









Academy of Finland project

- Automatically discovery language family relationships
- Analysis of language data
- Without relying on language-specific tools, resources, etc.





Academy of Finland project

- Automatically discovery language family relationships
- Analysis of language data
- Without relying on language-specific tools, resources, etc.





Academy of Finland project

- Automatically discovery language family relationships
- Analysis of language data
- Without relying on language-specific tools, resources, etc.





Academy of Finland project

- Automatically discovery language family relationships
- Analysis of language data
- Without relying on language-specific tools, resources, etc.





- Standard language processing tools not available
 - part-of-speech taggers, parsers, semantic role labellers
- Not much annotated data
 - E.g. large corpora of expert linguistic analyses
 - But raw data might be available.
- E.g. many Uralic languages
 - Not many speakers
 - E.g. Nganasan, ~100 speakers





- Standard language processing tools not available part-of-speech taggers, parsers, semantic role labellers
- Not much annotated data
 - E.g. large corpora of expert linguistic analyses
 - But raw data might be available.
- E.g. many Uralic languages
 - Not many speakers
 - E.g. Nganasan, ~100 speakers





- Standard language processing tools not available part-of-speech taggers, parsers, semantic role labellers
- Not much annotated data
 - E.g. large corpora of expert linguistic analyses
 - But raw data might be available
- E.g. many Uralic languages
 - Not many speakers
 - E.g. Nganasan, ~100 speakers





- Standard language processing tools not available part-of-speech taggers, parsers, semantic role labellers
- Not much annotated data
 - E.g. large corpora of expert linguistic analyses
 - But raw data might be available
- E.g. many Uralic languages
 - Not many speakers
 - E.g. Nganasan, ~100 speakers





- Standard language processing tools not available part-of-speech taggers, parsers, semantic role labellers
- Not much annotated data
 - E.g. large corpora of expert linguistic analyses
 - But raw data might be available
- E.g. many Uralic languages
 - Not many speakers

E.g. Nganasan, ~100 speakers





- Standard language processing tools not available part-of-speech taggers, parsers, semantic role labellers
- Not much annotated data
 - E.g. large corpora of expert linguistic analyses
 - But raw data might be available
- E.g. many Uralic languages
 - Not many speakers

E.g. Nganasan, ~ 100 speakers





Language Typology

- Grouping of languages according to their characteristics
- Mapping multi-dimensional space of:
 - similarities, differences
 - influence of contact
 - syntax, morphology
 - phonotactics, prosody





- Instead of labour-intensive study of languages, see what we can learn automatically from surface
- Text, speech recordings, transcribed speech
- Working with Martti Vainio's group on speech data
- Also looking at text input







- Instead of labour-intensive study of languages, see what we can learn automatically from surface
- Text, speech recordings, transcribed speech
- Working with Martti Vainio's group on speech data
- Also looking at text input







- Instead of labour-intensive study of languages, see what we can learn automatically from surface
- Text, speech recordings, transcribed speech
- Working with Martti Vainio's group on speech data
- Also looking at text input







- Instead of labour-intensive study of languages, see what we can learn automatically from surface
- Text, speech recordings, transcribed speech
- Working with Martti Vainio's group on speech data
- Also looking at text input







Question:

How much can we discover automatically about languages we know nothing about to start with?







Current work in Discovery Group Kemenkerabr Toha G Kemekerbige Macc Mracb BMTHegb 3-4

Find similarities between **characters/sounds** of different languages, in terms of how they are used







Look at contexts they appear in







Look at contexts they appear in







Difficulty: contexts specific to each language, can't compare







Novel **neural network** method finds common patterns of co-occurence Promising results







Can possibly use to identify related words where languages

- use different sounds to express them
- or write them differently (e.g. different scripts)