Deep Analysis

# Amazon (AWS) Textract

**Founded 2006 | HQ Seattle, WA | 25,000 employees (approx.) | $20B annual revenue (est.)**

*Textract may be one small piece of a big Amazon pie, but it is a pivotal piece. Textract is a catalyst for information management buyers and vendors to leverage Amazon Web Services (AWS) as a true platform rather than simply as a provider of cloud storage.*

## The Company

Amazon was founded in 1995 and launched Amazon Web Services (AWS), its subsidiary business, in 2006. AWS provides a wide range of on-demand cloud computing services, from storage to artificial intelligence. Amazon is headquartered in Seattle, Washington, and led by CEO Jeff Bezos. As of 2019, AWS had around $35 billion in annual revenue, with (our estimate) 25,000 employees. This report focuses on the AWS Textract service, which launched in May 2019.

## The Technology

Amazon describes its machine learning (ML) technology stack as a three-layered cake. At the bottom is a layer for advanced practitioners to work with deep learning frameworks such as Tensorflow or Pytorch. The next layer is for developers to work with Amazon SageMaker to build out managed ML capabilities. The top layer provides pre-built artificial intelligence

(AI) services that can be used immediately via an API integration point. Textract is one of these AI services offered by AWS. In simple terms, Textract provides data extraction from scanned documents. It does this by delivering ML-based optical character recognition (OCR) tools to read the document and accurately extract data and text.

OCR tools have been on the market for decades but came into the mainstream in the 1990s. However, though the concept of reading a letter or number in a scanned document is simple, it is extremely difficult to do with any degree of accuracy. The resolution quality of the scan and differing font sizes and types – not to mention handwriting quality – make the task of data extraction infinitely complex. No OCR does the job perfectly, including Textract; there will always be an error rate and exceptions to deal with. Yet even though Textract has been on the market less than a year, tests show its accuracy rate is as good as, and possibly better and more sophisticated than, many older systems on the market. And it should be noted that this accuracy rate will only improve as the system
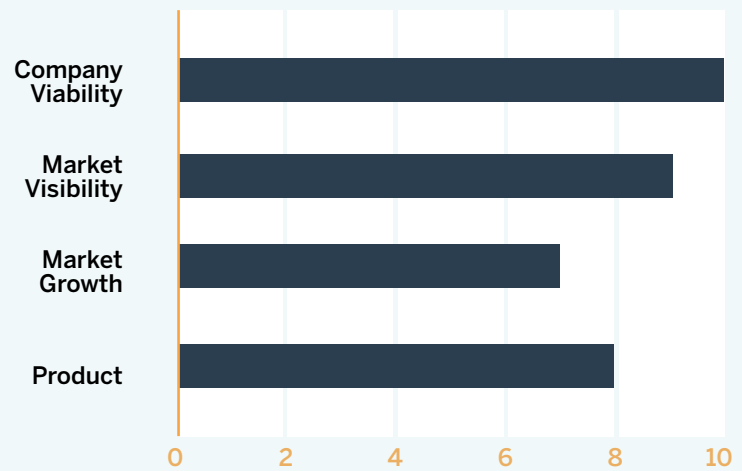
is used more extensively and the underlying machine learning has the opportunity to learn and improve.

As of today, Textract relies on well pre-trained models that, assuming your data is a match for those models, will produce excellent results. If not, then the results will not be as accurate. Where Textract goes further than many, if not most, traditional OCR systems is that it doesn't just capture raw text. It also has built-in form and entity extraction capabilities, and it can extract text from tables and images. So, for example, it can extract multiple tables within a document and recognize the keys and values to ensure the table data is moved intelligently.

Much of the market interest in Textract to date has focused on its performance versus traditional OCR systems. Though that is to be expected, we would argue that such a narrow focus misses the broader impact that Textract will likely have on the market. Textract is just one element in the AWS portfolio. Alongside it sit other AWS machine learning data extraction tools for video, images, and speech; tools for language translation, language comprehension, and text analytics. In turn, these tools are part of the broader AWS platform of on-demand computing and storage. Add to this a vast partner and developer ecosystem leveraging these combined capabilities, and the possibilities are enormous. The ambitions of AWS and the rapid rollout of their products represent a turning point for document and content management. That said, vast though these capabilities are, they are not without limitations.

Not every firm wants to use Amazon's cloud computing services, and many document process activities (particularly high-touch and highly sensitive activities) remain on-premises.



Figure 1
Amazon Textract Assessment

They will do so for the foreseeable future, and AWS appears to have no plans to cross the product chasm from cloud to on-premises. Just like any other AI or ML technology, the ones that AWS provides are only as good as the data they are trained on, and AWS's preference is to train the machines itself. However logical, this does place limitations on end users' ability to modify the system to their specific needs. Similarly, the quality of the scan or, for example, dodgy handwriting will challenge its accuracy rate.

However, in our analysis Textract is a compelling proposition when used as part of a broader application or in conjunction with other AWS services for digital transformation purposes. And this is where it gets interesting. Seen in isolation, Textract is just a modern twist on traditional OCR technology. But when Textract is seen as a tool to read a document, understand its structure, and extract information from it in the form in which it was meant to be read, that is surely just the starting point for transformation and automation.

## Our Opinion

Textract may be one small piece of a big Amazon pie, but it is a pivotal piece. Textract is a catalyst for information management buyers and vendors to leverage AWS as a true platform rather than simply as a provider of cloud storage. Though labeled as OCR, even a cursory look at Textract makes it clear that this is a fully-fledged content extraction tool, one designed to provide the data necessary for further automation. It's not perfect, particularly when business needs don't align with Amazon's pre-trained models, but it will improve with use over time.

The question is, where does Textract go from here? Clearly, AWS has other elements in its portfolio to complement Textract that can be used to build sophisticated applications: not only additional data analysis and extraction tools like Rekognition (video and images) and Transcribe (speech recognition), but also ElasticSearch, Blockchain, IoT, and even content management in the form of WorkDocs. What it lacks is the automation and process glue to pull all these elements together into a fully executable form. For now, the company seems happy to let its partner ecosystem do the work of building applications and taking them to market. But over time, as more considerable opportunities are identified, it's likely Amazon will take on more of that work, as they have already started to do in the call center world through Amazon Connect. What appears to be missing today are process and task automation tools in the AWS stack.

Figure 1 shows our assessment of Amazon Textract across four categories.

## Advice to Buyers

It's easier to state where Textract is not a fit than where it is a fit. If you run a complex capture operation in-house or outsourced, then Textract is not going to make a difference to you. The only impact it may have at some point is that one of your technology vendors may incorporate Textract into their software instead of some other option. But for areas of your business that are considering robotic process automation (RPA) or, for that matter, any other document-centric automation activity, Textract is well worth looking at and adding to your shortlist, as are some of the other AWS capture-related products. Today it is possible (in theory at least) to use Textract along with other AWS services to construct DIY content systems. For large volumes of legacy content, it's an interesting option to consider compared to traditional ECM and archiving technologies.

# 🔍 SOAR Analysis

## Strengths

→ Accurate, high-volume OCR capabilities

→ Part of a large family of related ML-based tools

## Aspirations

→ Be the leading capture technology in the market

→ Be a pacesetter in the content and data services industry

## Opportunities

→ Deliver more automation and process capabilities

→ Develop basic content-centric applications

## Results

→ Textract has improved notably in the months since its release

→ In a short time, Textract has become the No. 1 alternative to traditional OCR offerings

# About Deep Analysis

**Deep Analysis** is an advisory firm that helps organizations understand and address the challenges of innovative and disruptive technologies in the enterprise software marketplace.

Its work is built on decades of experience in advising and consulting to global technology firms large and small, from IBM, Oracle, and HP to countless start-ups.

Led by Alan Pelz-Sharpe, the firm focuses on Information Management and the business application of Cloud, Artificial Intelligence, and Blockchain. Deep Analysis recently published the book "Practical Artificial Intelligence: An Enterprise Playbook," co-authored by Alan and Kashyap Kompella, outlining strategies for organizations to avoid pitfalls and successfully deploy AI.

Deep Analysis works with technology vendors to improve their understanding and provide actionable guidance on current and future market opportunities.

Yet, unlike traditional analyst firms, Deep Analysis takes a buyer-centric approach to its research and understands real-world buyer and market needs versus the "echo chamber" of the technology industry.

## About the Author

Alan Pelz-Sharpe is the founder of Deep Analysis. He has over 25 years of experience in the IT industry, working with a wide variety of end-user organizations like FedEx, The Mayo Clinic, and Allstate, and vendors ranging from Oracle and IBM to start-ups around the world. Alan was formerly a Partner at The Real Story Group, Consulting Director at Indian Services firm Wipro, Research Director at 451, and VP for North America at industry analyst firm Ovum. He is regularly quoted in the press, including the *Wall Street Journal* and *The Guardian*, and has appeared on the BBC, CNBC, and ABC as an expert guest.

## Contact us:

info@deep-analysis.net

+1 978 877 7915