

automated speech and audio analysis for semantic access to multimedia



Franciska de Jong, Roeland Ordelman, Marijn Huijbregts
Human Media Interaction, University of Twente
The Netherlands

overview

1. variety of multimedia content and semantic annotation
2. automatic metadata generation by exploiting linguistic content:
 1. collateral data
 2. audio analysis:
 1. speech recognition
 2. audio classification
3. summary

content is king ...

- news (professionals & non-professionals)
- historical material (retrospective digitization)
 - interviews (oral history)
- meetings (governmental & corporate)
- presentations, lectures
- private material (lifelogs?)



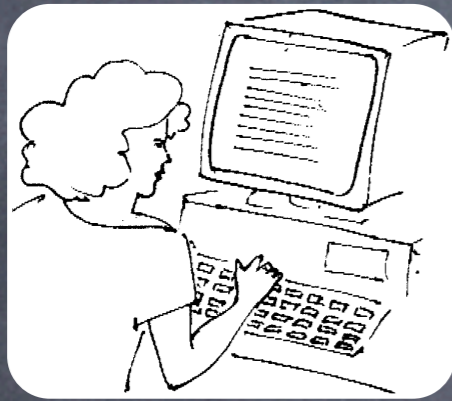
but metadata rules!

- ① networked electronic media and increasing size require automation of content-based extraction
- ① use linguistic content in multimedia archives: boost the accessibility
- ① Two examples:
 - ① 1995: subtitling for broadcast news retrieval
 - ① TRECVID: best performing systems exploit speech transcripts

semantic gap

- bridge gap between user needs and content features
- traditional approach: manual annotation + controlled vocabulary index terms
- enhance traditional approach: automatic metadata generation
- challenge: combine various types of metadata and exploit the added value

2. automatic metadata generation by exploiting linguistic content



manual
annotation



use collateral
data



audio
classification



speech
recognition

use audio
analysis



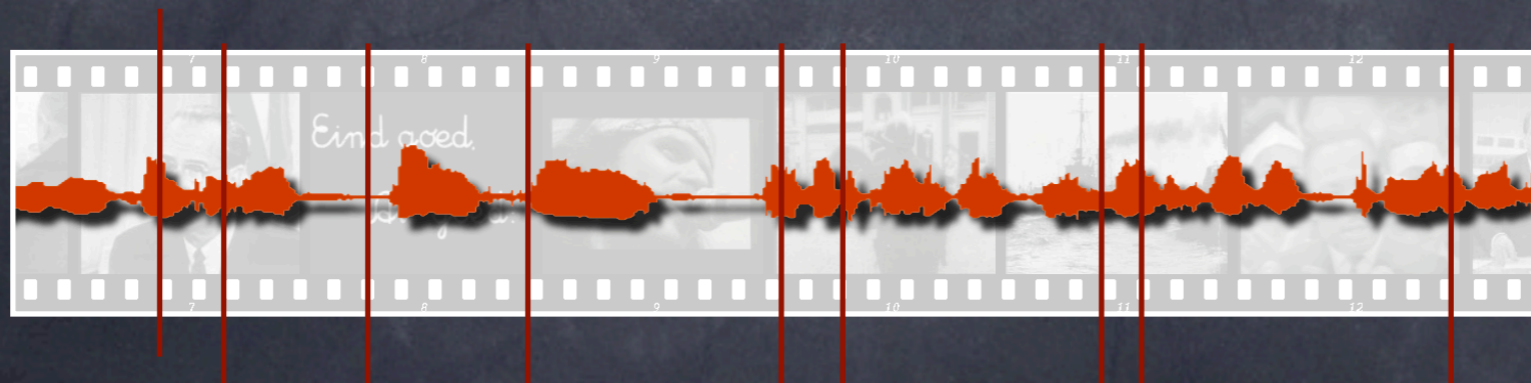
metadata for searching

exploiting collateral data

- broadcasts:
 - subtitling information for hearing impaired (time-codes included; Dutch subtitling has topic-boundaries)
 - teleprompter files, recording scripts
- other:
 - meeting minutes, notes (e.g., soccer matches, interviewer notes, presentation notes)

exploiting collateral data

- collateral data: locate documents
- searching within documents: alignment of text to multimedia document
 - label text with time-codes: speech recognition techniques (alignment)



alignment

- speech recognition training: forced alignment
- alignment: find most optimal distribution given audio + the words that are spoken
- caveats: large chunks, missing parts, model mismatch (e.g., low audio quality, old-fashioned speech)
- solutions: model adaptation, two-pass strategy

2.1 exploiting linguistic content: collateral data

Continuous Access To Cultural Heritage

Human Media Interaction

add a picture based on speech content

subtitling following speech

complete speech

relevant segment

search word

00:00.0 06:14.4

06:07.5

Terug

bevrijding van uw leven letterlijk en geestelijk

Browsing and searching WWII speeches of Dutch Queen Wilhelmina

cross-media mining

- not only synchronize audiovisual material that approximates the speech
- any kind of textual resource that is accessible: , open source titles and proprietary data
 - trusted webpages
 - newspaper articles
- shift focus from indexing individual documents to indexing multiple multimedia databases

cross-media news browsing

CV of principal person

Justitie.nl

drs. M.C.F. Verdonk, minister voor Vreemdelingenzaken en Integratie

Marie Cornelia Frederika (Rita) Verdonk werd op 18 oktober 1955 geboren te Utrecht.

Na het behalen van het diploma Atheneum aan het Niels Stensencollege te Utrecht studeerde zij sociologie (specialisatie organisatie/sociologie en criminologie) aan de Katholieke Universiteit Nijmegen (doctoraal examen 1983).

Merouze Verdonk was daarna tot 1996 werkzaam bij het ministerie van Justitie. Tot 1984 was zij directeur van de Dienst voor de Rechterlijke Gevoelenszaken. Zij

20060328-journaal-1958-069.smil

703Kbps

Verdonck heeft de leidinggevenden van de IND uit

Justitie

Ministerie van Justitie
Immigratie- en Naturalisatiedienst

IND, de toelatingsorganisatie van Nederland

Verblijfwijs

on topic governmental publication

20060328-journaal-1958-064.smil

702Kbps 0:24/0:29

toe die gesprekken lopen de asielzoekers mogelijk

related news segment

Volkskrant Den Haag

Kamer: IND altijd bij hoorzitting ex-asielzoekers

DEN HAAG - De Immigratie- en Naturalisatiedienst (IND) moet altijd aanwezig zijn bij gesprekken in Nederland tussen ex-asielzoekers en vertegenwoordigers van landen van herkomst. Dat stelde een meerderheid van de Tweede Kamer donderdag in een overleg met minister Verdonk (Vreemdelingenzaken).

on topic newspaper article

news segment (indexed on subtitles or speech)



Geef het onderwerp waar u op wilt zoeken:

Zoekterm(en)

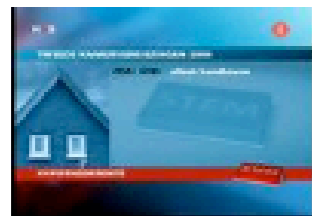
Informatie uit: Teletext Spraakherkenner

Sortering Kranten Datum Krant
 Score KrantSegment (1)
 Score JournaalSegment (2)

(1) hoe scoort de zoekterm in een krantenartikel
 (2) hoe scoort de tekst van een journaalsegment in een krantenartikel

[NB: Aantal krantenartikelen in database: 1419408]

[Klik op thumbnail om journaal te bekijken]



ZOEKRESULTATEN

Journaal maandag 13 november 2006 (00:00:41)

● Score: 11.

...De partijen willen dat er meer voor starters gebouwd wordt de bedrijven waar wij nu over Gunco bonus voor starters hoeven geen overdrachtsbelasting te betalen. Tella wil een stakers fonds voor goedkope leningen en de VVD denkt dat meer bouwen. Oplossing is dan de **hypotheekrente aftrek** twee jaar WW junior Freek Vossius handhaven de bedrijven daar bij de burens voor de nieuwe gevallen beperken tot twee en veertig procent nu strenge vijftig procent van de rente **aftrekbaar** ook GroenLinks wil de **aftrek** beperken voor alleen de lagere en middeninkomens en de SP legt de grens waarvoor bij hypotheek tot drie honderd vijftig duizend euro. ... >>

Krant	Datum	Score			Artikel titel
		Krant	Journaal		
de Volkskrant	6 jun 2003	8.897	42.79	✓	Huizenprijzen >>
 www.ad.nl	28 mrt 2002	10.93	38.42	✓	Heilig huisje >>
 www.ad.nl	2 mrt 2002	6.783	39.52	✓	PvdA: Maximale aftrek jonge huizenkoper >>
 www.ad.nl	12 dec 2001	12.24	49.70	✓	PvdA wil jonge huizenkopers bevoordelen >>
 Trouw	26 aug 2000	10.50	37.96	✓	Nederlands uniek, maar hoelang nog? >>

speech indexing

- speech recognition: full text annotations
- spoken document retrieval in the American-English broadcast news (BN) domain was declared "solved" with the NIST-sponsored TREC SDR track in 2000

speech indexing: main steps

- train acoustic models using speech that resembles speech in task domain
- train language models using text data that resembles speech & word usage in task domain

- de-multiplex, audio-conversion
- segmentation: speech/non-speech, speaker-diarization
- recognition (multi-pass: speaker/LM-adaptation)

- post-process speech transcripts (format, rich text)
- index & search

challenges

- rough baseline: at least 50% speech recognition accuracy required for IR
- speech type: planned/formal – spontaneous/informal – cross-talk
- speaker characteristics: non-native, dialect
- audio quality (recording conditions, acoustic environment)
- availability of training data
- vocabulary (OOV – QOV)

audio classification



- 👁️ speech (words)
- 👁️ structure (silence, music/jingles, speaker)
- 👁️ speaker:
 - 👁️ identity (gender, native, dialect, age)
 - 👁️ emotion

summary

- King-size multimedia content repositories
- require content-based extraction techniques
- exploit readily available collateral data and audio analysis techniques
- (multiple) linguistic annotations allow for cross-media browsing
- challenge: how to combine various types of metadata

PS: has this been recorded?