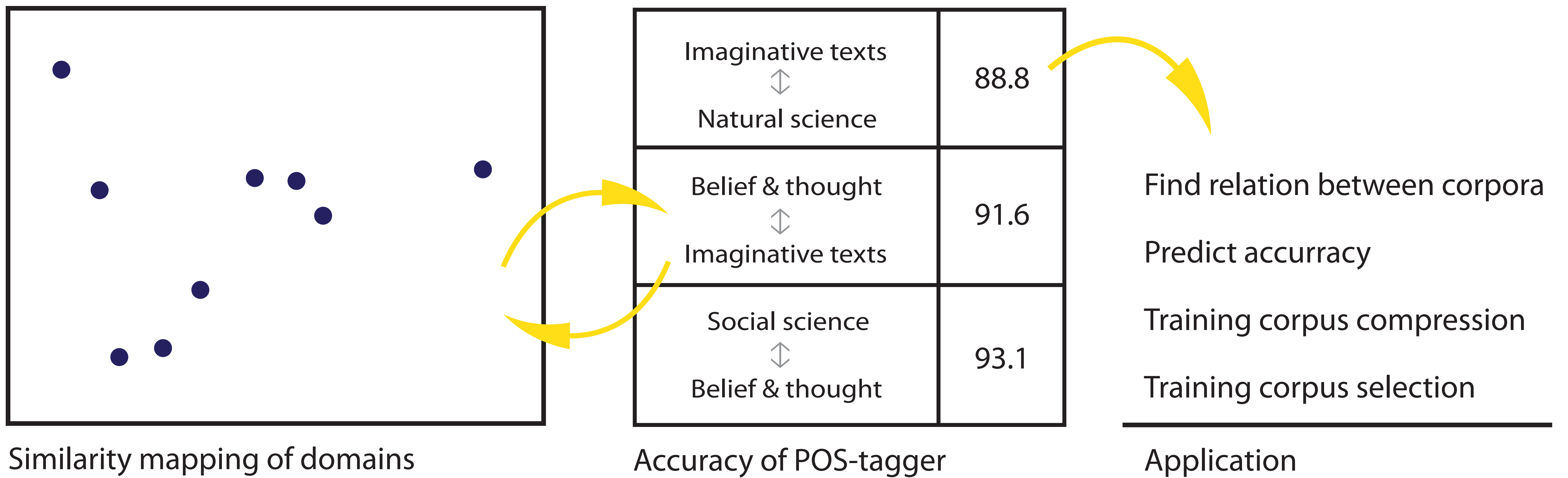


# USING DOMAIN SIMILARITY FOR PERFORMANCE ESTIMATION

VINCENT VAN ASCH  
WALTER DAELEMANS

CLiPS  
UNIVERSITY OF ANTWERP

OBJECTIVE: Trying to find a simple relation between token frequencies in corpora and the cross-domain accuracy of part-of-speech taggers.



Similarity mapping of domains

Accuracy of POS-tagger

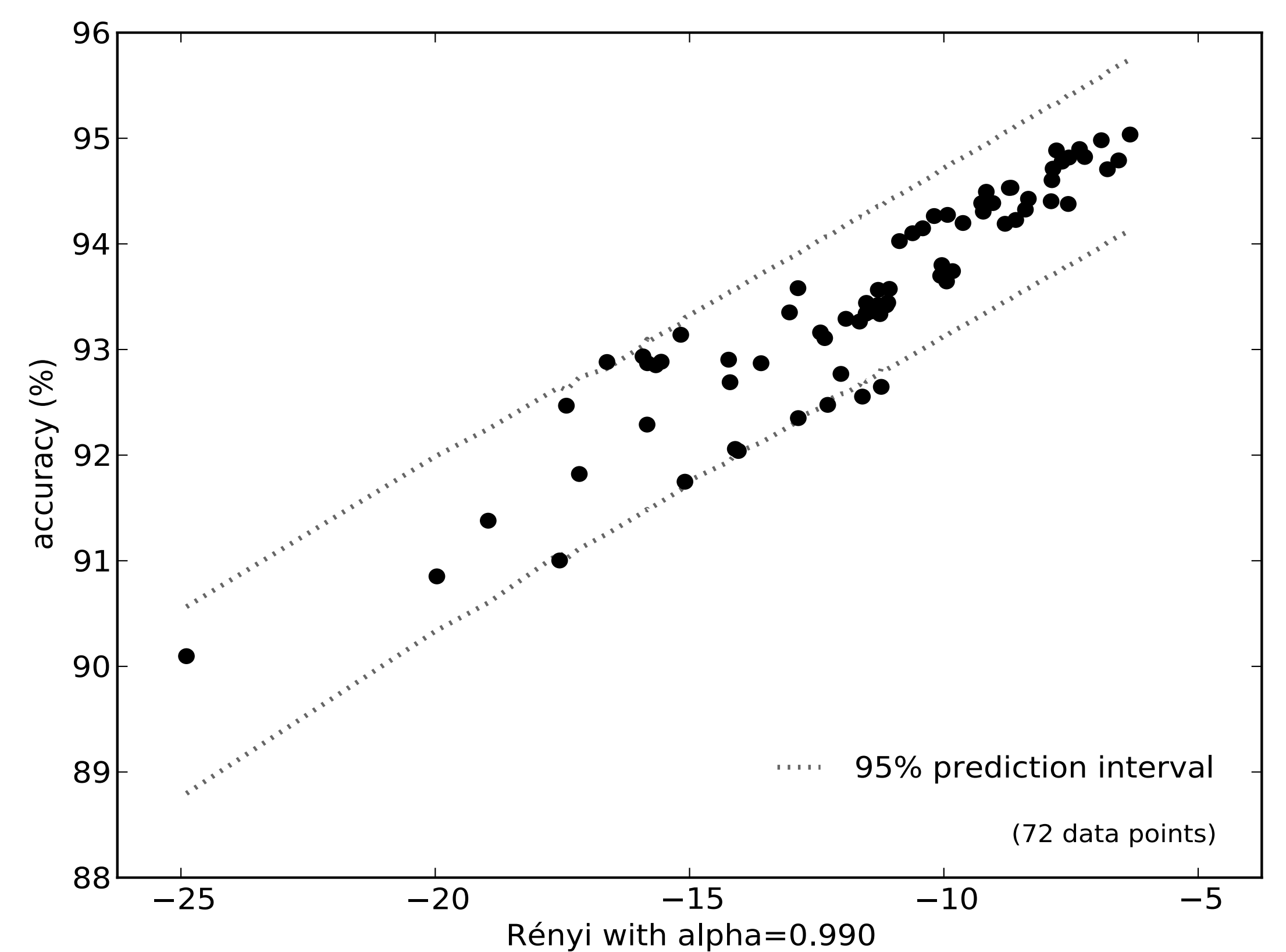
Application

Domains of BNC corpus	Arts	Applied science	Belief & thought	Commerce & finance	Leisure	Imaginative texts	Natural science	Social science	World affairs
Part-of-speech taggers	Rule base - majority			Memory-based - MBT			SVM-based - SVMTool		

Experimental design

	Imaginative texts	Natural science
smiled	203	1
sat	260	1
herself	363	1
development	1	161
DNA	1	184
data	1	233

Corpora as token frequencies



$$\frac{1}{\alpha-1} \log_2 \left( \sum_k p_k^{1-\alpha} q_k^\alpha \right)$$

Rényi-divergence with  $\alpha = 0.99$  gave the best linear correlations

Cross-validation experiments indicated:

- Linear relation between accuracy and metric
- Error reduction on accuracy prediction
- Linear relation for in-domain experiments