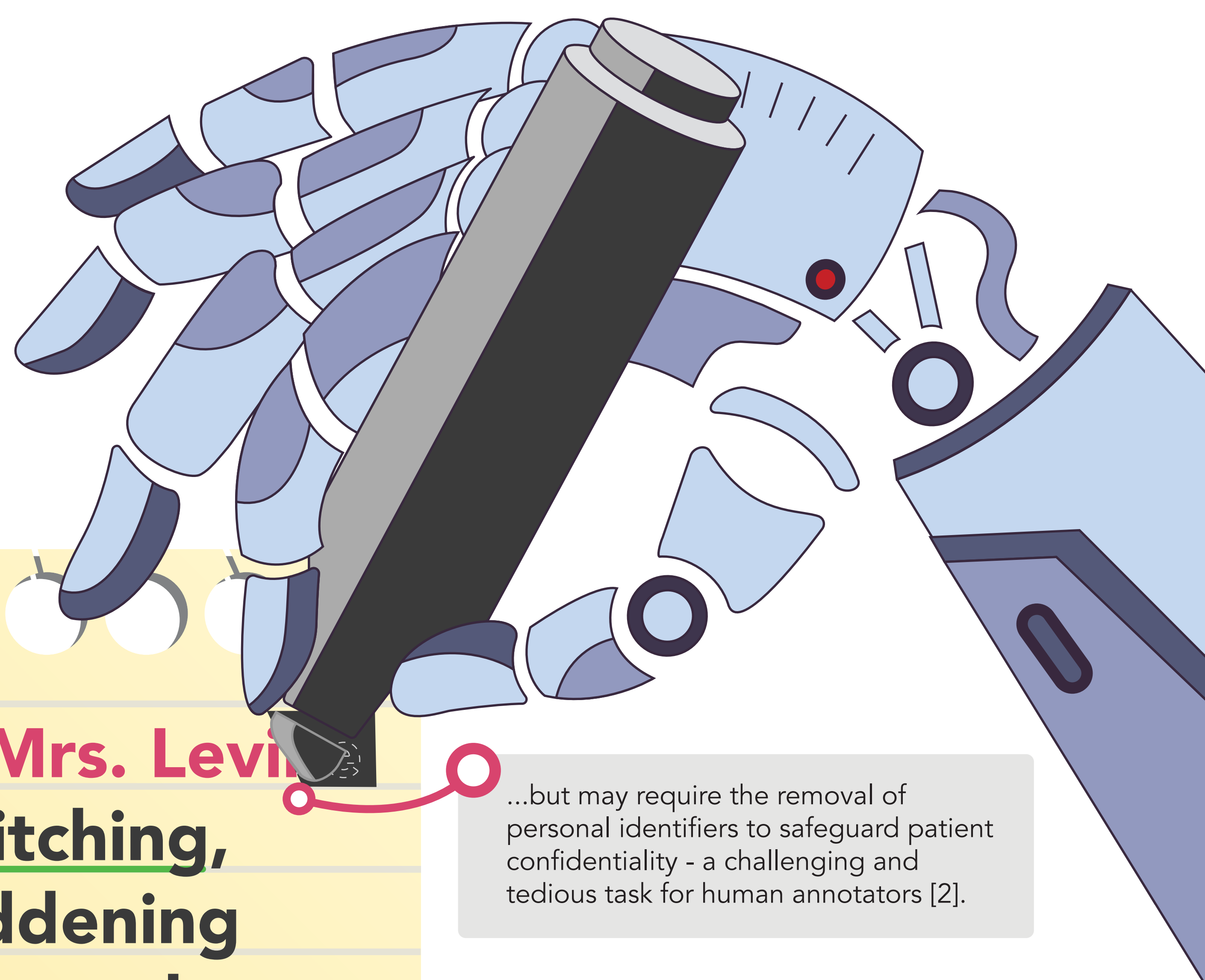


Automatic detection and removal of personal identifiers in case narratives using deep learning

Eva-Lisa Meldau, Sara Hedfors Vidlin, Lucie Gattepaille, Henric Taavola, Lovisa Sandberg, Yasunori Aoki, G. Niklas Norén **Uppsala Monitoring Centre, Uppsala, Sweden**



Case narratives can be crucial in the assessment of the causality and the clinical course of suspected adverse drug reactions...

...since they can provide information which may not be available in structured fields [1], for example about the severity of the reaction or the impact on the patient's quality of life.

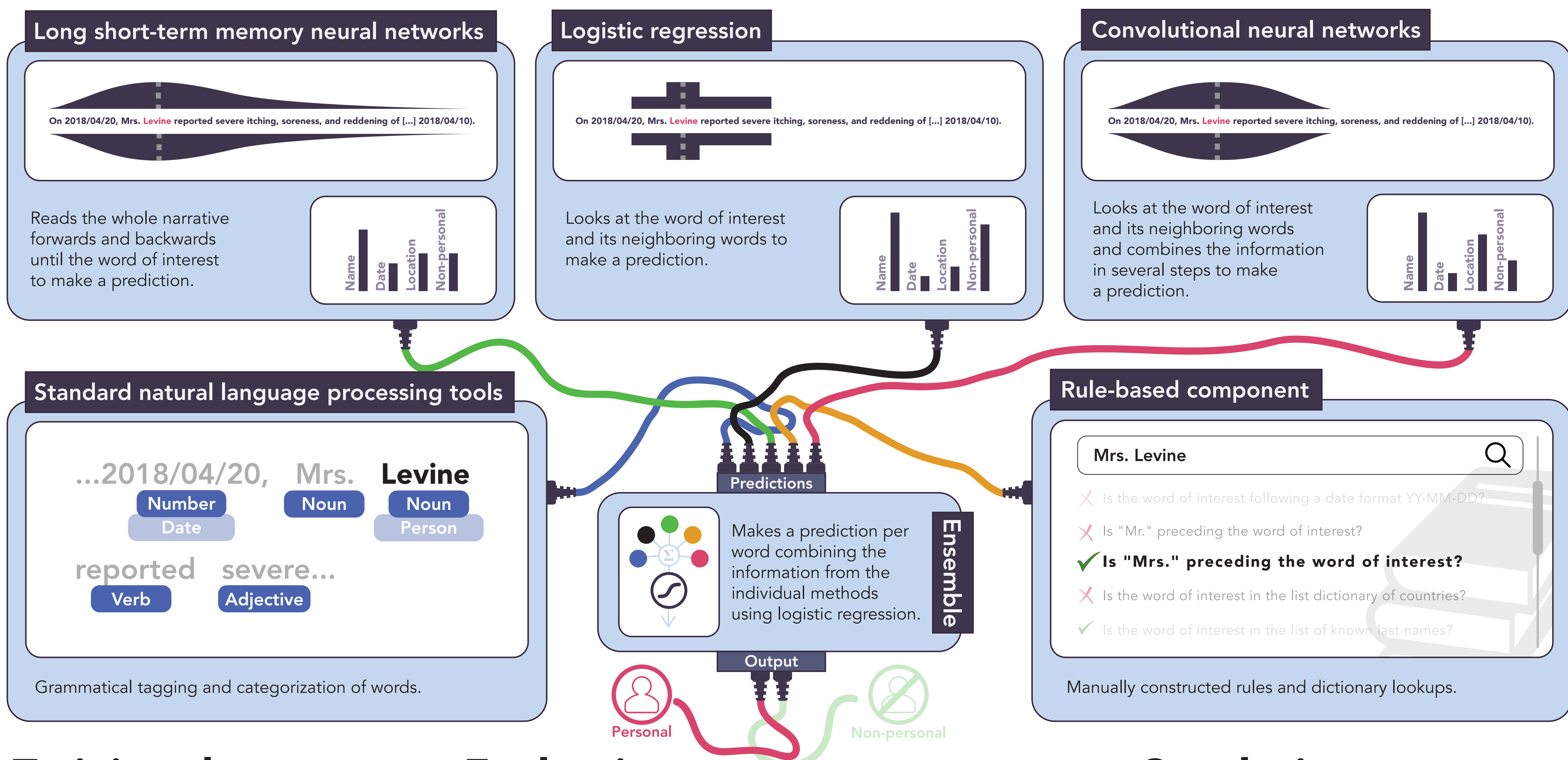
Sharing these narratives between organizations is important to allow for all involved to have the relevant information to fulfil their mission of protecting patients...

On **2018/04/20**, **Mrs. Levine** reported **severe itching, soreness, and reddening of the genital area** and an **inability to sit** while on therapy with **dapagliflozin** (started on **2018/04/10**).

...but may require the removal of personal identifiers to safeguard patient confidentiality - a challenging and tedious task for human annotators [2].

In this study, we aimed to develop an automatic de-identification method for narratives using deep learning (neural networks). We considered names, dates, and locations as personal identifiers, and all other words as 'non-personal' identifiers.

Methods:



Training data:

Individual algorithms: 521 medical records from the training set of the 2014 i2b2 de-identification challenge data set. [3]

Ensemble: 3/4 of the 269 records in the validation set from the same data set.

Evaluation:

i2b2

95% of the personal identifiers were removed
90% of the removed words were personal identifiers

On held-out 1/4 of the records in the i2b2 2014 validation data set.

VigiBase

90% of the personal identifiers were removed
45% of the removed words were personal identifiers

On 300 narratives from VigiBase, the WHO global database of individual case safety reports [4].

Conclusions:

The algorithm removed a greater proportion of the personal identifiers in i2b2 records than did human annotators [2] but at the expense of removing more of the other text. The performance on VigiBase narratives is promising considering that our method was only trained on medical records. With access to original narratives that can be annotated for further training and fine-tuning, we expect the performance to improve even further.

Acknowledgments:

The authors are indebted to the national centres who make up the WHO Programme for International Drug Monitoring and contribute reports to VigiBase. However, the opinions and conclusions of this study are not necessarily those of the various centres nor of WHO.

References:

- [1] Karimi, G, et al. Clinical stories are necessary for drug safety. Clin Med, 14(3): 326-327, 2014.
- [2] South, BR, et al. Evaluating the effects of machine pre-annotation and an interactive annotation interface on manual de-identification of clinical text. J Biomed Inform 2014; 50: 162-172.
- [3] Stubbs, A, Uzuner, O. Annotating longitudinal clinical narratives for de-identification: The 2014 i2b2/UTHealth corpus. J Biomed Inform, 58(Suppl): S20-29, 2015.
- [4] Lindquist, M. VigiBase, the WHO Global ICSR Database System: Basic Facts. Drug Inf J, 42(5): 409-19, 2008.

Access full resolution poster: bit.ly/umc-posters