# PARSEME

# SHARED TASK ANNOTATION TOOLS

Struga, 8 April 2016

# PHASE 2

# PHASE 2
# VALIDATION

# PHASE 2
# VALIDATION

## http://tiny.cc/parsemeValidate



**ParsemeValidatePhase2.jar**

`java -jar ParsemeValidatePhase2.jar Language`

# PHASE 2
# VALIDATION

**java -jar ParsemeValidatePhase2.jar Swedish**

```
Annotation with double annotate root type or not root
type?! or even using the same number for two different
annotations...
** -- File --> Pilot2 ST - Swedish - Annotator 1.tsv
** -- At line number --> 188
** -- The Token is --> 16   lägga   1  VPC
```



**Italian, Hebrew, Romanian, Swedish, Turkish,
English, Portuguese, Maltese, Spanish, Greek**

# PHASE 2
# VALIDATION

`java -jar ParsemeValidatePhase2.jar Polish`

All the annotations seem to be in good order!

**Polish, German, French, Hungarian**

# PHASE 2
# MERGING

# PHASE 2
# MERGING

| Rank | Token | No space (nsp) | MTW identifiers | First | |
|------|-------|----------------|-----------------|-------|---|
| | | | | MWE identifier | Type |
| 1 | W | | | | |
| 2 | stanie | | | | |
| 3 | obrzydzenia | | | | |
| 4 | przyprawiającego | | | [1, 1] | IPrepV, IPrepV |
| 5 | o | | | [1, 1] | |
| 6 | nowe | | | | |
| 7 | mdłości | | | | |
| 8 | nie | | | | |
| 9 | zauważył | | | | |
| 10 | nawet | nsp | | | |

| Rank | Token | No space (nsp) | MTW identifiers | First | |
|------|-------|----------------|-----------------|-------|---|
| | | | | MWE identifier | Type |
| 1 | Jakiś | | | | |
| 2 | czas | | | | |
| 3 | mierzyli | | | [null, 1] | [null, ID] |
| 4 | się | | | [null, 1] | |
| 5 | wzrokiem | nsp | | [null, 1] | |
| 6 | . | | | | |

# ANNOTATION TOOLS

# TOOL SELECTION PROCESS

BRAT
CAT
FLAT
CATMA
WebAnno
MAT
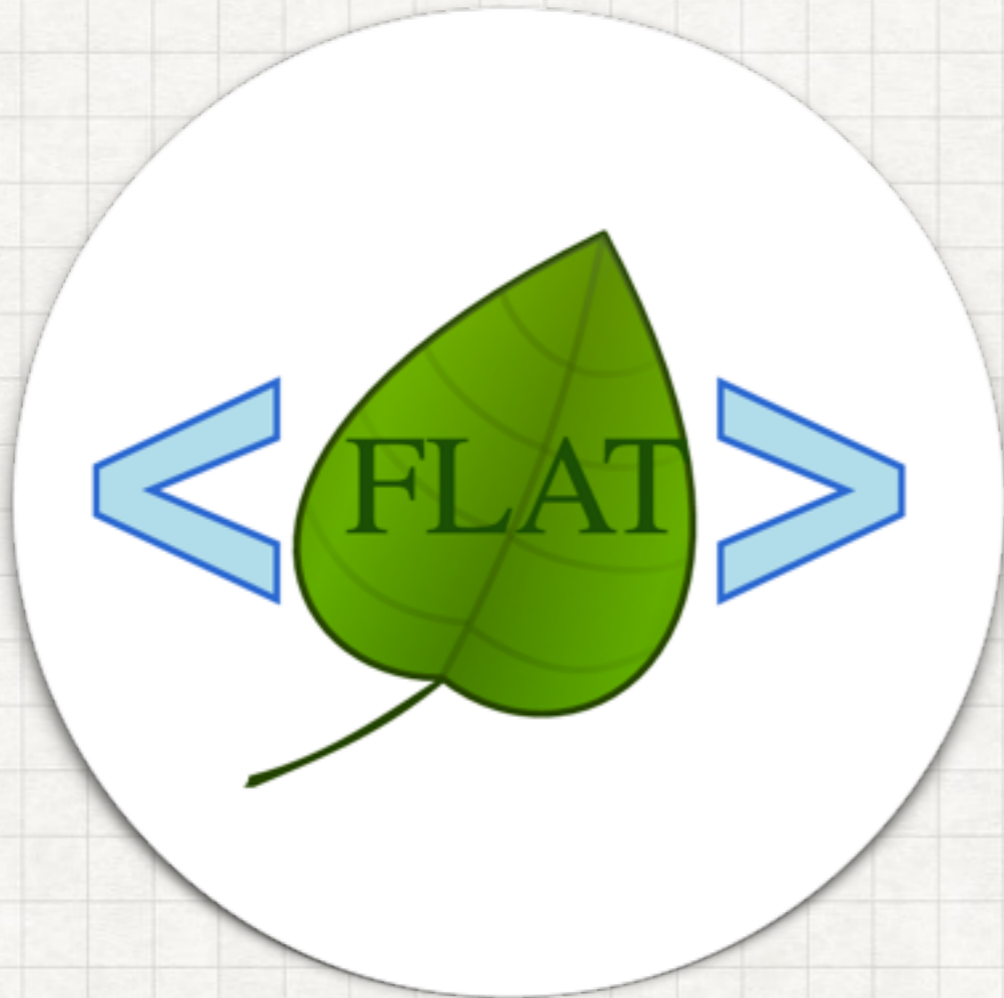Pub annotation
TELEGRAM-BOT

↓

BRAT
CAT
FLAT
TELEGRAM-BOT

# TOOL SELECTION PROCESS

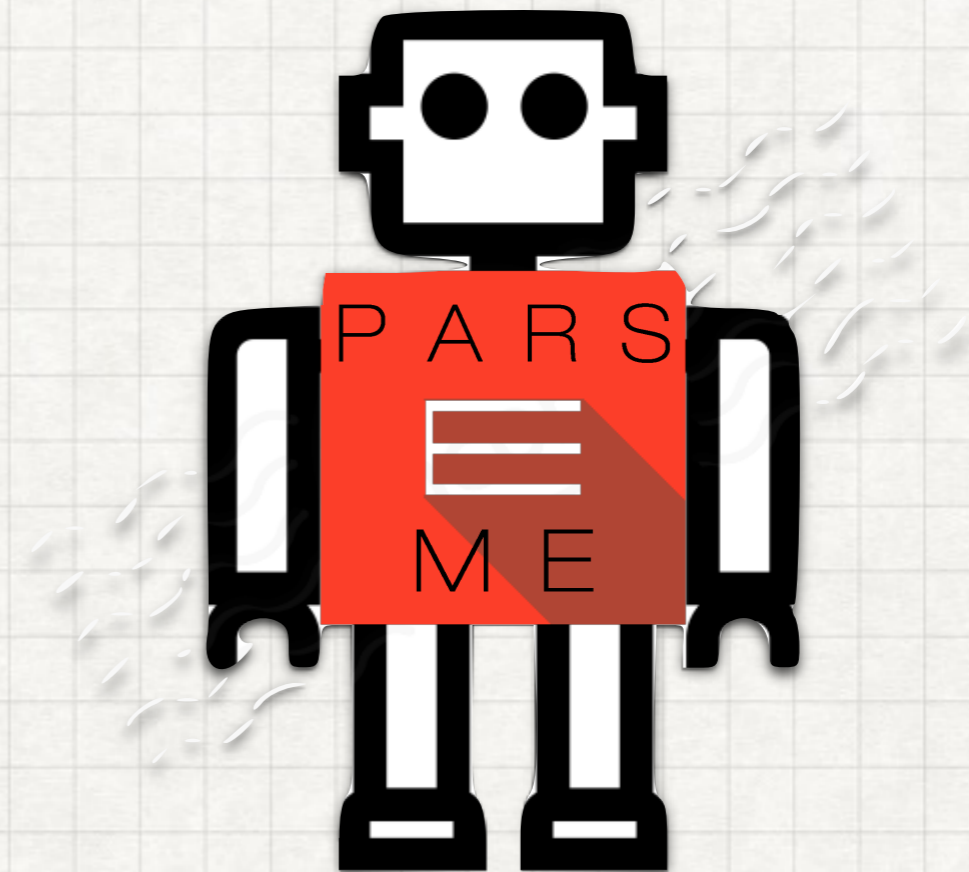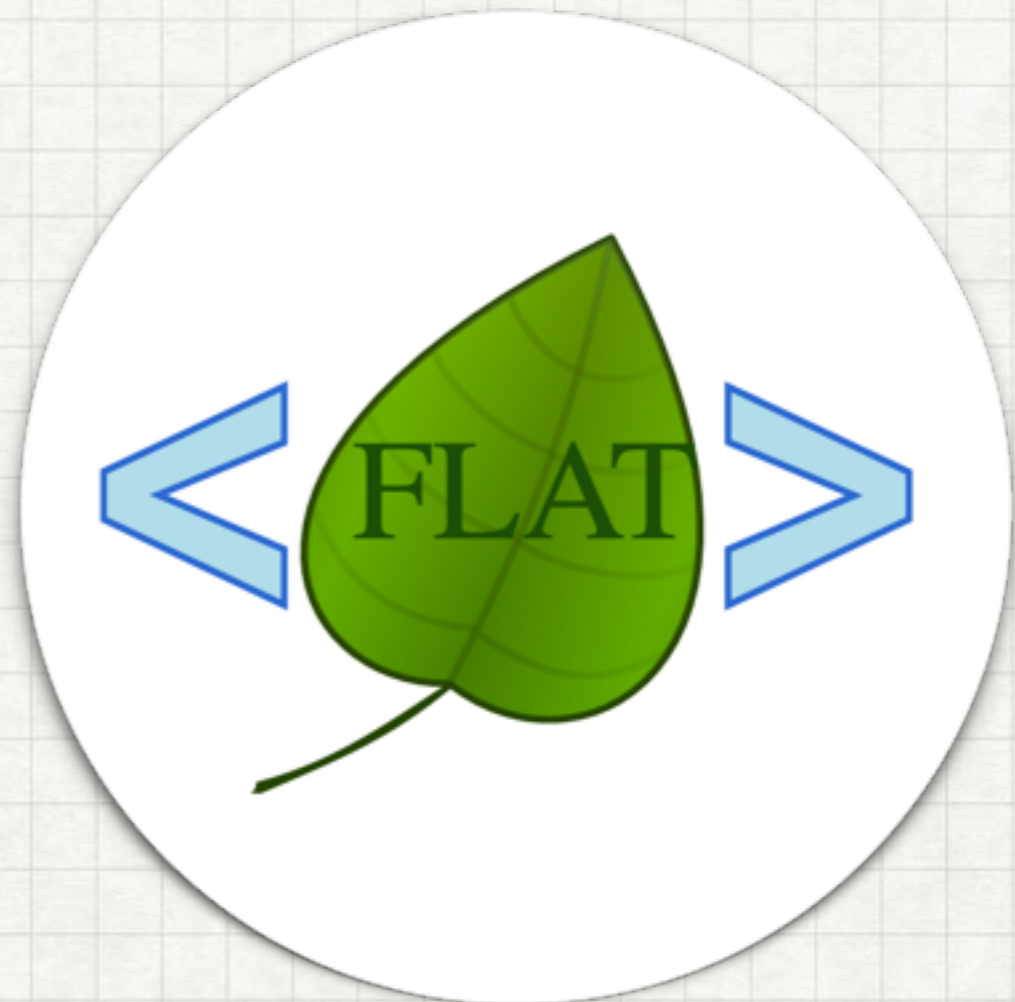| | BRAT | CAT | FLAT | TELEGRAM-BOT |
|---|---|---|---|---|
| **Tool:** | **BRAT** | **CAT** | **FLAT** | **TELEGRAM-BOT** |
| **WebSite:** | http://brat.nlplab.org/ | https://dh.fbk.eu/resources/cat-content-annotation-tool | https://github.com/proycon/flat/ | https://telegram.me/Parseme_Bot |
| **Demo:** | weaver.nlplab.org/~smp/brat-rtl/ | :8080/CAT_WEB_APP_i1 (username: sangati password: | http://flat.science.ru.nl/ | https://telegram.me/Parseme_Bot |
| **DECISION** | **Eliminated due to no support discontinuous overlapping MWEs, and heavy development effortneeded for customization** | **Eliminated due to no guarantee of support while the code is not open source.** | **Retained for a demo in Struga** | **Retained for a demo in Struga** |
| Is your tool **utf-8 compatible** (deals with different alphabets: Latin, Cyrillic, diacritics, etc.)? | Y | Y | Y | Y |
| Can it handle **right-to-left** languages? | Y (currently working on it, see Behrang conversation) | N | Not yet | Y |
| Server Needed | Y | N | Yes (central distribution repository - easy to upgrade when new features are available) | Y (via Google Application Engine - it should be free given that required resources are probably below standard quota but costs need to be verified) |
| Can it handle **predefined tokens**, that is to say, is there a way for the tool to accept a pre-tokenized sentence? This units will be the minimal-blocks the users can select, so ideally we don't want the selection to be on a sub-span of a token. | Pre-tokenisation is supported and the annotator can double-click on a token to select it in its entirety | Yes, but not below word level (spaces are inserted between tokens) | Y | Y (in some devices the highlighting of selected tokens may be limited) |
| Can it let the user assign to the selected sequence a category belonging to a set of **predefined categories**? | Y | Y | Y | Y |
| Can it let the user select a **discontinuous** sequence of tokens (non-adjacent tokens)? | Only via word relations ? | Y | Y | Y |
| Can it let the user select two or more **overlapping sequences** of tokens (and assign to each of them a separate category)? | Y | Yes (but rendering is not optimal) | Y | Y (not in the current demo, but can be implemented [SIMPLE]) |
| Does it have an **authentication procedure**, and can it handle 50-60 distinct users? | Y | Yes, but not yet tested with so many users | Yes, but not yet tested with so many users | Y |
| Can we define **user roles** (e.g., project managers, language managers, annotators)? | N | No (each user has a separate workspace) | Limited (only setting read or write permissions on workspaces) | Y (not in the current demo, but can be implemented [MODERATE]) |
| Does the tool provide the language managers with an **interface to upload a corpus and assign to each annotator a specific set of sentences**? Ideally we would need a sub-set of sentences to be annotated by at least two annotators to compute inter-annotation agreement. | No, though we have some functionality for comparisons. | N | Yes, but only manually as a preprocessing step | Y (in principle can be implemented [MODERATE+]: a user can submit a file to the bot, but might be wiser to do this manually) |
| Does it allow to **export the data** in a standard annotation format (xml, json, tsv)? | Y | Y | Y | Y |
| Does the tool include a **revision phase** where a subset of sentences annotated by two users goes to an **adjudicator** user (different from the previous two) who will need to select the correct annotation? | N | N | No, but willing to implement a separate user interface for that. | Y (not in the current demo, but can be implemented [MODERATE]) |
| Could we count on a reference person of the annotation tool who would be **willing to help** in the process of setting up the project and customizing the tool for our need. | We still respond to questions and feature requests. Most likely not be able to dedicate enough time for any major coding efforts. | Y/N (willing to answer questions but not solving issues or implementing new features) | Y (eager to implement new features and improve existing ones) | Y (eager to implement new features and improve existing ones) |
| Other comments | | - documents cannot be too big (need to split them) | - can be integrated with UCTO (tokenizer) to tokenize sentence or convert pre-tokenized sentence in FoLiA format<br>- All changes history saved, user can revert them<br>- document based (although you can activate sentence perspective where each sentence is separated from the next)<br>- documents cannot be too big (need to split them) | - the interface is rather modest (basically a messaging app with only text and buttons) and cannot be made fancier, but depending how you look at it, it may be an advantage as it imposes extra pressure to make the logic of the system tidy and as simple as possible. |

# TOOL SELECTION PROCESS
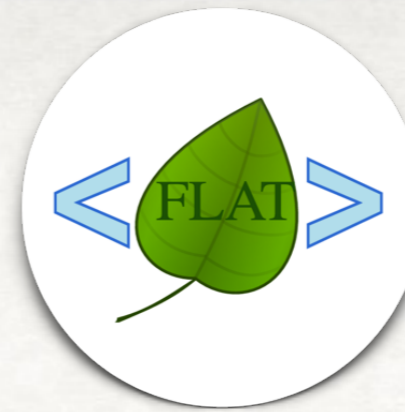
## FLAT

## TELEGRAM-BOT

# FLAT

- **FLAT**: web-application that offers an interface for the visualization and editing of FoLiA documents.

- **FoLiA**, a Format for Linguistic Annotation developed in the scope of the CLARIN-NL project and other projects.

- **Maarten van Gompel** (Centre for Language and Speech Technology Radboud Universiteit Nijmegen)

# FLAT

# FLAT

http://flat.science.ru.nl/register/

## Create an account

**Username:** [                    ]

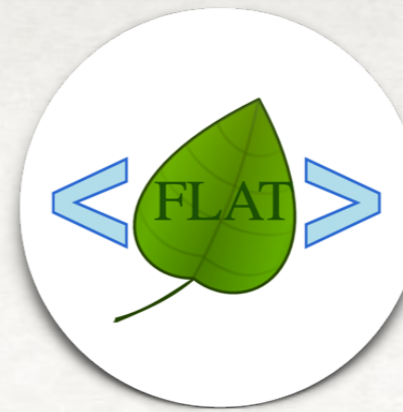Required. 30 characters or fewer. Letters, digits and @/./+/-/_ only.

**Password:** [                    ]

**Password confirmation:** [                    ]

Enter the same password as above, for verification.

Create the account

# FLAT

http://flat.science.ru.nl/login

## Log in

**Username:** kercos

**Password:** •••••••

**Configuration:** PARSEME Demo

Log in

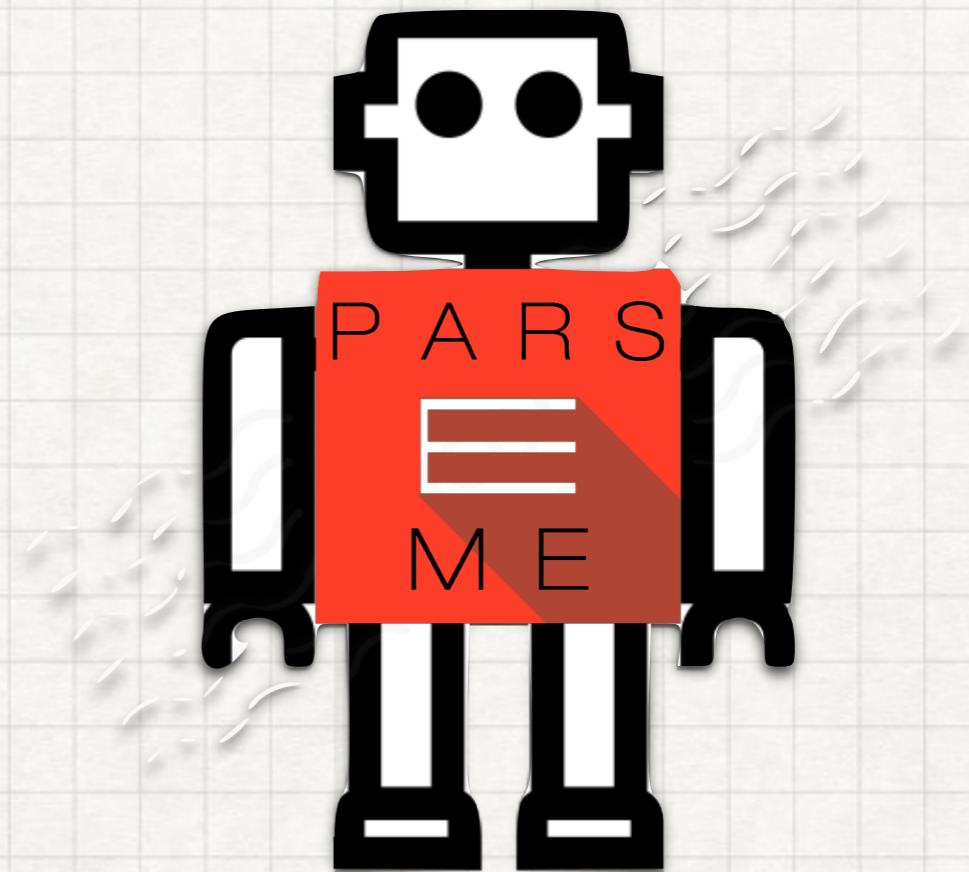(no account yet? Register here)

# FLAT

http://tiny.cc/flatCorpus

Pilot2_ST_English_FoLiA.zip

Extract it into a local folder

# TELEGRAM-BOT

- In house application (Federico with the help of Behrang)

- Based on Telegram platform (messaging application)

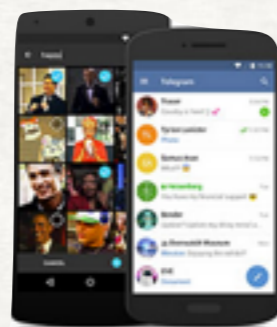- Text-based (simple interface with text and buttons)

# TELEGRAM-BOT

## web.telegram.org

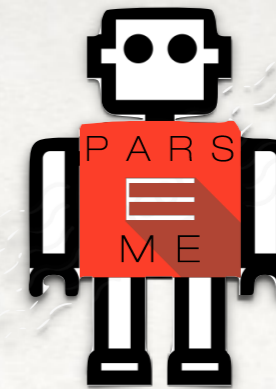or install **Telegram app** in your PC/tablet/phone



a native app for every platform

Telegram Web-version          Telegram for PC/Mac/Linux

# TELEGRAM-BOT

## web.telegram.org

# TELEGRAM-BOT

## web.telegram.org

# COMPARISON

| | FLAT | TELEGRAM-BOT |
|---|---|---|
| **CUSTOMIZATION** | FEATURES COMMON TO OTHER ANNOTATION TASKS (E.G., RIGHT TO LEFT LANGUAGES) | ANY FEATURE SPECIFIC TO PARSEME TASK (E.G., ROLES, TRAINING PHASE, VALIDATION) |
| **DEVICE COMPATIBILITY** | WEB INTERFACE: ACCESSIBLE VIA A BROWSER | VIA TELEGRAM APP: ACCESIBLE VIA A BROWSER OR DEDICATED APP (PC, TABLET or PHONE |
| **INTERFACE** | DOCUMENT BASED (MULTIPLE SENTENCES) | SENTENCE BASED |