# XCENTER DIGITAL

## Optical Character Recognition
## Best Practices

# Optical Character Recognition (OCR) – Overview

Although the scanning of paper-based invoices isn't considered e-invoicing, it is the natural first step for many organisations. Using OCR software, the data can be moved from a paper-based format to a digital format that can be entered in the Accounts Payable system. Outside of EDI and invoice portals, OCR has been a predominant tool of choice to enable the digitization of invoices.

Optical Character Recognition (OCR) technology is a hardware/software tool that takes a paper document, usually an invoice, scans then "reads" it and turns it into metadata that can be used to populate fields in a database.

From there the invoice can be brought into an electronic workflow for processing. Using OCR software, the data can be moved from a paper-based format to a digital format that can be entered in the AP system. OCR is the electronic conversion (through scanning) of invoices without extractable data (either paper or image files) into data that can be integrated directly (as an EDI or XML file) into a buyers Accounts Payable finance system for payment.

Whilst OCR solutions enables organisations to automate their AP processes to a certain extent, there are restrictions that are inherent to OCR technology, and which limit its impact beyond achieving a semi-automated state, where human intervention and errors are part and parcel of the technology in question. After all, we speak of "recognition" and not "extraction" when referring to OCR. Fundamentally, OCR solutions are all based on a similar probabilistic technology and methodology. For instance, the number "1" vs. lowercase letter "L", the number "0" vs. uppercase O, and so on. The latter is mitigated to some extent by the use of dictionaries (for example, "INVOICE" is more likely than "1NV0lCE"), but unfortunately invoice data such as the invoice number or the shipping reference, is usually not to be found in an OCR dictionary.

The challenge gets even more difficult when using OCR for invoice line item extraction. These inherent limitations of OCR result in varying accuracy recognition rates, which invariably requires human operators to check the results produced by OCR. Inaccuracies require manual intervention, leading to errors, long invoice processing time, and low percentage of "touchless" invoices or processed "straight-through".

# Optical Character Recognition (OCR) - Best Practices

### *Situational analysis*

OCR tools perform best with good quality typed documents. Therefore, if the scan quality is poor or the invoice is stained, wrinkled or the text is weakly printed then OCR will fail to successfully digitize such invoices. OCR tools also typically fail to read handwritten invoices, and therefore fonts that resemble handwriting create many errors during the OCR process. AP teams need to undertake a detailed analysis of the printed source material.

Specifics such as paper quality, characteristics such as language, font, and layout, and graphical elements are critical features that may affect the quality of data capture. This will provide information that will help to determine whether the data capture can be easily improved or not. For example, certain kind of historical documents may lack the lexical data that is required for OCR data entry and this can pose challenges. Image rich documents may need special measures to render them OCR compatible or may need improvised OCR data capture.

### *OCR process workflow*

It cannot be emphasized enough that having a well-delineated process workflow will determine the success or failure of your OCR data entry effort. A well-charted flow will ensure that your data capture and conversion succeeds as per your expectations.

### *OCR quality check procedures*

OCR data capture projects require quality assurance procedures in order to ensure a certain level of control. A quality assurance program will safeguard that the OCR results are on track and that the goals are achieved within the defined time.

Part of the procedure consists of the AP team carrying out a comprehensive review of the complete captured data set or a sectional review if historical reasons dictate this approach – for instance in the case of a spin-off. QA procedures also include tracking and modification of OCR errors. Ensure that your QA program is rigorously worked out and implemented and that the quality standards are adequately communicated to all staff involved.

### *Quality of master vendor data*

OCR solutions do not improve quality, accessibility and reliability of master data since there is a process gap between creation and digitization of invoices. It therefore helps to ensure accurate and up to date master vendor data.

### *Possible methods of processing*

a) In-house mailroom and small scan workplaces

b) Outsourcing (PO boxes need to be set up)

c) Combination of both methods (PO boxes + centralized scanning for the majority of documents coming from vendors, local scan workplaces with desktop-class scanners for documents received locally on branches – faster processing but still good quality)

### *Equipment - Production scanner(s), image enhancement technology (VRS) – examples:*

www.xcenter.digital
sales@xcenter.digital

Copyright 2021 XCENTER DIGITAL logo™

Background color removal, de-skew, de-speckle, automatic orientation, auto crop, edge clean-up, blank page removal, de-dithering…

*\* For small workplaces we recommend Kofax Express*

If small local scanning workplaces are not possible for budgetary reasons for example, we recommend to scan on MFP devices in full color and then use software binarization for image enhancement (explained later in "Email attachments" chapter )

### *Scanning parameters*

- Vertical resolution: 300DPI
- Horizontal resolution: 300DPI
- Color: black and white
- Size automatic cut according to the actual size of the document
- Output scanning: 1x multipage PDF file per document (incl appendixes if applicable)
- Image enhancement: Image enhancing techniques must be applied on document. Techniques such as: De-skew, Black border crop & remove, De-speckle / Noise reduction, Punch Hole fill, Blank page removal etc.