# Working with Multimedia Data in CMC Corpora

## International Conference on CMC and Social Media Corpora for the Humanities

**Dr. Thomas Schmidt**

linguisticbits.de | Musical Bits GmbH

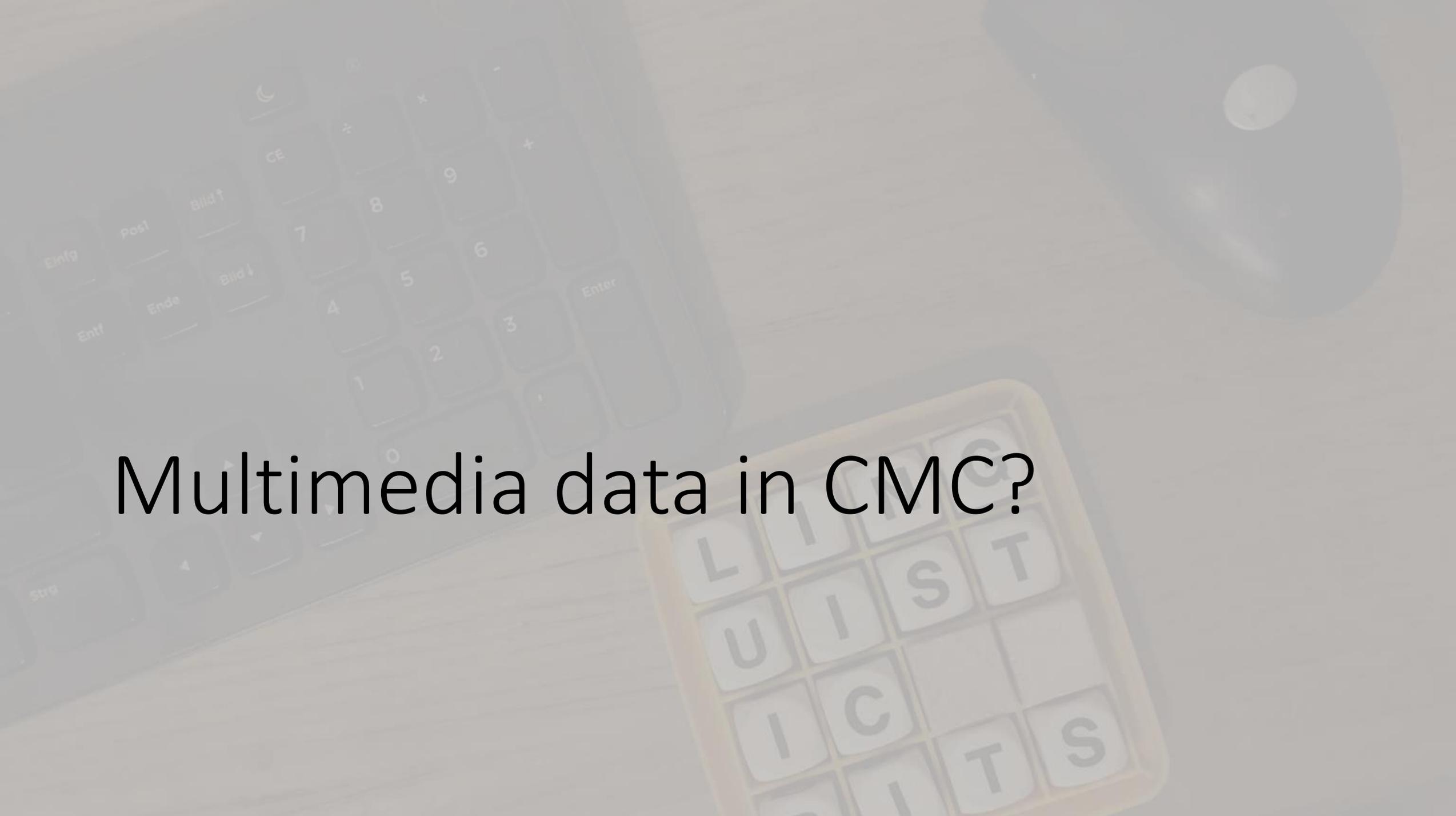thomas@linguisticbits.de

# Outline

1. Multimedia data in CMC?
2. Tools
   - Manual transcription and annotation
   - Automatic speech recognition
3. Formats, Standards, Interoperability
   - ISO/TEI standard
4. Q & A

# Multimedia data in CMC?

# Multimedia in CMC?

**Multimedia**
- Text
- Image
- Audio
- Video

**Multimodal**
- Written / Spoken / (Signed) – Alternative modes of language
- Speech + Gestures + Facial Expression (+ Body posture + …) – „Bodily communication"

**Multimodal CMC Corpora?**
- „Not-text" data
- Data from other modalities than writing
- Corpora taking these other media/modalities into account
- ➔ Represent (in writing) audio, video, (image) for corpus linguistic access
- ➔ Transcription

# Video conferencing

- Text (chat)

- Video incl. audio

- Speech in video

- Gesture + facial expression in video

- Emoticons in video

- Simultaneity of text and video

**Christopher** 08:45
Viel Erfolg bei den workshops

**Christopher** 15:39
Audio ▾

▶ ||||||||||||||||||||||||||||||||||||  0:34  ≡  1x  ⋮

| Hi. yeah, start Business Angels on V CS, et cetera. And, View transcript

**Transcript (auto-generated)** ✕
Christopher at 15:39

0:00  Hi.

0:05  yeah, start Business Angels on V CS, et cetera.

0:22  And,

**Thomas** 15:41
Nur LMU, von denen kommen einige Sachen, die ich jetzt auch nutze. Mit der TU hatte ich noch nix zu tun

# Slack channel

- Text

- Audio (voice messages)

- Derived text (auto-generated transcript)

- Alternation between text and audio

**Christopher** 08:45
Viel Erfolg bei den workshops

**Christopher** 15:39
Audio ▾

▶ ||||||||||||||||||||||||||||||||  0:34  ☰  1x  ⋮

Hi. yeah, start Business Angels on V CS, et cetera. And, View transcript

**Thomas** 15:41
Nur LMU, von denen kommen einige Sachen, die ich jetzt auch nutze. Mit der TU hatte ich noch nix zu tun

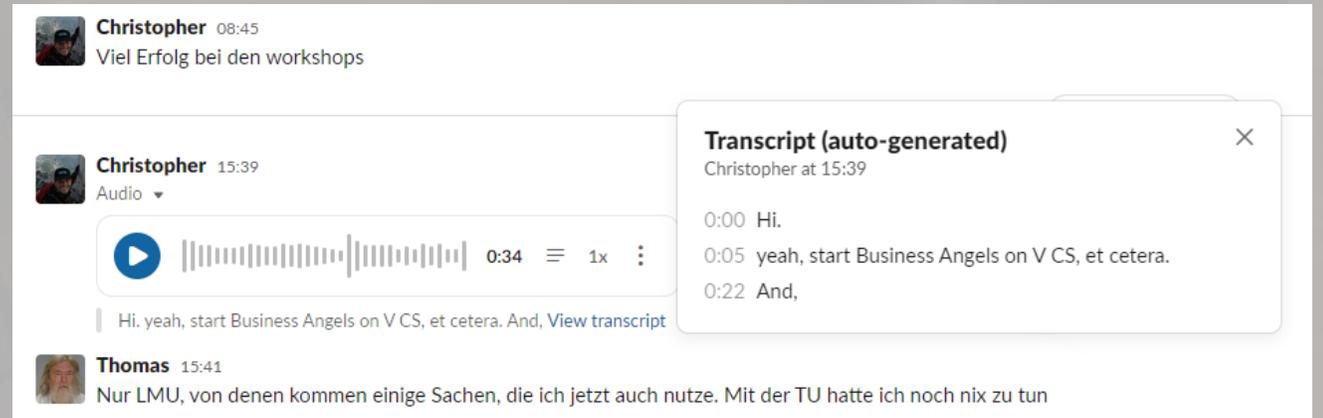**Transcript (auto-generated)**            ✕
Christopher at 15:39

0:00  Hi.
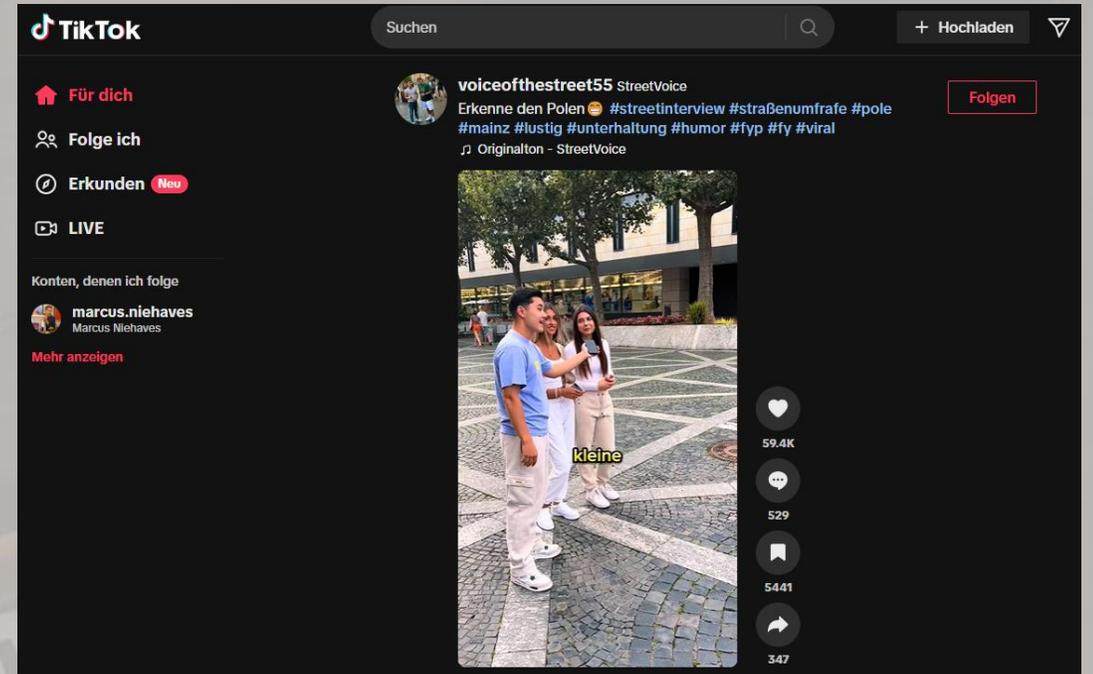0:05  yeah, start Business Angels on V CS, et cetera.
0:22  And,

# TikTok post

- Text

- Video

- Text refers to the video

# This conference

- Podcasts          [Babayode et al.]
- Audio/Video of Zoom and face-to-face-meetings (as comparison), interviews    [Steinsiek]
- Video-conferencing in Zoom          [Pabst et al.]
- Comments on Bilibili videos          [Zheng]
- Podcasts (vs. blog posts)    [Seemann et al.]
- Spoken corpora, gaze or walk annotations, kinect or motion capture data      [Ferger et al.]
- Online video film reviews  [Piroh]
- Short video data on TikTok          [Helenius]
- Videos of Authentic Social Interaction          [Krause et. al]
- Audio and video data from video sharing sites, streaming services and social media platforms [Coats]
- Multimodal WhatsApp discussion    [Mäkinen]

- Some mono-modal (podcasts)
- Audio/video of very different durations (seconds to hours)
- Different status of audio/video
- Comparison CMC
  <> Face-to-Face interaction

# Challenges

- Get your audio / video data transcribed and annotated

- Integrate it with text data

- Common basis for analysis

- In a FAIR-compliant manner (data sharing):
    - Interoperable (standardized, machine-readable formats)
    - Reusable (documented, conditions of use ➔ GDPR)

# Tools

**Sonntag**

Hello Schlotti, I wanna make an example for the CMC course. Can u hepl me? 😎 19:41 ✓✓

No idea what u r talking about but go ahead :p 19:42
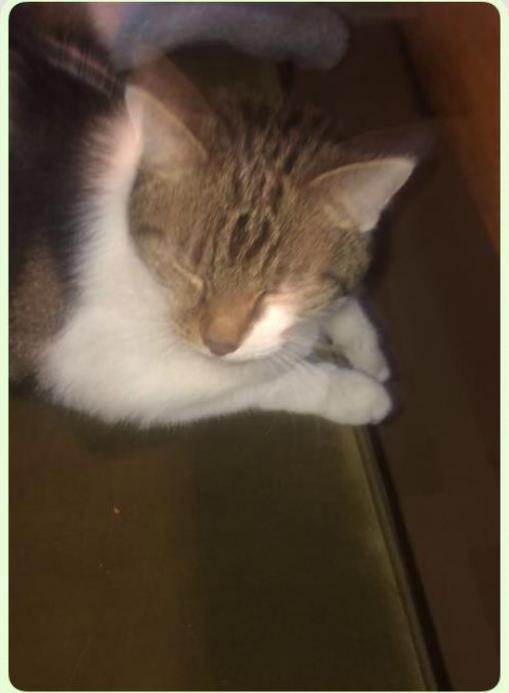
▶ ●━━━ 0:08 19:42 ✓✓

▶ ●━━━ 0:04 19:43

It's for the course. I said so!!! 19:43 ✓✓

▶ ●━━━ 0:06 19:44 ✓✓

I know I remember but I don't know what that is I SAID SO

know what that is I SAID SO 19:44
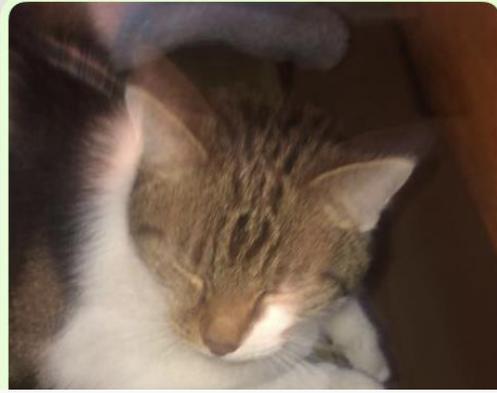


Thank u. Here is apicture of our cat. 19:45 ✓✓

▶ ●━━━ 0:03 19:45 ✓✓

See u 19:45

. Telekom.de 20:02  
< 11 Schlotti ███████  
zul. online heute um 19:30

. Telekom.de 20:03  
< 11 Schlotti ███████  
zul. online heute um 19:30

know what that is I SAID SO  
19:44

**Sonntag**

Hello Schlotti, I wanna make an example for the CMC course. Can u hepl me? 😎 19:41 ✓✓

No idea what u r talking about but go ahead :p 19:42



```
[11.09.23, 19:41:30] Thomas Schmidt: Hello Schlotti, I wanna make an example for the CMC course. Can u hepl me? 😎
[11.09.23, 19:42:40] Schlotti XXX: No idea what u r talking about but go ahead :p
[11.09.23, 19:42:53] Thomas Schmidt: <Anhang: 00000518-AUDIO-2023-09-11-19-42-53.opus>
[11.09.23, 19:43:34] Schlotti XXX: <Anhang: 00000519-AUDIO-2023-09-11-19-43-34.opus>
[11.09.23, 19:43:54] Thomas Schmidt: It's for the course. I said so!!!
[11.09.23, 19:44:04] Thomas Schmidt: <Anhang: 00000521-AUDIO-2023-09-11-19-44-04.opus>
[11.09.23, 19:44:40] Schlotti XXX: I know I remember but I don't know what that is I SAID SO
[11.09.23, 19:45:14] Thomas Schmidt: <Anhang: 00000523-PHOTO-2023-09-11-19-45-14.jpg>
[11.09.23, 19:45:20] Thomas Schmidt: <Anhang: 00000524-AUDIO-2023-09-11-19-45-20.opus>
[11.09.23, 19:45:29] Schlotti XXX: See u
```

0:06   19:44 ✓✓

0:03   19:45 ✓✓

I know I remember but I don't know what that is I SAID SO

See u 19:45

# Transcription tools

- Support for manual transcription
  - Alignment of transcript and audio/video
  - Structured data, ready for further processing
  - Further processing (annotation etc.) inside the tool

- Family of good practice tools:    **ELAN, EXMARaLDA, FOLKER, Praat,** Transcriber, CLAN

- Text editors, word processors    → No alignment, no structured data

- „Consumer tools": F4, inqScribe → deficits in interoperability

- QDA tools: MaxQDA, atlas.ti, NVivo → dito

- ASR tools: later

**ELAN (Eudico Linguistic Annotator)**

https://archive.mpi.nl/tla/elan

**EXMARaLDA Partitur-Editor**

https://www.exmaralda.org

FOLKER (FOLK-Editor)

https://www.exmaralda.org

**Praat**

https://www.fon.hum.uva.nl/praat/

- Multi-party conversation
- Video(s) (except Praat)

| ELAN | EXMARaLDA Partitur-Editor | FOLKER | Praat |
|---|---|---|---|
| Advanced video functionality | Part of a larger system: COMA (Corpus Manager), OrthoNormal (token annotation), EXAKT (Query) | | Advanced functionality for phonetic analysis |
| Complex annotation hierarchies | Direct support for ISO/TEI<br>Direct support for transcription systems (cGAT, HIAT) | | Scriptable |
| | Built-in support for tokenisation<br>Built-in support for masking / pseudonymization | | |
| **Complex** | **Not too simple ;-)** | **One tier per speaker** | **No video support** |

- Tools are interoperable (more later)
- Things to consider:
  - ✓ complexity of your data
  - ✓ workflow integration (other tasks)
  - ✓ expertise in your team / network
  - ✓ collaboration
  - ✓ recommendations by data centres

# Automatic Speech Recognition

- „Speech to text" – machine transcription
- Dramatic improvements in the last few years
- Potential to save lots (and lots) of effort (manual transcription 1:10 up to 1:100)
- Commercial companies: Amberscript, Trint, Google, SpeechMatics, …
- Big questions:
  - Quality and precision?
  - Data protection?

# ASR Example: Amberscript

- Upload to platform
- Pay (EUR 10 to 20 per hour)
- Submit for ASR
- Edit result online
- Download result

# ASR: Quality and precision

| Manual transcription | Amberscript |
|---|---|
| **[0.8]** Here we go. Here's my voice message. **[0.4]** So. **[0.1]** Why do you need that? **[1.4]** | Here we go. Here's my voice message. So why do you need that? |
| The thing is **((clears throat))**, I would need an **ehm ehm** what's it called voice message from you. Can you do that? | The thing is, I would need an what's it called voice message from you. Can you do that? |
| Oh, **m/ ((clears throat))** maybe I should have put that in a voice message. It's for the **course**. I said so. | Oh, maybe I should have put that in a voice message. It's for the **<span style="color:red">cause</span>**. I said so. |

- Word Error Rate (WER): <5% very good // <10% good // <20% acceptable // >20% ???
- Smoothing of „performance phenomena" – disfluencies, non-verbal
- Automatic Recognition → Manual correction : difficult with WER > 20%, difficult when high precision is required
- What WER to expect?

# ASR: Quality



- Ideal: Professional podcast – expect WER < 5%
- Voice message:
    + monologic
    +/- clean, standard, vocabulary
- Zoom conference
    + for clean surrounding
    +/- for all others

# Tasks

- Transcription

- Normalisation

- Masking (de-identification, ~~anonymization~~) in audio

- Masking (pseudonyms) in transcript text

- Lemmatisation

- POS tagging

- …

# Tasks and tools in workflows

- Methodological coherence: Data models, Annotation schemes
- Interoperability of tools: Getting data from A to B
- With a view to your own analytical demands
- With a view to re-usability
- From data acquistion to data preservation (and back: research data lifecycle)

# FOLK workflow



**Field access** → **Recording** → **Check consent** → **Check metadata** → **Edit recording**

→ **Transcription** → **Quality check** → **Normalisation** (haste hast Du) → **Lemmatisation** (hast haben) → **POS tagging** (haben VERB)

→ **Quality check** → **Database** → **Publication WWW** ← Browse ← Query ← Download **User**

# CMC workflow? Example WhatsApp voice messages

- Data collection [e.g. download chat: messages as txt, audio as *.opus / convert]

- ASR in Amberscript [download results as *.vtt]

- Manual correction and masking in FOLKER

- Normalisation in OrthoNormal

- Lemmatisation and POS-tagging with TreeTagger / STTS 2.0

- Query in EXAKT

- Export to ISO/TEI

# CMC workflow? WhatsApp voice messages

- Data collection
- ASR in Amberscript
- Manual **correction** and masking in FOLKER
- Normalisation in OrthoNormal
- Lemmatisation and POS-tagging with TreeTagger / STTS 2.0
- Query in EXAKT
- Export to ISO/TEI

| | Start | End | Speaker | Transcription text | Syntax | Time |
|---|---|---|---|---|---|---|
| 1 | 00:00.0 | 00:00.76 | X | (0.8) | ✓ | ✓ |
| 2 | 00:00.76 | 00:01.34 | X | Here we go. | ✓ | ✓ |
| 3 | 00:01.34 | 00:02.43 | X | Here_s my voice message. | ✓ | ✓ |
| 4 | 00:02.43 | 00:02.82 | X | [0.4] | ✗ | ✓ |

- cGAT syntax control
- automatic tokenisation

# CMC workflow? WhatsApp voice messages

- Data collection

- ASR in Amberscript

- Manual correction and **masking** in FOLKER

- Normalisation in OrthoNormal

- Lemmatisation and POS-tagging with TreeTagger / STTS 2.0

- Query in EXAKT

- Export to ISO/TEI



- Management of pseudonyms
- Selection of audio stretches to be masked
- Automatic insertion of noises into the audio

# CMC workflow? WhatsApp voice messages

- Data collection

- ASR in Amberscript

- Manual correction and masking in FOLKER

- **Normalisation** in OrthoNormal

- Lemmatisation and POS-tagging with TreeTagger / STTS 2.0

- Query in EXAKT

- Export to ISO/TEI



- Based on tokenisation
- Automatic for German (FOLK lexicon)

# CMC workflow? WhatsApp voice messages

- Data collection
- ASR in Amberscript
- Manual correction and masking in FOLKER
- Normalisation in OrthoNormal
- **Lemmatisation and POS-tagging with TreeTagger** / STTS 2.0
- Query in EXAKT
- Export to ISO/TEI



- Based on normalisation
- Can also be done in OrthoNormal
- OrthoNormal for manual correction

# CMC workflow? WhatsApp voice messages

**EXAKT search**

RegEx (Transcription) ⌄ | Search: | \b.o\b | ⌄ |

| # | S | Communication | Speaker | Left Context | Match | Right Context |
|---|---|---------------|---------|--------------|-------|---------------|
| 1 | ☑ | | X | [0.8] Here we | go | . Here's my voice message. [0.4] So. [0.1] Why do |
| 2 | ☑ | | X | [0.8] Here we go. Here's my voice message. [0.4] | So | . [0.1] Why do you need that? [1.4] |
| 3 | ☑ | | X | go. Here's my voice message. [0.4] So. [0.1] Why | do | you need that? [1.4] |

Types
Token
Select
Time:

[0.8] Here we go. Here's my voice message. [0.4] **So**. [0.1] Why do you need that? [1.4]
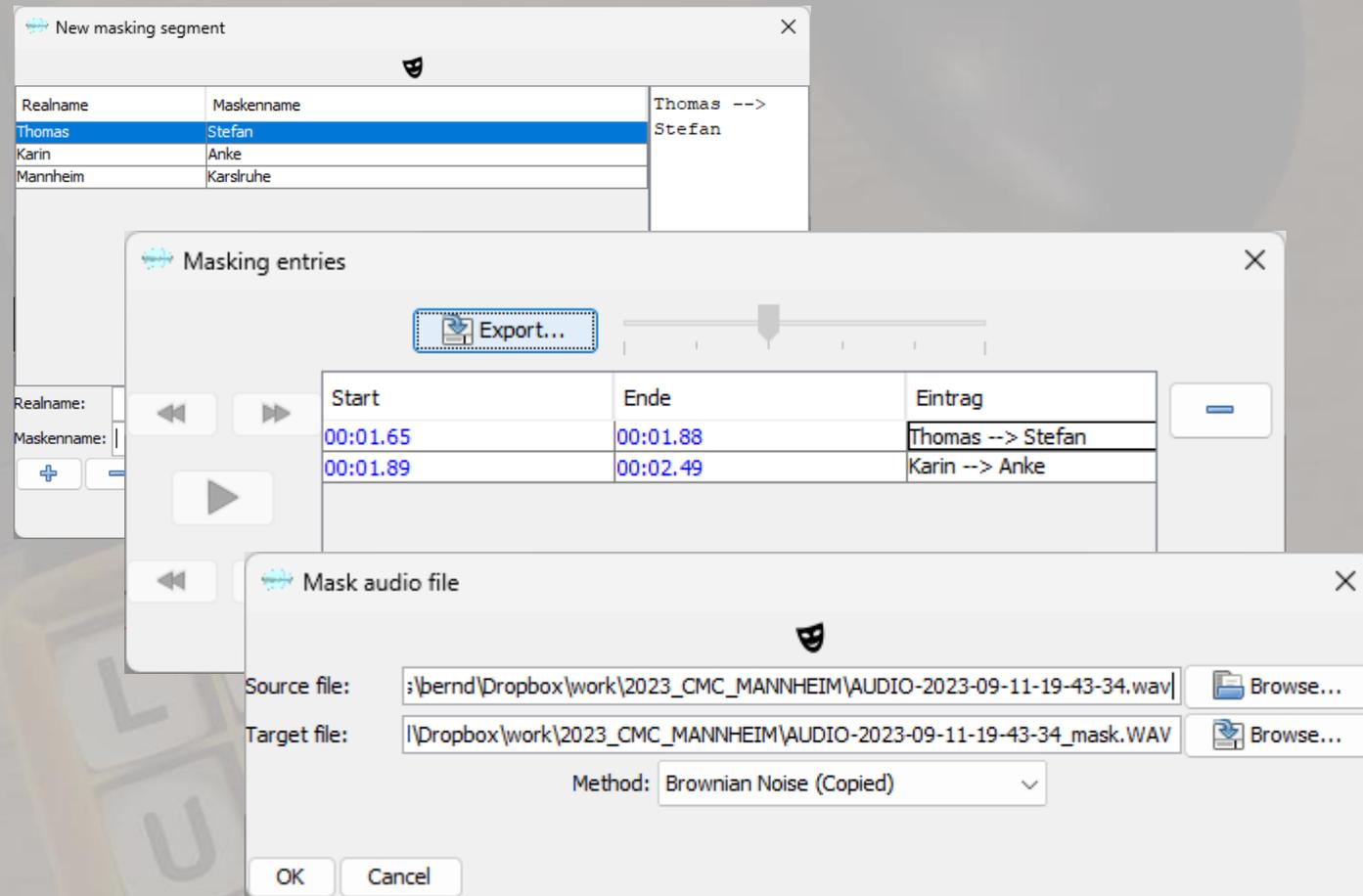
- Regular expression search
- Also on annotations (POS etc.)

- Data collection
- ASR in Amberscript
- Manual correction and masking in FOLKER
- Normalisation in OrthoNormal
- Lemmatisation and POS-tagging with TreeTagger / STTS 2.0
- **Query** in EXAKT
- Export to ISO/TEI

# CMC workflow? WhatsApp voice messages

- Data collection

- ASR in Amberscript

- Manual correction and masking in FOLKER

- Normalisation in OrthoNormal

- Lemmatisation and POS-tagging with TreeTagger / STTS 2.0

- Query in EXAKT

- Export to **ISO/TEI**



```
44    </teiHeader>
45  ▽ <text xml:lang="en">
46  ▽     <timeline unit="s">
47            <when xml:id="T0" interval="0.0" since="T0"/>
48            <when xml:id="T1" interval="0.0027775504057493327" since="T0"/>
49            <when xml:id="T2" interval="0.769381462392565" since="T0"/>
50            <when xml:id="T3" interval="1.3471119467884263" since="T0"/>
51            <when xml:id="T4" interval="2.433134155436415" since="T0"/>
```

- Directly from OrthoNormal
- Standard preservation format
- Compatible with further tools

```
59  ▽     <body>
60  ▽         <annotationBlock who="SPK0" start="T1" end="T9" xml:id="au1">
61  ▽             <u xml:id="u1"><pause dur="PT0.8S" xml:id="p1"/>
62                    <anchor synch="T2"/>
63                    <w xml:id="w1">Here</w>
64                    <w xml:id="w2">we</w>
65                    <w xml:id="w3">go</w>
66                    <pc xml:id="pc1">.</pc>
                      <anchor synch="T3"/>
                      <w xml:id="w4">Here</w>
                      <pc xml:id="pc2">'</pc>
                      <w xml:id="w5">s</w>
                      <w xml:id="w6">my</w>
                      <w xml:id="w7">voice</w>
                      <w xml:id="w8">message</w>
74                    <pc xml:id="pc3">.</pc>
75                    <anchor synch="T4"/>
76                    <pause dur="PT0.4S" xml:id="p2"/>
77                    <anchor synch="T5"/>
78                    <w xml:id="w9">So</w>
79                    <pc xml:id="pc4">.</pc>
80                    <anchor synch="T6"/>
81                    <pause dur="PT0.1S" xml:id="p3"/>
82                    <anchor synch="T7"/>
83                    <w xml:id="w10">Why</w>
```
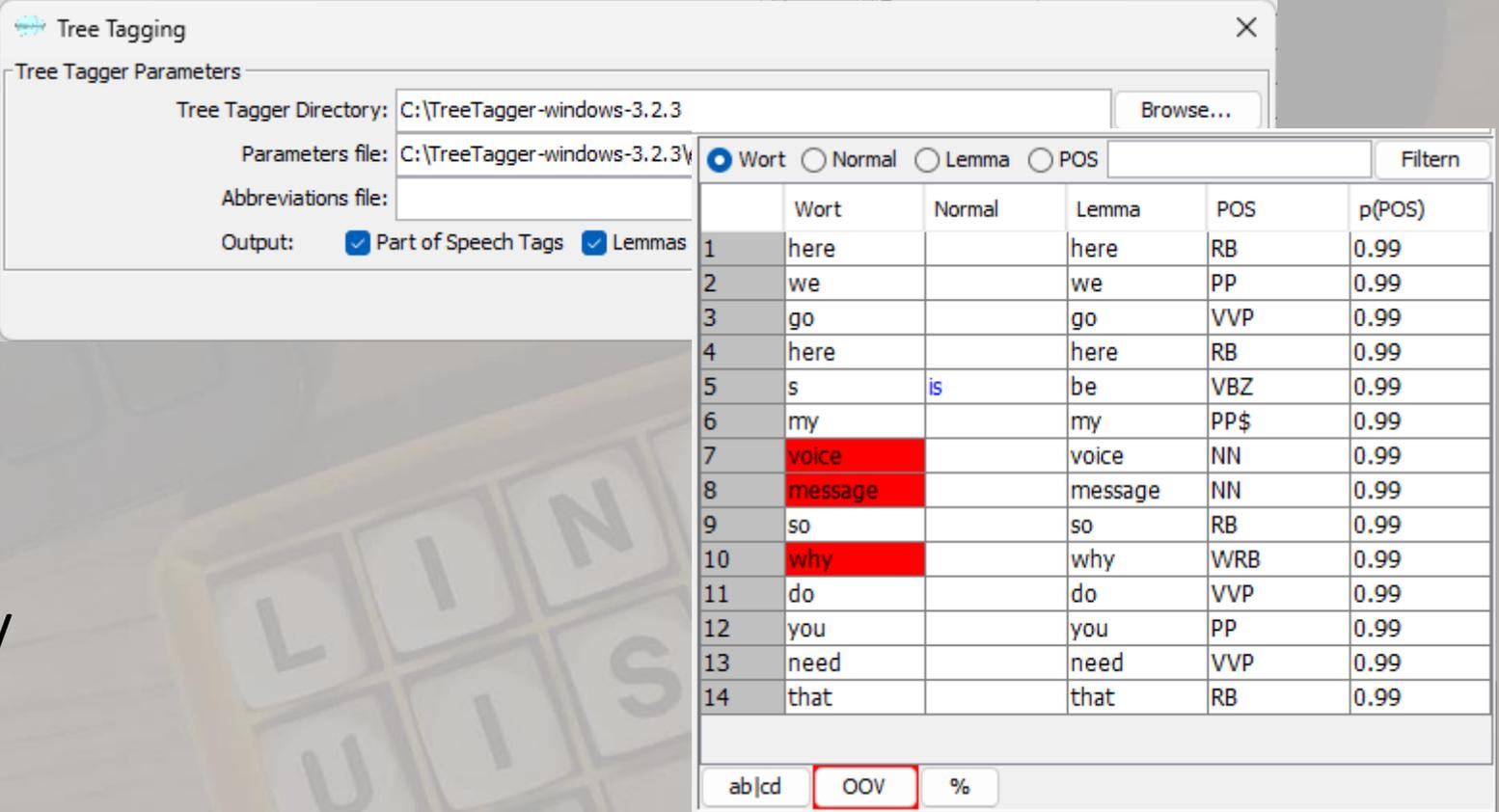
# Formats, Standards, Interoperability

# Interoperability

- Ability to exchange data
  - between tools, operating systems, etc.
  - between now and the future

- Minimum requirements
  - Structured data (Markup, CSV)
  - Documented
  - No proprietary, binary formats

- Ideally
  - Official standards
  - Semantic interoperability

# Formats

- ELAN, EXMARaLDA, FOLKER write XML formats
- Praat writes a well-defined text format, easily transformed to XML
- Very basic interoperability on the XML level
- Advanced interoperability via import and export filters in the tools
  - no information loss for simple data ➔ „round-tripping"
  - well-understood limits of interoperability ➔ ELAN > EXMARaLDA > FOLKER/Praat
- Tool formats are „de facto standards"

# Standards

- ISO 24624:2016 "Language resource management — Transcription of spoken language"
  - published by ISO in 2016, reviewed and confirmed in 2022
  - "endorsed" by the Text Encoding Initiative (TEI)
    - based on the TEI guidelines
    - TEI guidelines adapted to concepts needed for the standard
  - cross-relations to other parts of the guidelines
    - written text corpora
    - CMC corpora!
  - compatible with and supported directly or indirectly (via interoperability) by more than one tool
  - recommended / required by some CLARIN data centres

Michael Beißwenger, Harald Lüngen (2020): **CMC-core: a schema for the representation of CMC corpora in TEI.** *Corpus.*

Hedeland, Hanna / Schmidt, Thomas (2022): **The TEI-based ISO Standard 'Transcription of spoken language' as an Exchange Format within CLARIN and beyond.** *Selected Papers from the CLARIN Annual Conference 2021.*

# ISO/TEI Spoken vs. CMC TEI

**Transcript in ISO/TEI Spoken**

```
59 ▽        <body>
60 ▽            <annotationBlock who="SPK0" start="T1" end="T9" xml:id="au1">
61 ▽                <u xml:id="u1"><pause dur="PT0.8S" xml:id="p1"/>
62                     <anchor synch="T2"/>
63                     <w xml:id="w1">Here</w>
64                     <w xml:id="w2">we</w>
65                     <w xml:id="w3">go</w>
66                     <pc xml:id="pc1">.</pc>
67                     <anchor synch="T3"/>
68                     <w xml:id="w4">Here</w>
69                     <pc xml:id="pc2">'</pc>
70                     <w xml:id="w5">s</w>
71                     <w xml:id="w6">my</w>
72                     <w xml:id="w7">voice</w>
73                     <w xml:id="w8">message</w>
74                     <pc xml:id="pc3">.</pc>
75                     <anchor synch="T4"/>
76                     <pause dur="PT0.4S" xml:id="p2"/>
77                     <anchor synch="T5"/>
78                     <w xml:id="w9">So</w>
79                     <pc xml:id="pc4">.</pc>
80                     <anchor synch="T6"/>
81                     <pause dur="PT0.1S" xml:id="p3"/>
82                     <anchor synch="T7"/>
83                     <w xml:id="w10">Why</w>
```

**<post mode=„spoken"> in CMC-TEI**

```
<post mode="spoken" creation="human" synch="#t003" who="#A05"
    xml:id="m7"> Sagt Anne auch gerade. JA! Kann ich zustinmen. </post>
<post mode="written" creation="human" synch="#t003" who="#A02"
    xml:id="m8"> Da kostet ein Haarschnitt 50 € <figure type="emoji"
    creation="template">
    <desc type="meaning">face screaming in fear</desc>
    <desc type="unicode">U+1F631</desc></figure>
</post>
```

- Internal structure / Level of detail
  - „post"-internal time anchors
  - tokenisation
  - distinction words vs. non-words

# So what?

- What tool(s)? What workflow? What standard?
- It depends…
  - Status, amount and duration of audio/video in your CMC data
    - Sporadic and typically short (WhatsApp)
    - Main data type and longer (Zoom conference)
  - Envisaged processing of your corpus
    - „Plain text" database : Qualitative, example-based analysis
    - Detailed multi-level annotation: Corpus linguistics, quantification
  - Tool preferences
    - „End-user" tools
    - XML editors, scripts etc.
  - Your eco-system
    - Support by / requirements from a data center, colleagues and collaborators

Questions & attempts at Answers?

# Further advice and support

- CLARIN K-Center for CMC: Eurac Bozen / IJS Ljubljana / LLF France / IDS Mannheim
- Other centers in CLARIN (Europe) or NFDI (Germany)
- Good practice examples? Few with audio/video so far…
- Some free support for all tools presented here:
  - ELAN user forum / [support@exmaralda.org](mailto:support@exmaralda.org) / Praat mailing list
- Training courses for FOLKER (IDS) and EXMARaLDA (myself)
  - Next IDS course: October, 20th
- linguisticbits.de as a data management partner

Dr. Thomas Schmidt
https://linguisticbits.de
thomas@linguisticbits.de

| 21 | FOLK_E_00055_SE_01_T_01 | AM | hm (0.22) sehr liebenswürdig **dankeschön** (2.12) sehr nett (0.25) so … |
| 22 | FOLK_E_00055_SE_01_T_02 | AM | alle zusammen (( Lachansatz )) ja (.) **danke** ja (0.2) dir auch … |
| 23 | FOLK_E_00055_SE_01_T_03 | US | bandscheibenvorfall °h h° (( lacht )) gott sei **dank** (( Lachansatz )) °hh nein aber es … |
| 24 | FOLK_E_00055_SE_01_T_04 | US | sprechen (( Lachansatz )) (( Lachansatz )) (( lacht )) (( lacht )) (( lacht )) (0.8) °hh gott sei **dank** da bin ich … |
| 25 | FOLK_E_00055_SE_01_T_04 | US | … wurde gott sei **dank** is der nach … |
| 26 | FOLK_E_00055_SE_01_T_04 | US | äh gott sei **dank** (( Lachansatz )) (0.47) am sechzehnten dezember … |
| 27 | FOLK_E_00055_SE_01_T_05 | US | glas wein haben °h **danke** gerne (( Lachansatz )) wobei des … |
| 28 | FOLK_E_00055_SE_01_T_05 | US | prost prost prost **dankeschön** darauf dass wir … |
| 29 | FOLK_E_00057_SE_01_T_01 | ME | °h gut hh° (0.33) willkommen **danke** h° (0.34) äh h° wir haben … |
| 30 | FOLK_E_00057_SE_01_T_01 | ME | wieder rein °h jawohl **danke** … |
| 31 | FOLK_E_00058_SE_01_T_01 | HN | … jahrn °h gott sei **dank** noch mal im … |
| 32 | FOLK_E_00058_SE_01_T_01 | HN | okay (0.43) ganz herzlichen **dank** lara ich danke … |
| 33 | FOLK_E_00058_SE_01_T_01 | XL | dank lara ich **danke** auch wir bedanken … |
| 34 | FOLK_E_00059_SE_01_T_01 | HN | … wichtig °h ganz herzlichen **dank** hm (0.59) so °h |
| 35 | FOLK_E_00042_SE_01_T_01 | LP | da drauf dahin (0.61) **danke** na ja ich … |
| 36 | FOLK_E_00042_SE_01_T_02 | LK | … so ne art (.) **dankbarkeitsbewusstsein** von den mädels |
| 37 | FOLK_E_00042_SE_01_T_02 | LS | … bisschen mit der **dankbarkeit** dass sie dann … |
| 38 | FOLK_E_00042_SE_01_T_02 | AM | … das liegt an **dankbarkeit** dass die des … |
| 39 | FOLK_E_00042_SE_01_T_02 | LK | oder (0.82) ob des **dankbarkeit** is des will … |
| 40 | FOLK_E_00042_SE_01_T_02 | LP | ach also des **dankbarkeit** ich glaub ich … |