# Third Workshop on NLP for Similar Languages, Varieties and Dialects 2016 (VarDial 3)

Osaka, Japan
12 December 2016

**Additional copies of this publication are available from:**

# Table of Contents