**Lorentz Workshop:**

**Processing Ancient Text Corpora**

**17-21 February 2020**

**Part A: Corpora Primer**

# Syriac texts

*Wido van Peursen*

*Eep Talstra Centre for Bible and Computer*

**www.etcbc.nl**   **/etcbc**   **@etcbc_vu**   **company/etcbc/**

# Workshop Overview

1. **Corpora**
   A. Primer
   B. In-depth discussion

2. **Data structure and text model**

3. **Corpus analysis**

4. **Making knowledge accessible**

ETCBC

# Part A: Corpora Primer

"In the primer (Part A) various researchers will present a snapshot of their corpora, addressing questions such as: What challenges do scholars meet when dealing with the various scripts, languages and text types? How does text model and data structure interact with the scholarly traditions and conventions in the specific disciplines. How do we deal with multiple witnesses of the same texts?"

ETCBC

# Why study Syriac texts?

Syriac is the language of a long literary tradition including poems, ancient Bible translations, apocryphal stories, historiography, Bible commentaries, philosophical tractates, scientific works. It is the language of a neglected branch of early Christianity, and one of the most important languages to study the religious and cultural context in which Islam emerged. Syriac translations served as a bridge between Greeks and Arabs, and Syriac is still the liturgical language of many Christians in the Middle East and, increasingly, in the diaspora.

</> ETCBC

# Why study Syriac texts?

The Aramaic languages, in which Syriac takes pride of place, cover a period from the 10[th] cent. BC up to the present, and hence is one of the most well documented languages that allows for longitudinal linguistic analysis over three millennia.

ETCBC

# Annotated corpora

# Annotated corpora

- **Comprehensive Aramaic Lexicon (CAL)**
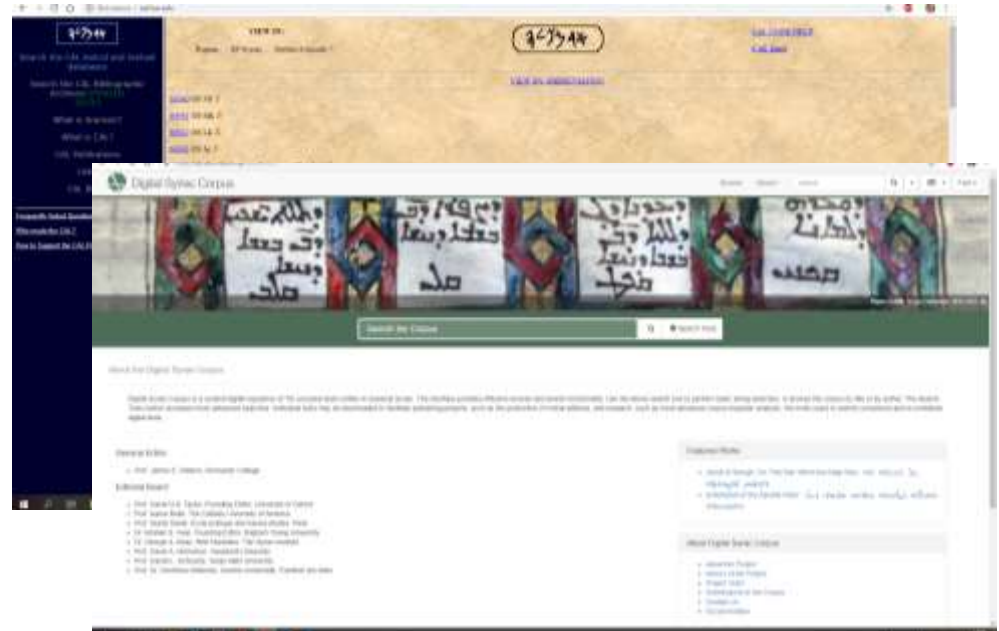- **Digital Syriac Corpus**
- **Text-Fabric (TF)**

ETCBC

# Annotated corpora

- **Comprehensive Aramaic Lexicon (CAL)**
- **Digital Syriac Corpus**
- **Text-Fabric (TF)**

# Annotated corpora

- **Comprehensive Aramaic Lexicon (CAL)**
- **Digital Syriac Corpus**
- **Text-Fabric (TF)**



ETCBC

# Annotated corpora

- **Comprehensive Aramaic Lexicon (CAL)**
- **Digital Syriac Corpus**
- **Text-Fabric (TF)**





</ETCBC

# Annotated corpora

- **Comprehensive Aramaic Lexicon (CAL)**

- **Digital Syriac Corpus**

- **Text-Fabric (TF)**

- **Parsing!**

# Structured data


ETCBC

# Structured data

- **Lexicographical resources**
- **Grammatical paradigms**
- **Liturgical data**
- **Geographical data**

ETCBC

# Structured data

- **Lexicographical resources**
- **Grammatical paradigms**
- **Liturgical data**
- **Geographical data**

# Structured data

- **Lexicographical resources**
- **Grammatical paradigms**
- **Liturgical data**
- **Geographical data**



```
ps: "person" =
   first: "first", second: "second",
   third: "third"
nu: "number" =
   sg: "singular", du: "dual",
   pl: "plural",
   unknown: "unknown"
gn: "gender" =
   f: "feminine", m: "masculine"
```

ETCBC

# Structured data

- **Lexicographical resources**
- **Grammatical paradigms**
- **Liturgical data**
- **Geographical data**

# Structured data

- **Lexicographical resources**
- **Grammatical paradigms**
- **Liturgical data**
- **Geographical data**

# Structured data

- **Lexicographical resources**

- **Grammatical paradigms**

- **Liturgical data**

- **Geographical data**

```
1    Ābājālūi (settlement)      [Syriac Not Available]  http://syriaca.org/place/881
2    Abba Hbisha (monastery)    [Syriac Not Available]  http://syriaca.org/place/2416
3    Abba Sahrowaĭ (monastery)          [Syriac Not Available]  http://syriaca.org/place/2417
4    Abba Ṣliba (monastery)     [Syriac Not Available]  http://syriaca.org/place/659
5    'Abdāsi (diocese)          [Syriac Not Available]  http://syriaca.org/place/660
6    Abdon (monastery)          [Syriac Not Available]  http://syriaca.org/place/2418
7    'Ābdûlâkandi (settlement)          [Syriac Not Available]  http://syriaca.org/place/934
8    Abnaye (diocese)           [Syriac Not Available]  http://syriaca.org/place/2276
9    Abnaye (settlement)        [Syriac Not Available]  http://syriaca.org/place/882
10   Abrō (settlement)          [Syriac Not Available]  http://syriaca.org/place/883
11   Acre (settlement)          [Syriac Not Available]  http://syriaca.org/place/14
12   'Ada (diocese)  [Syriac Not Available]  http://syriaca.org/place/2260
13   'Ādā (settlement)          [Syriac Not Available]  http://syriaca.org/place/935
14   Adana (settlement)         ܐܕܢܐ   http://syriaca.org/place/15
15   Adana (diocese) [Syriac Not Available]  http://syriaca.org/place/2277
16   'Adasiyya (settlement)     ܚܕܨܝܐ  http://syriaca.org/place/1492
17   Adeḥ (settlement)          [Syriac Not Available]  http://syriaca.org/place/884
18   Adharbayjān (region)       ܐܕܪܒܝܓܢ  http://syriaca.org/place/5
19   Adharbayjān (diocese)      [Syriac Not Available]  http://syriaca.org/place/2278
20   Adiabene (region)          [Syriac Not Available]  http://syriaca.org/place/993
21   Aegeae (settlement)        [Syriac Not Available]  http://syriaca.org/place/16
22   Aegyptus (province)        ܡܨܪܝܢ  http://syriaca.org/place/7
23   Afghanistan (region)       ܐܦܓܢܣܛܢ http://syriaca.org/place/1523
24   Aghjachă (settlement)      [Syriac Not Available]  http://syriaca.org/place/885
25   Aḥasîm (settlement)        [Syriac Not Available]  http://syriaca.org/place/886
26   al-Aḥmadiyya (settlement)          ܐܚܡܕܝܗ http://syriaca.org/place/1481
```

ETCBC

# Structured data

- **Lexicographical resources**

- **Grammatical paradigms**

- **Liturgical data**

- **Geographical data**



| | | | |
|---|---|---|---|
| 1 | Ābājālūi (settlement) | [Syriac Not Available] | http://syriaca.org/place/881 |
| 2 | Abba Hbisha (monastery) | [Syriac Not Available] | http://syriaca.org/place/2416 |
| 3 | Abba Sahrowaï (monastery) | [Syriac Not Available] | http://syriaca.org/place/2417 |
| 4 | Abba Şliba (monastery) | [Syriac Not Available] | http://syriaca.org/place/659 |
| 5 | 'Abdāsi (diocese) | [Syriac Not Available] | http://syriaca.org/place/660 |
| 6 | Abdon (monastery) | [Syriac Not Available] | http://syriaca.org/place/2418 |
| 7 | 'Ābdûlâkandi (settlement) | [Syriac Not Available] | http://syriaca.org/place/934 |
| 8 | Abnaye (diocese) | [Syriac Not Available] | http://syriaca.org/place/2276 |
| 9 | Abnaye (settlement) | [Syriac Not Available] | http://syriaca.org/place/882 |
| 10 | Abrö (settlement) | [Syriac Not Available] | http://syriaca.org/place/883 |
| 11 | Acre (settlement) | [Syriac Not Available] | http://syriaca.org/place/14 |
| 12 | 'Ada (diocese) | [Syriac Not Available] | http://syriaca.org/place/2260 |
| 13 | 'Ādā (settlement) | [Syriac Not Available] | http://syriaca.org/place/935 |
| 14 | Adana (settlement) | ܐܕܢܐ | http://syriaca.org/place/15 |
| 15 | Adana (diocese) | [Syriac Not Available] | http://syriaca.org/place/2277 |
| 16 | 'Adasivva (settlement) | ܐܕܣܝܐ | http://svriaca.org/place/1492 |

# Structured data

- **Lexicographical resources**
- **Grammatical paradigms**
- **Liturgical data**
- **Geographical data**

- **Combining textual data and structured data**

ETCBC

# Structured data

- **Lexicographical resources**
- **Grammatical paradigms**
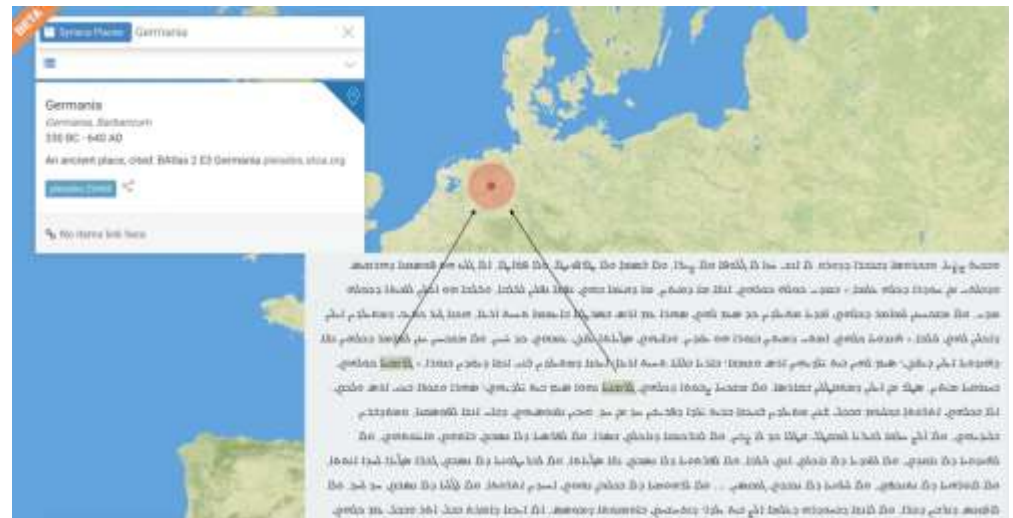- **Liturgical data**
- **Geographical data**



- **Combining textual data and structured data**

# Structured data

- **Lexicographical resources**

- **Grammatical paradigms**

- **Liturgical data**

- **Geographical data**

- **Combining textual data and structured data**

# Structured data

- **Lexicographi...
  resources**

- **Grammatica...**

- **Liturgical da...**

- **Geographica...**



- **Combining textual data
  and structured data**

- **Linked data**

ETCBC

# Text analytics

ETCBC

# Text analytics

- **Topic modelling**

- **Text clustering**

- **Etc, etc.**

# Text analytics

- **Topic modelling**
- **Text comparison**
- **Etc, etc.**

**1st CENTURY**

**2nd CENTURY**

Topics: Biblical heritage and Hellenistic culture

Topics: are Bardaisan's BLC and Ephrem's Refutations addressing same topics?

Bible

Bardaisan, Book of the Laws of the Countries

Ephrem the Syrian

**POLEMICS**

Lexemes
- Biblical vocabulary
- Greek loanwords
- Cosmological & astrological technical terms
- Named Entities

Lexicographical information SEDRA

LOD resources: syriac.org

Geography

eek Philosophy

# Desiderata

- **Annotated corpora**
  - **Syriac**
  - **Other forms of Aramaic, esp. Neo-Aramaic**
- **Analytical tools**
- **Human resources / Project funding**

ETCBC