AACL-IJCNLP 2022

**The 1st Workshop on
Information Extraction from Scientific Publications**

**Proceedings of the Workshop**

November 20, 2022

# Preface

The number of scientific papers published per year has exploded in recent years, strengthening its value as one of the main drivers for scientific progress. In astronomy alone, more than 41,000 new articles are published every year and the vast majority are available either via an open-access model or via pre-print services. Indexing the article's full-text in search engines helps discover and retrieve vital scientific information to continue building on the shoulders of giants, informing policy, and making evidence-based decisions. Nevertheless, it is difficult to navigate in this ocean of data; finding articles rely heavily on string matching searches and following citations/references. Still, new approaches are necessary to differentiate the signal from the noise more easily (e.g., finding the key articles about the medical condition we are interested in).

Simple string matching has substantial limitations, human language is ambiguous in nature, context matters, and we frequently use the same word and acronyms to represent a multitude of different meanings. Extracting structured and semantically relevant information from scientific publications (e.g., named-entity recognition, summarization, citation intention, linkage to knowledge graphs) allows better selection and filter articles.

The Workshop on Information Extraction from Scientific Publications (WIESP) is a forum to foster discussion and research using Natural Language Processing and Machine Learning. In this space, leading professionals, organizations, early career researchers and students can cooperate towards building the algorithms, models, and tools that will pave the way for machine comprehension of science in the future.

WIESP received 25 submissions, of which 16 were accepted (8 long papers, 4 short papers, and 4 shared task system papers).

WIESP was held on November 20th 2022.

# Organizing Committee

Tirthankar Ghosal, Charles University, CZ
Sergi Blanco-Cuaresma, Center for Astrophysics | Harvard & Smithsonian, USA
Alberto Accomazzi, Center for Astrophysics | Harvard & Smithsonian, USA
Robert M. Patton, Oak Ridge National Laboratory, USA
Felix Grezes, Center for Astrophysics | Harvard & Smithsonian, USA
Thomas Allen, Center for Astrophysics | Harvard & Smithsonian, USA

# Program Committee

Min-Yuh Day
Hen-Hsen Huang
Jheng-Long Wu
Daniel Acuna
Akiko Aizawa
Hamed Alhoori
Atilla Kaan Alkan
Thomas Allen
Hardik Arora
Premjith B
Partha Basuchowdhuri
Saprativa Bhattacharjee
Yimeng Dai
Xiang Dai
Vignesh Edithal
Sergey Feldman
Edward Fox
Zheng Gao
Daisuke Ikeda
Sarvnaz Karimi
Valia Kordoni
Sarvnaz Karimi
Valia Kordoni
Asheesh Kumar
Rishu Kumar
Sandeep Kumar
Xiangci Li
Shigeki Matsubara
Yoshitomo Matsubara
Sujit Pal
Rajesh Piryani
Trinita Roy
Neil Smalheiser
Wojtek Sylwestrzak
Rohan Tondulkar
George Tsatsaronis

Jan Philip Wahle

Ronin Wu

Wuhe Zou

# Table of Contents

# Conference Program