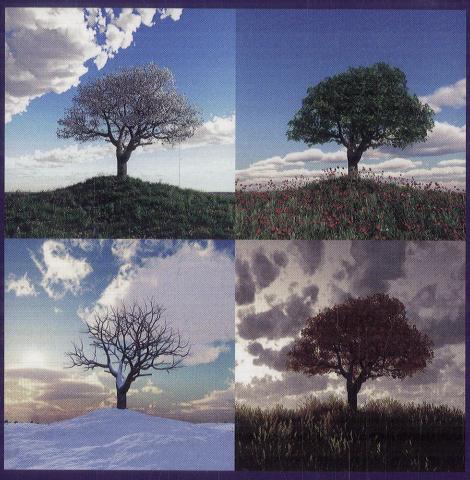# Temporal Data Mining

## Theophano Mitsa

# Table of Contents