

Challenges on the promising road to Automatic Speech Recognition of privacy- sensitive Dutch doctor-patient consultation recordings

Homo Medicinalis: <http://homed.ruhosting.nl>

Berrie van der Molen (b.j.vandermolen@uu.nl), **Cristian Tejedor-García** (<https://cristiantg.com>),
Henk van den Heuvel, Roeland Ordelman, Toine Pieters, Sandra van Dulmen, Arjan van Hessen



1: Medical domain adaptation in Automatic Speech Recognition

Automatic Speech Recognition (ASR) systems are made for “generic” Dutch.

How do we create an ASR infrastructure for the medical domain?

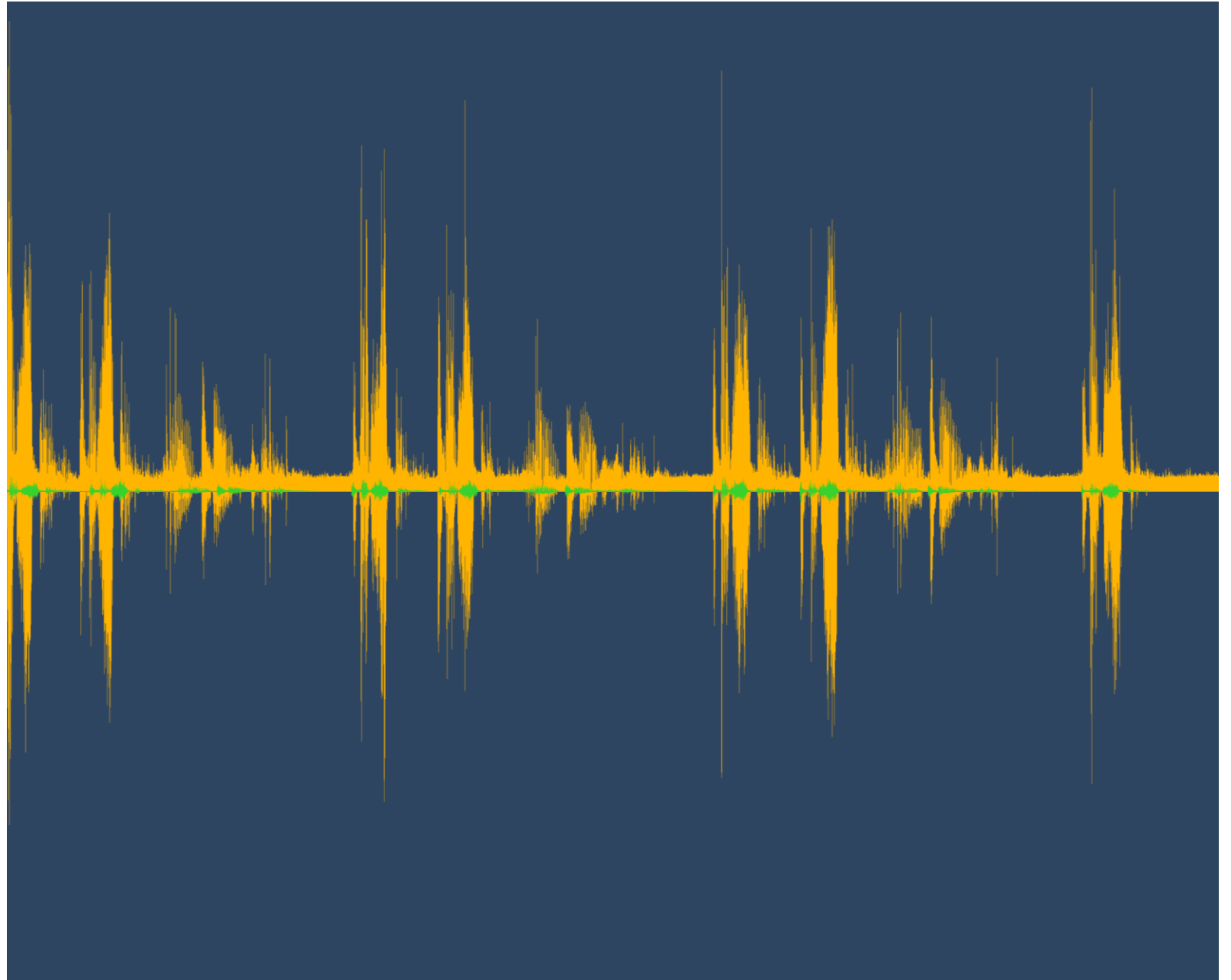
In this presentation

- CONTEXT 1. Medical domain adaptation in ASR
- 2. HoMed: medical domain adaptation
- CURRENT 3. Challenge: Training data with privacy-sensitive material
- 4. Challenge: Transcribing medical terms
- FUTURE 5. Challenge: Search and analysis of ASR-enriched doctor-patient consultations
- 5. Expected outcomes

2. Homo Medicinalis (HoMed): an infrastructure for ASR of Dutch doctor-patient consultation recordings

Stage 1: Language model
fine-tuning based on public
Medicijnjournaal episodes
and lists of medical terms

Stage 2: Semantic and
acoustic level adaptation
based on privacy-sensitive
doctor-patient consultation
recordings held at Nivel



2: Homo Medicinalis

DONE: Adding medical terms to general Dutch lexicon

DONE: Using alternative (non-privacy—sensitive) medically themed AV-recordings

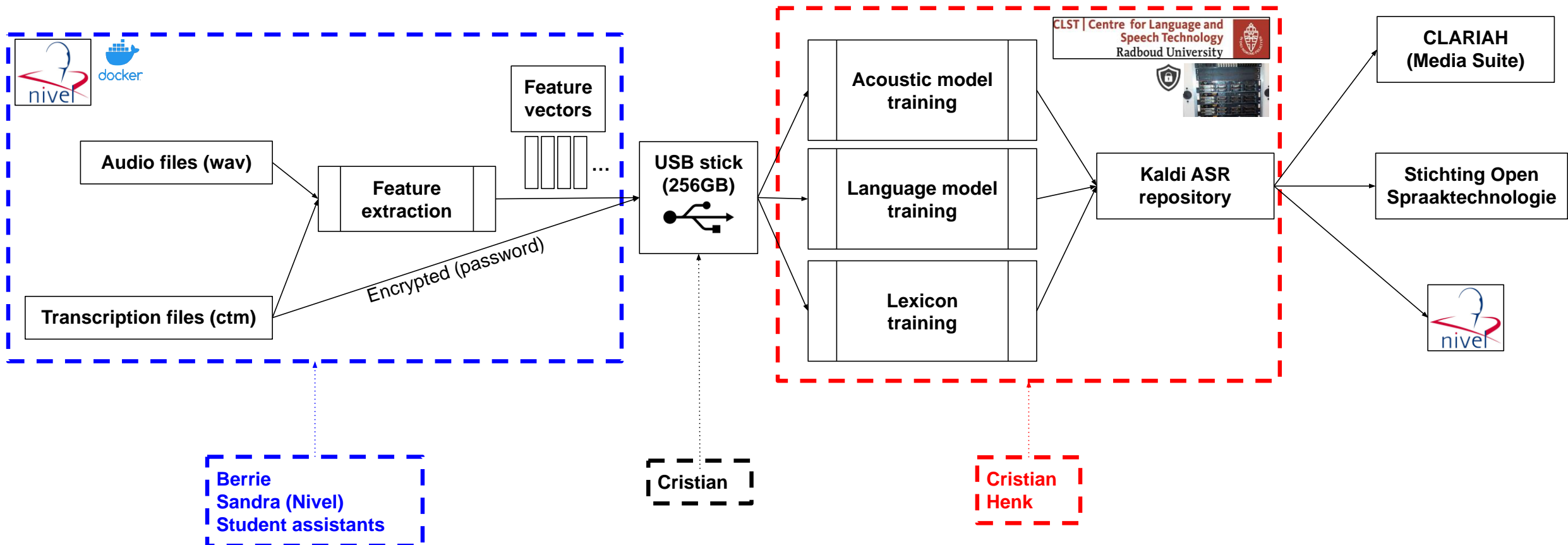
(*Medicijnjournaal*)

CURRENT: Using the privacy-sensitive recordings of Nivel

FUTURE: Developing distant and close reading strategy for ASR-enriched doctor patient consultation recordings



3: Doing research with privacy-sensitive recordings and transcripts



4a: Dealing with medical terms in ASR

Transcribing medical terms

- Mispronounced medical terms (solution: `mispronuciation[correct_pronunciation]` in transcript)
- More than one accepted spelling (we use [Farmacotherapeutisch Kompas](#))
- Digits in medical terms (e.g. "diabetes mellitus type 2"; "SGLT2-remmers")

Wees consequent in je schrijfwijze: als je eenmaal een bepaalde spelling gebruikt hebt, gebruik dan dezelfde wanneer de eigennaam nogmaals voorkomt in hetzelfde fragment.

- C2. **Gebruik hoofdletters voor titels van boeken, platen, rapporten, richtlijnen** etc. Als een titel uit meerdere woorden bestaat, worden alle woorden met een hoofdletter geschreven.

- ik heb Hugo Claus' boek Het Verdriet Van België gekocht.
- Heb jij De Naam Van De Roos al gelezen?
- De nieuwe richtlijn Cardiovasculair Risicomanagement is nu beschikbaar.
- NHG-Standaard⁷

- C3. **Gebruik hoofdletters voor gespelde letters.** Namen en kentekens worden regelmatig gespeld. Bij gebruik van een (quasi) pilotenalfabet moeten de betreffende woorden worden gebruikt; deze dienen wel te beginnen met een hoofdletter (zie ook F3-5).

- Ik sta op de sneekkade **S N E E K K A D E**
- Ik sta op de sneekkade **Sierra November Echo Echo Karel Karel Alpha Delta Echo**

- C8. **Gebruik xxx, -xxx, xxx- of -xxx- voor onverstanebare uitingen.** Gebruik deze code wanneer je een woord, een deel van een woord of een reeks woorden niet goed heeft verstaan en je er ook niks van kunt maken.

- ik heb al **xxx** bladzijden gelezen.

⁵ Let op: gebruik alleen een hoofdletter voor merknamen van geneesmiddelen, niet voor stoffen en generieke aanduidingen. Bijvoorbeeld: Het middel aciclovir wordt verkocht onder de merknaam Zovirax.

⁶ Met deze spelling wijken we dus nadrukkelijk af van de conventionele spelling; we doen dit om een zo consistent mogelijke transcriptie te verkrijgen en om eigennamen gemakkelijk als een geheel te herkennen.

⁷ Standaarden worden wel met hoofdletter geschreven, maar modules niet. Dus: NHG-Standaard maar FTO-module. Dit is omdat dit de meest gangbare manier van schrijven is, al wordt dit in de praktijk niet consequent gedaan.

4b: Dealing with medical terms in distant reading

Learning from close reading / learning from our transcripts.

Search and analysis strategies:

- *Symptoms / side effects*
- *Physical description of medicines*
- *Colloquial medicine descriptions*



5: Expected outcomes

- An infrastructure for ASR of doctor patient consultation recordings:
 - *In-house at Nivel*
 - *As component in the CLARIAH Media Suite*
 - *As code available via Stichting Open Spraaktechnologie*
- A methodology for training recordings in a privacy-sensitive institutional setting (e.g. justice, police, etc.)
- An adapted close and distant reading strategy for research of doctor patient consultations





**Universiteit
Utrecht**

Sharing science,
shaping tomorrow

Contact: Berrie van der Molen (b.j.vandermolen@uu.nl)

Homo Medicinalis is a collaborative PDI-SSH 2020 infrastructure project



NIVEL
Kennis voor betere zorg



Radboud Universiteit



**UNIVERSITY
OF TWENTE.**

