# Computer-assisted stylistic revision with incomplete and noisy feedback
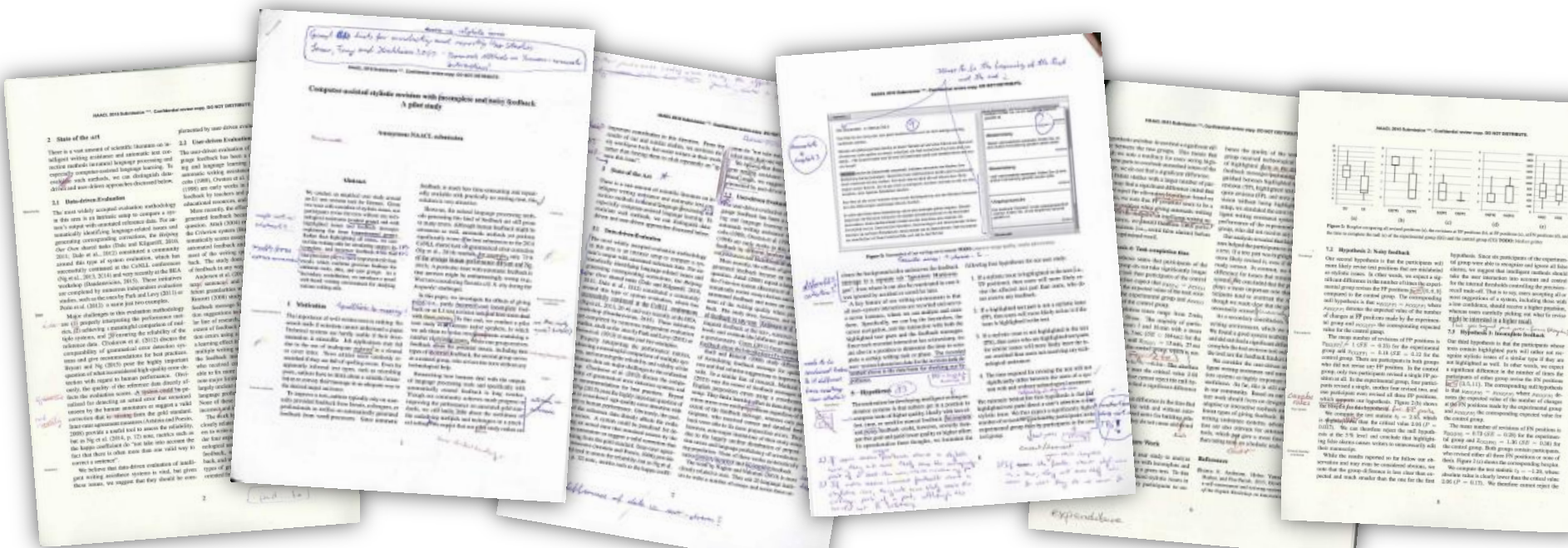## A pilot study

**Christian M. Meyer** and **Johann Frerik Koch**

The 11th Workshop on Innovative Use of NLP
for Building Educational Applications (BEA)
June 16, 2016. San Diego, CA, USA.

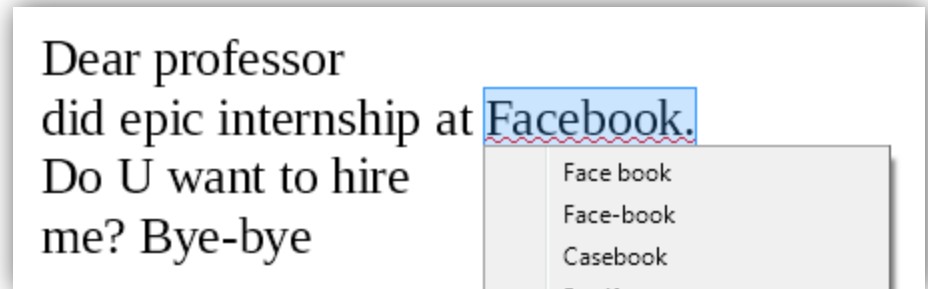TECHNISCHE
UNIVERSITÄT
DARMSTADT

UKP

# Vision

- Bryant & Ng (2015): best grammar correction software achieves only 73% of human performance

- **Our vision:** research new <u>useful</u> approaches to intelligent writing assistance with a focus on German native speakers

# Goal of this work

**There won't be perfect systems!** ☹

*How do users deal with incorrect and incomplete feedback?*

Dear professor
did epic internship at Facebook.
Do U want to hire
me? Bye-bye

- **Pilot user study**
- German L1 text revision task
- focus on stylistic issues
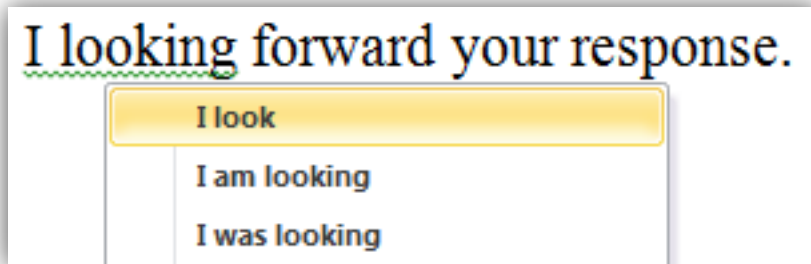
# Previous Work

**Data-driven evaluation**

- Shared tasks: HOO, CoNLL, BEA,…

- meaningful system comparison?

- interpretation of evaluation metrics?

- reliability of the reference data?

**User-driven evaluation**

- (Manual) feedback by teachers and peers

- Variation of feedback granularity, extent & formulation, time

- Nagata & Nakatani (2010): *"precision-oriented error detection is better than recall-oriented"*
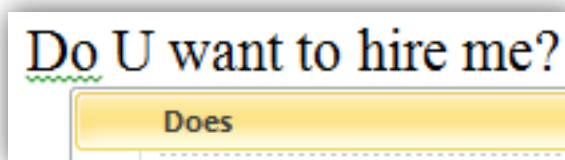
# Hypotheses

**H1**     If users receive **correct feedback**, they will more likely revise the corresponding section

# Hypotheses

**H1**  If users receive **correct feedback**, they will more likely revise the corresponding section

**H2**  If users receive **incorrect feedback**, they will more likely revise the corresponding section
– although it would not be necessary

# Hypotheses

**H1**   If users receive **correct feedback**, they will more likely revise the corresponding section

**H2**   If users receive **incorrect feedback**, they will more likely revise the corresponding section
– although it would not be necessary

**H3**   If users receive **incomplete feedback**, they will more likely miss issues not highlighted to them

# Hypotheses

**H1**     If users receive **correct feedback**, they will more likely revise the corresponding section

**H2**     If users receive **incorrect feedback**, they will more likely revise the corresponding section
– although it would not be necessary

**H3**     If users receive **incomplete feedback**, they will more likely miss issues not highlighted to them

**H4**     Providing automatic feedback does not affect the required **time to complete the task**
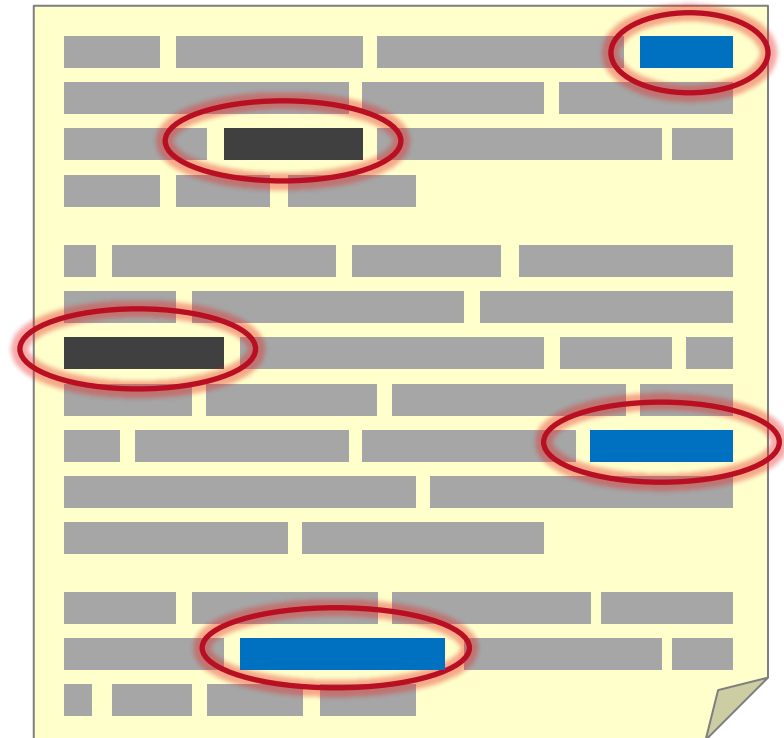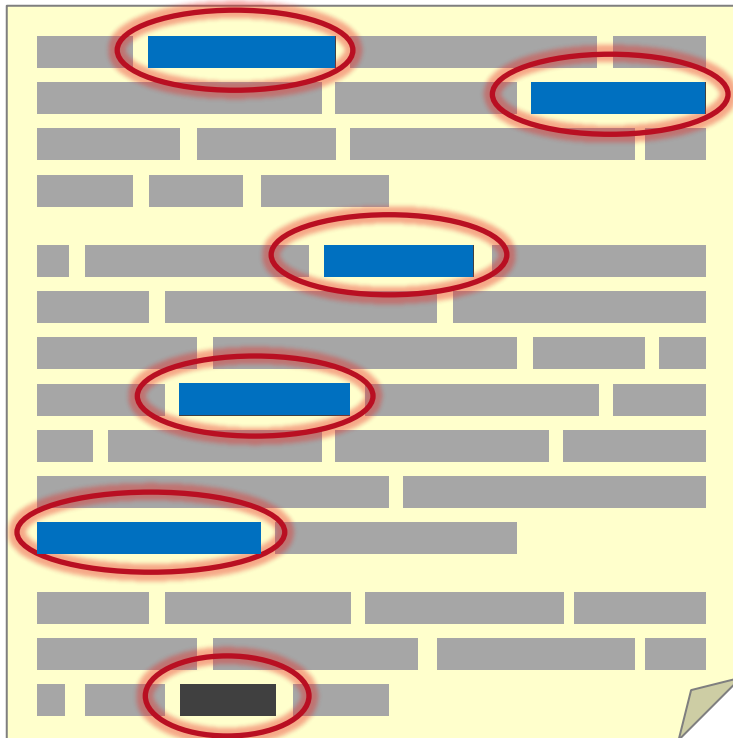
# Experimental Setup: Data



**$T_1$ News item**
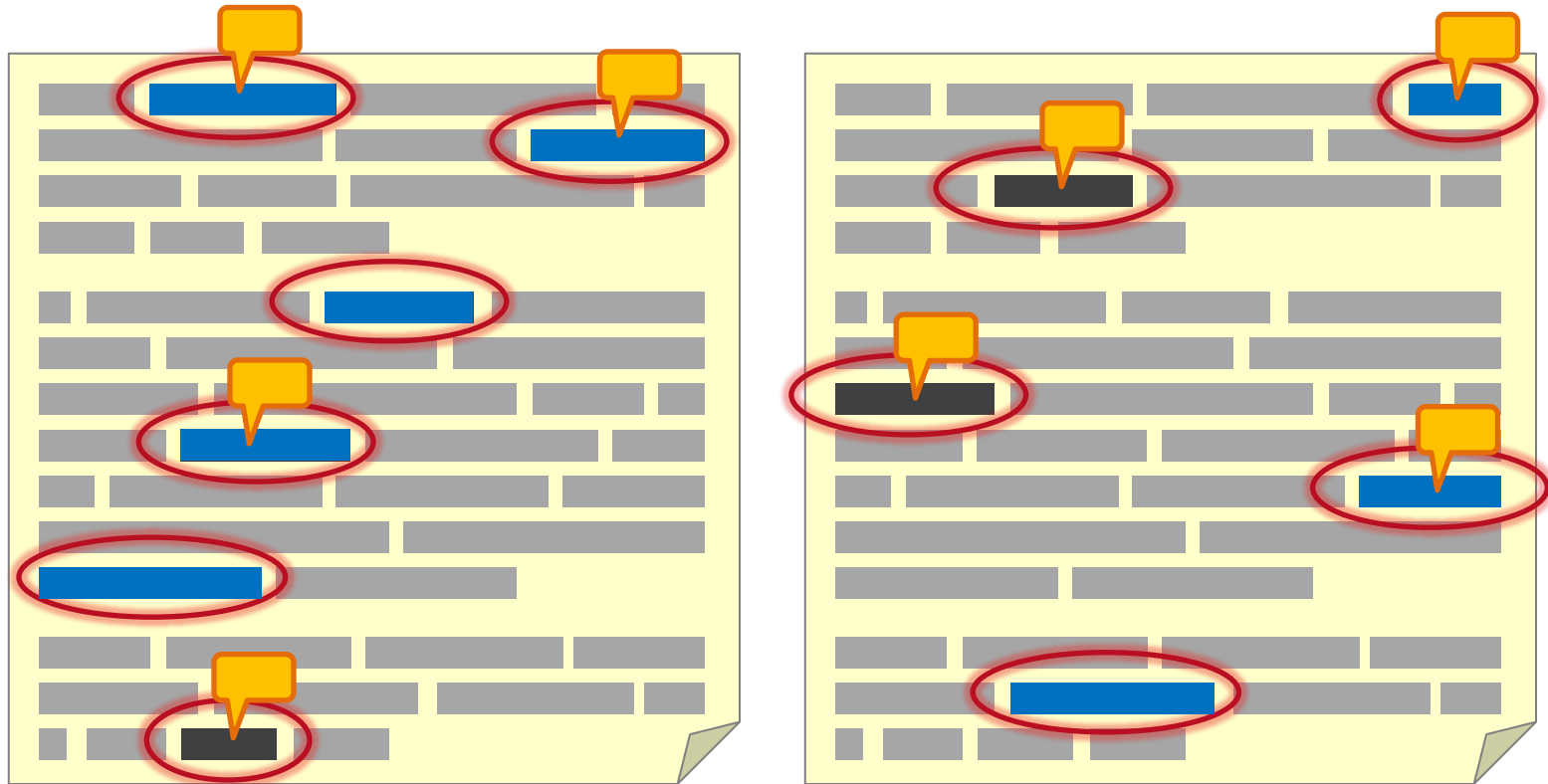206 words

**$T_2$ Wikipedia article**
183 words

# Experimental Setup: Data



**11 text positions**  **8 introduced issues**
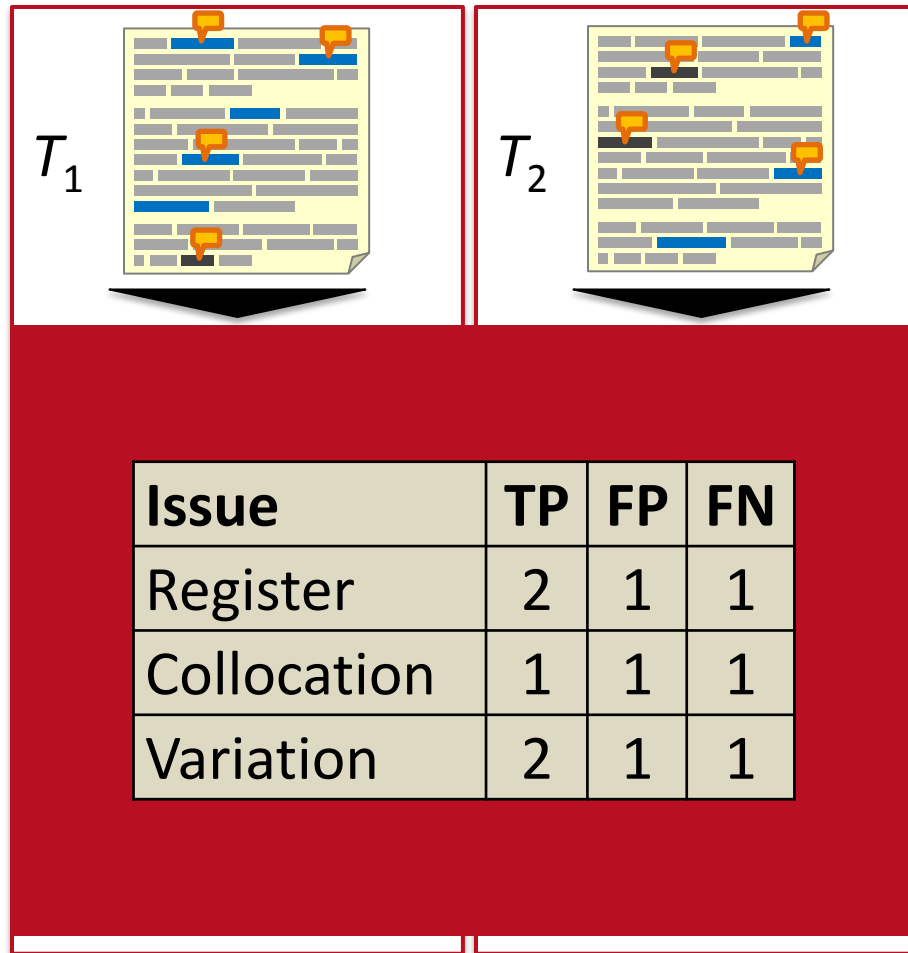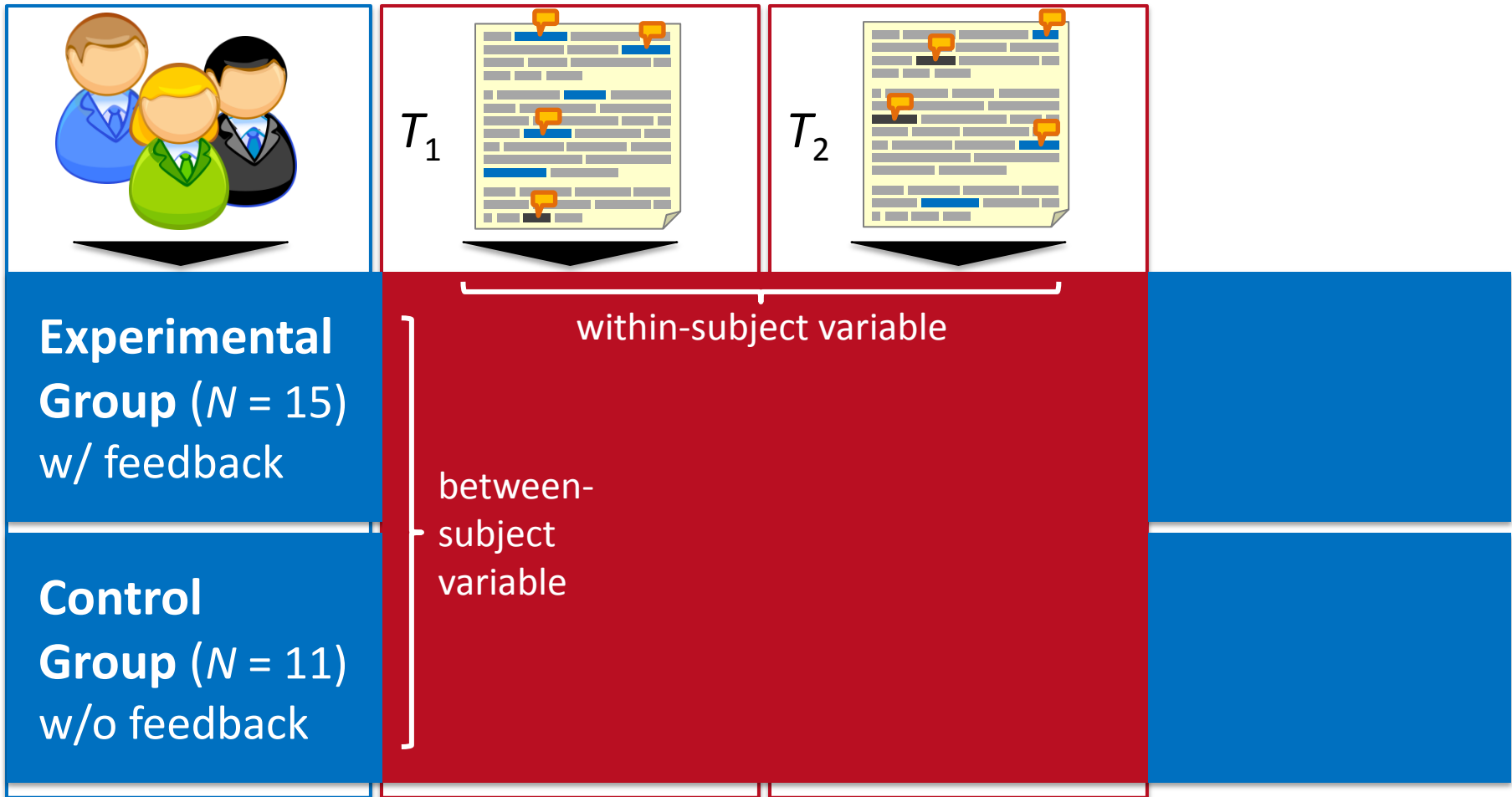
# Experimental Setup: Data



TP — correct feedback
FP — incorrect feedback
FN — incomplete feedback

# Experimental Setup: Data



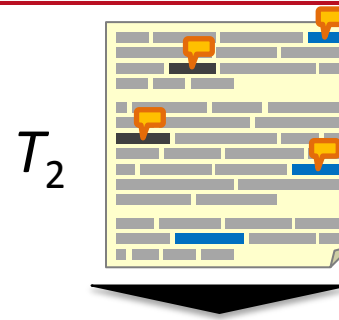| Issue | TP | FP | FN |
|---|---|---|---|
| Register | 2 | 1 | 1 |
| Collocation | 1 | 1 | 1 |
| Variation | 2 | 1 | 1 |

# Experimental Setup: Population



$T_1$ $T_2$

within-subject variable

**Experimental Group** ($N$ = 15) w/ feedback

**Control Group** ($N$ = 11) w/o feedback

between-subject variable

# Experimental Setup: Tool



$T_1$

$T_2$

**Experimental Group** ($N$ = 15) w/ feedback

**Control Group** ($N$ = 11) w/o feedback

# User Study

**New tool: InViEdit**
https://github.com/UKPLab/naacl-bea2016-writing-study

# Writing Assistance Software

https://github.com/UKPLab/naacl-bea2016-writing-study

**TECHNISCHE UNIVERSITÄT DARMSTADT**

Save Progress

Text editor

Selected highlight

Feedback messages

Discontinuous highlight

**System Usability Scale**
**SUS = 76.3**

> 68.0 "acceptable"

> 71.4 "good"

# Experimental Setup: Analysis



Experimental Group (*N* = 15) w/ feedback

Control Group (*N* = 11) w/o feedback

$T_1$

$T_2$

## User Study

### New tool: InViEdit

https://github.com/UKPLab/naacl-bea2016-writing-study

revised vs. not revised positions

revised vs. not revised positions

H1–H4

# Data Analysis

$$\begin{array}{rl}
11 & \text{positions (TP/FP/FN)} \\
\times \quad 26 & \text{participants} \\
\hline
= \quad 286 & \text{data points}
\end{array}$$

Data point $x$ = (revised vs. not revised)

|    | min(x) | $\overline{x}$ | SE | max(x) |
|----|-------|------|------|------|
| EG | 2     | 5.86 | 0.53 | 10   |
| CG | 0     | 3.18 | 0.74 | 8    |

Unpaired two sample Student's $t$ test
with significance level $\alpha = 0.05$  (P $\leq$ 0.05)

# H1: Correct Feedback helps

**Expectation**: $\mu_{EG(TP)} \neq \mu_{CG(TP)}$

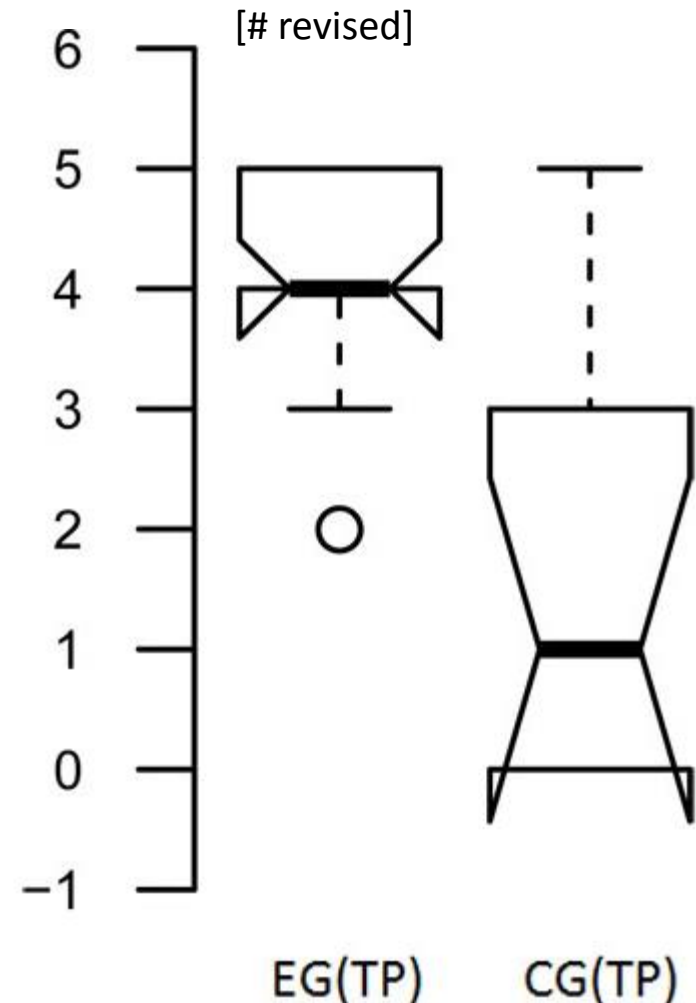**Arithmetic mean:**

$\overline{x}_{EG(TP)} = 4.13$  (SE = 0.23)

$\overline{x}_{CG(TP)} = 1.63$  (SE = 0.51)

**Test statistic:**

$t_{H1} = 4.85$

$|t_{H1}| > 2.06$   ($P < 0.0001$)

- reject null hypothesis at 5% level
- significant difference b/w groups ✓



[# revised]

EG(TP)     CG(TP)

# H2: Incorrect Feedback causes unnecessary revisions

**Expectation**: $\mu_{EG(FP)} \neq \mu_{CG(FP)}$

**Arithmetic mean:**
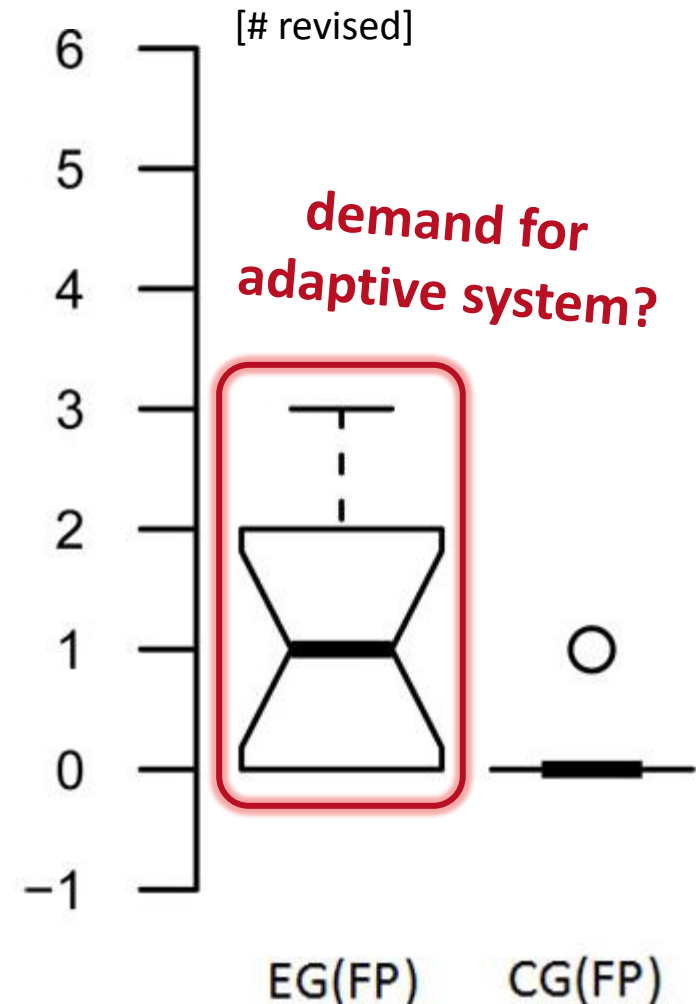
$\overline{x}_{EG(FP)} = 1$ $\quad$ (SE = 0.25)

$\overline{x}_{CG(FP)} = 0.18$ $\quad$ (SE = 0.12)

**Test statistic:**

$t_{H2} = 2.55$

$|t_{H2}| > 2.06$ $\quad$ ($P = 0.017$)

- reject null hypothesis at 5% level
- significant difference b/w groups ✓

[# revised]

demand for adaptive system?

EG(FP) $\quad$ CG(FP)

# H3: Incomplete Feedback causes users to miss similar issues

**Expectation**: $\mu_{EG(FN)} \neq \mu_{CG(FN)}$

**Arithmetic mean:**

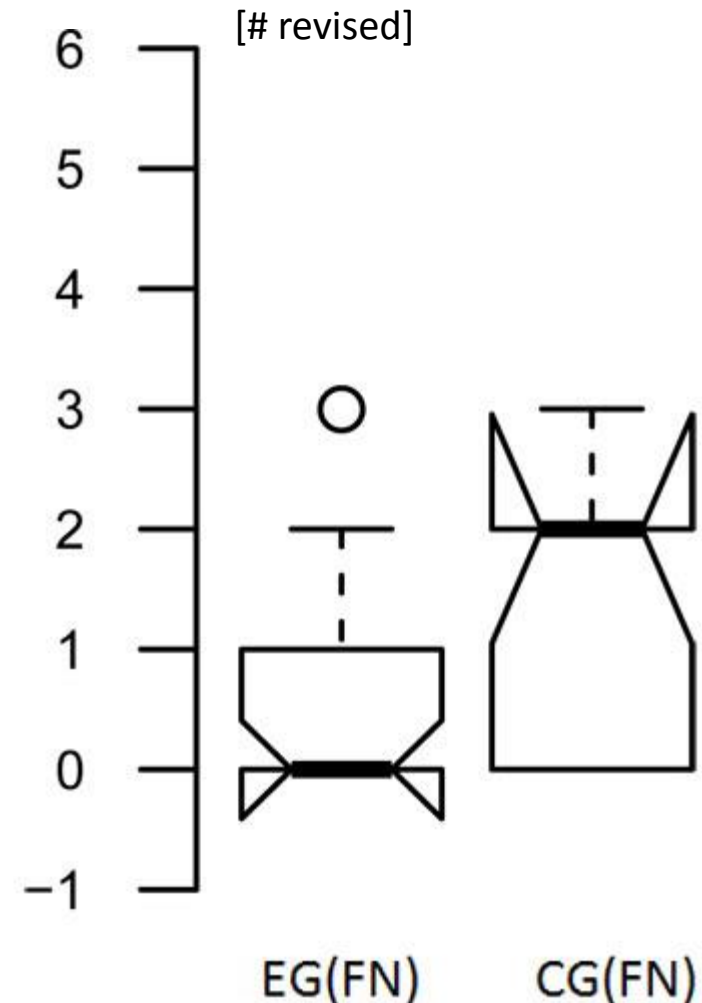$\overline{x}_{EG(FN)} = 0.73$  (SE = 0.28)

$\overline{x}_{CG(FN)} = 1.36$  (SE = 0.36)

**Test statistic:**

$t_{H3} = -1.39$

$|t_{H3}| < 2.06$   ($P = 0.17$)

- cannot reject null hypothesis
- cannot find significant difference ✖



[# revised]

EG(FN)     CG(FN)

# H4: Task completion time similar

**Expectation**: $\mu_{EG(\tau)} = \mu_{CG(\tau)}$

**Arithmetic mean:**

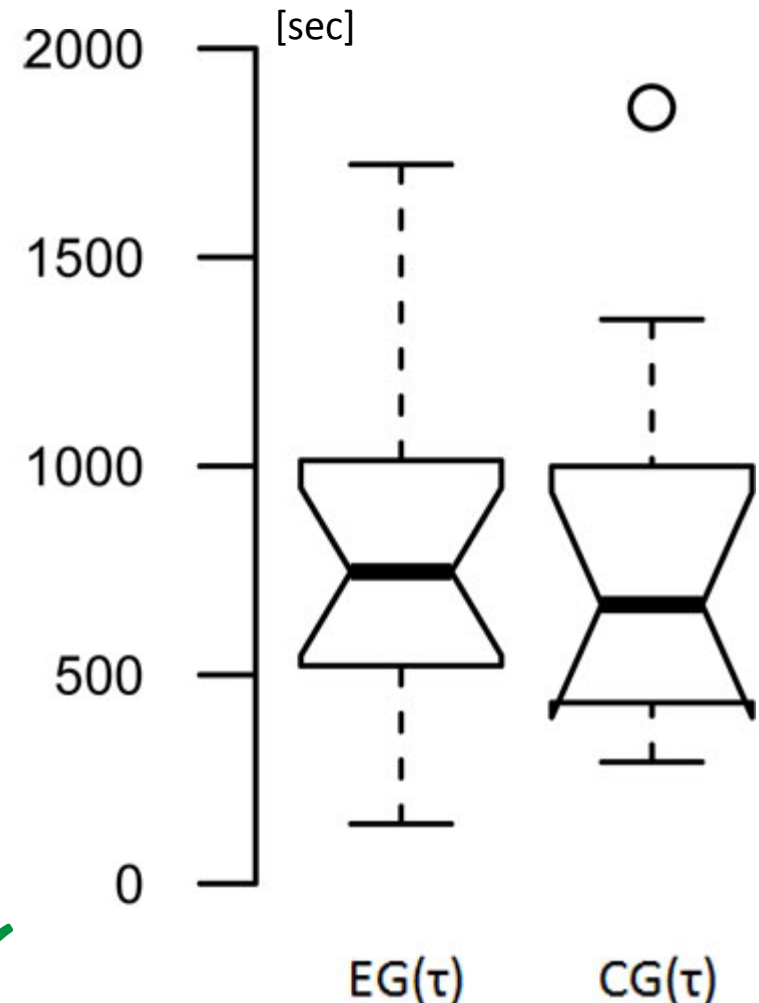$\overline{x}_{EG(\tau)}$ = 13 min,   3 sec  (SE = 104 sec)

$\overline{x}_{CG(\tau)}$ = 13 min, 27 sec  (SE = 144 sec)

**Test statistic:**

$t_{H4}$ = −0.14

$|t_{H4}|$ < 2.06   ($P$ = 0.89)

- cannot reject null hypothesis
- cannot find significant difference ✓

# Conclusion

**What we learned from this work:**

- correct feedback helps (H1)

- incorrect feedback problematic (H2) – overtrust?

  - but: demand for **adaptive systems**!

- tendency to miss FNs, but not significant (H3)

  - confirms "**precision** more important than recall"

- using feedback didn't take longer (H4)

**Software for evaluating writing assistance tools:**
https://github.com/UKPLab/naacl-bea2016-writing-study

# Kontakt / Contact

**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