# 7th Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial '2020)

Held online due to COVID-19

Barcelona, Spain
13 December 2020

# Table of Contents