

Stochastic grammars for transposon sequences

Ian Holmes

Abstract

TAGs and parsers for biological repeats.

Contents

1	Introduction	2
2	Definitions	3
2.1	Tree-Adjoining Grammars	3
2.1.1	Context-free grammars as a subset of TAGs	4
2.1.2	Derivation weights	5
2.1.3	Length distributions	6
2.2	A bottom-up (CYK) parsing algorithm	6
2.2.1	The dynamic programming matrix	7
2.2.2	The dynamic programming recursion	9
2.2.3	Fill order	9
2.3	A top-down parsing algorithm	10
3	A simple retrotransposon grammar	10
3.1	SCFG and LTR components	10
3.2	Developing the SCFG sub-grammar for transposon contents . . .	11
3.3	Supplying external hints	11
4	A simple DNA transposon grammar	12
5	Glossary of mathematical notation	13

6 Acknowledgments	14
--------------------------	-----------

7 References	14
---------------------	-----------

1 Introduction

Transposable elements (TEs), or *transposons*, are of great interest in molecular evolution [1], and an important aspect of genome annotation. There are several specializations in the overall task of transposon annotation: PILER [2] specializes in *de novo* transposon discovery, while REPCLASS specializes in classification of found transposons [3].

Many such programs classify TEs by their general structural features, particularly their terminal repeats: LTRs (Long Terminal Repeats) and TIRs (Terminal Inverted Repeats).

It is useful to build databases and profiles of known transposon families, for the purpose of classifying new ones. To date, the most comprehensive database of known transposons is REPBASE [4], whose profiles rely only on primary sequence homology models; that is, they do not make explicit use of terminal repeat structure.

A promising approach, that combines profile Hidden Markov Models (HMMs) of primary sequence homology (at the level of TE protein domains) with fast algorithms for detecting LTRs, is taken by LTRdigest [5]. The purpose of this paper is to represent the hybrid modeling approach of LTRdigest using formal grammars.

2 Definitions

2.1 Tree-Adjoining Grammars

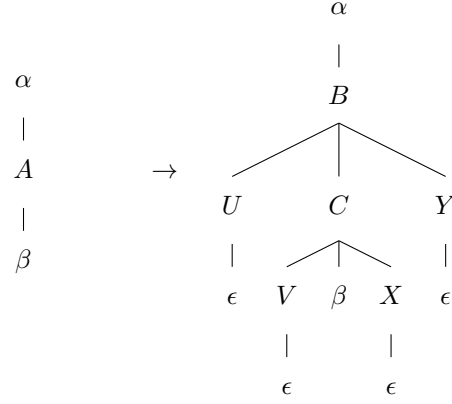
We define a minimal normal form of Tree-Adjoining Grammars (TAGs) suited to simplified presentation for biological sequence analysis, departing somewhat from the usual linguistic representation [6]. TAGs have previously been used in bioinformatics to model pseudoknots and other RNA structures [7, 8] and to model local duplications (Hickey and Blanchette, pers. comm.).

A TAG is a tuple $\mathcal{G} = (\mathcal{N}, \mathcal{T}, S, \mathcal{R}, \mathcal{W})$ where \mathcal{N} is a finite set of *node labels* (or *nonterminals*), \mathcal{T} is a finite set of *terminals* (disjoint from \mathcal{N}), $S \in \mathcal{N}$ is a distinguished *start label*, \mathcal{R} is a finite set of *transformation rules* and $\mathcal{W} : \mathcal{R} \rightarrow [0, \infty)$ is a *rule weight function*.

The process of generating an output sequence $Z \in \mathcal{T}^*$ using \mathcal{G} is referred to as a *derivation* of Z . The derivation consists of repeated local application of transformation rules to *intermediate trees*. We can commence a derivation from any *initial tree*, but ultimately we will be interested in derivations that commence with the following initial tree

$$\begin{array}{c} \epsilon \\ | \\ S \\ | \\ \epsilon \end{array}$$

Each “intermediate tree” is, formally, an ordered tree whose nodes are labeled from $(\mathcal{N} \cup \mathcal{T}^*)$. The transformation rules can take the various forms shown in Table 1. All the rules may be recognized as special cases of the general form



where $U, V, X, Y \in (\mathcal{N} \cup \mathcal{T}^*)$ and $B, C \in (\mathcal{N} \cup \{\epsilon\}^*)$. (The tree is ordered left-to-right.) Parsing this general form using a bottom-up, CYK approach would take $\mathcal{O}(|Z|^{10})$ time. If the rules are restricted to the types shown in Table 1, CYK parsing becomes $\mathcal{O}(|Z|^6)$ in time.

The derivation stops when no further transformations can be applied. The trees generated by this process have the property that every leaf node is labeled with a terminal sequence, while every internal node is labeled either with ϵ , or with a member of \mathcal{N} which never appears on the left-hand side of a rule in \mathcal{R} (and to which no transformations can therefore be applied).

To obtain the final output sequence, we read off the terminal sequences (at the leaves) in a depth-first traversal of the ordered tree (i.e. from just left of the root, moving anticlockwise around the tree, to just right of the root). The depth-first order of leaves in the general tree shown above is (U, V, β, X, Y) .

2.1.1 Context-free grammars as a subset of TAGs

We can further restrict the rules of Table 1 to recover the basic repertoire of a context-free grammar, or the weighted (stochastic) version, the SCFG. For example, SCFG rules of the form $A \rightarrow V X$, where A, V, X are nonterminals,

may be interpreted as bifurcation rules in Table 1, while SCFG rules of the form $A \rightarrow u \ C \ y$, where u, y are terminal strings and C is either a nonterminal or the empty string, may be interpreted as emission rules in Table 1.

2.1.2 Derivation weights

The *weight* of a derivation is the product of the weights of all rules used in the derivation. The weight of a given output sequence Z is the sum of the weights of all possible derivations of Z . Let $\mathcal{W}[\alpha \Rightarrow Z]$ denote the weight of output sequence Z , starting from initial tree α .

The final column of Table 1 shows a shorthand representation of each rule in Newick format, probably the most widely-understood bioinformatics format for representing tree structures. The initial tree, in this representation, is $((\epsilon)S)\epsilon$. In the table, we have omitted the placeholder ϵ 's at leaf nodes, so that (for example) the rule

$$((\beta)A)\alpha \rightarrow (V, \beta, X)\alpha$$

should strictly be read as

$$((\beta)A)\alpha \rightarrow ((\epsilon)V, \beta, (\epsilon)X)\alpha$$

Let $\mathcal{R}_n(A)$ denote the subset of rules of type n (according to Table 1) where the left-hand side is $((\beta)A)\alpha$ for a given $A \in \mathcal{N}$.

Note the special case of an emission rule where all the terminal strings are empty, $((\beta)A)\alpha \rightarrow ((\beta)C)\alpha$. Such a rule is referred to as a *transition* and may be written more compactly as $A \rightarrow C$.

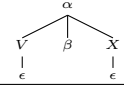

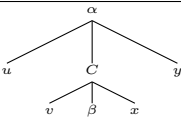
Type	From	To	Newick representation
(BIFURC)	$\begin{array}{c} \alpha \\ \\ A \\ \\ \beta \end{array}$ \rightarrow		$((\beta)A)\alpha \rightarrow (V, \beta, X)\alpha$
(ADJOIN)	$\begin{array}{c} \alpha \\ \\ A \\ \\ \beta \end{array}$ \rightarrow		$((\beta)A)\alpha \rightarrow (((\beta)C)B)\alpha$
(EMIT)	$\begin{array}{c} \alpha \\ \\ A \\ \\ \beta \end{array}$ \rightarrow		$((\beta)A)\alpha \rightarrow (u, ((v, \beta, x)C), y)\alpha$

Table 1: Types of transformation rule (i.e. tree adjunction rule) used in this paper. Here α, β represent any subtree; $A \in \mathcal{N}$ is the source node label; $B, C, V, X \in (\mathcal{N} \cup \{\epsilon\})$ are the destination node labels; ϵ is the empty string; and $u, v, x, y \in \mathcal{T}^*$ are (possibly empty) terminal strings.

2.1.3 Length distributions

To model genomic features, it is extremely useful to augment the TAG framework with length distributions; that is, to allow the rule weight $\mathcal{W}(\rho)$ to depend on distances such as $i, |Z| - l, L_0 = l - i, L_1 = k - j, j - i, l - k, m - i, j - m, n - k, l - n$ or $n - m$.

Note that the TAG framework implicitly allows some length distributions to be constructed, even without such augmentation; for example, the “waiting time” in a particular state is geometrically-distributed, and by arranging a series of such states, one can obtain other distributions, e.g. a negative binomial distribution. One can also obtain any distribution over a random variable $x \leq N$ by using N states, or approximate any distribution to some degree of precision with a finite number of states. However, it is sometimes more convenient simply to allow \mathcal{W} to depend directly on (i, j, k, l, m, n) .

2.2 A bottom-up (CYK) parsing algorithm

We can define a general parsing algorithm for TAGs that is the equivalent of the CYK (Cocke-Younger-Kasami) algorithm for SCFGs. The CYK algorithm is

closely related to the Inside algorithm, which computes the weight (probability) of a given output sequence. Here, we present the Inside version of the TAG algorithm; to obtain the CYK version, simply replace all summation operators (Σ) with max operators.

For the given output sequence $Z \in \mathcal{T}^*$, let $Z[i \dots j + 1]$ denote the substring from i through j inclusive, for $1 \leq i \leq j \leq |Z|$. Let $Z[i \dots i] = \epsilon$.

Define some indicator functions to match output substrings to emit rules

$$\begin{aligned}
 \Delta(i, j, x) &= \delta(Z[i \dots j] = x) \\
 \vec{\Delta}(i, x) &= \Delta(i, i + |x|, x) \\
 \overleftarrow{\Delta}(j, y) &= \Delta(j - |y|, j, y) \\
 \Delta_{iu} &= \vec{\Delta}(i, u) \\
 \Delta_{jv} &= \overleftarrow{\Delta}(j, v) \\
 \Delta_{kx} &= \vec{\Delta}(k, x) \\
 \Delta_{ly} &= \overleftarrow{\Delta}(l, y)
 \end{aligned}$$

2.2.1 The dynamic programming matrix

Dynamic programming is breaking down a problem recursively into a lot of smaller problems, whose solutions can be tabulated and re-used in computing the bigger problem. Here, the big problem is to find the probability (weight) that the initial symbol S will derive the output sequence Z . The smaller problems are the probabilities that a given symbol A will derive sequences $Z[i \dots j]$ to the left and $Z[k \dots l]$ to the right.

Introducing the additional terminal γ as a placeholder for the sequence

$Z[j \dots k]$, define

$$\begin{aligned} M(i, j, k, l, A) &= \mathcal{W}[(\gamma)A)\epsilon \Rightarrow Z[i \dots j]\gamma Z[k \dots l]] \\ M'(i, l, A) &= \mathcal{W}[(\epsilon)A)\epsilon \Rightarrow Z[i \dots l]] \end{aligned}$$

where $1 \leq i \leq j \leq k \leq l \leq |Z|$.

Thus $M(i, j, k, l, A)$ is the probability that the initial tree $((\gamma)A)\epsilon$ will generate $Z[i \dots j]$ to the left of γ and $Z[k \dots l]$ to the right of γ , whereas $M'(i, l, A)$ is the probability that the initial tree $((\epsilon)A)\epsilon$ will generate $Z[i \dots l]$. Note that

$$M'(i, l, A) = \sum_{k=i}^l M(i, k, k, l, A)$$

Note also the following boundary conditions:

$$\begin{aligned} M(i, i, k, k, \epsilon) &= 1 \\ M'(i, i, \epsilon) &= 1 \end{aligned}$$

The *inside sequence* refers to the sequences $Z[i \dots j]$ and $Z[k \dots l]$ (for M), or the sequence $Z[i \dots l]$ (for M').

The final derivation weight is

$$\mathcal{W}[S \Rightarrow Z] = M'(1, |Z|, S)$$

2.2.2 The dynamic programming recursion

The DP recursion relation is as follows (with B, C, V, X, u, v, x, y defined as in Table 1)

$$\begin{aligned}
M(i, j, k, l, A) = & \sum_{\rho \in \mathcal{R}_{BIF}(A)} W(\rho) M'(i, j, V) M'(k, l, X) \\
& + \sum_{\rho \in \mathcal{R}_{ADJ}(A)} \sum_{m=i}^j \sum_{n=k}^l W(\rho) M(i, m, n, l, B) M(m, j, k, n, C) \\
& + \sum_{\rho \in \mathcal{R}_{EMT}(A)} W(\rho) M(i + |u|, j - |v|, k + |x|, l - |y|, C) \Delta_{iu} \Delta_{jv} \Delta_{kx} \Delta_{ly}
\end{aligned}$$

Note that for M to be exactly computable, we require that there are no *null cycles* in the grammar. A null cycle is a series of transformations that, when applied consecutively to a given tree α , yield the original tree α again, or one that is trivially related to it (e.g. a tree that is identical to α after ϵ -labeled nodes of degree < 3 have been removed). Null cycles frequently arise via consecutive transitions of the form $A \rightarrow B$ followed by $B \rightarrow A$, but there are other possibilities too (for example, a bifurcation rule $((\beta)A)\alpha \rightarrow (V, \beta, X)\alpha$ followed by two transitions $V \rightarrow \epsilon$ and $X \rightarrow A$).

2.2.3 Fill order

The absence of null cycles implies that we can perform a *topological sort* (topo-sort) of \mathcal{N} based on the transition graph, yielding the order in which node labels must be visited in the Inside algorithm.

A suitable (but non-unique) Inside fill order is as follows:

- For $L_0 = 0$ to $|Z|$ (increasing), $i = 0$ to $|Z| - L_0$ (increasing):
 - Let $l = i + L_0$

- For $L_1 = L_0$ to 0 (decreasing), $j = i$ to $i + L_0 - L_1$ (increasing):
 - * Let $k = j + L_1$
 - * For $A \in \mathcal{N}$ (topo-sorted, transition destinations before sources):
 - Compute $M(i, j, k, l, A)$ and store
- For $A \in \mathcal{N}$ (topo-sorted, transition destinations before sources):
 - * Compute $M'(i, l, A)$ and store

2.3 A top-down parsing algorithm

The general idea is to try a top-down chart parser, such as the Earley parser [9] (a version of which exists for TAGs [10]) that can run very efficiently on preprocessed sequence, where the preprocessor augments the sequence with a table of tuples of the form (i, j, k, l, A, M) and the parser restricts the iteration to derivations that use some of these (i, j, k, l, A) co-ordinates.

[WRITE ME!]

3 A simple retrotransposon grammar

We are particularly interested in following class of grammars that describe specific arrangements of DNA-encoded protein domains flanked by LTRs (long terminal repeats).

3.1 SCFG and LTR components

The LTRs can be generated by emission rules (e.g. with $v = y = \epsilon$) followed either by a transition to the transposon interior, or (optionally) an adjunction rule, which may be an easier event on which to impose a hint-requirement constraint.

3.2 Developing the SCFG sub-grammar for transposon contents

- X generates a nucleotide sampled from the background distribution
- X_L generates ℓ background nucleotides, where $\ell \sim L$
- F_N samples a DNA sequence coding for family N from PFAM [11]
- I_L generates an intron of length $\ell \sim L$ (can be emitted by F_N)
- T_A generates a terminal inverted repeat, then transits to A

We can also make transitions back to S to generate a nested transposon insertion.

The grammar so described has some similarities to LTRdigest [5] and TENest [12].

3.3 Supplying external hints

The parsing algorithm uses $\mathcal{O}(|Z|^4)$ memory and $\mathcal{O}(|Z|^6)$ time. The hope is to accelerate it significantly by using externally-supplied “hints” as constraints on the locations of various features (especially the LTRs).

The hints file should include

- A set of tuples (i, j, k, l) indicating that $Z[i \dots j]$ and $Z[k \dots l]$ are (respectively) the 5' and 3' repeat regions of an LTR
- A set of tuples (i, j, N) indicating that $Z[i \dots j]$ is a match to PFAM family N

These hints can be generated by fast tools, e.g. suffix-tree based algorithms for finding LTRs [13], or GeneWise for finding DNA sequences that code for PFAM protein domains [14].

A very quick heuristic that might achieve most of the benefits of a more rigorous “hints” constraint would be to divide the genome into windows and only run the grammar on windows which contain K or more of the appropriate hints.

4 A simple DNA transposon grammar

[WRITE ME!]

5 Glossary of mathematical notation

Symbol	Meaning
\mathcal{G}	Grammar
\mathcal{N}	Set of node labels (nonterminals)
\mathcal{T}	Set of terminals
\mathcal{R}	Set of transformation rules
$\mathcal{R}_T(A)$	Set of type- T transformation rules (see Table 1) with $((\beta)A)\alpha$ on LHS
α, β	Arbitrary ordered node-labeled trees
\mathcal{W}	Rule weight function
$\mathcal{W}[\alpha \Rightarrow Z]$	Sum, over derivations of Z from α , of product over derivation rule weights
ϵ	The empty string
\mathcal{T}^*	Set of strings over \mathcal{T} , including the empty string
Z	Output sequence
$ Z $	Length of Z
i, j, k, l	Integer co-ordinates on Z , with $1 \leq i \leq j \leq k \leq l \leq Z $
m, n	Integer co-ordinates on Z , with $i \leq m \leq j$ and $k \leq n \leq l$
L_0, L_1	Integer lengths of Z -subsequences, with $0 \leq L_1 = (k - j) \leq L_0 = (l - i) \leq Z $
$Z[i \dots i]$	The empty string
$Z[i \dots j]$	(For $j > i$) Substring of Z from i to $j - 1$ inclusive
$\Delta(i, j, x)$	Indicates if rule string x matches output string $Z[i \dots j]$
$\vec{\Delta}(i, x)$	Indicates if rule string x matches output string Z , starting at position i
$\overleftarrow{\Delta}(j, y)$	Indicates if rule string y matches output string Z , ending before position j
$\Delta_{iu}, \Delta_{jv}, \Delta_{kx}, \Delta_{ly}$	Rule string matchers; Δ_{ns} matches string s ending/starting at position n
$M(i, j, k, l, A)$	Inside weight for $Z[i \dots j]$ (left) and $Z[k \dots l]$ (right) rooted at A
$M'(i, l, A)$	Inside weight for $Z[i \dots l]$ rooted at A
B, C, U, V, X, Y	Labels (nonterminal symbols, terminal strings, or ϵ) at nodes of the Newick tree $((U, ((V, X), Y)))$
u, v, x, y	Special symbols for U, V, X, Y when they are terminal strings

6 Acknowledgments

Sean Eddy co-organized an excellent workshop at U.Penn, “Language Modeling of Biological Data”, at which IH met Mark Steedman, Bonnie Webber and Aravind Joshi, and learned from them (and others) the secrets of TAGs, finite-state transducers, and mildly context-sensitive grammars. David Searls has been a long-time proponent of linguistics applied to bioinformatics. Dan Klein encouraged us to develop our own variations on grammar frameworks. Michael Souza introduced us to the pumping lemma and the equivalence of TAGs with various other frameworks. Lawrence Uricchio helped revise these ideas.

7 References

References

1. C. Feschotte. DNA transposons and the evolution of eukaryotic genomes. *Annual Review of Genetics*, 41:331–368, 2007.
2. R. C. Edgar and E. W. Myers. PILER: identification and classification of genomic repeats. *Bioinformatics*, 21 Suppl 1:i152–8, 2005.
3. C. Feschotte, U. Keswani, N. Ranganathan, M. L. Guibotsy, and D. Levine. Exploring repetitive DNA landscapes using REPCCLASS, a tool that automates the classification of transposable elements in eukaryotic genomes. *Genome Biol Evol*, 1:205–220, 2009.
4. V. V. Kapitonov and J. Jurka. A universal classification of eukaryotic transposable elements implemented in Repbase. *Nature reviews. Genetics*, 9(5):411–2; author reply 414, 2008.

5. S. Steinbiss, U. Willhoeft, G. Gremme, and S. Kurtz. Fine-grained annotation and classification of de novo predicted LTR retrotransposons. *Nucleic Acids Res.*, 37:7002–7013, Nov 2009.
6. A. Joshi and Y. Schabes. Tree-adjoining grammars, 1997.
7. H. Matsui, K. Sato, and Y. Sakakibara. Pair stochastic tree adjoining grammars for aligning and predicting pseudoknot RNA structures. *Bioinformatics*, 21:2611–2617, Jun 2005.
8. D. Chiang, A. K. Joshi, and D. B. Searls. Grammatical representations of macromolecular structure. *J. Comput. Biol.*, 13:1077–1100, Jun 2006.
9. J. Earley. An efficient context-free parsing algorithm. *Communications of the ACM*, 13(2):94–102, 1970.
10. L. Kallmeyer. *Parsing Beyond Context-Free Grammars*. Cognitive Technologies. Springer, 2010.
11. R. D. Finn, J. Tate, J. Mistry, P. C. Coghill, S. J. Sammut, Hans-Rudolf Hotz, G. Ceric, K. Forslund, S. R. Eddy, E. L. L. Sonnhammer, and A. Bateman. The Pfam protein families database. *Nucleic Acids Research*, 36(Database issue):D281–8, 2008.
12. B. A. Kronmiller and R. P. Wise. TEneST: automated chronological annotation and visualization of nested plant transposable elements. *Plant Physiol.*, 146:45–59, Jan 2008.
13. A. Kalyanaraman and S. Aluru. Efficient algorithms and software for detection of full-length LTR retrotransposons. *J Bioinform Comput Biol*, 4:197–216, Apr 2006.
14. E. Birney, M. Clamp, and R. Durbin. GeneWise and GenomeWise. *Genome Research*, 14(5):988–995, 2004.