# Size Matters: Finding the Most Informative Set of Window Lengths

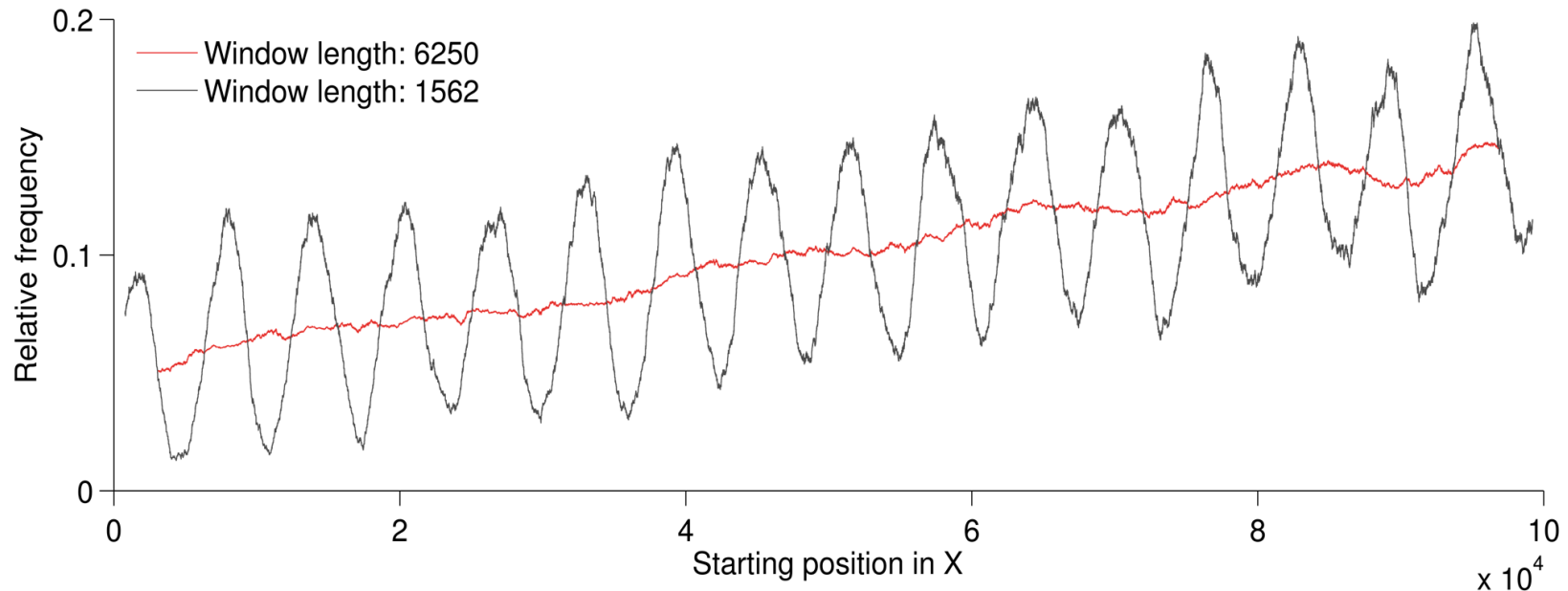Jefrey Lijffijt [1], Panagiotis Papapetrou [1,2], and
Kai Puolamäki [1]

[1] Aalto University, Finland
[2] Birkbeck, University of London, UK

# Summary

- Many sequence analysis algorithms use sliding windows

- **Problem: how to choose the length of the window**

- Novel problem setting and approach

- Solution: use several window lengths that can 'predict' all

- Solution can be computed efficiently

- Method works well

**Aalto University**
**School of Science**

**ECML-PKDD 2012**
**25/09/2012**
**2**

**Finding the most informative set of window lengths**
**Jefrey Lijffijt**

# Example



- Sequences often contain variability at different levels
  - E.g., multiple time scales, daily and weekly rhythm

- Statistic = relative event frequency

- Trends: slow increase and periodic increase/decrease

**Aalto University
School of Science**

**ECML-PKDD 2012
25/09/2012
3**

**Finding the most informative
set of window lengths
Jefrey Lijffijt**

# Related Work

- Solutions are wide-spread (for citations see the paper)

1. User has to choose
2. Optimize towards some objective
   - Fixed-length
   - Variable-length
   - Backing-off
   - Time-fading model (weighting)
   - Many others
3. Use all possible lengths

- None consider optimizing a set of window lengths

Aalto University
School of Science

Finding the most informative
set of window lengths
Jefrey Lijffijt

# Input Data

- *Event sequence X*
  $$X = x_1, ..., x_n, \; x_t \in \sigma$$
  - Fully ordered sequence of events
  - E.g., with four symbols: *ABBCDDAAADCACCABB*

- *Subsequence $X_j(i)$*
  $$X_j(i) = x_i, ..., x_{i+j-1}$$
  - Sequence of length *j* starting at index *i*

- *Statistic $f(X_j(i))$*
  - Measure of interest
    $$f(X_j(i)) = \# q \; occurs \; in \; X_j(i) \, / \, j$$
  - E.g., relative event frequency, or the type/token ratio
    $$f(X_j(i)) = \# types \; in \; X_j(i) \, / \, j$$
  - Can be any algorithm

**Aalto University**
**School of Science**

**ECML-PKDD 2012**
**25/09/2012**
**5**

**Finding the most informative**
**set of window lengths**
**Jefrey Lijffijt**

# Retain Information / Predict All

- Given a set of window sizes $\Omega$

$$\Omega = [\omega_1, \ldots, \omega_m]$$

- Goal: provide as much information wrt $f(X_\omega(i))$ for all $\omega$
  - Using $k$ window lengths

- That is, we want to predict $f(X_\omega(i))$ for all window lengths
  - Based on values $f(X_{\omega 1}(i)), \ldots, f(X_{\omega k}(i))$

**Aalto University**
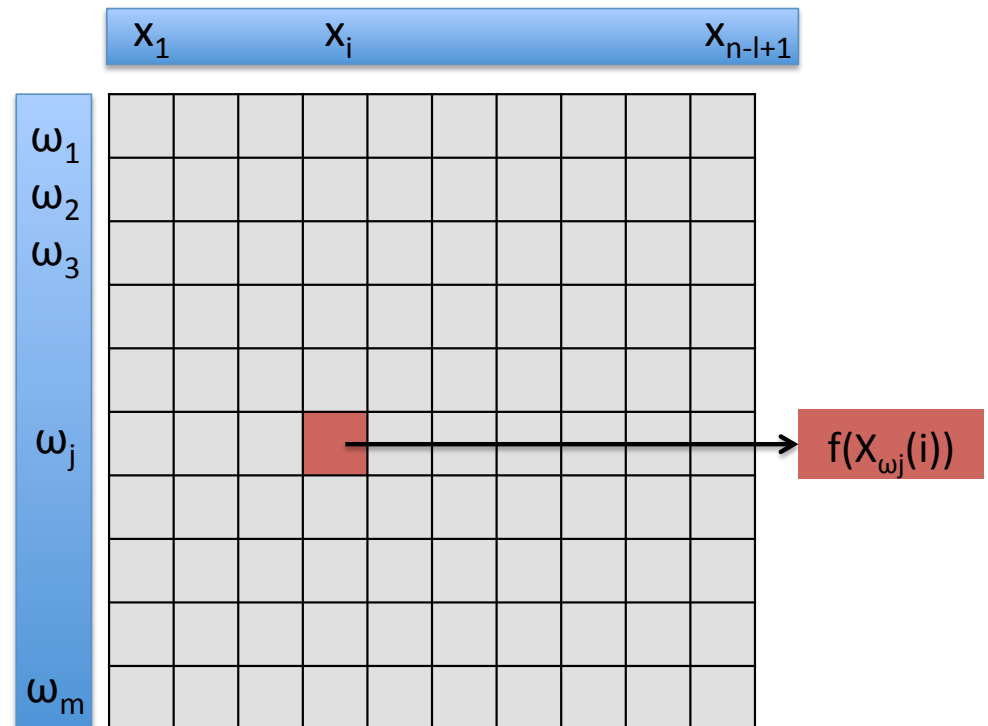**School of Science**

ECML-PKDD 2012
25/09/2012
6

Finding the most informative
set of window lengths
Jefrey Lijffijt

# Window-Trace Matrix

- Given a set of window sizes $\Omega$

$$\Omega = [\omega_1, ..., \omega_m]$$

- The Window-Trace matrix $T$ contains all $f(X_{\omega j}(i))$

$$T_{ji} = f(X_{\omega_j}(i))$$

- We compute only $N$ of the columns



$x_1$     $x_i$     $x_{n-l+1}$

$\omega_1$
$\omega_2$
$\omega_3$
$\omega_j$      $\longrightarrow$   $f(X_{\omega j}(i))$
$\omega_m$

Aalto University
School of Science

ECML-PKDD 2012
25/09/2012
7

Finding the most informative
set of window lengths
Jefrey Lijffijt

# Problem Statement

- *Problem 1 (Maximal variance).* Given a discrete sequence $X$, find a set $R = \{\omega_1, \ldots, \omega_k\}$ of $k$ window lengths that explain most of the variation in $X$, i.e., find a set $R$ that minimizes

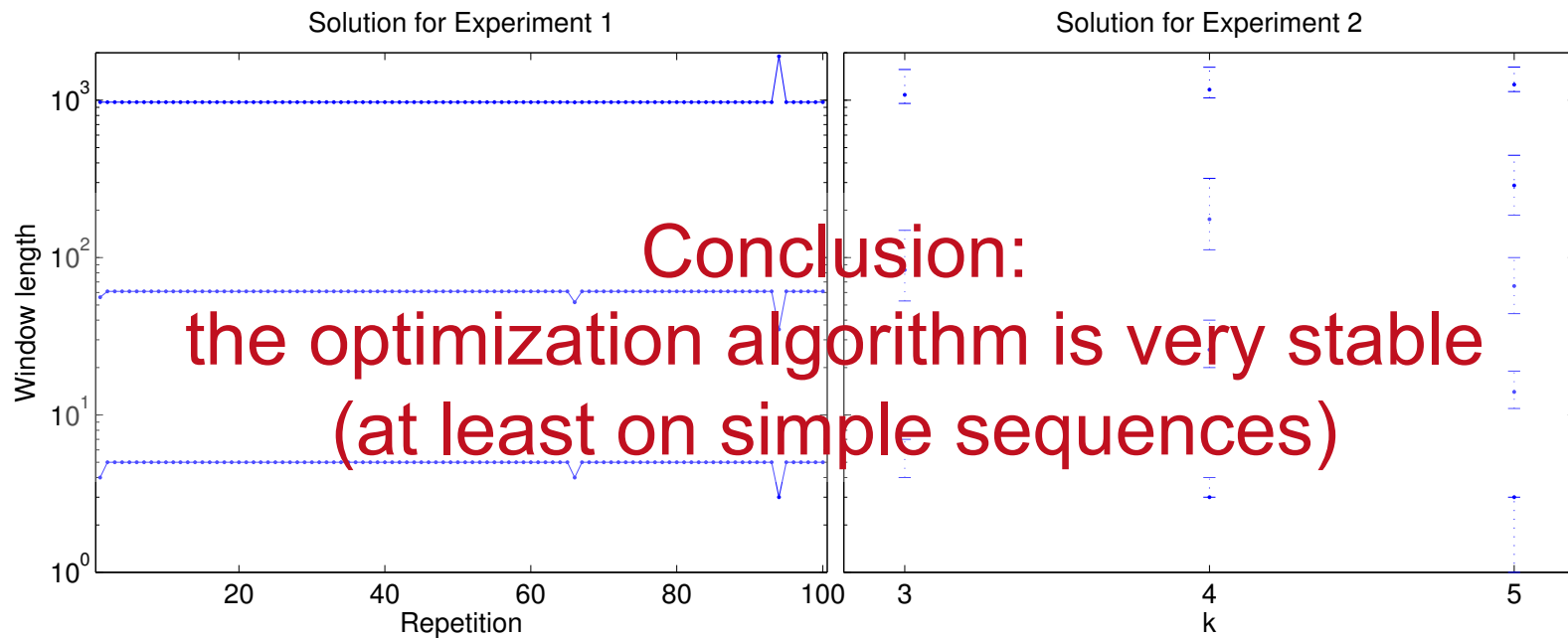$$\sum_{\omega_i \in \Omega} \min_{\omega_j \in R} d(\omega_i, \omega_j)$$

- The distance function that we use is squared error

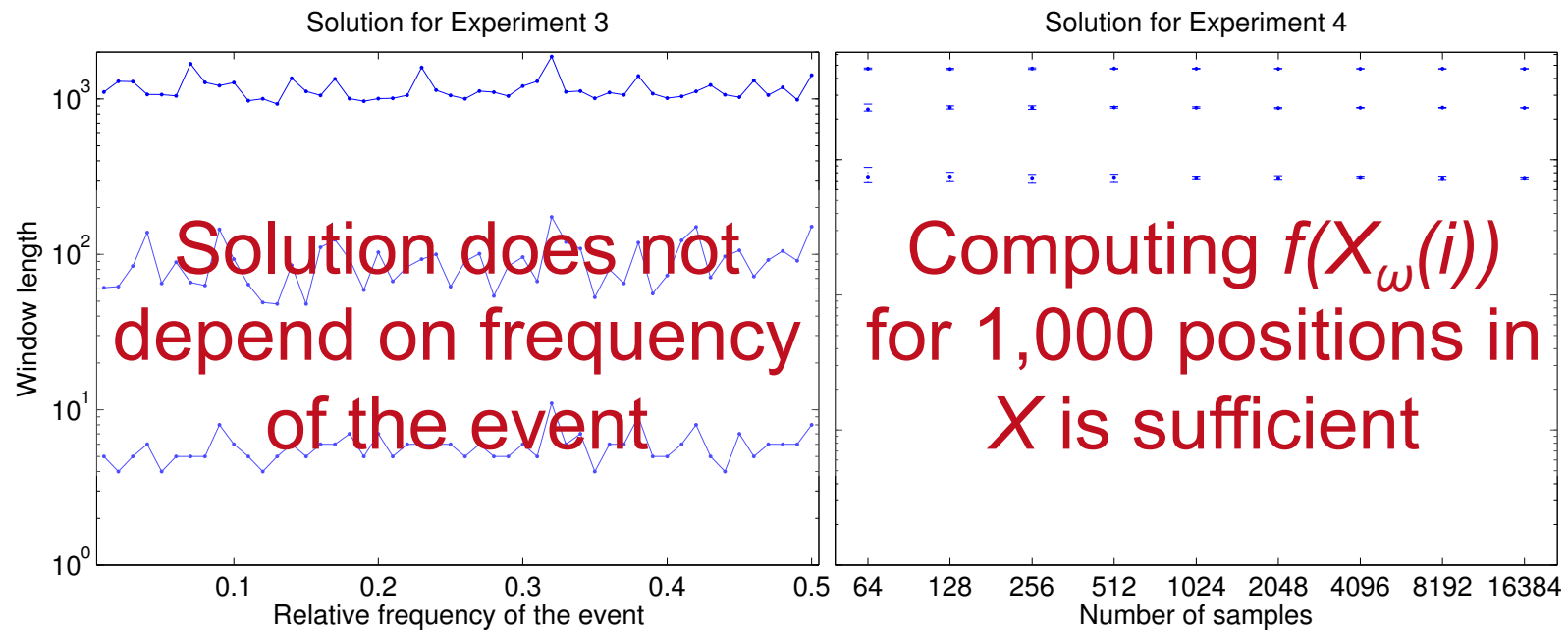$$d(\omega_i, \omega_j) = \sum_{k=1}^{n-l+1} (T_{ik} - T_{jk})^2$$

# Method

- Problem 1 is equivalent to the k-medoids problem

- NP-Hard (Aloise et al. 2009)

- Optimization algorithm:
  – Compute k-means clustering using Lloyd's algorithm
  – Include in $R$ the window lengths closest to each centroid
  – Repeat $r$ times and choose best solution (smallest error)

- Computational complexity: $O(r \cdot i \cdot k \cdot N \cdot m)$

**Aalto University**
**School of Science**

**ECML-PKDD 2012**
**25/09/2012**
**9**

**Finding the most informative**
**set of window lengths**
**Jefrey Lijffijt**

# Experiments on Synthetic Data (1/3)



Solution for Experiment 1         Solution for Experiment 2

Conclusion:
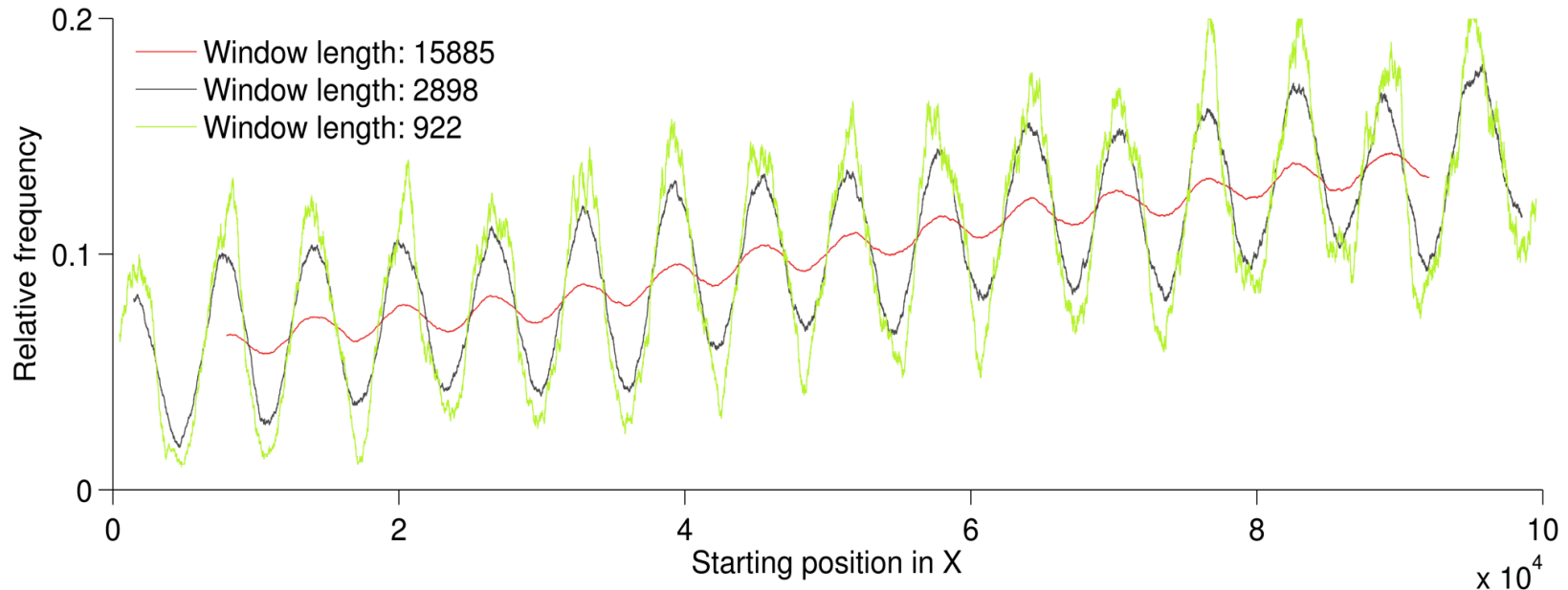the optimization algorithm is very stable
(at least on simple sequences)

- We studied the solution stability on Bernoulli sequences
  - Length = 10,000, p = 0.1
  1. Repeated runs on one sequence
  2. Repeatedly generate sequences

# Experiments on Synthetic Data (2/3)



Solution for Experiment 3

Solution for Experiment 4

Solution does not depend on frequency of the event

Computing $f(X_\omega(i))$ for 1,000 positions in $X$ is sufficient

- We studied the solution stability on Bernoulli sequences
  - Length = 10,000, p = 0.1
  3. Dependency on event frequency
  4. Dependency on number of samples (columns)

# **Experiments on Synthetic Data (3/3)**



- k = 3 solution for sequence shown at introduction

- Both trends clearly visible

- We can accurately estimate all other window lengths

**Aalto University
School of Science**

ECML-PKDD 2012
25/09/2012
12

**Finding the most informative
set of window lengths
Jefrey Lijffijt**

# Burstiness of words in Pride & Prejudice

- A word is *bursty* when it occurs in bursts and lulls
  - Areas with elevated and with lowered frequency

- Non-bursty words: the, and, a, in

- Bursty words: I, you, how, our

- We model burstiness using inter-arrival times
  - IAT: space between two consecutive occurrences of a word
  - Burstiness defined by MLE of Weibull $\beta$

**Aalto University**
**School of Science**

**ECML-PKDD 2012**
**25/09/2012**
**13**

**Finding the most informative set of window lengths**
**Jefrey Lijffijt**

# Burstiness of words in Pride & Prejudice



Bursty words give longer window lengths; they have a larger scale structure

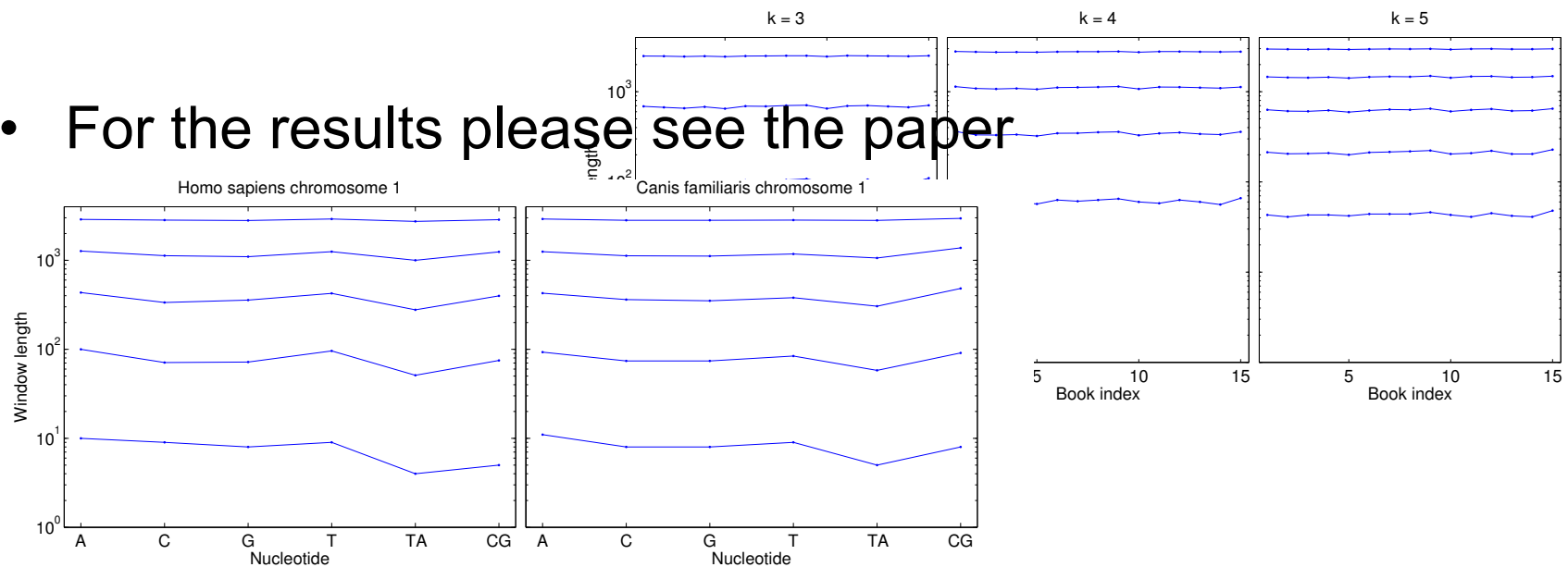| Freq. | Non-bursty | Bursty |
|---|---|---|
| Low | met, rest, right, help (1-4) | write, de, William, read (5-8) |
| Med | time, soon, other, only (9-12) | lady, has, can, may (13-16) |
| High | with, not, that, but (17-20) | you, is, my, his (21-24) |

**Aalto University**
**School of Science**

ECML-PKDD 2012
25/09/2012
14

Finding the most informative
set of window lengths
Jefrey Lijffijt

# Other Experiments

- Type/token ratio throughout several novels

- Frequency of (di-)nucleotides in DNA

- For the results please see the paper

# Summary

- Many sequence analysis algorithms use sliding windows

- **Problem: how to choose the length of the window**

- Novel problem setting and approach

- Solution: use several window lengths that can 'predict' all

- Solution can be computed efficiently

- Method works well

**Aalto University**
**School of Science**

**ECML-PKDD 2012**
**25/09/2012**
**16**

**Finding the most informative**
**set of window lengths**
**Jefrey Lijffijt**