

# Best Technologies for Text Extraction from Tender Documents

Extracting text accurately from tender documents is critical – missing even a small detail could mean losing a bid. Below, we provide a comprehensive comparison of leading OCR (Optical Character Recognition) technologies for this task, including cloud services and open-source tools. We consider document types (scanned images vs. digital PDFs), volume (a few thousand pages per month and growing), language support (primarily English with some Hindi), accuracy, cost, and deployment options.

## Requirements and Challenges

- **Document Formats:** Tender documents arrive as scanned images (needing OCR) or as digital PDFs (which may already contain a text layer). An optimal solution should handle both. For digital PDFs, text can be extracted directly (no OCR needed), preserving 100% accuracy. For scans, robust OCR is required.
- **Volume:** Currently ~20–30 tenders per month, each with 4–5 documents of ~20–30 pages. That's roughly 2,500–4,500 pages monthly (potentially ~5,000 as volume grows). The solution must scale to higher volumes as the bidding process speeds up.
- **Languages:** Documents are mostly in English, but some sections may be in Hindi or other regional languages. Multi-language support is important so that no text is missed.
- **Accuracy: No data can be missed.** The OCR must be highly accurate for printed text (and possibly occasional handwriting) to capture all project details, deadlines, requirements, etc. from tenders. High accuracy is paramount to avoid losing tenders due to missing information.
- **Budget and Deployment:** Open to paid cloud services (AWS, Google, Azure, etc.) as well as open-source or on-premise solutions. We must weigh the convenience and accuracy of cloud APIs against the control and privacy of local solutions. Data security is a consideration – ideally, avoid exposing sensitive tender data to external systems unless the benefit is clear.

## Cloud-Based OCR Services

Leading cloud OCR APIs offer excellent accuracy and scalability, with pay-as-you-go pricing. All three major providers (AWS, Google, and Azure) have **comparable base pricing (around \$1.50 per 1,000 pages)** for basic text extraction <sup>1</sup>. They also provide some free usage tier (e.g. first 1,000 pages free each month). Below we detail each:

### Amazon Textract (AWS)

Amazon Textract is a managed ML service for document text extraction and analysis:

- **Features:** It can do plain text OCR and also identify structured data like form fields and table cells. These advanced features ensure you don't miss data in forms or tables (e.g. key-value pairs, line items) within tenders <sup>2</sup> <sup>3</sup>. Textract can even handle handwriting to an extent <sup>2</sup>.

- **Accuracy:** Textract performs well, especially on **low-quality or messy scans**, where it has been noted to interpret words and layout effectively <sup>4</sup> . For clear printed text, its accuracy is on par with others. However, in at least one benchmark, Textract's overall accuracy slightly trailed the best performers (Google/Azure) on form-like documents <sup>5</sup> . In practice, Textract is highly reliable for English printed text, but it may occasionally miss subtle details that more refined engines catch.
- **Language Support: Limited** to mainly Latin-script languages. As of current documentation, Textract officially supports English, Spanish, Italian, Portuguese, French, and German <sup>6</sup> . It **does not reliably support Hindi** or other Indic scripts yet. If a tender has some Hindi text, Textract might return gibberish or nothing for those sections <sup>6</sup> . This is a drawback in multi-language scenarios.
- **Cost: Very low for our volume.** Basic text extraction (Detect Document Text API) costs about **\$0.0015 per page** (i.e. \$1.50 per 1,000 pages) <sup>7</sup> . At ~5,000 pages/month, that's only ~\$7.50. Advanced features are pricier: if extracting forms or tables using the Analyze Document API, AWS charges an additional ~\$0.05 per page for forms and ~\$0.015 per page for tables <sup>8</sup> <sup>9</sup> . These can stack if a page has both. For example, a page with a form and a table would cost \$0.065. In a worst case of all pages needing form & table analysis, 5,000 pages would cost ~\$325 <sup>10</sup> <sup>9</sup> – but if only some pages have tables/forms, the cost is proportionally lower. **Bottom line:** plain OCR is extremely affordable, and even with form/table analysis, the cost per tender is only a few cents.
- **Deployment & Security:** Textract is only offered as a cloud API (no on-premise version). You would need to upload documents to AWS for processing. AWS promises not to store or misuse the data, but it does mean tender documents leave your environment. If data residency or confidentiality is a concern, you'd need to consider encryption and AWS's compliance measures, or use another solution.

## Google Cloud Vision OCR / Document AI

Google offers OCR through its Cloud Vision API and a more specialized Document AI platform:

- **Features:** The **Cloud Vision API** provides OCR (called `DOCUMENT_TEXT_DETECTION`) that extracts all text and its layout positions <sup>11</sup> <sup>12</sup> . It's a straightforward service for getting text from images/PDFs. Google also has **Document AI** with pre-trained models for structured data extraction (invoices, forms, etc.), and a general Document OCR model that is optimized for documents. Document AI can classify and parse documents into fields, but using it is only necessary if you need a structured JSON output; for just text extraction, the base OCR is sufficient.
- **Accuracy:** Google's OCR is among the best in the industry. In diverse tests (including printed, handwritten, multilingual, and low-quality documents), **Google performed admirably on all document types** <sup>13</sup> . Its recognition quality is on par with Azure's – both tend to slightly outperform Textract in overall accuracy <sup>5</sup> . For clean documents, all three are almost equally accurate <sup>4</sup> . Google's OCR also excels at maintaining layout ordering. It handles handwriting reasonably but not perfectly (similar to others).
- **Language Support: Extensive.** Google's OCR supports over 100 languages, including Hindi and other Indic scripts. It can auto-detect multiple languages in one document. One review noted Google's language support is **on par with Microsoft's and Tesseract's – and better than Amazon's** <sup>14</sup> . This means if your tender has a mix of English and Hindi text, Google will likely capture both (where Textract would fail on Hindi).
- **Cost: Very low, similar to AWS.** Cloud Vision OCR is priced at **\$1.50 per 1,000 pages** (first 1,000 pages per month free) <sup>15</sup> <sup>16</sup> . So ~5,000 pages would be ~\$7.50, essentially the same ballpark as AWS. Document AI (the specialized processors) is higher – roughly \$0.01 per page for the general document model <sup>17</sup> (so ~\$50 for 5,000 pages) – but again, you might not need this unless you

require structured field extraction. The simpler OCR is both effective and cheapest. Google also offers \$300 free credit to new customers which can cover a lot of pages.

- **Deployment & Security:** Google's OCR is a cloud service only; no on-prem Google OCR is available (aside from using the open-source Tesseract engine that Google maintains – see Open Source section). Documents would be sent to Google's servers. Like AWS, Google assures data confidentiality and does not use the content for training by default. If cloud usage is acceptable, this is fine. If not, you'd use a local tool (Google doesn't provide an official local container as Azure does).

## Microsoft Azure Form Recognizer (Azure AI Document Intelligence)

Microsoft's Azure AI provides OCR and document processing under the **Document Intelligence/Form Recognizer** service:

- **Features:** Azure's **Read API** performs OCR on documents, and the service also offers **prebuilt models** for things like layout, tables, invoices, ID cards, etc., as well as **custom trainable models**. For tenders, the general **layout model** can extract text plus basic structure (tables, lines, paragraphs). Azure can also answer queries on documents and identify key-value pairs similar to Textract's forms feature. This is a very flexible platform – you can start with pure text extraction or move to structured extraction as needed.
- **Accuracy:** Azure's OCR is **top-tier**. Benchmarks show Azure and Google virtually tied at the top for accuracy on printed documents <sup>5</sup>. Azure's engine handled a variety of scripts (Latin, Chinese, Cyrillic) and even picked up some text hidden in document images in one test <sup>18</sup>. In practical terms, Azure's accuracy for English (and mixed-language) tender docs will be excellent – it is designed for “document-heavy” scenarios and often improves with each generation of the model. Any differences between Azure and Google in accuracy are minor and case-dependent. Notably, Azure and Google both significantly outperform a default open-source OCR like Tesseract on difficult inputs <sup>5</sup>.
- **Language Support: Extensive.** Azure's latest Read OCR model supports **164 languages**, including English, Hindi (Devanagari script), and other regional languages <sup>19</sup>. It supports mixed-language documents and will extract text in Hindi correctly (assuming the text is reasonably legible) <sup>19</sup>. This broad language support is a strong point (Azure expanded to many scripts including Devanagari, Cyrillic, Arabic, etc.). In short, Azure can handle the occasional Hindi sections in tenders without issue – a significant advantage over Textract.
- **Cost: On par with AWS/Google, if not slightly cheaper.** Azure's **Read OCR** is about **\$1.50 per 1,000 pages** for printed text (first 500 pages free per month on the free tier) <sup>20</sup>. That's ~\$7.50 for 5k pages. Azure's pricing is noted as **the most budget-friendly by a narrow margin** among the big three <sup>20</sup>. The difference isn't huge at our scale, but for completeness: after 1M pages, Azure's price drops (similar to others) – e.g. Azure and AWS drop to ~\$0.60/1000 beyond certain large volumes <sup>1</sup>. Specialized Azure models (layout, forms, etc.) might have slightly different pricing (e.g. an additional cost per page for form recognition or custom model), but for basic text it's the same low cost. Even using more advanced models, the cost per page is in the few-cents range. So cost will not be a barrier.
- **Deployment & Security:** Azure stands out for offering an **on-premises container** option. You can download the Azure OCR Docker container and run it in your environment, while still paying per page (it connects to Azure for billing/license) <sup>21</sup> <sup>22</sup>. This is great for data privacy – documents never leave your servers, satisfying strict security needs, yet you get the same OCR quality. Most customers use the cloud API for convenience <sup>23</sup>, but Azure's container is available if needed for compliance. In terms of ease of integration, Azure's service is known to be very developer-friendly with good documentation. One review highlighted that **Azure was the easiest to set up** among the

cloud OCRs, with minimal friction to get started <sup>24</sup>. On the downside, that same review noted Azure's OCR output sometimes included **erratic formatting**, like stray punctuation characters that weren't actually in the text <sup>25</sup>. This seems to be a minor quirk that required a little post-processing (e.g. filtering out lines of just dashes or commas). It wasn't observed with Google or AWS outputs. So Azure might need a tiny bit of cleanup logic, but it's not a deal-breaker.

## Other Commercial OCR Solutions

Aside from the "big three" cloud vendors, there are a few other notable OCR options:

- **ABBYY FineReader / ABBYY Cloud OCR:** ABBYY is an industry leader in OCR, long known for very high accuracy, especially on structured and printed documents. They offer SDKs and a cloud API. Accuracy is comparable to (or sometimes better than) the big cloud OCRs, and language support is broad. However, ABBYY's solutions are typically **licensed software (paid)** rather than on-demand pay-per-use. For example, ABBYY FineReader Engine can be deployed on-prem but has a licensing cost, and their cloud OCR SDK charges per page (pricing tends to be higher than the big three's \$1.5/1000 rate). ABBYY is worth considering if you need top accuracy and are willing to manage licensing or if you prefer an on-prem engine without the DIY of open-source.
- **IBM Watson / Others:** IBM's Watson Visual Recognition had OCR capabilities, and IBM also integrated ABBYY tech in some offerings. These are not as commonly used for standalone OCR now. There are also specialized APIs like **OCR.Space**, **Microsoft 365 (SharePoint) OCR**, etc., but they either use underlying engines from the above or have limitations. Given our needs, AWS/Google/Azure or ABBYY are the main commercial choices.

## Open-Source and Self-Hosted OCR Options

Open-source OCR engines can be run locally (or on your private cloud) to keep data in-house and avoid API costs. The trade-off is often lower accuracy out-of-the-box and more effort to set up and maintain. Key open-source tools include:

### Tesseract OCR

- **Overview:** Tesseract is a free, open-source OCR engine originally developed by HP and now maintained by Google <sup>26</sup>. It's command-line based and has wrappers for many languages (e.g. Python's `pytesseract`). It's widely used in many OCR applications (e.g. if you've used apps that OCR PDFs on your computer, they often use Tesseract under the hood).
- **Language Support:** Excellent. Tesseract comes with trained models for over 100 languages, including English and Hindi (Devanagari script) <sup>27</sup>. You can OCR multi-language documents by enabling multiple languages at once. This means Tesseract **will not miss Hindi text** from a language standpoint (assuming the model data is installed).
- **Accuracy: Good but not the very best.** For clean, high-quality scans or digital text, Tesseract can achieve high accuracy (near 98-99%). Its weaknesses show up with **noisy, low-resolution, or complex documents**. It **struggles with documents that are not clean, machine-printed text** – for example, scans with noise, uneven lighting, or any handwriting can reduce accuracy <sup>28</sup>. In comparisons, Tesseract usually scores below the cloud services on challenging documents <sup>5</sup>. It might misread characters (e.g. confuse O/0, I/1) or skip faded text that ML-based OCR could catch. That said, you can improve results by preprocessing images (deskewing, denoising, increasing

contrast, etc.). For our use (printed tender documents), Tesseract will likely do *fairly well* if the scans are decent. But if a tender PDF is a 10th-generation photocopy scan, Tesseract might miss some text that Azure/Google would catch.

- **Speed:** Tesseract is CPU-based and can process a page in maybe 1-2 seconds (depending on text amount and CPU power). 5,000 pages might take a few hours on one machine. This is usually fine for overnight processing, but not instantaneous. Cloud services internally parallelize across many machines, which is why they scale without trouble. With Tesseract, if volume grows, you may need to run OCR on multiple servers or threads.
- **Cost: Free license** (no per-page fees). The only costs are computing infrastructure. For ~5k pages/month, a single mid-range server (or even a desktop PC) could handle it. As volume increases, you might invest in more hardware or cloud VMs to run Tesseract. Still, compared to even \$50/month API cost, running your own server can be cost-effective if you have the expertise. Keep in mind the “cost” of developer time to integrate and maintain the OCR pipeline, and possibly the cost of errors: if Tesseract misses something and it’s not caught, the business impact could outweigh savings.
- **Integration:** Tesseract outputs plain text (or hOCR/pdf with layout). You would need to write code to ingest that text, then apply your own logic or AI to find the fields like project details, deadlines, etc. (This is true for any OCR though – raw text needs further processing for your use case). Tesseract is easy to call via scripts and well-documented <sup>26</sup>. Many PDF OCR tools (like **OCRmyPDF**) can use Tesseract to add text layer to scanned PDFs, which might be a quick way to make documents searchable.
- **When to use:** If you require data to stay on-prem and want to avoid any API usage, Tesseract is the go-to solution. It’s also a good fallback or baseline. You could run Tesseract first and then maybe only send to a cloud API those pages or lines that Tesseract found troublesome (advanced strategy to control costs/data sharing). Overall, Tesseract will **get the job done** for many cases, but you should budget for a robust QA process to ensure it truly didn’t miss critical data.

## Advanced Open-Source OCR (Deep Learning based)

In recent years, new open-source OCR engines powered by deep learning have emerged, often with superior performance to Tesseract on certain tasks:

- **docTR:** An end-to-end OCR library by Mindee, built on TensorFlow/PyTorch <sup>29</sup>. It uses modern text detection and recognition models. In a 2023 review, docTR was one of the top contenders among open-source tools, with accuracy approaching that of cloud APIs on printed text. It handles layout analysis as well. This could be a strong open option if you’re comfortable with Python and possibly using a GPU for inference.
- **EasyOCR:** A popular Python OCR package that is easy to use (hence the name). It’s based on deep learning and supports 80+ languages (including English and Hindi). EasyOCR often outperforms Tesseract on things like scene text or stylized fonts, but for standard documents its accuracy is in a similar ballpark. It’s worth testing – it’s pure Python and can run on CPU or GPU.
- **PaddleOCR:** An open-source OCR from Baidu’s PaddlePaddle framework. It has very accurate pre-trained models for multiple languages. PaddleOCR is known for good accuracy on Chinese and English text; it also provides layout and table recognition. It might handle noisy scans better than Tesseract due to the deep learning approach.
- **Others:** There are research models like Microsoft’s **TrOCR** (a transformer-based OCR model available on HuggingFace) which achieve excellent results on printed text. These require more ML know-how to use and possibly fine-tune. Also, frameworks like **LayoutLM** (for understanding document layout with text) might come into play after OCR to help categorize the extracted info.

Using these open-source DL models locally can potentially give you accuracy closer to Google/Azure, but the setup is more involved. You'll need to manage model dependencies, possibly have a GPU for speed, and handle memory/CPU usage for large documents. However, since your volume is not huge, even CPU might be okay with these models if moderately optimized.

**Summary of Open-Source vs Cloud:** Open-source tools cost nothing to use and keep data in-house, but typically require more effort to achieve the same accuracy. Cloud tools are practically plug-and-play with high accuracy and a very low cost at this scale (single-digit dollars). Many teams choose cloud OCR for convenience and reliability, unless data privacy rules mandate an on-prem solution.

## Accuracy and “No Data Missed” Considerations

Given that missing any tender detail is unacceptable, accuracy is the top priority. Here's how the options stack up and ways to ensure nothing is missed:

- **Cloud OCR Accuracy:** All three major cloud OCR engines have **very high accuracy on printed English text**. On clear documents, they will extract >99% of the text correctly (occasional minor OCR typos can occur, but rarely would they miss whole words or lines). Among them, **Google and Azure have a slight edge in independent tests** <sup>5</sup>. They consistently ranked at the top for accuracy and similarity to ground truth. **AWS Textract** is not far behind, and in certain cases (messy, low-quality scans) it might even do better at reading degraded text <sup>4</sup>. For multi-language content, Google and Azure clearly outperform Textract (since Textract might skip non-Latin text entirely). If your tender documents are standard printed material, you can expect all three to do a great job, with a slight confidence boost in favor of Azure/Google for edge cases.
- **Open-Source Accuracy:** Tesseract and similar open tools can approach the 90-98% accuracy range on printed docs, but the last few percent (where weird errors happen) can make a difference. For example, Tesseract might mis-parse a table or mix up letters in a way that requires manual correction. It's also more sensitive to image quality. **In one comparison, Tesseract's accuracy was measurably lower than Azure's and Google's on a forms dataset** <sup>5</sup>. That doesn't mean it will miss whole sections of text, but it might misread some numbers or letters. If a tender requirement is misread (e.g. “1%” read as “7%” or “2023” as “2028”), that could be problematic. Therefore, if using open source, you should implement a verification step (perhaps cross-check critical fields or have a human quickly review the extracted data).
- **Ensuring No Data is Missed:** Some strategies to maximize completeness:
- **Use Multi-Engine Cross Verification:** For extremely critical data, you could run two different OCR engines (say, Tesseract and Azure) and compare results. Any discrepancies could be flagged for review. This improves confidence but doubles processing.
- **Leverage Document Structure:** Since tenders have sections (Eligibility, Deadlines, etc.), you could parse the raw text and ensure each expected section is present. If, for example, the “Submission Deadline” date isn't found in the OCR text, that's a red flag to manually check that document.
- **Manual QA for Edge Cases:** Initially, you might have a person quickly skim the auto-extracted info vs. the original for a few tenders to gauge where the OCR might be slipping. This can help fine-tune your process (e.g. if tables are an issue, maybe use a table-specific extraction).
- **Use Layout/Forms extraction:** If tender documents come in structured forms or tables (sometimes government tenders have standard forms to fill), using Textract's form/table features or Azure's layout model could capture data in a structured way, potentially reducing errors (e.g. always picking the right field value rather than relying purely on text flow). This can complement raw text OCR.

In summary, to **absolutely minimize risk**, the safest route is to use one of the top cloud OCRs (Azure or Google) which have the highest accuracy and language coverage <sup>5</sup> <sup>14</sup> . Their cost is so low that it's worth it for peace of mind. If cloud is off-limits, using Tesseract or another open source engine with careful validation can work, but you'll need to be vigilant early on to ensure nothing is missed.

## Cost Comparison

One of the advantages of these services is that our scale of a few thousand pages is relatively small, so costs are very manageable. Here's a quick comparison for ~5,000 pages/month:

- **AWS Textract:** ~\$7.50 per month for text-only extraction (at \$1.5 per 1000 pages) <sup>7</sup> . If every page also needed form and table analysis, worst-case cost ~\$325 for 5k pages <sup>9</sup> – but that scenario is unlikely unless every page is densely structured. Realistically, you might selectively use form/table features on certain pages. Textract also has a 3-month free tier (1,000 pages text or 100 pages forms/tables free per month) <sup>30</sup> .
- **Google Cloud Vision:** ~\$7.50 per month for 5k pages (first 1k free, then \$1.5/1000) <sup>15</sup> <sup>16</sup> . Google's pricing remains \$1.5/1000 up to 5 million pages, then drops to \$0.60/1000 beyond that <sup>1</sup> . Document AI's \$0.01/page would be ~\$50/month if used <sup>17</sup> , but we likely don't need to incur that for just text extraction.
- **Microsoft Azure:** ~\$7.50 per month for 5k pages (first 500 pages free; \$1.5/1000 after) <sup>20</sup> . If scaled to huge numbers, Azure's price also eventually drops similarly. Azure slightly "edges out" others in cost by maybe a fraction of a dollar per thousand pages <sup>20</sup> – the difference at our scale is negligible.
- **Other fees:** All cloud providers charge by pages processed only; there are no separate compute charges (the pricing above is all-inclusive). There might be minor costs for storage if you upload files to cloud storage, but since you can send documents directly and our volumes are small, that's minimal.
- **Open Source:** No direct fees. Suppose you run Tesseract on an AWS EC2 instance or local server – the cost would be the server cost. For example, a medium cloud VM might be ~\$50/month, which is already higher than the API costs. But you likely have existing infrastructure or a lower-power machine could suffice. If the machine is doing other tasks (or if you only spin it up when needed), the marginal cost is low. In on-prem terms, the cost is electricity and wear on a machine you own. Essentially, open source is "free" but not truly zero cost if you factor in hardware and maintenance.

**Conclusion on cost:** All cloud solutions are inexpensive for the given volume – on the order of **a few dollars per month**. Even if the number of tenders doubles or triples, you're still well under \$50/month. This is likely trivial compared to the value of winning more tenders. Open-source is financially free but has an opportunity cost in development time and potentially slightly higher risk of errors. If budget is not extremely tight, cost should not be the deciding factor here – accuracy and reliability should be.

## Data Security Considerations

Since tender documents may be sensitive, it's worth noting how each approach handles security:

- **Cloud services** (AWS, Google, Azure) process data on their servers. All provide encryption in transit and at rest. They also have strict privacy policies – e.g., Google and Azure explicitly state they do not use your document content to train models and do not store it longer than needed to process. AWS

Textract likewise only processes data on the fly. All three are compliant with various standards (ISO, etc.). **However**, if your organization or client has a policy against uploading confidential documents to external services, this is a roadblock. Azure's on-prem container is a nice middle ground: you can keep data internal while still using their OCR engine <sup>22</sup>.

- **On-Premise/Open-Source:** Running OCR locally means the documents never leave your controlled environment. This is ideal for data privacy. You do need to ensure the system is secure (no unauthorized access to the OCR output, etc.). But essentially, you eliminate the exposure to third-party cloud providers. This may be necessary for certain government or highly confidential tenders.
- **Hybrid:** It's possible to design a system that uses open-source OCR for most of the document, and maybe only queries a cloud service for specific high-value fields or as a double-check. This limits what is sent out. For example, you might OCR everything with Tesseract, then if a crucial field (say "Bid Amount") is not confidently extracted, you send just that snippet image to an API for a second opinion. This is a complex approach and might not be needed given the strong privacy protections the big providers have, but it's an option if one must minimize external data exposure.
- **Compliance:** If you go the cloud route, ensure you use regions and settings that comply with data locality requirements (e.g. choose an India data center region for Azure or AWS if that's desired, so data doesn't leave the country). All providers let you specify region endpoints.

Given the note in your requirements about *not exposing internal data outside the portal*, if that is a hard rule, then the clear choice is an on-prem solution – either open-source OCR or Azure's container-based OCR – so that documents are processed within your infrastructure. If the business is open to cloud as long as it's secure (which they indicated they are open to using AWS/GCP), then leveraging those services can be done in a compliant way (for instance, many companies process documents in AWS/GCP after signing the appropriate agreements and using encryption).

## Recommendations

Considering all factors – accuracy, language needs, volume, cost, and security – here are our recommendations and alternatives:

- **Best Overall (Accuracy & Ease): Google Cloud Vision OCR or Azure Form Recognizer** are top choices. Both will reliably extract all text (English or Hindi) with very high accuracy <sup>14</sup> <sup>19</sup>. They are easy to integrate (just API calls) and cost only a few dollars a month at current volumes. Azure has an edge if you later want an on-prem deployment (container) or if you prefer slightly simpler setup. Google has equally strong OCR and might be preferred if you already use Google Cloud. Either of these will minimize the risk of missing data.
- **AWS Textract if using AWS stack:** Textract is a solid option if your infrastructure is AWS-centric or if you want the built-in form/table parsing features. It's almost as accurate on clean text, but **keep in mind language limitations** – if you expect non-English text (Hindi), Textract could miss those segments <sup>6</sup>. If that's rare or if you can live with that risk, Textract otherwise meets the needs well. The cost and speed are comparable, and it offers valuable form extraction capabilities which might be handy for automating tender forms.
- **Open-Source Local Solution:** If keeping data internal is paramount, you can start with **Tesseract OCR**. It's free and supports English/Hindi well <sup>27</sup>. However, plan for extra validation steps to ensure nothing was misread. You might also consider more advanced OS tools like **docTR or EasyOCR** to possibly improve accuracy on tough documents. This route requires more engineering: you'll need to set up an OCR server or service, manage performance (maybe batch OCR jobs with a queue), and



handle updates/tuning. Over time, you might incorporate machine learning to identify key fields from the text. This approach gives you full control and no external dependencies, at the cost of more initial effort.

- **Hybrid Approach:** You don't have to choose one exclusively. For example, you could use open-source OCR for initial processing (keeping data in-house) and then use a cloud OCR on the same documents as a "audit" occasionally to see if the results differ significantly, which could alert you to any missed data. Or use cloud OCR for the most critical documents while using local for others. Since volume is not huge, even double-processing some documents won't break the bank.
- **Leveraging PDF Text:** Regardless of OCR choice, ensure your pipeline first checks if a PDF has an embedded text layer (many digital PDFs do). If yes, **extract that text directly** (using a library like PDFBox or PyMuPDF) instead of OCR – this gives you 100% accurate text and is faster. Only pages that are pure images require OCR. This hybrid PDF handling will ensure you truly don't miss anything and speeds up processing of digital files.
- **Post-OCR Processing:** Remember that after text extraction, you plan to **categorize information** (deadlines, criteria, etc.) and auto-fill forms/emails. Achieving this will likely involve some combination of keyword matching, regex, or using an NLP/LLM to parse the text. For example, once you have the text, you could use rules to find "Submission Deadline: \_\_\_\_" or have a small ML model to pick out dates and requirements. The quality of OCR output will feed into this – fewer errors in OCR make the next steps easier. So choosing a high-accuracy OCR (cloud or a well-tuned model) will improve your success in the auto-filling stage.
- **Future Scaling:** If the number of tenders grows significantly (say 10x), all these solutions can handle it. Cloud APIs would charge more but remain affordable (e.g. 50k pages = ~\$75 with Google/Azure). Open-source would require scaling up servers or optimizing code, but could also handle it with investment in hardware. It's wise to design a pipeline that can run in parallel (process pages or documents concurrently) to speed up throughput if needed.

In conclusion, to **maximize accuracy and not miss any data**, a **cloud-based OCR solution (Azure or Google)** is highly recommended <sup>5</sup>. They offer proven accuracy, full support for English and Hindi text, and negligible cost at your scale. If cloud usage is acceptable, they provide the most "impactful" result with minimal data loss. If cloud is not an option, the open-source route with Tesseract or similar can be made to work, but be prepared to invest in validation and perhaps use multiple OCR engines to reach the same confidence level.

By implementing the right OCR technology now, coupled with AI for auto-filling and organization, you will greatly speed up tender processing and reduce the chance of human error – enabling your team to bid on more tenders each month with confidence.

#### Sources:

- AWS Textract pricing and features <sup>7 8 9</sup>
- AWS Textract supported languages (currently limited to a few Latin-based languages) <sup>6</sup>
- Google Cloud Vision OCR pricing and features <sup>15 16</sup>
- Google Document AI pricing (general document OCR) <sup>17</sup>
- Comparative OCR accuracy (Google & Azure leading, Amazon & Tesseract slightly behind) <sup>5</sup>
- OCR accuracy notes – Textract good for poor scans; all similar on clean text <sup>4</sup>
- Language support comparison (Google/Azure support Hindi and many languages; Amazon notably behind in this aspect) <sup>14 19</sup>
- Azure OCR ease-of-use and minor output quirks <sup>24 25</sup>

- Azure on-premise container availability for OCR <sup>21</sup> <sup>22</sup>
- Cost comparison of cloud OCR services (all around \$1.5 per 1000 pages) <sup>1</sup> <sup>20</sup>
- Open-source Tesseract overview and limitations <sup>27</sup> .

---

<sup>1</sup> <sup>13</sup> <sup>14</sup> <sup>18</sup> <sup>24</sup> <sup>25</sup> <sup>26</sup> <sup>27</sup> <sup>28</sup> <sup>29</sup> Our search for the best OCR tool in 2023, and what we found • MuckRock

<https://www.muckrock.com/news/archives/2023/oct/31/our-search-for-the-best-ocr-tool-in-2023-and-what-we-found/>

<sup>2</sup> <sup>3</sup> <sup>7</sup> <sup>8</sup> <sup>9</sup> <sup>10</sup> <sup>30</sup> Textract Pricing Page

<https://aws.amazon.com/textract/pricing/>

<sup>4</sup> <sup>20</sup> Decoding OCR: Your Guide to Picking the Right Text-Extraction Tool for Your Business

<https://www.velocis.in/blog/decoding-ocr>

<sup>5</sup> How to Compare OCR Tools: Tesseract OCR vs Amazon Textract vs Azure OCR vs Google OCR | by Federico Ricciuti | Medium

<https://federico-ricciuti.medium.com/how-to-compare-ocr-tools-tesseract-ocr-vs-amazon-textract-vs-azure-ocr-vs-google-ocr-ba3043b507c1>

<sup>6</sup> Is there a way to specify language while detecting text in AWS ...

<https://repost.aws/questions/QUVeFLYhy8RFWKRH0H94fXVA/is-there-a-way-to-specify-language-while-detecting-text-in-aws-textract>

<sup>11</sup> <sup>12</sup> <sup>15</sup> <sup>16</sup> Pricing | Cloud Vision API | Google Cloud

<https://cloud.google.com/vision/pricing>

<sup>17</sup> Pricing | Document AI | Google Cloud

<https://cloud.google.com/document-ai/pricing>

<sup>19</sup> <sup>21</sup> <sup>22</sup> <sup>23</sup> OCR - Optical Character Recognition - Azure AI services | Microsoft Learn

<https://learn.microsoft.com/en-us/azure/ai-services/computer-vision/overview-ocr>