

UNIVERSITY OF TARTU
FACULTY OF MATHEMATICS AND COMPUTER SCIENCE
INSTITUTE OF COMPUTER SCIENCE

Egon Elbre
Parallel Pattern Discovery
Master's Thesis Draft

Supervisor: J. Vilo, PhD

TARTU 2012

Contents

1	Introduction	4
1.1	Motivation and background	4
1.2	Pattern Discovery	4
1.3	Structure of the thesis	5
2	Definitions	6
2.1	Sequence	6
2.2	Pattern	6
2.3	Finite automaton	7
2.4	Dataset	7
2.5	Query	7
2.6	Query features	7
2.7	Pattern discovery problem	8
3	SPEXS	9
3.1	Algorithm	9
4	Abstractification	11
4.1	Algorithm	11
4.2	Pools	12
4.3	Filtering	12
4.4	Extending	13
4.4.1	Sequences	14

4.4.2	Groups and Stars	14
4.4.3	Optimizations	15
4.4.4	Other	15
5	Parallelization	16
5.1	Process	16
5.2	Adding Nodes	17
5.3	Extending	17
5.4	Dataset partitioning	17
5.5	Filtering	17
6	Implementation	19
6.1	Algorithm	19
6.2	Architecture	21
6.3	Configuration	21
6.4	Sets	22
6.5	Pools	22
6.6	Features and Filters	23
6.7	Debugging	24
7	Applications and experimental results	26
7.1	DNA sequences	26
7.2	Protein sequences	26
7.3	Text mining	26
7.4	Code mining	26
8	Conclusions	27
	Bibliography	27
A	Exhaustive Parallel Pattern Search	29

B	SPEXS2 Command Line Utility	30
B.1	Configuration	30
B.2	Running	30
C	SPEXS2 Source Code	31

Chapter 1

Introduction

1.1 Motivation and background

Collecting new data has been increasing more rapidly than algorithms and computer processing power. The average size of each dataset has also been increasing. This suggests that the only way to keep up with analysis is to parallelize algorithms.

One of main drivers of such large datasets is analysis of genomic and proteomic sequences. Regularities in such data can give new insights into how these patterns form and how they are related to the other features of the data.

More on importance

In this thesis we explore an algorithm for finding patterns and show how abstractions can make it scalable and flexible, and simpler both in theory and implementation compared with non-abstract version.

1.2 Pattern Discovery

Pattern discovery is a research area aiming to discover unknown patterns in a given set of data structures that are frequent and interesting according to

some measure.

Since the discovery algorithms are highly dependent on the data structures, that are being searched, the algorithm must be minimal in the requirements on the dataset to be applicable as wide range as possible. This also means that the patterns found must be dependent on the initial data.

1.3 Structure of the thesis

Basic idea is to select a basis algorithm. Abstrify away concreteness of implementation to support flexibility. Use the flexibility to substitute concreteness with parallel implementations.

Fix when structured

First we introduce definitions of our data and patterns in Chapter 2. In Chapter 3 we describe the general SPEXS algorithm as specified in Vilo et al. In Chapter 4 we generalize and abstrify the algorithm to get a more flexible and parallel algorithm. We will describe the implementation consideration of flexible algorithms in Chapter 5 using the abstract SPEXS algorithm as an example. Chapter 6 uses the implementation to show how it can be applied.

Chapter 2

Definitions

Pattern discovery is a problem of finding interesting patterns in some dataset. We discuss the algorithm in terms of DFAs and show how this can help to discover patterns in sequences.

2.1 Sequence

We use Σ to denote the set of tokens in the dataset, an *alphabet*. The *size* of the alphabet is $|\Sigma|$. *Tokens* can be numbers, letters, words or sentences - any symbol.

Any *sequence* $S = a_1a_2...a_n, \forall i \in \Sigma$ is called a *sequence* over the token set Σ . If the length of the string is 0, it is called an empty sequence or ϵ .

2.2 Pattern

A *pattern* is a structure that defines a set of structures Γ . We denote the set that a pattern defines as $x(\Gamma)$. If a structure α *matches* a pattern Γ , it means that $\alpha \in x(\Gamma)$.

better
defini-
tion

The most common form of such structures are regular expressions.

2.3 Finite automaton

A *deterministic finite automaton* (henceforth *DFA*) is an finite state machine (henceforth *FSM*) that accepts or rejects a sequence of tokens. _____

define

A *non-deterministic finite automaton* (henceforth *NFA*) is an FSM that accepts or rejects a sequence of tokens. _____

define

Both patterns and sequences can be represented by NFAs and hence is a good abstraction for both. This also means that if our algorithm works on NFAs it must work on any other dataset that is defined in terms of NFA.

2.4 Dataset

A *dataset* is a set of NFAs. In practice it may be more comfortable to view it as a set of sequence because of it's simpler structure.

example

2.5 Query

A *query* is a compound structure that has information about a pattern and it's matches in a dataset. The information about a match is the ending states where a pattern matches a DFA.

On sequences this means the ending positions of the pattern in the sequence.

example

2.6 Query features

Query feature is a function of the query. It gives information about the query such as the pattern representation, length, number of matches in the dataset.

example

Query interestingness is a function of the query whose results are well-ordered. This functions result allows to say whether one query is more interesting than the other. For convenience it is useful to represent that value in \mathbb{R} .

example

Query filter is a function of the query whose result is a boolean.

example

2.7 Pattern discovery problem

Pattern discovery problem is a process of finding most interesting queries according to some interestingness measure, and use the pattern they reflect.

Chapter 3

SPEXS

SPEXS is an pattern discovery algorithm described by Vilo et al. This algorithm finds patterns from a sequence. We take this as our basis for developing a new parallel algorithm. In this chapter we describe original algorithm so that we can later show the changes made to this algorithm.

3.1 Algorithm

We describe the general representation of the SPEXS algorithm. The original algorithm was as follows:

Algorithm 1 The SPEXS algorithm

Input: String S , pattern class \mathcal{P} , output criteria, search order, and fitness measure \mathcal{F}

Output: Patterns $\pi \in \mathcal{P}$ fulfilling all criteria, and output in the order of fitness \mathcal{F}

- 1: Convert input into sequences into a single sequence
 - 2: Initiate data structures
 - 3: $Root \leftarrow newnode$
 - 4: $Root.label \leftarrow \epsilon$
 - 5: $Root.pos \leftarrow (1, 2, \dots, n)$
 - 6: $enqueue(\mathcal{Q}, Root, order)$
 - 7: **while** $N \leftarrow dequeue(\mathcal{Q})$ **do**
 - 8: Create all possible extensions $P \in \mathcal{P}$ of N using $N.pos$ and S
 - 9: **for** extension P of N **do**
 - 10: **if** pattern P and position list $P.pos$ fulfill the criteria **then**
 - 11: $N.child \leftarrow P$
 - 12: calculate $\mathcal{F}(P, S)$
 - 13: $enqueue(\mathcal{Q}, P, order)$
 - 14: **if** P fulfills the output criteria **then**
 - 15: store P into output queue \mathcal{O}
 - 16: Report the list of top-ranking patterns from output queue \mathcal{O}
-

The main idea of the algorithm is that first we generate a pattern and a query that match all possible positions in the sequence. We then put this query into a queue for extending.

Extending a query means finding all queries whose patterns length is longer by 1. If any of the queries is fit by some criteria it will be put into the main queue, for further extension, and output queue for possible result.

Chapter 4

Abstraction

In this chapter we show how to make the algorithm more abstract by allowing flexibility through functions as parameters.

4.1 Algorithm

The algorithm in a more conventional view is:

Algorithm 2 The SPEXS2 algorithm

Input: *dataset*, *in* pool, *out* pool, *extender*, *extendable* filter, *outputtable* filter, *postprocess* function

Output: Patterns satisfying filters and *extender* are in *out* pool

```
1:  $\varepsilon \leftarrow NewEmptyQuery(dataset)$ 
2: in.put( $\varepsilon$ )
3: while q  $\leftarrow$  in.take() do
4:   extended  $\leftarrow$  extender(q, dataset)
5:   for qx  $\in$  extended do
6:     if extendable(qx) then
7:       in.put(qx)
8:     if outputtable(qx) then
9:       out.put(qx)
10:  postprocess(q)
```

When the algorithm starts we create an empty pattern query and put into the in pool. The in pool contains queries whose patterns should be further examined.

We pick a query from the in pool for extending. The extending means generating all queries whose pattern is larger by one. There can be several such queries.

If any of the queries should be further examined as defined by the extendable filter, it will be put into the in pool.

If the query is fit for output as defined by the outputtable filter, it will be put into the out pool.

If we extend each pattern at each step by one we guarantee that we examine the all patterns that conform to our criteria.

4.2 Pools

Pool is an abstract datatype for a collection of queries. The pool allows queries to be put into it and taken from it, also we can ask whether the pool is empty or not.

It has no guarantees on how the queries are stored internally and in which order they are taken out.

In practice this means we can use any collection such as list, set, queue as a pool. This gives us different performance characteristics.

This also means that pools can persist the queries if necessary.

4.3 Filtering

Filtering allows us to dramatically reduce the number of queries we have to look at. It also allows to select only interesting patterns by some criteria.

If we have interestingness measure we can create filter from it by defining

it's minimum or maximum value. One very usefule example would be a filter for limiting the pattern length.

By separating the extension and output filter, as opposed to SPEXS, we can still limit output without affecting the extension process. For example if we wish to see only patterns of length 3 we cannot do it with one filter. Since we need to extend patterns of length 0, 1 and 2.

4.4 Extending

The extending process is at the core of the algorithm. We shall look at how we can deal with different types.

The extending method is:

Algorithm 3 SPEXS2 extender

Input: q query, $next$ function

Output: result contains queries that have been extended by one

```

1:  $nexts \leftarrow new\ collection$ 
2: for  $pos \in matches(q)$  do
3:    $(token, pos) \leftarrow next(pos)$ 
4:    $nexts.put((token, pos))$ 
5:  $matches \leftarrow new\ map\ of\ token\ to\ position\ set$ 
6: for  $(token, pos) \in nexts$  do
7:   if  $matches[token]$  doesn't exist then
8:      $matches[token] \leftarrow \{\}$ 
9:    $matches[token] \leftarrow matches[token] + \{pos\}$ 
10:  $result \leftarrow new\ collection$ 
11: for  $(token, positions) \in matches$  do
12:    $qx \leftarrow new\ query$ 
13:    $qx.pattern \leftarrow q.pattern + token$ 
14:    $qx.positions \leftarrow positions$ 
15:    $result.add(qx)$ 

```

This extender depends on how $next$ is implemented, we shall look at different ways how to implement it.

4.4.1 Sequences

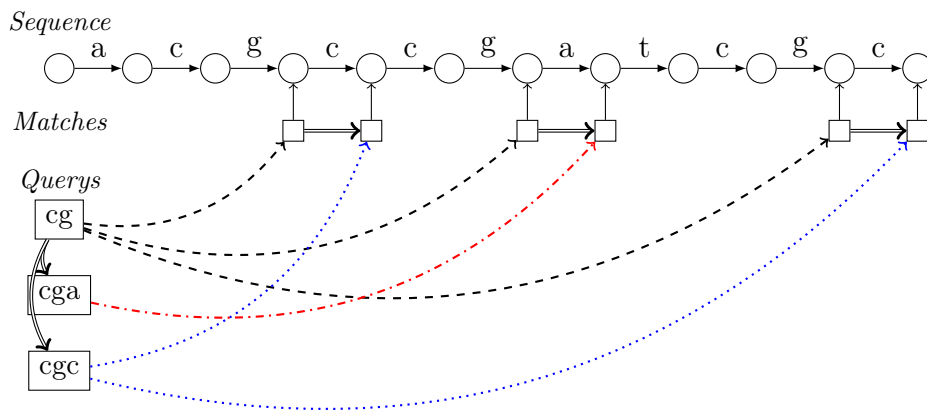
The simplest case how the next function behaves is when we are only looking for simple sequences.

Let's consider a sequence ACGCCGATCGC and a pattern CG.

Diagram of the ACG.CCG.ATCG.C and query for that pattern.

Next positions step is finding ACG.[C]CG.[A]TCG.[C].

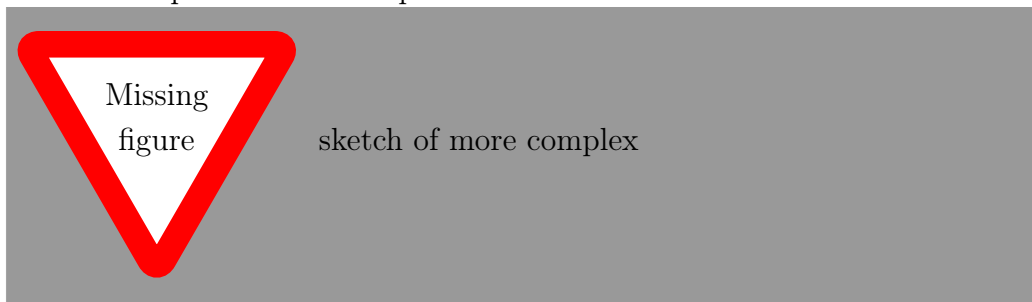
Grouping is [C] ==> CGC, [A] ==> CGA.



4.4.2 Groups and Stars

Since we want to find more interesting patterns we can add more information to the sequence DFA. Such as a star expression.

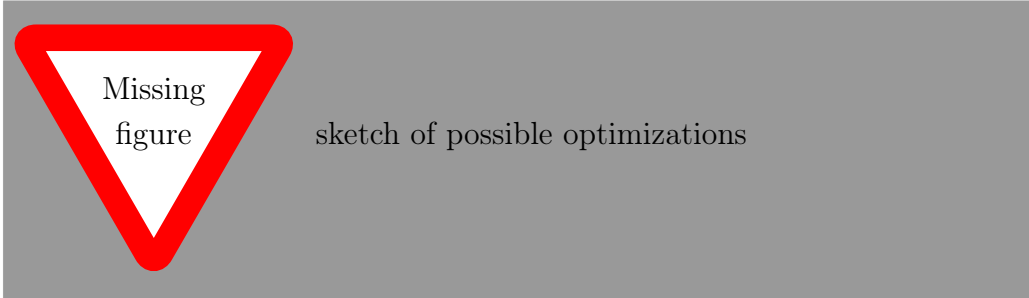
Same sequence with star paths.



4.4.3 Optimizations

Although we can add the group information to the DFA, it is more performant to use the information gathered from the extension of non-groups.

$$A[CG].Positions = AC.Positions \text{ union } AG.Positions$$

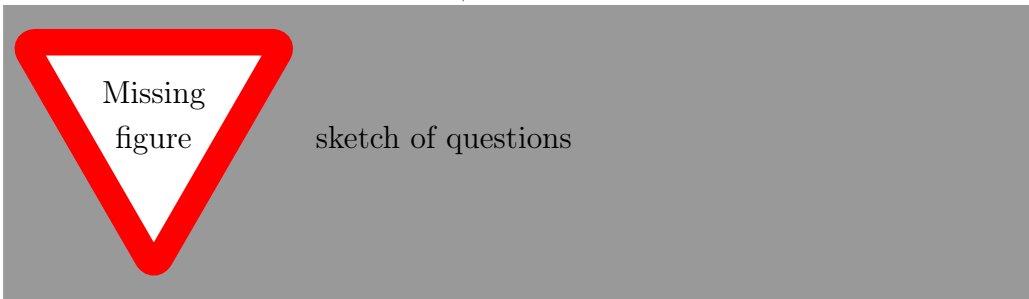


4.4.4 Other

There maybe several other extensions to the regular expressions.

Questionable position:

$$AC?.Positions = A.Positions + AC.Positions$$



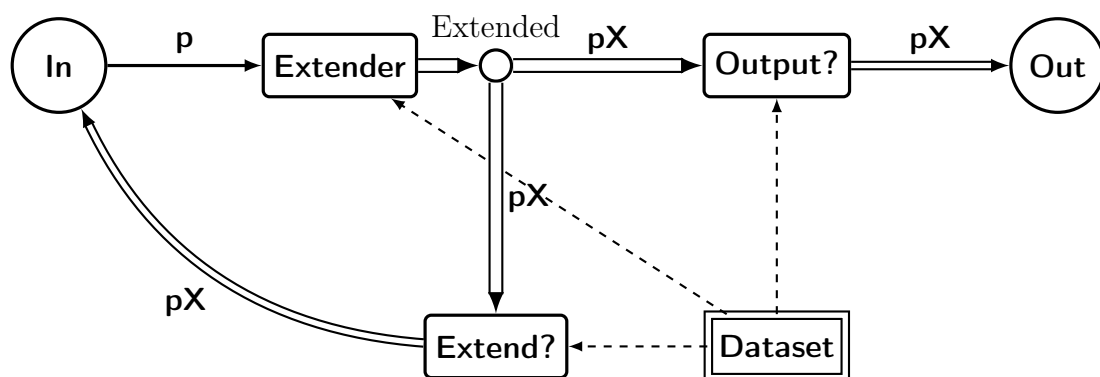
Chapter 5

Parallelization

Here we describe full parallelization of the previous algorithm. The fully parallelized code can be seen in appendix A.

5.1 Process

The main process of the algorithm as described by dataflow.[1, 2]



The easiest thing here to parallelize is the extender since it's interaction can be seen as a separate unit.

5.2 Adding Nodes

If instead of in/out pool we had several the algorithm can still work, if we have single take / put process to decide which actual pool to use.

We can use several extenders that work in parallel without problems since they do not need information about other queries nor input/output pools.

5.3 Extending

We still can also do the extending process in parallel:

```
func extender(query) {  
  // find all next positions  
  nexts = (map next query.matches)  
  matches = (group nexts by token)  
  result = (map newquery matches)  
  other = (map #(union (matches %1)) groups)  
  return result union other  
}
```

5.4 Dataset partitioning

If we look at where we need synchronization points if we partition the dataset. Whole parts of extension still apply for parts of dataset.

Synchronization is only needed for filtering. We don't need to know the whole query but just information about parts to make a decision whether to filter or not.

5.5 Filtering

This is the only place where we may need the whole information about the query.

Although many operations can be parallelized with map reduce.

example how to do counting

example what cannot be done easily on sharded data

Chapter 6

Implementation

Here we discuss a practical implementation, *spexs2*, for pattern discovery in sequences.

The actual implementation may need to diverge from the abstract definition for several reasons; mainly practicality, simplicity and performance. Many of the operations can be optimized for some particular type of datasets and configuration.

In this chapter we will discuss parts of program that the author considers non-trivial in it's design decisions. Information about the full source code is in the Appendix B and C.

6.1 Algorithm

The algorithm was implemented as:

```
type Setup struct {  
  Db  *Database  
  Out Pooler  
  In  Pooler  
  
  Extender Extender
```

```

Extendable Filter
Outputtable Filter

PostProcess PostProcess
}

func Run(s *Setup) {
prepareSpexs(s)
for {
p, valid := s.In.Take()
if !valid {
return
}

extensions := s.Extender(p)
for _, extended := range extensions {
if s.Extendable(extended) {
s.In.Put(extended)
}
if s.Outputtable(extended) {
s.Out.Put(extended)
}
}

if s.PostProcess(p) != nil {
break
}
}
}

```

Setup is a structure designed to hold the full setup of the algorithm. The algorithm is very similar to the algorithm described in the theoretical part, with the slight addition of PostProcess.

6.2 Architecture

The main criteria for designing program have been described in D. Parnas paper On the Decomposition of Programs. It suggests decomposing into isolated units and parts that are likely to change together. [3]

We chose the following decomposition:

```
Configuration
Dataset reader
Setup
Algorithm
Sets
Pools
Extenders
Features
Filters
Printer
```

This allows easily to add new data formats for input and output, add new query features and possible extenders.

6.3 Configuration

One problem with flexible algorithms is that they a lot of ways to be run. This often would need having tens or hundereds of program flags. To avoid this problem we decided to use a json file for the program configuration.

Example Configuration.

The problem with only using a json file is that when running from command line it may be more comfortable using flags. To solve this problem we added replacement strings into the json files that can be given in as a program argument.

```
"Datasets" : {
```

```
"fore" : { "File" : "$input$"  
...  
...
```

When using `spexs2 -conf conf.json input=filename` the input will be replaced by filename. Also there is optional default value if one wasn't given.

6.4 Sets

Since we need a collection how to store the matching positions it suggests the need for a set datatype.

If we have predictable distribution we can pack the sets better.

Although such optimizations can be avoided if during storing pools the sets are packed using some compression algorithm.

6.5 Pools

There can be different performance characteristics when using a particular implementation. Most importantly to support parallelism they should be ideally lock-free, but we can use locked version as well.

For the pools we have several choices: lifo, fifo or priority.

If we use a fifo queue as the in pool the algorithm does a breadth first search of patterns. This can be problematic since we would need a lot of memory to hold all the patterns in memory.

A lifo queue for in pool is a more reasonable choice for memory problems since we need to hold less patterns in memory.

A priority queue suits for the out pool since we can easily then use some feature to sort the queue and choose only a limited amount. We can use a simple trick to make limited size faster. Do a precheck against the worst element in the priority queue.

6.6 Features and Filters

One problem is that the amount of possible filters it's useful to construct them from some other features. Or if we wish to use multiple filters we can combine them.

We can use these features to find out something about the query. They each feature is defined as:

```
type Feature func(q *Query) (float64, string)
```

Most of features are in \mathfrak{R} , but for some there is some extended information that we may wish to know - hence the need for additional string value. One of the simplest is Pattern representation.

Also many of the features are defined in terms of multiple datasets. We can use a closure to easily define a more generic feature.

Example:

```
func Matches(dataset []int) Feature {
return func(q *Query) (float64, string) {
matches := countf(q.Pos, dataset)
return matches, ""
}
}
```

Here the Matches function creates a feature function defined for dataset.

We can use the name in the configuration file as "Matches(fore)".

If a feature returns a floating point value we can easily turn that into a filter by specifying a minimum or/and maximum value.

For example to give a lower and higher limits to some feature:

```
func featureFilter(feature Feature, min float64, max float64) Filter {
return func(q *Query) bool {
v, _ := feature(q)
return (min <= v) && (v <= max)
}
```



```
}  
}
```

Of course there are some filters that cannot be defined by features hence there is still possibility to make separate filters. Such as disallowing star symbol in the beginning of the pattern.

6.7 Debugging

Debugging is a important part of development process hence the need for more information how the algorithm is working.

Often this is resolved by adding some debug statments:

```
for i := 0; i < 100; i += 1{  
  printf("picking")  
  a := pick()  
  printf("picked %v", a)  
  
  printf("picking")  
  b := pick()  
  printf("picked %v", b)  
  
  out.put(a + b)  
}
```

This can be harmful to the readability of the code. To fix this debugging problem we use closures.

```
type PickFunc func()Thing  
func debuggable(fn PickFunc) PickFunc {  
  return func() Thing {  
    printf("start picking")  
    p := fn()  
    printf("picked %v" p)  
  }  
}
```

```

return p
}
}

pick = debuggable(pick)
for i := 0; i < 100; i += 1{
a := pick()
b := pick()
out.put(a + b)
}

```

As we can see the algorithm implementation is much more readable and we can inject different debugging statements without actually changing the algorithm.

Also we can now change the ways how to add debug info. One of would be a full stepwise debugger.

```

func debuggable(fn PickFunc) PickFunc {
return func() Thing {
mutex.lock()
p := fn()
while not continue
interact with user
mutex.unlock()
return p
}
}

```

Chapter 7

Applications and experimental results

Here we show examples for the program:

7.1 DNA sequences

make an example

7.2 Protein sequences

make an example

7.3 Text mining

make an example

7.4 Code mining

make an example

Chapter 8

Conclusions

In this thesis we showed how by making an algorithm more abstract and general we can also make it parallel. We showed that this algorithm can find patterns from sequences and also NFAs.

Although we showed that we can apply this algorithm on NFAs we did not analyze the performance characteristics. This suggests that this algorithm may be able to work on trees and graph, but would require slight modifications.

We also demonstrated how to make the algorithm more concrete and work well on some particular datasets. This also took into consideration further development and flexibility of the algorithm.

The SPEXS2 implementation currently is already used in biosequenceing and text mining.

Bibliography

- [1] G. Kahn, “The semantics of a simple language for parallel programming,” in *Information processing* (J. L. Rosenfeld, ed.), (Stockholm, Sweden), pp. 471–475, North Holland, Amsterdam, Aug 1974.
- [2] E. Lee and T. Parks, “Dataflow process networks,” *Proceedings of the IEEE*, vol. 83, pp. 773–801, may 1995.
- [3] D. L. Parnas, “On the criteria to be used in decomposing systems into modules,” *Commun. ACM*, vol. 15, pp. 1053–1058, Dec. 1972.
- [4] J. Vilo, *Pattern Discovery from Biosequences*. PhD thesis, University of Helsinki, 2002.
- [5] J. Dean and S. Ghemawat, “Mapreduce: simplified data processing on large clusters,” *Commun. ACM*, vol. 51, pp. 107–113, Jan. 2008.

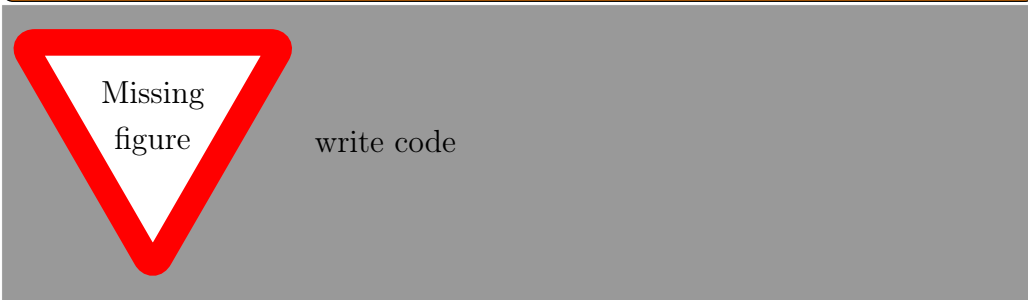
Appendix A

Exhaustive Parallel Pattern Search

This part assumes knowledge of *lisp* like languages. The code here is presented in Clojure and for an introduction see clojure.org/getting_started.

First we show reader macros that are different from other lisps. Then we present the algorithm with comments.

lists
sets
maps
lambdas
defn
map/reduce



Appendix B

SPEXS2 Command Line Utility

Write

B.1 Configuration

Write

B.2 Running

Write

Appendix C

SPEXS2 Source Code

The source code for the implementation of *SPEXS2* is available at github.com/egonelbre/spexs.

The source folder "src" has the following structure:

