# WG 2 "Retro-digitized dictionaries"

Progress Report

# WORKING GROUP 2 IN NUMBERS

- 50 members
  - Gender Balance: 22 female, 28 male, i.e. about 44 % of members are female
  - Early Career Investigators: 18 ESRs, i.e. about 36 % of members are ESRs
- Internationality: 20 countries (Austria (2), Belgium (1), Croatia (1), Denmark (3), Finland (3), France (4), Germany (6), Greece (2), Hungary (1), Israel (1), Italy (1), Macedonia (1), Netherlands (2), Poland (3), Portugal (4), Romania (4), Serbia (5), Slovakia (2), Switzerland (3), United Kingdom (1))

# OBJECTIVES OF WORKING GROUP 2

WG 2 will set up guidelines and standards for turning paper dictionaries into a digital format and develop common models in the field of e-lexicography for retro-digitised paper dictionaries already online or planning to go online.

To reach this goal WG 2 will

1. establish an overview of existing retro-digitised dictionaries and dictionaries which should be retro-digitised
2. create guidelines defining standard workflows for the digitisation of dictionaries
3. establish best practices for the encoding of information and the description of relevant information categories for paper dictionaries
4. establish best practices for dictionary enrichment and linking

# CONCRETE STEPS – BLOG *DIGILEX. LEGACY DICTIONARIES RELOADED*

# CONCRETE STEPS – BLOG *DIGILEX*

The blog posts reflect upon the central topics of WG 2 by

- describing the workflows for the digitisation of different dictionaries, stressing challenges, problems and solutions
- tackling encoding problems and thus exposing encoding standards
- dealing with issues of dictionary enrichment and dictionary linking.

https://digilex.hypotheses.org/



## Legacy Dictionaries Reloaded: Why Should We Bother?

The closest I've ever come to glimpsing hell was several years ago, reading an article in the New York Times, entitled "Justices Turning More Frequently to Dictionary, and Not Just for Big Words." The article cited the example of a certain Chief Justice John G. Roberts Jr. who had apparently parsed the meaning of a federal law by consulting the usual legal precedents (X vs Y) — and no less than five dictionaries. One of the words he looked up was "of". He discovered, lo and behold, that its meaning had something to do with belonging or possession.

# CONCRETE STEPS – BLOG *DIGILEX*

The blog will
- enhance the visibility of WG 2 and thus will contribute to the dissemination of the network's activities
- integrate the members in the working processes of WG 2 and thus help to achieve the defined goals of WG 2.
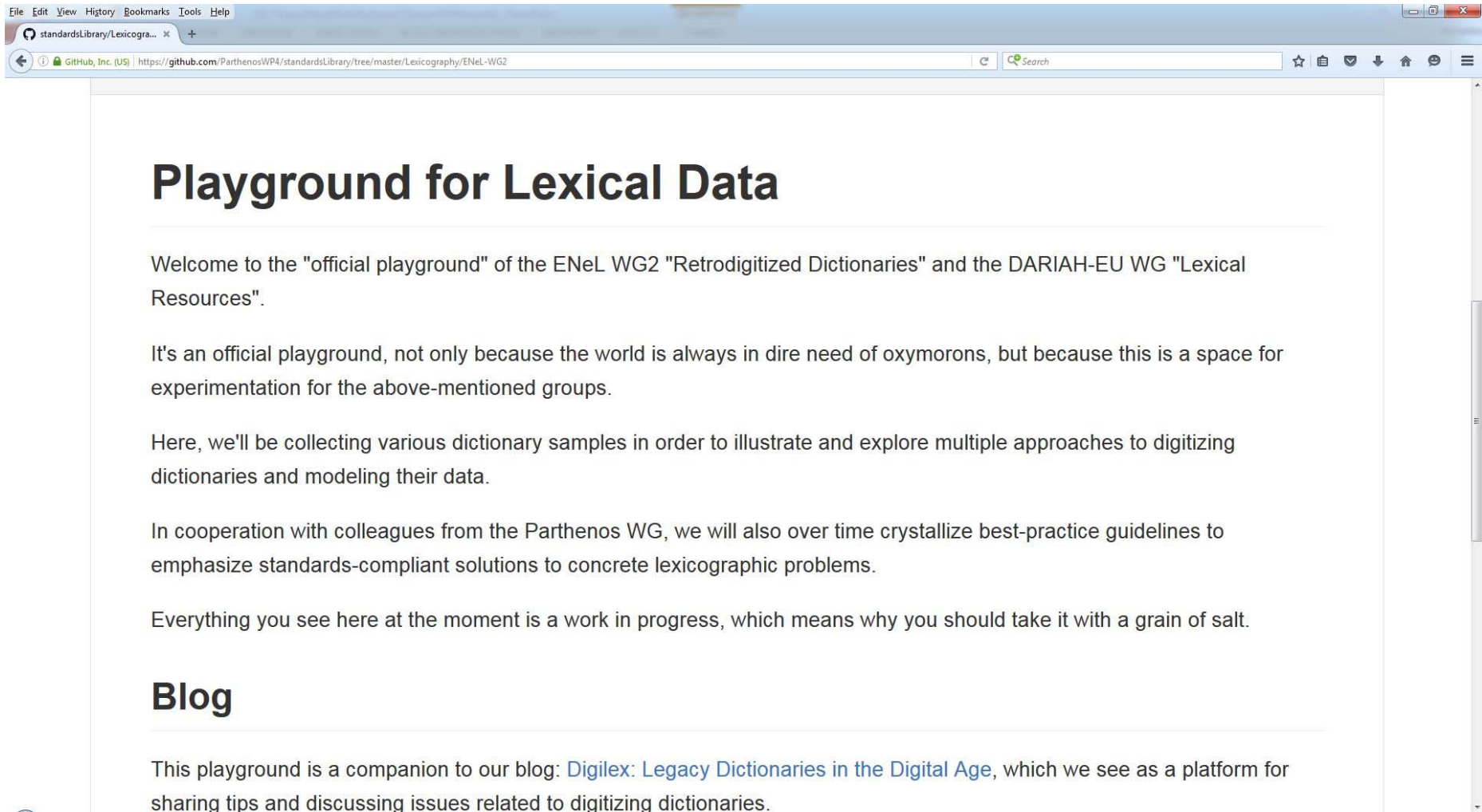
https://digilex.hypotheses.org/

## How Can I OCR My Dictionary?

One way to digitise a dictionary is using Optical Character Recognition or OCR. But is OCR feasible at all for my dictionary? And if so, which OCR program should I used, trainable or omnifont? And how about the workflow: should I train the OCR engine or not? And, finally, what should be the output format of my OCR? For those wanting to take the OCR adventure, here a very brief introduction.

## To OCR or not to OCR

Some texts are totally un-ocr-able:

# CONCRETE STEPS: PLAYGROUND ON GITHUB

# CONCRETE STEPS: COLLECTION OF DICTIONARY SAMPLES

# CONCRETE STEPS: CREATING AND DELIVERING GUIDELINES

https://github.com/ParthenosWP4/standardsLibrary/tree/master/Lexicography/ENeL-WG2

# CONCRETE STEPS: COOPERATION WITH #DARIAHTEACH

- #dariahTeach is producing an online platform for the delivery of high-quality Digital Humanities training materials.

- A separate online module on retro-digitizing dictionaries will be completed and published by the summer of 2017. It will be officially co-branded with ENeL.

# BARCELONA MEETING

- advanced follow-up to the Lisbon Training School in July 2015
- mixture of short presentations and hands-on sessions
  - offered additional instruction in more advanced topics that were not covered in Lisbon
  - addressed the issue of technical heterogeneity in lexical sources in a joint session with WG 4 chaired by Laurant Romary
  - worked on concrete encoding issues, challenges and problems

Keep doing what we are doing!