

Configuring Server-side Extraction



Applies To

Software	Version
Acrolinx Server	4.5, 4.6

Server-side extraction is the process of configuring the server to extract text instead of the Acrolinx client software. Server-side extraction is useful for customized Acrolinx integrations where there's no user interface for configuring text extraction settings. For example, Acrolinx integrations for rich text editors such as CKEditor use server-side extraction.

You define server-side extraction settings by adding context segmentation definition (CSD) files to your server installation and updating your language configuration file.

To configure server-side extraction, follow these steps:

1. [Create a server-side CSD file](#)
2. [Associate the server-side CSD file with a language or rule set](#)
3. *(Optional)* [Configure Acrolinx Server to generate an annotated copy of the checked document](#).
4. [Enable server-side extraction in your client software](#)

Restriction: Server-side extraction should not be used when the server is configured to save Acrolinx Scorecards on the client computer. For more information about this setting, see the [Acrolinx Server Administration Guide](#).

Context Segmentation Definitions (CSDs)

What is the Purpose of CSDs?

CSDs provide Acrolinx with instructions on how to get the text out of a certain kind of document. These instructions fall into two basic categories: segmentation and filtering.

When you check HTML or XML files, the settings in a CSD file determine the segmentation rules.

The CSD contains the following settings:

- The document types that are associated with the defined segmentation settings.
- Sentence-break elements.
- No-break elements.
- Elements to include or exclude when a check is performed.
- Attributes that have values which should be checked.

CSD File Naming and Location

CSDs are contained in files that have the following naming format:

<FILENAME>.properties.

Server-side CSD files usually are located in:

- %ACROLINX_CONFIGURATION_ROOT%\data\ on Windows.
- \$ACROLINX_CONFIGURATION_ROOT/data/ on Unix-based operating systems.

Creating and Modifying CSDs

To create or modify a CSD file, follow these steps:

1. If you don't have any CSD files, create one first

- a. Create a text file <FILENAME>.properties and save the file in your server configuration directory.

Example: %ACROLINX_CONFIGURATION_ROOT%\data\ <LANGUAGE_ID>

- b. Copy the following lines into the text file:

```
### General Information

csd_name=
version=1

### Matching Criteria

dtd.type=
dtd.source_reference=
dtd.client_signature=
dtd.public_id=
dtd.system_id=
dtd.root_element=
dtd.schema=

### Extraction Settings

# General Settings

ignoreCase=true
default_inclusion_mode=include

# Filtering
exclusion_elements=
inclusion_elements=
empty_elements=

# Break elements
no_break_elements=
sentence_break_elements=
token_break_elements=
default_break_level=

# Attributes
extract_attributes=
meta_information=

# Parenthetic elements
parenthetic_elements=
```

2. If you have an existing CSD that you'd like to edit, find the directory where your CSD profiles are stored and open the file:

Example: %ACROLINX_CONFIGURATION_ROOT%\data\ common\myCSD.properties

3. Update the values for each of the settings.

If you're not sure what some of the settings are for, see the [descriptions of each setting](#).

You can also use the following sample CSD to see how you could enter the different settings.

Note: We've provided a sample value for each of the matching criteria, but in reality you don't need to set all these criteria. The `dtd.type` is required but the others are optional.

```
### General Information

csd_name=My Test CSD
version=1

### Matching Criteria

dtd.type=XML
dtd.source_reference=.* /user-guides/. *.dita
dtd.language=en
dtd.client_signature=mF0YZzqQ2331Y2tlcg
dtd.public_id=--//OASIS//DTD DITA Concept//EN
dtd.system_id=concept.dtd
dtd.root_element=concept
dtd.schema=notes.xsd

### Extraction Settings

# General Settings
ignoreCase=true
default_inclusion_mode=include
mark_excluded_elements=false

# Filtering
exclusion_elements=metadata,cndname,codeph
inclusion_elements=ph,entry,title
empty_elements=img,xref

# Break elements
no_break_elements=sub,sup
sentence_break_elements=title,entry,li
token_break_elements=cr
default_break_level=sentence

# Parenthetic elements
parenthetic_elements=footnote,endnote

# Attributes
extract_attributes=img.?alt
meta_information=MyDocumentID:<book id>
```

4. Save your changes and if you haven't done so already, [associate the CSD with a rule set](#).

Configuring the Matching Criteria and General Settings for a CSD

Matching Criteria

The properties in a CSD file fall into two main categories: the matching criteria and the extraction settings themselves.

When you create a CSD, it's usually intended for a specific type of document. For example, you might create separate CSDs for help tasks and localization files. You use the matching criteria to match the appropriate documents with your CSD. Usually, you only need to define two or three matching criteria.

Use the following table to decide which matching criteria best suits your requirements.

Table 1. Matching Criteria and General Settings

Setting	csd_name
Description	A name that helps to differentiate the CSD from other CSDs in case you're reviewing the log files for CSD behavior.
Values	A unique name. You can give your CSD any name you like. <code>csd_name=CSD for Me</code>

Setting	version
Description	This is an internal version number that indicates how new the CSD is. The server uses this information to interpret the CSD as a legacy CSD type or a new CSD type.
Values	A number. In this case, the number 1. Just enter the number 1 as follows. <code>version=1</code> That's it, job done.
Setting	dtd.type
Description	Describes the type of file that was checked. The format could be text, html or xml.
Values	One of following values: HTML, XML, or WORD_XML To assign settings to XML files, you would enter the <code>dtd.type</code> identifier as follows: <code>dtd.type=XML</code>
Setting	dtd.source_reference
Description	The location of the file or files that the CSD should apply to. The server examines full path to the checked file. This path is sent by the Acrolinx plug-ins. You might use this property when there is no other data available to identify the file type.
Values	A file path or a path pattern entered as a regular expression. For example, to have the CSD apply to all DITA files in the "user-guides" directory, you would enter the property as follows: <code>dtd.source_reference=.*user-guides/*.dita</code> If you had multiple directories such as "user-guides" and "faqs" you would separate them with the OR operator " " like this: <code>dtd.source_reference=.*user-guides/*.dita .*faqs/*.dita</code> CAUTION: If you want to use more complex regular expressions, there are a few things you should look out for. See our detailed section on regular expressions in matching criteria .
Setting	dtd.language
Description	The language of the check as indicated by the Acrolinx plug-in. You can use this property to change the extraction settings based on the language of the text that's being checked. One scenario might be that you're checking a multilingual XML format such as TMX. You would want to include the parts of the file that match the checking language.
Values	A shortcode for the language which is usually two letters such "en" or "de". For example, to target the German segments in a TMX file, you would enter the property as follows: <code>dtd.language=de</code> In the same CSD, you might configure the inclusion elements property like this: <code>inclusion_elements=<tuv lang=de></code>
Setting	dtd.client_signature
Description	The signature that an Acrolinx plug-in sends when it authenticates with an Acrolinx server. You might use this identifier for customized Acrolinx integrations such as CMS integrations that don't send a plug-in short name. CAUTION: Only use this matching criterion if no other matching criteria work for you. Matching on the signature means that the CSD will only match when writers check from a certain editor. That's not really how CSDs are supposed to match but it can be useful if the host editor is the only way of determining the document type.
Values	A signature string. For example, to assign settings to an Acrolinx integration that uses the signature "xmF0YZzgQ233lY2tlcg", you would enter the <code>clientSignature</code> element as follows: <code>dtd.client_signature=xmF0YZzgQ233lY2tlcg</code>
Setting	dtd.public_id

Description	The public ID of an XML document type.
Values	A public ID string. For example, to assign settings to DITA concepts, you would enter the <code>publicId</code> identifier as follows: <pre>dtd.public_id=-//OASIS//DTD DITA Concept//EN<</pre>
Setting	<code>dtd.system_id</code>
Description	The system ID of an XML document type.
Values	A system ID string. For example, to assign settings to DITA concepts, you would enter the <code>systemId</code> identifier as follows: <pre>dtd.system_id=concept.dtd</pre>
Setting	<code>dtd.root_element</code>
Description	The root element of an XML or HTML document. You might use this identifier for simple XML documents that don't have a document type definition or schema.
Values	An element name. For example, to assign settings to XML files that begin with the element <code><product></code> , you might enter the <code>rootElement</code> element as follows: <pre>dtd.root_element=product</pre> You can also use this identifier to assign settings to HTML files. The root element for HTML documents is always <code>html</code> .
Setting	<code>dtd.schema</code>
Description	The name of an XML document schema.
Values	A schema name. For example, to assign settings to XML files that use the schema "notes", you would enter the <code>schemaName</code> identifier as follows: <pre>dtd.schema=notes.xsd</pre>

Table 1. Matching Criteria and General Settings

Setting	Description	Values
<code>csd_name</code>	A name that helps to differentiate the CSD from other CSDs in case you're reviewing the log files for CSD behavior.	A unique name. You can give your CSD any name you like. <pre>csd_name=CSD for Me</pre>
<code>version</code>	This is an internal version number that indicates how new the CSD is. The server uses this information to interpret the CSD as a legacy CSD type or a new CSD type.	A number. In this case, the number 1. Just enter the number 1 as follows. <pre>version=1</pre> That's it, job done.
<code>dtd.type</code>	Describes the type of file that was checked. The format could be text, html or xml.	One of following values: <code>HTML</code> , <code>XML</code> , or <code>WORD_XML</code> To assign settings to XML files, you would enter the <code>dtd.type</code> identifier as follows: <pre>dtd.type=XML</pre>

Setting	Description	Values
dtd.source_reference	<p>The location of the file or files that the CSD should apply to.</p> <p>The server examines full path to the checked file. This path is sent by the Acrolinx plug-ins.</p> <p>You might use this property when there is no other data available to identify the file type.</p>	<p>A file path or a path pattern entered as a regular expression.</p> <p>For example, to have the CSD apply to all DITA files in the "user-guides" directory, you would enter the property as follows:</p> <pre>dtd.source_reference=.* /user-guides/. *.dita</pre> <p>If you had multiple directories such as "user-guides" and "faqs" you would separate them with the OR operator " " like this:</p> <pre>dtd.source_reference=.* /user-guides/. *.dita .*/faqs/. *.dita</pre> <p>CAUTION: If you want to use more complex regular expressions, there are a few things you should look out for. See our detailed section on regular expressions in matching criteria.</p>
dtd.language	<p>The language of the check as indicated by the Acrolinx plug-in.</p> <p>You can use this property to change the extraction settings based on the language of the text that's being checked.</p> <p>One scenario might be that you're checking a multilingual XML format such as TMX. You would want to include the parts of the file that match the checking language.</p>	<p>A shortcode for the language which is usually two letters such as "en" or "de".</p> <p>For example, to target the German segments in a TMX file, you would enter to the property as follows:</p> <pre>dtd.language=de</pre> <p>In the same CSD, you might configure the inclusion elements property like this:</p> <pre>inclusion_elements=<tuv lang=de></pre>
dtd.client_signature	<p>The signature that an Acrolinx plug-in sends when it authenticates with an Acrolinx server.</p> <p>You might use this identifier for customized Acrolinx integrations such as CMS integrations that don't send a plug-in short name.</p> <p>CAUTION: Only use this matching criterion if no other matching criteria work for you.</p> <p>Matching on the signature means that the CSD will only match when writers check from a certain editor. That's not really how CSDs are supposed to match but it can be useful if the host editor is the only way of determining the document type.</p>	<p>A signature string.</p> <p>For example, to assign settings to an Acrolinx integration that uses the signature "xmF0YZzgQ233lY2tlcg", you would enter the <code>clientSignature</code> element as follows:</p> <pre>dtd.client_signature=xmF0YZzgQ233lY2tlcg</pre>
dtd.public_id	<p>The public ID of an XML document type.</p>	<p>A public ID string.</p> <p>For example, to assign settings to DITA concepts, you would enter the <code>publicId</code> identifier as follows:</p> <pre>dtd.public_id=-//OASIS//DTD DITA Concept//EN<</pre>
dtd.system_id	<p>The system ID of an XML document type.</p>	<p>A system ID string.</p> <p>For example, to assign settings to DITA concepts, you would enter the <code>systemId</code> identifier as follows:</p> <pre>dtd.system_id=concept.dtd</pre>
dtd.root_element	<p>The root element of an XML or HTML document.</p> <p>You might use this identifier for simple XML documents that don't have a document type definition or schema.</p>	<p>An element name.</p> <p>For example, to assign settings to XML files that begin with the element <code><product></code>, you might enter the <code>rootElement</code> element as follows:</p> <pre>dtd.root_element=product</pre> <p>You can also use this identifier to assign settings to HTML files. The root element for HTML documents is always <code>html</code>.</p>
dtd.schema	<p>The name of an XML document schema.</p>	<p>A schema name.</p> <p>For example, to assign settings to XML files that use the schema "notes", you would enter the <code>schemaName</code> identifier as follows:</p> <pre>dtd.schema=notes.xsd</pre>

You can enter a regular expression as the value for any matching criteria. Regular expressions are especially useful for the criterion `dtd.source_reference` because you can use them to define multiple file locations.

When you use regular expressions, bear in mind that you have to use the syntax for [Java-style regular expressions](#).

If you enter Windows-style paths, you must also escape the backslashes. In the Java-style regular expressions, you must escape backslashes twice. For example, suppose that your files are stored in the master directory `C:\DOC\user-guides\`. To match all DITA files in the directory you would enter the regular expression as follows:

```
dtd.source_reference=C:\\\\DOC\\\\user-guides\\\\.*.dita
```

Because of the idiosyncrasies of Java-style regular expressions, you have to replace each single backslash "\" with four backslashes "\\\\".

Configuring Extraction Settings in a CSD

After you're familiar with the different ways you can categorize elements in Acrolinx, you can start adding these elements to your extraction settings.

If you want to change the text extraction for elements that have specific attributes, you can add the attribute to your settings as well. For example, you could categorize `p` elements as excluded only when they have the value "internal" in the `audience` attribute. You add attributes to your CSD in a [specific syntax](#).

Use the following table as a guide to what you can configure in a CSD and how you add the extraction settings correctly:

Table 2. Extraction Settings

Setting	<code>ignoreCase</code>
Description	<p>If your element names appear in different cases you might want to turn this off.</p> <p>By default, Acrolinx ignores the casing of element names when assessing the extraction settings.</p> <p>For example, if you configure the element name <code>title</code> to be excluded, Acrolinx also matches the elements <code>TITLE</code> and <code>Title</code>.</p>
Values	<p>Boolean value.</p> <p><i>Example:</i> <code>true</code>, <code>false</code></p> <p>Suppose that you have an element such as <code>Data</code> and it means something different when it takes the uppercase form <code>DATA</code>.</p> <p>In this case, you would configure the extractions settings to be case-sensitive by entering the property as follows:</p> <pre>ignoreCase=false</pre>
Setting	<code>exclusion_elements</code>
Description	<p>If you don't want Acrolinx to check the text in an element at all, define it as excluded.</p>
Values	<p>Element names, element names with attributes.</p> <p><i>Example:</i> <code>metadata</code>, <code>cndname</code>, <code>codeph</code></p> <p>For example, you probably don't want to check the text in elements that contain code or identifiers.</p> <pre><metadata>document ID12345</metadata></pre> <p>Note: Excluded XML elements that are nested within included elements are still excluded. All other excluded items are ignored.</p> <p>In the following example, the <code><note></code> element is included and the <code><ph></code> element is excluded.</p> <pre><note>This text will be checked.<ph>This text will not be checked.</ph></note></pre> <p>The content in <code><ph></code> will still be excluded, even if it's contained within an included element.</p>
Setting	<code>inclusion_elements</code>
Description	<p>Use the property to define elements that are always included regardless of whether they're nested inside an excluded element.</p> <p>"Child" elements of inclusion elements are also always included.</p>
Values	<p>Element names, element names with attributes.</p> <p>For example, to always include <code>title</code> elements in a check, add it to your list of inclusion elements like this:</p> <pre>inclusion_elements=ph, entry, title</pre>
Setting	<code>sentence_break_elements</code>
Description	<p>If you have elements that contain sentences that don't end with a period, define them as sentence-break elements. This means that the end of the element should always be treated as a sentence end.</p>

Values	<p>Element names, element names with attributes. A very typical example is the title element:</p> <pre><title>This is a title</title></pre> <p>Unexpected <code>sentence_too_long</code> flags which actually consist of several sentences are a clear sign that one or more elements must be defined as a sentence break.</p> <p>In this case, you would add <code>title</code> to the list of sentence-break elements like this:</p> <pre>sentence_break_elements=title,entry,li</pre>
Setting	<code>no_break_elements</code>
Description	If you have elements that split a word into separate fragments, define it as a "no break" element.
Values	<p>Element names, element names with attributes. A typical example is the subscript element:</p> <pre>H<sub>2</sub>O concentration</pre> <p>If <code><sub></code> isn't classified, the incorrect text</p> <pre>H 2 O concentration</pre> <p>is sent to the server. If <code><sub></code> is a no break element, the correct text</p> <pre>H2O concentration</pre> <p>is sent to the server.</p> <p>In this case, you would add <code>sub</code> to the list of no break elements like this:</p> <pre>no_break_elements=sub,sup</pre>
Setting	<code>token_break_elements</code>
Description	<p>You might define elements as token breaks if you have words which aren't separated by a space.</p> <p>This setting defines the elements which should cause a token break. Adding a token break is basically adding a boundary between words.</p>
Values	<p>Element names, element names with attributes.</p> <p>Usually, you would add a space between words, but sometimes editors separate words with elements rather than spaces.</p> <p>For example, WordProcessingML uses <code><cr/></code> elements to indicate manually inserted line breaks, like this:</p> <pre>This sentence is<cr/>wrapped manually.</pre> <p>You wouldn't want the sentence to end at "is", or have the words interpreted as "iswrapped". So you would update the setting as follows:</p> <pre>token_break_elements=cr</pre>
Setting	<code>default_break_level</code>
Description	<p>This setting defines how Acrolinx should interpret elements that aren't already included in the "break element" settings or in any of the other settings.</p> <p>The following settings are "break element" settings:</p> <ul style="list-style-type: none"> <code>no_break_elements</code> <code>sentence_break_elements</code> <code>token_break_elements</code>
Values	<p>Enter one of the following values:</p> <ul style="list-style-type: none"> <code>none</code> Acrolinx doesn't add any break at all. <code>token</code> Acrolinx adds a token break. <code>sentence</code> Acrolinx adds a sentence break. <p>The default value for this setting is <code>token</code>.</p> <p>For example, if you wanted to have any undefined elements treated as sentence breaks, you would enter the property as follows:</p> <pre>default_break_level=sentence</pre>

Setting parenthetic_elements

Description Suppose that you have elements that surround fragments of text in the same way that you would put a side-note in parentheses in the middle of a sentence.

Example:

```
<p>This is the <footnote>This is a second sentence.</footnote> first sentence.</p>
```

This type of element is fairly rare and should be classified as a **parenthetical element** in the segmentation options.

Values Element names.

A typical example is a footnote in the middle of a sentence.

If `<footnote>` isn't classified as a **parenthetical element**, the incorrect text

```
This is the This is a secondsentence.first sentence.
```

is sent to the server. If `<footnote>` is classified as a parenthetical element, the correct text

```
This is the first sentence. This is a second sentence.
```

is sent to the server.

In this case, you would add `footnote` to the list of parenthetical elements like this:

```
parenthetic_elements=footnote,endnote
```

Setting empty_elements

Description If you have elements that don't contain any text such as text entities, but should still be treated as part of the sentence, define them as empty elements.

There are several things to consider when working with elements that don't contain any text:

- The plug-in replaces these elements with a placeholder character when processing the text for checking.

Sometimes elements that don't contain any text but are part of the sentence structure can cause style or grammar rules to flag incorrectly.

In the following example, the `xref` element contains no text but acts as a noun in the sentence.

```
<para>Slidethe<xref xidtype="IBSwitch" xrefid="ib001"></xref>firmlyintoplace.</para>
```

By default, the text "Slide the firmly into place." is sent to the server and the sentence is flagged for incorrect grammar. When you classify the elements as empty, you prevent this type of issue because empty elements are treated as nouns.

- Acrolinx always ignores the contents of empty elements.

This fact is important to consider if your empty element is usually empty but occasionally does contain text.

Considering the previous example, suppose that your xrefs are usually empty but you do have one xref that spans a fragment of your sentence.

```
<para>Slidethe<xref xidtype="IBSwitch" xrefid="ib001">first switch</xref>firmlyintoplace.</para>
```

The text that is within the `xref` tag will be ignored.

If you happened to misspell the word "first" as "frst", Acrolinx would not find this spelling issue because the text is contained within an empty element.

Values Element names.

For example, images are sometimes used in place of words like in the following excerpt:

```
Powered by <image href="acrolinx_logo.jpg"></image>.
```

In this case, you would add `img` to the list of empty elements like this:

```
empty_elements=img,xref
```

Setting default_inclusion_mode

Description When you run a check, this setting determines whether elements are included or excluded by default.

The following elements are *always* excluded or included, regardless of the `default_inclusion_mode` value:

- Elements that are defined as `exclusion_elements` or `inclusion_elements`
 - "Child" elements of exclusion or inclusion elements.
-

Values	<p>One of the following mode names:</p> <ul style="list-style-type: none"> • include • exclude <p>For example, if you would prefer to have all elements excluded by default unless you explicitly include them, you would enter the property as follows:</p> <pre>default_inclusion_mode=exclude</pre>
Setting	mark_excluded_elements
Description	<p>Configure the server to insert placeholders for excluded elements when processing the text.</p> <p>The default is false.</p> <p>This property prevents incorrect flags that are caused by excluded elements. For more information on this what this property does, see the Batch Checker documentation.</p>
Values	<p>Boolean value.</p> <p><i>Example:</i> true, false</p>
Setting	extract_attributes
Description	<p>Use this property to configure Acrolinx to include the values of attributes when sending text to the server for checking.</p> <p>You can use this feature in combination with customized style rules to validate the values of specific attributes. You can also use this feature to check attributes that contain sentences or keywords such as meta content. For more information on developing style rules to validate attribute values, contact your Acrolinx project consultant</p>
Values	<p>Attribute names, element names with attribute names. You can define how attributes are extracted in several ways:</p> <ul style="list-style-type: none"> • Extract the content of a specific attribute regardless of the element in which it appears. <p>Syntax: <code>.<AttributeName></code></p> <p>The following example shows how to configure the Batch Checker to extract and check the value of the <code>name</code> attribute:</p> <pre>extract_attributes=.name</pre> <p>This configuration would extract the text "Header" and "Checking Statistics" from the following code excerpt.</p> <pre> <tablename="Checking statistics"></pre> <ul style="list-style-type: none"> • Extract the content of a specific attribute when it's found in a specific element <p>Syntax: <code><ElementName>.<AttributeName></code></p> <p>The following example shows how to configure the Batch Checker to extract and check the value of the <code>alt</code> attribute within image elements:</p> <pre>extract_attributes=img.alt</pre> <p>This configuration would extract the text "Options Menu" and "Configuration Window" from the following code excerpt.</p> <pre><imgalt="Options Menu" src="13682.jpg"/> <imgalt="Configuration Window" src="3738.jpg"/></pre> <p>It's also possible to develop style rules to check for empty attributes so that you can check for images where the <code>alt</code> attribute exists but isn't defined.</p> <ul style="list-style-type: none"> • Extract the content of a specific attribute only when the element contains another attribute with a specific value. This syntax is useful for checking the content of meta tags <p>Syntax: <code><ElementName>.<AttributeName>#<restrictingAttribute>=<restrictingValue></code></p> <p>The following example shows how to configure the Batch Checker to extract and check the value of the <code>content</code> attribute from all <code>meta</code> elements as long as the <code>meta</code> element has the name "description".</p> <p><i>Example:</i> <code>extract_attributes=meta.content#name=description</code></p> <p>This configuration would extract the text "Topspin is about making your business better." from the following code excerpt.</p> <pre><meta name="Description" content="Topspin is about making your business better." /></pre>
Setting	meta_information

Description	<p>Some Acrolinx clients such as the Acrolinx Batch Checker can read the contents of certain elements or attributes within checked XML documents and write those values into the Scorecards as metadata.</p> <p>This functionality can be useful to process the Scorecards offline and to aggregate the report data by specific custom values. For more information about this feature, see the Batch Checker documentation.</p>
Values	<p>The input for this parameter is set with two values</p> <ul style="list-style-type: none"> The value Meta Attribute Name defines how a meta tag appears in the Scorecard. The Batch Checker writes the Meta Attribute Names as attributes into the <code><identifier></code> tag (prefixed with "meta_"). The Meta Attribute Name enables you to label your metadata. The value Source Element tells the Batch Checker which elements to look for in a checked XML document, and to insert the element values into new meta attributes of the <code><identifier></code> tag. <p>Enter the Meta Attribute Name value first, followed by a colon and then enter the Source Element Name surrounded by angled brackets: <code><...></code>.</p> <p><i>Example:</i> <code>meta_information=MetaTagName:<sourceElement></code></p> <p>The Source Element value can also reference Attributes of XML elements. To define an Attribute of a Source Element, enter the attribute name after the source element name separated by a space:</p> <p><i>Example:</i> <code>meta_information=MetaTagName:<sourceElementsourceAttribute></code></p> <p>The <code>meta_information</code> tag can also take multiple value pairs separated by a semicolon:</p> <p><i>Example:</i> <code>meta_information=MyDocumentID:<bookid>;MyRoleName:<bookelementrole></code></p>

Table 2. Extraction Settings

Setting	Description	Values
<code>ignoreCase</code>	<p>If your element names appear in different cases you might want to turn this off.</p> <p>By default, Acrolinx ignores the casing of element names when assessing the extraction settings.</p> <p>For example, if you configure the element name <code>title</code> to be excluded, Acrolinx also matches the elements <code>TITLE</code> and <code>Title</code>.</p>	<p>Boolean value.</p> <p><i>Example:</i> <code>true, false</code></p> <p>Suppose that you have an element such as <code>Data</code> and it means something different when it takes the uppercase form <code>DATA</code>.</p> <p>In this case, you would configure the extractions settings to be case-sensitive by entering the property as follows:</p> <pre>ignoreCase=false</pre>
<code>exclusion_elements</code>	<p>If you don't want Acrolinx to check the text in an element at all, define it as excluded.</p>	<p>Element names, element names with attributes.</p> <p><i>Example:</i> <code>metadata, cndname, codeph</code></p> <p>For example, you probably don't want to check the text in elements that contain code or identifiers.</p> <pre><metadata>document ID12345</metadata></pre> <p>Note: Excluded XML elements that are nested within included elements are still excluded. All other excluded items are ignored.</p> <p>In the following example, the <code><note></code> element is included and the <code><ph></code> element is excluded.</p> <pre><note>Thistextwillbechecked. <ph>Thistextwillnotbechecked.</ph></note></pre> <p>The content in <code><ph></code> will still be excluded, even if it's contained within an included element.</p>
<code>inclusion_elements</code>	<p>Use the property to define elements that are always included regardless of whether they're nested inside an excluded element.</p> <p>"Child" elements of inclusion elements are also always included.</p>	<p>Element names, element names with attributes.</p> <p>For example, to always include <code>title</code> elements in a check, add it to your list of inclusion elements like this:</p> <pre>inclusion_elements=ph, entry, title</pre>
<code>sentence_break_elements</code>	<p>If you have elements that contain sentences that don't end with a period, define the as sentence-break elements. This means that the end of the element should always be treated as a sentence end.</p>	<p>Element names, element names with attributes.</p> <p>A very typical example is the title element:</p> <pre><title>This is a title</title></pre> <p>Unexpected <code>sentence_too_long</code> flags which actually consist of several sentences are a clear sign that one or more elements must be defined as a sentence break.</p> <p>In this case, you would add <code>title</code> to the list of sentence-break elements like this:</p> <pre>sentence_break_elements=title,entry,li</pre>

Setting	Description	Values
no_break_elements	<p>If you have elements that split a word into separate fragments, define it as a "no break" element.</p>	<p>Element names, element names with attributes. A typical example is the subscript element:</p> <pre>H<sub>2</sub>O concentration</pre> <p>If <code><sub></code> isn't classified, the incorrect text</p> <pre>H 2 O concentration</pre> <p>is sent to the server. If <code><sub></code> is a no break element, the correct text</p> <pre>H2O concentration</pre> <p>is sent to the server.</p> <p>In this case, you would add <code>sub</code> to the list of no break elements like this:</p> <pre>no_break_elements=sub,sup</pre>
token_break_elements	<p>You might define elements as token breaks of you have words which aren't separated by a space.</p> <p>This setting defines the elements which should cause a token break. Adding a token break is basically adding a boundary between words.</p>	<p>Element names, element names with attributes.</p> <p>Usually, you would add a space between words, but sometimes editors separate words with elements rather than spaces.</p> <p>For example, WordProcessingML uses <code><cr/></code> elements to indicate manually inserted line breaks, like this:</p> <pre>This sentence is<cr/>wrapped manually.</pre> <p>You wouldn't want the sentence to end at "is", or have the words interpreted as "iswrapped". So you would update the setting as follows:</p> <pre>token_break_elements=cr</pre>
default_break_level	<p>This setting defines how Acrolinx should interpret elements that aren't already included in the "break element" settings or in any of the other settings.</p> <p>The following settings are "break element" settings:</p> <ul style="list-style-type: none"> <code>no_break_elements</code> <code>sentence_break_elements</code> <code>token_break_elements</code> 	<p>Enter one of the following values:</p> <ul style="list-style-type: none"> <code>none</code> Acrolinx doesn't add any break at all. <code>token</code> Acrolinx adds a token break. <code>sentence</code> Acrolinx adds a sentence break. <p>The default value for this setting is <code>token</code>.</p> <p>For example, if you wanted to have any undefined elements treated as sentence breaks, you would enter the property as follows:</p> <pre>default_break_level=sentence</pre>
parenthetic_elements	<p>Suppose that you have elements that surround fragments of text in the same way that you would put a side-note in parentheses in the middle of a sentence.</p> <p>Example:</p> <pre><p>This is the <footnote>This is a second sentence.</footnote> first sentence.</p></pre> <p>This type of element is fairly rare and should be classified as a parenthetical element in the segmentation options.</p>	<p>Element names.</p> <p>A typical example is a footnote in the middle of a sentence.</p> <p>If <code><footnote></code> isn't classified as a parenthetical element, the incorrect text</p> <pre>This is the This is a secondsentence.first sentence.</pre> <p>is sent to the server. If <code><footnote></code> is classified as a parenthetical element, the correct text</p> <pre>This is the first sentence. This is a second sentence.</pre> <p>is sent to the server.</p> <p>In this case, you would add <code>footnote</code> to the list of parenthetical elements like this:</p> <pre>parenthetic_elements=footnote,endnote</pre>
empty_elements	<p>If you have elements that don't contain any text such as text entities, but should still be treated as part of the sentence, define them as empty elements.</p> <p>There are several things to consider when working with elements that don't contain any text:</p>	<p>Element names.</p> <p>For example, images are sometimes used in place of words like in the following excerpt:</p> <pre>Powered by <image href="acrolinx_logo.jpg"></image>.</pre> <p>In this case, you would add <code>img</code> to the list of empty elements like this:</p>

Setting	Description	Values
	<ul style="list-style-type: none"> The plug-in replaces these elements with a placeholder character when processing the text for checking. <p>Sometimes elements that don't contain any text but are part of the sentence structure can cause style or grammar rules to flag incorrectly.</p> <p>In the following example, the <code>xref</code> element contains no text but acts as a noun in the sentence.</p> <pre><para>Slidethe<xref xidtype="IBSwitch" xrefid="ib001"> </xref>firmlyintoplace. </para></pre> <p>By default, the text "Slide the firmly into place." is sent to the server and the sentence is flagged for incorrect grammar. When you classify the elements as empty, you prevent this type of issue because empty elements are treated as nouns.</p> <ul style="list-style-type: none"> Acrolinx always ignores the contents of empty elements. <p>This fact is important to consider if your empty element is usually empty but occasionally does contain text.</p> <p>Considering the previous example, suppose that your xrefs are usually empty but you do have one xref that spans a fragment of your sentence.</p> <pre><para>Slidethe<xref xidtype="IBSwitch" xrefid="ib001">first switch</xref>firmlyintoplace. </para></pre> <p>The text that is within the <code>xref</code> tag will be ignored.</p> <p>If you happened to misspell the word "first" as "frst", Acrolinx would not find this spelling issue because the text is contained within an empty element.</p>	<code>empty_elements=img,xref</code>

Setting	Description	Values
default_inclusion_mode	<p>When you run a check, this setting determines whether elements are included or excluded by default.</p> <p>The following elements are <i>always</i> excluded or included, regardless of the <code>default_inclusion_mode</code> value:</p> <ul style="list-style-type: none"> • Elements that are defined as <code>exclusion_elements</code> or <code>inclusion_elements</code> • "Child" elements of exclusion or inclusion elements. 	<p>One of the following mode names:</p> <ul style="list-style-type: none"> • <code>include</code> • <code>exclude</code> <p>For example, if you would prefer to have all elements excluded by default unless you explicitly include them, you would enter the property as follows:</p> <pre>default_inclusion_mode=exclude</pre>
mark_excluded_elements	<p>Configure the server to insert placeholders for excluded elements when processing the text.</p> <p>The default is false.</p> <p>This property prevents incorrect flags that are caused by excluded elements. For more information on this what this property does, see the Batch Checker documentation.</p>	<p>Boolean value.</p> <p><i>Example:</i> <code>true, false</code></p>
extract_attributes	<p>Use this property to configure Acrolinx to include the values of attributes when sending text to the server for checking.</p> <p>You can use this feature in combination with customized style rules to validate the values of specific attributes. You can also use this feature to check attributes that contain sentences or keywords such as meta content. For more information on developing style rules to validate attribute values, contact your Acrolinx project consultant</p>	<p>Attribute names, element names with attribute names. You can define how attributes are extracted in several ways:</p> <ul style="list-style-type: none"> • Extract the content of a specific attribute regardless of the element in which it appears. <p>Syntax: <code>.<AttributeName></code></p> <p>The following example shows how to configure the Batch Checker to extract and check the value of the <code>name</code> attribute:</p> <pre>extract_attributes=.name</pre> <p>This configuration would extract the text "Header" and "Checking Statistics" from the following code excerpt.</p> <pre> <tablename="Checking statistics"></pre> • Extract the content of a specific attribute when it's found in a specific element <p>Syntax: <code><ElementName>.<AttributeName></code></p> <p>The following example shows how to configure the Batch Checker to extract and check the value of the <code>alt</code> attribute within image elements:</p> <pre>extract_attributes=img.alt</pre> <p>This configuration would extract the text "Options Menu" and "Configuration Window" from the following code excerpt.</p> <pre><imgalt="Options Menu" src="13682.jpg"/> <imgalt="Configuration Window" src="3738.jpg"/></pre> <p>It's also possible to develop style rules to check for empty attributes so that you can check for images where the <code>alt</code> attribute exists but isn't defined.</p>

Setting	Description	Values
		<p>Extract the content of a specific attribute only when the element contains another attribute with a specific value. This syntax is useful for checking the content of meta tags</p> <pre><ElementName>.<AttributeName># <restrictingAttribute>=<restrictingValue></pre> <p>The following example shows how to configure the Batch Checker to extract and check the value of the <code>content</code> attribute from all <code>meta</code> elements as long as the <code>meta</code> element has the name "description".</p> <p><i>Example:</i> <code>extract_attributes=meta.content#name=description</code></p> <p>This configuration would extract the text "Topspin is about making your business better." from the following code excerpt.</p> <pre><meta name="Description" content="Topspin is about making your business better." /></pre>
<code>meta_information</code>	<p>Some Acrolinx clients such as the Acrolinx Batch Checker can read the contents of certain elements or attributes within checked XML documents and write those values into the Scorecards as metadata.</p> <p>This functionality can be useful to process the Scorecards offline and to aggregate the report data by specific custom values. For more information about this feature, see the Batch Checker documentation.</p>	<p>The input for this parameter is set with two values</p> <ul style="list-style-type: none"> The value Meta Attribute Name defines how a meta tag appears in the Scorecard. The Batch Checker writes the Meta Attribute Names as attributes into the <code><identifier></code> tag (prefixed with "meta_"). The Meta Attribute Name enables you to label your metadata. The value Source Element tells the Batch Checker which elements to look for in a checked XML document, and to insert the element values into new meta attributes of the <code><identifier></code> tag. <p>Enter the Meta Attribute Name value first, followed by a colon and then enter the Source Element Name surrounded by angled brackets: <code><...></code>.</p> <p><i>Example:</i> <code>meta_information=MetaTagName:<sourceElement></code></p> <p>The Source Element value can also reference Attributes of XML elements. To define an Attribute of a Source Element, enter the attribute name after the source element name separated by a space:</p> <p><i>Example:</i> <code>meta_information=MetaTagName: <sourceElementsourceAttribute></code></p> <p>The <code>meta_information</code> tag can also take multiple value pairs separated by a semicolon:</p> <p><i>Example:</i> <code>meta_information=MyDocumentID: <bookid>;MyRoleName:<bookelementrole></code></p>

Syntax for Entering Attributes in a CSD

Most settings in a CSD file take the element names as their values, like this: `sentence_break_elements=codeph,cmdname`. However, you can also enter elements in combination with attributes. If you want to select elements based on their attributes, you enter the values in a specific syntax.

Syntax for Entering Attributes

Acrolinx can select elements for segmentation based on their attributes. You can then refine the segmentation behavior to respond to specific conditions. To enter an attribute as a value of a segmentation parameter, use the following syntax:

- `<elementNameattributeName>` - This value instructs Acrolinx to segment the element if it contains the specified attribute.
- `<elementNameattributeName="value">` - This value instructs Acrolinx to segment the element only if it contains the specified attribute with the specified value.

Remember: This syntax is case-sensitive. Ensure that the case of the elements, attributes, and attribute values match the case used in your XML documents. This syntax is used to segment text within elements only. To include the values of attributes for checking, you use a [different configuration](#).

The following example shows an exclusion element parameter set to exclude two elements based on their attributes.

```
exclusion_elements=<MODULEmodname="voicecontrol">,<PROCEDUREflag>
```

The first value tells Acrolinx to exclude text that is located in `module` elements that have a `modname` attribute set to `voicecontrol`. The second value tells Acrolinx to exclude text within `procedure` elements whenever the element contains the `flag` attribute, regardless of the attribute value.

To include the values of attributes for checking, you use a different syntax.

