

## A SYNTACTIC PATTERN RECOGNITION SYSTEM FOR DNA SEQUENCES

DAVID SEARLS and SHAN DONG

*Department of Genetics, University of Pennsylvania School of Medicine  
422 Curie Boulevard, Philadelphia, PA 19104-6145 USA*

### ABSTRACT

We review both theoretical and practical results of a linguistic approach to studying the structure of features of DNA sequences. Using generative grammars, complex assemblages can not only be described and analyzed abstractly, but also concretely, such that features can be searched for by a general-purpose parser. Our parser, called GENLANG, uses an extended logic grammar formalism and has found features as complex as tRNA genes, group I introns, and protein-encoding genes, within input sequences on a genomic scale.

### 1. Introduction

Recently a number of authors have adopted a “linguistic” view of DNA sequences. Most of their work has involved examinations of the occurrences of “words” in DNA in what is essentially an information-theoretic approach<sup>1,2</sup>, or the use of statistical analyses of vocabularies in the tradition of comparative linguistics<sup>3,4</sup>. These approaches, however, can be distinguished from another methodology in linguistics, pioneered by Noam Chomsky, which attempts to characterize higher-order phenomena in languages by way of generative *grammars*. Only a few authors have used this discipline in the realm of biological sequences<sup>5,6,7</sup>; we have found that such a linguistic approach proves useful not only in theoretical characterization of certain structural phenomena in sequences, but also in generalized pattern recognition in this domain, via *parsing*. In this paper we will review both these theoretical findings, and practical results using a logic-based parser for *syntactic pattern recognition*.

### 2. Background

We begin with the necessary background concerning the formalism of language theory, as well as some of its uses and implementations.

#### 2.1. Formal Language Theory

In formal language theory, a language is defined in terms of an *alphabet*, which is a finite set of *symbols*; in the case of DNA sequences, such symbols would be the nucleotide bases. A DNA molecule can then be represented as a finite *string* of symbols from this alphabet; a *language*, formally, is any set of such strings.

The concern of formal linguistics is the finite representation of languages which may themselves be infinite; the goal is an economy of expression, in an abstract representation, as an alternative to exhaustively enumerating all the allowable strings in a language. Such cogency may also have the benefit of capturing some kind of essential, clarifying generalization about the structure or *syntax* of a linguistic system, preferably related to the meaning or *semantics* of the language elements. For this purpose, language generators called *grammars* have proven extremely useful. Grammars specify languages through sets of *rules*, which achieve the desired succinctness largely by referring to each other and to themselves *recursively*. Perhaps the most important class of grammars is the *context-free grammars* (CFGs, which specify the *context-free languages*, CFLs). A CFG has a set of *terminals* or alphabetic elements, and an additional set of symbols called *nonterminals*; these symbols are used in a finite set of rules whose members are written as  $A \rightarrow u$  where  $A$  is a nonterminal and  $u$  is a string of terminals and nonterminals. A grammar generates the string elements of its language by taking a starting symbol  $S$  and rewriting it, by repeatedly finding a rule whose left-hand side matches some nonterminal in the current string, and substituting that rule's right-hand side, until the string contains all terminals. Such derivation steps are denoted by a double arrow,  $\implies$ , so that for the simple grammar with alphabet  $\{a, b, c\}$ , nonterminals  $S$  and  $X$ , and rules  $S \rightarrow aX$ ,  $X \rightarrow bX$ , and  $X \rightarrow c$ , one possible derivation is:

$$S \implies aX \implies abX \implies abbX \implies abbc$$

This grammar formalism would appear to be preferable to either trying to list the infinite set of strings  $\{ac, abc, abbc, abbbc, \dots\}$ , or using the informal description “all strings consisting of an  $a$  followed by any number of  $b$ 's followed by a  $c$ .” For one thing, it makes feasible the computational task of *parsing* a string to determine whether it is in the language specified by the grammar. A useful byproduct of parsing is the production of a *parse tree* reflecting the grammar rules applied and giving a kind of structural description of grammatical features in the input string—exactly the kind of output that is desired in describing certain biological sequence data.

CFGs have proven to be very useful in the field of compiler construction, where, in the form of BNF (Backus-Naur Form) descriptions, they are used to formally specify programming languages. An even more interesting application of computational linguistics, however, is in understanding natural language—a complex problem that has stimulated a large body of research. Although straightforward CFGs can be written that cover many aspects of natural language syntax, natural languages in their full generality are now thought to require greater than context-free power.

*Regular expressions*, using the basic operations of concatenation, disjunction (logical OR), and “star” (denoting any number of occurrences of its argument), also specify languages. However, the set of *regular* languages is strictly a subset of the CFLs, for no regular expression can specify certain *self-embedding* structures such as palindromes; computationally, these require a stack to store information about *de-*

*dependencies* between distant elements of the string. In fact, even the CFLs are a strict subset of the *context-sensitive* languages (CSLs), described by grammars that have more than one symbol on the LHS of rules. While CFLs are restricted to describing *nested* dependencies, CSLs can specify *crossing* dependencies, such as those found in *copy languages*, which contain duplicated strings of arbitrary extent. These language classes all take their place on the *Chomsky hierarchy* of languages, which categorizes the linguistic complexity of any given language, and which serves as the basis for analysis of the decidability and/or tractability of recognizing strings of any language with general-purpose parsers.  $O(n)$  parsers are easily designed for regular and *deterministic* CFLs—those that can be recognized without the need to backtrack on the input string—and  $O(n^3)$  parsers exist for *any* CFL. Certain well-defined characteristics of CFGs, such as *ambiguity* (the ability to derive the same strings via multiple distinct parses), determine their potential for more efficient general-purpose parsing, and for any *particular*, narrowly-defined language special-purpose linear-time recognizers can often be designed. Languages beyond context-free are increasingly more difficult to recognize by general-purpose parsers, and much effort has gone into defining language classes “slightly greater” than context-free that are adequate to a particular domain (such as natural language) yet can be parsed efficiently.

## 2.2. Syntactic Pattern Recognition

This methodology is not limited to computer languages or human natural languages, but can be extended to all manner of signals, images, or other data which have underlying structure. The field of Syntactic Pattern Recognition (SPR)<sup>8</sup> makes use of the tools and techniques of computational linguistics, such as grammars and parsers, to specify and search for patterns in data. Because grammars intrinsically promote the hierarchical abstraction of features, these can be built up to a very high level while maintaining a clear, modular “knowledge base.” Moreover, grammars by their nature detect individual features in this higher-level context, which creates a much greater degree of discrimination than isolated searches. SPR benefits from a strong formal foundation, but also incorporates features that extend the expressive power of grammars where necessary for the domain. For example, “noisy” signals can be dealt with by so-called stochastic grammars<sup>8</sup>, which incorporate probabilities into grammars in a natural way. SPR, in fact, has been classified as a form of *pattern-directed inference*, and indeed we have found that it provides an excellent framework for the incorporation of heuristics at many levels. SPR has been successfully applied to such problems as general signal processing, handwritten character recognition, and karyotype analysis by us<sup>9,10</sup> and many others<sup>8</sup>, and our initial results with SPR and the linguistic analysis of DNA (described below) suggest that it can address many of the complexities of this domain as well.

## 2.3. Logic Grammars

Prolog is a programming language that gives a procedural interpretation to a subset of first-order predicate logic. It uses a clausal form that lets programs be

written as databases of atomic predicates called *facts*, e.g., `protein(hemoglobin)`, and *rules* of the form `protease(X) :- protein(Y), degrades(X,Y)`. This can be read “X is a protease *if* there is a protein Y *and* X degrades Y.” Prolog’s rules and facts, together called *relations*, can be queried to perform inferences by backward-chaining proof, using a mechanism called *resolution*, and a powerful pattern-matching facility known as *unification*. The resulting system is able to perform computation as controlled deduction—in fact, a form of theorem proving.

Prolog’s history is closely linked with the formalism of Definite Clause Grammars (DCGs), and the notion that grammars can be expressed as rules of a Prolog program<sup>11</sup>. The process of parsing a string then becomes that of proving a theorem given that string as input and the “axioms” of a grammar. In practice, the previous example grammar would appear as in the code below.

```
s --> [a], x.                x --> [b], x | [c].
```

In Prolog, logical predicates begin with a lower-case letter, and in DCGs these correspond to nonterminals. (Prolog variables begin with upper-case letters.) Terminals are shown as Prolog list elements; lists generally appear within square brackets, with list elements separated by commas (e.g., `[a,b,c]`), but can also be double-quoted strings which correspond to lists of ASCII character codes, e.g., `"gattac"`. The vertical bar in the second rule is an “or” (disjunction), so that this rule for `x` actually represents the two rules given previously.

DCGs actually require a translation step to become Prolog clauses, because Prolog must have a mechanism for manipulating the input string, which it does by maintaining the string in “hidden” parameters of the nonterminals, called *difference lists*. In fact, the user may attach other parameters to nonterminals, as well as embed arbitrary code in rule bodies, all of which lends DCGs power far beyond CFGs. It is our belief that they offer a firm foundation for a language that can be a powerful descriptive tool for molecular genetic features, and also offers a flexible analytical and control mechanism through Prolog’s built-in DCG parser.

### 3. Theoretical Results

We suggested several years ago<sup>7</sup> that nucleic acids were beyond regular and at least context-free, based on the phenomenon of secondary structure: the stem portion of a stem-and-loop structure entails *nested dependencies* between base-paired residues, which are easily specified by a self-embedding CFG, but which in the general case are beyond the capabilities of any formal regular expression. It was also suggested<sup>7,12</sup> that DNA may be beyond context-free as well, due to the phenomenon of direct repeats, which constitute a *copy language* with crossing dependencies that cannot be described with essentially stack-based context-free formalisms. Subsequent work<sup>13</sup> formalized these conjectures.

Formal discussion of the linguistic status of DNA, like that of natural language, may be based on empirical phenomenology, but in the case of nucleic acids may also rest on the actual physical structure of the molecules, and in particular the ability to form secondary structure. We have offered formal proofs that idealized representations of such structure are indeed non-regular<sup>13</sup>. Beyond this, however, the simple existence of direct repeats is somewhat unsatisfying as evidence for the purely *formal* status of DNA, since it is only when direct repeats with no particular bounds on their extent can be shown to be *necessary* in a language that it can be said to be greater than context-free on that account. So as not to depend upon *ad hoc* phenomenology, we sought examples which, like that of inverted repeats, could be grounded in actual physical structures and processes arising in the molecules themselves. A series of such arguments were presented<sup>13</sup> based upon (1) the potential for circularization of DNA with terminal direct repeats, (2) unequal crossing-over in multiple tandem repeats, and, most importantly, (3) the existence of pseudoknots in structural RNA, which entail crossing dependencies between stems within each other's loops. It is interesting that dealing with pseudoknots has required major reimplementations of some RNA structure prediction programs<sup>14</sup>; the relative difficulty of this can in part be explained by the transition to greater-than-context-free recognition which can thus no longer be strictly stack-based. This points again to the utility of a solid formal linguistic characterization in the design of recognition algorithms in any given domain.

A number of additional formal results were given<sup>13</sup>, dealing with other linguistic attributes, again based on a somewhat idealized model of the structural characteristics of nucleic acids. Inverted repeats, for example, were shown to be *nondeterministic* languages, and the branching or *recursive* nature of secondary structures implies their language is *non-linear* (in this context, meaning that any grammar describing them must have a rule with more than one nonterminal on its right hand side). The *ambiguity* of general secondary structural grammars (that is, their ability to produce more than one essentially distinct parse for the same primary sequence) was also explored<sup>13</sup>, and it was shown that this grammatical ambiguity reflects alternative secondary structures in a biologically relevant way. (Since that time, we have proven that, while the most general language of ideal orthodox secondary structure is actually deterministic and thus unambiguous, biologically plausible secondary structure sublanguages are *inherently ambiguous*, i.e. impossible to describe by any unambiguous grammar). These results rule out the use of certain  $O(n^2)$  simplifications of general-purpose context-free parsers (which are otherwise  $O(n^3)$ ).

Given that biological sequences are beyond context-free, it is of interest to circumscribe their exact boundaries. We suggested<sup>7</sup> and subsequently demonstrated<sup>13</sup> that the language encompassing all of the phenomena described above in nucleic acid structure lies not only in the CSLs, but within a restricted subset known as the *indexed languages*. A subset of the indexed languages with a very perspicuous grammar formalism particularly well-suited to nucleic acids, known as *string variable grammar*, was developed by us<sup>12,13</sup>; examples are given in the next section.

We have also explored *closure* properties of the Chomsky hierarchy under biological operations—that is, whether language classes of interest, after undergoing certain biological processes, can be expected to remain at the same level in the Chomsky hierarchy or not. We have found<sup>13</sup> that regular languages and CFLs are closed under double-stranded replication and under simple recombinational events such as scission and ligation. Deterministic languages, however, are not closed under these operations, so that, for example, certain features can be recognized more efficiently on one strand of DNA than on the other, or in other than a leftmost fashion, suggesting the use of so-called *island parsing* strategies. With regard to *evolutionary* operations such as duplication, inversion, and transposition, it was shown that CFLs are *not* closed, suggesting that genomic rearrangements on an evolutionary scale may be responsible for increasing the mathematical complexity of the genetic language.

## 4. Practical Results

GENLANG, our logic-based syntactic pattern recognition system for DNA sequences, has been described previously<sup>7,10,12,13,15</sup>. Some of its newer developments will be reviewed here.

### 4.1. String Variables

One advantage of logic grammars lies in the rapid prototyping capabilities of Prolog, including the ability to easily add new syntactic constructs. For example, in GENLANG queries are of the form  $\langle pattern \rangle : \langle parse\ variable \rangle ==> \langle input \rangle$ , where new infix operators separate grammar elements: the *pattern* generally contains the top-level nonterminal in the grammar, the *parse variable* is a logic variable to which a parse tree will be bound, and the *input* is as described below. In the parse tree, nonterminals are typically adorned with information about their cost (in number of mismatches), their location in the input, the actual primary sequence recognized, etc. Similarly, *gaps* of either unbounded (...) or bounded (e.g. 19...27) extent were added to the language.

A more significant extension to DCGs, called *string variables*<sup>12,13</sup>, also benefits from the Prolog milieu. A string variable is a logic variable appearing in the body of a grammar rule, which stands for a string of arbitrary extent, and which may optionally have applied to it operators such as the tilde which denotes reverse complementarity. Such a feature in this domain makes it easy to specify even complex arrangements of direct and inverted repeats, such as are characteristic of secondary structure:

```
tandem_repeat ---> X, X.                stem_loop ---> X, ..., ~X.
pseudoknot ---> X, ..., Y, ~X, ..., ~Y.
attenuator ---> X, ..., ~X, ..., X.
```

These descriptions are at a much higher level than the corresponding context-free grammars<sup>12</sup>, and in fact most of them are even beyond context-free, yet despite their widely varying linguistic complexity they are expressed with comparable ease in

this formalism. Even so, these are relatively abstract descriptions and for purposes of parsing require “real world” constraints on their length and degree of mismatch allowed in stems; we next describe the mechanism for controlling these.

#### 4.2. Attributes

Objects in GENLANG can have attached to them an *attribute list* using the syntax,  $\langle object \rangle : [\langle attribute1 \rangle, \langle attribute2 \rangle, \dots, \langle attributeN \rangle]$ . There are four primary types of attributes, which generally entail operations on keywords: (1) *control* attributes which control the parse and the position on the input string at a meta level; (2) *constraint* attributes which impose constraints on (typically numerical) quantities denoted by keywords, such as the cost (normally, the number of mismatches) of a subtree; (3) *specification* attributes whereby the values of keywords such as ‘cost’ can be redefined with arbitrary expressions; and (4) *assignment* attributes which assign the values of keywords to logic variables which may then be carried through the parse and/or reported at the top level. For example:

```
foo:[cost=S+2*C] ---> "atg", ...:[step=3,S=size], bar:[size<50,C=cost].
```

Here, the control attribute `step=3` specifies that the gap (...) is to increase in increments of 3; the constraint attribute `size<50` keeps the span of the nonterminal `bar` under 50; the specification attribute `cost=S+2*C` redefines the cost of the nonterminal `foo` to an arithmetic function of the size of the gap and the cost of the `bar`; and the assignment attributes `S=size` and `C=cost` serve to bind those variables. (The default cost of `foo` would have been the number of mismatches in the "atg" plus those within `bar`.)

Attributes are managed by way of additional “hidden parameters” in the implementation of the grammar. Just as the difference lists serve to unburden the grammar designer of the low-level programming involved in input list management, a total of ten hidden parameters now hide such details as the accumulated cost of parse trees and the cost thresholds applied by the grammar, additional constraints on the size ranges of individual elements in the grammar, and the parse tree itself.

#### 4.3. Gaps

Gaps represent regions that are “skipped over”, but in fact it is the gaps that do the skipping—they constitute the multiple, embedded search engines of a typical grammar, and the source of most of its non-determinism—and so careful attention to their implementation has been necessary. One important feature is *delayed* (or *lazy*) *evaluation*: gaps encountered in the course of a parse are “packaged” and passed down the parse tree, and are not actually evaluated until they in turn encounter some feature with which they may combine for more efficient evaluation. For example, the combination of a lazy gap with a string of bases might under the right circumstances allow the string simply to be looked up in a hash table, rather than searched for in the primary sequence. GENLANG does, at the option of the user, hash its input into  $k$ -tuples of varying sizes, and permits many hundred-fold more efficient recognition of features such as direct repeats.

The implementation of lazy gaps also lends itself to finer control over the search strategy used by the parser. The logic-based parser is ordinarily breadth-first on the *input*, in the sense that all applicable rules will be tried at every position in a parse before moving on to the next position following a gap. However, a *lazy* gap will be passed to the first applicable rule, and that rule will be tried in every possible position permitted by the gap, before the gap is passed on to the next applicable rule—in other words, depth-first search on the input. The search style in GENLANG can be controlled either at a global level, or locally through the use of the attributes `deep`, `wide` or `best` (the latter performing best-first search within a defined range).

#### 4.4. Cost

A rule such as `foo ---> "gat" | "gaa" | "gta" | ...` would ordinarily succeed upon recognizing any of the disjuncts on the input string (or, any of the disjuncts with mismatches allowed up to the cost threshold established higher in the parse tree). However, with the addition of a `consensus` attribute, the disjuncts are treated as exemplars in the calculation, at compile time, of a weight matrix which now contributes the cost at this point of the parse. Moreover, the order in which the parser examines positions in the input string is reordered at compile time, so that the most “informative” positions—that is, those for which inappropriate input will most rapidly cause the nonterminal to exceed its cost threshold and thus fail—are examined first, for optimal efficiency. Base frequency data may also be entered into the grammar in tabular form, based on published data; however, compilation time even for large lists of exemplars is nearly as fast, and allows the user to enter new data freely, postulate classes by dividing exemplars among several nonterminals, etc. There are several methods available for calculating costs from base frequency data (e.g. the attribute `cost=neglog` uses a negative logarithm of base frequency), which are user-definable as well.

GENLANG also offers cost attributes that support rules concerning the local enrichment of features. For example, the current rule used for the pyrimidine-rich region between the branch and the 3' acceptor in a protein-encoding grammar is

```
pyrimidine_rich:[parse=[size,cost,list]] --->
  probably_pyrimidine, pyrimidine_rich | [].

probably_pyrimidine:[cost=(-5),once] ---> pyrimidine:[cost=0].
probably_pyrimidine:[cost=(3/5)] ---> base.
```

When a cost threshold is placed on this feature at a higher level in the grammar, the recursive rule finds the longest list it can of “probably\_pyrimidine’s” which are defined as either an authentic pyrimidine (its cost constrained to be zero), in which case the overall cost is decremented by 5, or as any other base, in which case the cost given as a ratio indicates that the cost is to be calculated as 3/5th of the way from the current cost to the threshold. This has the effect of giving “credit” for every hit (though the cost is not permitted to go below zero), and on every miss causing a



logarithmic decay in the cost until it intersects the threshold in a region that is too pyrimidine-poor. This paradigm is not limited to simple rules for single bases, but can be applied to repetitions of very high-level patterns as well.

Note also the definition constraint for `parse` in the first rule above; which attributes are reported in the parse tree (in this case, the feature's `size`, overall `cost`, and a `list` displaying its primary sequence), are user-definable at every level.

#### 4.5. *Input*

Prolog's linked lists proved to be an inappropriate data structure for the exceedingly long inputs found in the biological domain, and so GENLANG is actually a hybrid software system in which the input and certain other data structures, as well as some low-level procedures, are implemented in 'C'. Besides the hashing of the input described above, another important speedup is the implementation of *chart parsing*, a dynamic programming technique that prevents repeated recomputation of sub-trees of the parse and ameliorates much of the cost of non-deterministic parsing. This technique has been enhanced with a domain-specific feature that allows for rapid traversal of regions already parsed, using a byte-encoding of flags marking these regions for any particular nonterminal.

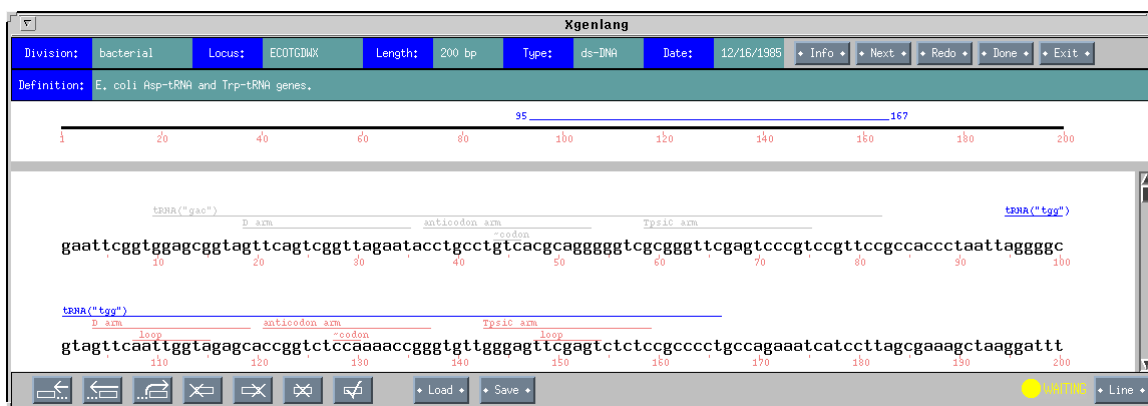
The form of the input presented to GENLANG can range from actual lists of bases, to sequence files, to directories of files, and even flat sequence database files, with individual entries handled sequentially, as if disjointed. This allows what amounts to grammar-based query of GenBank. The user may filter selections based on attributes of the database "object". Thus, a typical query is

```
(...,gene,...) ==> genbank:[division=rodent,locus='MUS*',  
                           definition='*complete cds*',date>1/1/92].
```

Text strings are given in single quotes and may include an asterisk "wild card". This query would search for a `gene` in the rodent file of genbank, selecting entries with locus ID's beginning with "MUS" (for mouse), whose definition contains the string "complete cds", which were entered after the beginning of this year. The processing of the input attributes is accomplished through the use of an additional logic grammar within GENLANG describing the structure of GenBank entries. Using this approach, new forms of data (e.g. ASN.1 or even relational databases) will be easily handled by simply adding new such attribute grammars, which can generally be developed in short order using DCGs.

#### 4.6. *Parse Visualization*

The current version of GENLANG has a graphical parse visualization and management tool, implemented in 'X' using Motif widgets and the Quintus Prolog interface. This tool provides a dynamic, real-time picture of the developing parse, so that the user is aware of the current state of computation at all times. This is important in the development of efficient grammars, since it can be made readily apparent to an observer when a highly non-deterministic feature is causing "thrashing" of the parse



**Figure 1.** The parse visualization tool, showing a parse of a pair of tRNA genes.

at some particular point in the input. It is also very useful in debugging of grammars, since it can be seen at what point a parse fails on test input that is known to be valid; this is especially important given the relative difficulty of debugging complex logic programs such as this, in the absence of the external, uniform point of reference provided by the input string, and a means of visualizing the activities of the logic programming relative to it. Finally, an unanticipated benefit of the parse visualization tool has been the capacity it affords for manual control over the parse at a meta level; a series of such controls are available from the graphical interface to allow sub-trees of the parse to be altered selectively, in a manner not otherwise available in a “flat” interface. The overall display for the parse visualization tool is shown (Figure 1).

At the top of this window are two rows giving salient header information extracted from the GenBank entry. At the upper right are five buttons controlling the top level of the input processing: *Exit* and *Done* buttons, a *Redo* button which causes the parse to start over on the current entry, a *Next* button which causes it to skip to the next entry, and an *Info* button which opens a scrolling window giving the entire GenBank header. The latter feature is most useful when comparing the results of a parse to information in the GenBank Features table. All of these controls have also been made available in the flat TTY interface to GENLANG, via modifications to the Prolog interrupt handler, so that any errant parse may be stopped by a control-C and redone, skipped over, compared with the Features table, etc. at any time.

The long narrow window near the top gives a picture of the overall input string, shown as a long line labelled with numerical positions. Immediately above this is a thin line extending over a portion of the entire sequence, with its bounds labelled numerically; this corresponds to the highest-level feature parsed.

The large main window portrays the primary sequence in a series of rows; the number of rows in the window and the number of bases placed in each row—either 100 or 200—are controlled by the user, as well as background colors and other graphical attributes. Above each row of primary sequence is drawn the parse tree as it develops. Lower-level nonterminals are first drawn immediately above the sequence as they are

```

tRNA("ctt"): span=90470/90583
  stem: list="ggttggtt"
  turn: list="tg"
  D arm: span=90479/90496
    stem: list="gcc"
    loop: list="gagcgggtctaa" cost=0
    ~stem: list="ggc" cost=0
  base: list="g"
  anticodon arm: span=90497/90546
    stem: list="cctgat"
    loop: cost=0
    pre_codon: list="tc"
    ~codon: list="aag"
    purine: list="a"
    intron: span=90509/90539
      ...: size=5
      "ctt": span=90514/90517
      ...: size=22
      base: list="a"
    ~stem: list="ctcagg" cost=1
  extra arm: span=90546/90559
    ...: size=13
  TpsiC arm: span=90559/90576
    stem: list="aagag"
    loop: list="ttcgaat" cost=6
    ~stem: list="ctctt" cost=0
    ~stem: list="agcaacc" cost=0.5

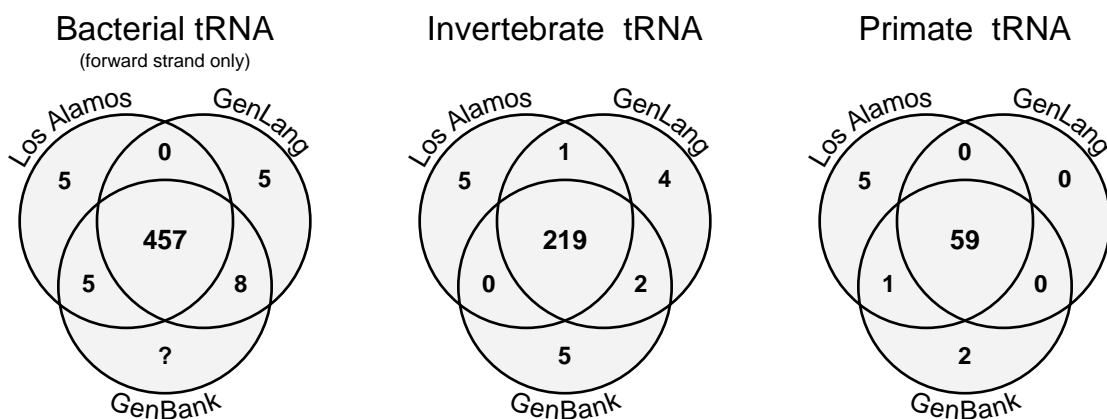
tRNA("aac"): span=127320/127393
  stem: list="gactcca"
  turn: list="tg"
  D arm: span=127329/127346
    stem: list="gcc"
    loop: list="aagttgggttaa" cost=0
    ~stem: list="ggc" cost=0
  base: list="g"
  anticodon arm: span=127347/127364
    stem: list="tgcga"
    loop: list="ctgttaa" cost=0
    pre_codon: list="ct"
    ~codon: list="gtt"
    purine: list="a"
    base: list="a"
    ~stem: list="tcgca" cost=0
  turn: list="agatc"
  TpsiC arm: span=127369/127386
    stem: list="gtgag"
    loop: list="ttcaacc" cost=73
    ~stem: list="ctcac" cost=0
    ~stem: list="tggggctc" cost=0.5

```

**Figure 2.** The first two parse trees returned by the tRNA grammar parsing chromosome III of *S. cerevisiae*. Note the intron and extra arm in the first, absent in the second.

recognized in the parse, as horizontal lines extending over the entire scope of the nonterminal and labelled by a tag associated with that nonterminal. Successively higher nonterminals are drawn in similar fashion above these as the parse progresses, in a bottom-up manner until a parse is completed. Which nonterminals are tagged and when they are drawn in the course of the parse are specified by the user with control attributes in the grammar; typically there is a tradeoff between the “informativeness” of the developing parse tree and the speed of the parse, and a useful development paradigm is to refine a grammar with a large amount of annotation and then to run the resulting specification with much less real-time annotation, or even batch-wise without benefit of the graphical interface.

The bottom left button causes the parser to backtrack and find another parse, and the adjacent buttons allow for variations on the usual backtracking scheme, as noted above. The user has the option of backtracking into not only the end, but also the beginning of any feature in the parse tree, so as to force it to move forward on the input string, and it can even be made to skip to any arbitrary position. The feature to be moved can be selected from a menu or by buttoning directly on the parse tree, and the desired position can be selected by buttoning on the sequence.

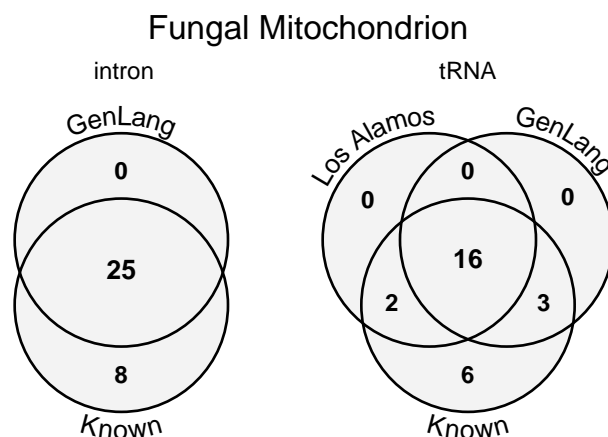


**Figure 3.** Results of batch runs of the tRNA grammar and the Los Alamos tRNA finder<sup>18</sup> against three GenBank (Release 71.0) divisions. The three-way Venn diagrams indicate the number of “hits” found by each program, and annotated in GenBank (except for the bacterial division, where only program hits were confirmed), with overlaps among the three sources indicated. Most “unannotated” hits were simply cases where two hits overlapped, and only one was correct, rather than outright false positives. At least one GENLANG-only hit appears to be an authentic unannotated tRNA, found on the opposite strand of an entry with two other confirmed tRNAs.

#### 4.7. Selected Results

Using a rudimentary protein-coding gene grammar, GENLANG was able to “discover” all five genes in the human beta globin cluster, in a region of 73,000 bp, with no false hits<sup>15</sup>; it also ignored the pseudogenes in this region, since they are syntactically invalid genes. Using relatively simple weight matrices for splice junctions and no compositional information about coding vs. non-coding regions, this grammar does not provide sufficient discrimination to reliably predict genes with more than a few introns, due to the combinatorial explosion of possibilities. Nevertheless, this nondeterminism of the parser proved to be useful in modeling thalassemia mutations in which alternative splicing results in a number of untranslatable messages also predicted by the parser<sup>15</sup>. Our current efforts in this area involve the incorporation of more sophisticated statistical and compositional measures into the syntactic framework of the grammar.

An *E. coli* tRNA grammar developed by us<sup>10</sup> was subsequently used by George Michaels at the NIH to parse all available *E. coli* sequence data (about 3,700,000 bp in both directions). After making one minor change to the grammar (loosening a constraint on a single base), the grammar found all 34 known tRNA genes in about three hours of total parse time, and was able to correct the standard *E. coli* map (EcoSeq<sup>16</sup> Version 5 in four cases (two sites wrong and two anticodons wrong) [G. Michaels, personal communication]. A generalized tRNA grammar, which was more loosely specified and which handled introns as well, was recently used by our group to parse all of yeast (*S. cerevisiae*) chromosome 3, over 315,000 bp, taking about eleven minutes in each direction (Figure 2). This parse found exactly the ten known



**Figure 4.** Results of parsing the mitochondrial genome of *P. anserina* with the Group I intron and tRNA grammars. While neither tRNA finder performed as well on the organelle division as on others, the GENLANG grammar was easily modified to detect five additional tRNAs on this genome (a total of 24/27), with no increase in false positives.

tRNA genes on this chromosome<sup>17</sup>; surprisingly, it was actually faster than a highly successful special-purpose tRNA finder from Los Alamos<sup>18</sup> (which however is faster with GenBank flat files). When run on a variety of divisions from GenBank, this grammar demonstrated both sensitivity and specificity well in excess of 95% when measured against known, annotated tRNA genes, and in several cases (like the Los Alamos program) apparently found unannotated tRNA genes (Figure 3).

The grammar is able to easily implement the constraint that tRNA introns contain within them a triplet which base-pairs with the anticodon, before splicing occurs. Although not done in the current grammar, it would also be relatively straightforward to incorporate postulated pseudoknot-like covariances between the D-arm and TΨC-arm loops, or the fact that the introns and extra arms are required to have extensive internal base-pairing. Obviously, for purposes of search it is also easy to add contextual clues, such as the fact that these genes tend to terminate in T-rich regions, or in yeast are found in proximity to certain transposable elements<sup>17</sup>.

Our most recent efforts in grammar design involve modelling of the self-splicing Group I introns, which compared to tRNA show a much more complex and variable secondary structure (including pseudoknots) with superimposed conserved regions<sup>19</sup>. This grammar was used to parse the 100,000 bp mitochondrial genome of the fungus *P. anserina*, with the results shown (Figure 4). We have parsed the same genome with a combination of the two grammars described here, plus a simple grammar for long open reading frames<sup>15</sup>, with results (Figure 5) that suggest that GENLANG may be a useful tool for assisting in the annotation of genomic sequence data.

**Figure 5 (overleaf).** Parse tree for the entire mitochondrial genome of *P. anserina* using a grammar composed of a disjunction of the Group I intron, tRNA, and ORF grammars. The parse took exactly two hours and produced most of the annotations found in published maps<sup>20</sup>.



## 5. Bibliography

1. W. Ebeling and M. A. Jimenez-Montano, *Math. Biosci.* **52** (1980) 53.
2. M. A. Jimenez-Montano, *Bull. Math. Biol.* **46(4)** (1980) 641.
3. V. Brendel, J. S. Beckmann, and E. N. Trifinov, *J. Biomol. Struct. Dynamics* **4** (1986) 11.
4. P. A. Pevzner, M. Y. Borodovsky, and A. A. Mironov, *J. Biomol. Struct. Dynamics* **6** (1989) 1013.
5. V. Brendel and H. G. Busse, *Nucleic Acids Res.* **12** (1984) 2561.
6. J. Collado-Vides, *Comput. Applic. in the Biosciences* **7(3)** (1991) 321.
7. D. B. Searls, in *Proceedings of the National Conference on Artificial Intelligence* (American Association for Artificial Intelligence, 1988), p. 386.
8. K. S. Fu, *Syntactic Pattern Recognition and Applications* (Prentice-Hall, Inc., Englewood Cliffs, NJ, 1982).
9. D. B. Searls, *Intelligent Systems Rev.* **1(4)** (1989) 67.
10. D. B. Searls and S. Liebowitz, in *Proceedings of the Workshop on Syntactic and Structural Pattern Recognition* (IAPR, 1990), p. 402.
11. F. C. N. Pereira and S. M. Shieber, *Prolog and Natural-Language Analysis* (Center for the Study of Language and Information, Stanford CA, 1987).
12. D. B. Searls, in *Logic Programming: Proceedings of the North American Conference*, eds. E. Lusk and R. Overbeek (MIT Press, 1989) p. 189.
13. D. B. Searls, in *Artificial Intelligence and Molecular Biology*, ed. L. Hunter (AAAI Press, 1992), p. 47.
14. J. P. Abrahams, M. van den Berg, E. van Batenburg, and C. Pleij, *Nucleic Acids Res.* **18** (1990) 3035.
15. D. B. Searls and M. O. Noordewier, in *Proceedings of the Conference on Artificial Intelligence Applications* (IEEE, 1991), p. 3.
16. K. E. Rudd, W. Miller, J. Ostell, and D. A. Benson, *Nucleic Acids Res.* **18** (1990) 313.
17. S. G. Oliver et al., *Nature* **357** (1992) 38.
18. G. A. Fichant and C. Burks, *J. Mol. Biol.* **220** (1991) 659.
19. T. R. Cech, *Gene* **73** (1988) 259.
20. D. J. Cummings, K. L. McNally, J. M. DoMenico, and E. T. Matsuura, *Curr. Genet.* **17** (1990) 375.