# Extracting facts from medical texts

**Robert Roosalu**, supervisor Sven Laur.

University of Tartu, Faculty of Science and Technology, Institute of Computer Science, MSc Computer Science

UNIVERSITY OF TARTU
Institute of Computer Science

STACC

Study IT in .ee
sponsored by Skype

## Introduction

Estonian *Geenivaramu* has built a large corpora of free-form texts from doctors, describing patient's diagnosis, treatment, etc. Successfully mining factual knowledge from this kind of data would lead us closer towards advancements in personal medicine, automated clinical trials, etc.

One chosen approach is via human defined patterns, such as regular expressions, or context free grammars. The results of the patterns are too numerous to be manually evaluated. To enable such evaluation, we have built PatternExaminer. The tool clusters the contexts of the extracted facts, which are then sampled for the final evaluation.

For the success of PatternExaminer, we need the highest quality text clustering capabilities.
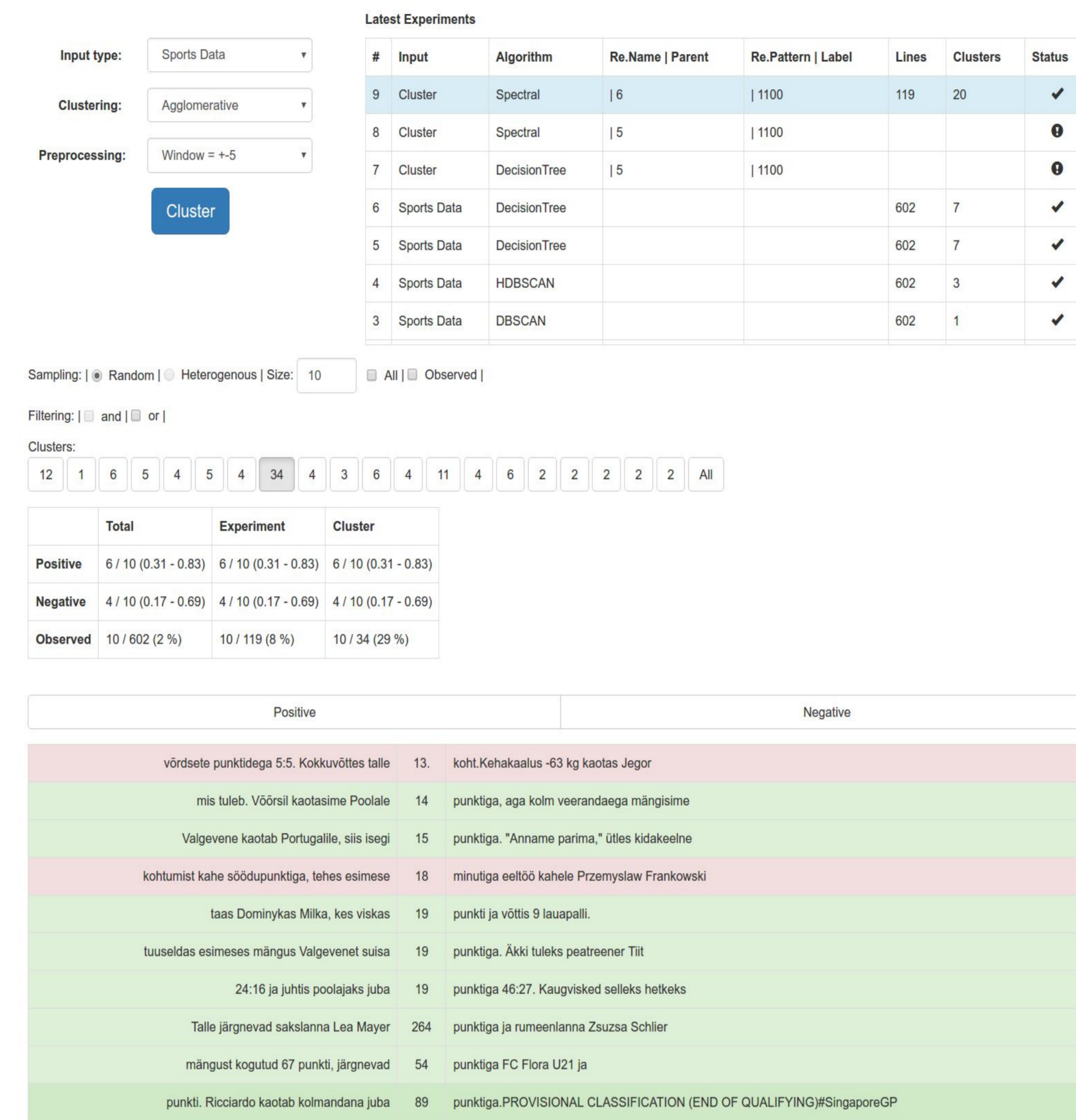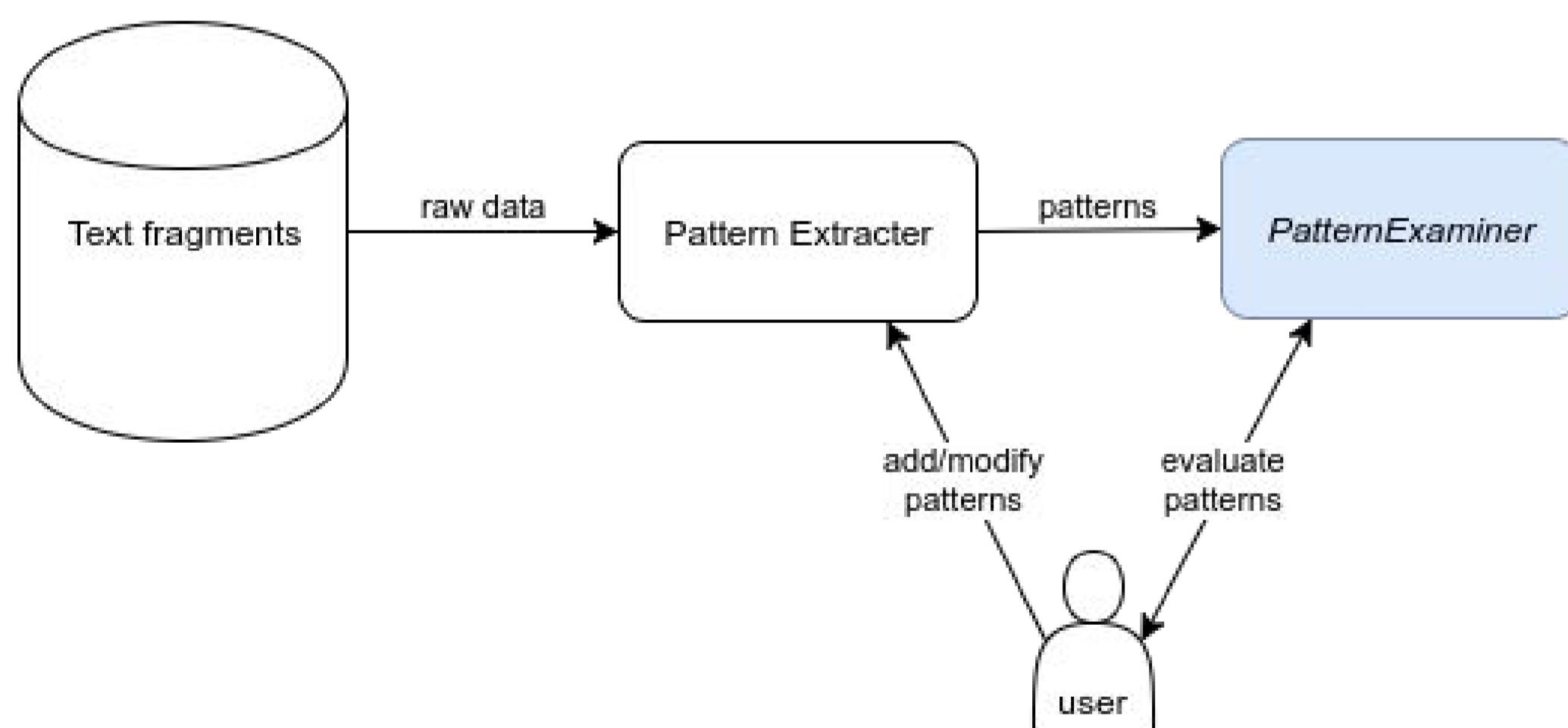
*Figure 2. User interface of PatternExaminer.*

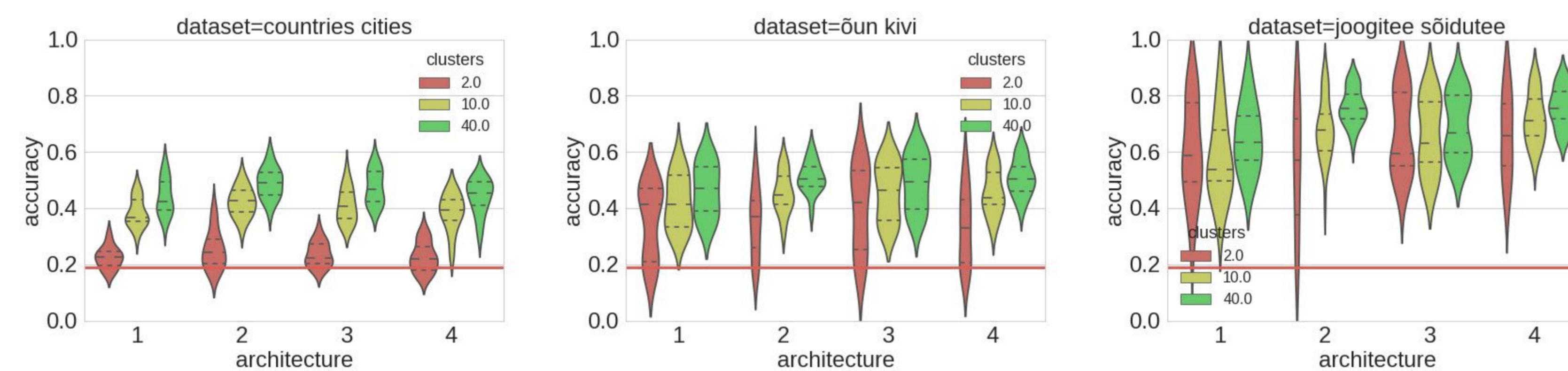*Figure 1. Workflow for validating the defined patterns.*

*Figure 3. Results of recurrent autoencoder experiments, Spectral clustering.*

## Experiments

For the dataset, eight different pairs of words were extracted from the Estonian Reference Corpus. Grid search was used over an array of various methods, parameters and hyperparameters, leading to a total of around 20000 different clustering experiments.

For evaluating the quality of clusterings, a custom set of metrics emulating the usage of PatternExaminer was employed.
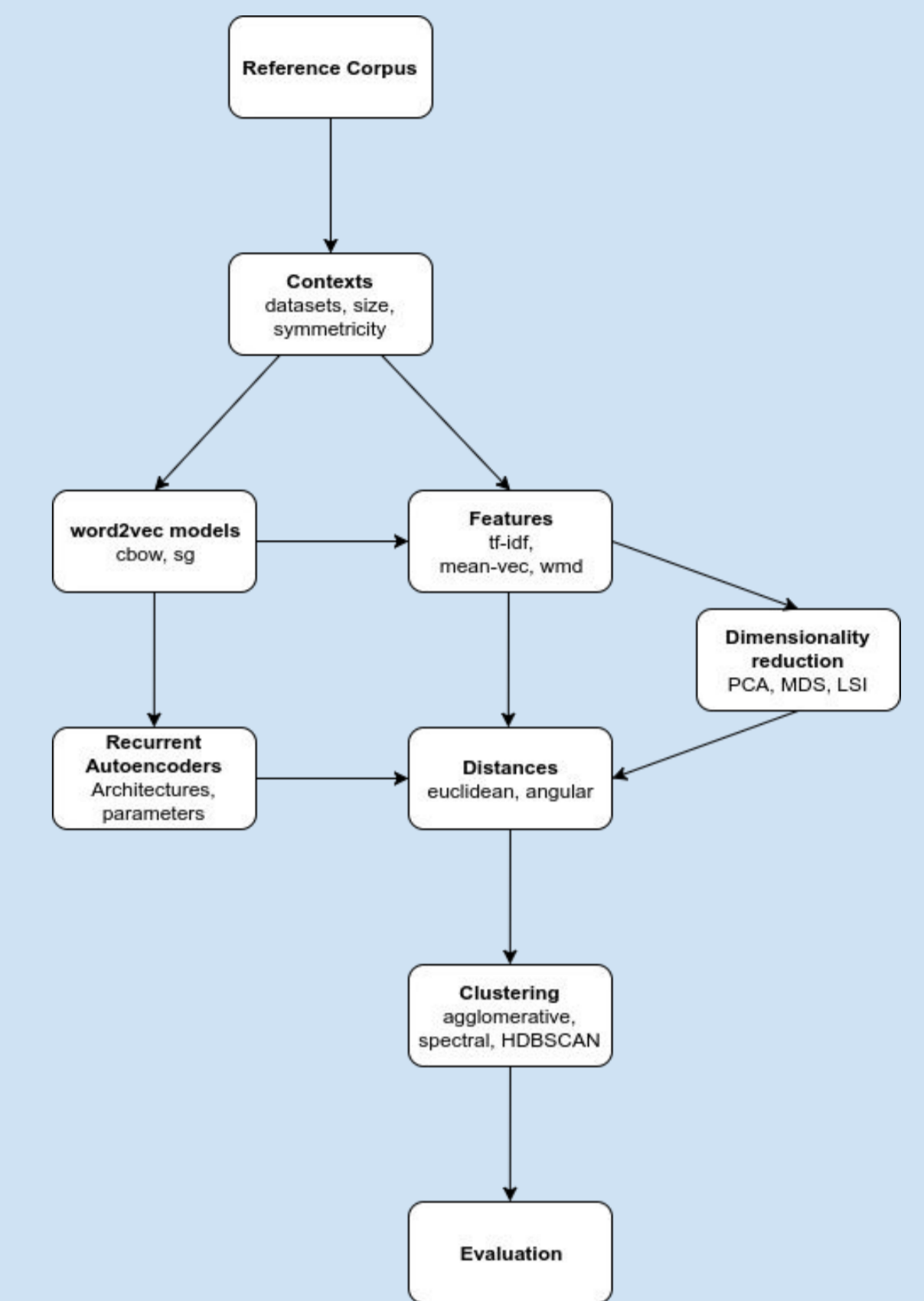
*Figure 5. Experiment pipeline.*

## Results

For analyzing such a large amount of experiment results, we used categorical faceted violin plots.

It was found, that spectral clustering is best for clustering textual data, in comparison with Agglomerative and HDBSCAN.

For the classical methods, a good context window is symmetric, but the size doesn't matter. Skip-gram word embeddings, averaged to the mean vectors and using angular distance yields the most separable embedding space. Dimensionality reduction did not prove useful in the experiments.

Simple recurrent neural network architectures achieved on-par or better results, than the mean-vector approach.

## Future work

Building a test set on the *Geenivaramu* dataset, to evaluate all following models.
Transfering the algorithms onto the *Geenivaramu* dataset, verifying if the current insights hold.
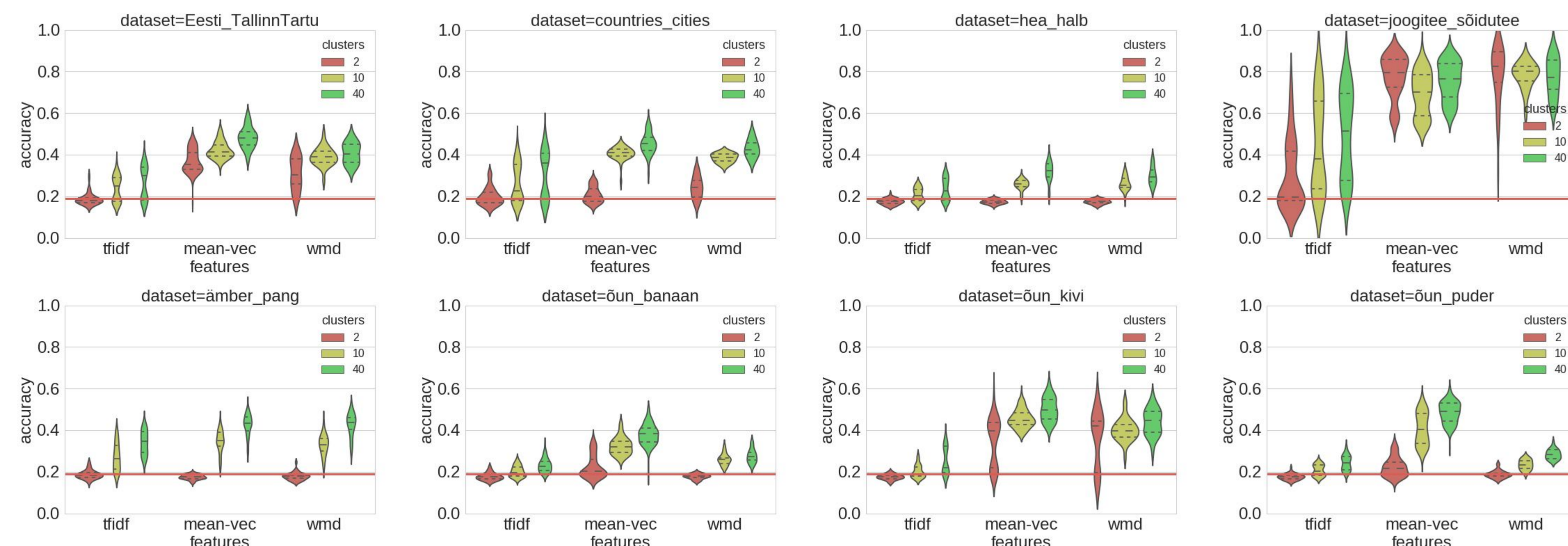Employing state-of-the art RNN architectures for the embedding.

*Figure 4. Various results, Spectral clustering.*