# Grooper™

# WHAT'S NEW IN 2.9

## UPCOMING FEATURES AND FUNCTIONALITY

28 July 2020

# WHAT'S NEW IN 2.9

- Data Integration Improvements
  - New MS Office Integration
  - New Box integration
  - Enhanced database export
  - New document viewer, supporting multiple document views
  - Data lookups can now be performed against CMIS sources
  - Improved OCR synthesis and paragraph detection
  - Classify, Extract, and Data review now support content type filtering
- Improvements for Architects
  - New Extractor weighting mechanisms
  - TF/IDF training builds training batch
  - Compile stats command for content containers
  - New TF/IDF variants for lexical classification
  - New weighted classification properties allow for new classification methodologies
  - New Inheritance Architecture
- Separation Improvements
  - Improved unstructured document separation training
  - Improved ESP auto-separation logic
  - New settings on document types to aid separation
  - New separation review UI
- New Data Annotation Option in data review
- Code expressions support LINQ and String Interpolation

# DATA INTEGRATION IMPROVEMENTS

# MS OFFICE INTEGRATION

Natively ingest multiple file formats generated by the MS office suite:
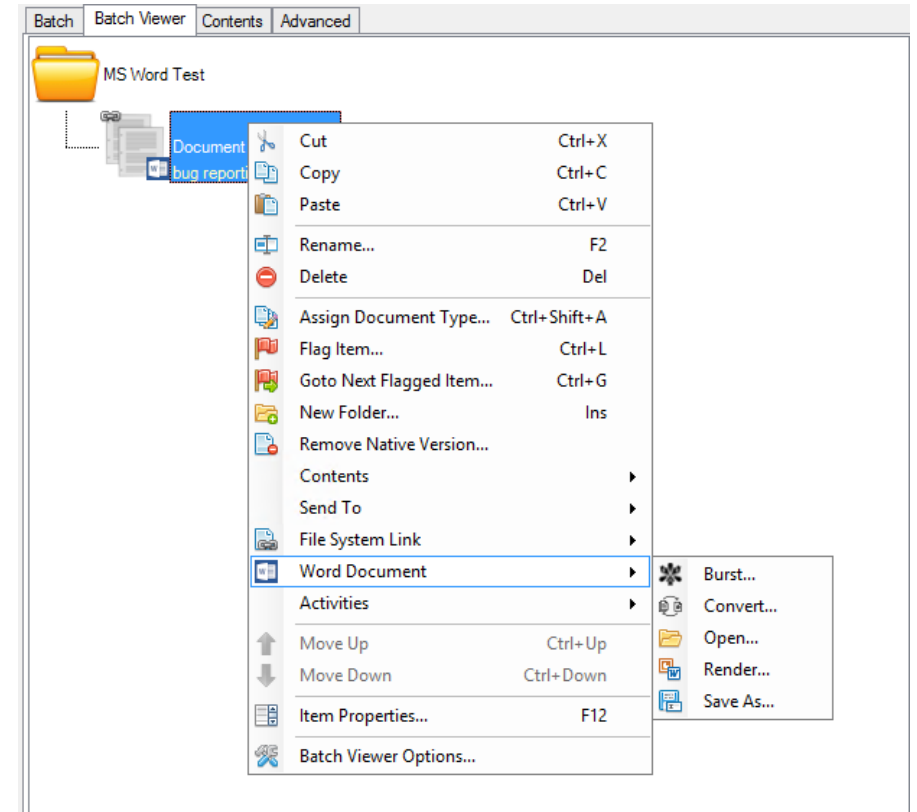
- Word document
- PowerPoint presentation
- Excel spreadsheet
- via CMIS, Outlook messages and mailboxes

Natively pull text from all supported filetypes

Removes the need for single-threaded rendering

MS office must be installed on machine running Design Studio

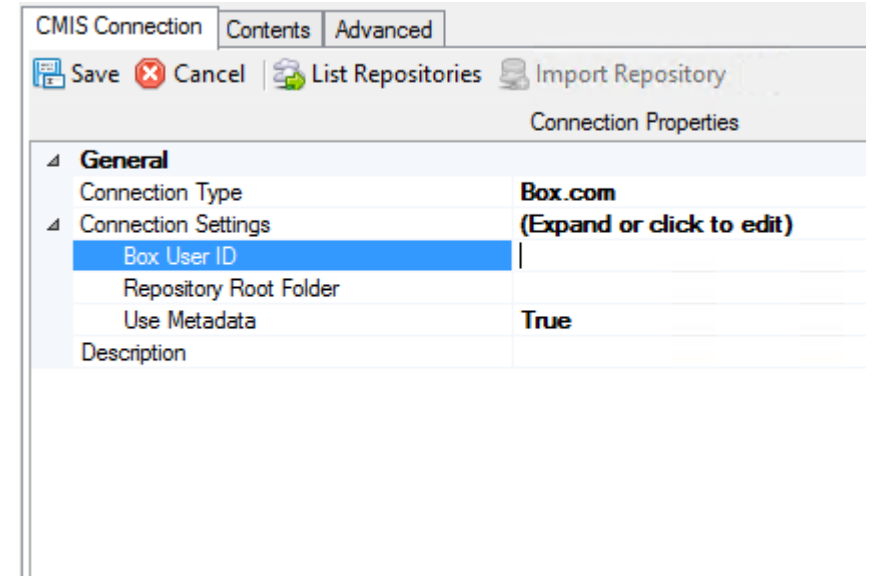Represents better integration with content production systems

# NEW BOX INTEGRATION

CONNECT DIRECTLY TO A BOX ACCOUNT TO IMPORT OR EXPORT FILES

WORKS WITH BOX BUSINESS OR ENTERPRISE

TAKES ADVANTAGE OF GROOPER'S EXISTING CMIS+ ARCHITECTURE, GIVING YOU ACCESS TO METADATA MAPPING AND EXPRESSIONS, CMIS LOOKUPS, ETC.

OUR FIRST-EVER INTEGRATION WITH A CLOUD ONLY PROVIDER; POSITIONS DEPLOYMENTS TO TAKE ADVANTAGE OF ONE OF THE LEADING STORAGE PROVIDERS

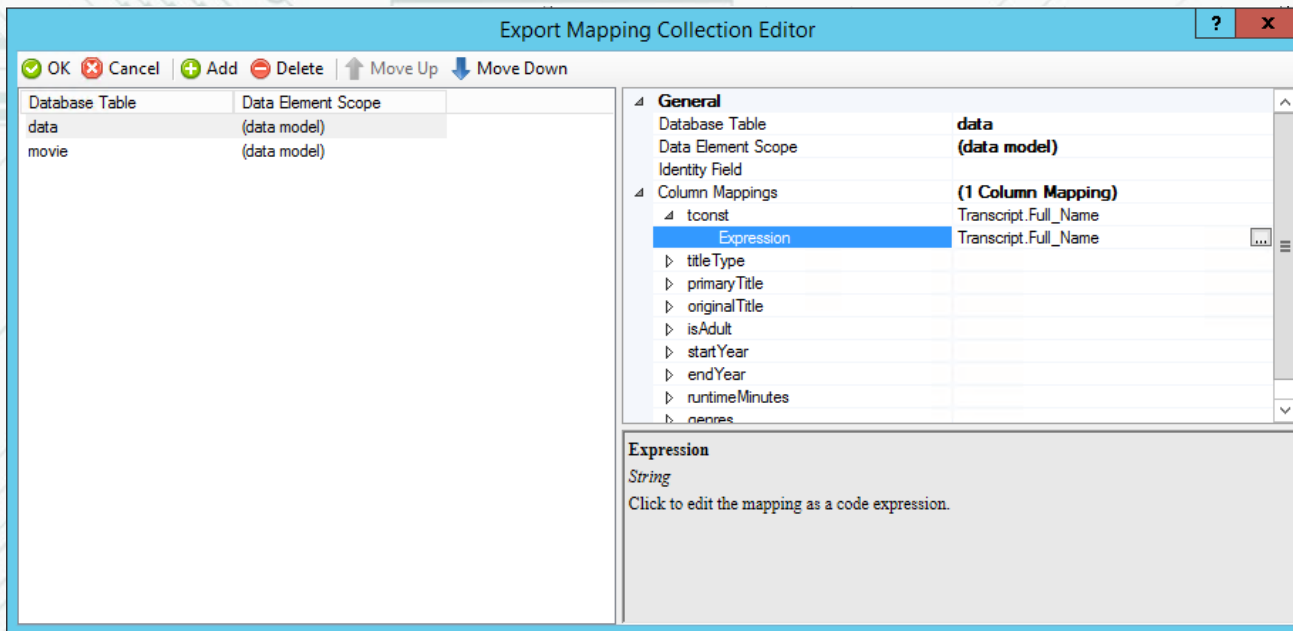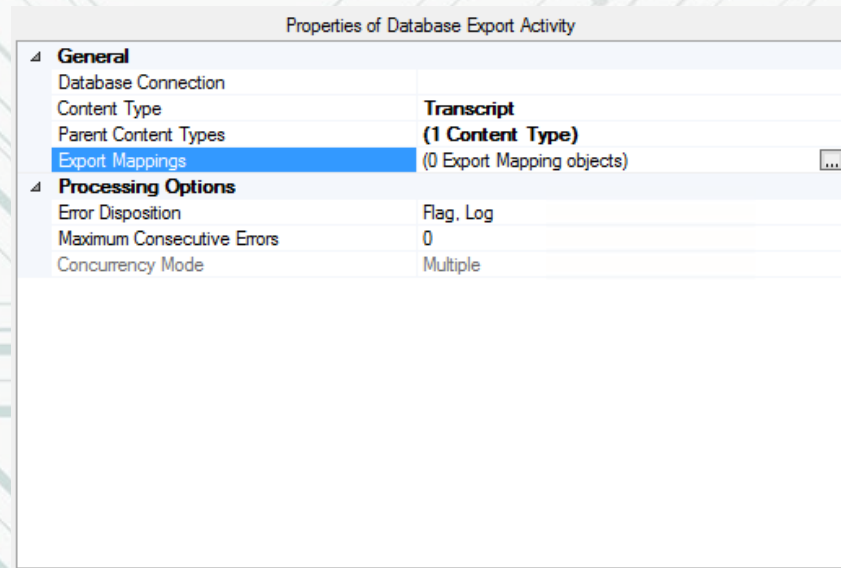CAN USE FOLDERS WITH GOVERNANCE LEVEL ASSIGNED TO TAKE ADVANTAGE OF GOVERNANCE FEATURES

# ENHANCED DATABASE EXPORT

Multiple exports can now be defined on a single database export step

- Must be within a single database; can span multiple tables
- Can choose separate data element scope on each table mapping object
- Greatly simplifies complicated or multipart database exports
- Supports SQL-server generated identity columns

## NEW DOCUMENT VIEWER

**ALLOWS FOR DISPLAY OF MULTIPLE DOCUMENT "RENDITIONS"**

Batch Folder | Contents | Advanced

Save | Cancel | Assign Document Type... | Flag Item... | Remove Native Version...

Native (Mail Message)

| General | |
| --- | --- |
| Content Type | |
| Flagged | False |
| Description | |
| **Attachment** | |
| File Name | Piers Plowman.eml |
| File Name Extension | .eml |
| MIME Type | message/rfc822 |
| File Type Description | Mail Message |
| MIME Type Info | Mail Message |
| Has Local Copy | True |
| Has PDF Version | False |
| Content Link | **CMIS Document Link** |

**Batch Folder**

A Batch Folder is a container object within the hierarchical structure of a Batch, which is used to represent logical folders and multipage documents.

**Remarks**

Batch Folders may contain Batch Page objects or other folders as children.

Batch Folder objects may be created in a variety of ways. If a Separation Profile is assigned in the Scan activity, then folders will be created in real time during scanning. The Separate

Piers Plowman.msg (2)
Piers Plowman.eml

**From:** Chris Dearner <cdearner@bisok.com>

**Sent:** Monday, October 7, 2019 2:09 PM

**To:** Chris Dearner <cdearner@bisok.com>

**Subject:** Piers Plowman

IN a somer seson,
Whan softe was the sonne
I shoop me into shroudes
As I a sheep weere
In habite as an heremite
Unholy of werkes,
Wente wide in this world
Wondres to here;
Ac on a May morwenynge
On Malverne hilles
Me bifel a ferly,
Of fairye me thoghte.
I was wery for-wandred,
And wente me to reste
Under a brood ban
By a bournes syde;
And as I lay and lenede,
And loked on the watres,
I slombred into a slepyng,
It sweyed so murye.

Thanne gan I meten
A merveillous swevene,
That I was in a wildernesse,
Wiste I nevere where,
And as I biheeld into the eest
An heigh to the sonne,
I seigh a tour on a toft
Trieliche y-maked,
A deep dale bynethe,
A dongeon therinne,
With depe diches and derke
And dredfulle of sighte.
A fair feeld ful of folk
Fond I ther bitwene,
Of alle manere of men,
The meene and the riche,
Werchynge and wandrynge,
As the world asketh.

# NEW VIEWER: ALTERNATIVE RENDITION

# CMIS DATA LOOKUPS

Functions similarly to a Database lookup and allows for population or validation of Grooper fields based on metadata located on CMIS objects

- Allows for better integration with CMIS data sources
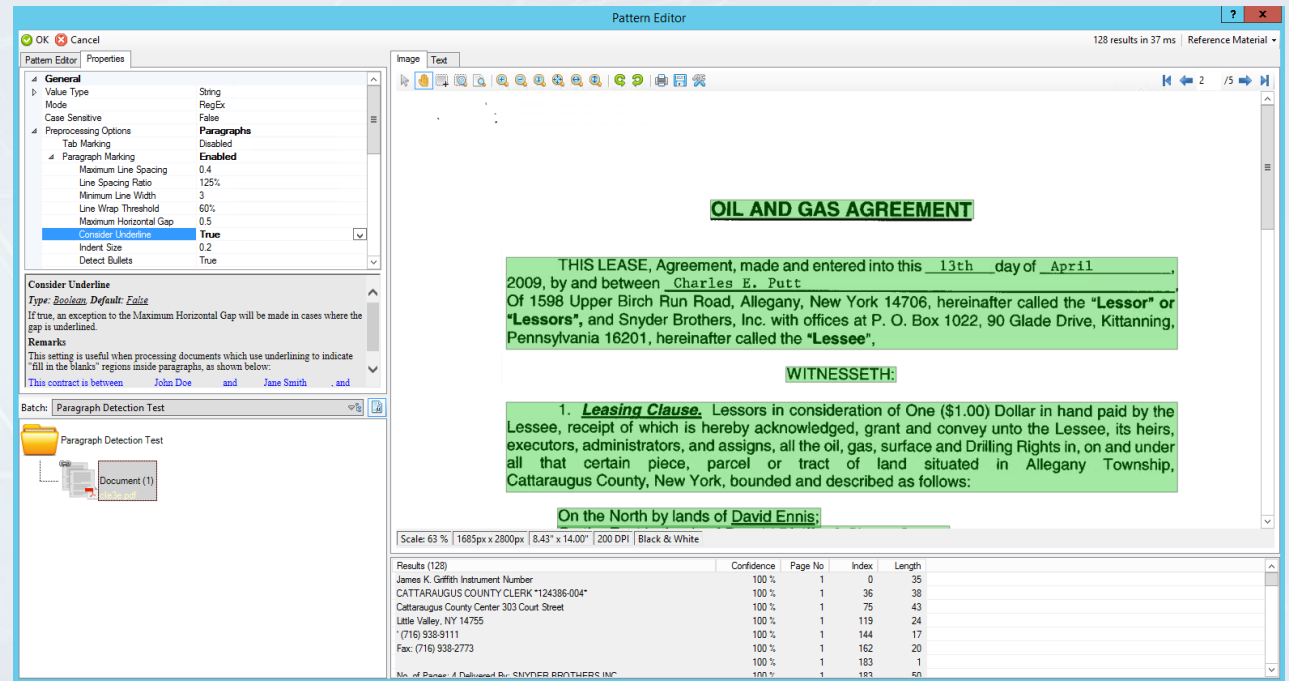- Allows for treating CMIS data sources as containing queryable information

# IMPROVED OCR SYNTHESIS AND PARAGRAPH DETECTION

Grooper can now detect paragraphs containing large sections of blank underlines correctly

There are additional various improvements to OCR Synthesis

IMPROVEMENTS FOR ARCHITECTS

# COMPILE STATS

COMPILE STATS PROVIDES COMPREHENSIVE STATISTICS ON CLASSIFICATION AND EXTRACTION ACROSS AN ENTIRE CONTENT MODEL

ASSISTS ADMINISTRATORS IN DEVELOPING AND TROUBLESHOOTING CONTENT MODELS, ALLOWING FOR MORE TARGETED WORK AND SHORTER DEVELOPMENT TIMES.

**Index Navigator Settings**

| | |
|---|---|
| **General** | |
| Flag Invalid Items | False |
| Processing Level | Level1 |
| Display Parent Folder(s) | False |
| Auto-Load Next Invalid Document | False |
| Content Type Filter | (0 Content Type objects) |
| Flag Messages | (empty) |

**Content Type Filter**

*Type: List of Content Type*

An optional list of Content Types to which the review process should be restricted.

**Remarks**

When this list is empty, all documents in scope will be included in the review process. If a list is provided, only document types matching one of the list entries will be included.

**Properties of Extract Activity**

| | |
|---|---|
| **General** | |
| Mode | Normal |
| Default Content Type | |
| Flag Invalid Items | False |
| Content Type Filter | (0 Content Type objects) |
| **Processing Options** | |
| Error Disposition | Flag, Log |
| Maximum Consecutive Errors | 0 |
| Concurrency Mode | Multiple |

**Properties of Classify Activity**

| | |
|---|---|
| **General** | |
| Content Model Scope | **Blump University** |
| Classification Level | DocType |
| Output Level | DocType |
| Apply To | All |
| Reclassify Mode | Overwrite |
| Supress Candidate List | False |
| Model Refresh Interval | 60 |
| **Processing Options** | |

# CONTENT TYPE FILTERING

**SUPPORTED IN EXTRACTION, CLASSIFICATION, AND DATA REVIEW ACTIVITIES**

- Appears in "index navigator settings" on Data Review Activity
- Allows for more focused or partial classification, extraction, and review
- Allows classification, extraction, and review to proceed in stages for larger or more complicated projects
- Supports targeted expansion or improvement
- Supports tiered classification decisioning
- Allows for classification to target only unclassified documents ("All," "Classified", or "Unclassified")

# NEW EXTRACTOR WEIGHTING MECHANISMS

**CONFIGURABLE OUTPUT CONFIDENCE AND OUTPUT MUTIPILIERS ON DATA TYPES**

- Manually set extraction confidence and multiply confidence by a set amount
- Allows for waterfall extraction techniques: most reliable methods can be backed up by less-reliable methods by using "order by confidence"

# WEIGHTED CLASSIFICATION PROPERTIES AND WATERFALL CLASSIFICATION

**ADVANCED CLASSIFICATION STRATEGIES FOR COMPLICATED DOCUMENT SETS**

- Positive extractors now classify documents with the confidence of their referenced extractor

- Using the confidence multiplier and output confidence configuration items on data types allows for the creation of waterfall classification strategies

- This allows for documents that are highly lexically dissimilar within a group or highly lexically similar across groups to be targeted with better accuracy (e.g. legal documents, oil and gas leases vs lease acknowledgements, etc.)

# AUTOBUILT TF/IDF TRAINING BATCH

TF/IDF training now automatically builds training batch, allowing for easier training after changes to feature collection or TF/IDF settings.

"Rebuild training" allows for tuning and A/B testing using identical training set and document training decisions

# NEW TF/IDF VARIANTS

- 'Use confidence' allows extraction confidence to factor into TF/IDF feature scoring

- 'Augmented' TF mode (in addition to logarithmic/'sublinear') allows for TF/IDF to partially account for document length (standard TF/IDF tends to favor longer documents)

- Smooth TF/IDF selectable under "Document Frequency Mode"

---

**Use Confidence**

*Type: Boolean, Default: False*

Specifies how the confidence of each feature instance should be calculated into the term frequency.

---

**Term Frequency Mode**

*Type: TfModes, Default: Normal*

Specifies how the term frequency (TF) of features should be calculated. Can be one of the following values:

- **Normal** - Term frequency is normalized to the size of the document. This mode provides matching which is completely independent of document size. In some classification scenarios, however, this mode can produce undesirable false positives, as it allows a very short document to match a very long document with a high degree of similarity.
- **Logarithmic** - Term frequency is scaled logarithmically. This mode is provided for backwards compatibility.
- **Augmented** - Term frequency is normalized to the frequency of the most-frequent feature. Provides an alternative to Normal mode where the size of the document plays more of a role in classification decisions. This mode also provides a Dampering Factor

---

**Document Frequency Mode**

*Type: IdfModes, Default: Normal*

Specifies how the inverse document frequency (IDF) of features should be calculated. Can be one of the following values:

- **Normal** - Standard IDF.
- **Smooth** - Smooth IDF.

# NEW INHERITANCE ARCHITECTURE

**DATA ELEMENT PROFILES ARE GONE; LONG LIVE DATA ELEMENT OVERRIDES**

- Data element profiles have been replaced by data element overrides
- Every property is now overridable
- Data Element Overrides tab includes built-in tester
- A fuller realization of Grooper's inheritance-based architecture; allows for more robust, more scalable data models.
- Saves time and reduces complexity when architecting large systems

# SEPARATION IMPROVEMENTS

# ESP AUTO-SEPARATION IMPROVEMENTS

**SEPARATION IS MORE ROBUST AND ACCURATE THAN EVER**

- Improved ESP auto-separation logic
- New Secondary Page Extractor
- Unstructured Documents now consider only first, middle, last pages; each n-page version is no longer considered a different formtype

New Separation review UI

SIMILAR TO ESP AUTO-SEPARATION TESTER

OTHER IMPROVEMENTS

# NEW DATA ANNOTATION OPTION IN DATA REVIEW

**DISPLAYS EXTRACTED DATA AT THE EXTRACTION LOCATION ON THE DOCUMENT**

- Greatly speeds up review
- Configurable properties include color, location
- Can disable for simple extraction cases or documents.

# CODE EXPRESSIONS: LINQ AND STRING INTERPOLATION SUPPORT IN EXPRESSIONS

**DO MORE WITH EXPRESSIONS**

- .NET string interpolation now supported in expressions
- LINQ now supported in expressions