# A scientist's everyday toolkit for transparent research

*…or improving research with OSF and GitHub!*

Georgia Loukatou

OSF

GitHub

Reproducible workflows*

- Personal profile, with file repositories for projects
- Time-stamping and contribution / version tracking
- Can be shared with collaborators
- Private or public
- Free

*Preserve current workflow of a project: all our material gathered in one place (open data and code), with documented procedure.
→ Make sure that scripts can be reused by us or others and results reproduced.

GitHub

# Reproducible Workflows

*Imagine having a hard time locating files of a project…*
*(data in DropBox, scripts in Documents, 10 versions of manuscript draft in Downloads)*
*…and can't find raw data anymore?*

- Make life easier
- Are proof of quality (which brings citations!)
- Give credit to those who actually did the work
- Advance research

# OSF : Open Science Framework

- [Homepage](#) (can be used as a personal page)

# OSF : Open Science Framework

- <u>Example1</u>: project in the making
  - Contributors: team members have access to everything all at once (no back and forth emails)
  - Components: data and code, with separate settings (i.e. if. database must stay private)
  - Privacy setting: supports both public and private (default) projects

# OSF : Open Science Framework

- Example1: project in the making
  - Size limit: practically none
  - **Add-on links: GitHub, Dropbox** …
  - Transparency: everything documented

# OSF : Open Science Framework

- Example 2: finished project
  - Reproducible pipeline: documentation, data and scripts, from cleaning data to results
  - Public

# OSF : Open Science Framework

- Example 2: finished project
  - File version

# View-only Link

When submitting a manuscript, create a view-only link leading to its OSF page:

- Link can be anonymous
- Proof of quality
- Appreciated by reviewers - one explicitly asked for it!
-

## Submitted manuscript:

Finally, Phonotactics from Utterances Determine Distributional Lexical Elements (PUDDLE) is an incremental alternative algorithm (Monaghan and Christiansen, 2010), where learners build a lexicon by entering every utterance that cannot be broken down further, and using such entries to find subparts in subsequent utterances.

WordSeg was used both for segmentation and evaluation. Each algorithm returns their input with spaces where the system hypothesizes a break. Evaluation is done with reference to word boundaries. Scripts used for corpus preprocessing and segmentation as well as results and supplementary material are available at https://osf.io/6q5e3/?view_only=d29bc605d45e4f be9a79508e456350e0.

man correlations (median=.42, range from -.15 to .98) suggested that there is a similar rank ordering of algorithm performance across languages. Inuktitut and Russian were the only languages not following the general ordering. The average correlation within our rough morphological groups (i.e., high-high, moderate-moderate, low-low) was .55, greater than across groups (.3).

### 4 Discussion

First, no algorithm performed systematically below chance level in our study. However, we cannot say that they all performed above chance for all languages either. This is mainly due to the good results in baseline p=0, especially salient for morphologically complex languages such as Inuktitut.

# Accepted version

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| MIr | 7/8 | 7/8 | 27 | 7 | Inu | 36 | Ind |
| FTPr | 7/8 | 7/8 | 25 | 11 | Inu | 30 | Rus |
| PUD | 6/8 | 6/8 | 22 | 7 | Ind | 34 | Ses |
| BTPa | 6/8 | 6/8 | 17 | 10 | Ses | 27 | Ind |
| MIa | 7/8 | 8/8 | 17 | 15 | Jap | 25 | Inu |
| BTPr | 6/8 | 5/8 | 14 | 9 | Inu | 22 | Yuc |
| Base0 | - | 1/8 | 13 | 6 | Tur | 35 | Inu |
| Base6 | 7/8 | - | 12 | 8 | Tur | 16 | Inu |

Table 2: Number of languages performing above baseline p=0 and p=1/6. Columns show the mean, the lowest and highest percentage of correctly segmented word tokens for each algorithm and the corresponding language. Languages are represented by the first three letters of their names. "PUD" stands for PUDDLE. "Base0" and "Base6" stand for baseline p=0 and p=1/6.

$xy$ divided by the product of the frequency of $x$ and that of $y$; the version in WordSeg draws from Saksida's implementation (Saksida et al., 2017). Whether to add a word boundary or not depends on a threshold, which can be based on a local comparison (*relative*, where one cuts if the TP or MI is lower than that for neighboring sequences); or a global comparison (*absolute*, where one cuts if the transition is lower than the average of all TP or MI over the sum of different phoneme bigrams). It should be noted that previous authors originally implemented TPs on syllables (Saksida

| | | | | | |
|---|---|---|---|---|---|
| Russian | 22 | 7 | AG | 31 | FTPa |
| Yucatec | 27 | 16 | MIa | 48 | AG |
| Sesotho | 24 | 9 | BTPr | 39 | AG |
| Indonesian | 29 | 7 | PUD | 65 | AG |
| Japanese | 26 | 14 | BTPa | 43 | AG |

Table 3: Mean percentage of correctly segmented word tokens for each language. Languages are listed in rough order of morphological complexity (see Table 1). Columns show the mean, lowest and highest percentage of correctly segmented word tokens per language, and the corresponding algorithm. "PUD" stands for PUDDLE.

subparts in subsequent utterances.

WordSeg was used both for segmentation and evaluation. Each algorithm returns their input with spaces where the system hypothesizes a break.[1] Evaluation is done with reference to orthographic word boundaries. Scripts used for corpus preprocessing and segmentation as well as results and supplementary material are available at https://osf.io/6q5e3/.

## 3 Results

Results are shown in Tables 2 (reporting on algorithms) and 3 (reporting on languages). Next, we address our research questions.

# It's easy and fast!

• Create a project,

• add contributors,

• make components,

• connect to GitHub, Dropbox and

• create a view-only link

• All these take 5 minutes! ☺

• TUTORIAL

# OSF Preprints

- complete manuscripts shared with a public audience without peer review.
  - Often, preprints are also submitted for peer review and publication in traditional scholarly journals.

- Why?
  - Paper too long to print, or the reviewing process never-ending, need fast feedback, or don't want to publish after all, but still get it out there!
  - **A preprint gives more exposure to your research (DOI, indexed by Google Scholar)**

- How?
  - OSF Preprints linked with communities such as PsyArXiv, BiorXiv, engrXiv…

  ⚠ just make sure before posting what is the policy of the journal you are submitting to/ were accepted to

# OSF Preregistrations

- support "registering" projects which freezes your project at a particular point in time and gives you an archival location for that version

## How to preregister

| Study Information | Sampling Plan | Variables | Study Design | Analysis plan |
|---|---|---|---|---|

| Study Information | Sampling Plan | Variables | Study Design | Analysis plan |
|---|---|---|---|---|
| • Title<br>• Authorship<br>• Research question<br>• Hypothesis | • Existing data<br>• Data collection procedure*<br>• Sample size rationale** | Manipulated or measured Variables | • Study type<br>• Blinding<br>• Randomization | • Statistical models (+interactions!)<br>• Follow-up analysis<br>• Inference criteria***<br>• Data exclusion<br>• Exploratory analysis |

***p-values, Bayes factors, specific model fit indices...

OSF, and many other templates out there, most minimal: AsPredicted.org Example

# Anything better than Preregistration?

- A **Registered Report** – your paper is accepted whatever the result!

- List of journals accepting RRs here: https://cos.io/rr/ e.g.: Brain and Behavior, Cognition and Emotion, Cortex, Nature Human Behavior, Psychological Science…

- Benefits:
  - reduced publication bias as negative results will not prevent publication.
  - authors receive constructive critical feedback prior to conducting the experiment.
  - enhances the credibility of the work.
  - [Reviewers more generous towards the work]

# Preregistered report

## How it works



Stage 1: REVIEW
- Authors submit a paper to review with no results and no discussion

In principle acceptance

Stage 2: REVIEW
- Authors add results and discussion

acceptance

In press
- Publish!

# GitHub

Is a repository hosting tool, that helps avoid this:

analyse_.sh
analyse_1!.sh
analyse_1.sh
analyse_2.sh
analyse_3.sh
analyse_4.sh
analyse_4last.sh
analyse_4lastTHISISTHEFINALONE.sh
analyse_final.sh
analyse_last.sh
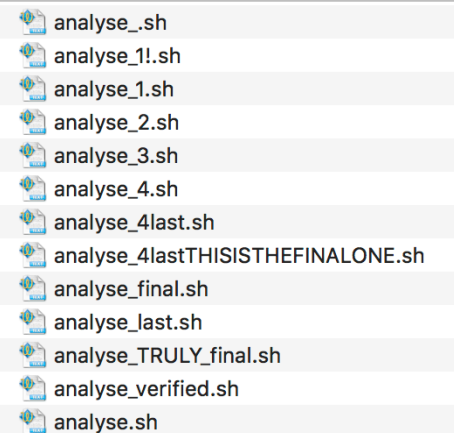analyse_TRULY_final.sh
analyse_verified.sh
analyse.sh

- When working on a project, it's hard to follow revisions (especially if you have collaborators!) -who changed what, when, and where those files are stored.
- GitHub keeps track of all the changes that have been pushed to the repository.
- **Version control** is the only reasonable way to keep track of changes in code, manuscripts, presentations, and data analysis projects.
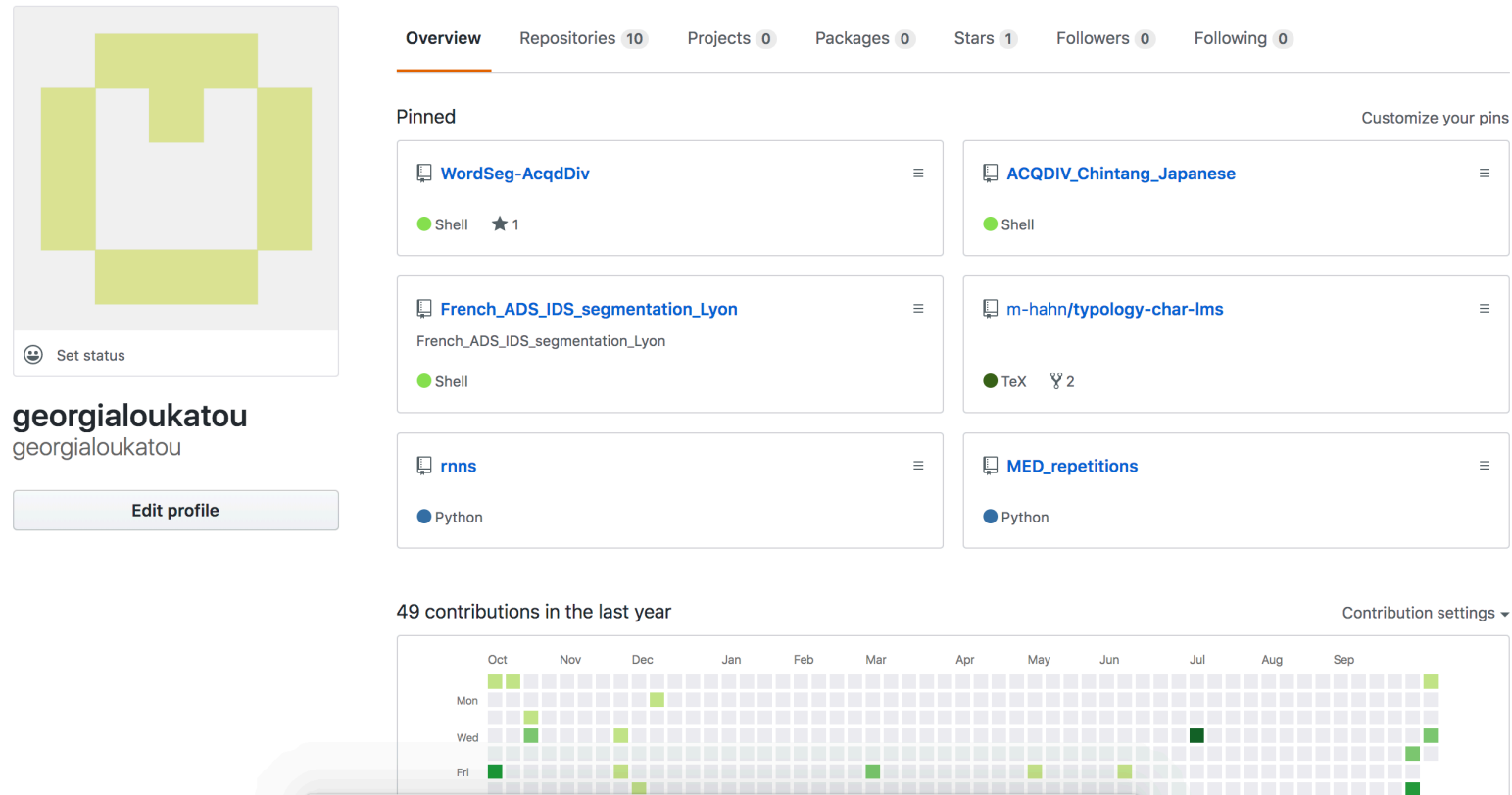
GitHub

# How does GitHub work?

Master branch (repository)

Create copy of project, new branch called *test*

Commit changes to *test*

Pull request: Propose to integrate changes in *test* to master

Review changes

Merge changes with Master

Even after merge, previous versions of a file can still be retrieved! –version history

Example: https://github.com/georgialoukatou

# GitHub example

# It's easy and fast!

- Create a repository

- Write a file

- Create a branch

- Commit a change

- Pull request and

- Merge

- All these take 5 minutes!  ☺

- TUTORIAL

# GitHub for programmers

- Moreover, it is one of the largest coding communities around, so using it can provide wide exposure for your project. *Github is like Facebook for programmers.*

- And it is command line friendly!

# Using the terminal

- git clone my-repo # clone a repository

- cd repo  # change into the `repo` directory and write a file, test.txt

- git add test.txt # git isn't aware of the file, stage it

- git commit –m "add FILE to initial commit" #take a snapshot of the staging area

- git push origin master # push changes to repository

- git branch my-branch # create new branch

- git checkout my-branch # switch to the branch and make changes

- git add FILE  # stage the changed file

- git commit –m "new change"

- git push origin my-branch # push changes to github

# Thank you!

georgialoukatou@gmail.com