# Techniques for Inferring Context-Free Lindenmayer Systems With Genetic Algorithm

Jason Bernard[a], Ian McQuillan[a]

[a]*Department of Computer Science, University of Saskatchewan, Saskatoon, Canada*

**Abstract**

Lindenmayer systems (L-systems) are a formal grammar system, where the most notable feature is a set of rewriting rules that are used to replace every symbol in a string in parallel; by repeating this process, a sequence of strings is produced. Some symbols in the strings may be interpreted as instructions for simulation software. Thus, the sequence can be used to model the steps of a process. Currently, creating an L-system for a specific process is done by hand by experts through much effort. The inductive inference problem attempts to infer an L-system from such a sequence of strings generated by an unknown system; this can be thought of as an intermediate step to inferring from a sequence of images. This paper evaluates and analyzes different genetic algorithm encoding schemes and mathematical properties for the L-system inductive inference problem. A new tool, the Plant Model Inference Tool for Deterministic Context-Free L-systems (PMIT-D0L) is implemented based on these techniques. PMIT-D0L is successfully evaluated on 28 known L-systems created by experts with alphabets up to 31 symbols, and PMIT-D0L can successfully infer even the largest of these L-systems in less than a few seconds. It is also evaluated and can correctly infer any system in a larger test set of algorithmically created L-systems with much larger alphabets.

*Keywords:* Lindenmayer systems, Plant modelling, Inductive inference, Genetic Algorithm

## 1. Introduction

Lindenmayer systems (L-systems), introduced in [1], are a bio-inspired grammar system that produces self-similar patterns that appear frequently in nature and especially in plants [2]. L-systems produce a sequence of strings, where a string is obtained by the parallel application of rewriting rules to the previous string. Certain symbols can be interpreted as instructions to create images, and therefore a sequence of strings can describe a temporal process, which can

(a) Fibonacci Bush after 7 generations as produced using vlab [2]

(b) Apple twig with blossoms as produced using vlab [2]

be visually simulated by software such as the "virtual laboratory" (vlab) [3]. Such simulations can incorporate different geometries [2, 4, 5], environmental factors [6, 7], and mechanistic controls [8, 4], and as such are useful for simulating plants. L-systems often consist of small textual descriptions that require little storage compared to real imagery. Certainly also, they have a low cost in currency, time, and labor to simulate *in silico* compared to actually growing a plant, and realistic imagery can be produced with a well-constructed L-system.

Formally, L-systems are described by an ordered tuple $G = (V, \omega, P)$ consisting of an alphabet $V$ (a finite set of allowed symbols), an axiom $\omega$ that is a string over $V$, and a finite set of productions, or rewriting rules, $P$. A deterministic context-free L-system (D0L-system) has rules of the form $A \to x$, where $A \in V$ is called the predecessor, $x$ is a string over V that is called the successor of $A$, with exactly one rule for each $A \in V$ as predecessor. Table 2 gives an example of a specific D0L-system created in [9]. Given a string $\omega_i = A_1 \cdots A_m$ where each $A_i \in V, 1 \le i \le m$, a derivation step $\Rightarrow$ is defined by $A_1 \cdots A_m \Rightarrow x_1 \cdots x_m$ where $A_i \to x_i$ is in $P$, for each $i$, $1 \le i \le m$. When this process is repeated $n - 1$ times starting with the axiom ($\omega = \omega_1 \Rightarrow \omega_2 \Rightarrow \cdots \Rightarrow \omega_n$), the sequence $(\omega_1, \ldots, \omega_n)$ is called the length-$n$ developmental sequence. By treating certain symbols as instructions for positioning and drawing in a 3D space (described in Section 2) temporal processes can be simulated using simulation software. Figures 1a and 1b show some structures built with D0L-systems in this fashion.

A difficult challenge is to determine an L-system that can accurately simulate a specific process, for example, modeling plant growth. In practice, this often involves manual measurements over time, scientific knowledge, and is done by hand by experts [2, 10, 5, 4]. Although this approach has been successful, it does have notable drawbacks. Producing a system manually requires an expert, who are in limited supply, and it does not scale to producing arbitrarily many (perhaps closely related) models. Indeed, the manual process currently has been described as requiring "tedious and intricate handwork" [11] that could be improved if an automatic approach could "infer rules and parameters automatically from real . . . images" [11]. Furthermore, when constructed manually, the more complex plant models require *a priori* knowledge of the underlying mechanics of the plant. In contrast, inferring L-systems automatically may be used to reveal scientific principles of the underlying process, or as stated by

2

Godin and Ferraro, automatic inference "could be further exploited in combination with investigations at a biomolecular level to better understand plant development." [5]

The ultimate goal of this line of research is to automatically determine an L-system from a sequence of plant images over time. There have been simplified variants of this problem that have been attempted thus far in the literature. One approach is to develop a tool as an aide for the expert to reduce the work load [12, 13]. With such approaches, the expert operator guides the tool towards a desirable model. Another approach is to build an automated method to convert a sequence of images directly to an L-system [14]. Yet another approach is to divide the problem into two separate steps, with the first step being an accurate segmentation of the images into a sequence of descriptive strings such that an L-system simulator would draw an approximation of the input images; and the second step is to infer the L-system from the sequence of strings. This second step essentially involves taking a sequence of strings as input, and determining the unknown L-system which could give this sequence as its developmental sequence. This is known as the inductive inference problem and it dates back to early work on L-systems studied from the perspective of decidability [15]. We will focus exclusively on this approach. For the first step, this type of plant image segmentation used on a temporal image sequence has been studied separately, e.g. [16, 17]. However, for inductive inference to be truly crucial in combination with image segmentation, it would need to be both fast and accurate so that it could be expanded to work with realistic complications such as noisy images, imperfect data, and other mechanisms built into L-systems. Existing work on inference and inductive inference of L-systems is described in Section 2.3. Most previous attempts to infer L-systems implemented in the literature have involved only a single string as input rather than a sequence of strings [18, 19]; with a single string, the goal is to find the correct L-system that generated this string at some point during its computation. Certainly, there are fewer L-systems that could possibly generate a given sequence of strings versus a given single string within the sequence. Therefore, it could be substantially easier both computationally and in terms of accuracy to infer the correct L-system from a sequence of strings, thereby motivating the study of this problem.

This paper creates the Plant Model Inference Tool for Deterministic Context-Free L-systems (PMIT-D0L) [20, 21] that aims to be an automated approach to solve the inductive inference problem for D0L-systems. Towards that goal, PMIT-D0L uses a genetic algorithm (GA) to search for an L-system that produces a sequence of strings provided as input. In general, GAs search solution spaces in accordance with the encoding scheme used for the problem, and to-date most existing approaches to L-system inference use similar encoding schemes. This paper presents and analyzes different encoding schemes, both existing and novel, to show which are most effective for inferring L-systems. Additionally, some mathematical properties are used to shrink the solution space.

Between the encoding schemes and the use of mathematical properties based on necessary conditions, PMIT-D0L is able to infer all L-systems in a test suite of 28 previously developed L-systems where the number of letters is up to 31

symbols; whereas, other approaches implemented in the literature are limited to 2 symbols as described in Section 2.3. The best encoding scheme, based on the novel approach of searching through allowable combinations of production *lengths* rather than productions directly, took no longer than 3.192 seconds for each L-system, and it took 0.391 seconds on average. All L-systems were inferred with 100% accuracy with all encoding schemes. This is notable as the GA is being used to essentially learn the simulations from data. In addition, to test how PMIT-D0L works on larger L-systems, we algorithmically create a large set of D0L-systems while varying alphabet size. And indeed, PMIT-D0L was able to infer all D0L-systems tested even with up to 134 letters in at most one minute, which is far larger than any L-system in the literature.

There are many future directions required in order to fully realize automatic inference of L-systems. Although many modern L-systems produced by experts use additional features, especially rules that have parameters (parameterized L-systems [2]), creating an inductive inference procedure for D0L-systems that is both fast and accurate is a big step forward. Firstly, it shows that the problem of inference from sequences of strings (and ultimately, images) has promise in the quest to determine a correct L-system versus other approaches. The additional data of a sequence of strings provides significantly more data than a single string to unambiguously recover a correct L-system in a fast way, and sequence data can be often practically obtained. Secondly, many uses of parameters in parameterized L-system rules behave like context-free L-systems during certain sections of their derivations (e.g. if the parameters are being used to incorporate a timing mechanism [2, 22]). The techniques developed in PMIT-D0L can be used for these sections, and can also be used to detect deviations corresponding to a change in the program via parameter. Therefore, studying D0L-system inference scientifically, and inferring D0L-systems in a fashion which can be extended into rules with parameters, is an important step towards the main long term objective. Indeed, PMIT-D0L provides both a fast and accurate implementation of inductive inference that is necessary for L-system inference from images, and is the first inductive inference implementation to do so.

The remainder of this paper is structured as follows. Section 2 will describe some existing automated approaches for inferring L-systems. Section 3 will discuss the different encoding schemes that can be used with PMIT-D0L, along with techniques for reducing the search space size, etc. Section 4 discusses the data set, performance metrics, and the results of the evaluation of PMIT-D0L. Finally, Section 5 concludes the work and discusses future directions.


## 2. Background

This section describes useful contextual and background information relevant to understanding this paper. It starts with describing some notation used. Since a GA is used as the search mechanism for this work, it contains a brief description of them. The section concludes with a discussion of some existing approaches to L-system inference.

### 2.1. Notation

An alphabet is a finite set of symbols. Given an alphabet $V$, a string (or word) over $V$ is any finite sequence of letters $A_1 A_2 \cdots A_n$, $A_i \in V, 1 \le i \le n$. The set of all words over $V$ is denoted by $V^*$, which contains the empty string denoted by $\lambda$. Given a word $x \in V^*$, $|x|$ is the length of $x$, and $|x|_B$ is the number of $B$'s in $x$, where $B \in V$. Given $V = \{B_1, \ldots, B_k\}$, the Parikh vector of a string $x \in V^*$ is $(|x|_{B_1}, \ldots, |x|_{B_k})$.

Given two words $x, y \in V^*$, $x$ is a subword of $y$ if $y = uxv$, for some $u, v \in V^*$ and in this case $y$ is said to be a superstring of $x$. Also, $x$ is a prefix of $y$ if $y = xv$ for some $v$, and $x$ is a suffix of $y$ if $y = ux$ for some $u$.

Given a D0L-system $G = (V, \omega, P)$ as defined in Section 1, the successor of $A$ is indicated by $succ(A)$. Given a rewriting rule $A \to succ(A)$, and $B \in V$, then the number of symbols $B$ in $succ(A)$ is called the growth of $B$ by $A$, denoted by $M(A, B)$. These values are stored in a $|V| \times |V|$ matrix called the growth matrix $M(G)$ [2]. Commonly, $V$ includes symbols to provide simple graphical instructions to simulation software (such as vlab [3]). One commonly used such instruction set is the "Turtle Graphics" [2]. It is imagined as manipulating a turtle through a 2D or 3D space with a pen on its back. The turtle has a state consisting of its position and orientation. The symbols "F" and "f" move the turtle forward along its current orientation with the pen on or off respectively. In 2D, the symbols "$+$" and "$-$" turn the turtle a predefined number of degrees left or right. In 3D, additional symbols are needed for pitch ("&" down and "^" up), and roll ("\" left and "/" right) [2]. For branching processes, the symbols "[" and "]" are used to start and stop a branch, which is implemented as pushing and popping the turtle's state on a stack and switching to it. It is usually the case that the symbols "[", "]", "$+$", "$-$" have identity productions. There are some instances where "F" may not have an identity production (e.g. some of the variants of "Fractal Plant" [2]). Given a sequence of $n$ words $\rho$ over $V$, $G$ is said to be compatible with $\rho$ if $\rho$ is $G$'s length-$n$ developmental sequence. To differentiate the turtle graphic symbols "$+$", "$-$" from the corresponding mathematical operators $+$ and $-$, the turtle graphics symbols will appear in bold as $\mathbf{+}$ and $\mathbf{-}$.

### 2.2. Background on Genetic Algorithm

The GA is described as follows by Bäck [23]. The GA is an optimization algorithm based on evolutionary principles used to efficiently search $N$-dimensional (usually) bounded spaces. An encoding scheme is applied to convert a problem's solution space into one describable by a virtual genome consisting of $N$ genes. Each gene can be either a binary, integer, or real value and represents, in a problem specific way, an element of the solution to the problem. While there exists several types of value encoding schemes, a literal encoding directly represents an element of the solution to a problem. An example of a literal encoding scheme uses gene values to represent the length of a successor, so a value of 3 indicates a successor length of 3. In contrast, a mapped encoding could instead use a real value from 0 to 1 subdivided into sections that represent the different possible solutions.

In evolutionary biology, increasingly fit offspring can be created over successive generations by intermixing the genes of parents. Similarly, a GA functions by iterating over the selection, crossover, mutation, and survival operators until at least one termination condition (e.g. a time limit) is met [23]. There exist different types of these operators; however, this paper will describe only those used here. The function of the GA is controlled by the parameters: population size ($P$), crossover weight ($C$), and mutation weight ($M$). Prior to the first iteration, a GA first produces an initial population of $P$ random solutions. The selection operator chooses some number of pairs of genomes from the population using a selection technique. One such technique, a roulette wheel, is one where the chance of any option being selected (in this case a genome) is proportional to an associated value (in this case, the genome's fitness). For each pair of genomes, the crossover operator swaps a random selection of genes between them, resulting in $P$ child genomes. The chance for any gene to be swapped is equal to $C$. The mutation operator changes a random selection of genes to a random valid value in each child genome. The chance of any individual gene being mutated is equal to $M$. The child genomes are added to the population, and the population is culled to size $P$, thereby keeping the most fit genomes (elite survival).

With PMIT-D0L, the following changes are made to the standard GA to encourage additional exploration. Although an individual genome may be selected for more than one pair, the same pair may not be selected more than once. If any genome has been modified by neither the crossover operator nor the mutation operator, then one gene is selected and mutated to ensure that at least one change has taken place. Where a mapped encoding is used, it is possible for two different genomes to map to the same solution. To prevent such solutions from dominating the population, genomes that map to the same solution are automatically culled (similarly during initialization, duplicated solutions are not permitted in the population).

*2.3. Existing Automated Approaches to L-system Inference*

Various approaches to L-system inference were surveyed in [24]. Here, only certain works most closely related to PMIT-D0L are described. There are several different broad approaches towards the problem: building by hand [2, 10, 4, 5], algebraic approaches [25, 18], using mathematical properties [18], and search approaches [19].

Inductive inference was studied theoretically (without implementation) by Hermann and Rozenberg [15], and Doucet [25]. In [15], the problem was studied from the perspective of decidability. In [25], Doucet defined a matrix equation to simplify the problem. Let $\rho = (\omega_1, \ldots, \omega_n)$ be input words over an $n-1$ letter alphabet, with the goal of finding a D0L-system that has $\rho$ as its length-$n$ developmental sequence. Let $Y$ be the $(n-1) \times (n-1)$ matrix such that row $i$ is the Parikh vector of $\omega_i$ for $i$, $1 \le i \le n-1$, and let $Z$ be the $(n-1) \times (n-1)$ matrix such that row $i$ is the Parikh vector of $\omega_{i+1}$ for $i$, $1 \le i \le n-1$. From the definition of matrix multiplication (also discussed in both [25, 26]), if a D0L-system $G$ has $\rho$ as its length-$n$ developmental sequence, then $YM(G) = Z$; that

is, the unknown growth matrix must be a solution to the equation $YX = Z$. That said, there can be solutions to this equation that are not growth matrices of D0L-systems generating $\rho$. However, it is possible to search only solutions to this equation rather than arbitrary successor lengths in order to find a D0L-system that generates $\rho$. Indeed, for a given solution to this equation, it is straightforward and efficient to find an L-system with it as growth matrix if it exists, as described below in Section 2.4. Recently, this theoretical algorithm was extended to work for context-sensitive L-systems [26]. We implement a similar approach here as one encoding scheme under the name *row reduced matrix encoding*. A somewhat similar approach [25] was implemented with a tool called LGIN [18] that infers L-systems from only a single string. LGIN looks exhaustively at the successor combinations, extracted from a single string in a developmental sequence that fulfills these equations. Since only a single string is used, the problem they are solving is more difficult than the inductive inference problem we are addressing, and indeed it does not guarantee to find a unique solution. LGIN only provides specific algorithms for one and two symbol alphabets (not including turtle graphic symbols), with larger alphabets described as "immensely complicated" [18] and they are not described algorithmically; however, LGIN was evaluated on six variants of "Fractal Plant" [2] and it was very fast having a peak time to find the L-system of less than one second for 5 of the 6 variants, and four seconds for the remaining variant.

Runqiang et al. [19] investigated inferring an L-system directly from a description of a single image (essentially one string) using a GA. In their approach, they encode each symbol within the successors as a gene. The fitness function attempts to match the candidate system to the input string. Their approach is limited to an alphabet size of 2 symbols and has a maximum combined length of all successors of 14. Their approach is 100% successful for variants of "Fractal Plant" [2] with $|V| = 1$, and a 66% success rate for variants of "Fractal Plant" as in Figure 2 with $|V| = 2$. Although they do not list any timings, their GA converged after a maximum of 97 generations, which suggests a short runtime. The main difference between their work and ours is they do not take a sequence of strings as input. They use an encoding scheme that we call *ordered sequence of symbols* (discussed further in Section 3.2), which we also implement and investigate here (on sequences of strings as input) which we compare to other encoding schemes.

*2.4. Scanning for Successors*

A technique for finding productions based on the successor lengths for every $A \in V$ was previously described in [26], which we call the *scanning process*. This technique is used extensively in this research, and is described as follows. With L-systems, although the symbols are replaced in parallel, the sequence of successors in the new word is unchanged from the sequence of the original symbols; i.e., if $\omega_i = A_1 A_2 \ldots A_m$ with $A_i \in V$, $1 \leq i \leq m$, then $\omega_{i+1} = succ(A_1)succc(A_2) \cdots succ(A_m)$. Consider the case of $A_1$ in $\omega_i$. To find $succ(A_1)$ it is only necessary to know the length of $succ(A_1)$. If $|succ(A_1)|$ is known (or different values for it are tested by searching), then the successor is the first
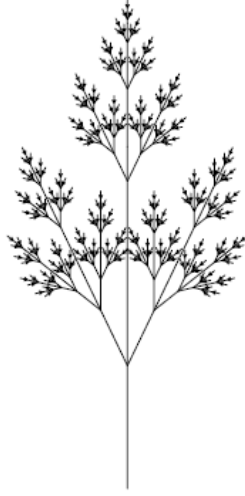
Figure 2: "Fractal Plant" variant #5 [2]

$|succ(A_1)|$ characters of $\omega_{i+1}$. The process of taking the next $l$ symbols (where $l$ is a hypothetical successor length) may be repeated for every new instance of a symbol encountered while scanning each word of a developmental sequence until every successor has been found, as described in [26]. With this fast algorithm, the goal is therefore to find a list of successor lengths that results in an L-system compatible with a developmental sequence; however, this may be done in a few ways. Most directly, a list of successor length combinations may be found by searching. Somewhat indirectly, the growth values for every $A, B \in V$ may instead be found, and a successor length for each $A \in V$ computed by summing the growth values for every $B \in V$.

## 3. Methodology

This section describes the design and procedure of PMIT-D0L. First, a high level overview of PMIT-D0L will be described, as this helps to contextualize the remainder of the section. Then, the techniques used to reduce the size of the defined search space are discussed. This is followed by a description of the different encoding schemes used to define and search the space to infer an L-system (the different encoding schemes are evaluated separately to see which work best). The final two sub-sections describe the process used to optimize the control parameters of the GA, and finally, the fitness function and termination conditions.

PMIT-D0L makes two assumptions regarding L-systems to be inferred: no successor is the empty word, and that for branching L-systems, the branching symbols [ and ] are properly nested within each production (this is a common assumption, e.g. [19] as improperly nested branching symbols are biologically meaningless). Most L-systems in the literature satisfy these assumptions.

8

As mentioned in Section 2.1, some symbols, such as the turtle interpretation symbols, are usually created to have identity productions. Hence, a set $C \subset V$ of *constants* is provided as input where it is assumed that all symbols in $C$ have identity productions. This separation is commonly done with inference procedures (e.g. [18] which separates out constants from the rest of the alphabet). In our experiments, we use the turtle interpretation symbols as the elements of $C$. Their known successors can speed up searches. In addition, let $\overline{V} = V - C$ (those symbols with unknown productions).

The pseudocode description of PMIT-D0L is provided in Algorithm 1. This algorithm has been implemented in C++ with CLR extensions. It takes as input a sequence of strings $\rho = (\omega_1, \ldots, \omega_n), \omega_i \in V^*, 1 \leq i \leq n$, an encoding scheme $En$ (the various encoding schemes are described in Section 3.2), and a list of constants $C$. For turtle interpretation, we use the fixed ordering of the elements of $C$: $[,], +, -$, followed by the symbols for yaw, pitch, roll, turn $180°$ (if they are present), then $F, f$. PMIT-D0L will return either a D0L-system compatible with $\rho$ or return that no D0L-system was found (which either means that one does not exist, or the GA terminated without finding a solution, as a GA cannot guarantee to find a solution) that is compatible with $\rho$.

At a high level, the algorithm first attempts to find partial solutions by removing all constants from $\rho$, and then it searches for solutions to this sequence of strings using the GA. For all solutions found by the GA with a fitness of 0, it adds them to a queue $Q$ (initialized in line 1) containing all partial solutions. Then, for each partial solution in the queue, it adds back in the first constant and tries determine, for each occurrence of this constant in $\rho$, which position of this constant in the next string must be produced from it (using a procedure described in detail in Section 3.1.3). Then, from each of these partial solutions, it uses the GA to search for new partial solutions extended from it. It proceeds similarly for the second constant, etc., until it has added back in all constants. It keeps track of the current constant to be added in to each partial solution with an index *currentConstant* into the list of constants $C$ (this value gets enqueued together with the partial solutions). One can see this basic procedure with the while loop on line 4, which removes all constants starting at index *currentConstant* from $\rho$ (line 5); uses the GA to search the search space (line 7); adds all new partial solutions found with a fitness of 0 to $Q$; takes one partial solution at a time from the queue (line 9); and if it has already included all constants, then it is a complete solution and the L-system has been found and is returned; otherwise, it adds in the next constant (line 9). The detailed procedure for adding one constant at a time to setup associations between consecutive strings is described in Section 3.1.3.

The GA is searched in accordance with the desired encoding scheme $En$ (described in Section 3.2). Before searching the search space with the GA, two approaches are used to reduce the size of the search space (executed in line 3 before the loop, and within the loop in line 6 each time before searching with the GA). As each of the two techniques can lead to reductions using the other technique, the two methods are executed alternatingly in a loop as part of this pseudocode line until there are no changes.

The first reduction technique is done by examining possible successor length combinations and growth matrix values. Reducing the range of possible production lengths is crucial, as is evident from the scanning process described in Section 2.4, which can infer an L-system quickly and correctly from only the successor lengths (and successor lengths can be calculated from growth matrix values). Hence, narrowing down the range of possible successor lengths or growth matrix values can significantly reduce the size of the search space. Thus, the first step is to initialize several programming variables in line 2 that keep track of lower and upper bounds on each of these, which get refined as additional data is scanned and deductions are made. For each $A \in V$, the current lower bound (respectively upper bound) on the successor length of $A$ is denoted by $A_{\min}$ (respectively $A_{\max}$), and $A_{\min}$ is initialized to be 1 as there are no erasing productions. The current upper and lower bounds for the growth of $B$ by $A$ for every pair of $A, B \in V$ is denoted by $(A, B)_{\min}$ and $(A, B)_{\max}$, respectively. Since all letters in $C$ have a known identity production, we initialize for each symbol $T \in C$, $T_{\min} = T_{\max} = 1$,$(T, T)_{\min} = (T, T)_{\max} = 1$ and $(T, A)_{\min} = (T, A)_{\max} = 0$ for every $A \in V, A \neq T$. The technique to reduce the search space by refining these values is described in detail below in Section 3.1.2.

The second reduction technique can be obtained using a concept called *successor fragments*. Given $A \in \overline{V}$, a word $\omega$ is an $A$-subword (respectively, $A$-prefix, $A$-suffix, $A$-superstring) if $\omega$ is a subword (respectively, prefix, suffix, superstring) of $succ(A)$. It is helpful to maintain, for each $A \in \overline{V}$, words that we have deduced must be an $A$-subword, and similarly for $A$-prefix, $A$-suffix, and $A$-superstring. It is evident that if there are multiple $A$-prefixes, then only the longest is of interest as the shorter ones are prefixes of the longest one, and similarly with the suffix. Hence, for each $A \in \overline{V}$, Algorithm 1 keeps strings $Sub_A$, $Pre_A$, and $Suf_A$ which stores the longest known $A$-subword, $A$-prefix, and $A$-suffix respectively, and $Sup_A$ is the shortest known $A$-superstring. Also in line 2, all strings are initialized to be empty. The detailed method to help determine these fragments is given in Section 3.1.1.

### 3.1. Search Space Reduction

In this section, the techniques that are used by PMIT-D0L to reduce the size of the solution space (with any of the encoding schemes) using mathematical properties of D0L-systems will be described. As previously mentioned, these are applied in Algorithm 1 in line 3 and 6 before each GA search, and the two techniques are alternatively executed in a loop until there are no more changes. Being based on necessary conditions guarantees that all valid solutions are in the remaining search space (if there is a D0L-system that can generate the input strings).

### 3.1.1. Refining Successor Relationships

As PMIT-D0L runs, it can determine additional successor fragments, which can help in turn to reduce growth bounds. Certain prefix and suffix fragments

**Data:** A sequence of strings $\rho$ over $V$, an encoding scheme $En$, and a
      list of symbols with known identity productions $C$

**Result:** D0L-system compatible with $\rho$, or that no compatible
      D0L-system is found

**1** $Q \longleftarrow \emptyset$ // initialize queue of solutions;

**2** // initialize state consisting of length and growth bound variables,
    fragment successors, and $currentConstant \longleftarrow 0$;

**3** refine length and growth bounds, and fragment successors;

**4 while** $currentConstant \leq |C|$ **do**

**5**      $\rho' \longleftarrow$ remove constants starting at $currentConstant$ from $\rho$;

**6**      refine length and growth bounds, and fragment successors;

**7**      search *space* based on $En$ (Sec. 3.2) and add each solution found
        with fitness 0 together with state to $Q$;

**8**      **if** $Q$ *is not empty* **then**

**9**          dequeue *solution* and state and make state current;

**10**         **if** $currentConstant = |C|$ **then**

**11**             return *solution*;

**12**         **else**

**13**             add letter $C[currentConstant]$ to *solution* (Sec. 3.1.3);

**14**             increase $currentConstant$;

**15**         **end**

**16**      **else**

**17**         return "none found";

**18**      **end**

**19 end**

**Algorithm 1:** High-level D0L-system inference procedure.

can be found for the first and last symbols in each input word by the following process. Consider two words such that $\omega_i \Rightarrow \omega_{i+1}$. It is possible to scan $\omega_i$ from left-to-right until the first symbol from $\overline{V}$ is scanned (say, $A$, where the word scanned is $\alpha A$). Then, in $\omega_{i+1}$, PMIT-D0L skips over the symbols of $C$ in $\alpha$ (since each symbol in $\alpha$ has a known identity production), and the next $A_{\min}$ symbols (the current lower bound for $|succ(A)|$), $\beta$ say, must be an A-prefix fragment. Furthermore, since branching symbols must be well-nested within a successor, if a [ symbol is met, the prefix fragment must also contain all symbols until the matching ] symbol is met. Similarly, an A-superstring fragment can be found by skipping $\alpha$, then taking the next $A_{\max}$ symbols from $\omega_{i+1}$ (the upper bound on $|succ(A)|$). If a superstring fragment contains a [ symbol without the matching ] symbol, then it is reduced to the symbol before the unmatched [ symbol. The lower and upper bounds $(A, B)_{\min}$ and $(A, B)_{\max}$ for each $B \in V$ can be then possibly improved by counting the number of $B$ symbols in any prefix and superstring fragments respectively. For a suffix fragment, the process is identical except the scan goes from right-to-left starting at the end of $\omega_i$.

**Example 1.** *Consider input strings:*

$$\omega_1 = +++A[-FF][+F]BF$$
$$\omega_2 = ++++A[-FF][-FF][+F][+F]BFF.$$

*Assume thus far PMIT-D0L has calculated that $A_{\min} = 2$ and $A_{\max} = 8$. It can scan $\omega_1$ until $A$ is found and record that $\alpha = +++$. An $A$-prefix fragment is $+A$ as those are the first two ($A_{\min}$) symbols of $\omega_2$ after skipping $\alpha$. An $A$-superstring fragment is $+A[-FF][$ as those are the first eight ($A_{\max}$) symbols of $\omega_2$ after skipping $\alpha$, which can be reduced to $+A[-FF]$ due to the unbalanced $[$ symbol. By counting within the prefix fragment, lower bounds on the growth for $A$ are $(A,+)_{\min} := 1$ and $(A,A)_{\min} := 1$, while upper bounds can be found from the superstring fragment to be $(A,+)_{\max} := 1, (A,-)_{\max} := 1, (A,A)_{\max} := 1, (A,[)_{\max} := 1, (A,])_{\max} := 1$ and $(A,F)_{\max} := 2$.*

If a symbol has a known successor, then it may be possible under certain circumstances to "line it up" with the symbol(s) it produces. In such circumstances, the derivation is sliced into two parts, and this reduces the possible productions for the symbols in each part. To illustrate the concept, consider the simple example:

$$\omega_1: A+B$$
$$\downarrow$$
$$\omega_2: \underbrace{ABA}_{\text{succ(A)}} + \underbrace{BBB}_{\text{succ(B)}}$$

As $+ \in C$, it must be that $+ \to +$ is a production, and the $+$ symbol in $\omega_1$ may only produce the $+$ in $\omega_2$. This splits the derivation such that everything to the left of the $+$ in $\omega_1$ must produce everything to the left of the $+$ in $\omega_2$, i.e. $A \to ABA$ must be a production. Similarly to the right of the $+$, $B \to BBB$ must be a production.

In practice, it is unusual for a single position of one string to be uniquely associated with a position in the next string, as in the example. More often, any individual position may be associated to multiple positions of the next word. However, a sequence of symbols, each with known successor, may be unique. For example, in the string $A[+B][-B]A[+[-C]][-[+D]]$, the individual symbols $[$, $]$, $+$, and $-$ alone might not uniquely associate to their successors; however, a sequence of symbols such as $][-$ or $]][-[+$ are potentially unique, and as such may be associated to a unique location. To make use of this observation, for each word, a list of possible associations between every position of a symbol in $C$ to positions of the next word is constructed. This list of associations is called a *marker map*. A marker map is constructed based on both individual symbols and sequences of symbols, which are referred to as *candidate markers*. Associating a candidate marker to potential successors takes into account that a number of symbols must be reserved for symbols that appear before and after

the candidate marker. For example, if $\omega_1 = A{+}BC{-}, \omega_2 = A{+}BC{+}C{-}$, then $+$ associates with both $+$'s in $\omega_2$, both are candidate markers. But since $B_{\min} + C_{\min} + -_{\min} \geq 3$ (as there are no empty word successors) the final 3 symbols of $\omega_1$ produce at least the $+C-$ of $\omega_2$. This eliminates the second $+$ in $\omega_2$ as being produced by the $+$ in $\omega_1$, and the $+$ in $\omega_1$ can only be associated to the first $+$. If, following the construction of the marker map, a candidate marker is not uniquely associated with its successor, then it is removed from the marker map.

Once the marker map has been calculated, it can help significantly to improve the length bounds and successor fragments. Consider a derivation $\omega_i \Rightarrow \omega_{i+1}$ expanded as

$$\omega_{i,1} A_1 \omega_{i,2} \cdots A_m \omega_{i,m+1} \Rightarrow \omega_{i+1,1} succ(A_1) \omega_{i+1,2} \cdots succ(A_m) \omega_{i+1,m+1},$$

where each $A_j$, $1 \leq j \leq m$ in $\omega_i$ has already been associated to the annotated successor in $\omega_{i+1}$ forming a marker. It follows that $\omega_{i,j} \Rightarrow \omega_{i+1,j}$ for all $j$, $1 \leq j \leq m + 1$. Indeed, from this, improved successor fragments, growth and length bounds may be found.

### 3.1.2. Refining Growth and Length Bounds

Here the bounds on $(A, B)_{\min}$ and $(A, B)_{\max}$ are improved. As all properties are run in a loop, these bounds are also influenced by successor fragments, and markers as described above, which can result in significantly improved bounds versus just examining what can be determined from Parikh vectors alone. A programming variable for the *accounted for growth* of a symbol $A \in V$ for $2 \leq i \leq n$, denoted as $G_{acc}(i, A)$ is:

$$G_{acc}(i, A) := \sum_{B \in V} (|\omega_{i-1}|_B \cdot (B, A)_{\min}).$$

The *unaccounted for growth* for a symbol $A$, denoted as $G_{ua}(i, A)$, is computed as $G_{ua}(i, A) := |\omega_i|_A - G_{acc}(i, A)$.

The unaccounted for growth can be used to improve the growth bounds. In particular, $(B, A)_{\max}$ is set (if it can be reduced) under the assumption that all unaccounted for $A$ symbols are produced by $B$ symbols. Furthermore, $(B, A)_{\max}$ is set to be the lowest such value computed for any word from 2 to $n$, where $B$ occurs, as any of the $n - 1$ words can be used to improve the maximum. And, $|succ(B)|_A$ must be less than or equal to $(B, A)_{\min}$ plus the additional unaccounted for growth of $A$ divided by the number of $B$ symbols (if there is at least one) in the previous word, as computed by

$$(B, A)_{\max} := \min_{\substack{2 \leq i \leq n, \\ |\omega_{i-1}|_B > 0}} \left( (B, A)_{\min} + \left\lfloor \frac{G_{ua}(i, A)}{|\omega_{i-1}|_B} \right\rfloor \right).$$

Indeed, the accounted for growth of $A$ is always updated whenever values of $(B, A)_{\min}$ change, and the floor function is used since $|succ(B)|_A$ is a non-negative integer. For example, if $\omega_{i-1} = ABA$, $\omega_i = ABABBBABA$, $(A, A)_{\min} =$

1, and $(B,A)_{\min} = 0$, then the accounted for growth of $A$ in $\omega_i$ is computed by $G_{acc}(i,A) = (A,A)_{\min} \cdot |\omega_{i-1}|_A + (B,A)_{\min} \cdot |\omega_{i-1}|_B = 1 \cdot 2 + 0 \cdot 1 = 2$. This leaves two $A$'s in $\omega_i$ unaccounted for. An upper bound on the value of $|succ(A)|_A$ is set when the $A$'s in $\omega_{i-1}$ produce all of the unaccounted for growth in $\omega_i$. So $A$ produces its minimum ($(A,A)_{\min} = 1$) plus the unaccounted for growth of $A$ in $\omega_i$ (2) divided by the number of $A$'s in $\omega_{i-1}$ ($|\omega_{i-1}|_A = 2$), hence $(A,A)_{\max} := 2$. Similarly, $(B,A)_{\max}$ is achieved when only $B$'s produce all unaccounted for growth of $A$; this sets $(B,A)_{\max}$ to $(B,A)_{\min} = 0$ plus the unaccounted for growth (2) divided by the number of $B$'s in $\omega_{i-1}$ (1), which is 2.

Once $(B,A)_{\max}$ has been determined for every $A, B \in V$, the observed words are re-processed to compute possibly improved values for $(B,A)_{\min}$. Indeed for each $(B,A)$, if $x := \sum_{\substack{C \in V \\ C \neq B}} (C,A)_{\max}$, and $x < |\omega_i|_A$, then this means that $|succ(B)|_A$ must be at least $\left\lceil \frac{|\omega_i|_A - x}{|\omega_{i-1}|_B} \right\rceil$, and then $(B,A)_{\min}$ can be set to this value if its bound is improved. For example, if $\omega_{i-1}$ has 2 $A$'s and 1 $B$, and $\omega_i$ has 10 $A$'s, and $(A,A)_{\max} = 4$, then at most two $A$'s produce eight $A$'s, thus one $B$ produces at least two $A$'s (10 total minus 8 produced at most by $A$), and $(B,A)_{\min}$ can be set to 2.

In a similar fashion, the length bounds $A_{\min}$ and $A_{\max}$ can be set using unaccounted for length.

### 3.1.3. Solution Projections

As previously defined, $C \subset V$ where all symbols in $C$ have a known identity production, and $\overline{V} = V - C$. Since a symbol in $\overline{V}$ cannot be produced by a symbol in $C$, in the while loop of Algorithm 1, it is possible to first infer an L-system over reduced alphabets with fewer constants, and add one constant at a time. For example, if $V = \{A, B, C, [, ], +, -\}$ and $C$ is a list of constants $([, ], +, -)$, then one can first find each successor of $A, B, C$ projected to $\{A, B, C\}$. After solving this initial problem, then a series of problems are solved adding in each symbol of $C$ to determine where it belongs in each successor. Note that "[" and "]" are completed together due to the assumption that they are properly nested. Overall, symbol filtering simplifies the inference problem by allowing for an iteration of lining up constants between consecutive words. Although more searches are needed, they are each in a smaller search space.

As described above, PMIT-D0L removes the symbols of $C$ temporarily by projecting $\rho$ onto a reduced alphabet, and then it iteratively re-adds each symbol of $C$ back into the problem one at a time. Let a solution to one of these reduced problems be called a *partial solution*, as it partly describes the final successors. The process for using the partial solutions towards the next partial solution by adding in the next symbol of $C$ is conceptually similar to that used for markers, as positions in pairs of consecutive words will be "lined up", and from there, successor relationships deduced. To describe this, some terminology for this process is provided. For every derivation step $\omega_i \Rightarrow \omega_{i+1}$, every position in $\omega_i$ is scanned, and associated to the possible locations in $\omega_{i+1}$ that it must produce. Let $\omega'_{i+1}$ be equal to $\omega_{i+1}$ but with all symbols of $C$ erased. As the algorithm proceeds, a position $j$ of $\omega_i$ is said to be *certain* if positions

$k, l$ have been determined such that letter $D$ at position $j$ of $\omega_i$ must produce exactly the symbols of $\omega_{i+1}$ between positions $j$ and $k$; and it is *uncertain* otherwise. If a position is certain, this means it has already been determined exactly the positions that this specific $D$ must produce in the next word (which also determines the successor of this letter). Note that if position $j$ is labelled by an element of $C$, then $k = l$ and position $k$ of $\omega_{i+1}$ must be $D$ since $D$ is a constant. This property is purely programmatic, and a position can change from uncertain to certain as the algorithm proceeds. Similarly, a position $j$ of $\omega_i$ labelled by some symbol $D \in \overline{V}$ is said to be *projected-certain* if positions $k, l$ have been determined such that letter $D$ at position $j$ of $\omega_i$ must produce exactly the symbols of $\omega'_{i+1}$ between positions $j$ and $k$; and it is *projected-uncertain* otherwise. This means, a position $j$ is projected-certain if it has been determined exactly which positions of $\omega_{i+1}$ can be produced by this position, ignoring adjacent positions labelled by symbols of $C$.

As we iterate the loop of Algorithm 1, every time we re-add a new symbol $D \in C$ in line 13, this involves scanning every position of each string of $\rho$ to try to determine whether they are certain. It is often possible to determine that certain positions are certain if adjacent positions are certain. We will describe the idea first with an example. Consider $\omega_1$ and $\omega_2$ in Equation 1 to 3, and assume that it has already been determined that the successor of $A$ projected to $\overline{V}$ is $BAB$, and the successor of $B$ projected to $\overline{V}$ is $AB$. That means that, by scanning $\omega_1$ from left-to-right, each position in $\omega_1$ labelled by a symbol of $\overline{V}$ is projected-certain, as the first $A$ must produce the first three symbols of $\omega'_2$, the $B$ must produce the next two symbols, and the second $A$ must produce the last three symbols. This also indicates that the first $A$ must produce $B[+A]B$ (labelled $\alpha$ in Equation 1), but it is not yet clear what adjacent symbols of $C$ can also be produced by this $A$ (so this $A$ is not yet certain). Similarly, the $B$ must also produce $A-[+B]$ of $\omega_2$, but it is unclear as to which adjacent symbols of $C$ are produced by $B$. Lastly, the final $A$ must produce the final $B[+A]B$ of $\omega_2$ but adjacent symbols of $C$ are unclear.

Next, the loop of Algorithm 1 re-adds the [ and ] symbols as the first two constants. As these symbols do not occur in $\omega_1$, there are no new symbols to line up and the procedure does not do anything. The next letter of $C$ it re-adds is the $+$ symbol. The $+$ symbol is uncertain at this stage, as the $+$ could produce either of the two annotated $+$'s in $\omega_2$ as shown in Equation 1; it cannot produce the first $+$ as the first $A$ must map to positions that include this symbol. Even though the $+$ is uncertain at this stage, the word $B[+A]B$ ($\alpha$ in Equation 1) can be declared an $A$-prefix. The next letter re-added is -, and it is concluded that the - is certain, as it has already been determined that the final $A$ maps to the final $B[+A]B$ with perhaps some adjacent symbols of $C$ to the left, and therefore, it cannot map to the second last -. Since - is now certain, this implies that the final $A$ is certain which indicates that the successor of $A$ is $B[+A]B$. This implies that the first $A$ must map to exactly the first $B[+A]B$ and the first $A$ is now certain, which also resolves the certainty of the first $+$, and lastly it implies that the $B$ is certain as well and produces exactly $+A-[+B]$. Thus,

the L-system is determined.

$$\omega_1 : A + B - A$$

$$\omega_2 : \underbrace{B[+A]B}_{\alpha} + + A - [+B] - B[+A]B \tag{1}$$

$$\omega_1 : A + B - A$$

$$\omega_2 : B[+A]B + + \underbrace{A - [+B]}_{\beta} - B[+A]B \tag{2}$$

$$\omega_1 : A + B - A$$

$$\omega_2 : B[+A]B + + A - [+B] - \underbrace{B[+A]B}_{\text{succ(A)}} \tag{3}$$

Hence, every time it executes line 13 of Algorithm 1, it scans each position from left-to-right of each string and assesses certainty of each position in this fashion; and it repeats this process in a loop until there are no more changes, so that it can fully consider how the certainty of positions can change as the certainty of other positions change.

### 3.2. Defining and Searching the Search Space

This section describes the different encoding schemes used in this research, and in some cases existing approaches to using encoding schemes [19, 25, 18], for inferring D0L-systems. Broadly, the encoding schemes can be broken down into three categories: ordered sequence of symbols (OSoS), growth-based, and length-based. The OSoS approaches take the viewpoint that a successor is an ordered sequence of unknown symbols, and so the search space is represented in this fashion. Another approach investigated in this research is to instead attempt to determine successor length combinations as the unknown, as an intermediate step, before determining the actual productions using the scanning process of Section 2.4. Similarly, the growth values may be inferred first, and then simply summed for each $A \in V$ to produce a successor length followed by the scanning process.

### 3.2.1. Ordered Sequence of Symbols Encoding

While the technique of building a search space based on the idea of searching for the symbol in each position of each successor has been previously investigated [19, 13], PMIT-D0L creates this search space with additional requirements.

16

For every $A \in V$, a successor may be encoded as follows. For the remainder of this section, we create a special symbol $\bar{\bar{\emptyset}}$. A number of genes equal to $A_{\max}$ is defined as this is the greatest number of symbols that may exist in the successor. The current values for $Pre_A$ and $Suf_A$ immediately identify a number of genes at the beginning and end equal to the length of the prefix and suffix respectively. For example, if $A_{\max} = 7$, $Pre_A = A$, and $Suf_A = BB$, then the genome would appear as follows A _ _ _ _ B B, where _ represents an unknown symbol, which could be set to $\bar{\bar{\emptyset}}$ if $|succ(A)| < 7$ (implying no symbol of $V$ exists in that position of this successor considered). Each of the genes are permitted to have a real-value from 0 to 1. Growth bounds can reduce the options however. To continue the previous example, if $(A, B)_{\max} = 2$, then since there are already two occurrences of $B$ symbols in the successor, $B$ needs not be a choice for any of the remaining genes. Minimum growth values can be similarly enforced. Further, if $(A, A)_{\min} = 3$, $(A, C)_{\min} = 2$, and the successor is A A _ _ _ B B, then the remaining genes must be either $A$ or $C$ regardless of how many symbols are in $V$. The lower bound on successor length is enforced by making $\bar{\bar{\emptyset}}$ unable to be selected until $A_{\min}$ symbols exist in the successor.

Furthermore, after a list of possible symbols for a gene is determined, instead of giving each symbol an equal chance of selection, the ranges can be improved based on frequency of letters (and letters with context) occurring in $\rho$. For example, if the choices for a gene are $A$ and $B$, then instead of setting the probabilities of $A$ and $B$ to 0.5 for each gene, it is weighted by the frequency with which $A$ and $B$ occur in $\rho$. It is also possible to incorporate a context (sliding a window) to improve the probabilities. For example, with the string $AABAACAAB$, if the successor state is $A\ A\ \_$, then using two symbols of context within the string shows $AAB$ occurs twice and $AAC$ once, and therefore $B$ is given a 2/3 probability and $A$ is given a 1/3 probability. Note, all strings of $\rho$ are used to compute the associated probability. This encoding scheme is called OSoS($N$), where $N$ is the length of the context. This paper evaluates both OSoS(1) and OSoS(2).

*3.2.2. Growth Encoding*

This approach, called PMIT-D0L(G), searches using the GA within the computed lower and upper bounds for $M(A, B)$, of which there are $|V|^2$ values. PMIT-D0L(G) uses a literal encoding scheme and is similar to those seen in [25, 18], where the correct value of a dimension is each value in $M(A, B)$. For each combination of growth matrix tested, the sum of each row is obtained to give a length and then the scanning process is used.

The implementation chosen for this encoding allows for the possibility that a candidate does not satisfy the property that, for each $B \in V$, and for the growth values currently being assessed in $M'$, $\sum_{A \in V}(|\omega_i|_A M'(A, B)) = |\omega_{i+1}|_B, 1 \leq i < n$ by choosing each gene independently from the values chosen for other genes. This avoids backtracking as the GA is free to select any values within the lower and upper bounds for each $M'(A, B)$. An alternate encoding scheme was also tested by adapting mapped ranges so that this length constraint equation needed to be satisfied. The results of an evaluation were found to be very similar to

those for the encoding scheme described above, and so results for this alternative approach are omitted.

### 3.2.3. Length Encoding

This approach, called PMIT-D0L(L), uses the scanning process that requires a successor length for each $A \in V$. Each dimension is mapped onto an integer value representing the length of a successor of a symbol in $V$. The values of each dimension represents the length of a successor, with the dimensions bounded by the computed upper and lower bounds for length. Compared to the growth-based approach, although the bounds on the individual dimensions are larger; i.e., the lower bound $A_{\min} \geq \sum_{B \in V}(A, B)_{\min}$, and the upper bound $A_{\max} \leq \sum_{B \in V}(A, B)_{\max}$, with the length-based approach. The number of dimensions in the search space is $|V|$ with the length-based approach.

As with the growth-based approach, an alternative encoding was also implemented that enforced the constraint, for lengths $x_A$ for $A \in V$ currently being assessed, $\sum_{A \in V}(|\omega_i|_A x_A) = |\omega_{i+1}|, 1 \leq i < n$. The evaluation showed that the results were also not significantly different overall, and so this approach is not discussed further.

### 3.2.4. Row-Reduced Matrix Encoding

Recall that Doucet [25] recognized that the productions could be represented as a matrix equation which is given in Section 2.3. A similar approach was implemented. Let $Y$ be the matrix where each row $i$, from 1 to $n - 1$ is the Parikh vector of $\omega_i$, and let $Z$ be the matrix where each row is a Parikh vector of $\omega_2$ to $\omega_n$. In this case, if $M$ is a growth matrix of an D0L-system with $\rho$ as its length-$n$ developmental sequence, then $YM = Z$ is true. In Doucet's original work, they proposed to solve for $M$, and where $Y$ is invertible, $Y^{-1}Z$ is a unique solution, and if the solution is not unique, to use linear Diophantine equations.

It is also possible to replace $M$ with the length of each production, called the *successor length matrix*, and $Z$ is replaced with the column vector consisting of the length of $\omega_2$ to $\omega_n$, and this modified equation $YX = Z$ must have the length of the productions of any D0L-system having $\rho$ as developmental sequence as a solution. We can try all solutions to this equation (of which there can be more solutions than the correct L-system) similarly to the growth matrix version of this equation, and test whether each solution gives a D0L-system compatible with $\rho$. The remainder of this discussion will be presented in the context of a length-based matrix; however, similar logic applies to a growth-based matrix by replacing growth values for successor lengths. Furthermore, Gaussian elimination can be applied to this equation in order to produce parameterized equations in terms of successor lengths. That is, after Gaussian elimination is applied, it results in a set of linear Diophantine equations, where the successor lengths are the variables, e.g. $5X_1 + 3X_2 = 24$. This would mean that 5 times the first successor length plus 3 times the second successor length must be 24. It is easier to search the space defined by solutions to these equations than searching all possible length combinations. Each successor length only gets substituted for variables that appear in exactly one equation. In these cases, there are an

18

| Parameter | PMIT-D0L | | | | | |
|---|---|---|---|---|---|---|
| | OSoS(1) | OSoS(2) | G | M+G | L | M+L |
| P | 110 | 105 | 90 | 95 | 105 | 100 |
| C | 0.80 | 0.80 | 0.85 | 0.90 | 0.85 | 0.85 |
| M | 0.17 | 0.14 | 0.07 | 0.09 | 0.10 | 0.10 |

Table 1: Optimized parameters for each variant of PMIT-D0L

infinite number of possible solutions over the integers, however when inferring L-systems, the successor lengths are constrained to be natural numbers and within the bounds on the lengths provided by the lengths of the words in $\rho$, and it is therefore finite. For each equation, the encoding scheme used to search for a solution has $N$ genes, where $N$ is the number of variables in an equation. The range of values for each gene is $A_{\min}$ to $A_{\max}$ for the symbol $A$ the gene is representing (which can be more restricted than solutions to Diophantine equations due to the additional mathematical properties in Section 3.1 used that takes the sequences of the words into account). This encoding scheme is designated as PMIT-D0L(M+L) to indicate the addition of the matrix operation. For the matrix based on growth values, it is designated as PMIT-D0L(M+G).

### 3.3. Parameter Optimization

As described in Section 2.2, the function of the GA is controlled by the population size ($P$), crossover weight ($C$), and mutation weight ($M$) parameters. Optimizing these parameters is difficult based on general principles, since the optimal settings will depend on the characteristics of the fitness landscape, which is problem specific [27]. As such, typically, the parameters are set by doing a *hyperparameter search*. Bergstra and Bengio [27] found that using a Random Search provides an effective means to optimize the GA's parameter settings. Using Random Search works as follows. A range of good values is selected for each control parameter. In this case, based on the work by Grefenstette [28], the ranges were set to $10 \leq P \leq 125$ in increments of 5, $0.6 \leq C \leq 0.95$ in increments of 0.05, and $0.01 \leq M \leq 0.20$ in increments of 0.01, with additional values of 0.001 and 0.0001 permitted. An initial mid-range value is selected for each parameter ($P = 60$, $C = 0.8$, $M = 0.10$), and this is considered the current parameter value set. Iteratively, sixteen trials of PMIT-D0L are executed with a random variant of the current parameter value set. Each parameter is randomly modified up or down by no more than two increments, i.e. $P$ may be modified by $-10$, $-5$, 0, $+5$, $+10$, while also remaining within the ranges above. The variant parameter set that provides the best fitness value is considered the new current parameter value set. In the case of a tie, which was quite common with PMIT-D0L, the fastest execution time is used. When none of the sixteen trials provide an improvement over the current parameter value set, the hyperparameter search terminates. The resulting parameter value sets for each variant of PMIT-D0L is shown in Table 1.

### 3.4. Fitness Function and Termination Conditions

After a candidate L-system $G$ is produced from the solution $S$, the following process is used to evaluate fitness. To begin, any $G$ which produces more than double the expected number of symbols is assigned the maximum fitness value so it (practically) guaranteed to be culled in the survival step. Starting with $\omega_1$, a developmental sequence of length $n$ is produced using $G$ denoted as $\overline{\rho}$. For each $\overline{\omega}_i \in \overline{\rho}$, or until it terminates (see below), the symbol in each position of $\omega_i$ is compared to the corresponding position in $\overline{\omega}_i$. An error is counted if the symbol does not match (like Hamming distance), or if there is no corresponding symbol (i.e., one of the strings is longer or shorter than the other). For example, when comparing $\omega_i = XYXXXY$ to $\overline{\omega}_i = XYYX$, there are four errors. The third and fourth symbols differ, and additionally $\omega_i$ has six symbols, while $\overline{\omega}_i$ has only four. This process terminates when the number of errors for the $i^{th}$ derivation is greater than zero, as any errors will cascade forward. The fitness value is computed as the number of errors divided by the number of expected symbols (e.g., $4/6$), plus the number of unchecked derivations $(n - i)$. This encourages the GA to find solutions that incrementally match $\rho$.

PMIT-D0L uses a three-part termination condition to determine when to stop running. Ideally, PMIT-D0L terminates when a solution is found with a fitness value of 0.0 as this corresponds to an L-system that produces $\rho$ as its length $n$ developmental sequence. PMIT-D0L will also terminate when the population is considered to have converged to prevent the GA from acting as a random search and skewing the results. First, the current generation $Gen_{best}$ is recorded whenever a new best solution is found. If an additional $Gen_{best}$ generations pass without finding a new best solution, the population is considered converged. To prevent random chance from causing early termination, PMIT-D0L must process at least $1,000$ generations. PMIT-D0L also terminates after a time limit is reached. For this paper, the time limit was set to four hours; however, this was mainly used to control the overall experimental time. In practice, a user may be willing to wait less or more time to find an L-system.

## 4. Evaluation

### 4.1. Data Set

To evaluate PMIT-D0L's ability to infer D0L-systems, ten fractals, plus the six plant-like fractal variants (one shown in Figure 2) inferred by the existing program LGIN [2, 18], and twelve other biological models were selected from the vlab online repository [3]. The biological models consist of ten algaes, apple twig with blossoms (shown in Figure 1b), and a "Fibonacci Bush" (shown in Figure 1a). The dataset compares favourably to similar studies where only some variants of one or two models are considered [18, 19]. The data set is also of greater complexity by considering models with alphabets from between 2 to 31 (excluding constants) symbols compared to two symbol alphabets [18, 19]. In all, 28 D0L-systems that have been created manually by experts were included. An example of a larger L-system is given in Table 2.

| Productions |
| --- |
| $c \rightarrow FFFz[+k][-r]FFfd$ |
| $z \rightarrow Fz$ |
| $k \rightarrow lmfF$ |
| $r \rightarrow stfF$ |
| $d \rightarrow FFFz[+k][-r]FFfe$ |
| $l \rightarrow fF$ |
| $m \rightarrow n$ |
| $s \rightarrow fF$ |
| $t \rightarrow u$ |
| $e \rightarrow FFFz[+fj]FFfg$ |
| $n \rightarrow fFF[-A]Fo$ |
| $u \rightarrow fFF[+A]Fv$ |
| $j \rightarrow abF$ |
| $g \rightarrow FFFz[+k][-r]FFfh$ |
| $A \rightarrow fFB$ |
| $o \rightarrow fFF[-B]Fp$ |
| $v \rightarrow fFF[+B]Fw$ |
| $a \rightarrow Ff$ |
| $b \rightarrow c$ |
| $h \rightarrow FFFz[+k][-r]FFfi$ |
| $B \rightarrow fFC$ |
| $p \rightarrow fFF[-C]Fq$ |
| $w \rightarrow fFF[+C]Fx$ |
| $i \rightarrow FFFz[-fj]FFfc$ |
| $C \rightarrow fFD$ |
| $q \rightarrow fFF[-D]F$ |
| $x \rightarrow fFF[+D]F$ |
| $D \rightarrow fF$ |

Table 2: L-system for *Dipterosiphonia* v1 [9].

In order to further examine the performance of PMIT-D0L using a simple GA beyond the test set above, an additional set of D0L-systems is algorithmically created for various alphabet sizes (ignoring constants) $|\overline{V}|$ starting with $|\overline{V}| = 2$. For each size $|\overline{V}|$, a set of 100 D0L-systems is procedurally generated.

The design of the L-system generator was intended to be simple, while still generating realistic successors. Towards this goal, the successors of the expertly created L-systems were examined. Many of the fractal L-systems have successors where production and graphical symbols alternate. For example, Dragon Curve's productions are $X \rightarrow X+YF+$ and $Y \rightarrow -F-Y$. Also, the successors typically consist predominantly of a few symbols of $\overline{V}$ followed by a constant. Alternatively, they have long sequences of constants which make the problem easier for PMIT-D0L to infer due to the ability to line up constants. For example, *Aphanocladia* has productions

$A \rightarrow BA, \qquad B \rightarrow U[-C]UU[+/C/]U, \qquad U \rightarrow ffFFFF,$
$C \rightarrow FFfFFfFFfFF[-FFFF]fFFfFF[+FFF]fFFfFF[-FF]fFFf.$

Similar successors are found in *Ditria reptans*, *Ditria zonaricola*, *Metamorphe* and others. Notice also that with the branching patterns, a directional symbol immediately follows a branch open symbol [, which makes sense otherwise it would continue inline with the preceding angle.

The following methodology is therefore used to procedurally generate D0L-systems. The axiom is created by concatenating up to 4 random symbols of $\overline{V}$. Each production is created by iteratively concatenating symbols until a randomly selected length is reached (with a caveat on length described below) where 10 is used as an upper bound on production length. This bound of 10 was chosen to be larger than the average production length (which is 7.8) for expertly created D0L-systems and because we probabilistically discourage long stretches of constants. For each position of a production, there is a base 80% chance of selecting a symbol of $\overline{V}$; this is reduced by 20% for every consecutive symbol of $\overline{V}$ (e.g. if $ABA$ has been selected, then there is a 20% chance that the next symbol will be a symbol from $\overline{V}$). Once it has chosen to use some constant (or some element of $\overline{V}$), the particular letter is chosen with equal probability. Branching patterns must comply with the following rules: 1) they must be properly nested within each production, and 2) enforcing that a direction change symbol occurs after a "[" (essentially this is a type of normal form). The length of the successor may be exceeded to enforce the branching rules, i.e., additional "]" symbols can be added to make it properly nested. The L-system is validated by generating $|\overline{V}| + 1$ strings for $\rho$ and confirming that every symbol in $\overline{V}$ occurs at least once within the first $\overline{V}$ strings of $\rho$; it is easy to show that this is equivalent to ensuring that every symbol of $\overline{V}$ could eventually be reached, and therefore the L-system does not contain useless symbols.

### 4.2. Performance Metrics

Two performance metrics are used to measure how well PMIT-D0L can infer D0L-systems. The first metric is *success rate* (SR) which is defined as
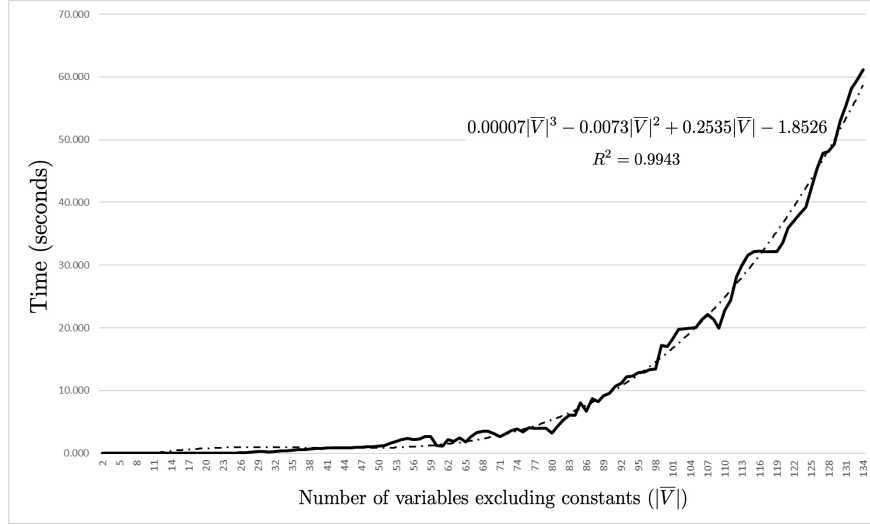
Figure 3: The solid line shows average MTTS across 100 procedurally generated D0L-systems for each number of non-constant variables ($\overline{V}$) from 1 to 134. The dash line shows the polynomial trend line for the MTTS results, and the equation is provided.

the percentage of times PMIT-D0L can find any L-system compatible with the input sequence. The second metric is *mean time to solve* (MTTS), in seconds, but measured to the millisecond level since some models can be determined in sub-second time. Time was measured using a single core of an Intel 4770 @ 3.4 GHz with 12 GB of RAM on Windows 10. These metrics are consistent with those found in literature [18, 19].

### 4.3. Results

The first set of results is MTTS results for the L-systems in the expert-created L-systems, shown in Table 3. The size of the variables (excluding constants) is given in the second column, and the MTTS using the various encoding schemes is given in columns 3 through 8. The SR was 100% for all L-systems and encoding schemes and is not shown. For PMIT-D0L(M+G) and PMIT-D0L(M+L), the systems where the matrix was invertible are marked with "*" as no searching was required. An average for each encoding technique is also provided. The lowest average was for PMIT-D0L(M+L).

Figure 3 shows $MTTS$ using PMIT-D0L(M+L) with the procedurally generated L-systems described in Section 4.1. The solid line gives the line graph of the average MTTS over 100 procedurally generated L-systems for each $|\overline{V}|$. This is tested for each size of $|\overline{V}|$ from 1 to 134, and all were inferred with 100% accuracy in less than one minute. The polynomial trend line is also provided.

### 4.4. Discussion

It is evident from Table 3 and the average row that OSoS(1) and OSoS(2) are performing worse than the other encoding schemes, and therefore using some

| Model | $|\overline{V}|$ | PMIT | | | | | |
|---|---|---|---|---|---|---|---|
| | | OSoS(1) | OSoS(2) | G | M+G | L | M+L |
| Algae | 2 | **0.001** | **0.001** | **0.001** | **0.001*** | **0.001** | **0.001*** |
| Cantor Dust | 2 | **0.001** | **0.001** | **0.001** | **0.001*** | **0.001** | **0.001*** |
| Dragon Curve | 2 | **0.001** | **0.001** | **0.001** | **0.001** | **0.001** | **0.001** |
| E-Curve | 2 | 24.312 | 23.075 | 0.026 | **0.025** | 0.088 | 0.029 |
| Fractal Plant v1 | 2 | **0.001** | **0.001** | **0.001** | **0.001*** | **0.001** | **0.001*** |
| Fractal Plant v2 | 2 | **0.001** | **0.001** | **0.001** | **0.001*** | **0.001** | **0.001*** |
| Fractal Plant v3 | 2 | **0.001** | **0.001** | **0.001** | **0.001*** | **0.001** | **0.001*** |
| Fractal Plant v4 | 2 | **0.001** | **0.001** | **0.001** | **0.001*** | 0.002 | **0.001*** |
| Fractal Plant v5 | 2 | **0.001** | **0.001** | **0.001** | **0.001*** | 0.003 | **0.001*** |
| Fractal Plant v6 | 2 | **0.001** | **0.001** | **0.001** | **0.001*** | 0.002 | **0.001*** |
| Gosper Curve | 2 | 0.022 | 0.023 | **0.001** | **0.001*** | 0.006 | **0.001*** |
| Koch Curve | 1 | **0.001** | **0.001** | **0.001** | **0.001*** | **0.001** | **0.001*** |
| Peano | 2 | 0.236 | 0.235 | **0.202** | 0.210 | 5.916 | 0.221 |
| Pythagoras Tree | 2 | **0.001** | **0.001** | **0.001** | **0.001*** | **0.001** | **0.001*** |
| Sierpenski Triangle v1 | 2 | **0.001** | **0.001** | **0.001** | **0.001*** | 0.010 | **0.001*** |
| Sierpenski Triangle v2 | 2 | **0.001** | **0.001** | **0.001** | **0.001*** | 0.003 | **0.001*** |
| *Aphanocladia* | 4 | 0.004 | 0.004 | 0.005 | **0.001** | 0.206 | 0.007 |
| *Dipterosiphonia* v1 | 28 | 11.885 | 10.871 | 123.008 | 3.820 | 178.718 | **1.639** |
| *Dipterosiphonia* v2 | 5 | 0.278 | **0.236** | 0.348 | 1.114 | 1.055 | 1.199 |
| *Ditria reptans* | 4 | 0.009 | 0.009 | 0.004 | **0.003** | 0.039 | **0.003** |
| *Ditria zonaricola* | 4 | 0.012 | 0.010 | 0.006 | **0.003** | 0.161 | 0.007 |
| *Herpopteros* | 4 | 0.019 | 0.017 | 0.007 | **0.004** | 0.070 | 0.006 |
| *Herposiphonia* | 5 | 0.026 | 0.022 | 0.016 | **0.013** | 0.190 | 0.015 |
| *Metamorphe* | 5 | 7732.255 | 512.040 | 1.381 | 3.793 | **0.632** | 2.387 |
| *Pterocladiella* | 30 | 22.631 | 8.805 | 0.944 | **0.881** | 4.120 | 3.192 |
| *Tenuissimum* | 31 | 0.851 | 0.871 | 0.603 | **0.520** | 120.619 | 1.141 |
| Apple Twig | 17 | 1.012 | 0.963 | **0.914** | **0.914** | 0.957 | 0.970 |
| Fibonacci Bush | 8 | 500.262 | 112.332 | 4.095 | 37.525 | 8.663 | **0.108** |
| Average | n/a | 296.208 | 23.912 | 4.699 | 1.744 | 11.481 | 0.391 |

Table 3: Comparison of MTTS in seconds for different encoding schemes for PMIT-D0L with the best MTTS bolded. $|\overline{V}|$ indicates the number of constant symbols in the L-system. SR is 100% for all executions. Results with "*" indicate an invertible matrix.

form of length with the scanning process seems to be best. However, OSoS(2) is faster overall than OSoS(1) (especially for *Metamorphe*), and therefore additional context is helping with OSoS. Overall, PMIT-D0L is 100% successful at inferring a diverse range of D0L-systems. The L-systems in the data set have different numbers of successors, with different lengths, and structures. With respect to using Doucet's [25] approach to find a unique solution, the matrix is found to be invertible for a little less than half of the L-systems in the test set, but never for any of the biological models. However, for both PMIT-D0L(M+G) and PMIT-D0L(M+L) the addition of the matrix operation to reduce the search space provides a benefit over PMIT-D0L(G) and PMIT-D0L(L) respectively. It is not so clear cut as to which encoding scheme is best, although PMIT-D0L(M+L) is the fastest overall, finishing in an average of 0.391 seconds. Certainly, it can be seen that PMIT-D0L(M+G) and PMIT-D0L(M+L) tend to be better than those without the matrix operations, but the timings tend to be

close. However, for Fibonacci Bush, PMIT-D0L(M+L) performed much better than PMIT-D0L(M+G). Overall, both are quite fast, and the same can be said for all of the growth-based and length-based encoding schemes. This leads perhaps to the conclusion that choosing between a growth-based or length-based encoding scheme is not so important, but rather the success of PMIT-D0L (of any type) is largely attributed to the space reduction techniques.

Since one benefit of automatic L-system inference is to infer L-systems that are not easily found by experts, the performance of PMIT-D0L was evaluated against procedurally generated D0L-systems of increasing complexity. It was found that PMIT-D0L seems to exhibit polynomial behaviour with respect to the number of successors with a trend line of $0.00007|\overline{V}|^3 - 0.0073|\overline{V}|^2 + 0.2535|\overline{V}| - 1.8526$. Using just a simple GA, PMIT-D0L was able to infer D0L-systems with $|\overline{V}| \leq 134$ in one minute or less. This is approximately five times larger than the largest D0L-system that could be found in the literature.

While all of the techniques are essential overall, the use of successor relationships extracts information by utilizing the sequence in which the symbols appear in $\rho$. This works by capitalizing on the fact that, even though the symbols are replaced in parallel, the order of the successors is the same as the symbols; i.e, if $\omega_i = A_1 A_2 \cdots A_m$ then $\omega_{i+1} = succ(A_1)succ(A_2) \cdots succ(A_m)$. Thus, if the relationship between $A_j$ and $succ(A_j)$ is known (or even partially known), much can be deduced about the location of $succ(A_j)$ in $\omega_{i+1}$ based on the location of $A_j$ in $\omega_i$. This in turn allows for the deduction of symbols near $A_j$. This is one of the main differences between PMIT-D0L and other existing approaches. Capturing information from the symbol sequence of strings in $\rho$ should provide guidance towards future investigations on inferring L-systems.

## 5. Conclusions

This paper presented an evaluation of different encoding schemes for the Plant Model Inference Tool for deterministic context-free L-systems (PMIT-D0L) to infer L-systems from a sequence of strings. Some of the encoding schemes are based on modifications of earlier works, while some are novel. The classical encoding schemes look at the problem of inferring successors by letting each position of each production be an unknown [19]. Here, we use a novel encoding scheme which considers the length of each production (or each Parikh vector) as an unknown, as it is straightforward to determine the L-system from the production lengths.

The evaluation of the different encoding schemes does not indicate a clear best encoding, however length-based approaches are the fastest for this test set. Much of this paper focused on techniques for reducing the search space size, using necessary conditions. Some of the techniques, such as setting up associations between constants (usually created from graphical symbols) are novel and particularly effective. The techniques are effective to the degree that the choice between growth-based and length-based is not particularly critical for this particular test set.

The inductive inference of L-systems from input strings allows for much more rapid development of models than the current approach of building models by hand [2, 10, 5, 4], which can take considerable effort. Additionally, by going directly from observation (strings) to a model, allows for mechanistic principles to be possibly revealed, as opposed to requiring expert knowledge to build the model.

Since PMIT-D0L seems capable of inferring L-systems with fairly large alphabets (at least 31 symbols in the test set in a fast manner, and much larger in the algorithmically generated L-systems), this work will be used as a base for investigating the inference of other, more complex, L-system extensions such as stochastic L-systems, parametric L-systems, and for using images as input. While inferring parametric L-systems is more complex than D0L-systems, some of the principles and techniques from this paper may be applicable, especially if these L-systems behave like D0L-systems in certain parts of their derivation, such as when parameters are used as a timing mechanism. Also, the use of markers and solution projection should still be applicable as well to parametric L-systems. As a final note, the speed of PMIT-D0L can be further enhanced using parallel processing, which will be explored in the future.

## 6. Acknowledgements

We thank the anonymous reviewers for valuable suggestions which considerably improved the presentation of the paper.

## References

[1] A. Lindenmayer, Mathematical models for cellular interaction in development, parts I and II, Journal of Theoretical Biology 18 (3) (1968) 280–315.

[2] P. Prusinkiewicz, A. Lindenmayer, The Algorithmic Beauty of Plants, Springer Science & Business Media, 2012.

[3] University of Calgary, Algorithmic Botany.
URL http://algorothmicbotany.org

[4] T. Nishida, K0L-system simulating almost but not exactly the same development-case of Japanese Cypress, Memoirs of the Faculty of Science, Kyoto University, Series B 8 (1) (1980) 97–122.

[5] C. Godin, P. Ferraro, Quantifying the degree of self-nestedness of trees: application to the structural analysis of plants, IEEE/ACM Transactions on Computational Biology and Bioinformatics 7 (4) (2010) 688–703.

[6] M. T. Allen, P. Prusinkiewicz, T. M. DeJong, Using L-systems for modeling source–sink interactions, architecture and physiology of growing trees: The L-PEACH model, New Phytologist 166 (3) (2005) 869–880.

[7] T. Watanabe, J. S. Hanan, P. M. Room, T. Hasegawa, H. Nakagawa, W. Takahashi, Rice morphogenesis and plant architecture: measurement, specification and the reconstruction of structural development by 3D architectural modelling, Annals of Botany 95 (7) (2005) 1131–1143.

[8] P. Prusinkiewicz, S. Crawford, R. S. Smith, K. Ljung, T. Bennett, V. Ongaro, O. Leyser, Control of bud activation by an auxin transport switch, Proceedings of the National Academy of Sciences 106 (41) (2009) 17431–17436.

[9] R. Morelli, R. Walde, E. Akstin, C. Schneider, L-system representation of speciation in the red algal genus Dipterosiphonia (Ceramiales, Rhodomelaceae), Journal of Theoretical Biology 149 (4) (1991) 453–465.

[10] P. Prusinkiewicz, L. Mündermann, R. Karwowski, B. Lane, The use of positional information in the modeling of plants, in: Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques, ACM, 2001, pp. 289–300.

[11] M. A. Galarreta-Valverde, M. M. Macedo, C. Mekkaoui, M. Jackowski, Three-dimensional synthetic blood vessel generation using stochastic L-systems, in: Medical Imaging: Image Processing, 2013, p. 86691I.

[12] C. Jacob, Genetic L-system programming: breeding and evolving artificial flowers with Mathematica, in: Proceedings of the First International Mathematica Symposium, 1995, pp. 215–222.

[13] K. J. Mock, Wildwood: The evolution of L-system plants for virtual environments, in: Proceedings of the 1998 IEEE World Congress on Computational Intelligence, IEEE, 1998, pp. 476–480.

[14] J. Guo, H. Jiang, B. Benes, O. Deussen, X. Zhang, D. Lischinski, H. Huang, Inverse procedural modeling of branching structures by inferring L-systems, ACM Transactions on Graphics 39 (5) (2020).

[15] G. Herman, G. Rozenberg, Developmental Systems and Languages, North-Holland, Amsterdam, 1975.

[16] N. Khan, O. Lyon, M. Eramian, I. McQuillan, A novel technique combining image processing, plant development properties, and the Hungarian algorithm, to improve leaf detection in maize, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW 2020), 2020, pp. 330–339.

[17] S. Das Choudhury, S. Bashyam, Y. Qiu, A. Samal, T. Awada, Holistic and component plant phenotyping using temporal image sequence, Plant Methods 14 (1) (2018) 35.

[18] R. Nakano, N. Yamada, Number theory-based induction of deterministic context-free L-system grammar, in: International Conference on Knowledge Discovery and Information Retrieval, SCITEPRESS, 2010, pp. 194–199.

[19] B. Runqiang, P. Chen, K. Burrage, J. Hanan, P. Room, J. Belward, Derivation of L-system models from measurements of biological branching structures using genetic algorithms, in: Proceedings of the International Conference on Industrial, Engineering and Other Applications of Applied Intelligent Systems, Springer, 2002, pp. 514–524.

[20] J. Bernard, I. McQuillan, New techniques for inferring L-systems using genetic algorithm, in: Proceedings of the 8th International Conference on Bioinspired Optimization Methods and Applications, Vol. 10835 of Lecture Notes in Computer Science, Springer, 2018, pp. 13–25.

[21] J. Bernard, I. McQuillan, A fast and reliable hybrid approach for inferring L-systems, in: Proceedings of the 2018 International Conference on Artificial Life, MIT Press, 2018, pp. 444–451.

[22] P. Prusinkiewicz, J. Hanan, Visualization of botanical structures and processes using parametric L-systems, in: Scientific visualization and graphics simulation, John Wiley & Sons, Inc., 1990, pp. 183–201.

[23] T. Bäck, Evolutionary Algorithms in Theory and Practice: Evolution Strategies, Evolutionary Orogramming, Genetic Algorithms, Oxford University Press, 1996.

[24] F. Ben-Naoum, A survey on L-system inference, INFOCOMP Journal of Computer Science 8 (3) (2009) 29–39.

[25] P. Doucet, The syntactic inference problem for D0L-sequences, L Systems (1974) 146–161.

[26] I. McQuillan, J. Bernard, P. Prusinkiewicz, Algorithms for inferring context-sensitive L-systems, in: 17th International Conference on Unconventional Computation and Natural Computation, Vol. 10867 of Lecture Notes in Computer Science, Springer, 2018, pp. 117–130.

[27] J. Bergstra, Y. Bengio, Random search for hyper-parameter optimization, Journal of Machine Learning Research 13 (2012) 281–305.

[28] J. J. Grefenstette, Optimization of control parameters for genetic algorithms, IEEE Transactions on Systems, Man and Cybernetics 16 (1) (1986) 122–128.