# 2021 IEEE Spoken Language Technology Workshop (SLT 2021)

Shenzhen, China
19 – 22 January 2021

Pages 1-566

**Additional Copies of This Publication Are Available From:**

Curran Associates, Inc
57 Morehouse Lane
Red Hook, NY  12571 USA
Phone:        (845) 758-0400
Fax:          (845) 758-2633
E-mail:       curran@proceedings.com
Web:          www.proceedings.com

CURRAN ASSOCIATES INC.
proceedings
.com

# TABLE OF CONTENTS

## SPEECH RECOGNITION: TRAINING METHODS

## SPEECH RECOGNITION: DATA AND FEATURE AUGMENTATION

# SPEECH RECOGNITION: GENERAL TOPICS

# SPEAKER RECOGNITION

## EMOTION RECOGNITION

# EXPRESSIVE SPEECH SYNTHESIS AND OTHERS

# NEURAL VOCODER AND OTHERS

## ANTI-SPOOFING AND DIARIZATION IN SPEAKER RECOGNITION

## MULTIMODAL PROCESSING

## SPEECH ANALYSIS AND ENHANCEMENT

## SPECIAL SESSION ON INTEGRATION OF SPEECH SEPARATION, RECOGNITION AND DIARIZATION TOWARDS REAL CONVERSATION PROCESSING 2

## INFORMATION RETRIEVAL FROM SPEECH AND SPEECH TRANSLATION

## NATURAL LANGUAGE PROCESSING

## SPOKEN LANGUAGE UNDERSTANDING AND SPOKEN DIALOGUE SYSTEMS

## RESOURCES AND EVALUATION

## ASSISTIVE TECHNOLOGIES

# SPECIAL SESSION ON INTEGRATION OF SPEECH SEPARATION, RECOGNITION AND DIARIZATION TOWARDS REAL CONVERSATION PROCESSING 1

# SPOKEN LANGUAGE UNDERSTANDING AND SOKEN DIALOGUE SYSTEMS

# HUMAN COMPUTER INTERACTION