

QUITA

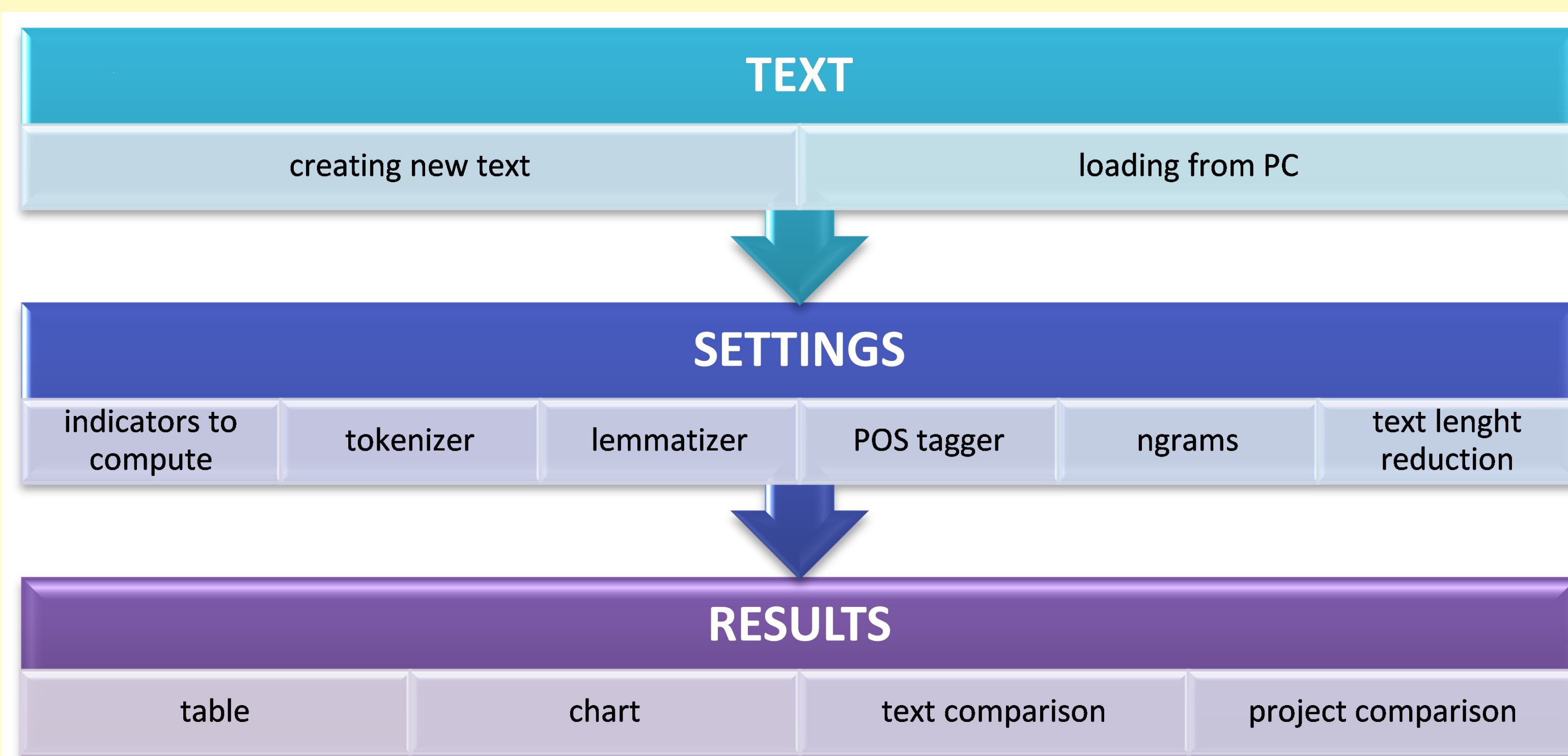
Quantitative Index Text Analyzer

Miroslav Kubát
Vladimír Matlach



Department of General Linguistic, Palacký University, Czech Republic

Our aim is to provide a user-friendly tool of quantitative text analysis for researchers from various disciplines (linguistics, criticism, history, sociology, psychology, politics, biology, etc.). QUITA combines all important parts of any quantitative research: obtaining results, statistical testing and graphical visualization. There is no need to use any additional software such as spreadsheet applications or special statistical programs.



STATISTICAL COMPARISON

Project	Tools	Settings
MTP-Multiple Text Project 1		
Project Settings Results Comparison: R1		
[czech] anglicke_jazy	[czech] fobyn	[czech] horadubal
[czech] fobyn	[czech] horadubal	[czech] krakati
[czech] horadubal	[czech] krakati	[czech] marketa_lazarova
[czech] krakati	[czech] marketa_lazarova	[czech] marketa
[czech] marketa	[czech] marketa_lazarova	[czech] obycejny_zivot
[czech] obycejny_zivot	[czech] pekar_jan	[czech] pekar_jan_mashout
[czech] pekar_jan	[czech] pekar_jan_mashout	[czech] pole_orna_a_valecna
[czech] pekar_jan_mashout	[czech] pole_orna_a_valecna	[czech] povetren
[czech] pole_orna_a_valecna	[czech] povetren	[czech] utek_do_budna
[czech] povetren	[czech] utek_do_budna	
[czech] utek_do_budna		
Average		
Standard deviation		

INDICATORS TO COMPUTE

Frequency Structure indicators

- Type-Token Ratio (TTR)
- h -point (h)
- Vocabulary Richness (R_v)
- Repeat Rate (RR)
- Relative Repeat Rate of McIntosh (RR_{mc})
- Hapax Legomenon Percentage (HL)
- Lambda (Λ)
- Gini Coefficient (G)
- Vocabulary Richness (R_4)
- Curve length (L)
- Curve length Indicator (R)
- Entropy (H)
- Adjusted Modulus (A)

Miscellaneous indicators

- Verb Distances (VD)
- Activity (Q) & Descriptivity (D)
- Writer's View (α)
- Average Tokens length (ATL)
- Thematic Concentration (TC)
- Secondary Thematic Concentration (STC)

TEXT-PROCESSING

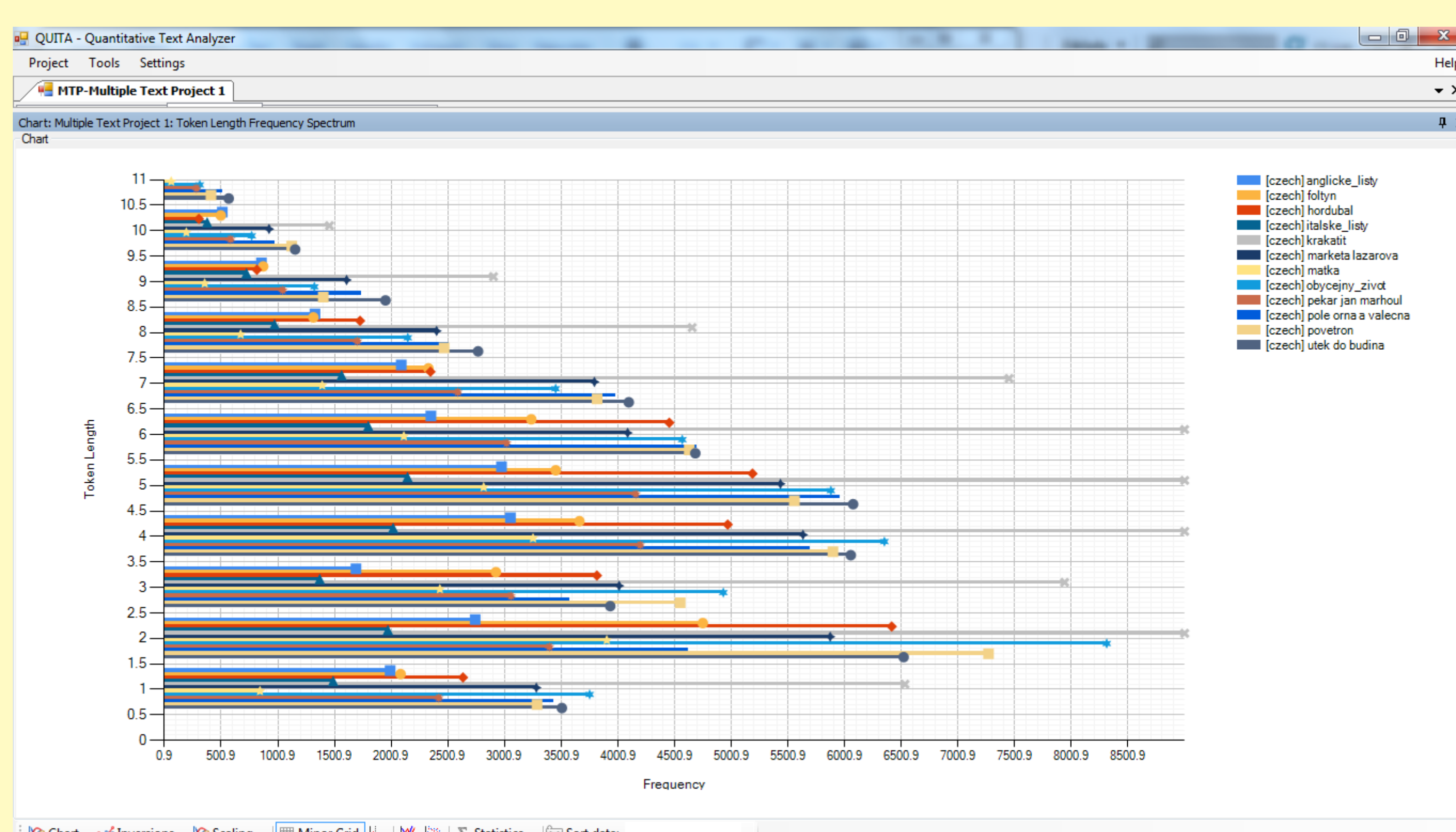
Pre-processing

- Tokenizer (word, line, char, DNA Triplet, DNA Nucleotide)
- Multilingual lemmatizer (AR, CZ, DE, DK, EN, ES, FI, FR, IT, NL, PT, RO, RU, SE)
- POS Tagger (It distinguishes parts of speech in a text)

Post-processing

- N-grams (QUITA enables creating char, word or whatever n-grams)
- Text length reduction

CREATING CHARTS



oltk.upol.cz/software

Acknowledgement: QUITA was supported by the student project IGA (no. FF_2013_031) of the Palacký University, Olomouc.