PrivateNLP2020

# The Second Workshop on Privacy in Natural Language Processing

Nov 20, 2020
Online

Order copies of this and other ACL proceedings from:

# Introduction

The PrivateNLP workshop aims to bring together practitioners and researchers from academia and industry to discuss the challenges and approaches to designing, building, verifying, and testing privacy preserving systems in the context of Natural Language Processing.

The workshop received 11 paper submissions this year and accepted 5 papers and one non-archival paper. These accepted papers cover membership inference attacks, text perturbation mechanisms, privacy mechanisms for language understanding tasks, automatic third-party identification in privacy policies, semantic similarity within privacy policies, and performance boosting of differentially private language models.

There are 5 invited speakers, Aaron Roth (University of Pennsylvania), Reza Shokri (National University of Singapore), Krishnaram Kenthapadi (Amazon AWS), Annabelle McIver (Macquarie University), and Mark Dras (Macquarie University). Each talk will cover different aspects of privacy-preserving NLP. We would like to thank the 33 Program Committee members who kindly reviewed the submissions, as well as the invited speakers, and the workshop co-organizers, Oluwaseyi Feyisetan (Amazon, USA), Sepideh Ghanavati (University of Maine, USA), Shervin Malmasi (Amazon, USA), and Patricia Thaine (University of Toronto, Canada).

**Organizers:**

Oluwaseyi Feyisetan, Amazon (USA)
Sepideh Ghanavati, University of Maine (USA)
Shervin Malmasi, Amazon (USA)
Patricia Thaine, University of Toronto (Canada)

**Program Committee:**

Aleksei Triastcyn, École Polytechnique Fédérale de Lausanne, (Switzerland)
Andreas Nautsch, EURECOM, (France)
Arne Köhn, Saarland University, (Germany)
Asma Eidhah Aloufi, Rochester Institute of Technology, (USA)
Balazs Pejo, Budapest University of Technology and Economics, (Hungary)
Benjamin Zi Hao Zhao, University of New South Wales, (Australia)
Briland Hitaj, SRI International, (USA)
Christian Weinert, Technische Universität Darmstadt, (Germany)
Congzheng Song, Cornell University, (USA)
Dinusha Vatsalan, Data61-CSIRO, (Australia)
Eleftheria Makri, Saxion University, (The Netherlands)
Elette Boyle, IDC Herzliya, (Israel)
Fang Liu, University of Notre Dame, (USA)
Gerald Penn, University of Toronto, (Canada)
Isar Nejadgholi, National Research Council, (Canada)
Jamie Hayes, University College London, (UK)
Jason Xue, University of Adelaide, (Australia)
Jaspreet Bhatia, (USA)
Julius Adebayom MIT, (USA)
Kambiz Ghazinour, State University of New York, (USA)
Ken Barker, University of Calgary, (Canada)
Liwei Song, Princeton University, (USA)
Luca Melis, Amazon, (USA)
Maximin Coavoux (University of Edinburgh), (UK)
Mitra Bokaei Hosseini, St. Mary's University, (USA)
Natasha Fernandes, Macquarie University, (Australia)
Nedelina Teneva, Amazon, (USA)
Peizhao Hu, Rochester Institute of Technology, (USA)
Sai Teja Peddinti, Google, (USA)
Shomir Wilson, Pennsylvania State University, (USA)
Tom Diethe, Amazon, (UK)
Travis Breaux, Carnegie Mellon University, (USA)
Xavier Ferrer, King's College London, (UK)

**Invited Speaker:**

Aaron Roth, University of Pennsylvania, (USA)
Reza Shokri, National University of Singapore, (Singapore)
Krishnaram Kenthapadi, Amazon AWS, (USA)
Annabelle McIver, Macquarie University, (Australia)
Mark Dras, Macquarie University, (Australia)

# Table of Contents

vii

# Conference Program

8:45–9:00      *Welcome*
Seyi Feyisetan

**9:00–10:15**      **Session 1**

**9:00–9:45**      ***Invited Talk - Aaron Roth***

9:45–10:15      *On Log-Loss Scores and (No) Privacy*
Abhinav Aggarwal, Zekun Xu, Oluwaseyi Feyisetan and Nathanael Teissier

**10:30–12:15**      **Session 2**

**10:30–11:15**      ***Invited Talk 2 - Reza Shokri***

11:15–11:45      *A Differentially Private Text Perturbation Method Using Regularized Mahalanobis Metric*
Zekun Xu, Abhinav Aggarwal, Oluwaseyi Feyisetan and Nathanael Teissier

11:45–12:15      *TextHide: Tackling Data Privacy in Language Understanding Tasks*
Yangsibo Huang, Zhao Song, Danqi Chen, Kai Li and Sanjeev Arora

**13:00–14:45**      **Session 3**

**13:00–13:45**      ***Invited Talk 3 - Mark Dras and Annabelle McIver***

13:45–14:15      *Identifying and Classifying Third-party Entities in Natural Language Privacy Policies*
Mitra Bokaie Hosseini, Pragyan K C, Irwin Reyes and Serge Egelman

14:15–14:45      *Surfacing Privacy Settings Using Semantic Matching*
Rishabh Khandelwal, Asmit Nayak, Yao Yao and Kassem Fawaz

**15:00–17:00   Session 4**

**15:00–15:45   *Invited Talk 4 - Krishnaram Kenthapadi***

15:45–16:15   *Differentially Private Language Models Benefit from Public Pre-training*
Gavin Kerrigan, Dylan Slack and Jens Tuyls

**16:15–16:45   *Finding Paper***

16:45–17:00   *Closing Remarks*
Seyi Feyisetan