

# Stochastic tree-adjoining grammars for modeling retrotransposons

Lawrence Uricchio, Ian Holmes

## Abstract

TAGs and parsers for biological repeats.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Definitions</b>	<b>2</b>
2.1	Tree-Adjoining Grammars . . . . .	2
2.2	A parsing algorithm . . . . .	5
2.2.1	The dynamic programming matrix . . . . .	5
<b>3</b>	<b>A simple retrotransposon grammar</b>	<b>5</b>
3.1	SCFG and LTR components . . . . .	6
3.2	Developing the SCFG sub-grammar for transposon contents . . .	6
3.3	Supplying external hints . . . . .	6
<b>4</b>	<b>Glossary of mathematical notation</b>	<b>7</b>
<b>5</b>	<b>Acknowledgments</b>	<b>8</b>
<b>6</b>	<b>References</b>	<b>8</b>

## 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 biological sequence analysis, as opposed to the linguistic representation elsewhere [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 set of *node labels*,  $\mathcal{T}$  is a set of *terminals* (disjoint from  $\mathcal{N}$ ),  $S \in \mathcal{N}$  is a distinguished *start label*,  $\mathcal{R}$  is a 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*, beginning with the 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.

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  $\epsilon$ , 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.

The final column of the above table 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  $\epsilon$ ’s at leaf nodes, so that (for example) the rule

$$((\beta)A)\alpha \rightarrow ((C, \beta)B)\alpha$$

Type	From		To	Newick representation
(1)	$\begin{array}{c} \alpha \\   \\ A \\   \\ \beta \end{array}$	$\rightarrow$	$\begin{array}{c} \alpha \\   \\ B \\ / \quad \backslash \\ C \quad \beta \\   \\ \epsilon \end{array}$	$((\beta)A)\alpha \rightarrow ((C, \beta)B)\alpha$
(2)	$\begin{array}{c} \alpha \\   \\ A \\   \\ \beta \end{array}$	$\rightarrow$	$\begin{array}{c} \alpha \\   \\ B \\ / \quad \backslash \\ \beta \quad C \\   \\ \epsilon \end{array}$	$((\beta)A)\alpha \rightarrow ((\beta, C)B)\alpha$
(3)	$\begin{array}{c} \alpha \\   \\ A \\   \\ \beta \end{array}$	$\rightarrow$	$\begin{array}{c} \alpha \\ / \quad \backslash \\ C \quad B \\   \quad   \\ \epsilon \quad \beta \end{array}$	$((\beta)A)\alpha \rightarrow (C, ((\beta)B))\alpha$
(4)	$\begin{array}{c} \alpha \\   \\ A \\   \\ \beta \end{array}$	$\rightarrow$	$\begin{array}{c} \alpha \\ / \quad \backslash \\ B \quad C \\   \quad   \\ \beta \quad \epsilon \end{array}$	$((\beta)A)\alpha \rightarrow (((\beta)B), C)\alpha$
(5)	$\begin{array}{c} \alpha \\   \\ A \\   \\ \beta \end{array}$	$\rightarrow$	$\begin{array}{c} \alpha \\   \\ B \\   \\ C \\   \\ \beta \end{array}$	$((\beta)A)\alpha \rightarrow (((\beta)C)B)\alpha$
(6)	$\begin{array}{c} \alpha \\   \\ A \\   \\ \beta \end{array}$	$\rightarrow$	$\begin{array}{c} \alpha \\ / \quad   \quad \backslash \\ u \quad B \quad x \\ / \quad   \quad \backslash \\ v \quad \beta \quad w \end{array}$	$((\beta)A)\alpha \rightarrow (u, ((v, \beta, w)B), x)\alpha$

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

should strictly be read as

$$((\beta)A)\alpha \rightarrow (((\epsilon)C, \beta)B)\alpha$$

Let  $\mathcal{R}_n \subseteq \mathcal{R}$  denote the subset of rules of type  $n$  according to Table 1.

## 2.2 A parsing algorithm

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

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 rules

$$\begin{aligned} \Delta(i, j, x) &= \delta(Z[i \dots j] = x) \\ \Delta(i, x) &= \Delta(i, i + |x|, x) \end{aligned}$$

### 2.2.1 The dynamic programming matrix

more to go here

## 3 A simple retrotransposon grammar

An example of an observed indexed grammar is the following class of grammars that can generate context-free structures flanked by LTRs (long terminal repeats).

### 3.1 SCFG and LTR components

### 3.2 Developing the SCFG sub-grammar for transposon contents

- $X$  generates a nucleotide sampled from the background distribution
- $X_L$  generates  $\ell$  background nucleotides, where  $\ell \sim L$
- $F_N$  samples a DNA sequence coding for family  $N$  from PFAM [9]
- $I_L$  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 [10].

Note a flaw of the framework is that, since the index string representing the “consensus” LTR must be directly observed at least once in the output (in fact, it corresponds to the 5’-most repeat in this grammar), we cannot allow a nested transposon insertion *within* that LTR.

### 3.3 Supplying external hints

The parsing algorithm uses  $\mathcal{O}(|Z|^4)$  memory and  $\mathcal{O}(|Z|^4)$  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 [11], or GeneWise for finding DNA sequences that code for PFAM protein domains [12].

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 Glossary of mathematical notation

Symbol	Meaning
$\mathcal{G}$	Grammar
$\mathcal{N}$	Set of node labels
$\mathcal{T}$	Set of terminals
$\mathcal{R}$	Set of transformation rules
$\mathcal{R}_n$	Set of transformation rules of type $n$ , according to Table 1
$\mathcal{W}$	Rule weight function
$\epsilon$	The empty string
$\mathcal{T}^*$	Set of strings over $\mathcal{T}$ , including the empty string
$Z$	Output sequence
$ Z $	Length of $Z$
$Z[i \dots j + 1]$	Substring of $Z$ from $i$ to $j$ inclusive ( $i$ starts at 1)
$Z[i \dots i]$	The empty string
$M(i, j, k, l, X)$	Max parse tree weight for $Z[k \dots l]$ rooted at $X[Z[i \dots j]]$
$\Delta(i, j, x)$	Indicates if rule string $x$ matches output string $Z[i \dots j]$
$\Delta(i, x)$	Indicates if rule string $x$ matches output string $Z$ , starting at position $i$

## 5 Acknowledgments

Our understanding of grammars applied to bioinformatics has been uplifted in conversations with Sean Eddy, Mark Steedman, Bonnie Webber, Aravind Joshi, David Searls, Dan Klein, and Michael Souza.

## 6 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.
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.
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. 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.
10. B. A. Kronmiller and R. P. Wise. TEneest: automated chronological annotation and visualization of nested plant transposable elements. *Plant Physiol.*, 146:45–59, Jan 2008.
11. 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.
12. E. Birney, M. Clamp, and R. Durbin. GeneWise and GenomeWise. *Genome Research*, 14(5):988–995, 2004.