# Explaining (how to improve) Diagnostic Reasoning

## — *The FAMULUS Project* —

Please don't hesitate to ask questions during the talk!

**Claudia Schulz**

15/01/2020   XAI Seminar

FAMULUS

# Explainable AI

data
sample

↓

**AI
System**

↓

task
solution
(output)



People with no idea about AI, telling me my AI will destroy the world

Me wondering why my neural network is classifying a cat as a dog..

Dog

# Explainable AI

# Learning Diagnostic Reasoning
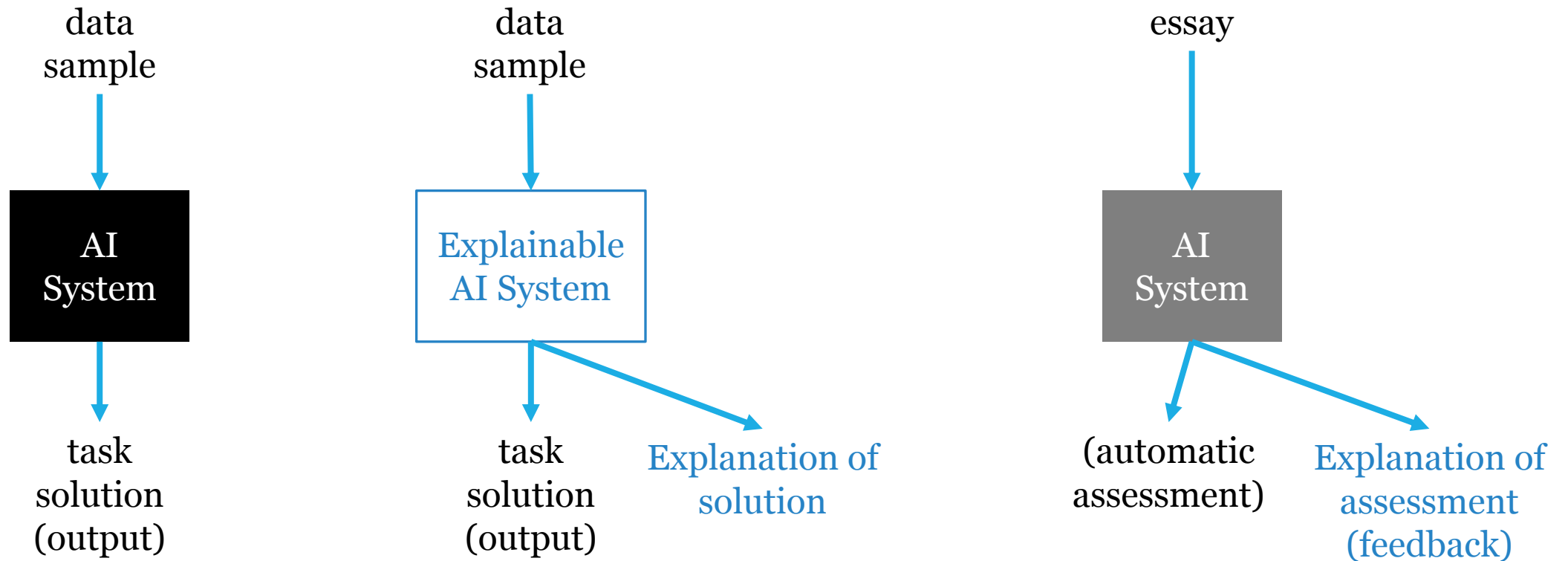
# Diagnostic Reasoning

# Case Simulation

# Individual Feedback

**Diagnostic Reasoning**

# Capture: Self-Explanation

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.

# Analyse: Reasoning Process



Diagnostic reasoning steps (epistemic activities) Fischer et al. 2014

Hypothesis Generation — *possible solutions*

Evidence Generation — *e.g. observations, deduction*

Evidence Evaluation — *evidence supports solution?*

Drawing Conclusions — *aggregate evidence to derive final solutions*

# Self-Explanation with Feedback

**Feedback**

Well done for thinking about different possible solutions, the generation of **hypotheses** is an important part of diagnosis.

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.

# Self-Explanation with Feedback

**Feedback**

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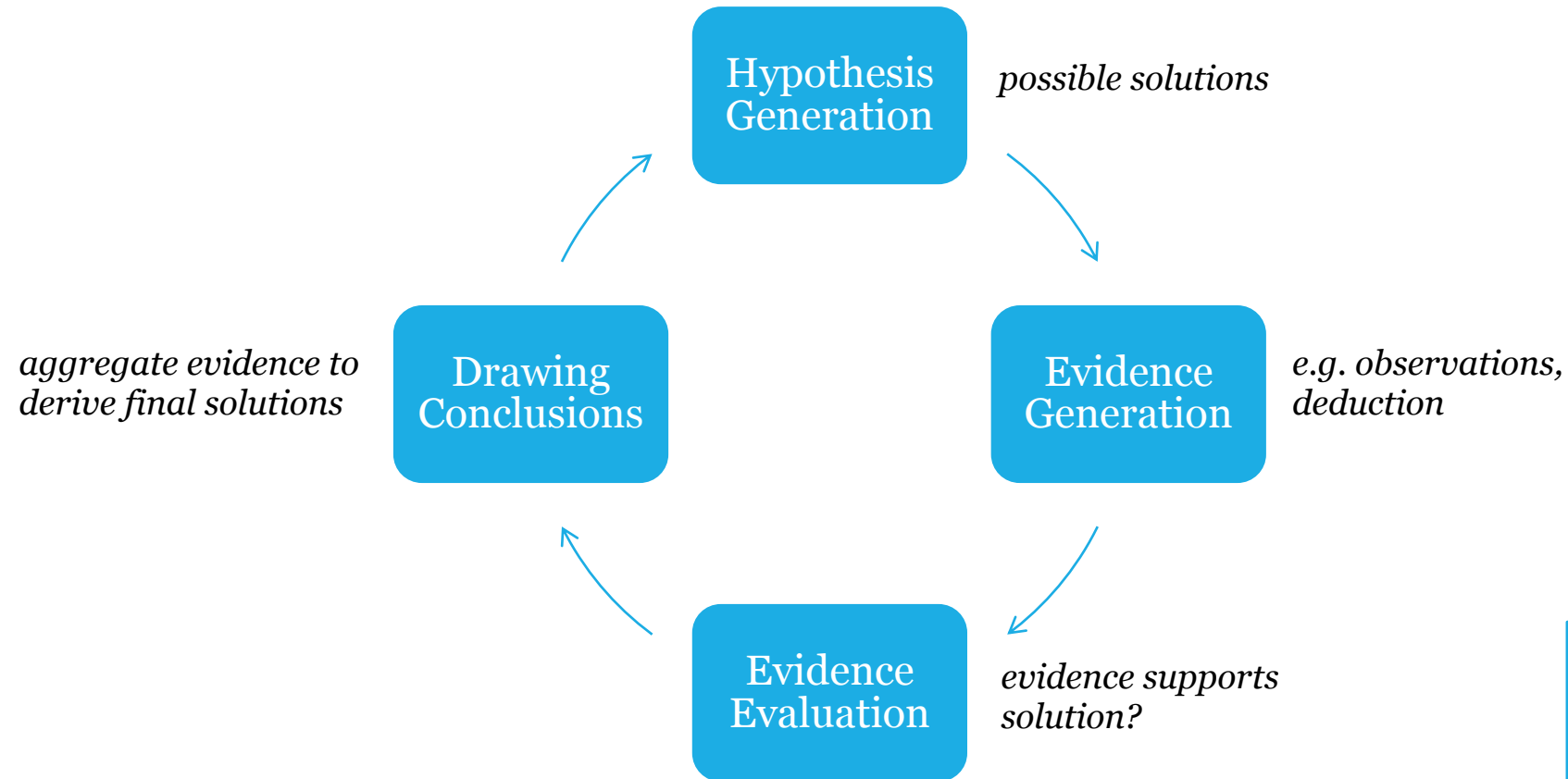
The patient reports to be lethargic and feverish.
From the anamnesis I learned that he had purulent tonsilitis and is still suffering from symptoms.
I first performed some laboratory tests and notice the decreased number of lymphocytes, which can be indicative of a bone marrow disease or an HIV infection.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

# Self-Explanation with Feedback

**Feedback**

After collecting 👎 and considering all evidence, you should decide on the most likely **diagnosis**. This is an important duty of a doctor.

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.

# Self-Explanation with Reasoning Steps

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    Evidence Generation
Evidence Evaluation      Drawing Conclusions

# Detecting Diagnostic Reasoning Steps

1) **Corpus Creation**
2) **Automatic Detection**

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    Evidence Generation
Evidence Evaluation       Drawing Conclusions

# Corpus Creation

Schulz, Meyer, Gurevych. "Challenges in the Automatic Analysis of Students' Diagnostic Reasoning."
*Proceedings of the 33rd AAAI Conference on Artificial Intelligence*. 2019.

# Corpus Creation

- Two domains:
  - Medicine Domain (MeD):        1131 self-explanations → 650 used
  - Teaching Domain (TeD):        976 self-explanations → 550 used

- (Domain) Expert annotators

- Cross-domain annotation scheme
  - Segmentation + classification
  - Easily adaptable to new domains

- German

# INCEpTION

https://inception-project.github.io/

# Corpus Creation

# Inter-Annotator Agreement

| Domain | $\alpha_U$ | $\alpha_U$-HG | $\alpha_U$-EG | $\alpha_U$-EE | $\alpha_U$-DC | $\alpha_U$-segment | $\uparrow\alpha_U$-pair | $\downarrow\alpha_U$-pair |
|---|---|---|---|---|---|---|---|---|
| medicine | 0.67 | 0.60 | 0.65 | 0.75 | 0.56 | 0.86 | 0.71 | 0.62 |
| teaching | 0.65 | 0.43 | 0.56 | 0.75 | 0.49 | 0.82 | 0.67 | 0.63 |

Table 1: Inter annotator agreement (IAA) in terms of Krippendorff's $\alpha_U$.

# Inter-Annotator Agreement

| Domain | $\alpha_U$ | $\alpha_U$-HG | $\alpha_U$-EG | $\alpha_U$-EE | $\alpha_U$-DC | $\alpha_U$-segment | $\uparrow\alpha_U$-pair | $\downarrow\alpha_U$-pair |
|--------|------------|---------------|---------------|---------------|---------------|--------------------|-------------------------|---------------------------|
| medicine | 0.67 | 0.60 | 0.65 | 0.75 | 0.56 | 0.86 | 0.71 | 0.62 |
| teaching | 0.65 | 0.43 | 0.56 | 0.75 | 0.49 | 0.82 | 0.67 | 0.63 |

Table 1: Inter annotator agreement (IAA) in terms of Krippendorff's $\alpha_U$.

# Inter-Annotator Agreement

| Domain | $\alpha_U$ | $\alpha_U$-HG | $\alpha_U$-EG | $\alpha_U$-EE | $\alpha_U$-DC | $\alpha_U$-segment | $\uparrow\alpha_U$-pair | $\downarrow\alpha_U$-pair |
|---|---|---|---|---|---|---|---|---|
| medicine | 0.67 | 0.60 | 0.65 | 0.75 | 0.56 | 0.86 | 0.71 | 0.62 |
| teaching | 0.65 | 0.43 | 0.56 | 0.75 | 0.49 | 0.82 | 0.67 | 0.63 |

Table 1: Inter annotator agreement (IAA) in terms of Krippendorff's $\alpha_U$.

# Inter-Annotator Agreement

| Domain | $\alpha_U$ | $\alpha_U$-HG | $\alpha_U$-EG | $\alpha_U$-EE | $\alpha_U$-DC | $\alpha_U$-segment | $\uparrow\alpha_U$-pair | $\downarrow\alpha_U$-pair |
|---|---|---|---|---|---|---|---|---|
| medicine | 0.67 | 0.60 | 0.65 | 0.75 | 0.56 | 0.86 | 0.71 | 0.62 |
| teaching | 0.65 | 0.43 | 0.56 | 0.75 | 0.49 | 0.82 | 0.67 | 0.63 |

Table 1: Inter annotator agreement (IAA) in terms of Krippendorff's $\alpha_U$.

# Inter-Annotator Agreement

| Domain | $\alpha_U$ | $\alpha_U$-HG | $\alpha_U$-EG | $\alpha_U$-EE | $\alpha_U$-DC | $\alpha_U$-segment | $\uparrow\alpha_U$-pair | $\downarrow\alpha_U$-pair |
|---|---|---|---|---|---|---|---|---|
| medicine | 0.67 | 0.60 | 0.65 | 0.75 | 0.56 | 0.86 | 0.71 | 0.62 |
| teaching | 0.65 | 0.43 | 0.56 | 0.75 | 0.49 | 0.82 | 0.67 | 0.63 |

Table 1: Inter annotator agreement (IAA) in terms of Krippendorff's $\alpha_U$.

| Domain | $\alpha_U$-HG&DC | $\alpha_U$-EE&DC | $\alpha_U$-HG&EE | $\alpha_U$-EG&EE | $\alpha_U$-EG&HG | $\alpha_U$-EG&DC |
|---|---|---|---|---|---|---|
| medicine | **0.71** | **0.85** | **0.78** | **0.78** | 0.61 | 0.56 |
| teaching | **0.62** | **0.81** | **0.77** | 0.72 | 0.47 | 0.48 |

Table 2: IAA ($\alpha_U$) when merging epistemic activities. Bold indicates a value higher than both single activities.

# Corpus Statistics

- majority vote (4/5, 3/4) + annotator meeting

- MeD av.  length:    63.8 tokens

- TeD av. Length:   100.2 tokens

|  |  | EG | EE | HG | DC |
|---|---|---|---|---|---|
| **MeD** | # | 219 | 2124 | 623 | 493 |
|  | av. # | 0.35 | 3.27 | 0.96 | 0.76 |
|  | av len. | 10.1 | 11.6 | 9.0 | 16.0 |
| **TeD** | # | 354 | 2671 | 311 | 444 |
|  | av. # | 0.64 | 4.86 | 0.57 | 0.81 |
|  | av. len. | 12.4 | 12.1 | 13.5 | 15.4 |

Table 3: Corpus statistics in terms of absolute number (#), average number per text (av. #), and average number of tokens (av. len), where EE/EG (and similar) denotes an overlap of an EG and EE segment.

# Corpus Statistics

- majority vote (4/5, 3/4) + annotator meeting

- MeD av.  length:    63.8 tokens

- TeD av. Length:   100.2 tokens

|     |         | EG   | EE   | HG   | DC   |
|-----|---------|------|------|------|------|
| MeD | #       | 219  | 2124 | 623  | 493  |
|     | av. #   | 0.35 | 3.27 | 0.96 | 0.76 |
|     | av len. | 10.1 | 11.6 | 9.0  | 16.0 |
| TeD | #       | 354  | 2671 | 311  | 444  |
|     | av. #   | 0.64 | 4.86 | 0.57 | 0.81 |
|     | av. len.| 12.4 | 12.1 | 13.5 | 15.4 |

Table 3: Corpus statistics in terms of absolute number (#), average number per text (av. #), and average number of tokens (av. len), where EE/EG (and similar) denotes an overlap of an EG and EE segment.

# Corpus Statistics

- majority vote (4/5, 3/4) + annotator meeting

- MeD av. length:    63.8 tokens

- TeD av. Length:   100.2 tokens

<div style="border:1px solid black;">

<mark>The x-ray and the subsequent</mark> MRI confirmed

a vertebral body fracture

</div>

|     |         | EG   | EE   | HG   | DC   | EG/EE | HG/DC | DC/EE | EG/HG | HG/EE | EG/DC |
|-----|---------|------|------|------|------|-------|-------|-------|-------|-------|-------|
| MeD | #       | 219  | 2124 | 623  | 493  | 5     | 4     | 342   | 0     | 12    | 4     |
|     | av. #   | 0.35 | 3.27 | 0.96 | 0.76 | –     | –     | –     | –     | –     | –     |
|     | av len. | 10.1 | 11.6 | 9.0  | 16.0 | 3.8   | 8.5   | 9.8   | –     | 5.7   | 6.8   |
| TeD | #       | 354  | 2671 | 311  | 444  | 8     | 2     | 143   | 3     | 8     | 3     |
|     | av. #   | 0.64 | 4.86 | 0.57 | 0.81 | –     | –     | –     | –     | –     | –     |
|     | av. len.| 12.4 | 12.1 | 13.5 | 15.4 | 7.9   | 22.0  | 10.9  | 6.0   | 11.1  | 11.7  |

Table 3: Corpus statistics in terms of absolute number (#), average number per text (av. #), and average number of tokens (av. len), where EE/EG (and similar) denotes an overlap of an EG and EE segment.

# Detecting Diagnostic Reasoning Steps

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.

1) **Corpus Creation**
2) **Automatic Detection**

Hypothesis Generation    Evidence Generation
Evidence Evaluation      Drawing Conclusions

# Automatic Detection of Diagnostic Reasoning Steps

Schulz, Meyer, Gurevych. "Challenges in the Automatic Analysis of Students' Diagnostic Reasoning." *Proceedings of the 33rd AAAI Conference on Artificial Intelligence*. 2019.

# Automatic Detection

**3 Challenges:**

1. segments of arbitrary length (C1),

2. distinguishing different epistemic activity types (C2)

3. overlapping epistemic activity segments (C3)

→ **multi-label problem:** $C' \subset C$

**Approach**: 3 problem transformations

**Multi-class sequence labelling**
$C = (\{B, I\} \times A) \cup \{O\}$
$A = \{HG, EG, EE, DC\}$

# Problem Transformations

*B/I/O*      *B/I/O*      *B/I/O*      *B/I/O*

| CRF | CRF | CRF | CRF |

| BiLSTM | BiLSTM | BiLSTM | BiLSTM |

**SEPARATE**:
multiple (single-label)
multi-class problems

*B/I/O – B/I/O – B/I/O – B/I/O*

| CRF |

| BiLSTM |

**CONCAT**:
unique (single-label)
multi-class problem

*B/I/O*      *B/I/O*      *B/I/O*      *B/I/O*

| CRF | CRF | CRF | CRF |

| BiLSTM |

**MULTI-OUTPUT**:
Multidimensional
classification problem

# Baseline Transformations

*B/I/O – EG/EE/HG/DC*

1. PREF-BASELINE: unique (single-label) multi-class problem
   - Without overlaps
   - Using preference order: DC > HG > EG > EE

2. MAJ-BASELNE: I-EE for all tokens

CRF

BiLSTM

# Evaluation Metrics

- Hamming Loss

$$HL = \frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} \frac{1}{|C|} \sum_{c \in C} \textbf{xor}(y_{x,c}, \hat{y}_{x,c})$$

$$y_{x,c} = \begin{cases} 1 & \textit{if token x has label c} \\ 0 & \textit{otherwise} \end{cases}$$

- C1 (Segmentation)

$$M_S(a) = macro\text{-}F1(C_a, \mathcal{X})$$

$$for\ a \in A = \{HG, EG, EE, DC\}$$

$$C_{HG} = \{B - HG, I - HG, O - HG\}$$

- C2 (Type Distinction)

$$M_A = macro\text{-}F1(\mathscr{P}(A), \mathcal{X})$$

- C3 (Overlaps)

$$M_O(a) = macro\text{-}F1(C_a, \mathcal{X}_{\text{overlap}})$$

# Automatic Detection: Results

$$HL = \frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} \frac{1}{|C|} \sum_{c \in C} \mathbf{xor}(y_{x,c}, \hat{y}_{x,c})$$

| | Architecture | $HL$ all | $M_S$ EG | EE | HG | DC | $M_A$ all | $M_O$ EG | EE | HG | DC |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **MeD** | MULTI-OUTPUT | 0.07 | 71.60 | 80.20[+] | 69.28 | 65.32 | 22.21[+] | 63.09 | 66.39[+] | 45.50 | 44.76 |
| | SEPARATE | 0.07 | 70.87 | 80.24[+] | 68.53 | 65.80 | 21.25[+] | 63.15 | 65.31[+] | 50.26 | 49.26 |
| | CONCAT | 0.06[+++] | 71.05 | 79.96[+] | 69.36 | 65.18 | 23.01[++] | 67.86 | 66.43[+] | 44.51 | 45.40 |
| | PREF-BASELINE | 0.07 | 70.02 | 75.46 | 69.32 | 65.74 | 19.77 | 52.91 | 38.87 | 46.34 | 49.03 |
| | MAJ-BA    | | | | | | | | | | 1.39 |
| | human u   | | | | | | | | | 8 | 76.50 |
| **TeD** | MULTI-OUTPUT | 0.07 | 78.55 | 78.87 | 57.18 | 61.77 | 19.96 | 58.42 | 71.98 | 32.61[+] | 47.10 |
| | SEPARATE | 0.07 | 76.38 | 79.47[+] | 57.05 | 57.52 | 18.34 | 54.68 | 78.89[+++] | 32.09 | 36.11 |
| | CONCAT | 0.06[++] | 78.71[+] | 79.07[+] | 57.12 | 62.53[+] | 21.68[+++] | 56.75 | 68.75[+] | 32.51 | 51.97[+] |
| | PREF-BASELINE | 0.06 | 77.60 | 77.21 | 55.67 | 61.02 | 18.93 | 57.25 | 45.15 | 36.62 | 49.71 |
| | MAJ-BASELINE | 0.11 | 31.75 | 23.11 | 32.03 | 30.97 | 4.42 | 31.21 | 30.75 | 32.61 | 6.28 |
| | human upper bound | 0.03 | 93.29 | 90.71 | 81.77 | 82.11 | 30.58 | 78.68 | 88.99 | 79.96 | 95.04 |

Conclusion: No distinction possible between neural architectures!

# Automatic Detection: Results

$$M_S(a) = macro\text{-}F1(C_a, \mathcal{X})$$

| | HL | $M_S$ | | | | $M_A$ | $M_O$ | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Architecture | all | EG | EE | HG | DC | all | EG | EE | HG | DC |
| **MeD** | | | | | | | | | | |
| Multi-Output | 0.07 | 71.60 | 80.20[+] | 69.28 | 65.32 | 22.21[+] | 63.09 | 66.39[+] | 45.50 | 44.76 |
| Separate | 0.07 | 70.87 | 80.24[+] | 68.53 | 65.80 | 21.25[+] | 63.15 | 65.31[+] | 50.26 | 49.26 |
| Concat | 0.06[+++] | 71.05 | 79.96[+] | 69.36 | 65.18 | 23.01[++] | 67.86 | 66.43[+] | 44.51 | 45.40 |
| Pref-Baseline | 0.07 | 70.02 | 75.46 | 69.32 | 65.74 | 19.77 | 52.91 | 38.87 | 46.34 | 49.03 |
| Maj-Baseline | 0.11 | 32.70 | 23.49 | 30.48 | 29.96 | 4.25 | 33.13 | 31.00 | 32.61 | 1.39 |
| human upper bound | 0.04 | 85.61 | 90.25 | 86.37 | 85.58 | 35.06 | 100.00 | 76.15 | 91.38 | 76.50 |
| **TeD** | | | | | | | | | | |
| Multi-Output | 0.07 | 78.53 | 78.87[+] | 57.16 | 61.77 | 19.96[+] | 58.42 | 71.98[+] | 32.61[+] | 47.10 |
| Separate | 0.07 | 76.38 | 79.47[+] | 57.05 | 57.52 | 18.34 | 54.68 | 78.89[+++] | 32.09 | 36.11 |
| Concat | 0.06[++] | 78.71[+] | 79.07[+] | 57.12 | 62.53[+] | 21.68[+++] | 56.75 | 68.75[+] | 32.51 | 51.97[+] |
| Pref-Baseline | 0.06 | 77.60 | 77.21 | 55.67 | 61.02 | 18.93 | 57.25 | 45.15 | 36.62 | 49.71 |
| Maj-Baseline | 0.11 | 31.75 | 23.11 | 32.03 | 30.97 | 4.42 | 31.21 | 30.75 | 32.61 | 6.28 |
| human upper bound | 0.03 | 93.29 | 90.71 | 81.77 | 82.11 | 30.58 | 78.68 | 88.99 | 79.96 | 95.04 |

# Automatic Detection: Results

$$M_S(a) = macro\text{-}F1(C_a, \mathcal{X})$$

| | HL | $M_S$ | | | | $M_A$ | $M_O$ | | | |
| Architecture | all | EG | EE | HG | DC | all | EG | EE | HG | DC |
|---|---|---|---|---|---|---|---|---|---|---|
| **MeD** | | | | | | | | | | |
| MULTI-OUTPUT | 0.07 | 71.60 | 80.20[+] | 69.28 | 65.32 | 22.21[+] | 63.09 | 66.39[+] | 45.50 | 44.76 |
| SEPARATE | 0.07 | 70.87 | 80.24[+] | 68.53 | 65.80 | 21.25[+] | 63.15 | 65.31[+] | 50.26 | 49.26 |
| CONCAT | 0.06[+++] | 71.05 | 79.96[+] | 69.36 | 65.18 | 23.01[++] | 67.86 | 66.43[+] | 44.51 | 45.40 |
| PREF-BASELINE | 0.07 | 70.02 | 75.46 | 69.32 | 65.74 | 19.77 | 52.91 | 38.87 | 46.34 | 49.03 |
| MAJ-BASELINE | 0.11 | 32.70 | 23.49 | 30.48 | 29.96 | 4.25 | 33.13 | 31.00 | 32.61 | 1.39 |
| human upper bound | 0.04 | 85.61 | 90.25 | 86.37 | 85.58 | 35.06 | 100.00 | 76.15 | 91.38 | 76.50 |
| **TeD** | | | | | | | | | | |
| MULTI-OUTPUT | 0.07 | 78.53 | 78.87[+] | 57.16 | 61.77 | 19.96[+] | 58.42 | 71.98[+] | 32.61[+] | 47.10 |
| SEPARATE | 0.07 | 76.38 | 79.47[+] | 57.05 | 57.52 | 18.34 | 54.68 | 78.89[+++] | 32.09 | 36.11 |
| CONCAT | 0.06[++] | 78.71[+] | 79.07[+] | 57.12 | 62.53[+] | 21.68[+++] | 56.75 | 68.75[+] | 32.51 | 51.97[+] |
| PREF-BASELINE | 0.06 | 77.60 | 77.21 | 55.67 | 61.02 | 18.93 | 57.25 | 45.15 | 36.62 | 49.71 |
| MAJ-BASELINE | 0.11 | 31.75 | 23.11 | 32.03 | 30.97 | 4.42 | 31.21 | 30.75 | 32.61 | 6.28 |
| human upper bound | 0.03 | 93.29 | 90.71 | 81.77 | 82.11 | 30.58 | 78.68 | 88.99 | 79.96 | 95.04 |

# Automatic Detection: Results

$$M_S(a) = macro\text{-}F1(C_a, \mathcal{X})$$

| | | HL | $M_S$ | | | | $M_A$ | $M_O$ | | | |
| | Architecture | all | EG | EE | HG | DC | all | EG | EE | HG | DC |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **MeD** | MULTI-OUTPUT | 0.07 | 71.60 | 80.20$^+$ | 69.28 | 65.32 | 22.21$^+$ | 63.09 | 66.39$^+$ | 45.50 | 44.76 |
| | SEPARATE | 0.07 | 70.87 | 80.24$^+$ | 68.53 | 65.80 | 21.25$^+$ | 63.15 | 65.31$^+$ | 50.26 | 49.26 |
| | CONCAT | 0.06$^{+++}$ | 71.05 | 79.96$^+$ | 69.36 | 65.18 | 23.01$^{++}$ | 67.86 | 66.43$^+$ | 44.51 | 45.40 |
| | PREF-BASELINE | 0.07 | 70.02 | 75.46 | 69.32 | 65.74 | 19.77 | 52.91 | 38.87 | 46.34 | 49.03 |
| | MAJ- | | | | | | | | | | 1.39 |
| | human | | | | | | | | | | 76.50 |
| **TeD** | MULTI | | | | | | | | | | 47.10 |
| | SEPARATE | 0.07 | 76.38 | 79.47$^+$ | 57.05 | 57.52 | 18.34 | 54.68 | 78.89$^{+++}$ | 32.09 | 36.11 |
| | CONCAT | 0.06$^{++}$ | 78.71$^+$ | 79.07$^+$ | 57.12 | 62.53$^+$ | 21.68$^{+++}$ | 56.75 | 68.75$^+$ | 32.51 | 51.97$^+$ |
| | PREF-BASELINE | 0.06 | 77.60 | 77.21 | 55.67 | 61.02 | 18.93 | 57.25 | 45.15 | 36.62 | 49.71 |
| | MAJ-BASELINE | 0.11 | 31.75 | 23.11 | 32.03 | 30.97 | 4.42 | 31.21 | 30.75 | 32.61 | 6.28 |
| | human upper bound | 0.03 | 93.29 | 90.71 | 81.77 | 82.11 | 30.58 | 78.68 | 88.99 | 79.96 | 95.04 |

Conclusion: Neural architectures perform segmentation reasonably well!

# Automatic Detection: Results

Upper Bound: 62.5

$$M_A = macro\text{-}F1(\mathscr{P}(A), \mathcal{X})$$

| | HL | $M_S$ | | | | $M_A$ | $M_O$ | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Architecture | all | EG | EE | HG | DC | all | EG | EE | HG | DC |
| **MeD** | | | | | | | | | | |
| Multi-Output | 0.07 | 71.60 | 80.20[+] | 69.28 | 65.32 | 22.21[+] | 63.09 | 66.39[+] | 45.50 | 44.76 |
| Separate | 0.07 | 70.87 | 80.24[+] | 68.53 | 65.80 | 21.25[+] | 63.15 | 65.31[+] | 50.26 | 49.26 |
| Concat | 0.06[+++] | 71.05 | 79.96[+] | 69.36 | 65.18 | 23.01[++] | 67.86 | 66.43[+] | 44.51 | 45.40 |
| Pref-Baseline | 0.07 | 70.02 | 75.46 | 69.32 | 65.74 | 19.77 | 52.91 | 38.87 | 46.34 | 49.03 |
| Maj-Baseline | | | | | | | | | | 1.39 |
| human upper bound | | | | | | | | | | 76.50 |
| **TeD** | | | | | | | | | | |
| Multi-Output | 0.07 | 78.35 | 78.87 | 57.16 | 61.77 | 19.96 | 58.42 | 71.98 | 32.01 | 47.10 |
| Separate | 0.07 | 76.38 | 79.47[+] | 57.05 | 57.52 | 18.34 | 54.68 | 78.89[+++] | 32.09 | 36.11 |
| Concat | 0.06[++] | 78.71[+] | 79.07[+] | 57.12 | 62.53[+] | 21.68[+++] | 56.75 | 68.75[+] | 32.51 | 51.97[+] |
| Pref-Baseline | 0.06 | 77.60 | 77.21 | 55.67 | 61.02 | 18.93 | 57.25 | 45.15 | 36.62 | 49.71 |
| Maj-Baseline | 0.11 | 31.75 | 23.11 | 32.03 | 30.97 | 4.42 | 31.21 | 30.75 | 32.61 | 6.28 |
| human upper bound | 0.03 | 93.29 | 90.71 | 81.77 | 82.11 | 30.58 | 78.68 | 88.99 | 79.96 | 95.04 |

**Conclusion: Distinction of different reasoning steps is highly challenging!**

# Automatic Detection: Results

*Upper Bound: 62.5*

$$M_A = macro\text{-}F1(\mathscr{P}(A), \mathcal{X})$$

|  | Architecture | HL | $M_S$ | | | | $M_A$ | $M_O$ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
|  |  | all | EG | EE | HG | DC | all | EG | EE | HG | DC |
| **MeD** | MULTI-OUTPUT | 0.07 | 71.60 | 80.20[+] | 69.28 | 65.32 | 22.21[+] | 63.09 | 66.39[+] | 45.50 | 44.76 |
|  | SEPARATE | 0.07 | 70.87 | 80.24[+] | 68.53 | 65.80 | 21.25[+] | 63.15 | 65.31[+] | 50.26 | 49.26 |
|  | CONCAT | 0.06[+++] | 71.05 | 79.96[+] | 69.36 | 65.18 | 23.01[++] | 67.86 | 66.43[+] | 44.51 | 45.40 |
|  | PREF-BASELINE | 0.07 | 70.02 | 75.46 | 69.32 | 65.74 | 19.77 | 52.91 | 38.87 | 46.34 | 49.03 |
|  | MAJ-BASELINE |  |  |  |  |  |  |  |  |  | 1.39 |
|  | human upper bound |  |  |  |  |  |  |  |  |  | 76.50 |
| **TeD** | MULTI-OUTPUT | 0.07 | 76.55 | 78.87 | 57.10 | 61.77 | 19.96 | 58.42 | 71.98 | 32.61[+] | 47.10 |
|  | SEPARATE | 0.07 | 76.38 | 79.47[+] | 57.05 | 57.52 | 18.34 | 54.68 | 78.89[+++] | 32.09 | 36.11 |
|  | CONCAT | 0.06[++] | 78.71[+] | 79.07[+] | 57.12 | 62.53[+] | 21.68[+++] | 56.75 | 68.75[+] | 32.51 | 51.97[+] |
|  | PREF-BASELINE | 0.06 | 77.60 | 77.21 | 55.67 | 61.02 | 18.93 | 57.25 | 45.15 | 36.62 | 49.71 |
|  | MAJ-BASELINE | 0.11 | 31.75 | 23.11 | 32.03 | 30.97 | 4.42 | 31.21 | 30.75 | 32.61 | 6.28 |
|  | human upper bound | 0.03 | 93.29 | 90.71 | 81.77 | 82.11 | 30.58 | 78.68 | 88.99 | 79.96 | 95.04 |

Conclusion: Overlapping segments are highly challenging!

# Automatic Detection: Results

Upper Bound: 62.5

$$M_A = macro\text{-}F1(\mathscr{P}(A), \mathcal{X})$$

| | | HL | $M_S$ | | | | $M_A$ | $M_O$ | | | |
| | Architecture | all | EG | EE | HG | DC | all | EG | EE | HG | DC |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **MeD** | MULTI-OUTPUT | 0.07 | 71.60 | 80.20+ | 69.28 | 65.32 | 22.21+ | 63.09 | 66.39+ | 45.50 | 44.76 |
| | SEPARATE | 0.07 | 70.87 | 80.24+ | 68.53 | 65.80 | 21.25+ | 63.15 | 65.31+ | 50.26 | 49.26 |
| | CONCAT | 0.06+++ | 71.05 | 79.96+ | 69.36 | 65.18 | 23.01++ | 67.86 | 66.43+ | 44.51 | 45.40 |
| | PREF-BASELINE | 0.07 | 70.02 | 75.46 | 69.32 | 65.74 | 19.77 | 52.91 | 38.87 | 46.34 | 49.03 |
| | MAJ-BASELINE | | | | | | | | | | 1.39 |
| | human upper | | | | | | | | | | 76.50 |
| **TeD** | MULTI-OUTPUT | 0.07 | 78.53 | 78.87+ | 57.16 | 61.77 | 19.96+ | 58.42 | 71.98+ | 32.61+ | 47.10 |
| | SEPARATE | 0.07 | 76.38 | 79.47+ | 57.05 | 57.52 | 18.34 | 54.68 | 78.89+++ | 32.09 | 36.11 |
| | CONCAT | 0.06++ | 78.71+ | 79.07+ | 57.12 | 62.53+ | 21.68+++ | 56.75 | 68.75+ | 32.51 | 51.97+ |
| | PREF-BASELINE | 0.06 | 77.60 | 77.21 | 55.67 | 61.02 | 18.93 | 57.25 | 45.15 | 36.62 | 49.71 |
| | MAJ-BASELINE | 0.11 | 31.75 | 23.11 | 32.03 | 30.97 | 4.42 | 31.21 | 30.75 | 32.61 | 6.28 |
| | human upper bound | 0.03 | 93.29 | 90.71 | 81.77 | 82.11 | 30.58 | 78.68 | 88.99 | 79.96 | 95.04 |

Conclusion: No architecture wins!

# Detecting Diagnostic Reasoning Steps

✓ 1) **Corpus Creation**
✓ 2) **Automatic Detection**

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          Evidence Generation
Evidence Evaluation            Drawing Conclusions
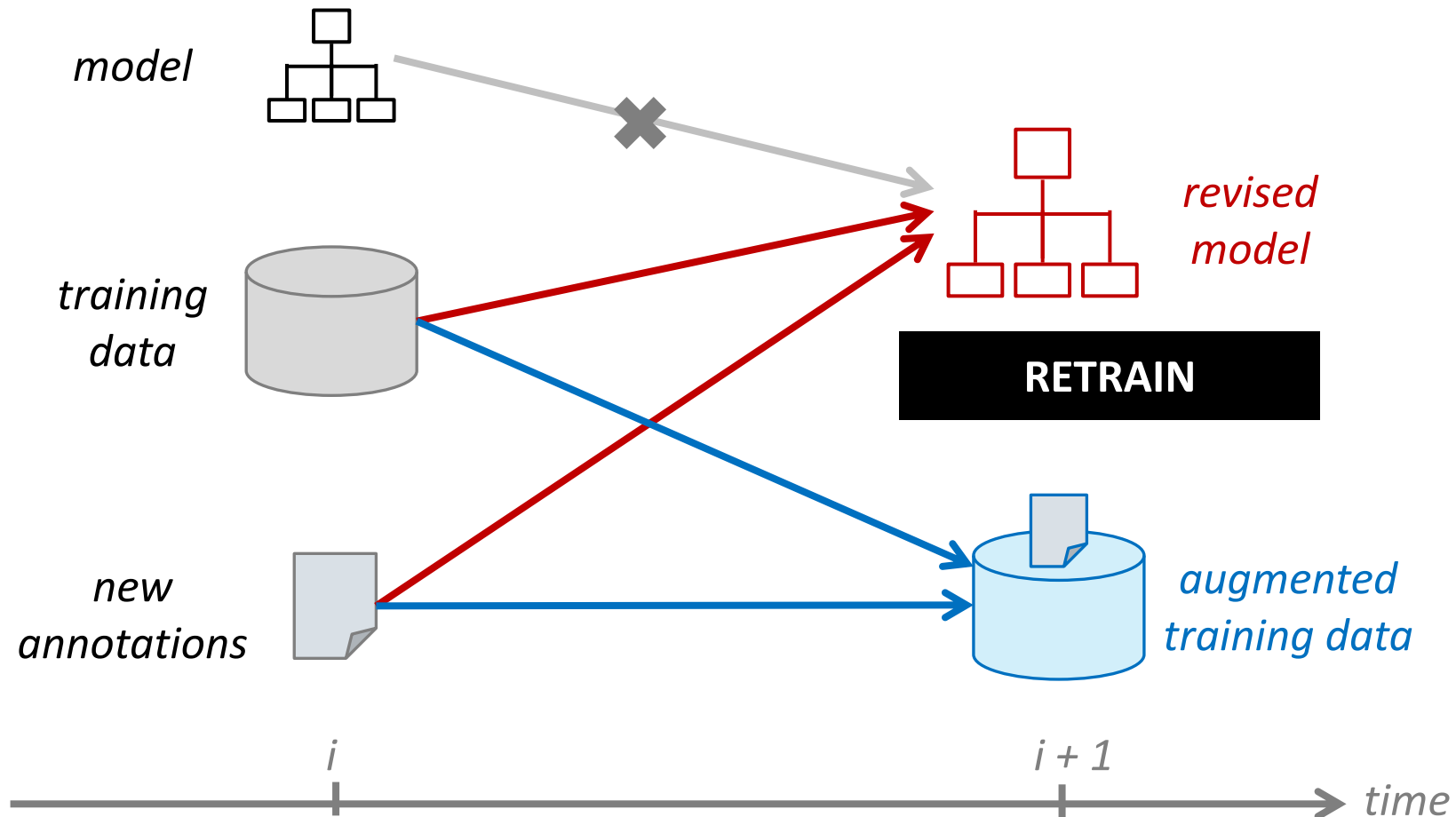
# Annotating the 2<sup>nd</sup> half of self-explanations

Schulz, et al. "Analysis of Automatic Annotation Suggestions for Hard Discourse-Level Tasks in Expert Domains."

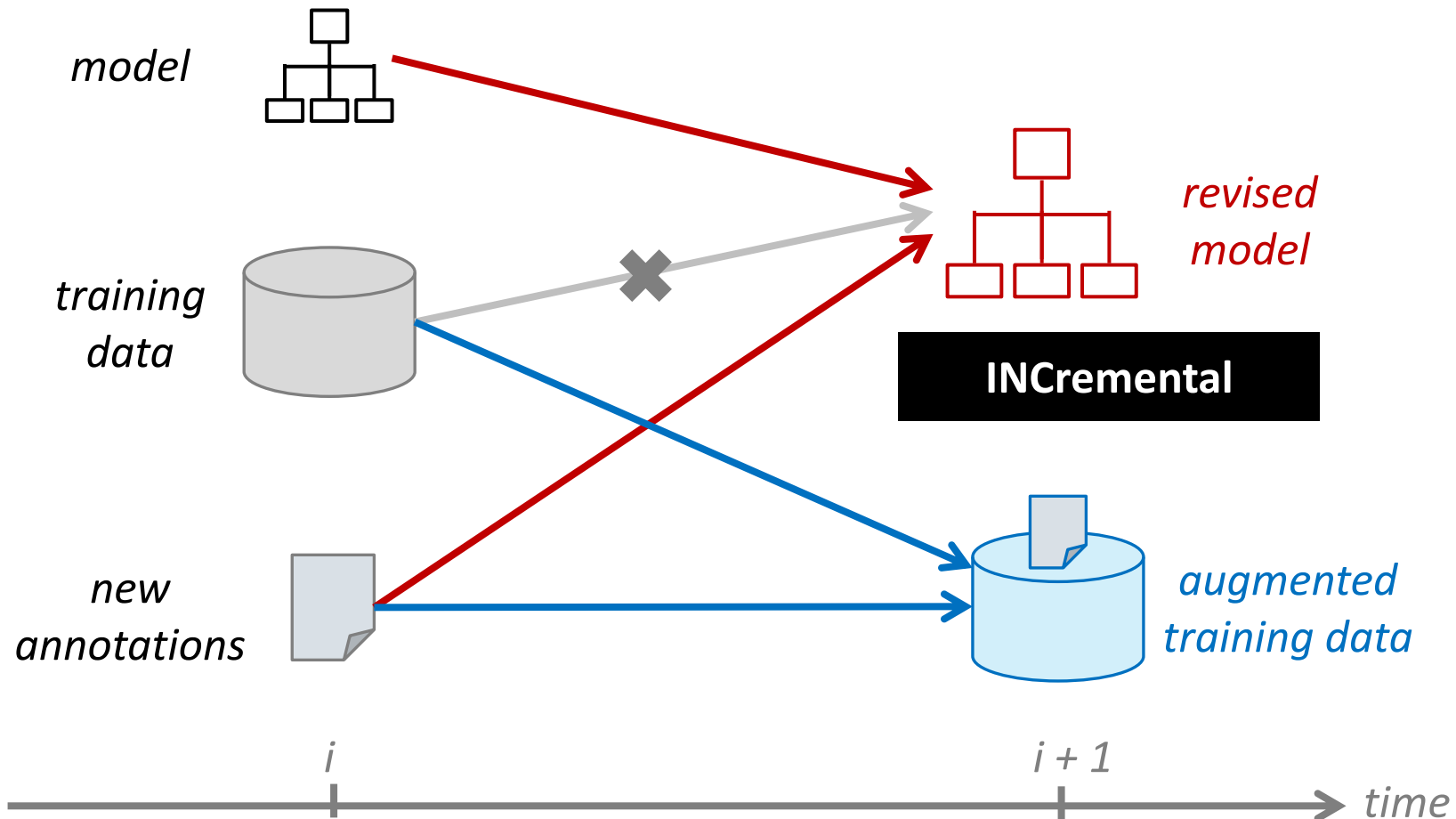*Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics.* 2019.
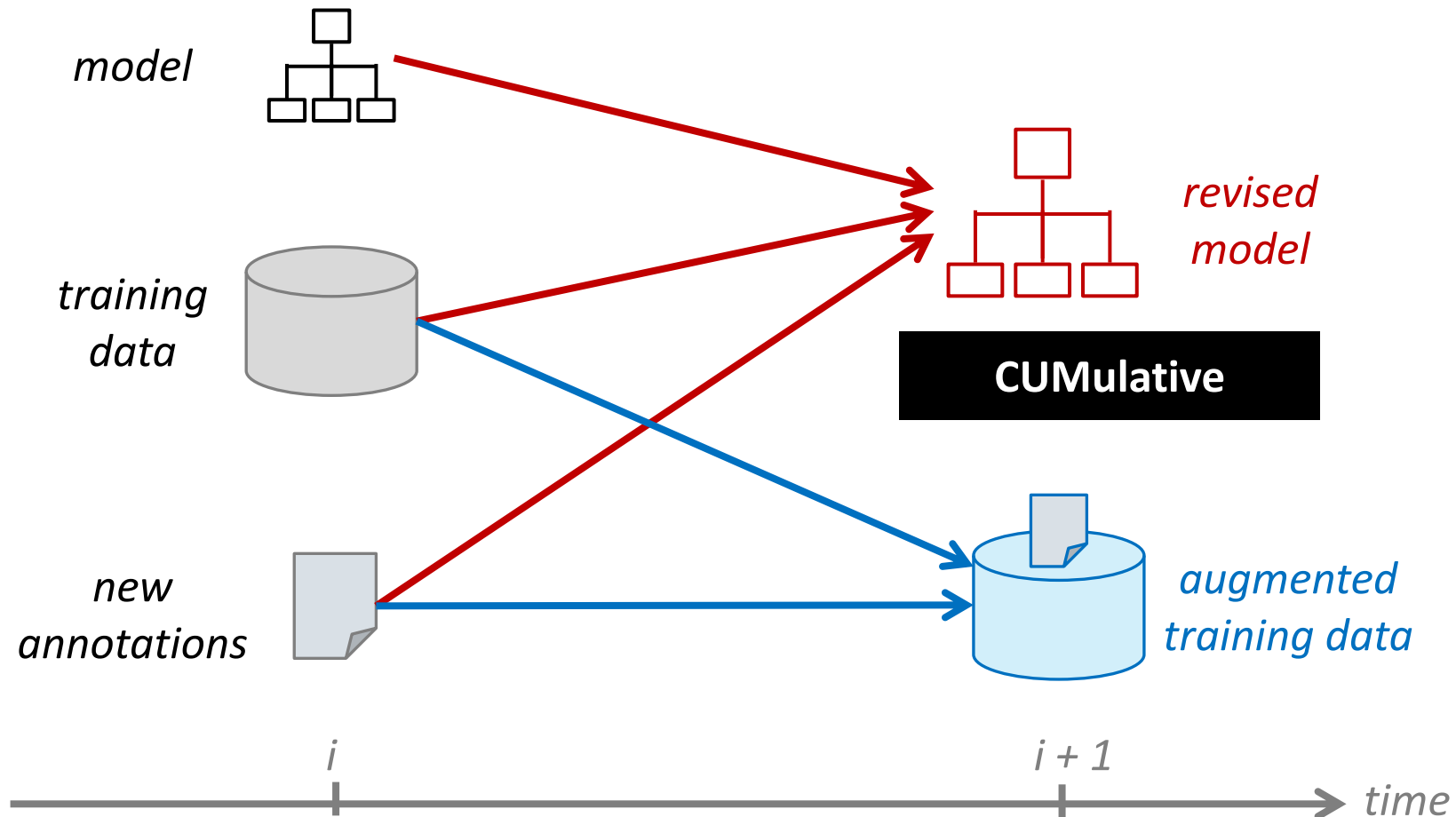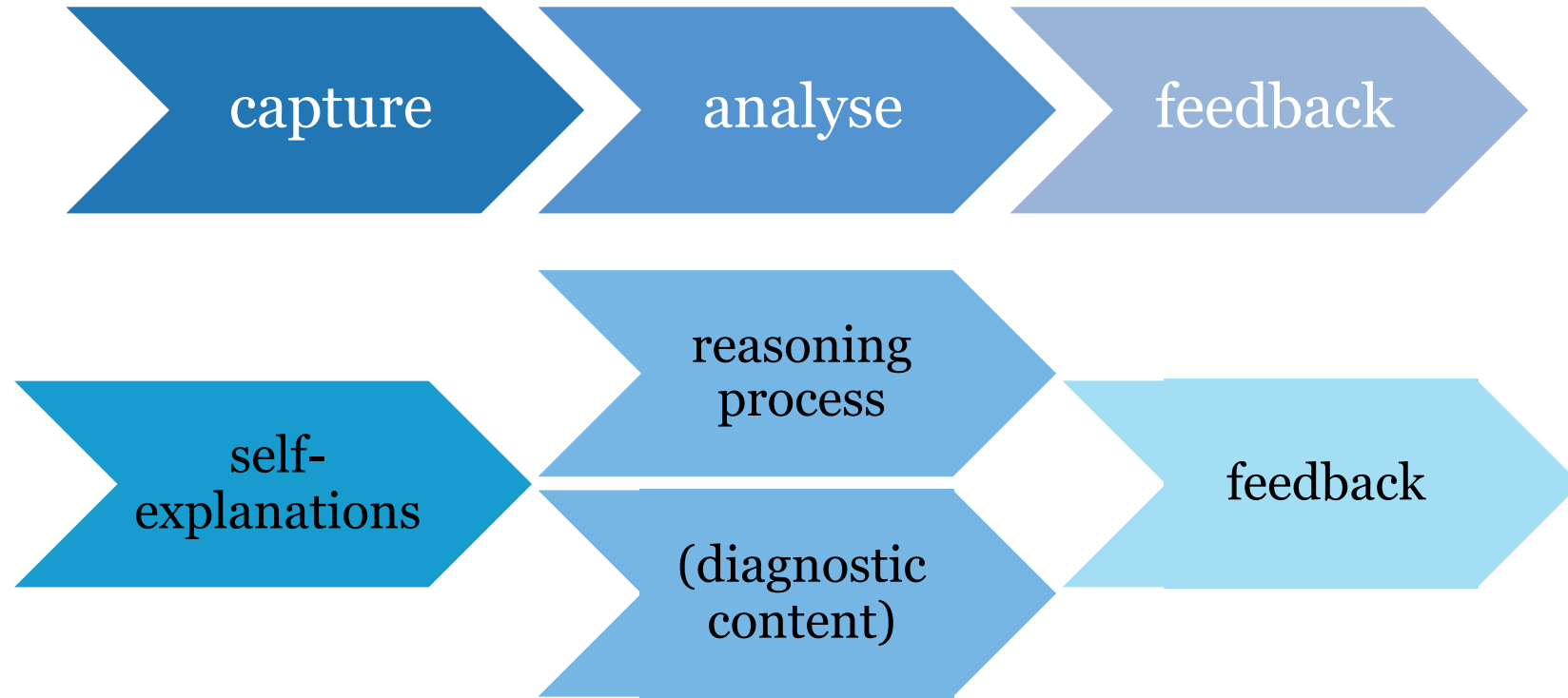
# Extending the Corpus



**INCEpTION**

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.
...d cultures are negative, so it is a virus, ...ction causing the symptoms.

**Suggestion:**
**Hypothesis generation (HG)**

✓ Accept     ✗ Reject

predictions
(PREF-BASELINE)

Hypothesis Generation     Evidence Generation
Evidence Evaluation        Drawing Conclusions

# Training Data and Suggestion Quality

**S1**    **S2**

**A1**

**A2**

**A3**

**A4**

**A5**

{EG, EE, HG, DC}

CRF

BiLSTM

**univ(ersal) model: $F_1 \approx .63$**
**pers(onalized) models: $F_1 \approx .55$**

# Annotation Suggestions in INCEpTION

# Annotation Suggestions - Setup

# Usefulness of Annotations

# Annotation Time



**[Minutes per text]**

# Reliability of Annotations



**[Krippendorff's $\alpha$]**

S1   S2   O1   O2   O3.1   O3.2   O4.1   O4.2

A1

**Conclusion: Annotation suggestions are helpful for experts and yield faster and more reliable annotations!**

A4
A5

**0.67**          **0.48**

# Intra-Annotator Consistency

# Human / Suggestion Model Agreement

**[Krippendorff's α]**



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

| | S1 | S2 | O1 | O2 | O3.1 | O3.2 | O4.1 | O4.2 |
|---|---|---|---|---|---|---|---|---|
| A1 | | | | | | | | |
| A4 | | | 0.56 | 0.48 | 0.55 | | 0.30 | |
| A5 | | | | | | | | |

# Iterative Model Training

Schulz, et al. "Analysis of Automatic Annotation Suggestions for Hard Discourse-Level Tasks in Expert Domains."

*Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 2019.

# Iterative Model Training



model

training
data

new
annotations

*revised
model*

*augmented
training data*

i

i + 1

time

# Iterative Model Training



model

training
data

new
annotations

revised
model

**RETRAIN**

augmented
training data

i

i + 1

time

# Iterative Model Training



model

training
data

new
annotations

revised
model

**INCremental**

augmented
training data

i

i + 1

time

# Iterative Model Training



model

training
data

new
annotations

*revised
model*

**CUMulative**

*augmented
training data*

i

i + 1

*time*

# Model Performance



Conclusion: INCremental Model Training yields good performance and allows for time-quality trade-offs

# Automatic Feedback on Diagnostic Reasoning

# Detecting Diagnostic Reasoning Steps

1) **Corpus Creation** ✓
2) **Automatic Detection** ✓

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 | Evidence Generation |
| Evidence Evaluation | Drawing Conclusions |

# Adaptive Feedback

**Diagnostic Reasoning**

# eLearning Platform

# eLearning Platform

# Automatic Feedback

*student's self-explanation*



☑ **Textaufgabe**

Die Körperliche Untersuchung war unauffällig. Allerdings waren im Labor die Entzündungswerte und Leberwerte auffällig. Der dicke Tropfen war negativ, daher war Malaria als Diagnose ausgeschlossen. Die Hepatitis Serologie war positiv und damit die Diagnose gesichert.

**Vielen Dank für Ihre Antwort!**

**Fallübersicht**

Die 36-jährige Frau Hoffmann stellt sich vor, mit einem seit einer Woche bestehenden grippalen Infekt. Als zusätzliche Symptome gibt sie Abgeschlagenheit, Appetitverlust, Übelkeit und Diarrhoe an. Sie war vor einem Monat ins Sansibar, vor der Reise wurde eine Gelbfieberimpfung durchgehführt.
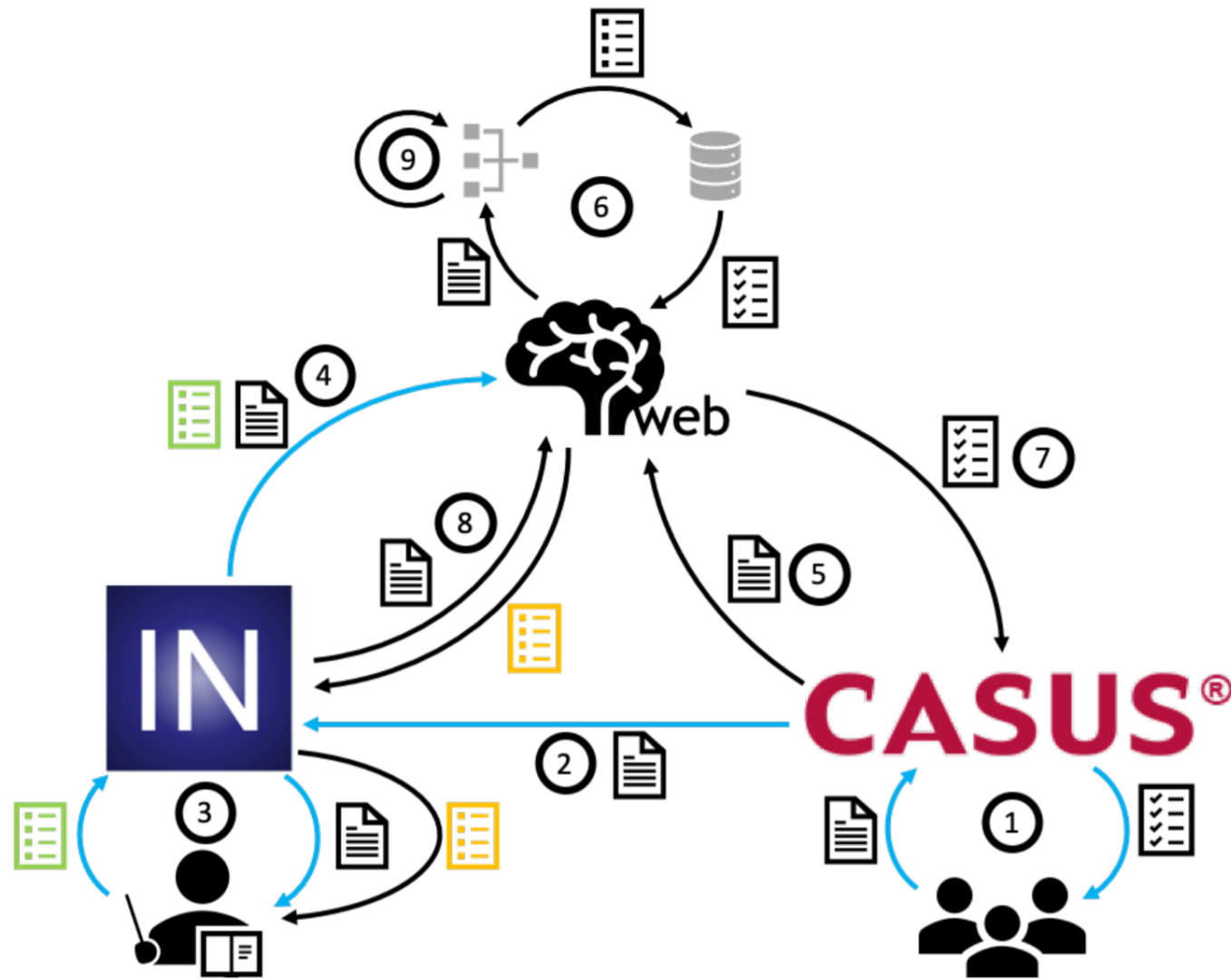
**Rückmeldung zu Differentialdiagnose**

💬 Bei einem einwöchigen grippalen Infekt mit Gliederschmerzen und Abgeschlagenheit wäre zunächst eine Influenza-Infektion denkbar gewesen. Für einen grippalen Infekt ist die Symptomatik allerdings zu langanhaltend, da dieser meist nach 3 Tagen abklingt.

💬 Bei einer Diarrhoe hättest du auch eine Darmerkrankung, wie die Gastroenteritis, vermuten sollen.

👁 Nicht schlecht, dass du eine Tropenkrankheit differentialdiagnostisch in Betracht gezogen hast. Möglich wären zB. Malaria, Dengue Fieber, Cholera etc.

*automatic adaptive feedback*

# Wrapping Up

# Explaining how to improve Diagnostic Reasoning

- Collect Self-Explanations

- Annotate Diagnostic Reasoning Steps
  - Train model for annotation suggestions → ease and speed-up

- Train Model for Detecting Reasoning Steps

- Use Model for Automatic Feedback

My publications
http://www.famulus-project.de/
https://inception-project.github.io/

For more questions, contact me:
**clauschulz1812@gmail.com**