

A framework for space-efficient string kernels

Djamal Belazzougui¹, Fabio Cunial²

(1) Department of Computer Science, University of Helsinki, Finland.(2) Max Planck Institute for Molecular Cell Biology and Genetics, Dresden, Germany.









In this work...

<i>k</i> -mers	Substrings
<i>k</i> -mer complexity	Substring complexity
<i>k</i> -mer kernel	Substring kernel
<i>k</i> -mer profiling	Length-weighted substring kernel
<i>k</i> -mer kernel profile	Character-weighted substring kernel
<i>k</i> -th order empirical entropy	Minimal absent words complexity
KL divergence profile	Minimal absent words distance
	Substring kernel with Markovian corrections
	KL divergence kernel

In this work...

<i>k</i> -mers	Substrings
<i>k</i> -mer complexity	Substring complexity
<i>k</i> -mer kernel	Substring kernel
<i>k</i> -mer profiling	Length-weighted substring kernel
<i>k</i> -mer kernel profile	Character-weighted substring kernel
<i>k</i> -th order empirical entropy	Minimal absent words complexity
KL divergence profile	Minimal absent words distance
	Substring kernel with Markovian corrections

KL divergence kernel

...in small space

Example applications

Alignment-free, whole-genome phylogeny reconstruction.

... for multiple values of *k* at once!

Calibrating the value of *k* for alignment-free distances, from a genome.

Calibrating the value of *k* for de Bruijn graphs, from a set of reads.



string

[1] Simon Gog. Compressed suffix trees: Design, construction, and applications. PhD thesis, University of Ulm, Germany, 2011.



[1] Simon Gog. Compressed suffix trees: Design, construction, and applications. PhD thesis, University of Ulm, Germany, 2011.

Randomized O(n) time from the string



[1] D. Belazzougui. Linear time construction of compressed text indices in compact space. STOC 2014.

[2] D. Belazzougui, G. Navarro, D. Valenzuela. Improved compressed indexes for full-text document retrieval. JDA, 18:3–13, January 2013.

[3] D. Belazzougui, F. Cunial, J. Kärkkäinen, V. Mäkinen. Versatile succinct representations of the bidirectional Burrows-Wheeler transform. ESA 2013.

rangeDistinct queries

 $\texttt{rangeDistinct}(i,j) = \{(c,\texttt{rank}(c,p_c),\texttt{rank}(c,q_c)): c \in \Sigma_{i,j}\}$

unordered

O(out) time, $\sigma \log n$ bits of working space using [1]

let's say O(out d) time in general

 $p_{\mathtt{G}}$

 $q_{\mathtt{G}}$

range(T) =
$$[12..26]$$

 $\Sigma_{12,26} = \{C, G, T, A, \$\}$

$$\texttt{rangeDistinct}(12,26) = \{(\texttt{G},2,4),\dots\}$$

[1] D. Belazzougui, G. Navarro, D. Valenzuela. Improved compressed indexes for full-text document retrieval. JDA, 18:3–13, January 2013.

	BW	Γ_T
	т	\$
	с	ATGTGTATTTGCT\$
	т	ATTTGCT\$
	т	CATGTGTATTTGCT\$
Ū	G	CT\$
	т	CTGTTCATGTGTATTTGCT\$
	т	CTTTCTGTTCATGTGTATTTGCT\$
	т	GCT\$
	т	GTATTTGCT\$
	т	GTGTATTTGCT\$
_	т	GTTCATGTGTATTTGCT\$
_	С	Τ\$
	G	TATTTGCT\$
	т	TCATGTGTATTTGCT\$
	т	TCTGTTCATGTGTATTTGCT\$
	т	TCTTTCTGTTCATGTGTATTTGCT\$
	т	T GCT\$
	G	T GTATTTGCT\$
	A	TGTGTATTTGCT\$
	С	TGTTCATGTGTATTTGCT\$
	G	TCATGTGTATTTGCT\$
	т	T TCTGTTCATGTGTATTTGCT\$
	\$	TTCTTTCTGTTCATGTGTATTTGCT\$
	т	TGCT\$
	с	TTTCTGTTCATGTGTATTTGCT\$
	A	TTTGCT\$

Enumerating nodes in no order 📁



ordered

	E	BWI	T
		т	\$
		с	ATGTGTATTTGCT\$
		т	ATTTGCT\$
		т	CATGTGTATTTGCT\$
		G	CT\$
		т	CTGTTCATGTGTATTTGCT\$
		т	CTTTCTGTTCATGTGTATTTGCT\$
		т	GCT\$
		т	GTATTTGCT\$
		т	GTGTATTTGCT \$
	_	т	GTTCATGTGTATTTGCT\$
\$	L	с	T \$
A	L	G	TATTTGCT\$
		т	T CATGTGTATTTGCT\$
с		т	TCTGTTCATGTGTATTTGCT\$
	L	т	TCTTTCTGTTCATGTGTATTTGCT\$
		т	T GCT\$
G		G	TGTATTTGCT\$
		A	TGTGTATTTGCT\$
	L	с	T GTTCATGTGTATTTGCT\$
		G	TTCATGTGTATTTGCT\$
		т	T TCTGTTCATGTGTATTTGCT\$
т		\$	T TCTTTCTGTTCATGTGTATTTGCT\$
		т	TIGCIŞ
		с	T TTCTGTTCATGTGTATTTGCT\$
	L	A	TTTGCT\$



k-mer complexity



k-mer kernel

$$\begin{aligned} \kappa(\mathbf{T}^1, \mathbf{T}^2) &= N/\sqrt{D^1 D^2} \in [-1..1] \\ N &= \sum_W \mathbf{T}^1[W] \mathbf{T}^2[W] \\ D^i &= \sum_W \mathbf{T}^i[W]^2 \end{aligned}$$



k-mer kernel

$$\kappa(\mathbf{T}^{1}, \mathbf{T}^{2}) = N/\sqrt{D^{1}D^{2}} \in [-1..1]$$
$$N = \sum_{W} \mathbf{T}^{1}[W]\mathbf{T}^{2}[W]$$
$$D^{i} = \sum_{W} \mathbf{T}^{i}[W]^{2}$$

"telescoping"



Rayan Chikhi and Paul Medvedev. Informed and automated k-mer size selection for genome assembly. Bioinformatics, 30(1):31–37, 2014.

Example application



[1] G.E. Sims, S.-R. Jun, G.A. Wu, S.-H. Kim. Alignment-free genome comparison with feature frequency profiles (FFP) and optimal resolutions. PNAS 106(8):2677–2682, 2009.

Similar problem: entropy



Simon Gog. Compressed suffix trees: Design, construction, and applications. PhD thesis, University of Ulm, Germany, 2011.

Similar problem: multiple kernels









 $O(nd + (k_2 - k_1)(f_2 - f_1))$ time, o(n) bits of space

In this work...

k-mers	Substrings
<i>k</i> -mer complexity	Substring complexity
<i>k</i> -mer kernel	Substring kernel
<i>k</i> -mer profiling	Length-weighted substring kernel
<i>k</i> -mer kernel profile	Character-weighted substring kernel
<i>k</i> -th order empirical entropy	Minimal absent words complexity
KL divergence profile	Minimal absent words distance
	Substring kernel with Markovian corrections
	KL divergence kernel







A.J. Smola, S.v.n. Vishwanathan. Fast kernels for string and tree matching. Advances in Neural Information Processing Systems 15, pages 585–592. 2003.





A.J. Smola, S.v.n. Vishwanathan. Fast kernels for string and tree matching. Advances in Neural Information Processing Systems 15, pages 585–592. 2003.



push on another stack the characters that form W

 $\lambda \log \sigma$ additional bits

G. Reinert, D. Chew, F. Sun, M.S. Waterman. Alignment-free sequence comparison (I): statistics and power. Journal of Computational Biology, 2009.

Minimal absent words complexity/Jaccard



S. Chairungsee, M. Crochemore. Using minimal absent words to build phylogeny. Theoretical Computer Science, 450:109–116, 2012.

$$\mathbb{P}(W) = \frac{\mathbb{P}(W[1..k-1]) \cdot \mathbb{P}(W[2..k])}{\mathbb{P}(W[2..k-1])}$$

if T was generated by a Markov process of order $\leq k - 2$

$$\tilde{p}_T(W) = \frac{p_T(W[1..k-1]) \cdot p_T(W[2..k])}{p_T(W[2..k-1])}$$

approximation with empirical probabilities

$$z_T(W) = \frac{p_T(W) - \tilde{p}_T(W)}{\tilde{p}_T(W)} \approx \frac{f_T(W) \cdot f_T(W[2..k-1])}{f_T(W[1..k-1]) \cdot f_T(W[2..k])} - 1 \quad \text{significance score}$$

J. Qi, B. Wang, B.-L. Hao. Whole proteome prokaryote phylogeny without sequence alignment: a *k*-string composition approach. J. Mol. Evol., 2004. A. Apostolico, O. Denas. Fast algorithms for computing sequence distances by exhaustive substring composition. Algorithms for Molecular Biology, 2008.

$$\mathbb{P}(W) = \frac{\mathbb{P}(W[1..k-1]) \cdot \mathbb{P}(W[2..k])}{\mathbb{P}(W[2..k-1])}$$

if T was generated by a Markov process of order $\leq k - 2$

$$\tilde{p}_T(W) = \frac{p_T(W[1..k-1]) \cdot p_T(W[2..k])}{p_T(W[2..k-1])}$$

approximation with empirical probabilities

$$z_T(W) = \frac{p_T(W) - \tilde{p}_T(W)}{\tilde{p}_T(W)} \approx \frac{f_T(W) \cdot f_T(W[2..k-1])}{f_T(W[1..k-1]) \cdot f_T(W[2..k])} - 1 \quad \text{significance score}$$
from repr

J. Qi, B. Wang, B.-L. Hao. Whole proteome prokaryote phylogeny without sequence alignment: a *k*-string composition approach. J. Mol. Evol., 2004. A. Apostolico, O. Denas. Fast algorithms for computing sequence distances by exhaustive substring composition. Algorithms for Molecular Biology, 2008.

$$\mathbb{P}(W) = \frac{\mathbb{P}(W[1..k-1]) \cdot \mathbb{P}(W[2..k])}{\mathbb{P}(W[2..k-1])}$$

if *T* was generated by a Markov process of order $\leq k - 2$

$$\tilde{p}_T(W) = \frac{p_T(W[1..k-1]) \cdot p_T(W[2..k])}{p_T(W[2..k-1])}$$

approximation with empirical probabilities

$$z_T(W) = \frac{p_T(W) - \tilde{p}_T(W)}{\tilde{p}_T(W)} \approx \frac{f_T(W) \cdot f_T(W[2..k-1])}{f_T(W[1..k-1]) \cdot f_T(W[2..k])} - 1 \quad \text{significance score}$$

$$W \text{ present}, W[2..k-1] \text{ not a maximal repeat} \Rightarrow z_T(W) = 0$$
$$W \text{ absent, but not minimal absent} \Rightarrow z_T(W) = 0$$
$$W \text{ minimal absent} \Rightarrow z_T = -1$$

J. Qi, B. Wang, B.-L. Hao. Whole proteome prokaryote phylogeny without sequence alignment: a *k*-string composition approach. J. Mol. Evol., 2004. A. Apostolico, O. Denas. Fast algorithms for computing sequence distances by exhaustive substring composition. Algorithms for Molecular Biology, 2008.



Similar problem

Compute the Kullback-Leibler divergence vector for all k in a range $[k_1..k_2]$:

$$\mathsf{KL}_k = \sum_{W \in [1..\sigma]^k} p_T[W] \cdot \left(\log(p_T[W]) - \log(\tilde{p}_T[W])\right)$$

Only *k*-mers whose infix is a maximal repeat can have nonzero contribution

Example application: find an upper bound on *k* for *k*-mer kernels.

G.E. Sims, S.-R. Jun, G.A Wu, S.-H. Kim. Alignment-free genome comparison with feature frequency profiles (FFP) and optimal resolutions. PNAS, 2009.

In this work...

Substrings k-mers Substring complexity *k*-mer complexity *k*-mer kernel Substring kernel Length-weighted substring kernel *k*-mer profiling Character-weighted substring kernel *k*-mer kernel profile *k*-th order empirical entropy Minimal absent words complexity Minimal absent words distance KL divergence profile Substring kernel with Markovian corrections

KL divergence kernel

all in O(nd + out) time and o(n) bits

In this work...

Substrings k-mers Substring complexity *k*-mer complexity *k*-mer kernel Substring kernel Length-weighted substring kernel *k*-mer profiling Character-weighted substring kernel *k*-mer kernel profile *k*-th order empirical entropy Minimal absent words complexity Minimal absent words distance KL divergence profile Substring kernel with Markovian corrections

KL divergence kernel

all in a single pass over the input

Algorithm design strategies

Independence on the order in which nodes of the suffix tree are enumerated

Telescoping

Batch processing to achieve output-sensitivity

Exploiting the traversal stack

Most "information" is "around" maximal repeats



A framework for space-efficient string kernels

Djamal Belazzougui¹, Fabio Cunial²

(1) Department of Computer Science, University of Helsinki, Finland.(2) Max Planck Institute for Molecular Cell Biology and Genetics, Dresden, Germany.