

Lernen zu Löschen: Vergessen digitaler Objekte als Gemeinschaftsaufgabe von Mensch und KI



Beitrag der Informatik

(Weiter-) Entwicklung eines partnerschaftlichen KI-Systems, das unterstützt, irrelevante Dateien zu identifizieren und zu löschen. Mit Induktiver Logischer Programmierung (ILP) werden wissensbasierte KI und maschinelles Lernen integriert. Lernen wird erklärungsorientiert und interaktiv umgesetzt, so dass der Mensch Systementscheidungen nachvollziehen und korrigieren kann.

Mensch-KI-Partnerschaft durch erklärbares und interaktives maschinelles Lernen

Name	Media Type	Size	Creation Date	Change Date	Access Date
LyfWtRmMK6.pdf	application/pdf	2995540	2005-06-14 07:24	2012-01-14 00:40	2014-05-06 11:07
RoMNTYSM.doc	application/msword	45705	2017-03-18 09:24	2019-02-26 08:56	2019-08-18 11:58
T0sWb5MTAE.doc	application/msword	54403	2013-08-28 12:34	2015-03-18 12:20	2021-03-03 15:27
U9ShWb5c.PDF	application/pdf	113533	2018-12-30 08:00	2019-05-16 10:13	2019-12-20 23:48
Vm49FtC5bEvVh.RTF	application/msword	47316	2018-12-28 20:21	2020-05-22 23:03	2020-09-03 20:23
gHIP6s.png	image/png	16042	2019-11-21 01:13	2020-06-22 08:38	2020-08-02 00:54
gYvZkGjereport - Copy (1).doc	application/msword	4367	2021-08-09 01:24	2021-08-19 07:42	2021-08-21 20:29
IGZTwTidX.png	image/png	10410	2017-08-22 04:58	2019-03-26 10:41	2019-08-07 09:40
r4009fmw.PNG	image/png	5236	2015-02-03 20:13	2017-05-16 08:00	2017-11-22 05:49
r4009fmw - Copy.PNG	image/png	5236	2017-08-19 18:53	2017-09-11 03:18	2017-10-16 15:44
r4009fmw - Copy (1).PNG	image/png	5236	2017-09-25 17:35	2017-10-09 07:13	2017-10-27 15:05

Abbildung: Datei r4009fmw.png wäre nicht als irrelevant klassifiziert, wenn sie im Vergleich zu r4009fmw-Copy.png neuer wäre.

Dare2Del nutzt

- Metainformationen von Dateien
- domänenspezifisches Wissen und Hintergrundtheorie zum Lernen von Irrelevanzentscheidungen

Dare2Del liefert Nutzenden Löschvorschläge

- die akzeptiert oder abgelehnt werden können

Dare2Del liefert Erklärungen für Irrelevanzentscheidungen

- die korrigierbar sind

Beitrag der Psychologie

Löschen soll so gestaltet und unterstützt werden, dass Personen sich tatsächlich damit auseinandersetzen. Dadurch sollen sich positive Konsequenzen sowohl für eigene als auch für organisationale Arbeitsabläufe ergeben.

Forschungsfragen Psychologie

- Welche Gestaltungsmerkmale (z.B. Erklärungen) muss das Assistenzsystem Dare2Del haben, um bei Löscheentscheidungen zu unterstützen?
- Inwieweit kann menschliches Vergessen durch Löschen unterstützt werden?
- Welche personenbezogenen Fähigkeiten (z.B. Fähigkeit zur Inhibition) gilt es zu beachten?
- Welche organisationalen Bedingungen (z.B. soziale Normen) haben einen Einfluss?

Ergebnisse der zweiten Projektphase

Psychologie

AP-P-1: Experimente zur Gestaltung der Löschvorschläge (mit vs. ohne Erklärungen)

- Erklärungen führen zu einer höheren Löschbereitschaft und helfen, das System glaubwürdiger und vertrauenswürdiger zu machen – vor allem, wenn sie überprüfbar sind
- Erklärungen reaktivieren kurzfristig aber auch Dokumentnamen und -inhalte (kein Vergessen)

AP-P-2: Experimente zu den Auswirkungen des Löschens mit und ohne Assistenzsystem

- Löschvorschläge des Assistenzsystems werden von etwa 50% der Nutzenden angenommen
- Die Vorgabe sozialer Normen zum Löscherhalten hat vergleichbare Effekte
- Löschen geht tendenziell mit geringerer Beanspruchung einher
- Individuelle Inhibitionsfähigkeit geht (unabhängig vom Löschen) mit höherer Aufgabenleistung einher

Informatik

Erklärbarkeit

- Verbale Erklärungen in unterschiedlichen Detaillierungsgraden
- Near-Miss Erklärungen

Interaktives Lernen

- Klassifikationsentscheidung und Erklärung durch Nutzer/innen modifizierbar
- Feste Regeln (Vorschriften) dürfen nicht durch Lernen änderbar sein
- 'Human-in-the-loop': Kompetenzerhalt, kein blindes Vertrauen

Ausgewählte Publikationen

Göbel, K., Niessen, C., Seufert, S., & Schmid, U. (2022). Explanatory machine learning for justified trust in human-AI collaboration: Experiments on file deletion recommendations. *Frontiers in Artificial Intelligence*, 5, Article 919534. <https://doi.org/10.3389/frai.2022.919534>

Göbel, K., Hensel, L., Schultheiss, O. C., & Niessen, C. (2022). Meta-analytic evidence shows no relationship between task-based and self-report measures of thought control. *Applied Cognitive Psychology*, 36(3), 659-672. <https://doi.org/10.1002/acp.3952>

Kiefer, S., Hoffmann, M., & Schmid, U. (2022). Semantic interactive learning for text classification: A constructive approach for contextual interactions. *Machine Learning and Knowledge Extraction*, 4(4), 994-1010.

Rabold, J., Siebers, M., & Schmid, U. (2022). Generating contrastive explanations for inductive logic programming based on a near miss approach. *Machine Learning*, 111(5), 1799-1820.

Ai, L., Muggleton, S.H., Hocquette, C., Gromowski, M., & Schmid, U. (2021). Beneficial and harmful explanatory machine learning. *Machine Learning*, 110(4), 695-721. <https://doi.org/10.1007/s10994-020-05941-0>

Göbel, K., & Niessen, C. (2021). Thought control in daily working life: How the ability to stop thoughts protects self-esteem. *Applied Cognitive Psychology*, 35(4), 1011 - 1022. <https://doi.org/10.1002/acp.3830>

Schmid, U. (2021). Interactive learning with mutual explanations in relational domains. In S. Muggleton and N. Charter, editors, *Human-like Machine Intelligence*, pp. 337-353. Oxford University Press. <https://doi.org/10.1093/oso/9780198862536.001.0001>

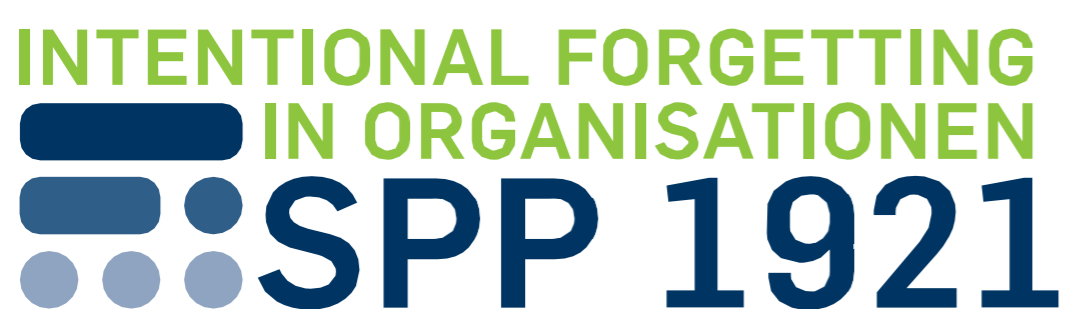
Gromowski, M., Siebers, M., & Schmid, U. (2020). A process framework for inducing and explaining Datalog theories. *Adv. Data Anal. Classif.* 14(4), 821-835. <https://doi.org/10.1007/s11634-020-00422-7>

Niessen, C., Göbel, K., Lang, J., & Schmid, U. (2020). Stop thinking: An experience sampling study on suppressing distractive thoughts at work. *Frontiers in Psychology*, 11, Article 1616. <https://doi.org/10.3389/fpsyg.2020.01616>

Niessen, C., & Lang, J. W. B. (2020). Cognitive control strategies and adaptive performance in a complex work task. *Journal of Applied Psychology*. Advance online publication. <https://doi.org/10.1037/apl0000830>

Niessen, C., Göbel, K., Siebers, M., & Schmid, U. (2020). Time to forget: Intentional forgetting in the digital world of work. *Zeitschrift für Arbeits- und Organisationspsychologie*, 64, 30-45. <https://doi.org/10.1026/0932-4089/a000308>

Siebers, M., & Schmid, U. (2019). Please delete that! Why should I? KI-Künstliche Intelligenz, 33(1), 35-44. <https://doi.org/10.1007/s13218-018-0565-5>



Gefördert durch



Lehrstuhl für Psychologie im Arbeitsleben
Prof. Dr. Cornelia Niessen

Projektmitarbeitende

Dr. Kyra Göbel¹
Angelina Olivia Wilczewski, B.Sc.²
Durgesh Nandini, M.Sc.²
(Projektmitglied 2020 – 2023)
Dipl.-Psych. Michael Siebers, B.Sc.²
(Projektmitglied 2016 – 2019)

¹ Friedrich-Alexander-Universität Erlangen-Nürnberg
² Otto-Friedrich-Universität Bamberg



Otto-Friedrich-Universität Bamberg



Professur Angewandte Informatik, insbes.
Kognitive Systeme
Prof. Dr. Ute Schmid