



Digital Business Ecosystem

Contract number° 507953

**Workpackage 8**  
**Population dynamics in the Evolutionary Environment**

**Deliverable 8.3**  
**Report on Population dynamics in EvE**



Project funded by the European Community under the "Information Society Technology" Programme

**Contract number:** 507953  
**Project acronym:** DBE  
**Title:** Digital Business Ecosystem

**Deliverable N°:** D8.3  
**Due date:** 05/2006  
**Delivery date:** 07/2006

**Short description:** The final report on population dynamics. Note that this task has been brought forward by five months.

**Author:** UBHAM  
**Partners contributed:** STU, LSE, HWU  
**Made available to:** Public

Versioning		
Version	Date	Author, Organisation
1.0	30/06/2006	J. E. Rowe, B. Mitavskiy, J. Woodward (UBHAM)
2.0	24/07/2006	J. E. Rowe, B. Mitavskiy, J. Woodward (UBHAM)

**Quality check**  
**1<sup>st</sup> internal reviewer:** Gerard Briscoe (HWU)  
**2<sup>nd</sup> internal reviewer:** Thomas Kurz (STU)



This work is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 2.5 License. To view a copy of this license, visit : <http://creativecommons.org/licenses/by-nc-sa/2.5/> or send a letter to Creative Commons, 543 Howard Street, 5th Floor, San Francisco, California, 94105, USA.



### Attribution-NonCommercial-ShareAlike 2.5

#### You are free:

- to copy, distribute, display, and perform the work
- to make derivative works

#### Under the following conditions:



**Attribution.** You must attribute the work in the manner specified by the author or licensor.



**Noncommercial.** You may not use this work for commercial purposes.



**Share Alike.** If you alter, transform, or build upon this work, you may distribute the resulting work only under a license identical to this one.

- For any reuse or distribution, you must make clear to others the license terms of this work.
- Any of these conditions can be waived if you get permission from the copyright holder.

**Your fair use and other rights are in no way affected by the above.**

# Contents

<b>1</b>	<b>Introduction and relationship to the project</b>	<b>6</b>
<b>2</b>	<b>An evolutionary approach to the dynamic set cover problem</b>	<b>8</b>
2.1	The dynamic set cover problem . . . . .	8
2.2	The UMDA algorithm . . . . .	9
2.3	Java code . . . . .	10
<b>3</b>	<b>Theoretical analysis of variable-length evolutionary systems</b>	<b>12</b>
<b>4</b>	<b>Some Results about the Markov Chains Associated to GPs and General EAs.</b>	<b>14</b>
4.1	Introduction . . . . .	14
4.2	Notation . . . . .	15
4.3	How Does a Heuristic Search Algorithm Work? . . . . .	19
4.4	The Markov Chain Associated to an Evolutionary Algorithm . . . . .	19
4.5	A Special Kind of Reproduction Steps and the Extended Geiringer Theorem . . . . .	20
4.6	Nonlinear Genetic Programming (GP) with Homologous Crossover. . . . .	24
4.7	The Statement of the Schema-Based Version of Geiringer's Theorem for Non-linear GP under Homologous Crossover. . . . .	29
4.8	How Do We Obtain Theorem 4.7.1 from Theorem 4.5.2? . . . . .	33
4.9	What Does Theorem 4.5.2 Tell Us in the Presence of Mutation for Non-linear GP? . . . . .	39
4.10	What Can Be Said in the Presence of Selection in the General Case? . . . . .	49
4.11	What are the relations $\triangleright$ and $\trianglerighteq$ for the case of fitness-proportional selection? . . . . .	53
4.12	What Can Be Said when the Last Elementary Step is Mutation? . . . . .	61
4.13	Conclusions . . . . .	64

# Executive Summary

This report is the final report from UBHAM on WP8 (sub-task S4, Population dynamics in the evolutionary environment) brought forward from month 36. It is a summary of work done at UBHAM during the period January–April 2006.

The report is in four main parts. The first part introduces the work and describes its relationship to the rest of the project. The second part outlines our work on extending evolutionary algorithms to cope with changing requirements (the dynamic set cover problem). The third part gives an overview of our continuing theoretical analysis of evolutionary systems with variable-length structures (including trees). The main theoretical results are presented in the fourth part, which is a technical paper, to be published in the journal *Theoretical Computer Science*.

The main “customers” within the DBE project of this work have been STU and, indirectly, Intel and Sun. The applications of our research have been of some assistance in the creation of the EvE architecture and algorithms. Most of the work, however, will be of long-term interest with an impact beyond the lifetime of the project.

# Chapter 1

## Introduction and relationship to the project

This deliverable is the final report for sub-task S4 (Population dynamics in the evolutionary environment), which is scheduled to run from month 19 to month 36, and builds on the work of S3, already reported [21] and the preliminary report from S4 [20]. The objectives of this sub-task are:

1. To investigate the effects of a changing environment on evolution. Such changes may come from external forces (e.g. changes in users' requirements) or internally (e.g. from the migration of information).
2. To inform the development of DBE Evolutionary Environment, by liaising with STU, Intel and Sun with regard to representations, operators and fitness definitions.
3. To study population dynamics of variable-sized structures. This will involve both empirical studies, as well as a theoretical extension of existing work on limit theorems for variable-sized strings.

In our previous reports [21, 20] we described a formalisation of the problem faced by the EvE as a (weighted) set covering problem, and showed that an evolutionary approach could be effective on such a problem. We then planned to build on this work in the following ways:

1. Working with STU to help with their initial implementation of the EvE.
2. Extending the evolutionary approach to deal with changing requirements.
3. Empirical and theoretical studies of the extended algorithm.
4. Continuing collaboration with STU (see [18]) and HWU (see [1]) in the design of a more sophisticated EvE (including visits to STU).

Unfortunately, due to personnel problems, UBHAM has had to withdraw from the project early. We have made some contribution to STU's development of the EvE, including collaborative visits. We have also begun work on designing the extension to deal with dynamic requests. This work is presented in chapter 2. However, the empirical and theoretical studies have not been completed, although the java code has been handed over to STU. We are also unable to assist STU with their further development.

We have made more progress, however, in our theoretical studies on variable-length evolution. Since it is envisioned that the extended algorithm may have to deal with such structures to represent SBVR statements (e.g. in tree form) this is a potentially rich area of research [19]. We summarise our results in chapter 3. Chapter 4 contains a copy of our most recent publication on this subject (to appear in *Theoretical Computer Science*). It is highly technical, and is included for completeness, since this is our final deliverable.

## Chapter 2

# An evolutionary approach to the dynamic set cover problem

### 2.1 The dynamic set cover problem

We briefly recall the set cover problem as a model for the DBE EvE. We begin by assuming a large collection of *features*. These are individual things that a consumer might want. We assume they are atomic. A *request* from a user is (in the simplest form) a list of desired features. Service providers make available services to users. A service comprises, amongst other things, a description of all the features which it can provide. So, from an abstract perspective, we think of a service as being a subset of features. The problem faced by the EvE is to find a collection of services so that, when taken as a whole, they provide all the features requested. We also want this to be done as cheaply as possible. We assume each service has an associated cost, and the user would like the cheapest collection of services that meets their request.

This problem can be solved effectively using an evolutionary approach, as we have seen in our previous reports. What we now consider is the *dynamic* set cover problem. Over time, a user may produce a sequence of different requests. Similarly, a set of users with a common (or similar) local service pool, may make different, but related requests. If the different requests were all unrelated to each other, then the most efficient thing to do would be to re-run the evolutionary algorithm from scratch on a random population. However, it is assumed that requests from related users will be somehow similar to each other. For example, there may be some features, or feature combinations, which frequently occur in requests from users of a certain type. We want to be able to exploit this commonality to make the evolutionary algorithm more efficient.

One way to approach this problem is to make sure the initial population used is biased towards services which occur frequently in response to user requests. We propose a probability vector  $v$  which assigns a certain probability to each service in the genepool. When a new request arrives, we generate the new initial population as follows. We suppose  $N$  is the population size, and there are  $S$  services in the genepool). Then to create a new member of the population, we create a binary string of length  $S$ ,



where the probability of setting the  $k^{th}$  bit to 1 is  $v_k$ . We repeat this  $N$  times to fill up the population. Each string is interpreted as a subset of services, where a 1 in position  $k$  indicates that service  $k$  is included.

When a new habitat is first created, it is given a genepool. It must also then get a vector  $v$  to help it construct populations. It could take this vector from a neighbouring habitat that corresponds to a similar user type. Or, if there is no suitable neighbour, it could start off by setting  $v_k$  to be a small probability such as  $1/N$ .

After a request comes from the user, the population is initialised using  $v$ . The evolutionary algorithm (see next section) is then run to produce a solution. That solution will be a collection of services that meets the request as cheaply as possible. We now want to update  $v$  so that the services that appeared in the solution become more likely to appear in the future. One scheme is to set

$$v'_k = \begin{cases} (1 - \alpha)v_k + \alpha & \text{if } k \text{ is in the solution} \\ \alpha v_k & \text{if } k \text{ is not in the solution} \end{cases}$$

where  $0 < \alpha < 1$  is a learning rate. When the next request comes in, the new version of  $v$  is used to create the initial population, which is now biased towards the more frequently used services. This form of selection is related to the population “thinning” algorithm proposed in [18], which acts as a pre-selection phase in the EvE. We propose that such a method be included in future versions.

## 2.2 The UMDA algorithm

The scheme described in the previous section incrementally modifies the frequencies of services in the initial population. We thought it might be appropriate, then, to investigate the use of an evolutionary algorithm which also operates on the frequencies of the services. The UMDA (‘Univariate Marginal Distribution Algorithm’) is such an evolutionary algorithm. It was invented by Heinz Muhlenbein [7] and has been well-analysed from a theoretical perspective [6, 23]. The algorithm follows three phases: selection, generation, mutation.

**Selection** A subset of individuals from the current population are selected depending on their fitness.

**Generation** For each service  $k$ , denote by  $p_k$  the frequency with which it appears in the selected population. Generate a new population using the probability vector  $p$ .

**Mutation** Mutate each bit of each string in the population with probability  $\mu$ .

There are several ways in which to perform selection. One may use the standard proportional or tournament selection schemes to select a number of individuals. Or one may select a certain fraction of the best individuals (e.g. take the best half of the population). The generation stage proceeds exactly as with the initial population, only now we use the vector  $p$ . The mutation rate is typically set to  $1/S$ , so that, on average, one service is included or deleted from a solution.

As a simple example, suppose that our population is as follows:

0	0	1	1	0	0	1	fitness = 100
1	0	1	1	0	0	0	fitness = 90
1	0	1	0	0	0	1	fitness = 75
0	1	0	1	0	0	1	fitness = 50
1	0	0	0	0	1	0	fitness = 30
0	1	1	1	0	0	1	fitness = 25

That is, the first string indicates a subset containing the third, fourth and seventh service, and so on. The population has been ranked in order of fitness. We now select the best half of the population, namely the first three strings. For each service (bit position) we consider the frequency with which it appears in these top three strings. For example, the first service occurs in the second and third string, but not the first string. Its frequency is therefore  $2/3$ . The frequencies for each service are:

<i>service</i>	1	2	3	4	5	6	7
<i>frequency</i>	$2/3$	0	1	$2/3$	0	0	$2/3$

We now generate a new population by sampling the services according to these probabilities. That is, every time we generate a new string, the first bit is set to a one with probability  $2/3$  and so on. So we might get:

1	0	1	1	0	0	1
1	0	1	0	0	0	1
1	0	1	1	0	0	0
0	0	1	1	0	0	1
0	0	1	1	0	0	1
1	0	1	0	0	0	0

We then apply mutation. That is, each bit has a probability  $\mu$  of changing. The value of  $\mu$  is usually set to be  $1/S$ , which in this case is  $1/7$  so that, on average, one bit is changed per string. The resulting population might be:

1	0	0	1	0	0	1
1	1	1	0	0	0	1
1	0	1	1	0	0	0
0	0	1	1	1	0	1
0	0	1	1	0	1	1
1	0	1	0	1	1	0

The fitnesses are then calculated, and the population ranked accordingly, and the next generation begins.

## 2.3 Java code

Java code implemented the UMDA algorithm can be found at

<http://www.cs.bham.ac.uk/~jer/scp.tar>.

After unpacking the tar file, the program must be compiled using

```
javac pop.java
```

which can then be run. Various parameters can be adjusted in the file

```
parameters.java.
```

This code is preliminary and for experimental purposes only.

## Chapter 3

# Theoretical analysis of variable-length evolutionary systems

We have previously reported several theoretical results in the analysis of evolutionary algorithms. These were generally in two areas:

- Results relating to variable-length structures (in the infinite population limit)
- Results relating to finite populations for strings.

As an example, we proposed and proved a generalisation of the classic Geiringer’s theorem for finite populations. Simply put, this states that the repeated action of crossover is to destroy the correlations between elements of solutions, so that a population tends to a state of “linkage equilibrium”. It is worth pointing out that the generation stage of the UMDA algorithm takes the population in a single step to this state, as it explicitly de-couples the linkage between services in a solution. UMDA can therefore be seen (at least theoretically) as applying repeated strong crossover at each generation.

In our current work, we have continued to look at these themes. In the following chapter, we present a technical paper containing our results, which we summarise here.

Firstly, we return again to Geiringer’s theorem [4], but now considering the evolution of variable size and shape individuals, namely trees. This, therefore, applies to standard Genetic Programming (in which we evolve programs in the form of trees). It also potentially applies to the DBE EvE if, as is envisaged, the types of requests and service combinations allowed become more sophisticated than simple lists and subsets. For example, a hierarchical workflow pattern could be represented as a tree structure.

We establish, for tree structures, the kinds of shapes (or *schemata* as they are technically known) which play a fundamental role in the analogue of Geiringer’s theorem, and show how crossover again has the effect of de-coupling the correlations that exist in the population.

On the issue of bloat (that is, uncontrolled growth of structures during evolution), there are a number of possible counter-measures that can be employed. The simplest

and strictest is to disallow items larger than a certain threshold. Less severe, but still effective, is to have a certain probability that items larger than the threshold will be killed off. Such a method (called the *Tarpeian method*) has been proposed and theoretically motivated by Riccardo Poli [10]. An even less strict method is to have a penalty function that penalises large structures, but this can be difficult to design correctly. See [8] for a survey of methods.

Secondly, we continue to look at general evolutionary algorithms from the Markov chain perspective, and begin to build a framework in which bounds on the *stationary distribution* can be proved. This distribution gives the probability that a given population (or a population with any desired property) will be found by the evolutionary algorithm in the long-run. The framework relies on the construction of a kind of *order* relation (technically a *pre-order*) between populations. Once such an order has been appropriately constructed, then one can use it to determine if one population is more likely than another to appear in the stationary distribution.

A simple application of this framework establishes, for example, that if we have an evolutionary algorithm which goes through the phases: mutation, selection, crossover (in that order), then it is impossible that the stationary distribution can be uniform. There are always some populations which are preferred to other. This is even true even if all individuals actually have the same fitness! We ascribe this to the implicit *biases* induced by crossover, which tends to “prefer” some types of population over others.

One of the practical consequences of this observation is that one has to be rather careful when designing crossover operators that one is aware that such biases are being introduced. For example, it is known that certain crossovers for variable-length string structures have biases towards strings of a certain length. That is, such strings will tend to be “over sampled” in the long run.

Although we have produced a considerable amount of theoretical work in this period, it is true to say that at the moment it is still largely of only theoretical interest. We include our paper in the next chapter for completeness, rather than expecting it to be of direct value to the project. We had hoped to use the remaining time of the project to use such results to help design appropriate crossover operators for structured services within the EvE but, with our withdrawal from the project, this must be left to others, or to some future work.

## Chapter 4

# Some Results about the Markov Chains Associated to GPs and General EAs.

*This chapter is a paper accepted for publication in Theoretical Computer Science.*

### 4.1 Introduction

Geiringer's classical theorem (see [4]) is an important part of GA theory. It has been cited in a number of papers: see, for instance, [12], [13], [17] and [22]. It deals with the limit of the sequence of population vectors obtained by repeatedly applying the crossover operator  $\mathcal{C}(p)_k = \sum_{i,j} p_i p_j r_{(i,j \rightarrow k)}$  where  $r_{(i,j \rightarrow k)}$  denotes the probability of obtaining the individual  $k$  from the parents  $i$  and  $j$  after crossover. In other words, it speaks to the limit of repeated crossover in the case of an infinite population. In [5], a new version of this result was proved for *finite* populations, addressing the limiting distribution of the associated Markov chain, as follows. Let  $\Omega = \prod_{i=1}^n A_i$  denote the search space of a given genetic algorithm (intuitively  $A_i$  is the set of alleles corresponding to the  $i^{\text{th}}$  gene and  $n$  is the chromosome length). Fix a population  $P$  consisting of  $m$  individuals with  $m$  being an even number.  $P$  can be thought of as an  $m$  by  $n$  matrix whose rows are the individuals of the population  $P$ . Write

$$P = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix}.$$

Notice that the elements of the  $i^{\text{th}}$  column of  $P$  are members of  $A_i$ . Continuing with the notation used in [12], denote by  $\Phi(h, P, i)$  where  $h \in A_i$  the proportion of rows, say  $j$ , of  $P$  for which  $a_{ji} = h$ . In other words, let  $R_h = \{j \mid 1 \leq j \leq m \text{ and } a_{ji} = h\}$ .

Now simply let  $\Phi(h, P, i) = \frac{|R_h|}{m}$ . The classical Geiringer theorem (see [4] or, [12] for modern notation) says that if one starts with a population  $P$  of individuals and runs a genetic algorithm (GA) in the absence of selection and mutation (crossover being the only operator involved) then, in the “long run”, the frequency of occurrence of the individual  $(h_1, h_2, \dots, h_n)$  before time  $t$ , call it  $\Phi(h_1, h_2, \dots, h_n, t)$ , approaches independence:

$$\lim_{t \rightarrow \infty} \Phi(h_1, h_2, \dots, h_n, t) = \prod_{i=1}^n \Phi(h, P, i).$$

Thereby, Geiringer’s theorem tells us something about the limiting frequency with which certain elements of the search space are sampled in the long run, provided one uses crossover alone. In [12] this theorem has been generalized to cover the cases of variable-length GA’s and homologous linear genetic programming (GP) crossover. The limiting distributions of the frequency of occurrence of individuals belonging to a certain schema under these algorithms have been computed. The special conditions under which such a limiting distribution exists for linear GP under homologous crossover have been established (see theorem 9 and section 4.2.1 of [12]). In [5] a rather powerful extension of the finite population version of Geiringer’s theorem has been established. In the current paper we shall use the recipe described in [5] to derive a version of Geiringer’s theorem for nonlinear GP with homologous crossover (see section 4.6 or [9] for a detailed description of how nonlinear GP with homologous crossover works) which is based on Poli hyperschemata (see section 4.6 or [9]). The first step in this procedure is to describe the search space and the appropriate family of reproduction transformations so that the resulting GP algorithm is bijective and self-transient in the sense of definition 5.2 of [5]. Then the generalized Geiringer theorem (theorem 5.2 of [5]) as well as corollaries 6.1 and 6.2 of [5] apply. The necessary details are summarized in the next few sections. A schema based version of Geiringer’s theorem for nonlinear GP applies even in the presence of “node-mutation” (see section 4.9).

The finite population Geiringer theorem established in [5] may completely describe the stationary distribution of the Markov chain associated to an evolutionary algorithm only in the absence of selection. In section 4.10 we introduce a pre-order relation on the states of a Markov chain associated to an evolutionary algorithm which is defined in terms of selection alone, and establish some general inequalities about the stationary distribution of this Markov chain when selection is the “last stage” in the cycle. In section 4.12 we demonstrate that the stationary distribution of the Markov chain associated to most evolutionary algorithms in the presence of selection can never be uniform when mutation rate is small enough, even if the fitness function is constant.

The material in sections 4.10, 4.11 and 4.12 is independent of the results in sections 5 - 9. Thus, the reader has an option of jumping to read section 4.10 right after section 4.

## 4.2 Notation

$\Omega$  is a finite set, called a *search space*.

$f : \Omega \rightarrow (0, \infty)$  is a function, called a *fitness function*. The goal is to find a maximum of the function  $f$ .

$\mathcal{F}_q$  is a collection of  $q$ -ary operations on  $\Omega$ . Intuitively  $\mathcal{F}_q$  can be thought of as the collection of reproduction operators: some  $q$  parents produce one offspring. In nature often  $q = 2$ , for every child has two parents, but in the artificial setting there seems to be no special reason to assume that every child has no more than two parents. When  $q = 1$ , the family  $\mathcal{F}_1$  can be thought of as asexual reproductions or mutations. The following definitions will be used in section 4.3 to describe the general evolutionary search algorithm. This approach makes it easy to state the Geiringer Theorem.

**Definition 4.2.1** A population  $P$  of size  $m$  is simply an element of  $\Omega^m$ . (Intuitively it is convenient to think of a population as a “column vector”.)

**Remark 4.2.1** There are 2 primary methods for representing populations: multi-sets and ordered multi-sets. Each has advantages, depending upon the particular analytical goals. Lothar Shmitt has published a number of papers which use the ordered multi-set representation to advantage (see, for instance, [15] and [16]). According to definition 4.2.1, in the current paper we continue the development of analysis based upon the presentation pioneered by Lothar Schmitt. The following example illustrates an aspect of the representation which the reader would do well to keep in mind:

**Example 4.2.1** Let  $\Omega = \{0, 1\}^3$ . Consider the populations  $\begin{pmatrix} 0 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$ ,  $\begin{pmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{pmatrix}$  and  $\begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 0 & 0 & 0 \end{pmatrix}$ . According to definition 4.2.1 (the ordered multi-set model which is exploited in the current paper) these are distinct populations despite the fact that they represent the same population under the multi-set model.

An *elementary step* is a probabilistic rule which takes one population as an input and produces another population of the same size as an output. For example, the following elementary step corresponds to the fitness-proportional selection which has been studied in detail by Wright and Fisher (see [24] and [3]).

**Definition 4.2.2** An elementary step of type 1 (alternatively, of type selection) takes

a given population  $P = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix}$  with  $x_i \in \Omega$  as an input. The individuals of  $P$  are evaluated:

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix} \rightarrow \begin{pmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_m) \end{pmatrix}$$



A new population

$$P' = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$$

is obtained where  $y_i$ 's are chosen independently  $m$  times from the individuals of  $P$  and  $y_i = x_j$  with probability  $\frac{f(x_j)}{\sum_{l=1}^m f(x_l)}$ .

In other words, all of the individuals of  $P'$  are among those of  $P$ , and the expectation of the number of occurrences of any individual of  $P$  in  $P'$  is proportional to the number of occurrences of that individual in  $P$  times the individual's fitness value. In particular, the fitter the individual is, the more copies of that individual are likely to be present in  $P'$ . On the other hand, the individuals having relatively small fitness value are not likely to enter into  $P'$  at all. This is designed to imitate the natural survival of the fittest principle.

Population  $P'$  is the output of this elementary step.

In order to define an elementary step of type 2 (reproduction) in a general setting which uses the ordered multi-set representation (see remark 4.2.1 and example 4.2.1) one needs to introduce the following definitions:

**Definition 4.2.3** Fix an ordered  $k$ -tuple of integers  $\mathbf{q} = (q_1, q_2, \dots, q_k)$ . Let  $K$  denote a partition of the set  $\{1, 2, \dots, m\}$  for some  $m \in \mathbb{N}$ . We say that partition  $K$  is  $\mathbf{q}$ -fit if every element of  $K$  consists of exactly  $q_i$  elements for some  $i$ . In logical symbols this means that if  $K = \{P_1, P_2, \dots, P_l\}$  then  $K$  is  $\mathbf{q}$ -fit if  $\forall 1 \leq j \leq l \exists 1 \leq i \leq k$  such that  $|P_j| = q_i$ . Denote by  $\mathcal{E}_{\mathbf{q}}^m$  the family of all  $\mathbf{q}$ -fit partitions of  $\{1, 2, \dots, m\}$  (i.e.  $\mathcal{E}_{\mathbf{q}}^m = \{K \mid K \text{ is a } \mathbf{q}\text{-fit partition of } \{1, 2, \dots, m\}\}$ ).

**Definition 4.2.4** Let  $\Omega$  be a set,  $\mathcal{F}_{q_1}, \mathcal{F}_{q_2}, \dots, \mathcal{F}_{q_k}$  be some fixed families of  $q_j$ -ary operations on  $\Omega$  ( $\mathcal{F}_{q_j}$  is simply a family of functions from  $\Omega^{q_j}$  into  $\Omega$ ), and  $p_1, p_2, \dots, p_k$  be probability distributions on  $(\mathcal{F}_{q_1})^{q_1}, (\mathcal{F}_{q_2})^{q_2}, \dots, (\mathcal{F}_{q_k})^{q_k}$  respectively. Let  $\mathbf{q} = (q_1, q_2, \dots, q_k)$ . Finally, let  $\wp_m$  be a probability distribution on the collection  $\mathcal{E}_{\mathbf{q}}^m$  of partitions of  $\{1, 2, \dots, m\}$  (see definition 4.2.3 above). We then say that the ordered  $2(k+1)$ -tuple  $(\Omega, \mathcal{F}_{q_1}, \mathcal{F}_{q_2}, \dots, \mathcal{F}_{q_k}, p_1, p_2, \dots, p_k, \wp_m)$  is a reproduction  $k$ -tuple of arity  $(q_1, q_2, \dots, q_k)$ .

The following definition of reproduction covers both, crossover and mutation. Definition 4.2.6 (see also remark 4.2.2) will make it possible to combine different reproduction operators in a simple and natural way.

**Definition 4.2.5** An elementary step of type 2 (alternatively, of type reproduction) associated to a given reproduction  $k$ -tuple  $(\Omega, \mathcal{F}_{q_1}, \mathcal{F}_{q_2}, \dots, \mathcal{F}_{q_k}, p_1, p_2, \dots, p_k, \wp_m)$

takes a given population  $P = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix}$  with  $x_i \in \Omega$  as an input.

The individuals of  $P$  are partitioned into pairwise disjoint tuples for mating according to the probability distribution  $\wp_m$ . For instance, if the partition selected according to  $\wp_m$  is  $K = \{(i_1^1, i_2^1, \dots, i_{q_1}^1), (i_1^2, i_2^2, \dots, i_{q_2}^2), \dots, (i_1^j, i_2^j, \dots, i_{q_j}^j) \dots\}$  the corresponding tuples are

$$Q_1 = \begin{pmatrix} x_{i_1^1} \\ x_{i_2^1} \\ \vdots \\ x_{i_{q_1}^1} \end{pmatrix} \quad Q_2 = \begin{pmatrix} x_{i_1^2} \\ x_{i_2^2} \\ \vdots \\ x_{i_{q_2}^2} \end{pmatrix} \quad \dots \quad Q_j = \begin{pmatrix} x_{i_1^j} \\ x_{i_2^j} \\ \vdots \\ x_{i_{q_j}^j} \end{pmatrix} \quad \dots$$

Having selected the partition, replace every one of the selected  $q_j$ -tuples  $Q_j = \begin{pmatrix} x_{i_1^j} \\ x_{i_2^j} \\ \vdots \\ x_{i_{q_j}^j} \end{pmatrix}$

with the  $q_j$ -tuples

$$Q' = \begin{pmatrix} T_1(x_{i_1^j}, x_{i_2^j}, \dots, x_{i_{q_j}^j}) \\ T_2(x_{i_1^j}, x_{i_2^j}, \dots, x_{i_{q_j}^j}) \\ \vdots \\ T_{q_j}(x_{i_1^j}, x_{i_2^j}, \dots, x_{i_{q_j}^j}) \end{pmatrix}$$

for a  $q_j$ -tuple of transformations  $(T_1, T_2, \dots, T_{q_j}) \in (\mathcal{F}_{q_j})^{q_j}$  selected randomly according to the probability distribution  $p_j$  on  $(\mathcal{F}_{q_j})^{q_j}$ . This gives us a new population

$$P' = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$$

which serves as the output of this elementary step.

Notice that a single child does not have to be produced by exactly two parents. It is possible that a child has more than two parents. Asexual reproduction (mutation) is also allowed.

**Definition 4.2.6** A cycle is a finite sequence of elementary steps, say  $\{s_n\}_{n=1}^j$ , which are either of type 1 or of type 2 and such that all of the steps in the sequence  $\{s_n\}_{n=1}^j$  have the same underlying search space and the same arity of input/output.

**Remark 4.2.2** Intuitively, these steps are linked together in such a way that the output of the step  $s_i$  is the input of the step  $s_{i+1}$ . This is why all of the steps in the same cycle must have the same underlying search space and the same arity of input/output (otherwise the input/output relationship does not make sense).

We are finally ready to describe a rather wide class of evolutionary heuristic search algorithms.

### 4.3 How Does a Heuristic Search Algorithm Work?

A general evolutionary search algorithm works as follows: Fix a *cycle*, say  $C = \{s_n\}_{n=1}^j$  (see definition 4.2.6). Now start the algorithm with an initial population

$$P = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix} \quad \text{The initial population } P \text{ may be selected completely randomly, or it}$$

may also be predetermined depending on the circumstances. The actual method of selecting the initial population  $P$  is irrelevant for the purposes of the current paper. To run the algorithm with cycle  $C = \{s_n\}$ , simply input  $P$  into  $s_1$ , run  $s_1$ , then input the output of  $s_1$  into  $s_2$  ... input the output of  $s_{j-1}$  into  $s_j$  and produce the new output, say  $P'$ . Now use  $P'$  as an initial population and run the cycle  $C$  again. Continue this loop finitely many times depending on the circumstances.

**Definition 4.3.1** *A sub-algorithm of a given evolutionary search algorithm defined by a cycle  $C = \{s_n\}_{n=1}^j$  is simply an evolutionary search algorithm defined by a subsequence  $\{s_{n_q}\}_{q=1}^l$  of the sequence  $C$  of elementary steps.*

*A recombination sub-algorithm is sub-algorithm defined by a sequence of elementary steps of type 2 (Reproduction) only.*

### 4.4 The Markov Chain Associated to an Evolutionary Algorithm

In [22] it has been pointed out that heuristic search algorithms give rise to the following Markov process<sup>1</sup> (see also [2], for instance): The state space of this Markov process is the set of all populations of a fixed size  $m$ . This set, in our notation, is simply  $\Omega^m$ . The transition probability  $p_{\mathbf{x}\mathbf{y}}$  is simply the probability that the population  $\mathbf{y} \in \Omega^m$  is obtained from the population  $\mathbf{x}$  by going through the cycle once (where the notion of a cycle is described in section 4.3: see definition 4.2.6 and remark 4.2.2). The aim of the current paper is to establish a few rather general properties of this Markov chain. In case when there are several algorithms present in our discussion we shall write  $\{p_{\mathbf{x}\mathbf{y}}^{\mathcal{A}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m}$  to denote the Markov transition matrix associated to the algorithm  $\mathcal{A}$  while  $\{p_{\mathbf{x}\mathbf{y}}^{\mathcal{B}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m}$  would denote the Markov transition matrix associated to the algorithm  $\mathcal{B}$ .

**Definition 4.4.1** *Fix an evolutionary search algorithm  $\mathcal{A}$ . Denote by  $p_{\mathbf{x}, \mathbf{y}}^n$  the probability that a population  $\mathbf{y}$  is obtained from the population  $\mathbf{x}$  upon the completion of  $n$  complete cycles (in the sense of definition 4.2.6 and remark 4.2.2) of the algorithm. We say that a population  $\mathbf{x}$  leads to a population  $\mathbf{y}$  under  $\mathcal{A}$  if and only if  $p_{\mathbf{x}, \mathbf{y}}^n > 0$  for some  $n$ . We also write  $\mathbf{x} \xrightarrow{\mathcal{A}} \mathbf{y}$  as a shorthand notation for  $\mathbf{x}$  leads to  $\mathbf{y}$ . (This terminology is adopted from [2].)*

<sup>1</sup>In the current paper the state space of this process is slightly modified for technical reasons which will be seen later.

## 4.5 A Special Kind of Reproduction Steps and the Extended Geiringer Theorem

To understand the intuitive meaning of the definition below, see sections 4.2 and 4.3.

**Definition 4.5.1** Given a set  $\Omega$  and a family of transformations  $\mathcal{F}_q$  from  $\Omega^q$  into  $\Omega$ , fix a  $q$ -tuple of transformations  $(T_1, T_2, \dots, T_q) \in (\mathcal{F}_q)^q$ . Now consider the transformation  $\langle T_1, T_2, \dots, T_q \rangle : \Omega^q \rightarrow \Omega^q$  sending any given element

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_q \end{pmatrix} \in \Omega^q \text{ into } \begin{pmatrix} T_1(x_1, x_2, \dots, x_q) \\ T_2(x_1, x_2, \dots, x_q) \\ \vdots \\ T_q(x_1, x_2, \dots, x_q) \end{pmatrix} \in \Omega^q$$

We say that the transformation  $\langle T_1, T_2, \dots, T_q \rangle$  is the tupling of the ordered  $q$ -tuple  $(T_1, T_2, \dots, T_q)$ .

**Definition 4.5.2** Given an elementary step of type 2 (reproduction) associated to the reproduction  $k$ -tuple  $\Omega = (\Omega, \mathcal{F}_{q_1}, \mathcal{F}_{q_2}, \dots, \mathcal{F}_{q_k}, p_1, p_2, \dots, p_k, \wp_m)$ , fix some index  $i$  with  $1 \leq i \leq k$  and denote by

$$\mathcal{G}(\Omega, q_i) = \{ \langle T_1, T_2, \dots, T_q \rangle : \Omega^{q_i} \rightarrow \Omega^{q_i} \mid T_j \in \mathcal{F}_{q_i}, p_i(T_1, T_2, \dots, T_{q_i}) > 0 \}$$

the family of all tuplings which have a positive probability of being selected.

**Remark 4.5.1** The family of tupling transformations  $\mathcal{G}(\Omega, q_i)$  described in definition 4.5.2 represents the family of  $q$  parents  $\rightarrow q$  children crossover transformations while the family  $\mathcal{F}_q$  represents the family of  $q$  parents  $\rightarrow 1$  child crossovers. Depending on the circumstances it may be more convenient to specify the family of  $q$  parents  $\rightarrow q$  children crossover transformations directly rather than specifying the families  $\mathcal{F}_q$  individually. We shall see an example of this situation in section 4.6. The family  $\mathcal{F}_q$  of  $q$  parents  $\rightarrow 1$  child crossovers can then be recovered from the family of  $q$  parents  $\rightarrow q$  children crossover transformations by using coordinate projections.

As mentioned in section 4.2, in nature often the arity of the reproduction transformations is 2 meaning that every child has 2 parents.

It turns out that quite many evolutionary algorithms, including the classical genetic algorithm and nonlinear (as well as linear) genetic programming are equipped with the reproduction steps having the following nice property which has been introduced and investigated in [5].

**Definition 4.5.3** A given elementary step of type 2 (reproduction) associated to the reproduction  $k$ -tuple  $(\Omega, \mathcal{F}_{q_1}, \mathcal{F}_{q_2}, \dots, \mathcal{F}_{q_k}, p_1, p_2, \dots, p_k, \wp_m)$  is said to be bijective (and self-transient) if it satisfies conditions 1 (and 2) stated below:

1.  $\forall 1 \leq i \leq k$  we have  $p_i(T_1, T_2, \dots, T_{q_i}) > 0 \implies \langle T_1, T_2, \dots, T_{q_i} \rangle$  (see definition 4.5.1 for the meaning of  $\langle T_1, T_2, \dots, T_{q_i} \rangle$ ) is a bijection (a one-to-one and onto map of  $\Omega^{q_i}$  onto itself).

2.  $\forall 1 \leq i \leq k \exists (T_1, T_2, \dots, T_{q_i}) \in (\mathcal{F}_{q_i})^{q_i}$  such that  $p_i(T_1, T_2, \dots, T_{q_i}) > 0$  and  $\langle T_1, T_2, \dots, T_{q_i} \rangle = \mathbf{1}$  where  $\mathbf{1} : \Omega^{q_i} \rightarrow \Omega^{q_i}$  denotes the identity map (i. e.  $\forall \mathbf{x} \in \Omega^{q_i}$  we have  $\langle T_1, T_2, \dots, T_{q_i} \rangle(\mathbf{x}) = \mathbf{x}$ ). We say that a recombination sub-algorithm (see definition 4.3.1) of a given evolutionary search algorithm is bijective (and self-transient) if every given term of the subsequence,  $s_{n_k}$  by which the sub-algorithm is defined is bijective (and self-transient).

**Remark 4.5.2** Notice that conditions 1 and 2 of definition 4.5.3 can be restated in terms of the family  $\mathcal{G}(\Omega, q_i)$  as follows:

1. Every transformation in the family of tuplings,  $\mathcal{G}(\Omega, q_i)$  is a bijection.
2.  $\mathbf{1} \in \mathcal{G}(\Omega, q_i)$  where  $\mathbf{1} : \Omega^{q_i} \rightarrow \Omega^{q_i}$  denotes the identity map.

In [5] the following nice facts have been established:

**Proposition 4.5.1** Let  $\mathcal{A}$  denote a bijective and self-transient algorithm (see definition 4.5.3). Then  $\xrightarrow{\mathcal{A}}$  is an equivalence relation.

Proposition 4.5.1 motivates the following definition:

**Definition 4.5.4** Given a bijective and self-transient algorithm  $\mathcal{A}$  and a population  $P \in \Omega^m$ , denote by  $[P]_{\mathcal{A}}$  the equivalence class of the population  $P$  under the equivalence relation  $\xrightarrow{\mathcal{A}}$ .

To alleviate the level of abstraction we illustrate proposition 4.5.1 and definition 4.5.4 with a couple of examples.

**Example 4.5.1** Consider a binary genetic algorithm over the search space  $\Omega = \{0, 1\}^n$  under the action of crossover alone. Let the population size be some even number  $m$ . Consider the following family of masked crossover transformations:  $\mathcal{F} = \{F_M \mid M \subseteq \{1, 2, \dots, n\}\}$  where each  $F_M$  is a binary operation (i. e. a function from  $\Omega^2$  into  $\Omega$ ) defined as follows: For every  $\mathbf{a} = (a_1, a_2, \dots, a_n)$  and  $\mathbf{b} = (b_1, b_2, \dots, b_n) \in \Omega^n$ ,  $F_M(\mathbf{a}, \mathbf{b}) = \mathbf{x} = (x_1, x_2, \dots, x_n) \in \Omega^n$  where

$$x_i = \begin{cases} a_i & \text{if } i \in M \\ b_i & \text{otherwise} \end{cases}$$

Let  $\mathcal{A}$  denote the evolutionary algorithm determined by a single elementary step of type 2 (crossover) which is associated to the reproduction  $\frac{m}{2}$ -tuple

$$\Omega = (\Omega, \mathcal{F}, \mathcal{F}, \dots, \mathcal{F}, p_1, p_2, \dots, p_{\frac{m}{2}}, \wp_m)$$

(see definitions 4.2.4 and 4.2.5) where the probability distributions  $p_i$  have the property that  $p_i(F_M, F_K) \neq 0$  only if  $K = M^c$  (here  $M^c$  denotes the complement of  $M$  in  $\{1, 2, \dots, n\}$ ). This assumption on the distributions  $p_i$  ensures that the elementary step of crossover associated to the reproduction  $\frac{m}{2}$ -tuple  $\Omega$  is bijective. Depending on the further properties of the distributions  $p_i$  and the distribution  $\wp_m$ , different types of equivalence relations  $\xrightarrow{\mathcal{A}}$  would be induced. Typically, in case of a classical GA crossover, the distributions  $p_i$  are all identical (i. e.  $p_1 = p_2 = \dots = p_{\frac{m}{2}} = p$ ) where

$p$  is the uniform distribution on  $\{1, 2, \dots, n\}$  and the distribution  $\wp_m$  is uniform over all partitions of  $\{1, 2, \dots, n\}$  into 2-element subsets. In such a case the equivalence relation  $\xrightarrow{\mathcal{A}}$  is determined by the numbers of 0's in the columns (or, equivalently, by the numbers of 1's in the columns). The reason this is so is that a population  $Q$  can be reached from a population  $P$  in by performing a sequence of crossover elementary steps only if it has the same amount of “genetic material” in every column since alleles are neither lost nor created during homologous crossover. Using the fact that every permutation can be obtained by performing enough transpositions, one can show the converse of this fact. This fact is a particular case of lemma 47 of [5]. For instance, if  $n = 5$  and  $m = 4$  we have

$$\begin{pmatrix} 0 & 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 \end{pmatrix} \xrightarrow{\mathcal{A}} \begin{pmatrix} 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 1 \end{pmatrix}.$$

Indeed, the number of 0's in both populations in the first column is 3, in the 2<sup>nd</sup>, 3<sup>rd</sup> and 4<sup>th</sup> columns is 2 and in the last column is 1. Thus the equivalence class corresponding to a given population  $P$  can be described by an ordered  $n$ -tuple  $[c]_P = (c_1, c_2, \dots, c_n)$  of numbers between 0 and  $m$  where  $c_i$  is the number of 0s in the  $i^{\text{th}}$  column of  $P$ . For example, if  $P$  is either one of the equivalent populations above then  $[c]_P = (3, 2, 2, 2, 1)$ .<sup>2</sup>

**Example 4.5.2** Continuing with example 4.5.1, consider the following family of mutation transformations  $\mathcal{M} = \{T_{\mathbf{u}} \mid \mathbf{u} \in \Omega\}$  where each transformation  $T_{\mathbf{u}}$  is defined as follows: Denote by  $+_2$  the addition modulo 2 ( $0 +_2 0 = 0$ ,  $1 +_2 0 = 1$ ,  $0 +_2 1 = 1$ ,  $1 +_2 1 = 0$ ). We then define  $T_{\mathbf{u}}$  to be the function from  $\Omega$  into itself which sends every  $\mathbf{a} = (a_1, a_2, \dots, a_n)$  to  $T_{\mathbf{u}}(\mathbf{a}) = \mathbf{a} \oplus \mathbf{u}$  where  $\oplus$  is componentwise addition modulo 2, i. e. given  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  and  $\mathbf{y} = (y_1, y_2, \dots, y_n) \in \Omega^n$ , the  $\oplus$  operation is defined as follows:  $\mathbf{x} \oplus \mathbf{y} = \mathbf{z}$  where  $\mathbf{z} = (z_1, z_2, \dots, z_n)$  with  $z_i = x_i +_2 y_i$ . Notice first that every transformation  $T_{\mathbf{u}}$  is bijective (in fact  $T_{\mathbf{u}} \circ T_{\mathbf{u}} = \mathbf{1}$  where  $\mathbf{1}$  is the identity map on  $\Omega$ ). Since every mutation transformation  $T_{\mathbf{u}}$  is uniquely determined by the element  $\mathbf{u} \in \Omega$ , defining a probability distribution on the family  $\mathcal{M}$  amounts to defining a probability distribution on  $\Omega = \{0, 1\}^n$ . To achieve a situation equivalent to the classical case where every bit is mutated independently with a small probability  $\epsilon > 0$  and remains unchanged with probability  $1 - \epsilon$ , we choose 1 with probability  $\epsilon$  and 0 with probability  $1 - \epsilon$  independently  $n$  times. Given a population of size  $m$  we let mutation be the elementary step associated to the reproduction  $m$ -tuple

$$\Omega_{\text{mutation}} = (\Omega, \mathcal{M}, \mathcal{M}, \dots, \mathcal{M}, p, p, \dots, p, \wp_m)$$

where  $p$  is the probability distribution on  $\mathcal{M}$  described above and  $\wp_m$  is the unique trivial probability distribution on the one-element set (since there is exactly one way to partition a given set into singleton subsets). Now let  $\mathcal{B}$  denote the algorithm determined by the elementary step of crossover as described in example 4.5.1 followed by the

---

<sup>2</sup>One point crossover under reasonable assumptions will produce the same equivalence relation.

elementary step of mutation as described above. Then the algorithm  $\mathcal{B}$  is ergodic in the sense of definition 58 of [5] which means that the equivalence relation  $\xrightarrow{\mathcal{B}}$  is trivial, i.e. there is only one equivalence class or, in other words, for any two populations  $P$  and  $Q$  we have  $P \xrightarrow{\mathcal{B}} Q$ . Indeed, thanks to the availability of mutation, any given population can be reached from any other given population in a single step with a small but a positive probability which means that any two given populations are equivalent under  $\xrightarrow{\mathcal{B}}$ .

The main result of [5] is the following fact:

**Theorem 4.5.2** *Let  $\mathcal{A}$  denote a bijective and self-transient algorithm. Then the Markov chain initiated at some population  $P \in \Omega^m$  is irreducible and its unique stationary distribution is the uniform distribution (on  $[P]_{\mathcal{A}}$ ).*

The classical versions of Geiringer theorem, such as the ones established in [4] and in [12] are stated in terms of the “limiting frequency of occurrence” of a certain element of the search space. The following definitions, which also appear in [5], make these notions precise in the finite population setting:

**Definition 4.5.5** *We define the characteristic function  $\mathcal{X} : \Omega^m \times \mathcal{P}(\Omega) \rightarrow \mathbb{N} \cup \{0\}$  as follows:  $\mathcal{X}(P, S) =$  the number of individuals of  $P$  which are the elements of  $S$ . (Recall that  $P \in \Omega^m$  is a population consisting of  $m$  individuals and  $S \in \mathcal{P}(\Omega)$  simply means that  $S \subseteq \Omega$ .)*

**Example 4.5.3** *For instance, suppose  $\Omega = \{0, 1\}^n$ ,  $P = \begin{pmatrix} 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 \end{pmatrix}$  and  $S \subseteq \Omega = \{0, 1\}^n$  is determined by the Holland schema  $(*, 1, *, 1, *)$ . Then  $\mathcal{X}(P, S) = 3$  because exactly three rows of  $P$ , the 1<sup>st</sup>, the 2<sup>nd</sup>, and the 5<sup>th</sup> are in  $S$ .*

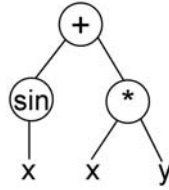
**Definition 4.5.6** *Fix an evolutionary algorithm  $\mathcal{A}$  and an initial population  $P \in \Omega^m$ . Let  $P(t)$  denote the population obtained upon the completion of  $t$  reproduction steps of the algorithm  $\mathcal{A}$  in the absence of selection and mutation. For instance,  $P(0) = P$ . Denote by  $\Phi(S, P, t)$  the proportion of individuals from the set  $S$  which occur before time  $t$ . That is,  $\Phi(S, P, t) = \frac{\sum_{s=1}^t \mathcal{X}(P(s), S)}{tm}$ . (Notice that  $tm$  is simply the total number of individuals encountered before time  $t$ . The same individual may be repeated more than once and the multiplicity contributes to  $\Phi$ .) Denote by  $\mathcal{X}(\square, S) : \Omega^m \rightarrow \mathbb{N}$  the restriction of the function  $\mathcal{X}$  when the set  $S$  is fixed (the notation suggests that one plugs a population  $P$  into the box).*

Intuitively,  $\Phi(S, P, t)$  is the frequency of encountering the individuals in  $S$  before time  $t$  when we run the algorithm starting with the initial population  $P$ .

## 4.6 Nonlinear Genetic Programming (GP) with Homologous Crossover.

In genetic programming, the search space,  $\Omega$ , consists of the parse trees which usually represent various computer programs.

**Example 4.6.1** A typical parse tree representing the program  $(+(\sin(x), *(x, y)))$  is drawn below:



Since computers have only a finite amount of memory, it is reasonable to assume that there are finitely many basic operations which can be used to construct programs and that every program tree has depth less than or equal to some integer  $L$ . Under these assumptions  $\Omega$  is a finite set. We may then define the search space as follows:

**Definition 4.6.1** Fix a signature  $\Sigma = (\Sigma_0, \Sigma_1, \Sigma_2, \dots, \Sigma_N)$  where  $\Sigma_i$ 's are finite sets.<sup>3</sup> We assume that  $\Sigma_0 \neq \emptyset$  and  $|\Sigma_j| \neq 1 \forall j$ <sup>4</sup>. The search space  $\Omega$  consists of all parse trees having depth at most  $L$ . Interior nodes having  $i$  children are labelled by the elements of  $\Sigma_i$ . The leaf nodes are labelled by the elements of  $\Sigma_0$ .

In order to study the appropriate family of reproduction (crossover) transformations with the aim of applying the generalized Geiringer theorem, it is most convenient to exploit Poli hyperschemata ([9] for a more detailed description).

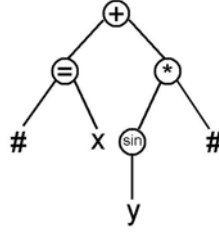
**Definition 4.6.2** A Poli hyperschema is a rooted parse tree which may have two additional labels for the nodes, namely  $\#$  and  $=$  signs (it is assumed, of course, that neither one of these denotes an operation). The  $=$  sign may label any interior node  $v$  of the tree. Since  $v$  does occur in the tree, we must have  $|\Sigma_i| > 0$ .) The  $\#$  sign can only label a leaf node. A given Poli hyperschema represents the set of all programs whose parse tree can be obtained by replacing the  $=$  signs with any operation of the appropriate arities and attaching any program trees in place of the  $\#$  signs. Different occurrences of  $\#$  or  $=$  may be replaced differently. We shall denote by  $S_t$  the set of programs represented by a hyperschema  $t$ .

Consider, for instance, the hyperschema  $t$  defined as  $(+(= (\#, x), *(\sin(y), \#)))$  which is pictured below:

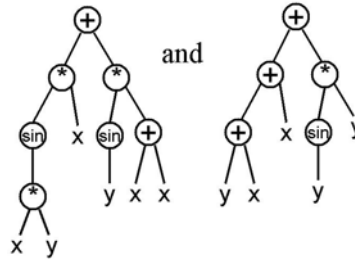
<sup>3</sup>Intuitively  $\Sigma_i$  is the set consisting of  $i$ -ary operations and  $\Sigma_0$  consists of the input variables. Formally this does not have to be the case though.

<sup>4</sup>The assumption that  $|\Sigma_j| \neq 1 \forall j$  does not cause any problems since we are free to select any elements from the search space that we want. On the other hand, this assumption helps us to avoid unnecessary complications when dealing with the poset of Poli hyperschemata later





A couple of programs fitting the hyperschema  $t$  are shown below:



In order to model the family of reproduction (crossover) transformation in a way which makes it obvious that GP is a bijective and self-transient algorithm, we shall introduce a partial order on the set of all Poli hyperschema so that every two elements have the least upper bound. The notion of the least upper bound will be also used to define the *common region* (see [11] for an alternative description of the notion of a common region).

**Definition 4.6.3** Denote by  $\mathcal{O}$  the set of all basic operations which can be used to construct the programs (i. e.  $\mathcal{O} = \Sigma_1 \cup \dots \cup \Sigma_N$ ) and by  $\mathcal{V}$  the set of all variables (i. e.  $\mathcal{V} = \Sigma_0$ ). Put the following partial order,  $\preceq$ , on the set  $\mathcal{O} \cup \mathcal{V} \cup \{=, \#\}$ :

1.  $\forall a, b \in \mathcal{O} \cup \mathcal{V}$  we have  $a \preceq b \iff a = b$ .
2.  $\forall a \in \mathcal{O}$  we have  $a \preceq =$ .
3.  $\forall a \in \mathcal{O} \cup \mathcal{V}$  we have  $a \preceq \#$ .
4.  $= \preceq =$ ,  $\# \preceq \#$  and  $= \preceq \#$ .

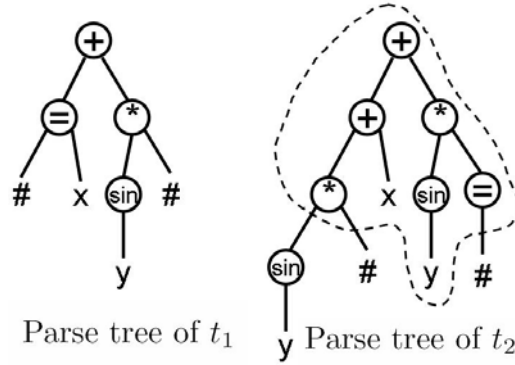
We shall also write  $a \succeq b$  to mean  $b \preceq a$ .

It is easy to see that  $\preceq$  is, indeed a partial order. Moreover, every collection of elements of  $\mathcal{O} \cup \mathcal{V} \cup \{=, \#\}$  has the least upper bound under  $\preceq$ . We are now ready to define the partial order relation on the set of all Poli hyperschemata:

**Definition 4.6.4** Let  $t_1$  and  $t_2$  denote two Poli hyperschemata. We say that  $t_1 \geq t_2$  if and only if the following two conditions are satisfied:

1. the tree corresponding to  $t_1$  when all of the labels are deleted is a subtree of the tree corresponding to  $t_2$  with all of the labels deleted.
2. Every one of the labels (which represents an operation or a variable) of  $t_1$  is  $\succeq$  the label of the node in the corresponding position of  $t_2$ .

**Example 4.6.2** For instance, the hyperschema  $t_1 = (+ (= (\#, x)), * (\sin(y), \#)) \geq t_2 = (+ (+ (* (\sin(x), y), x)), * (\sin(y), = (\#)))$ . Indeed, the parse trees of  $t_1$  and  $t_2$  appear on the picture below:



When all the labels in the dashed subtree of the parse tree of  $t_2$  are deleted one gets the tree isomorphic to that obtained from  $t_1$  by deleting all the labels. Thus condition 1 of definition 4.6.4 is satisfied. To see that condition 2 is fulfilled as well, we notice that the labels of  $t_1$  are  $\preceq$  to the corresponding labels of the dashed subtree of  $t_2$ : Indeed, we have  $+ \succeq +$ ,  $= \succeq +$ ,  $* \succeq *$ ,  $\# \succeq *$ ,  $x \succeq x$ ,  $\sin \succeq \sin$ ,  $\# \succeq =$  and  $y \succeq y$ .

Again it is easy to check that  $\geq$  is, indeed, a partial order relation on the collection of Poli hyperschemata. Proposition 4.6.1 below tells us even more:

**Proposition 4.6.1** Any given collection of Poli hyperschemata has the least upper bound under  $\geq$ .

*Proof:* Denote by  $\mathcal{S}$  a given collection of Poli hyperschemata. We provide an algorithm to construct the least upper bound of  $\mathcal{S}$  as follows: Copies of all the trees in  $\mathcal{S}$  are recursively jointly traversed starting from the root nodes to identify the parts with the same shape, i. e. the same arity in the nodes visited. Recursion is stopped as soon as an arity mismatch between corresponding nodes in some two trees from  $\mathcal{S}$  is present. All the nodes and links encountered are stored. This way we obtain a tree. It remains to stick in the labels. Each one of the interior nodes is labeled by the least upper bound of the corresponding labels of the trees in  $\mathcal{S}$ . The label of a leaf node is a variable, say  $x$ , if all the labels of the corresponding nodes of the trees in  $\mathcal{S}$  are  $x$  (which implies that they are leaf nodes themselves). In all other cases the label of the leaf node is the  $\#$  sign. It is not hard to see that this produces the least upper bound of the collection  $\mathcal{S}$  of parse trees.  $\square$

It was pointed out before, that programs themselves are Poli hyperschemata. The following fact is almost immediate from the explicit construction of the least upper bound carried out in the proof of proposition 4.6.1:

**Proposition 4.6.2** A given Poli hyperschema  $t$  is the least upper bound of the set  $S_t$  of programs determined by  $t$ .

From proposition 4.6.2 it follows easily that  $\geq$  is order isomorphic to the collection of subsets determined by the Poli hyperschemata:

**Proposition 4.6.3** *Let  $t$  and  $s$  denote Poli hyperschemata. Denote by  $S_t$  and  $S_s$  the subsets of the search space determined by the hyperschemata  $t$  and  $s$  respectively. Then  $t \geq s \iff S_t \supseteq S_s$ .*

There is another type of schemata which is useful to introduce in order to define the family of reproduction (crossover) transformations:

**Definition 4.6.5** *A shape schema is just a rooted ordered tree. If  $\tilde{t}$  is a given shape schema then  $S_{\tilde{t}}$  is just the set of all programs whose underlying tree when all the labels are deleted is precisely  $\tilde{t}$ . Given a Poli hyperschema  $s$ , we shall denote by  $\tilde{s}$  the underlying shape schema of  $s$ , i. e. the tree obtained by deleting all the labels in  $s$ .*

The notion of a common region which is equivalent to the one defined below also appears in [11]:

**Definition 4.6.6** *Given two Poli hyperschemata  $t$  and  $s$  we define their common region to be the underlying shape schema of the least upper bound of  $t$  and  $s$ .*

**Definition 4.6.7** *Fix a shape schema  $\tilde{t}$ . We shall say that the set  $C_{\tilde{t}} = \{(a, b) \mid a, b \text{ are program trees and } \tilde{t} \text{ is the common region of } a \text{ and } b\}$  is a component corresponding to the shape  $\tilde{t}$ .*

Notice that sets determined by the shape schemata partition the search space:

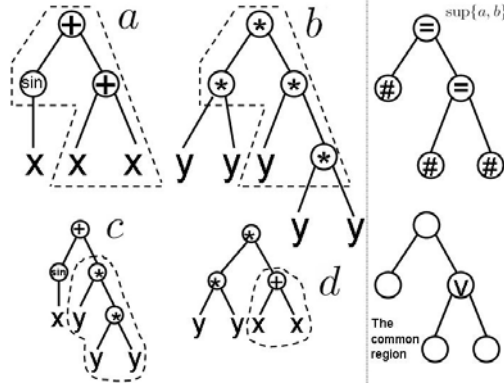
**Remark 4.6.1** *Notice that  $\Omega^2 = \bigcup_{\tilde{t} \text{ is a shape}} C_{\tilde{t}}$ . Moreover,  $C_{\tilde{t}} \cap C_{\tilde{s}} = \emptyset$  for  $\tilde{t} \neq \tilde{s}$ . (This is so because least upper bounds in a poset are uniquely determined and so the function sending  $(a, b) \rightarrow \sup(a, b) \rightarrow$  the underlying shape of  $\sup(a, b)$  is well defined. But then the sets  $C_{\tilde{t}}$  are simply the pre-images under this function of the individual shape schemata and, hence, form a partition of  $\Omega^2$ .)*

We now proceed to define the family of reproduction transformations. Our goal is to introduce a family of functions on  $\Omega^2$  in such a way that each one of them is easily seen to be bijective (see theorem 4.5.2, definition 4.5.2, definition 4.5.3 and remark 4.5.2). The idea is to define these transformations on each of the components first:

**Definition 4.6.8** *Fix a shape schema  $\tilde{t}$ . Fix a node,  $v$  of  $\tilde{t}$ . A one-point partial homologous crossover transformation  $T_v : C_{\tilde{t}} \rightarrow C_{\tilde{t}}$  is defined as follows: For given  $(a, b) \in C_{\tilde{t}}$  let  $T_v(a, b) = (c, d)$  where  $c$  and  $d$  are obtained from the program trees of  $a$  and  $b$  as follows: First identify the node  $v$  in the parse trees of  $a$  and  $b$  respectively. Now obtain the pair  $(c, d)$  by swapping the subtrees of  $a$  and  $b$  rooted at  $v$ . (This procedure is described in detail in [11] and it is also illustrated in the example below). Let  $\mathcal{G}_{\tilde{t}} = \{T_v \mid v \text{ is a node of } \tilde{t}\}$  denote the family of all partial homologous one-point crossover transformations associated to the shape  $\tilde{t}$ .*

The following example illustrates the concepts in definitions 4.6.5, 4.6.6 and 4.6.8:

**Example 4.6.3** In the upper left part of the picture parse trees of the two sample programs  $a$  and  $b$  are shown. Then on the upper right one can see the least upper bound of  $a$  and  $b$ . On the lower right the underlying tree of the least upper bound of  $a$  and  $b$  is drawn. According to definition 4.6.6, this tree is precisely the common region of the programs  $a$  and  $b$ . The isomorphic subtrees inside both,  $a$  and  $b$ , are emphasized inside the dashed areas:



A node  $v$  is selected inside the common region. The pair of children  $(c, d) = T_v(a, b)$  appears on the lower left of the picture above. The subtrees rooted at  $v$  which are swapped during crossover are emphasized inside the dashed area.

**Remark 4.6.2** One does need to show that for  $(a, b) \in C_{\tilde{t}}$  we have  $T_v(a, b) \in C_{\tilde{t}}$ . A rigorous argument can be given as follows: Clearly  $T_v : C_{\tilde{t}} \rightarrow \bigcup_{\tilde{t} \text{ is a shape}} C_{\tilde{t}}$  is a well-defined map. Moreover, since  $v$  is a node of the least upper bound of  $a$  and  $b$  and the pair  $(c, d)$  is obtained simply by swapping the corresponding subtrees rooted at  $v$ , we get  $s = \sup\{c, d\} \leq \sup\{a, b\}$ . Now consider the transformation  $F_v : C_{\tilde{s}} \rightarrow \bigcup_{\tilde{t} \text{ is a shape}} C_{\tilde{t}}$  and notice that, by definition, we have  $F_v(c, d) = (a, b)$ . But then, according to the reasoning above, we have  $\sup\{c, d\} \leq \sup\{a, b\}$ . Thereby, we get  $\sup\{c, d\} \leq \sup\{a, b\} \leq \sup\{c, d\} \implies \sup\{c, d\} = \sup\{a, b\} \implies \tilde{t} = \tilde{s}$ . This shows that  $T_v$  does, indeed, map into  $C_{\tilde{t}}$ . Moreover, in the process, we have also observed a couple of very important facts:

1.  $T_v \circ T_v = \mathbf{1}_{C_{\tilde{t}}}$  where  $\mathbf{1}_{C_{\tilde{t}}}$  denotes the identity map on  $C_{\tilde{t}}$ . This shows, in particular, that  $T_v$  is a bijection.
2.  $T_v$  preserves the least upper bounds:  $\sup\{a, b\} = \sup T_v(a, b)$ .

We are finally ready to define the family of reproduction transformations on the search space  $\Omega$  of all programs:

**Definition 4.6.9** For every shape schema  $\tilde{t}$  fix a node  $v_{\tilde{t}}$  of  $\tilde{t}$ . Define a one point crossover transformation  $T_{\{v_{\tilde{t}}\}_{\tilde{t} \text{ is a shape schema}}} : \Omega^2 \rightarrow \Omega^2$  to be the set-theoretic union of all partial crossover transformations of the form  $T_{v_{\tilde{t}}}$ . More explicitly, this means that whenever a given pair  $(a, b) \in \Omega^2$  we must have  $(a, b) \in C_{\tilde{s}}$  for a unique shape schema  $\tilde{s}$  (since, according to remark 4.6.1,  $\Omega^2$  is a disjoint union of components corresponding to various shapes). But then  $T_{\{v_{\tilde{t}}\}_{\tilde{t} \text{ is a shape schema}}}(a, b) = T_{v_{\tilde{s}}}(a, b)$ . Denote by

$\mathcal{G}$  the family of all crossover transformations together with the identity map on  $\Omega^2$ . For simplicity of notation we shall denote the transformations in  $\mathcal{G}$  by plain English letters:  $T$ ,  $F$  etc., keeping in mind that every such transformation is determined by making choices of partial crossover transformations on every one of the components.

**Remark 4.6.3** Thanks to remark 4.6.2, every one of the crossover transformations in the family  $\mathcal{G}$  is bijective (since it is a union of bijections on the pieces of a partition). It follows now that the generalized Geiringer theorem (theorem 4.5.2) applies to the case of homologous GP.

**Remark 4.6.4** It is also possible to model uniform GP crossover (this type of crossover is examined in detail in [11]) in the analogous manner. All of the results established in the current paper apply to this case without any modification.

## 4.7 The Statement of the Schema-Based Version of Geiringer's Theorem for Non-linear GP under Homologous Crossover.

As mentioned before, the schema-based version of Geiringer's theorem for non-linear GP is stated in terms of Poli hyperschemata.

**Definition 4.7.1** A Poli hyperschema of order  $i$  is a Poli hyperschema which has exactly  $i$  nodes whose label is not a  $\#$  or an  $=$  sign.

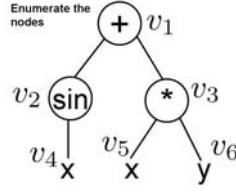
A configuration schema is a 0-order Poli hyperschema (i.e a hyperschema which has only the equal signs in the interior nodes and  $\#$  signs in the leaf nodes.)

An operation schema is a Poli hyperschema of order 1 (i. e. a hyperschema which has exactly one node whose label is not a  $\#$  or an  $=$  sign).

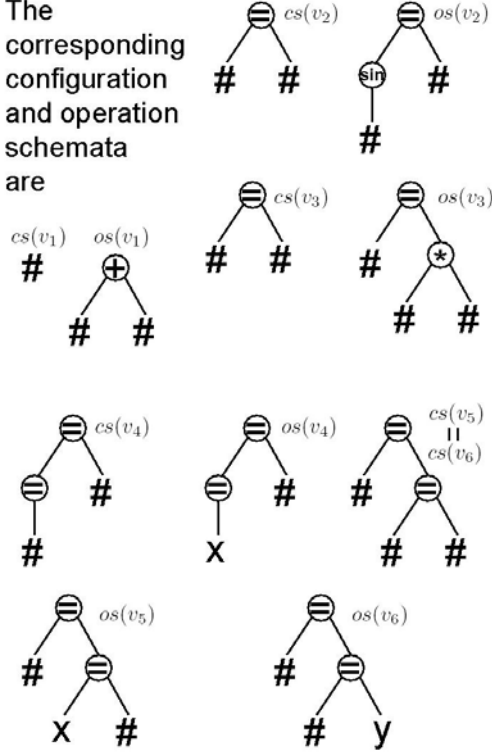
Fix an individual (a parse tree)  $\mathbf{u} \in \Omega$ . Let  $v$  denote any node of  $\mathbf{u}$ . Let  $B(v)$  denote the branch of the shape schema of  $\mathbf{u}$  from the root down to the node  $v$ . Let  $B^+(v) = B(v) \cup \{w \mid w \text{ is a child of some node } z \text{ of } B \text{ with } z \neq v\}$ . Now define  $cs(v)$  to be the configuration schema whose underlying shape schema is  $B^+(v)$ . Let  $o$  denote an operation or a variable (an element of  $\Sigma_i$  for some  $i$  between 0 and  $N$ ). Now obtain the operation schema  $os_o$  from  $cs(v)$  by attaching the node labelled by  $o$  in place of the  $\#$  sign at the node corresponding to  $v$  of  $cs(v)$ . Unless  $v$  is the leaf node of  $\mathbf{u}$ , all the children of this new node are the leaf nodes of  $os_o$  labelled by the  $\#$  sign. When  $o$  is the operation (or the variable) labelling the node  $v$  of  $\mathbf{u}$ , we shall write  $os(v)$  instead of  $os_o$ .

Notice that if  $v$  is a root node then  $cs(v)$  is just the schema which determines the entire search space, i. e. the parse tree consisting of a single node labelled by the  $\#$  sign. Example 4.7.1 illustrates definition 4.7.1.

**Example 4.7.1** Below we list all of the configuration schemata and operation schemata for the individual of example 4.6.1:



The corresponding configuration and operation schemata are



Recall from definition 4.5.5 that  $\mathcal{X}(P, S)$  denotes the number of individuals in the population  $P$  which are the elements of  $S \subseteq \Omega$ . The following definition makes it more convenient to state the schema-based version of Geiringer's theorem:

**Definition 4.7.2** Given a Poli hyperschema  $H$ , we shall write  $|H(P)|$  in place of  $\mathcal{X}(P, S_H)$  (see definition 4.6.2) to denote the number of individuals (counting repetitions) in the population  $P$  fitting the hyperschema  $H$ .

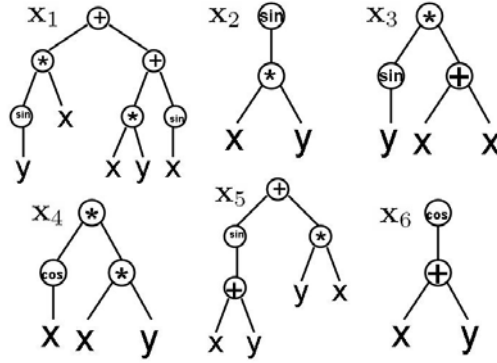
We can now finally state the Geiringer's theorem for non-linear GP under homologous crossover:

**Theorem 4.7.1** Fix an initial population  $P \in \Omega^m$  and an individual  $\mathbf{u} \in \Omega$ . Suppose every pair of individuals has a positive probability to be paired up for crossover and

every transformation in  $\mathcal{G}$  has a positive probability of being chosen<sup>5</sup>. Then the limiting frequency of occurrence of a given individual  $\mathbf{u}$ ,

$$\lim_{t \rightarrow \infty} \Phi(\mathbf{u}, P, t) = \prod_{v \text{ is a node of } \mathbf{u}} \frac{|os(v)(P)|}{|cs(v)(P)|}.$$

**Example 4.7.2** To illustrate how theorem 4.7.1 can be applied in practice, suppose we are interested in computing the frequency of encountering the individual  $\mathbf{u}$  from examples 4.6.1 and 4.7.1 when the initial population of 6 individuals pictured below is chosen:



The number of individuals in  $P$  fitting the operation schema  $os(v_1)$  is 2 (these are  $\mathbf{x}_1$  and  $\mathbf{x}_5$ ) while every individual fits the configuration schema  $cs(v_1)$ . Therefore  $\frac{|os(v_1)(P)|}{|cs(v_1)(P)|} = \frac{2}{6} = \frac{1}{3}$ . 4 individuals, namely  $\mathbf{x}_1$ ,  $\mathbf{x}_3$ ,  $\mathbf{x}_4$  and  $\mathbf{x}_5$  fit  $cs(v_2) = cs(v_3)$ , among these only 2 individuals, namely  $\mathbf{x}_3$  and  $\mathbf{x}_5$ , fit  $os(v_2)$  and 2 individuals,  $\mathbf{x}_4$  and  $\mathbf{x}_5$  fit  $os(v_3)$  so that  $\frac{|os(v_2)(P)|}{|cs(v_2)(P)|} = \frac{|os(v_3)(P)|}{|cs(v_3)(P)|} = \frac{2}{4} = \frac{1}{2}$ . Individuals  $\mathbf{x}_3$ ,  $\mathbf{x}_4$  and  $\mathbf{x}_5$  fit the configuration schema  $cs(v_4)$  while only  $\mathbf{x}_4$  fits the operation schema  $os(v_4)$  so that  $\frac{|os(v_4)(P)|}{|cs(v_4)(P)|} = \frac{1}{3}$ .  $\mathbf{x}_1$ ,  $\mathbf{x}_3$ ,  $\mathbf{x}_4$  and  $\mathbf{x}_5$  fit  $cs(v_5) = cs(v_6)$ . Among these only  $\mathbf{x}_3$  and  $\mathbf{x}_4$  fit  $os(v_5)$  while only  $\mathbf{x}_4$  fits  $os(v_6)$  so that  $\frac{|os(v_5)(P)|}{|cs(v_5)(P)|} = \frac{2}{4} = \frac{1}{2}$  and  $\frac{|os(v_6)(P)|}{|cs(v_6)(P)|} = \frac{1}{4}$ . Thereby, according to theorem 4.7.1, we obtain:

$$\begin{aligned} \lim_{t \rightarrow \infty} \Phi(\mathbf{u}, P, t) &= \prod_{i=1}^6 \frac{|os(v_i)(P)|}{|cs(v_i)(P)|} = \\ &= \frac{1}{3} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{3} \cdot \frac{1}{2} \cdot \frac{1}{4} = \frac{1}{288}. \end{aligned}$$

Roughly speaking, this means that if we run GP starting with the population  $P$  pictured above, in the absence of mutation and selection (crossover being the only step) for an infinitely long time, the individual  $\mathbf{u}$  will be encountered on average 1 out of 288 times.

<sup>5</sup>These conditions can be slightly relaxed, but we try to present the main idea only

**Example 4.7.3** Notice that linear GP (or, equivalently, variable length GA) as described in [12] is a special case of nonlinear GP when  $\forall i > 1 \Sigma_i = \emptyset$  and  $\Sigma_0$  and  $\Sigma_1 \neq \emptyset$ . Indeed, the elements of such a search space are parse trees such that every interior node has exactly one child and the depth of the tree is bounded by some integer  $N$ . One can think of such a tree as a sequence of labels  $(a_1, a_2, \dots, a_n)$ , the first label affiliated with the root node, second label with the child of the root node and so on. The label  $a_n$  is affiliated with the leaf node. This gives us a one-to-one correspondence, call it  $\phi$  between the search space for nonlinear GP in our specific case when  $\forall i > 1 \Sigma_i = \emptyset$  while  $\Sigma_0$  and  $\Sigma_1 \neq \emptyset$  and the search space for linear GP which preserves crossover. The following types of schemata have been introduced in [12]:

**Definition 4.7.3** The schema  $H = (*^{i-1}, h_i, \#)$  represents the subset  $S_H = \{\mathbf{x} = (x_1, x_2, \dots, x_l) \mid l > i \text{ and } x_i = h_i\}$ . In words,  $S_H$  is simply the set of all individuals whose length is at least  $i + 1$  and whose  $i^{\text{th}}$  allele is  $h_i$ .

**Definition 4.7.4** The schema  $H = (*^i, \#)$  represents the subset

$$S_H = \{\mathbf{x} = (x_1, x_2, \dots, x_l) \mid l > i\}.$$

In words,  $S_H$  is simply the subset of all individuals whose length is at least  $i + 1$ .

**Definition 4.7.5** The schema  $H = (*^{i-1}, h_i)$  represents the subset

$$S_H = \{\mathbf{x} = (x_1, x_2, \dots, x_i) \mid x_i = h_i\}$$

of the search space which is simply the set of all individuals of length exactly equal to  $i$  whose  $i^{\text{th}}$  (last) allele is  $h_i$ .

The reader may check that under the correspondence  $\phi$  the configuration schemata correspond to the schemata  $H_i = (*^i, \#)$  for  $i \geq 1$ , operation schemata correspond to the schemata of the form  $H = (*^{i-1}, h_i, \#)$  and of the form  $H = (*^{i-1}, h_i)$  for  $i > 1$ . Finally, the hyperschema  $t_{(1,1)}$  corresponds to the schema  $H = (h_1, \#)$ . Fix a population  $P \in \Omega^m$ . Recall that we denote by  $|H|$  the number of individuals in  $P$  which fit the schema  $H$  counting repetitions. Also recall from definition 4.5.6 that  $\Phi(S_H, P, 1) = \frac{|H|}{m}$  denotes the fraction of the number of individuals of  $P$  which fit the schema  $H$ . To abbreviate the notation we shall write  $\Phi(H, P, 1)$  instead of  $\Phi(S_H, P, 1)$ . Fix an individual  $\mathbf{u} = (h_1, h_2, \dots, h_n) \in \Omega$ . Theorem 4.7.1 tells us that

$$\begin{aligned} \lim_{t \rightarrow \infty} \Phi(\mathbf{u}, P, t) &= \frac{|(h_1, \#)|}{m} \cdot \left( \prod_{i=1}^{n-2} \frac{|(*^i, h_{i+1}, \#)|}{|(*^i, \#)|} \right) \cdot \frac{|(*^{n-1}, h_n)|}{|(*^{n-1}, \#)|} = \\ &= \frac{|(h_1, \#)|}{m} \cdot \left( \prod_{i=1}^{n-2} \frac{\frac{|(*^i, h_{i+1}, \#)|}{m}}{\frac{|(*^i, \#)|}{m}} \right) \cdot \frac{\frac{|(*^{n-1}, h_n)|}{m}}{\frac{|(*^{n-1}, \#)|}{m}} = \\ &= \Phi(h_1, \#) \cdot \left( \prod_{i=1}^{n-2} \frac{\Phi(*^i, h_{i+1}, \#)}{\Phi(*^i, \#)} \right) \cdot \frac{\Phi(*^{n-1}, h_n)}{\Phi(*^{n-1}, \#)} = \end{aligned}$$



$$= \Phi(*^{n-1}, h_n) \cdot \frac{\prod_{i=n-2}^0 \Phi(*^i, h_{i+1}, \#)}{\prod_{i=n-1}^1 \Phi(*^i, \#)} = \Phi(*^{n-1}, h_n) \cdot \prod_{i=n-1}^{i=1} \frac{\Phi(*^{i-1}, h_i, \#)}{\Phi(*^i, \#)}$$

which is precisely the formula obtained in [12].

## 4.8 How Do We Obtain Theorem 4.7.1 from Theorem 4.5.2?

The following couple of corollaries from [5] are useful in obtaining the schema-based versions of Geiringer theorem for various evolutionary algorithms. Throughout, we shall denote by  $\varrho_{[P]_{\mathcal{A}}}$  the uniform probability distribution on the set  $[P]_{\mathcal{A}}$  (see definition 4.5.4).

**Corollary 4.8.1** *Fix a bijective and self-transient algorithm  $\mathcal{A}$  and an initial population  $P \in \Omega^m$ . Fix a set  $S$  of individuals in  $\Omega$  ( $S \subseteq \Omega$ ). Then  $\lim_{t \rightarrow \infty} \Phi(S, P, t) = \frac{1}{m} E_{\varrho_{[P]_{\mathcal{A}}}}(\mathcal{X}(\square, S))$  (here  $E_{\varrho_{[P]_{\mathcal{A}}}}(f)$  denotes the expectation of the random variable  $f$  with respect to the uniform distribution on the set  $[P]_{\mathcal{A}}$ ).<sup>6</sup>*

To state the next corollary which brings us one step closer to deriving results similar in flavor to Geiringer's original theorem we need one more, purely formal, assumption about the algorithm:

**Definition 4.8.1** *We say that a given algorithm  $\mathcal{A}$  is regular if the following is true: for every population  $P = (x_1, x_2, \dots, x_m) \in \Omega^m$  and for every permutation  $\pi \in \mathcal{S}_m$ , the population obtained by permuting the elements of  $P$  by  $\pi$ , namely  $\pi(P) = (x_{\pi(1)}, x_{\pi(2)}, \dots, x_{\pi(m)}) \in [P]_{\mathcal{A}}$ . In words this says that the equivalence classes  $[P]_{\mathcal{A}}$  are permutation invariant.*

**Remark 4.8.1** *Definition 4.8.1 is only needed because our description of an evolutionary search algorithm uses the ordered multi-set model. This makes the generalized Geiringer theorem (theorem 4.5.2) look nice (the stationary distribution is uniform on  $[P]_{\mathcal{A}}$ ). A disadvantage of the multi-set model is that it allows algorithms which are not regular. If we were to use the model of [22] where the order of elements in a population is not taken into account (a reasonable assumption since most evolutionary algorithms used in practice are, indeed, regular) then the Generalized Geiringer theorem would have to be modified accordingly since the stationary distribution of the corresponding Markov process would be different from uniform (it is not difficult to compute it though since the corresponding Markov chain is just a "projection" of the one used in the current paper).*

**Corollary 4.8.2** *Fix a regular bijective and self-transient algorithm  $\mathcal{A}$  and an initial population  $P \in \Omega^m$ . Denote by  $\varrho_{[P]_{\mathcal{A}}}$  the uniform probability distribution on  $[P]_{\mathcal{A}}$  (see definition 4.5.4). Fix a set  $S$  of individuals in  $\Omega$  ( $S \subseteq \Omega$ ). Then we have  $\lim_{t \rightarrow \infty} \Phi(S, P, t) = \varrho_{[P]_{\mathcal{A}}}(\mathcal{V}_S)$  where*

$$\mathcal{V}_S = \{P \mid P \in [P]_{\mathcal{A}} \text{ and the } 1^{\text{st}} \text{ individual of } P \text{ is an element of } S\}.$$

---

<sup>6</sup>Throughout the paper, whenever a limit is involved, the equality is meant to hold for almost every infinite sequence of trials.

Corollaries 4.8.1 and 4.8.2 are proved in section 6 of [5]. When deriving schema-based versions of Geiringer theorem for a specific algorithm the following strategy may be implemented: Continuing with the notation in corollaries 4.8.1 and 4.8.2, suppose we are given a nested sequence of subsets of the search space:  $S_1 \supseteq S_2 \supseteq \dots \supseteq S_n$ . According to corollary 4.8.2,

$$\begin{aligned} \lim_{t \rightarrow \infty} \Phi(S_n, P, t) &= \varrho_{[P]_{\mathcal{A}}}(\mathcal{V}_{S_n}) = \frac{|\mathcal{V}_{S_n}|}{|[P]_{\mathcal{A}}|} = \frac{|\mathcal{V}_{S_n}|}{|\mathcal{V}_{S_{n-1}}|} \cdot \frac{|\mathcal{V}_{S_{n-1}}|}{|[P]_{\mathcal{A}}|} = \\ &= \frac{|\mathcal{V}_{S_n}|}{|\mathcal{V}_{S_{n-1}}|} \cdot \frac{|\mathcal{V}_{S_{n-1}}|}{|\mathcal{V}_{S_{n-2}}|} \cdot \dots \cdot \frac{|\mathcal{V}_{S_2}|}{|\mathcal{V}_{S_1}|} \cdot \frac{|\mathcal{V}_{S_1}|}{|[P]_{\mathcal{A}}|} = \\ &= \varrho_{[P]_{\mathcal{A}}}(\mathcal{V}_{S_1}) \cdot \prod_{j=0}^{n-2} \frac{|\mathcal{V}_{S_{n-j}}|}{|\mathcal{V}_{S_{n-j-1}}|} = \frac{1}{m} E_{\varrho_{[P]_{\mathcal{A}}}}(\mathcal{X}(\square, S)) \cdot \prod_{j=0}^{n-2} \frac{|\mathcal{V}_{S_{n-j}}|}{|\mathcal{V}_{S_{n-j-1}}|} \end{aligned}$$

Notice that  $\frac{|\mathcal{V}_{S_j}|}{|\mathcal{V}_{S_{j-1}}|}$  is just the proportion of populations in  $[P]_{\mathcal{A}}$  whose first individual is a member of  $S_j$  inside the set of populations in  $[P]_{\mathcal{A}}$  whose first individual is a member of  $S_{j-1}$ .

**Corollary 4.8.3** *Fix a regular, bijective and self-transient algorithm  $\mathcal{A}$  and an initial population  $P \in \Omega^m$ . Fix a nested sequence of subsets  $S_1 \supseteq S_2 \supseteq \dots \supseteq S_n$  of individuals in  $\Omega$  ( $S_1 \subseteq \Omega$ ). Then  $\lim_{t \rightarrow \infty} \Phi(S_n, P, t) = \frac{1}{m} E_{\varrho_{[P]_{\mathcal{A}}}}(\mathcal{X}(\square, S)) \cdot \prod_{j=0}^{n-2} \frac{|\mathcal{V}_{S_{n-j}}|}{|\mathcal{V}_{S_{n-j-1}}|}$  where, as before,  $\mathcal{V}_S$  denotes the set of all populations whose first individual is a member of  $S$  for a given subset  $S \subseteq \Omega$ .*

Denote by  $\mathcal{A}$  a given GP algorithm. Fix an individual  $\mathbf{u} \in \Omega$ . In order to apply corollary 4.8.3, we may choose a descending chain of Poli hyperschemata  $t_1 \geq t_2 \geq \dots \geq t_n = \mathbf{u}$ . Fix an initial population  $P$ . To avoid putting many subscripts, we shall write  $\mathcal{V}_t$  instead of  $\mathcal{V}_{S_t}$  for the set of all populations in  $[P]_{\mathcal{A}}$  (see definition 4.4.1) whose 1<sup>st</sup> individual is a member of  $S_t$  (the set of individuals determined by the hyperschema  $t$ ). In order to construct the desired sequence of nested hyperschemata, we assign the following numerical labelling to the nodes of the parse tree of  $\mathbf{u}$ : The nodes are labelled by the pairs of integer coordinates. The first coordinate shows the depth of the tree and the second coordinate shows how far to the right a given node at the depth specified by the first coordinate is located. Notice, for instance, that the root node is labelled by the coordinates (1, 1). We also introduce the following lexicographic linear ordering on the set of coordinate pairs:

**Definition 4.8.2**  $(a, b) \leq (c, d)$  if and only if either  $a \leq c$  or ( $a = c$  and  $b \leq d$ ).

It is well known and easy to verify that this defines a linear ordering.

**Definition 4.8.3** *Given a pair of coordinates  $(i, j)$ , denote by  $\uparrow(i, j)$  the immediate successor of  $(i, j)$  under the lexicographic ordering defined above. Explicitly,*

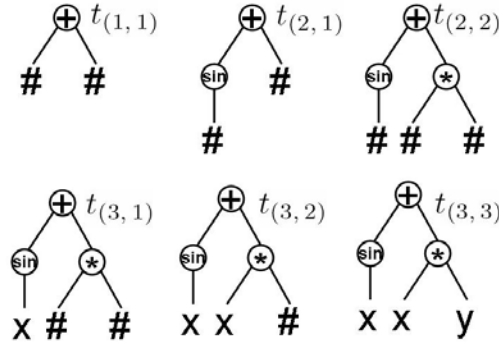
$$\uparrow(i, j) = \begin{cases} (i+1, 1) & \text{if } (i, j) \text{ labels the rightmost node of } \mathbf{u} \text{ at depth } i \\ (i, j+1) & \text{otherwise} \end{cases}$$

We obtain the desired nested sequence of hyperschemata for the given individual  $\mathbf{u}$  recursively in the following manner:

**Definition 4.8.4** Define  $t_{(1,1)}$  to be the hyperschema whose root node has the same label (operation) and arity as that of the root node of  $\mathbf{u}$ . All children of the root node are the leaf nodes labelled by the  $\#$  sign. Once the hyperschema  $t_{(i,j)}$  has been constructed, we obtain the hyperschema  $t_{\uparrow(i,j)}$  by attaching the node of  $\mathbf{u}$  with coordinate  $\uparrow(i,j)$  in place of the  $\#$  sign at coordinate  $\uparrow(i,j)$  to the parse tree of  $t_{(i,j)}$ . Unless this node, call it  $v$ , is a leaf node of  $\mathbf{u}$ , all children of this new node are the leaf nodes of  $t_{\uparrow(i,j)}$  labelled by the  $\#$  sign.

We illustrate the construction with an explicit example:

**Example 4.8.1** Below, the nested sequence  $t_{(1,1)} \geq t_{(2,1)} \geq t_{(2,2)} \geq t_{(3,1)} \geq t_{(3,2)} \geq t_{(3,3)}$  corresponding to the program of example 4.6.1 is drawn explicitly:



The formula for the limiting frequency of occurrence of a given program  $u$  in corollary 4.8.3 involves the ratios of the form  $\frac{\mathcal{V}_{t_{\uparrow(i,j)}}}{\mathcal{V}_{t_{(i,j)}}}$ . It turns out that these ratios can be expressed nicely in terms of the presence of certain configuration and operation schemata in the initial population  $P$ :

**Definition 4.8.5** Given a program tree  $\mathbf{u}$  and the corresponding nested sequence  $t_{(1,1)} \geq t_{(2,1)} \geq \dots \geq t_{(i,j)} \geq t_{\uparrow(i,j)} \geq \dots \geq t_{(l,k)} = \mathbf{u}$  of hyperschemata as in definition 4.8.4, for every  $(i,j) \neq (l,k)$ , denote by  $cs_{(i,j)}$  ( $os_{(i,j)}$ ) the configuration schema  $cs(v)$  (operation schema  $os(v)$ ) where  $v$  is the node of  $\mathbf{u}$  with coordinate  $\uparrow(i,j)$ .

**Example 4.8.2** Continuing with examples 4.6.1 and 4.7.1 notice that for the individual in these examples we have  $cs_{(1,1)} = cs_{(2,1)} = cs(v_2) = cs(v_3)$  while  $os_{(1,1)} = os(v_2)$  and  $os_{(2,1)} = os(v_3)$  (see example 4.7.1),  $cs_{(2,2)} = cs(v_4)$  while  $os_{(2,2)} = os(v_4)$  and  $cs_{(3,1)} = cs_{(3,2)} = cs(v_5) = cs(v_6)$  while  $os_{(3,1)} = os(v_5)$  and  $os_{(3,2)} = os(v_6)$ .

The following “orbit description” lemma is the reason for introducing configuration and operation schemata: We prove the lemma under the following special assumption:

**Definition 4.8.6** We say that a population  $P$  is special with respect to the individual  $\mathbf{u}$  if for every node  $v$  of  $\mathbf{u}$  and for every operation (or variable)  $o$  we have  $|os_o(P)| \leq 1$  where  $os_o$  is obtained from  $cs(v)$  by means of attaching the operation  $o$  at the leaf node of  $cs(v)$  corresponding to  $v$  as described in definition 4.7.1.

Definition 4.8.6 basically requires that no 2 operations (or variables) occurring in  $P$  at the specified location are the same. It turns out that the orbit description lemma stated below is a lot more convenient to prove under this special assumption. The general case will then follow by introducing enough extra labels for the operations and variables involved and then deleting the extra labels.

**Lemma 4.8.4** Fix an initial population  $P$  and a program  $\mathbf{u} \in \Omega$ . Assume that the population  $P$  is special with respect to the individual  $\mathbf{u}$ . Suppose every pair of individuals has a positive probability to be paired up for crossover and every transformation in  $\mathcal{G}$  has a positive probability of being chosen<sup>7</sup>. Consider the sequences of hyperschemata  $t_{(1,1)} \geq t_{(2,1)} \geq \dots \geq t_{(i,j)} \geq t_{\uparrow(i,j)} \geq \dots \geq t_{(l,k)} = \mathbf{u}$ ,  $\{cs_{(i,j)} \mid (i,j) \text{ is a coordinate of } \mathbf{u}, (i,j) \text{ is not the maximal coordinate}\}$  and  $\{os_{(i,j)} \mid (i,j) \text{ is a coordinate of } \mathbf{u}, (i,j) \text{ is not the maximal coordinate}\}$  corresponding to the individual  $\mathbf{u}$ . For a given hyperschema  $t$ , denote by  $|t(P)|$  the number of individuals in  $P$  which fit the hyperschema  $t$  counting repetitions. Suppose  $\forall$  non-maximal pairs of coordinates  $(i,j)$  we have  $|os_{(i,j)}(P)| \neq 0$  and  $|t_{(1,1)}(P)| \neq 0$ . Then it is true that 
$$\forall (i,j) \quad \frac{|\mathcal{V}_{t_{\uparrow(i,j)}}|}{|\mathcal{V}_{t_{(i,j)}}|} = \frac{1}{|cs_{(i,j)}(P)|}.$$

*Proof:* The key idea is to observe the following fact:

**Claim:** Fix a coordinate  $(i,j)$ . Fix any two operation schemata  $os_1$  and  $os_2$  which are obtained from  $cs_{(i,j)}$  by attaching either a variable or an operation at the node  $(i,j)$ . Suppose  $\exists$  individuals in  $P$  fitting both,  $os_1$  and  $os_2$ . Then  $|\mathcal{V}_{t_{(i,j)}} \cap \mathcal{V}_{os_1}| = |\mathcal{V}_{t_{(i,j)}} \cap \mathcal{V}_{os_2}|$ . *Proof:* Consider the map  $F : [P]_{\mathcal{A}} \rightarrow [P]_{\mathcal{A}}$  defined as follows: Given a population, say  $Q \in [P]_{\mathcal{A}}$ , notice that  $\exists$  an individual, say  $\mathbf{x}_1$ , in  $Q$  fitting the operation schema  $os_1$  (due to the way crossover is defined, the number of individuals fitting the operation schema  $os_1(Q)$  is the same in every population  $Q \in [P]_{\mathcal{A}}$ ). Moreover, such an individual is unique since we assumed that all operations appearing in the individuals of  $P$  are distinct. Likewise,  $\exists$  unique individual in  $Q$ , say  $\mathbf{x}_2$  fitting the operation schema  $os_2$ . Pair up individuals  $\mathbf{x}_1$  and  $\mathbf{x}_2$  and pair up the rest of the individuals arbitrarily for crossover. Select the crossover transformation  $T_v$  where  $v$  is the node with coordinate  $(i,j)$  for the pair  $(\mathbf{x}_1, \mathbf{x}_2)$  and choose the identity transformation for the rest of the pairs. Now let  $F(Q)$  be the population obtained upon the completion of the reproduction step described above (notice that  $F(Q) \in [P]_{\mathcal{A}}$  by definition of  $[P]_{\mathcal{A}}$ ). Notice also that  $F$  is its own inverse (i. e.  $F \circ F = \mathbf{1}_{[P]_{\mathcal{A}}}$ ). This tells us, in particular, that  $F$  is bijective. Moreover, it is clear from the definitions that  $F(\mathcal{V}_{t_{(i,j)}} \cap \mathcal{V}_{os_1}) \subseteq \mathcal{V}_{t_{(i,j)}} \cap \mathcal{V}_{os_2}$  and, likewise,  $F(\mathcal{V}_{t_{(i,j)}} \cap \mathcal{V}_{os_2}) \subseteq \mathcal{V}_{t_{(i,j)}} \cap \mathcal{V}_{os_1}$ . The desired conclusion follows at once.  $\square$

Now observe that  $t_{\uparrow(i,j)} = t_{(i,j)} \cap os_{(i,j)}$  so that  $\mathcal{V}_{t_{\uparrow(i,j)}} = \mathcal{V}_{t_{(i,j)}} \cap \mathcal{V}_{os_{(i,j)}}$  and  $t_{(i,j)} = \bigcup_{o \text{ is an operation or a variable}} (t_{(i,j)} \cap os_o)$  where  $os_o$  is obtained from  $cs_{(i,j)}$  by

<sup>7</sup>These conditions can be slightly relaxed, but we try to present the main idea only

attaching the operation (or variable)  $o$  at the node  $\uparrow(i, j)$ . Therefore we also have  $\mathcal{V}_{t(i, j)} = \bigcup_{o \text{ is an operation or a variable}} (\mathcal{V}_{t(i, j)} \cap \mathcal{V}_{os_o})$ . Since operations can not appear or disappear from a population during crossover,  $\mathcal{V}_{os_o} \neq \emptyset \implies \exists$  an individual in  $P$  fitting the operation schema  $os_o$ . Thus the only sets of the form  $\mathcal{V}_{t(i, j)} \cap \mathcal{V}_{os_o}$  which may possibly contribute to the union above are these for which  $\exists$  an individual in  $P$  fitting the operation schema  $os_o$ . According to the claim above, all such sets contribute exactly the same amount. Moreover, by assumption  $os(i, j)(P) \neq \emptyset$ , and so we have

$$|\mathcal{V}_{t(i, j)}| = n \cdot |\mathcal{V}_{t(i, j)} \cap \mathcal{V}_{os(i, j)}| = n \cdot |\mathcal{V}_{t(i, j)} \cap os(i, j)| = n \cdot |\mathcal{V}_{t_{\uparrow(i, j)}}| \implies \frac{|\mathcal{V}_{t_{\uparrow(i, j)}}|}{|\mathcal{V}_{t(i, j)}|} = \frac{1}{n}$$

where  $n$  is the number of operation schemata of the form  $os_o$  which are obtained from  $cs(i, j)$  by attaching a variable or an operation at the node with coordinate  $(i, j)$  and for which  $\exists$  an individual in  $P$  fitting the operation schema  $os_o$  and the last implication holds under the condition that  $|\mathcal{V}_{t(i, j)}| \neq 0$ . This condition is, indeed satisfied. (Suppose not. Let  $(a, b)$  denote the smallest coordinate such that  $|\mathcal{V}_{t(a, b)}| = 0$ . Notice that  $(a, b) \neq (1, 1)$  since  $|\mathcal{V}_{t(1, 1)}| \neq 0$ . (By assumption  $\exists$  an individual, say  $\mathbf{x}$ , in  $P$  fitting the hyperschema  $t_{(1, 1)}$ . Even if  $\mathbf{x}$  is not the 1<sup>st</sup> individual of  $P$ , by performing crossover of  $\mathbf{x}$  with the 1<sup>st</sup> individual of  $P$  at the root node one gets a population  $Q \in \mathcal{V}_{t(1, 1)}$ .) But then  $(a, b) = \uparrow(i, j)$  for some coordinate  $(i, j)$  and according to the equation above we have  $|\mathcal{V}_{t(i, j)}| = n \cdot |\mathcal{V}_{t_{\uparrow(i, j)}}| = n \cdot |\mathcal{V}_{t(a, b)}| = 0$  which contradicts the minimality of the coordinate  $(a, b)$ . So we conclude that  $|\mathcal{V}_{t(i, j)}| \neq 0$ ) Thereby

we have  $\frac{|\mathcal{V}_{t_{\uparrow(i, j)}}|}{|\mathcal{V}_{t(i, j)}|} = \frac{1}{n}$ . But  $cs(i, j) = \bigcup_{o \text{ is an operation or a variable}} os_o \implies cs(i, j)(P) = \bigcup_{o \text{ is an operation or a variable}} os_o(P)$ . Since we assumed that all of the operations and variables are distinct,  $\exists$  at most one individual in  $P$  fitting the operation schema  $os_o$  and it now follows that  $|cs(i, j)(P)| =$  the number of operation schemata of the form  $os_o$  such that  $os_o(P) \neq \emptyset$  which is precisely the number  $n$ . We finally obtain  $\frac{|\mathcal{V}_{t_{\uparrow(i, j)}}|}{|\mathcal{V}_{t(i, j)}|} = \frac{1}{|cs(i, j)|}$  which is precisely the conclusion of the lemma.  $\square$

**Remark 4.8.2** Given an individual  $\mathbf{u}$  and a population  $P$  consisting of  $m$  individuals, observe that the number of individuals fitting the hyperschema  $t_{(1, 1)}$  is the same in every population from  $[P]_{\mathcal{A}}$ , i. e.  $\forall Q \in [P]_{\mathcal{A}}$  we have  $|t_{(1, 1)}(Q)| = |t_{(1, 1)}(P)| = 1$ . It follows immediately now that  $\frac{1}{m} E_{\mathcal{Q}[P]_{\mathcal{A}}}(\mathcal{X}(\square, S_{t_{(1, 1)}})) = \frac{1}{m}$ .

We now combine corollary 4.8.3, remark 4.8.2 and lemma 4.8.4 to obtain the following special case of Geiringer theorem for nonlinear GP under homologous crossover in case when all of the operations appearing in the individuals of the initial population  $P$  are distinct:

$$\begin{aligned} \lim_{t \rightarrow \infty} \Phi(\mathbf{u}, P, t) &= \frac{1}{m} \cdot \prod_{(i, j) \text{ is not the maximal coordinate of } \mathbf{u}} \frac{1}{|cs(i, j)(P)|} = \\ &= \prod_{v \text{ is a node of } \mathbf{u}} \frac{1}{|cs(v)(P)|} \end{aligned}$$

(recall that when  $v$  is the root node of  $\mathbf{u}$ ,  $cs(v)$  determines the entire search space, and so  $\frac{1}{|cs(v)(P)|} = \frac{1}{m}$ ) To obtain the general case, suppose we are given an initial population  $P$ . For every node  $v$  of  $\mathbf{u}$  consider the set of operations  $\mathcal{O}(v) = \{o \mid |os_o(P)| \geq$

1 where  $os_o$  is obtained from  $cs(v)$  as in definition 4.7.1}. For every Moreover, for every operation (or variable)  $o \in \mathcal{O}(v)$  let  $\mathbf{x}_1^o, \mathbf{x}_2^o, \dots, \mathbf{x}_{|os_o(P)|}^o$  denote an enumeration of the individuals in  $P$  fitting the operation schema  $os_o(P)$ . Relabel the operation  $o$  occurring in the node  $v$  of  $\mathbf{x}_i^o$  by the formally different operation  $(o, i)$  (i. e. by the ordered pair  $(o, i)$  whose first element is the operation  $o$  itself and the second element is the index telling us in which individual of  $P$  the operation  $o$  labels the node  $v$ ). After all of the relabelling is complete we obtain a new population  $P'$  which is special with respect to the individual  $\mathbf{u}$  in the sense of definition 4.8.6. Formally speaking, we expand our signature  $\Sigma = (\Sigma_1, \Sigma_2, \dots, \Sigma_N)$  as in definition 4.6.1 by adding the operations (variables)  $(o, i)$  into  $\Sigma_j$  where  $j$  is the arity of the operation  $o$ . This gives us a new signature  $\Sigma^* = (\Sigma_1^*, \Sigma_2^*, \dots, \Sigma_N^*)$  where

$$\Sigma_j^* = \{o \mid o \in \Sigma_j \text{ and } o \notin \bigcup_{v \text{ is a node of } \mathbf{u}} \mathcal{O}(v)\} \cup \{(o, i) \mid o \in \mathcal{O}(v) \text{ for some } v \text{ and } 1 \leq i \leq |os_o(P)|\}.$$

Denote by  $\Omega^*$  the search space induced by the signature  $\Sigma^*$ . The natural projection maps  $p_j : \Sigma_j^* \rightarrow \Sigma_j$  sending  $0 \rightarrow o$  when  $o \notin \bigcup_{v \text{ is a node of } \mathbf{u}} \mathcal{O}(v)$  and  $(o, i) \rightarrow o$  when  $o \in \mathcal{O}(v)$  for some node  $v$  of  $\mathbf{u}$ , induce the natural “deletion of the extra labels” projection of the search spaces  $\varphi : \Omega^* \rightarrow \Omega$  where the individual  $\varphi(\mathbf{w}) \in \Omega$  is obtained from the individual  $\mathbf{w} \in \Omega^*$  by replacing the label of every node  $w$  of  $\mathbf{w}$  with  $p_j(w)$  where  $j$  is the arity of the node  $w$ . It is easily seen that the natural projection  $\varphi$  commutes with the crossover transformations in the sense that for any individuals  $\mathbf{x}, \mathbf{y} \in \Omega^*$  and for any crossover transformation  $T \in \mathcal{G}$  (see definition 4.6.9) we have  $\varphi(T(\mathbf{x}, \mathbf{y})) = T(\varphi(\mathbf{x}), \varphi(\mathbf{y}))$ .<sup>8</sup> Notice also that the population  $P$  can be obtained from the population  $P'$  by applying the natural projection  $\varphi$  to every individual of  $P'$ . Therefore, running the algorithm with the initial population  $P$  is the same thing as running the algorithm with the initial population  $P'$  and reading the output by applying the natural projection  $\varphi$ . The special case does apply to the population  $P'$ , as mentioned above, and so we have

$$\lim_{t \rightarrow \infty} \Phi(\mathbf{u}, P, t) = \sum_{\mathbf{w} \in \varphi^{-1}(\mathbf{u})} \lim_{t \rightarrow \infty} \Phi(\mathbf{w}, P, t) = \sum_{\mathbf{w} \in \varphi^{-1}(\mathbf{u})} \prod_{v \text{ is a node of } \mathbf{w}} \frac{1}{|cs(v)(P)|}.$$

Notice that  $\mathbf{w} \in \varphi^{-1}(\mathbf{u})$  precisely when the underlying shape schema of  $\mathbf{w}$  is the same as that of  $\mathbf{u}$ , call this shape schema  $t_{\mathbf{u}}$ , and the label of every node  $v$  of  $\mathbf{w}$  is  $(o, i)$  where  $o$  is the label of the node  $v$  of  $\mathbf{u}$ . According to the way the population  $P'$  was introduced, there are precisely  $|os(v)(P)|$  such labels (see also definition 4.7.1). We can then identify the preimage  $\varphi^{-1}(\mathbf{u})$  with the set  $\prod_{j=1}^K \{i \mid 1 \leq i \leq |os(v_j)|\}$  of ordered  $K$ -tuples of integers where  $K$  is the number of nodes in the parse tree of  $\mathbf{u}$  and  $v_1, v_2, \dots, v_K$  is any fixed enumeration of the nodes of  $\mathbf{u}$ , in the following manner: The identification map  $\iota : \prod_{j=1}^K \{i \mid 1 \leq i \leq |os(v_j)(P)|\} \rightarrow \varphi^{-1}(\mathbf{u})$  sends a given ordered  $K$ -tuple  $(i_1, i_2, \dots, i_K)$  into the tree  $\mathbf{w} = \iota((i_1, i_2, \dots, i_K))$  whose

<sup>8</sup>Of course, formally speaking, the two transformations  $T$  involved in the equation above are distinct, as they have different domains ( $\Omega^*$  and  $\Omega$  respectively), but they are determined by the same set of shape schemata and the same choice of nodes for crossover so we denote them by the same symbol.

underlying shape schema is  $t_{\mathbf{u}}$  and the label of a node  $v_j$  of  $\mathbf{w}$  is  $(o_j, i_j)$  where  $o_j$  is the label of the node  $v_j$  in the parse tree of  $\mathbf{u}$ . We finally obtain:

$$\begin{aligned}
\lim_{t \rightarrow \infty} \Phi(\mathbf{u}, P, t) &= \sum_{\mathbf{w} \in \varphi^{-1}(\mathbf{u})} \prod_{v \text{ is a node of } \mathbf{w}} \frac{1}{|cs(v)(P)|} = \\
&= \sum_{(i_1, i_2, \dots, i_K) \in \prod_{j=1}^K \{i \mid 1 \leq i \leq |os(v_j)|\}} \prod_{v \text{ is a node of } \mathbf{u}} \frac{1}{|cs(v)(P)|} = \\
&= \sum_{i_1=1}^{|os(v_1)(P)|} \sum_{i_2=1}^{|os(v_2)(P)|} \dots \sum_{i_K=1}^{|os(v_K)(P)|} \prod_{v \text{ is a node of } \mathbf{u}} \frac{1}{|cs(v)(P)|} \\
&= \prod_{j=1}^K \sum_{i_j=1}^{|os(v_j)(P)|} \frac{1}{|cs(v_j)(P)|} = \prod_{v \text{ is a node of } \mathbf{u}} \frac{|os(v)(P)|}{|cs(v)(P)|}
\end{aligned}$$

which is precisely the assertion of theorem 4.7.1.

## 4.9 What Does Theorem 4.5.2 Tell Us in the Presence of Mutation for Nonlinear GP?

In general, mutation is an elementary step of type 2 (see definition 4.2.5) which is determined by the reproduction 1-tuple of the form  $(\Omega, \mathcal{M}, p, \wp_m)$  where  $\mathcal{M}$  is a family of functions on  $\Omega$ . Notice that the set of partitions of the set of  $m$  elements into 1-element subsets consists of exactly one element – the partition into the singletons. This forces  $\wp_m$  to be the trivial probability distribution. We shall, therefore, omit it from writing:

**Definition 4.9.1** A mutation 1-tuple is a reproduction 1-tuple  $(\Omega, \mathcal{M}, p)$  where  $\mathcal{M}$  consists of functions on  $\Omega$  and  $1_\Omega \in \mathcal{M}$ . (Here  $1_\Omega : \Omega \rightarrow \Omega$  denotes the identity map.)

An ergodic mutation 1-tuple is a mutation 1-tuple  $(\Omega, \mathcal{M}, p)$  such that  $\forall x$ , and  $y \in \Omega \exists M \in \mathcal{M}$  with  $M(x) = y$  and  $p(M) > 0$ .

The following fact is a rather simple consequence of theorem 4.5.2 (see corollaries 7.1 and 7.2 of [5]):

**Corollary 4.9.1** Let  $\mathcal{A}$  denote a bijective and self-transient algorithm which involves an elementary step determined by an ergodic mutation. Then the Markov chain associated to the algorithm  $\mathcal{A}$  is irreducible and the unique stationary distribution of this Markov chain is uniform. In particular, the limiting frequency of occurrence of any given individual  $\mathbf{x}$  is  $\lim_{t \rightarrow \infty} \Phi(\{\mathbf{x}\}, P, t) = \frac{1}{|\Omega|}$  (see definition 4.5.6 for the meaning of  $\Phi(\{\mathbf{x}\}, P, t)$ ).

When dealing with nonlinear GP, depending on the circumstances, one may want to consider different types of mutation. Below we define one such possible mutation:



**Definition 4.9.2** Let  $\Omega$  denote the search space for nonlinear GP over the signature  $\Sigma = (\Sigma_0, \Sigma_1, \Sigma_2, \dots, \Sigma_N)$  where  $\Sigma_i$ 's are finite sets (see definition 4.6.1). Consider a configuration schema  $t$  and a node  $v$  of  $t$ . Let  $i$  denote the arity of the node  $v$  and let  $\pi$  denote a permutation of  $\Sigma_i$ . We define a node mutation transformation  $M_{t,v,\pi} : \Omega \rightarrow \Omega$  to be the function which sends a given program tree  $\mathbf{u}$  which fits the schema  $t$  to the program  $M_{t,v,\pi}(\mathbf{u})$  obtained from  $\mathbf{u}$  by replacing the label  $a \in \Sigma_i$  of the node  $v$  of  $\mathbf{u}$  with  $\pi(a)$  whenever  $\mathbf{u}$  fits the configuration schema  $t$  (if  $\mathbf{u}$  does not fit the schema  $t$  then  $M_{t,v,\pi}(\mathbf{u}) = \mathbf{u}$ ). We define the family of node mutations

$$\mathcal{M}_{node} = \{M_{t,v,\pi} : \Omega \rightarrow \Omega \mid \pi \in \mathcal{S}_{\Sigma_i} \text{ where } t \text{ is a configuration schema and } i \text{ is the arity of the node } v \text{ of } t\}.$$

As usual  $\mathcal{S}_X$  denotes the set of all permutations of the set  $X$ . Denote by  $\Omega_{NodeMut} = (\Omega, \mathcal{M}_{node}, p)$  the corresponding mutation 1-tuple.

Although node mutation described in definition 4.9.2 above is not ergodic in the sense of definition 4.9.1, it defines a bijective elementary step (see definition 4.5.3). Indeed, it is easy to see that the transformation  $M_{t,v,\pi^{-1}}$  is a 2-sided inverse of the transformation  $M_{t,v,\pi}$ . Thereby theorem 4.5.2 applies to nonlinear GP with homologous crossover and node mutation. It is also possible to derive a formula for the limiting frequency of occurrence of a given individual  $\mathbf{u}$ , namely  $\lim_{t \rightarrow \infty} \Phi(\mathbf{u}, P, t)$  much in the same way as in theorem 4.7.1. In order to state the corresponding result for nonlinear GP with homologous crossover and node mutation, it is convenient to introduce the following definitions first:

**Definition 4.9.3** Fix an individual  $\mathbf{u} \in \Omega$ . Let  $v$  denote a node of  $\mathbf{u}$  and consider the configuration schema  $cs(\uparrow v, i)$  obtained from the configuration schema  $cs(v)$  by attaching a node of degree  $i$  together with its  $i$  children in place of the  $\#$  sign at the node corresponding to  $v$  of  $cs(v)$ . The newly attached nodes are then labelled by the  $=$  and  $\#$  signs respectively. If the newly attached node is of arity 0 then it is a leaf node labelled by the  $=$  sign<sup>9</sup>. Furthermore, write  $cs(\uparrow v, \mathbf{u})$  in place of  $cs(\uparrow v, i)$  when  $i$  is the arity of the node  $v$  in  $\mathbf{u}$ . Also denote by  $os(v, o)$  the operation schema obtained from the configuration schema  $cs(v)$  by attaching a node labelled by the operation  $o$  together with its appropriate number of children in place of the  $\#$  sign. The children of the newly attached node (if there are any) are labelled by the  $\#$  signs.

**Definition 4.9.4** Given a mutation 1-tuple  $(\Omega, \mathcal{M}, p)$ , a configuration schema  $t$  (see definition 4.7.1), and a node  $v$  of  $t$  having arity  $i$ , denote by  $G(t, v)$  the group generated by all the permutations  $\pi \in \Sigma_i$  such that  $p(M_{s,v,\pi}) > 0$  for some configuration schema  $s$  such that the common region of  $s$  and  $t$  contains the node  $v$ . Fix an operation  $a \in \Sigma_i$ . Let

$$\mathcal{O}(t, v, a) = \{o \in \Sigma_i \mid \exists g \in G(t, v) \text{ with } g \cdot a = o\}$$

denote the orbit of the operation  $a$  under the action of the group  $G(t, v)$ .

---

<sup>9</sup>Formally speaking, according to definition 4.7.1,  $cs(\uparrow v, i)$  is not always a Poli hyperschema since it may contain a leaf node labelled by the  $=$  sign. However, such a schema also defines a subset of the search space  $\Omega$  in much the same way as Poli hyperschemata. The only difference is that a leaf node labelled by the  $=$  sign can be replaced by a variable only. One can not attach a nontrivial program tree to it.



Suppose we are given a population  $P$  of size  $m$  consisting of program trees from  $\Omega$ . Recall from definition 4.7.2 that we denote by  $|H(P)|$  the number of individuals in the population  $P$  fitting the schema  $H$ .

**Theorem 4.9.2** *Let  $\mathcal{A}$  denote an algorithm determined by 2 elementary steps of type 2 one of which is determined by the node mutation (see definition 4.9.2) and the other one by a homologous GP crossover. Suppose every one of the transformations in the family  $\mathcal{G}$  of GP homologous crossovers has a positive probability of being chosen.<sup>10</sup> Fix an individual (a program tree)  $\mathbf{u} \in \Omega$  and an initial population  $P$ . Let  $o(\mathbf{u}, v)$  denote the operation labelling the node  $v$  of the program tree  $\mathbf{u}$ . Denote by  $\hat{\mathbf{u}}$  the configuration schema obtained from the shape schema,  $\tilde{\mathbf{u}}$ , of  $\mathbf{u}$  (see definition 4.6.5) by labelling all the interior nodes of  $\mathbf{u}$  with the  $=$  signs and all the leaf nodes with the  $\#$  signs. Suppose that the probability distribution on the collection of node mutations is such that whenever  $v$  is a node of  $\mathbf{u}$  we have  $p(M_{s,v,\pi}) > 0 \implies s = cs(v)$  where as before  $cs(v)$  is the configuration schema of  $\mathbf{u}$  corresponding to the node  $v$ .<sup>11</sup> Then we have*

$$\lim_{t \rightarrow \infty} \Phi(\mathbf{u}, P, t) = \prod_{v \text{ is a node of } \mathbf{u}} \frac{\sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v, o(\mathbf{u}, v))} |os(v, o)(P)|}{\sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v, o)|}.$$

*Proof:* The proof of theorem 4.9.2 is very similar to the proof of theorem 4.7.1 and contains essentially no new ideas. We leave most of the details for the interested reader as an exercise and provide only the rough outline: Just like theorem 4.7.1, theorem 4.9.4 follows from corollary 4.8.3 by considering the nested sequence of hyperschemata  $t_{(1,1)} \geq t_{(2,1)} \geq \dots \geq t_{(i,j)} \geq t_{\uparrow(i,j)} \geq \dots \geq t_{(l,k)} = \mathbf{u}$  corresponding to the program  $\mathbf{u}$  (see definition 4.8.4). First we consider a special case when every set  $\Sigma_i$  consists of ordered pairs  $(l, o)$  where  $l$  is an integer, and mutation is allowed to modify only the operation  $o$  and is not allowed to change the integer  $l$ . We then prove theorem 4.9.2 in the special case when all the labels contained in the initial population  $P$  have distinct first coordinates. The general case then follows by introducing the extra integer labels for the first coordinate, applying the special case and then “erasing the integer part from the labels” in exactly the same way as it was done in the proof of theorem 4.7.1. The main difference lies in the claim proved inside lemma 4.8.4. The corresponding claim for the proof of theorem 4.9.2 then says the following:

**Lemma 4.9.3** *Fix a node  $v$  with coordinate  $(i, j)$ . Fix any two operation schemata  $os(a)$  and  $os(b)$  which are obtained from  $cs_{(i,j)}$  by attaching either a variable or an operation at the node with coordinate  $(i, j)$ . Suppose  $\exists$  individuals in  $P$  fitting both,  $os(c)$  and  $os(d)$  where  $a \in \mathcal{O}(cs(i, j), v, c)$  and  $b \in \mathcal{O}(cs(i, j), v, d)$ . Then  $|\mathcal{V}_{t_{(i,j)}} \cap \mathcal{V}_{os(a)}| = |\mathcal{V}_{t_{(i,j)}} \cap \mathcal{V}_{os(b)}|$ .*

Just like the claim inside of the lemma 4.8.4, lemma 4.9.3 is proved by constructing an explicit bijection between the sets  $\mathcal{V}_{t_{(i,j)}} \cap \mathcal{V}_{os(a)}$  and  $\mathcal{V}_{t_{(i,j)}} \cap \mathcal{V}_{os(b)}$ . The only difference is that these bijections make use of mutations as well as crossover.  $\square$

<sup>10</sup>Again we remark that this condition can be slightly relaxed but it does not introduce any new ideas of interest.

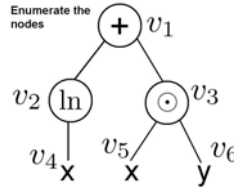
<sup>11</sup>Notice that this implies that  $\mathcal{O}(\hat{\mathbf{u}}, v, o(\mathbf{u}, v)) = \mathcal{O}(cs(v), v, o(\mathbf{u}, v))$

It may be worth mentioning that theorem 4.7.1 is a special case of theorem 4.9.2. Indeed, if the only mutation transformations chosen with positive probability are these which assign positive probability only to the mutations defined by the identity permutations, then every orbit  $\mathcal{O}(\hat{\mathbf{u}}, v, o)$  consists of exactly one element so that  $\forall t$  and  $v$  we have  $|\mathcal{O}(\hat{\mathbf{u}}, v, o)| = 1$ . To compress the language, we shall use  $\uplus$  to denote the union of *disjoint* sets. We then have  $\uplus_{o \in \mathcal{O}(\hat{\mathbf{u}}, v, o(\mathbf{u}, v))} os(v, o)(P) = os(v)(P)$  since  $o(\mathbf{u}, v)$  is the only operation inside of  $\mathcal{O}(\hat{\mathbf{u}}, v, o(\mathbf{u}, v))$  so that  $os(v)(P)$  is the only contributor to the disjoint union above. Moreover, we also have  $cs(v) = \uplus_{i=0}^N \uplus_{o \in \Sigma_i} os(v, o)$  so that we obtain

$$\sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v, o)| = \sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v, o)(P)| \cdot 1 = |cs(v)(P)|.$$

The formula in theorem 4.9.2 now simplifies to the one in theorem 4.7.1.

**Example 4.9.1** Continuing with example 4.7.2, suppose the signature  $\Sigma$  is defined as follows:  $\Sigma = (\Sigma_0, \Sigma_1, \Sigma_2)$  where  $\Sigma_0 = \{x, y, z, w\}$ ,  $\Sigma_1 = \{\sin, \cos, \tan, \cot, \ln\}$  and  $\Sigma_2 = \{+, -, *, \odot\}$  where  $\odot$  is a binary operation symbol different from  $+$ ,  $*$  and  $-$ . (Of course, the semantics of the binary operation  $\odot$  is irrelevant to the content of this example, but if the reader feels more comfortable with a concrete interpretation, they may assume, for instance, that  $x \odot y = \int_x^y e^{\xi^2} d\xi$ .) Now suppose the individual  $\mathbf{u}$  is the program  $(+(\ln(x), \odot(x, y)))$  pictured below (with nodes being enumerated just like in example 4.7.1) (it has the same shape schema as the individual in that example):



Notice that this individual has exactly the same set of configuration schemata as the corresponding set of configuration schemata in example 4.7.1 (the reader may see these configuration schemata pictured in that example). Suppose that the following mutation transformations are the only ones chosen with positive probability:

$$\begin{aligned} &M_{cs(v_1), v_1, (+, -)}, M_{cs(v_1), v_1, (*, \odot)}, M_{cs(v_2), v_2, (\ln, \sin, \cos)(\tan, \cot)}, \\ &M_{cs(v_2), v_2, (\tan, \cot)(\sin, \cos)}, M_{cs(v_2), v_2, (\ln, \sin)}, M_{cs(v_3), v_3, (+, \odot, *)}, \\ &M_{cs(v_3), v_3, (-, *)}, M_{cs(v_4), v_4, (x, w)(y, z)}, M_{cs(v_4), v_4, (+, *)(-, \odot)}, M_{cs(v_5), v_5, (x, y, z)}, \\ &M_{cs(v_5), v_5, (+, -, *, \odot)}, M_{cs(v_6), v_6, (x, w)(y, z)}, \text{ and} \\ &M_{cs(v_6), v_6, (x, w)}, M_{cs(v_6), v_6, (\sin, \cos, \tan, \cot)}, M_{cs(v_6), v_6, (+, -)}. \end{aligned}$$

Here we represent permutations in terms of their “disjoint cycle decompositions”: for example,  $(\ln, \sin, \cos)(\tan, \cot)$  represents the permutation on  $\Sigma_1$  which sends  $\ln$  into  $\sin$ ,  $\sin$  into  $\cos$  and  $\cos$  back into  $\ln$ . Likewise, it sends  $\tan$  into  $\cot$  and  $\cot$  back into

tan. If a cycle has length one (i.e. the element appearing in the cycle is a fixed point of the corresponding permutation) we omit that cycle from writing. For example,  $\odot$  and  $*$  are the fixed points of the permutation  $(+, -)$ . The corresponding permutation groups are:

$$\begin{aligned} G(cs(v_1), v_1) &= \langle (+, -), (*, \odot) \rangle, \\ G(cs(v_2), v_2) &= \langle (\ln, \sin, \cos)(\tan, \cot), (\tan, \cot)(\sin, \cos), (\ln, \sin) \rangle, \\ G(cs(v_3), v_3) &= \langle (+, \odot, *), (-, *) \rangle, \\ G(cs(v_4), v_4) &= \langle (x, w)(y, z), (+, *)(-, \odot) \rangle, \\ G(cs(v_5), v_5) &= \langle (x, y, z), (+, -, *, \odot) \rangle, \text{ and} \\ G(cs(v_6), v_6) &= \langle (x, w)(y, z), (x, w)(\sin, \cos, \tan, \cot), (+, -) \rangle. \end{aligned}$$

The cycle decomposition makes it easy to compute the corresponding orbits:

$$\begin{aligned} \mathcal{O}(cs(v_1), v_1, +) &= \{+, -\}, \mathcal{O}(cs(v_2), v_2, \ln) = \{\ln, \sin, \cos\}, \\ \mathcal{O}(cs(v_3), v_3, \odot) &= \{+, -, *, \odot\}, \mathcal{O}(cs(v_4), v_4, x) = \{x, w\}, \\ \mathcal{O}(cs(v_5), v_5, x) &= \{x, y, z\} \text{ and } \mathcal{O}(cs(v_6), v_6, y) = \{y, z\}. \end{aligned}$$

Now suppose the initial population is the same as in example 4.7.2. In order to apply theorem 4.9.2, for every node  $v$  of  $\mathbf{u}$  we need to compute the number  $|\text{os}(v, o)(P)|$ . Recall that the schema  $\text{os}(v, o)$  is obtained from the schema  $\text{os}(v)$  by attaching the operation  $o$  at the node  $v$  and labelling its children nodes by the  $\#$  signs. For the population  $P$  in example 4.7.2 it was already computed that  $|\text{os}(v_1, +)| = |\text{os}(v_1)| = 2$ . There are exactly 2 individuals (namely  $\mathbf{x}_3$  and  $\mathbf{x}_4$ ) fitting the schema  $\text{os}(v_1, *)$  and so  $|\text{os}(v_1, *)| = 2$ . Exactly one individual, namely  $\mathbf{x}_2$ , and one individual, namely  $\mathbf{x}_6$ , fit the schemata  $\text{os}(v_1, \sin)$  and  $\text{os}(v_2, \cos)$  respectively and so  $|\text{os}(v_1, \sin)| = |\text{os}(v_1, \cos)| = 1$ . For all the other operations  $o \in (\Sigma_0 \cup \Sigma_1 \cup \Sigma_2) - \{+, *, \sin, \cos\}$  there are no individuals in  $P$  fitting the schema  $\text{os}(v_1, o)$  and so we have  $|\text{os}(v_1, o)| = 0$ . There are no individuals in  $P$  fitting the schema  $\text{os}(v_2, \ln) = \text{os}(v_2)$  for the individual  $\mathbf{u}$  of the current example, and no individuals fitting the schemata of the form  $\text{os}(v_2, o)$  where  $o \notin \{*, \sin, \cos\}$  so that for such  $o$  we have  $|\text{os}(v_2, o)(P)| = 0$ . Moreover, there is exactly one individual, namely  $\mathbf{x}_1$  fitting the schema  $\text{os}(v_2, *)$  and exactly one, namely  $\mathbf{x}_4$  fitting the schema  $\text{os}(v_2, \cos)$  so that  $|\text{os}(v_2, *)| = |\text{os}(v_2, \cos)(P)| = 1$ ; exactly two individuals, namely  $\mathbf{x}_3$  and  $\mathbf{x}_5$  fit the schema  $\text{os}(v_2, \sin)$  so that  $|\text{os}(v_2, \sin)(P)| = 2$ . Continuing in this manner with the rest of the nodes of  $\mathbf{u}$  we obtain  $|\text{os}(v_3, o)(P)| = 0$  for  $o \notin \{+, *\}$ ;  $\mathbf{x}_1$  and  $\mathbf{x}_3$  fit  $\text{os}(v_3, +)$  while  $\mathbf{x}_4$  and  $\mathbf{x}_5$  fit  $\text{os}(v_3, *)$  and so  $|\text{os}(v_3, +)(P)| = |\text{os}(v_3, *)| = 2$ .  $|\text{os}(v_4, o)(P)| = 0$  for  $o \notin \{x, y, +\}$ ;  $\mathbf{x}_3$  is the only individual fitting the schema  $\text{os}(v_4, y)$ ,  $\mathbf{x}_4$  is the only individual fitting the schema  $\text{os}(v_4, x)$  and  $\mathbf{x}_5$  is the only individual fitting the schema  $\text{os}(v_4, +)$  and so we have  $|\text{os}(v_4, x)(P)| = |\text{os}(v_4, y)(P)| = |\text{os}(v_4, +)(P)| = 1$ .  $|\text{os}(v_5, o)(P)| = 0$  for  $o \notin \{*, x, y\}$ ;  $\mathbf{x}_1$  is the only individual fitting the schema  $\text{os}(v_5, *)$  and  $\mathbf{x}_5$  is the only individual fitting the schema  $\text{os}(v_5, y)$  while the individuals  $\mathbf{x}_3$  and  $\mathbf{x}_4$  are the only two which fit the schema  $\text{os}(v_5, x)$  so that we have  $|\text{os}(v_5, *)| = |\text{os}(v_5, y)(P)| = 1$  and  $|\text{os}(v_5, x)(P)| = 2$ .  $|\text{os}(v_6, o)(P)| =$

0 for  $o \notin \{\sin, x, y\}$ ; moreover,  $\mathbf{x}_1$  is the only individual fitting the schema  $os(v_6, \sin)$  and  $\mathbf{x}_4$  is the only individual fitting the schema  $os(v_6, y)$  while the individuals  $\mathbf{x}_3$  and  $\mathbf{x}_5$  are the only two which fit the schema  $os(v_6, x)$  so that we have  $|os(v_6, \sin)(P)| = |os(v_6, y)(P)| = 1$  and  $|os(v_6, x)(P)| = 2$ .

Finally for every node  $v$  and for every operation  $o \in \Sigma_0 \cup \Sigma_1 \cup \Sigma_2$  such that  $|os(v, o)(P)| \neq 0$  we need to compute  $|\mathcal{O}(\hat{\mathbf{u}}, v, o)|$ . For the node  $v_1$  these are  $\mathcal{O}(\hat{\mathbf{u}}, v_1, +) = \mathcal{O}(cs(v_1), v_1, +) = \{+, -\}$  so that  $|\mathcal{O}(cs(v_1), v_1, +)| = 2$ , and, likewise, from the description of the groups  $G(cs(v_i), v_i)$  given above, it is easy to compute that  $\mathcal{O}(cs(v_1), v_1, *) = \{*, \odot\}$ ,  $\mathcal{O}(cs(v_1), v_1, \sin) = \{\sin\}$  and  $\mathcal{O}(cs(v_1), v_1, \cos) = \{\cos\}$  so that

$$|\mathcal{O}(cs(v_1), v_1, \sin)| = |\mathcal{O}(cs(v_1), v_1, \cos)| = 1.$$

$\mathcal{O}(cs(v_2), v_2, *) = \{*\}$ ,  $\mathcal{O}(cs(v_2), v_2, \sin) = \mathcal{O}(cs(v_2), v_2, \cos) = \{\ln, \sin, \cos\}$  and so

$$|\mathcal{O}(cs(v_2), v_2, \sin)| = |\mathcal{O}(cs(v_2), v_2, \cos)| = 3;$$

$\mathcal{O}(cs(v_3), v_3, *) = \mathcal{O}(cs(v_3), v_3, +) = \Sigma_2$  so that

$$|\mathcal{O}(cs(v_3), v_3, *)| = |\mathcal{O}(cs(v_3), v_3, +)| = 4;$$

$\mathcal{O}(cs(v_4), v_4, x) = \{x, w\}$  and  $\mathcal{O}(cs(v_4), v_4, y) = \{y, z\}$  so that

$$|\mathcal{O}(cs(v_4), v_4, x)| = |\mathcal{O}(cs(v_4), v_4, y)| = 2;$$

$\mathcal{O}(cs(v_4), v_4, +) = \{+, *\}$  so that  $|\mathcal{O}(cs(v_4), v_4, +)| = 2$ ;  $\mathcal{O}(cs(v_5), v_5, *) = \Sigma_2$ ,  $\mathcal{O}(cs(v_5), v_5, x) = \mathcal{O}(cs(v_5), v_5, y) = \{x, y, z\}$  and so

$$|\mathcal{O}(cs(v_5), v_5, x)| = |\mathcal{O}(cs(v_5), v_5, y)| = 3;$$

$\mathcal{O}(cs(v_6), v_6, \sin) = \{\sin, \cos, \tan, \cot\}$ ,  $\mathcal{O}(cs(v_6), v_6, x) = \{x, w\}$  and, finally,  $\mathcal{O}(cs(v_6), v_6, y) = \{y, z\}$  so that

$$|\mathcal{O}(cs(v_6), v_6, x)| = |\mathcal{O}(cs(v_6), v_6, y)| = 2.$$

Now we are ready to compute the ratios of the form  $\frac{\sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v_j, o(\mathbf{u}, v_j))} |os(v_j, o)(P)|}{\sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v_j, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_j, o)|}$ . From the data computed above we have

$$\sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v_1, o(\mathbf{u}, v_1))} |os(v_1, o)(P)| = |os(v_1, +)(P)| + |os(v_1, -)(P)| = 2 + 0 = 2$$

( $\mathbf{x}_1$  and  $\mathbf{x}_5$  are the only two individuals in  $P$  which fit the schema  $os(v_1, +)$  while no individual in  $P$  fits  $os(v_1, -)$ );

$$\begin{aligned} \sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v_1, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_1, o)| &= |os(v_1, +)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_1, +)(P)| + \\ &+ |os(v_1, \sin)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_1, \sin)(P)| + |os(v_1, *) (P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_1, *) (P)| + \end{aligned}$$

$$+|os(v_1, \cos)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_1, \cos)(P)| = 2 \cdot 2 + 1 \cdot 1 + 2 \cdot 2 + 1 \cdot 1 = 10$$

and so

$$\frac{\sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v_1, o(\mathbf{u}, v_1))} |os(v_1, o)(P)|}{\sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v_1, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_1, o)|} = \frac{2}{10} = \frac{1}{5}.$$

Continuing in this manner, we obtain

$$\begin{aligned} \sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v_2, o(\mathbf{u}, v_2))} |os(v_2, o)(P)| &= |os(v_2, \ln)(P)| + |os(v_2, \sin)(P)| + \\ &+ |os(v_2, \cos)(P)| = 0 + 2 + 1 = 3 \end{aligned}$$

and

$$\begin{aligned} \sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v_2, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_2, o)| &= |os(v_2, *) (P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_2, *) (P)| + \\ &+ |os(v_2, \sin)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_2, \sin)(P)| + |os(v_2, \cos)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_2, \cos)(P)| = \\ &= 1 \cdot 1 + 2 \cdot 3 + 1 \cdot 3 = 10 \end{aligned}$$

so that

$$\begin{aligned} \frac{\sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v_2, o(\mathbf{u}, v_2))} |os(v_2, o)(P)|}{\sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v_2, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_2, o)|} &= \frac{3}{10}. \\ \sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v_3, o(\mathbf{u}, v_3))} |os(v_3, o)(P)| &= |os(v_3, \odot)(P)| + |os(v_3, +)(P)| + \\ &+ |os(v_3, -)(P)| + |os(v_3, *) (P)| = 0 + 2 + 0 + 2 = 4 \end{aligned}$$

and

$$\begin{aligned} \sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v_3, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_3, o)| &= |os(v_3, +)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_3, +)(P)| + \\ &+ |os(v_3, *) (P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_3, *) (P)| = 2 \cdot 4 + 2 \cdot 4 = 16 \end{aligned}$$

and so

$$\frac{\sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v_3, o(\mathbf{u}, v_3))} |os(v_3, o)(P)|}{\sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v_3, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_3, o)|} = \frac{4}{16} = \frac{1}{4}.$$

$$\sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v_4, o(\mathbf{u}, v_4))} |os(v_4, o)(P)| = |os(v_4, x)(P)| + |os(v_4, w)(P)| = 1 + 0 = 1$$

and

$$\begin{aligned} \sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v_4, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_4, o)| &= |os(v_4, y)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_4, y)| + \\ &+ |os(v_4, x)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_4, x)| + |os(v_4, +)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_4, +)| = 1 \cdot 2 + 1 \cdot 2 + 1 \cdot 2 = 6 \end{aligned}$$

and so

$$\frac{\sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v_4, o(\mathbf{u}, v_4))} |os(v_4, o)(P)|}{\sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v_4, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_4, o)|} = \frac{1}{6}.$$

$$\sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v_5, o(\mathbf{u}, v_5))} |os(v_5, o)(P)| = |os(v_5, x)(P)| + |os(v_5, y)(P)| +$$

$$+ |os(v_5, z)(P)| = 2 + 1 + 0 = 3$$

and

$$\sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v_5, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_5, o)| = |os(v_5, *) (P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_5, *)| +$$

$$+ |os(v_5, x)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_5, x)| + |os(v_5, y)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_5, y)| = 1 \cdot 4 + 2 \cdot 3 + 1 \cdot 3 = 13$$

and so

$$\frac{\sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v_5, o(\mathbf{u}, v_5))} |os(v_5, o)(P)|}{\sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v_5, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_5, o)|} = \frac{3}{13}.$$

Finally,

$$\sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v_6, o(\mathbf{u}, v_6))} |os(v_6, o)(P)| = |os(v_6, y)(P)| + |os(v_6, z)(P)| = 1 + 0 = 1$$

and

$$\sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v_6, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_6, o)| = |os(v_6, \sin)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_6, \sin)| +$$

$$+ |os(v_6, x)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_6, x)| + |os(v_6, y)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_6, y)| = 1 \cdot 4 + 2 \cdot 2 + 1 \cdot 2 = 10$$

so that

$$\frac{\sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v_6, o(\mathbf{u}, v_6))} |os(v_6, o)(P)|}{\sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v_6, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_6, o)|} = \frac{1}{10}.$$

Now we finally compute the product of these ratios and obtain:

$$\lim_{t \rightarrow \infty} \Phi(\mathbf{u}, P, t) = \prod_{i=1}^6 \frac{\sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v_i, o(\mathbf{u}, v_i))} |os(v_i, o)(P)|}{\sum_{i=0}^N \sum_{o \in \Sigma_i} |os(v_i, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v_i, o)|} =$$

$$= \frac{1}{5} \cdot \frac{3}{10} \cdot \frac{1}{4} \cdot \frac{1}{6} \cdot \frac{3}{13} \cdot \frac{1}{10} = \frac{3}{52000}.$$

At the opposite extreme is the case when every mutation transformation in the family  $\mathcal{M}_{\text{node}}$  has a positive probability of being chosen. In this case  $\mathcal{O}(\hat{\mathbf{u}}, v, o(\mathbf{u}, v)) = \Sigma_i$  where  $i$  is the arity of the operation  $o$ . In particular, for every operation  $o$  we have  $o \in \mathcal{O}(\hat{\mathbf{u}}, v, o(\mathbf{u}, v))$  if and only if  $o \in \Sigma_i$ . But then we have

$$\bigsqcup_{o \in \mathcal{O}(\hat{\mathbf{u}}, v, o(\mathbf{u}, v))} os(v, o) = cs(\uparrow v, \mathbf{u})$$

and so

$$\sum_{o \in \mathcal{O}(\hat{\mathbf{u}}, v, o(\mathbf{u}, v))} |os(v, o)(P)| = |cs(\uparrow v, \mathbf{u})(P)|.$$

We also have  $\biguplus_{o \in \Sigma_i} os(v, o) = cs(\uparrow v, i)$  so that

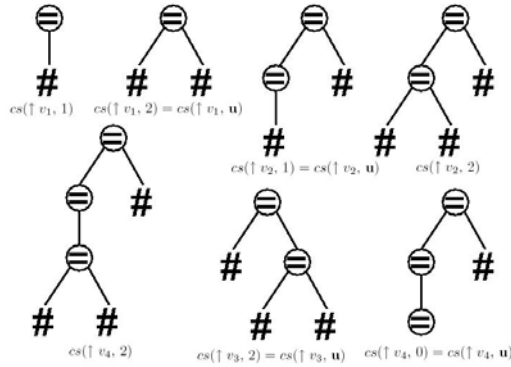
$$\sum_{o \in \Sigma_i} |os(v, o)(P)| \cdot |\mathcal{O}(\hat{\mathbf{u}}, v, o)| = |\Sigma_i| \cdot \sum_{o \in \Sigma_i} |os(v, o)(P)| = |cs(\uparrow v, i)(P)| \cdot |\Sigma_i|.$$

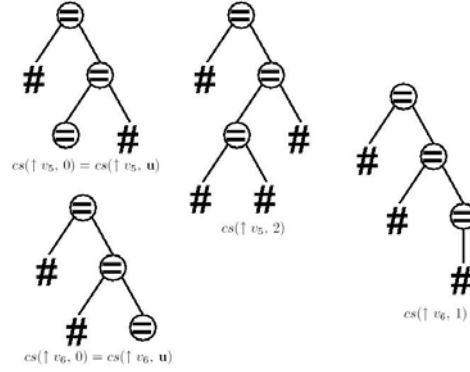
Combining these equations with theorem 4.9.2 we obtain:

**Corollary 4.9.4** *Let  $\mathcal{A}$  denote an algorithm determined by 2 elementary steps of type 2 one of which is determined by the node mutation (see definition 4.9.2) and the other one by a homologous GP crossover. Suppose every one of the transformations in the family  $\mathcal{G}$  of GP homologous crossovers has a positive probability of being chosen. Suppose also that for every node  $v$  of  $\mathbf{u}$  of arity  $i$  and for every permutation  $\pi$  of  $\Sigma_i$  we have  $p(M_{\hat{\mathbf{u}}, v, \pi}) = p(M_{cs(v), v, \pi}) > 0$ . Fix an individual (a program tree)  $\mathbf{u} \in \Omega$  and an initial population  $P$ . Then we have*

$$\lim_{t \rightarrow \infty} \Phi(\mathbf{u}, P, t) = \prod_{v \text{ is a node of } \mathbf{u}} \frac{|cs(\uparrow v, \mathbf{u})(P)|}{\sum_{i=0}^N |cs(\uparrow v, i)(P)| \cdot |\Sigma_i|}.$$

**Example 4.9.2** *Suppose the signature  $\Sigma = (\Sigma_0, \Sigma_1, \Sigma_2)$ , the initial population  $P$  and the individual  $\mathbf{u}$  are exactly as in example 4.9.1. Now suppose, (unlike in example 4.9.1, that for every permutation  $\pi$  of  $\Sigma_i$  we have  $p(M_{\hat{\mathbf{u}}, v, \pi}) > 0$ . Now corollary 4.9.4 applies and we can compute the frequency of occurrence of the individual  $\mathbf{u}$  according to the formula given there. To apply this formula we need to compute the numbers  $|cs(\uparrow v_j, i)(P)|$  where  $1 \leq j \leq 6$  and  $0 \leq i \leq 2$  (the numbers  $|cs(\uparrow v, \mathbf{u})(P)|$  are among these). The configuration schemata  $cs(v_i)$  for the individual  $\mathbf{u}$  are exactly the same as these for the individual of example 4.7.1 (since these two individuals have the same underlying shape schema) and they are pictured in that example. Below we display only these schemata  $|cs(\uparrow v_j, i)(P)|$  for which  $|cs(\uparrow v_j, i)(P)| \neq 0$ . According to definition 4.9.3 they are obtained from the corresponding schemata  $cs(v_i)$  by attaching a node which has  $i$  children (if  $i = 0$  it has no children) in place of the  $\#$  sign at the node  $v_i$  (which means that the  $\#$  sign can be replaced with an arbitrary variable but not with an operation symbol):*





There are exactly 2 individuals, namely  $\mathbf{x}_2$  and  $\mathbf{x}_6$  in  $P$  fitting the schema  $cs(\uparrow v_1, 1)$  so that  $|cs(\uparrow v_1, 1)(P)| = 2$ ;  $\mathbf{x}_1, \mathbf{x}_3, \mathbf{x}_4$  and  $\mathbf{x}_5$  are the only individuals in  $P$  which fit the schema  $cs(\uparrow v_1, 2) = cs(\uparrow v_1, \mathbf{u})$  so that  $|cs(\uparrow v_1, 2)(P)| = |cs(\uparrow v_1, \mathbf{u})(P)| = 4$ ;  $\mathbf{x}_3, \mathbf{x}_4$  and  $\mathbf{x}_5$  are the only individuals in  $P$  which fit the schema  $cs(\uparrow v_2, 1) = cs(\uparrow v_2, \mathbf{u})$  so that  $|cs(\uparrow v_2, 1)(P)| = |cs(\uparrow v_2, \mathbf{u})(P)| = 3$ ;  $\mathbf{x}_1$  is the only individual fitting the schema  $cs(\uparrow v_2, 2)$  so that  $|cs(\uparrow v_2, 2)(P)| = 1$ ;  $\mathbf{x}_1, \mathbf{x}_3, \mathbf{x}_4$  and  $\mathbf{x}_5$  are the only individuals in  $P$  which fit the schema  $cs(\uparrow v_3, 2) = cs(\uparrow v_3, \mathbf{u})$  so that  $|cs(\uparrow v_3, 2)(P)| = |cs(\uparrow v_3, \mathbf{u})(P)| = 4$ ;  $\mathbf{x}_3$  and  $\mathbf{x}_4$  are the only individuals in  $P$  which fit the schema  $cs(\uparrow v_4, 0) = cs(\uparrow v_4, \mathbf{u})$  so that  $|cs(\uparrow v_4, 0)(P)| = |cs(\uparrow v_4, \mathbf{u})(P)| = 2$ ;  $\mathbf{x}_5$  is the only individual fitting the schema  $cs(\uparrow v_4, 2)$  and so  $|cs(\uparrow v_4, 2)(P)| = 1$ ;  $\mathbf{x}_3, \mathbf{x}_4$  and  $\mathbf{x}_5$  are the only individuals in  $P$  fitting the schemata  $cs(\uparrow v_5, 0) = cs(\uparrow v_5, \mathbf{u})$  and/or  $cs(\uparrow v_6, 0) = cs(\uparrow v_6, \mathbf{u})$  and so we have  $|cs(\uparrow v_5, 0)| = |cs(\uparrow v_5, \mathbf{u})| = |cs(\uparrow v_6, 0)| = |cs(\uparrow v_6, \mathbf{u})| = 3$  and  $\mathbf{x}_1$  is the only individual fitting the either and both schemata  $cs(\uparrow v_5, 2)$  and/or  $cs(\uparrow v_6, 1)$  so that  $|cs(\uparrow v_5, 2)| = |cs(\uparrow v_6, 1)| = 1$ . From the definition of the signature  $\Sigma = (\Sigma_0, \Sigma_1, \Sigma_2)$  in example 4.9.1 we see that  $|\Sigma_0| = 4$ ,  $|\Sigma_1| = 5$  and  $|\Sigma_2| = 4$ . We are now ready to compute the ratios:

$$\begin{aligned} \frac{|cs(\uparrow v_1, \mathbf{u})(P)|}{\sum_{i=0}^3 |cs(\uparrow v_1, i)(P)| \cdot |\Sigma_i|} &= \frac{|cs(\uparrow v_1, \mathbf{u})(P)|}{|cs(\uparrow v_1, 2)(P)| \cdot |\Sigma_2| + |cs(\uparrow v_1, 1)(P)| \cdot |\Sigma_1|} = \\ &= \frac{4}{4 \cdot 4 + 2 \cdot 5} = \frac{2}{13}. \end{aligned}$$

$$\begin{aligned} \frac{|cs(\uparrow v_2, \mathbf{u})(P)|}{\sum_{i=0}^3 |cs(\uparrow v_2, i)(P)| \cdot |\Sigma_i|} &= \frac{|cs(\uparrow v_2, \mathbf{u})(P)|}{|cs(\uparrow v_2, 2)(P)| \cdot |\Sigma_2| + |cs(\uparrow v_2, 1)(P)| \cdot |\Sigma_1|} = \\ &= \frac{3}{1 \cdot 4 + 3 \cdot 5} = \frac{3}{19}. \end{aligned}$$

$$\begin{aligned} \frac{|cs(\uparrow v_3, \mathbf{u})(P)|}{\sum_{i=0}^3 |cs(\uparrow v_3, i)(P)| \cdot |\Sigma_i|} &= \frac{|cs(\uparrow v_3, \mathbf{u})(P)|}{|cs(\uparrow v_3, 2)(P)| \cdot |\Sigma_2|} = \\ &= \frac{4}{4 \cdot 4} = \frac{1}{4}. \end{aligned}$$



$$\begin{aligned} \frac{|cs(\uparrow v_4, \mathbf{u})(P)|}{\sum_{i=0}^3 |cs(\uparrow v_4, i)(P)| \cdot |\Sigma_i|} &= \frac{|cs(\uparrow v_4, \mathbf{u})(P)|}{|cs(\uparrow v_4, 0)(P)| \cdot |\Sigma_0| + |cs(\uparrow v_4, 2)(P)| \cdot |\Sigma_2|} = \\ &= \frac{2}{2 \cdot 4 + 1 \cdot 4} = \frac{1}{6}. \end{aligned}$$

$$\begin{aligned} \frac{|cs(\uparrow v_5, \mathbf{u})(P)|}{\sum_{i=0}^3 |cs(\uparrow v_5, i)(P)| \cdot |\Sigma_i|} &= \frac{|cs(\uparrow v_5, \mathbf{u})(P)|}{|cs(\uparrow v_5, 0)(P)| \cdot |\Sigma_0| + |cs(\uparrow v_5, 2)(P)| \cdot |\Sigma_2|} = \\ &= \frac{3}{3 \cdot 4 + 1 \cdot 4} = \frac{3}{16}. \end{aligned}$$

$$\begin{aligned} \frac{|cs(\uparrow v_6, \mathbf{u})(P)|}{\sum_{i=0}^3 |cs(\uparrow v_6, i)(P)| \cdot |\Sigma_i|} &= \frac{|cs(\uparrow v_6, \mathbf{u})(P)|}{|cs(\uparrow v_6, 0)(P)| \cdot |\Sigma_0| + |cs(\uparrow v_6, 1)(P)| \cdot |\Sigma_1|} = \\ &= \frac{3}{3 \cdot 4 + 1 \cdot 5} = \frac{3}{17}. \end{aligned}$$

And now corollary 4.9.4 tells us that

$$\lim_{t \rightarrow \infty} \Phi(\mathbf{u}, P, t) = \prod_{i=1}^6 \frac{|cs(\uparrow v_i, \mathbf{u})(P)|}{\sum_{j=0}^3 |cs(\uparrow v_i, j)(P)| \cdot |\Sigma_j|} = \frac{2}{13} \cdot \frac{3}{19} \cdot \frac{1}{4} \cdot \frac{1}{6} \cdot \frac{3}{16} \cdot \frac{3}{17} = \frac{9}{268736}.$$

It is possible to introduce mutation operators for nonlinear GP which are ergodic in the sense of definition 4.9.1, but the easiest thing to do is probably just to define the family  $\mathcal{M}_{erg}$  to be the family of all permutations of the search space  $\Omega$ . The probability distribution  $p$  must then be concentrated on any subset of  $\mathcal{M}$  which satisfies the ergodicity requirement of definition 4.9.1. This would ensure that corollary 4.9.1 applies.

## 4.10 What Can Be Said in the Presence of Selection in the General Case?

Theorem 4.5.2 established in [5] which allows us to deduce results such as theorems 4.7.1 and 4.9.2, applies only in the absence of selection. The theme of the remainder of the current paper is to establish a few basic properties of the Markov chains associated to evolutionary algorithms in the presence of fitness-proportional selection (as described in definition 4.2.2). Throughout the rest of the paper we shall break up our algorithm, call it  $\mathcal{A}$  into sub-algorithms and then consider their composition. This idea will be made clear below:

**Proposition 4.10.1** Denote by  $\mathcal{A}$  an evolutionary algorithm determined by the cycle  $(s_1, s_2, \dots, s_n)$ . Fix  $i$  with  $1 < i < n$  and let  $\mathcal{B}$  and  $\mathcal{C}$  denote the sub-algorithms determined by the cycles  $(s_1, s_2, \dots, s_i)$  and  $(s_{i+1}, s_2, \dots, s_n)$  respectively. Recall from section 4.5 that  $\{p_{\mathbf{x}\mathbf{y}}^{\mathcal{A}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m}$ ,  $\{p_{\mathbf{x}\mathbf{y}}^{\mathcal{B}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m}$  and  $\{p_{\mathbf{x}\mathbf{y}}^{\mathcal{C}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m}$  denote the Markov transition matrices associated to the algorithms  $\mathcal{A}$ ,  $\mathcal{B}$  and  $\mathcal{C}$  respectively. Then we have

$$\{p_{\mathbf{x}\mathbf{y}}^{\mathcal{A}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m} = \{p_{\mathbf{x}\mathbf{y}}^{\mathcal{C}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m} \cdot \{p_{\mathbf{x}\mathbf{y}}^{\mathcal{B}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m}$$

where  $\cdot$  denotes the usual matrix multiplication.

*Proof:* Denote by  $\lambda$  a probability distribution on  $\Omega^m$ . Completing a cycle of the algorithm  $\mathcal{A}$  amounts to completing a cycle of  $\mathcal{B}$  and then completing a cycle of  $\mathcal{C}$ . The next generation probability distribution upon the completion of the cycle of  $\mathcal{B}$  with the input distribution  $\lambda$  is  $\{p_{\mathbf{x}\mathbf{y}}^{\mathcal{B}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m} \cdot \lambda$ . Likewise the next generation distribution obtained upon the completion of a cycle of the algorithm  $\mathcal{C}$  with the input distribution  $\{p_{\mathbf{x}\mathbf{y}}^{\mathcal{B}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m} \cdot \lambda$  is just

$$\{p_{\mathbf{x}\mathbf{y}}^{\mathcal{C}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m} \cdot (\{p_{\mathbf{x}\mathbf{y}}^{\mathcal{B}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m} \cdot \lambda) = (\{p_{\mathbf{x}\mathbf{y}}^{\mathcal{C}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m} \cdot \{p_{\mathbf{x}\mathbf{y}}^{\mathcal{B}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m}) \cdot \lambda$$

which means that  $\{p_{\mathbf{x}\mathbf{y}}^{\mathcal{A}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m} = \{p_{\mathbf{x}\mathbf{y}}^{\mathcal{C}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m} \cdot \{p_{\mathbf{x}\mathbf{y}}^{\mathcal{B}}\}_{\mathbf{x}, \mathbf{y} \in \Omega^m}$  since the equation above holds for an arbitrary input distribution  $\lambda$ .  $\square$

We now proceed to study the effects of selection alone. First of all it is convenient to observe the following general fact:

**Definition 4.10.1** Let  $\{p_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}}$  denote a Markov transition matrix on a finite set  $\mathcal{X}$ . Fix  $\mathbf{x} \in \mathcal{X}$ . We define the transition support of  $\mathbf{x}$  to be the set  $S(\mathbf{x}) = \{\mathbf{z} \mid p_{\mathbf{z}\mathbf{x}} > 0\}$  of all states  $\mathbf{z}$  from which it is possible to get to  $\mathbf{x}$ .

**Definition 4.10.2** Let  $\{p_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}}$  denote a Markov transition matrix on a finite set  $\mathcal{X}$ . Fix  $\mathbf{x}$  and  $\mathbf{y} \in \mathcal{X}$ . We say that  $\mathbf{y} \supseteq \mathbf{x}$  if  $S(\mathbf{y}) \supseteq S(\mathbf{x})$  and  $\forall \mathbf{z} \in S(\mathbf{y})$  we have  $p_{\mathbf{z}\mathbf{y}} \geq p_{\mathbf{z}\mathbf{x}}$ . Moreover, if either  $S(\mathbf{y}) \supsetneq S(\mathbf{x})$  or  $p_{\mathbf{z}\mathbf{y}} > p_{\mathbf{z}\mathbf{x}}$  for some  $\mathbf{z} \in S(\mathbf{x})$  we write  $\mathbf{y} \supset \mathbf{x}$ .

Proposition 4.10.2 below provides the reason for definitions 4.10.1 and 4.10.2:

**Proposition 4.10.2** Let  $\{p_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}}$  denote a Markov transition matrix on a finite set  $\mathcal{X}$ . Fix  $\mathbf{u}$  and  $\mathbf{v} \in \mathcal{X}$  with  $\mathbf{u} \supseteq \mathbf{v}$  and an input probability distribution  $\lambda$  on  $\mathcal{X}$ . Denote by  $\rho$  the output distribution ( $\rho = \{p_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \lambda$ ). Then we have  $\rho(\mathbf{u}) \geq \rho(\mathbf{v})$ . Moreover, if  $\mathbf{u} \supset \mathbf{v}$  and  $\lambda(\mathbf{x}) > 0$  for every  $\mathbf{x} \in \mathcal{X}$  then  $\rho(\mathbf{u}) > \rho(\mathbf{v})$ .

*Proof:* This is a straightforward verification of the definitions:

$$\rho(\mathbf{u}) = \sum_{\mathbf{z} \in \mathcal{X}} \lambda(\mathbf{z}) p_{\mathbf{z}\mathbf{u}} = \sum_{\mathbf{z} \in S(\mathbf{u})} \lambda(\mathbf{z}) p_{\mathbf{z}\mathbf{u}} \succ \sum_{\mathbf{z} \in S(\mathbf{v})} \lambda(\mathbf{z}) p_{\mathbf{z}\mathbf{v}} = \sum_{\mathbf{z} \in \mathcal{X}} \lambda(\mathbf{z}) p_{\mathbf{z}\mathbf{v}} = \rho(\mathbf{v})$$

$$\text{where } \succ = \begin{cases} \geq & \text{if } \mathbf{u} \supseteq \mathbf{v} \\ > & \text{if } \mathbf{u} \supset \mathbf{v} \end{cases} . \quad \square$$

The following mild technical condition on a Markov transition matrix (which is easily satisfiable by most transition matrices modeling crossover and mutation) will extend proposition 4.10.2.

**Definition 4.10.3** We call a Markov transition matrix  $\{q_{xy}\}_{x,y \in \mathcal{X}}$  non-annihilating if  $\forall y \in \mathcal{X} \exists x \in \mathcal{X}$  such that  $q_{xy} > 0$ .

The main reason for introducing definition 4.10.3 is the following fact:

**Proposition 4.10.3** A given Markov transition matrix  $\{q_{xy}\}_{x,y \in \mathcal{X}}$  is non-annihilating if and only if for every input probability distribution  $\lambda$  on  $\mathcal{X}$  with  $\lambda(x) > 0$  for every  $x \in \mathcal{X}$ , the output distribution  $\rho = \{q_{xy}\}_{x,y \in \mathcal{X}} \cdot \lambda$  also satisfies the property that  $\rho(x) > 0$  for every  $x \in \mathcal{X}$ .

*Proof:* Given an input distribution  $\lambda$  with  $\lambda(x) > 0$  for every  $x \in \mathcal{X}$  and any state  $y \in \mathcal{X}$ , we have  $\rho(y) = \sum_{x \in \mathcal{X}} \lambda(x) \cdot q_{xy} > 0$  if and only if  $\lambda(z) \cdot q_{zy} > 0$  for some  $z \in \mathcal{X}$  if and only if  $q_{zy} > 0$  for some  $z \in \mathcal{X}$  (since  $\lambda(z) > 0$  automatically by assumption) if and only if the transition matrix  $\{q_{xy}\}_{x,y \in \mathcal{X}}$  is non-annihilating.  $\square$

Before proceeding any further it is worthwhile to mention that a product of non-annihilating transition matrices is non-annihilating:

**Corollary 4.10.4** Given non-annihilating Markov transition matrices  $\{q_{xy}\}_{x,y \in \mathcal{X}}$  and  $\{m_{xy}\}_{x,y \in \mathcal{X}}$ , the matrix  $\{r_{xy}\}_{x,y \in \mathcal{X}} = \{m_{xy}\}_{x,y \in \mathcal{X}} \cdot \{q_{xy}\}_{x,y \in \mathcal{X}}$  is non-annihilating as well.

*Proof:* Given an input distribution  $\lambda$  with  $\lambda(x) > 0$ , since  $\{q_{xy}\}_{x,y \in \mathcal{X}}$  is non-annihilating, the “intermediate” output distribution  $\mu = \{q_{xy}\}_{x,y \in \mathcal{X}}(\lambda)$  also has the property that  $\mu(x) > 0$  for all  $x \in \mathcal{X}$ . Now we have

$$\rho = \{r_{xy}\}_{x,y \in \mathcal{X}} \cdot \lambda = \{m_{xy}\}_{x,y \in \mathcal{X}} \cdot (\{q_{xy}\}_{x,y \in \mathcal{X}} \cdot \lambda) = \{m_{xy}\}_{x,y \in \mathcal{X}} \cdot \mu$$

also has the property that  $\rho(x) > 0$  for all  $x \in \mathcal{X}$  since it is an output of  $\mu$  under the non-annihilating transition matrix  $\{m_{xy}\}_{x,y \in \mathcal{X}}$ . By proposition 4.10.3, the transition matrix  $\{r_{xy}\}_{x,y \in \mathcal{X}}$  is non-annihilating as well.  $\square$

Though quite elementary, proposition 4.10.2 readily implies subtle and rather general inequalities about the stationary distributions of the Markov chains for which the last elementary step is selection:

**Corollary 4.10.5** Let  $\{q_{xy}\}_{x,y \in \mathcal{X}}$  and  $\{p_{xy}\}_{x,y \in \mathcal{X}}$  denote Markov transition matrices on a finite set  $\mathcal{X}$ . Fix  $\mathbf{u}$  and  $\mathbf{v} \in \mathcal{X}$  with  $\mathbf{u} \succeq \mathbf{v}$  where the  $\succeq$  and  $\succ$  relations are meant with respect to the matrix  $\{p_{xy}\}_{x,y \in \mathcal{X}}$  and an input probability distribution  $\lambda$  on  $\mathcal{X}$ . Denote by  $\rho$  the output distribution of the composed matrix  $\{p_{xy}\}_{x,y \in \mathcal{X}} \cdot \{q_{xy}\}_{x,y \in \mathcal{X}}$  ( $\rho = \{p_{xy}\}_{x,y \in \mathcal{X}} \cdot \{q_{xy}\}_{x,y \in \mathcal{X}} \cdot \lambda$ ). Then we have  $\rho(\mathbf{u}) \geq \rho(\mathbf{v})$ . Suppose, in addition, the matrix  $\{q_{xy}\}_{x,y \in \mathcal{X}}$  is non-annihilating. Now, if  $\mathbf{u} \succ \mathbf{v}$ , and  $\lambda(x) > 0$  for every  $x \in \mathcal{X}$  then  $\rho(\mathbf{u}) \succ \rho(\mathbf{v})$ .

*Proof:* Since

$$\{p_{xy}\}_{x,y \in \mathcal{X}} \cdot \{q_{xy}\}_{x,y \in \mathcal{X}} \cdot \lambda = \{p_{xy}\}_{x,y \in \mathcal{X}} \cdot (\{q_{xy}\}_{x,y \in \mathcal{X}} \cdot \lambda) = \{p_{xy}\}_{x,y \in \mathcal{X}} \cdot \mu$$

where  $\mu = \{q_{xy}\}_{x,y \in \mathcal{X}} \cdot \lambda$ , the desired conclusions follow by applying proposition 4.10.2 to the matrix  $\{p_{xy}\}_{x,y \in \mathcal{X}}$  and the input distribution  $\mu$ . For the second conclusion we use the assumption that  $\{q_{xy}\}_{x,y \in \mathcal{X}}$  is non-annihilating to deduce that  $\mu(\mathbf{x}) > 0$  for every  $\mathbf{x}$ .  $\square$

As an almost immediate consequence we deduce the following fact:

**Corollary 4.10.6** *Let  $\{q_{xy}\}_{x,y \in \mathcal{X}}$ ,  $\{m_{xy}\}_{x,y \in \mathcal{X}}$  and  $\{p_{xy}\}_{x,y \in \mathcal{X}}$  denote Markov transition matrices on a finite set  $\mathcal{X}$ . Suppose that the matrices  $\{q_{xy}\}_{x,y \in \mathcal{X}}$  and  $\{p_{xy}\}_{x,y \in \mathcal{X}}$  are non-annihilating while the matrix  $\{m_{xy}\}_{x,y \in \mathcal{X}}$  has the property that  $\forall \mathbf{x}, \mathbf{y} \in \mathcal{X} \ m_{xy} > 0$ . Fix  $\mathbf{u}$  and  $\mathbf{v} \in \mathcal{X}$  with  $\mathbf{u} \succeq \mathbf{v}$  where the  $\succeq$  and  $\succ$  relations are meant with respect to the matrix  $\{p_{xy}\}_{x,y \in \mathcal{X}}$ . Then the Markov chain determined by either one of the composed matrices  $\{p_{xy}\}_{x,y \in \mathcal{X}} \cdot \{m_{xy}\}_{x,y \in \mathcal{X}} \cdot \{q_{xy}\}_{x,y \in \mathcal{X}}$  or  $\{p_{xy}\}_{x,y \in \mathcal{X}} \cdot \{q_{xy}\}_{x,y \in \mathcal{X}} \cdot \{m_{xy}\}_{x,y \in \mathcal{X}}$  is irreducible. Let  $\pi$  denote the unique stationary distribution of the composed chain. Then we have  $\pi(\mathbf{u}) \geq \pi(\mathbf{v})$ . Suppose, in addition, the matrix  $\{q_{xy}\}_{x,y \in \mathcal{X}}$  is non-annihilating. Now, if  $\mathbf{u} \succ \mathbf{v}$ , and  $\lambda(\mathbf{x}) > 0$  for every  $\mathbf{x} \in \mathcal{X}$  then  $\pi(\mathbf{u}) > \pi(\mathbf{v})$ .*

*Proof:* The irreducibility of the composed chain is left as an exercise for the reader. As a hint, the reader may notice that from the assumptions that  $\{m_{xy}\}_{x,y \in \mathcal{X}}$  has all positive entries while  $\{p_{xy}\}_{x,y \in \mathcal{X}}$  and  $\{q_{xy}\}_{x,y \in \mathcal{X}}$  are non-annihilating, it follows that every one of the composed matrices has all positive entries and, hence, determines an irreducible Markov chain. The second conclusion follows from the fact that  $\{p_{xy}\}_{x,y \in \mathcal{X}}$  is the leftmost matrix in the composition by applying corollary 4.10.5 to the matrix  $\{p_{xy}\}_{x,y \in \mathcal{X}} \cdot \{a_{xy}\}_{x,y \in \mathcal{X}}$  where  $\{a_{xy}\}_{x,y \in \mathcal{X}} = \{m_{xy}\}_{x,y \in \mathcal{X}} \cdot \{q_{xy}\}_{x,y \in \mathcal{X}}$  or  $\{a_{xy}\}_{x,y \in \mathcal{X}} = \{q_{xy}\}_{x,y \in \mathcal{X}} \cdot \{m_{xy}\}_{x,y \in \mathcal{X}}$  with  $\pi$  being the input distribution which is then also the output distribution by stationarity. The condition of corollary 4.10.5 is satisfied thanks to corollary 4.10.4.  $\square$

When applying corollary 4.10.6 we have in mind that  $\{q_{xy}\}_{x,y \in \mathcal{X}}$  is the Markov transition matrix corresponding to recombination (i. e. a sub-algorithm determined by a single elementary step of type 2: see definitions 4.3.1 and 4.2.5),  $\{p_{xy}\}_{x,y \in \mathcal{X}}$  is the Markov transition matrix corresponding to selection (i. e. a sub-algorithm determined by a single elementary step of type 1: see definition 4.2.2) and  $\{m_{xy}\}_{x,y \in \mathcal{X}}$  is the Markov transition matrix corresponding to mutation. In order to apply proposition 4.10.2 to the case of fitness-proportional selection we need to determine the relation  $\succeq$  and  $\succ$  for this special case. Although this task is not difficult, it requires a careful “bookkeeping” analysis. This will be the subject of the next section. We end the current section with an immediate consequence (basically a restatement) of corollary 4.10.6:<sup>12</sup>

<sup>12</sup>It is worth mentioning that fitness-proportional selection is not the only possible type of selection. Other elementary steps of type 1 include, for instance, tournament selection and rank selection.

**Corollary 4.10.7** *Suppose we are given an evolutionary algorithm  $\mathcal{A}$  determined by the elementary steps  $s_1$ ,  $s_2$  and  $s_3$  where  $s_1$  and  $s_2$  are any elementary step (usually one of them is selection and the other is mutation) which define non-annihilating Markov transition matrices and such that one of these matrices has all positive entries, while  $s_3$  is the elementary step of type 1, i.e. selection. As before, let  $\{p_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}$  denote the Markov transition matrix determined by the elementary step  $s_3$  and  $\triangleright$  and  $\trianglerighteq$  are defined with respect to the transition matrix  $\{p_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}$ . Then the Markov chain determined by the algorithm  $\mathcal{A}$  with state space  $\mathcal{X} = \Omega^m$  is irreducible and its unique stationary distribution  $\pi$  satisfies  $\pi(\mathbf{x}) \geq \pi(\mathbf{y})$  and  $\pi(\mathbf{x}) > \pi(\mathbf{y})$  whenever  $\mathbf{x} \trianglerighteq \mathbf{y}$  and  $\mathbf{x} \triangleright \mathbf{y}$  respectively.*

## 4.11 What are the relations $\triangleright$ and $\trianglerighteq$ for the case of fitness-proportional selection?

This section is devoted to classifying the relations  $\triangleright$  and  $\trianglerighteq$  for the case of fitness-proportional selection. Although this task is not difficult, it requires a careful step-by-step analysis. The reader who is interested only in the net results can read only definitions 4.11.1 and 4.11.2, example 4.11.1, theorem 4.11.5 followed by examples 4.11.3, 4.11.4, 4.11.5 and 4.11.6 and theorem 4.11.6 which is illustrated by example 4.11.7. It is recommended (but not essential for understanding) that the reader does not omit the discussion between example 4.11.6 and theorem 4.11.6. We also strongly recommend that the reader makes him/herself familiar with lemma 4.11.1 since this fact is rather simple and reveals a very important step in the classification process.

**Definition 4.11.1** *Fix a population  $\mathbf{x} = (x_1, x_2, \dots, x_m) \in \Omega^m$  and denote by  $I(\mathbf{x}) = \{x \mid x = x_i \text{ for some } i \text{ with } 1 \leq i \leq m\}$  the set of all individuals in the population  $\mathbf{x}$ .*

**Lemma 4.11.1** *Given populations  $\mathbf{x}$  and  $\mathbf{y}$ , we have  $S(\mathbf{x}) \supseteq S(\mathbf{y})$  if and only if  $I(\mathbf{x}) \supseteq I(\mathbf{y})$ . In particular, a necessary condition for  $\mathbf{x} \trianglerighteq \mathbf{y}$  is that  $I(\mathbf{x}) \supseteq I(\mathbf{y})$ . Moreover, if  $I(\mathbf{x}) \subsetneq I(\mathbf{y})$  then  $\mathbf{x} \triangleright \mathbf{y}$ .*

*Proof:* Since individuals can only disappear (and new individuals can not appear) upon the completion of the elementary step of fitness-proportional selection (see definition 4.2.2) it follows immediately that for any populations  $\mathbf{z}$  and  $\mathbf{w}$  we have  $p_{\mathbf{z},\mathbf{w}} \geq 0$  if and only if  $I(\mathbf{z}) \supseteq I(\mathbf{w})$ . In other words  $S(\mathbf{w}) = \{\mathbf{z} \mid I(\mathbf{z}) \supseteq I(\mathbf{w})\}$ . It follows immediately now that if  $I(\mathbf{y}) \supseteq I(\mathbf{x})$  then  $S(\mathbf{x}) \supseteq S(\mathbf{y})$ . On the other hand, if  $S(\mathbf{x}) \supseteq S(\mathbf{y})$ , then, since  $\mathbf{y} \in S(\mathbf{y}) \subseteq S(\mathbf{x})$  we also have  $I(\mathbf{y}) \supseteq I(\mathbf{x})$  according to the characterization given above. We deduce now that  $I(\mathbf{y}) \supseteq I(\mathbf{x})$  if and only if  $S(\mathbf{x}) \supseteq S(\mathbf{y})$ . In particular, it follows immediately from the previous statement that  $I(\mathbf{y}) \supsetneq I(\mathbf{x})$  if and only if  $S(\mathbf{x}) \supsetneq S(\mathbf{y})$ . All of the remaining conclusions follow immediately from definition 4.10.2.  $\square$

**Definition 4.11.2** *Given a population  $\mathbf{x} = (x_1, x_2, \dots, x_m)$  and an individual  $x \in I(\mathbf{x})$ , denote by  $n(\mathbf{x}, x) = |\{i \mid x = x_i\}|$  the number of times  $x$  occurs in the population  $\mathbf{x}$ .*

**Example 4.11.1** Suppose

$$\mathbf{x} = (a, a, a, b, c, a, b, b, b) \text{ and } \mathbf{y} = (b, c, c, c, a, b, b, d, b).$$

Then  $I(\mathbf{x}) = \{a, b, c\}$  and  $I(\mathbf{y}) = \{a, b, c, d\}$ . We also have

$$n(\mathbf{x}, a) = n(\mathbf{x}, b) = 4 \text{ and } n(\mathbf{x}, c) = 1.$$

Likewise,

$$n(\mathbf{y}, a) = 1, n(\mathbf{x}, b) = 4, n(\mathbf{x}, c) = 3 \text{ and } n(\mathbf{x}, d) = 1.$$

According to definition 4.2.2, when performing fitness-proportional selection, the individuals are chosen independently with probability proportional to their fitness. Thus, if  $\mathbf{z} = (z_1, z_2, \dots, z_m)$  and  $\mathbf{x} = (x_1, x_2, \dots, x_m)$  is obtained from  $\mathbf{z}$  by performing fitness-proportional selection, then the probability that  $x_i = z$  for a given  $z \in I(\mathbf{z})$  is  $\frac{n(\mathbf{z}, z) \cdot f(z)}{\sum_{i=1}^m f(z_i)}$ . Thus  $p_{\mathbf{z}, \mathbf{x}} = \prod_{j=1}^m \frac{n(\mathbf{z}, x_j) \cdot f(x_j)}{\sum_{i=1}^m f(z_i)}$ . Moreover, every  $x \in I(\mathbf{x})$  occurs in the above product  $n(\mathbf{z}, x_j)$  times while every  $z \in I(\mathbf{z})$  occurs  $n(\mathbf{z}, z)$  times in the denominator sum of each of the multiples and so we deduce the following:

**Proposition 4.11.2** Given populations  $\mathbf{x} = (x_1, x_2, \dots, x_m)$  and  $\mathbf{z} = (z_1, z_2, \dots, z_m)$  we have

$$p_{\mathbf{z}, \mathbf{x}} = \begin{cases} \left( \frac{1}{\sum_{z \in I(\mathbf{z})} n(\mathbf{z}, z) \cdot f(z)} \right)^m \prod_{x \in I(\mathbf{x})} (n(\mathbf{z}, x))^{n(\mathbf{x}, x)} \cdot (f(x))^{n(\mathbf{x}, x)} & \text{if } I(\mathbf{x}) \subseteq I(\mathbf{z}) \\ 0 & \text{otherwise.} \end{cases}$$

In particular,  $p_{\mathbf{z}, \mathbf{x}}$  does not depend on the way the individuals in  $\mathbf{z}$  and in  $\mathbf{x}$  are ordered, but only depends on the sets  $I(\mathbf{x})$  and  $I(\mathbf{z})$  and the numbers  $n(\mathbf{x}, x)$  for  $x \in I(\mathbf{x})$  and  $n(\mathbf{z}, z)$  for  $z \in I(\mathbf{z})$ . In other words, if  $\sigma$  and  $\tau$  demote arbitrary permutations of the set  $\{1, 2, \dots, m\}$ , If  $\mathbf{x}_\sigma = (x_{\sigma(1)}, x_{\sigma(2)}, \dots, x_{\sigma(m)})$  and  $\mathbf{z}_\tau = (z_{\tau(1)}, z_{\tau(2)}, \dots, z_{\tau(m)})$  then  $p_{\mathbf{z}_\sigma, \mathbf{x}_\tau} = p_{\mathbf{z}, \mathbf{x}}$ .

In order to continue the investigation of the  $\supseteq$  relation for the case of fitness-proportional selection, it is convenient to introduce the following notions:

**Definition 4.11.3** Given populations  $\mathbf{x}$  and  $\mathbf{y}$  of size  $m$  let

$$I(\mathbf{y}|\mathbf{x}) = \{y \mid y \in I(\mathbf{y}), n(\mathbf{y}, y) > n(\mathbf{x}, y)\}$$

(if  $y \notin I(\mathbf{x})$  then  $n(\mathbf{x}, y) = 0$ ). Moreover, for  $y \in I(\mathbf{y}|\mathbf{x})$  let  $\kappa(\mathbf{y}|\mathbf{x}, y) = n(\mathbf{y}, y) - n(\mathbf{x}, y)$ .

**Example 4.11.2** Continuing with example 4.11.1, we have  $n(\mathbf{x}, x) > n(\mathbf{y}, x)$  if and only if  $x = a$  and so  $I(\mathbf{x}|\mathbf{y}) = \{a\}$ . Likewise,  $n(\mathbf{y}, y) < n(\mathbf{x}, y)$  if and only if  $y = c$  or  $y = d$  (since  $d \notin I(\mathbf{x})$  according to definition 4.11.3 we have  $n(\mathbf{x}, d) = 0 < 1 = n(\mathbf{x}, d)$ ) and so  $I(\mathbf{y}|\mathbf{x}) = \{c, d\}$ . Moreover, we also have  $\kappa(\mathbf{x}|\mathbf{y}, a) = 4 - 1 = 3$ ,  $\kappa(\mathbf{y}|\mathbf{x}, c) = 3 - 1 = 2$  and  $\kappa(\mathbf{y}|\mathbf{x}, d) = 1 - 0 = 1$ .

The sets  $I(\mathbf{y}|\mathbf{x})$  play a crucial role in discovering a sufficient and necessary condition for a population  $\mathbf{x} \supseteq \mathbf{y}$  in view of the fact below:

**Lemma 4.11.3** Given populations  $\mathbf{z}$ ,  $\mathbf{y}$  and  $\mathbf{x}$  with  $I(\mathbf{z}) \supseteq I(\mathbf{y}) \supseteq I(\mathbf{x})$ , have the following:

$$\sum_{x \in I(\mathbf{x}|\mathbf{y})} \kappa(\mathbf{x}|\mathbf{y}, x) = \sum_{y \in I(\mathbf{y}|\mathbf{x})} \kappa(\mathbf{y}|\mathbf{x}, y)$$

$$\frac{p_{\mathbf{z}\mathbf{x}}}{p_{\mathbf{z}\mathbf{y}}} = \frac{\prod_{x \in I(\mathbf{x}|\mathbf{y})} (n(\mathbf{z}, x))^{\kappa(\mathbf{x}|\mathbf{y}, x)} \cdot (f(x))^{\kappa(\mathbf{x}|\mathbf{y}, x)}}{\prod_{y \in I(\mathbf{y}|\mathbf{x})} (n(\mathbf{z}, y))^{\kappa(\mathbf{y}|\mathbf{x}, y)} \cdot (f(y))^{\kappa(\mathbf{y}|\mathbf{x}, y)}}.$$

*Proof:* Given populations  $\mathbf{z}$ ,  $\mathbf{y}$  and  $\mathbf{x}$ , from definition 4.11.3 it follows that for every  $x \in I(\mathbf{x})$  we have

$$n(\mathbf{x}, x) = \begin{cases} \min(n(\mathbf{x}, x), n(\mathbf{y}, x)) & \text{if } x \notin I(\mathbf{x}|\mathbf{y}) \\ \min(n(\mathbf{x}, x), n(\mathbf{y}, x)) + \kappa(\mathbf{x}|\mathbf{y}, x) & \text{if } x \in I(\mathbf{x}|\mathbf{y}) \end{cases}.$$

Likewise, for every  $y \in I(\mathbf{y})$  we have

$$n(\mathbf{y}, y) = \begin{cases} \min(n(\mathbf{x}, y), n(\mathbf{y}, y)) & \text{if } y \notin I(\mathbf{y}|\mathbf{x}) \text{ and } y \in I(\mathbf{x}) \\ \min(n(\mathbf{x}, y), n(\mathbf{y}, y)) + \kappa(\mathbf{y}|\mathbf{x}, y) & \text{if } y \in I(\mathbf{y}|\mathbf{x}) \text{ and } y \in I(\mathbf{x}) \\ \kappa(\mathbf{y}|\mathbf{x}, y) & \text{if } y \in I(\mathbf{y}|\mathbf{x}) \text{ and } y \notin I(\mathbf{x}) \end{cases}.$$

Since there are totally  $m$  elements in every population we must have

$$\sum_{y \in I(\mathbf{y})} n(\mathbf{y}, y) = \sum_{x \in I(\mathbf{x})} n(\mathbf{x}, x) = m.$$

Rearranging the terms in both sides of the last equation according to the observations made above, we obtain

$$\begin{aligned} & \sum_{x \in I(\mathbf{x})} \min(n(\mathbf{x}, x), n(\mathbf{y}, x)) + \sum_{x \in I(\mathbf{x}|\mathbf{y})} \kappa(\mathbf{x}|\mathbf{y}, x) = \\ & = \sum_{y \in I(\mathbf{x})} \min(n(\mathbf{x}, y), n(\mathbf{y}, y)) + \sum_{y \in I(\mathbf{y}|\mathbf{x})} \kappa(\mathbf{y}|\mathbf{x}, y) \end{aligned}$$

and the first desired equation follows by subtracting  $\sum_{x \in I(\mathbf{x})} \min(n(\mathbf{x}, x), n(\mathbf{y}, x))$  from both sides. The second equation follows by rearranging the multiples in the formula of proposition 4.11.2 according to the equation above and letting  $k(\mathbf{z}) = (\frac{1}{\sum_{z \in I(\mathbf{z})} n(\mathbf{z}, z) \cdot f(z)})^m$  so that we can write:

$$\begin{aligned} p_{\mathbf{z}\mathbf{x}} &= k(\mathbf{z}) \cdot \prod_{x \in I(\mathbf{x})} (n(\mathbf{z}, x))^{\min(n(\mathbf{x}, x), n(\mathbf{y}, x))} \cdot (f(x))^{\min(n(\mathbf{x}, x), n(\mathbf{y}, x))} \times \\ & \times \prod_{x \in I(\mathbf{x}|\mathbf{y})} (n(\mathbf{z}, x))^{\kappa(\mathbf{x}|\mathbf{y}, x)} \cdot (f(x))^{\kappa(\mathbf{x}|\mathbf{y}, x)} \end{aligned}$$

and, likewise,

$$p_{\mathbf{z}\mathbf{y}} = k(\mathbf{z}) \cdot \prod_{y \in I(\mathbf{y})} (n(\mathbf{z}, y))^{\min(n(\mathbf{x}, y), n(\mathbf{y}, y))} \cdot (f(y))^{\min(n(\mathbf{x}, y), n(\mathbf{y}, y))} \times$$

$$\times \prod_{y \in I(\mathbf{y}|\mathbf{x})} (n(\mathbf{z}, y))^{\kappa(\mathbf{y}|\mathbf{x}, y)} \cdot (f(y))^{\kappa(\mathbf{y}|\mathbf{x}, y)}.$$

Now the common factor

$$k(\mathbf{z}) \cdot \prod_{x \in I(\mathbf{x})} (n(\mathbf{z}, x))^{\min(n(\mathbf{x}, x), n(\mathbf{y}, x))} \cdot (f(x))^{\min(n(\mathbf{x}, x), n(\mathbf{y}, x))}$$

on the top and the bottom of the ratio  $\frac{p_{\mathbf{z}\mathbf{x}}}{p_{\mathbf{z}\mathbf{y}}}$  is canceled out and we obtain the desired formula.  $\square$

In fact, according to lemma 4.11.1, we have  $\mathbf{x} \supseteq \mathbf{y} \implies I(\mathbf{y}) \supseteq I(\mathbf{x})$ . Moreover, since new individuals can not appear as a result of selection, whenever  $\mathbf{z} \in S(\mathbf{y})$  (see definition 4.10.1 for the meaning of  $S(\mathbf{y})$ ) we must have  $I(\mathbf{z}) \supseteq I(\mathbf{y})$ . Therefore,  $\mathbf{x} \supseteq \mathbf{y}$  if and only if  $I(\mathbf{y}) \supseteq I(\mathbf{x})$  and  $\forall \mathbf{z}$  such that  $I(\mathbf{z}) \supseteq I(\mathbf{y})$  we have  $p_{\mathbf{z}\mathbf{x}} \geq p_{\mathbf{z}\mathbf{y}}$ . The condition  $p_{\mathbf{z}\mathbf{x}} \geq p_{\mathbf{z}\mathbf{y}}$  can be restated equivalently as  $\frac{p_{\mathbf{z}\mathbf{x}}}{p_{\mathbf{z}\mathbf{y}}} \geq 1$ . But, thanks to lemma 4.11.3, we have

$$\begin{aligned} \frac{p_{\mathbf{z}\mathbf{x}}}{p_{\mathbf{z}\mathbf{y}}} &= \frac{\prod_{x \in I(\mathbf{x}|\mathbf{y})} (n(\mathbf{z}, x))^{\kappa(\mathbf{x}|\mathbf{y}, x)} \cdot (f(x))^{\kappa(\mathbf{x}|\mathbf{y}, x)}}{\prod_{y \in I(\mathbf{y}|\mathbf{x})} (n(\mathbf{z}, y))^{\kappa(\mathbf{y}|\mathbf{x}, y)} \cdot (f(y))^{\kappa(\mathbf{y}|\mathbf{x}, y)}} = \\ &= \frac{\prod_{x \in I(\mathbf{x}|\mathbf{y})} (n(\mathbf{z}, x))^{\kappa(\mathbf{x}|\mathbf{y}, x)}}{\prod_{y \in I(\mathbf{y}|\mathbf{x})} (n(\mathbf{z}, y))^{\kappa(\mathbf{y}|\mathbf{x}, y)}} \cdot \frac{\prod_{x \in I(\mathbf{x}|\mathbf{y})} (f(x))^{\kappa(\mathbf{x}|\mathbf{y}, x)}}{\prod_{y \in I(\mathbf{y}|\mathbf{x})} (f(y))^{\kappa(\mathbf{y}|\mathbf{x}, y)}} \geq 1 \\ &\iff \frac{\prod_{x \in I(\mathbf{x}|\mathbf{y})} (n(\mathbf{z}, x))^{\kappa(\mathbf{x}|\mathbf{y}, x)}}{\prod_{y \in I(\mathbf{y}|\mathbf{x})} (n(\mathbf{z}, y))^{\kappa(\mathbf{y}|\mathbf{x}, y)}} \geq \frac{\prod_{y \in I(\mathbf{y}|\mathbf{x})} (f(y))^{\kappa(\mathbf{y}|\mathbf{x}, y)}}{\prod_{y \in I(\mathbf{x}|\mathbf{y})} (f(x))^{\kappa(\mathbf{x}|\mathbf{y}, x)}}. \end{aligned}$$

Observing that

$$\frac{\prod_{y \in I(\mathbf{y}|\mathbf{x})} (f(y))^{\kappa(\mathbf{y}|\mathbf{x}, y)}}{\prod_{y \in I(\mathbf{x}|\mathbf{y})} (f(x))^{\kappa(\mathbf{x}|\mathbf{y}, x)}}$$

does not depend on  $\mathbf{z}$  at all we deduce the following:

**Lemma 4.11.4** *Given populations  $\mathbf{x}$  and  $\mathbf{y}$  of size  $m$ , we have  $\mathbf{x} \supseteq \mathbf{y}$  if and only if  $I(\mathbf{y}) \supseteq I(\mathbf{x})$  and*

$$\min_{I(\mathbf{z}) \supseteq I(\mathbf{x})} \frac{\prod_{x \in I(\mathbf{x}|\mathbf{y})} (n(\mathbf{z}, x))^{\kappa(\mathbf{x}|\mathbf{y}, x)}}{\prod_{y \in I(\mathbf{y}|\mathbf{x})} (n(\mathbf{z}, y))^{\kappa(\mathbf{y}|\mathbf{x}, y)}} \geq \frac{\prod_{y \in I(\mathbf{y}|\mathbf{x})} (f(y))^{\kappa(\mathbf{y}|\mathbf{x}, y)}}{\prod_{y \in I(\mathbf{x}|\mathbf{y})} (f(x))^{\kappa(\mathbf{x}|\mathbf{y}, x)}}.$$

Thanks to lemma 4.11.4, the rest of our analysis boils down to constructing a population  $\mathbf{z}$  which minimizes the ratio over  $\mathbf{z}$  with  $I(\mathbf{z}) \supseteq I(\mathbf{y})$ . In view of proposition 4.11.2, without loss of generality, we can assume that the first  $|I(\mathbf{y})|$  individuals of  $\mathbf{z}$  enumerate the elements of  $I(\mathbf{y})$ , i.e.  $\mathbf{z} = \{y_1, y_2, \dots, y_{|I(\mathbf{y})|}, z_1, z_2, \dots, z_{m-|I(\mathbf{y})|}\}$ . Our goal is then to select  $z_1, z_2, \dots, z_{m-|I(\mathbf{y})|}$  in a way which minimizes this ratio. First, it is worth pointing out, that unless  $\mathbf{y} = \mathbf{x}_\sigma$  for some permutation  $\sigma$  of  $\{1, 2, \dots, m\}$  in the sense described in proposition 4.11.2 (in which case we trivially have  $\mathbf{x} \supseteq \mathbf{y}$  thanks to proposition 4.11.2), we can assume that  $I(\mathbf{y}|\mathbf{x}) \neq \emptyset$ . (Indeed,



according to lemma 4.11.3, we have  $\sum_{x \in I(\mathbf{x}|\mathbf{y})} \kappa(\mathbf{x}|\mathbf{y}, x) = \sum_{y \in I(\mathbf{y}|\mathbf{x})} \kappa(\mathbf{y}|\mathbf{x}, y)$ . If  $I(\mathbf{y}|\mathbf{x}) \neq \emptyset$  then  $\sum_{x \in I(\mathbf{x}|\mathbf{y})} \kappa(\mathbf{x}|\mathbf{y}, x) = \sum_{y \in I(\mathbf{y}|\mathbf{x})} \kappa(\mathbf{y}|\mathbf{x}, y) = 0$  which forces  $I(\mathbf{x}|\mathbf{y}) = \emptyset$ . But then for every  $y \in I(\mathbf{y})$  we have  $y \notin I(\mathbf{y}|\mathbf{x}) \implies I(\mathbf{y}) \subseteq I(\mathbf{x})$  and  $n(\mathbf{y}, y) \leq n(\mathbf{x}, y)$  and, since  $I(\mathbf{x}|\mathbf{y}) = \emptyset$ , also  $n(\mathbf{y}, y) \geq n(\mathbf{x}, y)$ . Summarizing, we obtain  $n(\mathbf{y}, y) = n(\mathbf{x}, y) \forall y \in I(\mathbf{y}) = I(\mathbf{y})$  which means that  $\mathbf{y} = \mathbf{x}_\sigma$ . Next, we observe that  $\forall i$  we have  $z_i \in I(\mathbf{y}|\mathbf{x})$ . (If not, then for some  $i$  we have  $z_i \notin I(\mathbf{y}|\mathbf{x})$ . In such a case, since replacing  $z_i$  with an element of  $I(\mathbf{y}|\mathbf{x}) \neq \emptyset$  will increase the denominator and either decrease (in case if  $z_i \in I(\mathbf{x}|\mathbf{y})$ ) or not influence the numerator in any way (since it is clear from definition 4.11.3 that  $I(\mathbf{y}|\mathbf{x}) \cap I(\mathbf{x}|\mathbf{y}) = \emptyset$ ) of the ratio on the L.H.S. of the inequality of lemma 4.11.4. This in turn would only decrease this ratio so that  $\mathbf{z}$  can not minimize it.) Since  $z_i$ s are chosen from the set  $I(\mathbf{y}|\mathbf{x})$ , and  $I(\mathbf{y}|\mathbf{x}) \cap I(\mathbf{x}|\mathbf{y}) = \emptyset$ , for every  $x \in I(\mathbf{x}|\mathbf{y})$   $x \neq z_i$  for  $1 \leq i \leq m - |I(\mathbf{y})|$  and  $x = y_j$  for a unique  $j$  with  $1 \leq j \leq |I(\mathbf{y})|$  (since  $I(\mathbf{y}) \supseteq I(\mathbf{x})$  and  $y_1, y_2, \dots, y_{|I(\mathbf{y})|}$  enumerate the elements of  $I(\mathbf{y})$ ) we have  $n(\mathbf{z}, x) = 1$  for every  $x \in I(\mathbf{x}|\mathbf{y})$ . But then the numerator of the ratios on the L.H.S. of lemma 4.11.4,  $\prod_{x \in I(\mathbf{x}|\mathbf{y})} (n(\mathbf{z}, x))^{\kappa(\mathbf{x}|\mathbf{y}, x)} = 1$  and the L.H.S. of the inequality in lemma 4.11.3 simplifies to

$$\min_{\mathbf{z} \in Q} \frac{1}{\prod_{y \in I(\mathbf{y}|\mathbf{x})} (n(\mathbf{z}, y))^{\kappa(\mathbf{y}|\mathbf{x}, y)}}$$

where  $Q = \{\mathbf{z} \mid \mathbf{z} = (y_1, y_2, \dots, y_{|I(\mathbf{y})|}, z_1, z_2, \dots, z_{m-|I(\mathbf{y})|}), y_1, y_2, \dots, y_{|I(\mathbf{y})|}$  enumerate  $I(\mathbf{y})$  and  $z_i \in I(\mathbf{y}|\mathbf{x})\}$ . Moreover, notice that  $z_i$ s can be chosen arbitrarily from the set  $I(\mathbf{y}|\mathbf{x})$  while for every  $y \in I(\mathbf{y}|\mathbf{x})$  there exists exactly one  $1 \leq j \leq |I(\mathbf{y})|$  with  $y = y_j$ . We then have  $n(\mathbf{z}, y) = 1 + |\{i \mid 1 \leq i \leq m - |I(\mathbf{y})|, z_i = y\}|$  and  $\sum_{y \in I(\mathbf{y}|\mathbf{x})} n(\mathbf{z}, y) = |I(\mathbf{y}|\mathbf{x})| + m - |I(\mathbf{y})|$ . On the other hand, given a finite sequence of natural numbers  $\{n_y\}_{y \in I(\mathbf{y}|\mathbf{x})}$  satisfying the constraint  $\sum_{y \in I(\mathbf{y}|\mathbf{x})} n_y = |I(\mathbf{y}|\mathbf{x})| + m - |I(\mathbf{y})|$ , we can construct  $\mathbf{z} = \{y_1, y_2, \dots, y_{|I(\mathbf{y})|}, z_1, z_2, \dots, z_{m-|I(\mathbf{y})|}\}$  with  $n(\mathbf{z}, y) = n_y$  by picking exactly  $n_y - 1$   $z_i$ s equaling to  $y$  for every  $y \in I(\mathbf{y}|\mathbf{x})$ . All of this is summarized in the main theorem below:

**Theorem 4.11.5** *Given populations  $\mathbf{x}$  and  $\mathbf{y}$  of size  $m$ , we have  $\mathbf{x} \supseteq \mathbf{y}$  with respect to fitness-proportional selection as described in definition 4.2.2, if and only if  $I(\mathbf{y}) \supseteq I(\mathbf{x})$  and*

$$\max_{\{n_y\}_{y \in I(\mathbf{y}|\mathbf{x})} \in Q(\mathbf{y}|\mathbf{x})} \prod_{y \in I(\mathbf{y}|\mathbf{x})} (n_y)^{\kappa(\mathbf{y}|\mathbf{x}, y)} \leq \frac{\prod_{x \in I(\mathbf{x}|\mathbf{y})} (f(x))^{\kappa(\mathbf{x}|\mathbf{y}, x)}}{\prod_{y \in I(\mathbf{y}|\mathbf{x})} (f(y))^{\kappa(\mathbf{y}|\mathbf{x}, y)}}$$

where

$$Q(\mathbf{y}|\mathbf{x}) = \{\{n_y\}_{y \in I(\mathbf{y}|\mathbf{x})} \mid \forall y \in I(\mathbf{y}|\mathbf{x}) n_y \in \mathbb{N}, \sum_{y \in I(\mathbf{y}|\mathbf{x})} n_y = |I(\mathbf{y}|\mathbf{x})| + m - |I(\mathbf{y})|\}.$$

Below we illustrate theorem 4.11.5 with a few simple examples:

**Example 4.11.3** *Continuing with examples 4.11.1 and 4.11.2, notice that we do have  $I(\mathbf{y}) \supsetneq I(\mathbf{x})$  and  $|I(\mathbf{y}|\mathbf{x})| + m - |I(\mathbf{y})| = 2 + 9 - 4 = 7$ . so according to theorem 4.11.5*

we have  $\mathbf{x} \triangleright \mathbf{y}$  if and only if

$$\frac{f(a)^3}{f(c)^2 \cdot f(d)} \geq \max_{n_c+n_d=7, n_c \geq 1 \text{ and } n_d \geq 1} n_c^2 \cdot n_d.$$

There are 6 possible pairs  $(n_c, n_d)$  over which we want to maximize the product  $n_c^2 \cdot n_d$ . These are (1, 6), (2, 5), (3, 4), (4, 3), (5, 2) and (6, 1). Moreover, by symmetry, since the power of the coefficient  $n_c$  in the product is bigger than that of  $n_d$  we only have to cheque 3 of these pairs: (4, 3), (5, 2) and (6, 1). The corresponding products are  $4^2 \cdot 3 = 48$ ,  $5^2 \cdot 2 = 50$  and  $6^2 \cdot 1 = 36$ . Then biggest one among these is 50 and so we deduce that  $\mathbf{x} \triangleright \mathbf{y}$  if and only if  $\frac{f(a)^3}{f(c)^2 \cdot f(d)} \geq 50$ .

**Example 4.11.4** Suppose  $\mathbf{y} = \{y_1, y_2, \dots, y_m\}$  with  $y_i \neq y_j$  for  $i \neq j$  (i.e.  $\mathbf{y}$  is a population consisting of distinct individuals). Now let  $\mathbf{x} = \{x_1, x_2, \dots, x_m\}$  with  $I(\mathbf{y}) \supseteq I(\mathbf{x})$ . Notice that in this example  $I(\mathbf{y}|\mathbf{x}) = I(\mathbf{y}) - I(\mathbf{x})$  while  $I(\mathbf{x}|\mathbf{y}) = \{x \mid \text{there is more than one } i \text{ s. t. } x = x_i\}$ . Moreover,  $\forall x \in I(\mathbf{x}|\mathbf{y})$  we have  $\kappa(\mathbf{x}|\mathbf{y}, x) = n(\mathbf{x}, x) - 1$  (since  $n(\mathbf{y}, x) = 1$ ) and  $\forall y \in I(\mathbf{y}|\mathbf{x})$  we have  $n(\mathbf{y}, y) = 1$  so that  $0 < \kappa(\mathbf{y}|\mathbf{x}, y) \leq n(\mathbf{y}, y)$  and we have  $\kappa(\mathbf{y}|\mathbf{x}, y) = 1$ . Finally, observe that the set  $Q(\mathbf{y}|\mathbf{x}) = \{\{1, 1, 1, \dots, 1\}\}$  since  $|I(y)| = m$  and we must have  $\sum_{y \in I(\mathbf{y}|\mathbf{x})} n_y = |I(\mathbf{y}|\mathbf{x})| + m - |I(y)| = |I(\mathbf{y}|\mathbf{x})|$  and  $n_y \geq 1$  which forces  $n_y = 1 \forall y \in I(\mathbf{y}|\mathbf{x})$ . According to theorem 4.11.5 we have  $\mathbf{x} \triangleright \mathbf{y}$  if and only if  $\frac{\prod_{x \in I(\mathbf{x}|\mathbf{y})} (f(x))^{\kappa(\mathbf{x}|\mathbf{y}, x)}}{\prod_{y \in I(\mathbf{y}|\mathbf{x})} (f(y))^{\kappa(\mathbf{y}|\mathbf{x}, y)}} \geq 1$  if and only if

$$\prod_{x \in I(\mathbf{x}|\mathbf{y})} (f(x))^{n(\mathbf{x}, x) - 1} \geq \prod_{y \in I(\mathbf{y}|\mathbf{x})} f(y).$$

**Example 4.11.5** Continuing with example 4.11.4, suppose, in addition, that there is exactly one individual in  $\mathbf{x}$  which occurs more than once in this population. That is, without loss of generality, let  $\mathbf{x} = (y_1, y_2, \dots, y_k, y_m, y_m, \dots, y_m)$  and  $\mathbf{y} = (y_1, y_2, \dots, y_k, y_{k+1}, \dots, y_m)$  where  $y_i \neq y_j$  for  $i \neq j$  and  $k < m$ . This is a special case of example 4.11.4 where  $I(\mathbf{x}|\mathbf{y}) = \{y_m\}$  and  $I(\mathbf{y}|\mathbf{x}) = \{y_{k+1}, y_{k+2}, \dots, y_{m-1}\}$ . According to the conclusion of example 4.11.4 we have  $\mathbf{x} \triangleright \mathbf{y}$  if and only if

$$(f(y_m))^{m-(k+1)} \geq \prod_{i=1}^{m-(k+1)} f(y_{k+i}) \text{ if and only if } f(y_m) \geq \sqrt[m-(k+1)]{\prod_{i=1}^{m-(k+1)} f(y_{k+i})}$$

which, in words, says that  $\mathbf{x} \triangleright \mathbf{y}$  if and only if the fitness of the unique repeated individual of  $\mathbf{x}$  is at least as large as the geometric mean of the fitness of all the individuals in  $\mathbf{y}$  which do not occur in  $\mathbf{x}$ . It is also worth pointing out that even if the inequality above is an equality we still have  $\mathbf{x} \triangleright \mathbf{y}$  since  $I(\mathbf{y}) \supsetneq I(\mathbf{x})$  and so  $S(\mathbf{x}) \supsetneq S(\mathbf{y})$ . In particular, even if the fitness function is flat, the relation  $\triangleright \neq \emptyset$  in case one uses fitness-proportional selection.

**Example 4.11.6** Now consider an “opposite extreme” to example 4.11.4 (in the sense of the diversity of elements in the population) where  $I(\mathbf{y}) = I(\mathbf{x}) = \{x, y\}$ . let’s say  $n(\mathbf{x}, x) = k$  (which implies that  $n(\mathbf{x}, y) = m - k$ ) and  $n(\mathbf{y}, x) = l < k$

(hence  $n(\mathbf{y}, y) = m - l$ ). It follows then that  $I(\mathbf{x}|\mathbf{y}) = \{x\}$  and  $I(\mathbf{y}|\mathbf{x}) = \{y\}$ . Moreover,  $\kappa(\mathbf{y}|\mathbf{x}, y) = \kappa(\mathbf{x}|\mathbf{y}, x) = k - l$ . Since  $|I(\mathbf{y}|\mathbf{x})| = |\{y\}|$ , it follows that  $Q(\mathbf{y}|\mathbf{x}) = \{\{k - l\}\}$  and which makes the maximization procedure trivial. According to theorem 4.11.5, we have  $\mathbf{x} \succeq \mathbf{y}$  if and only if  $(k - l)^{k-l} \leq \frac{f(x)^{k-l}}{f(y)^{k-l}}$  if and only if  $f(x) \geq (k - l) \cdot f(y)$ .

Theorem 4.11.5 tells us that in order to cheque if  $\mathbf{x} \succeq \mathbf{y}$  with respect to fitness-proportional selection we need to solve an integer optimization problem subject to linear constraints. Examples 4.11.4, 4.11.5 and 4.11.6 are particularly simple mainly because the sets  $Q(\mathbf{y}|\mathbf{x})$  were singletons so there was not much choice for the maximizing domain element. Although we do not intend to pursue studying this optimization problem in much detail since it is not the main subject of the current paper, it is worth mentioning that the method of Lagrange multipliers allows us to give an upper bound on the

$$\max_{\{n_y\}_{y \in I(\mathbf{y}|\mathbf{x})} \in Q(\mathbf{y}|\mathbf{x})} \prod_{y \in I(\mathbf{y}|\mathbf{x})} (n_y)^{\kappa(\mathbf{y}|\mathbf{x}, y)}$$

by letting  $n_y$ 's range over positive real numbers subject to the linear constraint  $\sum_{y \in I(\mathbf{y}|\mathbf{x})} n_y = |I(\mathbf{y}|\mathbf{x})| + m - |I(y)|$ . Moreover, if one wants an exact solution then the method allows to narrow down the choice of suitable integer sequences significantly: Indeed, according to the method of Lagrange multipliers, if any local maximum of a differentiable function  $f(\vec{n})$  on an open set  $D \subseteq \mathbb{R}^n$  subject to the constraint  $g(\vec{n}) = c$  where  $g$  is another differentiable function on  $D$  is achieved at a point  $q$ , then we must have  $\nabla f(q) = \lambda \cdot \nabla g(q)$  where  $\nabla$  denotes the gradient (derivative of a real-valued function) operator and  $\lambda$  is some constant proportionality coefficient (in other words, the gradients of  $f$  and  $g$  evaluated at the point  $q$  must be collinear vectors). In our case  $f, g : U \subseteq \mathbb{R}^{I(\mathbf{y}|\mathbf{x})} \rightarrow \mathbb{R}$  where  $U = \{\vec{n} \mid \vec{n} \in \mathbb{R}^{I(\mathbf{y}|\mathbf{x})} \text{ and } n_y \geq 0\}$  are defined according to the following formulas:  $f(\vec{n}) = \prod_{y \in I(\mathbf{y}|\mathbf{x})} (n_y)^{\kappa(\mathbf{y}|\mathbf{x}, y)}$  and  $g(\vec{n}) = \sum_{y \in I(\mathbf{y}|\mathbf{x})} n_y$ . Our goal is to maximize  $f$  subject to the constraint  $g(\vec{n}) = |I(\mathbf{y}|\mathbf{x})| + m - |I(y)|$  where  $\vec{n} = \{n_y\}_{y \in I(\mathbf{y}|\mathbf{x})}$ . Clearly  $\forall n_y$  we have  $\frac{\partial g}{\partial n_y} = 1$  and so the condition  $\nabla f(q) = \lambda \cdot \nabla g(q)$  boils down to the condition  $\frac{\partial f}{\partial n_u} = \frac{\partial f}{\partial n_v}$  for every  $u$  and  $v \in I(\mathbf{y}|\mathbf{x})$ . For any given  $w \in I(\mathbf{y}|\mathbf{x})$  we have

$$\frac{\partial f}{\partial n_w} = \kappa(\mathbf{y}|\mathbf{x}, w) \cdot (n_w)^{\kappa(\mathbf{y}|\mathbf{x}, w)-1} \cdot \prod_{y \in I(\mathbf{y}|\mathbf{x}), y \neq w} (n_y)^{\kappa(\mathbf{y}|\mathbf{x}, y)}.$$

Therefore, the equation  $\frac{\partial f}{\partial n_u} = \frac{\partial f}{\partial n_v}$  holds for every  $u$  and  $v \in I(\mathbf{y}|\mathbf{x})$  if and only if for every  $u$  and  $v \in I(\mathbf{y}|\mathbf{x})$  we have  $\kappa(\mathbf{y}|\mathbf{x}, u) \cdot n_v = \kappa(\mathbf{y}|\mathbf{x}, v) \cdot n_u$  if and only if  $\frac{n_u}{\kappa(\mathbf{y}|\mathbf{x}, u)} = \frac{n_v}{\kappa(\mathbf{y}|\mathbf{x}, v)}$  for all  $u$  and  $v \in I(\mathbf{y}|\mathbf{x})$ . In other words, the equation  $\frac{\partial f}{\partial n_u} = \frac{\partial f}{\partial n_v}$  holds for every  $u$  and  $v \in I(\mathbf{y}|\mathbf{x})$  if and only if the ratio  $\frac{n_y}{\kappa(\mathbf{y}|\mathbf{x}, y)} = \alpha$  is a constant independent of  $y \in I(\mathbf{y}|\mathbf{x})$ . Moreover, this also gives us  $n_u = \frac{n_v}{\kappa(\mathbf{y}|\mathbf{x}, v)} \cdot \kappa(\mathbf{y}|\mathbf{x}, u) = \alpha \cdot \kappa(\mathbf{y}|\mathbf{x}, u)$  for every  $u \in I(\mathbf{y}|\mathbf{x})$  and, according to the constraint, we also have

$$\sum_{y \in I(\mathbf{y}|\mathbf{x})} n_y = \alpha \cdot \sum_{y \in I(\mathbf{y}|\mathbf{x})} \kappa(\mathbf{y}|\mathbf{x}, y) = |I(\mathbf{y}|\mathbf{x})| + m - |I(y)|$$

which gives

$$\alpha = \frac{|I(\mathbf{y}|\mathbf{x})| + m - |I(y)|}{\sum_{y \in I(\mathbf{y}|\mathbf{x})} \kappa(\mathbf{y}|\mathbf{x}, y)}.$$

Notice that the point  $\vec{q}$  with coordinates  $n_u = \frac{|I(\mathbf{y}|\mathbf{x})| + m - |I(y)|}{\sum_{y \in I(\mathbf{y}|\mathbf{x})} \kappa(\mathbf{y}|\mathbf{x}, y)} \cdot \kappa(\mathbf{y}|\mathbf{x}, y)$  is the unique point which satisfies  $\nabla f(\vec{q}) = \lambda \cdot \nabla g(\vec{q})$ . We argue that this point must be the global maximum of the function  $f$  on the domain  $D = \{n_y \mid n_y \geq 0\} \cap g^{-1}(\{|I(\mathbf{y}|\mathbf{x})| + m - |I(y)|\})$  which is a closed and bounded subset of  $\mathbb{R}^{I(\mathbf{y}|\mathbf{x})}$  and, hence, is compact. Clearly the function  $f$  is continuous on  $D$  and, since  $D$  is compact it must achieve a minimum and a maximum on  $D$ . The only interesting case to consider is when  $|I(\mathbf{y}|\mathbf{x})| > 1$  (indeed, if  $|I(\mathbf{y}|\mathbf{x})| = 1$ , then  $D$  is a singleton set whose only point is  $\vec{q}$  so that it is trivially a global maximum). If maximum of  $f$  was not the point  $\vec{q}$  then it must be the point on the boundary of  $D$  (since it is the only interior point satisfying  $\nabla f(\vec{q}) = \lambda \cdot \nabla g(\vec{q})$ ). But every boundary point of  $D$  has at least one zero coordinate so that  $f(\vec{y}) = 0$  for every boundary point  $\vec{y}$  of  $D$ . On the other hand  $f(\vec{q}) > 0$ . Thus we deduce that  $f$  achieves a global maximum at the point  $\vec{q}$ . We then have the following sufficient condition for  $\mathbf{x} \succeq \mathbf{y}$ :

**Theorem 4.11.6** *Given populations  $\mathbf{x}$  and  $\mathbf{y}$  of size  $m$ , we have  $\mathbf{x} \succeq \mathbf{y}$  with respect to fitness-proportional selection as described in definition 4.2.2, if  $I(\mathbf{y}) \supseteq I(\mathbf{x})$  and*

$$\begin{aligned} & \frac{\prod_{x \in I(\mathbf{x}|\mathbf{y})} (f(x))^{\kappa(\mathbf{x}|\mathbf{y}, x)}}{\prod_{y \in I(\mathbf{y}|\mathbf{x})} (f(y))^{\kappa(\mathbf{y}|\mathbf{x}, y)}} \geq \\ & \geq \left( \frac{|I(\mathbf{y}|\mathbf{x})| + m - |I(y)|}{\sum_{y \in I(\mathbf{y}|\mathbf{x})} \kappa(\mathbf{y}|\mathbf{x}, y)} \right)^{\sum_{y \in I(\mathbf{y}|\mathbf{x})} \kappa(\mathbf{y}|\mathbf{x}, y)} \cdot \prod_{y \in I(\mathbf{y}|\mathbf{x})} \kappa(\mathbf{y}|\mathbf{x}, y)^{\kappa(\mathbf{y}|\mathbf{x}, y)} \end{aligned}$$

**Example 4.11.7** *Continuing with examples 4.11.1, 4.11.2 and 4.11.3, according to corollary 4.11.6 we have  $\mathbf{x} \succ \mathbf{y}$  (recall that we do have  $I(\mathbf{y}) \supseteq I(\mathbf{x})$ ) if*

$$\begin{aligned} & \frac{f(a)^3}{f(c)^2 \cdot f(d)} \geq \\ & \geq \left( \frac{|I(\mathbf{y}|\mathbf{x})| + m - |I(y)|}{\kappa(\mathbf{y}|\mathbf{x}, c) + \kappa(\mathbf{y}|\mathbf{x}, d)} \right)^{\kappa(\mathbf{y}|\mathbf{x}, c) + \kappa(\mathbf{y}|\mathbf{x}, d)} \cdot \kappa(\mathbf{y}|\mathbf{x}, c)^{\kappa(\mathbf{y}|\mathbf{x}, c)} \cdot \kappa(\mathbf{y}|\mathbf{x}, d)^{\kappa(\mathbf{y}|\mathbf{x}, d)} = \\ & = \left( \frac{7}{3} \right)^3 \cdot 2^2 \cdot 1^1 = 50. (814). \end{aligned}$$

Notice that the bound is only slightly larger than the exact one given in example 4.11.3. Moreover, although corollary 4.11.6 itself only provides a numerical bound, the method of Lagrange multipliers which was used to establish corollary 4.11.6, suggests how one can narrow down the search for the optimizing choice of coefficients by considering only the pairs  $(n_c, n_d)$  with integer coordinates which are closest to the point with coordinates  $x_u = \frac{|I(\mathbf{y}|\mathbf{x})| + m - |I(y)|}{\sum_{y \in I(\mathbf{y}|\mathbf{x})} \kappa(\mathbf{y}|\mathbf{x}, y)} \cdot \kappa(\mathbf{y}|\mathbf{x}, u)$  in “every direction”. In our specific example, this point is  $(\frac{7}{3} \cdot 2, \frac{7}{3} \cdot 1) = (4\frac{2}{3}, 2\frac{1}{3})$  and so the only potential candidates are  $(4, 3)$  and  $(5, 2)$ . We saw in 4.11.3 that the point  $(5, 2)$  is the winner. Of course, this “narrowing down” procedure is particularly useful for the cases when  $|I(\mathbf{y}|\mathbf{x})|$  is a large number.

## 4.12 What Can Be Said when the Last Elementary Step is Mutation?

Although not nearly as much can be said when the last elementary step is mutation, the following result is a rather general “anti-communism” theorem. It should be noted that a much stronger and more informative result which depends on the assumption that crossover is “pure” in the sense of [14] (meaning that identical pair of parents produce a pair of the same identical children) shall be established in a sequel paper.

**Theorem 4.12.1** *Let  $\{q_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}$  and  $\{p_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}$  denote Markov transition matrices on a finite set  $\mathcal{X}$ . Suppose  $\{q_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}$  is non-annihilating in the sense of definition 4.10.3. Also let  $\{\{m_{\mathbf{x}\mathbf{y}}^\delta\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}} \mid 0 < \delta < 1\}$  denote an indexed family of Markov transition matrices such that for every  $\epsilon > 0$  there exists  $r > 0$  such that for all  $\delta < r$  we have  $\|\{m_{\mathbf{x}\mathbf{y}}^\delta\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}} - I\| < \epsilon$  for some norm on the finite dimensional vector space of  $|\mathcal{X}| \times |\mathcal{X}|$  matrices<sup>13</sup>. Suppose also that for all  $\delta > 0$  with  $\delta < 1$  the composed Markov chain  $\mathcal{M}(\delta) = \{m_{\mathbf{x}\mathbf{y}}^\delta\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}} \cdot \{p_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}} \cdot \{q_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}$  is irreducible. Let  $\triangleright$  denote the relation associated with the Markov transition matrix  $\{p_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}$  (see definition 4.10.2). Finally, let  $\pi_\delta$  denote the unique stationary distribution of the Markov chain  $\mathcal{M}(\delta)$ . Then, for all small enough  $\delta$ , either there exists a state  $\mathbf{z} \in \mathcal{X}$  such that  $\pi_\delta(\mathbf{z}) < \frac{1}{|\mathcal{X}|}$  or, whenever  $\mathbf{x} \triangleright \mathbf{y}$ , we also have  $\pi_\delta(\mathbf{x}) > \pi_\delta(\mathbf{y})$ . In particular, as long as  $\triangleright \neq \emptyset$ , the stationary distribution of the Markov chain determined by the transition matrix  $\mathcal{M}(\delta)$  is never uniform for all sufficiently small “mutation rates”  $\delta$ .*

*Proof:* Denote by  $\Lambda = \{\{\lambda_{\mathbf{z}}\}_{\mathbf{z} \in \mathcal{X}} \mid \sum_{\mathbf{z} \in \mathcal{X}} \lambda_{\mathbf{z}} = 1, \lambda_{\mathbf{z}} \geq 0\}$  the probability simplex and let  $\Lambda_{\frac{1}{|\mathcal{X}|}} = \{\{\lambda_{\mathbf{z}}\}_{\mathbf{z} \in \mathcal{X}} \mid \sum_{\mathbf{z} \in \mathcal{X}} \lambda_{\mathbf{z}} = 1, \lambda_{\mathbf{z}} \geq \frac{1}{2|\mathcal{X}|}\}$ . For any given  $\mathbf{x} \triangleright \mathbf{y} \in \mathcal{X}$  consider a function  $f_{\mathbf{x},\mathbf{y}} : \Lambda_{\frac{1}{|\mathcal{X}|}} \rightarrow \mathbb{R}$  which sends a given  $\lambda \in \Lambda_{\frac{1}{|\mathcal{X}|}}$  to the number

$$\{p_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}} \cdot \{q_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}(\lambda)(\mathbf{x}) - \{p_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}} \cdot \{q_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}(\lambda)(\mathbf{y}) > 0$$

thanks to proposition 4.10.5. From basic point-set topology we know that the set  $\Lambda_{\frac{1}{|\mathcal{X}|}}$  is a compact topological space (it is a closed and bounded subset of  $\mathbb{R}^{|\mathcal{X}|}$  with  $|\mathcal{X}| < \infty$ ) and, moreover, the function  $f_{\mathbf{x},\mathbf{y}}$  is continuous (it is a restriction of a linear map). It follows then that the function  $f_{\mathbf{x},\mathbf{y}}$  achieves a minimum,  $\min(f_{\mathbf{x},\mathbf{y}})$ , on  $\Lambda_{\frac{1}{|\mathcal{X}|}}$ . Thanks to proposition 4.10.5 this minimum must be a positive number since the matrix  $\{q_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}$  is non-annihilating and every  $\lambda \in \Lambda_{\frac{1}{|\mathcal{X}|}}$  has the property that  $\lambda(\mathbf{z}) \geq \frac{1}{2|\mathcal{X}|} > 0$  for every  $\mathbf{z} \in \mathcal{X}$ . We now conclude that

$$\alpha = \min\{\min\{f_{\mathbf{x},\mathbf{y}}(\lambda) \mid \lambda \in \Lambda_{\frac{1}{|\mathcal{X}|}}\} \mid \mathbf{x} \triangleright \mathbf{y} \in \mathcal{X}\} > 0.$$

Now choose  $r > 0$  small enough so that whenever  $0 < \delta < r$  we have

$$\|\{m_{\mathbf{x}\mathbf{y}}^\delta\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}} - I\|_{\text{op}} < \min\{\frac{\alpha}{3}, \frac{1}{3|\mathcal{X}|}\}.$$

<sup>13</sup>It is a fact that all the norms on finite-dimensional vector spaces are equivalent. It is then irrelevant which norm we consider. For practical applications it is convenient to use  $\|\{a_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}\|_{\max} = \max\{|a_{\mathbf{x}\mathbf{y}}| \mid \mathbf{x}, \mathbf{y} \in \mathcal{X}\}$ . For the purpose of proving the theorem it seems most convenient to use the operator norm defined as  $\|\{a_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}\|_{\text{op}} = \sup\{\|\{a_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}(\vec{v})\| \mid \|\vec{v}\| = 1\}$  where  $\|(v_1, v_2, \dots, v_{|\mathcal{X}|})\| = \sum_{i=1}^{|\mathcal{X}|} |v_i|$ .

Choose any  $\delta$  satisfying  $0 < \delta < r$ . Now there are exactly two mutually exclusive and exhaustive cases:

**Case 1:**  $\exists n \in \mathbb{N}$  such that  $\{p_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \{q_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \mathcal{M}(\delta)^n(\Lambda) \subseteq \Lambda_{\frac{1}{|\mathcal{X}|}}$ .

In this case, let  $\gamma_\delta = \{p_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \{q_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \mathcal{M}(\delta)^n(\pi_\delta)$ . Since  $\pi_\delta$  is the stationary distribution of  $\mathcal{M}(\delta)$  (see the statement of the theorem), it is also the stationary distribution of  $\mathcal{M}(\delta)^{n+1}$  and it follows that

$$\begin{aligned} \{m_{\mathbf{xy}}^\delta\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}}(\gamma_\delta) &= \{m_{\mathbf{xy}}^\delta\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot (\{p_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \{q_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \mathcal{M}(\delta)^n(\pi_\delta)) = \\ &= (\{m_{\mathbf{xy}}^\delta\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \{p_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \{q_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}}) \cdot \mathcal{M}(\delta)^n(\pi_\delta) = \mathcal{M}(\delta)^{n+1}(\pi_\delta) = \pi_\delta \end{aligned}$$

and so

$$\|\gamma_\delta - \pi_\delta\| = \|(I - \{m_{\mathbf{xy}}^\delta\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}})(\gamma_\delta)\| \leq \|I - \{m_{\mathbf{xy}}^\delta\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}}\|_{\text{op}} < \frac{\alpha}{3}.$$

In particular,  $\forall \mathbf{x} \triangleright \mathbf{y} \in \mathcal{X}$  we have

$$|\gamma_\delta(\mathbf{x}) - \pi_\delta(\mathbf{x})| < \frac{\alpha}{3} \text{ and } |\gamma_\delta(\mathbf{y}) - \pi_\delta(\mathbf{y})| < \frac{\alpha}{3}$$

so that

$$\pi_\delta(\mathbf{x}) > \gamma_\delta(\mathbf{x}) - \frac{\alpha}{3} \text{ and } \pi_\delta(\mathbf{y}) < \gamma_\delta(\mathbf{y}) + \frac{\alpha}{3}$$

and, finally,

$$\pi_\delta(\mathbf{x}) - \pi_\delta(\mathbf{y}) > \gamma_\delta(\mathbf{x}) - \frac{\alpha}{3} - (\gamma_\delta(\mathbf{y}) + \frac{\alpha}{3}) = \gamma_\delta(\mathbf{x}) - \gamma_\delta(\mathbf{y}) - \frac{2\alpha}{3} > \frac{\alpha}{3} > 0$$

thanks to the choice of  $\alpha$ , and it follows, in this case, that  $\pi_\delta(\mathbf{x}) > \pi_\delta(\mathbf{y})$ .

**Case 2:**  $\forall n \in \mathbb{N}$  we have  $\{p_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \{q_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \mathcal{M}(\delta)^n(\Lambda) \not\subseteq \Lambda_{\frac{1}{|\mathcal{X}|}}$ .

In this case, first we claim that for every  $n \in \mathbb{N}$  there exists a distribution  $\gamma \in \mathcal{M}(\delta)^{n+1}(\Lambda)$  such that  $\gamma(\mathbf{z}) < \frac{5}{6|\mathcal{X}|}$  for some  $\mathbf{z} \in \mathcal{X}$ . Indeed, the assumption of case 2 says that there exists a distribution  $\lambda \in \{p_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \{q_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \mathcal{M}(\delta)^n(\Lambda)$  such that  $\lambda(\mathbf{z}) < \frac{1}{2|\mathcal{X}|}$ . But then  $\gamma = \{m_{\mathbf{xy}}^\delta\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}}(\lambda)$  is the distribution with the desired property. Indeed, since  $\lambda \in \{p_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \{q_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \mathcal{M}(\delta)^n(\Lambda)$ , it follows that  $\lambda = \{p_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \{q_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \mathcal{M}(\delta)^n(\eta)$  for some distribution  $\eta \in \Lambda$ . But then

$$\begin{aligned} \gamma &= \{m_{\mathbf{xy}}^\delta\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot (\{p_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \{q_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \mathcal{M}(\delta)^n(\eta)) = \\ &= (\{m_{\mathbf{xy}}^\delta\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \{p_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \cdot \{q_{\mathbf{xy}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}}) \cdot \mathcal{M}(\delta)^n(\eta) = \\ &= \mathcal{M}(\delta)^{n+1}(\eta) \in \mathcal{M}(\delta)^{n+1}(\Lambda). \end{aligned}$$

Moreover, since  $\delta < r$  we have

$$\begin{aligned} |\gamma(\mathbf{z}) - \lambda(\mathbf{z})| &\leq \|\gamma - \lambda\| = \|\{m_{\mathbf{xy}}^\delta\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}}(\lambda) - \lambda\| = \\ &= \|(\{m_{\mathbf{xy}}^\delta\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} - I)(\lambda)\| \leq \|(\{m_{\mathbf{xy}}^\delta\}_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} - I)\|_{\text{op}} < \frac{1}{3|\mathcal{X}|}. \end{aligned}$$

But then we also have  $\gamma(\mathbf{z}) \leq \lambda(\mathbf{z}) + \frac{1}{3|\mathcal{X}|} < \frac{1}{2|\mathcal{X}|} + \frac{1}{3|\mathcal{X}|} < \frac{5}{6|\mathcal{X}|}$  as desired. So we deduce every one of the sets  $\mathcal{M}(\delta)^{n+1}(\Lambda)$  contains a point  $\gamma_{n+1}$  with  $\gamma_{n+1}(\mathbf{z}) < \frac{5}{6|\mathcal{X}|}$  for some  $\mathbf{z} \in \mathcal{X}$ . It is well-known from Markov chain theory that the sequence of convex compact sets  $\{\mathcal{M}(\delta)^{n+1}(\Lambda)\}_{n=1}^{\infty}$  is nested ( $\mathcal{M}(\delta)^{n+1}(\Lambda) \supseteq \mathcal{M}(\delta)^{n+2}(\Lambda)$ ) and  $\bigcap_{n=1}^{\infty} \mathcal{M}(\delta)^{n+1}(\Lambda) = \{\pi_{\delta}\}$  where  $\pi_{\delta}$  is the unique stationary distribution of the Markov chain determined by the matrix  $\mathcal{M}(\delta)$ . Also, all the elements of the sequence  $\{\gamma_{n+1}\}_{n=1}^{\infty}$  inside of the compact set  $\Lambda$ , and, hence, the sequence  $\{\gamma_{n+1}\}_{n=1}^{\infty}$  has a convergent subsequence  $\{\gamma_{(n+1)_k}\}_{k=1}^{\infty}$ . But then  $\gamma_{(n+1)_k} \rightarrow \pi_{\delta}$  as  $k \rightarrow \infty$  (since the limit point must lie inside of every one of the compact sets  $\mathcal{M}(\delta)^{n+1}(\Lambda)$  and their intersection consists of a single point  $\pi_{\delta}$ ). Moreover, notice that since  $\mathcal{X}$  is a finite set while  $\{\gamma_{(n+1)_k}\}_{k=1}^{\infty}$  is an infinite sequence, according to the ‘‘pigeonhole principle’’ it follows that  $\exists \mathbf{z} \in \mathcal{X}$  such that infinitely many elements of the subsequence  $\{\gamma_{(n+1)_k}\}_{k=1}^{\infty}$  have the property that  $\{\gamma_{(n+1)_k}\}_{k=1}^{\infty}(\mathbf{z}) < \frac{5}{6|\mathcal{X}|}$ . In other words,  $\exists \mathbf{z} \in \mathcal{X}$  for which we can extract a subsequence  $\{\gamma_{(n+1)_{k_s}}\}_{s=1}^{\infty}$  of the convergent sequence  $\{\gamma_{(n+1)_k}\}_{k=1}^{\infty}$  with the property that  $\gamma_{(n+1)_{k_s}}(\mathbf{z}) < \frac{5}{6|\mathcal{X}|}$ . In particular,  $\gamma_{(n+1)_{k_s}}(\mathbf{z}) \rightarrow \pi_{\delta}(\mathbf{z})$  as  $s \rightarrow \infty$  and it follows that  $\pi_{\delta}(\mathbf{z}) \leq \frac{5}{6|\mathcal{X}|} < \frac{1}{|\mathcal{X}|}$  which is what we were after.  $\square$

When applying theorem 4.12.1 we have in mind that  $\{q_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}$  is the Markov transition matrix corresponding to recombination (i. e. a sub-algorithm determined by a single elementary step of type 2: see definitions 4.3.1 and 4.2.5),  $\{p_{\mathbf{x}\mathbf{y}}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}$  is the Markov transition matrix corresponding to selection (i. e. a sub-algorithm determined by a single elementary step of type 1: see definition 4.2.2) and  $\{m_{\mathbf{x}\mathbf{y}}^{\delta}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}$  is the Markov transition matrix corresponding to mutation with some ‘‘rate’’  $\delta$ . For the purpose of the current section, thanks to the generality of theorem 4.12.1, it is sufficient to assume only that  $m_{\mathbf{x}\mathbf{y}}^{\delta} > 0 \forall \mathbf{x}, \mathbf{y}$  and that  $\max(\{m_{\mathbf{x}\mathbf{y}}^{\delta} \mid \mathbf{x} \neq \mathbf{y} \in \mathcal{X}\}) \rightarrow 0$  as  $\delta \rightarrow 0$ . The following proposition tells us when mutation determined by the reproduction 1-tuple  $(\Omega, \mathcal{M}, p)$  satisfies conditions of theorem 4.12.1:

**Definition 4.12.1** *An ergodic family of mutations is an indexed family of ergodic mutation 1-tuples (see definition 4.9.1) of the form  $\{(\Omega, \mathcal{M}, p_{\delta})\}_{\delta \in (0,1)}$  where  $p_{\delta}(1_{\Omega}) \geq 1 - \delta$ .*

**Proposition 4.12.2** *Suppose  $\{(\Omega, \mathcal{M}, p_{\delta})\}_{\delta \in (0,1)}$  is an ergodic family of mutations as in definition 4.12.1. Then  $\forall \delta \in (0,1)$  the Markov transition matrix  $\{m_{\mathbf{x}\mathbf{y}}^{\delta}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}$  associated to the sub-algorithm determined by the mutation 1-tuple  $(\Omega, \mathcal{M}, p_{\delta})$  has the property that  $\|I - \{m_{\mathbf{x}\mathbf{y}}^{\delta}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}\| \rightarrow 0$  as  $\delta \rightarrow 0$ .*

*Proof:* Notice that  $\|I - \{m_{\mathbf{x}\mathbf{y}}^{\delta}\}_{\mathbf{x},\mathbf{y} \in \mathcal{X}}\| = \max\{m_{\mathbf{x}\mathbf{y}}^{\delta} \mid \mathbf{x} \neq \mathbf{y}\}$  and so it suffices to show that for every  $\mathbf{x} \neq \mathbf{y}$  we have  $m_{\mathbf{x}\mathbf{y}}^{\delta} \rightarrow 0$  as  $\delta \rightarrow 0$  (since the state space is finite). Notice also that  $\forall \mathbf{x} \neq \mathbf{y}$  we have  $0 < m_{\mathbf{x}\mathbf{y}}^{\delta} < 1 - m_{\mathbf{x}\mathbf{x}}^{\delta}$ . It now suffices to show only that  $\forall \mathbf{x} \in \mathcal{X} = \Omega^m$  we have  $1 - m_{\mathbf{x}\mathbf{x}}^{\delta} \rightarrow 0$  as  $\delta \rightarrow 0$ , or, equivalently, that  $\forall \mathbf{x} \in \mathcal{X} = \Omega^m$  we have  $m_{\mathbf{x}\mathbf{x}}^{\delta} \rightarrow 1$  as  $\delta \rightarrow 0$ . If we write  $\mathbf{x} = (x_1, x_2, \dots, x_m) \in \Omega^m$  then, since  $1_{\Omega}(x_i) = x_i$  we see that  $1 \geq m_{\mathbf{x}\mathbf{x}}^{\delta} \geq (p_{\delta}(1_{\Omega}))^m \geq (1 - \delta)^m \rightarrow 1$  as  $\delta \rightarrow 0$  and the desired conclusion follows.  $\square$



Combining theorem 4.12.1 with the conclusion of example 4.11.5 (saying that  $\triangleright \neq \emptyset$  for fitness-proportional selection) we deduce the following:

**Corollary 4.12.3** *Suppose for every  $0 < \delta < 1$  we are given an evolutionary algorithm  $\mathcal{A}_\delta$  determined by the cycle  $s_1, s_2, s_3^\delta$  where  $s_1$  is any elementary step (but usually it is an elementary step of type 2),  $s_2$  is the elementary step of type 1 (fitness-proportional selection as described in definition 4.2.2) and  $s_3^\delta$  is an elementary step of type 2 determined by an ergodic mutation 1-tuple chosen from an ergodic family of mutations (see definition 4.12.1). Then the Markov chain determined by the algorithm  $\mathcal{A}_\delta$  with state space  $\mathcal{X} = \Omega^m$  is irreducible and, for all small enough  $\delta$ , the unique stationary distribution of this Markov chain is not uniform.*

Corollary 4.12.3 tells us, in particular, that the stationary distribution of the Markov chain associated to an algorithm  $\mathcal{A}$  with the second elementary step being of type 1 (selection) is never uniform, even when the fitness function is flat. It is still reasonable to conjecture though, that in case of flat-fitness selection, under certain symmetry assumptions on recombination and mutation, everyone of the individuals in a given population is equally likely to occur “in the long run” in the sense of definition 4.5.6. Results of this nature (and even stronger) shall be established in the upcoming paper.

## 4.13 Conclusions

In the current paper we applied the methods developed in [5] to obtain a schema-based version of Geiringer’s theorem for non-linear GP with homologous crossover. The result enables us to calculate exactly the limiting distribution of the Markov chain associated with the evolution of a finite (fixed size) population under the action of repeated crossover, or the action of the mixture of crossover and mutation. This is an extension of the results for fixed and variable length strings given in [5] for finite populations.

The main result established in [5] applies only in the absence of selection and only when crossover and mutation are bijective (which is often, but not always the case). In the current paper we established a property of the stationary distribution of the Markov chain for a rather wide class of EAs. More specifically, we introduced a pre-order relation on the state space of a Markov chain which allows us to establish rather general inequalities concerning the stationary distribution of the Markov chain determined by an EA. This pre-order relation depends primarily on selection and not on the other stages determining an EA. In section 4.11 this partial order is completely classified for the case of fitness-proportional selection in section 4.11. More results on this issue, as well as some connection between the infinite and the finite population Geiringer theorems will appear in a forthcoming paper.



# Bibliography

- [1] G Briscoe and P De Wilde. D6.6 high-level design specification of the distributed intelligence system. *Digital Business Ecosystem*, Contract no 507953, 2006.
- [2] S. Coffey. An applied probabilist's guide to genetic algorithms. Master's thesis, The University of Dublin, 1999.
- [3] R. A. Fisher. *The genetical theory of natural selection*. Oxford: Clarendon Press, 1930.
- [4] H. Geiringer. On the probability of linkage in Mendelian heredity. *Annals of Mathematical Statistics*, 15:25–57, 1944.
- [5] B. Mitavskiy and J. Rowe. An extension of geiringer theorem for a wide class of evolutionary algorithms. *Evolutionary Computation*, 14, 2006.
- [6] Heinz Muhlenbein and Thilo Mahnig. Evolutionary algorithms and the Boltzmann distribution. In *Foundations of Genetic Algorithms (FOGA2002)*. Morgan Kaufmann Publishers, 2002.
- [7] Heinz Muhlenbein and Thilo Mahnig. Evolutionary computation and Wright's equation. *Theoretical Computer Science*, 2002.
- [8] Liviu Panait and Sean Luke. Alternative bloat control methods. In *Genetic and Evolutionary Computation – GECCO-2004*, pages 630–64. Springer-Verlag, 2004.
- [9] R. Poli. Hyperschema theory for GP with one-point crossover, building blocks, and some new results in GA theory. In R. Poli and W. Banzhaf *et al*, editors, *Genetic Programming, Proceedings of EuroGP'2000*, pages 163–180. Springer-Verlag, 2000.
- [10] R. Poli. A simple but theoretically-motivated method to control bloat in genetic programming. In *Genetic Programming, Proceedings of the 6th European Conference, EuroGP 2003*, pages 211–223. Springer-Verlag, 2003.
- [11] R. Poli and W. Langdon. On the search properties of different crossover operators in genetic programming. In *Proceedings of the Third Annual genetic programming conference*, pages 293–301, 1998.

- [12] R. Poli, C. Stephens, A. Wright, and J. Rowe. A schema-theory-based extension of geiringer’s theorem for linear GP and variable-length GAs under homologous crossover. In K. De Jong, R. Poli, and J. E. Rowe, editors, *Foundations of Genetic Algorithms*, volume 7, pages 45–62, 2002.
- [13] Y. Rabani, Y. Rabinovich, and A. Sinclair. A computational view of population genetics. In *Annual ACM Symposium on the Theory of Computing*, pages 83–92, 1995.
- [14] N. Radcliffe. The algebra of genetic algorithms. *Annals of Mathematics and Artificial Intelligence*, 10:339–384, 1994.
- [15] L. Schmitt. Theory of genetic algorithms. *Theoretical Computer Science*, 259:1–61, 2001.
- [16] L. Schmitt. Theory of genetic algorithms ii: models for genetic operators over the string-tensor representation of populations and convergence to global optima for arbitrary fitness function under scaling. *Theoretical Computer Science*, 310:181–231, 2004.
- [17] W. Spears. The equilibrium and transient behavior of mutation and recombination. In W. Martin and W. Spears, editors, *Foundations of genetic Algorithms*, volume 6, pages 241–260, 2000.
- [18] STU. Report on evolutionary and distributed fitness environment. *Internal*, 2006.
- [19] M De Tommasi. D15.3 BML framework 2nd release. *Digital Business Ecosystem*, 2005.
- [20] UBHAM. Report on population dynamics for variable-sized structures. DBE Deliverable 8.2.
- [21] UBHAM. Report on the evolution of high-level software components. DBE Deliverable 8.1.
- [22] M. Vose. *The simple genetic algorithm: foundations and theory*. MIT Press, 1999.
- [23] A. H. Wright, J. E. Rowe, R. Poli, and C. R. Stephens. A fixed point analysis of a genepool GA with mutation. In *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO 2002)*. Morgan Kaufmann Publishers, 2002.
- [24] S. Wright. Evolution in mendelian populations. *Genetics*, 16:97–159, 1931.