# TOPIC DETECTION AND TRACKING

## Event-based Information Organization

*edited by*

**James Allan**
*University of Massachusetts at Amherst*

# Contents