



HAL
open science

Unit testing, integration and deployment : dealing with diversity, interoperability and sustainability of digital corpora

Thibault Clérice, Bridget Almas, Marie-Claire Beaulieu, Stella Dee

► To cite this version:

Thibault Clérice, Bridget Almas, Marie-Claire Beaulieu, Stella Dee. Unit testing, integration and deployment : dealing with diversity, interoperability and sustainability of digital corpora. TEI Conference and Members' Meeting, Text Encoding Initiative, Oct 2015, Lyon, France. hal-04271323

HAL Id: hal-04271323

<https://hal.science/hal-04271323>

Submitted on 8 Nov 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Unit testing, integration and deployment : dealing with diversity, interoperability and sustainability of digital corpora

Thibault Clérice¹, Bridget Almas², Marie-Claire Beaulieu², and
Stella Dee²

¹Humboldt Chair of Digital Humanities, Universität Leipzig
²Tufts University

October 2015

The Open Philology Project (OPP) at Leipzig and its US affiliate, the Perseus Digital Library at Tufts (PDL), has years of experience developing extensive infrastructures for managing textual data for historical languages. With around 100 million words available on PDL, and millions more words coming through OPP, in a context of opening contributions from wide ranging communities of users, dealing with ingestion of new texts is a matter of security, flexibility and efficiency.

Over the last few years, PDL and OPP have been moving forward in implementing the Canonical Text Service URN norm and the Epidoc subset guidelines to allow for better interoperability and citability of its texts. We are now working towards supporting a scalable workflow centered on continuous curation of these texts, from both within and outside the PDL/OPP ecosystem. Key requirements for such a workflow are ease of maintenance and speedy deployment of texts for use by a wide variety of analytical services and user interfaces.

Drawing on software engineering best practices, we are building an architecture meant for continuous integrations¹: analogous to the way Travis² integrates with Github³, we are developing a customizable service that test individual files upon each contribution made to our public git repositories. The services can be configured to test and report status on a variety of checkpoints from schema compliance to CTS ready markup.

With a strong continuous integration service, we should be able to deal not only with a wide range of genres and languages, but also with a diversity of contributors. We can delegate the tedious tasks of checking markup to the

¹https://en.wikipedia.org/wiki/Continuous_integration

²<https://travis-ci.org>

³<https://github.com>

machine, leaving curators free to focus on the scholarship. We also expect that automating checks on the integrity and the adaptability of textual objects for specific frameworks can reduce the error rate and allow for shorter feedback loops to contributors and users of our corpora.