FinNLP 2020

# The Second Workshop on Financial Technology
# and Natural Language Processing
# in conjunction with IJCAI-PRICAI 2020

## Proceedings of the Workshop

January 5, 2021

Kyoto, Japan

# Sponsor

AIntu 科技部臺灣大學人工智慧研究中心

科技部人工智慧技術暨全幅健康照護聯合研究中心

MOST JOINT RESEARCH CENTER FOR AI TECHNOLOGY AND
ALL VISTA HEALTHCARE

# Preface

It is our great pleasure to welcome you to the Second Workshop on Financial Technology and Natural Language Processing (FinNLP).

The aim of FinNLP is to provide a forum for international participants to share knowledge on applying NLP to the FinTech domain. With the sharing of the researchers in this workshop, we hope that the challenging problems of blending FinTech and NLP will be identified, the future research directions will be shaped, and the scope of this interdisciplinary research area will be broadened.

Due to a great trial of all of us, the COVID-19 virus, the IJCAI-PRICAI conference is postponed to January 2021. In order to accelerate the development of this field, we decide to publish the proceedings in advance. This year, the participants of FinNLP still bring several novel ideas to this forum. We also cooperate with Fortia Financial Solutions to hold two shared tasks in FinNLP, including learning semantic representations (FinSim) and sentence boundary detection in PDF noisy text (FinSBD-2).

We have many people to thank. Dialekti Valsmou Stanislawski and Ismail El Maarouf lead their teams to hold the successful shared tasks and help review several submissions. All of program committee members work very hard to provide their insightful comments to the submissions, and help us select the suitable papers for FinNLP-2020. Many thanks to all participants for submitting their interesting works and sharing their ideas. Besides, we would like to express our gratitude to the MOST Joint Research Center of AI Technology and All Vista Healthcare for financial support.

Chung-Chi Chen, Hen-Hsen Huang, Hiroya Takamura, Hsin-Hsi Chen
FinNLP 2020 Organizers
July 2020

# Organizations

**Organizers**

Chung-Chi Chen, National Taiwan University

Hen-Hsen Huang, National Chengchi University

Hiroya Takamura, National Institute of Advanced Industrial Science and Technology

Hsin-Hsi Chen, National Taiwan University

**Shared Task Organizers**

Willy Au

Abderrahim Ait Azzi

Bianca Chong

Ismail El Maarouf

Youness Mansar

Virginie Mouilleron

Dialekti Valsamou-Stanislawski

**Program Committee**

Paulo Alves, Universidade Católica Portuguesa

Alexandra Balahur, European Commission's Joint Research Centre

Paul Buitelaar, Insight Centre for Data Analytics at NUIG

Damir Cavar, Indiana University

Pablo Duboue, Textualization Software Ltd.

Jinhua Du, American International Group, Inc.

Ismail El Maarouf, Fortia Financial Solutions

Sira Ferradans, Independent Researcher

Kiyoshi Izumi, The University of Tokyo

Changliang Li, Kingsoft Corporation

Nedim Lipka, Adobe Inc.

Hiroki Sakaji, The University of Tokyo

Dialekti Valsamou-Stanislawski, Fortia Financial Solutions

Chuan-Ju Wang, Academia Sinica

Annie T.T. Ying, EquitySim

Wlodek Zadrozny, University of North Carolina in Charlotte

# Table of Contents