# ELAC: Ensemble Learning for Anaphora- and Coreference-Resolution-Systems

Thomas Bögel, Lukas Funk, Andreas Kull
Softwareprojekt SS10

http://dakhma.net/elac

## Motivation & Architecture

### Motivation

Anaphora- and Coreference-Resolution-Systems (ACRS) are typically very specialized on a single phenomenon and failing on others, especially when it comes to domain adaptation.

Therefore we created a system, built around the machine learning tool WEKA, to combine several single ACR-Systems.
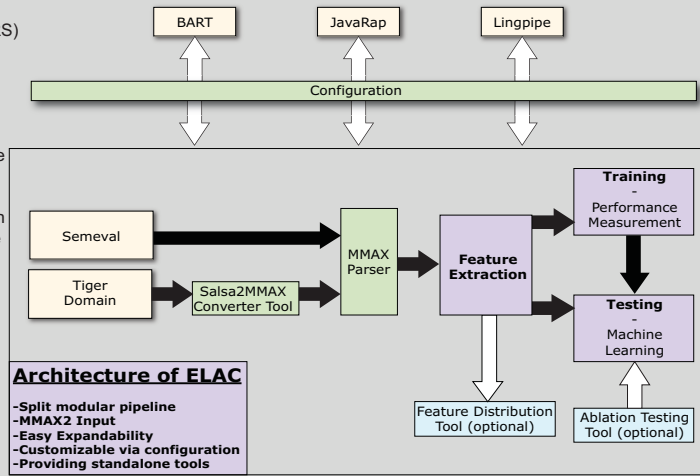
The goal was to test the ensemble learning approach in order to obtain improvements in comparison to a single system and on domain adaptation.
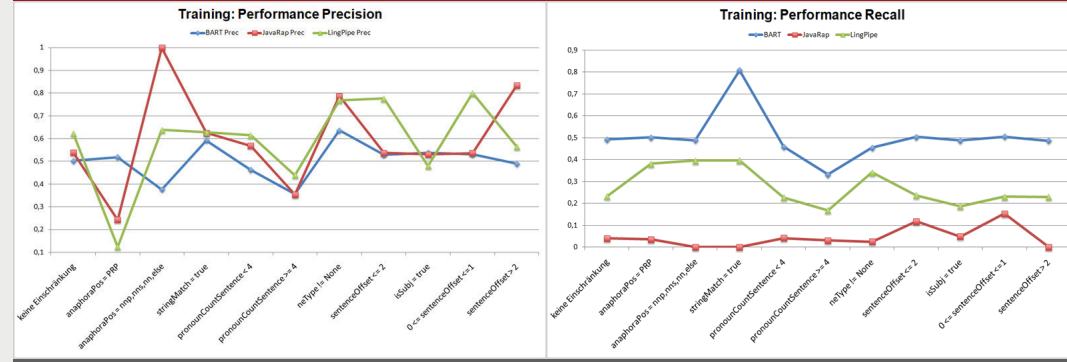
Our ressources:
– Training & Testing: Semeval2010 Task2
– Domain Adaptation: Sherlock Holmes
– MUC6 evaluation standard
– MMAX2 format

Our guidelines:
– Modularity
– Expandability
– Usability



**Architecture of ELAC**
-Split modular pipeline
-MMAX2 Input
-Easy Expandability
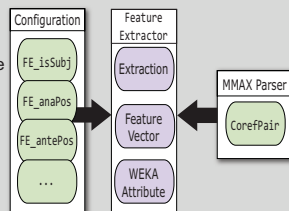-Customizable via configuration
-Providing standalone tools

## Training: Extraction & Measurement

### Feature Extractor

**Extraction:** The Coreference pairs and their features are extraced from the corpus.
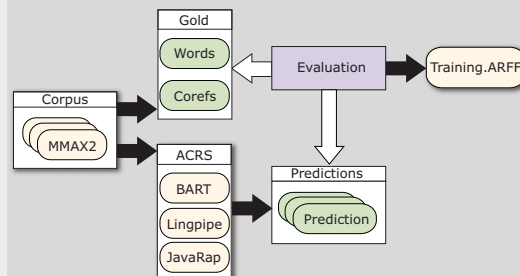
**Modularity:** New features can be added easily.



### Feature Dependent Performance Measurement

**Predictions of coreferences**, created by the ACR-Systems, are checked against the gold standard.
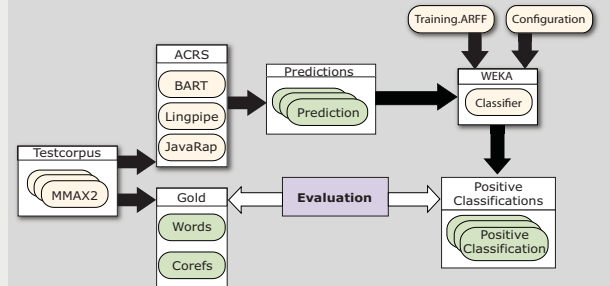
An **ARFF file for WEKA**, containing the positive and negative results of the evaluation, is created.



## Testing: Learning & Evaluation

### Machine Learning

**WEKA:** Classifier(s) are trained with an ARFF file
**ACR-Systems:** Predict anaphora and coreferences in test data
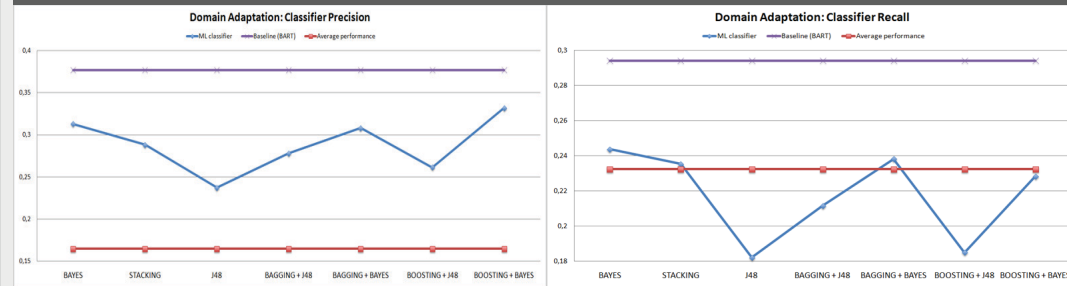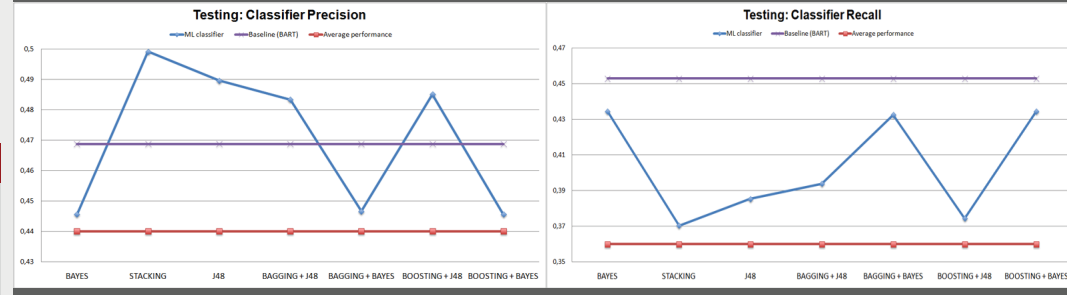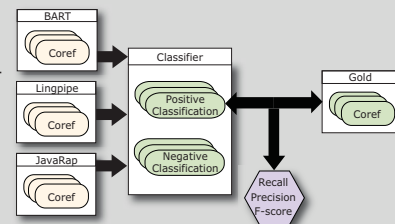**Classification:** Classifier classifies predictions as trustworthy or untrustworthy



### Evaluation

A **Positive list** with trust-worthy classifications is created by the classifier.

It is **compared to the gold standard** which is extracted from the test data.

**Recall, precision and F-score** are computed to evalute the Classifier(s)



## Findings: Data



Training: Performance Precision

Training: Performance Recall

Testing: Classifier Precision

Testing: Classifier Recall

Domain Adaptation: Classifier Precision

Domain Adaptation: Classifier Recall

## Findings: Conclusion

### Findings

**Training:**
– Performance of individual ACR-Systems with various features
– Based on 50% of Semeval2010 Task2 corpus
– MUC6 evaluation which causes problems for JavaRap
– BART is clearly dominating

**Testing:**
– Precision improved with an ensemble of classifiers
– Baseline's F-score not reached
– JavaRap made too many bad decisions

**Domain Adaptation:**
– Sherlock Holmes novel from Tiger/Salsa corpus
– BART has achieved the best F-score
– Baseline not reached
– ACR-Systems provided too much wrong data
– Better precision than the average of all systems

### Conclusion

– Ensemble Learning better than the average of all ACR-Systems
– BART is an overall robust system (only if trained properly)
– Best results probably with completely complementary systems