Chapter 6

Sets and Tables

In this chapter we consider abstract data types for sets and tables.

6.1 An Abstract Data Type for Sets

Sets undoubtedly play an important role in mathematics and are often needed in the implementation of various algorithms.

Question 6.1. Can you think of a reason why you would prefer using a set data type instead of a sequence?

Whereas a sequence is an ordered collection, a set is its *unordered* counterpart. In addition, a set has no duplicate element. In some algorithms and applications, ordering is not necessary and it may be important to make sure that there are no duplicates.

Specification of sets. We now define an abstract data type for sets. The definition follows the mathematical definition of sets from set theory and is purely functional. In particular, when updating a set (e.g. with insert or delete) it returns a new set rather than modifying the old set. The specification for sets is given in Abstract Data Type 6.2.

Question 6.3. Can you see a redundancy in this interface?

Note that the *bulk operations*, intersection, union, difference can be implemented in terms of the operations find, insert, delete. Indeed, they are the bulk versions of these operations.

• intersection — search for multiple elements instead of one.

Abstract Data Type 6.2 (Sets). For a universe of elements \mathbb{U} (e.g. the integers or strings), the SET abstract data type is a type \mathbb{S} representing the power set of \mathbb{U} (i.e., all subsets of \mathbb{U}) along with the following functions:

```
\begin{array}{llll} \text{empty} & : & \mathbb{S} & = & \emptyset \\ \text{size}(S) & : & \mathbb{S} \to \mathbb{N} & = & |S| \\ \text{singleton}(e) & : & \mathbb{U} \to \mathbb{S} & = & \{e\} \\ \text{filter}(f,S) & : & ((\mathbb{U} \to \mathbb{B}) \times \mathbb{S}) \to \mathbb{S} & = & \{s \in S \mid f(s)\} \\ \\ \text{find}(S,e) & : & \mathbb{S} \times \mathbb{U} \to \mathbb{B} & = & |\{s \in S \mid s = e\}| = 1 \\ \text{insert}(S,e) & : & \mathbb{S} \times \mathbb{U} \to \mathbb{S} & = & S \cup \{e\} \\ \text{delete}(S,e) & : & \mathbb{S} \times \mathbb{U} \to \mathbb{S} & = & S \setminus \{e\} \\ \\ \text{intersection}(S_1,S_2) & : & \mathbb{S} \times \mathbb{S} \to \mathbb{S} & = & S_1 \cap S_2 \\ \text{union}(S_1,S_2) & : & \mathbb{S} \times \mathbb{S} \to \mathbb{S} & = & S_1 \cup S_2 \\ \\ \text{difference}(S_1,S_2) & : & \mathbb{S} \times \mathbb{S} \to \mathbb{S} & = & S_1 \setminus S_2 \\ \end{array}
```

- where \mathbb{N} are the natural numbers (non-negative integers) and $\mathbb{B} = \{T, F\}$.
- union insert multiple elements.
- difference delete multiple elements.

Question 6.4. Can you implement the non-bulk operations in terms of the bulk operations?

We can implement find, insert, and delete in terms of the others.

```
 \begin{split} & \text{find}(S,e) = \text{size}(\text{intersection}(S, \text{singleton}(e))) = 1 \\ & \text{insert}(S,e) = \text{union}(S, \text{singleton}(e)) \\ & \text{delete}(S,e) = \text{difference}(S, \text{singleton}(e)) \end{split}
```

Question 6.5. Can you see a way to implement these bulk operations in terms of the non-bulk versions?

One simple idea would be to perform the non-bulk operations one by one, one after the other.

Question 6.6. Do you see a disadvantage to this approach?

The problem is that this approach is completely sequential.

It turns out these operations can be performed in parallel more efficiently than with using the non-bulk versions. We will not talk about the parallel-efficient implementations of these operations in this chapter, but will cover them in a later chapter. We will, however, talk about their cost specifications.

Remark 6.7. We write this definition to be generic and not specific to Standard ML. In our library, the type $\mathbb S$ is called set and the type $\mathbb U$ is called key, the arguments are not necessarily in the same order, and some of the functions are curried. For example, the interface for find is $find: set \to key \to bool$. Please refer to the documents for details. In the pseudocode, we will give in class and in the notes we will use standard set notation as in the right hand column of the table above.

Remark 6.8. You may notice that the interface does not contain a map function. If we try to generalize the notion of map from sequences, a map function does not make sense in the context of a set: if we interpret map to take in a collection, apply some function to each element and return a collection of the same structure. Consider a function that always returns 0. Mapping this over a set would return all zeros, which would then be collapsed into a singleton set, containing exactly 0. Therefore, such a map would allow reducing the set of arbitrary size to a singleton, which doesn't match the map paradigm (which traditionally preserves the structure and size).

Remark 6.9. Most programming languages either support sets directly (e.g., Python and Ruby) or have libraries that support them (e.g., in the C++ STL library and Java collections framework). They sometimes have more than one implementation of sets. For example, Java has sets based on hash tables and balanced trees. Unsurprisingly, the set interface in different libraries and languages differ in subtle ways. So, when using one of these interfaces you should always read the documentation carefully.

Cost specification for sets. So far, we have laid out a semantic interface, but before we can put it to use, we need to worry about the cost specification.

As we have discussed before, cost specifications depend on implementation.

Question 6.10. Can you think of a way to implement sets?

Sets can be implemented in several ways. The most common efficient ways used hashing or balanced trees. There are various tradeoffs in cost. For simplicity, we'll consider a cost model based on a balanced-tree implementation. We will cover how to implement these set operations when we talk about balanced trees later in the course. For now, a good intuition to have is that we use a comparison function to keep the elements in sorted order in a balanced tree.

Since a balanced tree implementation requires comparisons inside the various set operations, the cost of these comparisons affects the work and span. For this, we'll assume that compare has C_w work and C_s span.

	Work	Span
$size(S) \ singleton(e)$	O(1)	O(1)
filter(f,S)	$O\bigg(\sum_{e \in S} W(f(e))\bigg)$	$O\left(\log S + \max_{e \in S} S(f(e))\right)$
find(S,e) $insert(S,e)$ $delete(S,e)$	$O(C_w \cdot \log S)$	$O(C_s \cdot \log S)$
$intersection(S_1,S_2) \ union(S_1,S_2) \ difference(S_1,S_2)$	$O(C_w \cdot m \cdot \log(1 + \frac{n}{m}))$	$O(C_s \cdot \log(n+m))$

Question 6.12. Can you see the why the bulk operations are more work and span efficient?

If we perform for example $|S_2|$ applications of each of these operations, the work and span would be $|S_2| \log |S_1|$ (ignoring the cost of the comparison). This is significantly worse that what we would obtain with the bulk operations, which has has a vastly improved span (by a linear factor).

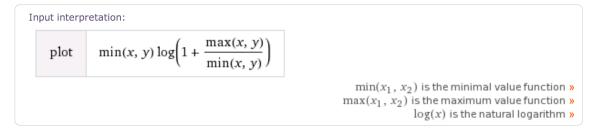
The work of the bulk operations might be somewhat more difficult to see. The key point to realize is that the bound is linear in the smaller of the two sets and logarithmic in the ratio of the two sets.

This means that if the ratio is very uneven, then it behaves very much like the naive algorithm that performs the non-bulk version repeatedly. For example, when one of the sets is a singleton, then the work is $O(\log n)$.

However, if the ratio is close to even, then it can reduce the total work to linear. When n=m, the work is simply

$$O(C_w \cdot m \cdot \log(1+1)) = O(C_w \cdot n).$$

This should not be surprising, because it corresponds to the cost of merging two approximately equal length sequences (effectively what these operations have to do).



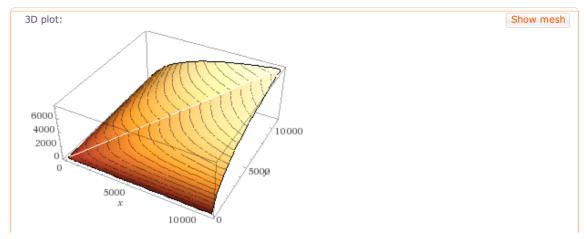


Figure 6.1: The work of the intersection, union, difference functions.

In fact, these bounds turn out to be both the theoretical upper and lower bounds for any comparison-based implementation of sets. We will get to this later. Figure 6.1 shows a plot of this function.

Exercise 6.13. Draw a two dimensional version of this figure by assuming one of S_1 to be of fixed size and varying the size of S_2 (or vice versa).

Consequently, in designing parallel algorithms it is good to think about how to use intersection, union, and difference instead of find, insert, and delete if possible.

Example 6.14. One way to convert a sequence to a set would be to insert the elements one by one, which can be coded as

```
1 fun fromSeq S =
2 Seq.iter Set.insert Set.empty S
```

However, the above is sequential. To do it in parallel we could instead do

```
\begin{array}{ll} 1 & \textit{fun fromSeq S} = \\ 2 & \textit{Seq.reduce Set.union Set.empty } \left\langle \left\{ x \right\} : x \in S \right\rangle \end{array}
```

Exercise 6.15. What is the work and span of the first version of from Seq.

Exercise 6.16. Show that on a sequence of length n the second version of from Seq does $O(C_w n \log n)$ work and $O(\log^2 n)$ span.

Summary 6.17. We talked about sets.

- Unordered collection.
- Unique elements.
- Supports efficient, find, insert, delete, operations serially and in parallel (bulk).

6.2 Tables: Assigning A Value to each Key

Since the elements of a set are unique (no duplicates), we can think of them as keys.

In many applications, it is important to be able to assign a value or data to each key. A *table* is an abstract data type that stores for each key data associated with it.

A table essentially stores *key-value* pairs in such a way that we can perform a range of operations quickly, e.g., finding the value for a key, inserting new key-value pairs, and deleting keys (and their values).

Question 6.18. Have you seen similar data structures before? Do you recall what they are called?

Tables are common data structures. They are also called dictionaries, associative arrays, maps, mappings, and functions (in set theory).

Most languages have tables either built in (e.g. dictionaries in Python, Perl, and Ruby), or have libraries to support them (e.g. map in the C++ STL library and the Java collections framework). We note that the interfaces for these languages and libraries have common features but typically differ in some important ways, so be warned. Most do not support the "parallel" operations we discuss.

Here we will define tables mathematically in terms of set theory before committing to a particular language.

109

Specifying tables. The specification of tables is quite similar to sets.

Definition 6.19. A table is set of key-value pairs where each key appears only once in the set.

We will use the following notation for a table

$$\{(k_1 \mapsto v_1), (k_2 \mapsto v_2), \dots, (k_n \mapsto v_n)\},\$$

where we have keys and values—and each key k_i mapped to a value v_i .

We choose this notation because that it makes clear that we are using tables rather than some other data structure such as sets.

Question 6.20. Can you see how we might represent tables as sets?

We can also represent a table as a set of key value pairs, e.g., $\{(k_1, v_1), (k_2, v_2), \dots, (k_n, v_n)\}$. Since keys are unique so are key-value pairs.

Question 6.21. Does this way of viewing maps remind you of a mathematical object?

Such sets are called *functions* in set theory since they map each key to a single value. We avoid this terminology so that we don't confuse it with functions in a programming language. However, note that the (find T) in the interface is precisely the "function" defined by the table T. In fact it is a *partial function* since the table might not contain all keys and therefore the function might not be defined on all inputs.

Abstract Data Type 6.22 defines tables.

Question 6.23. Can you see some differences between sets and tables?

There are several differences between sets and tables.

- The find function does not return a Boolean, but instead it returns the value associated with the key k. As it may not find the key in the table, its result may be bottom (\bot). For this reason, in the Table library, the interface for find is find: 'a table \rightarrow key \rightarrow 'a option, where 'a is the type of the values.
- When we insert a key-value pair, we can't simply ignore it if the key is already present, because the values might be different.

Abstract Data Type 6.22 (Tables). For a universe of keys \mathbb{K} , and a universe of values \mathbb{V} , the TABLE abstract data type is a type \mathbb{T} representing the power set of $\mathbb{K} \times \mathbb{V}$ restricted so that each key appears at most once (i.e., any set of key-value pairs where a key appears just once) along with the following functions:

```
\begin{array}{lll} \operatorname{empty} & : & \mathbb{T} & = & \emptyset \\ \operatorname{size}(T) & : & \mathbb{T} \to \mathbb{N} & = & |T| \\ \operatorname{singleton}(k,v) & : & \mathbb{K} \times \mathbb{V} \to \mathbb{T} & = & \{k \mapsto v\} \\ \operatorname{filter}(p,T) & : & (\mathbb{K} \times \mathbb{V} \to \mathbb{B}) \times \mathbb{T} \to \mathbb{T} & = & \{k \mapsto v\} \\ \operatorname{find}(T,k) & : & (\mathbb{V} \to \mathbb{V}) \times \mathbb{T} \to \mathbb{T} & = & \{k \mapsto f(v) : (k \mapsto v) \in T\} \\ \operatorname{find}(T,k) & : & \mathbb{T} \times \mathbb{K} \to (\mathbb{V} \cup \bot) & = & \begin{cases} v & (k \mapsto v) \in T \\ \bot & \text{otherwise} \end{cases} \\ \operatorname{insert}(f,T,(k,v)) & : & (\mathbb{V} \times \mathbb{V} \to \mathbb{V}) \times \mathbb{T} \times (\mathbb{K} \times \mathbb{V}) \to \mathbb{T} & = \\ \operatorname{merge}(f,T,\{k \mapsto v\}) \\ \operatorname{delete}(T,k)) & : & \mathbb{T} \times \mathbb{K} \to \mathbb{T} & = & \{(k' \mapsto v') \in T \mid k \neq k'\} \\ \operatorname{extract}(T,S) & : & \mathbb{T} \times \mathbb{S} \to \mathbb{T} & = & \{(k \mapsto v) \in T \mid k \in S\} \\ \operatorname{merge}(f,T_1,T_2) & : & (\mathbb{V} \times \mathbb{V} \to \mathbb{V}) \times \mathbb{T} \times \mathbb{T} \to \mathbb{T} & = \\ \operatorname{\forall} k \in \mathbb{K}, \begin{cases} k \mapsto f(v_1,v_2) & (k \mapsto v_1) \in T_1 \\ k \mapsto v_2 & (k \mapsto v_2) \in T_2 \\ k \mapsto v_2 & (k \mapsto v_2) \in T_2 \\ k \mapsto v_2 & (k \mapsto v_2) \in T_2 \end{cases} \\ \operatorname{erase}(T,S) & : & \mathbb{T} \times \mathbb{S} \to \mathbb{T} & = & \{(k \mapsto v) \in T \mid k \notin S\} \end{cases}
```

where \mathbb{S} is the power set of \mathbb{K} (i.e., any set of keys), \mathbb{N} are the natural numbers (non-negative integers), and $\mathbb{B} = \{T, F\}$.

For this reason, the insert function takes a function f as an argument,

$$f: \mathbb{V} \times \mathbb{V} \to \mathbb{V}$$
.

The purpose of f is to specify what to do if the key being inserted already exists in the table; f is applied to the two values. This function might simply return either its first or second argument, or it can be used, for example, to add the new value to the old one.

- The parallel counterpart of find is the extract function. The extract operation can be used to find a set of values in a table, returning just the table entries corresponding to elements in the set.
- The parallel counterpart of insert is the merge function, which takes a similar function to insert since it also has to consider the case that an element appears in both tables. The merge operation can add multiple values to a table in parallel by merging two tables.

- The parallel counterpart of delete is the erase function, The erase operation can delete multiple values from a table in parallel.
- We also introduce new specification-language notation for map and filter on tables:

$$\{k \mapsto f(v) : (k \mapsto v) \in T\}$$

is equivalent to map(f, T) and

$$\{(k \mapsto v) \in T \mid p(k, v)\}$$

is equivalent to filter(p, T).

Specifying the cost of tables. The costs of the table operations are very similar to sets.

	Work	Span
$size(T) \ singleton(k,v)$	O(1)	O(1)
filter(p,T)	$O\left(\sum_{(k\mapsto v)\in T}W(p(k,v))\right)$	$O\left(\log T + \max_{(k\mapsto v)\in T} S(f(k,v))\right)$
map(f,T)	$O\left(\sum_{(k\mapsto v)\in T}W(f(v))\right)$	$O\left(\log T + \max_{(k \mapsto v) \in T} S(f(v))\right)$
find(S,k) $insert(T,(k,v))$ $delete(T,k)$	$O(C_w \log T)$	$O(C_s \log T)$
extract (T_1,T_2) merge (T_1,T_2) erase (T_1,T_2)	$O(C_w m \log(1 + \frac{n}{m}))$	$O(C_s \log(n+m))$

As with sets there is a symmetry between the three operations extract, merge, and erase, and the three operations find, insert, and delete, respectively, where the prior three are effectively "parallel" versions of the earlier three.

Remark 6.25. We note that, in the SML Table library we supply, the functions are polymorphic (accept any type) over the values but not the keys. In particular the signature starts as:

```
signature TABLE =
   sig
3
    type 'a table
4
    type'a t = 'a table
     structure Key : EQKEY
    type key = Key.t
     structure Seq : SEQUENCE
7
     type a \ seq = 'a \ Seq.seq
9
     type set = unit table
10
    val find: 'a table \rightarrow key \rightarrow 'a option
11
12
```

The 'a in 'a table refers to the type of the value. The key type is fixed to be key. Therefore there are separate table structures for different keys (e.g. IntTable, StringTable). The reason to do this is because all the operations depend on the key type since keys need to be compared for a tree implementation, or hashed for a hash table implementation. Also note that the signature defines set to be a unit table. Indeed a set is just a special case of a table where there are no values.

Remark 6.26. In the SML Table library, we supply a collect operation that takes a sequence of key-value pairs and produces a table that maps every key in S to all the values associated with it in S, gathering all the values with the same key together in a sequence. This is equivalent to using a sequence collect followed by a Table. from Seq. Alternatively, it can be implemented as

```
1 function collect(S) =
2 let
3 val S' = \langle \{k \mapsto \langle v \rangle\} : (k, v) \in S \rangle
4 in
5 Seq.reduce (Table.merge Seq.append) \{\} S'
6 end
```

Exercise 6.27. *Figure out what this code does.*

113

Remark 6.28. Tables are similar to sets: they extend sets so that each key now carries a value. Their cost specification and implementations are also similar. In fact, in 15210 SML library, tables are implemented simply as sets where each key carries unit (constant) value.

6.3 Example: Bingle® It

Question 6.29. Can you think of some applications of sets and tables?

There are many. Here we consider an application of sets and tables to searching a corpus of documents.

In particular lets say one night, late, while avoiding doing your 210 homework you come up with a great idea: provide a service that indexes all the pages on the web so that people can search them by keywords. You figure a good name for such a service would be **Bingle**[®].

Question 6.30. Can you think of a simple algorithm to solve the problem.

Assuming that you have a copy of the Internet is some format, you can traverse it every time a query comes in.

Question 6.31. Going back to one of lectures on designing algorithms, can you see the problem with that?

The problem is redundancy. Every time a search comes is you have to traverse the whole data.

Instead, a better thing to do would be to build an *index* and use the index to avoid the redundant work. The idea would be for the index to organize the data in such a way to make searching efficient. You want queries to be fast since people will be running them all the time. On the other hand, the index can be slower to build, because you will run it every once in a while to keep your data up to date.

Question 6.32. Can you thing of the types of queries that you would wish to provide for?

You may want to support are logical queries on words involving And, Or, and AndNot. For example a query might look like

```
"CMU" And "fun" And ("courses" Or "clubs")
```

and it would return a list of web pages that match the query (*i.e.*, contain the words "CMU", "fun" and either "courses" or "clubs"). This list would include the 15-210 home page, of course.

Remark 6.33. This idea has been thought of before. Indeed these kinds of searchable indexes date back to the 1970s with systems such as Lexis for searching law documents. Today, beyond web searches, searchable indices are an integral part of most mailers and operating systems. The different indices support somewhat different types of queries. For example, by default Google supports queries with And and adjacent to but with their advanced search you can search with Or, AndNot as well as other types of searches.

```
Example 6.34. As a simple set of documents, we can consider the tweets made by some of your friends yesterday.
```

where the identifiers are the names, and the contents is the tweet.

On this set of documents, searching for "fun" and "club" would return "jack", "mary",

"sue", "peter", and "john"; lots of fun. (Note that our search returned "john" as well, even though he wasn't having that much fun.)

You can imagine that you would want to support an interface such as the following.

The input to makeIndex is a sequence of pairs each consisting of a document identifier (e.g. the URL) and the contents of the document as a single text string.

Example 6.35. Continuing on our example, we can use the interface to make an index of these tweets:

```
val f = (find (makeIndex(T))) : word \rightarrow docs
```

In addition to making the index, we partially apply find on the index. This makes it possible to use find with the index.

For example, the code,

```
toSeq(And(f "fun", Or(f "class", f "club"))) \\ \Rightarrow \langle "jack", "mary", "sue", "john" \rangle
```

returns all the documents (tweets) that contain "fun" and either "class" or "club". The code,

```
size(AndNot(f "fun", f "tiddlywinks")) \Rightarrow 4
```

returns the number of documents that contain "fun" and not "tiddlywinks".

Question 6.36. Can you think of a way to implement this interface? Let's start with the makeIndex function.

We can implement this interface using sets and tables. The *makeIndex* function can be implemented as follows.

```
1 function makeIndex(docs) =
2 let
3 function tagWords(id, str) = \langle (w, id) : w \in tokens(str) \rangle
4 val Pairs = flatten \langle tagWords(d) : d \in docs \rangle
5 val Words = Table.collect(Pairs)
6 in
7 \{w \mapsto Set.fromSeq(d) : (w \mapsto d) \in Words\}
8 end
```

The tagWords function takes a document as a pair consisting of the document identifier and contents, breaks the string into tokens (words) and tags each token with the identifier returning a sequence of these pairs.

Example 6.37. *Here is an example of how tagWords works:*

```
tagWords("jack", "chess club was fun")

⇒ ⟨("chess", "jack"),("club", "jack"), ("was", "jack"), ("fun", "jack")⟩
```

To build the index, we apply tagWords to all documents, and flatten the result to a single sequence.

In our example the result would start as:

```
Pairs = \(("chess", "jack"), ("club", "jack"), ("was", "jack"), ("fun", "jack"), ("I", "mary"), ("had", "mary"), ("fun", "mary"), ...
```

Using Table.collect, we then collect the entries by word creating a sequence of matching documents.

In our example it would start:

Finally, for each word the sequences of document identifiers is converted to a set. Note the notation that is used to express a map over the elements of a table.

Question 6.38. Do you see how we might implement the rest of the interface, which includes functionality for performing searches?

The rest of the interface can be implemented as follows:

```
function find T v = Table.find <math>T v function And(s_1,s_2) = s_1 \cap s_2 function Or(s_1,s_2) = s_1 \cup s_2 function AndNot(s_1,s_2) = s_1 \setminus s_2 function size(s) = |s| function toSeq(s) = Set.toSeq(s)
```

Cost. Assuming that all tokens have a length upper bounded by a constant, the cost of makeIndex is dominated by the collect, which is basically a sort. The work is therefore $O(n \log n)$ and the span is $O(\log^2 n)$, assuming the words have constant length.

117

Note that if we do a size(f "red") the cost is only $O(\log n)$ work and span. It just involves a search and then a length.

If we do And (f "fun", Or (f "courses", f "classes")) the worst case work and span are at most:

$$W = O(|f("fun")| + |f("courses")| + |f("classes")|)$$

 $S = O(\log |index|)$

The sum of sizes is to account for the cost of the And and Or. The actual cost could be significantly less especially if one of the sets is very small.

.