

# Iterating on multiple collections in synchrony

STEFANO PERNA

Department of Computer Science,  
National University of Singapore  
(e-mail: dcsstef@nus.edu.sg)

VAL TANNEN

Department of Computer and Information Science,  
University of Pennsylvania  
(e-mail: val@cis.upenn.edu)

LIMSOON WONG

Department of Computer Science,  
National University of Singapore  
(e-mail: wongls@comp.nus.edu.sg)

---

## Abstract

Modern programming languages typically provide some form of comprehension syntax which renders programs manipulating collection types more readable and understandable. However, comprehension syntax corresponds to nested loops in general. There is no simple way of using it to express efficient general synchronized iterations on multiple ordered collections, such as linear-time algorithms for low-selectivity database joins. *Synchrony fold* is proposed here as a novel characterization of synchronized iteration. Central to this characterization is a *monotonic isBefore predicate* for relating the orderings on the two collections being iterated on, and an *antimonotonic canSee predicate* for identifying matching pairs in the two collections to synchronize and act on.

A restriction is then placed on Synchrony fold, cutting its extensional expressive power to match that of comprehension syntax, giving us *Synchrony generator*. Synchrony generator retains sufficient intensional expressive power for expressing efficient synchronized iteration on ordered collections. In particular, it is proved to be a natural generalization of the database merge join algorithm, extending the latter to more general database joins. Finally, *Synchrony iterator* is derived from Synchrony generator as a novel form of iterator. While Synchrony iterator has the same extensional and intensional expressive power as Synchrony generator, the former is better dovetailed with comprehension syntax. Thereby, algorithms requiring synchronized iterations on multiple ordered collections, including those for efficient general database joins become expressible naturally in comprehension syntax.

---

## 1 Introduction

Comprehension syntax, together with simple appeals to library functions, usually provides clear, understandable and short programs. Such a programming style for collection manipulation avoids loops and recursion as these are regarded as harder to understand and more error-prone. However, current collection-type function libraries appear lacking direct support that takes effective advantage of a *linear ordering* on collections for programming

47 in the comprehension style, even when such an ordering can often be made available by  
48 sorting the collections at a linearithmic overhead. We do not argue that these libraries lack  
49 expressive power *extensionally*, as the *functions* that interest us on ordered collections are  
50 easily expressible in the comprehension style. However, they are expressed *inefficiently*  
51 in such a style. We give a practical example in Section 2. We sketch in Section 3 proofs  
52 that, under a suitable formal definition of the restriction that gives the comprehension style,  
53 efficient algorithms for low-selectivity database joins, for example, cannot be expressed.  
54 Moreover, in this setting, these algorithms remain inexpressible even when access is given  
55 to any single library function such as `foldLeft`, `takeWhile`, `dropWhile`, and `zip`.

56 We proceed to fill this gap through several results on the design of a suitable collection-  
57 type function. We notice that most functions in these libraries are defined on one collection.  
58 There is no notion of any form of general synchronized traversal of two or more collections  
59 other than `zip`-like mechanical lock-step traversal. This seems like a design gap: synchron-  
60 ized traversals are often needed in real-life applications and, for an average programmer,  
61 efficient synchronized traversals can be hard to implement correctly.

62 Intuitively, a “synchronized traversal” of two collections is an iteration on two collec-  
63 tions where the “moves” on the two collections are coordinated, so that the current position  
64 in one collection is not too far from the current position in the other collection; i.e., from  
65 the current position in one collection, one “can see” the current position in the other collec-  
66 tion. However, defining the idea of “position” based on physical position, as in `zip`, seems  
67 restrictive. So, a more flexible notion of position is desirable. A natural and logical choice  
68 is that of a linear ordering relating items in the two collections; i.e. a linear ordering on the  
69 union of the two ordered collections. Also, given two collections which are sorted accord-  
70 ing to the linear orderings on their respective items, a reasonable new linear ordering on  
71 the union should respect the two linear orderings on the two original collections; i.e. given  
72 two items in an original collection where the first “is before” the second in the original  
73 collection, then the first should be before the second in the linear ordering defined on the  
74 union of the two collections.

75 Combining the two motivations above, our main approach to reducing the complexity  
76 of the expressed algorithms is to traverse two or more sorted collections in a *synchronized*  
77 manner, taking advantage of relationships between the linear orders on these collections.  
78 The following summarizes our results.

79 An addition to the design of collection-type function libraries is proposed in Section 4.  
80 It is called *Synchrony fold*. Some theoretical conditions, viz. monotonicity and antimono-  
81 tonicity, that characterize efficient synchronized iteration on a pair of ordered collections  
82 are presented. These conditions ensure the correct use of *Synchrony fold*. *Synchrony fold*  
83 is then shown to address the intensional expressive power gap articulated above.

84 *Synchrony fold* has the same extensional expressive power as `foldLeft`; it thus captures  
85 functions expressible by comprehension syntax augmented with typical collection-type  
86 function libraries. Because of this, *Synchrony fold* is not sufficiently precisely filling the  
87 intensional expressive power gap for comprehension syntax sans library function. A restric-  
88 tion to *Synchrony fold* is proposed in Section 5. This restricted form is called *Synchrony*  
89 *generator*. It has exactly the same extensional expressive power as comprehension syntax  
90 without any library function, but it has the intensional expressive power to express efficient  
91  
92

algorithms for low-selectivity database joins. Synchrony generator is further shown to correspond to a rather natural generalization of the database merge join algorithm (Blasgen & Eswaran, 1977; Mishra & Eich, 1992). The merge join was proposed half a century ago, and has remained as a backbone algorithm in modern database systems for processing equijoin and some limited form of non-equijoin (Silberschatz *et al.*, 2016), especially when the result has to be outputted in a specified order. Synchrony generator generalizes it to the class of non-equijoin whose join predicate satisfies certain antimonotonicity conditions.

Previous works have proposed alternative ways for compiling comprehension syntax, to enrich the repertoire of algorithms expressible in the comprehension style. For example, Wadler & Peyton Jones (2007) and Gibbons (2016), have enabled many relational database queries to be expressed efficiently under these refinements to comprehension syntax. However, these refinements only took equijoin into consideration; non-equijoin remains inefficient in comprehension syntax. In view of this and other issues, an iterator form—called *Synchrony iterator*—is derived from Synchrony generator in Section 6. While Synchrony iterator has the same extensional and intensional expressive power as Synchrony generator, it is more suitable for use in synergy with comprehension syntax. Specifically, Synchrony iterator makes efficient algorithms for simultaneous synchronized iteration on multiple ordered collections expressible in comprehension syntax.

Last but not least, Synchrony fold, Synchrony generator, and Synchrony iterator have an additional merit compared to other codes for synchronized traversal of multiple sorted collections. Specifically, they decompose such synchronized iterations into three orthogonal aspects, viz. relating the ordering on the two collections, identifying matching pairs, and acting on matching pairs. This orthogonality arguably makes for a more concise and precise understanding and specification of programs, hence improved reliability, as articulated by Schmidt (1986), Sebesta (2010) and Hunt & Thomas (2000).

## 2 Motivating example

Let us first define events as a data type.<sup>1</sup> An event has a *start* and an *end* point, where *start* < *end*, and typically has some additional attributes (e.g., an *id*) which do not concern us for now; cf. Figure 1. Events are ordered lexicographically by their *start* and *end* point: If an event *y* starts before an event *x*, then the event *y* is ordered before the event *x*; and when both events start together, the event which ends earlier is ordered before the event which ends later. Some predicates can be defined on events; e.g., in Figure 1, *isBefore*(*y*, *x*) says event *y* is ordered before event *x*, and *overlap*(*y*, *x*) says events *y* and *x* overlap each other.

Consider two collections of events, *xs*: *Vec*[*Event*] and *ys*: *Vec*[*Event*], where *Vec*[*.*] denotes a generic collection type, e.g., a vector.<sup>2</sup> The function *ov1*(*xs*, *ys*), defined in

<sup>1</sup> Scala (Odersky *et al.*, 2019) is used in this paper as the ambient language for a concrete discussion.

<sup>2</sup> In this paper, for convenience, *Vec*[*.*] is taken as the Scala *Vector*[*.*]. This allows us to assume postpends *:+* and *:++* are constant/linear time in their right argument, and prepends *+*: and *++*: are constant/linear time in their left argument. It is fine to take *Vec*[*.*] as *List*[*.*]; in this case, the postpends should be swapped by prepends, and some *reverse* has to be inserted into some of the codes. These list-specific details are not germane to understanding the key ideas of this paper. Hence, we adopt vectors as our generic collection type in general. Nonetheless, when we reach the final description of our last result, Synchrony iterator, at the end of Section 6.1, we will use concrete collection types, including an instance of list.

```

139 case class Event(start: Int, end: Int, id: String)
140 // Constraint: start < end
141
142 val isBefore = (y: Event, x: Event) => {
143   (y.start < x.start) ||
144   (y.start == x.start && y.end < x.end)
145 }
146
147 val overlap = (y: Event, x: Event) => {
148   (x.start < y.end && y.start < x.end)
149 }
150
151 def ov1(xs: Vec[Event], ys: Vec[Event]) = {
152   for (x <- xs; y <- ys; if overlap(y, x)) yield (x, y)
153 }
154
155 def ov2(xs: Vec[Event], ys: Vec[Event]) = {
156   // Requires: xs and ys sorted lexicographically by (start, end).
157   def aux(
158     xs: Vec[Event], ys: Vec[Event],
159     zs: Vec[Event], acc: Vec[(Event, Event)])
160   : Vec[(Event, Event)] =
161     // Key Invariant: aux(xs, ys, Vec(), acc) = acc ++ ov1(xs, ys)
162     if (xs.isEmpty) acc
163     else if (ys.isEmpty && zs.isEmpty) acc
164     else if (ys.isEmpty) aux(xs.tail, zs, Vec(), acc)
165     else {
166       val (x, y) = (xs.head, ys.head)
167       (isBefore(y, x), overlap(y, x)) match {
168         case (true, false) => aux(xs, ys.tail, zs, acc)
169         case (false, false) => aux(xs.tail, zs ++: ys, Vec(), acc)
170         case (_, true) => aux(xs, ys.tail, zs :+ y, acc :+ (x, y))
171       }
172     }
173   aux(xs, ys, Vec(), Vec())
174 }

```

Fig. 1. A motivating example. The functions  $ov1(xs, ys)$  and  $ov2(xs, ys)$  are equal on inputs  $xs$  and  $ys$  which are sorted lexicographically by their start and end point. While  $ov1(xs, ys)$  has quadratic time complexity  $O(|xs| \cdot |ys|)$ ,  $ov2(xs, ys)$  has time complexity  $O(|xs| + k|ys|)$  when each event in  $ys$  overlaps fewer than  $k$  events in  $xs$ .

Figure 1, retrieves the events in  $xs$  and  $ys$  that overlap each other. Although this comprehension syntax-based definition has the important virtue of being clear and succinct, it has quadratic time complexity  $O(|xs| \cdot |ys|)$ . An alternative implementation  $ov2(xs, ys)$  is given in Figure 1 as well. On  $xs$  and  $ys$  which are sorted lexicographically by  $(start, end)$ ,  $ov1(xs, ys) = ov2(xs, ys)$ . Notably, the time complexity of  $ov2(xs, ys)$  is  $O(|xs| + k|ys|)$ , provided each event in  $ys$  overlaps fewer than  $k$  events in  $xs$ . The proofs for these claims will become obvious later, from Theorem 4.4.

The function  $ov1(xs, ys)$  exemplifies a database join, and the join predicate is  $overlap(y, x)$ . Joins are ubiquitous in database queries. Sometimes, a join predicate is a conjunction of equality tests; this is called an equijoin. However, when a join predicate comprises entirely of inequality tests, it is called a non-equijoin;  $overlap(y, x)$  is a special form of non-equijoin which is sometimes called an interval join. Non-equijoin is quite common in practical applications. For example, given a database of taxpayers, a query

retrieving all pairs of taxpayers where the first earns more but pays less tax than the second is an interval join. As another example, given a database of mobile phones and their prices, a query retrieving all pairs of competing phone models (i.e. the two phone models in a pair are priced close to each other) is another special form of non-equijoin called a band join.

Returning to  $ov1(xs, ys)$  and  $ov2(xs, ys)$ , the upper bound  $k$  on the number of events in  $xs$  that an event in  $ys$  can overlap with is called the selectivity of the join.<sup>3</sup> When restricted to  $xs$  and  $ys$  which are sorted by their `start` and `end` point,  $ov1(xs, ys)$  and  $ov2(xs, ys)$  define the same function. However, their time complexity is completely different. The time complexity of  $ov1(xs, ys)$  is quadratic. On the other hand, the time complexity of  $ov2(xs, ys)$  is a continuum from linear to quadratic, depending on the selectivity  $k$ . In a real-life database query,  $k$  is often a very small number, relative to the number of entries being processed. So, in practice,  $ov2(xs, ys)$  is linear.

The definition  $ov1(xs, ys)$  has the advantage of being obviously correct, due to its being expressed using easy-to-understand comprehension syntax. Whereas,  $ov2(xs, ys)$  is likely to take even a skilled programmer much more effort to get right. This example is one of many functions having the following two characteristics. Firstly, these functions are easily expressible in a modern programming language using only comprehension syntax. However, this usually results in a quadratic or higher time complexity. Secondly, there are linear-time algorithms for these functions. Yet, there is no straightforward way to provide linear-time implementation for these functions using comprehension syntax without using more sophisticated features of the programming language and its collection-type libraries.

The proof for this intensional expressiveness gap in a simplified theoretical setting is outlined in the next section and is shown in full in a companion paper (Wong, 2021). It is the main objective of this paper to fill this gap as simply as possible.

### 3 Intensional expressiveness gap

As alluded to earlier, what we call *intensional expressive power* in this paper refers to the class of mappings from input to output that can be expressed, as in Fortune *et al.* (1983) and Felleisen (1991). In particular, so long as two programs in a language  $\mathcal{L}$  produce the same output given the same input, even when these two programs differ greatly in terms of time complexity, they are regarded as expressing (implementing) the same function  $f$ , and are thus equivalent and mutually substitutable.

However, we focus here on improving the ability to express algorithms, that is, on *intensional expressive power*. Specifically, as in many past works (Abiteboul & Vianu, 1991; Biskup *et al.*, 2004; Van den Bussche, 2001; Colson, 1991; Suciu & Paredaens, 1997; Suciu & Wong, 1995; Wong, 2013), we approach this in a coarse-grained manner by considering the time complexity of programs. In particular, an algorithm which implements a function  $f$  in  $\mathcal{L}$  is considered inexpressible in a specific setting if every program implementing  $f$  in  $\mathcal{L}$  under that setting has a time complexity higher than this algorithm.

Since Scala and other general programming languages are Turing complete, in order to capture the class of programs that we want to study with greater clarity, a restriction

<sup>3</sup> The results in this work remain valid when  $k$  is defined instead as the average number (rounded up to a whole number) of events in  $xs$  that an event in  $ys$  overlaps with.

needs to be imposed. Informally, user-programmers<sup>4</sup> are allowed to use comprehension syntax, collections, and tuples; but they are not allowed to use while-loops, recursion, and nested collections; and they are not allowed to define new data types and new higher-order functions (a higher-order function is a function whose result is another function or is a nested collection.) They are also not allowed to call functions in the collection-type function libraries of these programming languages, unless specifically permitted.

This is called the “first-order restriction.” Under this restriction, following Suciu & Wong (1995), when a user-programmer is allowed to use a higher-order function from a collection-type library, e.g., `foldLeft(e)(f)`, the function `f` which is user-defined can only be a first-order function. A way to think about this restriction is to treat higher-order library functions as a part of the syntax of the language, rather than as higher-order functions.

Under such a restriction, some functions may become inexpressible; and even when a function is expressible, its expression may correspond to a drastically inefficient algorithm (Biskup *et al.*, 2004; Suciu & Paredaens, 1997; Wong, 2013). In terms of extensional expressive power, the first-order restriction of our ambient language, Scala, is easily seen to be complete with respect to flat relational queries (Buneman *et al.*, 1995; Libkin & Wong, 1997), which are queries that a relational database system supports (such as joins). The situation is less clear from the perspective of intensional expressive power.

This section outlines results suggesting that Scala under the first-order restriction cannot express efficient algorithms for low-selectivity joins, and that this remains so even when a programmer is permitted to access some functions in Scala’s collection-type libraries. Formal proofs are provided in a companion paper (Wong, 2021).

To capture the first-order restriction on Scala, consider the nested relational calculus  $\mathcal{NRC}$  of Wong (1996).  $\mathcal{NRC}$  is a simply-typed lambda calculus with Boolean and tuple types and their usual associated operations; set types, and primitives for empty set, forming singleton sets, union of sets, and `flatMap` for iterating on sets;<sup>5</sup> and base types with equality tests and comparison tests. Replace its set type with linearly ordered set types, and assume that ordered sets, for computational complexity purposes, are traversed in order in linear time; this way, ordered sets can be thought of as lists. The replacement of set types by linearly ordered set types does not change the nature of  $\mathcal{NRC}$  in any drastic way, because  $\mathcal{NRC}$  can express a linear ordering on any arbitrarily deeply nested combinations of tuple and set types given any linear orderings on base types; cf. Libkin & Wong (1994). Next, restrict the language to its flat fragment; i.e., nested sets are not allowed. This restriction has no impact on the extensional expressive power of  $\mathcal{NRC}$  with respect to functions on non-nested sets, as shown by Wong (1996). Denote this language as  $\mathcal{NRC}_1(\leq)$ , where the permitted extra primitives are listed explicitly between the brackets.

<sup>4</sup> In this paper, we separate an implementer-programmer who implements programming constructs and library functions from a user-programmer who uses these. The former has access to all features of the programming language. The latter, in the context of this paper, is restricted to Scala under the first-order constraint plus specifically permitted library functions which the former provides. When we say a programmer, we refer to either programmer. So, in this paper, the implementer-programmer is the one implementing the proposed Synchrony fold, Synchrony generator, and Synchrony iterator. And the user-programmer is the one implementing the examples `ovi`, `ovCount`, `mtgi`, etc.

<sup>5</sup> `flatMap` is Scala’s terminology. It is also known as `bind` in the Haskell parlance.

Some terminologies are needed for stating the results. To begin, by an object, we mean the value of any combination of base types, tuples, and sets that is constructible in  $\mathcal{NRC}_1(\leq)$ .

A level-0 atom of an object  $C$  is a constant  $c$  which has at least one occurrence in  $C$  that is not inside any set in  $C$ . A level-1 atom of an object  $C$  is a constant  $c$  which has at least one occurrence in  $C$  that is inside a set. The notations  $atom^0(C)$ ,  $atom^1(C)$ , and  $atom^{\leq 1}(C)$  respectively denote the set of level-0 atoms of  $C$ , the set of level-1 atoms of  $C$ , and their union. The level-0 molecules of an object  $C$  are the sets in  $C$ . The notation  $molecule^0(C)$  denotes the set of level-0 molecules of  $C$ . E.g., suppose  $C = (c_1, c_2, \{(c_3, c_4), (c_5, c_6)\})$ ; then  $atom^0(C) = \{c_1, c_2\}$ ,  $atom^1(C) = \{c_3, c_4, c_5, c_6\}$ ,  $atom^{\leq 1}(C) = \{c_1, c_2, c_3, c_4, c_5, c_6\}$ , and  $molecule^0(C) = \{\{(c_3, c_4), (c_5, c_6)\}\}$ .

The level-0 Gaifman graph of an object  $C$  is defined as an undirected graph  $gaifman^0(C)$  whose nodes are the level-0 atoms of  $C$ , and edges are all the pairs of level-0 atoms of  $C$ . The level-1 Gaifman graph of an object  $C$  is defined as an undirected graph  $gaifman^1(C)$  whose nodes are the level-1 atoms of  $C$ , and the edges are defined as follow: If  $C = \{C_1, \dots, C_n\}$ , the edges are pairs  $(x, y)$  such that  $x$  and  $y$  are in the same  $atom^0(C_i)$  for some  $1 \leq i \leq n$ ; if  $C = (C_1, \dots, C_n)$ , the edges are pairs  $(x, y) \in gaifman^1(C_i)$  for some  $1 \leq i \leq n$ ; and there are no other edges. The Gaifman graph of an object  $C$  is defined as  $gaifman(C) = gaifman^0(C) \cup gaifman^1(C)$ ; cf. Gaifman (1982).

Let  $e(\vec{X})$  be an expression  $e$  whose free variables are  $\vec{X}$ . Let  $e[\vec{C}/\vec{X}]$  denote the closed expression obtained by replacing the free variables  $\vec{X}$  by the corresponding objects  $\vec{C}$ . Let  $e[\vec{C}/\vec{X}] \Downarrow C'$  mean the closed expression  $e[\vec{C}/\vec{X}]$  evaluates to the object  $C'$ ; the evaluation is performed according to a typical call-by-value operational semantics (Wong, 2021).

It is shown by Wong (2021) that  $\mathcal{NRC}_1(\leq)$  expressions can only manipulate their input in highly restricted local manners. In particular, expressions which have at most linear time complexity are able to mix level-0 atoms with level-0 and level-1 atoms, but are unable to mix level-1 atoms with themselves.

**Lemma 3.1** (Wong (2021), Lemma 3.1). *Let  $e(\vec{X})$  be an expression in  $\mathcal{NRC}_1(\leq)$ . Let objects  $\vec{C}$  have the same types as  $\vec{X}$ , and  $e[\vec{C}/\vec{X}] \Downarrow C'$ . Suppose  $e(\vec{X})$  has at most linear time complexity with respect to the size of  $\vec{X}$ .<sup>6</sup> Then for each  $(u, v) \in gaifman(C')$ , either  $(u, v) \in gaifman(\vec{C})$ , or  $u \in atom^0(\vec{C})$  and  $v \in atom^1(\vec{C})$ , or  $u \in atom^1(\vec{C})$  and  $v \in atom^0(\vec{C})$ .*

Here is a grossly simplified informal argument to provide some insight on this “limited-mixing” lemma. Consider an expression  $X \text{ f1atMap } f$ , where  $X$  is the variable representing the input collection and  $f$  is a function to be performed on each element of  $X$  in the usual manner of  $\text{f1atMap}$ . Then, the time complexity of this expression is  $O(n \cdot \hat{f})$ , where  $n$  is the number of items in  $X$  and  $O(\hat{f})$  is the time complexity of  $f$ . Clearly,  $O(n \cdot \hat{f})$  can be linear only when  $O(\hat{f}) = O(1)$ . Intuitively, this means  $f$  cannot have a subexpression of the form  $X \text{ f1atMap } g$ . Since  $\text{f1atMap}$  is the sole construct in  $\mathcal{NRC}_1(\leq)$  for accessing and manipulating the elements of a collection, when  $f$  is passed an element of  $X$ , there is no way for it to access a different element of  $X$  if  $f$  does not have a subexpression of the form  $X \text{ f1atMap } g$ . So, it is not possible for  $f$  to mix the components from two different elements

<sup>6</sup> In this work, all mentions of time complexity are with respect to input size.

of  $X$ . Unavoidably, many details are swept under the carpet in this informal argument, but are taken care of by Wong (2021).

This limited-mixing handicap remains when the language is further augmented with typical functions—such as `dropWhile`, `takeWhile`, and `foldLeft`—in collection-type libraries of modern programming languages. Even the presence of a fictitious operator, `sort`, for instantaneous sorting, cannot rescue the language from this handicap.

**Lemma 3.2** (Wong (2021), Lemma 5.1). *Let  $e(\vec{X})$  be an expression in  $\mathcal{NRC}_1(\leq, \text{takeWhile}, \text{dropWhile}, \text{sort})$ . Let objects  $\vec{C}$  have the same types as  $\vec{X}$ , and  $e(\vec{C}/\vec{X}) \Downarrow C'$ . Suppose  $e(\vec{X})$  has at most linear time complexity. Then there is a number  $k$  that depends only on  $e(\vec{X})$  but not on  $\vec{C}$ , and a set  $A \subseteq \text{atom}^{\leq 1}(\vec{C})$  where  $|A| \leq k$ , and for each  $(u, v) \in \text{gaifman}(C')$ , either  $(u, v) \in \text{gaifman}(\vec{C})$ , or  $u \in \text{atom}^0(\vec{C})$  and  $v \in \text{atom}^1(\vec{C})$ , or  $u \in \text{atom}^1(\vec{C})$  and  $v \in \text{atom}^0(\vec{C})$ , or  $u \in A$  and  $v \in \text{atom}^1(\vec{C})$ ,  $u \in \text{atom}^1(\vec{C})$  and  $v \in A$ .*

**Lemma 3.3** (Wong (2021), Lemma 5.4). *Let  $e(\vec{X})$  be an expression in  $\mathcal{NRC}_1(\leq, \text{foldLeft}, \text{sort})$ . Let objects  $\vec{C}$  have the same types as  $\vec{X}$ , and  $e(\vec{C}/\vec{X}) \Downarrow C'$ . Suppose  $e(\vec{X})$  has at most linear time complexity. Then there is a number  $k$  that depends only on  $e(\vec{X})$  but not on  $\vec{C}$ , and a set  $A \subseteq \text{atom}^{\leq 1}(\vec{C})$  where  $|A| \leq k$ , such that for each  $(u, v) \in \text{gaifman}(C')$ , either  $(u, v) \in \text{gaifman}(\vec{C})$ , or  $u \in \text{atom}^0(\vec{C})$  and  $v \in \text{atom}^1(\vec{C})$ , or  $u \in \text{atom}^1(\vec{C})$  and  $v \in \text{atom}^0(\vec{C})$ , or  $u \in A$  and  $v \in \text{atom}^1(\vec{C})$ , or  $v \in A$  and  $u \in \text{atom}^1(\vec{C})$ .*

The inexpressibility of efficient algorithms for low-selectivity joins in  $\mathcal{NRC}_1(\leq, \text{takeWhile}, \text{dropWhile}, \text{sort})$ , and  $\mathcal{NRC}_1(\leq, \text{foldLeft}, \text{sort})$  can be deduced from Lemmas 3.1, 3.2, and 3.3. The argument for  $\mathcal{NRC}_1(\leq, \text{foldLeft}, \text{sort})$  is provided here as an illustration. Let `zip(xs, ys)` be the query that pairs the  $i$ th element in `xs` with the  $i$ th element in `ys`, assuming that the two input collections `xs` and `ys` are sorted and have the same length. Without loss of generality, suppose the  $i$ th element in `xs` has the form  $(o_i, x_i)$  and that in `ys` has the form  $(o_i, y_i)$ ; suppose also that each  $o_i$  occurs only once in `xs` and once in `ys`, each  $x_i$  does not appear in `ys`, and each  $y_i$  does not appear in `xs`. Clearly, `zip(xs, ys)` is a low-selectivity join; in fact, its selectivity is precisely one. Let  $\text{xs} = \{(o_1, u_1), \dots, (o_n, u_n)\}$  and  $\text{ys} = \{(o_1, v_1), \dots, (o_n, v_n)\}$ . Let  $C' = \{(u_1, o_1, o_1, v_1), \dots, (u_n, o_n, o_n, v_n)\}$ . Then `zip(xs, ys) = C'`. Then  $\text{gaifman}(C') = \{(u_1, v_1), \dots, (u_n, v_n)\} \cup \Delta$ , where  $\Delta$  are the edges involving the  $o_j$ 's in  $\text{gaifman}(C')$ . Clearly, for  $1 \leq i \leq n$ ,  $(u_i, v_i) \in \text{gaifman}(C')$  but  $(u_i, v_i) \notin \text{gaifman}(\text{xs}, \text{ys}) = \text{xs} \cup \text{ys}$ . Now, for a contradiction, suppose  $\mathcal{NRC}_1(\leq, \text{foldLeft}, \text{sort})$  has a linear-time implementation for `zip`. Then, by Lemma 3.3, either  $u_i \in \text{atom}^0(\text{xs}, \text{ys})$ , or  $v_i \in \text{atom}^0(\text{xs}, \text{ys})$ , or  $u_i \in A$ , or  $v_i \in A$  for some  $A$  whose size is independent of `xs` and `ys`. However, `xs` and `ys` are both sets; thus,  $\text{atom}^0(\text{xs}, \text{ys}) = \{\}$ . This means  $A$  has to contain every  $u_i$  or  $v_i$ . So,  $|A| \geq n = |\text{xs}| = |\text{ys}|$  cannot be independent of `xs` and `ys`. This contradiction implies there is no linear-time implementation of `zip` in  $\mathcal{NRC}_1(\leq, \text{foldLeft}, \text{sort})$ .

A careful reader may realise that  $\mathcal{NRC}_1(\leq)$  does not have the `head` and `tail` primitives commonly provided for collection types in programming languages. However, the absence of `head` and `tail` in  $\mathcal{NRC}_1(\leq)$  is irrelevant in the context of this paper. To see this, consider these two functions: `taken(xs)` which returns in  $O(n)$  time the first  $n$  elements of `xs`, and `dropn(xs)` which drops in  $O(n)$  time the first  $n$  elements of `xs`, when `xs` is ordered.

So,  $\text{head}(xs) = \text{take}_1(xs)$  and  $\text{tail}(xs) = \text{drop}_1(xs)$ . The proof given by Wong (2021) for Lemma 3.2 can be copied almost verbatim to obtain an analogous limited-mixing result for  $\mathcal{NRC}_1(\leq, \text{take}_n, \text{drop}_n, \text{sort})$ .

Since `zip` is a manifestation of the intensional expressive power gap of  $\mathcal{NRC}_1(\leq)$  and its extensions above, one might try to augment the language with `zip` as a primitive. This makes it trivial to supply an efficient implementation of `zip`. Unfortunately, this does not escape the limited-mixing handicap either.

**Lemma 3.4** (Wong (2021), Lemma 5.7). *Let  $e(\vec{X})$  be an expression in  $\mathcal{NRC}_1(\leq, \text{zip}, \text{sort})$ . Let objects  $\vec{C}$  have the same types as  $\vec{X}$ , and  $e(\vec{C}/\vec{X}) \Downarrow C'$ . Suppose  $e(\vec{X})$  has at most linear time complexity. Then there is a number  $k$  that depends only on  $e(\vec{X})$  but not on  $\vec{C}$ , and an undirected graph  $K$  where the nodes are a subset of  $\text{atom}^{\leq 1}(\vec{C})$  and each node  $w$  of  $K$  has degree at most  $nk$ ,  $n$  is the number of times  $w$  appears in  $\vec{C}$ , such that for each  $(u, v) \in \text{gaifman}(C')$ , either  $(u, v) \in \text{gaifman}(\vec{C}) \cup K$ , or  $u \in \text{atom}^0(\vec{C})$  and  $v \in \text{atom}^1(\vec{C})$ , or  $u \in \text{atom}^1(\vec{C})$  and  $v \in \text{atom}^0(\vec{C})$ .*

It follows from Lemma 3.4 that there is no linear-time implementation of `ov1(xs, ys)` in  $\mathcal{NRC}_1(\leq, \text{zip}, \text{sort})$ . To see this, suppose for a contradiction that there is an expression  $f(xs, ys)$  in  $\mathcal{NRC}_1(\leq, \text{zip}, \text{sort})$  that implements `ov1(xs, ys)` with time complexity  $O(|xs| + h|ys|)$  when each event in `ys` overlaps fewer than  $h$  events in `xs`. Let  $k_0$  be the  $k$  induced by Lemma 3.4 on  $f$ . Suppose without loss of generality that no start and end points in `xs` appears in `ys`, and vice versa. Then setting  $h > k_0$  produces the desired contradiction.

$\mathcal{NRC}_1(\leq)$  is designed to express the same functions and algorithms that first-order restricted Scala is able to express. A bare-bone fragment of Scala that corresponds to  $\mathcal{NRC}_1(\leq)$  can be described as follows. In terms of data types: Base types such as `Boolean`, `Int`, and `String` are included. The operators on base types are restricted to `=` and `<=` tests. Other operators on base types (e.g., functions from base types to base types) can generally be included without affecting the limited-mixing lemmas. Tuple types over base types (i.e., all tuple components are base types) are included. The operators on tuple types are the tuple constructor and the tuple projection. A collection type is included, and the Scala `Vector[.]` is a convenient choice as a generic collection type; however, only vectors of base types and vectors of tuples of base types are included. The operators on vectors are the vector constructor, the `flatMap` on vectors, the vector append `++`, and the vector emptiness test; when restricted to these operators, vectors essentially behave as sets. It is also possible to use other Scala collection types—e.g., `List[.]`—instead of `Vector[.]`, so long as the operators are restricted to a constructor, append `++`, and emptiness test. Some other common operators on collection types, e.g., `head` and `tail`, can also be included, though adding these would make the language correspond to  $\mathcal{NRC}_1(\leq, \text{take}_1, \text{drop}_1)$  instead of  $\mathcal{NRC}_1(\leq)$  and, as explained earlier, this does not impact the limited-mixing lemmas. In terms of general programming constructs: Defining functions whose return types are any of the data types above (i.e. return types are not allowed to be function types), making function calls, and using comprehension syntax and `if-then-else` are all permitted. Although pared to such a bare bone, this highly restricted form of Scala retains sufficient expressive power; e.g., all flat relational queries can be easily expressed using it.

Thus, Lemma 3.1 implies there is no efficient implementation of low-selectivity joins, including `ov1(xs, ys)`, in first-order restricted Scala. Lemma 3.2 implies there is no efficient implementation of low-selectivity joins in first-order restricted Scala even when the programmer is given access to `takeWhile` and `dropWhile`. Lemma 3.3 implies there is no efficient implementation of low-selectivity joins in first-order restricted Scala even when the programmer is given access to `foldLeft`. Lemma 3.4 implies there is no efficient implementation of low-selectivity joins in first-order restricted Scala even when the programmer is given access to `zip`. Moreover, these limitations remain even when the programmer is further given the magical ability to do sorting infinitely fast.

## 4 Synchrony fold

Comprehension syntax is typically translated into nested `flatMap`'s, each `flatMap` iterating independently on a single collection. Consequently, like any function defined using comprehension syntax, the function `ov1(xs, ys)` in Figure 1 is forced to use nested loops to process its input. While it is able to return correct results even for an unsorted input, it overkills and overpays a price in its quadratic time complexity when its input is already appropriately sorted. In fact, `ov1(xs, ys)` is still overpaying the quadratic-time complexity price when its input is unsorted, because sorting can always be performed when needed for a relatively affordable linearithmic overhead.

In contrast, the function `ov2(xs, ys)` in Figure 1 is linear in time complexity when selectivity is low, which is much more efficient than `ov1(xs, ys)`. There is one fundamental explanation for this efficiency: The input `xs` and `ys` are sorted and `ov2(xs, ys)` directly exploits this sortedness to iterate on `xs` and `ys` in “synchrony,” i.e. in a coordinated manner, akin to the merge step in a merge sort (Knuth, 1973) or a merge join (Blasgen & Eswaran, 1977; Mishra & Eich, 1992). However, its codes are harder to understand and to get right.

It is desirable to have an easy-to-understand-and-check linear-time implementation that is as efficient as `ov2(xs, ys)` but using only comprehension syntax, without the acrobatics of recursive functions, while-loops, etc. This leads us to the concepts of *Synchrony fold*, *Synchrony generator*, and *Synchrony iterator*. Synchrony fold is presented in this section. Synchrony generator and iterator are presented later in Sections 5 and 6 respectively.

### 4.1 Theory of Synchrony fold

The function `ov2(xs, ys)` exploits the sortedness and the relationship between the orderings of `xs` and `ys`. In Scala's collection-type function libraries, functions such as `foldLeft` are also able to exploit the sortedness of their input. Yet there is no way of individually using `foldLeft` and other collection-type library functions mentioned earlier—as suggested by Lemma 3.2, Lemma 3.3, and Lemma 3.4—to obtain linear-time implementation of low-selectivity joins, without defining recursive functions, while-loops, etc. The main reason is that these library functions are mostly defined on a single input collection. Hence, it is hard for them to exploit the relationship between the orderings on two collections. And there is no obvious way to process two collections using any one of these library functions alone, other than in a nested-loop manner, unless the ambient programming language has more

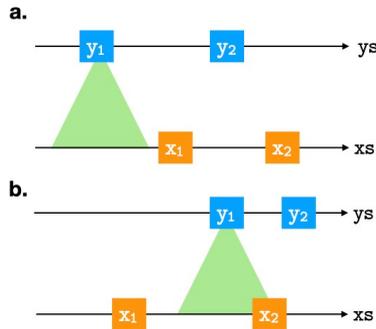


Fig. 2. Visualization of monotonicity and antimonotonicity. Two collections  $xs$  and  $ys$  are sorted according to some orderings, as denoted by the two arrows. The `isBefore` predicate is represented by the relative horizontal positions of items  $x_i$  and  $y_j$ ; i.e., if  $y_j$  is before  $x_i$ , then  $y_j$  is before  $x_i$ . The `canSee` predicate is represented by the shaded green areas. **a.** If  $y_1$  is before  $x_1$  and cannot see  $x_1$ , then  $y_1$  is also before and cannot see any  $x_2$  which comes after  $x_1$ . So, every  $x_i$  that matches  $y_1$  has been seen; it is safe to move forward to  $y_2$ . **b.** If  $y_1$  is not before  $x_1$  and cannot see  $x_1$ , then any  $y_2$  which comes after  $y_1$  is also not before and cannot see  $x_1$ . So, every  $y_j$  that matches  $x_1$  has been seen; it is safe to move forward to  $x_2$ .

sophisticated ways to compile comprehensions (Marlow *et al.*, 2016; Wadler & Peyton Jones, 2007), or unless multiple library functions are used together.

Scala’s collection-type libraries do provide the function `zip` which pairs up elements of two collections according to their physical position in the two collections, viz. first with first, second with second, and so on. However, by Lemma 3.4, this mechanical pairing by `zip` cannot be used to implement efficient low-selectivity joins, which require more general notions of pairing where pairs can form from different positions in the two collections.

So, we propose `syncFold`, a generalization of `foldLeft` that iterates on two collections in a more flexible and synchronized manner. For this, we need to relate positions in two collections by introducing two logical predicates `isBefore(y, x)` and `canSee(y, x)`, which are supplied to `syncFold` as two of its arguments. Informally, `isBefore(y, x)` means that, when we are iterating on two collections  $xs$  and  $ys$ , in a synchronized manner, we should encounter the item  $y$  in  $ys$  before we encounter the item  $x$  in  $xs$ . And `canSee(y, x)` means that the item  $y$  in  $ys$  corresponds to or matches the item  $x$  in  $xs$ ; in other words,  $x$  and  $y$  form a pair which is of interest. Note that an item  $y$  “corresponds to or matches” an item  $x$  does not necessarily mean the two items are the same. For example, when items are events as defined in Section 2, in the context of `ov1` and `ov2`, an event  $y$  corresponds to or matches an event  $x$  means the two events overlap each other. Obviously, an item does not need to be an atomic object; it can be a tuple or an object having a more complex type.

The `isBefore(y, z)` and `canSee(y, x)` predicates are characterized respectively by the monotonicity and antimonotonicity conditions defined below and depicted in Figure 2. To provide formal definitions, let the notation  $(x \ll y \mid zs)$  mean “an occurrence of  $x$  appears physically before an occurrence of  $y$  in the collection  $zs$ .” That is,  $(x \ll y \mid zs)$  if and only if there are  $i < j$  such that  $x = z_i$  and  $y = z_j$ , where  $z_1, z_2, \dots, z_n$  are the items in  $zs$  listed in their order of appearance in  $zs$ . Note that  $(x \ll x \mid zs)$  if and only if  $x$  occurs at least twice in  $zs$ .

Also, a sorting key of a collection  $zs$  is a function  $\phi(\cdot)$  with an associated linear ordering  $<_\phi$  on its codomain such that, for every pair of items  $x$  and  $y$  in  $zs$  where  $\phi(x) \neq \phi(y)$ , it is the case that  $(x \ll y \mid zs)$  if and only if  $\phi(x) <_\phi \phi(y)$ . Note that a collection may have zero, one, or more sorting keys. Two sorting keys  $\phi(\cdot)$  and  $\psi(\cdot)$  are said to have comparable codomains if their associated linear orderings are identical; i.e. for every  $z$  and  $z'$ ,  $z <_\phi z'$  if and only if  $z <_\psi z'$ . For convenience, in this situation, we write  $<$  to refer to  $<_\phi$  and  $<_\psi$ .

**Definition 4.1** (Monotonicity of `isBefore`). *An `isBefore` predicate is monotonic with respect to two collections  $(xs, ys)$ , which are not necessarily of the same type, if it satisfies the conditions below.*

1. If  $(x \ll x' \mid xs)$ , then for all  $y$  in  $ys$ : `isBefore`( $y, x$ ) implies `isBefore`( $y, x'$ ).
2. If  $(y' \ll y \mid ys)$ , then for all  $x$  in  $xs$ : `isBefore`( $y, x$ ) implies `isBefore`( $y', x$ ).

**Definition 4.2** (Antimonotonicity of `canSee`). *Let `isBefore` be monotonic with respect to  $(xs, ys)$ . A `canSee` predicate is antimonotonic with respect to `isBefore` if it satisfies the conditions below.*

1. If  $(x \ll x' \mid xs)$ , then for all  $y$  in  $ys$ : `isBefore`( $y, x$ ) and not `canSee`( $y, x$ ) implies not `canSee`( $y, x'$ ).
2. If  $(y \ll y' \mid ys)$ , then for all  $x$  in  $xs$ : not `isBefore`( $y, x$ ) and not `canSee`( $y, x$ ) implies not `canSee`( $y', x$ ).

To appreciate the monotonicity conditions, imagine two collections  $xs$  and  $ys$  are being merged without duplicate elimination into a combined list  $zs$ , in a manner that is consistent with the `isBefore` predicate and the physical order of appearance in  $xs$  and  $ys$ . To do this, let  $xs$  comprises  $x_1, x_2, \dots, x_m$  as its elements and  $(x_1 \ll x_2 \ll \dots \ll x_m \mid xs)$ ; let  $ys$  comprises  $y_1, y_2, \dots, y_n$  as its elements and  $(y_1 \ll y_2 \ll \dots \ll y_n \mid ys)$ ; and let  $z_i$  denote the  $i$ th element of  $zs$ . As there is no duplicate elimination, each  $z_i$  is necessarily a choice between some element  $x_j$  in  $xs$  and  $y_k$  in  $ys$ , and  $i = j + k - 1$ , unless all elements of  $xs$  or  $ys$  have already been chosen earlier. Let  $\alpha(i)$  be the index of the element  $x_j$ , i.e.  $j$ ; and  $\beta(i)$  be the index of the element  $y_k$ , i.e.  $k$ . Obviously,  $\alpha(1) = \beta(1) = 1$ . And  $zs$  is necessarily constructed as follows: If  $\alpha(i) > m$  or `isBefore`( $y_{\beta(i)}, x_{\alpha(i)}$ ), then  $z_i = y_{\beta(i)}$ ,  $\alpha(i+1) = \alpha(i)$ , and  $\beta(i+1) = \beta(i) + 1$ ; otherwise,  $z_i = x_{\alpha(i)}$ ,  $\alpha(i+1) = \alpha(i) + 1$ , and  $\beta(i+1) = \beta(i)$ .

Notice that in constructing  $zs$  above, only the `isBefore` predicate is used. The existence of a monotonic predicate `isBefore` with respect to  $(xs, ys)$  does not require  $xs$  and  $ys$  to be ordered by any sorting keys. For example, an “always true” `isBefore` predicate simply puts all elements of  $ys$  before all elements of  $xs$  when merging them into  $zs$  as described above. However, such trivial `isBefore` predicates have limited use.

When  $xs$  and  $ys$  are ordered by some sorting keys, more useful monotonic `isBefore` predicates are definable. For example, as an easy corollary of the construction of  $zs$  above, if  $xs$  and  $ys$  are ordered according to some sorting keys  $\phi(\cdot)$  and  $\psi(\cdot)$  with comparable codomains (i.e.,  $<_\phi$  and  $<_\psi$  are identical and thus can be denoted simply as  $<$ ), then a predicate defined as `isBefore`( $y, x$ ) =  $\psi(y) < \phi(x)$  is guaranteed monotonic with respect to  $(xs, ys)$ . To see this, without loss of generality, suppose for a contradiction that  $(x_i \ll x_j \mid xs)$ ,  $\phi(x_i) \neq \phi(x_j)$ , and `isBefore`( $y, x_i$ ), but not `isBefore`( $y, x_j$ ). This means  $\phi(x_i) < \phi(x_j)$ ,

553  $\psi(y) < \phi(x_i)$ , but  $\psi(y) \not< \phi(x_j)$ . This gives the desired contradiction that  $\phi(x_j) < \phi(x_i)$ .  
 554 This  $\text{isBefore}(y, x) = \psi(y) < \phi(x)$  is a natural bridge between the two sorted collec-  
 555 tions. Specifically, define  $\omega(i) = \phi(z_i)$  if  $z_i$  is from  $x_s$  and  $\omega(i) = \psi(z_i)$  if  $z_i$  is from  $y_s$ ;  
 556 and let  $\omega(z_s)$  denote the collection comprising  $\omega(1), \dots, \omega(n+m)$  in this order. Then,  
 557  $(\omega(i) \ll \omega(j) \mid \omega(z_s))$  implies  $\omega(i) \leq \omega(j)$ . That is,  $\omega(z_s)$  is linearly ordered by  $<$ , the  
 558 associated linear ordering shared by the two sorting keys  $\phi(\cdot)$  and  $\psi(\cdot)$  of  $x_s$  and  $y_s$ .

559 Next, to appreciate the antimonicity conditions, one may eliminate the double nega-  
 560 tives and read these antimonicity conditions as: (1) If  $\text{isBefore}(y, x)$  and  $(x \ll x' \mid x_s)$ ,  
 561 then  $\text{canSee}(y, x')$  implies  $\text{canSee}(y, x)$ ; and (2) If not  $\text{isBefore}(y, x)$  and  $(y \ll y' \mid y_s)$ ,  
 562 then  $\text{canSee}(y', x)$  implies  $\text{canSee}(y, x)$ . Imagine that the  $x$ 's and  $y$ 's are placed on the same  
 563 straight line, from left to right, in a manner consistent with  $\text{isBefore}$  (e.g., as explained  
 564 above). Then, if  $\text{canSee}$  is antimonic to  $\text{isBefore}$ , its antimonicity implies a “right-  
 565 sided” convexity. That is, if  $y$  can see an item  $x$  of  $x_s$  to its right, then it can see all  $x_s$  items  
 566 between itself and this  $x$ . Similarly, if  $x$  can be seen by an item  $y$  of  $y_s$  to its right, then it  
 567 can be seen by all  $y_s$  items between itself and this  $y$ . No “left-sided” convexity is required  
 568 or implied however.

569 It follows that any  $\text{canSee}$  predicate which is reflexive and convex always satisfies the  
 570 antimonicity conditions when  $\text{isBefore}$  satisfies the monotonicity conditions. So, we  
 571 can try checking convexity and reflexivity of  $\text{canSee}$  first, which is a more intuitive task.  
 572 Moreover, though this will not be discussed here, certain optimizations—which are use-  
 573 ful in a parallel distributed setting—are enabled when  $\text{canSee}$  is reflexive and convex.  
 574 Nonetheless, we must stress that the converse is not true. That is, an antimonic  $\text{canSee}$   
 575 predicate needs not be reflexive or convex; e.g., the  $\text{overlap}(y, x)$  predicate from Figure 1  
 576 is an example of a nonconvex antimonic predicate, and the inequality  $m < n$  of two  
 577 integers is an example of a nonreflexive convex antimonic predicate.

578 **Proposition 4.3** (Reflexivity and convexity imply antimonicity). *Let  $x_s$  and  $y_s$  be two*  
 579 *collections, which are not necessarily of the same type. Let  $z_s$  be a collection of some*  
 580 *arbitrary type. Let  $\phi : x_s \rightarrow z_s$  be a sorting key of  $x_s$  and  $\psi : y_s \rightarrow z_s$  be a sorting key of*  
 581  *$y_s$ . Then  $\text{isBefore}$  is monotonic with respect to  $(x_s, y_s)$ , and  $\text{canSee}$  is antimonic with*  
 582 *respect to  $\text{isBefore}$ , if there are predicates  $<_{z_s}$  and  $\triangleleft_{z_s}$  such that all the conditions below*  
 583 *are satisfied.*

- 585 1.  $\phi$  preserves order:  $(x \ll x' \mid x_s)$  implies  $(\phi(x) \ll \phi(x') \mid z_s)$
- 586 2.  $\psi$  preserves order:  $(y \ll y' \mid y_s)$  implies  $(\psi(y) \ll \psi(y') \mid z_s)$
- 587 3.  $<_{z_s}$  preserves  $\text{isBefore}$ :  $\text{isBefore}(y, x)$  if and only if  $\psi(y) <_{z_s} \phi(x)$
- 588 4.  $<_{z_s}$  is monotonic with respect to  $(z_s, z_s)$
- 589 5.  $\triangleleft_{z_s}$  preserves  $\text{canSee}$ :  $\text{canSee}(y, x)$  if and only if  $\psi(y) \triangleleft_{z_s} \phi(x)$
- 590 6.  $\triangleleft_{z_s}$  is reflexive: for all  $z$  in  $z_s$ ,  $z \triangleleft_{z_s} z'$
- 591 7.  $\triangleleft_{z_s}$  is convex: for all  $z_0$  in  $z_s$  and  $(z \ll z' \ll z'' \mid z_s)$ ,  $z \triangleleft_{z_s} z_0$  and  $z'' \triangleleft_{z_s} z_0$  implies  
 592  $z' \triangleleft_{z_s} z_0$ ; and  $z_0 \triangleleft_{z_s} z$  and  $z_0 \triangleleft_{z_s} z''$  implies  $z_0 \triangleleft_{z_s} z'$

593 *In particular, when  $x_s = y_s = z_s$ , and  $\text{isBefore}$  is monotonic with respect to  $(x_s, y_s)$  and*  
 594 *thus  $(z_s, z_s)$ , conditions 1 to 5 above are trivially satisfied by setting the identity function*  
 595 *as  $\phi$  and  $\psi$ ,  $\text{isBefore}$  as  $<_{z_s}$ , and  $\text{canSee}$  as  $\triangleleft_{z_s}$ . Thus, a reflexive and convex  $\text{canSee}$  is*  
 596 *also antimonic.*

The antimonotonicity conditions provide us with two rules for moving on to the next  $x$  or the next  $y$ ; cf. Figure 2. Specifically, according to Antimonotonicity Condition 1, when the current  $y$  in  $ys$  is before the current  $x$  in  $xs$ , and this  $y$  cannot “see” (i.e. does not match) this  $x$ , then this  $y$  cannot see any of the following items in  $xs$  either. Therefore, it is not necessary to try matching the current  $y$  to the rest of the items in  $xs$ , and we can move on to the next item in  $ys$ . On the other hand, according to Antimonotonicity Condition 2, when the current  $y$  in  $ys$  is not before the current  $x$  in  $xs$ , and this  $y$  cannot see this  $x$ , then all subsequent items in  $ys$  cannot see this  $x$  either. Therefore, it is not necessary to try matching the current  $x$  to the rest of the items in  $ys$ , and we can safely move on to the next item in  $xs$ .

When neither rule is triggered, regardless of whether the current  $y$  in  $ys$  is or is not before the current  $x$  in  $xs$ , this  $y$  can see this  $x$ . That is, we have a matching pair of  $x$  and  $y$  to perform some specified actions on. After the actions are performed, we can choose to move on to the next item in  $xs$  or in  $ys$ . In this work, we decide to keep the collection  $xs$  as the reference and to move on to the next item in the collection  $ys$ . Since the next item in  $xs$  may be an item that the current  $y$  can see, before moving on to the next item in  $ys$ , we should also “save” the current  $y$ ; when we eventually move on to the next item in  $xs$ , we must remember to “rewind” our position in  $ys$  back to all these  $y$ ’s saved during the processing of the current  $x$ .

Together, these conditions lead to what we call a *Synchrony fold*—the `syncFold` function defined in Figure 3—which iterates on two collections in synchrony.

What does `syncFold(f, e, bf, cs)(xs, ys)` do? To answer this question, consider the function `slowFold(f, e, cs)(xs, ys)` which is also defined in Figure 3. The function `slowFold(f, e, cs)(xs, ys)` first initializes an internal variable `acc` to `e`; then iterates through every pair of  $x$  in  $xs$  and  $y$  in  $ys$ , and updates `acc` to `f(x, y, acc)` whenever `cs(y, x)`; at the end of the iteration, it outputs the value of `acc`.

Remarkably, when `bf` is monotonic with respect to  $(xs, ys)$  and `cs` is antimonotonic with respect to `bf`, `syncFold(f, e, bf, cs)(xs, ys)` computes the same result as `slowFold(f, e, cs)(xs, ys)`. Furthermore, `syncFold` has a potentially linear complexity  $O(|xs| + k|ys|)$  in terms of number of calls to the function  $f$ , when `cs` has degree  $< k$  in the sense that  $|\{x \in xs \text{ such that } cs(y, x)\}| < k$  for each  $y$  in  $ys$ . Whereas, `slowFold` has quadratic complexity  $O(|xs| \cdot |ys|)$ .

**Theorem 4.4** (Synchrony fold). *Suppose `isBefore` is monotonic with respect to  $(xs, ys)$  and `canSee` is antimonotonic with respect to `isBefore`.*

1. `syncFold(f, e, isBefore, canSee)(xs, ys) = slowFold(f, e, canSee)(xs, ys)`.
2. `slowFold(f, e, canSee)(xs, ys)` calls the function  $f$  a total of  $|xs| \cdot |ys|$  number of times.
3. `syncFold(f, e, isBefore, canSee)(xs, ys)` calls the function  $f$  at most  $|xs| + k|y|$  number of times, if `canSee` has degree  $< k$  with respect to  $(xs, ys)$ .

**Proof** For Part 1, consider the function `aux(xs, ys, zs, acc)` in `syncFold`. Suppose `isBefore` is monotonic with respect to  $(xs, zs ++ ys)$ , and `canSee` is antimonotonic with

```

645 def syncFold[A,B,C]
646   (f: (A,B,C) => C, e: C, bf: (B,A) => Boolean, cs: (B,A) => Boolean)
647   (xs: Vec[A], ys: Vec[B])
648 : C = {
649   // Requires: bf monotonic wrt (xs,ys); cs antimonotonic wrt bf.
650   // Assumes: isEmpty, head, tail are constant time;
651   //           prepend (++) is linear in its left argument;
652   //           single-item postpend (:+) is constant time.
653
654   def aux(xs: Vec[A], ys: Vec[B], zs: Vec[B], acc: C): C =
655     if (xs.isEmpty) acc
656     else if (ys.isEmpty && zs.isEmpty) acc
657     else if (ys.isEmpty) aux(xs.tail, zs, Vec(), acc)
658     else {
659       val (x, y) = (xs.head, ys.head)
660       (bf(y, x), cs(y, x)) match {
661         case (true, false) =>
662           // Antimonotonicity Condition 1:
663           // bf(y,x) & !cs(y,x) => all x' after x: !cs(y,x')
664           // So, y can be discarded safely; move on to next y.
665           aux(xs, ys.tail, zs, acc)
666
667         case (false, false) =>
668           // Antimonotonicity Condition 2:
669           // !bf(y,x) & !cs(y,x) => all y' after y: !cs(y',x)
670           // So x can be discarded safely. But the next x may
671           // still be able to see some y saved earlier in zs.
672           aux(xs.tail, zs ++: ys, Vec(), acc)
673
674         case (_, true) =>
675           // At this point, cs(y,x); so process (x,y) using f.
676           // Save this y as it may see next x; move on to next y.
677           aux(xs, ys.tail, zs :+ y, f(x, y, acc))
678       }
679     }
680   aux(xs, ys, Vec(), e)
681 }
682
683 def slowFold[A,B,C]
684   (f: (A,B,C) => C, e: C, cs: (B,A) => Boolean)
685   (xs: Vec[A], ys: Vec[B])
686 : C = {
687   var acc: C = e
688   for (x <- xs; y <- ys; if cs(y, x)) { acc = f(x, y, acc) }
689   return acc
690 }

```

Fig. 3. Definitions of `syncFold` and `slowFold`. These two programs compute the same results when `bf` is monotonic with respect to  $(xs, ys)$  and `cs` is antimonotonic with respect to `bf`. However, `syncFold` is more efficient than `slowFold`.

respect to `isBefore`. If `xs` is non-empty, let `x` be `xs.head`, and  $z_1, \dots, z_n$  be the items in `zs` such that `canSee(z1, x)`, ..., `canSee(zn, x)`, and  $acc = f(x, z_n, \dots f(x, z_1, e) \dots)$ . If `xs` is empty, let  $acc = e$ . Then, an induction on  $(|xs|, |ys|)$  shows that

```

686   aux(xs, ys, zs, acc)
687 = slowFold(f, e, canSee)(xs, zs ++: ys)

```

So,

```

    syncFold(f, e, isBefore, canSee)(xs, ys)
691 = aux(xs, ys, Vec(), e)
692 = slowFold(f, e, canSee)(xs, ys)
693

```

694 For Part 2, it is obvious that `slowFold(f, e, canSee)(xs, ys)` calls the function  $f$  a total  
695 of  $|xs| \cdot |ys|$  number of times.

696 For Part 3, on each call to `aux` in `syncFold`, either  $xs$  or  $ys$  is shortened by 1 item. This  
697 gives  $|xs| + |ys|$  calls to `aux`. In some calls,  $ys$  is prepended with  $zs$ . Recall the assumption  
698 that `canSee` has degree  $< k$ . Thus, each item in  $ys$  can see fewer than  $k$  items in  $xs$ . So,  
699 the total size of  $zs$  summed over all the calls to `aux` is at most  $(k - 1)|ys|$ ; these are the  
700 maximum number of additional calls to `aux`. Therefore, the total number of calls to `aux`,  
701 and thus to  $f$ , is at most  $|xs| + k|ys|$ . ■

## 703 4.2 Second Synchrony fold

704 `SyncFold(f, e, bf, cs)(xs, ys)` discards items in  $xs$  that no item in  $ys$  sees. This may not  
705 be desired in some situations, e.g., when someone actually wants to retrieve those items  
706 in  $xs$  that no item in  $ys$  sees. Also, `syncFold` pairs up each  $x$  in  $xs$  with each  $y$  in  $ys$  that  
707 sees it, and applies the function  $f$  on these pairs one by one. This may not be convenient  
708 in some situations; e.g., when someone wants to count the number of  $y$ 's that see an  $x$ .  
709 Hence, it might be useful to also provide a second Synchrony fold function `syncFoldGrp`  
710 which processes, as a group, those  $y$ 's that see an  $x$ .

711 An astute reader might have realised that, in the definition provided in Figure 3, `syncFold`  
712 keeps the  $y$ 's that can see the current  $x$  in the collection  $zs$ . So, as defined in Figure 4,  
713 `syncFoldGrp(f, e, bf, cs)(xs, ys)` is just `syncFold` with  $f$  applied to  $(x, zs, acc)$  instead  
714 of  $(x, y, acc)$ . The function `syncFoldGrp(f, e, bf, cs)(xs, ys)` computes the same result  
715 as `slowFoldGrp(f, e, cs)(xs, ys)`, which is also defined in Figure 4 and is much easier  
716 to understand. However, while the former can be linear in time complexity, the latter is  
717 quadratic.  
718  
719  
720

721 **Theorem 4.5** (Second Synchrony fold). *Suppose `isBefore` is monotonic with respect*  
722 *to  $(xs, ys)$  and `canSee` is antimonotonic with respect to `isBefore`. Then,*

- 723 1. `syncFoldGrp(f, e, isBefore, canSee)(xs, ys) =`  
724 `slowFoldGrp(f, e, canSee)(xs, ys).`

725 *Suppose further that `canSee` has degree  $< k$  with respect to  $(xs, ys)$ , and  $f$  has linear time*  
726 *complexity in its second argument, and other arguments have negligible influence on  $f$ 's time*  
727 *complexity. Then,*

- 729 2. `slowFoldGrp(f, e, canSee)(xs, ys)` *has time complexity*  $O((|xs| + k)|ys|)$ .
- 730 3. `syncFoldGrp(f, e, isBefore, canSee)(xs, ys)` *has time complexity*  $O(|xs| +$   
731  $2k|ys|)$ .

732  
733  
734 **Proof** For Part 1, consider the function `aux(xs, ys, zs, acc)` in `syncFoldGrp`. Suppose  
735 `isBefore` is monotonic with respect to  $(xs, zs ++: ys)$ , and `canSee` is antimonotonic with  
736

```

737 def syncFoldGrp[A,B,C]
738   (f: (A,Vec[B],C)=>C, e: C, bf: (B,A)=>Boolean, cs: (B,A)=>Boolean)
739   (xs: Vec[A], ys: Vec[B])
740 : C = {
741   // Requires: bf monotonic wrt (xs,ys) & cs antimonotonic wrt bf.
742
743   def aux(xs: Vec[A], ys: Vec[B], zs: Vec[B], acc: C): C =
744     if (xs.isEmpty) acc
745     else if (ys.isEmpty && zs.isEmpty) acc
746     else if (ys.isEmpty) aux(xs.tail, zs, Vec(), f(xs.head, zs, acc))
747     else {
748       val (x,y) = (xs.head, ys.head)
749       (bf(y, x), cs(y, x)) match {
750         case (true, false) =>
751           // Antimonotonicity Condition 1:
752           // bf(y,x) & !cs(y,x) => all x' after x: !cs(y,x')
753           // So, y can be discarded safely; move on to next y.
754           aux(xs, ys.tail, zs, acc)
755
756         case (false, false) =>
757           // Antimonotonicity Condition 2:
758           // !bf(y,x) & !cs(y,x) => all y' after y: !cs(y',x)
759           // So, x can be discarded. And the y accumulated in zs
760           // should now be processed by f in one go. Note: the
761           // next x may be able to see some y accumulated in zs.
762           aux(xs.tail, zs ++: ys, Vec(), f(x, zs, acc))
763
764         case (_, true) =>
765           // At this point, cs(y,x).
766           // Accumulate this y in zs; move on to next y.
767           aux(xs, ys.tail, zs :+ y, acc)
768       }
769     }
770
771   aux(xs, ys, Vec(), e)
772 }
773
774 def slowFoldGrp[A,B,C]
775   (f: (A,Vec[B],C) => C, e: C, cs: (B,A) => Boolean)
776   (xs: Vec[A], ys: Vec[B])
777 : C = {
778   var acc: C = e
779   for (x <- xs; zs = for (y <- ys; if cs(y, x)) yield y) {
780     acc = f(x, zs, acc)
781   }
782   return acc
783 }

```

Fig. 4. Definitions of `syncFoldGrp` and `slowFoldGrp`. They compute the same results when `bf` is monotonic with respect to  $(xs, ys)$  and `cs` is antimonotonic with respect to `bf`. However, `syncFoldGrp` is more efficient than `slowFoldGrp`.

respect to `isBefore`. Suppose also that `canSee(z, x)` for each `z` in `zs`, when `xs` is non-empty and `x` is `xs.head`. Then, an induction on  $(|xs|, |ys|)$  shows that

```

778   aux(xs, ys, zs, acc)
779 = slowFoldGrp(f, acc, canSee)(xs, zs ++: ys)

```

So,

```

783   syncFoldGrp(f, e, isBefore, canSee)(xs, ys)
784   = aux(xs, ys, Vec(), e)
785   = slowFoldGrp(f, e, canSee)(xs, ys)

```

786 For Part 2, the theorem assumes that `canSee` has degree  $< k$ , and `f` has time complexity  
787 linear in its second argument and independent of its other arguments. The first assumption  
788 implies that the total size of `zs` over all the calls to `f` is at most  $k|ys|$ . The second assumption  
789 implies that the total time complexity due to calls to `f` is  $O(k|ys|)$ . The nested loops of  
790 `slowFoldGrp`, excluding calls to `f`, has  $O(|xs| \cdot |ys|)$  time complexity. Thus, summing these  
791 two components gives a quadratic time complexity,  $O((|xs| + k)|ys|)$ .

792 For Part 3, again recall the two assumptions of the theorem, viz. `canSee` has degree  $<$   
793  $k$ , and `f` has time complexity linear in its second argument and independent of its other  
794 arguments. The first assumption implies that the total size of `zs` over all the calls to `f` is  
795 at most  $k|ys|$ . The second assumption implies that the total time complexity due to calls  
796 to `f` is  $O(k|ys|)$ . In addition, as in `syncFold`, there are at most  $|xs| + k|ys|$  calls to `aux` in  
797 `syncFoldGrp`. Summing these gives a linear time complexity,  $O(|xs| + 2k|ys|)$ . ■

798 Now, let `snoc(x, zs, a) = a :+ (x, zs)` add `(x, zs)` to the end of a collection `a`. Then,

```

799
800   slowFoldGrp(snoc, Vec(), cs)(Vec(x), ys)
801   = Vec((x, for (y <- ys; if cs(y, x)) yield y))
802
803   slowFoldGrp(snoc, Vec(), cs)(xs, ys)
804   = for (x <- xs; (x', zs) <- slowFoldGrp(snoc, Vec(), cs)(Vec(x), ys)) yield (x', zs)

```

805 The corollary below now follows from Theorem 4.5. This corollary is helpful for a deeper  
806 understanding of `syncFoldGrp`, leading later to the design of Synchrony iterator in Section 6.  
807

808 **Corollary 4.6.** *Let `isBefore` be monotonic with respect to  $(xs, ys)$ , and `canSee` be*  
809 *antimonotonic with respect to `isBefore`. Let `snoc(x, zs, a) = a :+ (x, zs)`. Then,*

```

810
811   syncFoldGrp(snoc, Vec(), isBefore, canSee)(xs, ys)
812   = for (x <- xs; (x', zs) <- syncFoldGrp(snoc, Vec(), isBefore, canSee)(Vec(x), ys)
813       yield (x', zs)

```

### 814 4.3 Synchrony fold vs foldLeft

815 As mentioned earlier, `syncFold` and `syncFoldGrp` are generalizations of `foldLeft`. In particu-  
816 lar, as shown below, `foldLeft` is definable via either of them.

```

817
818   xs.foldLeft(e)(g)
819   = syncFold((x, _, a) => g(a, x), e, (_, _) => true, (_, _) => true)(xs, Vec())
820   = syncFoldGrp((x, _, a) => g(a, x), e, (_, _) => true, (_, _) => true)(xs, Vec())

```

821 Furthermore, both definitions are as efficient as the implementation of `foldLeft` in  
822 collection-type libraries; e.g., if the function `g` above has  $O(1)$  time complexity, then both  
823 implementations of `foldLeft` above have  $O(|xs|)$  time complexity, same as any typical  
824 implementation of `foldLeft` in collection-type libraries of modern programming languages.  
825  
826  
827  
828

At the same time, functions expressible by `syncFold` and `syncFoldGrp` are also expressible in first-order restricted Scala when `foldLeft` is available. Let `isBefore` be monotonic with respect to  $(xs, ys)$ , and `canSee` be antimonotonic with respect to `isBefore`. Then,

```

829
830
831
832   syncFold(f, e, isBefore, canSee)(xs, ys)
833 = xs.foldLeft(e)((a, x) =>
834   ys.foldLeft(a)((a', y) => if (canSee(y, x)) f(x, y, a') else a'))
835
836   syncFoldGrp(f', e, isBefore, canSee)(xs, ys)
837 = xs.foldLeft(e)((a, x) =>
838   f'(x, for (y <- ys; if canSee(y, x)) yield y, a))

```

These implementations of `syncFold` and `syncFoldGrp` in terms of `foldLeft` are quadratic in time complexity. They are also somewhat more convoluted than the implementations of `foldLeft` in terms of `syncFold` and `syncFoldGrp`. Perhaps more ingenious programmers can find some simpler ways of implementing `syncFold` and `syncFoldGrp` solely in terms of `foldLeft`. Unfortunately, due to Lemma 3.3, there is no way they can find an efficient linear-time implementation of either one using `foldLeft` alone under the first-order restriction.

**Proposition 4.7** (*SyncFold and syncFoldGrp are conservative extensions of foldLeft*). *The extensional expressive power of Scala under the first-order restriction, when foldLeft is available, is the same with or without syncFold and syncFoldGrp. However, more efficient algorithms for some functions (e.g., a linear-time algorithm for low-selectivity join) can be defined using syncFold and syncFoldGrp than using foldLeft in Scala under the first-order restriction.*

It is worth noting that `syncFold(f, e, isBefore, canSee)(xs, ys)` and the expression `syncFoldGrp((x, zs, a) => zs.foldLeft(a)((a', z) => f(z, a')), e, isBefore, canSee)(xs, ys)` compute the same function at comparable time complexity. So, `syncFold` can be defined efficiently using `syncFoldGrp`. Similarly, `SyncFoldGrp` can also be implemented at comparable time complexity using `syncFold`. To wit, as presented later in Section 5.1, efficient `takeWhile` and `dropWhile` are definable by `syncFold`; in turn, efficient `syncFoldGrp` needs only a straightforward modification to the implementation—using `foldLeft`, `takeWhile`, and `dropWhile`—of the function `groups2` shown in Figure 13 of Section 8.1.

#### 4.4 Synchrony fold in action

Linear time complexity for the example from Section 2, `ov1(xs, ys)`, can be achieved using `syncFold`. The codes for `ov3(xs, ys)` below shows that a user-programmer only has to provide straightforward definitions for the `isBefore` and `canSee` predicates; for this example, these are the `isBefore` and `overlap` predicates defined earlier in Figure 1. There is no worry about getting the “synchronized” iteration of `xs` and `ys` right, as `syncFold` takes care of this already. The linear time complexity is easily appreciated using Theorem 4.4 when `overlap` has a low degree with respect to  $(xs, ys)$ , i.e. each event in `ys` overlaps few events in `xs`.

```

871 def ov3(xs: Vec[Event], ys: Vec[Event]) = {
872   // Requires: xs and ys are sorted lexicographically by (start, end).
873   // Note: isBefore and overlap are as defined in Figure 1.

```

```

875   def f(x: Event, y: Event, acc: Vec[(Event, Event)]) = acc :+ (x, y)
876   syncFold(f, Vec(), isBefore, overlap)(xs, ys)
877 }

```

There is a loose end to be tied up in the example above, viz. verifying that `isBefore` is monotonic with respect to  $(xs, ys)$  and `overlap` is antimonotonic with respect to `isBefore`. This is omitted here, as it is straightforward under the assumption that  $(xs, ys)$  are lexicographically ordered by the `start` and `end` point of their events.

As an example of `syncFoldGrp`, it is used below to count in potentially linear time the number of events in `ys` that each event in `xs` overlaps with. The linear time complexity follows from Theorem 4.5.

```

885 def ovCount(xs: Vec[Event], ys: Vec[Event]): Vec[(Event, Int)] = {
886   // Requires: xs and ys are sorted lexicographically by (start, end).
887   // Note: isBefore and overlap are as defined in Figure 1.
888   def f(x: Event, zs: Vec[Event], acc: Vec[(Event, Int)]) = {
889     acc :+ (x, zs.length)
890   }
891   syncFoldGrp(f, Vec(), isBefore, overlap)(xs, ys)
892 }

```

Comparing `ov3(xs, ys)` above and `ov2(xs, ys)` from Figure 1, the design of Synchrony fold makes clear three orthogonal aspects of the event-overlap example: connecting the orderings on the two collections, identifying matching pairs, and acting on matching pairs. With regard to connecting the orderings on the two collections, the “navigation” is captured by the `isBefore` predicate. With regard to identification of matching pairs, it is captured by the `canSee` predicate (i.e. `overlap`). Finally, with regard to action on matching pairs, this is captured by the function `f`. Making these three orthogonal aspects explicit brings about a more concise and precise understanding (Schmidt, 1986; Hunt & Thomas, 2000; Sebesta, 2010). For example, assuming `isBefore` is monotonic with respect to  $(xs, ys)$  and `canSee` is antimonotonic with respect to `isBefore`, one can read `syncFold(f, e, isBefore, canSee)(xs, ys)` simply as “for each pair in  $(xs, ys)$  satisfying `canSee`, do `f` on it.” Hopefully, this clarity makes it easier to see mistakes, and thus easier to write programs correctly.

This simple way to read Synchrony fold programs was in fact formalized earlier via Theorem 4.4 and 4.5. These two theorems reveal the extensional equivalence of `syncFold` and `slowFold`, and of `syncFoldGrp` and `slowFoldGrp`. While `slowFold` and `slowFoldGrp` are intuitive, they use some local side effects. Now, comparing `ov3(xs, ys)` and `ov1(xs, ys)`, a straightforward relationship between a restricted form of `syncFold` and `syncFoldGrp` and comprehension syntax can be further discerned below; this time without side effects. This relationship also shows that any join whose predicate  $p(y, x)$  can be decomposed into an antimonotonic predicate `canSee(y, x)` and a residual predicate  $h(y, x)$ , can be implemented using `syncFold` and `syncFoldGrp` efficiently.

**Proposition 4.8** (Comprehending `syncFold` and `syncFoldGrp`). *Suppose  $xs$  and  $ys$  are two collections, `isBefore` is monotonic with respect to  $(xs, ys)$ , and `canSee` is antimonotonic with respect to `isBefore`. Then, these three Scala programs express the same function:*

1. `for (x <- xs; y <- ys; if canSee(y, x) && h(y, x)) yield g(x, y)`

2. `syncFold(f, Vec(), isBefore, canSee)(xs, ys)`, where  
`f(x, y, acc) = if (h(y,x)) { acc :+ g(x, y) } else acc`
3. `syncFoldGrp(f', Vec(), isBefore, canSee)(xs, ys)`, where  
`f'(x, zs, acc) = acc :++ for (z <- zs; if h(z, x)) yield g(x, z)`

However, when `canSee` has a low degree and `g` and `h` have  $O(1)$  time complexity, the first program is quadratic while the second and third programs are linear in their time complexity with respect to  $|xs|$  and  $|ys|$ .

## 5 Synchrony generator

Lemma 3.1 indicates that an intensional expressiveness gap already exists in first-order restricted Scala sans library functions. And Lemma 3.3 further indicates that this same gap exists practically unmitigated when first-order restricted Scala is augmented with `foldLeft`. On the one hand, Proposition 4.7 shows that the two Synchrony folds are conservative extensions of first-order restricted Scala augmented with `foldLeft`, and significantly increases the algorithmic richness of this fragment of Scala. On the other hand, Proposition 4.7 also means that Synchrony fold is an overkill as a solution for this gap which originated at the level of first-order restricted Scala without library functions, since Synchrony fold adds much extra extensional expressive power to this fragment of Scala while fixing its intensional expressive power gap.

This section identifies a restriction on Synchrony fold to fix this gap at its root, i.e. at the level of unaugmented first-order restricted Scala. The significance of this restricted form, viz. *Synchrony generator*, in the context of database joins is also discussed.

### 5.1 Deriving Synchrony generator

As mentioned, we wish to identify some restriction on Synchrony fold to cut its extensional expressive power to that of first-order restricted Scala sans library functions. Proposition 4.8 suggests the two solutions `syncMap` and `syncFlatMap`, shown in Figure 5.

Ignoring efficiency issues, the functions expressible by `syncMap` and `syncFlatMap` are already expressible just using comprehension syntax, when `isBefore` is monotonic with respect to  $(xs, ys)$  and `canSee` is antimotonic with respect to `isBefore`. Specifically,

```

syncMap(f, isBefore, canSee)(xs, ys)
= for (x <- xs; y <- ys; if canSee(y, x)) yield f(x, y)

syncFlatMap(f, isBefore, canSee)(xs, ys)
= for (x <- xs; z <- f(x, for (y <- ys; if canSee(y, x)) yield y))
  yield z

```

Thus, `syncMap` and `syncFlatMap` do not add extensional expressive power to first-order restricted Scala sans library functions, but add to it sufficient algorithmic power to implement efficient low-selectivity joins.

In fact, an even more stringent restriction, the *Synchrony generators*, `syncGen` and `syncGenGrp`, also depicted in Figure 5, can provide the same extra intensional expressive power as `syncMap` and `syncFlatMap`. This is because

```

967 def syncMap[A,B,C]
968   (f: (A,B)=>C, bf: (B,A)=>Boolean, cs: (B,A)=>Boolean)
969   (xs: Vec[A], ys: Vec[B])
970 : Vec[C] = {
971   // Requires: bf monotonic wrt (xs, ys); cs antimonotonic wrt bf.
972   val step = (x: A, y: B, acc: Vec[C]) => acc :+ f(x,y)
973   syncFold(step, Vec(), bf, cs)(xs, ys)
974 }
975
976 def syncFlatMap[A,B,C]
977   (f: (A,Vec[B])=>Vec[C], bf: (B,A)=>Boolean, cs: (B,A)=>Boolean)
978   (xs: Vec[A], ys: Vec[B])
979 : Vec[C] = {
980   // Requires: bf monotonic wrt (xs, ys); cs antimonotonic wrt bf.
981   val step = (x: A, zs: Vec[B], acc: Vec[C]) => acc :++ f(x,zs)
982   syncFoldGrp(step, Vec(), bf, cs)(xs, ys)
983 }
984
985 def syncGen[A,B]
986   (isBefore: (B,A) => Boolean, canSee: (B,A) => Boolean)
987   (xs: Vec[A], ys: Vec[B])
988 : Vec[(A,B)] = {
989   // Requires: bf monotonic wrt (xs, ys); cs antimonotonic wrt bf.
990   val step = (x: A, y: B, acc: Vec[(A,B)]) => acc :+ (x, y)
991   val e: Vec[(A,B)] = Vec()
992   syncFold(step, e, isBefore, canSee)(xs, ys)
993 }
994
995 def syncGenGrp[A,B]
996   (bf: (B,A) => Boolean, cs: (B,A) => Boolean)
997   (xs: Vec[A], ys: Vec[B])
998 : Vec[(A,Vec[B])] = {
999   // Requires: bf monotonic wrt (xs, ys); cs antimonotonic wrt bf.
1000   val step = (x: A, zs: Vec[B], acc: Vec[(A,Vec[B])]) => acc :+ (x,zs)
1001   val e: Vec[(A,Vec[B])] = Vec()
1002   syncFoldGrp(step, e, bf, cs)(xs, ys)
1003 }

```

Fig. 5. Definitions of syncMap, syncFlatMap, syncGen and SyncGenGrp.

```

995 syncFlatMap(f, isBefore, canSee)(xs, ys)
996 = for ((x, zs) <- syncGenGrp(isBefore, canSee)(xs, ys); z <- f(x, zs)) yield z
997
998 syncMap(f, isBefore, canSee)(xs, ys)
999 = for ((x, y) <- syncGen(isBefore, canSee)(xs, ys)) yield f(x, y)

```

Strictly speaking, syncGenGrp is not first-order restricted as it returns a nested collection. However, let us constrain it to be used strictly as a generator  $(x, zs) \leftarrow \text{syncGenGrp}(bf, cs)(xs, ys)$  in a comprehension construct, with the understanding that  $\text{for } ((x, zs) \leftarrow \text{syncGenGrp}(bf, cs)(xs, ys); \dots) \text{ yield } e$  is desugared to  $\text{syncFlatMap}((x, zs) \Rightarrow \text{for } (\dots) \text{ yield } e, bf, cs)(xs, ys)$ . With this constraint, syncGenGrp can justifiably be viewed as a first-order construct, as it becomes mere syntactic sugar which gets desugared into a first-order construct.

As shown earlier, syncMap and syncFlatMap are expressible as functions in comprehension syntax. And syncGen and syncGenGrp are desugared into syncMap and syncFlatMap. So, the theorem below follows.

1010  
1011  
1012

**Theorem 5.1.** *The extensional expressive power of Scala under the first-order restriction, is the same with or without any of `syncMap`, `syncFlatMap`, `syncGen`, and `syncGenGrp`. However, more efficient algorithms for some functions (e.g., a linear-time algorithm for low-selectivity join) can be defined when any of `syncMap`, `syncFlatMap`, `syncGen`, and `syncGenGrp` is made available in this fragment of Scala.*

For illustration, the function `ov1(xs, ys)` from Figure 1 is expressed below using `syncGen`. This version, `ov4(xs, ys)`, as with `ov3(xs, ys)` in Section 4.4, has linear time complexity when selectivity is low.

```

1022 def ov4(xs: Vec[Event], ys: Vec[Event]): Vec[(Event, Event)] = {
1023   // Requires: xs and ys sorted lexicographically by (start, end).
1024   // Note: isBefore and overlap are as defined in Figure 1.
1025   syncGen(isBefore, overlap)(xs, ys)
1026 }

```

Recall also Lemma 3.2 that  $\mathcal{NRC}_1(\leq, \text{takeWhile}, \text{dropWhile}, \text{sort})$  cannot realise efficient low-selectivity joins. Therefore, first-order restricted Scala augmented with `takeWhile` and `dropWhile`, by themselves, cannot implement `syncGen` and `syncGenGrp` efficiently. On the other hand, both `takeWhile` and `dropWhile` can be realised quite efficiently and succinctly using either Synchrony generator. For example,

```

1032 for (x <- xs.takeWhile(p)) yield f(x)
1033 = for ((_, x) <- syncGen((_,_)=>false, (x,_)=>p(x))(Vec(()), xs)) yield f(x)
1034
1035 for (x <- xs.dropWhile(p)) yield f(x)
1036 = for ((x,_) <- syncGen((y,x)=>!y && !p(x), (y,_)=>y)(xs, Vec(false,true))) yield f(x)

```

## 5.2 Synchrony generator vs database merge join

As mentioned in Section 2, a database join having a join predicate comprising entirely of equality tests is an equijoin, and those comprising entirely of inequality tests is a non-equijoin. A relational database system executes joins using a variety of strategies (Blasgen & Eswaran, 1977; Mishra & Eich, 1992; Silberschatz *et al.*, 2016). Where possible, a relational database system decomposes a join predicate into an equijoin part and a residual part; it then executes the equijoin part using either an index join (if suitable indices are available), or a merge join (if indices are not available but the relations are already appropriately sorted), or a sort-merge join or a hash join (if indices are not available and the relations are not already sorted); finally, it executes the residual part as a selection predicate on the result of the equijoin. So, the time complexity is always linear or at worst linearithmic for a join which has an equijoin part that has low selectivity. For a non-equijoin, most relational database systems execute it using nested loops, which have quadratic time complexity. However, some relational database systems can execute some restricted forms of non-equijoin, such as a band join  $x.a \leq y.b \leq x.c$ , more efficiently (e.g., in linear time, when the band join predicate  $x.a \leq y.b \leq x.c$  has low selectivity.)

The Synchrony generator `syncGen(isBefore, canSee)` is closely related to, and is an elegant generalization of, the merge join used in relational database systems. In relational database systems, the merge join is always applied on a pair of relational tables `xs` and `ys`

1059 which are sorted according to some sorting keys  $\phi(\cdot)$  and  $\psi(\cdot)$  with comparable codomains.  
 1060 This induces a linear ordering  $\text{isBefore}(y, x) = \psi(y) < \phi(x)$  on items in the two tables.  
 1061 So, by construction, this  $\text{isBefore}$  predicate is monotonic with respect to  $(x_s, y_s)$ .

1062 For the standard merge join (Silberschatz *et al.*, 2016), the join predicate  $\text{canSee}$  must  
 1063 comprise entirely of equality tests (i.e. an equijoin). So,  $\text{canSee}$  is reflexive and con-  
 1064 vex; thus, by Proposition 4.3, it is antimonotonic with respect to  $\text{isBefore}$ . All relational  
 1065 database systems also support a variant of the merge join where the join predicate is a  
 1066 single inequality like  $\text{canSee}(y, x) = x.a < y.b$  or  $\text{canSee}(y, x) = x.a \leq y.b$ . The for-  
 1067 mer is antimonotonic and convex but not reflexive, the latter is convex and reflexive and  
 1068 thus also antimonotonic. Some database systems support a range join predicate of the  
 1069 form  $\text{canSee}(y, x) = x.a - \epsilon \leq y.b \leq x.a + \epsilon$ ; cf. DeWitt *et al.* (1991). This is a reflex-  
 1070 ive and convex predicate. Thus, by Proposition 4.3, it is antimonotonic with respect to  
 1071  $\text{isBefore}$ . Newer database systems support a band join predicate of the form  $\text{canSee}(y, x) =$   
 1072  $x.a \leq y.b \leq x.c$ . This predicate is antimonotonic but not convex. Some database systems  
 1073 support an interval join predicate of the form  $\text{canSee}(y, x) = x.a \leq y.b \ \&\& \ y.c \leq x.d$  on  
 1074 some special data types, such as those associated with time periods, where the constraints  
 1075  $x.a \leq x.d$  and  $y.c \leq y.b$  are known or enforced by these database systems; cf. Piatov *et al.*  
 1076 (2016) and Dignoes *et al.* (2021). This predicate, taking into account the two associated  
 1077 constraints, is antimonotonic but not convex.

1078 As all these  $\text{canSee}$  predicates are antimonotonic, they can be used legitimately in  
 1079  $\text{syncGen}(\text{isBefore}, \text{canSee})(x_s, y_s)$ . And this computes the corresponding equijoin, single-  
 1080 inequality merge join, range join, band join, and interval join. Clearly, the monotonicity of  
 1081 the  $\text{isBefore}$  predicate and the antimonotonicity of the  $\text{canSee}$  predicate constitute a more  
 1082 general and more elegant condition for correctness than the adhoc syntactic forms required  
 1083 by current formulations of equijoin, single-inequality merge join, range join, band join,  
 1084 and interval join.

1085 There have been many works introducing join algorithms in the database community  
 1086 to handle non-equijoin, from early studies by DeWitt *et al.* (1991) to recent studies by  
 1087 Piatov *et al.* (2016) and Dignoes *et al.* (2021). These works generally require a combina-  
 1088 tion of new data structures, new evaluation techniques, and even exploitation of hardware  
 1089 features of modern CPU architectures. These are tools which are not part of the repertoire  
 1090 of an average programmer. Moreover, these works consider only some syntactic forms. In  
 1091 contrast, Synchrony generator efficiently and uniformly implements a more general class  
 1092 of non-equijoin without requiring any of these. This makes Synchrony generator rather  
 1093 appealing as an addition to collection-type function libraries of programming languages.

1094 Moreover, by Theorem 4.4, the time complexity of  $\text{syncGen}(\text{isBefore}, \text{canSee})(x_s, y_s)$   
 1095 is  $O(|x_s| + k|y_s|)$  where  $k$  is the degree of the  $\text{canSee}$  predicate. It is worth noting that  
 1096 the size of the result of  $\text{syncGen}(\text{isBefore}, \text{canSee})(x_s, y_s)$  is also  $O(|x_s| + k|y_s|)$ , which  
 1097 obviously constitutes a lowerbound on the efficiency of any algorithm for computing the  
 1098 same join. So, despite Synchrony generator being much simpler and more general than ear-  
 1099 lier algorithms for various more restricted forms of non-equijoin, the time complexity of  
 1100  $\text{syncGen}(\text{isBefore}, \text{canSee})(x_s, y_s)$  is already asymptotically optimal. Even more impres-  
 1101 sive, it does this while staying strictly within the extensional expressive power of first-order  
 1102 restricted Scala unaugmented with any library function.

```

1105 case class Event(start: Int, end: Int, id: String)
1106 // Constraint: start < end
1107
1107 val isBeforeWithId = (y: Event, x: Event) => {
1108     (y.id < x.id) ||
1109     (y.id == x.id && y.start < x.start) ||
1110     (y.id == x.id && y.start == x.start && y.end < x.end)
1111 }
1112
1112 val overlapWithId(y: Event, x: Event) => {
1113     (y.id == x.id) &&
1114     (x.start < y.end && y.start < x.end)
1115 }
1116
1116 def ovWithId(xs: Vec[Event], ys: Vec[Event]) = {
1117     // Requires: xs and ys are sorted by (id, start, end)
1118     syncGen(isBeforeWithId, overlapWithId)(xs, ys)
1119 }
1120
1120 // Query1: ovWithId directly translated into SQL
1121 SELECT x.*, y.*
1122 FROM xs AS x, ys AS y
1123 WHERE y.id = x.id AND x.start < y.end AND y.start < x.end
1124
1124 // Query2: ovWithId implemented as "Union of band joins" in SQL
1125 SELECT x.*, y.*
1126 FROM xs AS x JOIN ys AS y ON y.start < x.start AND x.start < y.end
1127 WHERE y.id = x.id
1128 UNION ALL
1129 SELECT x.*, y.*
1130 FROM xs AS x JOIN ys AS y ON x.start <= y.start AND y.start < x.end
1131 WHERE y.id = x.id

```

Fig. 6. A variation of the event-overlap example. `ovWithId(xs, ys)` computes the same function as the two SQL queries on inputs `xs` and `ys` which are sorted lexicographically by `(id, start, end)`.

Therefore, the formulation of Synchrony generator and the monotonicity and anti-monotonicity conditions on the associated `isBefore` and `canSee` predicates add conceptual elegance and algorithmic clarity in characterizing and generalizing the merge join.

For a further appreciation of what this brings, consider a slight variation of the event-overlap example. As shown in Figure 6, this time, events are categorized by their `id` attribute (where there can be many events of each category), and are ordered lexicographically by their `id`, `start`, and `end` attributes. An event `y` is now considered before another event `x` either when `y` has a smaller category `id` than `x`, or they have the same category `id` and `y` starts before `x`, or they have the same category `id` and start together but `y` ends earlier than `x`, as defined by the function `isBeforeWithId` in the figure. Similarly, two events now are considered overlapping only when they have the same category `id` and they overlap in time, as defined by the function `overlapWithId` in the figure. The function `ovWithId(xs, ys)` returns the overlapping same-category events in `xs` and `ys`. It has time complexity  $O(|xs| + k|ys|)$  if each event in `ys` overlaps fewer than  $k$  same-category events in `xs`, as per the time complexity of Synchrony generator.

The direct translation of `ovWithId(xs, ys)` into SQL is given as `query1` in Figure 6. Notice that it does not meet the syntactic requirement of a band join. So, a relational database system has to execute it using nested loops, resulting in  $O(|xs| \cdot |ys|)$  time complexity. If

the relational database system supports a “single inequality” variant of the merge join, it can cut the time complexity by half; but this is still quadratic.

As  $\text{start} < \text{end}$  holds for any event, it can be shown that `overlapWithId(y, x)` if and only if  $y.\text{id} = x.\text{id}$  and either  $y.\text{start} < x.\text{start} < y.\text{end}$  or  $x.\text{start} \leq y.\text{start} < x.\text{end}$ . So, an alternative SQL query can use the “union of two band joins” idea of Dignoes *et al.* (2021) to implement the same-category event-overlap function. This is `query2` in Figure 6.

While there are different implementations of band join, their time complexity is lower bounded by output size. Thus, optimistically, the time complexity of each of the two band join is  $O(|x_s| + k|y_s|)$  if each event in  $y_s$  overlaps fewer than  $k$  events in  $x_s$ . As there are two band joins and one union, the time complexity is  $O(2|x_s| + 2k|y_s|)$ , assuming the result of the second band join is directly concatenated to the first. Some implementations of band join do not support equality predicate; e.g., Dignoes *et al.* (2021) had to modify Postgres to make its band join support an equality predicate. In this case,  $x_s$  and  $y_s$  have to be re-sorted using only their `start` and `end` attributes and the selectivity  $k'$  must now include overlaps of events in different categories (so,  $k' > k$ ). Then the time complexity becomes worse.

It is worth remarking that, as demonstrated by Dignoes *et al.* (2021), the “union of two band joins” idea is the current state of art in implementing interval join in relational database systems research. The Synchrony generator implementation `ovWithId(x_s, y_s)` has time complexity  $O(|x_s| + k|y_s|)$  when the selectivity is  $k$ , which compares favourably to  $O(2|x_s| + 2k|y_s|)$ . Importantly, it works directly on the `overlapWithId(y, x)` predicate (and any other antimonotonic predicates); whereas, for a relational database system, a user-programmer has to be skilled enough to recast an interval join to the more optimizer-friendly “union of two band joins.” Another useful virtue is that the result of `syncGen(x_s, y_s)` is in the same order as  $x_s$ , while the result produced by the “union of two joins” has lost this ordering. Thus, if the result is to be used as an input to a subsequent query (see the arranging-meeting example in Figure 7), the former might be usable directly; whereas, the latter might require extra sorting effort.

### 5.3 Synchronized iteration on multiple collections

Synchrony fold and derivatives described earlier are synchronizing iterations on two collections. How about synchronizing iterations on three or more collections using these functions? Consider a user-programmer writing a program `mtg0(ws, xs, ys, zs)` for finding the common overlaps between four collections of events. If  $w_s$ ,  $x_s$ ,  $y_s$ , and  $z_s$  are the available time slots of four people, then `mtg0(ws, xs, ys, zs)` are the time slots they are available to meet together.

A naive definition for `mtg0` in comprehension syntax, aiming at clarity, is given first in Figure 7. While `mtg0` is easy to understand, its quartic time complexity begs for improvement. A quick improvement is to insert some `overlap` predicates to eliminate nonoverlapping time slots as early as possible, as done by `mtg1` in Figure 7. If each available time slot of a person overlaps fewer than  $k$  time slots of another person, the time complexity of `mtg1` is quadratic, viz.  $O(|w_s|(|x_s| + k|y_s| + k^2|z_s| + k^3))$ . This is still not very efficient. So, `mtg2` in Figure 7 is an attempt using Synchrony generator to obtain a more efficient implementation. It also makes use of a nice idea on parallel comprehension-cum-monadic zip (Gibbons, 2016).

```

1197 def mtg0(ws: Vec[Event], xs: Vec[Event], ys: Vec[Event], zs: Vec[Event]
1198 ): Vec[Event] =
1199   for (
1200     w <- ws; x <- xs; y <- ys; z <- zs;
1201     s = max(w.start, x.start, y.start, z.start);
1202     e = min(w.end, x.end, y.end, z.end);
1203     if s < e
1204   ) yield Event(start = s, end = e, id = w.id + x.id + y.id + z.id)
1205
1206 def mtg1(ws: Vec[Event], xs: Vec[Event], ys: Vec[Event], zs: Vec[Event]
1207 ): Vec[Event] =
1208   for (
1209     w <- ws;
1210     x <- xs; if overlap(x, w);
1211     y <- ys; if overlap(y, w);
1212     z <- zs; if overlap(z, w);
1213     s = max(w.start, x.start, y.start, z.start);
1214     e = min(w.end, x.end, y.end, z.end);
1215     if s < e
1216   ) yield Event(start = s, end = e, id = w.id + x.id + y.id + z.id)
1217
1218 def mtg2(ws: Vec[Event], xs: Vec[Event], ys: Vec[Event], zs: Vec[Event]
1219 ): Vec[Event] = {
1220   // Requires: ws, xs, ys, and zs sorted by (start, end).
1221   // Issue: The first four lines of codes below
1222   //         breaks the first-order restriction.
1223   val wxss = syncGenGrp(isBefore, overlap)(ws, xs)
1224   val wyss = syncGenGrp(isBefore, overlap)(ws, ys)
1225   val wzss = syncGenGrp(isBefore, overlap)(ws, zs)
1226   val wxyzs = zip3(wxss, wyss, wzss)
1227
1228   for (
1229     (wxs, wys, wzs) <- wxyzs;
1230     (w, xss) = wxs;
1231     (_, yss) = wys;
1232     (_, zss) = wzs;
1233     x <- xss; y <- yss; z <- zss;
1234     s = max(w.start, x.start, y.start, z.start);
1235     e = min(w.end, x.end, y.end, z.end);
1236     if s < e
1237   ) yield Event(start = s, end = e, id = w.id + x.id + y.id + z.id)
1238 }

```

Fig. 7. The arranging-meeting example.

Assuming  $ws$ ,  $xs$ ,  $ys$ , and  $zs$  are sorted lexicographically based on start and end point of events, and all events overlap fewer than  $k$  other events, the time complexity of `mtg2(ws, xs, ys, zs)` is linear,  $O(|ws| + 2k|xs| + |ws| + 2k|ys| + |ws| + 2k|zs| + 2|ws| + k^3|ws|) = O((k^3 + 5)|ws| + 2k(|xs| + |ys| + |zs|))$ . Note that the  $5|ws|$  overheads are due to (1) zipping `wxss`, `wyss`, and `wzss`; (2) scanning `wxyzs`; and (3) scanning `ws` three times when synchronizing with `xs`, `ys`, and `zs`. Nonetheless, this is much better than the quartic time complexity of `mtg0` and quadratic time complexity of `mtg1`, albeit it will be further improved in the next Section when Synchrony iterator is introduced.

Associated with the  $5|ws|$  overheads is also the need to construct and store the intermediate collections `wxss`, `wyss`, `wzss`, and `wxyzs`, as Scala constructs these eagerly. Moreover, `wxss`, `wyss`, `wzss`, and `wxyzs` are nested collections; this breaks the first-order restriction. In

1243 addition, the need to scan `ws` three times may be an issue when `ws` is a large data stream, as  
 1244 it implies needing to buffer the whole stream in memory. In a lazy programming language,  
 1245 this issue may go away, depending on how clever its garbage collector is.

1246 Another somewhat unsatisfactory issue is the need for defining the function `zip3` to com-  
 1247 bine the three sequences of synchronizations of `xs`, `ys`, and `zs` to `ws`. We already know from  
 1248 the limited-mixing lemmas of Section 3 that `zip` is not efficiently definable under the first-  
 1249 order restriction, even though it is a straightforward two-line recursive function. And when  
 1250 the calendars of more people have to be synchronized, a zoo of `zip4`, `zip5`, etc. have to be  
 1251 written as well.<sup>7</sup> This issue is attributable to Scala’s unsophisticated treatment of compre-  
 1252 hension syntax. In a programming language (e.g., Haskell) which has more powerful ways  
 1253 to compile comprehension syntax using alternative binding semantics as well as enhance-  
 1254 ments to comprehension syntax design (Marlow *et al.*, 2016; Wadler & Peyton Jones,  
 1255 2007; Gibbons, 2016; Lindley *et al.*, 2011), this issue of breaking the first-order restriction  
 1256 will likely disappear, though the  $5|ws|$  time-complexity overheads highlighted earlier will  
 1257 likely remain.

## 1258 6 Synchrony iterator

1259 Recall again the motivating example from Figure 1, `ov1(xs, ys) = for (x <- xs; y <- ys;`  
 1260 `if overlap(x, y) yield (x, y)`. Besides its poor quadratic time complexity which has  
 1261 already been highlighted, it suffers from another problem. If `vec[.]` is a streaming data  
 1262 type, i.e. `xs` and `ys` are dynamic data generated continuously as events in them take place,  
 1263 then `ov1(xs, ys)` has to buffer all of `ys` and cannot move on to the second item in `xs` until  
 1264 the data stream `ys` is finished.

1265 Our `syncFold` and `syncFoldGrp`—and thus, `syncMap`, `syncFlatMap`, `syncGen`, and  
 1266 `syncGenGrp`—do not suffer this same problem because, by Antimonotonicity Condition 2,  
 1267 they can move on to the next item in `xs` as soon as the current item in `ys` is after the current  
 1268 item in `xs` and cannot see the item. So, Synchrony fold and its derivatives do not have to  
 1269 buffer for all of `ys`. Nonetheless, the definitions of Synchrony fold in Section 4.1 and 4.2  
 1270 do not produce any output until all items in `xs` and `ys` have been processed. As the actual  
 1271 processing by Synchrony fold only needs to see a small chunk of the two data streams at  
 1272 a time to compute the result for this small chunk, it is desirable to be able to return results  
 1273 incrementally in an on-demand manner.

1274 A possible solution is using the relationship between `foldLeft` and `foldRight` to derive,  
 1275 from `syncGenGrp`, a lazy version `syncGenGrpLazy` where its result type is a `LazyList`.<sup>8</sup> While  
 1276 this is sufficient for getting a streaming version of Synchrony generator, `syncGenGrpLazy`  
 1277 has similar issues as `syncGenGrp` when it comes to synchronizing multiple collections. The

1278 <sup>7</sup> The reader may find this `zip` issue confusing. Are we not already using recursion and other features when we  
 1279 define Synchrony fold and generator? Why are we complaining about having to define `zip3` in `mtg2`? Recall,  
 1280 in this paper, we separate an implementer-programmer who implements programming constructs and library  
 1281 functions from a user-programmer who uses these. The former has access to all features of Scala. The latter, in  
 1282 the context of this paper, is restricted to first-order Scala plus specifically permitted library functions which the  
 1283 former provides. As `mtg2` is an example of how to use Synchrony generator, it is expected to be written by the  
 1284 user-programmer. The user-programmer, being restricted to first-order Scala, thus cannot define an efficient  
 1285 `zip`. So, this user-programmer will have to write a much clumsier-looking program than `mtg2` for efficiency’s  
 1286 sake; the clumsier-looking but efficient program is such an eyesore that we decided to omit it from this paper.

1287 <sup>8</sup> In Scala, items in a `LazyList` are computed only when they are needed and are memoized.

arranging-meeting example of Figure 7, for instance, would have exactly the same implementation using `syncGenGrpLazy` as the version using `syncGenGrp`, viz. `mtg2`, but with every occurrence of `syncGenGrp` replaced by `syncGenGrpLazy`. In particular, a user-programmer implementing it would also be required to write the functions for `LazyList` version of `zip3`, `zip4`, etc. depending on the number of people required for the meeting, as a generic `zip` function for zipping an arbitrary number of collections of arbitrary different types cannot be assigned a valid type in current strongly typed programming languages.

As another mechanism for incrementally producing items on demand, an iterator comes to mind. A normal Scala iterator `yi` provides a `yi.next()` method which on-the-fly computes and produces the next item in the iteration. Although this is simple, we decided against it. The reason is `yi.next()` is iterating only on one collection. A user-programmer would thus be forced to organize the synchronization with the other collections using some additional mechanism, and we would be back to square one.

These two problems—viz. streaming and multi-collection synchronized iteration—are addressed in this section. In particular, *Synchrony iterator* is conceived in this section. A Synchrony iterator `yi = new EIterator(ys, isBefore, canSee)` provides a `yi.syncedWith(x)` method that on-the-fly computes and produces the items in the iteration on `ys` that should be synchronized to (i.e. can see) the item `x`, under the assumption that successive invocations of `syncedWith` are given, as the values of `x`, successive items of a collection `xs` ordered such that `isBefore` is monotonic with respect to  $(xs, ys)$  and `canSee` is antimonotonic with respect to `isBefore`.

This design of Synchrony iterator has two advantages over the standard iterator. Firstly, a nice byproduct of this design is that the same `x` can be used to synchronize multiple Synchrony iterators simultaneously. Using this alternative way to express multi-collection synchronized iteration avoids the `zip` issue mentioned in the discussion on `mtg2`. Secondly, like iterators in general, Synchrony iterator requires side effects. However, unlike the standard iterator, some safe-use conditions can be provided on Synchrony iterator. These conditions isolate these side effects and are sufficient for restoring transparent equational reasoning for programs involving Synchrony iterator.

### 6.1 Designing Synchrony iterator

Deriving a version of Synchrony generator that incrementally computes and returns its result is a fairly typical programming problem. So, we give it a try first by looking at the definition of `syncGenGrp`. The codes for the Synchrony generator `syncGenGrp`, after unfolding through `syncGenGrp(bf, cs)(xs, ys) = syncFoldGrp((x, zs, a) => a :+ (x, zs), Vec(), bf, cs)(xs, ys)`, are shown in the top half of Figure 8.

It is quite apparent that a simple rearrangement of the `aux` function used in defining `syncGenGrp` is sufficient to make it return one element of the result at a time. This is shown in the bottom half of Figure 8. In this rearrangement, the `EIterator` class is introduced. Objects of this class are called *e iterators* (pronounced “iterators.”) An eiterator `yi = new EIterator(ys, isBefore, canSee)` can be regarded as an *enhanced* iterator on the collection `ys`. The eiterator is characterized by a method `yi.syncedWith(x)`, which is the rearranged `aux` function from `syncGenGrp`.

```

1335 // The codes for syncGenGrp after unfolding through
1336 //   syncGenGrp(bf, cs)(xs, ys) =
1337 //   syncFoldGrp((x, zs, a) => a :+ (x, zs), Vec(), bf, cs)(xs, ys).
1338 def syncGenGrp[A,B]
1339   (bf: (B,A) => Boolean, cs: (B,A) => Boolean)
1340   (xs: Vec[A], ys: Vec[B])
1341 : Vec[(A,Vec[B])] = {
1342
1343   def aux(xs: Vec[A], ys: Vec[B], zs: Vec[B], acc: Vec[(A,Vec[B])])
1344   : Vec[(A,Vec[B])] = {
1345     if (xs.isEmpty) acc
1346     else if (ys.isEmpty && zs.isEmpty) acc
1347     else if (ys.isEmpty) aux(xs.tail, zs, Vec(), acc :+ (xs.head, zs))
1348     else {
1349       val (x,y) = (xs.head, ys.head)
1350       (bf(y, x), cs(y, x)) match {
1351         case (true, false) => aux(xs, ys.tail, zs, acc)
1352         case (false, false) => aux(xs.tail, zs ++: ys, Vec(), acc :+ (x,zs))
1353         case (_, true) => aux(xs, ys.tail, zs :+ y, acc)
1354       }
1355     }
1356   }
1357
1358   aux(xs, ys, Vec(), Vec())
1359 }
1360
1361 // Rearranging syncGenGrp's aux function to return one element
1362 // of the result at a time. This provides a preliminary
1363 // implementation of Synchrony iterator.
1364 class EIterator[A,B](
1365   elems: Vec[B],
1366   bf: (B,A)=>Boolean, cs:(B,A)=>Boolean) {
1367
1368   private var es = elems
1369
1370   def syncedWith(x: A): Vec[B] = {
1371     def aux(zs: Vec[B]): Vec[B] = {
1372       if (es.isEmpty && zs.isEmpty) zs
1373       else if (es.isEmpty) { es = zs; zs }
1374       else {
1375         val y = es.head
1376         (bf(y, x), cs(y, x)) match {
1377           case (true, false) => { es = es.tail; aux(zs) }
1378           case (false, false) => { es = zs ++: es; zs }
1379           case (_, true) => { es = es.tail; aux(zs :+ y) }
1380         }
1381       }
1382     }
1383     aux(Vec())
1384   }
1385 }

```

Fig. 8. Preliminary definition of EIterator, shown along side the unfolded definition of syncGenGrp. The syncedWith method of the former is derived from the aux function of the latter.

The theorem below shows that when  $yi = \text{new EIterator}(ys, \text{isBefore}, \text{canSee})$  is a fresh iterator on  $ys$ , calling  $yi.\text{syncedWith}(x)$  on each successive item  $x$  in  $xs$ , returns the corresponding successive item  $(x, zs)$  in  $\text{syncGenGrp}(\text{isBefore}, \text{canSee})(xs, ys)$ . Furthermore, the total time complexity is the same. Thus, an iterator generates—at the same efficiency—the same items produced by a Synchrony generator, and it produces these

items one at a time when its `syncedWith` method is called iteratively. For this reason, an iterator is called a *Synchrony iterator*.

**Theorem 6.1.** *Suppose `isBefore` is monotonic with respect to  $(xs, ys)$ , and `canSee` is antimonotonic with respect to `isBefore`. Then, the following two programs define the same function.*

1. `syncGenGrp(isBefore, canSee)(xs, ys)`
2. `val yi = new EIterator(ys, isBefore, canSee)`  
`for (x <- xs; zs <- yi.syncedWith(x)) yield (x, zs)`

*Both programs have time complexity  $O(|xs| + 2k|ys|)$ , assuming `isBefore` and `canSee` have constant time complexity and each item in `ys` can see fewer than  $k$  items in `xs`.*

**Proof** When an iterator `yi = new EIterator(ys, isBefore, canSee)` is freshly created, its internal variable `es` is initialized to the collection `ys`. Let the items in `xs` be  $x_1, \dots, x_n$ , in this order. Suppose `isBefore` is monotonic with respect to  $(xs, ys)$ , and `canSee` is antimonotonic with respect to `isBefore`. Suppose `yi.syncedWith(x1)`, ..., `yi.syncedWith(xn)` are called in this order. Let `zs1`, ..., `zsn` be the corresponding results. Let `es1`, ..., `esn` be the value of the internal variable `es` at the end of each of these calls. And let `e0 = ys`.

By construction, for each `y` in `ys`, such that `isBefore(y, xj)`, it is the case `canSee(y, xj)` if and only if `y` is in `zsj` and `esj`. Also, by construction, for each `y` in `ys`, such that not `isBefore(y, xj)`, it is the case that `y` is in `esj`; and by Antimonotonicity Condition 2, `y` is in `zsj` if and only if `canSee(y, xj)`. Thus, for `y` in `ys`, `y` is in `zsj` if and only if `canSee(y, xj)`. So,

```
zsj = for (y <- ys; if canSee(y, xj)) yield y
```

```
Vec((xj, zsj)) = syncGenGrp(isBefore, canSee)(Vec(xj), ys)
```

Then, the first part of the theorem follows from Corollary 4.6,

```
syncGenGrp(isBefore, canSee)(xs, ys)
= for (x <- xs; zs <- yi.syncedWith(x)) yield (x, zs)
```

Looking at the definition of the function `aux` in `syncedWith`, when processing `yi.syncedWith(xj)`, each time `aux` is called, it reduces the number of items in `esj-1` (and thus `ys`) by 1; or it increases the number of items by  $|zs<sub>j</sub>|$  exactly once, when it returns. As mentioned earlier, the items in `zsj` are those `y` that can see `xj`. Thus, the total number of times `aux` gets called when processing `yi.syncedWith(x1)`, ..., `yi.syncedWith(xn)`, is  $|ys| + \sum_j |zs<sub>j</sub>|$ . By assumption of the theorem, each item in `ys` can see fewer than  $k$  items in `xs`. So, each item in `y` appears in fewer than  $k$  distinct `zsj`. Thus,  $\sum_j |zs<sub>j</sub>| < k|ys|$ . Also, the prepend operator `zs ++: es` is linear in  $|zs|$ ; these add an overhead of  $\sum_j |zs<sub>j</sub>| < k|ys|$ . So, `for (x <- xs; zs <- yi.syncedWith(x)) yield (x, zs)` has time complexity  $O(|xs| + 2k|ys|)$ . This proves the second part of the theorem. ■

The following useful details can also be extracted from the proof above.

**Proposition 6.2.** *Suppose `isBefore` is monotonic with respect to  $(xs, ys)$ , and `canSee` is antimonotonic with respect to `isBefore`. Let the iterator `yi = new EIterator(ys, isBefore, canSee)` be freshly created. Let  $x_1, \dots, x_n$  be some*

```

1427 def mtg3(
1428   ws: Vec[Event], xs: Vec[Event], ys: Vec[Event], zs: Vec[Event]
1429 ): Vec[Event] = {
1430   // Requires: ws, xs, ys, zs sorted lexicographically by (start, end).
1431   // Note: isBefore and overlap are as defined in Figure 1.
1432   val xi = new EIterator(xs, isBefore, overlap);
1433   val yi = new EIterator(ys, isBefore, overlap);
1434   val zi = new EIterator(zs, isBefore, overlap);
1435   for (
1436     w <- ws;
1437     x <- xi.syncedWith(w);
1438     y <- yi.syncedWith(w);
1439     z <- zi.syncedWith(w);
1440     s = max(w.start, x.start, y.start, z.start);
1441     e = min(w.end, x.end, y.end, z.end);
1442     if s < e
1443   ) yield Event(start = s, end = e, id = w.id + x.id + y.id + z.id)
1444 }

```

Fig. 9. The arranging-meeting example expressed using Synchrony iterator.

of the items in  $xs$ , with possible repetitions and omissions of items in  $xs$ , such that for  $1 \leq j < j' \leq n$  where  $x_j \neq x_{j'}$ , it is the case that  $(x_j \ll x_{j'} \mid xs)$ . Suppose  $y_i.\text{syncedWith}(x_j)$  is called in sequence for each  $x_j$ . Let  $zs_j$  be the corresponding result and  $es_j$  be the value of the internal variable  $es$  of the iterator  $y_i$  at the end of each of these calls. And let  $e_0 = ys$ . Then,

$$\begin{aligned}
 zs_j &= \{ \text{new EIterator}(es_{j-1}, \text{isBefore}, \text{canSee}) \}.\text{syncedWith}(x_j) \\
 &= \{ \text{new EIterator}(ys, \text{isBefore}, \text{canSee}) \}.\text{syncedWith}(x_j) \\
 &= \text{for } (y \leftarrow ys; \text{if } \text{canSee}(y, x_j)) \text{ yield } y
 \end{aligned}$$

That is, only the ordering of  $x_j$  matters when calling `syncedWith`; repetitions and omissions of items in  $xs$  have no impact.

The design of Synchrony iterator thus meets the objective of incrementally computing and producing synchronized items in a collection  $ys$  to items in a collection  $xs$ .

Fortuitously, the synchronization provided by Synchrony iterator is specified via  $y_i.\text{syncedWith}(x)$ ; i.e.,  $xs$  does not need to be given as part of the specification. This design facilitates the simultaneous synchronization of items in multiple collections to the same item  $x$ . In particular, let  $y_{i_1}, \dots, y_{i_n}$  be iterators on the  $n$  collections  $ys_1, \dots, ys_n$  that are to be synchronized to items in  $xs$ . Then for each item  $x$  in  $xs$ , the methods  $y_{i_1}.\text{syncedWith}(x)$ ,  $\dots$ ,  $y_{i_n}.\text{syncedWith}(x)$  are called to achieve simultaneous synchronized iteration on the  $n$  collections to the collection  $xs$ , like this:

```

1463 val yi1 = new EIterator(ys1, bf1 cs1); ...; val yin = new EIterator(ysn, bfn csn);
1464 for (x <- xs; y1 <- yi1.syncedWith(x); ...; yn <- yin.syncedWith(x)) yield ...

```

The function `mtg3` in Figure 9, which revisits the arranging-meeting example from Section 5.3, illustrates this simultaneous synchronization. Notice that `mtg3` is first order. And, in contrast to the approach adopted earlier by `mtg2` from Figure 7, `mtg3` dispenses with the need to define a zoo of `zip`'s to structure synchronized iteration in multiple collections.

However, the time complexity of `mtg3` may not be as good as `mtg2`. Suppose all events overlap fewer than  $k$  other events. The time complexity of `mtg3` is  $O(|ws| +$

1473  $2k|xs| + 2k^2|ys| + 2k^3|zs| + k^3|ws|$ ), because  $x \leftarrow xi.syncedWith(w)$  is called once for each  
 1474  $w$ ,  $y \leftarrow yi.syncedWith(w)$  is called once for each  $x$ , and  $z \leftarrow zi.syncedWith(w)$  is called once  
 1475 for each  $y$ .

1476 Fortunately, although  $y \leftarrow yi.syncedWith(w)$  and  $z \leftarrow zi.syncedWith(w)$  are called mul-  
 1477 tiple times for different  $x$ 's and  $y$ 's respectively, these calls depend on  $w$  and not on  $x$  and  $y$ .  
 1478 So, the problem is easy to solve by making `syncedWith` remember its immediate last result.  
 1479 Figure 10 shows the revised `EIterator` class incorporating this simple solution. We have so  
 1480 far used `Vec[.]` to denote a collection type. Perhaps this gives the appearance that  $xs$  and  $ys$   
 1481 are collection types of the same kind, i.e., both are vectors, both are lists, etc. Actually, this  
 1482 does not need to be the case. So, we show this in this revised version of `EIterator` as well.  
 1483 Specifically, `yi = new EIterator(ys, bf, cs)` now constructs an iterator for any kind of  
 1484 iterable object  $ys$ , e.g., a `LazyList[B]`, which is Scala's preferred data type for representing  
 1485 data streams. Also, the method `yi.syncedWith(x)` now returns a `List[B]`. And the collec-  
 1486 tion  $xs$ , where  $x$  comes from, can be yet another kind of collection type, e.g., a `Vector[A]`.  
 1487 Incidentally, the `prepend ++:` and `postpend :+` operations on vectors, are not needed in this  
 1488 version of `EIterator`.

1489 With this simple modification to Synchrony iterator, the time complexity of `mtg3`  
 1490 becomes  $O((k^3 + 1)|ws| + 2k(|xs| + |ys| + |zs|))$ . Now, `mtg3` is even more efficient than  
 1491 `mtg2`, successfully reducing the latter's  $5|ws|$  overheads to  $|ws|$ , as well as avoiding the con-  
 1492 struction of several large intermediate collections. Moreover, even when  $ws$ ,  $xs$ ,  $ys$ , and  $zs$   
 1493 are large dynamic data streams, `mtg3` can produce their common time slots incrementally  
 1494 as overlapping events arrive.

1495 It is worth diving deeper into the details of the revised definition of `EIterator` in  
 1496 Figure 10. `EIterator` memoizes the previous result in `ores` and the previous value of  $x$   
 1497 in `ox`. If the next value of  $x$  is same as the one memoized earlier in `ox`, the result memoized  
 1498 earlier in `ores` is returned immediately. Otherwise, the synchronized iteration resumes from  
 1499 `ores` and continues onward to `es`. This actually kills a second bird with the same stone: In  
 1500 the earlier definition of Synchrony iterator in Figure 8, when both `bf(y, x)` and `cs(y, x)`  
 1501 are false,  $zs$  must be prepended back to `es` before returning  $zs$  as the result. This prepending  
 1502 step is dispensed with in this revised definition of `EIterator` as the result is now already  
 1503 memoized in `ores` and the iteration in response to the next call value  $x$  resumes from `ores`  
 1504 before continuing onward to `es`.

1505 Under the hood in Scala, being a `List[.]`, `ores` is a “boxed value”; i.e., it is a pointer.  
 1506 Thus, if there are multiple consecutive  $x$ 's which have the same value, the corresponding  
 1507 `yi.syncedWith(x)` results are exactly the same pointer. This has a rather nice practical impli-  
 1508 cation, akin to factorized databases (Olteanu & Schleich, 2016). As an illustration, let the  
 1509 collection  $xs$  be just a sequence repeating the same value  $u$ , and the collection  $ys$  be just a  
 1510 sequence repeating the same value  $v$ . Suppose also that `cs(v, u)` is true. Then, it does not  
 1511 matter whether `bf(v, u)` is true, `for (x <- xs; zs = yi.syncedWith(x)) yield (x, zs)` has  
 1512 linear physical size  $O(|xs| + |ys|)$ , even though—semantically—there are  $|xs| \cdot |ys|$  number  
 1513 of items from  $xs$  and  $ys$  in it. Although not explored here, this property may be further  
 1514 exploited for designing more efficient algorithms, e.g., for database query processing,  
 1515 perhaps in the manner of Henglein & Larsen (2010) and Olteanu & Schleich (2016).

1516 Also, in practice, `isBefore(y, x)` and `canSee(y, x)` predicates do not use all the infor-  
 1517 mation in  $y$  and  $x$ . In fact, they often have a form like `bf(y, x) = bfk( $\psi(y)$ ,  $\phi(x)$ )` and

1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511  
1512  
1513  
1514  
1515  
1516  
1517  
1518

```

1519 class EIterator[A,B](
1520   elems: Iterable[B],
1521   bf: (B,A)=>Boolean, cs: (B,A)=>Boolean)
1522 {
1523   private var es: Iterable[B] = elems
1524   private var ores: List[B] = List() // last result
1525   private var ox: Option[A] = None // last x
1526
1527   // When iterating, use items in ores before items in es.
1528   private def empty = es.isEmpty && ores.isEmpty
1529   private def hd = if (ores.isEmpty) es.head else ores.head
1530   private def nx() = if (ores.isEmpty) { es = es.tail }
1531                     else { ores = ores.tail }
1532
1533   def syncedWith(x: A): List[B] = {
1534     def aux(zs: List[B]): List[B] =
1535       if (empty) { zs }
1536       else {
1537         val y = hd
1538         (bf(y, x), cs(y, x)) match {
1539           case (true, false) => { nx(); aux(zs) }
1540           case (false, false) => { zs }
1541           case (_, true) => { nx(); aux(y +: zs) }
1542         }
1543       }
1544     // Use the last result if this x is same as the last x
1545     if (ox == Some(x)) { ores }
1546     else { ox = Some(x); ores = aux(List()).reverse; ores }
1547   }
1548 }
1549
1550 class EIteratorWithKey[KA,KB,A,B](
1551   keya: A => KA, keyb: B => KB,
1552   elems: Iterable[B],
1553   bfk: (KB,KA)=>Boolean, csk: (KB,KA)=>Boolean)
1554 extends EIterator[A,B](null, null, null)
1555 {
1556   // EIterator ek for synchronizing elems to keya(x) instead of x.
1557   private val ek: EIterator[KA,B] = {
1558     val bf = (y: B, kx: KA) => bfk(keyb(y), kx)
1559     val cs = (y: B, kx: KA) => csk(keyb(y), kx)
1560     new EIterator(elems, bf, cs)
1561   }
1562
1563   // Override syncedWith(x) by ek.syncedWith(keya(x)).
1564   // This is equivalent to defining the isBefore and canSee
1565   // predicates for EIteratorWithKey as:
1566   //   bf(y, x) = bfk(keyb(y), keya(x))
1567   //   cs(y, x) = csk(keyb(y), keya(x))
1568   override def syncedWith(x: A): List[B] = ek.syncedWith(keya(x))
1569 }

```

Fig. 10. Revised definition of Synchrony iterator `EIterator` and its derivative `EIteratorWithKey`, whose `isBefore` predicate (`bfk`) and `canSee` predicate (`csk`) are defined using sorting keys (`keya`, `keyb`).

$cs(y, x) = csk(\psi(y), \phi(x))$ , where  $\psi(\cdot)$  and  $\phi(\cdot)$  are some sorting keys of  $ys$  and  $xs$  respectively. As  $xs$  is sorted by  $\phi(\cdot)$ , for any  $(x_i \ll x_j \mid xs)$  where  $\phi(x_i) = \phi(x_j)$ , it is the case that  $\phi(x_i) = \phi(x_k) = \phi(x_j)$  for all  $(x_i \ll x_k \ll x_j \mid xs)$ ; this is so even when  $x_i \neq x_k \neq x_j$ . This means  $yi.syncedWith(x_i) = yi.syncedWith(x_k) = yi.syncedWith(x_j)$ ,

1565 assuming `yi.syncedWith(xi)`, `yi.syncedWith(xk)`, and `yi.syncedWith(xj)` are called in  
 1566 this sequence. However, when  $x_i \neq x_k \neq x_j$ , `yi.syncedWith(xi)`, `yi.syncedWith(xk)`, and  
 1567 `yi.syncedWith(xj)` would be pointers to three separate physical lists comprising exactly  
 1568 the same sequence of items from `ys`. To avoid this situation, instead of memoizing  
 1569 the argument `x`, the `syncedWith` method of `EIterator` should memoize  $\phi(x)$ . In Scala,  
 1570 this can be accomplished by defining a subclass `EIteratorWithKey` of `EIterator`, where  
 1571 `EIteratorWithKey` redefines `syncedWith(x)` to `syncedWith( $\phi(x)$ )`, as shown in Figure 10. Then,  
 1572 instead of creating an iterator by `yi = new EIterator(ys, bf, cs)`, it can be created as  
 1573 `yi = new EIteratorWithKey( $\phi(\cdot)$ ,  $\psi(\cdot)$ , ys, bfk, csk)`.

## 1574 6.2 Safe use of Synchrony iterator

1575 Synchrony iterator is defined using side effects. Each time `syncedWith` is invoked on an  
 1576 iterator, the local variable `es` and its local result cache `ores` and `ox` are updated. This can  
 1577 make a program difficult to understand when Synchrony iterator is used in an undisciplined  
 1578 way. Therefore, the following conditions are imposed to ensure better discipline in using  
 1579 Synchrony iterator. To specify these conditions, the notation  $\mathcal{F}[\cdot]$  denotes an expression  
 1580 with a “hole”—called a “context”—and  $\mathcal{F}[e]$  denotes the same expression but with the  
 1581 expression `e` substituted into the hole.  
 1582

1583  
 1584 **Definition 6.3** (Safe-use). *The following conditions are presumed to hold on a program  
 1585 for each expression `yi.syncedWith(x)` that appears in the program.*

- 1586 1. *There is a collection `xs`, and `x` takes successive values in `xs`. That is, the expression  
 1587 `yi.syncedWith(x)` appears in an enclosing expression that binds `x` to the collec-  
 1588 tion `xs`. In general, the enclosing expression  $\mathcal{C}[\text{yi.syncedWith}(x)]$  looks like, or gets  
 1589 desugared into, one of these forms:*

1590 `xs flatMap (x =>  $\mathcal{F}[\text{yi.syncedWith}(x)]$ )`

1591 `xs map (x =>  $\mathcal{F}[\text{yi.syncedWith}(x)]$ )`

1592 `xs filter (x =>  $\mathcal{F}[\text{yi.syncedWith}(x)]$ )`

- 1593 2. *`yi` is an iterator on some collection `ys`.*
- 1594 3. *`isBefore` is monotonic with respect to  $(xs, ys)$ .*
- 1595 4. *`canSee` is antimonotonic with respect to `isBefore`.*
- 1596 5. *`yi.syncedWith(x)` produces the same value as `ys filter (y => canSee(y, x))`, though  
 1600 not necessarily with the same efficiency, in the context of this program. That is,  
 1601  $\mathcal{C}[\text{yi.syncedWith}(x)] = \mathcal{C}[\text{ys filter}(y \Rightarrow \text{canSee}(y, x))]$ .*

1602  
 1603  
 1604 It may seem onerous to programmers to have these conditions imposed on them. In  
 1605 reality, they only need to take responsibility for Safe-use Conditions 3 and 4, as these are  
 1606 non-trivial for the compiler to verify automatically in some cases; nonetheless, they are  
 1607 often easy to achieve. The other safe-use conditions are easy for a compiler to check or to  
 1608 enforce in pragmatically, as explained below, or to train programmers to comply with.  
 1609  
 1610

Safe-use Condition 1 is trivial, and can be easily checked and enforced by the compiler. It simply says a Synchrony iterator on a collection  $ys$  should always be used inside the scope of the generator that  $ys$  is synchronized to.

Safe-use Condition 2 is also trivial. It is just standard type checking.

Safe-use Condition 5, though seems non-trivial at first sight, can be achieved in a pragmatic way which can be enforced by the compiler. In fact, only two basic rules are needed. First, if there is another expression  $yi.syncedWith(x')$  on the same iterator  $yi$ , we must have  $x == x'$ . That is, all occurrences of the iterator  $yi$  are identical; i.e. synchronized to the same  $x$  in  $xs$ . Or, better still, insist on  $yi$  to occur only twice in the program, once when the iterator is being created (i.e.  $yi = \text{new EIterator}(ys, \text{isBefore}, \text{canSee})$ ), and once when the iterator is being used for the only time (i.e.  $yi.syncedWith(x)$ ). Second, the iterator  $yi$  should be constructed immediately before the generator of  $xs$ . That is, programmers should always use  $yi.syncedWith(x)$  inside an enclosing expression that looks like, or gets desugared to, one of these forms:

```
val yi = new EIterator(ys, isBefore, canSee)
xs flatMap (x =>  $\mathcal{F}[yi.syncedWith(x)]$ )
```

```
val yi = new EIterator(ys, isBefore, canSee)
xs map (x =>  $\mathcal{F}[yi.syncedWith(x)]$ )
```

```
val yi = new EIterator(ys, isBefore, canSee)
xs filter (x =>  $\mathcal{F}[yi.syncedWith(x)]$ )
```

Even though iterators have side effects that change their state (viz. their local variables  $es$ ,  $ores$ , and  $ox$ ), the two rules under Safe-use Condition 5 isolate these side effects. Without loss of generality, by the second rule, suppose an iterator appears like this:

```
val yi = new EIterator(ys, isBefore, canSee)
xs flatMap (x =>  $\mathcal{F}[yi.syncedWith(x)]$ )
```

By the first rule,  $yi$  appears only in the exact form  $yi.syncedWith(x)$ , whose value depends only on  $x$ . Being in a comprehension,  $x$  takes successive values  $x_j$  of  $xs$ . So, by Proposition 6.2, it is guaranteed that  $yi.syncedWith(x_j) = ys \text{ filter } (y => \text{canSee}(y, x_j))$ . That is,

```
{ val yi = new EIterator(ys, isBefore, canSee);
  xs flatMap (x =>  $\mathcal{F}[yi.syncedWith(x)]$ ) }
= xs flatMap (x =>  $\mathcal{F}[ys \text{ filter } (y => \text{if canSee}(y, x))]$ )
```

In other words, it permits the left-hand-side (which has side effects) to be replaced by the right-hand-side (which has no side effects) when one is reasoning extensionally. Thus, despite its side effects, under the safe-use conditions, one might justifiably claim that Synchrony iterator is a purer programming paradigm than a standard iterator.

Incidentally, the equivalence highlighted above also implies that Synchrony iterator, under the safe-use conditions, is a conservative extension of first-order restricted Scala sans library functions.

**Theorem 6.4.** *The extensional expressive power of Scala under the first-order restriction, is the same with or without Synchrony iterator under the safe-use conditions. However, more efficient algorithms for some functions (e.g., a linear-time algorithm for low-selectivity join) can be defined when Synchrony iterator is made available in this fragment of Scala.*

### 6.3 Referential transparency of Synchrony iterator

The Safe-use conditions of Synchrony iterator assure its referential transparency. A rather attractive implication is that equational reasoning that holds for standard collection types, but fails on standard iterators, holds for Synchrony iterator. This is a direct consequence of Safe-use Condition 5. We illustrate this with the equations for code motion, redundant-code elimination, and parallelism.

#### Code motion

The “code-motion” equation below is valid for standard collection types, provided the free variables of  $e_2$  are a subset of the free variables of  $u \Rightarrow \mathcal{F}[e_2]$ , and  $e_2$  has no observable side effects.

$$e_1 \text{ flatMap } (u \Rightarrow \mathcal{F}[e_2]) \\ = \{ \text{val } v = e_2; e_1 \text{ flatMap } (u \Rightarrow \mathcal{F}[v]) \}$$

The code-motion equation is inapplicable to standard iterators. In contrast, it is applicable to Synchrony iterator under safe-use conditions. Specifically, the following holds:

$$e_1 \text{ flatMap } (u \Rightarrow \mathcal{F}[yi.\text{synchronized}(x)]) \\ = \{ \text{val } v = yi.\text{synchronized}(x); e_1 \text{ flatMap } (u \Rightarrow \mathcal{F}[v]) \}$$

The validity of this code-motion equation is a consequence of Safe-use Condition 5. To wit, assume  $yi = \text{new EIterator}(ys, bf, cs)$  for some  $ys$ ,  $bf$ , and  $cs$ , then proceed as follow.

$$e_1 \text{ flatMap } (u \Rightarrow \mathcal{F}[yi.\text{synchronized}(x)]) \\ = e_1 \text{ flatMap } (u \Rightarrow \mathcal{F}[ys \text{ filter } (y \Rightarrow cs(y, x))]) \\ = \{ \text{val } v = ys \text{ filter } (y \Rightarrow cs(y, x)); e_1 \text{ flatMap } (u \Rightarrow \mathcal{F}[v]) \} \\ = \{ \text{val } v = yi.\text{synchronized}(x); e_1 \text{ flatMap } (u \Rightarrow \mathcal{F}[v]) \}$$

#### Redundant-code elimination

The “redundant-code elimination” equation below is valid for standard collection types, provided the expression  $e$  has no observable side effects.

$$(e \text{ flatMap } f) ++ (e \text{ flatMap } g) \\ = \{ \text{val } v = e; (v \text{ flatMap } f) ++ (v \text{ flatMap } g) \}$$

This redundant-code elimination equation is inapplicable when  $e$  is an expression having an iterator type. In contrast, under safe-use conditions, it is applicable to Synchrony iterator despite its having side effects. Specifically, the following holds:

```

1703     (yi.syncedWith(x) flatMap f) ++ (yi.syncedWith(x) flatMap g)
1704 = { val v = yi.syncedWith(x); (v flatMap f) ++ (v flatMap g) }

```

1705 The validity of this redundant-code elimination is again a consequence of Safe-use  
1706 Condition 5. As before, assume `yi = new EIterator(ys, bf, cs)` for some `ys`, `bf`, and `cs`,  
1707 and proceed as follow.

```

1708     (yi.syncedWith(x) flatMap f) ++ (yi.syncedWith(x) flatMap g)
1709 = (ys filter (y => cs(y, x) flatMap f)) ++ (ys filter (y => cs(y, x) flatMap g))
1710 = { val v = ys filter (y => cs(y, x)); (v flatMap f) ++ (v flatMap g) }
1711 = { val v = yi.syncedWith(x); (v flatMap f) ++ (v flatMap g) }
1712

```

### 1713 *Homomorphism over flatMap*

1714 The “homomorphism” equation below is valid for standard collection types, provided  
1715 expressions `e`, `f`, and `g` have no observable side effects. This equation is the basis for  
1716 parallelization of `flatMap` in, e.g., Hadoop-like platforms.

```

1718     (e ++ f) flatMap g
1719 = (e flatMap g) ++ (f flatMap g)

```

1720 A similar homomorphism equation holds for `syncedWith`. Suppose  
1721 `val yi = new EIterator(us ++ vs, bf, cs)` for some `us`, `vs`, `bf`, and `cs`. Then, after replac-  
1722 ing `val yi = new EIterator(us ++ vs, bf, cs)` by `{ val ui = new EIterator(us, bf, cs);`  
1723 `val vi = new EIterator(vs, bf, cs) }`, the equation below holds.

```

1724     yi.syncedWith(x)
1725 = ui.syncedWith(x) ++ vi.syncedWith(x)

```

1726 This equation follows because

```

1727     yi.syncedWith(x)
1728 = (us ++ vs) filter (y => cs(y, x))
1729 = (us filter (y => cs(y, x))) ++ (vs filter (y => cs(y, x)))
1730 = ui.syncedWith(x) ++ vi.syncedWith(x)

```

1731 This equation offers a simple way to parallelize `syncedWith`.

## 1732 *6.4 Possible syntax for Synchrony iterator*

1733 Considering the safe-use conditions, it is perhaps pertinent to suggest a syntax for  
1734 Synchrony iterator that automatically enforces all the safe-use conditions, apart from the  
1735 safe-use conditions on monotonicity and antimonotonicity. One possibility is to introduce  
1736 the following generator pattern into comprehension syntax:

```

1737     (x, zs1, ..., zsn) <- xs syncWith(ys1, bf1, cs1) ...
1738                               syncWith(ysn, bfn, csn)

```

1739 This way, the `EIterator` class can be hidden from user-programmers, and they can be told  
1740 that `zsj = ys.filter((y) => csj(y, x))` at all times in terms of value, as per Proposition 6.2,  
1741 but is obtained very efficiently.

1742

This generator pattern is compiled by desugaring it to

```

1749
1750 (x, zs1, ..., zsn) <- {
1751     val yi1 = new EIterator(ys1, bf1, cs1); ...;
1752     val yin = new EIterator(ysn, bfn, csn)
1753     for (
1754         x <- xs;
1755         zs1 = yi1.syncedWith(x); ...;
1756         zsn = yin.syncedWith(x);
1757     ) yield (x, zs1, ..., zsn) }

```

Usual “deforestation” rules (Wadler, 1990) should be able to optimize this further to remove the intermediate collection introduced by this desugaring. If not, the generator pattern can also be desugared into the chain of generator and assignment patterns below:

```

1761     yi1 = new EIterator(ys1, bf1, cs1); ...;
1762     yin = new EIterator(ysn, bfn, csn)
1763     x <- xs;
1764     zs1 = yi1.syncedWith(x); ...;
1765     zsn = yin.syncedWith(x);
1766

```

The program `mtg4` in Figure 11 is a rewrite of the program `mtg3` from Figure 9 using this suggested syntax for Synchrony iterator. As can be seen, using this syntax, the three Synchrony iterators `xi`, `yi`, and `zi` that earlier appeared explicitly in `mtg3` are now tucked away from sight. The user-programmer is thus presented with a pure functional comprehension syntax which uses a slightly enhanced generator form.<sup>9</sup>

## 7 Some use-cases and a stress test

Synchrony fold and Synchrony iterator for querying relational databases in general, genomic datasets in particular, and timestamped data streams is discussed here. A stress test on using them on genomic datasets is also presented.

### 7.1 Relational database queries

The use-case of Synchrony fold and Synchrony iterator in the context of relational database querying should be quite clear already. Further technical and theoretical details are given in a companion paper (Wong, 2021). So, here, we just point out that only one extra function is needed to make all relational database queries (including group-by, order-by, and aggregate functions) efficiently implementable in first-order restricted Scala endowed with

<sup>9</sup> This tantalizing syntax is used for illustrative purpose later in Section 7.4. However, in the rest of this work, we eschew using it in favour of the plain `yi = new EIterator(ys, bf, cs)` and `yi.syncedWith(x)` as our iterator constructs. The idea and design of Synchrony fold, Synchrony generator, and Synchrony iterator are partly driven by our desire in suggesting a small set of library functions for general synchronized iteration. The `EIterator` and `EIteratorWithKey` classes and the functions defining Synchrony fold and Synchrony generator have the crucial advantage of being readily copied and adopted for a wide variety of programming languages without modifying any of their compilers. Whereas, introducing new syntax into any programming language faces the obstacle of modifying its compiler, which requires significantly more technical effort (perhaps also requires lots of lobbying); it is thus an unlikely scenario for most programming languages.

```

1795 def mtg4(
1796   ws: Vec[Event], xs: Vec[Event], ys: Vec[Event], zs: Vec[Event]
1797 ): Vec[Event] = {
1798   // Requires: ws, xs, ys, zs sorted lexicographically by (start, end).
1799   // Note: isBefore and overlap are as defined in Figure 1.
1800   for (
1801     (w, wxs, wys, wzs) <- ws syncWith(xs, isBefore, overlap)
1802     syncWith(ys, isBefore, overlap)
1803     syncWith(zs, isBefore, overlap)
1804     x <- wxs; y <- wys; z <- wzs;
1805     s = max(w.start, x.start, y.start, z.start);
1806     e = min(w.end, x.end, y.end, z.end);
1807     if s < e
1808   ) yield Event(start = s, end = e, id = w.id + x.id + y.id + z.id)
1809 }

```

Fig. 11. The arranging-meeting example revisited again. The program `mtg4` is a rewrite of the program `mtg3` from Figure 9 using the generator syntax suggested for Synchrony iterator.

Synchrony fold and Synchrony iterator. That extra function is `sortWith(f)(xs)` which sorts the collection `xs: Vec[A]` using the ordering function `f: (A,A) => Boolean`. This is because Synchrony fold and Synchrony iterator require their input to be suitably sorted beforehand.

Actually, sorting at quadratic time complexity is already expressible in first-order restricted Scala endowed with Synchrony fold. Theorem 4.5 implies that all functions defined using first-order restricted Scala with Synchrony fold has time complexity of the form  $O(m^n)$  where  $m$  is input size. Since efficient sorting requires  $\Omega(m \log(m))$  time in general, this means sorting takes  $\Theta(m^2)$  time in first-order-restricted Scala with Synchrony fold. Thus, it is necessary to provide an efficient `sortWith` sorting function to user-programmers, to ensure that they are able to implement any relational database queries efficiently in this framework.

It is worth remarking that on-disk sorting is much easier to implement than on-disk indexed tables; cf. Silberschatz *et al.* (2016). It is thus a virtue of Synchrony iterator, which needs only the former when processing very large collections, relative to approaches that try to compile join-expressing comprehensions into indexed tables.

## 7.2 Genometric queries

A second use-case is genometric queries on genomic datasets. BEDOPS (Neph *et al.*, 2012) and GMQL (Masseroli *et al.*, 2019) are two notable toolkits for processing these datasets, and support similar query operations. The former via unix-style commands. The latter via a specialized GenoMetric Query Language. The data model is highly constrained in such domain-specific toolkits. GMQL is used here for illustration, modulo some liberty taken with GMQL's syntax.

There are only a few main object types. The first main object type is the genomic region; this is `Bed(chrom: String, start: Int, end: Int, ...)` for a region located on chromosome `chrom`, beginning at position `start`, ending at position `end`, plus some other pieces of information which are omitted here. Regions on the same chromosome are ordered by their `start` and `end` point lexicographically; regions on different chromosomes are ordered by their `chrom` value. This ordering  $<_{\text{Bed}}$  defines the default `isBefore` predicate on regions, viz. `isBefore(y, x)` if and only if  $y <_{\text{Bed}} x$ . The next main object is the genome, which

is a BED file; it is a large text file in the BED format (Neph *et al.*, 2012), the de facto format for this kind of information in the bioinformatics community. A BED file is just a collection of regions, abstracted here as `Vec[Bed]`. The next main object type is the sample, `Sample(bedFile: Vec[Bed], meta: ...)`, which is just a BED file and its associated meta-data. The last main object is the sample database, which is just a collection of samples; it is abstracted here as `Vec[Sample]`.

Queries at the level of samples mainly select samples from the sample database to analyze. Queries at the level of BED files mainly extract and process regions of interest. The first kind of queries are basically simplified relational database queries. The second kind of queries are the specialized ones that a relational database cannot handle efficiently. The reason is that these queries invariably have a join predicate which is a conjunction of “genometric” predicates. The GMQL “genometric” predicates in essence are:  $DL(n)(y, x)$ , meaning the regions overlap or their nearest points are less than  $n$  bases apart;  $DG(n)(y, x)$ , meaning the regions do not overlap and their nearest points are more than  $n$  bases apart; and a few other similar ones. GMQL also implicitly imposes, on genometric predicates, the constraint that  $y$  and  $x$  are no further apart than a system-fixed number of bases (e.g., 200,000 bases.) For a reader who is unfamiliar with genomics, “bases,” or “bp,” is the unit used for describing distance on a genome.

GMQL queries can be easily modeled and efficiently implemented in our Synchrony iterator framework. Let  $xs: Vec[Bed]$  and  $ys: Vec[Bed]$  be two BED files sorted in accordance to  $<_{Bed}$ . Then, `isBefore` is monotonic with respect to  $(xs, ys)$ . Genometric predicates such as  $DL(n)$  are antimonotonic with respect to `isBefore`. Genometric predicates such as  $GL(n)$  are not antimonotonic with respect to `isBefore`. As GMQL automatically inserts  $DL(200000)$  as an additional genometric predicate into a query, a query has at least one antimonotonic genometric predicate.

The implementation of GMQL using Synchrony iterator is described in a companion paper (Perna *et al.*, 2021). Here, we just briefly describe a more complex GMQL query operator,  $JOIN(g_1, \dots, g_n; f, h, j)(xss, yss)$ . This GMQL query finds all pairs of samples  $xs$  in  $xss$  and  $ys$  in  $yss$  satisfying the join predicate  $j(xs, ys)$  on samples. Then for each such pair of samples  $xs$  and  $ys$ , for each pair of regions  $x$  in  $xs.bedFile$  and  $y$  in  $ys.bedFile$  satisfying all the specified genometric predicates  $g_1(y, x), \dots, g_n(y, x)$ , it builds a new region  $f(x, y)$ ; these new regions are put into a new BED file  $xys$ ; finally, a new sample having BED file  $xys$  and metadata  $h(x.meta, y.meta)$  is produced.

$JOIN(g_1, \dots, g_n; f, h, j)(xss, yss)$  is naturally and efficiently embedded into first-order Scala via comprehension syntax and Synchrony iterator. To wit, it is realized by

```

1876 for (xs <- xss; ys <- yss; if j(x, y))
1877   yield {
1878     val yi = new EIterator(ys.bedFile, isBefore, p)
1879     val xys = for (x <- xs.bedFile; y <- yi.syncedWith(x); if q(y,x))
1880       yield f(x,y)
1881     Sample(bedFile = xys, meta = h(xs.meta, ys.meta))
1882   }

```

where  $p$  is the conjunction of all the antimonotonic predicates among  $g_1 \dots, g_n$  and  $q$  is the conjunction of all the remaining predicates among  $g_1 \dots, g_n$ . In fact, our Synchrony-based GMQL implementation does this decomposition of the list of input genometric predicates into  $p$  and  $q$  automatically.

### 7.3 A stress test

We stress-tested Synchrony iterator by re-implementing GMQL using Synchrony iterator. The GMQL engine (Masseroli *et al.*, 2019) is a state-of-the-art purpose-built system for querying genomic datasets. GMQL is optimized for sample databases containing many samples, with each sample having a large BED file (Neph *et al.*, 2012) containing tens of thousands to hundreds of thousands of genomic regions. GMQL achieves high performance by binning the genome into chunks and comparing different bins concurrently (Gulino *et al.*, 2018).

As the GMQL is based on Scala, we reimplemented it using Synchrony iterator in Scala; this way, the influence of programming language and compiler differences is eliminated. The Synchrony implementation comes with a sequential mode (samples and their BED files are processed in a strictly sequential manner) and a sample-parallel mode (BED files of different samples are processed in parallel but regions in a BED file are processed in a sequential manner.) This reimplementation comprises circa 4,000 lines of Scala codes as counted by `cloc`; and makes use of some Scala function libraries. In contrast, the original GMQL engine comprises circa 24,000 lines of Scala codes and uses more Scala function libraries and also Spark function libraries. This comparison reveals the merit of Synchrony iterator in enabling complex algorithms to be expressed in a succinct high-level manner.

For benchmarking, we deployed the GMQL engine on a local installation of Apache Spark, which simulates a small cluster on a single multicore machine. We refer to this as the GMQL *command-line interface*, or CLI. The machine is a laptop with 2.6 GHz 6-Core i7, 16 GB 2667 MHz DDR4, 500 GB SSD. Despite the simplicity of our implementation, it significantly outperforms GMQL CLI on essentially all test queries and on the full range of dataset sizes and equals GMQL CLI on the largest-size datasets. This is a strong testimony to Synchrony iterator as an elegant idea for expressing efficient synchronized iterations on multiple collections in a succinct and easy-to-understand manner. The implementation and detailed evaluation are presented in a companion paper (Perna *et al.*, 2021). The implementation is available at <https://www.comp.nus.edu.sg/~wong1s/projects/synchrony>.

We present below some comparison results on a simple region `MAP` query. The GMQL `MAP` query takes two sample databases `xss` and `yss` and produces for each pair of BED files `xs.bedFile` in `xss` and `ys.bedFile` in `yss`, and each region `x` in `xs.bedFile`, the number of regions in `ys.bedFile` that it overlaps with. GMQL executes its `MAP` operator in a four-level deeply nested loop, in a brute-force parallel manner; i.e., all BED file pairs are analyzed in parallel. For each BED file pair, the BED files are chopped into bins; the bins are paired; and all bin pairs are analyzed in parallel. Ignoring parallelism, the complexity is  $O(n^2m^2)$  assuming both `xss` and `yss` contain  $n$  BED files and each BED file contains  $m \gg n$  regions. The Synchrony iterator version uses a two-level nested loop to pair up the BED files, but each pair of BED files is analyzed using a Synchrony iterator:

```

1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
for (xs <- xss; ys <- yss)
yield {
  val yi = new EIterator(ys.bedFile, isBefore, DL(0))
  for (x <- xs.bedFile; r = yi.syncedWith(x))
  yield (x, r.length)
}

```

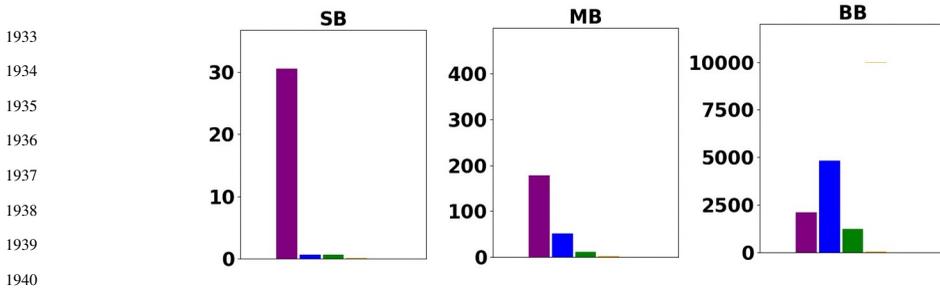


Fig. 12. Performance of GMQL CLI and Synchrony emulation on simple region MAP. Time in seconds, average of 30 runs for SB and MB, and 5 runs for BB. *Purple*: GMQL CLI. *Blue*: Sequential Synchrony emulation. *Green*: Sample-parallel Synchrony emulation.

The sample-parallel version runs the two-level nested loop in parallel but the Synchrony iterator sequentially. The sequential version does everything sequentially. Ignoring parallelism, the complexity is  $O((2k + 1)mn^2)$  where  $k$  is a small number corresponding to the maximum number of overlaps a region can have with other regions.

For this paper, the three versions are run on three input settings (SB, MB, BB) containing varying number of BED files, where each BED file has 100,000 regions. The setting SB means both  $x_{ss}$  and  $y_{ss}$  contain exactly one BED file; thus, there is exactly one BED file pair to analyze. The setting MB means both  $x_{ss}$  and  $y_{ss}$  contain exactly ten BED files; thus, there are 100 BED file pairs to analyze. The setting BB means both  $x_{ss}$  and  $y_{ss}$  contain exactly one hundred BED files; thus, there are 10,000 BED file pairs to analyze. Roughly  $x_{ss}$  and  $y_{ss}$  are each of size circa 10MB, 96MB, and 696MB on disk in settings SB, MB, and BB. The timings are shown in Figure 12. It is clear that the two Synchrony iterator versions are far more efficient than GMQL CLI. Only in the BB setting, GMQL CLI is able to beat the strict sequential Synchrony iterator. But GMQL CLI's brute-force parallelism is still no match to sample-parallel Synchrony iterator for these and other settings considered.

#### 7.4 Stream queries

As a last use-case, we model timestamped data streams. Two kinds of objects are considered for this purpose. The first kind is called observations. An observation  $x: \text{Obs}(\text{at}: \text{Int}, \text{id}: \dots)$  has a timestamp  $x.\text{at}$ , which is the time the observation is obtained, and has some other pieces of information that are irrelevant for our purpose. All timestamps are the number of nanoseconds that have elapsed since a fixed reference time point. The second kind is called observation streams. An observation stream is just a collection of observations,  $xs: \text{Vec}[\text{Obs}]$ .

Observations are intrinsically ordered by their timestamps. Thus, it is natural to define the following as the `isBefore` predicate on observation streams:

$$\text{bf}(y, x) = y.\text{at} < x.\text{at}$$

and it is also natural to assume that observation streams are sorted by timestamps by default. A variety of `canSee` predicates can be easily defined, such as:

$$\text{notAfter}(y, x) = ! (y.\text{at} > x.\text{at})$$

```
1979   within(n)(y, x) = abs(y.at - x.at) <= n
```

1980 For convenience, let `clk` be a stream of observations representing regular clock ticks at  
 1981 intervals of 1 milliseconds. Also, let `xss`, `yss`, and `zss` be several streams of observations.  
 1982 Then a variety of observation processing can be easily and efficiently expressed. Just for  
 1983 practice, to see how it looks, the suggested syntax for abstracting away Synchrony iterator  
 1984 from Section 6.4 is used here; the programs below are not legitimate Scala.

- 1985 • `cartesian(f)(clk, xss, yss, zss)` groups observations in `xss`, `yss`, and `zss` into 1  
 1986 millisecond time-synchronized blocks; applies `f` to each block to generate a new  
 1987 observation stream.  
 1988

```
1989   cartesian(f)(clk, xss, yss, zss)
1990   = { val cs = (u:Obs, v:Obs) => within(1000)(u,v) && notAfter(u,v)
1991       for ((c, xs, ys, zs) <- clk syncWith(xss, bf, cs)
1992           syncWith(yss, bf, cs)
1993           syncWith(zss, bf, cs)
1994       ) yield f(c, xs, ys, zs) }
```

- 1995 • `mostRecent(f)(clk, xss, yss, zss)`, applies `f` to the last observation in `xss`, `yss`,  
 1996 and `zss` within each 1 millisecond block. Skips a block if any stream contains no  
 1997 observation in that block of time.  
 1998

```
1999   mostRecent(f)(clk, xss, yss, zss)
2000   = { val cs = (u:Obs, v:Obs) => within(1000)(u,v) && notAfter(u,v)
2001       for ((c, xs, ys, zs) <- clk syncWith(xss, bf, cs)
2002           syncWith(yss, bf, cs)
2003           syncWith(zss, bf, cs);
2004       if !(xs.isEmpty || ys.isEmpty || zs.isEmpty)
2005       ) yield f(c, xs.last, ys.last, zs.last) }
```

- 2006 • `affineMostRecent(f)(clk, xss, yss, zss)`, applies `f` to the last observation in `xss`,  
 2007 `yss`, and `zss` within each 1 millisecond block. If a block has a stream which contains  
 2008 no observation in this block of time, keep observations in this block and consider  
 2009 them with the next block.  
 2010

```
2011   affineMostRecent(f)(clk, xss, yss, zss)
2012   = { val cs = (u:Obs, v:Obs) => within(2000)(u,v) && notAfter(u,v)
2013       def nd(us: Vec[Obs], t: Int) = us.filter(_.at <= t).isEmpty
2014       for ((c, xs, ys, zs) <- clk syncWith(xss, bf, cs)
2015           syncWith(yss, bf, cs)
2016           syncWith(zss, bf, cs);
2017       if !(xs.isEmpty || ys.isEmpty || zs.isEmpty);
2018       (lx, ly, lz, oc) = (xs.last, ys.last, zs.last, c.at - 999);
2019       if (lx.at > oc && ly > oc && lz > oc) ||
2020           ((lx.at > oc || ly > oc || lz > oc) &&
2021            (nd(xs, oc) || nd(ys, oc) || nd(zs, oc)))
2022       ) yield f(c, lx, ly, lz) }
```

As can be seen, a variety of temporal stream processing and synchronization operators, akin to those in Bracevac *et al.* (2018), can be implemented in comprehension syntax using Synchrony iterator. Notably, provided there are not too many events within each 1 millisecond block and  $f$  has at most linear time complexity, all of these examples have linear time complexity. If the observation type has a structure that carries more information (e.g., length of observation, if observation extends over a period of time), an even richer variety of antimonotonic predicates can be defined and used in specifying stream synchronization.

To some extent, this use-case illustrates that Synchrony iterator is not restricted to database query processing. Rather, it is capturing and generalizing common patterns and forms of synchronized iterations, such as those found in database query processing and in stream event processing.

## 8 Other possibilities

### 8.1 Grouping

It is instructive to look at the function `groups` in the top half of Figure 13. It was suggested by a reviewer as an approach to efficient implementation of relational joins. The idea is based on grouping. This suggestion inspired us to introduce the Synchrony generator `syncGenGrp`, which we did not describe in the initial draft of this paper mainly because it returns a nested collection and thus, strictly speaking, does not meet the first-order restriction requirement.

The reviewer probably had in mind a function like `syncGenGrp` and provided the function `groups` in Figure 13 as the implementation. However, this only works correctly when the join predicate `cs` is a conjunction of equality tests, i.e. an equijoin. Here is an example to show that it does not correctly implement a join in general. Let us regard `xs: Vec[Event]` and `ys: Vec[Event]` as lists of line segments sorted by  $(start, end)$ . Consider the following line segments.

```
a = Event(start = 10, end = 70, id = "a")
b = Event(start = 20, end = 30, id = "b")
c = Event(start = 40, end = 80, id = "c")
d = Event(start = 60, end = 90, id = "d")
```

Let `isBefore` and `overlap` be as defined in Figure 1. Let `xs` be a singleton containing the line segment `d` and `ys` comprises the line segments `a`, `b` and `c` in this order. Then `ovl(xs, ys)` evaluates to exactly the two pairs  $(d, a)$  and  $(d, c)$ . In agreement with `ovl(xs, ys)`, `syncGenGrp(isBefore, overlap)(xs, ys)` evaluates to the singleton  $(d, \text{Vec}(a, c))$ . Whereas, `groups(isBefore, overlap)(xs, ys)` incorrectly evaluates to an empty collection.

Perhaps instead of `val yt = ys.dropWhile(y => bf(y, x))`, the reviewer meant `val yt = ys.dropWhile(y => bf(y, x) && !cs(y, x))`. This revised `groups(bf, cs)(xs, ys)` works correctly when `bf` is monotonic with respect to  $(xs, ys)$  and `cs` is reflexive and convex with respect to `bf`. It does not work as expected when `cs` is antimonotonic but not convex. We mentioned earlier that predicates which are reflexive and convex are also antimonotonic, and that the converse is not true. The `overlap` predicate on events is such

```

2071 def groups[A,B]
2072   (bf: (B,A) => Boolean, cs: (B,A) => Boolean)
2073   (xs: Vec[A], ys: Vec[B])
2074 : Vec[(A,Vec[B])] = {
2075   def step(acc: (Vec[(A,Vec[B])], Vec[B]), x: A)
2076   : (Vec[(A, Vec[B])], Vec[B]) = {
2077     val (xzss, ys) = acc
2078     // this works only for equijoin cs:
2079     val yt = ys.dropWhile(y => bf(y, x))
2080     // this works for convex cs:
2081     // val yt = ys.dropWhile(y => bf(y, x) &&& ! cs(y, x))
2082     val zs = yt.takeWhile(y => cs(y, x))
2083     (xzss :+ (x, zs), yt)
2084   }
2085   val e: (Vec[(A,Vec[B])], Vec[B]) = (Vec(), ys)
2086   val (xzss, _) = xs.foldLeft(e)(step _)
2087   return xzss
2088 }
2089
2090 def groups2[A,B]
2091   (bf: (B,A) => Boolean, cs: (B,A) => Boolean)
2092   (xs: Vec[A], ys: Vec[B])
2093 : Vec[(A,Vec[B])] = {
2094   // Requires: bf monotonic wrt (xs, ys); cs antimonotonic wrt bf.
2095   val step = (acc: (Vec[(A,Vec[B])], Vec[B]), x: A) => {
2096     val (xzss, ys) = acc
2097     val maybes = ys.takeWhile(y => bf(y, x) || cs(y, x))
2098     val yes = maybes.filter(y => cs(y, x))
2099     val nos = ys.dropWhile(y => bf(y, x) || cs(y, x))
2100     (xzss :+ (x, yes), yes ++: nos)
2101   }
2102   val e: (Vec[(A,Vec[B])], Vec[B]) = (Vec(), ys)
2103   val (xzss, _) = xs.foldLeft(e)(step)
2104   return xzss
2105 }

```

Fig. 13. Alternative attempts to define `syncGenGrp`. The function `groups` is only correct when `cs` is an equijoin predicate. The function `groups2` is equivalent to `syncGenGrp` and has comparable efficiency.

an example; it is antimonotonic but not convex. This can be seen using the line segments given above: `overlap(a, d)` and `overlap(c, d)` but not `overlap(b, d)`. Indeed, this revised `groups` function returns the singleton `(d, Vec(a))`, which is still incorrect.

In order to get the correct semantics as `syncGenGrp`, the definition of `groups` must be modified to account for antimonotonicity. This can be done as in `groups2`, depicted in the bottom half of Figure 13. In `groups2`, for each `x` in `xs`, `step` iterates on the current copy of `ys`, to divide it using `takeWhile` and `dropWhile`. The function `dropWhile` stops the iteration on `ys` as soon as an item `y` in `ys` is encountered such that both `bf(y, x)` and `cs(y, x)` are false, and yields the remainder `nos`. This early stopping is correct due to Antimonotonicity Condition 2. The function `takeWhile` copies items on `ys` until an item `y` in `ys` is encountered such that both `bf(y, x)` and `cs(y, x)` are false, obtaining the prefix `maybes`. Those items in `maybes` that can see `x` are extracted into `yes` and returned as the result for this `x`. The function `step` also updates `ys` to `yes ++: nos`, “rewinding” it and setting it up for the next item in `xs`. Thus, for the next `x`, the iteration on `ys` effectively skips over all the items in `ys` that are before it and cannot see it. This skipping is correct due to Antimonotonicity Condition 1.

Group2 and syncGenGrp can be shown to define the same function and have similar time complexity. The concepts of monotonicity and antimonotonicity seems fundamental to achieving efficient synchronized iteration. In particular, groups2 needs both of these to ensure correctness. Their use in groups2 has further clarified how these concepts interact to synchronize iteration. Specifically, synchronized iteration on two collections, which are already sorted in a comparable ordering, is characterized by knowing precisely when to start, stop, and “rewind.” Hence, by explicitly parameterizing its iteration control on monotonicity and antimonotonicity, Synchrony fold can perhaps be regarded as a programming construct that characterizes efficient synchronized iteration on ordered collections.

## 8.2 Indexed tables

A reviewer also introduced us to the recent works of Gibbons (2016) and Gibbons *et al.* (2018), which investigate programming language embedding of joins that avoid naive evaluation strategies. These works provide an elegant mathematical foundation for indexing and grouping, leading to efficient implementation of equijoins. Underlying the theoretical perspective of these works is the use of indexed tables.

Indexed tables are part of the repertoire of collection-type libraries of modern programming languages. For example, the collection-type libraries of Scala provide a groupBy method. For a collection  $ys: \text{Vec}[A]$ , and an indexing function  $f: A \Rightarrow K$ ,  $ys.groupBy(f)$  builds and returns an indexed table  $ms: \text{Map}[K, \text{Vec}[A]]$ , where  $ms(ky)$  is the precomputed result of `for (y <- ys; if f(y) == ky) yield y`. Assuming  $f$  is a function of  $O(1)$  time complexity, the indexed table is constructed in  $O(|ys|)$  time, and accessing  $ms(f(y))$  takes  $O(1)$  time.

This groupBy function is useful for implementing efficient equijoin by a user-programmer, if we disregard the fact that it returns a nested collection and is thus not first order. For a direct example, assuming  $f$  and  $g$  are constant-time functions, the equijoin `for (x <- xs; y <- ys; if g(x) == f(y)) yield (x, y)` can be computed in  $O(|xs| + |ys|)$  time using an indexed table as `{ val ms = ys.groupBy(f); for (x <- xs; y <- ms(g(x)) yield (x, y) }`. This corresponds to an index-look join strategy, which a database system usually uses when it expects there are not many  $y$  that matches any  $x$  at all. As another example, Gibbons (2016) describes an interesting perspective where collections are viewed as indexed tables, and derives a zip-parallel comprehension for joining them. This strategy corresponds to the index-scan join strategy, which a database system uses when it expects most  $x$  matches at least one  $y$  and vice versa.

Using an indexed-table approach to implement an equijoin has a key advantage that the input collections do not need to be sorted to begin with. Provided all indexed tables which are needed can fit into memory, implementing an equijoin using an indexed-table approach is superior to using Synchrony iterator. If one or more of the input collections are unsorted or are in some unsuitable orderings, these inputs have to be sorted before Synchrony iterator can be used to process them; this can be a significant overhead when the input collections which require sorting are large.

On the other hand, there are several limitations with indexed tables, especially when we are not operating on an actual database system. Firstly, this means the indexed table

is an in-memory structure; so, it is not suitable for very large collections.<sup>10</sup> Secondly, the indexed table has to be completely constructed before it is used; so, it is not suitable for data streams. Lastly, and crucially, as an indexed table relies on exact equality to directly retrieve entries, it can easily implement efficient equijoin but it cannot implement non-equijoin such as band join and interval join.

In contrast, Synchrony iterator does not suffer these limitations. In other words, Synchrony iterator is a more general and more uniform approach for realizing efficient equijoin and a large class of non-equijoin. Synchrony iterator is thus justifiably appealing.

## 9 Concluding remarks

Modern programming languages typically provide some form of comprehension syntax for manipulating collection types. In this regard, comprehension syntax does not add extensional expressive power, but it makes programs much more readable (Trinder, 1991; Buneman *et al.*, 1994). Comprehensions typically correspond to nested loops. So, it is difficult to use comprehension syntax to express efficient algorithms for, e.g., database joins. This has partly motivated developments that introduced alternative binding semantics for comprehension syntax, so that some comprehensions are not compiled into nested loops. For example, parallel and grouping comprehension were introduced to enable implementation of efficient database queries in the style of comprehension syntax (Wadler & Peyton Jones, 2007; Gibbons, 2016; Gibbons *et al.*, 2018). Nonetheless, it has not been formally demonstrated that efficient algorithms for, e.g., equijoin cannot be implemented without making such refinements to comprehension syntax.

The first contribution of this paper is to highlight, in a precise sense, comprehension syntax suffers a limited-mixing handicap. In particular, this formally confirms that efficient algorithms for low-selectivity database joins—and this includes equijoin—cannot be implemented using comprehension syntax in the first-order setting (i.e., first-order restricted Scala in the context of this paper.) This justifies, from the intensional expressive power point of view, that these interesting works are necessary.

Although there is no efficient implementation for low-selectivity database joins in the first-order setting, they are nonetheless expressible as functions in the first-order setting. Therefore, the gap is purely in the intensional expressive power of comprehension syntax. So, we considered whether any function commonly provided in the collection-type libraries of modern programming languages is able to fix this gap. The limited-mixing handicap of comprehension syntax in the first-order setting remains even after adding any one of `foldLeft`, `takeWhile`, `dropWhile`, and `zip`; and most common functions in these libraries are derivatives of `foldLeft`.

<sup>10</sup> Even when we are operating on an actual database system, this fits-into-memory issue can be a problem in a situation where the database system has to—as is often the case—process many queries concurrently. While hash tables needed by a query may fit into memory, this may prevent hash tables needed by other queries to fit, thereby affecting the overall performance of the system. This was a reason that even though the hash join (Silberschatz *et al.*, 2016), which is based on dynamically constructing a hash table, has been known and implemented in database systems a long time ago, its use was discouraged (e.g., hash join was routinely disabled in Oracle 11g systems) until recent times when systems with very large memory have become common.

2209 The second contribution of this paper is to identify and propose a candidate library  
2210 function which fills this gap. We noticed that, apart from `zip`, the notion of general synchrony-  
2211 zed iteration on multiple collections is conspicuously absent from current collection-type  
2212 libraries. Hence, to kill two birds with one stone, we looked for a function that encapsulates  
2213 a common pattern of general synchronized iteration on multiple collection. Arguably,  
2214 `foldLeft` is the most powerful function in collection-type libraries of modern programming  
2215 languages, as most other commonly found functions in these libraries are extensionally  
2216 expressible using `foldLeft`. So, as an upperbound, we identified Synchrony fold, which is a  
2217 novel synchronized iteration function that expresses the same functions as `foldLeft` in the  
2218 first-order setting and yet expresses more algorithms, including efficient low-selectivity  
2219 database joins. Furthermore, just as a simple restriction can be imposed on `foldLeft` to  
2220 cut its extensional expressive power to precisely match comprehension syntax, a similar  
2221 restriction can be imposed on Synchrony fold to cut its extensional expressive power  
2222 to precisely match comprehension syntax. This restricted form is Synchrony generator.  
2223 Synchrony generator expresses exactly the same functions as comprehension syntax in  
2224 the first-order setting; but it expresses a richer repertoire of algorithms, including efficient  
2225 low-selectivity database joins. Hence, Synchrony generator is a conservative extension of  
2226 comprehension syntax that precisely fills its intensional expressiveness gap.

2227 Synchrony generator is nonetheless not well dovetailed with comprehension syntax in  
2228 the first-order setting. In particular, synchronized iteration over multiple ordered collec-  
2229 tions simultaneously apparently can only be expressed using Synchrony generator in an  
2230 aesthetically clumsy manner in the first-order setting. When a function `zipn` for simultane-  
2231 ously zipping  $n$  collections is available, efficient synchronized iteration over  $n$  collections  
2232 can be succinctly and elegantly expressed using Synchrony generator and this function.  
2233 However, `zipn` is outside the first-order setting. Moreover, this approach carries overheads  
2234 of  $n$  extra scans of at least one dataset. Another limitation of this approach is that it is  
2235 not user-programmer friendly: A zoo of `zip3`, `zip4`, etc. have to be provided, as a single  
2236 `zip` function for zipping an arbitrary number of collections of different types cannot be  
2237 assigned a valid type in strongly typed programming languages.

2238 The third contribution of this paper is Synchrony iterator. We found that Synchrony  
2239 generator is algorithmically equivalent to iterating on the items in a first collection, and  
2240 invoking Synchrony iterator on each of these items to efficiently return matching items  
2241 in a second collection. Synchrony iterator thus smoothly dovetails with comprehension  
2242 syntax. More importantly, it enables efficient synchronized iteration on multiple collections  
2243 to be simply expressed in comprehension syntax in a first-order setting and without the  $n$   
2244 extra-scan overheads.

2245 Synchrony fold, Synchrony generator, and Synchrony iterator can be regarded as capturing  
2246 an intuitive pattern of efficient synchronized iteration on ordered collections. They  
2247 suggest that efficient synchronized iteration on ordered collections are characterized by a  
2248 monotonic `isBefore` predicate that relates the orderings of the input collections, and an anti-  
2249 monotonic `canSee` predicate that identifies matching pairs to act on. The antimonotonicity  
2250 conditions on `canSee` further informs that the efficiency of the synchronization arises from  
2251 exploiting “right-sided convexity” of the matching items. Indeed, together, these predi-  
2252 cates make explicit where to start, stop, and rewind an iteration on two collections, thereby  
2253 achieving efficient synchronization.  
2254

2255 The fourth contribution of this paper is the revelation that efficient synchronized iter-  
2256 ation on ordered collections is captured by such a pattern which is characterized by the  
2257 monotonic `isBefore` and antimonotonic `canSee` predicates. A corollary of this fourth con-  
2258 tribution is the result that Synchrony generator (and thus Synchrony iterator) is a natural  
2259 generalization of the merge join algorithm (Blasgen & Eswaran, 1977; Mishra & Eich,  
2260 1992) widely used in database systems for decades for realizing efficient equijoins. With  
2261 a simple modification implied by Synchrony generator, the modified merge join algorithm  
2262 can work as long as the join predicate is antimonotonic with respect to the sort order of the  
2263 relations being joined.

2264 Lastly, we briefly described using Synchrony iterator to re-implement GMQL (Masseroli  
2265 *et al.*, 2019), which is a state-of-the-art query system for large genomic datasets. The  
2266 Synchrony-based re-implementation is more efficient than GMQL, and is also six-fold  
2267 shorter in terms of number of lines of codes, thereby validating the theory and design of  
2268 Synchrony fold and Synchrony iterator.

2269 This paper primarily illustrates Synchrony fold, Synchrony generator, and Synchrony  
2270 iterator using examples based on low-selectivity database joins. Nonetheless, Section 7.4  
2271 briefly showcases using Synchrony iterator to specify event stream processing operators.  
2272 This suggests Synchrony fold and Synchrony iterator capture patterns of efficient synchro-  
2273 nized iteration, showing that they can be parameterized by a pair of monotonic `isBefore`  
2274 and antimonotonic `canSee` predicates. However, our notion of synchronized iteration, as  
2275 encapsulated by Synchrony fold, generator, and iterator, is quite constrained. It maybe a  
2276 worthwhile future work to understand what interesting yet common patterns of efficient  
2277 synchronized iteration are not encapsulated by Synchrony fold and Synchrony iterator.

## 2278 **Conflicts of Interest**

2281 None

## 2284 **Acknowledgements**

2285 Stefano Ceri invited us to the *GeCo Workshop on Challenges in Data-Driven Genomic*  
2286 *Computing*, held in Como, Italy, in March 2019. This work evolved from the talk given  
2287 by LW at the workshop and the ensuing interesting discussions with VT and SP. We thank  
2288 Stefano for his invitation and surreptitious seeding of this work.

2290 Jeremy Gibbons and the reviewers provided very useful suggestions on this paper. They  
2291 also brought many relevant works and ideas to our attention. Their comments greatly  
2292 enriched our perspective in this work. We thank them for their invaluable contribution  
2293 in helping us improve this work.

2294 This work was supported by National Research Foundation, Singapore, under its  
2295 Synthetic Biology Research and Development Programme (Award No: SBP-P3); and by  
2296 Ministry of Education, Singapore, Academic Research Fund Tier-1 (Award No: MOE T1  
2297 251RES1725). In addition, VT was supported in part by a Kwan Im Thong Hood Cho  
2298 Temple Visiting Professorship, and LW was supported in part by a Kwan Im Thong Hood  
2299

2301 Cho Temple Chair Professorship. Any opinions, findings, and recommendations expressed  
2302 herein are those of the authors, and do not reflect the views of these grantors.

## 2305 References

- 2306 Abiteboul, S. and Vianu, V. (1991) Generic computation and its complexity. *Proceedings of 23rd*  
2307 *ACM Symposium on the Theory of Computing* pp. 209–219.
- 2308 Biskup, J., Paredaens, J., Schwentick, T. and den Bussche, J. V. (2004) Solving equations in the  
2309 relational algebra. *SIAM Journal on Computing* **33**(5):1052–1066.
- 2310 Blasgen, M. and Eswaran, K. (1977) Storage and access in relational databases. *IBM Systems Journal*  
2311 **16**(4):363–377.
- 2312 Bracevac, O., Amin, N., Salvaneschi, G., Erdweg, S., Eugster, P. and Mezini, M. (2018) Versatile  
2313 event correlation with algebraic effects. *Proceedings of ACM on Programming Languages*  
2314 **2**(ICFP):67.
- 2315 Buneman, P., Libkin, L., Suciuc, D., Tannen, V. and Wong, L. (1994) Comprehension syntax.  
2316 *SIGMOD Record* **23**(1):87–96.
- 2317 Buneman, P., Naqvi, S., Tannen, V. and Wong, L. (1995) Principles of programming with complex  
2318 objects and collection types. *Theoretical Computer Science* **149**(1):3–48.
- 2319 Colson, L. (1991) About primitive recursive algorithms. *Theoretical Computer Science* **83**:57–69.
- 2320 DeWitt, D. J., Naughton, J. F., Schneider, D. A. (1991) An evaluation of non-equijoin algorithms.  
2321 *Proceedings of 17th International Conference on Very Large Data Bases* pp. 443–452.
- 2322 Dignoes, A., Boehlen, M. H., Gamper, J., Jensen, C. S. and Moser, P. (2021) Leveraging range  
2323 joins for the computation of overlap joins. *The VLDB Journal*, [https://doi.org/10.1007/  
2324 s00778-021-00692-3](https://doi.org/10.1007/s00778-021-00692-3).
- 2325 Felleisen, M. (1991) On the expressive power of programming languages. *Science of Computer*  
2326 *Programming* **17**:35–75.
- 2327 Fortune, S., Leivant, D. and O’Donnell, M. (1983) The expressiveness of simple and second-order  
2328 type structures. *Journal of the ACM* **30**(1):151–185.
- 2329 Gaifman, H. (1982) On local and non-local properties. *Proceedings of the Herbrand Symposium,*  
2330 *Logic Colloquium ’81* pp. 105–135. North Holland.
- 2331 Gibbons, J. (2016) Comprehending ringads. Lindley, S., McBride, C., Trinder, P. and Sannella, D.  
2332 (eds), *A List of Successes That Can Change the World* pp. 132–151.
- 2333 Gibbons, J., Henglein, F., Hinze, R. and Wu, N. (2018) Relational algebra by way of adjunctions.  
2334 *Proceedings of the ACM on Programming Languages* **2**(ICFP):86.
- 2335 Gulino, A., Kaitoua, A. and Ceri, S. (2018) Optimal binning for genomics. *IEEE Transactions on*  
2336 *Computers* **68**(1):125–138.
- 2337 Henglein, F. and Larsen, K. F. (2010) Generic multiset programming with discrimination-based joins  
2338 and symbolic Cartesian products. *Higher-Order and Symbolic Computation* **23**(3):337–370.
- 2339 Hunt, A. and Thomas, D. (2000) *The Pragmatic Programmer: From Journeyman to Master*.  
2340 Addison-Wesley.
- 2341 Knuth, D. E. (1973) *The Art of Computer Programming: Sorting and Searching*. Addison Wesley.
- 2342 Libkin, L. and Wong, L. (1994) Aggregate functions, conservative extension, and linear orders. Beerl,  
2343 C., Ogori, A. and Shasha, D. E. (eds), *Proceedings of 4th International Workshop on Database*  
2344 *Programming Languages, New York, August 1993* pp. 282–294. Springer-Verlag. See also UPenn  
2345 Technical Report MS-CIS-93-36.
- 2346 Libkin, L. and Wong, L. (1997) Query languages for bags and aggregate functions. *Journal of*  
*Computer and System Sciences* **55**(2):241–272.
- Lindley, S., Wadler, P. and Yallop, J. (2011) Idioms are oblivious, arrows are meticulous, monads  
and promiscuous. *Electronic Notes in Theoretical Computer Science* **229**(5):97–117.
- Marlow, S., Peyton-Jones, S., Kmett, E. and Mokhov, A. (2016) Desugaring Haskell’s do-notation  
into applicative operations. *ACM SIGPLAN Notices* **51**(12):92–104.

- 2347 Masseroli, M., Canakoglu, A., Pinoli, P., Kaitoua, A., *et al.* (2019) Processing of big heteroge-  
2348 nous genomic datasets for tertiary analysis of next generation sequencing data. *Bioinformatics*  
2349 **35**(5):729–736.
- 2350 Mishra, P. and Eich, M. H. (1992) Join processing in relational databases. *ACM Computing Surveys*  
2351 **24**(1):63–113.
- 2352 Neph, S., Kuehn, M. S., Reynolds, A. P., Haugen, E., Thurman, R. E., Johnson, A. K., Rynes, E.,  
2353 Maurano, M. T., Vierstra, J., Thomas, S., Sandstorm, R., Humbert, R. and Stamatoyannopoulos,  
2354 J. A. (2012) BEDOPS: High-performance genomic feature operations. *Bioinformatics*  
2355 **28**(4):1919–1920.
- 2356 Odersky, M., Spoon, L. and Venners, B. (2019) *Programming in Scala: A Comprehensive Step-by-*  
2357 *Step Guide*. Artima Inc.
- 2358 Olteanu, D. and Schleich, M. (2016) Factorized databases. *ACM SIGMOD Record* **45**(2):5–16.
- 2359 Perna, S., Pinoli, P., Tannen, V., Ceri, S. and Wong, L. (2021) *Synchronized iteration for genomic*  
2360 *data processing*. Manuscript available from [https://www.comp.nus.edu.sg/~wongls/  
2361 projects/synchrony/synchrony-gmql-v12.pdf](https://www.comp.nus.edu.sg/~wongls/projects/synchrony/synchrony-gmql-v12.pdf).
- 2362 Piatov, D., Helmer, S. and Dignoes, A. (2016) An interval join optimized for modern hardware.  
2363 *Proceedings of 32nd IEEE International Conference on Data Engineering* pp. 1098–1109.
- 2364 Schmidt, D. A. (1986) *Denotational Semantics: A Methodology For Language Development*. Allyn  
2365 and Bacon.
- 2366 Sebesta, R. W. (2010) *Concepts of Programming Languages*. Addison-Wesley.
- 2367 Silberschatz, A., Korth, H. F. and Sudarshan, S. (2016) *Database System Concepts*. 7th edn.  
2368 McGraw-Hill.
- 2369 Suci, D. and Paredaens, J. (1997) The complexity of the evaluation of complex algebra expressions.  
2370 *Journal of Computer and Systems Sciences* **55**(2):322–343.
- 2371 Suci, D. and Wong, L. (1995) On two forms of structural recursion. *LNCS 893: Proceedings of 5th*  
2372 *International Conference on Database Theory* pp. 111–124. Springer-Verlag.
- 2373 Trinder, P. W. (1991) Comprehensions, a query notation for DBPLs. *Proceedings of 3rd*  
2374 *International Workshop on Database Programming Languages, Nahplion, Greece* pp. 49–62.  
2375 Morgan Kaufmann.
- 2376 Van den Bussche, J. (2001) Simulation of the nested relational algebra by the flat relational algebra,  
2377 with an application to the complexity of evaluating powerset algebra expressions. *Theoretical*  
2378 *Computer Science* **254**(1–2):363–377.
- 2379 Wadler, P. (1990) Deforestation: Transforming programs to eliminate trees. *Theoretical Computer*  
2380 *Science* **73**:231–248.
- 2381 Wadler, P. and Peyton Jones, S. (2007) Comprehensive comprehension: Comprehensions with ‘order  
2382 by’ and ‘group by’. *Haskell ’07: Proceedings of ACM SIGPLAN Workshop on Haskell* pp. 61–72.
- 2383 Wong, L. (1996) Normal forms and conservative extension properties for query languages over  
2384 collection types. *Journal of Computer and System Sciences* **52**(3):495–505.
- 2385 Wong, L. (2013) A dichotomy in the intensional expressive power of nested relational calculi aug-  
2386 mented with aggregate functions and a powerset operator. *Proceedings of 32nd ACM Symposium*  
2387 *on Principles of Database Systems* pp. 285–295.
- 2388 Wong, L. (2021) *Addressing an intensional expressiveness gap of comprehension syn-*  
2389 *tax*. Manuscript available from [https://www.comp.nus.edu.sg/~wongls/projects/  
2390 synchrony/v5-wls-natural2021.pdf](https://www.comp.nus.edu.sg/~wongls/projects/synchrony/v5-wls-natural2021.pdf)