# 9th Web as Corpus Workshop

# (WaC-9 2014)

**Held at the 14th Conference of the European Chapter
of the Association for Computational Linguistics**

# Gothenburg, Sweden
# 26 April 2014

# Table of Contents