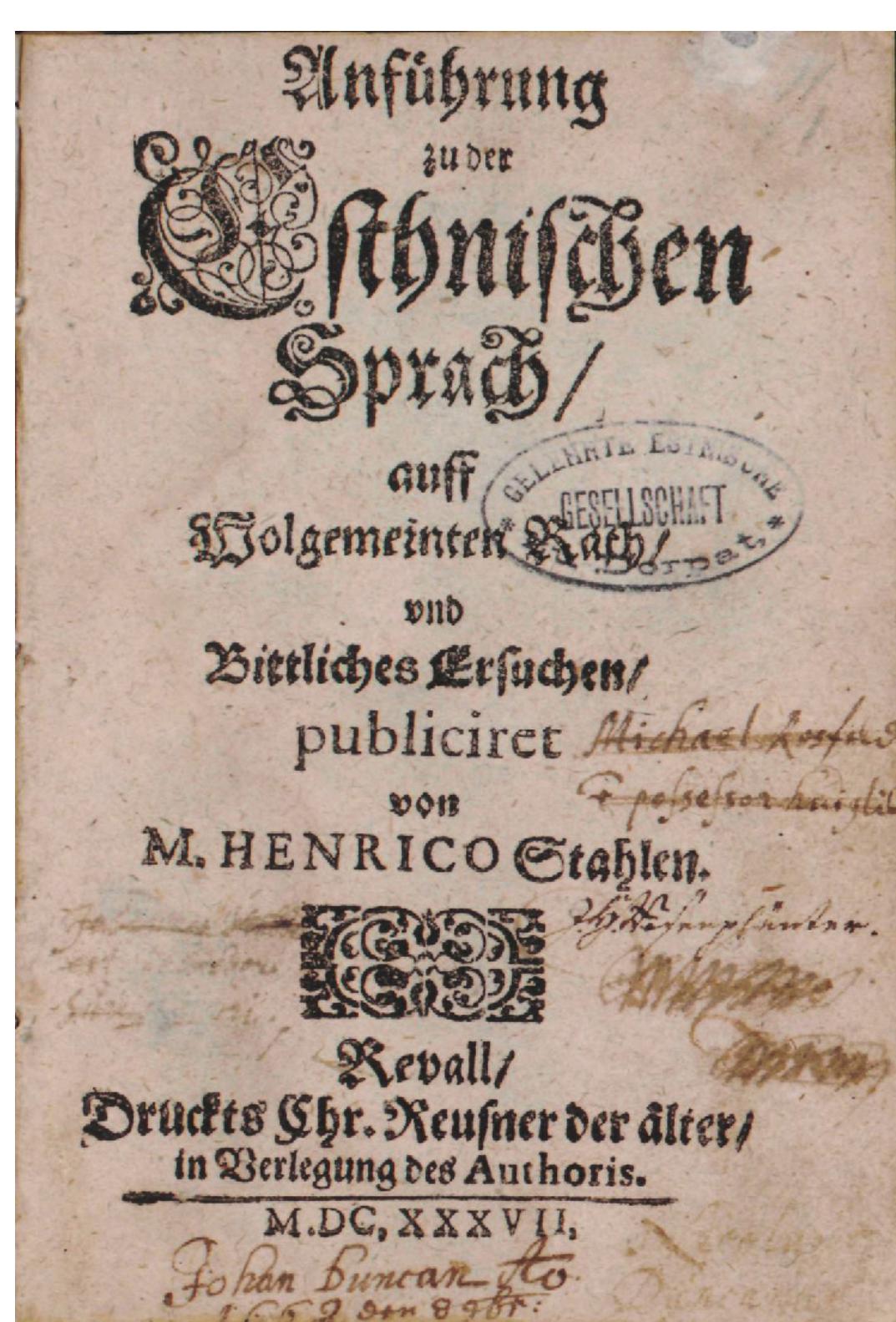


BRIDGING 17TH AND 18TH CENTURY ESTONIAN AND LLMs: UNLOCKING HISTORICAL DICTIONARIES

Madis Jürviste – lexicographer, junior researcher – Institute of the Estonian Language & University of Tartu
 Tiina Paet – researcher, Institute of the Estonian Language
 Sven-Erik Soosaar – senior researcher, Institute of the Estonian Language

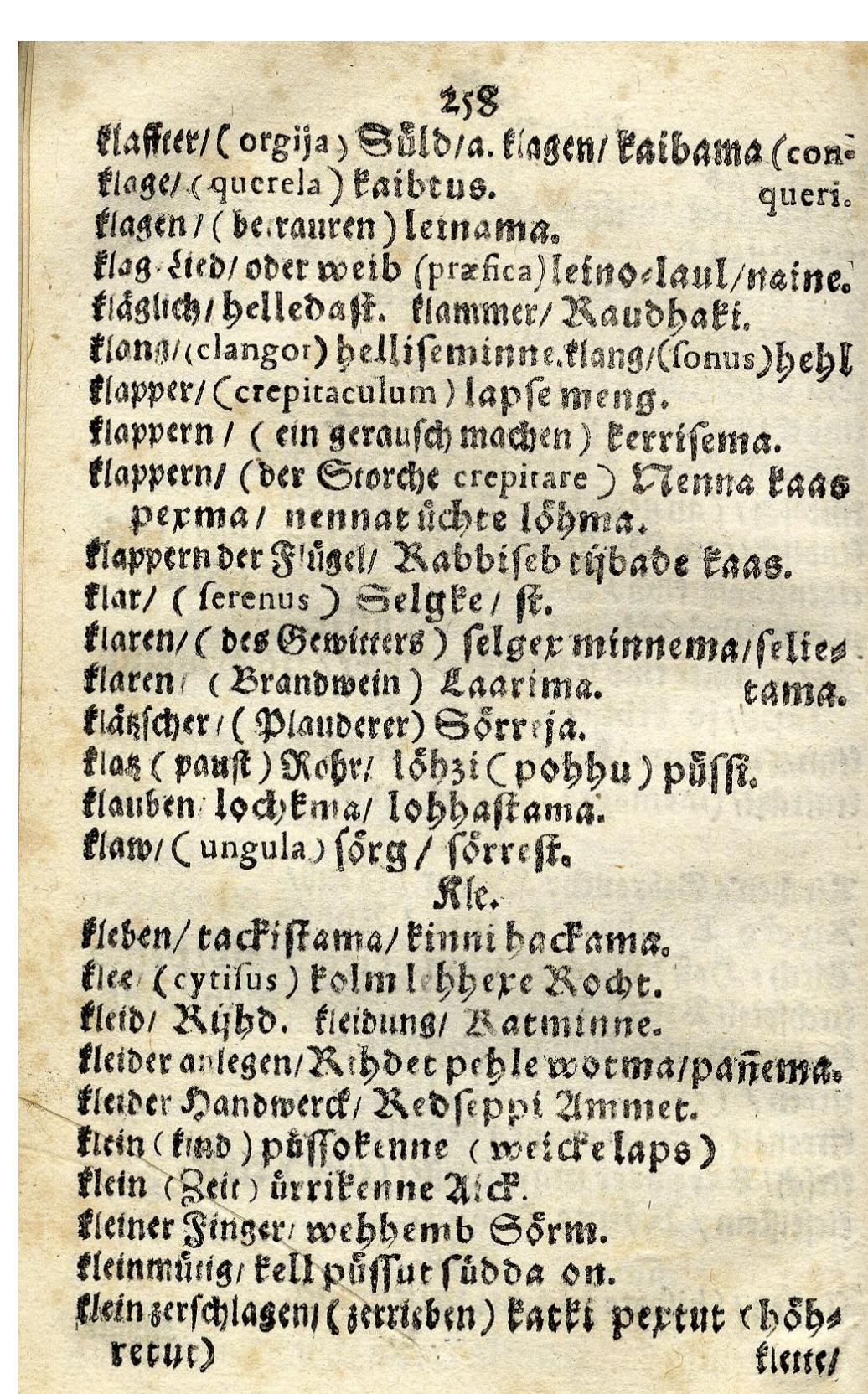
* This study was funded through the R&D project 'Applying large language models to lexicography: new opportunities and challenges' (EKKD-III1).



Heinrich Stahl, "Anführung zu der Esthnischen Sprach ..." (1637)

Rie 'cloth, garment' Ger: Kleid

Stahl 1637: **rihd**
 Göseken 1660: **Rijhd**
 Vestring ~1720: **Rie**
 Helle 1732: **rie**
 Hupel 1780: **rie**
 Hupel 1818: **riid, ride, rie**
 Wiedemann 1893: **rie**
 1918...2025: **riie**



Heinrich Göseken, "Manuductio ..." (1660), p. 258

Dictionaries studied

- Stahl, Heinrich 1637. Anführung zu der Esthnischen Sprach.
- Gutslaff, Johannes 1648. Observationes Grammaticae circa lingam Esthonicam.
- Göseken, Heinrich 1660. Manuductio ad Linguam Oesthonicam. Anführung zur Öhstnischen Sprache.
- Vestring, Salomo Heinrich 2000 [~1720]. Lexicon Esthonica Germanicum.
- Helle, Anton Thor 1732. Kurtzgefazzte Anweisung Zur Ehstnischen Sprache.
- Hupel, August Wilhelm 1780. Ehstnische Sprachlehre für beide Hauptdialekte den revalschen und dörptschen; nebst einem vollständigen Wörterbuch.

Research

Overview and initial experiments

Our research explores the significant potential of Large Language Models (LLMs) in historical lexicography. The larger aim of the project is to create a historical layer using LLM technology (see example on the left, *rie*.)

Material

Six German-Estonian and Estonian-German dictionaries, spanning from 1637-1780.

LLMs tested

GPT-4o, Claude 3 Opus, Gemini 1.5 Pro.

Aim

Explore the potential use of LLMs in studying 17th and 18th century lexicography, as well as finding the most suitable model for analysing historical Estonian lexical data.

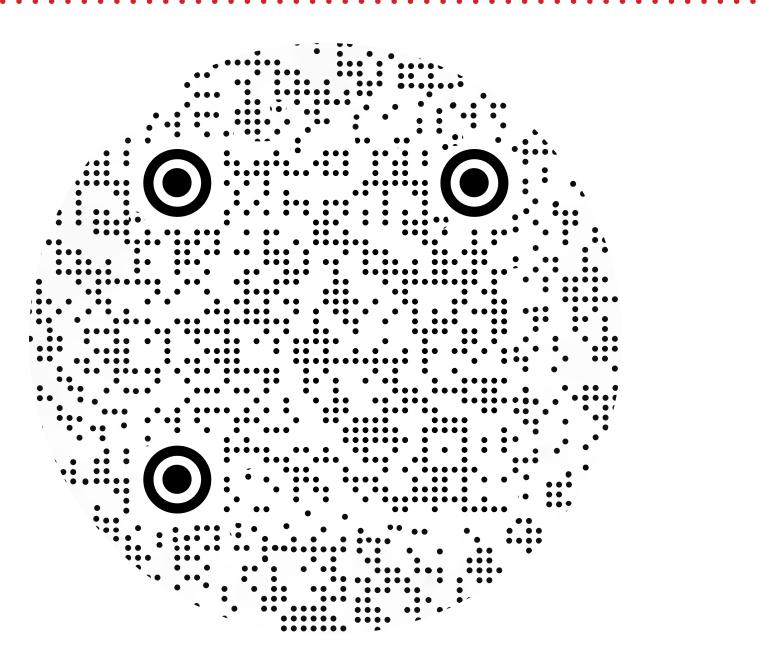
Experiments

- Experiment 1: 30 professional titles (*kohtunik / sundija* 'judge', *varas / kühnames* 'thief') and societal roles across all the six sources.
- Experiment 2: 54 words from Gutslaff (1648) (South Estonian). Three categories tested: identical forms (*kirp* 'flea'), minor orthographic variations (*Hawd / haud* 'grave'), completely different lexemes (*ramato nachk* lit. 'book skin' / *paber* 'paper').
- Experiment 3: 20 loanwords without German equivalents across 3 sources. Three categories tested: minimal differences (*pipar / pippar* / *pepper* 'pepper'); somewhat different forms with overlapping letters (*äädikas / Ettickas* 'vinegar'); radical divergences (*väävel / koirasitt* 'sulfur').

Success rates

- E1: 90% accuracy (Claude 3 Opus).
- E2: 76-81% accuracy (Gemini 1.5 Pro, Claude 3 Opus).
- E3: 60-90% accuracy rates (Claude 3 Opus).

References ->



Case-study: Gutslaff 1648. Enriching historical dictionaries with modern forms and meanings

Challenge and approach: Historical dictionaries contain entries (*karrakaro, padia, särge sap*) that are cryptic for modern readers, especially non-linguists. Envisaged solution: provide modern word forms and meanings for historical lexicographical data.

- Aim:** Establishing replicable workflows applicable to other historical dictionaries and language pairs (Gutslaff 1648 has undergone manual analysis already: Lepajõe 1998, Viitkar 2005).

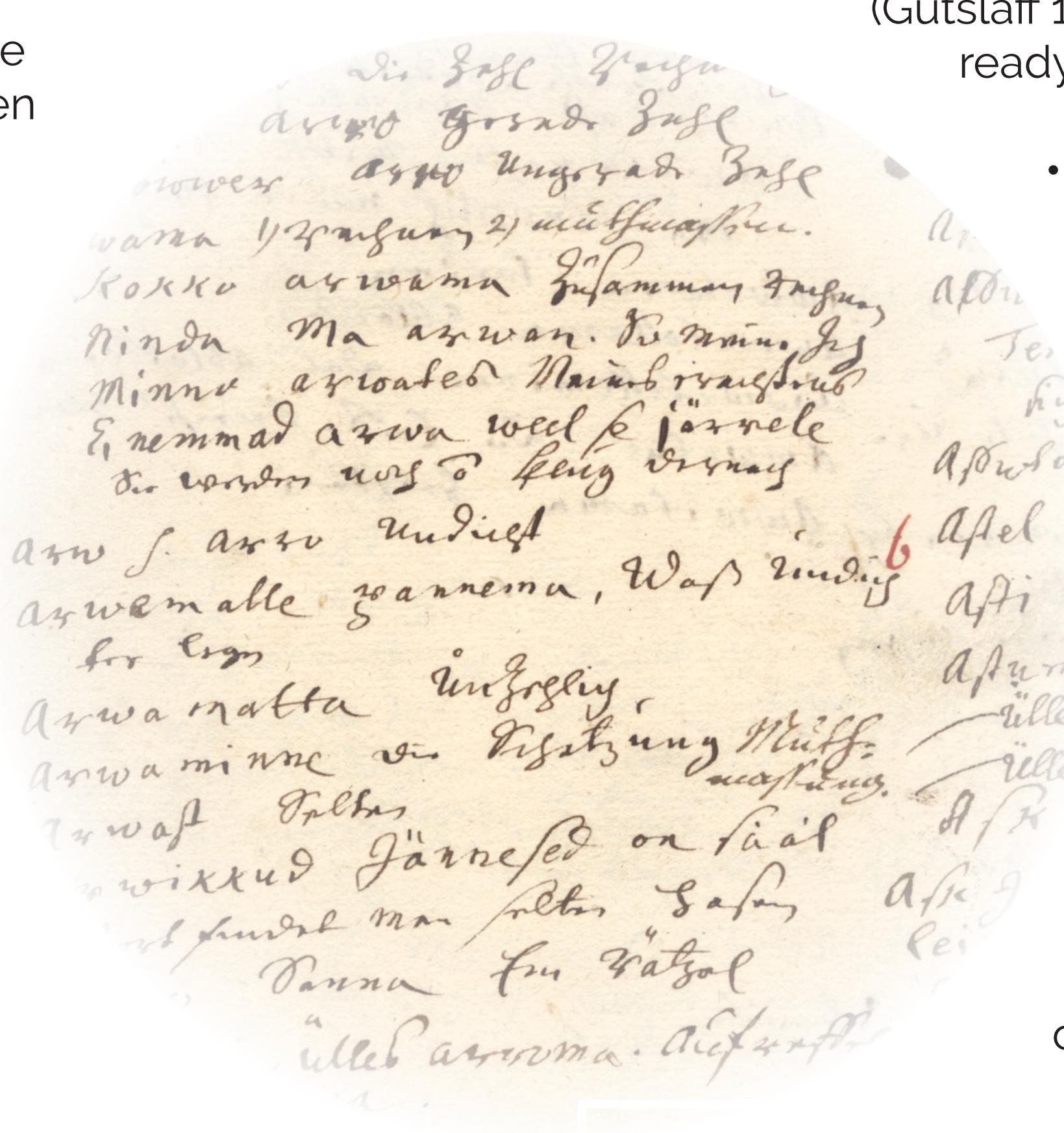
- Example :** [et] Ahwoneck – [de] Barß > mod. ahven, 'perch'.

Method: multi-stage LLM-assisted analysis (cross-source mapping), followed by expert (human) validation and oversight.

Tested LLM: Claude 4 Opus (API).

Test sample: 342 entries out of 1714 (20%).

Success rate: 81% of entries (278 out of 342) required no correction. 11% needed additional review. 8% of entries required complete manual override.



```
"Püsse", "püss", "püss", "Büchse", "Büchse", "Püss" on püss"
"Püssohand", "pisuhänd", "kratt", "koletis", "Drache", "Püssohand" on pisuhänd, kratt"
"Putsa", "sulg", "sulg", "Feder", "Feder", "Putsa" on sulg"
"Putt", "puder", "puder", "Brey", "Brey", "Putt" on puder"
"Raip", "raibe", "surukeha", "Aass", "Aas", "Raip" on raibe"
"Ramat", "raamat", "kiri", "Brief", "Brief", "Ramat" on raamat kui kiri"
"Ramat", "raamat", "Buch", "Buch", "Ramat" on raamat"
"Ramm", "surupaarid", "kanderam", "Baare", "Bahre", "Ramm" on surupaarid"
"Rasa", "raas", "Leivatükk", "Brocke Brods", "Brotstück", "Rasa" on raas, leivatükk"
"Raw", "rasv", "rasv", "Fett", "Fett", "Raw" on rasv"
"Reiss", "libejää", "sile jää", "glat Eis", "glattes Eis", "Reiss" on libejää"
"Rind", "ind", "ind", "Brust", "Brust", "Rind" on rind"
"Risti", "rist", "rist", "Kreutz", "Kreuz", "Risti" on rist"
"Risti inniminni", "ristinimene", "Kristlane", "Christ", "Risti inniminni" on ristiinnimene"
"röhmus", "röhmus", "röhmus", "annutig", "annutig", "röhmus" on röhmus"
"röhckma", "ruttama", "kiirustama", "jövödma", "eilen", "röhckma" on ruttama"
"rummal", "rumal", "rumal", "dumm", "dumm", "rummal" on rumal"
"Russick", "rusik", "rusik", "Faust", "Faust", "Russick" on rusik"
"saddama", "sadama", "sadama", "langem", "fallen", "fallen/regnen", "saddama" on sadama"
"salvama", "salvama", "hammustama", "beissen", "salvama" on salvama"
"Sanna", "saun", "saun", "Badstube", "Badstube", "Sanna" on saun"
"Sant", "sant", "vaene", "Bettler", "arn", "Sant" on sant, vaene"
"Sarah", "sarapuu", "sarapuu", "Aschenbau(m)", "Haseistrach", "Sarah" on sarapuu"
"Savissiwick", "vaskuss", "vaskuss", "Blindschleich", "Blindschleiche", "Savissiwick" on vaskuss"
"selge", "selge", "selge", "Clarh", "Klar", "selge" on selge"
"Seng", "säng", "voori", "Bette", "Bett", "Seng" on säng, voori"
"sick", "sikk", "sokk", "Bock", "Bock/Ziegenbock", "sick" on sikk"
```

Excerpt from an API-request (CSV-output): enriching Gutslaff's 1648 dictionary with modern word forms and meanings. Data points:
 >>> [original et], [modern et], [meaning], [original del], [modern del], [LLM comment]

Conclusions

- LLMs are capable of historical lexicographical analysis.
- Optimal resource efficiency in lexicographical analysis.
- Cross-linguistic potential: adaptability to other language pairs.
- Scalability: Gutslaff 1648 enrichment experiment is applicable to other historical dictionaries.
- Expert validation is essential: LLM analysis combined with human expert validation (manual review) yields best results.

Contacts

- madis.juryviste@eki.ee
- tiina.paet@eki.ee
- sveneri.koosaar@eki.ee

