# OASIS ebXML Registry

# Proposal: Content-based Discovery

# Category: New feature

# Date: August 15, 2002

# Version 0.2

## Authors: Farrukh Najmi, Nikola Stojanovic, Matt MacKenzie

Table of Contents

## Status of this Document

This note describes the initial proposal for the `Content-based Discovery` work item for OASIS ebXML Registry V3.0. It is expected that the Query sub-team of the OASIS ebXML Registry TC will improve upon this initial proposal and then submit it for consideration by ebXML Registry TC at large.

## 1  Abstract

This document proposes a new feature being added to the query capabilities of the OASIS ebXML Registry targeted for version 3.0. The `Content-based Discovery` feature enables the registry to handle queries that can predicate their results on the content defined by the repository items in addition to predicating on the metadata defined by the RegistryObjects. The feature enables client to discover repository items based upon specific criteria matching their content.

## 2  Motivation

The following motivations drive this proposal:

1. Radically improve the discovery capabilities of the OASIS ebXML Registry.
2. Enable typical business use-cases that require discovery of content based upon the data within the content.

## 3  Assumptions

The following assumptions are made in this proposal:

1. The design center of the proposal will focus on content-based discovery of repository items that are XML documents.

   The proposal will allow for supporting content-based discovery on other types of content besides XML documents.

## 4  Design Goals

The following design goals have been pursued in this proposal:

Require no new interfaces to allow this feature to be implemented in a V2.x registry.

## 5  External Dependencies

This proposal depends upon the following external artifacts and event:
   o   No external dependencies

# 6   Use Cases

There are many scenarios where content-based discovery is necessary.

## 6.1  Find All CPPs Where Role is "Buyer"

A company that sells a product using the RosettaNet PIP3A4 Purchase Order process
wants to find CPPs for other companies where the Role element of the CPP is that of
"Buyer".

## 6.2  Find All XML Schema's That Use Specified Namespace

A client may wish to discover all XML Schema documents in the registry that use an
XML namespace containing the word "oasis".

## 6.3  Find All WSDL Descriptions with a SOAP Binding

An ebXML registry client is attempting to discover all repository items that are WSDL
descriptions that have a SOAP binding defined. Note that SOAP binding related
information is content within the WSDL document and not metadata.

# 7   Content-based Discovery

```
[Note]The following will be a chapter in ebRS
       immediately following Chapter 8 on Queries.
```

This chapter describes the Content-based discovery facility of the ebXML Registry. This
facility enables clients to discover repository items based upon the content contained
within the repository item.

*The essence of the content-based discovery features is based upon the ability to
selectively convert repository content into RegistryObject metadata.*

A registry uses one or more content indexing services to automatically index repository
items when they are submitted to the registry. Indexing a repository item creates
RegistryObject metadata such as Classification instances. The indexed metadata enables
clients to discover the repository item using existing query capabilities of the registry.

```
[Note]The term index is used to refer to
       RegistryObject Metadata generated from
       selective repository item content. It should
       not be confused with databases indexes. It is
       named such because it is similar in concept to
       database indexes, which are metadata generated
       from content.
```

## 7.1 Content Indexing Service

Figure 1 shows that conceptually, a content indexing service (or indexer) accepts as input a repository item and generates as output one or more Classification instances that are used to classify the ExtrinsicObject for that repository item. In addition an indexer accepts as control input an index definition file, which is also a repository item.

Index Definition File

Indexable Content → Content Indexing Service → Index Metadata

☞☞**Figure 1: Abstract Content Indexing Service: Inputs and Outputs**

### 7.1.1 Illustrative Example

Figure 2 shows a UML instance diagram to illustrate how a Content Indexing Service is used. The content indexing service is the normative Default XML Indexing Service described in section 7.9.

- o In the center we see a Content Indexing Service name defaultXMLIndexer.
- o On the left side we see a CPP repository item and its ExtrinsicObject inputExtObjForCPP being input as Indexable Content to the defaultXMLIndexer.
- o On top we see an XSLT style sheet repository item and its ExtrinsicObject being sent as an Index Definition File to the defaultXMLIndexer.
- o On the right we see the outputExtObjForCPP, which is the modified ExtrinsicObject for the CPP. We also see a Classification roleClassification, which classifies the CPP by the Role element within the CPP. These are the Index Metadata generated as a result of the indexer indexing the CPP.

126
127 ☞☞**Figure 2: Example of CPP indexing using Default XML Indexer**

## 128  **7.2  Index Definition File**

129 The Index Definition File describes what information should the indexer extract from the
130 repository item and subsequently map it to the generated Classification(s). This
131 specification does not define the format of the Index Definition File. Each indexer is free
132 to define its own Index Definition File format in an indexer specific manner. The only
133 constraint in this specification is that the index definition file must be a repository item.

## 134  **7.3  Index-able Content**

135 The index-able content is the content that the client wishes to be indexed by the Content
136 Indexing Service. As such it is the subject of the content indexing action.
137 This specification does not define the format of index-able content. This specification
138 describes how a client may register arbitrary indexers for indexing arbitrary content
139 types.
140 The most common use case for an indexer is for indexing XML documents. Therefore,
141 this specification also provides a normative definition for a specialized XML Content
142 Indexer in section 7.9.
143 An ebXML Registry must provide native built-in support for the normative XML Content
144 Indexer.
145 In addition, an ebXML Registry must allow clients to register arbitrary indexers for
146 arbitrary content. In either case the registry must use the appropriate indexer if one exists,
147 to index a repository item when it is submitted.

## 148  **7.4  Index Metadata**

149 A content indexing service indexes a repository item by processing it and extracting
150 specific information content as specified by the Index Definition File. The content
151 indexing service must map the extracted content to index metadata in form of instances of
152 RIM classes.

153 For example, the index metadata may consists of:

154 Classification instances

155     o   ExternalIdentifier instances

156     o   ExternalLink instances

157     o   The name attribute for the ExtrinsicObject for the index-able content

158     o   The description attribute for the ExtrinsicObject for the able-able content

159 A content indexing service is free to generate any class defined by RIM as index
160 metadata in an application specific manner.


## 161   7.5  Content Indexing Protocol

162 The interface of the content indexing service must implement a single method called
163 indexContent. The indexContent method accepts an IndexContentRequest as parameter
164 and returns an IndexContentResponse as its response if there are no errors.

165 The IndexContentRequest contains repository items that need to be indexed. The
166 resulting IndexContentResponse contains the metadata that gets generated by the Content
167 Indexing Service as a result of indexing the specified repository items.

168 The content indexing protocol is abstract and does not specify the interface or behavior of
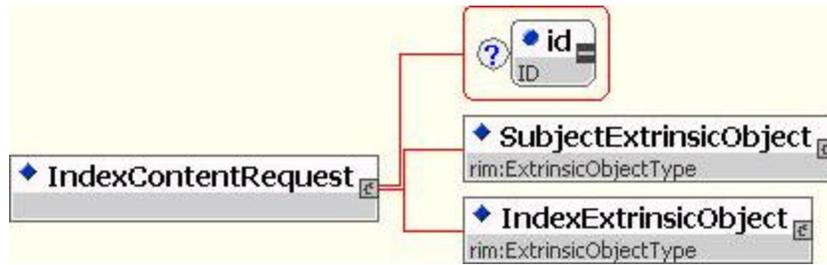169 any specific Content Indexing Service.



170
171 ☞☞**Figure 3: Content Indexing Protocol**


### 172   7.5.1  IndexContentRequest

173 The IndexContentRequest is used to submit repository items to a Content Indexing
174 Service so that it can create index metadata for the specified repository items.

175 **7.5.1.1 Syntax:**



176
177 ☞☞**Figure 4: IndexContentRequest Syntax**

178 **7.5.1.2 Parameters:**

179   ✍ *id*: Inherited request id attribute common to all requests.

180   ✍ *IndexExtrinsicObject*: This parameter specifies the ExtrinsicObject for the
181    repository item that the caller wishes to specify as the Index Definition
182    file. This specification does not specify the format of this repository item.
183    There must a corresponding repository item as an attachment to this
184    request. The corresponding repository item should follow the same rules
185    as attachments in SubmitObjectsRequest.

186   ✍ *SubjectExtrinsicObject:* This parameter specifies the ExtrinsicObject for
187    the repository item that the caller wishes to be indexed. This specification
188    does not specify the format of this repository item. There must a
189    corresponding repository item as an attachment to this request. The
190    corresponding repository item should follow the same rules as attachments
191    in SubmitObjectsRequest.

192

193 **7.5.1.3 Returns:**

194 This request returns an IndexContentResponse upon success. See section 7.5.2 for details.

195 **7.5.1.4 Exceptions:**

196 In addition to the exceptions common to all requests, the following exceptions may be
197 returned:

198   ✍ MissingRepositoryItemException: signifies that the caller did not provide
199    a required repository item as an attachment to this request.

200   ✍ UnsupportedIndexException: signifies that this Content Indexing Service
201    did not support the IndexExtrinsicObject provided by the client.

202   ✍ UnsupportedSubjectException: signifies that this Content Indexing
203    Service did not support the SubjectExtrinsicObject provided by the client.
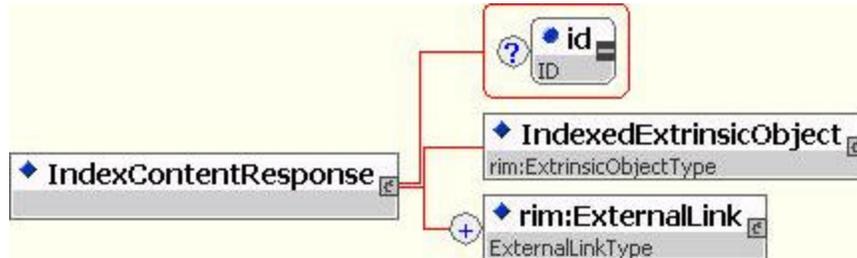
204

205 **7.5.2   IndexContentResponse**

206 The IndexContentRequest is sent by the Content Indexing Service as a response to an
207 IndexContentRequest.

208

209 **7.5.2.1   Syntax:**



210
211 ☞☞**Figure 5: IndexContentResponse Syntax**

212 **7.5.2.2   Parameters:**

213 ✍ *ExternalLink:*  This parameter specifies one or more ExternalLink
214 elements that may be generated as index metadata during the indexing of
215 the repository item.

216 ✍ *id*: id attribute inherited from RegistryResponseType.

217 ✍ *IndexedExtrinsicObject:*  This parameter specifies the modified
218 ExtrinsicObject for the repository item that has been indexed by the
219 Content Indexing Service. The Content Indexing Service may add
220 metadata such as Classifications, ExternalIdentifiers, name, description
221 etc. to the IndexedExtrinsicObject element. There must not be an
222 accompanying repository item as an attachment to this request.

223

224 # 7.6  Publishing a Content Indexing Service

225 Any publisher may publish an arbitrary content indexing service to an ebXML Registry.
226 The content indexing service must be published using the existing LifeCycleManager
227 interface. The publisher must use the existing SubmitObjectsRequest to publish:

228 A Service instances with two required slots named supportedObjectType and
229 supportedMimeType. The values of these slots are explained in section 7.7.1.

230   o  A ServiceBinding instance contained within the Service instance
231   o  An ExternalLink instance on the ServiceBinding that must be resolvable to a web
232      page describing:
233        o   The format of the supported Index-able Content
234        o   The format of the supported Index Definition File

235 Note that that no SpecificationLink is required since this specification is implicit for
236 Content Indexing Services.

237 A content indexing service must be published with a default index definition file that
238 must be an ExtrinsicObject and repository item pair. The ExtrinsicObject for the index
239 definition must have two required slots named supportedObjectType and
240 supportedMimeType. The values of these slots are explained in section 7.7.2.

### 7.6.1  Multiple Indexers and Index Definition Files

242 This specification allows clients to submit multiple indexers and index definition files for
243 the same mimeType/objectType. How a registry handles multiple indexer and index
244 definition file submission for the same type of content is a matter of registry specific
245 policy. If a registry does not allow this then it must send an InvalidRequestException
246 with a reason, when a duplicate indexer or index def is submitted. If a registry allows this
247 then it must provide a conflict resolution mechanism to select the appropriate indexer and
248 index definition file in some registry specific manner.

### 7.6.2  Restrictions On Publishing Content Indexing Services

250 A client may submit any content indexing service or index definition file. A registry may
251 use registry specific policies to determine whether a client submitted content indexing
252 service or index definition file are acceptable. For example a registry may require that the
253 content indexing service or index definition file does not create excessive metadata. A
254 registry may reject a SubmitObjectRequest with an InvalidRequestException and give a
255 reason why the request was rejected, upon receiving requests publishing Content
256 Indexing Service or Index Definition File that is unreasonable.

## 7.7  Dynamic Content Indexing

258 Some time during or after a publisher submits a repository item, the registry must check
259 to see if there is a Content Indexing Service and index definition file registered for that
260 type of repository item. This is referred to as Content Indexing Service resolution as
261 described in section 7.7.1 and index definition file resolution as described in section
262 7.7.2.
263 If a Content Indexing Service and index definition file are found then the registry must
264 invoke that service using the Content Indexing Protocol. In the invocation, it gives a
265 repository item as Index-able Content and a repository item as Index Definition File
266 within an IndexContentRequest. The Content Indexing Service must index the content
267 and return the modified ExtrinsicObject for the Index-able Content such that it has index
268 metadata generated from relevant portions of the Index-able Content.
269 The registry must store the repository item along with the modified ExtrinsicObject
270 annotated with the index metadata once the Content Indexing Protocol is completed.
271 Note that a registry may do dynamic content indexing synchronous with the original
272 SubmitObjectRequest request or it may do so asynchronously sometime after the request
273 is committed. It is suggested that asynchronous indexing latency should be no more than
274 24 hours.
275 The result of dynamic content indexing is that index-able content gets indexed
276 dynamically when it is submitted. Once indexed it is possible to use the index metadata to
277 do dynamic content-based discovery of the index-able content.

### 7.7.1   Content Indexing Service Resolution Algorithm

When a registry receives a submission of an ExtrinsicObject EO1 and repository item pair, it must use the following algorithm to determine or resolve the content indexing service to be used to index that content:

1. Check if a Service instance for the indexer exists that has a slot named supportedObjectType with value matching the objectType of the ExtrinsicObject EO1, AND a slot named supportedMimeType with value matching the mimeType of the ExtrinsicObject. If so use that indexer.

2. Otherwise, Check if a Service instance for the indexer exists that has a slot named supportedObjectType with value matching the objectType of the ExtrinsicObject EO1. If so use that indexer.

3. Check if a Service instance for the indexer exists that has a slot named supportedMimeType with value matching the mimeType of the ExtrinsicObject EO1. If so use that indexer.

If no indexer is found then content should not be indexed. If an indexer is found then the registry must resolve an index definition file as defined next.

### 7.7.2   Index Definition File Resolution Algorithm

When a registry receives a submission of an ExtrinsicObject instance EO1 and repository item pair, it must first resolve a content indexing service as described in section 7.7.1.

If a content indexing service has been resolved then the registry must use the following algorithm to determine or resolve the index definition file to be used to index that content:

1. Check if an ExtrinsicObject instance for the index definition file exists that has a slot named supportedObjectType with value matching the objectType of the ExtrinsicObject EO1, AND a slot named supportedMimeType with value matching the mimeType of the ExtrinsicObject. If so use that index definition file.

2. Otherwise, Check if an ExtrinsicObject instance for the index definition file exists that has a slot named supportedObjectType with value matching the objectType of the ExtrinsicObject EO1. If so use that index definition file.

3. Check if an ExtrinsicObject instance for the index definition file exists that has a slot named supportedMimeType with value matching the mimeType of the ExtrinsicObject EO1. If so use that index definition file.

If no index definition file is found then content should not be indexed.

## 7.8   Dynamic Content-based Discovery

As described earlier, index-able content is automatically indexed when it is submitted to the registry. This content may subsequently be dynamically discovered using the index metadata within existing AdhocQueryRequest. Because the index metadata is based upon index-able content, an AdhocQueryRequest can perform dynamic content- based

317   discovery.

## 7.9  Default XML Content Indexer

319   An ebXML Registry must provide the XML Content Indexing Service natively as a built-
320   in service. The XML content indexing service accepts an XML instance document as its
321   input and it accepts an XSLT Style sheet as a Content Definition File. Each type of
322   content should have its own unique XSLT style sheet. For example and ebXML CPP
323   document should have a specialize ebXML CPP index definition style sheet. The XML
324   content indexing service must apply the XSLT style sheet to the XML instance document
325   input to generate the index metadata. Since a single style sheet must be applied to both
326   the ExtrinsicObject and the Index-able Content, we must assume the two documents to be
327   composed within a single virtual document the schema for which is as follows:
328
329   Do we need the outer <somename> tag at all?
330   <somename>
331          <ExtrinsicObject/>
332          <CPP/>
333   </somename>

### 7.9.1  Publishing of Default XML Content Indexer

335   The default XML Content Indexing Service need not be explicitly published to an
336   ebXML Registry. An ebXML Registry must provide the XML Content Indexing Service
337   natively as a built-in service. This built-in service must be published to the registry as
338   part of the intrinsic bootstrapping of required data within the registry.

## 8   Notes

340   These notes are here to not lose the thought and will be merged into the proposal later.

341      o  Need replacement term for index. Choices suggested so far are: promotion,

342      o  Need illustrative example to include sample Index Def file, sample CPP and
343         sample output.

344      o  Do we need to replace repository item with ExtrinsicObject for inedex-able
345         content? Reason, EO can be without repository item. [FN] No I don't think so.
346         Nikola do you still want this?

347      o  How to establish ClassificationScheme, identificationScheme etc.

348      o  Instead of using slots maybe we can do Classifications to pre-defined
349         schemes/nodes. Nikola please explain this. Do we still need it?

350      o  Same client goes to 2 registries that have 2 different indexer. User may be
351         confused by different output for same input and control.

352      o  IndexContentResponse: how to handle any non-composed metadata such as
353         ExternalIdentifier, Package etc.?

354      o  How to track generated metadata separate from submitted metadata? Should we
355          also log which indexer and index file created it?

356      o  Should we allow a way for client to override default index def file and/or default
357          indexer?

358

359

360