

Metacat - Bug #5696

pathQuery returns eml docs which have no public access granted

08/24/2012 12:06 PM - gastil gastil

<b>Status:</b>	Resolved	<b>Start date:</b>	08/24/2012
<b>Priority:</b>	Normal	<b>Due date:</b>	
<b>Assignee:</b>	ben leinfelder	<b>% Done:</b>	0%
<b>Category:</b>	metacat	<b>Estimated time:</b>	0.00 hour
<b>Target version:</b>	2.0.5	<b>Spent time:</b>	0.00 hour
<b>Bugzilla-Id:</b>	5696		
<b>Description</b>			
<p>As far as I remember, non-public eml docs did not used to be returned in pathQuery result sets in earlier versions of metacat. This is with <a href="http://metacat.lternet.edu/knb/metacat?action=getversion">http://metacat.lternet.edu/knb/metacat?action=getversion</a></p> <p>&lt;version&gt;2.0.3&lt;/version&gt;</p> <p>A pathQuery returns an eml doc which does not have public read access. Example: knb-lter-sev.389.3</p> <p>with</p> <pre>&lt;access authSystem="knb" order="denyFirst" scope="document"&gt; &lt;allow&gt; &lt;principal&gt;uid=SEV, o=lter, dc=ecoinformatics, dc=org&lt;/principal&gt; &lt;permission&gt;all&lt;/permission&gt; &lt;/allow&gt; &lt;/access&gt;</pre> <p>A pathQuery returned this in its result set:</p> <pre>&lt;document&gt; &lt;docid&gt;knb-lter-sev.389.3&lt;/docid&gt; &lt;docname&gt;eml&lt;/docname&gt; &lt;doctype&gt;eml://ecoinformatics.org/eml-2.0.1&lt;/doctype&gt; &lt;createdate&gt;2005-07-29&lt;/createdate&gt; &lt;updatedate&gt;2012-08-22&lt;/updatedate&gt; &lt;param name="@packageId"&gt;sev.00389.1&lt;/param&gt; &lt;param name="dataset/title"&gt;Lightning Strike Data for New Mexico, 1989&lt;/param&gt; &lt;/document&gt;</pre> <p>This may be related in part to bug <a href="#">#5553</a> (not sure). The denyFirst may be part of the problem. The older revisions also had denyFirst.</p>			

History

#1 - 09/03/2012 11:20 AM - ben leinfelder

I've recreated the scenario on my local Metacat installation -- it appears the permissions for the older revision are still being applied to the newer revision in the search results. I suspect this is related to how the index is purged, or not, as the case seems to indicate.

#2 - 09/04/2012 08:43 AM - ben leinfelder

Metacat v2.0.4 will include a fix for this issue.

#3 - 09/28/2012 08:11 AM - ben leinfelder

Seems this query can be quite expensive when the DB has a large number of documents. Re-working to remove the max(rev) condition - hoping that it does not require a massive overhaul of the QuerySpecification->SQL code.

#4 - 09/28/2012 09:20 AM - ben leinfelder

With the join to the xml\_documents table, the response is better - but not that great (3 minutes for "tree" keyword search:

MetacatHandler.handleSQuery - squery:  
<pathquery version="1.2">  
<querytitle>Advanced Search</querytitle>

```
<returnfield>keyword</returnfield>
<returndoctype>eml://ecoinformatics.org/eml-2.1.0</returndoctype>
<returndoctype>eml://ecoinformatics.org/eml-2.0.1</returndoctype>
<returndoctype>eml://ecoinformatics.org/eml-2.0.0</returndoctype>
<querygroup operator="UNION">
<queryterm searchmode="contains" casesensitive="false">
<value>tree</value>
<pathexpr>keyword</pathexpr>
</queryterm>
</querygroup>
</pathquery>
ran in 179648 ms [edu.ucsb.nceas.metacat.MetacatHandler]
```

#### #5 - 09/28/2012 09:31 AM - ben leinfelder

The expensive part seems to be the subqueries of all public-read docs and all public-read-deny docs:

```
SELECT docid,docname,doctype,date_created, date_updated, rev FROM xml_documents WHERE docid IN ((SELECT DISTINCT docid FROM
xml_path_index WHERE ((UPPER LIKE TREE AND path LIKE keyword) ))) AND (docid IN (SELECT id.docid from xml_access xa, identifier id,
xml_documents xmld WHERE id.guid = xa.guid AND id.docid = xmld.docid AND id.rev = xmld.rev AND ( (lower(principal_name) = 'public') AND
perm_type = 'allow' AND permission > 3)) AND docid NOT IN (SELECT id.docid from xml_access xa, identifier id, xml_documents xmld WHERE
id.guid = xa.guid AND id.docid = xmld.docid AND id.rev = xmld.rev AND ( (lower(principal_name) = 'public') AND perm_type = 'deny' AND perm_order
='allowFirst' AND permission > 3) ))
```

#### #6 - 10/18/2012 04:40 PM - ben leinfelder

I've reworked how access was being checked -- now we have a simpler clause there and the current revision handling is done "higher up" in the query -- this saves us a lot of time when we come to the access control clauses.

#### #7 - 03/27/2013 02:31 PM - Redmine Admin

Original Bugzilla ID was 5696