# Multiple Sequence Alignment in Historical Linguistics

A Sound Class Based Approach

Johann-Mattis List[*]

[*]Institute for Romance Languages and Literature
Heinrich Heine University Düsseldorf

2011/06/25

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

## Structure of the Talk

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

# Sequences

**Definition 1**

## Sequences

### Definition 1

Given an *alphabet* (a non-empty finite set, whose elements are called *characters*), a *sequence* is an ordered list of characters drawn from the alphabet. The elements of sequences are called *segments*. (cf. Böckenbauer & Bongartz 2003: 30f)

# Sequences

# Sequences

# Sequences

# Sequences

# Sequences

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Pairwise Sequence Alignment
Multiple Sequence Alignment

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Pairwise Sequence Alignment
Multiple Sequence Alignment

# Alignment Analyses

## Definition 2

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Pairwise Sequence Alignment
Multiple Sequence Alignment

## Alignment Analyses

**Definition 2**

An *alignment* of two sequences *s* and *t* is a two-row matrix in which both sequences are aranged in such a way that all matching and mismatching segments occur in the same column, while empty cells, resulting from empty matches, are filled with gap symbols. (cf. Kruskal 1983)

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Pairwise Sequence Alignment
Multiple Sequence Alignment

# Alignment Analyses

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Pairwise Sequence Alignment
Multiple Sequence Alignment

# Alignment Analyses

**Sequences**
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Pairwise Sequence Alignment
Multiple Sequence Alignment

# Alignment Analyses

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

## Pairwise Sequence Alignment

→ Construct a matrix in which all segments of two sequences are confronted with each other and with gap characters.

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

## Pairwise Sequence Alignment

→ Construct a matrix in which all segments of two sequences are confronted with each other and with gap characters.

→ Calculate the score of all subsequences recursively by filling the matrix from left to right and from top to bottom.

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Pairwise Sequence Alignment**
**Multiple Sequence Alignment**

## Pairwise Sequence Alignment

→ Construct a matrix in which all segments of two sequences are confronted with each other and with gap characters.

→ Calculate the score of all subsequences recursively by filling the matrix from left to right and from top to bottom.

→ Employ a scoring function in each recursion step, which evaluates, whether the characters in each cell should be matched with themselves or with gap characters.

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

## Pairwise Sequence Alignment

→ Construct a matrix in which all segments of two sequences are confronted with each other and with gap characters.

→ Calculate the score of all subsequences recursively by filling the matrix from left to right and from top to bottom.

→ Employ a scoring function in each recursion step, which evaluates, whether the characters in each cell should be matched with themselves or with gap characters.

→ Retrieve the alignment by applying a traceback function which reconstructs the 'path of choices' (Durbin 2002) which led to the final value.

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

| T | E | S | T | - | - | - | - |
|---|---|---|---|---|---|---|---|
| - | - | - | - | T | E | S | T |

**8**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

| T | E | S | T | - | - | - |
|---|---|---|---|---|---|---|
| - | - | - | T | E | S | T |

**6**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

| T | E | S | - | T | - | - | - |
|---|---|---|---|---|---|---|---|
| - | - | - | T | - | E | S | T |

**8**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

| T | E | S | T | - | - | - |
|---|---|---|---|---|---|---|
| - | - | T | - | E | S | T |

**7**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

| T | E | - | S | T | - | - | - |
|---|---|---|---|---|---|---|---|
| - | - | T | - | - | E | S | T |

**8**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

## Pairwise Sequence Alignment

| T | E | S | T | - | - | - |
|---|---|---|---|---|---|---|
| - | T | - | - | E | S | T |

**7**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

| T | - | E | S | T | - | - | - |
|---|---|---|---|---|---|---|---|
| - | T | - | - | - | E | S | T |

**8**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

| T | E | S | T | - | - | - |
|---|---|---|---|---|---|---|
| T | - | - | - | E | S | T |

**6**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

| T | E | S | T | - | - |
|---|---|---|---|---|---|
| T | - | - | E | S | T |

**5**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

|   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| T | E | S | - | T | - | - |
| T | - | - | E | - | S | T |

**6**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

```
T   E   S   T   -   -
T   -   E   -   S   T
```

**5**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

| T | E | - | S | T | - | - |
|---|---|---|---|---|---|---|
| T | - | E | - | - | S | T |

**6**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

```
T    E    S    T    -    -
T    E    -    -    S    T
```

**4**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

T     E     S     T     -
T     E     -     S     T

**3**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

```
T    E    -    S    T    -
T    E    S    -    -    T
```

**4**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

| T | E | S | T | - |
|---|---|---|---|---|
| T | E | S | - | T |

**2**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

**Pairwise Sequence Alignment**
Multiple Sequence Alignment

# Pairwise Sequence Alignment

```
T    E    S    T
T    E    S    T
```

**0**

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Pairwise Sequence Alignment
**Multiple Sequence Alignment**

# Multiple Sequence Alignment

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Pairwise Sequence Alignment
**Multiple Sequence Alignment**

## Multiple Sequence Alignment

→    Align all sequences pairwise and store the scores in a matrix.

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Pairwise Sequence Alignment
**Multiple Sequence Alignment**

## Multiple Sequence Alignment

→ Align all sequences pairwise and store the scores in a matrix.

→ Construct a guide tree from the matrix.

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
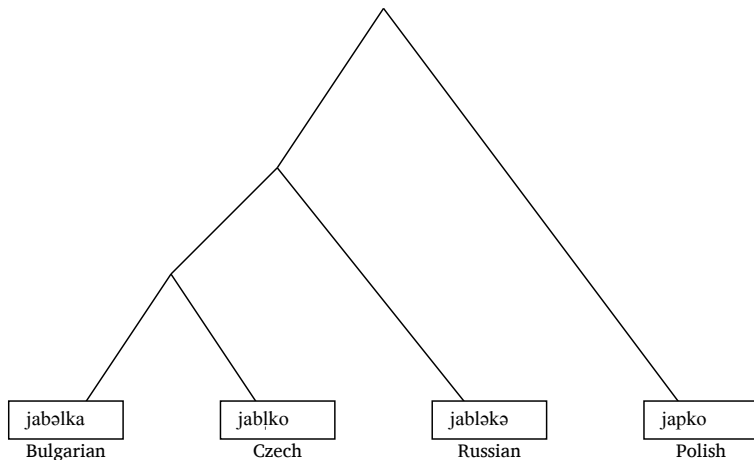**Performance of the Method**

Pairwise Sequence Alignment
**Multiple Sequence Alignment**

## Multiple Sequence Alignment

→ Align all sequences pairwise and store the scores in a matrix.
→ Construct a guide tree from the matrix.
→ Align all sequences along the guide tree, going from its leaves to its root.

**Sequences**
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
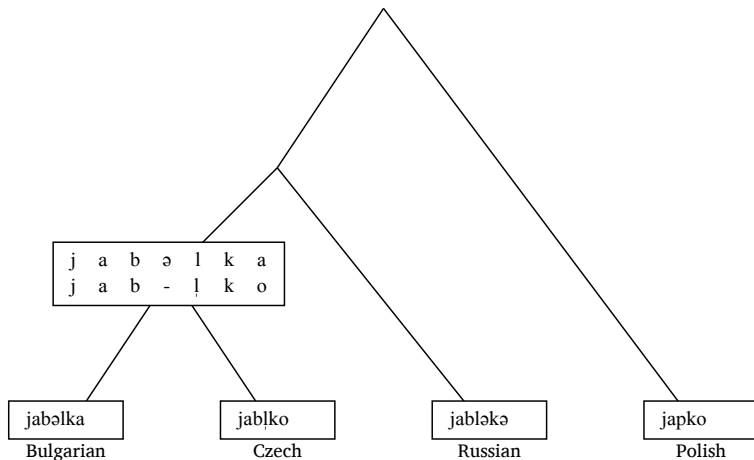A New Method for Multiple Sequence Alignment
Performance of the Method

Pairwise Sequence Alignment
**Multiple Sequence Alignment**

# Multiple Sequence Alignment



*jabǎlka*
'apple'
Bulgarian

*jablko*
'apple'
Czech

*jabloko*
'apple'
Russian

*jabłko*
'apple
Polish

**Sequences**
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Pairwise Sequence Alignment
**Multiple Sequence Alignment**

## Multiple Sequence Alignment

| jabəlka | jabl̩ko | jabləkə | japko |
|---------|---------|---------|-------|
| Bulgarian | Czech | Russian | Polish |

**Sequences**
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Pairwise Sequence Alignment
**Multiple Sequence Alignment**

# Multiple Sequence Alignment

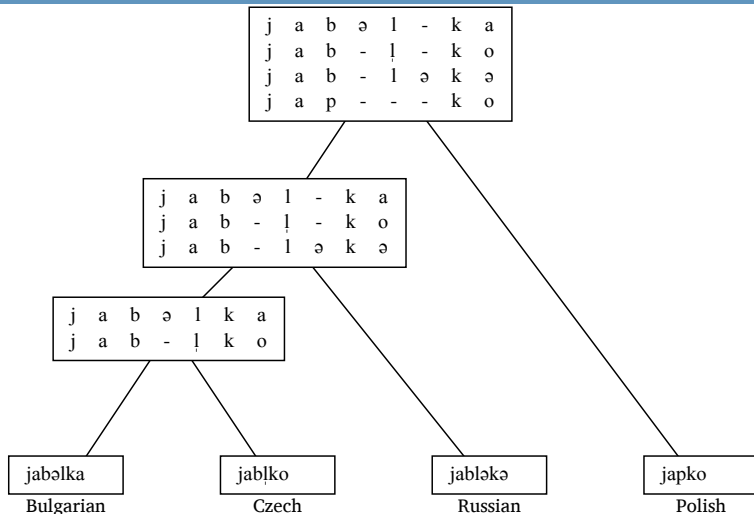**Sequences**
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Pairwise Sequence Alignment
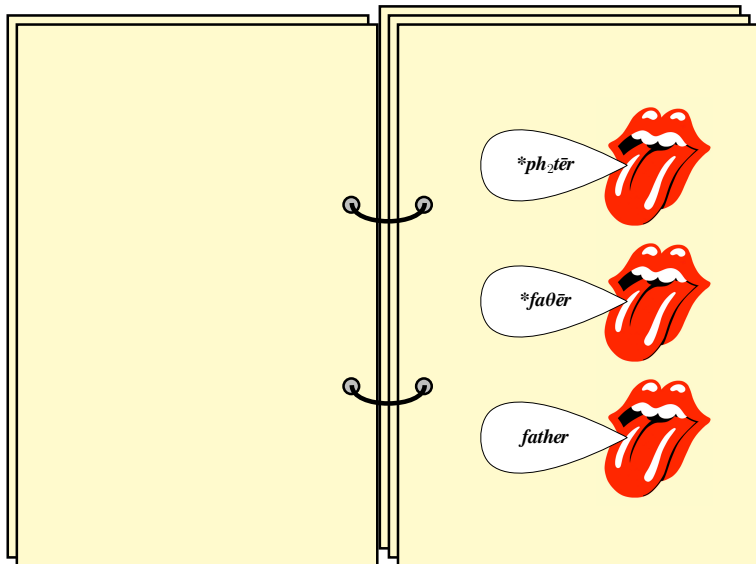**Multiple Sequence Alignment**

# Multiple Sequence Alignment



| jabəlka | jabḷko | jabləkə | japko |
| Bulgarian | Czech | Russian | Polish |

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Pairwise Sequence Alignment
**Multiple Sequence Alignment**

# Multiple Sequence Alignment

Sequences
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Pairwise Sequence Alignment
**Multiple Sequence Alignment**

# Multiple Sequence Alignment

**Sequences**
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
**A New Method for Multiple Sequence Alignment**
Performance of the Method

Pairwise Sequence Alignment
**Multiple Sequence Alignment**

# Multiple Sequence Alignment

**Sequences**
**Alignment Analyses**
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Pairwise Sequence Alignment
**Multiple Sequence Alignment**

# Multiple Sequence Alignment

Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

Sequence Similarity
Sound Classes

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Sequence Similarity
Sound Classes

# Sequence Comparison in Historical Linguistics

*German*

| tsʰ | aː | n |
|---|---|---|

*English*

| t | ʊː | θ |
|---|---|---|

*Italian*

| d | ɛ | n | t | e |
|---|---|---|---|---|

*French*

| d | ã |
|---|---|

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Sequence Similarity
Sound Classes

# Sequence Comparison in Historical Linguistics

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Sequence Similarity
Sound Classes

# Sequence Comparison in Historical Linguistics

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Sequence Similarity
Sound Classes

# Sequence Comparison in Historical Linguistics

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Sequence Similarity
Sound Classes

# Sequence Comparison in Historical Linguistics

*Proto-Germanic*    | t | | a | | n | | θ | | - |

*Proto-Romance*    | d | | e | | n | | t | | e |

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Sequence Similarity
Sound Classes

# Sequence Comparison in Historical Linguistics



*Proto-Germanic* | t | a | n | θ | - |

*Proto-Indo-European* | d | e | n | t | - |

*Proto-Romance* | d | e | n | t | e |

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Sequence Similarity
Sound Classes

## Sequence Comparison in Historical Linguistics

*Proto-Indo-European*

| d | e | n | t |
|---|---|---|---|

Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

Sequence Similarity
Sound Classes

# Sequence Comparison in Historical Linguistics

Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

**Sequence Similarity**
Sound Classes

# Sequence Similarity

**Synchronic Sequence Similarity**

**Diachronic Sequence Similarity**

Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

**Sequence Similarity**
Sound Classes

## Sequence Similarity

**Synchronic Sequence Similarity**

Sequences are judged to be similar if the segments of the sequences are phonetically similar ('phenotypic resemblence', Lass 1997).

**Diachronic Sequence Similarity**

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Sequence Similarity**
Sound Classes

## Sequence Similarity

**Synchronic Sequence Similarity**

Sequences are judged to be similar if the segments of the sequences are phonetically similar ('phenotypic resemblence', Lass 1997).

**Diachronic Sequence Similarity**

Sequences are judged to be similar if the segments of the sequences correspond *systematically* ('genotypic resemblence', Lass 1997).

Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

**Sequence Similarity**
Sound Classes

## Sequence Similarity

**Synchronic Sequence Similarity**

**Diachronic Sequence Similarity**

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Sequence Similarity**
**Sound Classes**

## Sequence Similarity

**Synchronic Sequence Similarity**

| Greek | mati | 'eye' | ≈ | Malay | mata | 'eye' |
|-------|------|-------|---|-------|------|-------|
| Greek | θεɔs | 'god' | ≈ | Spanish | diɔs | 'god' |

**Diachronic Sequence Similarity**

Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

**Sequence Similarity**
Sound Classes

## Sequence Similarity

**Synchronic Sequence Similarity**

| Greek | mati | 'eye' | ≈ | Malay | mata | 'eye' |
|-------|------|-------|---|-------|------|-------|
| Greek | θɛɔs | 'god' | ≈ | Spanish | diɔs | 'god' |

**Diachronic Sequence Similarity**

| German | tsʰaːn | 'tooth' | ≈ | English | tʊːθ | 'tooth' |
|--------|--------|---------|---|---------|------|---------|
| Spanish | etʃo | 'fact' | ≈ | French | fɛ | 'fact' |

Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

Sequence Similarity
**Sound Classes**

# Sound Classes

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Sequence Similarity**
**Sound Classes**

## Sound Classes

→ '[Even] the most divergent languages show examples of phonetic change which are remarkably similar' (Arlotto 1972: 77).

Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

Sequence Similarity
**Sound Classes**

## Sound Classes

→  '[Even] the most divergent languages show examples of phonetic change which are remarkably similar' (Arlotto 1972: 77).

→  Sounds which often occur in correspondence relations in genetically related languages can be clustered into *classes* (cf. Dolgopolsky 1986, Burlak & Starostin 2005).

Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

Sequence Similarity
**Sound Classes**

## Sound Classes

→ '[Even] the most divergent languages show examples of phonetic change which are remarkably similar' (Arlotto 1972: 77).

→ Sounds which often occur in correspondence relations in genetically related languages can be clustered into *classes* (cf. Dolgopolsky 1986, Burlak & Starostin 2005).

→ In contrast to the pure notion of synchronic and diachronic similarity, *sound classes* incorporate phonetic detail **and** systematic correspondence patterns within a probabilistic framework.
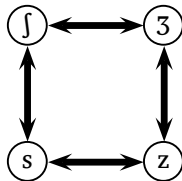
Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

Sequence Similarity
**Sound Classes**

# Sound Classes

Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

Sequence Similarity
**Sound Classes**

# Sound Classes

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

Sequence Similarity
**Sound Classes**

# Sound Classes

Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

Sequence Similarity
**Sound Classes**

# Sound Classes

Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

Sequence Similarity
**Sound Classes**

# Current Sound Class Approaches

**Dolgopolsky (1986)**

**Holman et al. (2011)**

Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

Sequence Similarity
**Sound Classes**

# Current Sound Class Approaches

**Dolgopolsky (1986)**

Based on an empirical basis, speech sounds are divided into ten types, distinguished 'in such a way that phonetic correspondences inside a "type" are more regular than those between different types' (Dolgopolsky 1986: 35).

**Holman et al. (2011)**

Sequences
Alignment Analyses
**Sequence Comparison in Historical Linguistics**
A New Method for Multiple Sequence Alignment
Performance of the Method

Sequence Similarity
**Sound Classes**

## Current Sound Class Approaches

**Dolgopolsky (1986)**

Based on an empirical basis, speech sounds are divided into ten types, distinguished 'in such a way that phonetic correspondences inside a "type" are more regular than those between different types' (Dolgopolsky 1986: 35).

**Holman et al. (2011)**

Sounds are represented in ASJP code (Brown et al. 2008), a transcription system, which reduces the full range of the IPA alphabet to 41 symbols. An automatic approach for the calculation of the frequency of sound correspondences is applied to a large database (Wichmann et al. 2010), giving the possibility to determine transition probabilities between the 41 sound classes.

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
**A New Method for Multiple Sequence Alignment**
Performance of the Method

Main Ideas
Working Procedure
Implementation

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
**A New Method for Multiple Sequence Alignment**
Performance of the Method

**Main Ideas**
Working Procedure
Implementation

# Main Ideas

**Sound Classes**

**Scoring Functions**

**Position-Specific Scoring**

**Swap Check**

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
**Working Procedure**
**Implementation**

## Main Ideas

**Sound Classes**

Phonetic sequences are internally represented as sound classes.

**Scoring Functions**

**Position-Specific Scoring**

**Swap Check**

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
Working Procedure
Implementation

## Main Ideas

**Sound Classes**

Phonetic sequences are internally represented as sound classes.

**Scoring Functions**

Scoring functions define specific transition probabilities among sound classes.

**Position-Specific Scoring**

**Swap Check**

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
Working Procedure
Implementation

## Main Ideas

**Sound Classes**

Phonetic sequences are internally represented as sound classes.

**Scoring Functions**

Scoring functions define specific transition probabilities among sound classes.

**Position-Specific Scoring**

Substitution scores vary according to prosodic context.

**Swap Check**

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
Working Procedure
Implementation

## Main Ideas

**Sound Classes**

Phonetic sequences are internally represented as sound classes.

**Scoring Functions**

Scoring functions define specific transition probabilities among sound classes.

**Position-Specific Scoring**

Substitution scores vary according to prosodic context.

**Swap Check**

Alignments are automatically searched for swapped sites.

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
**A New Method for Multiple Sequence Alignment**
Performance of the Method

**Main Ideas**
Working Procedure
Implementation

# Position-Specific Scoring

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
Working Procedure
Implementation

## Position-Specific Scoring

→    It is assumed that sound change occurs more frequently in prosod-
ically *weak* positions of phonetic sequences (cf. Geisler 1992).

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
Working Procedure
Implementation

## Position-Specific Scoring

→   It is assumed that sound change occurs more frequently in prosodically *weak* positions of phonetic sequences (cf. Geisler 1992).

→   Given the sonority structure of a phonetic sequence, one can, apart from the *initial* and *final* positions, distinguish positions of *ascending*, *maximum* and *descending* sonority.

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
**Working Procedure**
**Implementation**

## Position-Specific Scoring

→ It is assumed that sound change occurs more frequently in prosodically *weak* positions of phonetic sequences (cf. Geisler 1992).

→ Given the sonority structure of a phonetic sequence, one can, apart from the *initial* and *final* positions, distinguish positions of *ascending*, *maximum* and *descending* sonority.

→ These positions can be ordered in a *hierarchy of strength* (initial > ascending > descending > maximum > final).

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
Working Procedure
Implementation

## Position-Specific Scoring

→   It is assumed that sound change occurs more frequently in prosod-ically *weak* positions of phonetic sequences (cf. Geisler 1992).

→   Given the sonority structure of a phonetic sequence, one can, apart from the *initial* and *final* positions, distinguish positions of *ascending*, *maximum* and *descending* sonority.

→   These positions can be ordered in a *hierarchy of strength* (initial > ascending > descending > maximum > final).

→   Based on the relative strength of all sites in a phonetic sequence, substitution scores and gap penalties are modified by scaling fac-tors, favoring changes in weaker positions and aggravating them in stronger positions.

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
Working Procedure
Implementation

# Position-Specific Scoring

**j     a     b     ə     l     k     a**

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
Working Procedure
Implementation

# Position-Specific Scoring



sonority
increases

**j a b ə l k a**

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
Working Procedure
Implementation

## Position-Specific Scoring

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
**Working Procedure**
**Implementation**

# Position-Specific Scoring

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Main Ideas
Working Procedure
Implementation

## Working Procedure

INPUT SEQUENCES

jabl̩ko
jabəlka
jabləkə
japkɔ

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
**A New Method for Multiple Sequence Alignment**
Performance of the Method

Main Ideas
**Working Procedure**
Implementation

# Working Procedure

## SOUND CLASS CONVERSION

| | | |
|---|---|---|
| jabłko | → | yablko |
| jabəlka | → | yab3lka |
| jabləkə | → | yabl3k3 |
| japkɔ | → | yapko |

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Main Ideas
Working Procedure
Implementation

# Working Procedure

```
┌─────────────────────────────────┐
│  DISTANCE CALCULATION           │
└─────────────────────────────────┘

yablko   0.00  0.14  0.34  0.12
yab3lka  0.14  0.00  0.46  0.28
yabl3k3  0.34  0.46  0.00  0.44
yapko    0.12  0.28  0.44  0.00
```

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Main Ideas
**Working Procedure**
Implementation

## Working Procedure



CLUSTER ANALYSIS

```
yablko
yab3lka
yabl3k3
yapko
```

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
**Working Procedure**
Implementation

## Working Procedure

PROGRESSIVE ALIGNMENT

```
yablko
yab3lka
yabl3k3
yapko
```

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Main Ideas
Working Procedure
Implementation

## Working Procedure

PROGRESSIVE ALIGNMENT

```
y   a   b   –   l   k   o
y   a   b   3   l   k   a
yabl3k3
yapko
```

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Main Ideas
**Working Procedure**
Implementation

## Working Procedure

```
┌─────────────────────────────────┐
│  PROGRESSIVE ALIGNMENT          │
└─────────────────────────────────┘

y    a    b    –    l    –    k    o
y    a    b    3    l    –    k    a
y    a    b    –    l    3    k    3
yapko
```

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
Performance of the Method

Main Ideas
Working Procedure
Implementation

# Working Procedure

```
PROGRESSIVE ALIGNMENT

y   a   b   –   l   –   k   o
y   a   b   3   l   –   k   a
y   a   b   –   l   3   k   3
y   a   p   –   –   –   k   o
```

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
**Working Procedure**
Implementation

# Working Procedure



SWAP CHECK

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| y | a | b | – | l | – | k | o |
| y | a | b | 3 | l | – | k | a |
| y | a | b | – | l | 3 | k | 3 |
| y | a | p | – | – | – | k | o |

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
**Working Procedure**
**Implementation**

## Working Procedure

| IPA CONVERSION |
|---|

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| y | a | b | ... | → | j | a | b | ... |
| y | a | b | ... | → | j | a | b | ... |
| y | a | b | ... | → | j | a | b | ... |
| y | a | p | ... | → | j | a | p | ... |

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Main Ideas
**Working Procedure**
Implementation

## Working Procedure

| OUTPUT MSA |
|---|

| j | a | b | - | l | k | o |
|---|---|---|---|---|---|---|
| j | a | b | ə | l | k | a |
| j | a | b | l | ə | k | ə |
| j | a | p | - | - | k | ɔ |

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Main Ideas**
**Working Procedure**
**Implementation**

## Implementation

The algorithm is implemented as part of the LingPy library (List 2011, see `http://lingulist.de/lingpy/`). LingPy is a suite of open source Python modules for sequence comparison, distance analyses, data operations and visualization methods in quantitative historical linguistics.

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Evaluation
Results

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Evaluation**
**Results**

# Evaluation

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Evaluation**
Results

## Evaluation

→ In biological analyses, the performance of alignment algorithms is traditionally tested by comparing manually edited alignments (*reference alignment*) with those produced by the respective algorithms (*test alignment*).

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Evaluation**
Results

## Evaluation

→ In biological analyses, the performance of alignment algorithms is traditionally tested by comparing manually edited alignments (*reference alignment*) with those produced by the respective algorithms (*test alignment*).

→ There exist different evaluation measures for determining the goodness of an algorithm.

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
**Performance of the Method**

**Evaluation**
Results

# Evaluation

→  In biological analyses, the performance of alignment algorithms is traditionally tested by comparing manually edited alignments (*reference alignment*) with those produced by the respective algorithms (*test alignment*).

→  There exist different evaluation measures for determining the goodness of an algorithm.

**(a)**  The *percentage of identical columns* score (PIC) calculates how many columns match in the reference and the test alignment.

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Evaluation**
Results

## Evaluation

→ In biological analyses, the performance of alignment algorithms is traditionally tested by comparing manually edited alignments (*reference alignment*) with those produced by the respective algorithms (*test alignment*).

→ There exist different evaluation measures for determining the goodness of an algorithm.

**(a)** The *percentage of identical columns* score (PIC) calculates how many columns match in the reference and the test alignment.

**(b)** The *percentage of identical rows* score (PIR) calculates how many rows match in the reference and the test alignment.

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
**Performance of the Method**

**Evaluation**
Results

## Evaluation

→ In biological analyses, the performance of alignment algorithms is traditionally tested by comparing manually edited alignments (*reference alignment*) with those produced by the respective algorithms (*test alignment*).

→ There exist different evaluation measures for determining the goodness of an algorithm.

**(a)** The *percentage of identical columns* score (PIC) calculates how many columns match in the reference and the test alignment.

**(b)** The *percentage of identical rows* score (PIR) calculates how many rows match in the reference and the test alignment.

**(c)** The *sum-of-pairs* score (SOP) calculates the size of the intersection of aligned pairs of residues in the reference and the test alignment divided by the size of aligned pairs of residues in the reference alignment (cf. Thompson et al. 1999).

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Evaluation**
**Results**

## Evaluation

→    In biological analyses, the performance of alignment algorithms is traditionally tested by comparing manually edited alignments (*reference alignment*) with those produced by the respective algorithms (*test alignment*).

→    There exist different evaluation measures for determining the goodness of an algorithm.

**(a)**   The *percentage of identical columns* score (PIC) calculates how many columns match in the reference and the test alignment.

**(b)**   The *percentage of identical rows* score (PIR) calculates how many rows match in the reference and the test alignment.

**(c)**   The *sum-of-pairs* score (SOP) calculates the size of the intersection of aligned pairs of residues in the reference and the test alignment divided by the size of aligned pairs of residues in the reference alignment (cf. Thompson et al. 1999).

**(d)**   The *modified rand index* (MRI) checks 'whether the same elements are together in the [test] alignment and the [reference] alignment' (Prokić et al. 2009: 21).

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Evaluation**
Results

# The BulDial Gold Standard (Prokić et al. 2009)

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Evaluation**
Results

# The BulDial Gold Standard (Prokić et al. 2009)

→    152 manually edited MSAs

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Evaluation**
Results

# The BulDial Gold Standard (Prokić et al. 2009)

→ 152 manually edited MSAs
→ 192 taxa (dialect points)

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Evaluation**
Results

## The BulDial Gold Standard (Prokić et al. 2009)

→ 152 manually edited MSAs

→ 192 taxa (dialect points)

→ ca. 30,000 sequences

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
**Performance of the Method**

**Evaluation**
Results

## Two Test Models

|                     | **DOLGO**          | **ASJP**           |
| ------------------- | ------------------ | ------------------ |
| **Sound Classes**   | Dolgopolsky (1986) | ASJP-Code          |
| **No. of Symbols**  | 11                 | 41                 |
| **Scoring Function**| simple matching    | Brown et al. (2011)|
| **Default Gap Pen.**| -6                 | -12                |
| **Main Algorithm**  | LingPy             | LingPy             |

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

**Evaluation**
**Results**

## Two Test Models

|                      | **ALPHA (modif. by Prokić et al. 2009)** |
| -------------------- | ---------------------------------------- |
| **Sound Classes**    | no                                       |
| **No. of Symbols**   | full IPA                                 |
| **Scoring Function** | CV distinction                           |
| **Default Gap Pen.** | unknown                                  |
| **Main Algorithm**   | ALPHAMALIG                               |

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
**Performance of the Method**

Evaluation
**Results**

## Results

|  | ASJP | DOLGO | ALPHA |
|---|---|---|---|
| **Perfect Alignments** | 132 (87%) | 123 (81%) | 103 (69%) |
| **Perc. of Ident. Col.** | 0.9313 | 0.8952 | 0.8409 |
| **Per. of Ident. Rows** | 0.9531 | 0.9043 | 0.7632 |
| **Sum of Pairs** | 0.9901 | 0.9855 | 0.9825 |
| **Modified Rand Index** | 0.9902 | 0.9844 | 0.9824 |

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
**Performance of the Method**

Evaluation
**Results**

# Results

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
**Performance of the Method**

Evaluation
**Results**

# Advantages of LingPy

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Evaluation
**Results**

## Advantages of LingPy

→ The method for swap detection identifies 19 of 21 swapped sites in LingPy-ASJP, 13 of them are aligned properly.

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Evaluation
**Results**

## Advantages of LingPy

→ The method for swap detection identifies 19 of 21 swapped sites in LingPy-ASJP, 13 of them are aligned properly.

→ The scoring function in LingPy-ASJP is based on a large empirical basis, allowing fine distinctions along with the extended model of sound classes.

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
**Performance of the Method**

Evaluation
**Results**

## Advantages of LingPy

→ The method for swap detection identifies 19 of 21 swapped sites in LingPy-ASJP, 13 of them are aligned properly.

→ The scoring function in LingPy-ASJP is based on a large empirical basis, allowing fine distinctions along with the extended model of sound classes.

→ The application of prosodic profiles enhances both the calculation of the guide tree and the alignment process.

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
**Performance of the Method**

Evaluation
**Results**

# Examples

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Evaluation
**Results**

# Examples

**ALPHAMALIG**

| Aldomirovci | v | - | r | ɑ | - | ʧ | ɑ | m |
| Asparuhovo | v | - | r | ɣ | ʃ | t | ɑ | m |
| Panagjurishte | v | ɣ | r | - | ʃ | t | ə | m |
| Rakovica | v | - | r̩ | - | ʃ | t | ɑ | m |
| Stambolovo | v | - | r | ɣ | ç | t | ə | m |

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Evaluation
**Results**

# Examples

### ALPHAMALIG

| Aldomirovci | v | - | r | ɑ | - | tʃ | ɑ | m |
|---|---|---|---|---|---|---|---|---|
| Asparuhovo | v | - | r | ɣ | ʃ | t | ɑ | m |
| Panagjurishte | v | ɣ | r | - | ʃ | t | ə | m |
| Rakovica | v | - | r̩ | - | ʃ | t | ɑ | m |
| Stambolovo | v | - | r | ɣ | ç | t | ə | m |

### LingPy-ASJP

| Aldomirovci | v | r | ɑ | - | tʃ | ɑ | m |
|---|---|---|---|---|---|---|---|
| Asparuhovo | v | r | ɣ | ʃ | t | ɑ | m |
| Panagjurishte | v | ɣ | r | ʃ | t | ə | m |
| Rakovica | v | r̩ | - | ʃ | t | ɑ | m |
| Stambolovo | v | r | ɣ | ç | t | ə | m |

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

Sequences
Alignment Analyses
Sequence Comparison in Historical Linguistics
A New Method for Multiple Sequence Alignment
**Performance of the Method**

Evaluation
**Results**

# Examples

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Evaluation
**Results**

## Examples

### LingPy-DOLGO

| | | | | | | |
|---|---|---|---|---|---|---|
| Aldomirovci | u | n | e | t | r | e |
| Asparuhovo | - | v | ɣ | t | rʲ | ə |
| Babjak | f | n | e | t | r | e |
| Bachkovo | - | v | ɑ | t | rʲ | ə |
| Bagrenci | u | n | e | t | r | e |

**Sequences**
**Alignment Analyses**
**Sequence Comparison in Historical Linguistics**
**A New Method for Multiple Sequence Alignment**
**Performance of the Method**

Evaluation
**Results**

# Examples

## LingPy-DOLGO

| | | | | | | |
|---|---|---|---|---|---|---|
| Aldomirovci | u | n | e | t | r | e |
| Asparuhovo | - | v | ɤ | t | rʲ | ə |
| Babjak | f | n | e | t | r | e |
| Bachkovo | - | v | ɑ | t | rʲ | ə |
| Bagrenci | u | n | e | t | r | e |

## LingPy-ASJP

| | | | | | | |
|---|---|---|---|---|---|---|
| Aldomirovci | u | n | e | t | r | e |
| Asparuhovo | v | - | ɤ | t | rʲ | ə |
| Babjak | f | n | e | t | r | e |
| Bachkovo | v | - | ɑ | t | rʲ | ə |
| Bagrenci | u | n | e | t | r | e |

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF