

Efficient Optimally Lazy Algorithms for Minimal-Interval Semantics*

Paolo Boldi[†]

Sebastiano Vigna

Dipartimento di Scienze dell'Informazione, Università degli Studi di Milano

Abstract

Minimal-interval semantics [5] associates with each query over a document a set of intervals, called *witnesses*, that are incomparable with respect to inclusion (i.e., they form an antichain): witnesses define the minimal regions of the document satisfying the query. Minimal-interval semantics makes it easy to define and compute several sophisticated proximity operators, provides snippets for user presentation, and can be used to rank documents. In this paper we provide algorithms for computing conjunction and disjunction that are linear in the number of intervals and logarithmic in the number of operands; for additional operators, such as ordered conjunction and Brouwerian difference, we provide linear algorithms. In all cases, space is linear in the number of operands. More importantly, we define a formal notion of *optimal laziness*, and either prove it, or prove its impossibility, for each algorithm. Optimal laziness implies that the algorithms do not assume random access to the input intervals, and read as little input as possible to produce a certain output. We cast our results in the general framework of antichain completions of interval orders, making our algorithms directly applicable to other domains.

Contents

1	Introduction	2
2	Minimal-interval semantics	3
3	Operators	5
4	Lazy evaluation	6
4.1	Optimal laziness	7
5	General remarks	8
6	Algorithms based on indirect queues	9
6.1	Basic comparators	10
6.2	The OR operator	10
6.3	The AND operator	12

*A preliminary version of some of the results in this paper appeared in [3].

[†]This work is partially supported by the EC Project DELIS.

7 Greedy algorithms	15
7.1 The BLOCK operator	15
7.2 The AND _{<} operator	17
7.3 Brouwerian difference	21
8 Previous work	22
9 Conclusions	22

1 Introduction

Search engines are a popular way to retrieve information in the web. However, the classical problem studied by the theory of information retrieval, that of answering a query by returning the set of documents that match the information provided by the user, is complicated by the huge number of documents to be taken into consideration. On the web retrieving many relevant documents is usually not a problem—the documents are simply too many already. Rather than recall, precision (in particular, precision in the first 10–20 results) is the main issue.

A first possibility for extending the user capabilities is *query expansion*, an automatic or semi-automatic mechanism that aims at enriching a given query, by using for example some semantics extracted from the context, or by asking directly the user what is the intended meaning of his/her query. In this case, we start from a very simple query, perhaps expressed in some natural language and finally produce a richer (hopefully, more specific) query that is to be submitted to the search engine.

A different, complementary approach is that of providing the user with more powerful (but understandable) operators, which however requires to depart from the Boolean model. In this paper we pursue this path, focusing on *minimal-interval semantics*, a semantic model that uses *antichains of intervals of natural numbers* to represent the semantics of a query; this is the natural framework in which operators such as ordered conjunction, proximity restriction, etc., can be defined and combined freely. Each interval is a *witness* of the satisfiability of the query, and defines a region of the document that the query satisfies (words in the document are numbered starting from 0, so regions of text are identified with intervals of integers). For instance, a query formed by the conjunction of two terms is satisfied by the minimal intervals of the document containing both terms.

This approach has been defined and studied in full extent by Clarke, Cormack and Burkowski in their seminal paper [5]. They showed that antichains have a natural lattice structure that can be used to interpret conjunctions and disjunctions in queries. Moreover, it is possible to define several additional operators (proximity, followed-by, and so on) directly on the antichains. The authors have also described families of successful ranking schemes based on the number and length of the intervals involved [4].

The main feature of minimal-interval semantics is that, by its very definition, an antichain of intervals cannot contain more than w intervals, where w is the number of words in the document. Thus, it is in principle possible to compute all minimal-interval operators in time linear in the document size. This is not true, for instance, if we consider different interval-semantics approaches in which *all* intervals are retained and indexed (e.g., the PAT system [7] or the `sgrep` tool [10]), as the overall number of regions is quadratic in the document size.

In this paper, we attack the problem of providing efficient algorithms for the computation of such operators. As a subproblem, we can compute the proximity of a set of terms, and indeed we are partly inspired by previous work on proximity [16, 14]. Our algorithms are linear in the number of input intervals. For conjunction and disjunction, there is also a multiplicative logarithmic factor in the number of input antichains, which however can be shown to be essentially unavoidable in the disjunctive case. The space used by all algorithms is linear in the number of input antichains (in

fact, we need to store just one interval per antichain), so they are a very extreme case of stream transformation algorithms [1, 9]. Moreover, our algorithms are (with one exception, for which we prove an impossibility result) *optimally lazy*, that is, while building their results they do not advance the input lists more than necessary.¹

We believe that the existence of (almost) linear lazy algorithms for minimal-interval semantics makes it the natural candidate for advancing web search engines beyond a purely Boolean model: in particular, the possibility of limiting the interval width has a very natural interpretation for the user in terms of proximity, and ordered conjunction has obvious applications.

Minimal intervals can also be used together with other standard information-retrieval techniques. For instance, the Indri search engine [15] expands a query into a number of subqueries, many of which are interval-based, and combines the results.

In Section 2 we briefly introduce minimal-interval semantics, referring to the original paper for examples and motivations. The presentation is rather algebraic, and uses standard terms from mathematics and order theory (e.g., “interval” instead of “extent” as in [5]). The resulting structure is essentially identical to that described in the original paper, but our systematic approach makes good use of well-known results from order theory, making the introduction self-contained. For some mathematical background, see, for instance, Birkhoff’s classic [2].

Another advantage of our approach is that by representing abstractly regions of text as intervals of natural numbers we can easily highlight connections with other areas of computer science: antichains of intervals have been used for role-based access control [6], and for testing distributed computations [11]. The problem of computing operators on antichains has thus an intrinsic interest that goes beyond the problems of information retrieval. This is the reason why we cast all our results in the general framework of antichain completion of intervals on arbitrary (totally) ordered finite sets.

Finally, we present our algorithms. First we discuss algorithms based on queues, and then greedy algorithms.²

2 Minimal-interval semantics

Given a finite totally ordered set O , let us denote with \mathcal{I}_O the set of intervals of O (a subset X of O is an *interval* if $x, y \in X$ and $x < z < y$ imply $z \in X$; note that $\emptyset \in \mathcal{I}_O$) ordered by inclusion. Our working example will always be $w = \{0, 1, \dots, w-1\}$, where w represents the number of words in a document, numbered starting from 0 (see Figure 1); elements of \mathcal{I}_w can be thought of as regions of text.

Given intervals I and J , the interval *spanned* by I and J is the least interval containing I and J (in fact, their least upper bound in \mathcal{I}_O). Nonempty intervals will be denoted by $[\ell . . r]$, where ℓ is the left extreme and r is the right extreme (i.e, the smallest and largest element in the interval). Intervals are ordered by containment: when we want to order them by *reverse* containment instead, we shall write $\mathcal{I}_O^{\text{op}}$ (“op” stands for “opposite”).

The idea behind minimal-interval semantics [5] is that every interval in \mathcal{I}_w is a *witness* that a given query is satisfied by a document made of w words. *Smaller witnesses imply a better match, or more information*; in particular, if an interval is a witness any containing interval is a witness. We also expect that *more witnesses imply more information*. Thus, when expressing the semantics of a query, we discard non-minimal intervals, as there are intervals that provide more relevant information. As a result, minimal-interval semantics associates with each query an *antichain*³ of intervals. For instance,

¹In fact, the algorithms presented here differ significantly from those presented in [3] precisely because of the quest for optimal laziness.

²A free implementation of all algorithms described in this paper is available as a part of MG4J (<http://mg4j.dsi.unimi.it/>).

³An *antichain* of a partial order is a subset of elements that are pairwise incomparable.

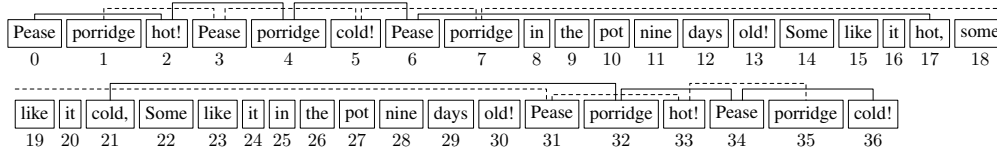


Figure 1: A sample text; the intervals corresponding to the semantics of the query “(hot OR cold) AND porridge AND pease” are shown. For easier reading, every other interval is dashed.

in Figure 1 we see a short passage of text, and the antichain of intervals corresponding to a query. Note that, for instance, the interval $[0..3]$ is not included because it is not minimal.

It is however more convenient to start from an algebraic viewpoint. An *order ideal* X (henceforth called just an *ideal*) is a subset of a partial order that is closed downwards: if $y \leq x$ and $x \in X$, then $y \in X$. The *ideal completion* of an order P is a distributive lattice whose elements are the ideals of P ordered by inclusion. We are interested in computing operators on the ideal completion of $\mathcal{I}_O^{\text{op}}$, which will be the base of our semantics:

$$\mathcal{E}_O = \{ X \subseteq \mathcal{I}_O^{\text{op}} \mid X \text{ is an ideal} \}.$$

It is known that an ideal over a finite partial order is uniquely represented by the antichain of its maximal elements. Intuitively, the antichain of maximal elements is the “upper border” of the ideal. Because of this bijection, antichains of intervals are endowed with a partial order, and with the algebraic structure of a distributive lattice, which turns out to be a very handy representation of \mathcal{E}_O .

The lattice of antichains \mathcal{E}_w thus defined is essentially the classic Clarke–Cormack–Burkowski minimal-interval lattice, with the important difference that since we allow the empty interval, we have a top element that has the empty interval only as a witness. For the purposes of this paper, the difference is immaterial, though.

To make the reader grasp more easily the meaning of \mathcal{E}_O , we now describe in an elementary way its order and its lattice operations (note that we are not giving a definition: the operations are simply the reflection on the set of antichains of those of \mathcal{E}_O). Given antichains A and B , we have

$$A \leq B \iff \forall I \in A \quad \exists J \in B \quad J \subseteq I.$$

Intuitively, $A \leq B$ if every witness I in A (an interval) can be substituted by a better (or equal) witness J in B , where “better” means that the new witness J is contained in I .

Correspondingly, the \vee of two antichains A and B is given by the union of the intervals in A and B from which non-minimal intervals have been eliminated. Finally, the \wedge of A and B is given by the set of all intervals spanned by a pair of intervals $I \in A$ and $J \in B$, from which non-minimal intervals have been eliminated. It is this very natural algebraic structure that has led to the definition of the Clarke–Cormack–Burkowski lattice.

For instance, consider from Figure 1 the positions of “porridge” (1,4,7,32,35), “pease” (0,3,6,31,34) and “hot OR cold” (2,5,17,21,33,36), seen as sets of singleton intervals; by picking one interval from each of the three sets, we generate a large number of spanned intervals, but the minimal ones are just

$$\{ [0..2], [1..3], [2..4], [3..5], [4..6], [5..7], [6..17], [7..31], \\ [21..32], [31..33], [32..34], [33..35], [34..36] \}.$$

A simple snippet extraction algorithm would compute greedily the first k smallest nonoverlapping intervals of the antichain, which would yield, for $k = 3$, the intervals $[0..2]$, $[3..5]$, $[31..33]$, that is, “Pease porridge hot!”, “Pease porridge cold!”, and, again, “Pease porridge hot!”. A ranking

scheme such as those proposed in [4] would use the number and the length of these intervals to assign a score to the document with respect to the query.

Finally, we remark that the intervals in an antichain can be ordered in principle either by left or by right extreme, but these orders can be easily shown to be the same, so we can say that the intervals in an antichain are naturally linearly ordered by their extremes.

3 Operators

For the rest of the paper, we assume that we are operating on antichains based on an unknown total order O for which we just have a comparison operator. We use $\pm\infty$ to denote a special element that is strictly smaller/larger than all elements in O . Before getting to the core of the paper, however, we highlight the connection with query resolution in a search engine.

Search engines use inverted lists to index their document collections [19]. The algorithms described in this paper assume that, besides the documents in which a term appear, the index makes available the *positions* of all occurrences of a term *in increasing order* (this is a standard assumption, as it is necessary to perform gap-encoding).

Given a query (that we shall not define formally: the syntax is implied by our choice of operators), we first obtain the list of documents that could possibly satisfy the query; this is a routine process that involves merging and intersecting lists. Once we know that a certain document might satisfy the query, we want to find its witnesses, if any. To do so, we interpret the terms appearing in the query as lists of singleton intervals (the term positions), and apply in turn each operator appearing in the query. The resulting antichain represents the minimal-interval semantics (i.e., the set of witnesses) of the query with respect to the document.

For completeness, we define explicitly the operators⁴ AND and OR, which are applied to a list of input antichains A_0, A_1, \dots, A_{m-1} , resulting in the \wedge and \vee , respectively, of the antichains A_0, A_1, \dots, A_{m-1} . Besides, we consider other useful operators that can be defined directly on the antichain representation [5]. With this aim, let us introduce a relation \ll between intervals: $I \ll J$ iff $x < y$ for all $x \in I$ and $y \in J$.

1. (“*disjunction operator*”) OR, given input antichains A_0, A_1, \dots, A_{m-1} , returns the set of minimal intervals among those in $A_0 \cup A_1 \cup \dots \cup A_{m-1}$.
2. (“*conjunction operator*”) AND, given input antichains A_0, A_1, \dots, A_{m-1} , returns the set of minimal intervals among those spanned by the tuples in $A_0 \times A_1 \times \dots \times A_{m-1}$.
3. (“*phrasal operator*”) BLOCK, given input antichains A_0, A_1, \dots, A_{m-1} , returns the set of intervals of the form $I_0 \cup I_1 \cup \dots \cup I_{m-1}$ with $I_i \in A_i$ ($0 \leq i < m$) and $I_{i-1} \ll I_i$ ($0 < i < m$).
4. (“*ordered non-overlapping conjunction operator*”) AND $_{<}$, given input antichains A_0, A_1, \dots, A_{m-1} , returns the set of minimal intervals among those spanned by the tuples $\langle I_0, I_1, \dots, I_{m-1} \rangle \in A_0 \times A_1 \times \dots \times A_{m-1}$ satisfying $I_{i-1} \ll I_i$.
5. (“*low-pass operator*”) LOWPASS $_k$, given an input antichain A , returns the set of intervals from A not longer than k .

⁴The reader might be slightly confused by the fact that we are using \wedge and AND to denote essentially the same thing (similarly for \vee and OR). The difference is that \wedge is a binary operator, whereas AND has variable arity. Even if the evaluation of AND could be reduced, by associativity, to a composition of \wedge s, from the viewpoint of the computational effort things are quite different.

6. (“Brouwerian difference⁵ operator”) Given two antichains A (the minuend) and B (the subtrahend), the difference $A - B$ is the set of intervals $I \in A$ for which there is no $J \in B$ such that $J \subseteq I$.

More informally, given input antichains A_0, A_1, \dots, A_{m-1} , the operator BLOCK builds sequences of *consecutive* intervals, each of which is taken from a different antichain, in the given order. It can be used, for instance, to implement a phrase operator. The $\text{AND}_{<}$ operator is an ordered-AND operator that returns intervals spanned by intervals coming from the A_i , much like the AND operator. However, in the case of $\text{AND}_{<}$ the left extremes of the intervals must be nondecreasing, and the intervals must be nonoverlapping. This operator can be used, for instance, to search for terms that must appear in a specified order. LOWPASS_k restricts the result to intervals shorter than a given threshold, and be easily combined with AND or $\text{AND}_{<}$ to implement searches for terms that must not be too far apart, and possibly appear in a given order. Finally, the Brouwerian difference considers the interval in the subtrahend as “poison” and returns only those intervals in the minuend that are not poisoned by any interval in the subtrahend; this operator finds useful applications, for example, in the case of passage search if the poisoning intervals are taken to be small (possibly singleton) intervals around the passage separators (e.g., end-of-paragraph, end-of-sentence, etc.).

Note that the natural lattice operators AND and OR cannot return the empty antichain when all their inputs are nonempty. This is not true of the above operators: for instance, BLOCK might fail to find a sequence of consecutive intervals even if all its inputs are nonempty.

Finally, we remark that all intervals satisfying the definition of the BLOCK operator are minimal. Indeed, assume by contradiction that for two concatenations of minimal intervals we have $[\ell . . r] \subset [\ell' . . r']$ (which implies either $\ell' < \ell$ or $r < r'$). Assume that $\ell' < \ell$ (the case $r < r'$ is similar), and note that removing the first component interval from both concatenations we still get intervals strictly containing one another. We iterate the process, obtaining two intervals of A_{m-1} strictly containing one another.

4 Lazy evaluation

The main point of this paper is that algorithms for computing operators on antichain of intervals should be always *lazy* and *linear in the input intervals*: if an algorithm is lazy, when only a small number of intervals is needed (e.g., for presenting snippets) the computational cost is significantly reduced. Linearity in the input intervals is the best possible result for a lazy algorithm, as input must be read at some point. All algorithms described in this paper satisfy this property, albeit in the case of AND and OR there is also a logarithmic factor in the number of input antichains.

Note that if the inverted index provides random-access lists of term positions, algorithms such as those proposed in [5] might be more appropriate for first-level operators (e.g., logical operators computed directly on lists of term positions), as by accessing directly the term positions they achieve complexity proportional to $ms \log n$, where n is the overall number of intervals in the input antichains, m is the number of antichains, and s is the number of results. Nonetheless, as soon as one combines several operators, the advantage of a lazy linear implementation is again evident. Moreover, s is in principle bounded only by n , and the estimate above hides the fact that the input antichains must have been computed, with a time cost and space occupancy $\Omega(n)$.

The logarithmic factor in the number of antichains can be easily proved to be unavoidable for the OR operator in a model in which intervals can be handled just by comparing their extremes:

Theorem 1 Every algorithm to compute OR that is only allowed to compare interval extremes requires $\Omega(n \log n)$ comparisons for n input intervals.

⁵This operator satisfies the property that $A - B \leq C$ iff $A \leq B \vee C$; it is sometimes called *pseudo-difference*, and its definition is dual to that of relative pseudo-complement [2].

Proof. It is possible to sort n distinct integers by computing the OR of n antichains, each made by just one singleton interval containing one of the integers to be sorted. The resulting antichain is exactly the list of sorted integers. By an application of the $\Omega(n \log n)$ lower bound for sorting in this model, we get to the result. ■

4.1 Optimal laziness

The term “lazy” is usually quoted informally, in particular in the context of functional or declarative programming. In this paper we consider algorithms that access input antichains under the form of lists that return the corresponding intervals in their natural order. We want to define formally a notion of laziness that makes it possible proving rigorously optimality results. We restrict to algorithms that read their inputs from an array of lists. Each list is accessible via a “next” function that returns the next element from the list, and when a list is empty it returns **null**. Analogously, each algorithm has a “next” function that returns the next output, and when the output is over it returns **null**. So such algorithms can be thought of as producing an output list.

Given an algorithm \mathcal{A} , an input I (i.e., an array of lists), let us write $\rho_i^{\mathcal{A}}(I, p)$ for the number of elements (including possibly **null**) read by \mathcal{A} from the i -th list of the input array I when the p -th output is produced (sometimes, we will omit \mathcal{A} , I or p when they are clear from the context); when writing $\rho_i^{\mathcal{A}}(I, p)$ we shall always assume that the $0 \leq i < m$ (where m is the number of input lists) and that the output of \mathcal{A} on input I contains at least p intervals.

A first property that we would like our algorithms to feature is that there is no algorithm that uses strictly less inputs:

Definition 1 Two algorithms are *functionally equivalent* iff they produce the same output list when they are given the same input lists. An algorithm \mathcal{A} is *minimally lazy* if, for every functionally equivalent algorithm \mathcal{B} such that

$$\rho_i^{\mathcal{B}}(I, p) \leq \rho_i^{\mathcal{A}}(I, p)$$

for all I in the set of inputs and all p , we actually have

$$\rho_i^{\mathcal{B}}(I, p) = \rho_i^{\mathcal{A}}(I, p).$$

In fact, for most of our algorithms we will be able to prove a more interesting property:

Definition 2 An algorithm \mathcal{A} is *k-lazy* iff for every functionally equivalent algorithm \mathcal{B} , and for all input I , and all i and p we have

$$\rho_i^{\mathcal{A}}(I, p) \leq \rho_i^{\mathcal{B}}(I, p) + k.$$

An algorithm \mathcal{A} is *optimally lazy* if it is k -lazy for some k , and there exist no functionally equivalent $(k - 1)$ -lazy algorithm.

Optimally lazy algorithms advance their inputs as little as possible when emitting an output. Minimally optimally lazy algorithm have the further property that no improvement can be obtained on a particular input without getting a worse result on some other input. Note that since by definition there are no k -lazy algorithms when k is negative, a 0-lazy algorithm is always minimally and optimally lazy.

There is a subtlety in Definition 1 and 2 that is worth remarking. By requiring that the parameter p is never greater than the number of intervals in the output, we are not considering how many elements are read from the input lists to emit the final **null**. In principle, this choice implies that even minimally optimally lazy algorithms may consume useless input elements to emit their final **null**. A more thorough analysis would be required to include also this case, but it would yield a further

subdivision of the above taxonomy of optimality: indeed, for some problems we consider it is easy to show there is no **null**-optimal solution. We think that such an analysis would add little value to the present work, as behaving lazily on non-**null** outputs is a sufficiently strong property by itself.

5 General remarks

In the description and in the proofs of our algorithms, we use interchangeably A_i to denote the i -th input *antichain* and the *list* returning its intervals in their natural order (and, ultimately, **null**). This ambiguity should cause no difficulty to the reader.

To simplify the exposition, in the pseudocode we often test whether a list is empty. Of course, this is not allowed by our model, but in all such cases the following instruction retrieves the next interval from the same list. Thus, the test can be replaced by a call that retrieves the next interval and tests for **null**. Finally, we can assume that after the function “next” has returned **null** once, it will keep returning **null** thereafter: this behaviour can be obtained by using an extra flag that avoids entering the function altogether.

In all our algorithms, we do not consider the case of inputs equal to the top of the lattice (the antichain formed by the empty interval). For all our operators, the top either determines entirely the output (e.g., OR) or it is irrelevant (e.g., AND). Analogously, we do not consider the case of inputs equal to the bottom of the lattice (the empty antichain), which can be handled by a test on the first input read.

More generally, when proving optimal laziness, it is common to meet situations in which an initial check is necessary to rule out obvious outputs. The initial check can make the algorithm analysis more complicated, as its logic could be wildly different from the true algorithm behaviour. To simplify this kind of analysis, we prove the following metatheorem, which covers the cases just described; in the statement of the theorem, \mathcal{A} represent the algorithm performing the initial check, whereas \mathcal{B} does the real job:

Theorem 2 Let \mathcal{B} be an algorithm defined on a set of inputs B , and \mathcal{A} be defined on a larger set of inputs $A \supseteq B$, and such that

- on all inputs $I \in B$, \mathcal{A} outputs a one-element list containing a special element, say \perp , and
- for all $I \in B$ and all i , $\rho_i^{\mathcal{A}}(I, 1) \leq \rho_i^{\mathcal{B}}(I, 1)$.

Then, there exists an algorithm, denoted by $\mathcal{A} \star \mathcal{B}$, such that

- $\mathcal{A} \star \mathcal{B}$ is functionally equivalent to \mathcal{B} on B ;
- $\mathcal{A} \star \mathcal{B}$ is functionally equivalent to \mathcal{A} on $A \setminus B$;
- if \mathcal{A} and \mathcal{B} are (minimally) optimally lazy on $A \setminus B$ and B , respectively, then $\mathcal{A} \star \mathcal{B}$ is (minimally) optimally lazy on A .

Proof. Algorithm $\mathcal{A} \star \mathcal{B}$ simulates algorithm \mathcal{A} and caches the input read so far. If \mathcal{A} emits any element different from \perp , the simulation goes on until \mathcal{A} is done, without caching the input any longer; otherwise, $\mathcal{A} \star \mathcal{B}$ starts executing \mathcal{B} on the cached input and possibly on the remaining part of the input until \mathcal{B} is done.

It is immediate to check that $\mathcal{A} \star \mathcal{B}$ is indeed functionally equivalent to \mathcal{A} and \mathcal{B} on $A \setminus B$ and B , respectively, and moreover

$$\rho_i^{\mathcal{A} \star \mathcal{B}}(I, p) = \begin{cases} \rho_i^{\mathcal{A}}(I, p) & \text{if } I \in A \setminus B \\ \rho_i^{\mathcal{B}}(I, p) & \text{if } I \in B. \end{cases}$$

enqueue(Q, x)	insert item with index x in the queue
topIndex(Q)	returns the index of the top item
top(Q)	returns the top item
dequeue(Q)	returns the index of the top item and deletes it from the queue
change(Q)	signals that the top item has changed
size(Q)	returns the number of indices currently in the queue

Table 1: The operations available for an indirect priority queue.

Suppose now that \mathcal{A} is a -lazy and \mathcal{B} is b -lazy for some minimal a and b , and let $c = \max\{a, b\}$. For every algorithm \mathcal{C} that is functionally equivalent to $\mathcal{A} \star \mathcal{B}$, we have that $\rho_i^{\mathcal{C}}(I, p) \leq \rho_i^{\mathcal{B}}(I, p) + b$ for all $I \in B$, and $\rho_i^{\mathcal{C}}(I, p) \leq \rho_i^{\mathcal{A}}(I, p) + a$ for all $I \in A \setminus B$. But then, using the observation above, $\rho_i^{\mathcal{C}}(I, p) \leq \rho_i^{\mathcal{A} \star \mathcal{B}}(I, p) + c$ for all $I \in A$, so $\mathcal{A} \star \mathcal{B}$ is c -lazy.

Suppose now that \mathcal{C} is functionally equivalent to $\mathcal{A} \star \mathcal{B}$ but that it is $(c - 1)$ -lazy, and assume that $c = b$ (the other case is analogous). Then, for all $I \in B$, $\rho_i^{\mathcal{C}}(I, p) \leq \rho_i^{\mathcal{A} \star \mathcal{B}}(I, p) + c - 1 = \rho_i^{\mathcal{B}}(I, p) + b - 1$; but since \mathcal{C} is also functionally equivalent to \mathcal{B} on B , the latter inequality contradicts the minimality of b .

For minimal laziness, suppose that \mathcal{C} is functionally equivalent to $\mathcal{A} \star \mathcal{B}$ and such that $\rho_i^{\mathcal{C}}(I, p) \leq \rho_i^{\mathcal{A} \star \mathcal{B}}(I, p)$ for all $I \in A$. In particular, this means that $\rho_i^{\mathcal{C}}(I, p) \leq \rho_i^{\mathcal{A}}(I, p)$ for all $I \in A \setminus B$, and $\rho_i^{\mathcal{C}}(I, p) \leq \rho_i^{\mathcal{B}}(I, p)$ for all $I \in B$. The minimal laziness of \mathcal{A} and \mathcal{B} imply that $\rho_i^{\mathcal{C}}(I, p) = \rho_i^{\mathcal{A}}(I, p)$ for all $I \in A \setminus B$ and $\rho_i^{\mathcal{C}}(I, p) = \rho_i^{\mathcal{B}}(I, p)$ for all $I \in B$, hence $\rho_i^{\mathcal{C}}(I, p) = \rho_i^{\mathcal{A} \star \mathcal{B}}(I, p)$ for all $I \in A$. ■

Incidentally, we observe that $\mathcal{A} \star \mathcal{B}$ requires in general more space than \mathcal{A} or \mathcal{B} , because of caching; nonetheless, in all our applications we will need to cache just one item per input list.

6 Algorithms based on indirect queues

The algorithms we provide for AND and OR are inspired by the plane-sweeping technique used in [16] for their proximity algorithm, which is on its own right a variant of the standard sorted-list merge. The algorithms are implemented using an *indirect priority queue*.

An indirect priority queue Q is a data structure based on an array (called the *reference array*), which is managed outside the queue itself, and a priority order that compares items from the reference array. At each time, the queue contains a set of indices into the reference array (initially, a specified set, possibly empty). An array index x can be added to the queue calling the function enqueue(Q, x).

The *index* of the least item in the reference array with respect to the priority order can be accessed by invoking the function topIndex(Q). The index of the least item with respect to the priority order is also returned by dequeue(Q), which additionally removes the index from Q . Analogously, top(Q) return the least *item* in the reference array with respect to the priority order.

The data structure assumes that the only item of the reference array that might change its value is the top item. Such a change must be communicated immediately to the queue by calling the function change(Q). Table 1 summarises the operations available on an indirect priority queue.⁶

A trivial array-based implementation requires linear space (in the number of input lists) and has constant cost for all operations modifying the queue, whereas retrieving the top requires linear time. A better implementation uses a priority queue (e.g., based on a heap) with linear space and logarithmic time complexity for all operations modifying the queue. Sophisticated heaps with linear costs

⁶Actually, a more appropriate name would be *semi-indirect* queue: an indirect queue has a change operation that restores the correct state after a change in the value associated to any index.

for several operations do not modify significantly the overall behaviour, as each time the queue is advanced the interval corresponding to the top index becomes greater: there are data structures that make it possible to *decrease* in constant time the top, but not *increase it* (otherwise we could sort in linear time by comparison).

All algorithms based on indirect priority queues have time complexity $O(n \log m)$ if the input is formed by m antichains containing n intervals overall, and use $O(m)$ space. This is immediate, as all loops contain exactly one queue advancement.

6.1 Basic comparators

Our algorithms will be based on two priority orders. The first one, denoted by \trianglelefteq , is defined by

$$[\ell \dots r] \trianglelefteq [\ell' \dots r'] \iff r < r' \text{ or } r = r' \text{ and } \ell \geq \ell'.$$

In other words, $[\ell \dots r] \trianglelefteq [\ell' \dots r']$ if $[\ell \dots r]$ ends before or is a suffix of $[\ell' \dots r']$. Note in particular that (somewhat counterintuitively) $[\ell \dots r] \trianglelefteq [\ell' \dots r']$ iff $\ell \geq \ell'$.

The second order, denoted by \preceq , is defined by

$$[\ell \dots r] \preceq [\ell' \dots r'] \iff \ell < \ell' \text{ or } \ell = \ell' \text{ and } r \geq r'.$$

In other words, $[\ell \dots r] \preceq [\ell' \dots r']$ if $[\ell \dots r]$ starts before or prolongs $[\ell' \dots r']$. Note in particular that $[\ell \dots r] \preceq [\ell' \dots r']$ iff $r \geq r'$, and that the following implication holds:

$$[\ell \dots r] \subseteq [\ell' \dots r'] \implies [\ell \dots r] \trianglelefteq [\ell' \dots r'] \text{ and } [\ell' \dots r'] \preceq [\ell \dots r]$$

The algorithms for AND/OR use an indirect priority queue with priority order \preceq or \trianglelefteq . The reference array underlying the queue contains one interval per input antichain. In the initialisation phase, the reference array is filled with the first interval from each antichain, and the queue contains all indices.

To simplify the description, we define a procedure $\text{advance}(Q)$ that updates with the next interval the list associated with the top index and notifies the queue of the change. If the update cannot be performed because the list is empty, the top index is dequeued. The function is described in pseudocode in Algorithm 1.

Algorithm 1 The advance function.

```

0  procedure advance( $Q$ ) begin
1     $i \leftarrow \text{topIndex}(Q)$ ;
2    if  $A_i$  is not empty then
3       $[\ell_i \dots r_i] \leftarrow \text{next}(A_i)$ ;
4      change( $Q$ )
5    else
6      dequeue( $Q$ )
7    end;
8  end;
```

6.2 The OR operator

We start with the simplest nontrivial operator. To compute the OR of the antichains A_0, A_1, \dots, A_{m-1} , we merge them using an indirect priority queue Q with priority order \trianglelefteq .

We keep track of the last interval c returned (initially, $c = [-\infty \dots -\infty]$). When we want to compute the next interval, we advance Q as long as the top interval contains c , and then if the queue is not empty we return the top. The algorithm⁷ is described in pseudocode in Algorithm 2.

Theorem 3 Algorithm 2 for OR is correct.

Proof. First of all, note that all intervals in A_0, A_1, \dots, A_{m-1} are assigned to c at some point, unless they contain a previously returned interval. Thus, we just have to prove that only minimal intervals are returned.

Let $[\ell \dots r]$ be a non-minimal element of $A_0 \cup A_1 \cup \dots \cup A_{m-1}$, and $[\ell' \dots r']$ the *largest* (according to \preceq) minimal interval contained in $[\ell \dots r]$. After returning $[\ell' \dots r']$ (which certainly appears at the top of the queue *before* $[\ell \dots r]$ due to the fact that \subseteq implies \preceq), all intervals in the queue have a right extreme larger than or equal to r' . When we advance the queue, and until we get past $[\ell \dots r]$, the top interval will always contain $[\ell' \dots r']$, for otherwise there would be a minimal interval with right extreme between r' and r , and $[\ell' \dots r']$ would not be largest. Thus, the while loop will remove at some point $[\ell \dots r]$.

To prove that all returned intervals are unique, we just have to note that when I is returned, all other copies of I are in the reference array. Thus, at the next call the while loop will be repeated until all remaining copies are discarded. ■

Algorithm 2 The algorithm for the OR operator. Note that the second part of the while condition is actually equivalent to “ $\text{left}(\text{top}(Q)) \leq \text{left}(c)$ ” due to the monotonicity of the top-interval right extreme.

```

0  Initially  $c \leftarrow [-\infty \dots -\infty]$  and  $Q$  contains one interval from each  $A_i$ .
1  function next begin
2    while  $Q$  is not empty and  $c \subseteq \text{top}(Q)$  do
3      advance( $Q$ )
4    end;
5    if  $Q$  is empty then return null;
6     $c \leftarrow \text{top}(Q)$ ;
7    return  $c$ 
8  end;

```

Theorem 4 Algorithm 2 for OR is minimally and optimally lazy.

Proof. We show that the algorithm (let us call it \mathcal{A}) is 0-lazy. The first output of the algorithm requires reading exactly one interval from each list. No correct algorithm can emit the first output without this data.

Suppose now that for an algorithm \mathcal{A}^* it happens that

$$\rho_i^{\mathcal{A}^*}(I, p) < \rho_i^{\mathcal{A}}(I, p)$$

for some input I and some i and p . Since upon returning the p -th output $[\ell \dots r]$ the reference array contains the least interval (w.r.t. \preceq) after $[\ell \dots r]$ from each list, this means that \mathcal{A}^* emits $[\ell \dots r]$ having read from the i -th input list an interval $[\ell' \dots r']$ strictly smaller than $[\ell \dots r]$ according to \preceq ; this means that either $r' < r$, or $r' = r$ and $\ell < \ell'$, but the latter case is ruled out by minimality

⁷Note that this algorithm, as discussed in Section 8, can be derived from the dominance algorithms presented in [12].

of $[\ell \dots r]$. Thus, $r' < r$, and \mathcal{A}^* would return an incorrect output if the i -th input list would return $[s \dots s]$ as next input, with $r' < s < r$. ■

Note that for the last proof the genericity of the underlying order is essential: if we knew that there are no elements between r' and r we could not obtain the contradiction.

6.3 The AND operator

Then AND operator is much more subtle. The priority order of Q is \preceq , and additionally the queue keeps track of the largest right extreme of intervals in the reference array, which will we call *the right extreme of Q* (we just need a variable that is maximised with the right extreme of each new input interval, as at the first dequeuing we shall return **null**). We say that Q is *full* if it contains exactly m indices.

At any time, the interval *spanned* by Q is the interval defined by the left extreme of the top interval and the right extreme of Q : it will be denoted by $\text{span}(Q)$. Clearly, it is the minimum interval containing all intervals currently in the queue.

We keep track of the last interval c returned (initially, $c = [-\infty \dots -\infty]$). When we want to compute the next interval, we first advance Q until the spanned interval does not contain c , and in case Q is no longer full we return **null**. Then, we store the interval $[\ell \dots r]$ currently spanned by Q as a candidate and advance Q . If the new interval spanned by Q is included in $[\ell \dots r]$ we repeat the operation, updating the candidate. Otherwise (or if Q is no longer full) we just return the candidate. The algorithm is described in pseudocode in Algorithm 3.

Algorithm 3 The algorithm for the AND operator. Note that the second part of the first while condition can be substituted with “ $\text{left}(\text{top}(Q)) = \text{left}(c)$ ” because of monotonicity of the largest right extreme, and that the second part of the second while condition can be substituted with “ $\text{right}(c) = \text{right}(Q)$ ” by monotonicity of the top-interval left extreme.

```

0  Initially  $c \leftarrow [-\infty \dots -\infty]$  and  $Q$  contains one interval from every  $A_i$ .
1  function next begin
2    while  $Q$  is full and  $c \subseteq \text{span}(Q)$  do
3      advance( $Q$ )
4    end;
5    if  $Q$  is not full then return null;
6    do
7       $c \leftarrow \text{span}(Q)$ ;
8      if  $c = \text{top}(Q)$  then return  $c$ ;
9      advance( $Q$ )
10   while  $Q$  is full and  $\text{span}(Q) \subseteq c$ ;
11   return  $c$ 
12 end;

```

Theorem 5 Algorithm 3 for AND is correct.

Proof. We say that a queue configuration is *complete* if it contains all copies of the top interval from all lists that contain it. Now observe that *every complete configuration of an indirect priority queue is entirely defined by its top interval*. More precisely, if the top is an interval I from list i , then for every other list j the corresponding interval J in the queue is the minimum interval in A_j larger than

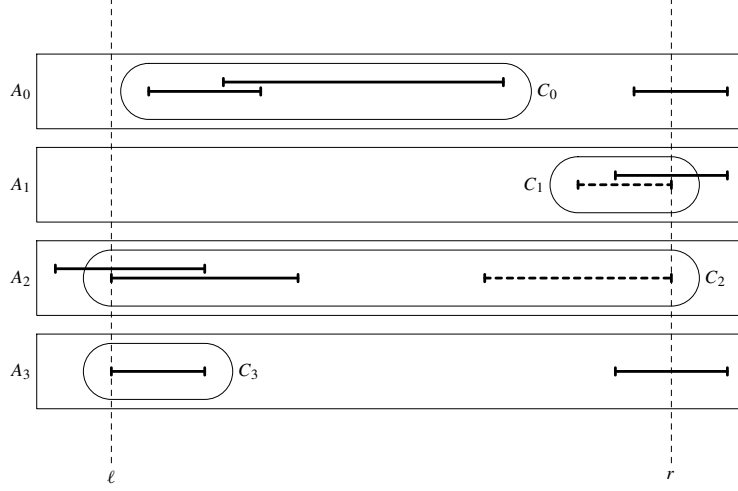


Figure 2: A sample configuration found in the proofs of Theorems 5 and 6. The dashed intervals are right delimiters. The first two input lists are in the inner set; the last two input lists are in the conflict set; the last input list is also in the resolution set.

or equal to I (according to \preceq). Indeed, suppose by contradiction that there is another interval K from A_j satisfying

$$I \preceq K < J.$$

Then, at some point K must have entered the queue, and when it has been dequeued the top must have become some interval $I' \preceq I$, so we get

$$K \preceq I' \preceq I \preceq K,$$

which yields $K = I$: a contradiction, as we assumed the configuration of the queue to be complete.

We now show that for every minimal interval $[\ell \dots r]$ in the AND of A_0, A_1, \dots, A_{m-1} there is a complete configuration of Q spanning $[\ell \dots r]$. Consider for each i the set C_i of intervals of A_i contained in $[\ell \dots r]$. At least one of these sets must contain a (necessarily unique) *right delimiter*, that is, an interval of the form $[\ell' \dots r]$ (see Figure 2). Moreover, at least one of the sets containing a delimiter must be a singleton. Indeed, if every C_i containing a right delimiter would also contain some other interval, the right extreme of that interval would clearly be smaller than r : the maximum of such right extremes, say $r' < r$, would define a spanned interval $[\ell \dots r']$ showing that $[\ell \dots r]$ was not minimal. We conclude that at least one C_i , say $C_{\bar{i}}$, is a singleton containing a right delimiter.

Let I_i be the leftmost interval in each C_i ; these intervals are a complete configuration of Q : if $I_i = [\ell \dots r']$ is the \preceq -smallest among such intervals and if $I_i \in A_j$ necessarily $I_i = I_j$, because A_j cannot contain two intervals with the same left extreme. The set of intervals also spans $[\ell \dots r]$ (because the right extreme of $I_{\bar{i}}$ is r , and the left extreme of the \preceq -least interval I_i is ℓ). We conclude that all minimal intervals in the output are eventually spanned by Q .

However, no minimal interval can be spanned during the first while loop, unless it has been already returned, as all intervals spanned in that loop contain a previously returned interval (notice that at the first call the loop is skipped altogether). Finally, if an interval is spanned in the second while loop and we do not get out of the loop, the next candidate interval will be smaller or equal. We conclude that sooner or later all minimal intervals cause an interruption of the second while loop, and are thus returned.

We are left to prove that if an interval is returned, it is necessary minimal. If we exit the loop using the check on the top interval, the returned interval is necessary minimal. Otherwise, assume that the interval $[\ell \dots r]$ spanned by Q at the start of the second while loop is not minimal, so $[\bar{\ell} \dots r] \subset [\ell \dots r]$, for some minimal interval $[\bar{\ell} \dots r]$ that will be necessarily spanned later (as we already proved that all minimal intervals are returned). Since the right extreme of Q is nondecreasing, the second while loop will pass through intervals of the form $[\ell' \dots r]$, with $\ell < \ell' < \bar{\ell}$, until we exit the loop.

Finally, we remark the uniqueness of all returned intervals is guaranteed by the first while loop. ■

Note that our algorithm for AND *cannot* be 0-lazy, because the choices made by the queue for equal intervals cause different behaviours. For instance, on the input lists $\{[0 \dots 0], [2 \dots 2]\}$, $\{[1 \dots 1]\}$, $\{[0 \dots 0], [2 \dots 2]\}$ the algorithm advances the last list before returning $[0 \dots 1]$, but there is a variant of the same algorithm that keeps intervals sorted lexicographically by \leq and by input list index, and this variant would advance the first list instead.

Nonetheless:

Theorem 6 Algorithm 3 for AND is minimally and optimally lazy.

Proof. We denote Algorithm 3 with \mathcal{A} , and let \mathcal{A}^* be a functionally equivalent algorithm. Let us number the intervals appearing in a certain input $I = A_0, A_1, \dots, A_{m-1}$: in particular, let $[\ell_i^j \dots r_i^j]$ be the j -th interval appearing in A_i . For sake of simplicity, let us identify the **null** returned as last element by the input lists with the interval $[\infty \dots \infty]$ (it is immediate to see that \mathcal{A} behaves identically). Let us write ρ_i (respectively, ρ_i^*) for $\rho_i^{\mathcal{A}}(I, p)$ (respectively, $\rho_i^{\mathcal{A}^*}(I, p)$), and $[\ell \dots r]$ be the p -th output interval; let also s_i be the index of the first interval in list A_i that is included in $[\ell \dots r]$.

We divide the indices of the input lists in two sets: the *inner set* is the set of indices i for which $\ell < \ell_i^{s_i}$ (that is, the first interval of A_i included in $[\ell \dots r]$ has left extreme larger than ℓ); the *conflict set* is the set of indices i for which $\ell = \ell_i^{s_i}$ (that is, the first interval of A_i included in $[\ell \dots r]$ has left extreme equal to ℓ). Finally, the *resolution set* is a subset of the conflict set containing those indices i for which $r_i^{s_i+1} > r$ (that is, the successor of the first interval of A_i included in $[\ell \dots r]$ is no longer contained in $[\ell \dots r]$). Note that the resolution set is always nonempty, or otherwise $[\ell \dots r]$ would not be minimal (recall that we substituted **null** with $[\infty \dots \infty]$). The situation is depicted in Figure 2.

We remark the following facts:

- (i). for all i , $\rho_i^* \geq s_i$; that is, when \mathcal{A}^* outputs $[\ell \dots r]$ it has read at least the first interval of the antichain with left extreme larger than or equal to ℓ ; otherwise, \mathcal{A}^* would emit $[\ell \dots r]$ even on a modified input in which A_i has no intervals contained in $[\ell \dots r]$ (such intervals have index equal to or greater than s_i , so they have not been seen by \mathcal{A}^* , yet);
- (ii). for all i in the inner set, $\rho_i = s_i \leq \rho_i^*$;
- (iii). for all i in the conflict set, $\rho_i \in \{s_i, s_i + 1\}$; that is, in the case an antichain *does* contain an interval J with left extreme ℓ , either the last interval read by \mathcal{A} when $[\ell \dots r]$ is output is exactly J , or it is the interval just after J ;
- (iv). if for some i we have $[\ell_i^{s_i} \dots r_i^{s_i}] = [\ell \dots r]$, then $\rho_j = s_j$ for all j , because we exit the second while loop at line 8;
- (v). otherwise, there is a unique index \bar{i} in the resolution set such that $\rho_{\bar{i}} = s_{\bar{i}} + 1$ (i.e., $r_{\bar{i}}^{\rho_{\bar{i}}} > r$), and for all other resolution indices i we have $\rho_i = s_i$ (i.e., $r_i^{\rho_i} \leq r$); this happens because we interrupt the second while loop when we see the first interval whose right extreme exceeds r (at line 10).

Let us first prove that \mathcal{A} is 1-lazy by showing that $\rho_i \leq \rho_i^* + 1$: this is true for all indices in the inner set because of (ii), and for all indices in the conflict set because of $\rho_i \leq s_i + 1 \leq \rho_i^* + 1$ (by (iii) and (i)).

Now, let us show that \mathcal{A}^* cannot be 0-lazy. Suppose it is such; then, in particular, $\rho_i^* \leq \rho_i$ for all indices i , and we can assume w.l.o.g. that $\rho_i^* < \rho_i$ for some i (if for all inputs, all output prefixes and all i we had $\rho_i^* = \rho_i$, then we would conclude that \mathcal{A} is 0-lazy as well, contradicting the observation made before this theorem).

Note that we can also assume w.l.o.g. not to be in case (iv) (as in that case $\rho_i = \rho_i^*$ for all i), which also implies that $\ell \neq r$. Thus, the unique index \bar{i} of (v) is also the only index in the resolution set such that $\rho_{\bar{i}}^* = s_{\bar{i}} + 1$ (\mathcal{A}^* must advance some list in the resolution set, or it would emit a wrong output on a modified input in which the $(s_i + 1)$ -th interval of A_i is $[r \dots r]$ for all i in the conflict set).

Let i_0, i_1, \dots, i_{t-1} be the indices in the conflict set for which $\rho_{i_p} = s_{i_p} + 1$, in the order in which they are accessed from the corresponding lists by \mathcal{A} : clearly $i_{t-1} = \bar{i}$ is the only resolution index in this sequence, by (v). Let j_0, j_1, \dots, j_{u-1} be the indices in the conflict set for which $\rho_{j_p}^* = s_{j_p} + 1$, in the order in which they are accessed from the corresponding lists by \mathcal{A}^* . Necessarily, $\{j_0, j_1, \dots, j_{u-1}\} \subseteq \{i_0, i_1, \dots, i_{t-1}\}$ (because $s_{j_p} + 1 = \rho_{j_p}^* \leq \rho_{j_p} \leq s_{j_p} + 1$, hence $\rho_{j_p} = s_{j_p} + 1$) and inclusion is strict (because, for some index i , $\rho_i^* < \rho_i$, hence $s_i \leq \rho_i^* < \rho_i \leq s_i + 1$, which implies that $i = i_v$ for some v , whereas $i \neq j_v$ for all v).

Let p be the first position that \mathcal{A} and \mathcal{A}^* choose differently, that is, $i_p \neq j_p$ (this happens at least at the position of j_0, j_1, \dots, j_{u-1} where \bar{i} appears). We build a new input similar to A_0, A_1, \dots, A_{m-1} , except for A_{i_p} and A_{j_p} , which are identical up to their interval of left extreme ℓ ; then, A_{i_p} continues with $[r' \dots r']$ for some $r' > r$ (so i_p is in the resolution set), whereas A_{j_p} continues with $[r \dots r]$ (so j_p is in the inner set). On this input, to output $[\ell \dots r]$ \mathcal{A} advances the input list A_{j_p} strictly less than \mathcal{A}^* , which contradicts the assumption on \mathcal{A}^* . ■

7 Greedy algorithms

The remaining operators admit greedy algorithms: they advance the input lists in a specified order until some condition becomes true. The case of LOWPASS_k is of course trivial, and the algorithm for BLOCK is essentially a restatement in terms of intervals of the folklore algorithm for phrasal queries. They are both optimally and minimally lazy. The case of $\text{AND}_<$ and Brouwerian difference are more interesting: $\text{AND}_<$ is the only algorithm for which we prove the impossibility of an optimally or minimally lazy implementation in the general case.

All greedy algorithms have time complexity $O(n)$ if the input is formed by m antichains containing n intervals overall, and use $O(m)$ space. This is immediate, as all loops advance at least one input list.

7.1 The BLOCK operator

The BLOCK operator is the only one that can be implemented exclusively if the underlying total order is discrete, that is, if it admits a notion of successor. In discussing this algorithm, we shall assume that every element $x \in O$ has a successor, denoted by $x + 1$, satisfying $x < x + 1$ and $x \leq y \leq x + 1 \implies x = y$ or $y = x + 1$.

We keep track of a current interval for all lists A_0, A_1, \dots, A_{m-1} ; initially, these intervals are set to $[-\infty \dots -\infty]$. When we want to compute the next interval, we update the interval associated to the first list. Then, we try to fix index i (initially, $i = 1$). To do so, we advance the list A_i until the returned interval has left extreme larger than the right extreme of the current interval for A_{i-1} . If we go too far, we just advance the first list, reset i to 1 and restart the process, otherwise we increment i . When we find an interval for A_{m-1} we return the interval spanned by all current intervals. The algorithm is described in pseudocode in Algorithm 4.

Theorem 7 Algorithm 4 for BLOCK is correct.

Algorithm 4 The algorithm for the BLOCK operator.

```

0  Initially  $[\ell_k \dots r_k] \leftarrow [-\infty \dots -\infty]$  for all  $0 \leq k < m$ .
1  function next begin
2    if  $A_0$  is empty then return null;
3     $[\ell_0 \dots r_0] \leftarrow \text{next}(A_0)$ ;
4     $i \leftarrow 1$ ;
5    while  $i < m$  do
6      while  $\ell_i \leq r_{i-1}$  do
7        if  $A_i$  is empty then return null;
8         $[\ell_i \dots r_i] \leftarrow \text{next}(A_i)$ 
9      end;
10     if  $\ell_i = r_{i-1} + 1$  then  $i \leftarrow i + 1$ 
11     else begin
12       if  $A_0$  is empty then return null;
13        $[\ell_0 \dots r_0] \leftarrow \text{next}(A_0)$ ;
14        $i \leftarrow 1$ 
15     end
16   end;
17   return  $[\ell_0 \dots r_{m-1}]$ 
18 end;

```

Proof. At the start of an iteration of the external while loop (line 5) with a certain index i we clearly have $r_k + 1 = \ell_{k+1}$ for $k = 0, 1, \dots, i - 2$. Thus, if we complete the execution of the loop we certainly return a correct interval.

To complete the proof, we start by proving the following invariant property: at line 5, for all $0 < j < m$ there are no intervals in A_j with left extreme in $[r_{j-1} + 1 \dots \ell_j - 1]$. In other words, the j -th current interval $[\ell_j \dots r_j]$ has either left extreme smaller than or equal to r_{j-1} , or it is the first interval in A_j whose left extreme is larger than r_{j-1} . The property is trivially true at the beginning, and advancing $[\ell_0 \dots r_0]$ cannot change this fact. We are left to prove that the execution of the internal while loop (line 6) cannot either.

During the execution of the loop at line 6, only $[\ell_i \dots r_i]$ can change. This affects the invariant because it modifies the intervals $[r_{i-1} + 1 \dots \ell_i - 1]$ and $[r_i + 1 \dots \ell_{i+1} - 1]$, but in the second case the interval is made smaller, so the invariant is *a fortiori* true. In the first case, at the beginning of the execution of the internal while loop either $r_{i-1} + 1 \leq \ell_i - 1$, that is, $r_{i-1} < \ell_i$, so the loop is not executed at all and the invariant cannot change, or $r_{i-1} + 1 > \ell_i - 1$, which means that the interval $[r_{i-1} + 1 \dots \ell_i - 1]$ is empty, and the loop will advance $[\ell_i \dots r_i]$ up to the first interval in A_i with a left extreme larger than r_{i-1} , making again the invariant true.

Suppose now that there are $[\bar{\ell}_0 \dots \bar{r}_0], [\bar{\ell}_1 \dots \bar{r}_1], \dots, [\bar{\ell}_k \dots \bar{r}_k]$ satisfying $r_i + 1 = \ell_{i+1}$ for some $k > 0$ and $0 \leq i < k$. We prove by induction on k that at some point during the execution of the algorithm we will be at the start of the external while loop with $i = k$ and $[\ell_j \dots r_j] = [\bar{\ell}_j \dots \bar{r}_j]$ for $j = 0, 1, \dots, k$. The thesis is trivially true for $k = 0$. Assume the thesis for $k - 1$, so we are at the start of the external while loop with $i = k - 1$ and $\ell_j = \bar{\ell}_j, r_j = \bar{r}_j$ for $j = 0, 1, \dots, k - 1$. Because of the invariant, either $[\ell_k \dots r_k] = [\bar{\ell}_k \dots \bar{r}_k]$ or $[\ell_k \dots r_k]$ will be advanced by the execution of the internal while loop up to $[\bar{\ell}_k \dots \bar{r}_k]$. Thus, at the end of the external while loop the thesis will be true for k . We conclude that all concatenations of intervals from A_0, A_1, \dots, A_{m-1} are returned.

We note that all intervals returned are unique (minimality has been already discussed in Section 3), as $[\ell_0 \dots r_0]$ is advanced at each call, so a duplicate returned interval would imply the existence of two

comparable intervals in A_0 . ■

Theorem 8 Algorithm 4 for BLOCK is minimally and optimally lazy.

Proof. The algorithm is trivially 0-lazy, as all outputs are uniquely determined by a tuple of intervals from the inputs. An algorithm \mathcal{A}^* advancing an input list A_i less than Algorithm 4 for some output $[\ell \dots r]$ would emit $[\ell \dots r]$ even if we truncated A_i after the last interval read by \mathcal{A}^* . ■

Practical remarks. In the case of intervals of integers, the advancement of the first list at the end of the outer loop can actually be iterated until $r_0 \geq \ell_i - i$. This change does not affect the complexity of the algorithm, but it may reduce number of iterations of the outer loop. In case the input antichains are entirely formed by singletons⁸, a folklore algorithm aligns the singletons circularly rather than starting from the first one (since they are singletons, once the position of an interval is fixed all the remaining ones are, too). The main advantage is that of avoiding to resolve several alignments if the first few terms appear often consecutively, but not followed by the remaining ones.

7.2 The $\text{AND}_{<}$ operator

The algorithm for computing this operator is a medley of the algorithms for AND and for BLOCK: as in the case of AND, we must check that future intervals are not smaller than our current candidate $[\ell' \dots r']$; as in the case of BLOCK, there is no queue and the lists A_0, A_1, \dots, A_{m-1} are advanced greedily. Again, we keep track of a current interval $[\ell_i \dots r_i]$ for every list A_i ; initially, these intervals are $[-\infty \dots -\infty]$, except for the first one, which is taken from the first list. The algorithm is described in pseudocode in Algorithm 5; an informal description follows.

The core of the algorithm is in the loop starting at line 8: this loop tries to align the i -th interval, that is, advance it until $[\ell_{i-1} \dots r_{i-1}] \ll [\ell_i \dots r_i]$. The loop starting at line 5 aims at aligning all intervals; note that we assume as an invariant that, after the first execution, every time we discover that the i -th interval is already aligned we can conclude that also the remaining intervals (the ones with index larger than i) are aligned as well (second condition at line 7).

The loop at line 5 can be interrupted as soon as, trying to align the i -th interval, we exhaust the i -th list or find an interval whose right extremes exceeds b , the left extreme of the $(m-1)$ -th interval forming the current candidate alignment. If any such condition is satisfied, the current candidate is certainly minimal and can thus be returned.

Upon a successful alignment (line 14), we have a new candidate: note that either this is the first candidate (i.e., $[\ell' \dots r'] = [\infty \dots \infty]$ before the assignment), or its right extreme coincides with the one of the previous candidate (i.e., $r' = r_{m-1}$ before the assignment), whereas its left extreme is certainly strictly larger. In either case, we try to see if we can advance the first interval and find a new, smaller candidate with a new alignment: this should explain the outer loop.

Theorem 9 Algorithm 5 for $\text{AND}_{<}$ is correct.

Proof. Let us say that a sequence $[\ell'_h \dots r'_h] \ll [\ell'_{h+1} \dots r'_{h+1}] \ll \dots \ll [\ell'_{k-1} \dots r'_{k-1}]$ of intervals ($h < k \leq m$), one from each list $A_h, A_{h+1}, \dots, A_{k-1}$, is *leftmost* if, for all $h < j < k$, there are no intervals in A_j with left extreme in $(r'_{j-1} \dots \ell'_j)$: such a sequence is uniquely determined by k and by $[\ell'_h \dots r'_h]$. Let $[\bar{\ell} \dots \bar{r}]$ be the interval returned at the last call (initially, $[\bar{\ell} \dots \bar{r}] = [\infty \dots \infty]$). Then, the following invariant holds at the start of the loop at line 5:

1. $[\ell_0 \dots r_0] \ll [\ell_1 \dots r_1] \ll \dots \ll [\ell_{i-1} \dots r_{i-1}]$ is leftmost;

⁸We emphasize this case because this is what happens with phrasal queries all of whose subqueries are simple terms; implementation may treat this special case differently to obtain further optimization, for instance using *ad hoc* indices [17].

Algorithm 5 The algorithm for the $\text{AND}_{<}$ operator. For sake of simplicity, we use the convention that returning $[\infty \dots \infty]$ means returning **null**, and that if one of the input lists is exhausted the function returns **null**.

```

0  Initially  $[\ell_0 \dots r_0] \leftarrow \text{next}(A_0)$ ,  $[\ell_k \dots r_k] \leftarrow [-\infty \dots -\infty]$  for all  $0 < k < m$  and  $i \leftarrow 1$ .
1  function next begin
2     $[\ell' \dots r'] \leftarrow [\infty \dots \infty]$ ;
3     $b \leftarrow \infty$ ;
4    forever
5      forever
6        if  $r_{i-1} \geq b$  then return  $[\ell' \dots r']$ ;
7        if  $i = m$  or  $\ell_i > r_{i-1}$  then break;
8        do
9          if  $r_i \geq b$  or  $A_i$  is empty then return  $[\ell' \dots r']$ ;
10          $[\ell_i \dots r_i] \leftarrow \text{next}(A_i)$ 
11         while  $\ell_i \leq r_{i-1}$ ;
12          $i \leftarrow i + 1$ ;
13       end;
14      $[\ell' \dots r'] \leftarrow [\ell_0 \dots r_{m-1}]$ ;
15      $b \leftarrow \ell_{m-1}$ ;
16      $i \leftarrow 1$ ;
17     if  $A_0$  is empty then return  $[\ell' \dots r']$ ;
18      $[\ell_0 \dots r_0] \leftarrow \text{next}(A_0)$ 
19   end;
20 end;

```

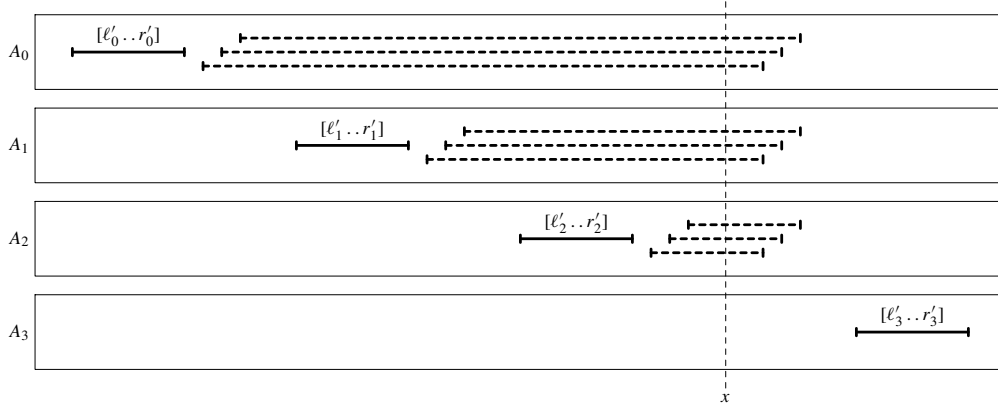


Figure 3: A sample configuration found in the proof of Theorem 10. In this case, $m = 4$ and $k = 1$. The dashed intervals are those of the form $[u_i^j \dots v^j]$: while reading such intervals it is impossible to decide whether the continuous intervals span an element of the output.

2. if $\ell_i \neq -\infty$ also $[\ell_i \dots r_i] \ll [\ell_{i+1} \dots r_{i+1}] \ll \dots \ll [\ell_{m-1} \dots r_{m-1}]$ is leftmost;
3. if $[\ell_{i-1} \dots r_{i-1}] \ll [\ell_i \dots r_i]$ then this pair is leftmost.

The fact that this invariant holds is easy to check; in particular, see the inner while loop at line 8 and the exit at line 7.

We now show that each output interval $[\bar{\ell} \dots \bar{r}]$ is at some time assigned to $[\ell' \dots r']$. Note that $i > 0$ at all times, so $[\ell_0 \dots r_0]$ is assigned only at the end of the infinite loop. This means that $[\ell_0 \dots r_0]$ runs through the whole first input list.

Thus, as soon as $\ell_0 = \bar{\ell}$ the inner loop will either compute the leftmost representation of $[\bar{\ell} \dots \bar{r}]$, or exit prematurely. In the second case, the function will necessarily complete the leftmost representation at the next call. We conclude that leftmost representations of all output intervals are assigned to $[\ell' \dots r']$ eventually: since $[\bar{\ell} \dots \bar{r}]$ is minimal, it will be emitted before $[\ell' \dots r']$ is assigned again. Uniqueness follows by uniqueness of leftmost representations. ■

It is not difficult to see that there is *no algorithm* for $\text{AND}_{<}$ that is k -lazy for any k , except for the case $m = 2$; indeed:

Theorem 10 If $m > 2$, there exist no optimally lazy algorithm for $\text{AND}_{<}$.

Proof. By contradiction, let \mathcal{B} be k -lazy, and observe that, on any given input I , every algorithm for $\text{AND}_{<}$, before emitting its p -th output $[\ell \dots r]$, must have reached at least the leftmost sequence $[\ell'_0 \dots r'_0] \ll [\ell'_1 \dots r'_1] \ll \dots \ll [\ell'_{m-1} \dots r'_{m-1}]$ spanning it. Now, choose any $x \in (r'_{m-2} \dots \ell'_{m-1})$ and, for all $i = 0, \dots, m-2$, take an arbitrary sequence $u_i^0 < u_i^1 < u_i^2 < \dots < u_i^{k+1} \in (r'_i \dots \min\{\ell'_{i+1}, x\})$; also choose an arbitrary sequence $v^0 < v^1 < v^2 < \dots < v^{k+1} \in (x \dots \ell'_{m-1})$. Run \mathcal{B} on a different input J , obtained as follows: whenever \mathcal{B} asks for an input from list $i < m-1$, we use the original intervals from I only up to $[\ell'_i \dots r'_i]$, and then we do the following: if $i < m-1$, we start offering $[u_i^0 \dots v^0]$, $[u_i^1 \dots v^1]$ and so on; as far as the last input list is concerned, we do not make any change. An example of this construction is given in Figure 3.

Note that the intervals are chosen so that \mathcal{B} cannot yet emit $[\ell \dots r]$, because there is always some chance for it not to be minimal. We stop testing \mathcal{B} as soon as, for some \bar{i} , \mathcal{B} has read at least $k+2$ inputs after $[\ell'_{\bar{i}} \dots r'_{\bar{i}}]$ from the \bar{i} -th list for some $\bar{i} < m-1$; let J' be the portion of J read by

\mathcal{B} so far, let \bar{j} any index different from \bar{i} and from $m - 1$, and let \mathcal{A} be an algorithm for $\text{AND}_{<}$ obtained from \mathcal{B} by modifying its behavior on the input as follows: when faced with an input that coincides with I up to $[\ell'_0 \dots r'_0] \ll [\ell'_1 \dots r'_1] \ll \dots \ll [\ell'_{m-1} \dots r'_{m-1}]$ inclusive, it then reads one more interval for each list and, if all these intervals contain any common point, say z , it starts reading from list \bar{j} until an interval not including z is reached, or until the \bar{j} -th list ends, in which case it emits $[\ell \dots r]$. Note that this modification does not harm the correctness of the algorithm, but now $\rho_i^{\mathcal{A}}(J', p) + k + 1 = \rho_i^{\mathcal{B}}(J', p)$ which contradicts the k -laziness of \mathcal{B} . ■

Hence, for $\text{AND}_{<}$, there is no hope for our algorithm to be optimally lazy in the general case; yet, it enjoys three interesting properties:

Theorem 11 Let \mathcal{A} be Algorithm 5 for $\text{AND}_{<}$.

1. \mathcal{A} is minimally lazy;
2. \mathcal{A} is minimally and optimally lazy when $m = 2$;
3. for any functionally equivalent algorithm \mathcal{B} , $\rho_i^{\mathcal{A}}(I, p) \leq \rho_i^{\mathcal{B}}(I, p + 1)$; that is, our algorithm, to produce any output, never reads more input than \mathcal{B} needs to produce its next output.

Proof. (1) Suppose that \mathcal{B} is functionally equivalent to \mathcal{A} and $\rho_j^{\mathcal{B}}(I, p) \leq \rho_j^{\mathcal{A}}(I, p)$ for every j , I and p , and $\rho_{\bar{j}}^{\mathcal{B}}(\bar{I}, \bar{p}) < \rho_{\bar{j}}^{\mathcal{A}}(\bar{I}, \bar{p})$ for some specific \bar{j} , \bar{I} and \bar{p} . Let $[\ell \dots r]$ be the \bar{p} -th output on input \bar{I} , and $[\ell'_0 \dots r'_0] \ll \dots \ll [\ell'_{m-1} \dots r'_{m-1}]$ be its leftmost spanning sequence (Figure 4 displays an example); when \mathcal{A} outputs $[\ell \dots r]$, we have that $[\ell_j \dots r_j] = [\ell'_j \dots r'_j]$ for all $j > i$, whereas the i -th list is over or is such that $r_i \geq \ell'_{m-1}$ (with leftmost r_i), $[\ell_0 \dots r_0] \ll \dots \ll [\ell_{i-1} \dots r_{i-1}]$ is leftmost and $\ell < \ell_0$. Since no correct algorithm can emit $[\ell \dots r]$ before scanning its input up to the leftmost spanning sequence, necessarily $\bar{j} \leq i$.

Moreover, necessarily $\bar{j} \neq i$: otherwise, we could modify the inputs by substituting the unread intervals of the lists $A_i, A_{i+1}, \dots, A_{m-2}$ with a suitable sequence of aligned intervals which, together with the remaining ones, would span $[\ell_0 \dots r]$; this would make $[\ell \dots r]$ non minimal.

Now, suppose that J is an input equal to \bar{I} but modified so that the \bar{j} -th list ends immediately after the last interval read by \mathcal{B} : on input J , algorithm \mathcal{A} does not read a single interval from list i beyond $[\ell'_i \dots r'_i]$, because it emits $[\ell \dots r]$ as soon as the test for emptiness of $A_{\bar{j}}$ is performed. So $\rho_i^{\mathcal{A}}(J, \bar{p}) < \rho_i^{\mathcal{B}}(J, \bar{p})$, a contradiction.

(2) We prove that \mathcal{A} is 0-lazy in that case. Indeed, when a certain output $[\ell \dots r]$ is ready to be produced, \mathcal{A} tries to read one more interval $[\ell_0 \dots r_0]$ from the first list, and this is unavoidable (any other algorithm must do this, or otherwise we might modify the next interval so that $[\ell \dots r]$ is not minimal). This interval has a right extreme larger than or equal to ℓ_1 , or otherwise $[\ell \dots r]$ would not be minimal: \mathcal{A} exits at this point, so it is 0-lazy.

(3) This is trivial: when \mathcal{A} outputs an interval, it has not yet reached (or, it has just reached) the leftmost sequence spanning the following output, and no correct algorithm could ever emit the next output before that point. ■

Practical remarks. In the case of intervals of integers, the check for $r_i \geq b$ can be replaced by $r_i \geq b - (m - i - 2)$, and the check for $r_{i-1} \geq b$ by $r_i \geq b - (m - i - 1)$, obtaining in some case faster detection of minimality. If the input antichains are entirely formed by singletons, the check $r_i \geq b$ can be removed altogether, as in that case we know that $r_i = \ell_i \leq r_{i-1} < b$.

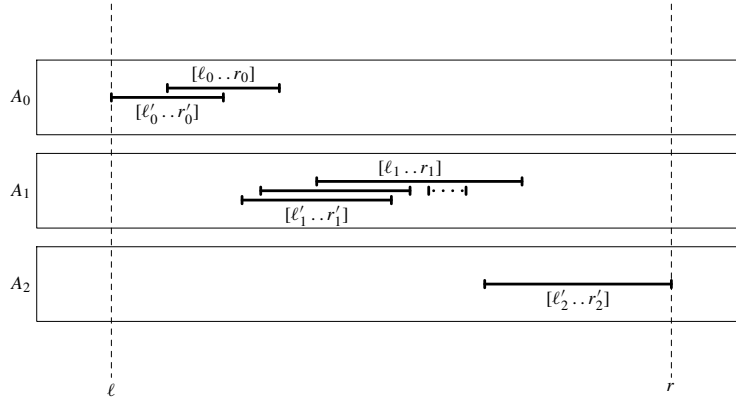


Figure 4: A sample configuration found in the proof of Theorem 11. The intervals $[\ell'_i \dots r_i]$ form a leftmost spanning sequence, and $i = 1$, so $\bar{j} = 0$. Note that no algorithm can avoid reading $[\ell_1 \dots r_1]$, or it would fail if we replaced it with the dotted interval.

7.3 Brouwerian difference

The Brouwerian difference $M - S$ between antichains M (the minuend) and S (the subtrahend) can be computed by searching greedily, for each interval $[\ell \dots r]$ in M , the first interval $[\ell' \dots r']$ in S for which $\ell' \geq \ell$ or $r' \geq r$. We keep track of the last interval $[\ell' \dots r']$ read from the input list S (initially, $[\ell' \dots r'] = [-\infty \dots -\infty]$) and update it until $\ell' \geq \ell$ or $r' \geq r$. At that point, if we did not exhaust S and $[\ell' \dots r'] \subseteq [\ell \dots r]$ (in which case $[\ell \dots r]$ should not be output) we continue scanning M ; otherwise, we return $[\ell \dots r]$. The algorithm is described in pseudocode in Algorithm 5.

Algorithm 6 The algorithm for Brouwerian difference.

```

0  Initially  $[\ell' \dots r'] \leftarrow [-\infty \dots -\infty]$ .
1  function next begin
2    while  $M$  is not empty do
3       $[\ell \dots r] \leftarrow \text{next}(M)$ ;
4      while  $\ell' < \ell$  and  $r' < r$  and  $S$  is not empty do
5         $[\ell' \dots r'] \leftarrow \text{next}(S)$ 
6      end;
7      if  $S$  is empty or  $[\ell' \dots r'] \not\subseteq [\ell \dots r]$  then return  $[\ell \dots r]$ 
8      end;
9      return null
10 end;

```

Theorem 12 Algorithm 6 for Brouwerian difference is correct.

Proof. Note that at the start of the inner while loop (line 4) $[\ell' \dots r']$ contains either the leftmost interval of S such that $\ell' \geq \ell$ or $r' \geq r$, or some interval preceding it. This is certainly true at the first call, and remains true after the execution of the inner while loop because of the first part of its exit condition (line 4). Finally, advancing the list of M cannot make the invariant false.

Given the invariant, at the end of the inner loop $[\ell' . . r']$ contains the leftmost interval of S such that $\ell' \geq \ell$ or $r' \geq r$, if such an interval exists. Note that if $[\ell' . . r']$ is not contained in $[\ell . . r]$, then no other interval of S is. Indeed, if $\ell' < \ell$ this means that $r' \geq r$, so all preceding intervals have too small left extremes, and all following intervals have too large right extremes (the same happens *a fortiori* if $\ell' \geq \ell$). Thus, the test at line 7 will emit $[\ell . . r]$ if and only if it belongs to the output. ■

Theorem 13 Algorithm 6 for Brouwerian difference is minimally and optimally lazy.

Proof. The algorithm (let us call it \mathcal{A}) is trivially 0-lazy: when $[\ell . . r]$ is output, \mathcal{A} has read just $[\ell . . r]$ from M and the first element $[\ell' . . r']$ of S such that $\ell' \geq \ell$ or $r' \geq r$. If either interval has not been read by some other algorithm \mathcal{A}^* , \mathcal{A}^* would fail if we removed altogether $[\ell . . r]$ from M or if we substituted $[\ell' . . r']$ with $[\ell . . r]$ and deleted all following intervals in S . ■

8 Previous work

The only attempt at linear lazy algorithms for minimal-interval region algebras we are aware of is the work of Young–Lai and Tompa on *structure selection queries* [18], a special type of expressions built on the primitives “contained-in”, “overlaps”, and so on, that can be evaluated lazily in linear time. Their motivations are similar to ours — application of region algebras to very large text collections. Similarly, Navarro and Baeza–Yates [13] propose a class of algorithms that using tree-traversals are able to compute efficiently several operations on overlapping regions. Their motivations are efficient implementation of structured query languages that permit such regions. Albeit similar in spirit, they do not provide algorithms for any of the operators we consider, and they do not provide a formal proof of laziness.

The manipulation of antichain of intervals can be translated into manipulation of points in the plane compared by *dominance* — coordinatewise ordering. Indeed, $[\ell . . r] \supseteq [\ell' . . r']$ iff the point $(\ell, -r)$ is dominated by the point $(\ell', -r')$. Dominance problems have been studied for a long time in computational geometry: for instance, [12] presents an algorithm to compute the maximal elements w.r.t. dominance. This method can be turned into an algorithm for antichains of intervals by coupling it with a simple (right-extreme based) merge to produce an algorithm for the OR operator. One has just to notice that since dominance is symmetric in the extremes, the mapping $[\ell . . r] \mapsto (-r, \ell)$ turns minimal intervals (by containment) into maximal points (by dominance). The algorithm described in [12] assume a decreasing first-coordinate order of the points, which however is an *increasing* ordering by right extreme on the original intervals. After some cleanup, the algorithm turns out to be identical to our algorithm for OR (albeit the authors do not study its laziness).

The other operators have no significant geometric meaning, and to the best of our knowledge there is no algorithm in computational geometry that computes them.

Lazy evaluation is a by-now classical topic in the theory of computation, dating back to the mid-70s [8], originally introduced for expressing the semantics of call-by-need in functional languages. However, the notion of lazy optimality used in this paper is new, and we believe that it captures as precisely as possible the idea of optimality in accessing sequentially multiple lists of inputs in a lazy fashion.

9 Conclusions

We have provided efficient algorithms for the computation of several operators on the lattice of interval antichains. The algorithms for lattice operations require time $O(n \log m)$ for m input antichains containing n intervals overall, whereas the remaining algorithms are linear in n . In particular, the

algorithm for OR has been proved to be optimal in a comparison-based model. Moreover, the algorithms are minimally and optimally lazy (with the exception of $\text{AND}_<$ when $m > 2$, in which case we prove an impossibility result) and use space linear in the number of input antichains. Our algorithms compare favourably with previously known techniques [5], which in particular required random access to the inputs.

An interesting open problem is that of providing a matching lower bound for the AND operator, at least for a comparison-based computational model.

References

- [1] Ziv Bar-Yossef, Ravi Kumar, and D. Sivakumar. Reductions in streaming algorithms, with an application to counting triangles in graphs. In *SODA '02: Proceedings of the thirteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 623–632, Philadelphia, PA, USA, 2002. Society for Industrial and Applied Mathematics.
- [2] Garrett Birkhoff. *Lattice Theory*, volume XXV of *AMS Colloquium Publications*. American Mathematical Society, third (new) edition, 1970.
- [3] Paolo Boldi and Sebastiano Vigna. Efficient lazy algorithms for minimal-interval semantics. In Fabio Crestani, Paolo Ferragina, and Mark Sanderson, editors, *Proc. SPIRE 2006*, number 4209 in *Lecture Notes in Computer Science*, pages 134–149. Springer–Verlag, 2006.
- [4] Charles L. A. Clarke and Gordon V. Cormack. Shortest-substring retrieval and ranking. *ACM Trans. Inf. Syst.*, 18(1):44–78, 2000.
- [5] Charles L. A. Clarke, Gordon V. Cormack, and Forbes J. Burkowski. An algebra for structured text search and a framework for its implementation. *Comput. J.*, 38(1):43–56, 1995.
- [6] J. Crampton and G. Loizou. The completion of a poset in a lattice of antichains. *International Mathematical Journal*, 1(3):223–238, 2001.
- [7] G. H. Gonnet. PAT 3.1: An efficient text searching system. User’s manual. Technical report, Center for the New Oxford English Dictionary. University of Waterloo, Waterloo, Canada, 1987.
- [8] Peter Henderson and Jr. James H. Morris. A lazy evaluator. In *POPL '76: Proceedings of the 3rd ACM SIGACT-SIGPLAN symposium on Principles on programming languages*, pages 95–103, New York, NY, USA, 1976. ACM Press.
- [9] Monika Henzinger, Prabhakar Raghavan, and Sridhar Rajagopalan. Computing on data streams, 1998.
- [10] Jani Jaakkola and Pekka Kilpeläinen. Nested text-region algebra. Technical Report C-1999-2, Department of Computer Science, University of Helsinki, 1999.
- [11] Guy-Vincent Jourdan, Jean-Xavier Rampon, and Claude Jard. Computing on-line the lattice of maximal antichains of posets. *Order*, 11(3):197–210, 1994.
- [12] H. T. Kung, Fabrizio Luccio, and Franco P. Preparata. On finding the maxima of a set of vectors. *J. Assoc. Comput. Mach.*, 22(4):469–476, 1975.
- [13] Gonzalo Navarro and Ricardo Baeza-Yates. A class of linear algorithms to process sets of segments. In *Proc. CLEI'96*, volume 2, pages 671–682, 1996.
- [14] Jürg Nievergelt and Franco P. Preparata. Plane-sweep algorithms for intersecting geometric figures. *Comm. ACM*, 25(10):739–747, 1982.

- [15] The Lemur Project. Indri. <http://www.lemurproject.org/indri/>.
- [16] Kunihiro Sadakane and Hiroshi Imai. Fast algorithms for k -word proximity search. *IEICE Trans. Fundamentals*, E84-A(9), September 2001.
- [17] Hugh E. Williams, Justin Zobel, and Dirk Bahle. Fast phrase querying with combined indexes. *ACM Trans. Inf. Syst.*, 22(4):573–594, 2004.
- [18] Matthew Young-Lai and Frank Wm. Tompa. One-pass evaluation of region algebra expressions. *Inf. Syst.*, 28(3):159–168, 2003.
- [19] Justin Zobel and Alistair Moffat. Inverted files for text search engines. *ACM Comput. Surv.*, 38(2), 2006.