

Tutorial accompanying the study "Computer-Assisted Language Comparison: State of the Art"

This tutorial supplements the study "Computer-Assisted Language Comparison: State of the Art". In this tutorial, we explain in detail, how our workflow can be tested and applied.

The workflow consists of several Python libraries that interact, one producing the data that can be used by the other. Since the data is available in different stages, each stage allows us to intervene by correcting errors manually that were made by the automated approach.

For users who are interested in testing our workflow on their local machine or further applying it in their own research, some basic knowledge of the Python programming language and the commandline will be required. All the software offered here is available in the form of free software. For more information on LingPy, the main programming library used here, we recommend users to check the tutorial¹ accompanying the study "Sequence comparison in computational historical linguistics"² by List et al. (2018)[1].

1. Code Ocean Capsule

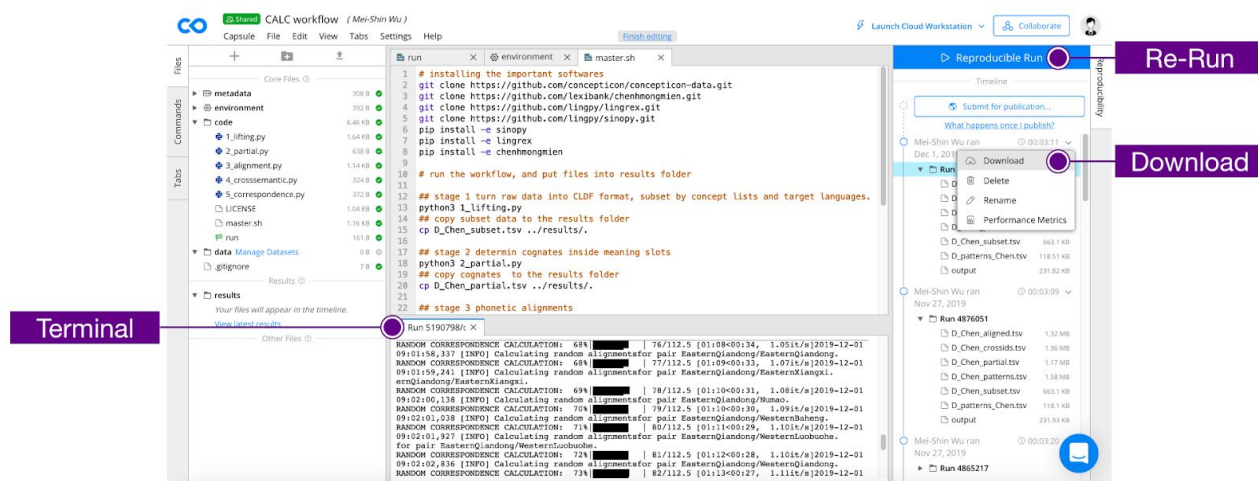
In order to facilitate it for users to quickly test our workflows without installing the software, we have set up a Code Ocean Capsule which users can use to run the code remotely. Code Ocean is an open access platform which enables researchers to reproduce their or others' experiments. For a detailed introduction to the Code Ocean platform³, please refer to the website. To see how our experiments can be run from within the Code Ocean Capsule, follow the following steps:

- a) Navigate to the capsule: <https://codeocean.com/capsule/8178287/tree/v2>
- b) Press the **Re-Run** button to reproduce the results.
- c) View the progression in the **Terminal** panel.
- d) Download all results and unzip the .zip file for further inspection on EDICTOR.

¹ <https://github.com/lingpy/lingpy-tutorial>

² <https://academic.oup.com/jole/article/3/2/130/5050100>

³ <https://codeocean.com/>



The following files can be found in the downloaded file:

File	Stage	Section
D_Chen_subset.tsv	From raw data to tokenized data	3.1
D_Chen_partial.tsv	From Tokenized Data to Cognate Sets	3.2
D_Chen_aligned.tsv	From Cognate Sets to Alignments	3.3
D_Chen_crosssids.tsv	From Alignments to Cross-Semantic Cognates	3.4
D_Chen_patterns.tsv	From Cross-Semantic Cognates to Sound Correspondence	3.5
D_Chen_distance.dst	Validation	4.2, 4.3
D_Chen_tree.tre	Validation	4.2, 4.3

2. Installation Instructions

We assume that users who are interested in running the workflow on their local machine are familiar with the essentials of command-line operations and system administration on either Unix-like systems (such as Linux and MacOS) or Windows systems. Also, users should have Python⁴ installed, including the package manager `pip`. Additionally, the version control system⁵

⁴ <https://www.python.org/>, Version 3.5 or higher

⁵ <https://git-scm.com/>

`git` will be required. We strongly encourage users to run this code in a virtual environment. A virtual environment is a practical solution for creating independent configurations for testing and experimenting, with no interference on the system-wide installation and without requiring complex virtualization or containerization solutions. The Python Packaging User Guide⁶ gives clear instructions on setting up a virtual environment on Windows, Linux and macOS.

We start by installing the dependencies from the commandline. In order to do so, we first download the code that we will use with help of `git`.

```
$ git clone https://github.com/lingpy/workflow-paper.git
$ cd workflow-paper
```

Now that we have done this, we can install all the packages we will need with help of `pip`.

```
$ pip install -r requirements.txt
```

Now that this has been done, we need to configure the access to reference catalogs, such as Concepticon⁷ and CLTS⁸ in order to make sure that they can be accessed readily by the code. This can be done with help of the `catconfig` argument submitted with the `cldfbench` package which organizes the linguistic datasets.

```
$ cldfbench catconfig
```

You will be prompted to ask if you want to clone actual versions of Concepticon, Glottolog, and CLTS, and the easiest way to deal with this is to agree and type “y” in all cases.

3. Getting Started

There are two basic ways in which you can run our workflow:

1. You can run it by downloading a set of Python scripts and running them directly on your computer.

⁶ <https://packaging.python.org/guides/installing-using-pip-and-virtual-environments/>

⁷ <https://github.com/concepticon/concepticon-data>

⁸ <https://github.com/cldf-clts/clts/>

2. You can use the `cldfbench` package to run the commands via the commandline, without downloading the data directly.

The advantage of solution 2 is that you do not have to download extra data, since we have integrated the code directly in the `lexibank` version of the dataset of Hmong Mien languages by Chén (2012)[2]. Once this dataset has been installed (and this is the first package we have installed in the previous section as part of all dependencies needed), you can type commands on your commandline, and the code will be carried out. The disadvantage is that the code example itself is not that easy to process for people less experienced with Python. For this reason, we will only note the commands in each of the steps we discuss in the following, and not explain them in more detail.

3.1 From Raw Data to Tokenized Data

The first script essentially loads the data from the repository and creates a wordlist that contains a subselection of all the data that was used. Some aspects of the more difficult “lifting” of data have already been done and distributed along with the original data package⁹, which specifically also contains the orthography profile in the file `etc/orthography.tsv` and can be automatically applied with help of the `cldfbench` package.

```
$ cldfbench lexibank.makecldf chenhmongmien
```

But since the data is available in the form of a `cldf` package with the original orthography already tokenized to the formats we need, you can also skip this step and convert the data to the wordlist format required by the `lingpy` package.

```
$ python 1_select.py
```

If you want to test the version from the CLDF-repository directly with `cldfbench`, you can type:

```
$ cldfbench chenhmongmien.wf_select.
```

This will select a part of the languages and a part of the concepts, as indicated in the main study and write them to a file `D_Chen_subsets.tsv`. Additionally, you will see some statistics on the

⁹ <https://github.com/lexibank/chenhmongmien>

terminal, specifically a table indicating the coverage for each language. If you want to select all languages, and not just a subset, type:

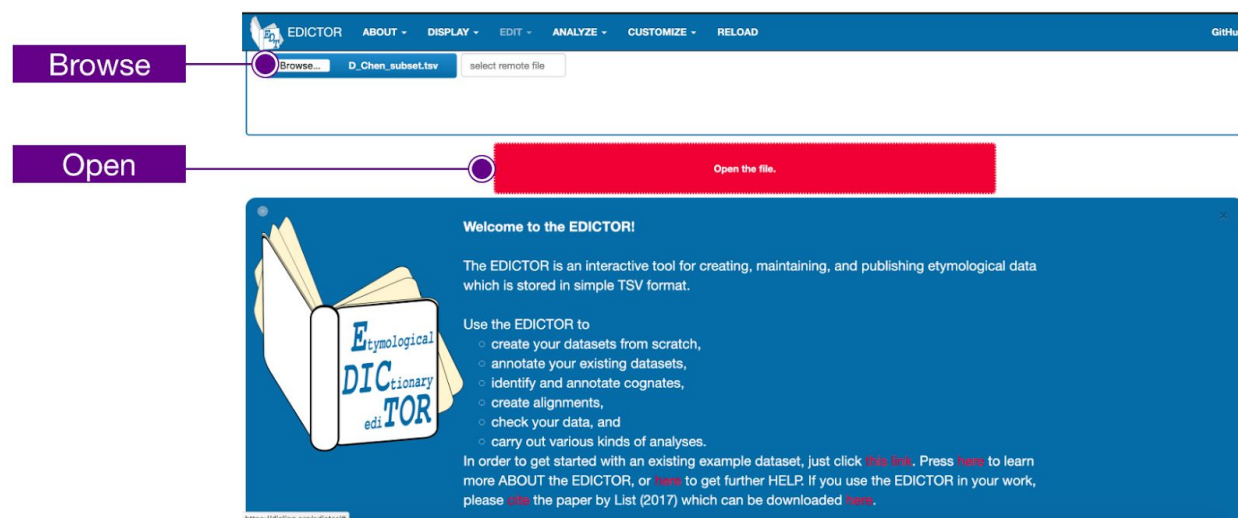
```
$ python 1_select.py all
```

The output `A_Chen_subset.tsv` is generated due to the argument `all` is used. Once the argument `all` is used in the first stage, it has to be added to the rest of stages to ensure that the workflows process the correct files.

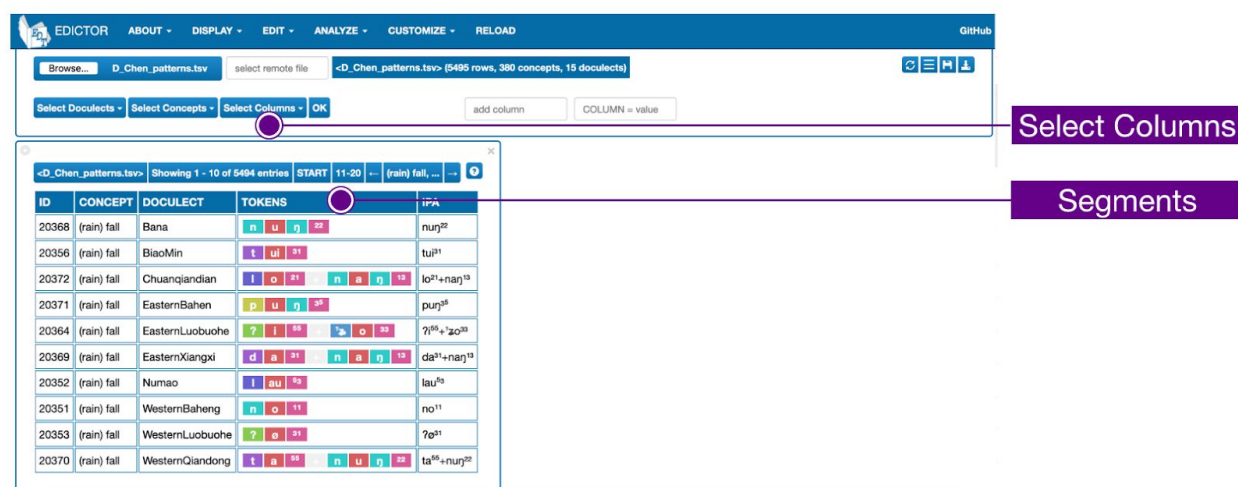
Doculect	Words	Coverage
Bana	502	1.00
BiaoMin	488	0.97
CentralGuizhouChuanqiandian	454	0.90
Chuanqiandian	501	1.00
EasternBahen	492	0.98
EasternLuobuohe	499	0.99
EasternQiandong	442	0.88
EasternXiangxi	492	0.98
Numao	490	0.98
WesternBaheng	500	1.00
WesternLuobuohe	488	0.97
WesternQiandong	494	0.98
WesternXiangxi	502	1.00
Younuo	500	1.00
ZaoMin	455	0.91

Already now you can inspect the data with the help of the [EDICTOR](https://digling.org/edictor/) tool. In order to do so, open the tool's website at <https://digling.org/edictor/> and wait until the page is loaded (note that we recommend to browse EDICTOR in Firefox, but GoogleChrome should also not cause further problems).

The data is in the file `D_Chen_subset.tsv`, in order to load it to the tool, press the **Browse** button and select the file. Once this has been done, press the **Open the file** button to examine the data, as illustrated in the following figure.



The segmented strings are displayed in the TOKENS column. Press **Select Columns** to inspect the raw forms and other aspects of the data, as shown in the following figure.



In order to save data to your computer, after you have manually edited them, you need to “download” them. This may be a bit surprising, since effectively, you do not download the data, but since the EDICTOR is working on a browser, it does not have any access to the data on your computer, and **download** is the only way to communicate with your machine. Thus, in order to save your data and load it to your machine, you first have to press the **save** icon at the top-right corner in order to store the edited data in the web browser. When now pressing the

download icon at the top-right, your browser will either directly download the data and store them in your download folder, or it will ask you to specify a specific file destination.

The screenshot shows the EDICTOR web interface. At the top, there is a navigation bar with links: EDITOR, ABOUT, DISPLAY, EDIT, ANALYZE, CUSTOMIZE, and RELOAD. Below this, a file browser shows the selected file 'D_Chen_patterns.tsv' with a 'Download' button. The main area displays a table with columns: ID, CONCEPT, DOCULECT, TOKENS, and IPA. The table contains 10 rows of data, including entries for 'Bana', 'BiaoMin', 'Chuanqiandian', 'EasternBahen', 'EasternLuobuohe', 'EasternXiangxi', 'Numao', 'WesternBaheng', 'WesternLuobuohe', and 'WesternQiandong'. To the right of the table, there is a 'Save' button.

ID	CONCEPT	DOCULECT	TOKENS	IPA
20368	(rain) fall	Bana	n u ŋ ʔ	nun ²²
20356	(rain) fall	BiaoMin	t u i ʔ	tu ²¹
20372	(rain) fall	Chuanqiandian	l o ʔ i a ŋ ʔ	lo ²¹ +nan ¹³
20371	(rain) fall	EasternBahen	p u ŋ ʔ	puŋ ³⁸
20364	(rain) fall	EasternLuobuohe	ʔ i ʔ ʔ ʔ ʔ ʔ	ʔi ¹⁰ +ʔo ²⁰
20369	(rain) fall	EasternXiangxi	d a ʔ i a ŋ ʔ	da ²¹ +nan ¹³
20352	(rain) fall	Numao	l au ʔ	lau ²⁵
20351	(rain) fall	WesternBaheng	n o ʔ	no ¹¹
20353	(rain) fall	WesternLuobuohe	ʔ a ʔ	ʔa ²¹
20370	(rain) fall	WesternQiandong	t a ʔ i a ŋ ʔ	ta ¹⁰ +nun ²²

Be careful when editing data in the EDICTOR without saving and downloading them. If you close your browser, all the edits you made will be lost, so you should regularly save and download your data when working with the EDICTOR. As a shortcut, you can also type CONTROL+S to save and CONTROL+E to “export” the data (i.e., to download them).

3.2 From Tokenized Data to Cognate Sets

Partial cognate detection is an important task, specifically when working with Southeast Asian language data. The algorithm we use for this task was first proposed in the study “Using Sequence Similarity Networks to Identify Partial Cognates in Multilingual Wordlists” by List et al. (2016)[3], where the algorithm is described in due detail.

To illustrate how the algorithm works, we provide an example with four words for ‘moon’ in the Eastern Baheng, Eastern Qiandong, Bana and Biao Min language varieties.

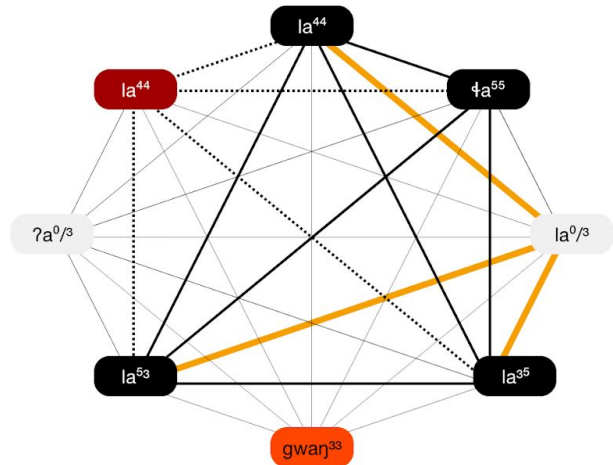
The major steps of the algorithm are the following:

1. Calculate the distances of all morpheme pairs.
2. Create a fully connected network from the distance scores.
3. Filter the network by deleting edges in the following fashion:
 - A. Two morphemes in the same word should not be linked (see the dashed lines in the following figure).

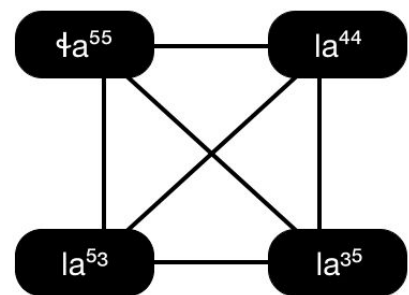
B. A morpheme in a word should not be linked to two morphemes in another word (see the yellow edges in the figure).

- Remove the edges with similarity scores below a given threshold.

DOCULECT	IPA	
EasternBaheng	la ^{0/3}	ɬa ⁵⁵
EasternQiandong	la ⁴⁴	la ⁴⁴
Bana	la ^{0/4}	la ³⁵
BiaoMin	la ⁵³	gwan ³³



Once this has been done, an algorithm for Community Detection in networks[4] is used to partition the network into “communities”, with each community representing one partial cognate set.



In order to calculate partial cognates, we use the algorithm as provided by the `lingpy` software package and apply it to our subselection of languages.

```
$ python 2_partial.py
```

If you want to test the version from the CLDF-repository directly with `cldfbench`, you can type:

```
$ cldfbench chenhmongmien.wf_partial.
```

This will take some time when you run it the first time. The data can be found in the file `D_Chen_partial.tsv`.

To inspect the data with EDICTOR, load `D_Chen_partial.tsv` as shown before. Then press **DISPLAY** to select **SETTINGS** in the drop-down menu. Select **PARTIAL** in the **Morphology and Colexification Mode** entry. Press the **Refresh** button.

The screenshot shows the EDICTOR web application interface. On the left, a table displays data from the file `D_Chen_patterns.tsv`. The table has columns: ID, CONCEPT, DOCULECT, TOKENS, IPA, and COGIDS. The data shows various linguistic entries with their corresponding tokens and COGIDS values.

On the right, the 'Settings of the Editor' panel is open. It contains several settings: 'Preview' set to 10, 'Cognate IDs' set to 'PARAMETERID', 'Partial Cognates' set to 'COGIDS', 'Tokens' set to 'tokens column', 'Alignments' set to 'ALIGNMENT', and 'Morphology and Colexification Mode' set to 'PARTIAL'. A 'Refresh' button is located next to the 'Preview' setting.

Three purple callout boxes with white text are overlaid on the image:

- 'Display' points to the 'DISPLAY' button in the top navigation bar.
- 'Refresh' points to the 'Refresh' button in the 'Settings of the Editor' panel.
- 'Partial' points to the 'PARTIAL' radio button in the 'Morphology and Colexification Mode' section.

In order to investigate the partial cognates, you need to select the column which stores the identifiers. To do so, press **Select Columns** and select **COGIDS** in the drop-down menu.

If you right-click on any number in the “COGIDS” column, a pop-up window will open and show all the cognate sets for a given word form in the form of an alignment. Since we have not yet aligned the data, the alignment will be wrong at this point.

The screenshot shows the EDICTOR application interface. The top bar includes menu items: EDITOR, ABOUT, DISPLAY, EDIT, ANALYZE, CUSTOMIZE, and RELOAD. Below the bar, there's a file browser showing 'D_Chen_partial.tsv' with 5502 rows, 360 concepts, and 15 doculects. A toolbar below the browser has buttons for 'Select Doculects', 'Select Concepts', 'Select Columns', and 'OK', along with 'add column' and 'COLUMN = value' options.

The main window displays a table with columns: ID, CONCEPT, DOCULECT, TOKENS, IPA, and COGIDS. The table shows 15 entries. A right-click context menu is visible over the table, with a callout 'Right click!'. A detailed view of the entry for 'Bana' (COGIDS: 6542, 6541) is shown on the right, with a callout 'Alignment' pointing to the alignment grid.

ID	CONCEPT	DOCULECT	TOKENS	IPA	COGIDS
31	sun	Bana	l a n	la ³¹ +ni ³¹	6542 6541
43	sun	BiaoMin	p i t au	pi ³¹ +tau ³¹	6541 6548
49	sun	CentralGuizhouChuanqiandian	q e	qe ³⁴	6541
40	sun	Chuanqiandian	q o	qo ³⁴	6541
28	sun	EasternBahen	l a o e	la ³¹ +ue ³⁵	6542 6541
44	sun	EasternLuobuohe	q o p n a </td <td>qo³¹+na³¹</td> <td>6547 6541</td>	qo ³¹ +na ³¹	6547 6541
47	sun	EasternQiangdong	q a n ei	q'an ³⁵ +nei ³⁴	6546 6541
33	sun	EasternXiangxi	q a ei	q'a ³⁵ +ei ³⁴	6543 6541
48	sun	Numao	m a q a j	ma ³² +qa ³¹	6545 6541
29	sun	WesternBaheng	q o n a j ei	q'on ³¹ +?an ³⁵ +ei ³⁴	6546 6550 6541

3.3 From Cognate Sets to Alignments

To align the data, we use the new procedure for template-based alignment, which is available from the `lingrex` package which we have installed as one of the requirements of our workflow, and the `sinopy` package, which helps us to compute syllable templates from all morphemes in the data. Running the code is again straightforward.

```
$ python 3_alignment.py
```

If you want to test the version from the CLDF-repository directly with `cldfbench`, you can type:

```
$ cldfbench chenhmongmien.wf_alignment
```

The aligned data will be stored in the file `D_Chen_aligned.tsv`. To inspect the alignments in EDICTOR, load this file and follow the previous steps we mentioned in Section 3.2. In addition to selecting the **COGIDS** column now, we also select the **STRUCTURE** column, since this column provides the templates for each morpheme, which we have automatically added to the data with help of `sinopy`.

EDICTOR ABOUT DISPLAY EDIT ANALYZE CUSTOMIZE RELOAD GitHub

Browse: D_Chen_aligned.tsv select remote file <D_Chen_aligned.tsv> (10980 rows, 360 concepts, 15 doculects)

Select Doculects Select Concepts Select Columns OK add column COLUMN = value

ID	CONCEPT	DOCULECT	TOKENS	IPA	COGIDS	STRUCTURE
31	sun	Bana	l a n	la ^h +n ¹³	6542 6541	int+int
45	sun	Bana	l a n	la ^h +n ¹³	2279 2272	int+int
45	sun	Bana	l a n	la ^h +n ¹³	2279 2272	int+int
32	sun	BiaoMin	g l i t au n	ŋ ¹³ +tau ³¹	2272 2276	int+int
43	sun	BiaoMin	g l i t au n	ŋ ¹³ +tau ³¹	6541 6548	int+int
39	sun	CentralGuizhouChuanqiandian	g e	ŋe ⁴³	2272	int
50	sun	Chuanqiandian	g e	ŋe ⁴³	2272	int
49	sun	EasternBahen	l a e	la ^h /a+pe ³¹	2279 2272	int+int
49	sun	EasternBahen	l a e	la ^h /a+pe ³¹	2279 2272	int+int
41	sun	EasternLuobuohe	g o e n a n	qo ^h /e+na ³¹	2281 2272	int+int

As we already mentioned, if you right-click on any number in the “COGIDS” column, a pop-up window will show the alignment. Click on the `=` sign to modify the alignment. The modification itself is very straightforward: just click on a sound segment to move it to the right, and click on a gap segment to delete this segment.

Bana «sun» (COGIDS: 2279,2272)

DOCULECTS	CONCEPTS	ID: 2279	2272
WesternBaheng	sun	g e i	35
Numao	sun	g a n	33
WesternLuobuohe	sun	g a n	31
BiaoMin	sun	g i n	44
ZaoMin	sun	g ai n	44
Younuo	sun	g n o	34
CentralGuizhouChuanqiandian	sun	g e n	35
WesternXiangxi	sun	g e n	31
EasternLuobuohe	sun	g a n	13
Bana	sun	l a e	53
EasternXiangxi	sun	g e i	34
EasternQiandong	sun	g e i	44
WesternQiandong	sun	g e i	35
EasternBahen	sun	l a e	43
Chuanqiandian	sun	g e i	31
Bana	firewood	ts ^h e n	313
ZaoMin	firewood	h o n	53
Younuo	firewood	ts ^h i n	13
BiaoMin	firewood	ts a n	31

Cognate set "2279" links the following 6 entries:

	l	a	e	n	
Bana	l	a	e	n	313
EasternBahen	l	a	e	n	31
Bana	ts ^h	e	n		53
ZaoMin	h	o	n		13
Younuo	ts ^h	i	n		31
BiaoMin	ts	a	n		31

EDIT ALIGN CLOSE

3.4 From Alignments to Cross-Semantic Cognates

The algorithm for cross-semantic cognate detection as we propose it here is illustrated in more detail in the main study. It is implemented as part of the `lingrex` package. Again, it is straightforward to run the code.

```
$ python 4_crosssemantic.py
```

If you want to test the version from the CLDF-repository directly with `cldfbench`, you can type:

```
$ cldfbench chenhmongmien.wf_crosssemantic
```

The output file is `D_Chen_crossids.tsv`, and we load it into the EDICTOR tool, just as we did before, but when checking the **SETTINGS** in the menu this time, we need to specify that the column “CROSSIDS” holds the partial cognates. To do so, just type in **CROSSIDS** in the text field **Partial Cognates** in the settings menu and then press the **refresh** button.

Settings of the Editor:

Use this menu to select your settings for an ongoing EDICTOR session. In order to change the general sessions before loading the EDICTOR, you need to specify a customer URL in the **Custom Settings**. Refer to the **Help** panel to understand more about the general ideas behind the EDICTOR. Alternatively, just hover with the mouse pointer over the respective settings and read the instructions which will appear.

Preview: 15

Cognate IDs: PARAMETERID

Partial Cognates: CROSSIDS

Tokens: tokens column

Alignments: ALIGNMENT

Morphology and Colexification Mode: ☐ FULL ☒ PARTIAL

Morphemes: morpheme annotation column

Sources: source

To inspect the distribution of partial cognates, press **ANALYZE** in the top-level menu and select **Cognate sets** in the drop-down menu.

Analyze

Cognate Sets

As a result, a new panel will open and show the distribution of all cognate sets across the different language varieties. Pressing the red button with the cognate set identifier on the left will open the alignment. Pressing the yellow buttons with the word identifiers will show you the original morpheme. On the right, in the column **CONCEPTS**, you will find those cognate sets which are attested for more than one concept as separated by a comma. Clicking on this field will modify the main wordlist panel in such a way that only the selected concepts will appear.

Investigate cognate sets in the data

Select Sets: 1 of 14 Sets

CROSSIDS	Wes	Eas	Wes	Ban	Eas	Wes	Zas	Chu	You	Bis	Eas	Wes	Num	Eas	Can	CONCEPTS
112	409	409	409	409	409	409	409	409	409	409	409	409	409	409	409	leaf, back, branch, help, road
241	2934	2935	2940	2941	2944	2947	2949	2950	2951	2952	2955	2954	2956			blood
281	4925	4926	4927	4929	4931	4932	4938	4940	4942	4943	4946	4945	4947			branch, twig
388	18655	19291	19292	19294	18662	19300	18668	19306	18672	18673	18676	18675	18677			climb (a tree), go upstairs
481	18124	18126	18128	18130	18131	18134	18137	18139	18140	18141	18142	18144	18146			cry
575	17303	17304	17305	17307	17310	17313	17316	17318	17320	17321	17324	17323	17325			drink
644	10529	10530	10533	10535	10539	10542	10545	10546	10547	10550	10549	10551				eight, egg
681	11413	11415	11417	11418	11421	11424	11426	11427	11428	11429	11432	11431	11433			everyday
893	2089	2090	2091	2095	2099	2102	2104	2105	2106	2107	2110	2108	2111			eye
724	6551	6552	6553	6555	6557	6558	6561	6564	6566	6568	6569	6571	6572			fat meat
830	16750	16751	16752	16754	16756	16757	16759	16767	16768	16771	16770	16772				fly
901	5048	5049	5050	5052	5054	5055	5058	5063	5064	5065	5066	5069	5070			fruit
962	14439	14440	14441	14445	14446	14449	14452	14435	14456	14457	14460	14459	14461			good
2447	7642	7643	7648	7648	7649	7652	7657	7658	7659	7660	7663	7710	7664			house, thatched cottage, ill...
1172	956	957	958	960	962	963	968	971	973	974	977	976	978			human being
1191	20894	20895	20896	20898	20900	20901	20904	20907	20909	20910	20911	20914	20916			illuminate
1282	20992	20993	20994	20996	20998	21002	21005	21007	21008	21009	21013	21012	21014			know
1294	20440	20447	20448	20452	20453	20456	20459	20461	20463	20464	20467	20466	20468			laugh (smile)

CROSSIDS

Comma!

Click to Expand

3.5 From Cross-Semantic Cognates to Sound Correspondence Patterns

As a final step, we will try to infer the major correspondence patterns in the data, using the algorithm by List (2019)[5] which is available from the `lingrex` package. Running the code is straightforward, as before.

```
$ python 5_correspondence.py
```

If you want to test the version from the CLDF-repository directly with `cldfbench`, you can type:

```
$ cldfbench chenhmongmien.wf_correspondence
```

This creates two output files. One, called `D_Chen_patterns.tsv` is the file without wordlist that can be loaded by EDICTOR and inspected, and one file contains the patterns that have been inferred alone, called `D_patterns_Chen.tsv`. In order to inspect the patterns, we recommend to use the EDICTOR tool, which requires the same steps that we already applied when loading our cross-semantic cognates. Once this has been done, press the **ANALYZE** button in the top menu and select **CORRESPONDENCE PATTERNS** in the drop-down menu.

The screenshot shows the PHOENIX software interface. The 'ANALYZE' menu is open, displaying options: PHONOLOGY, MORPHOLOGY, CORRESPONDENCES, COGNATE SETS, and CORRESPONDENCE PATTERNS. The 'CORRESPONDENCE PATTERNS' option is selected, and a sub-menu is visible with 'START' and 'END' buttons. Below the menu, a table displays data for the file 'D_Chen_patterns.tav'. The table has columns: ID, CONCEPT, DOCULECT, TOKENS, and IPA. The data rows show various linguistic entries with their corresponding tokens and IPA transcriptions.

ID	CONCEPT	DOCULECT	TOKENS	IPA
20368	(train) fall	Bana	n u n	nur ²²
20366	(train) fall	BiaoMin	t u l	tu ²¹
20372	(train) fall	ChuangJiandian	t o l	lo ²¹ +na ²¹
20371	(train) fall	EasternBahen	p u n	pu ²²
20364	(train) fall	EasternLuobuohu	t i l	ʔi ²¹ +ʔo ²¹
20369	(train) fall	EasternXiangel	d i l	da ²¹ +na ²¹
20352	(train) fall	Numao	t a u	tau ²²
20351	(train) fall	WesternBaheng	n o l	no ²¹
20353	(train) fall	WesternLuobuohu	t i l	ʔi ²¹
20370	(train) fall	WesternQiangdong	t a n	ta ²¹ +nur ²²

In order to allow for a good display, the doculect names are all abbreviated. Hovering the mouse cursor on an abbreviation will show you the full name.

The screenshot shows the WordNet Explorer application interface. At the top, there's a navigation bar with tabs: EDITOR, ABOUT, DISPLAY, EDIT, ANALYZE, CUSTOMIZE, and RELOAD. Below this is a search bar with the text "O_Chen_patterns.tsv" and a "select remote file" button. To the right of the search bar, it says "<O_Chen_patterns.tsv> (5466 rows, 380 concepts, 16 doculects)". Below the search bar are buttons for "Select Doculects", "Select Concepts", "Select Columns", and "OK". There's also a "add column" button and a "COLUMN = value" input field.

The main content area is titled "Investigate correspondence patterns in the data". Below the title are buttons for "Select Sets", "THR", "38", "PREV", "30", "OK", and "1-4 of 6 Sites". Below these buttons is a table with columns: COGNATES, INDEX, PATTERN, CONCEPTS, Num, Wes, Bis, Zao, You, Wes, Ess, Ban, Ess, Wes, Ess, Chu, Wes, Cen, Ess. The table contains several rows of data, each representing a crossword puzzle. The first row is for crossword 1939, pattern 1, with the pattern "da / 30" and the concept "salt". The second row is for crossword 2043, pattern 1, with the pattern "da / 50" and the concept "sharp". The third row is for crossword 246, pattern 1, with the pattern "da / 823" and the concept "blood". The fourth row is for crossword 738, pattern 1, with the pattern "da / 823" and the concept "fear (be afraid)". The fifth row is for crossword 1567, pattern 1, with the pattern "da / 412" and the concept "mushroom". The sixth row is for crossword 365, pattern 1, with the pattern "da / 541" and the concept "climb (a tree), go upstairs".

Annotations in the image include a purple circle and line pointing to the "CROSSIDS" label, another purple circle and line pointing to the "Doculect" label, and a purple circle and line pointing to the "Correspondence Patterns" label.

Clicking on a cell in the correspondence pattern panel will allow you to see not only the sound in question, but the full morpheme in which this sound occurs.

Investigate correspondence patterns in the data

Select Sets ▾ THR 38 PREV 30 OK -- 1-6 of 6 Sites -- ↻

COGNATES	INDEX	PATTERN	CONCEPTS	Num	Wes	Bia	Zao	You	Wes	Eas	Ban	Eas	Wes	Eas	Chu	Wes	Cen	Eas	SIZE
1939	1	ds / 50	salt	ts ei 13	ts i 33	dz n 43	Ø	Ø	ts ts ds	z	a	ts ts	Ø	Ø	Ø	Ø	Ø	Ø	1.07 / 2
2043	1	ds / 50	sharp	ts e 44	Ø	Ø	Ø	Ø	Ø	ts	Ø	Ø	a	ts	Ø	ts	n	Ø	1.07 / 2
246	1	ds / ...	blood	ts a n 13	ts a n 44	a n 34	z a m 24	ts u n 33	ts ts	Ø	ts	a	ts ts	Ø	ts	a	Ø	1.60 / 2	
738	1	ds / ...	fear (be afraid)	ts	ts e 34	Ø	z a 42	ts i 44	ts ts ds	Ø	a	Ø	ts	Ø	ts	a	Ø	1.60 / 2	
1567	1	ds / ...	mushroom	ts ei 33	ts i 31	ts au 33	k i u 44	Ø	k	ts ds	g	ts	k	ts	Ø	ts	ts	0.00 / 1	
395	1	ds / ...	climb (a tree),go upstairs	ts ei 44	ts i n 35	ts a 34	h o 42	Ø	ts ts ds	Ø	ts	ts ts	ts	ts	ts	ts	ts	0.00 / 1	

4. Validation

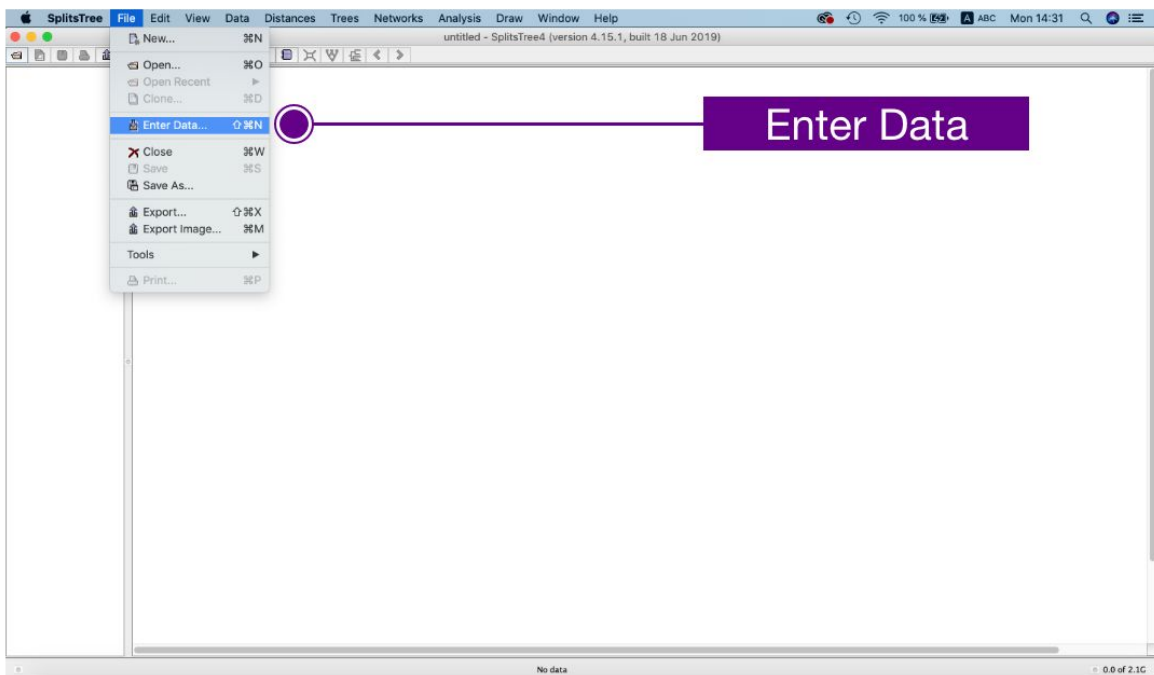
We calculate the shared cognates between language pairs and output the scores in the form of a pairwise distance matrix. The script `6_phylogeny.py` gives two documents, a distance matrix (`A_Chen_distance.dst` or `D_Chen_distance.dst`) and a tree file, based on a Neighbor-Joining analysis (`A_Chen_tree.tre` or `D_Chen_tree.tre`).

There are many ways to work with the distance matrix, here, we give one of the approaches to visualize the matrix as a neighbor-net network with the help of SplitsTree.

To get started, first make sure to install SplitsTree¹⁰ [6] and follow the installation instructions. In order to compute the distance matrix with our code, use the command line (here we compute it for the entire dataset, so we run it with the keyword `all`)

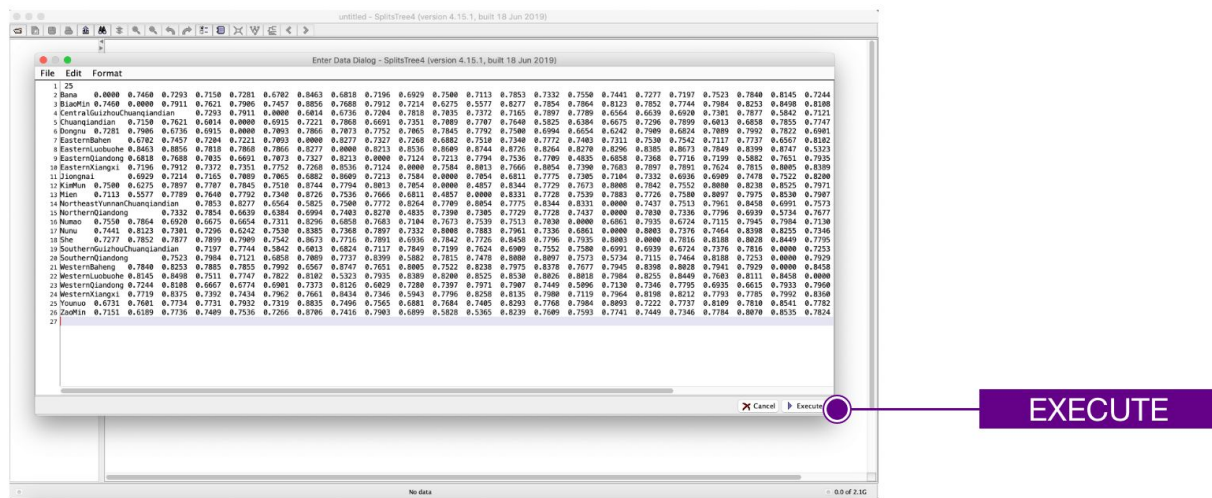
```
$ python 6_phylogeny.py all
```

To generate a Neighbor-Net from the distance matrix, open the file `A_Chen_distance.dst` or `D_Chen_distance.dst` with any plain text editor and start the SplitsTree software. Then click on **File** and **Enter Data**, as shown in the image below.

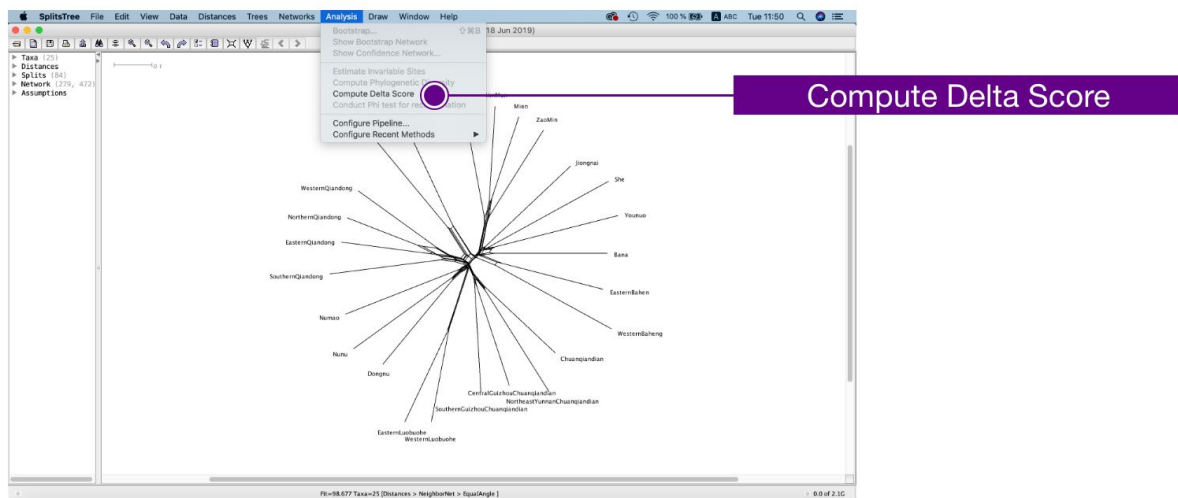


¹⁰ <https://software-ab.informatik.uni-tuebingen.de/download/splitstree4/welcome.html>

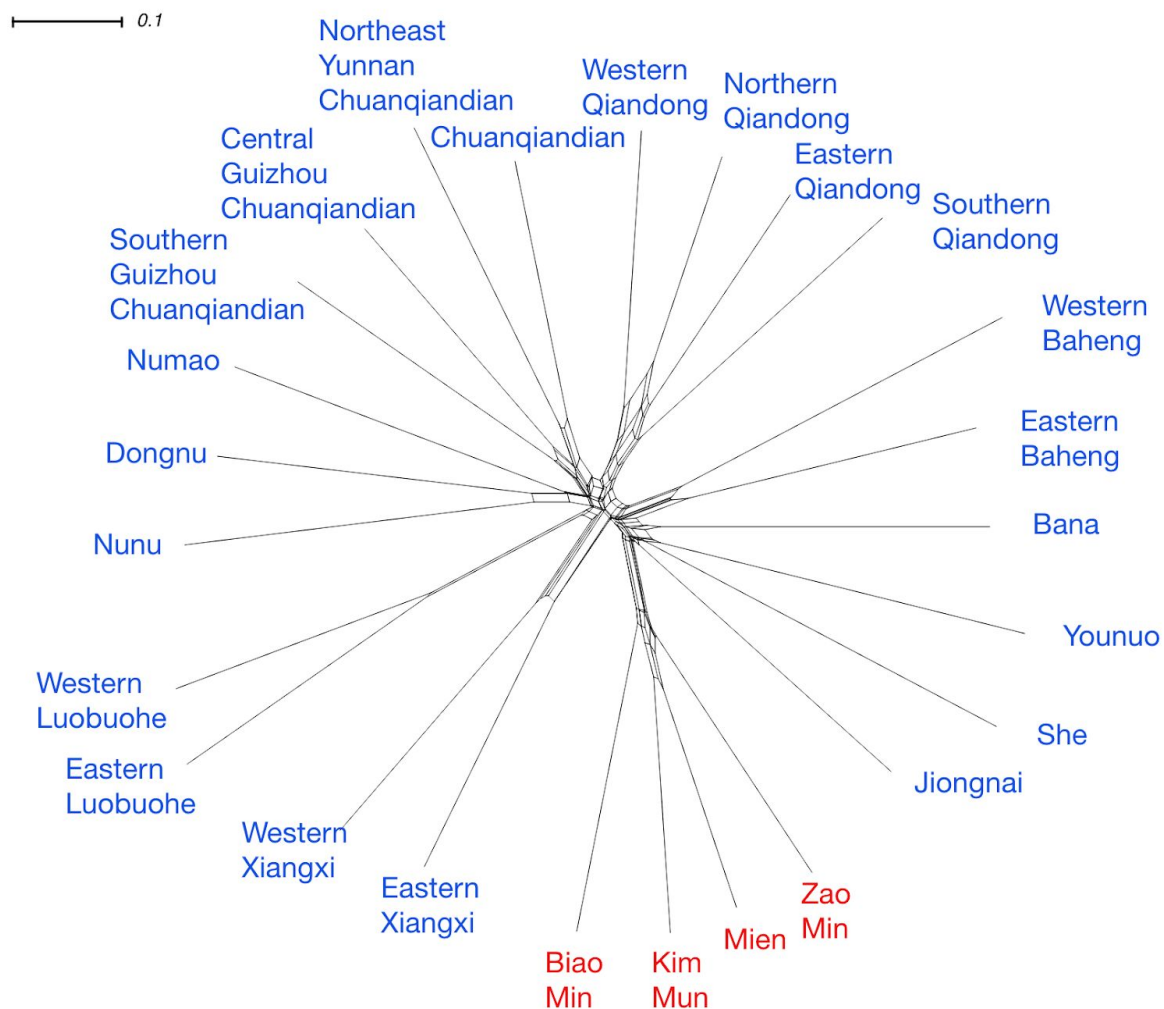
Then copy the distance matrix and paste it into the **Enter Data Dialog**, and press **Execute**.



You can now inspect the network. To analyze the data further, you can compute the delta scores, showing the degree of reticulation in the data, by pressing **Analysis** and then **Compute Delta Score**, as shown below.



The resulting Neighbor-Net is shown in the following figure. For the purpose of illustration, the Mienic language varieties are colored in red, the Hmongic group is highlighted in blue.



The following table shows the delta scores we computed from the data.

Taxon	Delta score
Bana	0.34706
Biao Min	0.27289
Central Guizhou Chuanqiandian	0.29924
Chuanqiandian	0.29172
Dongnu	0.32416

Eastern Baheng	0.32056
Eastern Luobuohe	0.33529
Eastern Qiangong	0.32083
Eastern Xiangxi	0.33736
Jiongnai	0.32644
Kim Mun	0.26992
Mien	0.25672
Northeast Yunnan Chuanqiandian	0.29748
Northern Qiandong	0.28447
Numao	0.34185
Nunu	0.32375
She	0.31671
Southern Guizhou Chuanqiandian	0.34376
Southern Qiandong	0.30988
Western Baheng	0.35259
Western Luobuohe	0.3211
Western Qiandong	0.31137
Western Xiangxi	0.35174
Younuo	0.2996
Zao Min	0.26797

The average delta score is 0.313. As mentioned before, the distances between taxa are calculated via shared cognates. The shorter the distances between two taxa, the higher the similarities between them. If the taxa share cognates not only within their group but also outside their groups, the network finds it challenging to determine the best cluster for them. The larger the reticular structure, or the less tree-like the data is, the higher is the delta score. For one

particular language variety's delta score this means that this specific language contributes to a certain amount of conflict in the data.

5. Conclusion

In this tutorial, we provided details of how to execute our workflow for Computer-Assisted Language comparison, using the scripts we wrote, while at the same time illustrating how the results can be manually inspected and modified. We have not discussed the details of the code we wrote, but we recommend users proficient in Python to have a look.

6. References

1. List J-M, Walworth M, Greenhill SJ, Tresoldi T, Forkel R. Sequence comparison in computational historical linguistics. *Journal of Language Evolution* [Internet]. 2018;3(2):130–44. Available from: <https://academic.oup.com/jole/article/3/2/130/5050100?guestAccessKey=cf8fe64e-3996-4cb1-ba2c-317a7cd81bf4>
2. 陳其光CQ. Miàoyáo yǔwén [Internet]. Běijīng: Zhōngyāng Mínzú Dàxué 中央民族大学 [Central Institute of Minorities]; 2012. Available from: https://en.wiktionary.org/wiki/Appendix:Hmong-Mien_comparative_vocabulary_list
3. List J-M, Lopez P, Baptiste E. Using sequence similarity networks to identify partial cognates in multilingual wordlists. In: *Proceedings of the Association of Computational Linguistics 2016 (Volume 2: Short Papers)* [Internet]. Berlin: Association of Computational Linguistics; 2016. pp. 599–605. Available from: <http://anthology.aclweb.org/P16-2097>
4. Rosvall M, Bergstrom CT. Maps of random walks on complex networks reveal community structure. *Proc Natl Acad Sci USA*. 2008;105(4):1118–23.
5. List J-M. Automatic inference of sound correspondence patterns across multiple languages. *Computational Linguistics* [Internet]. 2019;1(45):137–61. Available from: https://www.mitpressjournals.org/doi/full/10.1162/coli_a_00344
6. Huson DH. SplitsTree: Analyzing and visualizing evolutionary data. *Bioinformatics*. 1998;14(1):68–73.