# 5th Workshop on Vision and Language (VL'16)

Held at ACL 2016

Berlin, Germany
12 August 2016

**proceedings**
.com

# Table of Contents