

Scalable Sequence Similarity Search and Join in Main Memory on Multi-cores

Astrid Rheinländer and Ulf Leser

Humboldt-Universität zu Berlin, Department of Computer Science
Berlin, Germany

Abstract. Similarity-based queries play an important role in many large scale applications. In bioinformatics, DNA sequencing produces huge collections of strings, that need to be compared and merged. We present PeARL, a data structure and algorithms for similarity-based queries on many-core servers. PeARL indexes large string collections in compressed tries which are entirely held in main memory. Parallelization of searches and joins is performed using MapReduce as the underlying execution paradigm. We show that our data structure is capable of performing many real-world applications in sequence comparisons in main memory. Our evaluation reveals that PeARL reaches a significant performance gain compared to single-threaded solutions. However, the evaluation also shows that scalability should be further improved, e.g., by reducing sequential parts of the algorithms.

1 Introduction

Similarity-based searches and joins are important for many applications such as document clustering or plagiarism detection [7,16]. In bioinformatics, similarity-based queries are used for sequence read alignment or for finding homologous sequences between different species. In recent years, much effort has been spent on developing tools to speed up similarity-based queries on sequences. Many prominent tools use sophisticated index structures and filter techniques that enable significant runtime improvements [2,8,9].

A challenge arises from the immense growth of sequence databases in the past few years. For example, the number of sequences stored in EMBL grows exponentially every year and sums up to more than 300 billion nucleotides as of May 2011. One strategy to deal with this huge amount of data is to divide it into smaller parts and perform analyses partition-wise in parallel. For this scenario, Google developed the programming paradigm MapReduce to enable a massively-parallel processing of huge data sets in large distributed systems of commodity hardware [4]. However, the main bottleneck of distributed MapReduce is network bandwidth and disk I/O. Therefore, another option is to design data structures and algorithms that adapt the MapReduce paradigm for many-core servers [11]. We argue that modern many-core servers, combined with the constantly falling prices for main memory, are perfectly suited to perform many

real-world applications in sequence analysis. Such settings are much easier to maintain and do not suffer from bandwidth problems.

In this paper, we challenge the current opinion that problems in sequence analysis already have grown so big that distributed systems are the only solution. We present PEARL, a main-memory data structure and parallel algorithms for similarity-based search and join operations on sequence data. In particular, our data structure uses compressed tries. In tries, the complexity for exact searches only depends on string lengths and not on the number of stored strings [14]. This allows an efficient execution of exact searches even in large tries. In order to retain these advantages for similarity-based queries, we store additional information at each node that enable early pruning of whole subtrees. Previously, we demonstrated that these strategies effectively speed up similarity-based queries in PETER [12], a disk-based index structure and predecessor of PeARL.

A crucial aspect in designing data structures for similarity based queries that interact with MapReduce is to support proper data partitioning. Specifically, we show how tries on top of large string collections can be compressed and partitioned for enabling in-memory MapReduce based search and join operations. To our knowledge, this is the first work that parallelizes similarity-based string searches and joins in tries. Our evaluation reveals that PeARL’s similarity-based algorithms scale well.

The rest of this paper is organized as follows: Section 2 introduces basic concepts needed for the design of our data structure and algorithms. We describe design principles of PeARL and algorithms for similarity search and join, as well as our parallelization strategy in Sect. 3. We evaluate our tool in Sect. 4 and discuss related work in Sect. 5. Finally, we conclude our paper with an outlook to future work.

2 Preliminaries

Let Σ^* be the set of all strings of any finite length over an alphabet Σ . The length of a string $s \in \Sigma^*$ is denoted by $|s|$. A substring $s[i \dots j]$ of s starts at position i and ends at position j , ($1 \leq i \leq j \leq |s|$). Any substring of length $q \in \mathbb{N}$ is called q -gram. Conceptually, we will ground our algorithms on operators for similarity search and similarity join, which are defined as follows:

Let s be a string, R a bag of strings, d a distance function, and k a threshold. The similarity-based search operator is defined as $\text{sim}_{\text{search}}(s, R, k) = \{r | d(r, s) \leq k, r \in R\}$. Similarly, for two bags of strings R, S , the similarity-based join operator is defined as $\text{sim}_{\text{join}}(R, S, k) = \{(r, s) | d(r, s) \leq k, r \in R, s \in S\}$.

In PeARL, we support Hamming and edit distance as similarity measure. We focus on edit distance based operations in this paper, but see [12] for the key ideas on queries using Hamming distance. In general, the edit distance of r and s is computed in $O(|r| * |s|)$ using dynamic programming. As we are mostly interested in finding highly similar strings within a previously defined distance threshold k , we use the k -banded alignment algorithm [5] with time complexity $O(k * \max\{|r|, |s|\})$ instead.

Our parallelization strategy is inspired by the well-known programming model MapReduce, a two-step approach that consists of a map and a reduce phase [4]. Essentially, data is stored in **<key, value>**-pairs and partitioned into several subsets. In the map step, a user-defined function is applied to each input item $\langle k_i, v_i \rangle$ and an intermediate list of $\langle k_j, v_j \rangle$ pairs is emitted. All intermediate items generated by map are grouped on the basis of the keys and finally, the user-defined reduce function is applied to each group in order to assemble the final result set.

3 Data Structure and Algorithms

In this section, we introduce our data structure PeARL together with algorithms for executing similarity string searches and joins in parallel.

Conceptually, a PeARL index (see Fig. 1) is based on radix trees [10] and defined as follows: Let R be a bag of strings. A PeARL index P_R for R consists of a set of rooted, compressed tries T_R , a sequence string seq , and a **<key, value>** data structure *StringIDMap* and meets the following conditions:

- (*Identification of strings*) The string seq is a concatenation of all $r \in R$. We assign a unique ID to each r , assembled from a serial number, the length of r , and the start position of r in seq .
- (*Node types*) We distinguish between infix nodes and string nodes. An infix node is a node that represents some substring r_l of r , $|r_l| \geq 1$. Each node u represents a sequence of characters of length $l \geq 1$. The labels of any two children v, w of u start with different characters. Every r maps to exactly one node $x \in T_R$ such that the concatenation of all labels from T_R 's root to x exactly is r . Such nodes x are called string node. We store a pair that consists of the node ID of x and the UID of r in the *StringIDMap*. If R contains multiple copies of r , all corresponding UIDs are assigned to x .
- (*Storing infixes*) Node labels are not stored directly in node u , but retrieved via lookups in seq . Thus, u stores length and start position of the represented infix in seq .
- (*Additional information*) Each node u stores additional attributes, namely the minimum (*min*) and maximum (*max*) lengths of strings stored in the subtree starting at u , a character frequency vector fv and a bit-string qGr . The character frequency vector $fv(u)$ consists of $|\Sigma|$ components and counts the number of occurrences of $c_i \in \Sigma$ in the prefix represented by u in component i . Similarly, a bit in qGr at position i represents the i th string of all strings over Σ of length q in lexicographical order. Bit i is set to 1, if the prefix represented by node u contains the corresponding q -gram.
- (*Trie partitioning*) For very large string collections, we expect the upper levels of a trie to be completely filled. Therefore, we partition a single PeARL trie into multiple tries on the basis of shared prefixes. Each partition is identified by the prefix which was used for partitioning (see Fig. 1). The prefix length used for partitioning is user-defined.

Figure 1 displays a PeARL index for strings over $\Sigma = \{A, C, G, T\}$. Grey nodes are string nodes, white nodes are infix nodes. Edge labels are not stored in the index itself, but are displayed for better comprehensibility only. Displayed q -gram sets indicate which bits in qGr are set.

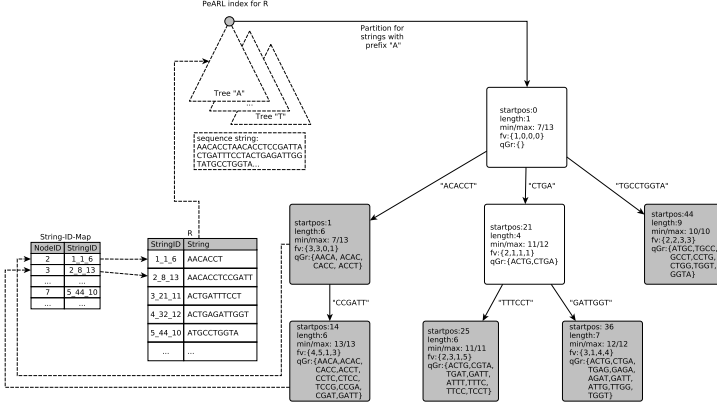


Fig. 1. PeARL index structure

3.1 Algorithms

Building the PeARL index for a set of strings R works as follows: In a first step, R is sorted lexicographically, UIDs are assembled, and R is split into multiple partitions based on shared prefixes. For each partition $R_i \subseteq R$, we start with an empty trie T_{R_i} and iteratively insert each string contained in R_i using preorder DFS traversal. After all strings from R_i have been inserted, we iterate once over the whole trie and update the information min/max , fv and qGr .

Similar to indexing, our algorithms for similarity-based searches and joins are also grounded on preorder DFS traversal of all trie partitions. Each algorithm is equipped with filtering strategies. These filters, namely prefix and edit distance pruning [14], character frequency pruning [1], and q -gram filtering [6], have been introduced in slightly different contexts before. Their concrete usage and efficiency for trie-based search and join queries is shown in [12]. Therefore, we only briefly summarize our search and join strategies in the following and concentrate on our novel parallelization scheme later.

Similarity search starts with a given search string q and traverses each trie partition in a PeARL index starting at root. Whenever a new child of the current node is reached, we first check whether we can prune this node (see [12] for details on filtering). If all filters have been passed successfully, we compute the edit distance between the query and the prefix of the node. If the distance exceeds a threshold k , we start a backtracking routine and traverse the remaining, not yet examined paths in the trie. Otherwise we descend forward to the leaves. When a string node x is reached and $d(q, x) \leq k$ holds, we report a match.

Similarity join for two sets R, S takes two PeARL indices P_R, P_S as input. Each trie partition T_{R_i} is joined with each partition T_{S_j} . Recall that both tries are partitioned by prefixes. We first check the partition prefixes on edit distance and it might happen that k is already exceeded. In this case, we skip the corresponding trie pair. Otherwise, we compute the similarity-based intersection of both partitions. As for search, we start at the root nodes and traverse both tries concurrently. When unseen nodes are reached, we check all filters and prune, if possible. Whenever two string nodes $x \in T_{R_i}, y \in T_{S_j}$ are reached, and given that $d(x, y) \leq k$ holds, we report a match.

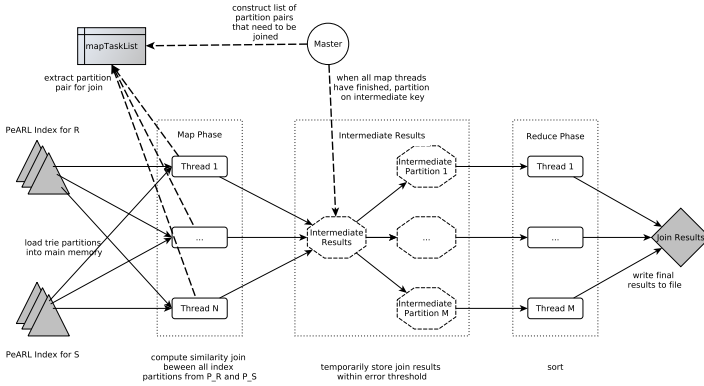


Fig. 2. MapReduce workflow of similarity joins in PeARL

3.2 Parallelization with In-memory MapReduce

We use MapReduce to parallelize PeARL for an execution on multi-core servers. However, a usage in distributed scenarios is conceptually also possible as PeARL trie partitions could as well be spread over nodes in a distributed file system.

Recall that a user-defined function is applied to each input item $\langle k_i, v_i \rangle$ in the map phase. Depending on the specific task, we use the map phase to either execute the similarity join of any two PeARL partitions or, to search a certain string in each partition of a PeARL index. Reduce phases are typically used to compute aggregates of intermediate results. Figure 2 shows the workflow for parallelizing similarity joins in PeARL with MapReduce. A master routine takes two PeARL indices P_R, P_S as input, together with an error threshold k , and a number of available threads t . As string collections stored with PeARL are already partitioned into multiple tries, we get a natural data partitioning for the map phase. The master generates a set of map tasks (stored in a FIFO data structure `mapTaskList`), such that each trie partition $T_{R_i} \in P_R$ is joined with each trie partition $T_{S_j} \in P_S$ and starts the map phase.

Each map thread has access to `mapTaskList` and extracts one task (T_{R_i}, T_{S_j}) out of this list. After some initialization steps, map calls the join routine, that executes the similarity join of $T_{R_i} \bowtie_k T_{S_j}$ and returns the set of all similar string pairs contained in (T_{R_i}, T_{S_j}) within the given distance k . These items are inserted into an intermediate data structure. For each similar string pair (r, s) , an intermediate key is set to the UID of r . When one map iteration has finished and as long as `mapTaskList` is not empty, the map thread extracts the next (T_{R_i}, T_{S_j}) pair out of this list and again computes the similarity join.

When all map tasks have been processed, the master partitions all intermediate data on the basis of intermediate keys and passes each partition to a separate reduce thread. This ensures that all similar string pairs which involve r are assigned to the same intermediate partition. Finally, reduce sorts all (r, s) pairs based by edit distance. Optionally, reduce can also emit the number of similar strings found in S for each r , or filter the results found for r on best score.

Parallelizing similarity searches is analog to the parallelization of similarity joins. The main difference is that P_S is replaced with one or a list of search sequences. If not existent, each search pattern is assigned a unique ID. For searches, the `mapTaskList` contains $\langle k_i, v_i \rangle$ pairs where k_i is a partition prefix of and v_i consists of T_{R_i} and the search sequence(s). As for join, similarity search is performed in the map phase.

4 Evaluation

We evaluated the performance of PeARL on a NUMA server with 24 cores and 256 GB main memory available. All experiments were executed with `numactl -localalloc` to control the memory accession strategy and thread placement. Test data sets (see Table 1) are extracted from dbEST¹ as of March 7th, 2011 for the organism mouse. Indexing is linear in the number of indexed strings [12] and is not included in the reported measurements. In terms of memory consumption, PeARL needs roughly 20 GB of main memory to index all infixes of length 2,000 bp in the *C. elegans* genome (roughly 100M strings). For computing the similarity join $\text{III} \bowtie_k \text{IV}$, PeARL needs approx. 8 GB of main memory.

Table 1. Data sets extracted from dbEST

Set	# strings	∅ length (min / max)	# characters
I	10,000	511.99 (49 / 1,190)	5,120,495
II	240,000	455.94 (18 / 2,160)	109,425,487
III	300,000	446.74 (18 / 2,160)	134,023,819
IV	1,000,000	512.12 (7 / 3,920)	512,123,043

¹ www.ncbi.nlm.nih.gov/dbEST/

4.1 Performance of Similarity-Based Operators

First, we compared the performance of all similarity-based operations in PeARL with its predecessor PETER in single-threaded mode. The main difference of both tools is that in PeARL, all parts of the index are kept in main memory whereas in PETER, disk I/O was necessary during search and join. Another difference is that q -gram sets in PeARL are stored persistently in the index whereas previously, q -grams were computed on the fly. Trie partitioning and parallelization was also not present in the predecessor. Overall, we observed that these improvements increased the efficiency of our filters. Whereas in PETER, filtering lead to runtime improvements of up to 80% compared to the baseline with no filters enabled, we now achieve runtime improvements of up to 99% caused by filtering (data not shown).

We evaluated the runtime of similarity search and measured 10,000 individual searches of non-indexed patterns from set I in the PeARL index for set IV, see Fig. 3 (left). In single-threaded mode, searches in PETER ran significantly faster than in PeARL (up to factor 10 on $k = 2$). This is not surprising, as there is some overhead introduced in PeARL by the added functionality for MapReduce based parallelization, which is also present in single-threaded searches. However, we will see in the following section that this overhead pays out for multi-threaded similarity searches and joins. We also compared PeARL to Flamingo, a library for string searching developed at UC Irvine². As displayed in Fig. 3(left), PeARL outperforms Flamingo for search in single threaded mode for small thresholds (factor 20 for $k = 1$ and factor 3 for $k = 2$). For larger k , Flamingo begins to outperform PeARL.

For evaluating the runtime of similarity joins in PETER and PeARL, we computed the join between set IV and varying subsets taken from set II. As shown in Fig. 3 (right), similarity joins in PeARL are computed considerably faster than in PETER. For example, we reached an improvement of factor 3 on $k = 2$ at a join cardinality of $2e+11$. Generally speaking, the implemented improvements in PeARL are the more profitable when indexed string sets grow large. We could not compare PeARL to Flamingo for joins, since no reference implementation was available.

4.2 Scalability of PeARL

We compared the multi-threaded execution of 10,000 individual searches of patterns from set I in set IV with PeARL (24 threads) to a single-threaded execution with PeARL and Flamingo. As displayed in Fig. 4 (left), the multi-threaded execution in PeARL outperforms the single-threaded execution in Flamingo with factors in the range of 6 ($k = 3$) to 57 ($k = 1$). We also observed that the 24-threaded outperforms the single-threaded execution in PeARL with factors in the range of 5.5 ($k = 1$) to 6.2 ($k = 3$). For similarity joins, we could only compare the 24-threaded to the single-threaded execution in PeARL since no

² <http://flamingo.ics.uci.edu/>

external reference implementation was available. Thus, we measured the execution times of $\text{III} \bowtie_{k \in \{1,2,3\}} \text{IV}$. As displayed in Fig. 4 (right), we measured a runtime improvement of factors in the range of 4.2 ($k = 2$) to 4.9 ($k = 1$).

When analyzing the parallelized search and join algorithms in terms of speed-up, the first step is to estimate the fractions of parallelizable and non-parallelizable parts in our algorithms. In general, the parallelizable fractions dominate, since only reading the indices into main memory, extracting tasks from `mapJoinList`, sorting intermediate partitions before executing reduce, and writing the final output to file is performed in serial. We estimated the size of the parallelizable fraction based on the measured speed-up using $N = 24$ CPU cores. According to this, 10 % of our search and 20 % of our join algorithm remain serial.

Figure 5 (left) displays the speed-up of searches of all ESTs from set I in the indexed set IV with regard to the number of CPU cores. We observed that the speed-up for measured runtimes almost perfectly fits the theoretical curve of Amdahl’s law for $P = 0.90$. Similarly, we observed for joins that the measured speed-up fits well to Amdahl’s law for $P = 0.80$ (see Fig. 5 (right)). This indicates that estimating the non-parallelizable fraction with 10 % for searches and 20 % for joins is sound. Using 24 CPU cores with 24 map and reduce workers, we achieve a speed-up of our join algorithm of 4.3. According to that, the maximal speed-up for join is 4.9 using $\geq 1,000$ cores. This indicates that executing the current implementation of PeARL is limited by the serial parts contained in our algorithms.

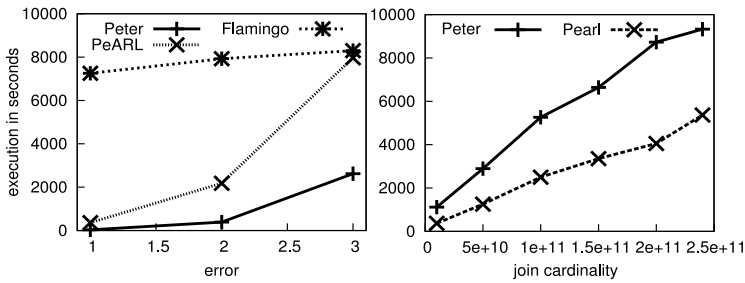


Fig. 3. Performance of single-threaded similarity operations. Left: search of 10,000 patterns from set I in set IV. Right: join $\text{IV} \bowtie_{k=2} \text{II}$ on subsets of II.

5 Related Work

Morrison [10] introduced prefix trees as an index structure for storing strings and exact string matching. Shang et al. [14] extended prefix trees with dynamic programming techniques to perform inexact matching. Prefix pruning was studied in [14] and is based on the observation that edit distance can only grow with prefix length. Aghili et al. [1] proposed character frequency distance based filtering to reduce candidate sets for similarity-based string searches. Indexing methods

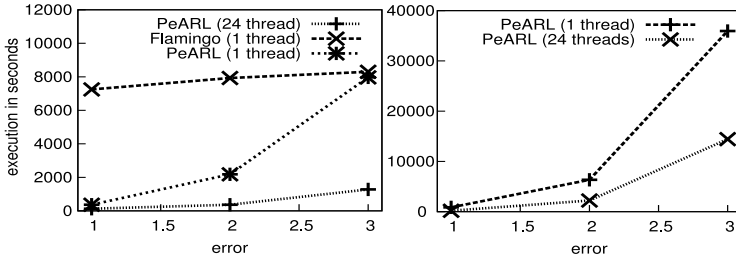


Fig. 4. Performance of parallelization. Left: search. Right: join.

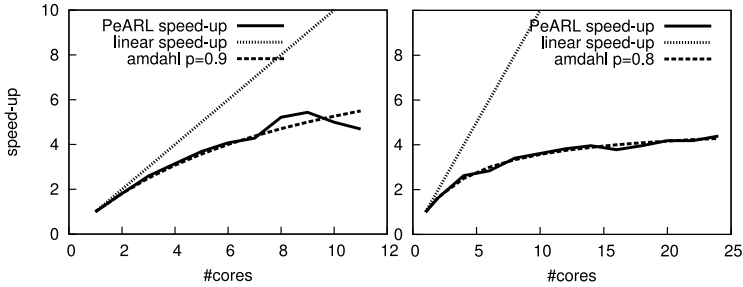


Fig. 5. PeARL speed-up. Left: sim. search on $k = 2$. Right: sim. join on $k = 3$.

based on q -grams restrict search spaces efficiently for edit distance based operations. They take advantage of the observation that two strings are within a small edit distance iff they share a large number of q -grams [15].

The MapReduce programming model for parallel data analysis was initially proposed by Dean and Ghemawat [4]. Vernica et al. [17] present an algorithm set-similarity string joins with distributed MapReduce. We could not compare to their solution, since no in-memory version was available. Ranger et al. [11] developed a MapReduce based programming framework for shared-memory multi-core servers with a scalability almost reaching hand-coded solutions.

A main application for similarity-based string searches and joins in bioinformatics is read alignment. Almost all tools follow the seed-and-extend approach. BLAST [2] seeds the alignment with hash-table indices and extend the initially ungapped seeds with a banded local alignment algorithm. However, algorithms that use only ungapped seeds might miss some valuable alignments. BWA-SW [8] is one tool that allows gap and mismatches in the seeds. We also applied PeARL for read alignment and compared the execution times to BWA-SW. BWA-SW significantly outperforms PeARL (data not shown), but it must be noted that it is a heuristic that misses solutions, while PEARL solves the alignment problem exactly. CloudBurst [13] is another another tool for read alignment using MapReduce on top of Hadoop [3]. A comparison between PEARL and CloudBurst is pending.

6 Conclusions and Future Work

In this paper, we presented PeARL, a data structure and parallel algorithms for similarity-based search and join operations in compressed tries. PeARL is parallelized in main memory with MapReduce on a multi-core server. Our evaluation revealed that the speed-up of our search and join algorithms executed on multi-core servers cannot grow infinitely large due to the serial parts contained in our workflow. We are currently working on reducing these bottlenecks and on performing a detailed comparison between PeARL and CloudBurst.

References

1. Aghili, S.A., Agrawal, D., El Abbadi, A.: BFT: Bit Filtration Technique for Approximate String Join in Biological Databases. In: Nascimento, M.A., de Moura, E.S., Oliveira, A.L. (eds.) SPIRE 2003. LNCS, vol. 2857, pp. 326–340. Springer, Heidelberg (2003)
2. Altschul, S.F., Gish, W., Miller, W., Myers, E.W., Lipman, D.J.: Basic local alignment search tool. *J. Molecular Biology* 215(3), 403–410 (1990)
3. Bialecki, A., Cafarella, M., Cutting, D., O’Malley, O.: Hadoop, <http://hadoop.apache.org/>
4. Dean, J., Ghemawat, S.: MapReduce: simplified data processing on large clusters. *Communications of the ACM* 51(1), 107 (2008)
5. Fickett, J.W.: Fast optimal alignment. *Nucl. Acids Res.* 12(1Part1), 175–179 (1984)
6. Gravano, L., Ipeirotis, P.G., Jagadish, H.V., Koudas, N., Muthukrishnan, S., Srivastava, D.: Approximate string joins in a database (Almost) for free. In: *Proc. VLDB*, pp. 491–500 (2001)
7. Hoad, T.C., Zobel, J.: Methods for identifying versioned and plagiarized documents. *J. American Society for Information Science and Technology* 54, 203–215 (2003)
8. Li, H., Durbin, R.: Fast and accurate long-read alignment with Burrows–Wheeler transform. *Bioinformatics* 26(5), 589–595 (2010)
9. Li, H., Homer, N.: A survey of sequence alignment algorithms for next-generation sequencing. *Briefings in Bioinformatics* 11(5), 473–483 (2010)
10. Morrison, D.R.: PATRICIA – practical algorithm to retrieve information coded in alphanumeric. *Journal of the ACM* 15(4), 514–534 (1968)
11. Ranger, C., Raghuraman, R., Penmetsa, A., Bradski, G.R., Kozyrakis, C.: Evaluating mapreduce for multicore and multiprocessor systems. In: *Proc. HPCA*, pp. 13–24 (2007)
12. Rheinländer, A., Knobloch, M., Hochmuth, N., Leser, U.: Prefix Tree Indexing for Similarity Search and Similarity Joins on Genomic Data. In: Gertz, M., Ludäscher, B. (eds.) SSDBM 2010. LNCS, vol. 6187, pp. 519–536. Springer, Heidelberg (2010)
13. Schatz, M.C.: Cloudburst. *Bioinformatics* 25, 1363–1369 (2009)
14. Shang, H., Merrett, T.H.: Tries for approximate string matching. *IEEE TKDE* 8(4), 540–547 (1996)
15. Sutinen, E., Tarhio, J.: Filtration with q-Samples in Approximate String Matching. In: Hirschberg, D.S., Meyers, G. (eds.) CPM 1996. LNCS, vol. 1075, pp. 50–63. Springer, Heidelberg (1996)
16. Vakali, A., Pokorný, J., Dalamagas, T.: An Overview of Web Data Clustering Practices. In: Lindner, W., Fischer, F., Türker, C., Tzitzikas, Y., Vakali, A.I. (eds.) EDBT 2004. LNCS, vol. 3268, pp. 597–606. Springer, Heidelberg (2004)
17. Vernica, R., Carey, M.J., Li, C.: Efficient parallel set-similarity joins using mapreduce. In: *Proc. SIGMOD*, pp. 495–506 (2010)