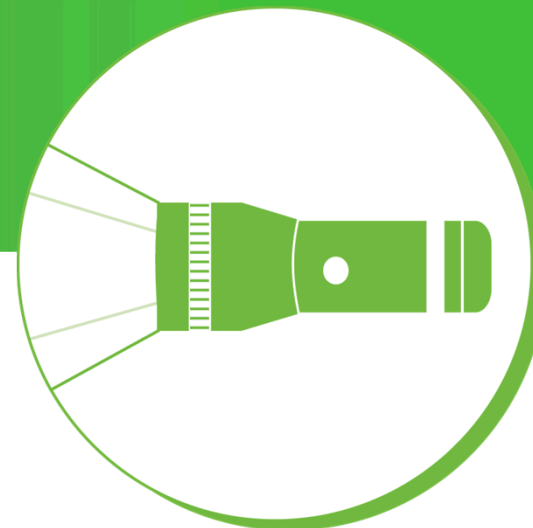


ECM and Dark Data

Turn on the light to improve compliance, cut storage, and leverage document assets

Peter Duff
CEO, Adlib
Vice Chair, PDF Association

Roger Beharry Lall, ecm^P
Director, Product Marketing, Adlib



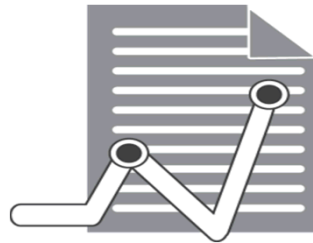
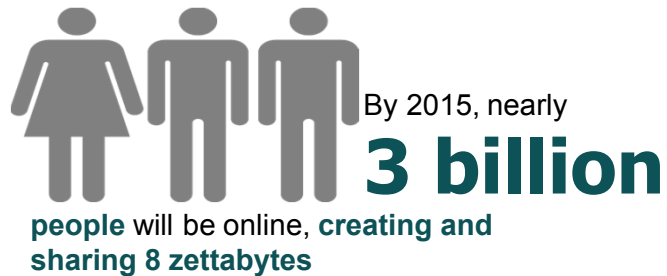
The Data Explosion

INDIVIDUALS
ARE CREATING VAST
AMOUNTS
OF DATA.



By 2020, **B2B transactions**
on the internet will reach **450**
billion
per day

WHAT'S CONTRIBUTING TO THE EXPLOSION?



File Type Growth Rate of
Consumer Internet Traffic

File Sharing **23%**
Data **29%**
(CAGR 2010-2015)

Enterprise data will grow

650%,

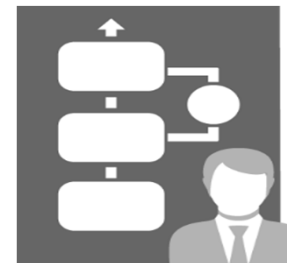
partially due to regulations like the
Sarbanes-Oxley Act requiring
companies to store financial records



DATA...
WE HAVE A
PROBLEM.



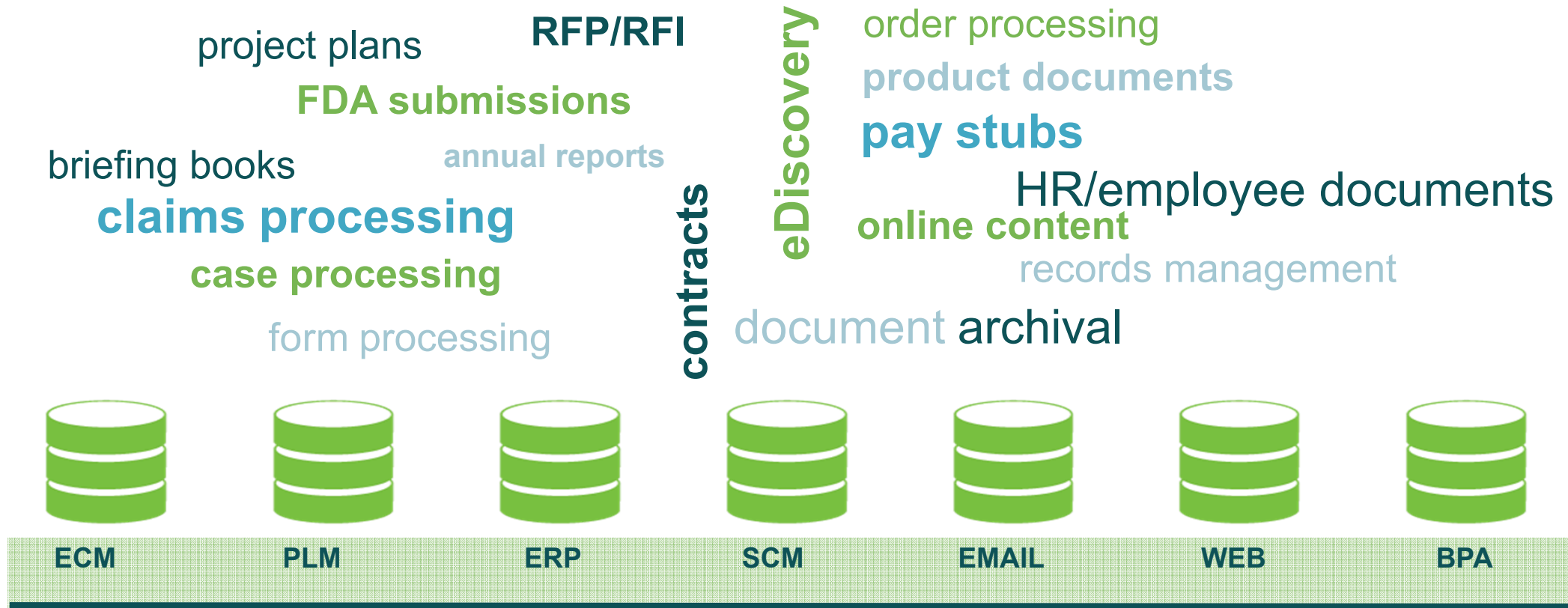
35% of the
digital universe is **subject**
to compliance and
regulations



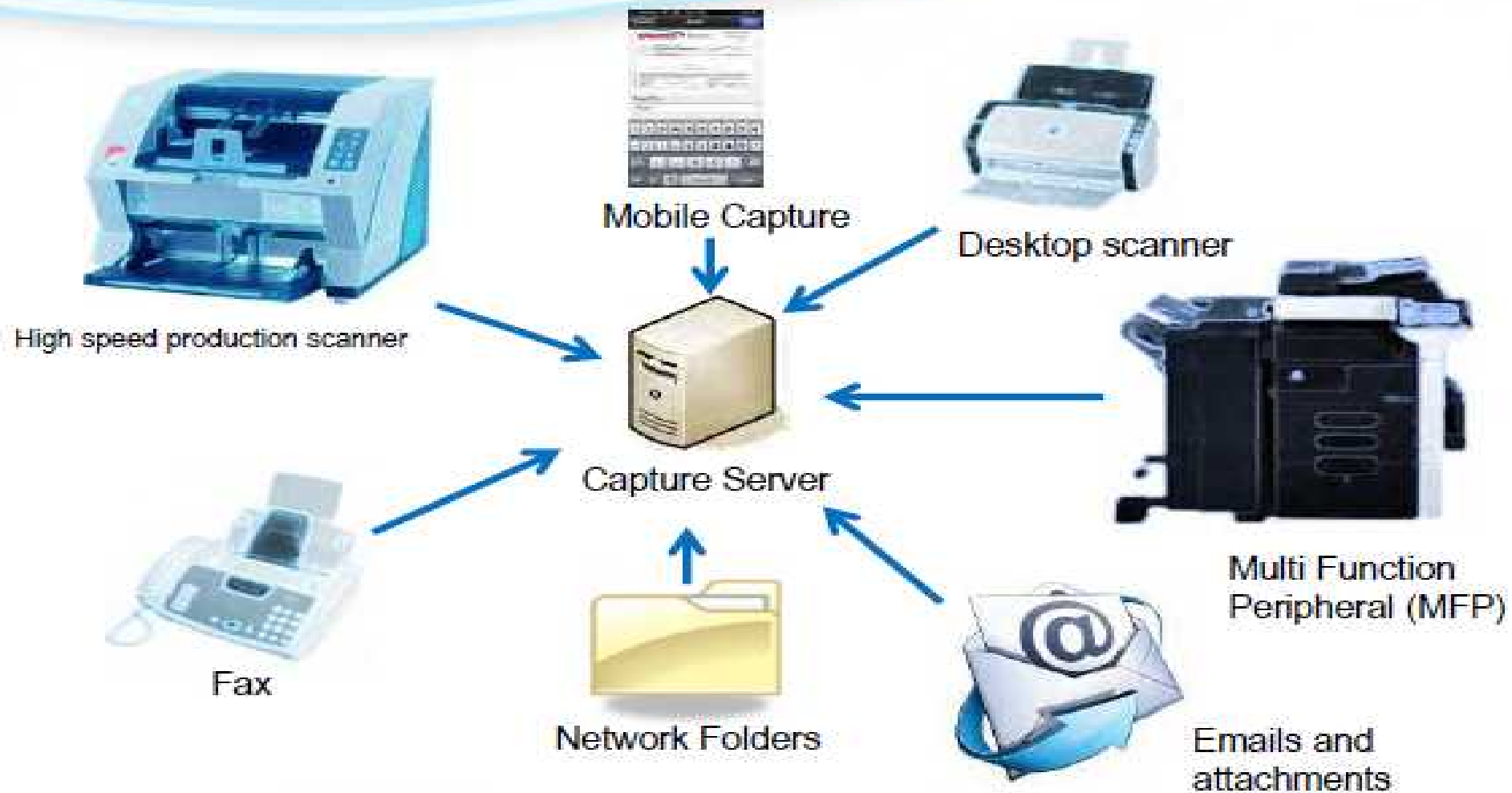
In the next decade, the number
of files will grow by a factor of
75 while **IT professions**
will grow by **less** than
a factor of **1.5**

Document Complexity

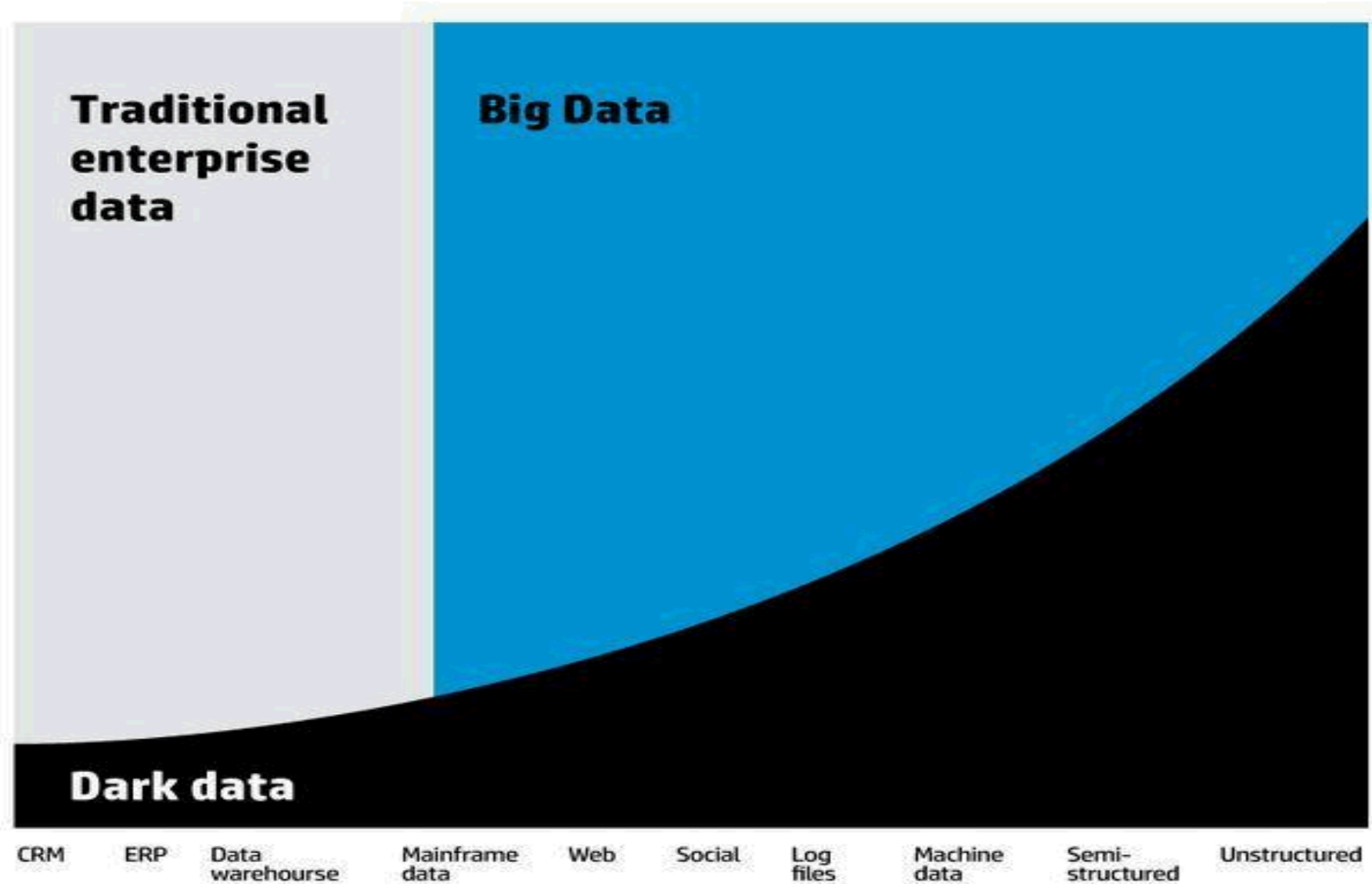
Variety Of Systems, Processes And Formats



Multi-Channel Input / Batch Creation



Growth of Dark Data





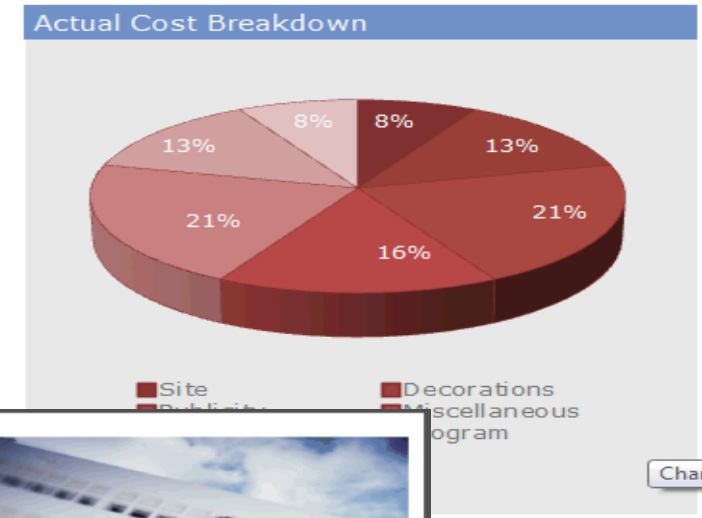
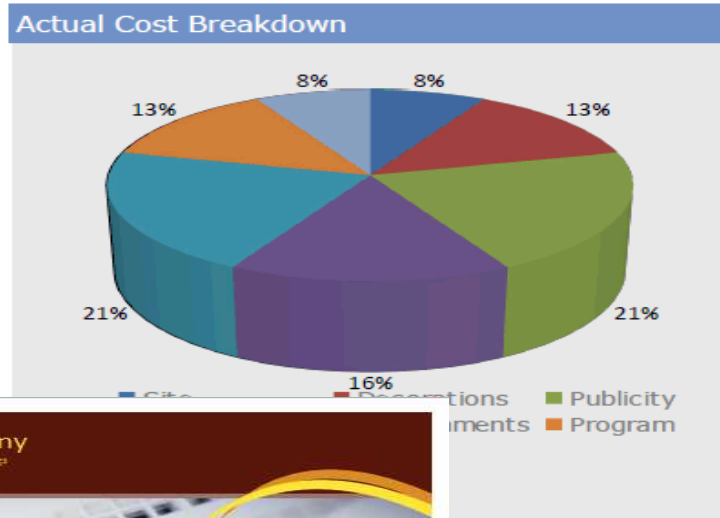
Gartner

The information assets that organizations collect, process and store during regular business activities, but generally fail to use for other purposes.

IDC

Up to 90% of Big Data is Dark Data.

Dirty Data. Dark Data. No Data!

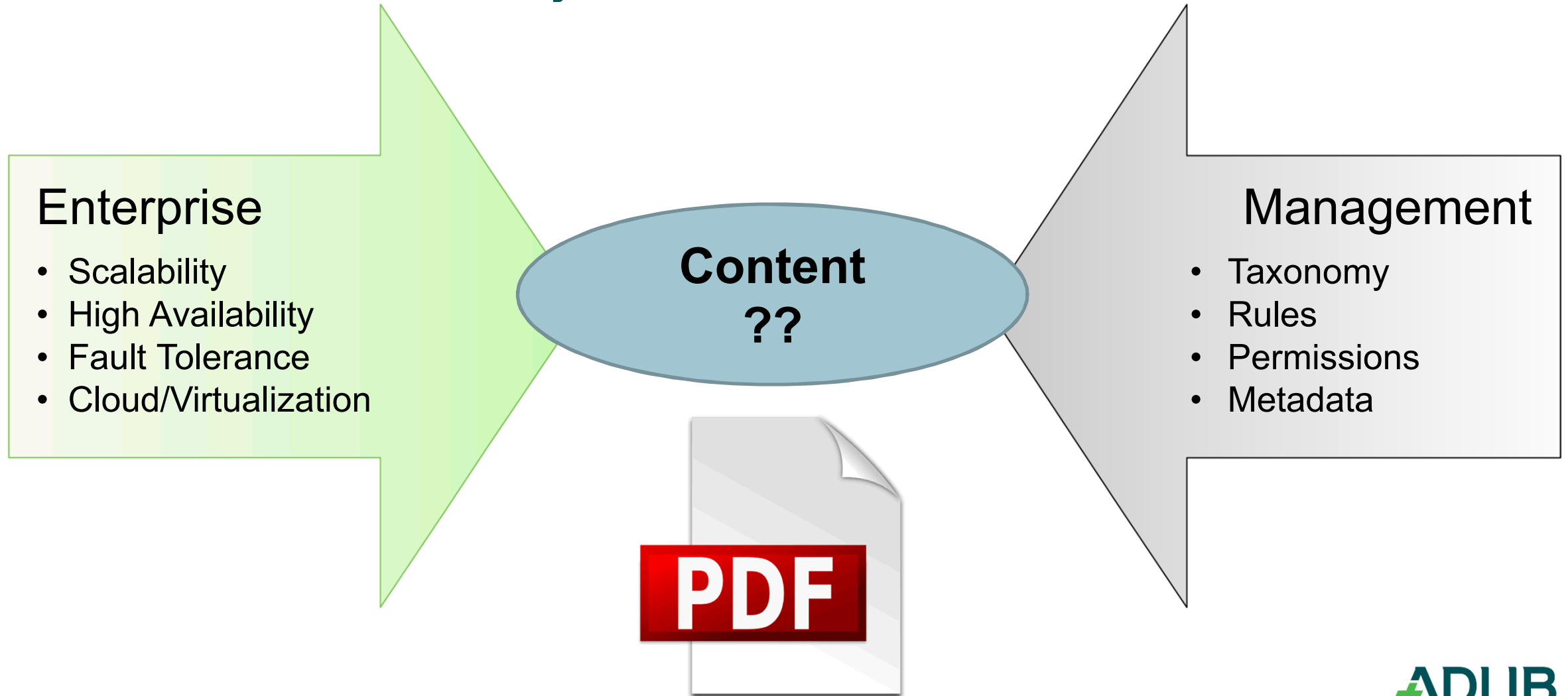


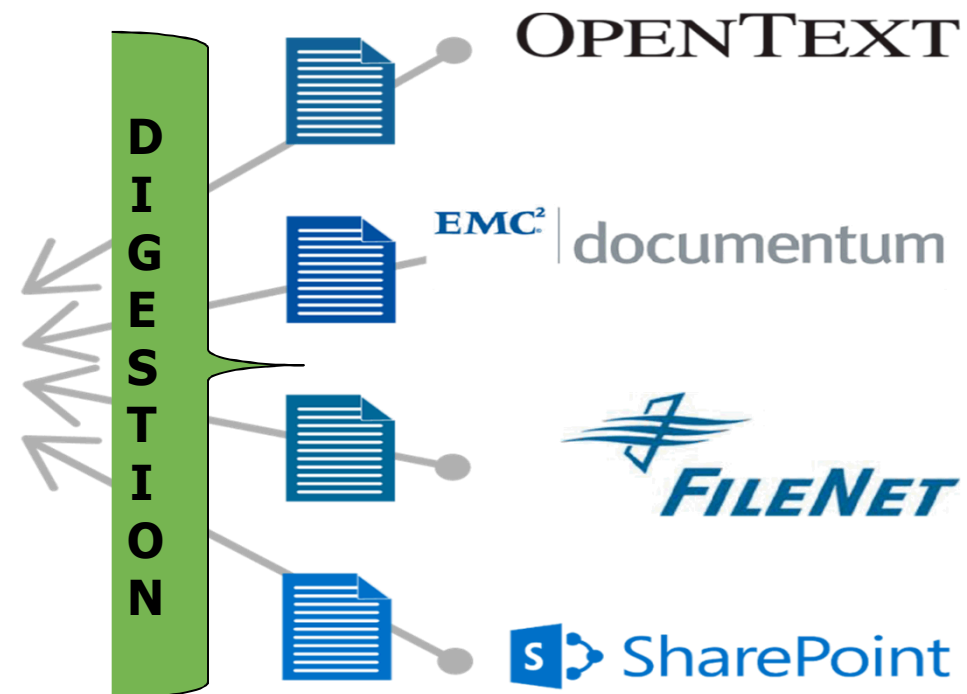
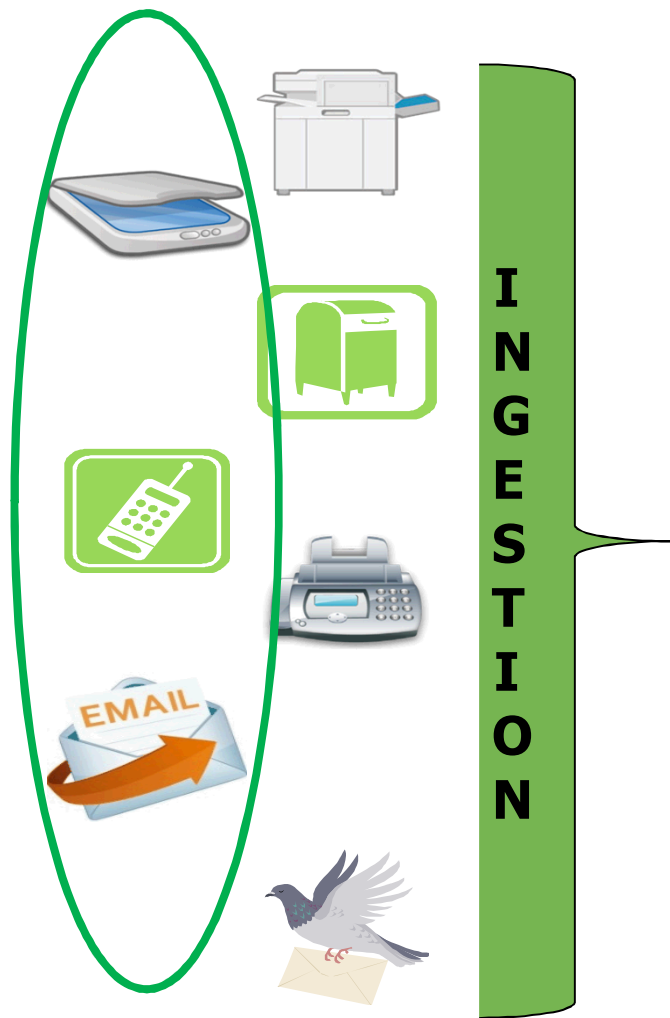
Original Source



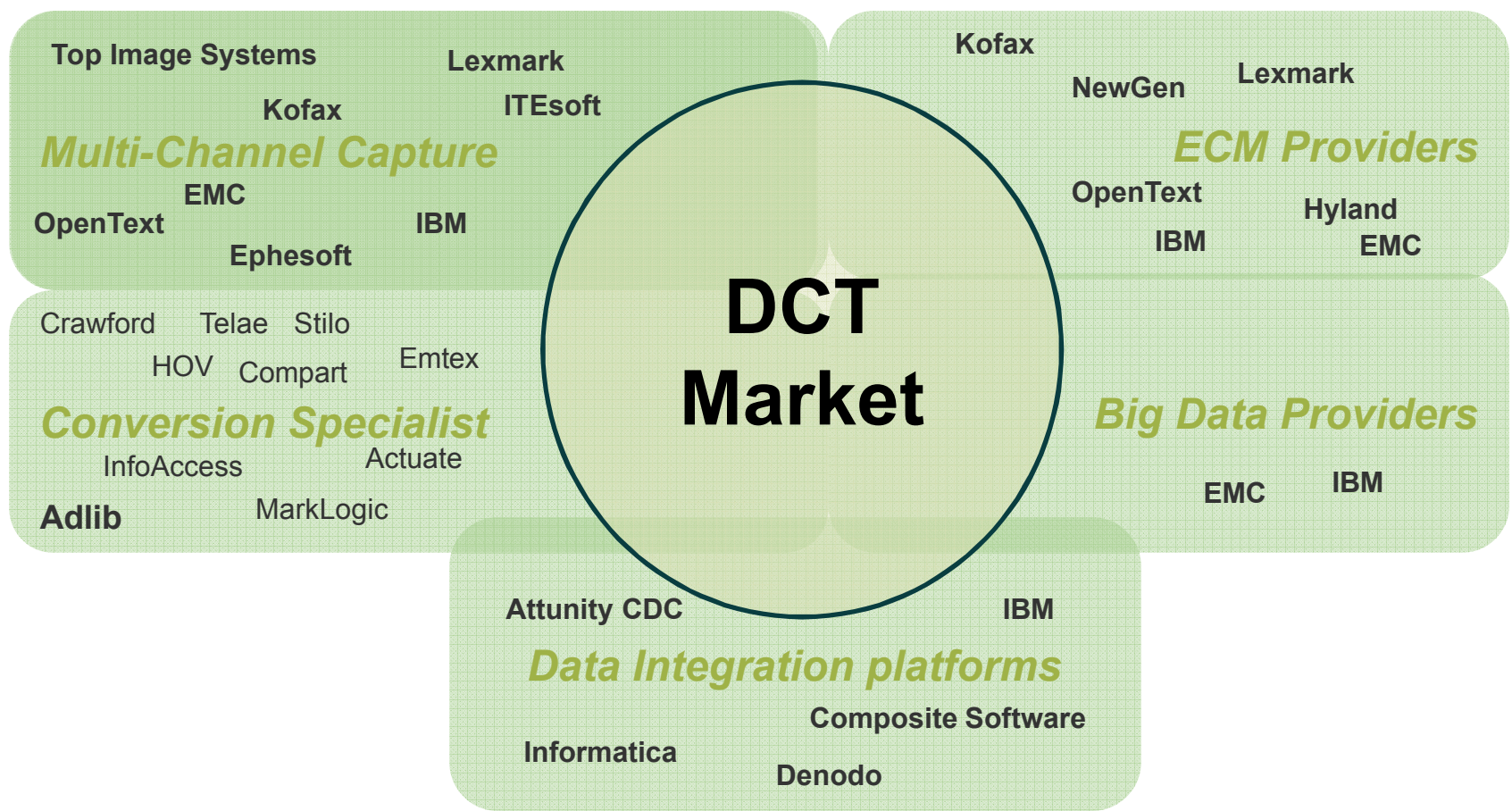
Rendered Image

Enterprise Content Management: Only Part of the Answer

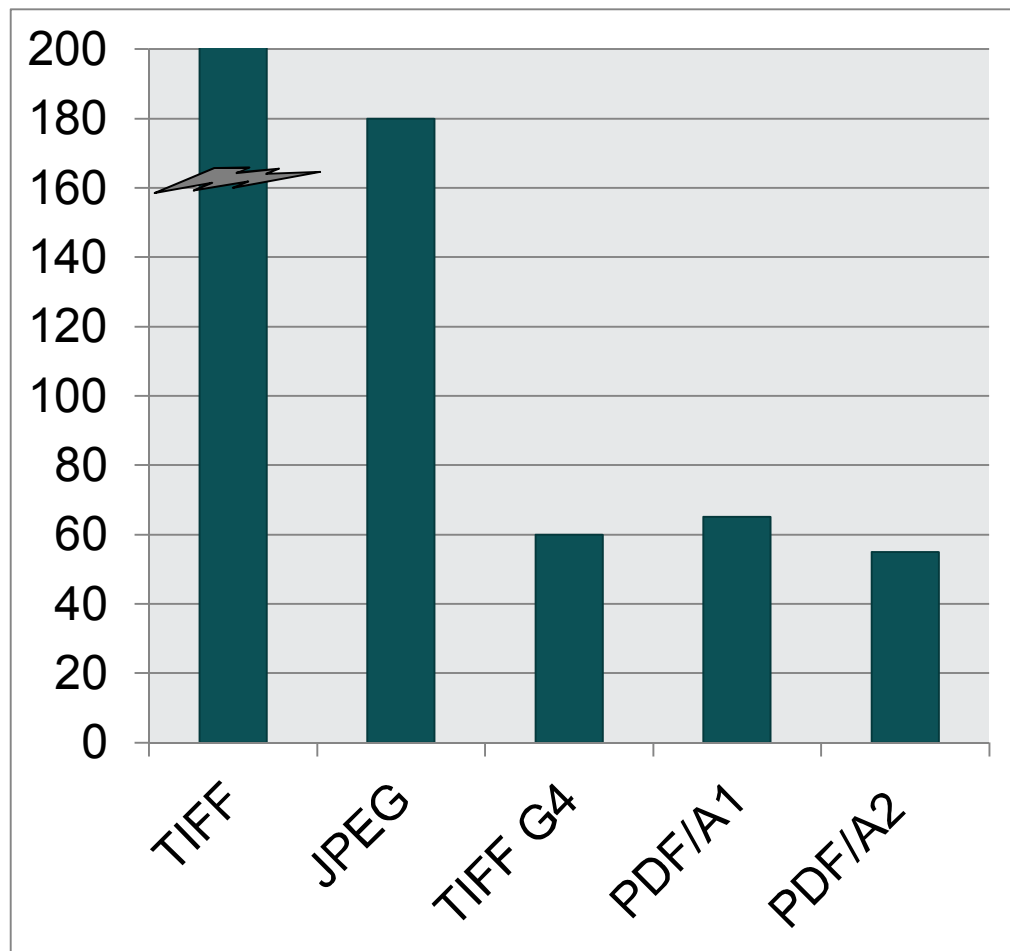




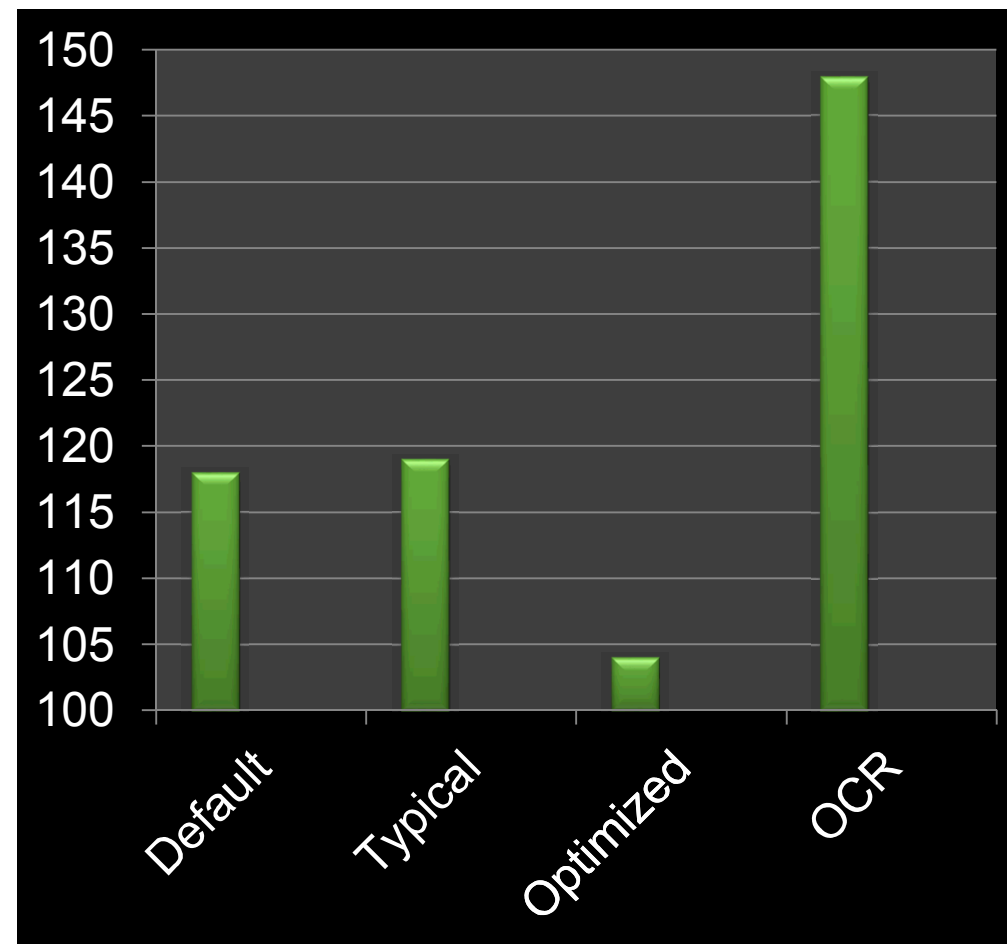
Document Content Transformation



File Size Optimization for Storage Reduction

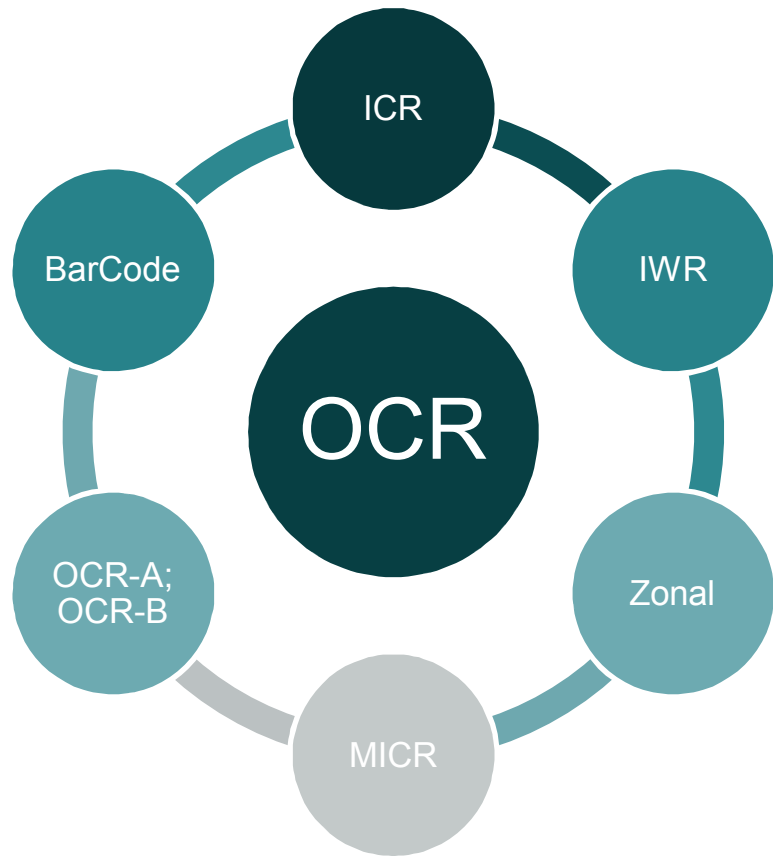


OPTIMAL FORMAT



OPTIMAL PDF

Optical Character Recognition (OCR)



Converting printed or written text characters—captured as images during scanning—into computer-based, encoded text.

Benefits of OCR Capabilities

- Liberating information for electronic searches
- Delivering industry-leading accuracy
- Supporting regulatory mandates
- **Make content immediately findable from the moment of capture.**

Search Enabled Documents

ECT
123 Compliance Street
San Francisco, CA 94114, USA
Phone: 415 555 12125

P.O. DATE	REQUISITIONER

QTY	DESCRIPTION
22	Installation

Adobe Acrobat

Acrobat has finished searching the document. No matches were found.

John Doe
Authorized by

Site Actions

ADLIB Sales ▶ Site Search Results

Adlib Demo Sales Human Resources Executive Research and Development Nintex

This Site: Sales

Try searching again in [All Sites](#).

1-2 of 2 results

No refinements available

- [Signature](#)
Authors: S
<http://sp20>
- [Quotation](#)
REQUISITIONER
...
Authors: S
<http://sp20>

ECT
123 Compliance Street
San Francisco, CA 94114, USA
Phone: 415 555 12125

Ship To: Ministry of Defence
ATTN: Mr. John Doe

Quotation Number
#: 000F1806042011

P.O. DATE	REQUISITIONER	SHIPPED VIA	F.O.B. POINT	TERMS

QTY	DESCRIPTION	UNIT PRICE	TOTAL
22	Installation		\$125,000.00
22	Quality Assurance - Air Stress Test and Diagnostic		\$25,000.00
22	Annual Maintenance Agreement		\$750,000.00
		SUBTOTAL	\$900,000.00
		HST	\$117,000.00
		TOTAL (US)	\$1,017,000.00

John Doe
Authorized by

06/10/2011
Date

Zonal OCR

```
<?xml version="1.0" encoding="UTF-16" ?>
<PDFINFO>
- <OCRZONES>
  <OCRZONE NAME="BARCODE1" PAGE="1" LEFT="150" TOP="365" WIDTH="383" HEIGHT="42" TYPE="Barcode">THISISABARCODE</OCRZONE>
  <OCRZONE NAME="BARCODE2" PAGE="1" LEFT="750" TOP="365" WIDTH="311" HEIGHT="43" TYPE="Barcode">18669911704</OCRZONE>
  <OCRZONE NAME="OMR1" PAGE="1" LEFT="335" TOP="596" WIDTH="64" HEIGHT="60" TYPE="Optical Mark">1</OCRZONE>
  <OCRZONE NAME="OMR2" PAGE="1" LEFT="811" TOP="596" WIDTH="64" HEIGHT="60" TYPE="Optical Mark">0</OCRZONE>
  <OCRZONE NAME="CASE_NUMBER" PAGE="1" LEFT="1010" TOP="882" WIDTH="81" HEIGHT="16" TYPE="Text">9873425</OCRZONE>
  <OCRZONE NAME="ShippingAddress" PAGE="1" LEFT="150" TOP="986" WIDTH="271" HEIGHT="101" TYPE="Text">Adlib Software, Inc. 215-
  3228 South Service Road Burlington, ON L7N 3H8</OCRZONE>
  <OCRZONE NAME="MICR1" PAGE="1" LEFT="191" TOP="1310" WIDTH="202" HEIGHT="19" TYPE="Magnetic Ink
  Character">041000124</OCRZONE>
  <OCRZONE NAME="MICR2" PAGE="1" LEFT="529" TOP="1310" WIDTH="183" HEIGHT="19" TYPE="Magnetic Ink
  Character">0000000000</OCRZONE>
  <OCRZONE NAME="MICR3" PAGE="1" LEFT="848" TOP="1311" WIDTH="220" HEIGHT="18" TYPE="Magnetic Ink
  Character">0000010000</OCRZONE>
</OCRZONES>
</PDFINFO>
```

Zonal OCR Demo Document

Barcode Sample:



OMR Sample:



OMR1 - Checked



OMR2 - Not Checked

TEXT Samples

Shipping Address:

Adlib Software, Inc.
215-3228 South Service Road
Burlington, ON
L7N 3H8

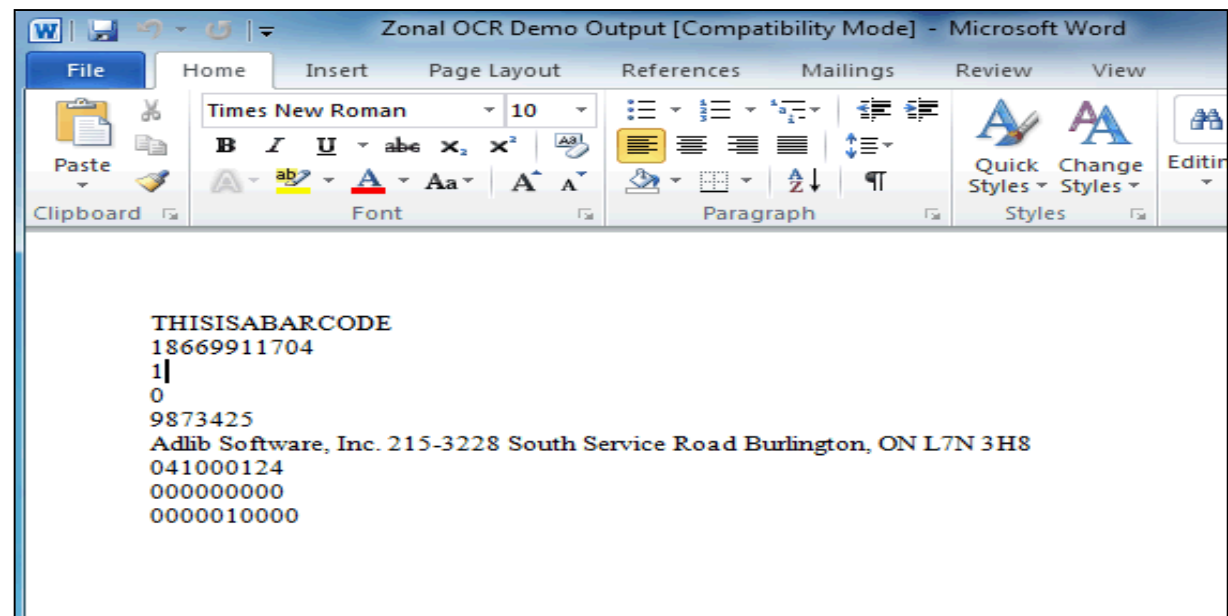
CASE NUMBER: 9873425

MICR Sample (Check Font)

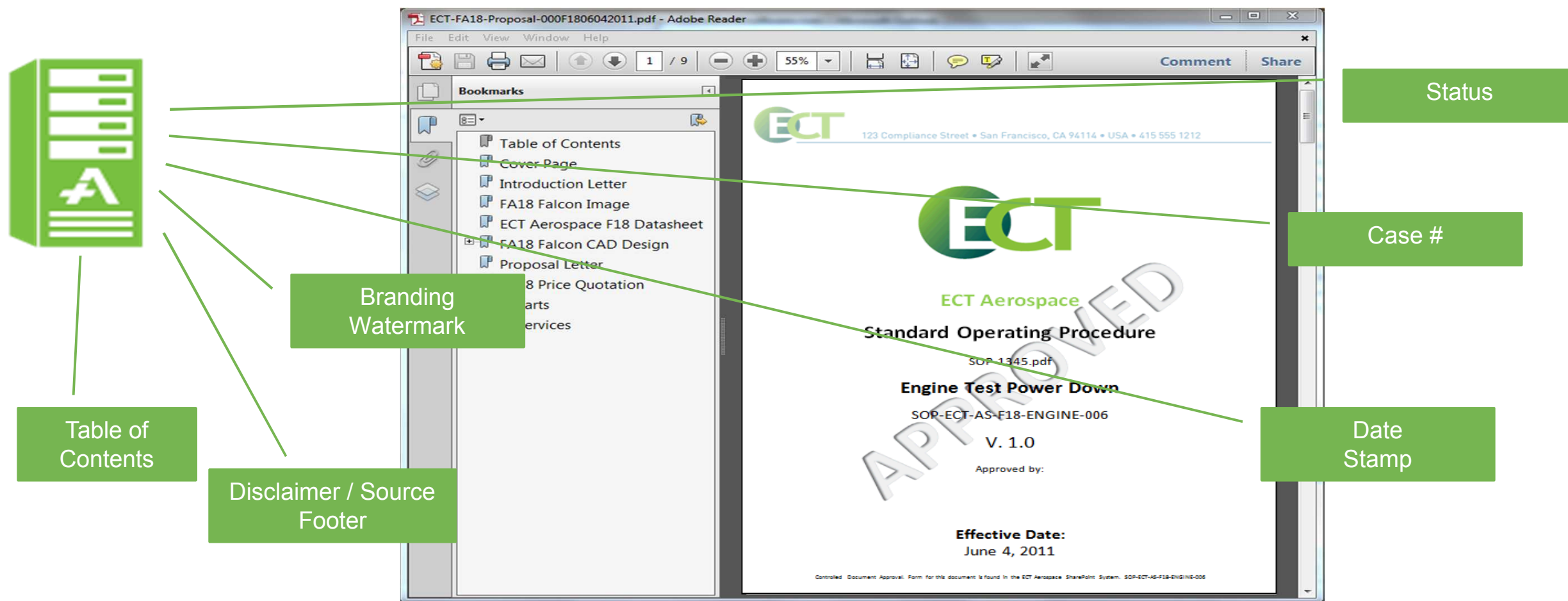
041000124

0000000000

0000010000



Metadata Extraction



Content Comparisson

[illegible]

The image shows a highly degraded and noisy scan of a document page. The text is mostly illegible due to heavy noise and low contrast. Three red circles are drawn around specific characters or words that are partially legible:

- A red circle at the top center highlights the word "WINONA".
- A red circle in the lower-left quadrant highlights the letter "U".
- A red circle in the lower-right quadrant highlights the letter "W".

The background is filled with a dense, chaotic pattern of black and white speckles and noise, making the rest of the text unreadable.

[illegible]

Applications of Image Analysis:

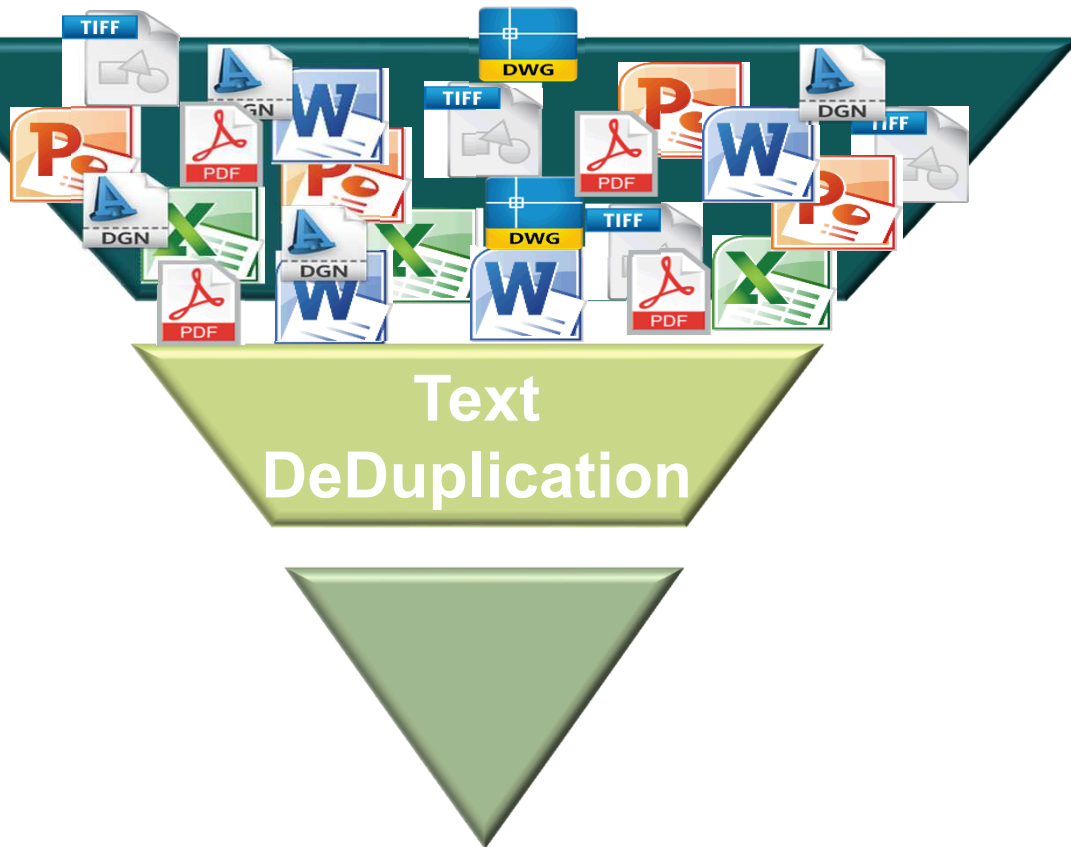
[illegible]

- * XML extractions
- * De-Duplication
- * Auto classification
- * Signature detection
- * Contract comparison
- * Revisions/versioning
- * Expiration management
- * Template confirmations

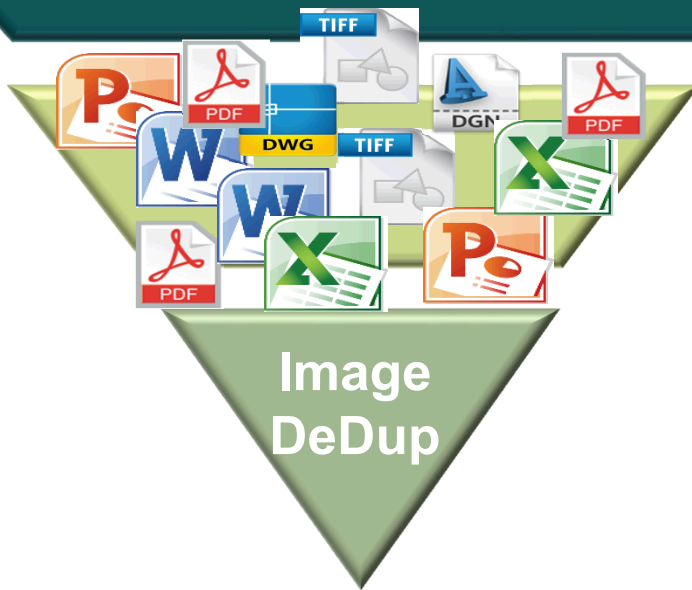


Hash DeDuplication

- Analyzes hash values of all files
- Duplicates identified & removed



- Compares text (natively)
- OCR Image only content
- Duplicates identified based on threshold & removed



- Pixel by pixel comparison
- Duplicates identified based on threshold & removed



Hash DeDuplication

Text
DeDuplication

Image
DeDup



Powered By:

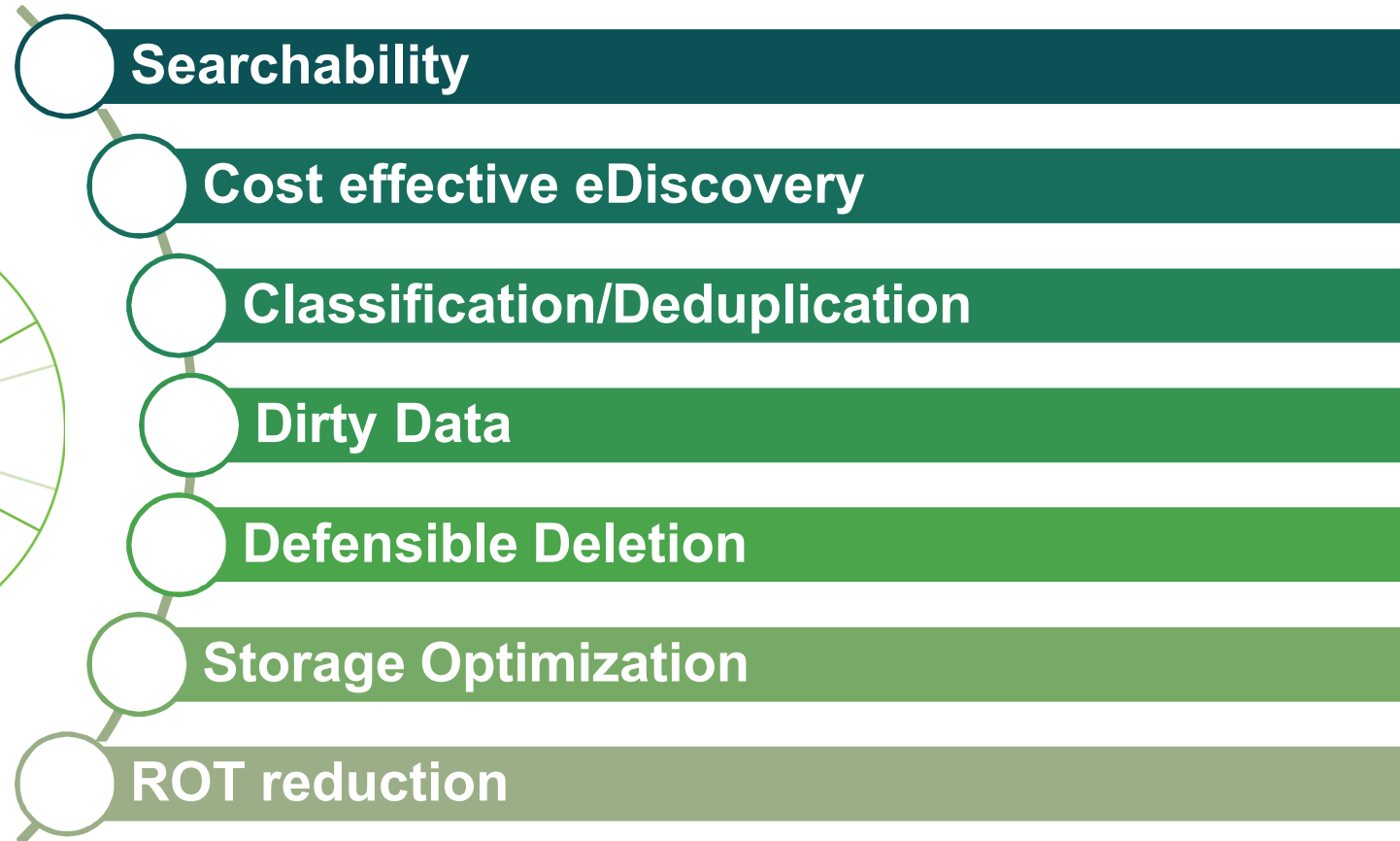
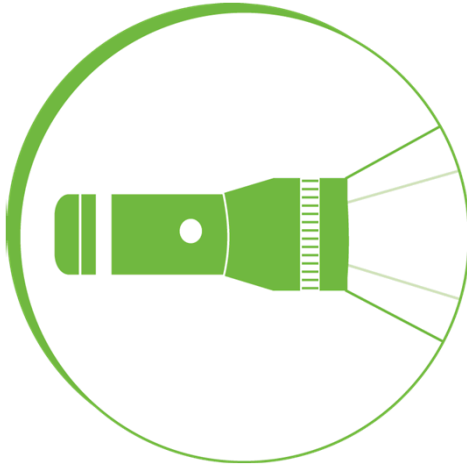


Document Lifecycle Overview





Leveraging the PDF Standard to Understand Dark Data and Improve Document Processes



A stage spotlight shining from the top left corner onto a dark, textured floor, creating a bright, circular pool of light. The background is dark and indistinct.

Thank You

Questions:
www.adlibsoftware.com