COLING 2012

# 24th International Conference on Computational Linguistics

# Proceedings of COLING 2012: Demonstration Papers

**Program chairs:**
**Martin Kay and Christian Boitet**

**8-15 December 2012**
**Mumbai, India**

**Diamond sponsors**

Tata Consultancy Services
Linguistic Data Consortium for Indian Languages (LDC-IL)

**Gold Sponsors**

Microsoft Research
Beijing Baidu Netcon Science Technology Co. Ltd.

**Silver sponsors**

IBM, India Private Limited
Crimson Interactive Pvt. Ltd.
Yahoo
Easy Transcription & Software Pvt. Ltd.

# Table of Contents

iv