

Searching and Indexing Genomic Databases via Kernelization

Travis Gagie and Simon J. Puglisi

December 4, 2014

Abstract

The rapid advance of DNA sequencing technologies has yielded databases of thousands of genomes. To search and index these databases effectively, it is important that we take advantage of the similarity between those genomes. Several authors have recently suggested searching or indexing only one reference genome and the parts of the other genomes where they differ. In this paper we survey the twenty-year history of this idea and discuss its relation to kernelization in parameterized complexity.

1 Introduction

The Human Genome Project took thirteen years and three billion dollars to sequence a human genome, but the latest next-generation sequencing methods take only a few days and a few thousand dollars. With these methods, initiatives such as the 1000 Genomes Project and the 100 000 Genomes Project are now feasible. Advances in sequencing have far outstripped advances in computer processors and random-access memory, however, so it is increasingly challenging to make use of the data available. For example, while modern aligners can easily hold in memory the index for approximate pattern matching on a single human genome, they cannot handle thousands of human genomes. Schneeberger et al. [29] proposed that we index the common parts of the genomes only once for them all, but we index the parts near variation sites for each genome. Ferrada [12] suggested indexing the parts of all the genomes near boundaries between phrases in the LZ77 parse of the database. This is more general and may give better compression but requires the LZ77 parse, which is difficult to compute when the database does not fit in memory. Wandelt et al. [34] proposed using a modified parse in which phrases must occur in a reference genome, which is easier to compute. (When papers have appeared in journals we cite those versions, although their chronological order may differ from that of previous versions.) Danek et al. [7] recently showed that with this general approach we can store an index for approximate pattern matching on the database from the 1000 Genomes Project, in the memory of a commodity personal computer.

This has so far not been possible with competing approaches, as surveyed by Vyverman et al. [32].

When we are not given an upper bound on the pattern length, we can use one of the competing indexes that does not require such a bound or we can scan, with an online pattern-matching algorithm, the reference genome and the parts of the other genomes near phrase boundaries. Wandelt and Leser [33] and Rahn et al. [27] proposed the latter idea specifically for approximate pattern matching in genomic databases, but the general approach has a twenty-year history in the field of compressed pattern matching. In this paper we survey that history and relate it to current research: in Section 2 we discuss some relevant data compression schemes and how they have been augmented to support fast random-access reading; in Section 3 we discuss how they have been used to speed up pattern-matching; in Section 4 we discuss how they have been used in compressed indexing. While writing this survey, we realized that scanning or indexing only parts of the database and then mapping the solution for those parts onto a solution for the whole database, is like kernelization in parameterized complexity. (We note that kernels in parameterized complexity bear no relation to operating system kernels nor to kernels in machine learning.) We emphasize this perspective because we feel that computing a pattern-matching kernel is an interesting problem in itself, regardless of how we process it later, and deserving of further study. Of course, the nature and even the existence of the kernel depend on the problem we are trying to solve.

2 Compression with Random-Access Reading

In general, the best compression of highly repetitive datasets is achieved with the LZ77 algorithm by Ziv and Lempel [35]. Suppose $S[1..n]$ is a string with $S[n] = \$$, which is an end-of-file symbol that does not occur elsewhere in S . LZ77 works by parsing S into phrases such that, for each phrase $S[i..j]$, $S[i..j-1]$ occurs in $S[1..j-2]$ but $S[i..j]$ does not occur in $S[1..j-1]$; that phrase is stored as a triple consisting of a pointer to $S[i..j]$'s first occurrence in S (which is called the phrase's source), $j-i$, and $S[j]$. The LZ77 encoding of S takes $\mathcal{O}(z \log n)$ bits, where z is the number of phrases in the parse. For example, in the following verses vertical lines indicate phrase boundaries:

```
9|9-|b|o|t|t|e|s|-o|f|-b|e|e|r|-o|n|-t|h|e|-w|a|l|l|-9|9-bottles-of-beer-
I|f-o|n|e-o|f-t|h|o|s|e|-b|o|t|t|l|e|-s|h|o|u|l|d|-h|a|p|p|e|n|-t|o|-f|a|l|l-
98|-bottles-of-beer-on-the-wall-

98|-bottles-of-beer-on-the-wall-98-bottles-of-beer-
I|f-o|n|e-o|f-t|h|o|s|e|-b|o|t|t|l|e|-s|h|o|u|l|d|-h|a|p|p|e|n|-t|o|-f|a|l|l-
97|-bottles-of-beer-on-the-wall. . .
```

(We have displayed the verses with linebreaks to increase readability, but we have not considered them while computing the parse.) Although these verses may be annoyingly similar by the standards of natural language, they are far

less similar than human genomes. Indeed, most repetitive biological datasets are much too similar (as well as much too large) for us to use them as informative examples.

One drawback of LZ77 compression is that reading a character in a compressed string can be very slow. Rytter [28] and Charikar et al. [4] showed how we can turn that parse into a balanced straight-line program (SLP) for S with $\mathcal{O}(z \log n)$ rules. An SLP for S is a context-free grammar in Chomsky normal form that generates S and only S ; it is balanced if the height of each subtree in the parse tree is logarithmic in that subtree's size. It follows from Rytter's and Charikar et al.'s results that we can store S in $\mathcal{O}(z \log^2 n)$ bits and support random-access reading of any substring of S with length ℓ in $\mathcal{O}(\log n + \ell)$ time. Verbin and Yiu [31] showed that this is nearly optimal in the worst case. Bille et al. [3] showed how, given even an unbalanced SLP for S with r rules, we can store S in $\mathcal{O}(r \log n)$ bits and support random-access reading in $\mathcal{O}(\log n + \ell)$ time. Rytter's, Charikar et al.'s and Bille et al.'s constructions are not practical, but there are practical grammar-based compressors, such as those by Larsson and Moffat [24] and Maruyama and Tabei [26]. As far as we know, block graphs by Gagie et al. [14, 16] are the most practical grammar-like representations for random-access reading. The LZ78 algorithm by Ziv and Lempel [36] does not compress repetitive datasets as well as LZ77, but the LZ78 encoding of S can easily be augmented to support random-access reading in $\mathcal{O}(\log \log n + \ell)$ time. LZ78 also works by parsing S into phrases but then each phrase must extend a previous phrase plus one character. Because of this property, the LZ78 encoding of S has $\Omega(\sqrt{n})$ phrases, even when $S = a^n$.

In the example above, the first verse contains many phrase boundaries but the second verse contains only three. Kuruppu et al. [22] proposed that, given a set of similar strings (or one string that can easily be divided into similar substrings), we store the first string in plain text as a reference and compress the others with a version of LZ77 that restricts phrases' sources to occur in the reference. They called this scheme Relative Lempel-Ziv (RLZ) and showed it compresses genomic databases very well in practice (although it too uses $\Omega(\sqrt{n})$ phrases, even when $S = a^n$) and there are several implementations of this approach, such as those by Deorowicz and Grabowski [8], Kuruppu et al. [21] and Ferrada et al. [11]. Even when there is no obvious reference, Kuruppu et al. [23] showed we can often build one by sampling the dataset: intuitively, if a substring is common then it is likely to appear in our sample, and if it is not then we lose little by not compressing it well; this can be formalized using results about SLPs.

3 Searching

Farach and Thorup [10] observed that the first occurrence of any pattern $P[1..m]$ in S must cross or end at a phrase boundary in the LZ77 parse. Kärkkäinen and Ukkonen [18] showed how, if we already know the locations of P 's occurrences in S that cross or end at phrase boundaries, then we can deduce the locations

of all its other occurrences from the structure of the parse. By the same arguments, LZ78 also has these properties and Karpinski, Rytter and Shinohara [19] simultaneously proved similar results for SLPs. Bille et al. [2] observed that any substring of S within edit distance k of P (i.e., any of P 's approximate matches) has length at most $m + k$, and any such substring that does not cross or end at an LZ78 phrase boundary must be an exact copy of an earlier one that does. They gave an algorithm for approximate pattern matching in LZ78 strings that works by extracting the $m + k$ and $m + k - 1$ characters before and after each LZ78 phrase boundary, respectively, using a technique similar to those discussed in Section 2; scanning the resulting substrings with any online algorithm for approximate pattern matching in uncompressed strings; and then deducing the locations of the other approximate matches from the structure of the parse.

Bille et al. [3] extended this approach to show how we can find all P 's approximate matches in S from an SLP for S . Recently, Gagie et al. [15] extended it further to show how we can preprocess the LZ77 parse of S in $\mathcal{O}(z \log n)$ time such that later, given P and k , we can find all P 's occ approximate matches in $\mathcal{O}(z \min(mk, m + k^4) + occ)$ time. Their algorithm works by extracting the $m + k$ and $m + k - 1$ characters before and after each LZ77 phrase boundary, respectively, and then continuing as with the algorithm by Bille et al. [2]. The set of substrings we extract is like a kernel in parameterized complexity: the total length of the substrings can be much smaller than n , but a solution on them can quickly be mapped to a solution on all of S . For our example from Section 2 with $m = 4$ and $k = 1$, the kernel is

```
99-bottles-of-beer-on-the-wall-99-bo
eer-If-one-of-those-bottles-should-happen-to-fall-98-bot
ll-98-bot
eer-If-on
ll-97-bot...
```

If we want a kernel consisting of only a single string, we can concatenate the substrings with $k + 1$ copies of \$ between each consecutive pair. Notice that if we are careful, we can avoid scanning the fourth substring “eer-If-on”, since it occurs in the second substring.

We do not wish to leave the impression that kernelization is the only approach used in compressed pattern matching, nor even that the papers mentioned above are the only ones that use it. We have focused on those papers because we feel they are the most relevant to the practical bioinformatics papers by Wandelt and Leser [33] and Rahn et al. [27] mentioned in Section 1. Those authors were apparently unaware of the field of compressed pattern matching and re-invented kernelization specifically for approximate pattern matching in genomic databases, with kernels based on RLZ instead of LZ77, LZ78 or SLPs. This may be because the earlier researchers using kernelization for pattern matching did not publicize their ideas in interdisciplinary forums or implement their ideas in tools usable by other scientists.

4 Indexing

Kärkkäinen and Ukkonen [18] gave the first LZ-based index, which supported exact pattern matching and stored S separately and uncompressed. They used Patricia trees and range reporting to find a set of candidate matches crossing or ending at LZ77 phrase boundaries; verified them by checking S ; and then used more range reporting to find the other matches. We can obtain various time-space tradeoffs by compressing S and use the methods discussed in Section 1 to extract the characters needed to verify candidate matches. Claude and Navarro [5] modified Kärkkäinen and Ukkonen’s index to use a grammar-compressed encoding of S , and Kreft and Navarro [20] modified it to use the encoding of S produced by a version of LZ77 they called LZ-End, which supports fast random-access reads starting at phrase boundaries. Arroyuelo et al. [1] and Do et al. [9] gave indexes based on LZ78 and RLZ, respectively, and Maruyama et al. [25] and Takabatake et al. [30] gave indexes based on the edit-sensitive parsing by Cormode and Muthukrishnan [6]. Gagie et al. [13] recently gave a version of Kärkkäinen and Ukkonen’s index that uses a total of $\mathcal{O}(z \log^2 n)$ bits and returns the locations of all P ’s *occ* occurrences in S in $\mathcal{O}(m \log m + \text{occ} \log \log n)$ time. These indexes require no assumptions about the pattern.

Kärkkäinen and Sutinen [17] gave an index based on a version of LZ77 that allows phrases to overlap by $q - 1$ characters, where q is a parameter. If P has length exactly q , then their index returns the locations of all P ’s occurrences in S in optimal $\mathcal{O}(m + \text{occ})$ time. If we are given an upper bound M on the pattern length at construction time, then even with Kärkkäinen and Ukkonen’s original version, we need keep only a kernel of the text and can use $\mathcal{O}(z \log n + zM \log \sigma)$ bits in total, where σ is the size of the alphabet. We suspect this escaped investigation for so long because it seemed too obvious and inelegant to be theoretically interesting, and the need to index massive, highly repetitive datasets in practice has become pressing only since the development of next-generation sequencing methods.

The use of kernelization for indexing was eventually investigated by Schneeberger et al. [29], although they did not present kernelization as a separate process because their work was application-driven. As noted in Section 1, they proposed that, given a database of genomes from the same species, we index the common parts of the genomes only once for them all, but we index the parts near variation sites for each genome. Wandelt et al. [34] and Danek et al. [7] gave similar results, essentially using a kernel based on the RLZ parse. Like Schneeberger et al., these authors indexed the kernels using specific methods based on q -grams or seeds. Danek et al.’s index for the database for the 1000 Genomes Project is the first one to fit in the memory of a commodity personal computer. Ferrada et al. [12] emphasized kernelization (albeit not under that name) in terms of the LZ77 parse, which is more general and may give better compression, and pointed out that we can use any index for approximate pattern matching to store the kernel. One point they did not comment on, and which we hope to have clarified in this paper, is that we can consider kernels based

on LZ77, LZ78, RLZ, other compression schemes, or possibly other algorithms entirely. These kernels may be easier to compute when the database does not fit in memory, or have other useful properties that make them preferable in some situations. One interesting problem is how we can best maintain a dynamic kernel for an expanding database. This could allow us to align reads against a genomic database and then add the newly-assembled genome, which could be useful when dealing with mutating cancer genomes or changing strains of a disease during an outbreak.

Acknowledgement

Many thanks to Fabio Cunial, Paweł Gawrychowski, Simon Grabowski, Juha Kärkkäinen, Veli Mäkinen, Gonzalo Navarro, Esa Pitkänen, Yasuo Tabei and Niko Välimäki, for helpful discussions.

References

- [1] D. Arroyuelo, G. Navarro, and K. Sadakane. Stronger Lempel-Ziv based compressed text indexing. *Algorithmica*, 62:54–101, 2012.
- [2] P. Bille, R. Fagerberg, and I. L. Gørtz. Improved approximate string matching and regular expression matching on Ziv-Lempel compressed texts. *ACM Transactions on Algorithms*, 6, 2009.
- [3] P. Bille, G. M. Landau, R. Raman, K. Sadakane, S. R. Satti, and O. Weimann. Random access to grammar-compressed strings. In *Proceedings of the 22nd Symposium on Discrete Algorithms (SODA)*, pages 373–389, 2011.
- [4] M. Charikar, E. Lehman, D. Liu, R. Panigrahy, M. Prabhakaran, A. Sahai, and A. Shelat. The smallest grammar problem. *IEEE Transactions on Information Theory*, 51:2554–2576, 2005.
- [5] F. Claude and G. Navarro. Improved grammar-based compressed indexes. In *Proceedings of the 19th Symposium on String Processing and Information Retrieval (SPIRE)*, pages 180–192, 2012.
- [6] G. Cormode and S. Muthukrishnan. The string edit distance matching problem with moves. *ACM Transactions on Algorithms*, 3, 2007.
- [7] D. A. Danek, S. Deorowicz, and S. Grabowski. Indexes of large genome collections on a PC. *PLoS ONE*, 9, 2014. e109384.
- [8] S. Deorowicz and S. Grabowski. Robust relative compression of genomes with random access. *Bioinformatics*, 27:2979–2986, 2011.

- [9] H. H. Do, J. Jansson, K. Sadakane, and W.-K. Sung. Fast relative Lempel-Ziv self-index for similar sequences. *Theoretical Computer Science*, 532:14–30, 2014.
- [10] M. Farach and M. Thorup. String matching in Lempel-Ziv compressed strings. *Algorithmica*, 20:388–404, 1998.
- [11] H. Ferrada, T. Gagie, S. Gog, and S. J. Puglisi. Relative Lempel-Ziv with constant-time random access. In *Proceedings of the 21st Symposium on String Processing and Information Retrieval (SPIRE)*, pages 13–17, 2014.
- [12] H. Ferrada, T. Gagie, T. Hirvola, and S. J. Puglisi. Hybrid indexes for repetitive datasets. *Philosophical Transactions of the Royal Society A*, 327, 2014. Article no. 2016.
- [13] T. Gagie, P. Gawrychowski, J. Kärkkäinen, Y. Nekrich, and S. J. Puglisi. LZ77-based self-indexing with faster pattern matching. In *Proceedings of the 11th Latin American Symposium on Theoretical Informatics (LATIN)*, pages 731–742, 2014.
- [14] T. Gagie, P. Gawrychowski, and S. J. Puglisi. Faster approximate pattern matching in compressed repetitive texts. In *Proceedings of the 22nd International Symposium on Algorithms and Computation (ISAAC)*, pages 653–662, 2011.
- [15] T. Gagie, P. Gawrychowski, and S. J. Puglisi. Faster approximate pattern matching in compressed repetitive texts. *Journal of Discrete Algorithms*, 2014. in press.
- [16] T. Gagie, C. Hoobin, and S. J. Puglisi. Block graphs in practice. In *Proceedings of the 2nd International Conference on Algorithms for Big Data (ICABD)*, pages 30–36, 2014.
- [17] J. Kärkkäinen and E. Sutinen. Lempel-Ziv index for q -grams. *Algorithmica*, 21:137–154, 1998.
- [18] J. Kärkkäinen and E. Ukkonen. Lempel-Ziv parsing and sublinear-size index structures for string matching. In *Proceedings of the 3rd South American Workshop on String Processing (WSP)*, pages 141–155, 1996.
- [19] M. Karpinski, W. Rytter, and A. Shinohara. An efficient pattern-matching algorithm for strings with short descriptions. *Nordic Journal of Computing*, 4:172–186, 1997.
- [20] S. Kreft and G. Navarro. On compressing and indexing repetitive sequences. *Theoretical Computer Science*, 483:115–133, 2013.
- [21] S. Kuruppu, B. Beresford-Smith, T. C. Conway, and J. Zobel. Iterative dictionary construction for compression of large DNA data sets. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 9:137–149, 2012.

- [22] S. Kuruppu, S. J. Puglisi, and J. Zobel. Relative Lempel-Ziv compression of genomes for large-scale storage and retrieval. In *Proceedings of the 17th Symposium on String Processing and Information Retrieval (SPIRE)*, pages 201–206, 2010.
- [23] S. Kuruppu, S. J. Puglisi, and J. Zobel. Reference sequence construction for relative compression of genomes. In *Proceedings of the 18th Symposium on String Processing and Information Retrieval (SPIRE)*, pages 420–425, 2011.
- [24] N. J. Larsson and A. Moffat. Offline dictionary-based compression. In *Proceedings of the Data Compression Conference (DCC)*, pages 296–305, 1999.
- [25] S. Maruyama, M. Nakahara, N. Kishiue, and H. Sakamoto. ESP-index: A compressed index based on edit-sensitive parsing. *Journal of Discrete Algorithms*, 18:100–112, 2013.
- [26] S. Maruyama and Y. Tabei. Fully online grammar compression in constant space. In *Proceedings of the Data Compression Conference (DCC)*, pages 173–182, 2014.
- [27] R. Rahn, D. Weese, and K. Reinert. Journalized string tree — a scalable data structure for analyzing thousands of similar genomes on your laptop. *Bioinformatics*, 2014. in press.
- [28] W. Rytter. Application of Lempel-Ziv factorization to the approximation of grammar-based compression. *Theoretical Computer Science*, 302(1–3):211–222, 2003.
- [29] K. Schneeberger, J. Hagmann, S. Ossowski, N. Warthmann, S. Gesing, O. Kohlbacher, D., and Weigel. Simultaneous alignment of short reads against multiple genomes. *Genome Biology*, 10, 2009. Article no. R98.
- [30] Y. Takabatake, Y. Tabei, and H. Sakamoto. Improved ESP-index: A practical self-index for highly repetitive texts. In *Proceedings of the 13th Symposium on Experimental Algorithms (SEA)*, pages 338–350, 2014.
- [31] E. Verbin and W. Yu. Data structure lower bounds on random access to grammar-compressed strings. In *Proceedings of the 24th Symposium on Combinatorial Pattern Matching (CPM)*, pages 247–258, 2013.
- [32] M. Vyverman, B. De Baets, V. Fack, and P. Dawyndt. Prospects and limitations of full-text index structures in genome analysis. *Nucleic Acids Research*, 40:6993–7015, 2012.
- [33] S. Wandelt and U. Leser. String searching in referentially compressed genomes. In *Proceedings of the Conference on Knowledge Discovery and Information Retrieval (KDIR)*, pages 95–102, 2012.

- [34] S. Wandelt, J. Starlinger, M. Bux, and U. Leser. RCSI: scalable similarity search in thousand(s) of genomes. *Proceedings of the VLDB Endowment*, 6:1534–1545, 2013.
- [35] J. Ziv and A. Lempel. A universal algorithm for sequential data compression. *IEEE Transactions on Information Theory*, 23(3):337–343, 1977.
- [36] J. Ziv and A. Lempel. Compression of individual sequences via variable-rate coding. *IEEE Transactions on Information Theory*, 24(5):530–536, 1978.