RELISH Symposium
*Rendering Endangered Lexicons Interoperable through Standards Harmonization: RELISH meets LOEWE*
Frankfurt/Main, October 10, 2011

GOETHE
UNIVERSITÄT
FRANKFURT AM MAIN

# SLIEKKAS | KALT

**S**enosios **LIE**tuvių **K**albos **K**orpus**AS**   **K**orpus **A**lt**L**i**T**auisch

Jolanta Gelumbeckaitė
Institute of Empirical Linguistics
gelumbeckaite@em.uni-frankfurt.de

# Aim of SLIEKKAS

➢ to create a qualitative multilevel electronic retrieval engine for multilateral linguistic research of Old Lithuanian,

➢ to allow for reliable results of diachronic Lithuanian language studies,

➢ to enable the implementation of the two biggest desiderata of Baltic linguistics, the **Old Lithuanian grammar** and the **historic dictionary of Lithuanian**.
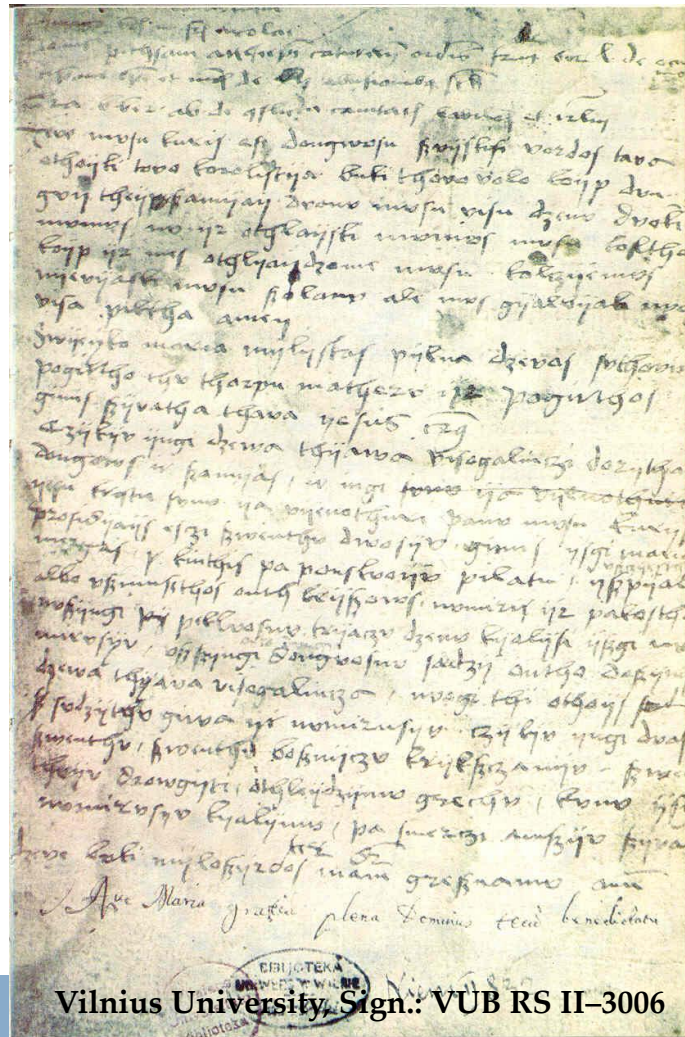
### What is SLIEKKAS doing?

✓ Joining the current process of contemporary historical corpora, and not copying them or reinventing the wheel.

✓ Consolidating already existing databases of Old Lithuanian and using already existing resources optimally.

✓ Developing a concept of the Old Lithuanian reference corpus as well as its scientific and technical base.
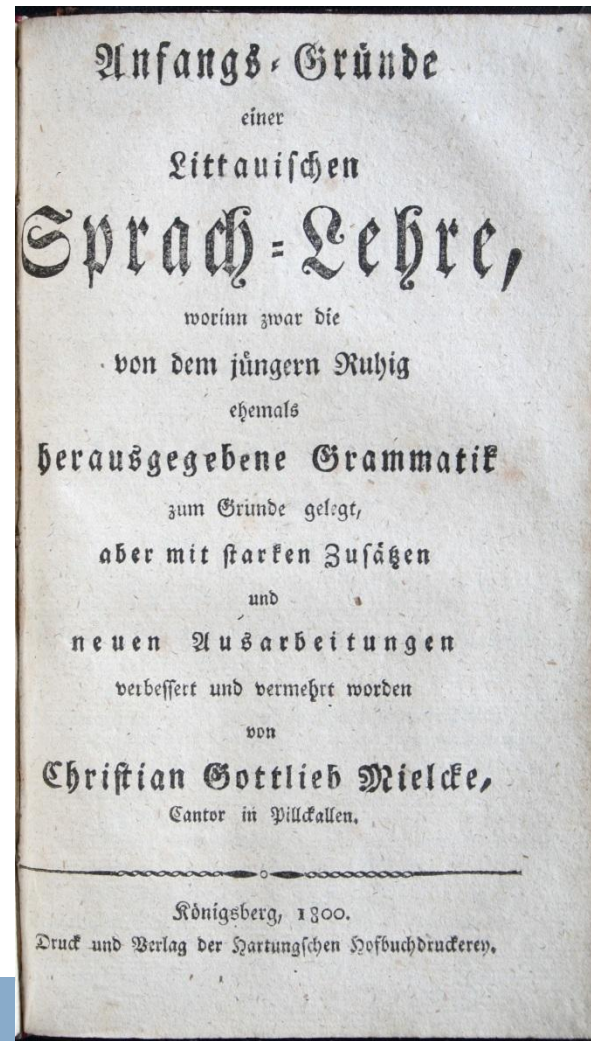
# Arrangement of the Corpus

o Digitisation according to the principles of the diplomatic literal transcription and conversion to the XML format;

o Multi-layer-stand-off annotation:

- Header-Informations,

- Structural, palaeographic, and textological annotations,

- Lemmatising of word forms (according to the language period), glossing in Standard Lithuanian, and translating into Modern Lithuanian,

- POS-Tagging. Parts of speech and grammatical forms (morphology and in certain cases syntactic morphology);

o Alignment of the Lithuanian texts with the digital facsimiles (autograph, print) and with such parallel corpora as with the annotated Latin, German, Polish etc. source texts.

# Old Lithuanian (1500–1800), ca. 10 million words

The oldest known Lithuanian text:
*Pater noster, Ave Maria* and *Credo*
in the *Tractatus sacerdotalis* (Straßburg: Martin Flach, 1503) by Nicolaus de Blony

Lithuanian Grammar by
Christian Gottlieb Mielcke (1732–1807)
*Anfangs=Gründe einer Littauischen Sprach=Lehre* (Königsberg, 1800)

Anfangs-Gründe
einer
Littauischen
Sprach=Lehre,
worinn zwar die
von dem jüngern Ruhig
ehemals
herausgegebene Grammatik
zum Grunde gelegt,
aber mit starken Zusätzen
und
neuen Ausarbeitungen
verbessert und vermehrt worden
von
Christian Gottlieb Mielcke,
Cantor in Pillckallen.

Königsberg, 1800.
Druck und Verlag der Hartungschen Hofbuchdruckerey.

# Referenzkorpus Altdeutsch

Scientifically as well as technologically SLIEKKAS leans onto the Old German Reference Corpus *DDD – Referenzkorpus Altdeutsch (750–1050)*

# Senieji raštai

## www.lki.lt/seniejirastai

Institute of Lithuanian Language, Vilnius. The most comprehensive and constantly increasing database of Old Lithuanian Writings. The digital versions of the texts are prepared according to the principles of the diplomatic literal transcription. Out of 72 texts (1573–1816, over 3 m. text words) 40 texts along with their KWIC concordances are currently accessible to external users online and for downloading (*Word-doc* format).

Deficiencies:

- no special search engine,
- no clear solution for the lemmatising and the grammatical tagging.

WP 1r,6     Ir kad prifſiartinaĳa meſtap ‹←meſtep› Je=
WP 1r,7     ruſalemJr ‹←Je=ruſalenn› atteĳa kiemap Beth= φ.

WP 1r,8     Schas dienas Egłai malanauſiei ‹←malanauſe› krikßa=
WP 1r,9     nis ir krikßankas bralei ir ſeſeris ing Cħun Je=
WP 1r,10    ſų turreſime tris wetas ‹←wetos›.

**A**

**ak** (1) interj.: *Ak* 16₁₄

**akis** (1) sf.: gen. pl. *Akiu* 14₄

**akmenėlis** (2) sm. demin.: acc. sg. *Akmeneli* 5₁₀ * gen. pl. *Akmenelů* 9₂₃

**akmuo** (1) sm.: gen. sg. *akmens* 21₆

**ale** 'bet' (2) conj.: *ałe* 5₁₆ 11₁₆

**anas** (1) pron. demonstr.: nom. pl. *anie* 12₁₄

The Lithuanian material consists of ca. 200.000 words.

Every character (resp. every word) is aligned with the index databank. This makes it possible to register all parallel positions of the words in the corpus as well as with the facsimile of the autograph (folio to folio so far).

Single corpora as well as the reference corpus can be searched by using special database query forms.

The texts are annotated merely with regard to the lemmatisation.

Deficiency:

– no grammatical annotations.



02.11.2011

# ePaveldas
## www.epaveldas.lt

# Vilnius University Library
## www.mb.vu.lt/sk-kolekcijos

GOETHE
UNIVERSITÄT
FRANKFURT AM MAIN

The databank of the virtual cultural heritage system contains 314 digitised facsimiles of handwritten and printed Lithuanian texts from the period 1547–1863.

VUL digital colletions comprise among others unique copies of Old Lithuanian grammars, dictionaries, and catechisms.

Digital facsimiles in the databases of such single libraries as the Duke August Library in Wolfenbüttel.
This library both preserves the original and places at the disposal a digital facsimile of the oldest Lithuanian Codex from 1573.

1.  3 prayers (*Pater noster, Ave Maria, Credo*), ca. 1520, 168 text words.

2.  Martynas Mažvydas, *Catechismusa prasty žadei* (*MžK* 1547), 7 000 text words.

3.  Martynas Mažvydas, *Giesme S. Ambrasziejaus* (*MžGA* 1549), 2 418 text words.

4.  Wolfenbüttel Postilla (*WP* 1573), 130 559 text words.

5.  Baltramiejus Vilentas, *Enchiridion* (*VE* 1579), 10 199 text words.

6.  Mikalojus Daukša, *Kathechismas* (*DK* 1595), 13 702 text words.

7.  Heinrich Johann Lysius, *Mažas Katgismas* (*LyK* 1719), 6 136 text words.

8.  Gabrielis Engelis, *Mažas Katgismas* (*EnK* 1722), ca. 10 000 text words.

9.  Kristijonas Donelaitis, *Metai* (ca. 1765–1775), 17 291 text words.

# Grammatical annotation

- *communis opinio* – any grammatical tagging of a text is a "perilous activity, because the text thereby loses [its] integrity" and becomes contaminated through various views of different linguistic schools (John M. Sinclair, *Current Issues in Corpus Linguistics*, 2004);

- no in any way grammatically annotated contemporary Lithuanian text-corpus;

- no standards for tagging of the Lithuanian texts;

- Lithuanian corpora as merely raw text accumulations;

- The Contemporary Lithuanian Corpus (donelaitis.vdu.lt), over 100 million words.

# Grammatical annotation of SLIEKKAS

o   predominantly restricted to morphology,

o   POS-Tagging and the morphological description of single word forms,

o   hierarchical: starting with the grammatical class of the lemma and going to individual unchangeable, and changeable morphological categories of the lemma as well as of the individual word form attested.

–   The difficulties begin in the establishing of the parts of speech already ("There are ten parts of speech, and they are all troublesome", Mark Twain, *The Awful German language* ). The Lithuanian grammar traditionally distinguishes 11 parts of speech, since the beginning of the 20th century setting apart interjections and onomatopoeica. The POS of the SLIEKKAS consits although of 10 more or less "troublesome" grammatical classes:

1. **ADJ** – adjective,

2. **ADV** – adverb,

3. **AP** – adposition,

4. **CARD** – cardinal number, **ADJO** ordinal number,

5. **ITJ** – interjection + onomatopoeic,

6. **KO** – conjunction,

7. **N** –noun: **NA** common noun, **NT** proper noun,

8. **P** – pronoun: **PD** demonstrative, **PI** indefinite, **PK** interrogative, **PPER** personal,

9. **PTK** – particle,

10. **V** – verb: **VA** auxiliary, **VV** main.

The grammatical class is indicated according to the lemma and to the word form attested. This distinction enables us to indicate:

- changes in the grammatical classes (nominalisation of the adjectives, adjectivisation of the participles, adverbialisation of the nouns, adjectives or participles, and turning of some nouns into adpositions),

- ambivalent issues as adverbialisation of the oblique case forms:

| atėjo | vidunakčiu |
|---|---|
| *come_Pret_3* | *midnight_Sg_Inst* |
| | lemma: NA |
| | attested form: NAA (common noun, adverbialised) |

*(he, she, it, they) came in the midnight*

- distinction of prepositions and postpositions within the class of the adposition,

- a clear definition of such non-finite forms of the verb as infinitives, participles, semiparticiples, gerunds, gerundives, supines, and adverbialised infinitives.

# Morphology and syntactic morphology (1)

- analytic (compound) tense forms, consisting of a copula and a participial (present or past, active or passive) predicative:

  - distinguishing between an auxiliary and a main verb already in the definition of the grammatical class of the lemma,

  - between predicative, attributive or adverbial usage of the participles in the characterisation of the word forms attested.

- compound nominal predicates, consisting of a copula and of a nominal predicative:

| Jonas | buvo | pranašas | / | pranašu |
|---|---|---|---|---|
| *John*_Sg_Nom | *be*_Pret_3 | *prophet*_Sg_**Nom** | / | *prophet*_Sg_**Inst** |
| | | lemma: NA | | |
| | | attested form: NAP (common noun, predicative) | | |

*John was a prophet*

# Morphology and syntactic morphology (2)

- secondary predicatives, expressed by nouns, adjectives, pronouns, numerals or participles:

| gandras | parlėkė | linksmas |
|---|---|---|
| stork_Sg_Nom | fly back home_Pret_3 | cheerful_Masc_Sg_Nom |
| | | lemma: ADJ |
| | | attested form: ADJP (adjective, predicative) |

*The stork came back home being cheerful*

- ✓ cf. the adverb *linksmai* in this position:

| gandras | parlėkė | linksmai |
|---|---|---|
| stork_Sg_Nom | fly back home_Pret_3 | cheerfully |

*The stork came back home cheerfully*

# Morphological annotation

The further morphological annotation is also hierarchical. It consists of 3 layers, which include the information about:

1) the unchangeable morphological categories of the lemma (like gender and flexional class),

2) the unchangeable morphological categories of the word form attested: they can differ from these of the lemma,

3) the flexional morphological categories or the word form attested.

# Data analysis tools (Toolbox)

There are no specialised diachronic dictionaries of Old Lithuanian, which we could apply for the (semi)automatisation of the annotation process. Therefore we have to either perform the morphological annotation manually or to choose such data analysis tools like the Toolbox, which enables semiautomatisations of interlinear annotations.

Through recognising and defining the problems we search for the best solutions. Going through the 'earthworm' limb by limb.

# Thank you for your attention!