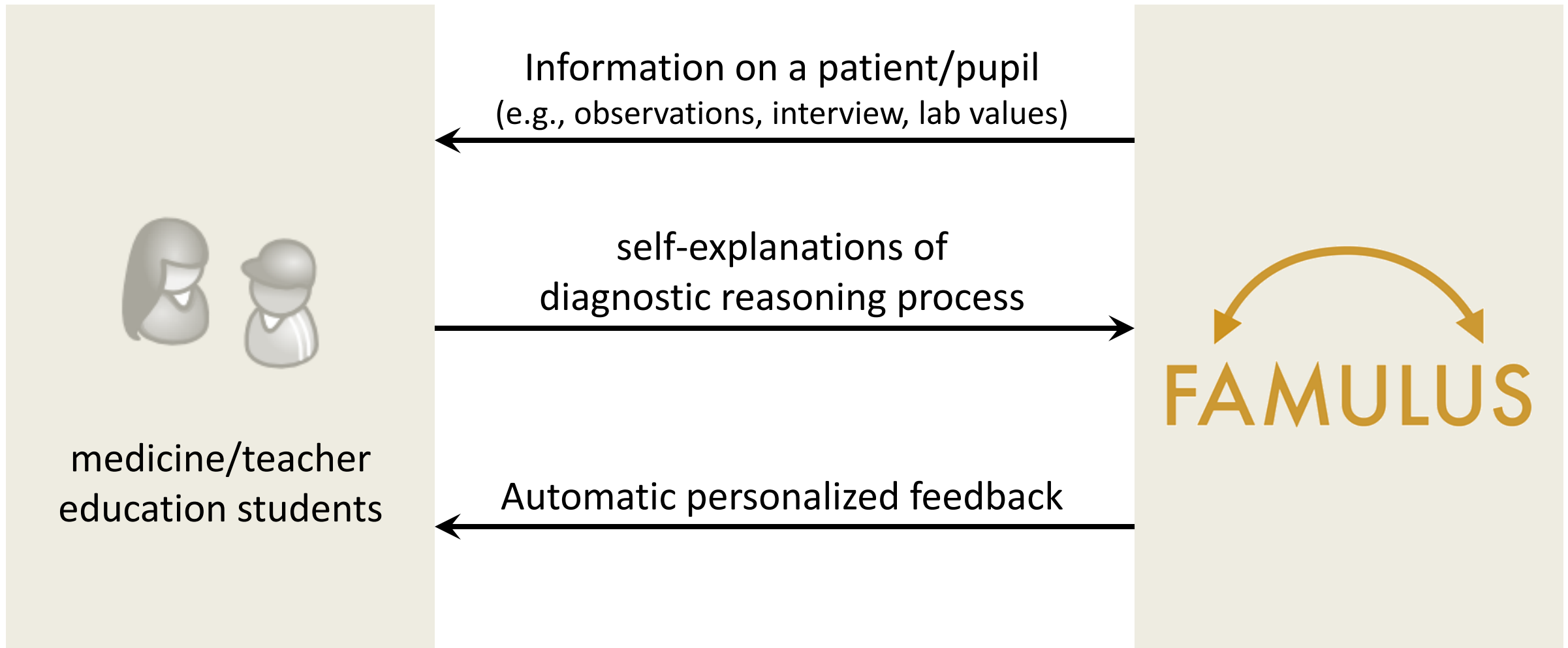# Analysis of Automatic Annotation Suggestions for Hard Discourse-Level Tasks in Expert Domains

**Claudia Schulz, Christian M. Meyer, Jan Kiesewetter, Michael Sailer,
Elisabeth Bauer, Martin R. Fischer, Frank Fischer, and Iryna Gurevych**

FAMULUS

http://famulus-project.de

UKP

# Learning to Diagnose

Information on a patient/pupil
(e.g., observations, interview, lab values)

self-explanations of
diagnostic reasoning process

medicine/teacher
education students

Automatic personalized feedback

FAMULUS

# Diagnostic Reasoning

The patient reports to be lethargic and feverish. From the anamnesis I learned that he had purulent tonsilitis and is still suffering from symptoms. I first performed some laboratory tests and notice the decreased number of lymphocytes, which can be indicative of a bone marrow disease or an HIV infection. The HIV test is positive. However, the results from the blood cultures are negative, so it is a virus, parasite, or a fungal infection causing the symptoms.

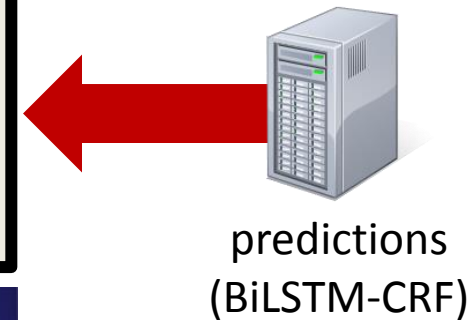| | |
|---|---|
| Hypothesis generation (HG) | Evidence evaluation (EE) |
| Drawing conclusions (DC) | Evidence generation (EG) |

# Research Question

The patient reports to be lethargic and feverish. From the anamnesis I learned that he had purulent tonsilitis and is still suffering from symptoms. I first performed some laboratory tests and notice the decreased number of lymphocytes, which can be indicative of a bone marrow disease or an HIV infection. The HIV

**How can we improve this hard and time-consuming annotation task?**

the symptoms.

| | |
|---|---|
| Hypothesis generation (HG) | Evidence evaluation (EE) |
| Drawing conclusions (DC) | Evidence generation (EG) |

# Annotation Suggestions

The patient reports to be lethargic and feverish. From the anamnesis I learned that he had purulent tonsilitis and is still suffering from symptoms. I first performed some laboratory tests and notice the decreased number of lymphocytes, which can be indicative of a bone marrow disease or an HIV infection. The HIV test is positive. However, the results from the blood cultures are negative, the symptoms... infection causing the sympt...

Suggestion:
**Hypothesis generation (HG)**

✓ Accept     ✗ Reject

INCEpTION annotation platform
https://inception-project.github.io

INCEpTION

predictions
(BiLSTM-CRF)

UKP

# Training Data and Suggestion Quality

# Training Data and Suggestion Quality

S1    S2

A1

A2

A3

A4

A5

{EG, EE, HG, DC}

↑

CRF

↑

BiLSTM

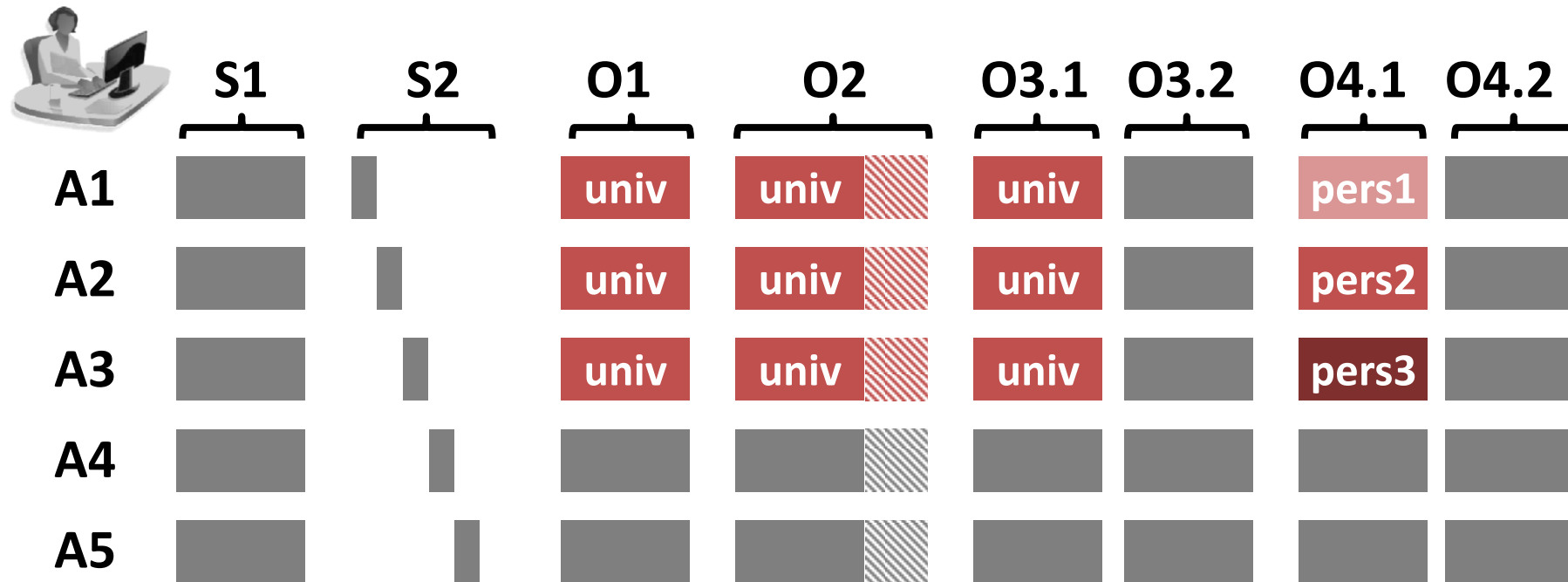$\Longrightarrow$

**univ(ersal) model: $F_1 \approx .63$**
**pers(onalized) models: $F_1 \approx .55$**
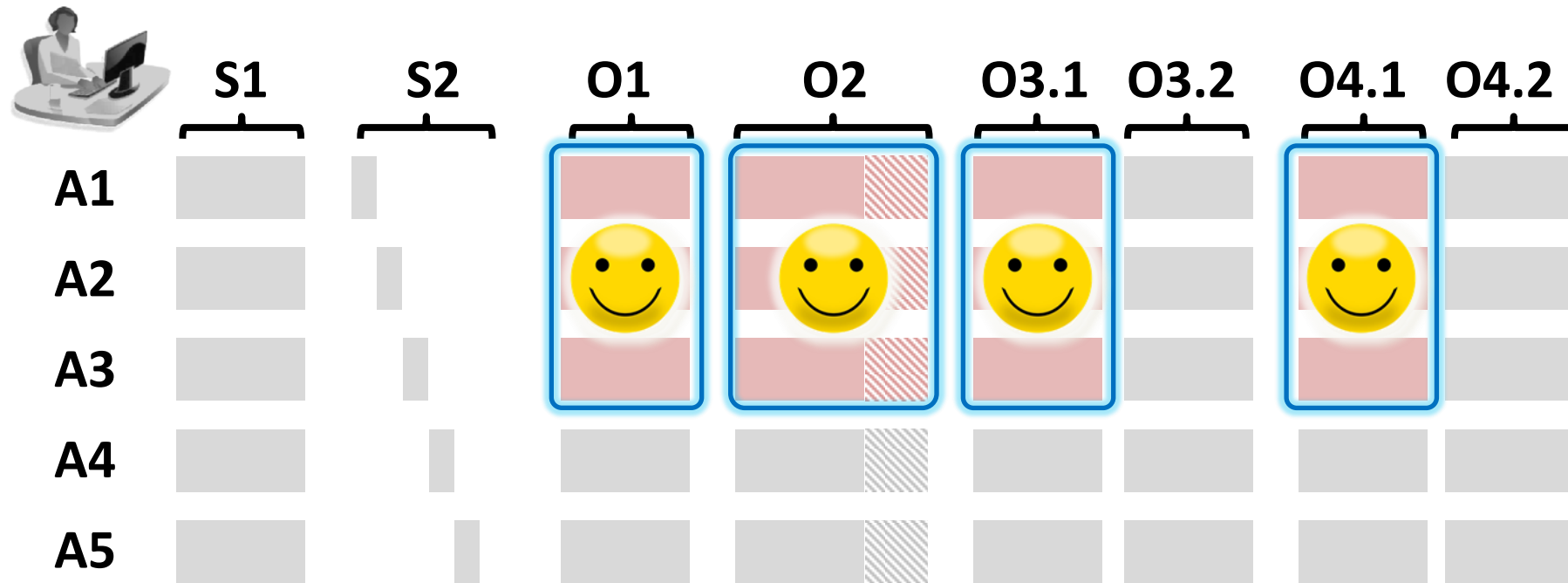
# Effectiveness of Annotation Suggestions

# Effectiveness of Annotation Suggestions



*presentation: focus on medical use case*
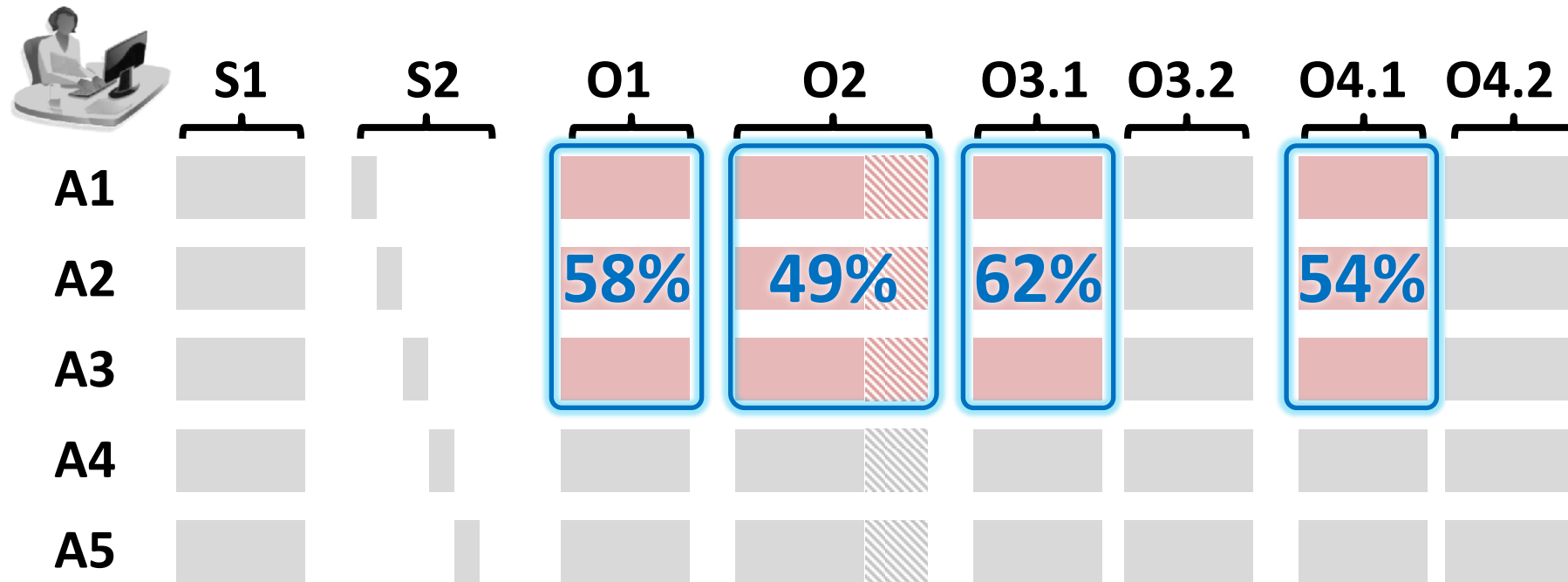*paper: results for same setup in teacher education!*

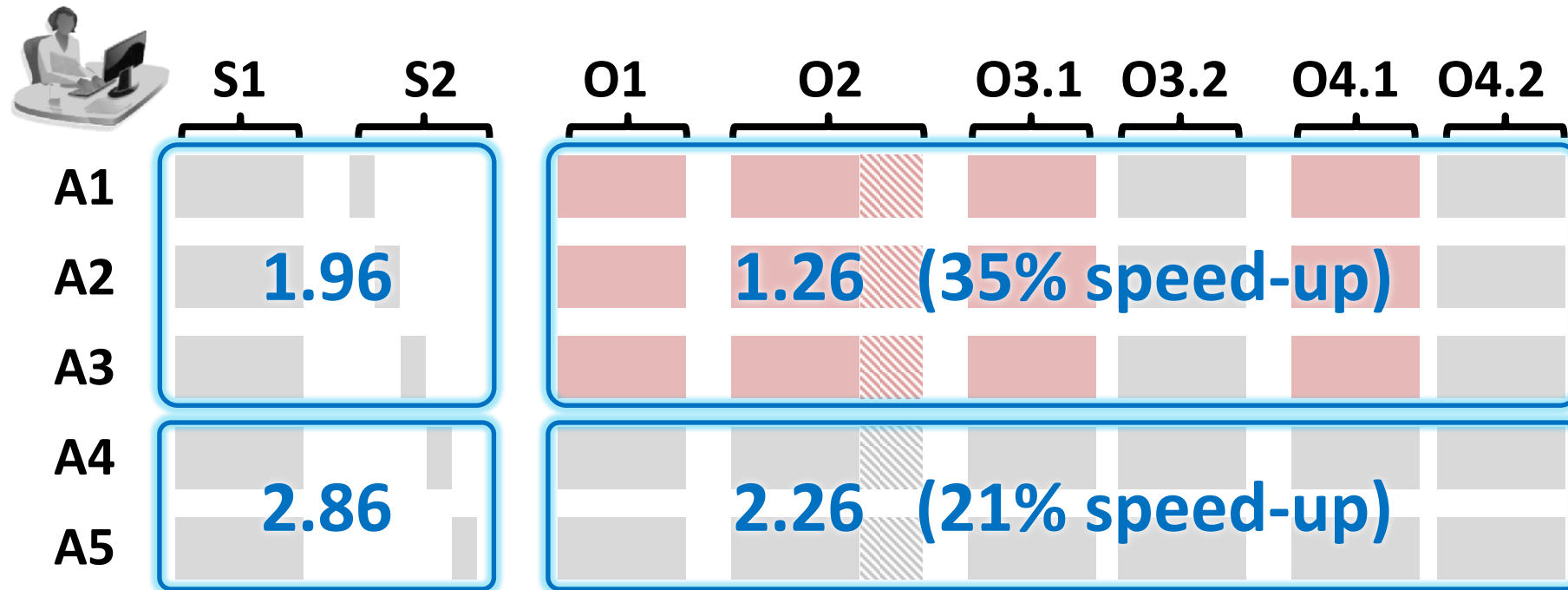# Usefulness of Annotations

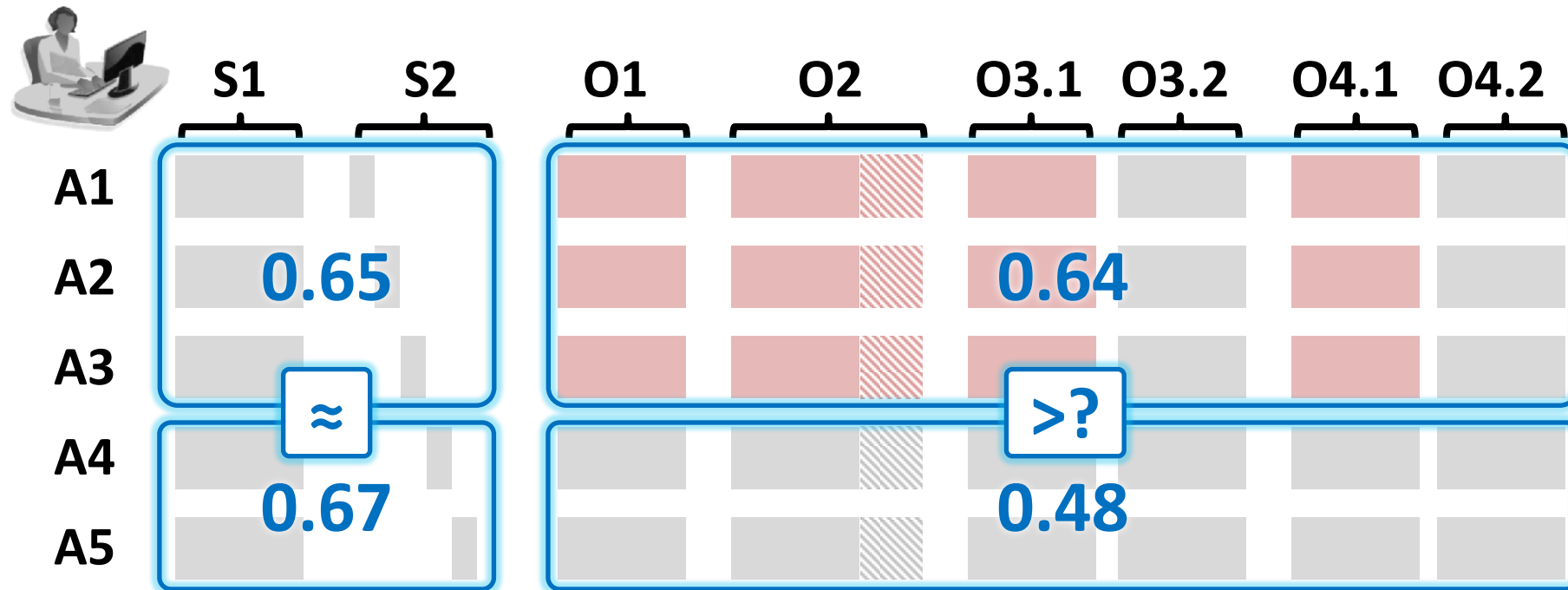**[Annotator happiness]**

# Usefulness of Annotations

# Annotation Time



[Minutes per text]

|  | S1 | S2 | | O1 | O2 | O3.1 | O3.2 | O4.1 | O4.2 |
|---|---|---|---|---|---|---|---|---|---|
| A1 | | | | | | | | | |
| A2 | 1.96 | | | 1.26 (35% speed-up) | | | | | |
| A3 | | | | | | | | | |
| A4 | | | | | | | | | |
| A5 | 2.86 | | | 2.26 (21% speed-up) | | | | | |

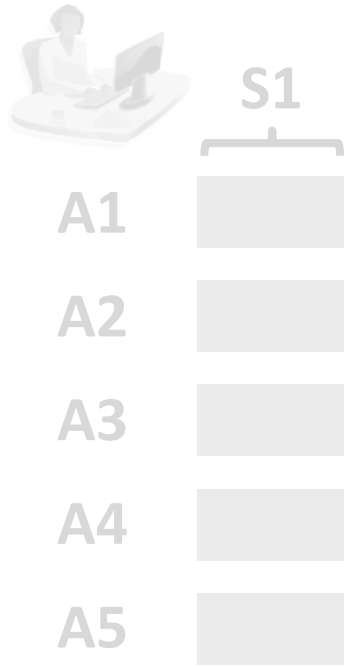# Reliability of Annotations

[Krippendorff's α]

# Effects of Annotation Suggestions



Conclusion 1: Annotation suggestions are helpful for experts
and yield faster and (maybe) more reliable annotations!

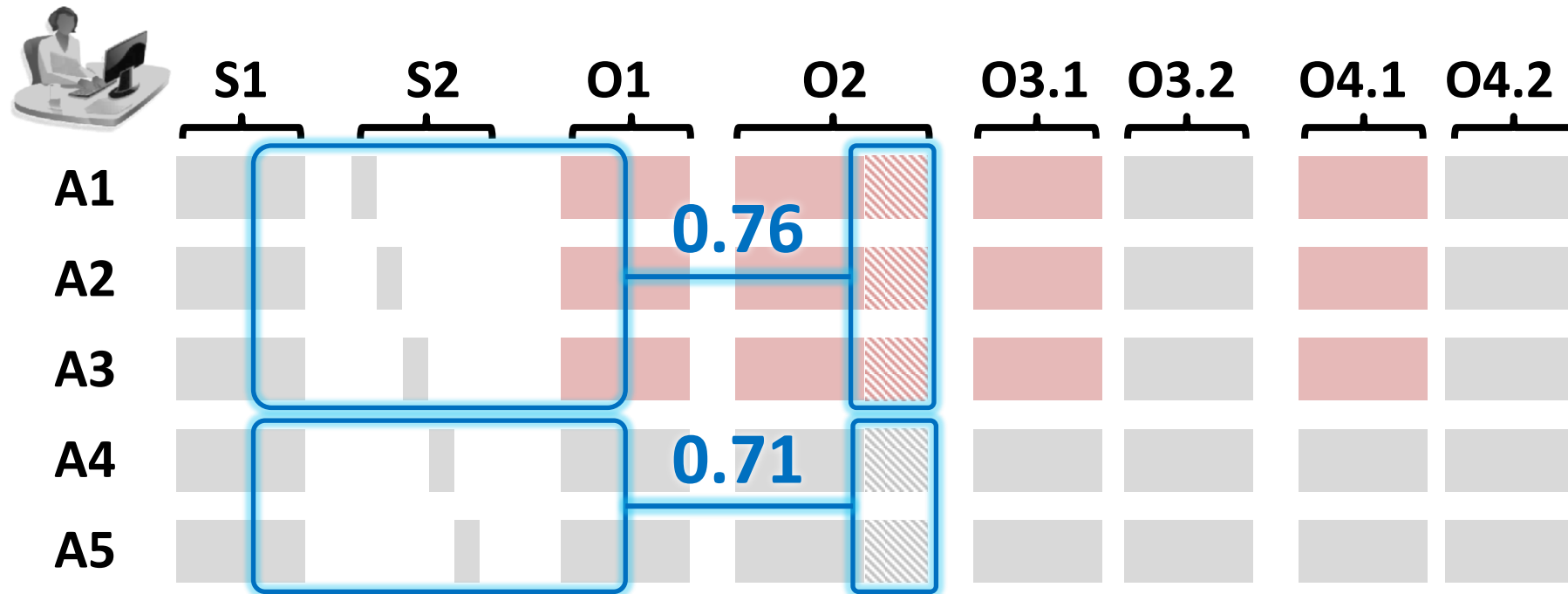# But: Do predictions bias the decisions?

# Intra-Annotator Consistency

[Krippendorff's $\alpha$]

# Human–Machine Agreement

# Human–Machine Agreement

|  | S1 | S2 | O1 | O2 | O3.1 | O3.2 | O4.1 | O4.2 |
|---|---|---|---|---|---|---|---|---|
| A1 |  |  |  |  |  |  |  |  |
| A2 | 0.65 |  | 0.64 | 0.55 | 0.69 | 0.52 | 0.42 | 0.41 |
| A3 |  |  |  |  |  |  |  |  |
| A4 |  |  |  |  |  |  |  |  |
| A5 | 0.67 |  | 0.56 | 0.48 | 0.55 | 0.45 | 0.30 | 0.45 |

# Further Analysis of Annotation Bias

- **Pairwise agreement between the A1–A3 and the A4–A5 groups**
  - → A1–A3 do not behave differently than A4–A5

- **Distribution of labels**
  - → no systematic difference

# Further Analysis of Annotation Bias

- **Pairwise agreement between the A1–A3 and the A4–A5 groups**
  - → A1–A3 do not behave differently than A4–A5

- **Distribution of labels**
  - → no systematic difference

- **Distribution of disagreements**
  - → only small differences

|    | EG | EE | DC | HG |
|----|----|----|----|----|
| EG | -  | 7% | 1% | 0% |
| EE | 7% | -  | 22%| 13%|
| DC | 1% | 22%| -  | 7% |
| HG | 0% | 13%| 7% | -  |

*with suggestions A1–A3*

|    | EG | EE | DC | HG |
|----|----|----|----|----|
| EG | -  | 5% | 1% | 2% |
| EE | 5% | -  | 21%| 14%|
| DC | 1% | 21%| -  | 8% |
| HG | 2% | 14%| 8% | -  |

*without suggestions A4–A5*

# Effects of Annotation Suggestions



|  | S1 | S2 | O1 | O2 | | O3.1 | O3.2 | O4.1 | O4.2 |
|----|----|----|----|----|----|----|----|----|----|
| A1 |  |  | univ | univ |  | univ |  | pers1 |  |
| A2 |  |  | univ | univ |  | univ |  | pers2 |  |
| A3 |  |  | univ | univ |  | univ |  | pers3 |  |
| A4 |  |  |  |  |  |  |  |  |  |
| A5 |  |  |  |  |  |  |  |  |  |

**Conclusion 2: Some evidence for annotation bias, but negligible, as no systematic discrepancy compared to the control setup!**

# Interactive Model Training

*model*

*training data*

*i*

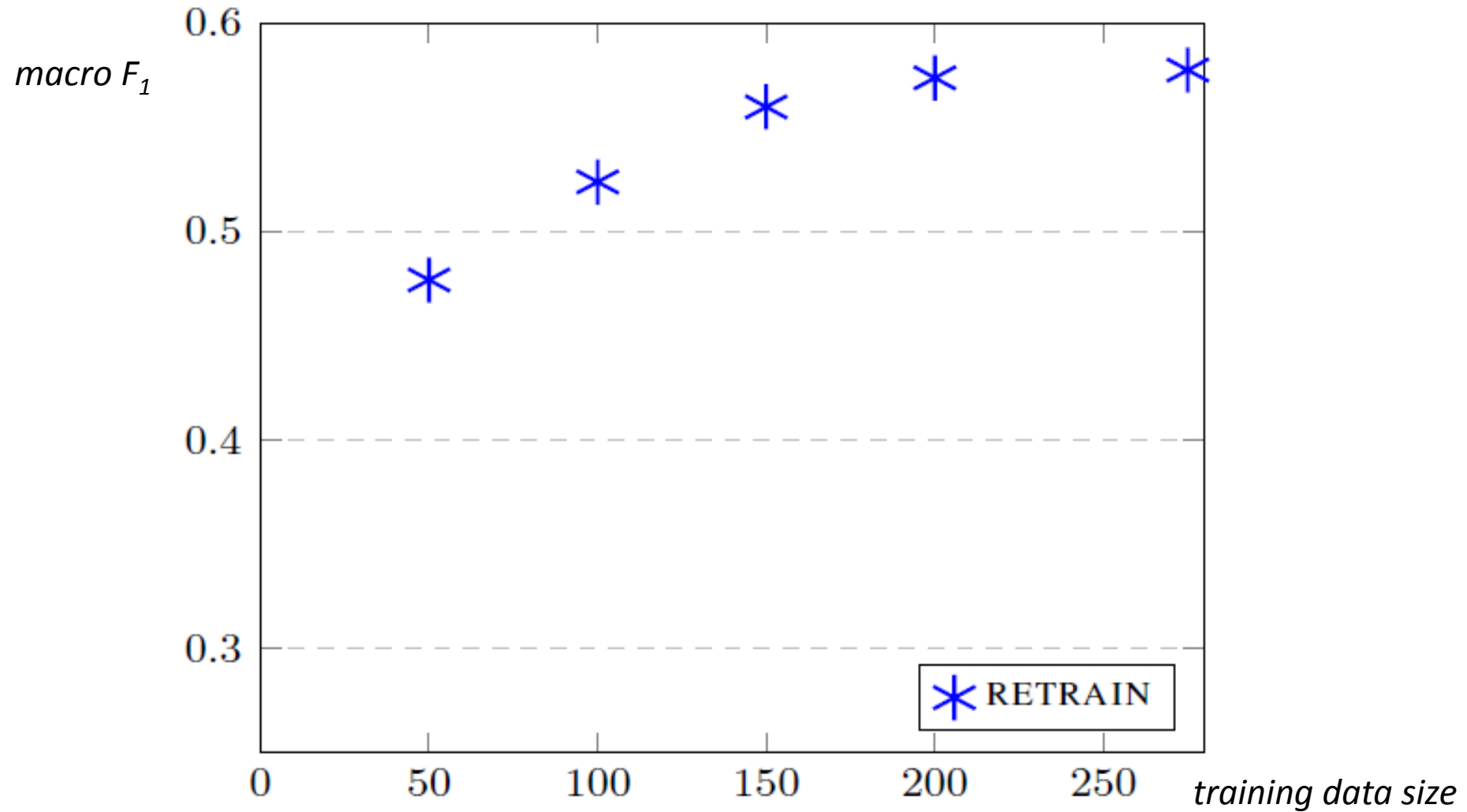*time*

# Interactive Model Training
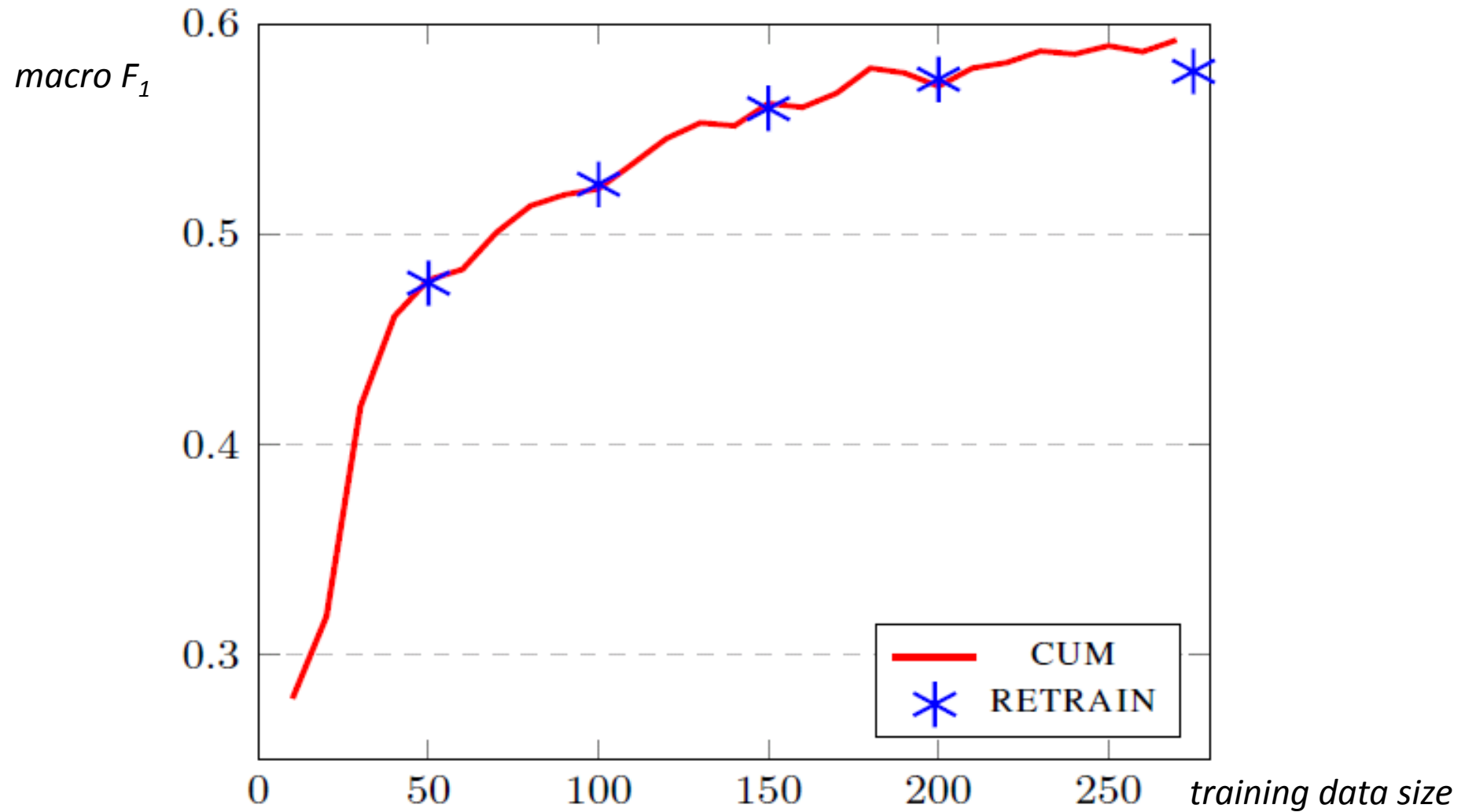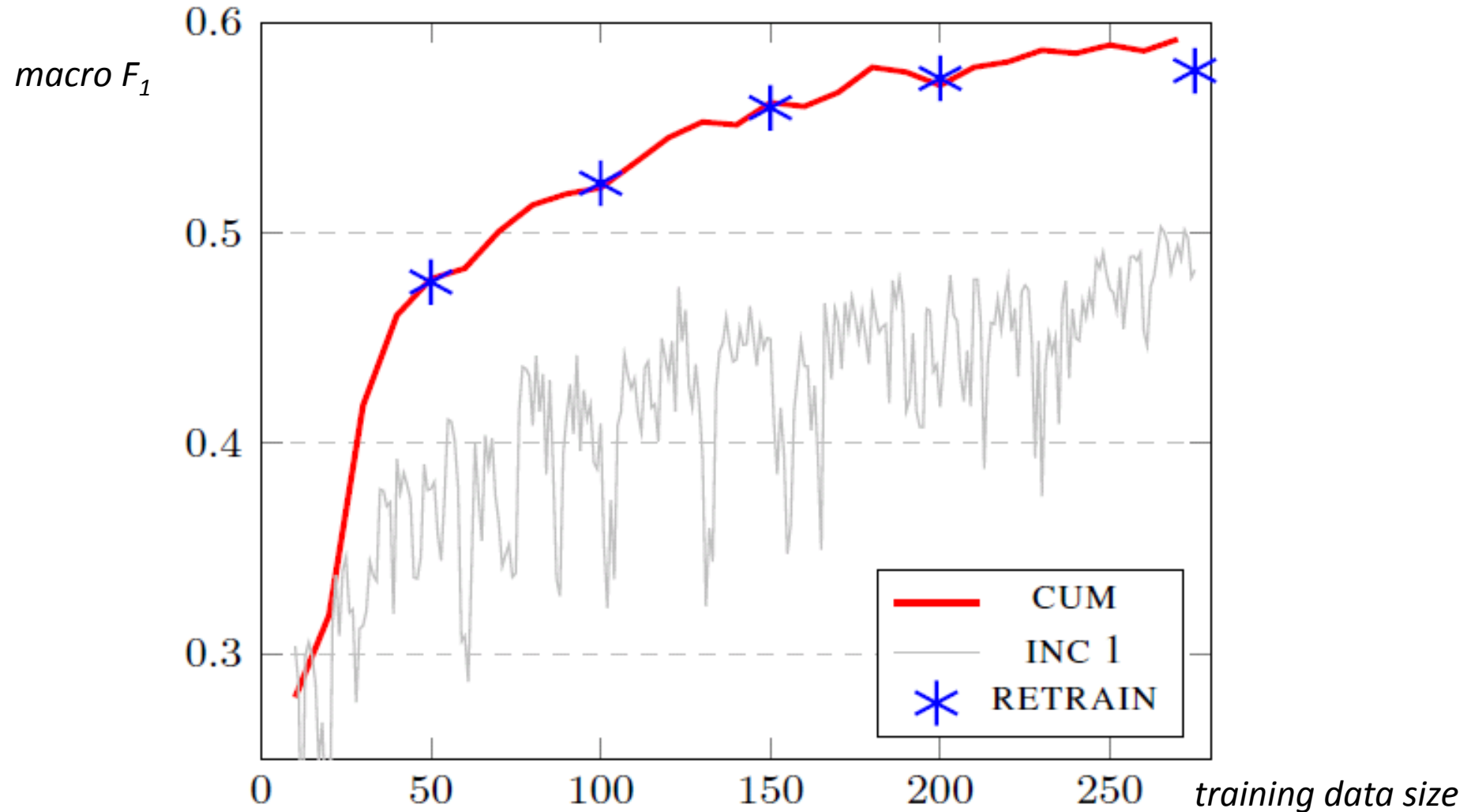
# Interactive Model Training

# Interactive Model Training



model

training
data

new
annotations

revised
model

**INCremental**

augmented
training data

i

i + 1

time

# Interactive Model Training



model

training
data

new
annotations

revised
model

**CUMulative**

augmented
training data

*i*

*i + 1*

time

# Model Performance

# Model Performance

# Model Performance

# Model Performance

# Interactively Trained Suggestions



Conclusion 3: Interactive Model Training yields good performance
and allows for time–quality trade-offs!

# Summary

Conclusion 1: Annotation suggestions are helpful for experts
and yield faster and (maybe) more reliable annotations!

Conclusion 2: Some evidence for annotation bias, but negligible,
as no systematic discrepancy compared to the control setup!

Conclusion 3: Interactive Model Training yields good performance
and allows for time–quality trade-offs!

**Reproducibility**
data: https://tudatalib.ulb.tu-darmstadt.de/handle/tudatalib/2001
model: https://github.com/UKPLab/aaai19-diagnostic-reasoning

# Thank you for your attention!

Conclusion 1: Annotation suggestions are helpful for experts
and yield faster and (maybe) more reliable annotations!

Conclusion 2: Some evidence for annotation bias, but negligible,
as no systematic discrepancy compared to the control setup!

Conclusion 3: Interactive Model Training yields good performance
and allows for time–quality trade-offs!

**Reproducibility**
data: https://tudatalib.ulb.tu-darmstadt.de/handle/tudatalib/2001
model: https://github.com/UKPLab/aaai19-diagnostic-reasoning

# Kontakt / Contact

**Dr. Christian M. Meyer**
Technische Universität Darmstadt
Ubiquitous Knowledge Processing Lab

Hochschulstr. 10, 64289 Darmstadt, Germany
+49 (0)6151 16–25293
+49 (0)6151 16–25295
meyer (at) ukp.informatik.tu-darmstadt.de