



Jan Gorisch, Department of Pragmatics; Thomas Schmidt, linguisticbits.de

EVALUATING WORKFLOWS FOR CREATING ORTHOGRAPHIC TRANSCRIPTS FOR ORAL CORPORA BY TRANSCRIBING FROM SCRATCH OR CORRECTING ASR-OUTPUT

LREC-COLING 2024, Turin

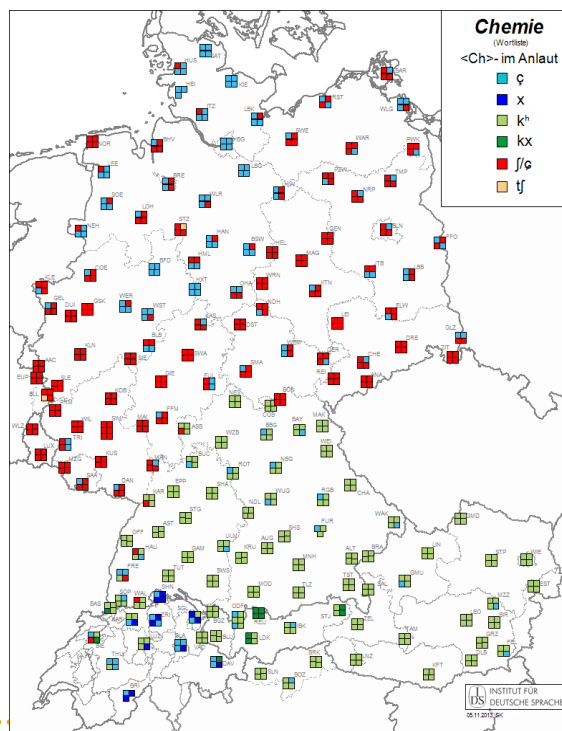
INTRODUCTION

The AGD* hosts...

- Conversation Corpora
- Variation Corpora
- Extra-territorial Varieties



also historic data from
the 1950s and 60s



IDS

LEIBNIZ-INSTITUT FÜR
DEUTSCHE SPRACHE

```
Sprecher Ortsdaten Koordinaten Geocode Geografische_Länge
33 <Sigle_in_Transkripten>AAC1</Sigle_in_Transkripten>
34 <Anmerkungen/>
35 </Basisdaten>
36 <Ortsdaten Typ="Geburtsort">
37 <Land Kürzel="DE">Deutschland</Land>
38 <Region>Nordrhein-Westfalen</Region>
39 <Kreis>Nicht vorhanden</Kreis>
40 <Ortsname>Aachen</Ortsname>
41 <Koordinaten>
42 <Geocode>
43 <Geografische_Breite>50.776207</Geografische_Breite>
44 <Geografische_Länge>6.083788</Geografische_Länge>
45 </Geocode>
46 <Planquadrat>Nicht vorhanden</Planquadrat>
47 <Anmerkungen>Nicht dokumentiert</Anmerkungen>
48 </Koordinaten>
49 <Ortsteil>Nicht dokumentiert</Ortsteil>
50 <Ortsbeschreibung>Nicht vorhanden</Ortsbeschreibung>
51 <Aufenthaltsdauer>1988-2006</Aufenthaltsdauer>
52 <KB-Link>
53 <Anmerkungen>18 Jahre</Anmerkungen>
54 </Ortsdaten>
55 <Ortsdaten Typ="Schulort">
56 <Ortsdaten Typ="Wohnort">
57 <Sprachdaten>
58 <Sprachkenntnisse Sprachname="Englisch">
59 <Sprachkenntnisse Sprachname="Französisch">
60 <Sprachkenntnisse Sprachname="Griechisch">
61 <Sprachproduktion>
62 <Sprachgebrauch Domäne="In der Familie">
63 <Sprachgebrauch Domäne="Mit Freunden">
64 <Sprachgebrauch Domäne="Mit Mitschülern">
65 <Sprachgebrauch Domäne="Im Schulunterricht">
66 <Sprachgebrauch Domäne="In formellen Situationen (Vortrag/Rede)">
67 </Sprachdaten>
68 <Beziehung_zu_anderem_Sprecher Kennung_anderer_Sprecher="DH--_S_00">
69 <Typ_der_Bezugsperson>Klassenkamerad</Typ_der_Bezugsperson>
70 <Dauer_der_Beziehung>Nicht dokumentiert</Dauer_der_Beziehung>
```

*AGD: Archive for Spoken German (<https://agd.ids-mannheim.de>)

TRANSCRIPTS (STATE 2018)

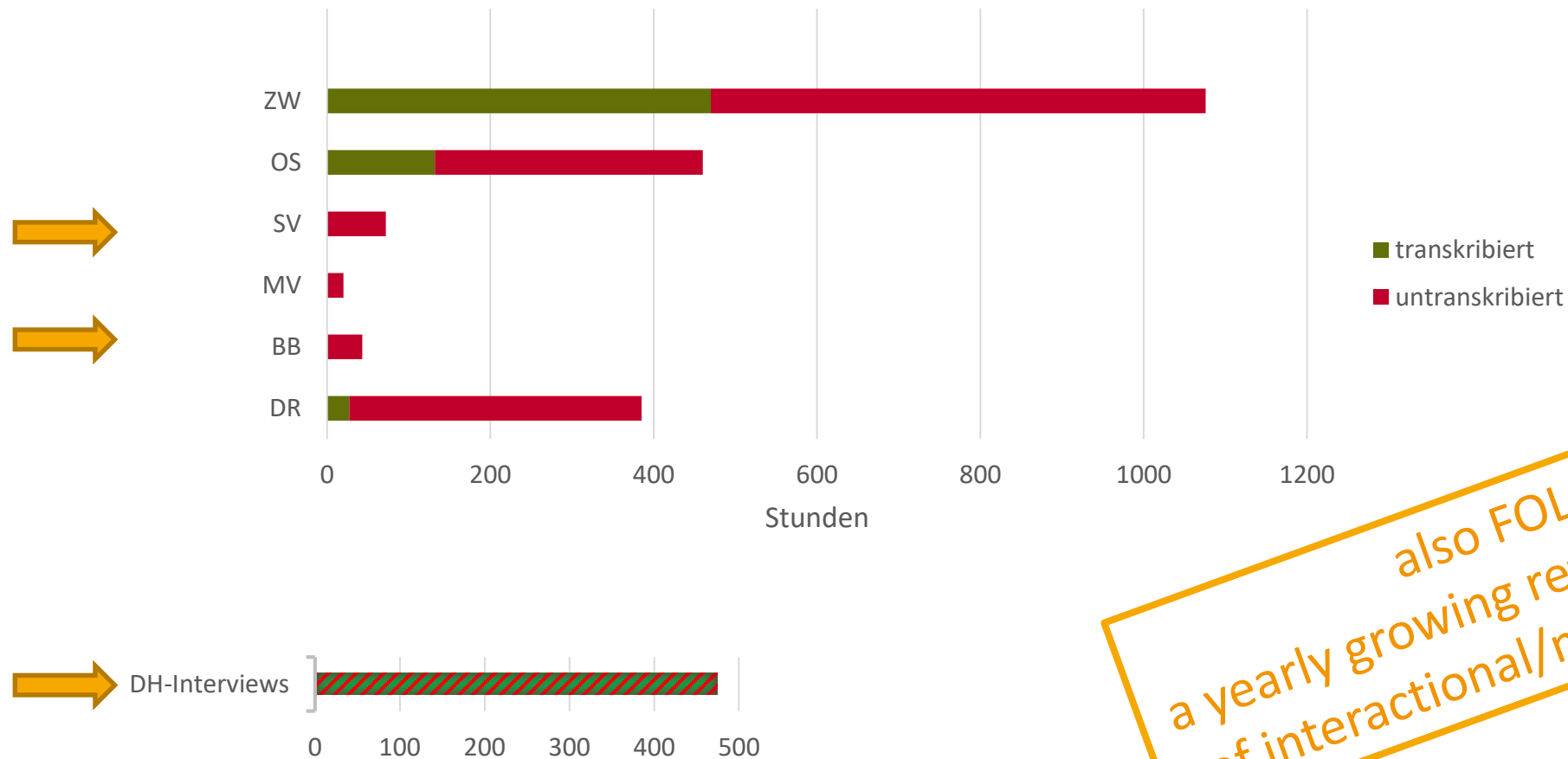
| CORPUS (selection) | # record. total | # record. transcribed | % transcribed | hours transcribed | hours untranscribed |
|---------------------------|--------------------|--------------------------|------------------|----------------------|------------------------|
| ZW – Zwirner | 5796 | 2495 | 43% | 470 | 606 |
| OS – ehem. Dt. Ostgebiete | 981 | 280 | 29% | 132 | 328 |
| SV – Südwest. u. Vorarlb. | 242 | 0 | 0% | 0 | 72 |
| MV – Varia | 72 | 0 | 0% | 0 | 20 |
| BB – Böblingen | 73 | 2 | 3% | 1 | 42 |
| DR – dt. Mundart. DDR | 444 | 33 | 7% | 27 | 358 |
| total | 7608 | 2810 | 37% | 630 | 1426 |

| | | | | | |
|-------------------------------|-----|---------|------|-------------------|-------------------|
| DH-Interviews – Deutsch Heute | 688 | 688 / 2 | ~50% | (interviewee) 237 | (interviewer) 237 |
|-------------------------------|-----|---------|------|-------------------|-------------------|

TRANSCRIPTS (STATE 2018)

IDS

LEIBNIZ-INSTITUT FÜR
DEUTSCHE SPRACHE



also FOLK:
a yearly growing reference corpus
of interactional/multimodal Talk

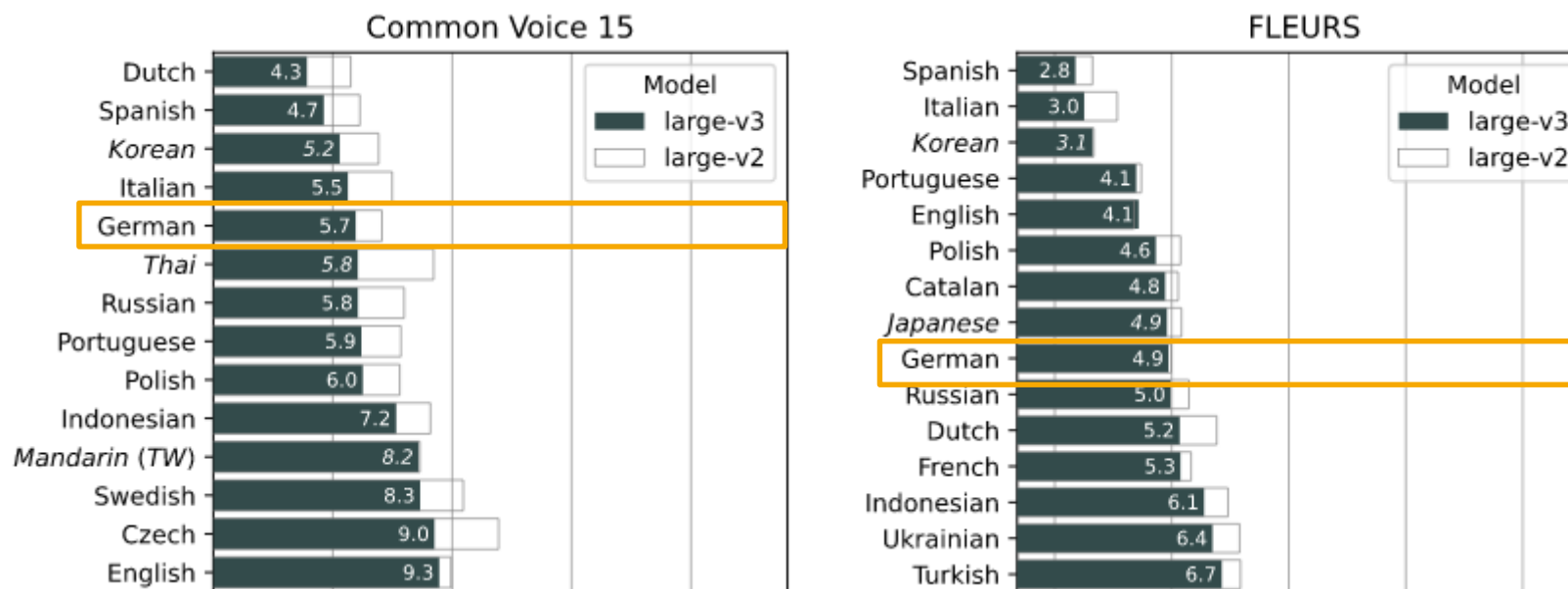
- Estimated need for transcription: 75.000h
 - e.g. 15 students for 10 years
 - + Overhead for management, quality control, technical supervision, documentation
 - Advanced dialect-competence(s) necessary – Alsatian, Low-German, Silesian, (data also contain Frisian, Sorbian, Dutch)
- Central, manual transcription in principle too expensive, not organisable
- Alternatives (cf. Brinckmann 2009)
 - Outsourcing
 - „Crowd“-Sourcing
 - (partial) automation

ASR



OPENAI'S WHISPER MODEL

- Word-Error-Rate (WER)



As discussed in the accompanying paper, we see that performance on transcription in a given language is directly correlated with the amount of training data we employ in that language.

TODAY'S RESEARCH QUESTIONS

Huge improvements in ASR-quality in recent years (months)

- When is ASR-output good enough? (e.g. for training oral corpora)
- If ASR-output needs to be better, are there more efficient ways of transcribing from scratch?

We might now be at a point where correcting ASR-output is easier than transcribing from scratch

Is ASR the solution to the problem?

Or do we need a different problem to the solution?

MATERIAL FOR THE EXPERIMENT

- Speech-biographic interviews of the DH-Korpus („deutsch heute“)
- Hannover (speaker HAN1, HAN2, ... HAN4 and interviewer NL)
- Innsbruck (speaker IBK1, IBK2, ..., IBK4 and interviewer MF)
- Later: Recordings from DH-Zurich, corpora: BB, SV,

ASR: OpenAI-Whisper

medium model, later: large-v1

Speaker diarization acc

pyannote.audio

Conversion to EXML

editor format (*.exb)

! Limitations !

The AGD is an operating Archive:
Besides evaluation, we try to produce transcripts and
provide annotators with the best input possible

EXAMPLE & INTERFACE



IDS

LEIBNIZ-INSTITUT FÜR
DEUTSCHE SPRACHE

EXMARaLDA Partitur-Editor 1.7 [N:\Workspace\ASR\Korrektur_ASR_Transkripte\PaulRoelle\Korrektur_Abschnitt5\DH--_E_00062_SE_31_A_01_DF_01_Abschnitt5_T_01_DF_01.exb]

File Edit View Transcription Tier Event Timeline Format Web Services Legacy Help

Hast du das Gefühl, die sprechen genauso wie du oder leicht anders?

02:25.02 13.92 02:38.94

02:28 02:29 02:30 02:31 02:32 02:33 02:34 02:35

+ Add event... Append interval 1 [*] [*] [*] [*] [*] [P]

| | 38 [02:23.0] | 39 [02:25.0] | 40 [02:31.0] | 41 |
|------------|--------------------------------|---|---|----|
| MF | die? | Hast du das Gefühl, die sprechen genauso wie du oder leicht anders? | | H |
| IBK3_ortho | einundzwanzigsterund vierzehn. | | Ich glaube, meine Geschwister sind mehr so, reden mehr so wie meine Eltern, glaube ich. | |
| None | | | | |
| [segments] | | | | |

Done.

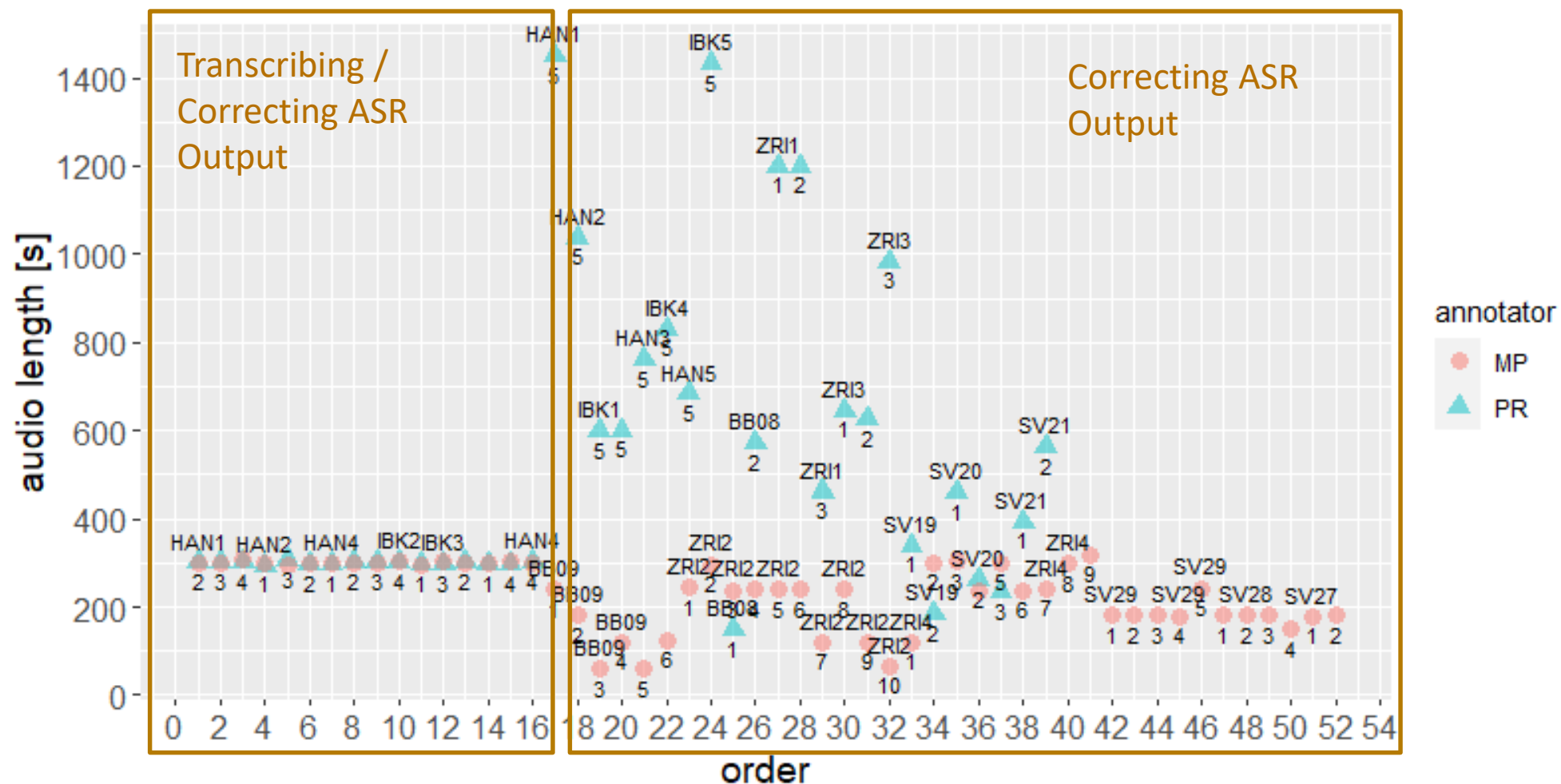
[16:34:12] Partitur-Editor started

Segmentation: **GENERIC** Player: **JDS-Player**

EXPERIMENTAL DESIGN / TASK(S)

IDS

LEIBNIZ-INSTITUT FÜR
DEUTSCHE SPRACHE



INSTRUCTIONS AND METRIC

“Create orthographic transcripts”

- Alignment
- Mask tier

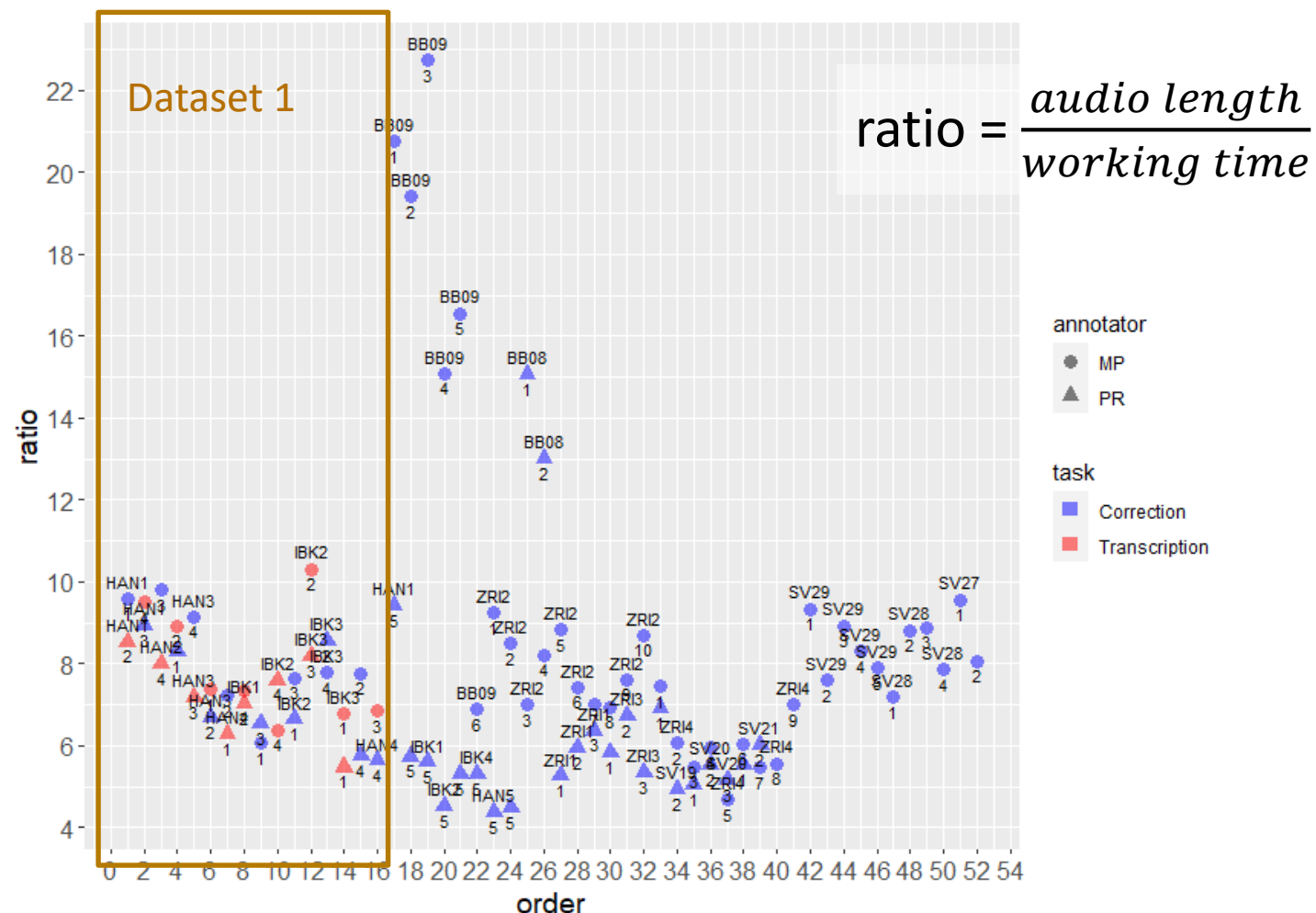
“Note the (working) time spent per stretch of recording”

Metric: ratio of $\frac{\text{audio length}}{\text{working time}}$

RESULTS

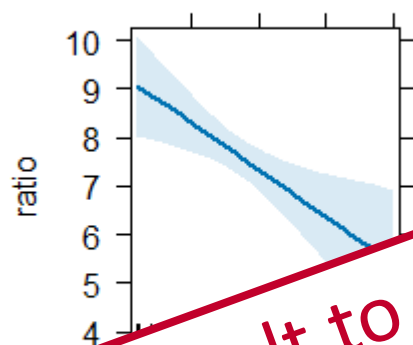
Linear regression model in R

- Criterion variable:
 - ratio
- Predictors:
 - Task (transcribing vs. correcting)
 - Corpus
 - Recording place
 - Annotator
 - Order



TRANSCRIBING VS. CORRECTING

order effect plot

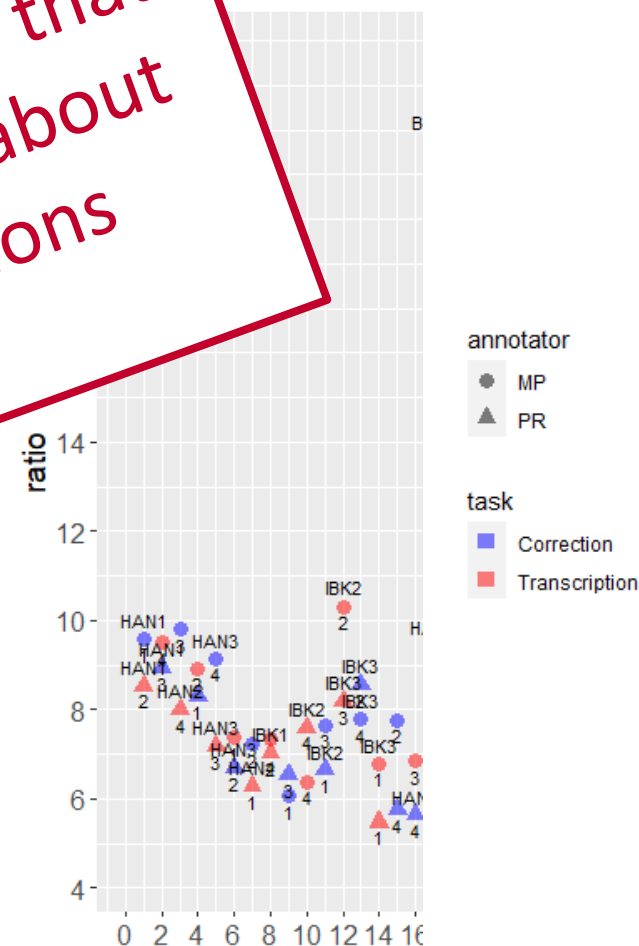


annotator effect plot



Difficult to report a non-effect, but it seems that:
Correcting and editing ASR-output takes about
the same effort as creating transcriptions
“manually” from scratch.

| Factor | | z | p | |
|-------------|---------|---------|---------|---------|
| (Intercept) | 19.240 | < 2e-16 | *** | |
| task | 0.3705 | -0.298 | 0.76780 | |
| order | 0.0618 | -3.115 | 0.0043 | ** |
| placeIBK | 0.5344 | 0.5685 | 0.3555 | |
| annotatorPR | -0.8195 | 0.3692 | -2.220 | 0.035 * |

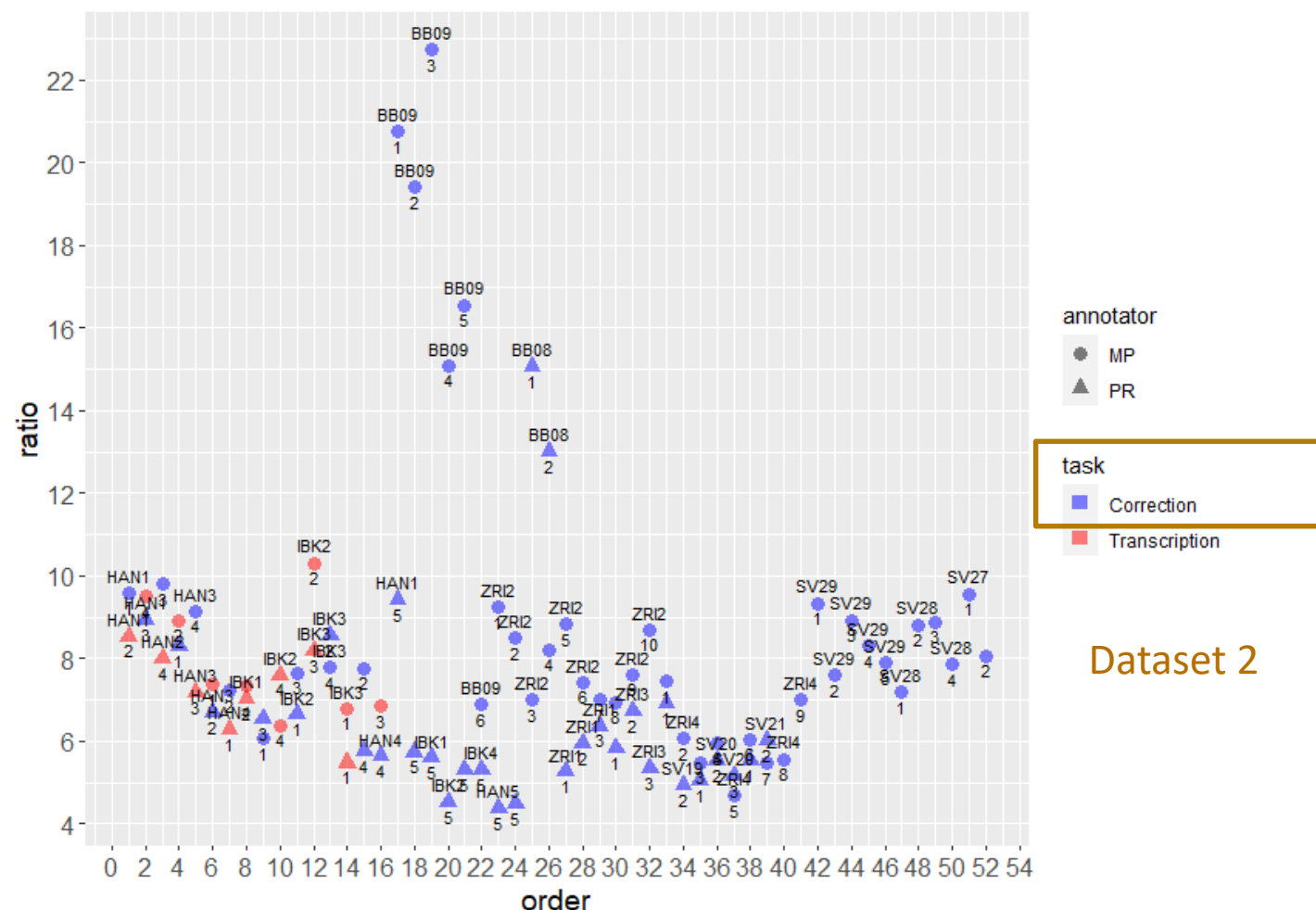


Adjusted R-squared = 0.33

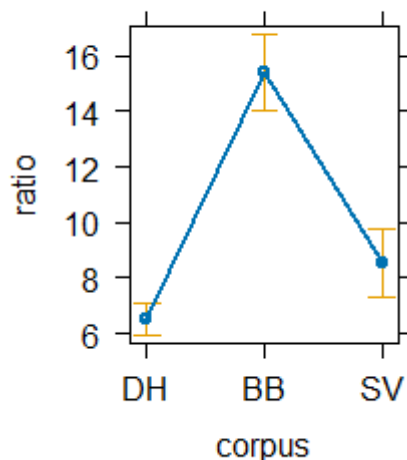
CORRECTING ONLY

Linear regression model

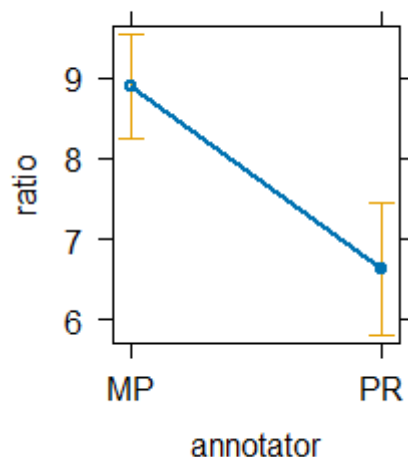
- Criterion variable:
 - ratio
- Predictors:
 - Corpus
 - Recording place
 - Annotator
 - Order



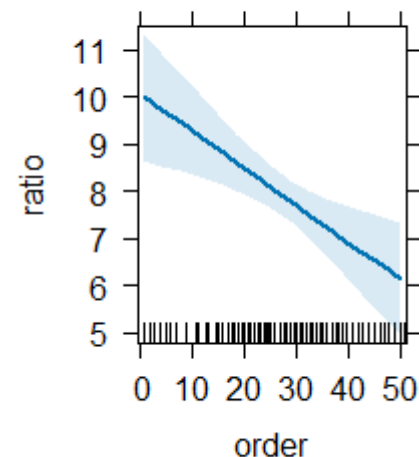
corpus effect plot



annotator effect plot



order effect plot



| Factor | Estimate | Std. Error | z | p | |
|-------------------|----------|------------|--------|---------|-----|
| (Intercept) | 9.451 | 0.691 | 13.679 | < 2e-16 | *** |
| corpusBB | 8.888 | 0.7327 | 12.131 | < 2e-16 | *** |
| corpusSV | 2.0189 | 0.7552 | 2.637 | 0.0094 | ** |
| annotatorPR | −2.2656 | 0.5945 | −3.811 | 0.0003 | *** |
| order | −0.0787 | 0.0244 | −3.232 | 0.0019 | ** |
| audio_length_in_s | −0.0004 | 0.00098 | 0.374 | 0.7094 | |

Adjusted R-squared = 0.75

QUALITATIVE RESULTS

COMMENTS FROM ANNOTATORS

| Stretch | Comment |
|---------|---|
| HAN2-1 | Speakers are extremely often on the wrong tier. Much overlap. |
| HAN3-2 | Again often the wrong tier. The interviewee speaks rather clear. |
| HAN4-1 | Little overlap. Missing uhm and yes, as always. |
| HAN4-4 | Extremely often wrong tier. |
| IBK1-2 | Extremely often overlap. |
| IBK1-3 | No opportunity is missed to overlap. On the contrary, the dialect is surprisingly well transcribed. |
| IBK2-4 | As above; and the interviewee speaks extremely unclear and quiet. |
| IBK2-1 | Again missing uhm and yes in overlap. |
| IBK3-3 | The interviewee speaks a lot. |
| IBK3-2 | Very often wrong tier; repetitions are almost never detected. |
| IBK4-1 | Little talk and little overlap. |
| IBK1-5 | The alignment of the last minute was completely off; the first time this kind of error. |
| HAN3-5 | More often the wrong tier than in all previous transcripts. |
| HAN4-5 | Both speak very clear and in longer sequences. |
| BB08 | Correction might have taken longer as more listening was necessary. |
| SV19-1 | Despite strong dialect, easy to transcribe because of little overlap. |
| SV22-1 | One speaker speaks very unclearly and the interviewer is extremely quiet. |
| SV26-1 | Possibly shorter correction time, as much was simply incomprehensible. |

Example from Südwestdeutschland and Vorarlberg

Ref: wir sind zuerst natürlich selbstständiger Konsum gewesen

Ref: we have first of course autonomous cooperative been

'Ref: at the beginning we have been an autonomous cooperative of course'

Hyp: Wir waren zuerst natürlich selbstständiger Konsum

Hyp: we were first of course autonomous cooperative

'Hyp: first we were an autonomous cooperative of course'

Cut Excerpt:

SV--_E_00019 from 08:25.34 to 08:28.70

And add 1 s silence at the beginning

ANECDOTES II

ON AMBERSCRIPT

Despite doing speaker recognition, Amberscript does not seem to be able to deal with:
Genus-accord.

Ségolène Royale (feminin) said:

- “je ne suis pas sortie”

Amberscript transcribes (masculin):

- “je ne suis pas sorti”



CONCLUSIONS & FUTURE WORK

Problem: Transcription bottleneck

Solution: ASR

Improve and optimize the solution:

- Exploit prompting mechanism
- Avoid some normalizations
- Different post-processing
- Adapt system to data (fine-tuning)

Adjust or redefine the problem:

- Analyze the imperfect output
 - Use additional ASR information, e.g. probability of word-detection
1. Generic query on ASR-transcribed corpus
 2. Refine manually the results

Record audio as close to the (individual) speaker as possible

fin



EVOLVING MODELS... (LARGE) V1, V2, V3 AND AMOUNTS OF TRAINING DATA

- large-v1: Training Data: 13344 hours for German

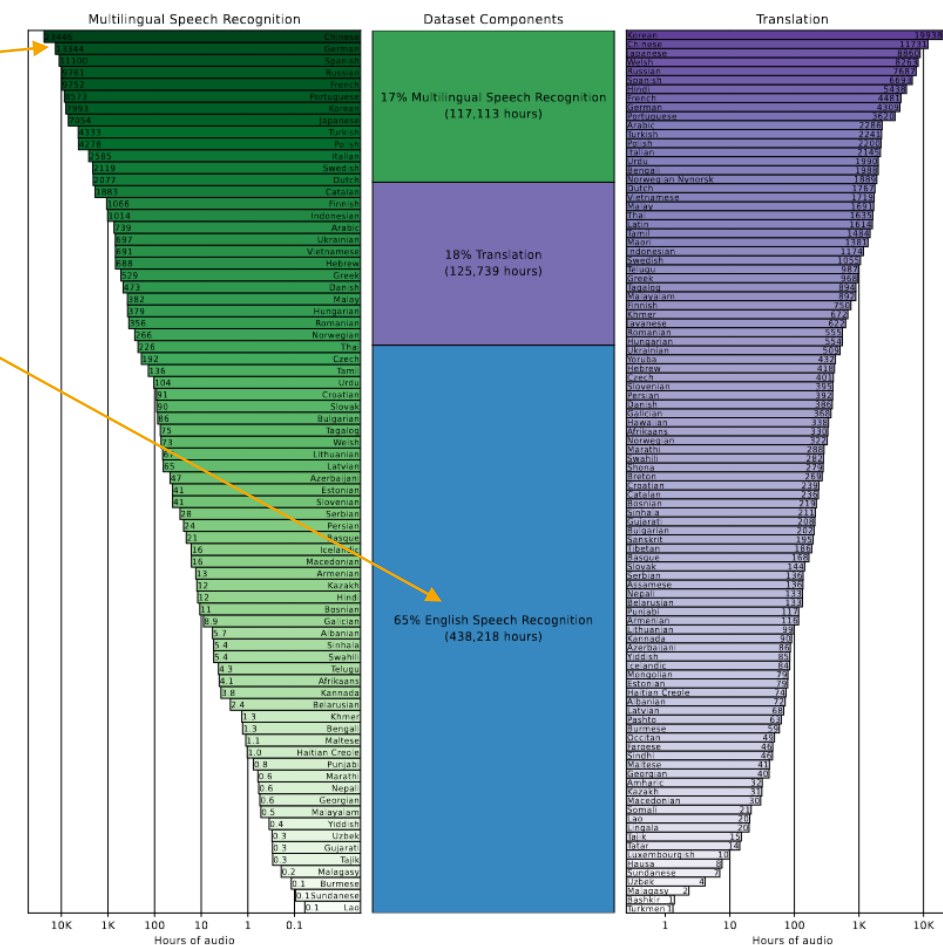
13344

German

65% English Speech Recognition
(438,218 hours)

- large-v2: Compared to the Whisper large model, the large-v2 model is trained for 2.5x more epochs with added regularization
- large-v3:
 - The input uses 128 Mel frequency bins instead of 80
 - 1 million hours of weakly labeled audio and 4 million hours of pseudolabeled audio collected using large-v2.

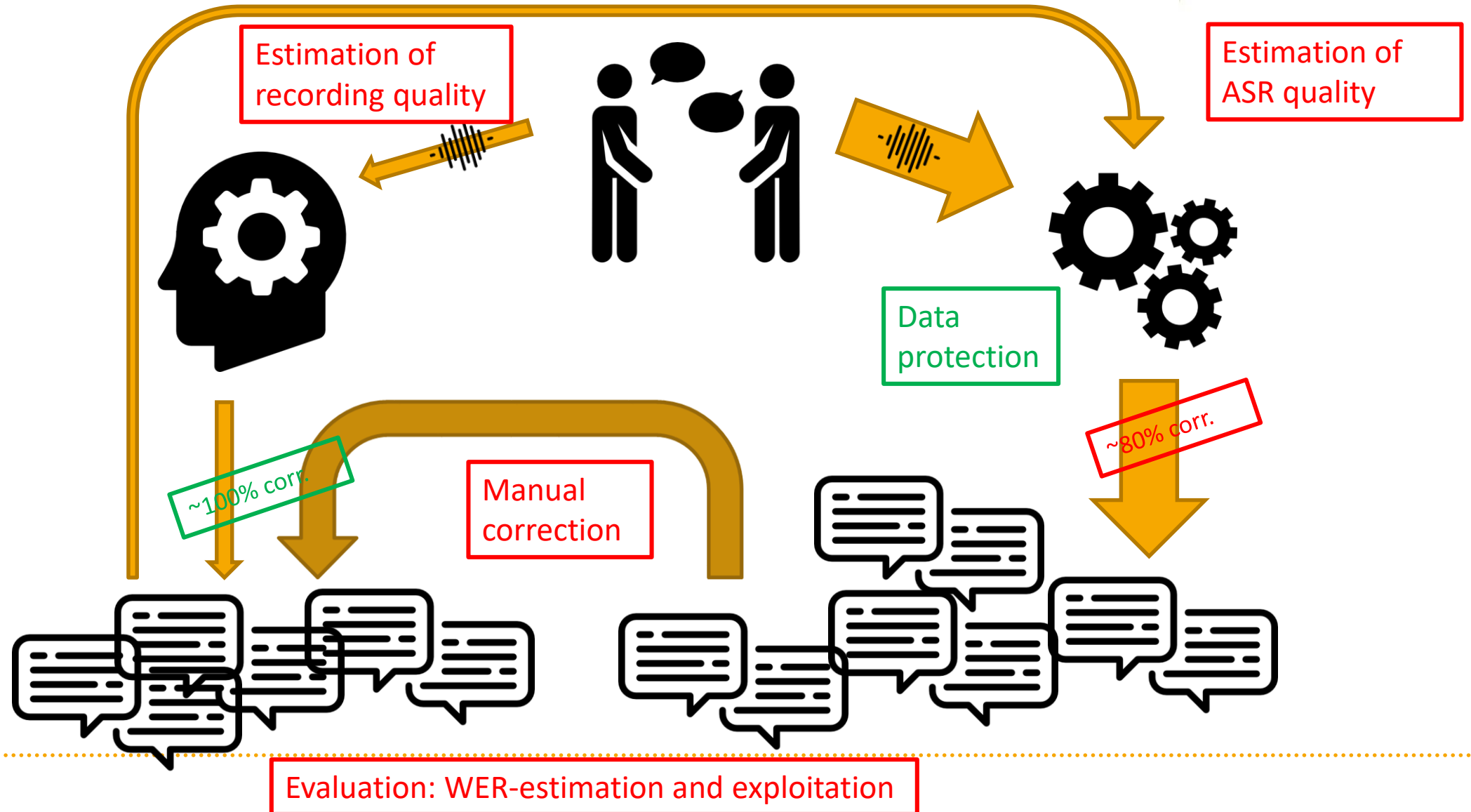
E. Training Dataset Statistics



CHALLENGES / IDEAS

IDS

LEIBNIZ-INSTITUT FÜR
DEUTSCHE SPRACHE



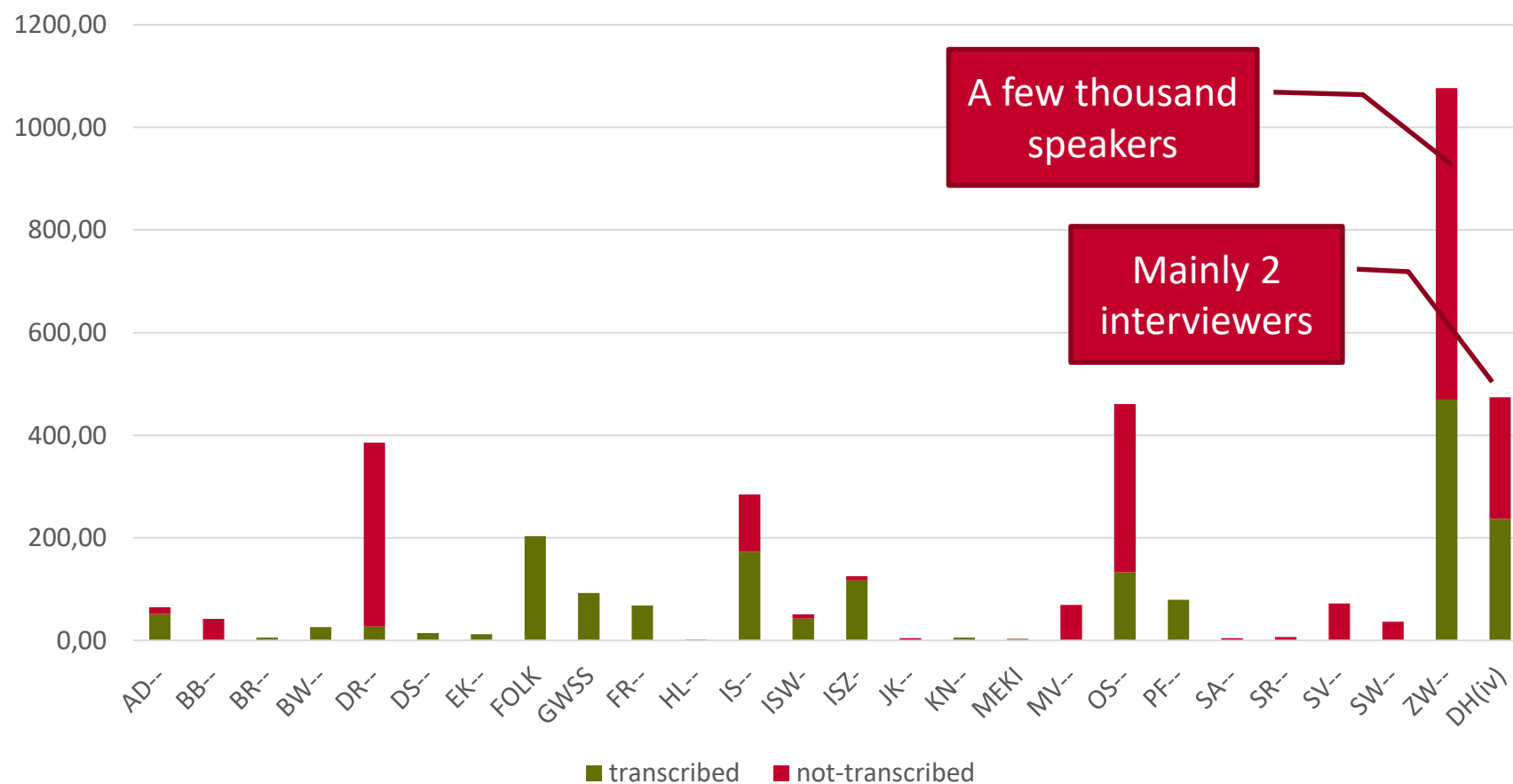
- Who knows of literature on “average time spent for (orthographic) transcription”
 - According to types of transcription conventions (cGAT, normalized)
 - Is there any literature on correcting ASR-output? Any tools for that specific purpose?
 - Are there open (speech) recognizers around that can also recognize/output ...
 - Backchannels
 - Hesitations
 - Repetitions
 - Annotations of non-verbal sounds / sounds of background noise
 - Dialectal speech from old speakers recorded in the 1960s in rural areas?
 - Which tools can be used for evaluating/benchmarking ASR-output (WER, Insertions, Deletions, Substitutions)
 - Kaldi, HResult in HTK-3, ..., ??diff in Gitlab??
-

REFERENCE(S)

Brinckmann, Caren (2009): Transcription Bottleneck of Speech Corpus Exploitation. In: Lyding, Verena (Ed.): *LULCL II 2008 - Proceedings of the Second Colloquium on Lesser Used Languages and Computer Linguistics*. Bozen-Bolzano, 13th-14th November 2008. Bozen-Bolzano: EURAC. 165-179.



transcription status per corpus



METHOD FOR THE

- 5-minute chunks
- Two student helpers (Maja Peer & Paul Rölle)
- Two tasks (T & K)
 - T(ranscribing) from scratch
 - C(orrecting) ASR-output
 - ✓ speaker association
 - ✓ content
 - ✓ alignment
- Training (30 Min. for C; 20 Min. for T)
- Measure: working time (stopwatch)
- Metric: ratio of $\frac{\text{audio length}}{\text{working time}}$

| task | 1 st chunk | 2 nd chunk | 3 rd chunk | 4 th chunk | 5 th chunk | interviewer | interviewee/ rec. name |
|------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-------------|---------------------------|
| C | Maja | | Paul | | | NL | HAN1 |
| T | | Paul | | Maja | | | |
| C | Paul | | Maja | | | NL | HAN2 |
| T | | Maja | | Paul | | | |
| C | | Paul | | Maja | | NL | HAN3 |
| T | Maja | | Paul | | | | |
| C | | Maja | | Paul | | NL | HAN4 |
| T | Paul | | Maja | | | | |

Hannover

| task | 1 st chunk | 2 nd chunk | 3 rd chunk | 4 th chunk | 5 th chunk | interviewer | interviewee/ rec. name |
|------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-------------|---------------------------|
| C | Maja | | Paul | | | MF | IBK1 |
| T | | Paul | | Maja | | | |
| C | Paul | | Maja | | | MF | IBK2 |
| T | | Maja | | Paul | | | |
| C | | Paul | | Maja | | MF | IBK3 |
| T | Maja | | Paul | | | | |
| C | | Maja | | Paul | | MF | IBK4 |
| T | Paul | | Maja | | | | |

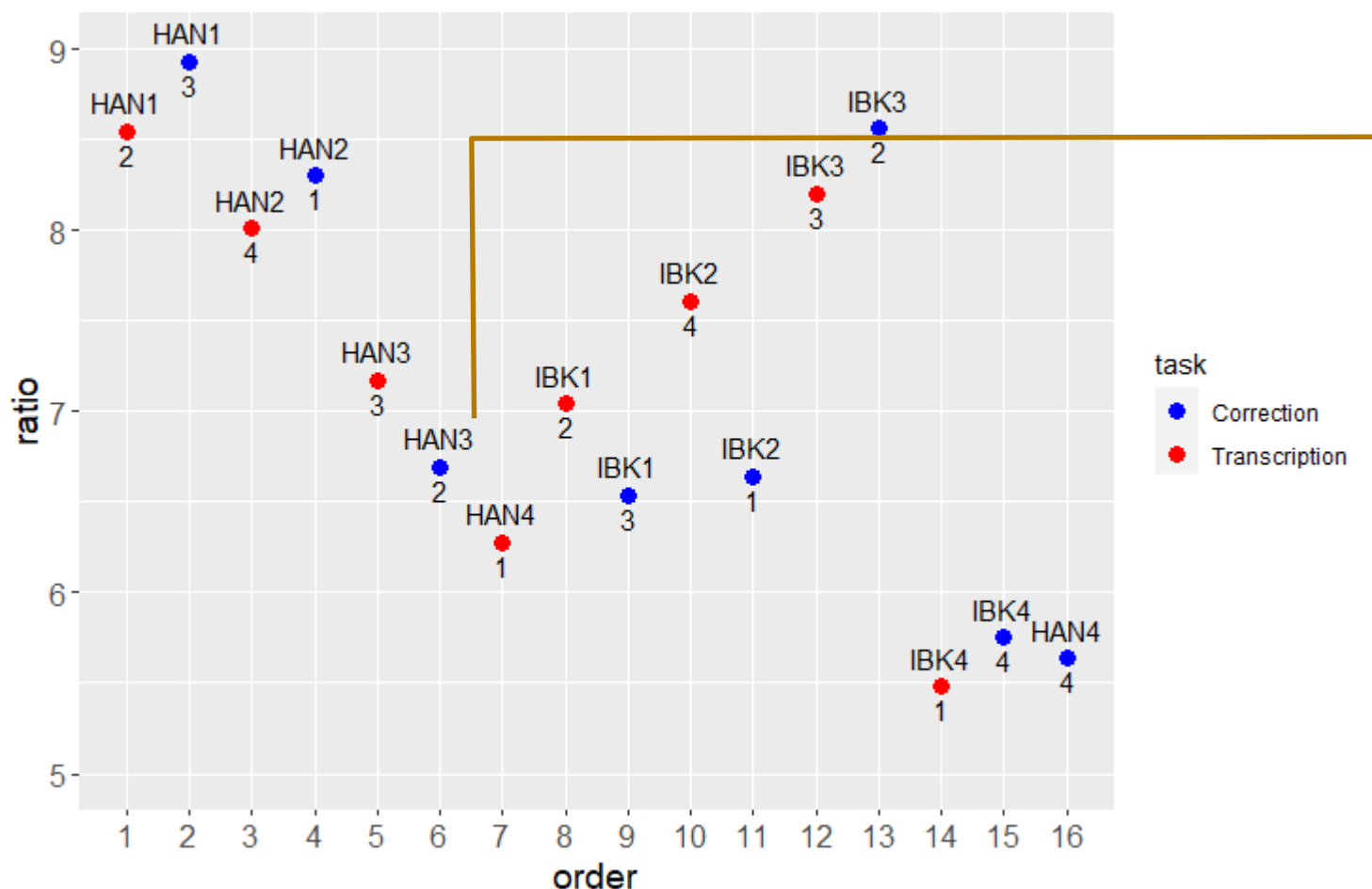
Innsbruck

COMMENTS FROM ANNOTATOR

Paul (after training):

- “**Transcribing** is much more pleasant.”
- “it is easier to get into a rhythm/flow.”
- “At the **correcting** task one has to concentrate on several things at once.”
- “The alignment is often slightly off”
- “In most cases, backchannels and hesitation markers are missing”

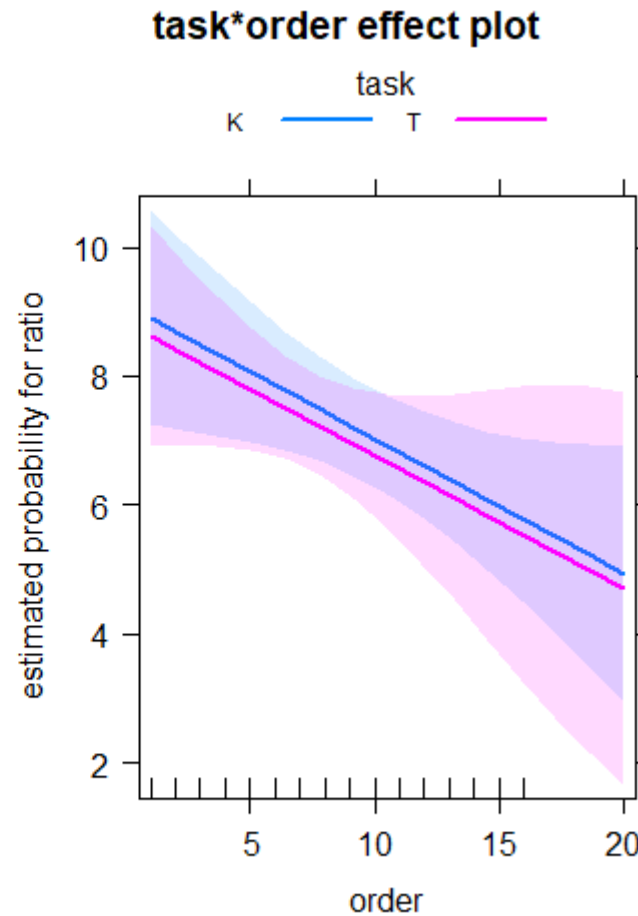
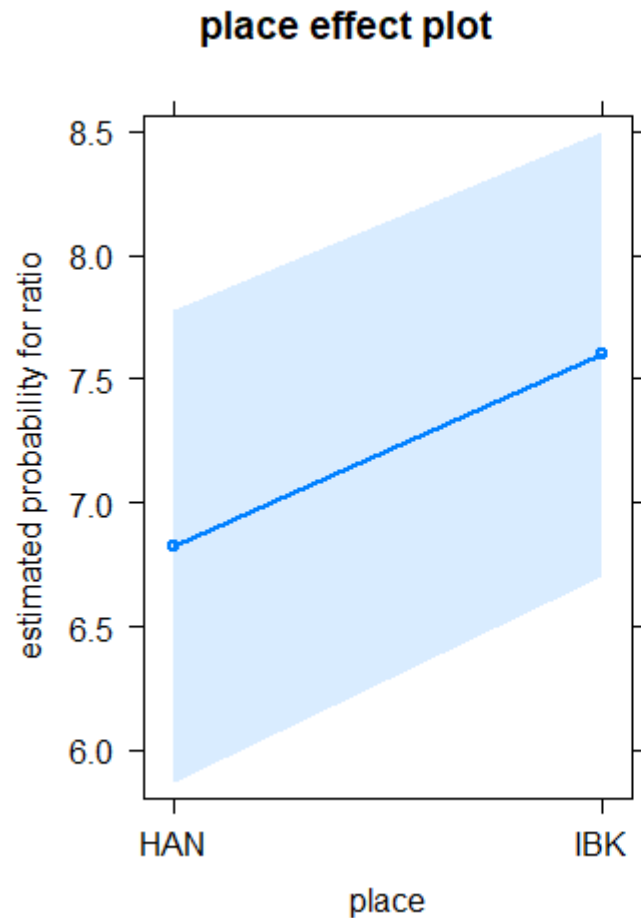
RESULTS FROM ANNOTATOR PAUL



Paul (after 6 files): “in the meantime, I also got into a flow with **correcting**. I now know what to look for and thus have a better rhythm.”

RESULTS

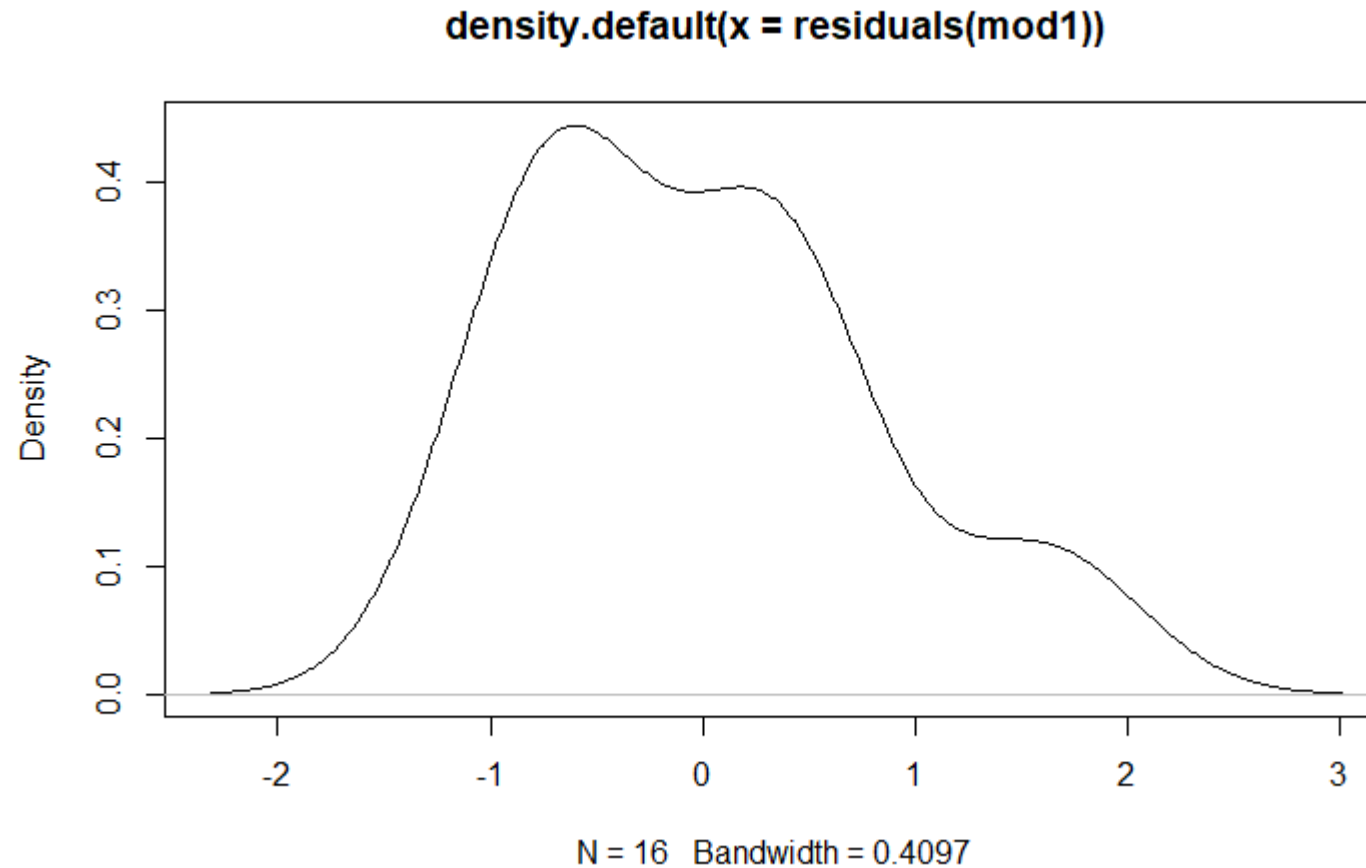
```
mod <- lm(ratio ~ task * order + place, data = df)
```



- Both **transcribing** and **correcting** seem to get faster over time
- Independent from the recording place (Hannover or Innsbruck)

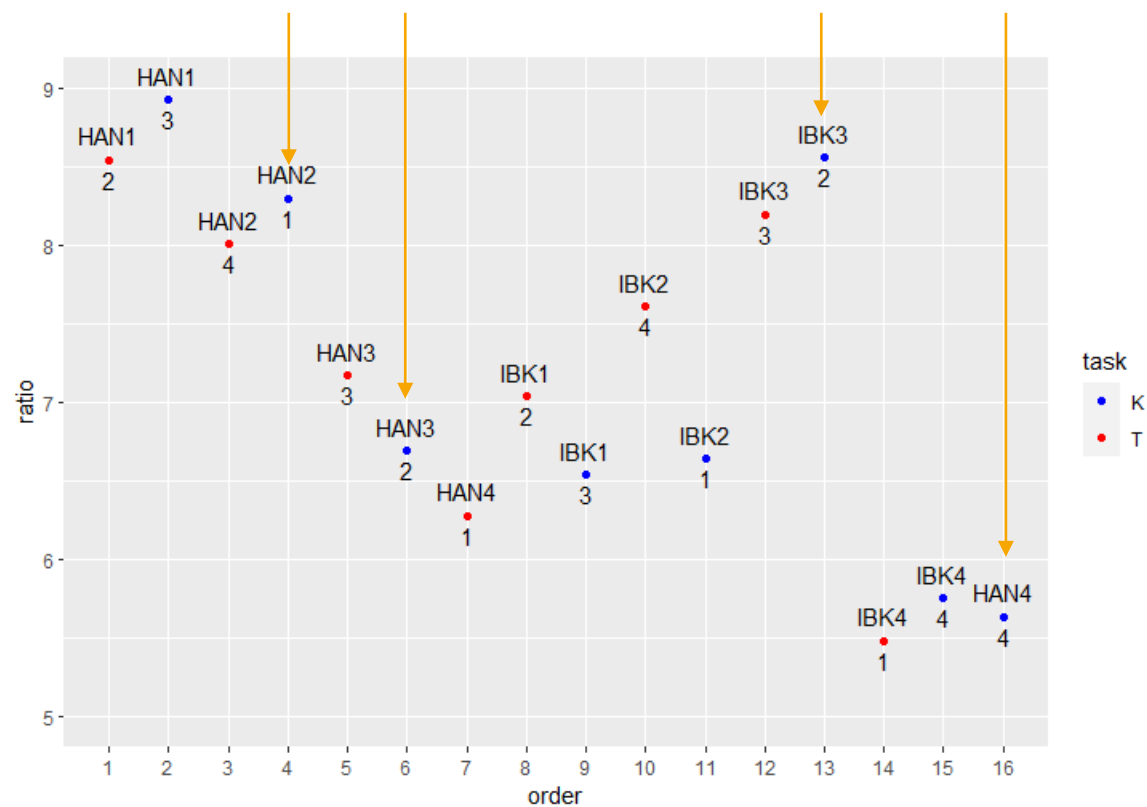
DISTRIBUTION OF RESIDUALS

```
plot(density(residuals(mod1)))
```

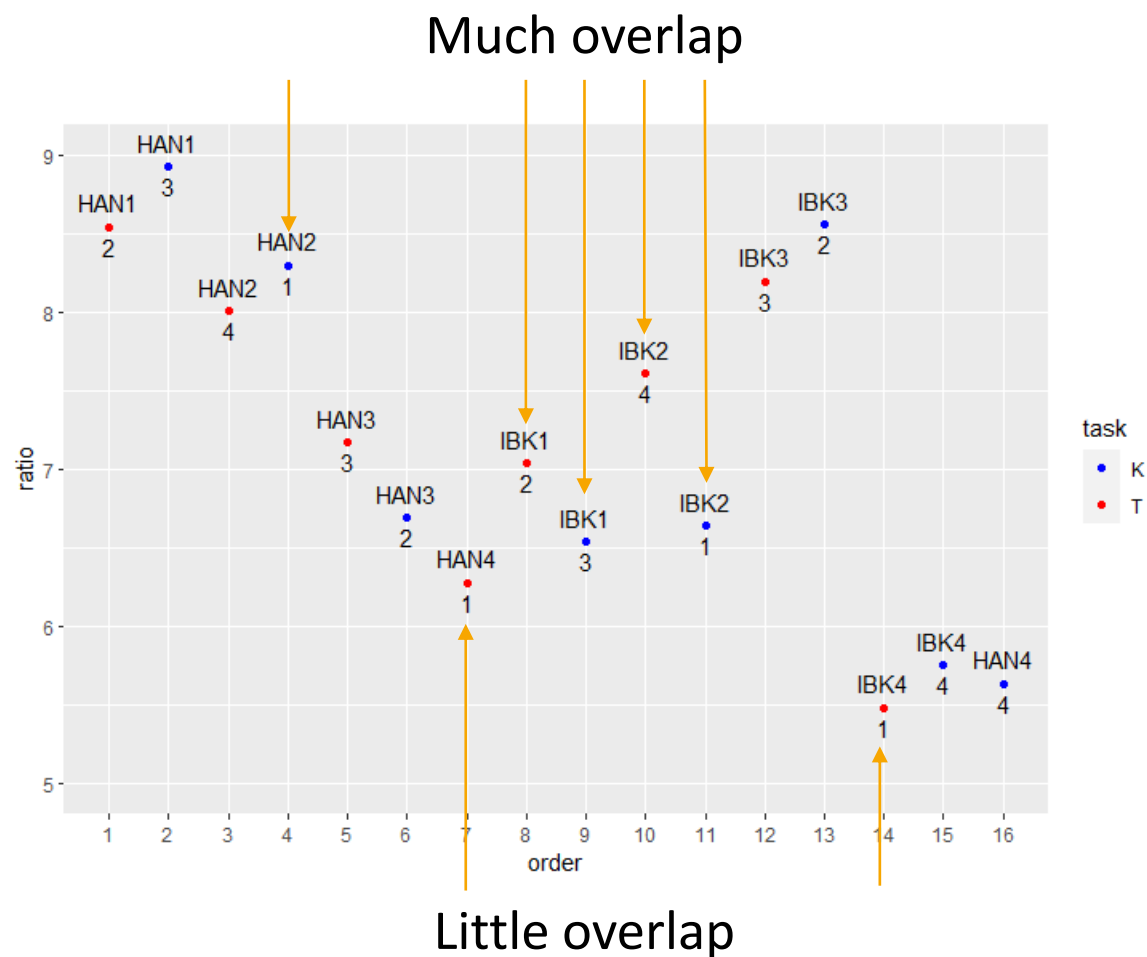


- Not really normally distributed
- !!! only 16 data points
- ??? we'll see with the next 16

Speakers are very often on the wrong tier

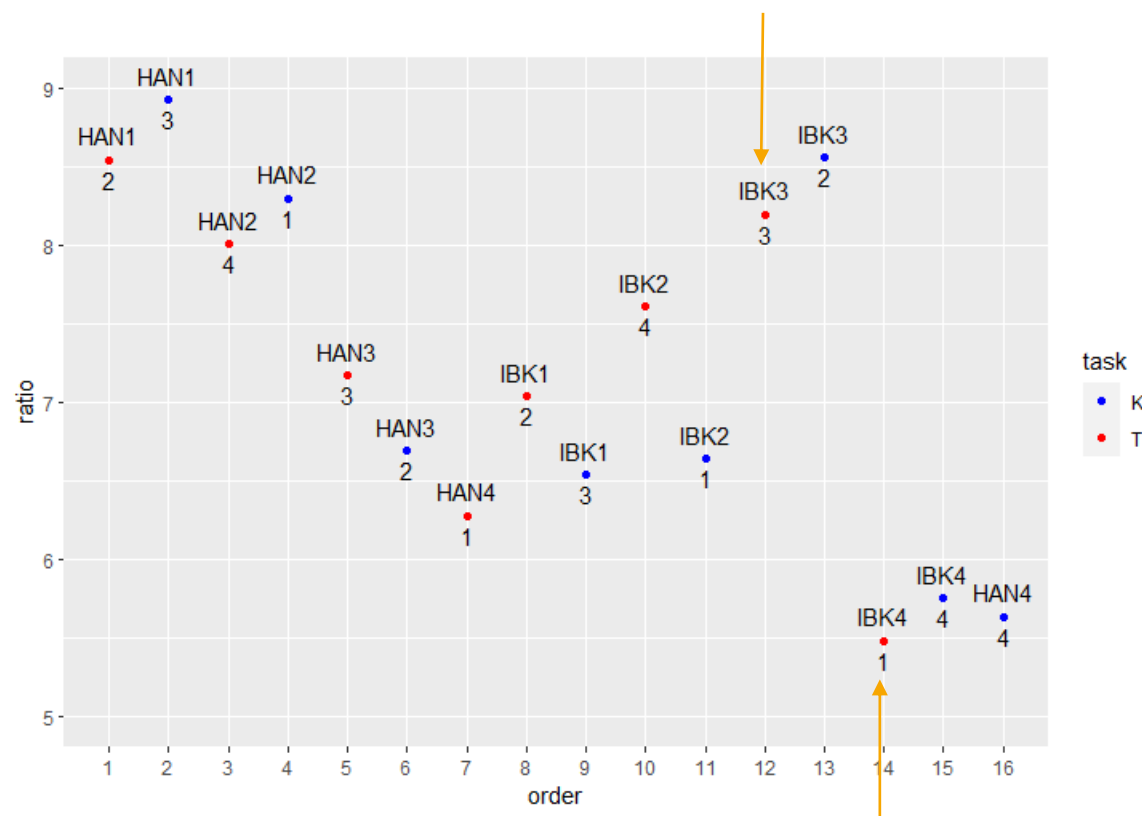


COMMENTS FROM ANNOTATOR



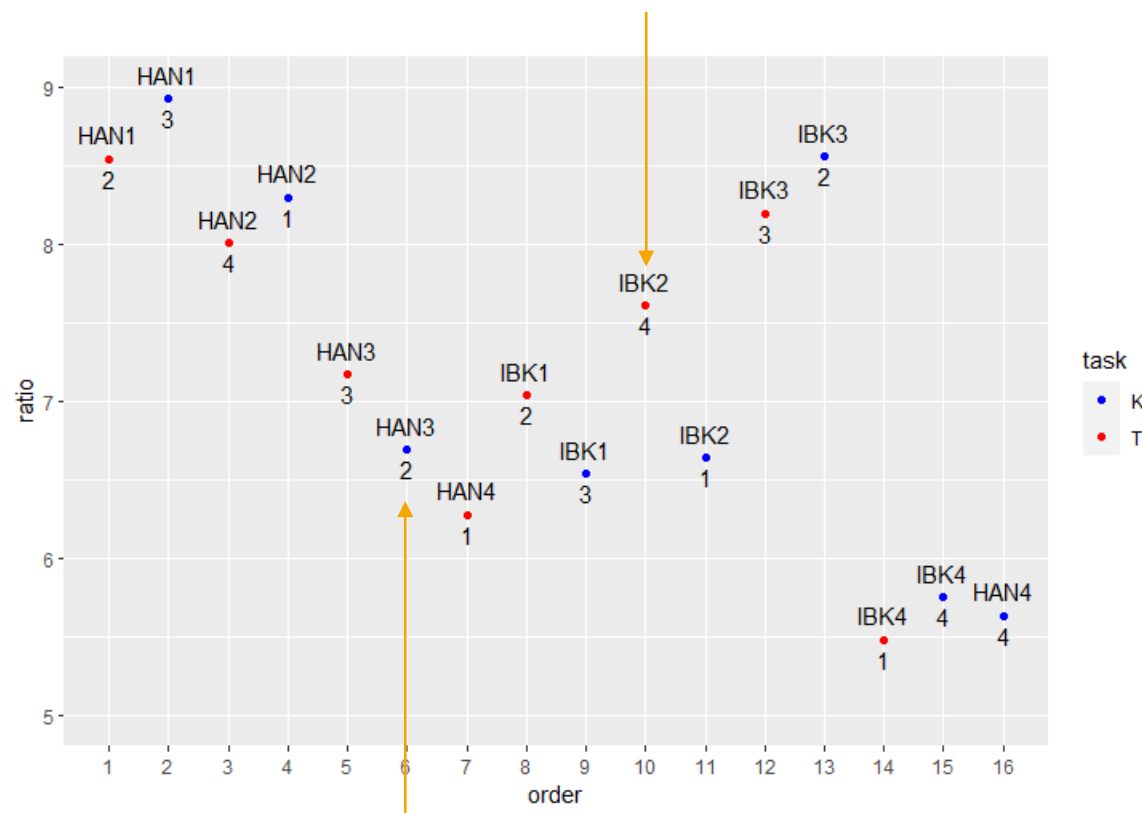
COMMENTS FROM ANNOTATOR

The interviewee speaks a lot



Speakers don't speak much

The interviewee speaks extremely unclear and quiet



The interviewee speaks rather clear

Word repetitions are almost never recognized by the ASR system

