# Deliverable D 3.5.1

# Automatic Test Set

| | |
|---|---|
| **Author(s):** | Juli Bakagianni, Sokratis Sofianopoulos, Stelios Piperidis (ILSP), Arle Lommel, Aljoscha Burchardt (DFKI) |
| **Dissemination Level:** | Public |
| **Date:** | 11.01.2014 (revised 10.12.2014) |

## D 3.5.1: Automatic Test Set

| Grant agreement no. | 296347 |
|---|---|
| Project acronym | QTLaunchPad |
| Project full title | Preparation and Launch of a Large-scale Action for Quality Translation Technology |
| Funding scheme | Coordination and Support Action |
| Coordinator | Prof. Hans Uszkoreit (DFKI) |
| Start date, duration | 1 July 2012, 24 months |
| Distribution | Restricted |
| Contractual date of delivery | December 2013 |
| Actual date of delivery | January 2014 (revised December 2014) |
| Deliverable number | D3.5.1 |
| Deliverable title | Automatic Test Set |
| Type | Other – Report |
| Status and version | Final v1.0 |
| Number of pages | 42 |
| Contributing partners | ILSP, DFKI |
| WP leader | DFKI |
| Task leader | ILSP |
| Authors | Juli Bakagianni, Sokratis Sofianopoulos, Stelios Piperidis, Arle Lommel, Aljoscha Burchardt |
| EC project officer | Kimmo Rossi |
| The partners in QTLaunchPad are: | Deutsches Forschungszentrum für Künstliche Intelligenz (DFKI), Germany |
| | Dublin City University (DCU), Ireland |
| | Institute for Language and Speech Processing, R.C. "Athena" (ILSP/ATHENA RC), Greece |
| | The University of Sheffield (USFD), United Kingdom |

**D 3.5.1: Automatic Test Set**

# Contents

## Update

This deliverable has been updated to clarify a few points and to correct some minor formatting issues.

**D 3.5.1: Automatic Test Set**

# 1    Introduction

This report presents the different types of processes defined and implemented for testing the layers of the MT-infrastructure prototype developed in QTLaunchPad. As reported in Deliverables 3.4.1 and 7.1.3, two layers are currently under development: one layer that caters for the data infrastructure and data (pre)processing workflows, and another one that supports MT quality evaluation workflows.

The data infrastructure and data (pre)processing workflows are supported by the META-SHARE infrastructure and its extensions for the purposes of QTLaunchpad. To this end, we have started designing and implementing at a prototype level the extension of the current META-SHARE functionalities by providing an additional language processing mechanism for processing language datasets (essentially language corpora) with appropriate natural language tools. Language processing tools are provided as SOAP web services. According to the current design, when the user selects to process a dataset, a list of all available annotation services for each relevant annotation level (e.g. tokenization & sentence splitting, pos tagging, lemmatization, alignment) are provided for the given language, and resource type. As soon as the user selects a tool, the server invokes a service that dispatches the corpus to the specific web service for processing. The system informs the user about the requested job via the messaging service of the platform. When the processing has been completed, the new (annotated) dataset is automatically stored and indexed in the repository, and the user is appropriately informed. If the user, for any reason, requests to process a dataset with a specific tool, and this dataset has already been processed by the specific tool, then the system will just forward the user to the processed dataset that has been created and stored in the repository. Section 2 reports on the automatic procedures that are being followed for testing this data documentation, storage and data processing layer of the infrastructure prototype under development.

The quality evaluation layer of the infrastructure prototype (Deliverable 3.4.1) takes care of and facilitates building quality evaluation workflows involving human and/or automatic assessment methods. It consists in a data model together with data collections, DB user interfaces and data maintenance tools for the management of MT-related data collections. This application is used to administer existing data sets, e.g. the existing WMT data, and new data, such as the data sets compiled for and by the QTLaunchPad shared task in 2014.The prototype demonstrates an entire workflow that includes at least the following steps: 1) selection of source text to be translated, 2) routing through one or more MT engines, 3) sending to Quest for quality estimation (which can then be passed to translate5), 4) error markup of all or part of the source and/or target text in translate5 based on the results of QE, 5) generation of quality scores in translate5. Section 3 presents the steps that can be followed for a simple user test of the translate5-based translation quality evaluation layer that verifies that the related markup can be added and safely retrieved.
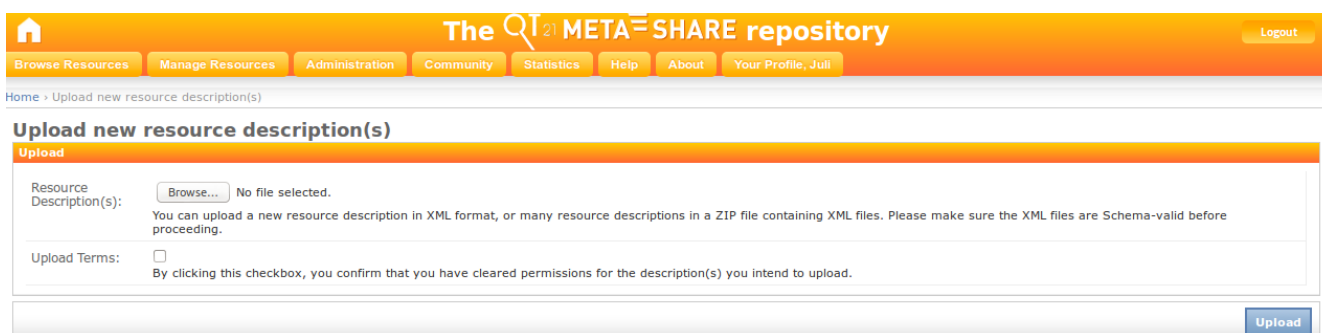
## 2. Testing the QT21 Resources Repository and its language processing extensions

Tests are created for the entire processing procedure, that is, from importing a resource description and resource actual data to processing a resource by services of various annotation levels and to the integrity of the derived resource. We have run the tests periodically during the implementation in order to guarantee the correctness of the processing mechanism.

Section 2.1 defines the tests that verify the uploading of resource descriptions, while Section 2.2 presents the tests that assert the uploading of resource actual data to the QT21 repository. Section 2.3 defines the tests that validate the stability and correctness of the processing mechanism.

### 2.1 Upload resource description

In the QT21 repository, providers can upload metadata-based resource descriptions by issuing a command to the command line interface or by using the META-SHARE graphical interface (see Figure 1). Moreover, providers can upload a single archive or a batch of archives by compressing them to a zip archive. A resource description would be successfully imported to the QT21 repository, if it is valid according to the META-SHARE metadata schema.

Figure 1

The respective tests to validate the correctness and the stability of uploading xml archives are given in Appendix 1, from line 1 to line 76. To validate that the code works as expected,
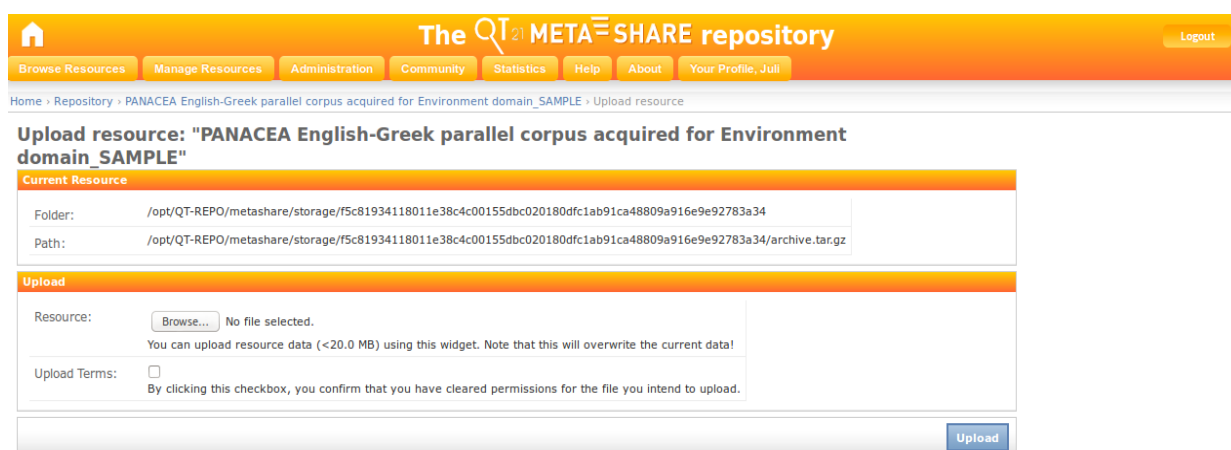
## D 3.5.1: Automatic Test Set

this collections contains tests that verify the following cases:

- import all xml files in the given directory (line 38 – line 42)

- import valid xml file (line 44 – line 47)

- import invalid xml file (line 59 – line 63)

- import valid xml files that are compressed to zip archive (line 65 – line 69)

- import invalid xml files that are compressed to zip archive (line 71 – line 76)

### 2.2. Upload data

Providers can upload the actual data of either their or their editor group resources, when the administrator can upload the actual data of any resource. The data must be compressed to an archive that belongs to one of the following file types: zip, tar.gz, gz, tgz, tar or bzip2. Currently, the maximum size of an uploaded file is 10MB. The providers can upload the data from the GUI as shown in Figure 2.

Figure 2

The tests that validate the method of uploading the actual resource data works as expected are shown in Appendix 1, from line 78 to line 267. These tests assert the following cases:

- admin user can upload data to any resource (line 154 – line 178)

- editor user can upload actual resource data  to his/her owned resources or to resources that belong to his/her editor group (line 180 – line 224)

- editor user can update actual resource data  to his/her owned resources or to

resources that belong to his/her editor group by re-uploading (line 226 – line 255)

- editor user cannot upload actual resource data to resources that do not belong to his/her editor group (line 257 – line 267)

## 2.3. Processing

In the QT21 repository, editors can process stored resources with linguistic services, which are also stored in the repository. The services are distinguished in the following annotation levels: tokenization, sentence splitting, part of speech tagging, lemmatization, dependency parsing, chunking and sentence alignment.

Currently, the QT21 repository provides processing and annotation services for Greek and English, by deploying the services developed in the framework of the PANACEA project ( http://panacea-lr.eu ). The services that are involved in the processing procedure are the ILSP NLP services and the DCU NLP services. The first collection of services is used for processing Greek text corpora complying with the xces data format, whereas the second collection is used for processing English text corpora in the txt data format.

By selecting to process a dataset with a service (see Figure 3), a list of services shows up, from which the user can choose his/her preferred one.



Figure 3

The services that are provided for annotation are relevant to the language, the data format

**D 3.5.1: Automatic Test Set**

and the annotation level of the input resource. The mapping between the specifications of an input resource and the adequate services for processing this input are given below:

- when a resource is a Greek monolingual text corpus having xces data format, then the relevant services are given in Figure 4.



Figure 4

- when a resource is an English monolingual text corpus having txt data format, then the relevant services are given in Figure 5, and

## D 3.5.1: Automatic Test Set



Figure 5

- when a resource is an English – Greek parallel corpus, then the selection of the preferred services is substantiated in two steps:

  - the selection of the annotation level (see Figure 6) and



Figure 6

  - the selection of the services that correspond to the selected annotation level (see

**D 3.5.1: Automatic Test Set**

Figure 7)

Figure 7

The tests that verify that the adequate services can process a resource are given in Appendix 1, from line 352 to line 406. Specifically, these tests assert the following cases:

- the editor selects to process a Greek monolingual text corpus, which comprises data in xces format. The services that are shown in Figure 4 must be provided to the user for selection (line 352 – line 369)

- the editor selects to process a English monolingual text corpus, which comprises data in txt format. The services that are presented in Figure 5 must be provided to the user for selection (line 371 – line 387)

- the editor selects to process an English - Greek parallel text corpus. The user selects Chunk annotation and the services that are shown in Figure 7 are provided to the user for selection (line 389 – line 406)

As soon as the user selects a tool, the server will invoke the equivalent service that will dispatch the corpus for processing. In the case of higher-level annotation tools, such as chunkers, processing is implemented using workflows. A workflow consists of a series of services, where the output of a lower annotation level service is the input of a next service of higher level. For example, the service for tagging/lemmatizing a Greek corpus is implemented as a workflow by invoking the processing chain shown in Figure 8.

Figure 8

When the processing of the resource has completed successfully, the system will:

- upload the output resource (metadata and data) to the QT21 repository (see Figure 9)



Figure 9

- inform the user about the requested job via the messaging service of the platform. The email will contain the link of the output resource on the QT21 web interface (see Figure 10).

## D 3.5.1: Automatic Test Set



Figure 10

If the processing has failed (for example the service is unavailable for a long time), the user will be informed via email message.

All tests that refer to processing procedures can be found in Appendix 1, from line 408 to line 536 and from line 668 to line 710. These tests verify the cases listed below and they are validated by the following criteria: a) an email has been sent to the user and b) the output resource has been stored to the QT21 repository.

- process Greek monolingual  text corpus with sentence splitter and tokenizer (line 408 – line 419)
- process Greek monolingual  text corpus with chunker (line 421 – line 432)
- process Greek monolingual  text corpus with lemmatizer (line 435 – line 446)
- process English monolingual  text corpus with tokenizer (line 448 – line 459)
- process English monolingual  text corpus with Part of Speech tagger (line 462 – line 473)
- process English monolingual  text corpus with parser (line 476 – line 487)
- process English – Greek parallel  text corpus with sentence splitter and tokenizer (line 490 – line 504)
- process English – Greek parallel  text corpus with Part of Speech tagger (line 507 – line 521)
- process English – Greek parallel  text corpus with sentence aligner (line 524 – line 536)
- Make concurrent calls (namely 1, 3, 5, 10 and 20 calls) to a Greek service

**D 3.5.1: Automatic Test Set**

> (tagger/lemmatizer) (line 672 – line 683)

- Make concurrent calls (namely 1, 3, 5, 10 and 20 calls) to an English service (treetagger) (line 685 – line 696)
- Make concurrent calls (namely 1, 3, 5, 10 and 20 calls) to a sentence alignment service (Berkeley aligner) (line 698 – line 710)

The annotated resource, which derives from a processing procedure, is stored automatically in the QT21 repository, while the equivalent metadata are also produced. The data of the output resource are stored following the same directory structure as the initial input data. Furthermore, each data filename takes an extension according to the annotation level.

All tests regarding the integrity of the output resource data (all files of the input resource have to be processed and all output files must have the appropriate extension) can be found in Appendix 1, from line 538 to line 632. These tests assert the following cases:

- a sentence-splitter and tokenizer processes a Greek monolingual text corpus and the derived resource data files are as many as the input data files and are named with the appropriate extension ('.tok') (line 538 – line 578)
- Part of Speech taggers-lemmatisers process a Greek - English parallel text corpus and the derived resource data files are as many as the input data files and are named with the appropriate extension ('.taglem') (line 580 – line 632)

If the user requests to process a resource with a specific web service, which has already been provided for another user, then the system will just forward the user to the processed resource that has been created in the repository (see Figure 11).



Figure 11

The test that verifies the above procedure is presented in Appendix 1, from line 634 to line 666.

# 3. User Testing of the translation quality evaluation layer

NOTE (December 2014): The following instructions refer to an older version of translate5 and are not compatible with current versions, although they can be adapted to newer versions. In addition the user credentials are no longer valid. This test set will not be maintained. Note as well that translate5 is a UI-driven tool. At the time this test set was created it was not possible to create an automated test. For more information on translate5 and annotation, please visit http://www.qt21.eu/launchpad/translate5.

1. Go to http://www.translate5.net

2. Click on Starting the demo version:

**Starting the demo version**

3. Login with the following credentials:

User name: manager
Password: asdfasdf

**translate5 Login**

**Login**

User

manager

Password

••••••••

Login

Reset password

5. Click on "Add Task" in the *Task Overview* panel:

**D 3.5.1: Automatic Test Set**



6. In the *Create task* window provide a title, specify the source (English) and target (Spanish) languages, and select the file translate5test.zip (available from http://www.qt21.eu/downloads/translate5test.zip) and then click *Add Task*.



7. Click *Add User*



8. Select "translator1" as a user and click *Add User*:

**D 3.5.1: Automatic Test Set**



9. Click *Save changes*:



10. Click *Logout* in the *Task Overview* window



11. Click *Starting the demo version* once more and repeat the login process using username "translator1" and "asdfasdf" as the password.

**D 3.5.1: Automatic Test Set**

12. When you reach the *Task Overview* window you should see a listing of available projects. One of them should have the title you selected in Step 6. Click on the "Edit task" icon (second from the left) to open the project for editing:

| | open | translate5 test | _ | 0 | Englisch (en) | Spanisch (es) |
|---|---|---|---|---|---|---|

13. In the *Segment list and editor* window perform the following edits by double clicking on segments in the *Target text reworked* column, highlighting text, and then selecting issues by clicking on the *Add QM Subsegment* button in the *QM Subsegments* pane and clicking Save after making the noted annotations:

- Segment 1:
    - Mark both instances of "the" with *Fluency:Grammar:Function words.*
- Segment 2:
    - Mark "has nothing" with *Accuracy:Mistranslation*.
- Segment 3:
    - Select the space between *both* and *lived* and then add "of us" to the comment field in the *QM Subsegments* pane and then mark the space with *Accuracy:Omission.*
- Segment 4:
    - Select "is entertaining" and mark it as *Fluency:Grammar:Tense/aspect/mood*
- Segment 5:
    - Select "either the charisma of the fashion designer", then add "the charisma of the fashion designer either" to the comment field in the *QM Subsegments* pane and then mark the text with *Fluency:Grammar:Word order.*
- Segment 6:
    - Select "Sharon Stone and John Kennedy Jr were" and mark it as *Fluency:Grammar:Word order*.

**D 3.5.1: Automatic Test Set**



*Marking issues in translate5*

The results after making all edits should like the following (the sort order of the segments is insignificant):



14. Select *Tasks: Finish task*



15. Logout from the *Task Overview* page and log back in as "manager".

**D 3.5.1: Automatic Test Set**

16. In the Task Overview page, click on the export icon (second from right, depicting a round object with a green arrow) for the project created in steps 6–9:



and select "Export"



17. After completing Step 16 a zip file will download. After unzipping this file, you should have a directory containing a single CSV file. Opening the CSV file will reveal that it contains XML markup containing the MQM issue tagging. Converting it into tabular format should reveal text equivalent to the following (allowing for whitespace differences and differences in ID values):

| mid | quelle | ziel |
|---|---|---|
| abces/2012/12/01/239729-1_es_MT2 | "Valentino prefiere la elegancia a la notoriedad" | "Valentino prefers <mqm:startIssue type="Function words" severity="critical" note="" agent="Trans Lator One" id="73"/>the <mqm:endIssue id="73"/>elegance to <mqm:startIssue type="Function words" severity="critical" note="" agent="Trans Lator One" id="74"/>the <mqm:endIssue id="74"/>fame" |
| abces/2012/12/01/239729-10_es_MT1 | Estar enamorado de la realeza no tiene nada de malo. | To be in love with royalty <mqm:startIssue type="Terminology" severity="critical" note="" agent="Trans Lator One" id="77"/>has |

**D 3.5.1: Automatic Test Set**

| | | |
|---|---|---|
| | | nothing<mqm:endIssue id="77"/> wrong. |
| abces/2012/12/01/239729-12_es_MT1 | En los años 60 y 70, ambos vivíamos en los Alpes y éramos buenos amigos. | In the 1960s and 1970s, both<mqm:startIssue type="Omission" severity="critical" note="of us" agent="Trans Lator One" id="75"/> <mqm:endIssue id="75"/>lived in the Alps and we were good friends. |
| abces/2012/12/01/239729-13_es_MT2 | Valentino es un anfitrión espectacular y entretiene con generosidad y elegancia. | Valentino is a spectacular host and <mqm:startIssue type="Tense/aspect/mood" severity="critical" note="" agent="Trans Lator One" id="78"/>is entertaining<mqm:endIssue id="78"/> with generosity and elegance. |
| abces/2012/12/01/239729-18_es_MT2 | La modelo argentina Valeria Mazza tampoco olvida el carisma del modisto. | The Argentinian model Valeria Mazza does not forget <mqm:startIssue type="Word order" severity="critical" note="the charisma of the fashion designer either." agent="Trans Lator One" id="76"/>either the charisma of the fashion designer.<mqm:endIssue id="76"/> |
| abces/2012/12/01/239729-20_es_MT2 | Éramos veinte personas, entre las que estaban Sharon Stone y John Kennedy Jr. | We were twenty people, among whom <mqm:startIssue type="Word order" severity="critical" note="" agent="Trans Lator One" id="79"/>Sharon Stone and John Kennedy Jr were<mqm:endIssue id="79"/>. |

18. Verify that the MQM tagging in the downloaded file matches the markup highlighted in red above.

**D 3.5.1: Automatic Test Set**

# Appendix I

```python
1.  class ImportTest(TestCase):
2.      """
3.      Tests the import procedure for resources
4.      """
5.
6.      test_editor_group = None
7.      super_user = None
8.
9.      @classmethod
10.     def setUpClass(cls):
11.         LOGGER.info("running '{}' tests...".format(cls.__name__))
12.         test_utils.set_index_active(False)
13.
14.     @classmethod
15.     def tearDownClass(cls):
16.         test_utils.set_index_active(True)
17.         LOGGER.info("finished '{}' tests".format(cls.__name__))
18.
19.     def setUp(self):
20.         """
21.         Set up the import test
22.         """
23.         test_utils.setup_test_storage()
24.
25.         ImportTest.test_editor_group = EditorGroup.objects.create(
26.                                 name='test_editor_group')
27.
28.         ImportTest.super_user = User.objects.create_superuser('superuser',
    'su@example.com', 'secret')
29.
30.     def tearDown(self):
31.         """
32.         Clean up the test
33.         """
34.         test_utils.clean_resources_db()
```

**D 3.5.1: Automatic Test Set**

```
35.        test_utils.clean_storage()

36.        test_utils.clean_user_db()

37.

38.    def _test_import_dir(self, path):

39.      """Asserts that all XML files in the given directory can be imported."""

40.      _files = os.listdir(path)

41.      for _file in _files:

42.        self._test_import_xml_file("%s%s" % (path, _file))

43.

44.    def test_import_xml(self):

45.      """Asserts that a single XML file can be imported."""

46.      self._test_import_xml_file('{}/repository/fixtures/testfixture.xml'

47.                    .format(ROOT_PATH))

48.

49.    def _test_import_xml_file(self, xml_path):

50.      """Asserts that the given XML file can be imported."""

51.      successes, failures = test_utils.import_xml_or_zip(xml_path)

52.      self.assertEqual(1, len(successes),

53.         'Could not import file {} -- successes is {}, failures is {}'

54.            .format(xml_path, successes, failures))

55.      self.assertEqual(0, len(failures),

56.         'Could not import file {} -- successes is {}, failures is {}'

57.            .format(xml_path, successes, failures))

58.

59.    def test_broken_xml(self):

60.      _currfile = '{}/repository/fixtures/broken.xml'.format(ROOT_PATH)

61.      successes, failures = test_utils.import_xml_or_zip(_currfile)

62.      self.assertEqual(0, len(successes), 'Should not have been able to import file
       {}'.format(_currfile))

63.      self.assertEqual(1, len(failures), 'Should not have been able to import file
       {}'.format(_currfile))

64.

65.    def test_import_zip(self):

66.      _currfile = '{}/repository/fixtures/tworesources.zip'.format(ROOT_PATH)

67.      successes, failures = test_utils.import_xml_or_zip(_currfile)

68.      self.assertEqual(2, len(successes), 'Could not import file {}'.format(_currfile))

69.      self.assertEqual(0, len(failures), 'Could not import file {}'.format(_currfile))

70.
```

**D 3.5.1: Automatic Test Set**

```
71.    def test_import_broken_zip(self):

72.        _currfile = '{}/repository/fixtures/onegood_onebroken.zip'.format(ROOT_PATH)

73.        successes, failures = test_utils.import_xml_or_zip(_currfile)

74.        self.assertEqual(1, len(successes), 'Could not import file {} -- successes is {},
       failures is {}'.format(_currfile, successes, failures))

75.        self.assertEqual(1, len(failures), 'Could not import file {} -- successes is {},
       failures is {}'.format(_currfile, successes, failures))

76.        self.assertEquals('broken.xml', failures[0][0])

77.

78. class DataUploadTests(TestCase):

79.    """

80.    Tests related to uploading actual resource data for a pure resource

81.    description.

82.    """

83.    # static variables to be initialized in setUpClass():

84.    test_editor_group = None

85.    editor_user = None

86.    editor_login = None

87.    superuser_login = None

88.    testfixture = None

89.    testfixture2 = None

90.    testfixture3 = None

91.    testfixture4 = None

92.

93.    @classmethod

94.    def setUpClass(cls):

95.        LOGGER.info("running '{}' tests...".format(cls.__name__))

96.        test_utils.set_index_active(False)

97.

98.        DataUploadTests.test_editor_group = \

99.            EditorGroup.objects.create(name='test_editor_group')

100.

101.            # test user accounts

102.            DataUploadTests.editor_user = test_utils.create_editor_user(

103.                'editoruser', 'editor@example.com', 'secret',

104.                (DataUploadTests.test_editor_group,))

105.            superuser = User.objects.create_superuser(

106.                'superuser', 'su@example.com', 'secret')
```

## D 3.5.1: Automatic Test Set

```
107.
108.         # login POST dicts
109.         DataUploadTests.editor_login = {
110.            REDIRECT_FIELD_NAME: ADMINROOT,
111.            LOGIN_FORM_KEY: 1,
112.            'username': DataUploadTests.editor_user.username,
113.            'password': 'secret',
114.         }
115.         DataUploadTests.superuser_login = {
116.            REDIRECT_FIELD_NAME: ADMINROOT,
117.            LOGIN_FORM_KEY: 1,
118.            'username': superuser.username,
119.            'password': 'secret',
120.         }
121.
122.     @classmethod
123.     def tearDownClass(cls):
124.        test_utils.clean_user_db()
125.        test_utils.set_index_active(True)
126.        LOGGER.info("finished '{}' tests".format(cls.__name__))
127.
128.     def setUp(self):
129.        test_utils.setup_test_storage()
130.        # first test resource
131.        DataUploadTests.testfixture = \
132.            _import_test_resource(DataUploadTests.test_editor_group)
133.        # second test resource which is only visible by the superuser
134.        DataUploadTests.testfixture2 = \
135.            _import_test_resource(None, TESTFIXTURE2_XML)
136.        # third test resource which is owned by the editoruser
137.        DataUploadTests.testfixture3 = \
138.            _import_test_resource(None, TESTFIXTURE3_XML,
      pub_status=PUBLISHED)
139.        DataUploadTests.testfixture3.owners.add(DataUploadTests.editor_user)
140.        DataUploadTests.testfixture3.save()
141.        # fourth test resource which is owned by the editor user
142.        DataUploadTests.testfixture4 = _import_test_resource(
```

**D 3.5.1: Automatic Test Set**

```
143.                DataUploadTests.test_editor_group, TESTFIXTURE4_XML,
144.                pub_status=INTERNAL)
145.          DataUploadTests.testfixture4.owners.add(DataUploadTests.editor_user)
146.          DataUploadTests.testfixture4.save()
147.
148.        def tearDown(self):
149.            # we have to clean up both the resources DB and the storage folder as
150.            # uploads change the storage folder
151.            test_utils.clean_resources_db()
152.            test_utils.clean_storage()
153.
154.        def test_superuser_can_always_upload_data(self):
155.            """
156.            Verifies that a superuser may upload actual resource data to any
157.            resource.
158.            """
159.            client = test_utils.get_client_with_user_logged_in(
160.                DataUploadTests.superuser_login)
161.            for res in (DataUploadTests.testfixture, DataUploadTests.testfixture2,
162.                    DataUploadTests.testfixture3, DataUploadTests.testfixture4):
163.                self.assertFalse(res.storage_object.get_download(),
164.                  "the test LR must not have a local download before the test")
165.                _data_upload_url = "{}repository/resourceinfotype_model/{}/" \
166.                    "upload-data/".format(ADMINROOT, res.id)
167.                response = client.get(_data_upload_url)
168.                self.assertEquals(200, response.status_code, "expected that a " \
169.                        "superuser may always upload resource data")
170.                with open(DATA_UPLOAD_ZIP, 'rb') as _fhandle:
171.                    response = client.post(_data_upload_url,
172.                        {'uploadTerms': 'on', 'resource': _fhandle}, follow=True)
173.                    self.assertContains(response, "was changed successfully.")
174.                # refetch the storage object from the data base to make sure we have
175.                # the latest version
176.                _so = StorageObject.objects.get(pk=res.storage_object.pk)
177.                self.assertTrue(_so.get_download(),
178.                  "uploaded LR data must be available as a local download")
179.
```

25

**D 3.5.1: Automatic Test Set**

```
180.        def test_editor_can_upload_data(self):
181.            """
182.            Verifies that an editor user may upload actual resource data to owned
183.            resources and to resources that are in her editor group.
184.            """
185.            client = test_utils.get_client_with_user_logged_in(
186.                DataUploadTests.editor_login)
187.            # make sure data can be uploaded to owned resources:
188.            self.assertFalse(DataUploadTests.testfixture3.storage_object \
189.                  .get_download(),
190.                "the test LR must not have a local download before the test")
191.            _data_upload_url = "{}repository/resourceinfotype_model/{}/" \
192.                "upload-data/".format(ADMINROOT, DataUploadTests.testfixture3.id)
193.            response = client.get(_data_upload_url)
194.            self.assertEquals(200, response.status_code, "expected that an " \
195.                        "editor may upload resource data to owned resources")
196.            with open(DATA_UPLOAD_ZIP, 'rb') as _fhandle:
197.                response = client.post(_data_upload_url,
198.                    {'uploadTerms': 'on', 'resource': _fhandle}, follow=True)
199.                self.assertContains(response, "was changed successfully.")
200.            # refetch the storage object from the data base to make sure we have
201.            # the latest version
202.            _so = StorageObject.objects.get(
203.                pk=DataUploadTests.testfixture3.storage_object.pk)
204.            self.assertTrue(_so.get_download(),
205.                "expected that uploaded LR data is available as a local download")
206.            # make sure data can be uploaded to editable resources:
207.            self.assertFalse(
208.                DataUploadTests.testfixture.storage_object.get_download(),
209.                "the test LR must not have a local download before the test")
210.            _data_upload_url = "{}repository/resourceinfotype_model/{}/" \
211.                "upload-data/".format(ADMINROOT, DataUploadTests.testfixture.id)
212.            response = client.get(_data_upload_url)
213.            self.assertEquals(200, response.status_code, "expected that an " \
214.                "editor may upload resource data to resources of her editor group")
215.            with open(DATA_UPLOAD_ZIP, 'rb') as _fhandle:
216.                response = client.post(_data_upload_url,
```

## D 3.5.1: Automatic Test Set

```
217.                   {'uploadTerms': 'on', 'resource': _fhandle}, follow=True)
218.              self.assertContains(response, "was changed successfully.")
219.              # refetch the storage object from the data base to make sure we have
220.              # the latest version
221.              _so = StorageObject.objects.get(
222.                  pk=DataUploadTests.testfixture.storage_object.pk)
223.              self.assertTrue(_so.get_download(),
224.                  "expected that uploaded LR data is available as a local download")
225.
226.          def test_editor_can_update_own_upload_data(self):
227.              """
228.              Verifies that an editor user may update actual resource data by
229.              re-uploading to owned resources.
230.              """
231.              client = test_utils.get_client_with_user_logged_in(
232.                  DataUploadTests.editor_login)
233.              # add resource data first which can be updated later:
234.              shutil.copy(DATA_UPLOAD_ZIP, u'{}/archive.zip'.format(
235.                  DataUploadTests.testfixture3.storage_object._storage_folder()))
236.              DataUploadTests.testfixture3.storage_object.compute_checksum()
237.              DataUploadTests.testfixture3.storage_object.save()
238.              # make sure that there is really resource data stored now:
239.              _so = StorageObject.objects.get(
240.                  pk=DataUploadTests.testfixture3.storage_object.pk)
241.              _old_checksum = _so.checksum
242.              self.assertNotEqual(None, _old_checksum)
243.              # vreify that the resource data update works:
244.              _data_upload_url = "{}repository/resourceinfotype_model/{}/" \
245.                  "upload-data/".format(ADMINROOT, DataUploadTests.testfixture3.id)
246.              response = client.get(_data_upload_url)
247.              self.assertContains(response,
248.                  DataUploadTests.testfixture3.storage_object.get_download())
249.              with open(DATA_UPLOAD_ZIP_2, 'rb') as _fhandle:
250.                  response = client.post(_data_upload_url,
251.                      {'uploadTerms': 'on', 'resource': _fhandle}, follow=True)
252.                  self.assertContains(response, "was changed successfully.")
253.              _so = StorageObject.objects.get(
```

## D 3.5.1: Automatic Test Set

```
254.                pk=DataUploadTests.testfixture3.storage_object.pk)
255.            self.assertNotEqual(_so.checksum, _old_checksum)
256.
257.        def test_editor_cannot_upload_data_to_invisible_resources(self):
258.            """
259.            Verifies that an editor user must not upload actual resource data to
260.            resources that do not belong to her editor group and which are not hers.
261.            """
262.            client = test_utils.get_client_with_user_logged_in(
263.                DataUploadTests.editor_login)
264.            response = client.get("{}repository/resourceinfotype_model/{}/" \
265.                "upload-data/".format(ADMINROOT, DataUploadTests.testfixture2.id))
266.            self.assertIn(response.status_code, (403, 404), "expected that an " \
267.                "editor must not upload resource data to invisible resources")
268.
269.    class ProcessingTests(TestCase):
270.
271.        @classmethod
272.        def setUpClass(cls):
273.            LOGGER.info("running '{}' tests...".format(cls.__name__))
274.            test_utils.set_index_active(False)
275.
276.        @classmethod
277.        def tearDownClass(cls):
278.            test_utils.set_index_active(True)
279.            LOGGER.info("finished '{}' tests".format(cls.__name__))
280.
281.        def setUp(self):
282.            """
283.            Creates a new storage object instance for testing.
284.            """
285.            test_utils.setup_test_storage()
286.
287.            self.user = User.objects.create_superuser("superuser",
     "su@example.com", "secret")
288.            self.client = Client()
289.            self.client.login(username=self.user.username, password='secret')
```

## D 3.5.1: Automatic Test Set

```
290.
291.        def tearDown(self):
292.            """
293.            Removes the storage object instance from the database after testing.
294.            """
295.            test_utils.clean_resources_db()
296.            test_utils.clean_user_db()
297.            test_utils.clean_storage()
298.
299.        def add_corpus(self, xmlfile, archive):
300.            '''adds corpus (metadata and data)'''
301.            res = test_utils.import_xml(
302.                '{0}/processing/fixtures/{1}'.format(settings.ROOT_PATH, xmlfile))
303.            res.storage_object.publication_status = PUBLISHED
304.            res.storage_object.update_storage()
305.
306.            data_upload_url = "{}repository/resourceinfotype_model/{}/" \
307.                    "upload-data/".format(ADMINROOT, res.id)
308.            response = self.client.get(data_upload_url)
309.
310.            with open('{0}/processing/fixtures/{1}'.format(settings.ROOT_PATH,
        archive), 'rb') as fhandle:
311.                response = self.client.post(data_upload_url,
312.                    {'uploadTerms': 'on', 'resource': fhandle}, follow=True)
313.            return res
314.
315.        def add_service(self, xmlfile):
316.            '''adds service's metadata '''
317.            res = test_utils.import_xml(
318.                '{0}/processing/fixtures/{1}'.format(settings.ROOT_PATH, xmlfile))
319.            res.storage_object.publication_status = PUBLISHED
320.            res.storage_object.update_storage()
321.            return res
322.
323.        def assert_monolingual_corpus_processing(self, _corpusxmlfile,
        _corpusarchive, _servicefile):
324.            _corpus = self.add_corpus(_corpusxmlfile, _corpusarchive)
325.            _service = self.add_service(_servicefile)
```

**D 3.5.1: Automatic Test Set**

```
326.
327.            success = ProcessingTask.delay(self.user, _corpus, _service)
328.            mail_sent, output_resource_url = success.get()
329.            response = self.client.get(output_resource_url)
330.
331.            self.assertTrue(success.successful())
332.            self.assertEquals(len(mail_sent.outbox), 1)
333.            self.assertEquals(mail_sent.outbox[0].subject, 'QT-LAUNCHPAD
       annotated resource is ready')
334.            self.assertEqual(response.status_code, 200)
335.
336.        def assert_parallel_corpus_processing(self, _corpusxmlfile,
       _corpusarchive, _firstservicefile, _secondservicefile=None):
337.            _corpus = self.add_corpus(_corpusxmlfile, _corpusarchive)
338.            _firstservice = self.add_service(_firstservicefile)
339.            _secondservice = None
340.            if not _secondservicefile is None:
341.                _secondservice = self.add_service(_secondservicefile)
342.
343.            success = ParallelProcessingTask.delay(self.user, _corpus,
       _firstservice, _secondservice)
344.            mail_sent, output_resource_url = success.get()
345.            response = self.client.get(output_resource_url)
346.
347.            self.assertTrue(success.successful())
348.            self.assertEquals(len(mail_sent.outbox), 1)
349.            self.assertEquals(mail_sent.outbox[0].subject, 'QT-LAUNCHPAD
       annotated resource is ready')
350.            self.assertEqual(response.status_code, 200)
351.
352.        def test_select_right_tools1(self):
353.            '''
354.            Verifies that the relevant services, regarding the language,
355.            the format and the annotation level of the input resource,
356.            show up
357.            '''
358.            #add PANACEA Environment Greek monolingual corpus_SAMPLE
359.            _corpus = self.add_corpus('el_corpus.xml', 'ENV_EL.tar.gz')
```

## D 3.5.1: Automatic Test Set

```
360.
361.            _corpus_id = _corpus.storage_object.identifier
362.
363.            response =
      self.client.post('{0}/repository/processing/{1}/'.format(DJANGO_URL, _corpus_id))
364.
365.            ILSP_services = ['ILSP Sentence splitting and tokenization','ILSP FBT
      Parser', 'ILSP Lemmatization',
366.                    'ILSP NLP services', 'ILSP Chunking','ILSP Dependency
      Parsing',
367.                    'ILSP Named Entity Recognition']
368.          for service in ILSP_services:
369.              self.assertContains(response, service, status_code=200)
370.
371.        def test_select_right_tools2(self):
372.            '''
373.            Verifies that the relevant services, regarding the language,
374.            the format and the annotation level of the input resource,
375.            show up
376.            '''
377.            #add PANACEA Labour English monolingual corpus_SAMPLE
378.            _corpus = self.add_corpus('en_corpus.xml', 'LAB_EN.tar.gz')
379.
380.            _corpus_id = _corpus.storage_object.identifier
381.
382.            response =
      self.client.post('{0}/repository/processing/{1}/'.format(DJANGO_URL, _corpus_id))
383.
384.            en_services = ['europarl_tokeniser', 'europarl_sentence_splitter',
      'berkeley_tagger',
385.                    'treetagger', 'berkeley_parser']
386.          for service in en_services:
387.              self.assertContains(response, service, status_code=200)
388.
389.        def test_select_right_tools3(self):
390.            '''
391.            Verifies that the relevant services, regarding the language,
392.            the format and the annotation level of the input resource,
393.            show up
```

## D 3.5.1: Automatic Test Set

394.      '''

395.      #add PANACEA English-Greek parallel corpus acquired

396.      #for Labour Legislation domain_SAMPLE

397.      _corpus = self.add_corpus('en_el_corpus.xml', 'LAB_EN_EL.tar.gz')

398.

399.      _corpus_id = _corpus.storage_object.identifier

400.

401.      response = self.client.post('{0}/repository/processing/annotation-levels/{1}/'.format(DJANGO_URL, _corpus_id),

402.      {'annotation_level' : 'chunking'}, follow=True)

403.

404.      chunkers = ['ILSP Chunker', 'berkeley_parser']

405.      for service in chunkers:

406.      self.assertContains(response, service, status_code=200)

407.

408.      def test_processing_el_corpus_with_ilsp_sst(self):

409.      '''

410.      Verifies the processing of the PANACEA Environment Greek monolingual corpus_SAMPLE

411.      by the ILSP Sentence splitter and tokenizer

412.      '''

413.      #add PANACEA Environment Greek monolingual corpus_SAMPLE

414.      _corpus_xmlfile = 'el_corpus.xml'

415.      _corpus_archive = 'ENV_EL.tar.gz'

416.      #add ILSP Sentence splitter and tokenizer

417.      _ilsp_sst = 'ilsp_sst.xml'

418.

419.      self.assert_monolingual_corpus_processing(_corpus_xmlfile, _corpus_archive, _ilsp_sst)

420.

421.      def test_processing_el_corpus_with_ilsp_chunker(self):

422.      '''

423.      Verifies the processing of the PANACEA Environment Greek monolingual corpus_SAMPLE

424.      by the ILSP Chunker

425.      '''

426.      #add PANACEA Environment Greek monolingual corpus_SAMPLE

427.      _corpus_xmlfile = 'el_corpus.xml'

**D 3.5.1: Automatic Test Set**

```
428.            _corpus_archive = 'ENV_EL.tar.gz'

429.            #add ILSP Chunker

430.            _ilsp_chunker = 'ilsp_chunker.xml'

431.

432.            self.assert_monolingual_corpus_processing(_corpus_xmlfile,
        _corpus_archive, _ilsp_chunker)

433.

434.

435.        def test_processing_el_corpus_with_ilsp_lemmatizer(self):

436.            '''

437.            Verifies the processing of the PANACEA Environment Greek
        monolingual corpus_SAMPLE

438.            by the ILSP Lemmatizer

439.            '''

440.            #add PANACEA Environment Greek monolingual corpus_SAMPLE

441.            _corpus_xmlfile = 'el_corpus.xml'

442.            _corpus_archive = 'ENV_EL.tar.gz'

443.            #add ILSP Lemmatizer

444.            _ilsp_lemmatizer = 'ilsp_lemmatizer.xml'

445.

446.            self.assert_monolingual_corpus_processing(_corpus_xmlfile,
        _corpus_archive, _ilsp_lemmatizer)

447.

448.        def test_processing_en_corpus_with_europarl_tokeniser(self):

449.            '''

450.            Verifies the processing of the PANACEA Labour English monolingual
        corpus_SAMPLE

451.            by the europarl tokeniser

452.            '''

453.            #add PANACEA Labour English monolingual corpus_SAMPLE

454.            _corpus_xmlfile = 'en_corpus.xml'

455.            _corpus_archive = 'LAB_EN.tar.gz'

456.            #add europarl tokeniser

457.            _europarl_tokeniser = 'europarl_tokeniser.xml'

458.

459.            self.assert_monolingual_corpus_processing(_corpus_xmlfile,
        _corpus_archive, _europarl_tokeniser)

460.

461.
```

## D 3.5.1: Automatic Test Set

```
462.          def test_processing_en_corpus_with_treetagger(self):

463.              '''

464.              Verifies the processing of the PANACEA Labour English monolingual
     corpus_SAMPLE

465.              by the treetagger

466.              '''

467.              #add PANACEA Labour English monolingual corpus_SAMPLE

468.              _corpus_xmlfile = 'en_corpus.xml'

469.              _corpus_archive = 'LAB_EN.tar.gz'

470.              #add treetagger

471.              _treetagger = 'treetagger.xml'

472.

473.              self.assert_monolingual_corpus_processing(_corpus_xmlfile,
     _corpus_archive, _treetagger)

474.

475.

476.          def test_processing_en_corpus_with_berkeley_parser(self):

477.              '''

478.              Verifies the processing of the PANACEA Labour English monolingual
     corpus_SAMPLE

479.              by the berkeley parser

480.              '''

481.              #add PANACEA Labour English monolingual corpus_SAMPLE

482.              _corpus_xmlfile = 'en_corpus.xml'

483.              _corpus_archive = 'LAB_EN.tar.gz'

484.              #add berkeley parser

485.              _berkeley_parser = 'berkeley_parser.xml'

486.

487.              self.assert_monolingual_corpus_processing(_corpus_xmlfile,
     _corpus_archive, _berkeley_parser)

488.

489.

490.          def test_processing_en_el_corpus_with_sentence_splitters(self):

491.              '''

492.              Verifies the processing of the PANACEA English-Greek parallel corpus
     acquired

493.              for Labour Legislation domain_SAMPLE by sentence splitters

494.              '''

495.              #add PANACEA English-Greek parallel corpus acquired
```

## D 3.5.1: Automatic Test Set

```python
496.            #for Labour Legislation domain_SAMPLE
497.            _corpus_xmlfile = 'en_el_corpus.xml'
498.            _corpus_archive = 'LAB_EN_EL.tar.gz'
499.            #add ILSP Sentence splitter and tokenizer
500.            _ilsp_sst = 'ilsp_sst.xml'
501.            #add europarl sentence splitter
502.            _europarl_sentence_splitter = 'europarl_sentence_splitter.xml'
503.
504.            self.assert_parallel_corpus_processing(_corpus_xmlfile,
       _corpus_archive, _ilsp_sst, _europarl_sentence_splitter)
505.
506.
507.        def test_processing_en_el_corpus_with_pos_taggers(self):
508.            '''
509.            Verifies the processing of the PANACEA English-Greek parallel corpus
       acquired
510.            for Labour Legislation domain_SAMPLE by POS taggers/lemmatizers
511.            '''
512.            #add PANACEA English-Greek parallel corpus acquired
513.            #for Labour Legislation domain_SAMPLE
514.            _corpus_xmlfile = 'en_el_corpus.xml'
515.            _corpus_archive = 'LAB_EN_EL.tar.gz'
516.            #add ILSP NLP services
517.            _ilsp_nlp = 'ilsp_nlp.xml'
518.            #add treetagger
519.            _treetagger = 'treetagger.xml'
520.
521.            self.assert_parallel_corpus_processing(_corpus_xmlfile,
       _corpus_archive, _ilsp_nlp, _treetagger)
522.
523.
524.        def test_processing_en_el_corpus_with_sentence_aligner(self):
525.            '''
526.            Verifies the processing of the PANACEA English-Greek parallel corpus
       acquired
527.            for Labour Legislation domain_SAMPLE by sentence aligner
528.            '''
529.            #add PANACEA English-Greek parallel corpus acquired
```

**D 3.5.1: Automatic Test Set**

```
530.            #for Labour Legislation domain_SAMPLE

531.            _corpus_xmlfile = 'en_el_corpus.xml'

532.            _corpus_archive = 'LAB_EN_EL.tar.gz'

533.            #add hunalign

534.            _hunalign = 'hunalign.xml'

535.

536.            self.assert_parallel_corpus_processing(_corpus_xmlfile,
        _corpus_archive, _hunalign)

537.

538.        def test_output_files1(self):

539.            '''

540.            Verifies that all the input resource's data has been annotated and that

541.            the output resource filenames have the right extension

542.            '''

543.            #add PANACEA Environment Greek monolingual corpus_SAMPLE

544.            _corpus = self.add_corpus('el_corpus.xml', 'ENV_EL.tar.gz')

545.            _corpus_data_path = _corpus.storage_object.get_download()

546.            #add ILSP sentence splitter and tokenizer

547.            _ilsp_sst = self.add_service('ilsp_sst.xml')

548.

549.            success = ProcessingTask.delay(self.user, _corpus, _ilsp_sst)

550.            self.assertTrue(success.successful())

551.

552.            processing = Processing.objects.get(input_resource=_corpus,
        first_tool=_ilsp_sst)

553.            output_resource = processing.output_resource

554.            output_data_path = output_resource.storage_object.get_download()

555.            input_archive = tarfile.open(_corpus_data_path, "r:gz")

556.            input_filenames = []

557.            for filename in input_archive.getnames():

558.                if len(filename.split('.')) == 1:

559.                    continue

560.                input_filenames.append(os.path.basename(filename))

561.            input_archive.close()

562.

563.            output_archive = tarfile.open(output_data_path, "r:gz")

564.            output_filenames = []

565.            _extension  = 'tok'
```

**D 3.5.1: Automatic Test Set**

```
566.            for filename in output_archive.getnames():
567.                if len(filename.split('.')) == 1:
568.                    continue
569.                output_filenames.append(os.path.basename(filename))
570.            output_archive.close()
571.
572.            self.assertEqual(len(input_filenames), len(output_filenames))
573.            #verifies that with the deletion of the output file extension
574.            #the filenames of the input and output resource are the same
575.            output_filenames_noextension = []
576.            for filename in output_filenames:
577.
       output_filenames_noextension.append(filename.replace('.{0}'.format(_extension),'') )
578.            self.assertSetEqual(set(input_filenames),
       set(output_filenames_noextension))
579.
580.        def test_output_files2(self):
581.            '''
582.            Verifies that all the input resource's data has been annotated and that
583.            the output resource filenames have the right extension
584.            '''
585.            #add PANACEA Environment Greek monolingual corpus_SAMPLE
586.            _corpus = self.add_corpus('en_el_corpus.xml', 'LAB_EN_EL.tar.gz')
587.            _corpus_data_path = _corpus.storage_object.get_download()
588.            #add ILSP NLP services
589.            _ilsp_nlp = self.add_service('ilsp_nlp.xml')
590.            #add treetagger
591.            _treetagger = self.add_service('treetagger.xml')
592.
593.            success = ParallelProcessingTask.delay(self.user, _corpus, _ilsp_nlp,
       _treetagger)
594.            self.assertTrue(success.successful())
595.
596.            processing = Processing.objects.get(input_resource=_corpus,
       first_tool=_ilsp_nlp, second_tool=_treetagger)
597.            output_resource = processing.output_resource
598.            output_data_path = output_resource.storage_object.get_download()
599.
```

## D 3.5.1: Automatic Test Set

```
600.            input_archive = tarfile.open(_corpus_data_path)
601.            input_filenames = []
602.            for filename in input_archive.getnames():
603.                if len(filename.split('.')) == 1:
604.                    continue
605.                input_filenames.append(os.path.basename(filename))
606.            input_archive.close()
607.
608.            output_archive = tarfile.open(output_data_path,"r:gz")
609.            output_filenames = []
610.            output_en_filenames = []
611.            output_el_filenames = []
612.            _extension  = 'taglem'
613.            for filename in output_archive.getnames():
614.                if len(filename.split('.')) == 1:
615.                    continue
616.                elif 'en' in filename:
617.                    output_en_filenames.append(os.path.basename(filename))
618.                elif 'el' in filename:
619.                    output_el_filenames.append(os.path.basename(filename))
620.                else:
621.                    output_filenames.append(os.path.basename(filename))
622.            output_archive.close()
623.
624.            self.assertEqual(len(input_filenames), len(output_filenames))
625.            self.assertEqual(len(input_filenames), len(output_en_filenames))
626.            self.assertEqual(len(input_filenames), len(output_el_filenames))
627.            #verifies that with the deletion of the output file extension
628.            #the filenames of the input and output resource are the same
629.            output_filenames_noextension = []
630.            for filename in output_filenames:
631.
    output_filenames_noextension.append(filename.replace('.{0}'.format(_extension),'') )
632.            self.assertSetEqual(set(input_filenames),
    set(output_filenames_noextension))
633.
634.        def test_re_processing(self):
635.            '''
```

## D 3.5.1: Automatic Test Set

636.        Verifies that when requested processing has already been executed,

637.        the user will be informed and redirected to the output resource

638.        '''

639.        #add PANACEA Environment Greek monolingual corpus_SAMPLE

640.        _corpus = self.add_corpus('el_corpus.xml', 'ENV_EL.tar.gz')

641.        #add ILSP Feature-based multi-tiered POS Tagger

642.        _ilsp_fbt = self.add_service('ilsp_fbt.xml')

643.

644.        success = ProcessingTask.delay(self.user, _corpus, _ilsp_fbt)

645.        self.assertTrue(success.successful())

646.

647.        _corpus_id = _corpus.storage_object.identifier

648.        _corpus_name = _corpus.identificationInfo.get_default_resourceName()

649.        _ilsp_fbt_name = _ilsp_fbt.identificationInfo.get_default_resourceName()

650.

651.        processing = Processing.objects.get(input_resource=_corpus,
first_tool=_ilsp_fbt)

652.        output_resource = processing.output_resource

653.        output_resource_name =
output_resource.identificationInfo.get_default_resourceName()

654.        output_data_path = output_resource.storage_object.get_download()

655.

656.        self.assertEqual(output_resource_name, '{0} annotated by
{1}'.format(_corpus_name, _ilsp_fbt_name))

657.        self.assertTrue(os.path.exists(output_data_path))

658.

659.

660.        response =
self.client.post('{0}/repository/processing/{1}/'.format(DJANGO_URL, _corpus_id),

661.                                {'tool' : _ilsp_fbt_name}, follow=True)

662.        text = 'The {0} has already been annotated by the
{1}.'.format(_corpus_name,

663.                                                _ilsp_fbt_name)

664.        self.assertEqual('repository/processing_message.html',

665.         response.templates[0].name)

666.        self.assertContains(response, text,status_code=200)

667.

668.    class processCorpusTest(unittest.TestCase):

## D 3.5.1: Automatic Test Set

```
669.            files =
   ['1.xml','2.xml','3.xml','5.xml','6.xml','8.xml','9.xml','10.xml','11.xml','13.xml','14.xml','15.
   xml','16.xml','17.xml',

670.                                '18.xml','19.xml','20.xml','22.xml','23.xml','27.xml']

671.

672.        def testEN(self):

673.            process = ProcessMonolingual()

674.                result=0

675.                homeFolder="/home/QTLP_shared/LAB_en"

676.                url="http://www.cngl.ie/panacea-soaplab2-
   axis/services/panacea.treetagger"

677.            with concurrent.futures.ThreadPoolExecutor(max_workers=10) as e:

678.            future_to_res = {e.submit(process.process(homeFolder, file, url)): file
   for file in files}

679.                for future in concurrent.futures.as_completed(future_to_res):

680.                        res = future_to_res[future]

681.                        result+=1

682.        print result

683.        self.assertEqual(result, 20, "Failed Test")

684.

685.        def testGR(self):

686.        process = ProcessMonolingual()

687.        result=0

688.                homeFolder="/home/QTLP_shared/LAB_el"

689.                url="http://nlp.ilsp.gr/soaplab2-
   axis/typed/services/getstarted.ilsp_nlp?wsdl"

690.        with concurrent.futures.ThreadPoolExecutor(max_workers=10) as e:

691.            future_to_res = {e.submit(process.process(homeFolder, file, url)): file
   for file in files}

692.            for future in concurrent.futures.as_completed(future_to_res):

693.                        res = future_to_res[future]

694.                        result+=1

695.        print result

696.        self.assertEqual(result, 20, "Failed Test")

697.

698.        def testEL-EN_sentalign(self):

699.        process = ProcessAlignment()

700.        result=0

701.                homeFolder="/home/QTLP_shared/el_en"
```

## D 3.5.1: Automatic Test Set

702.                         sentsplit_url="http://www.cngl.ie/panacea-soaplab2-axis/services/panacea.europarl_sentence_splitter"

703.                         align_url="http://www.cngl.ie/panacea-soaplab2-axis/services/panacea.berkeley_aligner"

704.            with concurrent.futures.ThreadPoolExecutor(max_workers=10) as e:

705.                future_to_res = {e.submit(process.process(homeFolder, file, sentsplit_url, align_url)): file for file in files}

706.                for future in concurrent.futures.as_completed(future_to_res):

707.                            res = future_to_res[future]

708.                            result+=1

709.        print result

710.        self.assertEqual(result, 20, "Failed Test")