# Paragraph Extraction

**Paragraph Extraction** can be used to extract a value that lies in a particular paragraph present over the document.

In Paragraph Extraction first a paragraph is identified and then from within that paragraph a particular value matching with a given regex pattern is extracted.

Paragraphs are identified on the basis of the following conditions:

- Regex match for start pattern is treated as the start of the paragraph only if there is no span(word) present to the left of found Regex Match.
- Ephesoft takes the average white space between lines and segregates the text body on the basis of white space being larger than the average space.
- If any line ends with the End pattern if defined, then it takes priority over the line spacing mechanism and the paragraphs end on that line even if the next lines satisfy the spacing condition.
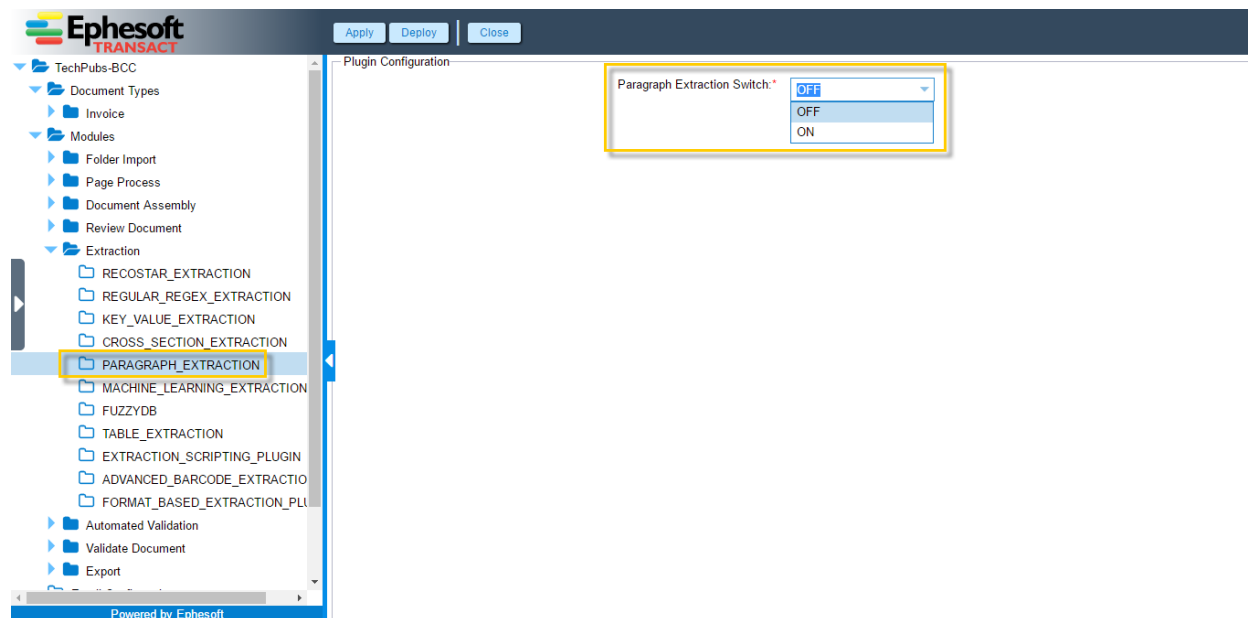
The start pattern for a paragraph can be a title of the paragraph or starting words of the paragraph. You can configure the extraction rule accordingly. During extraction, paragraph wrapping is handled by default while using **Paragraph Extraction Rule**.

This functionality enables you, as an administrator of batch classes to configure extraction rules for index fields.

# Configuration

The index field values for which **Paragraph Extraction** is configured are extracted using a plugin.

**PARAGRAPH_EXTRATION** Plugin governs the extraction of configured index field while using **Paragraph Extraction**.



This plugin has only one configuration which is a switch. If the value of the switch is set to **ON**, the configured index filed is extracted, else not. By default, the switch is set to **OFF**.

# Configuring Paragraph Extraction Rule for an Index Field

Configuration for **Paragraph Extraction Rule** is similar to any other extraction rule such as the KV Extraction Rule.
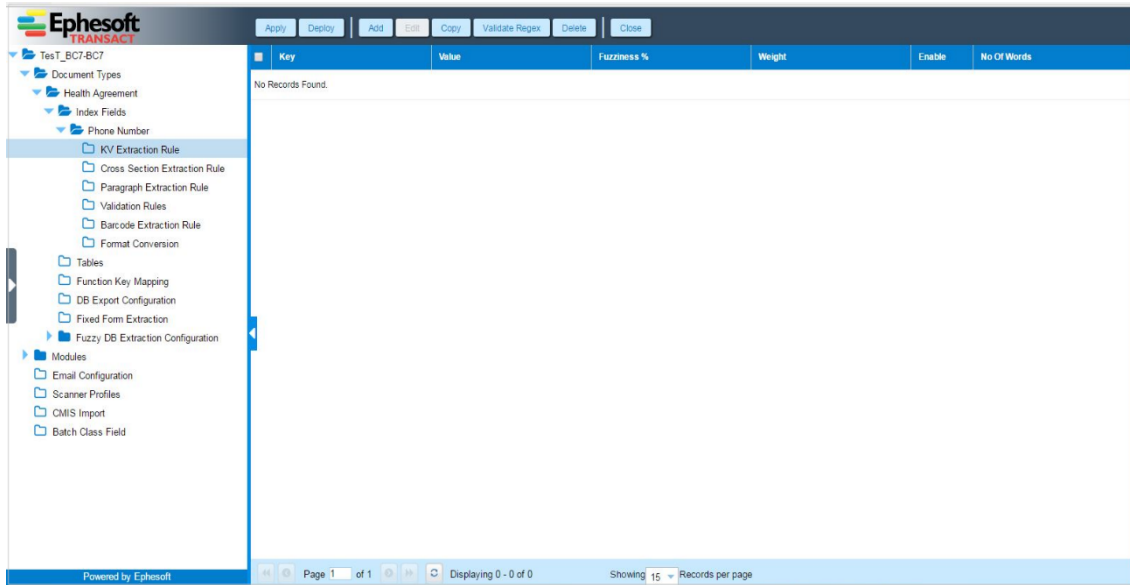
**To configure Paragraph Extraction Rule for an index field**

1. From the DCMA Home page, click **ADMINSITRATOR** and select **BATCH CLASS MANAGEMENT.**

   The Ephesoft Enterprise **Login** page displays.

2. Enter valid credentials to login.

   The **Batch Class Management** screen displays.

3. Select the batch class from the list in the **Batch Class Management** screen and click **OPEN.**

   The batch class opens with **Document Types** node selected by default.

4. Select the document type from the list and click **OPEN.**

The document type node expands displaying a list of index fields.
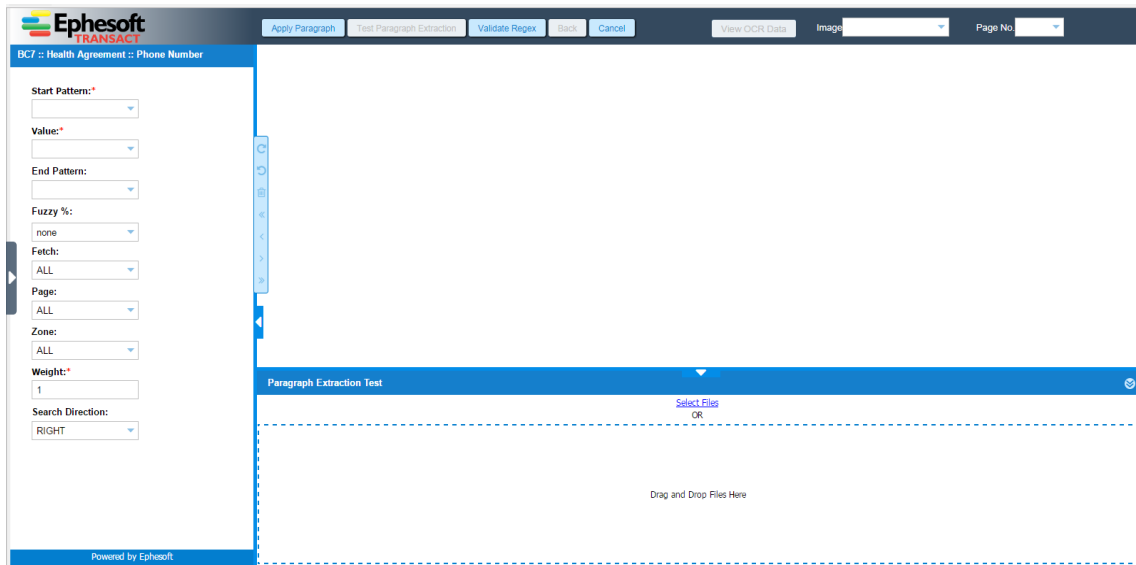
5. Select the index field from the list and click **OPEN.**

   The index field node expands displaying all the available extraction rules in the left navigation pane and **KV Extraction Rule** selected by default as shown in the image below.
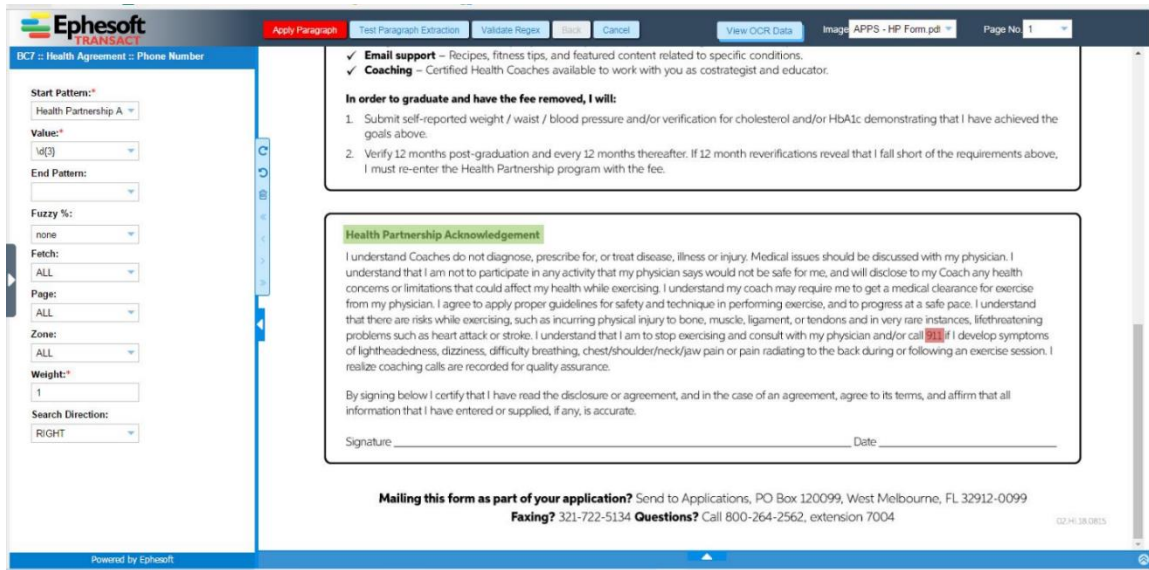


6. Select **Paragraph Extraction Rule** from the navigation pane and click **ADD.**

   The following screen displays.

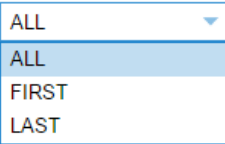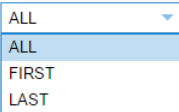7. Click **Select Files** link from **Paragraph Extraction Test** section and upload a image file.

The uploaded image is displayed in the image view pane.
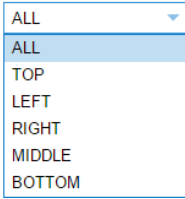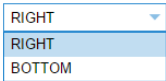


> While using **Paragraph Extraction Rule**, image overlays don't have any significance in terms of extraction. You can use them to create the regular expressions for **Start Pattern**, **Value and End Pattern**.

8. Enter the relevant configuration details as described in the table below:

| Component | Description |
|-----------|-------------|
| **Start Pattern** | This parameter is to configure a start regex pattern for paragraphs. Start of Paragraph is identified using this regex pattern.<br><br>You can enter a regular expression or use **Regex Builder**/**Regex Pool** options to enter a search pattern. |
| **Value** | This parameter is to find values within the identified paragraph.<br><br>You can enter a regular expression or use **Regex Builder**/**Regex Pool** options to enter a search pattern. |
| **End Pattern** | This is an optional parameter, and can be used to end a paragraph if any line of the paragraph ends with this pattern. |

| Component | Description |
|---|---|
| | **End Pattern:** <br> Regex Builder <br> Regex Pool <br><br> You can enter a regular expression or use **Regex Builder**/**Regex Pool** options to enter a search pattern. |
| **Fuzzy %** | You can use this parameter to do a fuzzy search while searching for Paragraphs using the Start Pattern. <br><br> **Fuzzy %:** <br> none <br> none <br> 10% <br> 20% <br> 30% |
| **Fetch** | Depending on the value you select for this parameter, the Value tag in the batch xml is populated with ALL/FIRST/LAST Regex matches in the paragraph. <br><br> **Fetch:** <br> ALL <br> ALL <br> FIRST <br> LAST <br><br> This parameter has the following options: |

| | | |
|---|---|---|
| **ALL** | All the values matching the regex are updated in batch.xml with space as a delimiter. | |
| **FIRST** | First value matching the regex going from left to right in the paragraph is updated in batch.xml. | |
| **LAST** | Last value matching the regex going from left to right in the paragraph is updated in batch.xml. | |

| Component | Description |
|---|---|
| **Page** | You have three options available to choose from for this parameter: ALL, FIRST, and LAST. <br><br> **Page:** <br> ALL <br> ALL <br> FIRST <br> LAST <br><br> Depending on the selected value, the extraction algorithm runs on ALL/FIRST/LAST Page(s) of the document. |
| **Zone** | Every page in divided into 5 zones: TOP, MIDDLE, BOTTOM, LEFT and RIGHT along with the default option of ALL. |

| Component | Description |
|---|---|
| | **Zone:**<br>ALL<br><br>ALL<br>TOP<br>LEFT<br>RIGHT<br>MIDDLE<br>BOTTOM<br><br>You can use this parameter to specify the portion of the page where the algorithm searches for start value of paragraph to extract it.<br><br>For example, if you configure this parameter value as BOTTOM, the start pattern of the paragraph is searched only in the BOTTOM zone. |
| **Weight** | You can use this parameter to implement weighted confidence values.<br><br>**Weight:***<br><br>1<br><br>This is used to give bias/weightage to a particular extraction rule. |
| **Search Direction** | You can use this option to select where to look for the paragraph i.e., on the RIGHT of the start pattern or BELOW the start pattern.<br><br>**Search Direction:**<br>RIGHT<br><br>RIGHT<br>BOTTOM<br><br>The default search direction is RIGHT. |

9. Click **Test Paragraph Extraction** from the toolbar on top of the page.

   The extraction result is highlighted on the image as an overlay and are also displayed in the

   **Paragraph Extraction Test** grid as shown in the image below.

10. Click **Apply Paragraph** to apply the rule to the index field.

---

If you click **Cancel** without saving changes, the following confirmation message displays.



Click **Save** to save changes or click **Discard** to discard any configuration changes and navigate to the **Paragraph Extraction Rule** screen.

---

# Property File Configuration

The property file for paragraph extraction is dcma-paragraph-extraction.properties and is available at:

**Ephesoft\Application\WEB-INF\classes\META-INF\dcma-paragraph-extraction\dcma-paragraph-extraction.properties**

| Configurable property | Type of value | Value options | Description |
|---|---|---|---|
| **paragraph.endParagrpahSpacesSwitch** | String | ON\|OFF<br>**Default: ON** | This switch governs whether to end the paragraph at current line, if the current line length is less than threshold*average line length of the current paragraph. |
| **paragraph.endParagrpahSpacesThreshold** | Float | Between 0 and 1<br>**Default: .85** | Threshold to be used for endParagraphLineSwitch. |