, including handling complex tables, challenging layouts, mixed fonts, and more. Otherwise, things become a circus at the AI layer.

5.2 Native Multimodal Vision: The "End-to-End" Paradigm

A newer, more powerful approach uses native multimodal vision models (e.g., Gemini 3, GPT-5, dots.OCR) that treat the document as an **image first**, not just as extracted text. Instead of running a fragile OCR pipeline and then asking an LLM to “fix” the results, these models reason directly over the visual layout of the page. This allows them to jointly understand content and structure in one pass, which is especially important for tables, forms, and other layout-heavy documents where meaning is encoded spatially.

The Vision Workflow:

- The model renders the PDF pages as high-resolution images.
- The model's vision encoder analyzes the image directly. It "sees" the table grid, the bold headers, and the spatial relationships.
- The model generates the answer or structured output directly from the visual input. This means that the model can **interpret** that a number is a "Total" because it is **visually**

located at the bottom right of a grid, even if there are no explicit gridlines. This approach is termed "**End-to-End**" (**E2E**) parsing because it unifies detection and recognition into a single inference pass.

5.3 The Hallucination Risk

While Native Vision models are superior for layout understanding, they introduce the risk of **hallucination**. Because they are generative (predicting the next token), there is a non-zero probability that the model might "correct" a messy number into something it thinks looks plausible but is factually wrong.

In Research Administration, where a single digit error in a budget is unacceptable, this necessitates a "**Human-in-the-Loop**" workflow, where the AI's output is treated as a *proposal* to be validated, not a *fact* to be blindly accepted.

6 Review of OCR Solutions for Research Administrators

We will now evaluate ten prominent solutions against the specific needs of Research Administration: accuracy on NoAs, table handling, cost, and data privacy. We categorize them into **Legacy/Open Source**, **Commercial Cloud**, and **Vision-Language Models**.

6.1 Legacy and Open-Source Solutions

These tools are free or low-cost and run locally, ensuring data privacy, but often lack the sophistication to handle complex federal forms without extensive custom coding.



A. Tesseract (Google)

- **Overview:** The grandfather of open-source OCR. Originally developed by HP in the 1980s, it is now maintained by Google. It uses an LSTM-based neural network for character recognition.
- **Strengths:** Free, fast, purely local (runs on-device), and ubiquitous. It supports over 100 languages.
- **Weaknesses: Layout Unaware.** Tesseract processes text line-by-line. It *destroys* the structure of multi-column documents or tables unless heavily pre-processed. It always outputs a "text soup" rather than a structured document.
- **RA Verdict:** Obsolete for primary document parsing in 2026. It serves best as a fallback engine or for simple, single-column text extraction.



B. PaddleOCR (Baidu)

- **Overview:** A rich, industrial-grade OCR toolkit developed by Baidu. It is highly optimized for speed and supports a vast array of languages.
- **Strengths:** Extremely fast and lightweight. It includes specialized tools for table extraction (PP-Structure). It performs better than Tesseract on "in-the-wild" images (e.g., photos of receipts).
- **Weaknesses:** It remains a **pipeline** system (Detection → Recognition). Errors in the detection phase propagate. While efficient, it lacks the semantic reasoning of a Large Language Model to interpret ambiguous grant terms.
- **RA Verdict:** A strong contender for high-volume indexing where speed is critical, but it struggles with the logical complexity of dense grant awards compared to newer VLMs.

6.2 Commercial Cloud Giants

These solutions offer high ease of use and powerful pre-trained models but come with recurring costs and **significant data security concerns** (sending unredacted grant data to the cloud).



C. Amazon Textract (AWS)

- **Overview:** A fully managed machine learning service from AWS. It features a specialized "AnalyzeExpense" API designed for invoices and receipts.
- **Strengths:** Deep integration with the AWS ecosystem (S3, Lambda). It offers a "Query" system where users can ask questions like "What is the Period of Performance?" directly to the document.
- **Weaknesses:** **Table accuracy** is often cited as lower (~78-82%) compared to newer models, particularly on complex, non-standard layouts found in government awards. It can be expensive (\$1.50 per 1,000 pages) for large archival projects.²⁴
- **RA Verdict:** Good for institutions already heavily invested in AWS infrastructure, but arguably "last generation" capability for complex table reconstruction.



D. Google Document AI

- **Overview:** Google Cloud's enterprise solution. It uses "processors" specialized for different document types (e.g., Invoice Processor, Form Parser).
- **Strengths:** Leverages Google's massive training datasets. It creates a "knowledge graph" of the document, understanding relationships between entities.
- **Weaknesses:** It functions as a "black box": it is difficult to tune if it fails on a specific NSF layout. Pricing is complex and metered. Data privacy is a concern for grants involving Controlled Unclassified Information (CUI).
- **RA Verdict:** Powerful, but potentially overkill and too opaque for the specific, bespoke needs of postaward parsing.



E. Mistral OCR (Mistral AI)

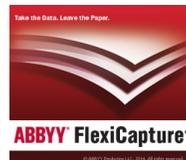
- **Overview:** Mistral OCR is an advanced Optical Character Recognition API built to convert documents (PDFs, scanned pages, images) into structured, machine-readable outputs. It goes beyond simple text extraction, preserving hierarchy, tables, and context so that documents are easier to ingest into downstream systems or analytics workflows.
- **Strengths:** Strong document understanding with high accuracy across text, tables, and mixed content. Benchmarks show competitive performance on multilingual and complex layouts, and it can produce structured outputs like Markdown or JSON directly from input documents. It also supports high throughput (thousands of pages per minute) and offers self-hosted options for institutions with strict data privacy requirements.
- **Weaknesses:** Like other generative OCR services, it may still struggle with deeply nested or borderless financial tables compared to highly specialized table-parsing systems. Hallucination of numeric values remains a risk when relied on without human validation. The self-hosted option is available with a special commercial license rather than as fully open weights anyone can download.
- **RA Verdict:** A powerful, next-generation OCR option for research offices needing to extract and structure text and tables from diverse document types.

Best suited as a primary OCR engine where high-fidelity structured output is needed but still benefits from human-in-the-loop validation for critical numeric fields.



F. Microsoft Azure AI Document Intelligence

- **Overview:** Formerly known as "Form Recognizer." It uses advanced layout analysis to map document structure.
- **Strengths:** Generally rated higher than AWS or Google for **layout retention** and table structure accuracy (87-93%). It integrates seamlessly with Microsoft Power Automate, a tool widely used in university administration workflows.
- **Weaknesses:** It has a steeper learning curve for configuration and custom model training compared to some competitors.
- **RA Verdict:** The strongest of the "Big Three" cloud providers for table-heavy grant documents, especially for institutions on the Microsoft 365 stack.



G. ABBYY FlexiCapture

- **Overview:** The long-standing incumbent in the enterprise OCR market. It combines OCR with rigid, rule-based templates.
- **Strengths:** Robust on-premises deployment options (solving the privacy issue). Highly customizable templates allow for precise data extraction if the document format is constant.
- **Weaknesses:** *rigidity*: It relies heavily on templates. If the NSF changes their award format slightly, the template breaks and requires manual reconfiguration by an IT specialist. Licensing costs are high.
- **RA Verdict:** Reliable but rigid. It represents the "old way" of doing high-volume scanning: effective for standardized forms but brittle for variable correspondence.

H. Rossum

- **Overview:** A "*cognitive data capture*" platform that focuses heavily on financial documents (invoices).
- **Strengths:** Features a unique User Interface (UI) designed for human validation. The system "learns" from user corrections, improving over time. High accuracy on invoices (95%+).
- **Weaknesses: Cost** is significant (**starts ~\$15k/year**). It is primarily tuned for Accounts Payable (AP) workflows, not necessarily the text-heavy, multi-page layouts of Grant Notices.
- **RA Verdict:** Excellent for the post-award financial side (processing sub-award invoices), but potentially ill-suited for the initial award setup and NoA parsing. Too expensive for some institutions.

I. Nanonets Nanonets

- **Overview:** A "no-code" AI platform that allows users to *train* custom models with small datasets!
- **Strengths:** Extremely flexible. An RA office could upload 50 NSF awards and train a specific "NSF Model" without coding. It handles tables well and offers workflow automation.
- **Weaknesses:** Pricing for "Pro" tiers can be high (**\$999/mo/model**). While "low code," it still requires someone in the office to manage the training sets and model lifecycle, which introduces technical and logistical burdens.
- **RA Verdict:** High potential for bespoke needs, but the cost model may be prohibitive for smaller research offices.

6.3 Vision-Language Models (VLMs)

This category represents the *cutting edge*: models that *unify vision and language* to achieve state-of-the-art accuracy while being efficient enough to run locally.



J. DeepSeek-OCR

- **Overview:** A newly released 3-billion parameter model utilizing "Contextual Optical Compression." It focuses on compressing visual information into tokens for efficient processing.
- **Strengths:** Extremely efficient token usage; capable of processing very long documents by compressing images. It outputs structured Markdown directly.
- **Weaknesses:** Benchmarks suggest DeepSeek trades some structural precision for its extreme compression. In direct comparisons on the OlmOCR benchmark, it scores ~75.4, whereas other specialized models score higher.³¹ It also has reported issues with complex table reconstruction compared to models specifically tuned for layout. Also, many state and federal organizations have banned the use of DeepSeek products.
- **RA Verdict:** A *strong contender*, but slightly less precise on the dense, complex layouts typical of federal forms compared to the leader in this category.



K. dots.OCR (*The Preferred Choice*)

- **Overview:** A 1.7-billion parameter unified Vision-Language Model developed by RedNote HiLab. It uses a Qwen2.5-VL base architecture.
- **Mechanism:** It unifies layout detection, text recognition, and table parsing into a single "end-to-end" generation pass. It treats the document processing task as a single conversation: "Here is an image, describe its structure in JSON."
- **Strengths:** State-of-the-art (SOTA) performance on tables and formulas. Its compact size allows it to run on consumer-grade GPUs (e.g., NVIDIA RTX 4090) or even Mac Silicon (with optimizations). It supports prompt-based task switching (e.g., "Extract Table" vs. "Read Text").
- **Weaknesses:** As a Transformer-based model, it is computationally heavier (slower) than pipeline tools like PaddleOCR, though faster than massive general VLMs.
- **RA Verdict:** The "*Goldilocks*" solution: smart enough to understand layout like GPT-5, small enough to run locally like Tesseract, and specifically tuned for the structured chaos of complex documents.

7. Why **dots.OCR** is the Preferred OCR Choice for Research Administrators in 2026

Based on a synthesis of architectural capabilities, benchmark performance, and the specific constraints of Research Administration (data privacy, budget, and document complexity), **dots.OCR** emerges as the superior strategic choice in early 2026.

This conclusion is driven by four key pillars: **Unified Architecture, Table Precision, Efficiency/Privacy Balance, and Adaptability.**

7.1 The Power of Unified Vision-Language Architecture

The fundamental flaw of traditional solutions (like Textract or ABBYY) is the separation of layout analysis and text recognition. If the layout engine fails to see a "ghost line" in a borderless budget table, the text recognition engine reads across the row, merging the "Salaries" column with the "Fringe Benefits" column. This error is fatal for financial data entry.

dots.OCR avoids this by using a single [Transformer model](#) to process the entire page context simultaneously.³⁷

- **Impact on NoAs:** It understands that the text "Year 1" is structurally related to the column of numbers below it, not because it drew a bounding box based on pixel density, but because it *understands the visual semantic relationship*.
- **Result:** It generates structured JSON or Markdown that preserves the logical hierarchy of the Notice of Award. This creates a direct map: `{"Budget": {"Year 1": {"Direct": 50000}}}`. This structural integrity significantly reduces the risk of data entry errors where dates or dollars are assigned to the wrong category.

7.2 State-of-the-Art Table & Layout Parsing

In Research Administration, the hardest information to extract is rarely narrative text. It is structured data embedded in layout. Nearly every high-value field that an SPO or post-award office cares about lives inside a table: line-item budgets, cost share breakdowns, effort commitments, reporting schedules, subrecipient lists, period-of-performance summaries, and compliance certifications. These are not just "text blocks" but spatial data structures where meaning is encoded in rows, columns, headers, and visual grouping.

Traditional OCR pipelines treat tables as an inconvenient special case. They flatten a two-dimensional grid into a one-dimensional stream of tokens, often destroying the very structure that defines what each number means. Borderless tables, spanning headers, nested footnotes, and multi-level column groupings common in federal forms routinely break heuristic table detectors, leading to outputs that look machine-readable but are semantically wrong. Once that structure is lost, downstream LLMs cannot reliably reconstruct it.

State-of-the-art document AI systems therefore live or die by their table and layout parsing performance. High accuracy here is not a “nice to have.” It directly determines whether budget totals are attributed to the correct categories, whether indirect costs are distinguished from direct costs, and whether reporting obligations are captured correctly. In practice, this is the difference between automation that quietly saves staff hours and automation that creates new compliance risk by producing plausible but incorrect structured data.

- **Evidence:** On the [OmniDocBench](#), a rigorous industry benchmark for document parsing, dots.OCR achieves a **Table TEDS (Tree Edit Distance-based Similarity) score of 88.6%**. This score statistically outperforms Google’s Gemini 2.5 Pro (85.8%) and other competitors.
- **Significance:** This means dots.OCR is objectively better at reconstructing the complex, nested, and often borderless tables found in federal grant attachments than one of the world’s most advanced proprietary models. For an RA office, this translates to fewer hours spent manually correcting “broken” spreadsheet exports and higher confidence in the data pushed to the ERP.

7.3 The Efficiency and Privacy Sweet Spot

Research Administrators deal with sensitive data. Grant proposals contain proprietary intellectual property (IP); NoAs contain salary information which is Personally Identifiable Information (PII). Grants from the Department of Defense (DOD) or Department of Energy (DOE) often come with **Controlled Unclassified Information (CUI)** restrictions that strictly regulate data handling.

- **The Cloud Risk:** Sending this data to OpenAI (GPT-5) or Google (DocAI) introduces privacy, security, and potential export control complications. It requires complex legal vetting and contractual agreements with providers.
- **The dots.OCR Advantage:** At only **1.7 billion parameters**, dots.OCR is remarkably compact. It is small enough to run on a single standard workstation.

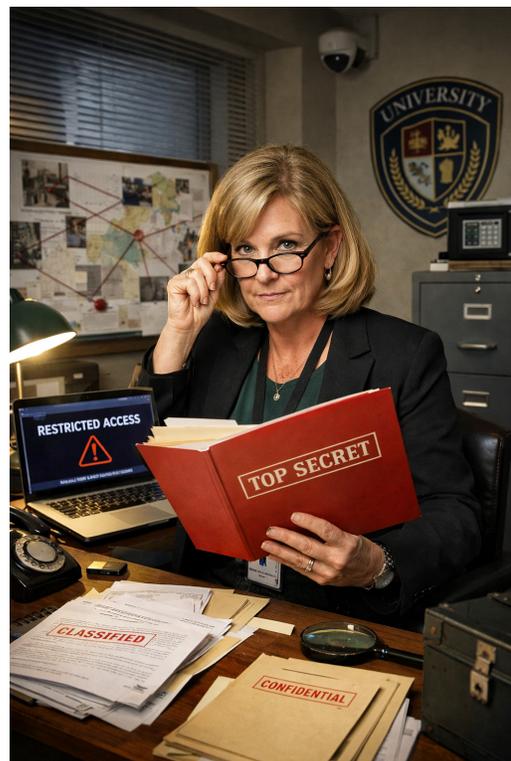


Figure 9 - Research Administrators have to handle sensitive data such as PII, FCI, and even CUI

- **Operational Benefit:** An institution can deploy dots.OCR *on-premise* or within a private Virtual Private Cloud (VPC). The data never leaves the university’s control.

This capability is relevant for compliance with strict federal data security mandates.

7.4 Adaptability via Prompting

Unlike ABBYY, which requires rigid templates that break when a format changes, or Tesseract, which has no control interface, dots.OCR is ***promptable***.

- **Use Case:** An administrator can feed a NoA into the system with specific prompts.
 - **Prompt:** "Extract the layout only" → To analyze the document structure.
 - **Prompt:** "Grounding OCR" → To find the specific coordinates of the "Total Award Amount" for audit trails.
- **Flexibility:** This allows the model to adapt to new funding agency formats without retraining. If the NIH releases a new "Policy Page" format next year, the prompt-based nature of the VLM allows it to adapt to the visual context naturally in a future-proofing way.

7.5 Comparative Summary: dots.OCR vs. The Field

Feature	Legacy (Tesseract/ABBYY)	Cloud Giants (Textract/Google)	General VLMs (GPT5/Gemini)	dots.OCR
Architecture	Pipeline / Template	Proprietary Pipeline	Large VLM	Unified Compact VLM
Table Accuracy	Low / Brittle	Moderate (Costly)	High	SOTA (88.6% TEDS)
Data Privacy	High (Local)	Low (Cloud Egress)	Low (Cloud Egress)	High (Local Deployment)
Cost Model	License / Free	Per Page Metered	Per Token (Expensive)	Open Source / Compute Only
Setup	High (Templates)	Low (API)	Low (API)	Moderate (Self-Host)
Hardware	Low (CPU)	None (SaaS)	None (SaaS)	Moderate (Consumer GPU)

Conclusion: dots.OCR is the current preferred choice because it democratizes "Big Tech" document intelligence. It gives the Research Administration office the table-parsing power of a Gemini or GPT-class frontier model without the exorbitant per-page costs or data privacy compromises, all in a package that fits on local hardware.

8. Implementation Guide for Research Administrators

Implementing a solution like ***dots.OCR*** moves the RA office from a paradigm of ***manual data entry and cut-and-paste*** to one of ***automated data extraction + validation***. This shift requires not just software, but a change in workflow.

8.1 Deployment Architecture & Hardware

To run dots.OCR effectively, an institution does not need a supercomputer.

- **Hardware:** A dedicated workstation or small server is sufficient. The recommended specification is an NVIDIA GPU with at least 16GB of VRAM (e.g., RTX 3090/4090) to handle batch processing of multi-page PDFs. For Mac-based offices, recent community efforts have enabled dots.OCR to run on Apple Silicon (M1/M2/M3) chips, making it accessible even on standard laptops.
- **Integration:** The model can be wrapped in a simple Python API (using frameworks like [Flask](#) or [FastAPI](#)).



Data scientists at [AI4RA at the University of Idaho](#) developed and disseminate [dots_ocr_api](#), a simple open source API wrapper around dots.OCR. Once deployed and configured on your server or workstation at your institution, you can easily submit a PDF and retrieve markdown.

Example Automation Workflow:

- NoA PDF arrives via email from funding agency.
- Script automatically sends PDF to the dots.OCR local server.
- dots.OCR returns structured Markdown.
- An LLM extraction script maps the JSON to ERP fields ("Award Number," "Dates," "Budget").
- Data is pushed to a "Staging Area" in the ERP system, ready for human-in-the loop validation before committing.

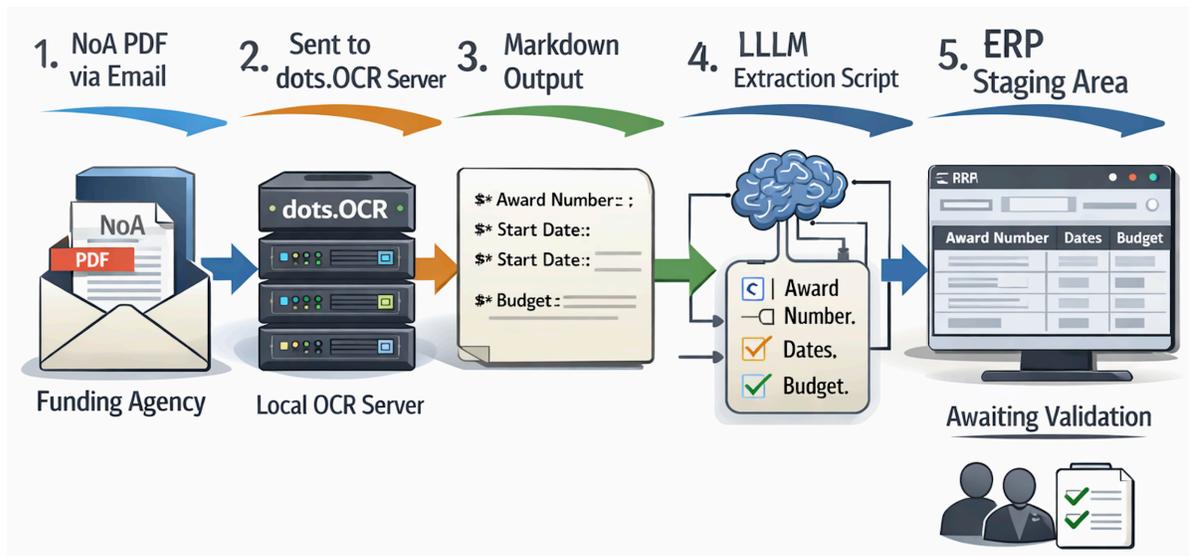


Figure 10. An example RA workflow using dots.OCR to convert agency correspondence into properly ingested institutional data.

8.2 Bridging the DATA Act Gap

While waiting for federal agencies to fully comply with **DATA Act** machine-readable standards, dots.OCR acts as the essential bridge. It essentially "upgrades" legacy PDF NoAs into the machine-readable format that agencies *should* be providing. This empowers institutions to perform analytics on their grant portfolios. As an example instantly querying "Show me all awards received in 2024 that contain a 'Data Sharing' clause," a task that would previously require opening and reading hundreds of files.

Convenience tools like **Vandalizer**, currently under development by the University of Idaho through funding by NSF GRANTED, lower the technical barriers so RAs can perform OCR, extraction, and reporting workflows using only a browser.



Figure 11 - OCR is an essential first step for document loading into user-friendly RA tools like Vandalizer.

8.3 Managing Expectations and Human-in-the-Loop

While dots.OCR is State-of-the-Art, it is not magic.

- **Handwriting:** While capable, it may struggle with scribbled notes in margins compared to specialized handwriting models.
- **Resolution:** It performs best at 200 DPI. Low-quality faxes of sub-award agreements may still require manual review.
- **Hallucination:** As a generative model, there is a risk of "hallucination." However, dots.OCR's grounding mechanisms minimize this. The "Human-in-the-loop" validation step, where an administrator reviews the extracted data before final commit, remains an essential component of the workflow.



Figure 12 - Humans must remain in the loop.

9. Conclusion

The administrative burden of federal research compliance is effectively a tax on scientific discovery. Every hour a research administrator spends retyping a budget table from a PDF is an hour not spent guiding a faculty member through a complex proposal, finishing an award setup, or managing a critical compliance risk. The sheer volume of data, combined with the rigidity of legacy documentation formats, has created a sustainability crisis in Sponsored Programs Offices.

The evolution of Optical Character Recognition, from the pixel-counting machines of the 1920s to the vision-language intelligence of today, offers a way out of this "swivel chair" paradigm. While commercial cloud solutions offer power, they introduce data sovereignty and security risks and recurring costs that strain university budgets.

dots.OCR (and systems like it) represent a watershed moment for this industry. By packaging state-of-the-art visual understanding, table reconstruction, and prompt-based flexibility into a compact, open-source model you can run on your campus, it provides Research Administration offices with a tool that is:

- **Technically Superior** for the specific, table-heavy layouts of grant awards.
- **Fiscally Responsible** by eliminating per-page Cloud Service fees.
- **Operationally Secure** by allowing sensitive award data to remain entirely in your control behind an institutional firewall.

For the Research Administrator seeking to modernize their office's data intake infrastructure in 2026, dots.OCR is a strategic preference. It reliably transforms documents such as the Notice of Award from a static image into a dynamic data asset, helping fulfill the digital promise of modern research administration.