```
*********************************************************
   12650 Tue Sep  9 16:17:33 2008
new/src/modules/server/catalog.py
3166 feed generation needs performance improvement
3306 feed returns invalid last-modified header
*********************************************************
  1 #!/usr/bin/python2.4
  2 #
  3 # CDDL HEADER START
  4 #
  5 # The contents of this file are subject to the terms of the
  6 # Common Development and Distribution License (the "License").
  7 # You may not use this file except in compliance with the License.
  8 #
  9 # You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
 10 # or http://www.opensolaris.org/os/licensing.
 11 # See the License for the specific language governing permissions
 12 # and limitations under the License.
 13 #
 14 # When distributing Covered Code, include this CDDL HEADER in each
 15 # file and include the License file at usr/src/OPENSOLARIS.LICENSE.
 16 # If applicable, add the following below this CDDL HEADER, with the
 17 # fields enclosed by brackets "[]" replaced with your own identifying
 18 # information: Portions Copyright [yyyy] [name of copyright owner]
 19 #
 20 # CDDL HEADER END
 21 #
 22 # Copyright 2008 Sun Microsystems, Inc.  All rights reserved.
 23 # Use is subject to license terms.

 25 import subprocess
 26 import threading
 27 import signal
 28 import os
 29 import sys
 30 import cherrypy

 32 import pkg.catalog as catalog
 33 import pkg.fmri as fmri
 34 #endif /* ! codereview */
 35 import pkg.indexer as indexer
 36 import pkg.server.query_engine as query_e

 38 from pkg.misc import SERVER_DEFAULT_MEM_USE_KB
 39 from pkg.misc import emsg

 41 class ServerCatalog(catalog.Catalog):
 42         """The catalog information which is only needed by the server."""

 44         def __init__(self, cat_root, authority = None, pkg_root = None,
 45             read_only = False, index_root = None, repo_root = None,
 46             rebuild = True):

 48                 self.index_root = index_root
 49                 self.repo_root = repo_root

 51                 # The update_handle lock protects the update_handle variable.
 52                 # This allows update_handle to be checked and acted on in a
 53                 # consistent step, preventing the dropping of needed updates.
 54                 # The check at the top of refresh index should always be done
 55                 # prior to deciding to spin off a process for indexing as it
 56                 # prevents more than one indexing process being run at the same
 57                 # time.
 58                 self.searchdb_update_handle_lock = threading.Lock()

 60                 if self.index_root:
```

```
 61                         self.query_engine = \
 62                                 query_e.ServerQueryEngine(self.index_root)

 64                 if os.name == 'posix':
 65                         try:
 66                                 signal.signal(signal.SIGCHLD,
 67                                     self.child_handler)
 68                         except ValueError:
 69                                 emsg("Tried to create signal handler in "
 70                                     "a thread other than the main thread")

 72                 self.searchdb_update_handle = None
 73                 self._search_available = False
 74                 self.deferred_searchdb_updates = []
 75                 self.deferred_searchdb_updates_lock = threading.Lock()

 77                 self.refresh_again = False

 79                 catalog.Catalog.__init__(self, cat_root, authority, pkg_root,
 80                     read_only, rebuild)

 82                 if not self._search_available:
 83                         self._check_search()

 85         def whence(self, cmd):
 86                 if cmd[0] != '/':
 87                         tmp_cmd = cmd
 88                         cmd = None
 89                         path = os.environ['PATH'].split(':')
 90                         path.append(os.environ['PWD'])
 91                         for p in path:
 92                                 if os.path.exists(os.path.join(p, tmp_cmd)):
 93                                         cmd = os.path.join(p, tmp_cmd)
 94                                         break
 95                         assert cmd
 96                 return cmd

 98         def refresh_index(self):
 99                 """ This function refreshes the search indexes if there any new
100                 packages. It starts a subprocess which results in a call to
101                 run_update_index (see below) which does the actual update.
102                 """

104                 self.searchdb_update_handle_lock.acquire()

106                 if self.searchdb_update_handle:
107                         self.refresh_again = True
108                         self.searchdb_update_handle_lock.release()
109                         return

111                 try:
112                         fmris_to_index = set(self.fmris())

114                         indexer.Indexer.check_for_updates(self.index_root,
115                             fmris_to_index)

117                         if fmris_to_index:
118                                 if os.name == 'posix':
119                                         cmd = self.whence(sys.argv[0])
120                                         args = (cmd, "--refresh-index", "-d",
121                                             self.repo_root)
122                                         try:
123                                                 self.searchdb_update_handle = \
124                                                     subprocess.Popen(args,
125                                                         stderr = \
126                                                         subprocess.STDOUT)
```

```
127                                    except Exception, e:
128                                            emsg("Starting the indexing "
129                                                "process failed")
130                                            raise
131                            else:
132                                    self.run_update_index()
133                    else:
134                            # Since there is nothing to index, setup
135                            # the index and declare search available.
136                            # We only log this if this represents
137                            # a change in status of the server.
138                            ind = indexer.Indexer(self.index_root,
139                                SERVER_DEFAULT_MEM_USE_KB)
140                            ind.setup()
141                            if not self._search_available:
142                                    cherrypy.log("Search Available",
143                                        "INDEX")
144                            self._search_available = True
145            finally:
146                    self.searchdb_update_handle_lock.release()

148    def run_update_index(self):
149            """ Determines which fmris need to be indexed and passes them
150            to the indexer.

152            Note: Only one instance of this method should be running.
153            External locking is expected to ensure this behavior. Calling
154            refresh index is the preferred method to use to reindex.
155            """
156            fmris_to_index = set(self.fmris())

158            indexer.Indexer.check_for_updates(self.index_root,
159                fmris_to_index)

161            if fmris_to_index:
162                    self.__update_searchdb_unlocked(fmris_to_index)
163            else:
164                    ind = indexer.Indexer(self.index_root,
165                        SERVER_DEFAULT_MEM_USE_KB)
166                    ind.setup()

168    def _check_search(self):
169            ind = indexer.Indexer(self.index_root,
170                SERVER_DEFAULT_MEM_USE_KB)
171            if ind.check_index_existence():
172                    self._search_available = True
173                    cherrypy.log("Search Available", "INDEX")

175    def build_catalog(self):
176            """ Creates an Indexer instance and after building the
177            catalog, refreshes the index.
178            """
179            self._check_search()
180            catalog.Catalog.build_catalog(self)
181            # refresh_index doesn't use file modification times
182            # to determine which packages need to be indexed, so use
183            # it to reindex if it's needed.
184            self.refresh_index()

186    def child_handler(self, sig, frame):
187            """ Handler method for the SIGCLD signal.  Checks to see if the
188            search database update child has finished, and enables searching
189            if it finished successfully, or logs an error if it didn't.
190            """
191            try:
192                    signal.signal(signal.SIGCHLD, self.child_handler)
```

```
193                    except ValueError:
194                            emsg("Tried to create signal handler in "
195                                "a thread other than the main thread")
196            # If there's no update_handle, then another subprocess was
197            # spun off and that was what finished. If the poll() returns
198            # None, then while the indexer was running, another process
199            # that was spun off finished.
200            rc = None
201            if not self.searchdb_update_handle:
202                    return
203            rc = self.searchdb_update_handle.poll()
204            if rc == None:
205                    return

207            if rc == 0:
208                    self._search_available = True
209                    cherrypy.log("Search indexes updated and available.",
210                        "INDEX")
211                    # Need to acquire this lock to prevent the possibility
212                    # of a race condition with refresh_index where a needed
213                    # refresh is dropped. It is possible that an extra
214                    # refresh will be done with this code, but that refresh
215                    # should be very quick to finish.
216                    self.searchdb_update_handle_lock.acquire()
217                    self.searchdb_update_handle = None
218                    self.searchdb_update_handle_lock.release()

220                    if self.refresh_again:
221                            self.refresh_again = False
222                            self.refresh_index()
223            elif rc > 0:
224                    # XXX This should be logged instead
225                    # If the refresh of the index failed, defensively
226                    # declare that search is unavailable.
227                    self._search_available = False
228                    emsg(_("ERROR building search database, rc: %s"))
229                    emsg(_(self.searchdb_update_handle.stderr.read()))

231    def __update_searchdb_unlocked(self, fmri_list):
232            """ Takes a fmri_list and calls the indexer with a list of fmri
233            and manifest file path pairs. It assumes that all needed
234            locking has already occurred.
235            """
236            assert self.index_root
237            fmri_manifest_list = []

239            # Rather than storing those, simply pass along the
240            # file and have the indexer take care of opening and
241            # reading the manifest file. Since the indexer
242            # processes and discards the manifest structure (and its
243            # search dictionary for that matter) this
244            # is much more memory efficient.

246            for f in fmri_list:
247                    mfst_path = os.path.join(self.pkg_root,
248                                        f.get_dir_path())
249                    fmri_manifest_list.append((f, mfst_path))

251            if fmri_manifest_list:
252                    index_inst = indexer.Indexer(self.index_root,
253                        SERVER_DEFAULT_MEM_USE_KB)
254                    index_inst.server_update_index(fmri_manifest_list)

256    def search(self, token):
257            """Search through the search database for 'token'.  Return a
258            list of token type / fmri pairs."""
```

```
259                         assert self.index_root
260                         if not self.query_engine:
261                                 self.query_engine = \
262                                         query_e.ServerQueryEngine(self.index_root)
263                         query = query_e.Query(token, case_sensitive=False)
264                         return self.query_engine.search(query)

266                 @staticmethod
267                 def read_catalog(catalog, dir, auth=None):
268                         """Read the catalog file in "dir" and combine it with the
269                         existing data in "catalog"."""

271                         catf = file(os.path.join(dir, "catalog"))
272                         for line in catf:
273                                 if not line.startswith("V pkg") and \
274                                     not line.startswith("C pkg"):
275                                         continue

277                                 f = fmri.PkgFmri(line[7:])
278                                 ServerCatalog.cache_fmri(catalog, f, auth)

280                         catf.close()

282 #endif /* ! codereview */
```

```
**********************************************************
    16866 Tue Sep  9 16:17:34 2008
new/src/modules/server/feed.py
3166 feed generation needs performance improvement
3306 feed returns invalid last-modified header
**********************************************************
    1 #!/usr/bin/python2.4
    2 #
    3 # CDDL HEADER START
    4 #
    5 # The contents of this file are subject to the terms of the
    6 # Common Development and Distribution License (the "License").
    7 # You may not use this file except in compliance with the License.
    8 #
    9 # You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
   10 # or http://www.opensolaris.org/os/licensing.
   11 # See the License for the specific language governing permissions
   12 # and limitations under the License.
   13 #
   14 # When distributing Covered Code, include this CDDL HEADER in each
   15 # file and include the License file at usr/src/OPENSOLARIS.LICENSE.
   16 # If applicable, add the following below this CDDL HEADER, with the
   17 # fields enclosed by brackets "[]" replaced with your own identifying
   18 # information: Portions Copyright [yyyy] [name of copyright owner]
   19 #
   20 # CDDL HEADER END
   21 #
   22 # Copyright 2008 Sun Microsystems, Inc.  All rights reserved.
   23 # Use is subject to license terms.

   25 """feed - routines for generating RFC 4287 Atom feeds for packaging server

   27    At present, the pkg.server.feed module provides a set of routines that, from
   28    a catalog, allow the construction of a feed representing the activity within
   29    a given time period."""

   31 import cherrypy
   32 from cherrypy.lib.static import serve_file
   33 import cStringIO
   34 import datetime
   35 import httplib
   36 import os
   37 import rfc822
   38 import sys
   39 #endif /* ! codereview */
   40 import time
   41 import urllib
   42 import xml.dom.minidom as xmini

   44 from pkg.misc import get_rel_path, get_res_path
   45 import pkg.server.catalog as catalog
   38 import pkg.catalog as catalog
   46 import pkg.fmri as fmri
   47 import pkg.Uuid25 as uuid

   49 MIME_TYPE = 'application/atom+xml'
   50 CACHE_FILENAME = "feed.xml"
   51 RFC3339_FMT = "%Y-%m-%dT%H:%M:%SZ"

   53 def dt_to_rfc3339_str(ts):
   54         """Returns a string representing a datetime object formatted according
   55         to RFC 3339.
   56         """
   57         return ts.strftime(RFC3339_FMT)

   59 def rfc3339_str_to_ts(ts_str):
```

```
   60         """Returns a timestamp representing 'ts_str', which should be in the
   61         format specified by RFC 3339.
   62         """
   63         return time.mktime(time.strptime(ts_str, RFC3339_FMT))

   65 def rfc3339_str_to_dt(ts_str):
   66         """Returns a datetime object representing 'ts_str', which should be in
   67         the format specified by RFC 3339.
   68         """
   69         return datetime.datetime(*time.strptime(ts_str, RFC3339_FMT)[0:6])

   71 def ults_to_ts(ts_str):
   72         """Returns a timestamp representing 'ts_str', which should be in
   73         updatelog format.
   74         """
   75         # Python doesn't support fractional seconds for strptime.
   76         ts_str = ts_str.split('.')[0]
   77         # Currently, updatelog entries are in local time, not UTC.
   78         return time.mktime(time.strptime(ts_str, "%Y-%m-%dT%H:%M:%S"))

   80 def ults_to_rfc3339_str(ts_str):
   81         """Returns a timestamp representing 'ts_str', which should be in
   82         updatelog format.
   83         """
   84         ltime = ults_to_ts(ts_str)
   85         # Currently, updatelog entries are in local time, not UTC.
   86         return dt_to_rfc3339_str(datetime.datetime(
   87             *time.gmtime(ltime)[0:6]))

   89 def fmri_to_taguri(rcfg, f):
   90         """Generates a 'tag' uri compliant with RFC 4151.  Visit
   91         http://www.taguri.org/ for more information.
   92         """
   93         return "tag:%s,%s:%s" % (rcfg.get_attribute("feed",
   94             "authority"), f.get_timestamp().strftime("%Y-%m-%d"),
   95             urllib.unquote(f.get_url_path()))

   97 def init(scfg, rcfg):
   98         """This function performs general initialization work that is needed
   99         for feeds to work correctly.
  100         """

  102         if not scfg.is_read_only():
  103                 # RSS/Atom feeds require a unique identifier, so
  104                 # generate one if isn't defined already.  This
  105                 # needs to be a persistent value, so we only
  106                 # generate this if we can save the configuration.
  107                 fid = rcfg.get_attribute("feed", "id")
  108                 if not fid:
  109                         # Create a random UUID (type 4).
  110                         rcfg._set_attribute("feed", "id", uuid.uuid4())

  112                 # Ensure any configuration changes are reflected in the feed.
  113                 __clear_cache(scfg)

  115 def set_title(request, rcfg, doc, feed, update_ts):
  116         """This function attaches the necessary RSS/Atom feed elements needed
  117         to provide title, author and contact information to the provided
  118         xmini document object using the provided feed object and update
  119         time.
  120         """

  122         t = doc.createElement("title")
  123         ti = xmini.Text()
  124         ti.replaceWholeText(rcfg.get_attribute("feed", "name"))
  125         t.appendChild(ti)
```

```
126             feed.appendChild(t)

128             l = doc.createElement("link")
129             l.setAttribute("href", cherrypy.url())
130             l.setAttribute("rel", "self")
131             feed.appendChild(l)

133             # Atom requires each feed to have a permanent, universally unique
134             # identifier.
135             i = doc.createElement("id")
136             it = xmini.Text()
137             it.replaceWholeText("urn:uuid:%s" % rcfg.get_attribute("feed", "id"))
138             i.appendChild(it)
139             feed.appendChild(i)

141             # Indicate when the feed was last updated.
142             u = doc.createElement("updated")
143             ut = xmini.Text()
144             ut.replaceWholeText(dt_to_rfc3339_str(update_ts))
145             u.appendChild(ut)
146             feed.appendChild(u)

148             # Add our icon.
149             i = doc.createElement("icon")
150             it = xmini.Text()
151             it.replaceWholeText(get_res_path(request, rcfg.get_attribute(
152                 "feed", "icon")))
153             i.appendChild(it)
154             feed.appendChild(i)

156             # Add our logo.
157             l = doc.createElement("logo")
158             lt = xmini.Text()
159             lt.replaceWholeText(get_res_path(request, rcfg.get_attribute(
160                 "feed", "logo")))
161             l.appendChild(lt)
162             feed.appendChild(l)

164             maintainer = rcfg.get_attribute("repository", "maintainer")
165             # The author information isn't required, but can be useful.
166             if maintainer:
167                     name, email = rfc822.AddressList(maintainer).addresslist[0]

169                     if email and not name:
170                             # If we got an email address, but no name, then
171                             # the name was likely parsed as a local address. In
172                             # that case, assume the whole string is the name.
173                             name = maintainer
174                             email = None

176                     a = doc.createElement("author")

178                     # First we have to add a name element. This is required if an
179                     # author element exists.
180                     n = doc.createElement("name")
181                     nt = xmini.Text()
182                     nt.replaceWholeText(name)
183                     n.appendChild(nt)
184                     a.appendChild(n)

186                     if email:
187                             # If we were able to extract an email address from the
188                             # maintainer information, add the optional email
189                             # element to provide a point of communication.
190                             e = doc.createElement("email")
191                             et = xmini.Text()
```

```
192                             et.replaceWholeText(email)
193                             e.appendChild(et)
194                             a.appendChild(e)

196                     # Done with the author.
197                     feed.appendChild(a)

199 operations = {
200         "+": ["Added", "%s was added to the repository."],
201         "-": ["Removed", "%s was removed from the repository."],
202         "U": ["Updated", "%s, an update to an existing package, was added to "
203             "the repository."]
204 }

206 def add_transaction(request, scfg, rcfg, doc, feed, txn, fmris):
199 def add_transaction(request, scfg, rcfg, doc, feed, txn):
207         """Each transaction is an entry.  We have non-trivial content, so we
208         can omit summary elements.
209         """

211         e = doc.createElement("entry")

213         tag, fmri_str = txn["catalog"].split()
214         f = fmri.PkgFmri(fmri_str)

216         # Generate a 'tag' uri, to uniquely identify the entry, using the fmri.
217         i = xmini.Text()
218         i.replaceWholeText(fmri_to_taguri(rcfg, f))
219         eid = doc.createElement("id")
220         eid.appendChild(i)
221         e.appendChild(eid)

223         # Attempt to determine the operation that was performed and generate
224         # the entry title and content.
225         if txn["operation"] in operations:
226                 op_title, op_content = operations[txn["operation"]]
227         else:
228                 # XXX Better way to reflect an error?  (Aborting will make a
229                 # non-well-formed document.)
230                 op_title = "Unknown Operation"
231                 op_content = "%s was changed in the repository."

233         if txn["operation"] == "+":
227                 c = scfg.updatelog.catalog
234                 # Get all FMRIs matching the current FMRI's package name.
235                 matches = fmris[f.pkg_name]
236                 if len(matches["versions"]) > 1:
237                         # Get the oldest fmri.
238                         of = matches[str(matches["versions"][0])][0]
229                 matches = catalog.extract_matching_fmris(c.fmris(),
230                     f.get_name(), matcher=fmri.exact_name_match)

232                 if len(matches) > 1:
233                         # Get the oldest fmri (it's the last entry).
234                         of = matches[-1]

240                         # If the current fmri isn't the oldest one, then this
241                         # is an update to the package.
242                         if f != of:
243                                 # If there is more than one matching FMRI, and
244                                 # it isn't the same version as the oldest one,
245                                 # we can assume that this is an update to an
246                                 # existing package.
247                                 op_title, op_content = operations["U"]

249         # Now add a title for our entry.
```

```
250          etitle = doc.createElement("title")
251          ti = xmini.Text()
252          ti.replaceWholeText(" ".join([op_title, fmri_str]))
253          etitle.appendChild(ti)
254          e.appendChild(etitle)

256          # Indicate when the entry was last updated (in this case, when the
257          # package was added).
258          eu = doc.createElement("updated")
259          ut = xmini.Text()
260          ut.replaceWholeText(ults_to_rfc3339_str(txn["timestamp"]))
261          eu.appendChild(ut)
262          e.appendChild(eu)

264          # Link to the info output for the given package FMRI.
265          e_uri = get_rel_path(request, 'info/0/%s' % f.get_url_path())

267          l = doc.createElement("link")
268          l.setAttribute("rel", "alternate")
269          l.setAttribute("href", e_uri)
270          e.appendChild(l)

272          # Using the description for the operation performed, add the FMRI and
273          # tag information.
274          content_text = op_content % fmri_str
275          if tag == "C":
276                  content_text += "  This version is tagged as critical."

278          co = xmini.Text()
279          co.replaceWholeText(content_text)
280          ec = doc.createElement("content")
281          ec.appendChild(co)
282          e.appendChild(ec)

284          feed.appendChild(e)

286 def update(request, scfg, rcfg, t, cf):
287          """Generate new Atom document for current updates.  The cached feed
288          file is written to scfg.repo_root/CACHE_FILENAME.
289          """

291          # Our configuration is stored in hours, convert it to seconds.
292          window_seconds = rcfg.get_attribute("feed", "window") * 60 * 60
293          feed_ts = datetime.datetime.fromtimestamp(t - window_seconds)

295          d = xmini.Document()

297          feed = d.createElementNS("http://www.w3.org/2005/Atom", "feed")
298          feed.setAttribute("xmlns", "http://www.w3.org/2005/Atom")

300          set_title(request, rcfg, d, feed, scfg.updatelog.last_update)

302          d.appendChild(feed)

304          # The feed should be presented in reverse chronological order.
305          def compare_ul_entries(a, b):
306                  return cmp(ults_to_ts(a["timestamp"]),
307                      ults_to_ts(b["timestamp"]))

309          # Get the entire catalog in the format returned by catalog.cache_fmri,
310          # so that we don't have to keep looking for possible matches.
311          fmris = {}
312          catalog.ServerCatalog.read_catalog(fmris,
313              scfg.updatelog.catalog.catalog_root)

315 #endif /* ! codereview */
```

```
316          for txn in sorted(scfg.updatelog.gen_updates_as_dictionaries(feed_ts),
317              cmp=compare_ul_entries, reverse=True):
318                  add_transaction(request, scfg, rcfg, d, feed, txn, fmris)
305                  add_transaction(request, scfg, rcfg, d, feed, txn)

320          d.writexml(cf)

322 def __get_cache_pathname(scfg):
323          return os.path.join(scfg.repo_root, CACHE_FILENAME)

325 def __clear_cache(scfg):
326          if scfg.is_read_only():
327                  # Ignore the request due to server configuration.
328                  return

330          pathname = __get_cache_pathname(scfg)
331          try:
332                  if os.path.exists(pathname):
333                          os.remove(pathname)
334          except IOError:
335                  raise cherrypy.HTTPError(
336                      httplib.INTERNAL_SERVER_ERROR,
337                      "Unable to clear feed cache.")

339 def __cache_needs_update(scfg):
340          """Checks to see if the feed cache file exists and if it is still
341          valid.  Returns False, None if the cache is valid or True, last
342          where last is a timestamp representing when the cache was
343          generated.
344          """
345          cfpath = __get_cache_pathname(scfg)
346          last = None
347          need_update = True
348          if os.path.isfile(cfpath):
349                  # Attempt to parse the cached copy.  If we can't, for any
350                  # reason, assume we need to remove it and start over.
351                  try:
352                          d = xmini.parse(cfpath)
353                  except Exception:
354                          d = None
355                          __clear_cache(scfg)

357                  # Get the feed element and attempt to get the time we last
358                  # generated the feed to determine whether we need to regenerate
359                  # it.  If for some reason we can't get that information, assume
360                  # the cache is invalid, clear it, and force regeneration.
361                  fe = None
362                  if d:
363                          fe = d.childNodes[0]

365                  if fe:
366                          utn = None
367                          for cnode in fe.childNodes:
368                                  if cnode.nodeName == "updated":
369                                          utn = cnode.childNodes[0]
370                                          break

372                          if utn:
373                                  last_ts = rfc3339_str_to_dt(utn.nodeValue)

375                                  # Since our feed cache and updatelog might have
376                                  # been created within the same second, we need
377                                  # to ignore small variances when determining
378                                  # whether to update the feed cache.
379                                  update_ts = scfg.updatelog.last_update.replace(
380                                      microsecond=0)
```

```
382                                if last_ts >= update_ts:
383                                        need_update = False
384                                else:
385                                        last = rfc3339_str_to_ts(utn.nodeValue)
386                        else:
387                                __clear_cache(scfg)
388                else:
389                        __clear_cache(scfg)

391        return need_update, last

393 def handle(scfg, rcfg, request, response):
394        """If there have been package updates since we last generated the feed,
395        update the feed and send it to the client.  Otherwise, send them the
396        cached copy if it is available.
397        """

399        cfpath = __get_cache_pathname(scfg)

401        # First check to see if we already have a valid cache of the feed.
402        need_update, last = __cache_needs_update(scfg)

404        if need_update:
405                # Update always looks at feed.window seconds before the last
406                # update until "now."  If last is none, we want it to use "now"
407                # as its starting point.
408                if last is None:
409                        last = time.time()

411                if scfg.is_read_only():
412                        # If the server is operating in readonly mode, the
413                        # feed will have to be generated every time.
414                        cf = cStringIO.StringIO()
415                        update(request, scfg, rcfg, last, cf)
416                        cf.seek(0)
417                        buf = cf.read()
418                        cf.close()

420                        # Now that the feed has been generated, set the headers
421                        # correctly and return it.
422                        response.headers['Content-type'] = MIME_TYPE
```

```
424                        # Return the current time and date in GMT.
425                        response.headers['Last-Modified'] = rfc822.formatdate()
```

```
410                        response.headers['Last-Modified'] = \
411                            datetime.datetime.now().isoformat()
427                        response.headers['Content-length'] = len(buf)
428                        return buf
429                else:
430                        # If the server isn't operating in readonly mode, the
431                        # feed can be generated and cached in inst_dir.
432                        cf = file(cfpath, "w")
433                        update(request, scfg, rcfg, last, cf)
434                        cf.close()

436        return serve_file(cfpath, MIME_TYPE)
```