

# Automating information extraction from documents

July 4, 2019

*Nobody said it was easy*

*No one ever said it would be this hard*

Truer words haven't been spoken (or sang really) when it comes to extracting data from piles and piles of documents. As much as we like to believe otherwise, an enormous number of company transactions still involve paper work instead of digital means. Accountants and finance departments understand the pain and are desperately looking for an end-to-end solution which automates and digitalises the accounting workflows. Book-keepers spend long hours on manual data entry from invoices, receipts and other documents costing meaningless effort, money and time. Good grief! why are still struggling to find a reliable solution that can put an end to this everlasting drain of resources. The complexity of automating this problem revolves around the variability of documents' format and structure. An invoice or a receipt can be created in millions of different styles which makes it impossible to take a 'template' based approach or 'fingerprinting' to sound fancier. Adding to the misery, reading text via traditional OCR (optical character recognition) is enigmatic and prone to error because it largely depends on the picture quality, shadows, and perspective from where the picture was captured. Enter the new kid on the block, Staple AI, which uses computer vision and machine learning to imitate the human approach to extract the information from the semi-structured documents.

So Staple AI imitates human behaviour but why exactly? After all, they make errors, are really slow, and they demand (lots of) money. However, they are pretty good when it comes to finding information on a piece of paper just through a short glimpse. For a document such as an invoice, a human doesn't need to read the entire text to understand and pick out the key-field values. Instead, they look for the structural and contextual cues, the general layout, and the formatting of the text. For example, it's fairly easy for a human to extract the 'Invoice number' or the 'Total amount' in the image below. No where in the image is it written: 'The invoice number is WN201901002' but we all understand it almost instantaneously. But why exactly is it so simple ? Is it the text in bold with the number on the right? Or is it the assignment via a colon that assigns a number to the word 'invoice number'? Probably the structure and the localisation of the text (as compared to a big paragraph like this) in a certain way helps. For machines, semi-structured information like this is difficult to interpret with traditional NLP (Natural Language Processing) tools because mostly there is no grammar or proper sentences involved which can correlate a value to field. So indeed copying a human sounds great!

Staple extracts meaningful information from documents such as invoices, receipts, and purchase orders and transfers the digital data directly to your integrated accounting software or just to you in a clean JSON format (if that's what you are into). The Staple platform hunts for the key-information required for scanning the document such as the supplier's name, tax, total amount,

BILL TO	Invoice Number: WN201901002		
[REDACTED] Pte Ltd	Invoice Date: January 9, 2019		
32 Carpenter Street #03-01	Payment Due: January 9, 2019		
059911			
Singapore			
[REDACTED]@[REDACTED].com	Amount Due (SGD): \$0.00		
<b>Services</b>	<b>Quantity</b>	<b>Price</b>	<b>Amount</b>
Deposit - Provision of Employee Full Stack Developer	1	\$300.00	\$300.00
Deposit - Provision of Employee Front-End Developer	1	\$300.00	\$300.00
	<b>Total:</b>	<b>\$600.00</b>	
Payment on January 11, 2019 using a bank payment:			\$600.00
	<b>Amount Due (SGD):</b>	<b>\$0.00</b>	

invoice number, and most importantly, the line-items. The document can be in a pdf or an image format. For pdfs, the text is unskewed, dewarped, and without any washed off ink which makes life easier. Images are rather tricky because people are lazy and not everyone has the photography skills of a sixteen year old. Other problems such as shadows, local curling of the paper and washed off ink are extremely common. Thankfully there are sophisticated computer vision tools available that can get rid of the messy background and dewarp the picture.

In the first step, Staple cleans the image and prepares it for an optimal text detection to attain high reading accuracy. In the second step, the processed image is passed through an OCR engine where the text is extracted conserving the spatial information of the text on the page. After all, that's exactly how a human understands the context with the help of the structure and location of the text. After text detection, Staple's document scanning platform runs the detected text and its spatial information on the page through its proprietary model and outputs a JSON file with the required fields and their values, the confidence of the values found, and the bounding box of the value with respect to the top left corner of the page. Staple understands that every document is beautiful and unique in its own way (just like every human) and no 'template' based society can contain their beauty. Staple celebrates this uniqueness and has built the invoice scanning platform such that it is capable of scanning documents independent of their structure, layout and design.

To visualise the capability of Staple's document scanning platform, Staple has built a free API and a live demo which can be used to scan documents such as invoices and receipts. At the moment, pdfs and images of invoices and receipts are accepted. The API is deployed via Amazon web services using Elastic beanstalk which dockerises the Flask application. Elastic Beanstalk maintains the deployment, manages the load and auto-scales with the help of a load balancer. The auto-scaling and load-capacity varies from 1 to 100,000 documents being scanned in a day. One can use the API separately or have Staple integrate your Xero accounting software where your details of the document would automatically be posted. Moreover, Staple has built a Reconciler to further make the life of accountants easier by reconciling payments via Stripe. Enough said! Is it any good? Let's find out!

So let's free humans from all the mundane accounting tasks and leave them alone to what they do best, netflix and chill!

<p><b>INVOICE</b> Reg. No. <b>Address</b> Singapore URL www.  Address BILL TO Singapore Email com</p> <p>Invoice Number: WN201901002 Invoice Date: January 11, 2019 Payment Due: January 9, 2019 Amount Due SGD \$0.00</p> <p><b>Services</b></p> <table border="1"> <tr> <th>Quantity</th> <th>Price</th> <th>Amount</th> </tr> <tr> <td>1</td> <td>\$300.00</td> <td>\$300.00</td> </tr> <tr> <td>1</td> <td>\$300.00</td> <td>\$300.00</td> </tr> </table> <p>Total: SGD \$600.00 Payment on January 11, 2019 using a bank payment: SGD \$600.00</p> <p><b>Currency</b></p> <p>Amount Due (SGD): SGD \$0.00</p> <p><b>Notes</b> Account Holder: Account name Bank Name: DBS Account number SGD Account Number: SGD Account number Bank Code: Branch Code:</p> <p>Thank you for your business! Help us spread the word and you'll receive a referral fee!</p>	Quantity	Price	Amount	1	\$300.00	\$300.00	1	\$300.00	\$300.00	 <p>玲记成记食供有限公司 INV NO: 1001723 DATE: 22/3/2019</p> <p><b>BILL TO</b> Address SINGAPORE (T20)</p> <p><b>DELIVER TO</b> Address</p> <table border="1"> <thead> <tr> <th>OUR REFERENCE PO number</th> <th>CUSTOMER REFERENCE 3365</th> <th>TERMS 15</th> <th>PAGE 1</th> </tr> <tr> <th>ITEM CODE</th> <th>DESCRIPTION</th> <th>QUANTITY</th> <th>UNIT PRICE</th> <th>AMOUNT</th> </tr> </thead> <tbody> <tr><td>YEC-067</td><td>TAIWAN GREEN / SIO PAK CHOW 小白菜</td><td>1.00 KG</td><td>2.70</td><td>2.70</td></tr> <tr><td>YEC-221</td><td>CARROT WHITE / RADDISH 白萝卜</td><td>5.00 KG</td><td>1.55</td><td>8.65</td></tr> <tr><td>YEC-070</td><td>ONION SPRING (LOCAL) 生葱</td><td>1.00 KG</td><td>4.90</td><td>4.90</td></tr> <tr><td>YEC-325</td><td>SWEET POTATOES (JAPANESE) 日本番薯</td><td>5.00 KG</td><td>3.00</td><td>15.00</td></tr> <tr><td>YEC-147</td><td>ONION LANCE WHITE (PEELED) 白大葱肉</td><td>4.00 KG</td><td>1.50</td><td>6.00</td></tr> <tr><td>YEC-203</td><td>BLUEBERRY FRESH (225gm) 蓝莓</td><td>5.00 PKT</td><td>5.80</td><td>29.00</td></tr> <tr><td>YEC-326</td><td>HERB FRESH MINT (50gm) 香草</td><td>2.00 BLD</td><td>2.20</td><td>4.40</td></tr> <tr><td>YEC-029</td><td>ASPARAGUS FRESH (AUS) USA 美芦笋</td><td>1.00 KG</td><td>15.00</td><td>15.00</td></tr> <tr><td>YEC-174</td><td>POPCORN (LOCAL) (100GM) 爆米花</td><td>5.00 BLD</td><td>1.50</td><td>7.50</td></tr> <tr><td>YEC-170</td><td>BBQ FRESH OLIVES (50gm) ISABEL 烟香草</td><td>3.00 BUNDLE</td><td>2.50</td><td>7.50</td></tr> <tr><td>YRT-009</td><td>APPLE GREEN 苹果</td><td>10.00 KGS</td><td>0.39</td><td>3.90</td></tr> <tr><td>YRT-019</td><td>TOFU ROC SILKEN (LOCAL) 椰豆腐</td><td>4.00 PKT</td><td>0.75</td><td>3.00</td></tr> <tr><td>YEC-135</td><td>BEAN FINE (KENYA / FRANCE) 小毛豆</td><td>0.60 KG</td><td>7.00</td><td>4.20</td></tr> <tr><td>YRT-022</td><td>ORANGES 橙</td><td>5.00 KGS</td><td>0.35</td><td>1.75</td></tr> <tr><td>YRT-002</td><td>RAVANA PHILIPPINES 香蕉</td><td>4.00 KG</td><td>1.90</td><td>7.60</td></tr> <tr><td>YEC-156</td><td>CORN ON COB FRESH 玉米棒</td><td>20.00 KGS</td><td>0.60</td><td>12.00</td></tr> <tr><td>YEC-014</td><td>CELERI IMPORTED 芥芹</td><td>1.00 KG</td><td>2.00</td><td>2.00</td></tr> <tr><td>YEC-298</td><td>LITCHIE GREEN CORAL 椴珊瑚生菜</td><td>1.00 KG</td><td>4.50</td><td>4.50</td></tr> <tr><td>YRT-029</td><td>FIREAPPLE (HONEY)</td><td>2.00 KGS</td><td>1.80</td><td>3.60</td></tr> <tr><td>YEC-013</td><td>CAULIFLOWER 花椰菜</td><td>5.00 KG</td><td>2.30</td><td>11.50</td></tr> <tr><td>YEC-045</td><td>BAKED CHICKEN 鸡块</td><td>1.00 KG</td><td>2.30</td><td>2.30</td></tr> <tr><td>YEC-397</td><td>CHERRY RED (SWEETNESS) 无籽葡萄</td><td>1.00 KG</td><td>6.50</td><td>6.50</td></tr> <tr><td>YEC-075</td><td>TOOMATES (LOCAL) 香茄</td><td>6.00 KG</td><td>1.90</td><td>11.40</td></tr> <tr><td>YEC-059</td><td>PUMPKIN (LOCAL) 南瓜</td><td>3.50 KG</td><td>1.20</td><td>4.20</td></tr> <tr><td>YRT-018</td><td>WATER MELON RED (SWEETNESS) 西瓜</td><td>5.50 KG</td><td>1.20</td><td>6.60</td></tr> <tr><td>YEC-158</td><td>PETAL FRESH PEELD CHAMPUANO 菠萝</td><td>1.00 PKT</td><td>18.00</td><td>18.00</td></tr> <tr><td>YRT-032</td><td>ROCK MELON / CANTALOUP</td><td>2.00 KG</td><td>4.00</td><td>8.00</td></tr> <tr><td>YRT-014</td><td>STRAWBERRY FRESH 草莓</td><td>1.50 KG</td><td>15.80</td><td>23.70</td></tr> <tr><td>YRT-037</td><td>LEMON YELLOW 黄柠檬</td><td>6.00 KG</td><td>3.30</td><td>19.80</td></tr> </tbody> </table> <p>Retrieved the above goods in good and fit condition. A Member of D2P Credit Business - Your pro-active credit profile to your company It's contribute towards building Reconcile &amp; Accepted by _____ _____ Sub Total: SGD 204.19 Tax total: SGD 16.83 Amount Due: SGD 220.02</p>	OUR REFERENCE PO number	CUSTOMER REFERENCE 3365	TERMS 15	PAGE 1	ITEM CODE	DESCRIPTION	QUANTITY	UNIT PRICE	AMOUNT	YEC-067	TAIWAN GREEN / SIO PAK CHOW 小白菜	1.00 KG	2.70	2.70	YEC-221	CARROT WHITE / RADDISH 白萝卜	5.00 KG	1.55	8.65	YEC-070	ONION SPRING (LOCAL) 生葱	1.00 KG	4.90	4.90	YEC-325	SWEET POTATOES (JAPANESE) 日本番薯	5.00 KG	3.00	15.00	YEC-147	ONION LANCE WHITE (PEELED) 白大葱肉	4.00 KG	1.50	6.00	YEC-203	BLUEBERRY FRESH (225gm) 蓝莓	5.00 PKT	5.80	29.00	YEC-326	HERB FRESH MINT (50gm) 香草	2.00 BLD	2.20	4.40	YEC-029	ASPARAGUS FRESH (AUS) USA 美芦笋	1.00 KG	15.00	15.00	YEC-174	POPCORN (LOCAL) (100GM) 爆米花	5.00 BLD	1.50	7.50	YEC-170	BBQ FRESH OLIVES (50gm) ISABEL 烟香草	3.00 BUNDLE	2.50	7.50	YRT-009	APPLE GREEN 苹果	10.00 KGS	0.39	3.90	YRT-019	TOFU ROC SILKEN (LOCAL) 椰豆腐	4.00 PKT	0.75	3.00	YEC-135	BEAN FINE (KENYA / FRANCE) 小毛豆	0.60 KG	7.00	4.20	YRT-022	ORANGES 橙	5.00 KGS	0.35	1.75	YRT-002	RAVANA PHILIPPINES 香蕉	4.00 KG	1.90	7.60	YEC-156	CORN ON COB FRESH 玉米棒	20.00 KGS	0.60	12.00	YEC-014	CELERI IMPORTED 芥芹	1.00 KG	2.00	2.00	YEC-298	LITCHIE GREEN CORAL 椴珊瑚生菜	1.00 KG	4.50	4.50	YRT-029	FIREAPPLE (HONEY)	2.00 KGS	1.80	3.60	YEC-013	CAULIFLOWER 花椰菜	5.00 KG	2.30	11.50	YEC-045	BAKED CHICKEN 鸡块	1.00 KG	2.30	2.30	YEC-397	CHERRY RED (SWEETNESS) 无籽葡萄	1.00 KG	6.50	6.50	YEC-075	TOOMATES (LOCAL) 香茄	6.00 KG	1.90	11.40	YEC-059	PUMPKIN (LOCAL) 南瓜	3.50 KG	1.20	4.20	YRT-018	WATER MELON RED (SWEETNESS) 西瓜	5.50 KG	1.20	6.60	YEC-158	PETAL FRESH PEELD CHAMPUANO 菠萝	1.00 PKT	18.00	18.00	YRT-032	ROCK MELON / CANTALOUP	2.00 KG	4.00	8.00	YRT-014	STRAWBERRY FRESH 草莓	1.50 KG	15.80	23.70	YRT-037	LEMON YELLOW 黄柠檬	6.00 KG	3.30	19.80
Quantity	Price	Amount																																																																																																																																																																		
1	\$300.00	\$300.00																																																																																																																																																																		
1	\$300.00	\$300.00																																																																																																																																																																		
OUR REFERENCE PO number	CUSTOMER REFERENCE 3365	TERMS 15	PAGE 1																																																																																																																																																																	
ITEM CODE	DESCRIPTION	QUANTITY	UNIT PRICE	AMOUNT																																																																																																																																																																
YEC-067	TAIWAN GREEN / SIO PAK CHOW 小白菜	1.00 KG	2.70	2.70																																																																																																																																																																
YEC-221	CARROT WHITE / RADDISH 白萝卜	5.00 KG	1.55	8.65																																																																																																																																																																
YEC-070	ONION SPRING (LOCAL) 生葱	1.00 KG	4.90	4.90																																																																																																																																																																
YEC-325	SWEET POTATOES (JAPANESE) 日本番薯	5.00 KG	3.00	15.00																																																																																																																																																																
YEC-147	ONION LANCE WHITE (PEELED) 白大葱肉	4.00 KG	1.50	6.00																																																																																																																																																																
YEC-203	BLUEBERRY FRESH (225gm) 蓝莓	5.00 PKT	5.80	29.00																																																																																																																																																																
YEC-326	HERB FRESH MINT (50gm) 香草	2.00 BLD	2.20	4.40																																																																																																																																																																
YEC-029	ASPARAGUS FRESH (AUS) USA 美芦笋	1.00 KG	15.00	15.00																																																																																																																																																																
YEC-174	POPCORN (LOCAL) (100GM) 爆米花	5.00 BLD	1.50	7.50																																																																																																																																																																
YEC-170	BBQ FRESH OLIVES (50gm) ISABEL 烟香草	3.00 BUNDLE	2.50	7.50																																																																																																																																																																
YRT-009	APPLE GREEN 苹果	10.00 KGS	0.39	3.90																																																																																																																																																																
YRT-019	TOFU ROC SILKEN (LOCAL) 椰豆腐	4.00 PKT	0.75	3.00																																																																																																																																																																
YEC-135	BEAN FINE (KENYA / FRANCE) 小毛豆	0.60 KG	7.00	4.20																																																																																																																																																																
YRT-022	ORANGES 橙	5.00 KGS	0.35	1.75																																																																																																																																																																
YRT-002	RAVANA PHILIPPINES 香蕉	4.00 KG	1.90	7.60																																																																																																																																																																
YEC-156	CORN ON COB FRESH 玉米棒	20.00 KGS	0.60	12.00																																																																																																																																																																
YEC-014	CELERI IMPORTED 芥芹	1.00 KG	2.00	2.00																																																																																																																																																																
YEC-298	LITCHIE GREEN CORAL 椴珊瑚生菜	1.00 KG	4.50	4.50																																																																																																																																																																
YRT-029	FIREAPPLE (HONEY)	2.00 KGS	1.80	3.60																																																																																																																																																																
YEC-013	CAULIFLOWER 花椰菜	5.00 KG	2.30	11.50																																																																																																																																																																
YEC-045	BAKED CHICKEN 鸡块	1.00 KG	2.30	2.30																																																																																																																																																																
YEC-397	CHERRY RED (SWEETNESS) 无籽葡萄	1.00 KG	6.50	6.50																																																																																																																																																																
YEC-075	TOOMATES (LOCAL) 香茄	6.00 KG	1.90	11.40																																																																																																																																																																
YEC-059	PUMPKIN (LOCAL) 南瓜	3.50 KG	1.20	4.20																																																																																																																																																																
YRT-018	WATER MELON RED (SWEETNESS) 西瓜	5.50 KG	1.20	6.60																																																																																																																																																																
YEC-158	PETAL FRESH PEELD CHAMPUANO 菠萝	1.00 PKT	18.00	18.00																																																																																																																																																																
YRT-032	ROCK MELON / CANTALOUP	2.00 KG	4.00	8.00																																																																																																																																																																
YRT-014	STRAWBERRY FRESH 草莓	1.50 KG	15.80	23.70																																																																																																																																																																
YRT-037	LEMON YELLOW 黄柠檬	6.00 KG	3.30	19.80																																																																																																																																																																

