

AcTo: How to Build a Network of Integrated Projects for Medieval Occitan

Why Medieval Occitan?

While being the ancestor of Modern Occitan, a modern language spoken by minorities in France, Spain (Catalonia), and Italy, Medieval Occitan is also, and crucially, the language of a corpus of texts fundamental for the pre-modern cultural history of Europe. Indeed the corpus of Old Occitan literature, and especially the texts of the Troubadours have had a great influence in the development of Modern European literature and beyond.

“Per solatz revelhar,
qu’es tròp endormitz,
et per Pretz, qu’es faiditz,
acolhir e tornar,
me cudèi trabalhar”
Guiraud de Bornèlh

AcTo <https://acto.hypotheses.org/>

AcTo stands for *Acolhir e Tornar*, “to welcome and return”, and is a network of data and resource centers for the study of Medieval Occitan headed by Université Paul-Valéry in Montpellier, France, which gathers together projects from different countries (France, Italy, Spain, Germany, UK).

The aim of the project is to federate existing resources (**digital editions, lexicons, but also tools**), harmonising the data and metadata encoding within the projects as well as with international standards

List of AcTo projects



AcTo objectives

- Alignment of metadata and documentation in collaboration with CIRDOC and Occitanica.eu
- Annotation and referencing of place and persons’ names in digital edition
- Copyright and legal issues
- Orthographic normalisation and lemmatisation across projects**

Current work on Orthographic Normalisation and Lemmatisation

Medieval Occitan orthography was not standardised. Digital editions, while preserving the verbatim transcription, should also allow for search by normalised forms. Harmonising the normalisation as well as the lemmatisation choices is a crucial prerequisite for a federated search throughout all existing corpora.

Current experiment on automatically analysing the text based on

- Rule based methods
- Machine learning methods re-training models derived from Old French



Link each token in the TEI digital edition to the reference lexicon, Dictionnaire d’Occitan Médiéval (DOM).



Pyrrha Dashboard

Thalamus.test

Corpus Thalamus.test - List of tokens

Id	Form	Lemma	POS	Morph	Context	Similar	Save
1	L	lo2	DETdef	NOMB=>[GENRE=m][CAS=r]	L' an MLXXXVII, los crestians prezeron Barasona. L'		Save
2	an	an	NOMcom	NOMB=>[GENRE=m][CAS=r]	L' an MLXXXVII, los crestians prezeron Barasona. L' an		Save
3	MLXXXVII	lbrum8	ADUcar	NOMB=>[GENRE=m][CAS=r]	L' an MLXXXVII, los crestians prezeron Barasona. L' an MCLXXXII		Save
4	-	-	PONtbl	MORPH=empty	L' an MLXXXVII, los crestians prezeron Barasona. L' an MCLXXXII		Save
5	los	lo2	DETdef	NOMB=>[GENRE=m][CAS=n]	L' an MLXXXVII, los crestians prezeron Barasona. L' an MCLXXXII		Save
6	crestians	crestian	NOMcom	NOMB=>[GENRE=m][CAS=n]	L' an MLXXXVII, los crestians prezeron Barasona. L' an MCLXXXII		Save
7	prezeron	prendre	VERc3g	MODE=ind[TEMPS=pp][PERS=3][NOMB=pp]	L' an MLXXXVII, los crestians prezeron Barasona. L' an MCLXXXII		Save

OMÉLiE project
Outils et méthodes pour l’édition linguistique enrichie

- Use of a post-correction environment to create a training and test corpus

AcTo and



- Better metadata
- Visibility of resources in the VLO
- Federated content search
- Use of best practices and standards (TEI, LMF, ...)
- Collaboration for the development of NLP pipelines

Gilda Caïti-Russo Laboratoire LLACS Univ Paul-Valéry Montpellier
Jean-Baptiste Camps Centre Jean-Mabillon École nationale des chartes Université PSL, Paris
Gilles Couffignal Université Paris-Sorbonne
Francesca Frontini Laboratoire PRAXILING Univ Paul-Valéry Montpellier
Hervé Lieutard Laboratoire LLACS Univ Paul-Valéry Montpellier
Elisabeth Reichle Ludwig Maximilian University of Munich
Maria Selig Universität Regensburg