# 18th Nordic Conference of Computational Linguistics (NODALIDA 2011)

Riga, Latvia
11 – 13 May 2011

**Editors:**

**Bolette Sandford Pedersen**       **Inguna Skadina**
**Gunta Nespore**

# Contents

## IV   Student papers                                                   327