

Improving Accuracy & Efficiency in Document Handling for Business Processes

Alison Major
Computer Science
Lewis University
Romeoville, Illinois, USA
AlisonMMajor@lewisu.edu

Abstract—Data is found in a variety of digital formats: email messages, spreadsheets, images, sound files. Many tools have been designed to assist in automatically extracting data from certain formats. Optical Character Recognition (known as OCR) can be used to convert images into text, though not always in a structured and useful format. Natural-language processing (known as NLP) allows computers to extract information from text written in the natural language form. Is there a method to allow multiple forms of information to be received by a system and automatically pull the desired data from those different formats? For example, if a customer submits a purchase order as a scanned image that we want to automatically pull item numbers from, what is the ideal method for processing that file and information? Similarly, if a customer submits an order in an email message in prose or with an embedded table, what are our options for extracting the desired information? This paper will review technologies that use OCR, NLP, and artificial intelligence to find optimal options for extracting data from a variety of inputs.

Index Terms—optical character recognition (OCR), artificial intelligence (AI), machine learning (ML), natural language processing (NLP), robotic process automation (RPA)

I. INTRODUCTION

Many companies receive input from users in a variety of ways. Some may email, some send attached files. There may be a web application used to submit forms of information or capture a photo of a document on their phone. Technology has provided us with ample options to submit information.

For a human to run submitted information through a business process, they need to read the data in whichever format it is received, then understand it enough to enter it into the appropriate system. This basic decision-making is generally not a difficult task, but when needed at a large scale, can become not only tiresome but prone to entry errors.

During a controlled 2009 study at UNLV, 215 students were provided with 30 datasheets that contained six types of data to process. When the students only used visual confirmation of the correctness of the data, they made an average of 10.23 errors. Progressive steps to automate confirming information improved the error count [1].

Employees with mountains of data to enter may not have time for more than a quick visual confirmation, resulting in more errors. While data can be processed cheaply by hiring large numbers of low-cost employees (for example, Amazon

Mechanical Turk), more can be done to provide efficient and accurate methods for data handling.

To demonstrate challenges and considerations in how to process a variety of inputs, we'll review the challenges and some data behind the manual entry of information (Section II). With the awareness of the problem at hand, we will define rules-based approaches and optical character recognition (OCR), with the challenges involved with these methods (Section III). Built on the foundation of the business flow for manual data entry and the use of rules and OCR, we will present a flow diagram for processing data (Section IV). This section will break down these steps for processing, define data classification types, and review the steps needed to handle the data in a repeatable way.

We'll focus on receiving orders and adding them to a company's ordering system. Ideal situations would involve an integrated system, where customers enter an order in an application (like amazon.com) and the order is then directly passed into the company's Enterprise Resource Planning (ERP) system. This type of system is used to provide the company with insights and internal controls, provide a central database, is used by accounting, supply chain, sales, and many other departments within a business. Most situations are not ideal, and orders are received in other methods that must then be entered into the ERP manually. When there are a large number of customers, or when customers carry a lot of influence in how they place orders (using a method that is convenient for them and the systems they use), there may not be opportunities to attempt direct integration.

Many solutions to parts of this problem are already available. We will then explore combining this flow for several document types into a single pipeline, review how several companies successfully used this method (Section V), and explore how we might do likewise (Section VI).

II. MANUAL DATA ENTRY

In 2018, Goldman Sachs reported that the direct and indirect cost of manual data entry for global businesses was estimated to be about \$2.7 trillion [2]. Gartner's study in 2019 found that avoidable rework in accounting departments amounted to 30% of a full-time employees' time; for an accounting staff of 40 full-time employees, this amounts to 25,000 hours per year and about \$878,000 [3]!

While modern applications provide us with many conveniences, many companies or industries have legacy processes in place. People deal with paper forms and documents every day, which then must be typed into digital systems. Employees will receive emails, faxes, and files that they must then determine which type of document it is and then process it accordingly for that customer, usually by typing the data contained into another system.

This manual data entry leaves room for error, as Gartner found. Additionally, manual entry is very time-consuming; talented employees could better spend their time on tasks that provide real value to building the customer relationship and improving the customer experience, rather than focusing on wearying data entry. These studies reiterate that people make mistakes; it is costly for businesses and worth the time to consider and address.

III. UNDERSTANDING EXTRACTION METHODS

A. Rules-Based Approach

Files that are already machine-readable are a good place to start in understanding how we can extract information. Documents that are already stored digitally can be accessed with simple computer programs. Once the file has been read by the program, it can find values and information in the contents using regular expressions and based on expected positions.

A2 -> PO#: "1234567"

D-Column (Item Numbers)
D9 -> Item #: "124"

PO# 1234567		Ship To Location		XYZ Back Door		
Ship Date 09/07/21				789 Main Ave		
Ship To XYZ Back Door				Anytown, USA		
Bill To XYZ Company						
Line #	Quantity	Unit of Measure	Item #	Description	Unit Price	Charge
1	7	CASE	123	Green Balls 16" Diameter	\$20.00	\$140.00
2	6	CASE	122	Red Balls 12" Diameter	\$18.00	\$108.00
3	8	CASE	124	Blue Balls 24" Diameter	\$25.00	\$200.00
4	2	CASE	143	Orange Balls 6" Diameter	\$15.00	\$30.00
Total Quantity		23	Total Cost		\$478.00	

Some fine print here. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Integer at enim tempor eros aliquet risusmod et in felis. Duis id risusmod magna. Lorem ipsum dolor sit amet, consectetur adipiscing elit.

Nullam eget suscipit oeci, sit amet tristique justo. Maecenas aliquam consectetur tellus ac fermentum. Nulla varius at dolor eget suscipit. Proin metus metus, condimentum eget diam ut, tristique dictum nisi.

G12 -> Total Cost: "\$478.00"

Fig. 1. We can programmatically find values in spreadsheets using positions and rules.

A document like that in "Fig. 1" has values stored in specific cells within the document. If we know where to look, we can tell a computer program to extract that value and store it for us. For example, in our figure, we can see that cell A2 contains the value for our purchase order number, the D column (in cells 7-10) has item numbers, and G12 has the total cost. We can map these locations to variables where we want to store this information.

B. Optical Character Recognition

Optical Character Recognition, more commonly known as OCR, is a technique that results in converting images and files into machine-readable data (information in a form that a computer can process). The OCR technology reviews images

for fonts and shapes, matching information with text that can be reviewed and stored.

While providing actionable information from different types of digital files, there are still many challenges with OCR. Images and pages may be received with incorrect orientation or skewed at an angle (imagine the paper is fed crookedly into a fax or scanner, or a picture taken with a phone that is not perfectly aligned). There may be noise or distortion on the page, or handwriting and other marks obscuring the printed text. If background colors or images are present, these can also impact the technology's ability to get a clear read on the information.

Some of these challenges can be overcome through the use of preprocessing. Pages can be scaled and adjusted to be the correct orientation. Filters can be applied to clean up noise. Newer iterations of OCR technologies also apply machine learning (ML) that can be used to better identify text characters when the image is unclear.

"PURCHASE ORDER"

PURCHASE ORDER

INVOICE: XYZ Company
123 State Street
Anytown, USA

XYZ Back Door
789 Main Ave
Anytown, USA

#1234567

DATE: 08/07/2021

DATE REQUIRED: 09/01/2021

PH: (111) 111-1111

FX: (222) 222-2222

COUNT	ITEM NUM	DESCRIPTION	UNIT PRICE	COST
7	123	Green Balls 16" Diameter	\$20.00	\$140.00
6	122	Red Balls 12" Diameter	\$18.00	\$108.00
8	124	Blue Balls 24" Diameter	\$25.00	\$200.00
2	143	Orange Balls 6" Diameter	\$15.00	\$30.00

124

TOTAL COUNT: 23

TOTAL COST: \$478.00

"\$478.00"

Fig. 2. Using OCR, we can translate the image into text, then we can programmatically find values for that text using positions and rules.

Introducing OCR to the process of receiving files from customers gives us a digital version of the data. Rather than an employee typing the information character-for-character, output from an OCR engine provides the employee with data that can be copy-pasted with minor corrections. Combining OCR with the rules-based approach, we can gather data from an image (see "Fig. 2"). While small, this technique introduces some new efficiencies that we can leverage in an improved data

processing flow.

IV. PROCESSING DATA

Let's say we run a ball factory that produces balls of different colors and sizes. We receive orders from our customers in many ways. One of our customers submits orders using a templated spreadsheet generated from their internal system. Another customer faxes their (different) templated order along, which is then received digitally as a scanned image. Yet another sends an email message containing the information for their order. With more customers, there can be a wide variety of ways that orders are submitted!

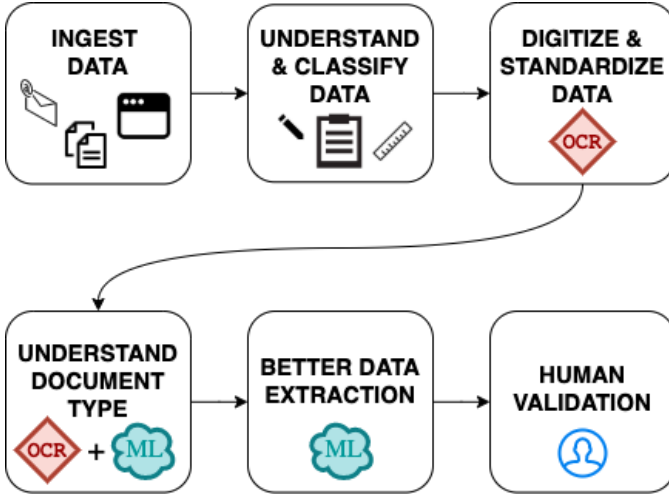


Fig. 3. Document flow in a generic business process for extracting data from different inputs.

Enforcing a specific business process will standardize the format and methods that customers use to place orders. We then could build an uncomplicated process to ingest orders using simple scripts and rules. We could even create macros in our spreadsheets to process the information and make it easy to enter into our ordering system.

Not all companies will be able to expect a standard input on all customers; instead, they must consider how to handle the different types of documents by creating a business flow like the one in “Fig. 3”. The features in this general business flow are detailed below.

A. Ingest Data

The first step in our flow involves ingesting data. This means that we are receiving it somehow. The intake step may be an employee receiving a spreadsheet or in the body of the message.

Some companies have different ingestion methods that allow for integrated flow to their ordering system, like having users place orders by logging into a website and entering their orders directly. As our initially defined ideal scenario, integrated systems allow for little human error once the order has been entered by the customer. New businesses may have this option, but those that have been around longer may have customers

that are used to a certain method or have their legacy ERP systems or large volumes of data that make manual entry into a web application difficult.

To continue in our journey through this business process, we'll assume that we are ingesting documents that come in various types and formats and are not part of an integrated system. We will assume that our customer orders for our ball factory are entered manually.

B. Understand & Classify Data

Since we are receiving orders in several file types (spreadsheet, pdf, email message, etc) that could be handled in different ways, we must find a method for understanding and classifying the data into one of three primary categories: Structured Data, Semi-Structured Data, and Unstructured Data. Once we have organized our documents into these groups, we can take the appropriate next action.

1) *Structured Data*: Often in the format of comma-separated values (CSV) and spreadsheets, structured data is formatted in a consistent template. The file contains tables with values in predictable positions. These documents could also be pdfs or other file types; the common thread is that the information is in the same layout for every document received.

Because of its predictability, this type of data can be gathered through simple data extraction scripts or template-based OCR. Rule-based approaches work well because we always know where to expect information that can be linked to a key-value pair (in the form of “key: value” like “date: 8/14/2021” or “PO: 1234567” where “date” and “PO” are keys).

PO# 1234567				Ship To Location		XYZ Back Door	
Ship Date 09/01/21						789 Main Ave	
Ship To XYZ Back Door						Anytown, USA	
Bill To XYZ Company							
Line #	Quantity	Unit of Measure	Item #	Description	Unit Price	Charge	
1	7	CASE	123	Green Balls 16" Diameter	\$20.00	\$140.00	
2	6	CASE	122	Red Balls 12" Diameter	\$18.00	\$108.00	
3	8	CASE	124	Blue Balls 24" Diameter	\$25.00	\$200.00	
4	2	CASE	143	Orange Balls 6" Diameter	\$15.00	\$30.00	
Total Quantity		23	Total Cost			\$478.00	
Some fine print here: Lorem ipsum dolor sit amet, consectetur adipiscing elit. Integer sit enim tempor eros aliquet euismod et in felis. Duis id euismod magna. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Nullam eget suscipit orci, sit amet tristique justo. Maecenas aliquam consectetur tellus ac fermentum. Nulla varius ut dolor eget suscipit. Proin metus metus, condimentum eget diam ut, tristique dictum nisi.							

Fig. 4. A templated spreadsheet for a purchase order.

Using our ball factory example, the spreadsheet in “Fig. 4” could be an example of a structured document, if all spreadsheets received follow the same layout template. This structured data allows for scripts to be written that can map information by using the expected position of a value or with regular expressions to collect the value from the file. We’ve proven a sample of this in the simple Python code in Appendix A (for related unit tests, see Appendix B).

In “Fig. 5” we can imagine a straightforward set of rules that determines where information exists if all documents follow the same layout. When data is structured, we have a predictable place to look for the values we need, and gathering information in an automated method can be quite simple.



Fig. 5. All documents that match “Layout 1” can be passed through the same set of rules, “Rules 1,” extracting the values for each key piece of data in the document and storing in a standardized format.

LINE #	QTY.	ITEM #	DESCRIPTION	UNIT PRICE	COST
1	7	123	Green Balls 16" Diameter	\$20.00	\$140.00
2	6	122	Red Balls 12" Diameter	\$18.00	\$108.00
3	8	124	Blue Balls 24" Diameter	\$25.00	\$200.00
4	2	143	Orange Balls 6" Diameter	\$15.00	\$30.00
TOTAL QTY.				23	
TOTAL COST					\$478.00

Fig. 6. A templated scanned image for a purchase order.

Similarly, we might receive a scanned image of a purchase order that equally follows a template (like that in “Fig. 6”), but must be read using OCR. Processing the file through a basic OCR reader could provide us with a data structure (see snippets of the type of data an OCR engine can produce in Appendix C), which can then be organized through additional scripts, though a bit more cumbersome to handle. Assuming all scanned images received are provided in this same layout, we can still treat the documents received as structured data.

2) *Semi-Structured Data*: If the same layout template of the structured data isn’t followed, the changes or variations require new rules. New designs, layouts, or columns in an item table might be a result of different regulations for a region requiring additional table columns in a purchase order, or larger customers using different systems that generate the files being provided.

Documents that contain the same type of information but in different designs and layouts fall into the category of semi-structured data. The keys in key-value pairs are the same across all documents, but the position on the page might vary, as seen in “Fig. 7”. Similarly, we might see scanned images like “Fig. 8” with varied layouts.

In this case, our rule-based approach gives us good accuracy for a document with one layout (ex: “Layout 1”), but low accuracy for a document with another layout (ex: “Layout 2”), like in “Fig. 9”. Appendix D provides additional code that can be applied to the same methods used in Appendix A that will work for the spreadsheet in “Fig. 7”, providing an example of

PURCHASE ORDER 1234567		SHIP TO XYZ Back Door 789 Main Ave Anytown, USA		SHIP DATE 09/01/21	
DESC	PRICE	QTY	U/M	ITEM NUM	COST
Green Balls 16" Diameter	\$20.00	7	CASE	123	\$140.00
Red Balls 12" Diameter	\$18.00	6	CASE	122	\$108.00
Blue Balls 24" Diameter	\$25.00	8	CASE	124	\$200.00
Orange Balls 6" Diameter	\$15.00	2	CASE	143	\$30.00
Total Quantity		23		Total Cost	\$478.00

Some fine print here. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Integer at enim tempor eros aliquet euismod et in felis. Duis id euismod magna. Lorem ipsum dolor sit amet, consectetur adipiscing elit.

Nullam eget suscipit orci, sit amet tristique justo. Maecenas aliquam consectetur tellus ac fermentum. Nulla varius at dolor eget suscipit. Proin metus metus, condimentum eget diam ut, tristique dictum nisi.

Fig. 7. A templated spreadsheet for a purchase order containing similar information to “Fig. 4”, but with a different layout.

COUNT	ITEM NUM	DESCRIPTION	UNIT PRICE	COST
7	123	Green Balls 16" Diameter	\$20.00	\$140.00
6	122	Red Balls 12" Diameter	\$18.00	\$108.00
8	124	Blue Balls 24" Diameter	\$25.00	\$200.00
2	143	Orange Balls 6" Diameter	\$15.00	\$30.00
TOTAL COUNT:			23	
TOTAL COST:				\$478.00

Fig. 8. A templated scanned image for a purchase order containing similar information to “Fig. 6”, but with a different layout.

creating rules for each layout.

Writing a unique rule for each type is manageable when there are only a few variations in the layout templates. This approach of creating rules for every layout variation will have difficulty scaling for large volumes of semi-structured data (see “Fig. 10”). Companies will have to devise a template creation and management process, adding and adjusting scripts and rules with each change and new layout.

Introducing artificial intelligence (AI) for growing volumes of semi-structured data can make it easier to scale at this stage of document classification and data extraction, avoiding template management. The field of AI combines computer science and hearty datasets to automate problem-solving [4]. Machine learning (a specific form of AI) uses models to automatically detect desired information from input documents. Machine learning (ML) applies algorithms to “learn” and improves accuracy in its program by receiving more samples to train on and improve its models. Vihar Kurama suggests the use of several pre-trained ML models (FastRCNN, Attention

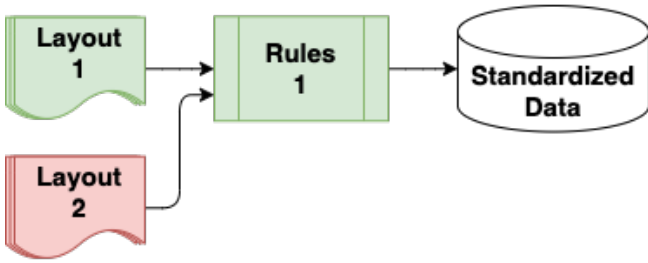


Fig. 9. Documents with “Layout 1” applied to “Rules 1” will extract values with good accuracy, but documents with “Layout 2” applied to the same rules will extract values with low accuracy.

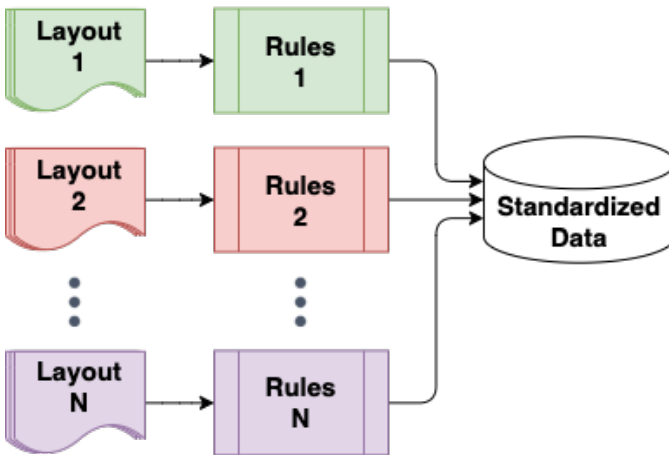


Fig. 10. For a large number of layouts, rule and template management becomes a job in itself and does not scale well.

OCR, and Graph Convolution) for a hybrid solution of both rules-based data extraction and ML when working with semi-structured data such as our examples [5].

Important to note is that the reliance on automated extraction should always be kept in check by a regular human review of metrics, accuracy, and confidence scores. Thankfully the processing of more samples to a machine learning model results in higher accuracy of the data retrieved.

3) *Unstructured Data*: The final classification for documents received would be unstructured data. This document type has no key-value pairs, may contain handwritten text, is free-flowing, and lacks consistency. Likely there are no tables (or the tables lack column headers) or a prescriptive way to automate the identification of which values to assign to which purposes. The message in “Fig. 11” is an example of unstructured information.

Scripts and ML are not enough to understand and classify this data. Instead, we may need to consider other techniques, such as natural language processing (NLP). Another branch of artificial intelligence, NLP provides computers with the ability to understand and translate the human language text of

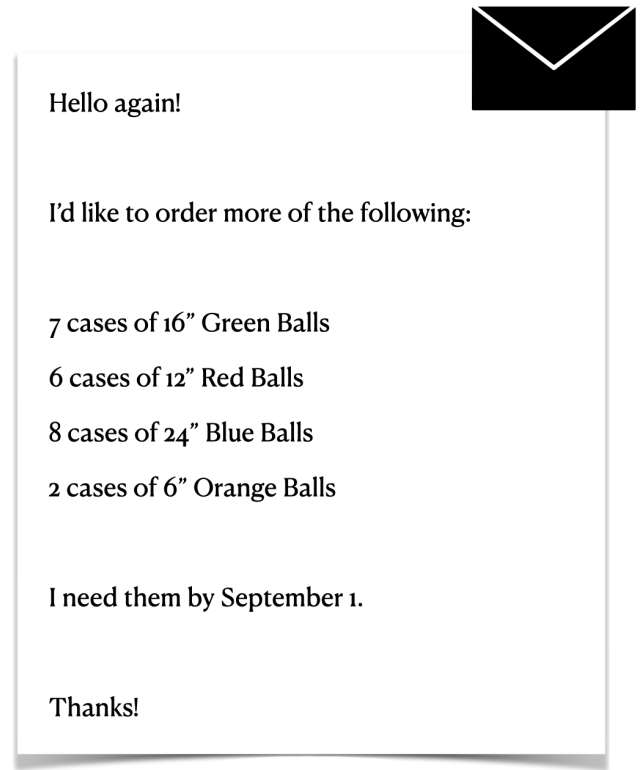


Fig. 11. Messages like those received in emails are easy to understand from a human, but need to be translated to something a computer can understand using natural language processing (NLP).

a document into something the machine can make decisions from [6].

C. Digitize & Standardize Data

With the process of understanding and classifying our data as structured, semi-structured, or unstructured, we now have some digital information that we can put to use. Our first pass at the documents might not have been deep enough to understand the document type; our business process may be more efficient at this point to gather just enough information to determine where to feed the document. This may mean that we’re capturing a page header, footer, or image with specific information to determine what type of document we’re working with.



Fig. 12. Capturing information from the document header may help us categorize the document by type.

D. Understand Document Types

With the standardized data from the previous step, we can determine if the document is an invoice, a purchase order, an application form, and so on. Documents that contain a “purchase order” header like in “Fig. 12” would get routed to the process for the ordering system, those with invoice headers might be sent to the Accounts Payable team, and application forms could be routed to the appropriate department, etc.

As a human, understanding the document types only takes a glance. If we consider automating this step for our pipeline, we need enough information to make the decision. If our data is structured, a rules-based system makes this very easy. Semi-structured and unstructured data can make this more difficult, especially if many document types are ingested into the same system.

Depending on the complexity of the business process involved, we may consider introducing machine learning or natural language processing in this step. Or, we improve our ingestion step to separate document types by directing users to submit different types of documents to the appropriate email address, shared folder, web portal, etc. The latter option may have some friction at the beginning but could be a very simple condition to overcome with a little customer training.

E. Better Data Extraction

Once we know our document type (invoice vs. purchase order, etc.) we can now apply more specific data extraction methods. This may be more heavily focused on machine learning models if the documents are semi-structured or natural language processing if they are unstructured.

Because we’ll assume that our data is semi-structured for our ball factory, this stage is a good placement for intelligent OCR, that is, OCR that involves ML or some form of AI. Model-based OCR engines can improve the process of digitizing the documents at scale while continuing to reduce errors with every training.

F. Human & Automated Validation

Nearing the end of our process flow, we’ve now reached the point where our data has been ingested and completely digitized into a standardized format. While it’s true that with machine learning we may need fewer human checks as we train with more samples, a certain level of human-in-the-loop is necessary (and a good idea) when automating document handling. Confirming accurate information has been gathered from the ingested documents is vital for clean data and correct orders being created.

Thankfully, the continuous feed of more samples over time improves the ML model and builds confidence in the process. Regardless, there should always be a human check to verify incoming information.

Beyond the human check, we can also add automatic validation by cross-referencing the data we’ve gathered against information in our system(s). Do the item numbers and customers exist in our system? Are the items in the order out of stock? If using machine learning or OCR, we can also flag low-confidence scores from the data extraction for human review.

V. SINGLE PIPELINE

All of this has helped us understand the problem and an approach to handling documents. We have an idea of which types of tools might be useful and each particular step, but now we must understand how to string it all together into a single, cohesive pipeline, since finding a single solution is our goal!

A business process might begin as a simple macro or script, improving the time and effort needed for small tasks. From there, the process may see continued improvements through IT or business solutions. At some point, additional automation can be considered to raise the bar to another level: robotic process automation.

A. Robotic Process Automation

Robotic process automation, more commonly referred to as RPA, is the method of automating repetitive, often administrative, tasks. Generally speaking, virtual workers created for RPA will take advantage of a user interface on a screen, much like a human worker would do. These digital assistants will capture data and manipulate applications based on programmed decision-making. Similar to writing a programming script, RPA bots are made to mimic the same steps a human would do in particularly repetitive tasks, freeing up the human workers to focus on more valuable tasks.

While a powerful tool, the category of RPA is still relatively new in the technology field. In a survey of 100 document

The image shows a sample purchase order form. Red annotations highlight specific areas: a red circle around the shipping address, a red circle around the item details table, and a red circle around the total cost. Red arrows point from text labels to these circles.

Shipping Address

Item Details

Total Cost

QUANTITY	ITEM NUM	DESCRIPTION	UNIT PRICE	TOTAL
7	123	Green Balls 16" Diameter	\$20.00	\$140.00
5	122	Red Balls 12" Diameter	\$18.00	\$108.00
8	124	Blue Balls 24" Diameter	\$25.00	\$200.00
2	143	Orange Balls 6" Diameter	\$15.00	\$30.00

TOTAL COUNT: 23
TOTAL COST: \$478.00

Fig. 13. Knowing our document type, we can do better data extraction. If our document type is a purchase order, we may know to extract the shipping address, item details, and total cost.

automation leaders by the Association of Intelligent Information Management, 18% of the respondents use RPA and 12% weren't sure of what RPA was. RPA with advanced capture has a place in the plans of 25% of the respondents, but 42% don't plan to make use of RPA. [7]

Beyond standard RPA, where we simply automate the flow itself, we can level up our process some more by using intelligent RPA, which leverages forms of artificial intelligence, as we've already discussed, for document handling.

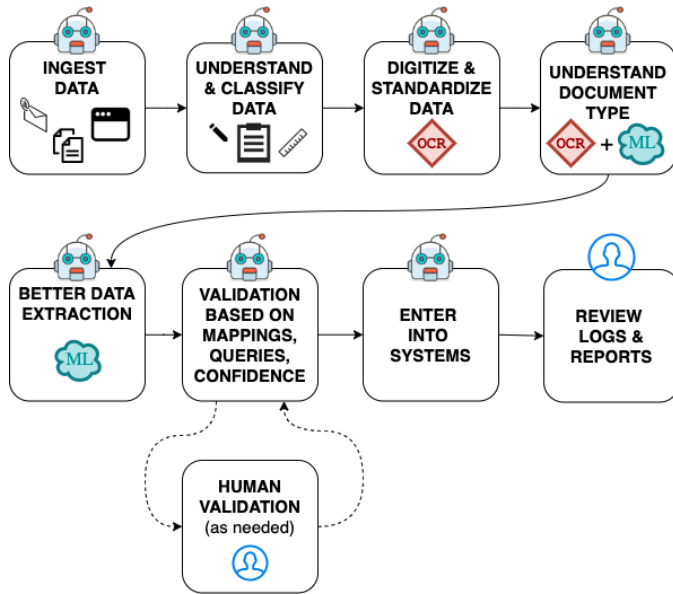


Fig. 14. Document Flow in a Business Process with RPA Bot and Human Worker ownership for each step.

The flow in “Fig. 14” is similar to the business process we began within “Fig. 3.” However, we can see that with the use of RPA, we can assign many of previous stages to a virtual worker; a human does not need to move the order document from one stage to the next. With the rules and AI in place, our digital assistant can not only ingest the data, but also classify, organize by type, standardize the data, validate it by cross-checking against existing table mappings (and alerting a human employee if further inspection is needed), but then also enter the data into the appropriate system (perhaps an ERP). Legacy systems may have no convenient integrations; RPA can interact with the user interface just like a human employee would, faster and with the cost-saving accuracy that a business requires.

This intelligent enhancement to a previously cumbersome and very manual flow frees up the human employees to turn their focus back to what matters: the value of customer relationships and understanding their needs to grow sales.

B. Case Study: PepsiCo

PepsiCo faced the challenge of a document processing flow that was slow, prone to errors, and very manual. When they applied a single process flow and involved the ABBYY FlexiCapture solution to their documents, they found that

they had lower error rates, reduced processing times, and their Accounts Payable staff were able to work much more efficiently [8].

PepsiCo’s new flow is as follows:

- 1) Paper invoices (in any of 5 languages) are scanned locally
- 2) Scanned invoices are...
 - a) sent to the POWER Project
 - b) identified with the correct entity
 - c) recognized by FlexiCapture
- 3) New digital documents and data are sent to verification stations
- 4) Data is verified
- 5) Data is forwarded to PepsiCo’s Image Vision for approval
- 6) Data is exported in XML
- 7) Data is processed in PepsiCo’s SAP ERP solution

“FlexiCapture’s accuracy is also proving stable under a heavy workload. In its first three months the POWER Project solution has processed two thousand batches, over 21,000 documents and nearly 40,000 pages without issues – in a mix of five languages.” [8]

In a review done by PricewaterhouseCoopers (PwC), ABBYY FlexiCapture was rated best in class in data validation and data classification, as well as data extraction in both structured and semi-structured documents. It also ranked well for multi-lingual capabilities, and the ability to read barcodes and tables. Google Tesseract and Microsoft’s Modi ranked better for screenscraping in desktop, web, and documents. PwC goes on to note that “the integration of the ABBYY OCR engine not only enhances automation for rules-driven processes, but also adds the flavour of NLP and widens the scope of automation.” [9]

This single pipeline solution has proven that it can scale and scale well! Powerful tools that fit the specific process need, like ABBYY FlexiCapture, combined with the power of automation, enabled PepsiCo to not only improve the intake process for documents at one location but to create a solution that accurately processes the information and can be scaled across the entire enterprise.

C. Case Study: Sysco

With a similar need and desire to improve the business processes, Sysco set out to also improve the accuracy and efficiency of many tasks across the company. By creating a team focused on RPA solutions, Sysco sought to create a new culture around automation enterprise-wide.

The Sysco RPA team has provided more than 200,000 hours of productivity given back to the business and has handled over 5.1 million transactions over the first 3.5 years of the team’s existence. Over 150 automations have been created and span across almost all business units. Their expansion looks to include specialty companies outside of the United States and includes a growing pipeline of over 200 business processes that are automation candidates. [10]

With at least one automation that involved OCR tools used to parse digital faxes, the Sysco team explored several other solutions varying from simple rules-based automations to more complicated and impactful processes [11]. Interestingly, this team not only standardizes and automates processes for other business units but also created a standardized business process and pod structure to improve their own efficiency and accuracy in determining which automations to create next and in the bot creation itself.

D. General Application

While the case studies presented above provide compelling narratives of the power of RPA in combination with OCR and AI, it may be difficult to understand where to begin. By identifying business cases that are repetitive and have a volume of effort that could be reduced through the use of scripts or automations, we have a place to start.

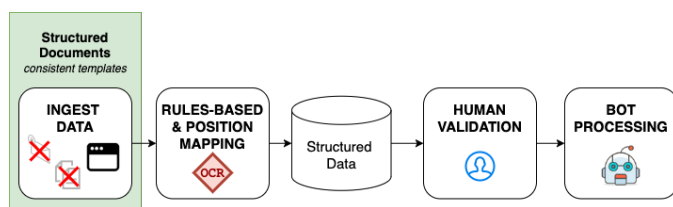


Fig. 15. Starting with an automated flow for only structured documents allows for full understanding of the process and gets employees familiar with using an RPA virtual assistant. Limiting the input to a single layout also eliminates needing to classify the data and determine the document type.

Businesses may experience a temptation to view automation as an all-or-nothing scenario. By starting with a small section of automation first, like “Fig. 15,” we can provide benefits quite quickly. Automating small steps in a process provides an opportunity to work out issues (not only bugs in the programming but also finding flaws or forgotten exceptions in the business process itself). When chosen correctly, automating one step in a business process can deliver a return on investment as soon as the feature is being adopted by users. This smaller section also helps us to understand nuances and challenges before we move on to more complicated stages.

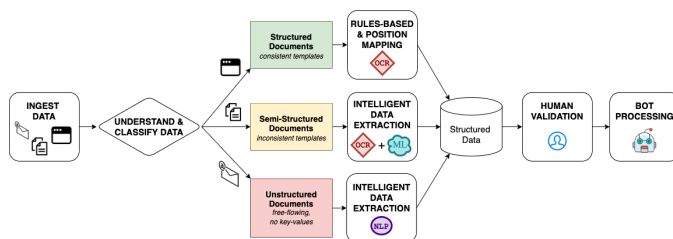


Fig. 16. By starting with one flow and adding in more entry points that lead to the same structured data, we can expand our automated process in an iterative way.

By each iteration of our process flow, we can provide new input formats that are handled by different rules and technologies (“Fig. 16”). We can add handling for structured data that is received through our web application, semi-structured

data received as files on a shared drive, and unstructured data captured in an email message. We can continue to expand for all data classifications while still leveraging the automated stages we built in the first round once the data is stored in a standardized format.

With many powerful solutions out there, advantages can be found in small steps first. Are employees typing in data from paper forms regularly? First, determine if the business process itself can change to an integrated solution. If it is not feasible to have the user change their input method within a reasonable timeframe, continue with the existing process and begin improving it by ensuring that those documents have been scanned into a system. Once digitally saved, run an OCR engine against the document, providing your employees with something they can at least copy-paste. Or, perhaps you already have the information in structured spreadsheets. Write a script that automatically captures the key-value pairs and table data, storing the information in a standardized format like a database. Once the process flow is working smoothly to the point of consistent and standardized information, consider adding some type of automated solution that can enter the data directly into the system, rather than an employee copying the information in.

By taking small, iterative steps, any company, large or small, can work its way into an automated flow. While some solutions like ABBYY FlexiCapture can provide top-of-the-line solutions, the cost may be prohibitive to some businesses. Start with the standard, rules-based methods, and enhance the areas that will still provide enough return-on-investment to justify the effort. The beautiful thing about the single-process pipeline idea is that more document types could be ingested over time, training the process with where to send each type. Some stages might continue to be human-handled, but there may be small improvements that allow the human to perform that step faster, like a dashboard that gives the employee a preview of the document. These types of process improvements might enable an employee to quickly glance and classify the information manually. Process improvement may save on the implementation of intelligent OCR or ML while still providing time-saving benefits.

VI. SOLUTIONS & COST

There are many tools out there and it can quickly become overwhelming to narrow them down. If you’re entering into the field of RPA, you may choose Blue Prism’s Decipher or Automation Hero’s solution. Perhaps you work with spreadsheets and gravitate towards Monarch. Microsoft has an entire suite of options with the Microsoft Power Apps. And we can’t forget the solutions that Amazon provides.

Each tool has its strengths and weaknesses, in addition to varying levels of cost. Understand the types of documents you’ll be working with and determine your needs. If you’ve started with automating structured documents, and are ready to include other classifications that require intelligent sorting, consider which models and training might be necessary.

I've heard the phrases, "the best exercise is the one you'll do," or "the best camera is the one in your hand." This idea can apply here as well: the best tool to improve your process is whichever one you start using. One technology may not have huge advantages over another. Some options may be tightly coupled to their own branded integrations. Others might have a nicer user interface. Becoming paralyzed by the number of solutions available is a possibility; give your research period a time-box, then choose whichever option(s) best suit the business needs and the skills of your technical team.

VII. CONCLUSION

Manual data entry is time-consuming and error-prone. The introduction of OCR (combined with ML) or NLP can result in fewer typos and opens the door to the improved process flow. Beyond that, considering RPA within the business process can create further efficiencies, giving skilled employees time back to focus on more valuable tasks that have a higher impact on the company's profitability.

Many well-tested solutions are available, with pre-trained models available to implement from the start, reducing the number of required samples to get a machine learning solution off the ground quickly. However, the strongest solutions can be cost-prohibitive. Not to be deterred, we can still provide improvements by automating smaller sections of the process flow or enhancing the human step (thus bypassing the need for intelligent automation) to make it simpler to perform.

While still a young category of technology, RPA provides further enhancements by not only stringing together the steps into a single process flow but also moving into working with legacy systems without the need to wait for modern integrations.

Depending on the business needs, more research can (and should) be done on which tools may provide the best solutions. A quick search on the web can provide appropriate overviews of the most popular tools. For example, Nanonets provides a list of pros and cons for the best OCR tools of 2021 [12]. Depending on the needs and time available (for development and potential training of models) for each business case, the right tool for the job may change.

Iterating through development is key for a quick return on investment and limiting rework. Starting with rule-based solutions can provide a solution that is faster to deliver and provides insights into nuances and challenges within the process flow. Expanding into RPA in the next iteration builds confidence and trust in the automation, freeing up more time for human employees to focus on valuable tasks. Entering into intelligent RPA opens the floodgates to the types of documents that can be handled automatically, with human review.

The problem at hand is that data can be received in many formats; manually extracting that information can be tedious at best and riddled with errors at its worst. Is there a single-process solution that businesses can implement to improve this process? Based on what we've covered, the answer is yes and no. Each business process has some unique nuances that require specific solutions. Some very powerful options

are out there that provide less training and setup. However, all companies can start with small solutions (even improving the business process itself!), then iterate and build on those solutions creating a robust, single-process built of many small steps.

REFERENCES

- [1] Matt Harris. When good info goes bad: The real cost of human data errors, 2014. <https://ungerboeck.com/resources/when-good-info-goes-bad-the-real-cost-of-human-data-errors> [Online; accessed 15-August2021].
- [2] James Schneider Ph.D., Bill Schultz, Julia McCrimlisk, Jesse Hulsing, and Ryan M. Nash CFA. B2B: How the next payments frontier will unleash small business. EQUITY RESEARCH, 2018. https://known-production.s3.amazonaws.com/uploads/attachment/file/3460/Global%2BTechnology_%2BB2B_%2BHow%2Bthe%2Bnext%2Bpayments%2Bfrontier%2Bwill%2Bunleash%2Bsmall%2Bbusiness%2B.pdf [Online; accessed 13-August2021].
- [3] Justin Lavelle. Gartner says robotic process automation can save finance departments 25,000 hours of avoidable work annually, 2019. <https://www.gartner.com/en/newsroom/press-releases/2019-10-02-gartner-says-robotic-process-automation-can-save-fina> [Online; accessed 13-August2021].
- [4] IBM Cloud Education. Artificial intelligence (AI), 2020. <https://www.ibm.com/cloud/learn/what-is-artificial-intelligence> [Online; accessed 15-August2021].
- [5] Vihar Kurama. A comprehensive guide to OCR with RPA and document understanding, 2021. <https://nanonets.com/blog/ocr-with-rpa-and-document-understanding-uipath/> [Online; accessed 13-August2021].
- [6] IBM Cloud Education. Natural language processing (NLP), 2020. <https://www.ibm.com/cloud/learn/natural-language-processing> [Online; accessed 14-August2021].
- [7] Garrett Hollander. [survey results] intelligent document automation trends, 2019. <https://resources.m-files.com/blog/survey-results-intelligent-document-automation-trends-3> [Online; accessed 13-August2021].
- [8] PepsiCo automates invoice processing with ABBYY FlexiCapture, 2021. <https://www.abbyy.com/customer-stories/pepsico-automates-invoice-processing-with-abbyy-flexicapture/> [Online; accessed 13-August2021].
- [9] PwC. Robotic process automation and intelligent character recognition: Smart data capture. PricewaterhouseCoopers Private Limited, 2018. <http://www.pwc.in> [Online; accessed 13-August2021].
- [10] Alison Major and Kim Meredith. Sysco: Scaling an rpa program without compromise, 2021. https://www.brighttalk.com/resource/core/349721/blue-prism-cafe-webinar-with-sysco_764477.pdf [Online; accessed 15-August2021].
- [11] Blue Prism Café. Sysco: Scaling an intelligent automation program without compromise, 2021. <https://www.brighttalk.com/webcast/17260/476962> [Webinar; accessed 15-August2021].
- [12] Prithiv S. Best ocr software of 2021, 2021. <https://nanonets.com/blog/ocr-software-best-ocr-software/> [Online; accessed 15-August2021].
- [13] OCR.SPACE is a service of a9t9 software GmbH. OCRSpace, 2021. <https://ocr.space/> [Online; accessed 8-August2021].

APPENDIX A
PYTHON METHODS TO COLLECT DATA FROM A SPREADSHEET

A. *spreadsheets/methods.py*

```
1 import string
2
3
4 def get_po_number(workbook, cells):
5     """Returns the value for the Purchase Order Number as a string"""
6     value = get_cell_value(workbook, cells["poNum"])
7     value_as_int = int(value)
8     value_as_string = str(value_as_int)
9     return value_as_string
10
11
12 def get_bill_to_name(workbook, cells):
13     """Returns the value for the Bill-To Name"""
14     return get_cell_value(workbook, cells["billTo"])
15
16
17 def get_ship_to_name(workbook, cells):
18     """Returns the value for the Ship-To Name (only first line)"""
19     value = get_cell_value(workbook, cells["shipTo"])
20     value_after_first_line = value.split("\n")
21     return value_after_first_line[0]
22
23
24 def get_ship_date(workbook, cells):
25     """Returns the value for the Shipping Date"""
26     return get_cell_value(workbook, cells["shipDate"])
27
28
29 def get_total_quantity(workbook, cells):
30     """Returns the value for the Total Quantity as an integer"""
31     value = get_cell_value(workbook, cells["totalQty"])
32     value_as_int = int(value)
33     return value_as_int
34
35
36 def get_total_cost(workbook, cells):
37     """Returns the value for the Total Cost"""
38     return get_cell_value(workbook, cells["totalCost"])
39
40
41 def get_cell_value(workbook, cell):
42     """Gets the value of the cell (ex: 'B1') from a workbook"""
43     if cell == "":
44         """No cell was provided."""
45         return "[cell not mapped]"
46
47     column = get_column_value(cell)
48     row = get_row_value(cell)
49     return workbook.loc[row].at[column]
50
51
52 def get_column_value(cell):
53     """Gets the 'loc' value from a cell"""
54     alphabet = dict()
55     for index, letter in enumerate(string.ascii_lowercase):
56         alphabet[letter] = index
57
58     value = cell[0].lower()
59     return alphabet[value]
60
61
62 def get_row_value(cell):
```

```

63     """Gets the 'at' value from a cell"""
64     value = cell[1:]
65     return int(value) - 1
66
67
68 def get_item_line_number(item_table, item_row):
69     """Gets the item's line number from the table"""
70     try:
71         return item_table.at[item_row, 'LINE_NUM']
72     except KeyError:
73         return "[not available]"
74
75
76 def get_item_quantity(item_table, item_row):
77     """Gets the item's quantity from the table"""
78     try:
79         return item_table.at[item_row, 'QTY']
80     except KeyError:
81         return "[not available]"
82
83
84 def get_item_unit_of_measure(item_table, item_row):
85     """Gets the item's unit of measure from the table"""
86     try:
87         return item_table.at[item_row, 'UOM']
88     except KeyError:
89         return "[not available]"
90
91
92 def get_item_number(item_table, item_row):
93     """Gets the item's identification number from the table"""
94     try:
95         value = item_table.at[item_row, 'ITEM_NUM']
96         value_as_int = int(value)
97         return str(value_as_int)
98     except KeyError:
99         return "[not available]"
100
101
102 def get_item_description(item_table, item_row):
103     """Gets the item's description from the table"""
104     try:
105         return item_table.at[item_row, 'DESC']
106     except KeyError:
107         return "[not available]"
108
109
110 def get_item_unit_price(item_table, item_row):
111     """Gets the item's unit price from the table"""
112     try:
113         return item_table.at[item_row, 'UNIT_PRICE']
114     except KeyError:
115         return "[not available]"
116
117
118 def get_item_cost(item_table, item_row):
119     """Gets the item's total order cost from the table"""
120     try:
121         return item_table.at[item_row, 'COST']
122     except KeyError:
123         return "[not available]"

```

B. *spreadsheets/excel_order_one.py*

These mappings are specific to the spreadsheet represented in “Fig. 4”.

```

1 import pandas as pd
2

```

```

3 FILE_PATH = '../documents/ORDERS.xlsx'
4 SHEET_NAME = 'Order Sample 1'
5
6 """The workbook to pull values from"""
7 WORKBOOK = pd.read_excel(
8     FILE_PATH,
9     sheet_name=SHEET_NAME,
10    header=None
11 )
12
13 """The table for item information"""
14 ITEM_TABLE = pd.read_excel(
15     FILE_PATH,
16     sheet_name=SHEET_NAME,
17     usecols="A:G",
18     skiprows=5,
19     nrows=4
20 )
21 """Mapping column names for item table"""
22 ITEM_TABLE.columns = [
23     'LINE_NUM',      # Line #
24     'QTY',           # Quantity
25     'UOM',           # Unit of Measure
26     'ITEM_NUM',      # Item #
27     'DESC',          # Description
28     'UNIT_PRICE',    # Unit Price
29     'COST'           # Charge
30 ]
31
32 """Mapping of the values to cells"""
33 ORDER_ONE_CELLS = {
34     "poNum": "B1",
35     "billTo": "B4",
36     "shipTo": "B3",
37     "shipDate": "B2",
38     "totalQty": "B12",
39     "totalCost": "G12"
40 }

```


APPENDIX B
UNIT TESTS FOR PYTHON METHODS IN APPENDIX A

A. *tests/test_spreadsheets_methods.py*

```
1 import datetime
2 import unittest
3 import pandas as pd
4 from parameterized import parameterized
5 from spreadsheets.methods import get_column_value, get_row_value, get_po_number, get_total_cost, get_to
6     get_ship_date, get_ship_to_name, get_bill_to_name, get_cell_value, get_item_line_number, get_item_q
7     get_item_unit_of_measure, get_item_number, get_item_description, get_item_unit_price, get_item_cost
8
9 FILE_PATH = '../documents/ORDERS.xlsx'
10 SHEET_NAME = 'Order Sample 1'
11
12 """The workbook to pull values from"""
13 WORKBOOK = pd.read_excel(
14     FILE_PATH,
15     sheet_name=SHEET_NAME,
16     header=None
17 )
18
19 """Mapping of the values to cells"""
20 CELLS = {
21     "poNum": "B1",
22     "billTo": "B4",
23     "shipTo": "B3",
24     "shipDate": "B2",
25     "totalQty": "B12",
26     "totalCost": "G12"
27 }
28
29 """Fake item table for testing"""
30 ITEM_TABLE = pd.DataFrame(
31     [
32         [1, 7, 'CASE', '123', 'Item 1', 20, 140],
33         [2, 6, 'PALLET', 124.0, 'Item 2', 18, 108]
34     ],
35     columns=[
36         'LINE_NUM',
37         'QTY',
38         'UOM',
39         'ITEM_NUM',
40         'DESC',
41         'UNIT_PRICE',
42         'COST'
43     ]
44 )
45
46
47 class SpreadsheetsMethods(unittest.TestCase):
48     @parameterized.expand([
49         ["B1 gives 1", "B1", 1],
50         ["B2 gives 1", "B2", 1],
51         ["A1 gives 0", "A1", 0],
52         ["Z9 gives 25", "Z9", 25]
53     ])
54     def test_get_column_value(self, name, cell, expected):
55         actual = get_column_value(cell)
56         self.assertEqual(expected, actual)
57
58     @parameterized.expand([
59         ["B1 gives 1", "B1", 0],
60         ["B2 gives 1", "B2", 1],
61         ["A1 gives 0", "A1", 0],
62         ["Z9 gives 25", "Z9", 8],
```

```

63         ["F13 gives 12", "F13", 12],
64         ["D101 gives 100", "D101", 100]
65     ])
66     def test_get_row_value(self, name, cell, expected):
67         actual = get_row_value(cell)
68         self.assertEqual(expected, actual)
69
70     @parameterized.expand([
71         ["A1 gives 'PO#'", "A1", "PO#"],
72         ["blank gives '[cell not mapped]'", "", "[cell not mapped]"]
73     ])
74     def test_get_cell_value(self, name, cell, expected):
75         actual = get_cell_value(WORKBOOK, cell)
76         self.assertEqual(expected, actual)
77
78     def test_get_po_number(self):
79         expected = "1234567"
80         actual = get_po_number(WORKBOOK, CELLS)
81         self.assertEqual(expected, actual)
82
83     def test_get_bill_to_name(self):
84         expected = "XYZ Company"
85         actual = get_bill_to_name(WORKBOOK, CELLS)
86         self.assertEqual(expected, actual)
87
88     def test_get_ship_to_name(self):
89         expected = "XYZ Back Door"
90         actual = get_ship_to_name(WORKBOOK, CELLS)
91         self.assertEqual(expected, actual)
92
93     def test_get_ship_date(self):
94         expected = datetime.datetime(2021, 9, 1, 0, 0) # 09/01/21
95         actual = get_ship_date(WORKBOOK, CELLS)
96         self.assertEqual(expected, actual)
97
98     def test_get_total_quantity(self):
99         expected = 23
100         actual = get_total_quantity(WORKBOOK, CELLS)
101         self.assertEqual(expected, actual)
102
103     def test_get_total_cost(self):
104         expected = 478.00
105         actual = get_total_cost(WORKBOOK, CELLS)
106         self.assertEqual(expected, actual)
107
108     @parameterized.expand([
109         ["Row 1 gives line '1'", ITEM_TABLE, 0, 1],
110         ["Row 2 gives line '2'", ITEM_TABLE, 1, 2],
111         ["Row 3 gives '[not available]'", ITEM_TABLE, 2, '[not available]']
112     ])
113     def test_get_item_line_number(self, name, table, row, expected):
114         actual = get_item_line_number(table, row)
115         self.assertEqual(expected, actual)
116
117     @parameterized.expand([
118         ["Row 1 gives quantity 7", ITEM_TABLE, 0, 7],
119         ["Row 2 gives quantity 6", ITEM_TABLE, 1, 6],
120         ["Row 3 gives '[not available]'", ITEM_TABLE, 2, '[not available]']
121     ])
122     def test_get_item_quantity(self, name, table, row, expected):
123         actual = get_item_quantity(table, row)
124         self.assertEqual(expected, actual)
125
126     @parameterized.expand([
127         ["Row 1 gives 'CASE'", ITEM_TABLE, 0, 'CASE'],
128         ["Row 2 gives 'PALLET'", ITEM_TABLE, 1, 'PALLET'],
129         ["Row 3 gives '[not available]'", ITEM_TABLE, 2, '[not available]']

```

```

130     ])
131     def test_get_item_unit_of_measure(self, name, table, row, expected):
132         actual = get_item_unit_of_measure(table, row)
133         self.assertEqual(expected, actual)
134
135     @parameterized.expand([
136         ["Row 1 gives item number '123'", ITEM_TABLE, 0, '123'],
137         ["Row 2 gives item number '124'", ITEM_TABLE, 1, '124'],
138         ["Row 3 gives '[not available]'", ITEM_TABLE, 2, '[not available]']
139     ])
140     def test_get_item_number(self, name, table, row, expected):
141         actual = get_item_number(table, row)
142         self.assertEqual(expected, actual)
143
144     @parameterized.expand([
145         ["Row 1 gives 'Item 1'", ITEM_TABLE, 0, 'Item 1'],
146         ["Row 2 gives 'Item 2'", ITEM_TABLE, 1, 'Item 2'],
147         ["Row 3 gives '[not available]'", ITEM_TABLE, 2, '[not available]']
148     ])
149     def test_get_item_description(self, name, table, row, expected):
150         actual = get_item_description(table, row)
151         self.assertEqual(expected, actual)
152
153     @parameterized.expand([
154         ["Row 1 gives 20.00", ITEM_TABLE, 0, 20.00],
155         ["Row 2 gives 18.00", ITEM_TABLE, 1, 18.00],
156         ["Row 3 gives '[not available]'", ITEM_TABLE, 2, '[not available]']
157     ])
158     def test_get_item_unit_price(self, name, table, row, expected):
159         actual = get_item_unit_price(table, row)
160         self.assertEqual(expected, actual)
161
162     @parameterized.expand([
163         ["Row 1 gives 140.00", ITEM_TABLE, 0, 140.00],
164         ["Row 2 gives 108.00", ITEM_TABLE, 1, 108.00],
165         ["Row 3 gives '[not available]'", ITEM_TABLE, 2, '[not available]']
166     ])
167     def test_get_item_cost(self, name, table, row, expected):
168         actual = get_item_cost(table, row)
169         self.assertEqual(expected, actual)
170
171
172 if __name__ == '__main__':
173     unittest.main()

```

B. tests/test_excel_order_one.py

```
1 import datetime
2 import unittest
3 from parameterized import parameterized
4 from spreadsheets.excel_order_one import WORKBOOK, ORDER_ONE_CELLS, ITEM_TABLE
5 from spreadsheets.methods import get_po_number, get_bill_to_name, get_ship_to_name, get_ship_date, get_
6     get_total_cost, get_item_line_number, get_item_quantity, get_item_unit_of_measure, get_item_number,
7     get_item_description, get_item_unit_price, get_item_cost
8
9
10 class ExcelOrderOne(unittest.TestCase):
11     def test_get_po_number(self):
12         expected = "1234567"
13         actual = get_po_number(WORKBOOK, ORDER_ONE_CELLS)
14         self.assertEqual(expected, actual)
15
16     def test_get_bill_to_name(self):
17         expected = "XYZ Company"
18         actual = get_bill_to_name(WORKBOOK, ORDER_ONE_CELLS)
19         self.assertEqual(expected, actual)
20
21     def test_get_ship_to_name(self):
22         expected = "XYZ Back Door"
23         actual = get_ship_to_name(WORKBOOK, ORDER_ONE_CELLS)
24         self.assertEqual(expected, actual)
25
26     def test_get_ship_date(self):
27         expected = datetime.datetime(2021, 9, 1, 0, 0) # 09/01/21
28         actual = get_ship_date(WORKBOOK, ORDER_ONE_CELLS)
29         self.assertEqual(expected, actual)
30
31     def test_get_total_quantity(self):
32         expected = 23
33         actual = get_total_quantity(WORKBOOK, ORDER_ONE_CELLS)
34         self.assertEqual(expected, actual)
35
36     def test_get_total_cost(self):
37         expected = 478.00
38         actual = get_total_cost(WORKBOOK, ORDER_ONE_CELLS)
39         self.assertEqual(expected, actual)
40
41     @parameterized.expand([
42         ["Row 1 gives line '1'", ITEM_TABLE, 0, 1],
43         ["Row 2 gives line '2'", ITEM_TABLE, 1, 2],
44         ["Row 3 gives line '3'", ITEM_TABLE, 2, 3],
45         ["Row 4 gives line '4'", ITEM_TABLE, 3, 4]
46     ])
47     def test_get_item_line_number(self, name, table, row, expected):
48         actual = get_item_line_number(table, row)
49         self.assertEqual(expected, actual)
50
51     @parameterized.expand([
52         ["Row 1 gives quantity 7", ITEM_TABLE, 0, 7],
53         ["Row 2 gives quantity 6", ITEM_TABLE, 1, 6],
54         ["Row 3 gives quantity 8", ITEM_TABLE, 2, 8],
55         ["Row 4 gives quantity 2", ITEM_TABLE, 3, 2]
56     ])
57     def test_get_item_quantity(self, name, table, row, expected):
58         actual = get_item_quantity(table, row)
59         self.assertEqual(expected, actual)
60
61     @parameterized.expand([
62         ["Row 1 gives 'CASE'", ITEM_TABLE, 0, 'CASE'],
63         ["Row 2 gives 'CASE'", ITEM_TABLE, 1, 'CASE'],
64         ["Row 3 gives 'CASE'", ITEM_TABLE, 2, 'CASE'],
65         ["Row 4 gives 'CASE'", ITEM_TABLE, 3, 'CASE']
```



```

66     ])
67     def test_get_item_unit_of_measure(self, name, table, row, expected):
68         actual = get_item_unit_of_measure(table, row)
69         self.assertEqual(expected, actual)
70
71     @parameterized.expand([
72         ["Row 1 gives item number '123'", ITEM_TABLE, 0, '123'],
73         ["Row 2 gives item number '124'", ITEM_TABLE, 1, '122'],
74         ["Row 3 gives item number '123'", ITEM_TABLE, 2, '124'],
75         ["Row 4 gives item number '123'", ITEM_TABLE, 3, '143']
76     ])
77     def test_get_item_number(self, name, table, row, expected):
78         actual = get_item_number(table, row)
79         self.assertEqual(expected, actual)
80
81     @parameterized.expand([
82         ["Row 1 gives the description", ITEM_TABLE, 0, 'Green Balls\n16" Diameter'],
83         ["Row 2 gives the description", ITEM_TABLE, 1, 'Red Balls\n12" Diameter'],
84         ["Row 3 gives the description", ITEM_TABLE, 2, 'Blue Balls\n24" Diameter'],
85         ["Row 4 gives the description", ITEM_TABLE, 3, 'Orange Balls\n6" Diameter']
86     ])
87     def test_get_item_description(self, name, table, row, expected):
88         actual = get_item_description(table, row)
89         self.assertEqual(expected, actual)
90
91     @parameterized.expand([
92         ["Row 1 gives 20.00", ITEM_TABLE, 0, 20.00],
93         ["Row 2 gives 18.00", ITEM_TABLE, 1, 18.00],
94         ["Row 3 gives 25.00", ITEM_TABLE, 2, 25.00],
95         ["Row 4 gives 15.00", ITEM_TABLE, 3, 15.00]
96     ])
97     def test_get_item_unit_price(self, name, table, row, expected):
98         actual = get_item_unit_price(table, row)
99         self.assertEqual(expected, actual)
100
101     @parameterized.expand([
102         ["Row 1 gives 140.00", ITEM_TABLE, 0, 140.00],
103         ["Row 2 gives 108.00", ITEM_TABLE, 1, 108.00],
104         ["Row 3 gives 200.00", ITEM_TABLE, 2, 200.00],
105         ["Row 4 gives 30.00", ITEM_TABLE, 3, 30.00]
106     ])
107     def test_get_item_cost(self, name, table, row, expected):
108         actual = get_item_cost(table, row)
109         self.assertEqual(expected, actual)
110
111
112 if __name__ == '__main__':
113     unittest.main()

```

APPENDIX C

DATA COLLECTED AS JSON FROM A SCANNED IMAGE

A full sample of the JSON for the scanned image from Order 1 (“Fig. 6”), generated by an OCR scanning tool (OCRSpace [13]), can be found stored on GitHub: https://github.com/amajor/research-in-cs-course-project/Scripts/documents/OCRSpace_Order1.json

A full sample of the JSON for the scanned image from Order 2 (“Fig. 8”), generated by an OCR scanning tool (OCRSpace [13]), can be found stored on GitHub: https://github.com/amajor/research-in-cs-course-project/Scripts/documents/OCRSpace_Order2.json

Below is a very small sampling of some lines at the beginning of a JSON file (lines 1-40) containing the data recognized from the scanned image (“Fig. 6”):

```
1 {
2   "ParsedResults": [
3     {
4       "Overlay": {
5         "Lines": [
6           {
7             "LineText": "INVOICE TO: xyz Company",
8             "Words": [
9               {
10                "WordText": "INVOICE",
11                "Left": 74.5,
12                "Top": 105,
13                "Height": 13,
14                "Width": 64.5
15              },
16              {
17                "WordText": "TO:",
18                "Left": 149.5,
19                "Top": 105.5,
20                "Height": 12.5,
21                "Width": 25.5
22              },
23              {
24                "WordText": "xyz",
25                "Left": 182,
26                "Top": 105.5,
27                "Height": 13.5,
28                "Width": 27
29              },
30              {
31                "WordText": "Company",
32                "Left": 219.5,
33                "Top": 105.5,
34                "Height": 17,
35                "Width": 66
36              }
37            ],
38            "MaxHeight": 17,
39            "MinTop": 105
40          },

```

APPENDIX D
PYTHON METHODS TO COLLECT DATA FROM A SPREADSHEET

A. *spreadsheets/excel_order_two.py*

These mappings are specific to the spreadsheet represented in “Fig. 7”.

```
1 import pandas as pd
2
3 FILE_PATH = '../documents/ORDERS.xlsx'
4 SHEET_NAME = 'Order Sample 2'
5
6 """The workbook to pull values from"""
7 WORKBOOK = pd.read_excel(
8     FILE_PATH,
9     sheet_name=SHEET_NAME,
10    header=None
11 )
12
13 """The table for item information"""
14 ITEM_TABLE = pd.read_excel(
15     FILE_PATH,
16     sheet_name=SHEET_NAME,
17     usecols="A:F",
18     skiprows=3,
19     nrows=4
20 )
21 """Mapping column names for item table"""
22 ITEM_TABLE.columns = [
23     'DESC',          # DESC
24     'UNIT_PRICE',    # PRICE
25     'QTY',           # QTY
26     'UOM',           # U/M
27     'ITEM_NUM',      # ITEM_NUM
28     'COST'           # COST
29 ]
30
31 """Mapping of the values to cells"""
32 ORDER_TWO_CELLS = {
33     "poNum": "A2",
34     "billTo": "",
35     "shipTo": "C2",
36     "shipDate": "F2",
37     "totalQty": "C10",
38     "totalCost": "F10"
39 }
```

APPENDIX E
UNIT TESTS FOR PYTHON MAPPINGS FOR SPREADSHEET ORDER TWO

A. *tests/test_excel_order_two.py*

```
1 import datetime
2 import unittest
3 from parameterized import parameterized
4 from spreadsheets.excel_order_two import WORKBOOK, ORDER_TWO_CELLS, ITEM_TABLE
5 from spreadsheets.methods import get_po_number, get_bill_to_name, get_ship_to_name, get_ship_date, get_
6     get_total_cost, get_item_line_number, get_item_quantity, get_item_unit_of_measure, get_item_number,
7     get_item_description, get_item_unit_price, get_item_cost
8
9
10 class ExcelOrderTwo(unittest.TestCase):
11     def test_get_po_number(self):
12         expected = "1234567"
13         actual = get_po_number(WORKBOOK, ORDER_TWO_CELLS)
14         self.assertEqual(expected, actual)
15
16     def test_get_bill_to_name(self):
17         expected = "[cell not mapped]"
18         actual = get_bill_to_name(WORKBOOK, ORDER_TWO_CELLS)
19         self.assertEqual(expected, actual)
20
21     def test_get_ship_to_name(self):
22         expected = "XYZ Back Door"
23         actual = get_ship_to_name(WORKBOOK, ORDER_TWO_CELLS)
24         self.assertEqual(expected, actual)
25
26     def test_get_ship_date(self):
27         expected = datetime.datetime(2021, 9, 1, 0, 0) # 09/01/21
28         actual = get_ship_date(WORKBOOK, ORDER_TWO_CELLS)
29         self.assertEqual(expected, actual)
30
31     def test_get_total_quantity(self):
32         expected = 23
33         actual = get_total_quantity(WORKBOOK, ORDER_TWO_CELLS)
34         self.assertEqual(expected, actual)
35
36     def test_get_total_cost(self):
37         expected = 478.00
38         actual = get_total_cost(WORKBOOK, ORDER_TWO_CELLS)
39         self.assertEqual(expected, actual)
40
41     @parameterized.expand([
42         ["Row 1 gives [not available]", ITEM_TABLE, 0, "[not available]"],
43         ["Row 2 gives [not available]", ITEM_TABLE, 1, "[not available]"],
44         ["Row 3 gives [not available]", ITEM_TABLE, 2, "[not available]"],
45         ["Row 4 gives [not available]", ITEM_TABLE, 3, "[not available]"]
46     ])
47     def test_get_item_line_number(self, name, table, row, expected):
48         actual = get_item_line_number(table, row)
49         self.assertEqual(expected, actual)
50
51     @parameterized.expand([
52         ["Row 1 gives quantity 7", ITEM_TABLE, 0, 7],
53         ["Row 2 gives quantity 6", ITEM_TABLE, 1, 6],
54         ["Row 3 gives quantity 8", ITEM_TABLE, 2, 8],
55         ["Row 4 gives quantity 2", ITEM_TABLE, 3, 2]
56     ])
57     def test_get_item_quantity(self, name, table, row, expected):
58         actual = get_item_quantity(table, row)
59         self.assertEqual(expected, actual)
60
61     @parameterized.expand([
62         ["Row 1 gives 'CASE'", ITEM_TABLE, 0, 'CASE'],
```



```

63         ["Row 2 gives 'CASE'", ITEM_TABLE, 1, 'CASE'],
64         ["Row 3 gives 'CASE'", ITEM_TABLE, 2, 'CASE'],
65         ["Row 4 gives 'CASE'", ITEM_TABLE, 3, 'CASE']
66     ])
67     def test_get_item_unit_of_measure(self, name, table, row, expected):
68         actual = get_item_unit_of_measure(table, row)
69         self.assertEqual(expected, actual)
70
71     @parameterized.expand([
72         ["Row 1 gives item number '123'", ITEM_TABLE, 0, '123'],
73         ["Row 2 gives item number '124'", ITEM_TABLE, 1, '122'],
74         ["Row 3 gives item number '123'", ITEM_TABLE, 2, '124'],
75         ["Row 4 gives item number '123'", ITEM_TABLE, 3, '143']
76     ])
77     def test_get_item_number(self, name, table, row, expected):
78         actual = get_item_number(table, row)
79         self.assertEqual(expected, actual)
80
81     @parameterized.expand([
82         ["Row 1 gives the description", ITEM_TABLE, 0, 'Green Balls\n16" Diameter'],
83         ["Row 2 gives the description", ITEM_TABLE, 1, 'Red Balls\n12" Diameter'],
84         ["Row 3 gives the description", ITEM_TABLE, 2, 'Blue Balls\n24" Diameter'],
85         ["Row 4 gives the description", ITEM_TABLE, 3, 'Orange Balls\n6" Diameter']
86     ])
87     def test_get_item_description(self, name, table, row, expected):
88         actual = get_item_description(table, row)
89         self.assertEqual(expected, actual)
90
91     @parameterized.expand([
92         ["Row 1 gives 20.00", ITEM_TABLE, 0, 20.00],
93         ["Row 2 gives 18.00", ITEM_TABLE, 1, 18.00],
94         ["Row 3 gives 25.00", ITEM_TABLE, 2, 25.00],
95         ["Row 4 gives 15.00", ITEM_TABLE, 3, 15.00]
96     ])
97     def test_get_item_unit_price(self, name, table, row, expected):
98         actual = get_item_unit_price(table, row)
99         self.assertEqual(expected, actual)
100
101     @parameterized.expand([
102         ["Row 1 gives 140.00", ITEM_TABLE, 0, 140.00],
103         ["Row 2 gives 108.00", ITEM_TABLE, 1, 108.00],
104         ["Row 3 gives 200.00", ITEM_TABLE, 2, 200.00],
105         ["Row 4 gives 30.00", ITEM_TABLE, 3, 30.00]
106     ])
107     def test_get_item_cost(self, name, table, row, expected):
108         actual = get_item_cost(table, row)
109         self.assertEqual(expected, actual)
110
111
112 if __name__ == '__main__':
113     unittest.main()

```