# Evolving state grammar for modeling DNA and RNA structures

**Ajay Kumar \*, Nidhi Kalra, Sunita Garhwal**

*CSED, Thapar University, Patiala, India,147004*
*\*Corresponding author E-mail: ajayloura@gmail.com*

### Abstract

In this paper, we represent bio-molecular structures (Attenuator, Extended Pseudoknot Structure, Kissing Hairpin, Simple H-type structure, Recursive Pseudoknot and Three-knot Structure) using state grammar. These representations will be measured using descriptional complexity point of views. Results indicate that the proposed approach is more succinct in terms of production rules and variables over the existing approaches. Another major advantage of the proposed approach is state grammar can be represented by deep pushdown automata, whereas no such automaton exists for matrix ins-del system.

*Keywords*: Descriptional Complexity; Kissing Hairpin; Pseudoknot; Ribonucleic Acid; Simple H-Type; State Grammar; Three-Knot.

## 1. Introduction

Natural computing combines the formal models and algorithmic techniques to solve the problem inspired by nature with the inclusion of natural materials (i.e. molecules). Deoxyribonucleic Acids (DNA) computing also comes under the umbrella of natural computing. Nowadays, a major concern of the bioinformatics is the analysis of DNA, Ribonucleic Acids (RNA) and protein sequences. DNA and RNA molecules form a complimentary pair which results in a pattern formation in the sequence. The grammatical formalism of these biological sequences is used in solving many bioinformatics problems such as multiple alignment calculations, classification, and prediction of primary and secondary structures. DNA and RNA are responsible for the development, growth and functioning of all known living organisms and viruses.

Chomsky grammar systems [1] are found to be ideal for representing the interactions of nucleotides as there is a similarity between formal languages and bio-molecules language. Some sequences like a hairpin (the language of palindrome) can be represented by a context-free grammar, whereas other sequences like attenuator $\{u\bar{u}^R u\bar{u}^R \mid u \in \Sigma^*_{DNA}\}$ cannot be represented by a context-free grammar. Similarly, other biological structures such as extended pseudoknot, recursive pseudoknot, simple-H type, kissing hairpin, Three-knot includes cross-dependency, and they require a higher class of formal language than context-free. In this paper, we will represent these bio-molecular structures using state grammar (a type of regulated grammar). The regulated grammar consists of production rules similar as context-free grammar, but certain restrictions are imposed on these grammars to represent the cross dependencies. State grammar is a rule-based regulated grammar in which restrictions are imposed in terms of states.

Prior Work: The concept of state grammar was introduced by Kasai [2]. A state grammar is a rule-based regulated grammar in which restrictions are imposed in terms of states. Various representations of DNA and RNA sequences using formal grammar and automata have been found in the literature. Sung [3] represented RNA secondary structure loops such as a hairpin loop, an internal loop, bulge loop and double helix using context-sensitive grammar.

Sakakibara et al. [4] modeled RNA structure loops using stochastic context-free grammar. Further, Sakakibara [5] modeled RNA structure loops using pair hidden Markov models. Yuki and Kasami [6] modeled RNA structure loops using stochastic multiple context-free grammars. Brown and Wilson [7] modeled RNA pseudoknot structure using the intersection of stochastic context-free grammars.

Rivas and Eddy [8] used cross-interaction grammar for representing RNA secondary structure. Searls [9] used indexed grammar to represent DNA and RNA sequences such as tandem repeat, inverted repeat and pseudoknot. Searls [10] also represented DNA sequences using string variable grammar. Mizoguchi et al. [11] used stochastic multiple context-free grammars to represent various classes of pseudoknots. Cai et al. [1] represented RNA pseudoknot structure using parallel communicating grammar systems. Kuppusamy et al. [12] represented DNA and RNA secondary structures using matrix insertion-deletion system. Kalra and Kumar [3] represented tandem repeat, inverted repeat and pseudoknot using state grammar.

In this paper, we analyze the representation of bio-molecular structures with the basic descriptional complexity in terms of a number of production rules and number of variables. Results are compared with the matrix insertion-deletion system for the similar representations. After introducing some preliminary concepts in Section 2, we represent attenuator, extended pseudoknot, H-type, three-knot structure, recursive pseudoknot structure and kissing hairpin using state grammar. Section 4 consists of results and discussion.

## 2. Preliminaries

In this section, some basic notations and definition are discussed. $\Sigma_D = \{g, c, a, t\}$ and $\Sigma_R = \{g, c, a, u\}$ denote DNA and RNA alphabet respectively. $\Sigma_D^*$ denotes the free monoid generated by $\Sigma_D$. $\lambda$ denotes empty string or null string. In DNA and RNA, pairing occurs between complement pair in purines and pyrimidine. The complement of a symbol $d$ is denoted by $\bar{d}$. Purines are

classified into adenine ($a$) and guanine ($g$), while pyrimidine is classified into cytosine ($c$) thymine ($t$) and uracil ($u$).

In DNA:

$$\bar{a} = t$$

$$\bar{t} = a$$

$$\bar{g} = c$$

$$\bar{c} = g$$

In DNA:

$$\bar{a} = u$$

$$\bar{u} = a$$

$$\bar{g} = c$$

$$\bar{c} = g$$

Watson-Crick pairing occurs in RNA. DNA and RNA are important macromolecules that exist in every form of life. They are made from monomers known as nucleotides. Each nucleotide consists of a pentose carbon sugar, a phosphate group, and a nitrogenous base. If the sugar is ribose, then the polymer is RNA. If the sugar is deoxyribose, then the polymer is DNA.

Def. 2.1: Regulated grammar [14-16] is quintuple $(N, \Sigma, S, P, RG)$ where $N$ is a set of non-terminals, $\Sigma$ is an alphabet, S is the start symbol, P is the set of production rules, and RG is the restriction applied on the derivations of strings, and it depends on the type of regulated grammar.

Regulated grammar is classified into rule-based and context-based grammatical regulation [14].

Def. 2.2 [13]: A state grammar is a quintuple $G(V, Q, \Sigma, P, S)$, where $V$ is a finite set of symbols, $Q$ is a finite set of states such that $V \cap Q = \phi$, $\Sigma \subseteq V$ is an alphabet of terminals, $P \subseteq (Q \times (V - \Sigma) \times (Q \times V^+))$ is a finite relation over the productions, and $S \in V - \Sigma$ is the start symbol.

Example 1: Consider the state grammar $G_s = (\{S, A, B, 0, 1, 2\}, \{p_0, p_1, p_2\}, \{0, 1, 2\}, P, S)$ where the production rules are

$$(p_0, S) \rightarrow (p_0, AB) \quad (p_0, A) \rightarrow (p_1, 0A\ 1) \quad (p_1, B) \rightarrow (p_0, 2B)$$

$$(p_0, A) \rightarrow (p_2, 01) \quad (p_2, B) \rightarrow (p_2, 2)$$

Consider the string $s = 000111222$

$$S_{p0} \rightarrow AB_{p0} \rightarrow 0A1B_{p1} \rightarrow 0A12B_{p0} \rightarrow 00A112B_{p1} \rightarrow$$

$$00A1122B_{p0} \rightarrow 000111122B_{p2} \rightarrow 000111222_{p2}$$

The non-context-free language generated by $G_s$ is $L(G_s) = \{0^n 1^n 2^n \mid n \geq 1\}$.

Def. 2.3 [17]: Matrix insertion-deletion system $I(V, \Sigma, A, R)$ where $V$ is a finite set of symbols, $\Sigma \subseteq V$ is an alphabet of terminals, A is a finite language over $V$, $R$ is a finite set of triple in a matrix format $[(u_1, \alpha_1 \mid \beta_1, v_1)...(u_n, \alpha_n \mid \beta_n, v_n)]$, $(u_i, v_i) \in V^* \times V^*$ and $(\alpha_i, \beta_i) \in (V^+ \times \{\lambda\}) \bigcup (\{\lambda\} \times V^+)$.

Def. 2.4 [17]: Given a matrix insertion-deletion system $I(V, \Sigma, A, R)$. Descriptional complexity measure of $I$ in terms of variables and production is defined by

$$prod(I) = \left( \sum_{m \in R} m + |R| \right). |A|,$$

$$var(I) = |V| - |T|$$

Here $|R|$ denote a total number of rules in a matrix insertion-deletion system.

Def. 2.5: Descriptional complexity of a state grammar $G(V, Q, \Sigma, P, S)$ is defined by

$$prod(G) = |P| + |Q|$$

$$var(I) = |V| - |T|$$

The descriptional complexity of example 1 is

$$prod(G) = |P| + |Q| = 5 + 3 = 8$$

And

$$var(I) = |V| - |T| = |\{S, A, B, 0, 1, 2\}| - |\{0, 1, 2\}| = 5 - 3 = 2.$$

Recently, various researchers had researched different direction of automata theory, especially focusing on deep pushdown automata and state grammar (See [13, 18-22] for more details)

## 3. State grammar for bio-molecular structures

This section describes the state grammar for bio-molecular sequences found in DNA and RNA.

Proposition 1. The attenuator language $L_A = \{u\bar{u}^R u\bar{u}^R \mid u \in \Sigma_{DNA}^*\}$ can be generated by state grammar. Fig 1. represents attenuator structure of the sequence *gtcgacagcgct*.



$u\bar{u}_R$                                        $u\bar{u}_R$

**Fig. 1:** Attenuator Structure.

Proof: Grammar $G_1 = \{(S, A, B, u, \bar{u}), (q_0, q_1, q_2), (u, \bar{u}), P, S\}$, where $u \in \{a, g, c, u\}$ and $\bar{u}$ is the complement of $u$. State grammar productions are defined as follows:

$$(q_0, S) \rightarrow (q_0, AB)$$

$$(q_0, A) \rightarrow (q_1, uA\bar{u})$$

$$(q_1, B) \rightarrow (q_0, uB\bar{u})$$

$$(q_0, A) \rightarrow (q_2, \lambda)$$

$$(q_2, B) \rightarrow (q_2, \lambda)$$

Derivation for input string $w = gtatacgtatac$

$$S_{q_0} \rightarrow AB_{q_0} \rightarrow gAcB_{q_1} \rightarrow gAcgBc_{q_0} \rightarrow gtAacgBc_{q_1} \rightarrow gtAacgtBac_{q_0}$$
$$\rightarrow gtaAtacgtBac_{q_1} \rightarrow gtaAtacgtaBtac_{q_0} \rightarrow gtatacgtaBtac_{q_2}$$
$$\rightarrow gtatacgtatac_{q_2}$$

**Proposition 2.** The extended pseudoknot language $L_{ExP} = \{u_1 v_1 \bar{u}_1^R u_2 \bar{v}_1^R \bar{u}_2^R \mid u_1, u_2, v_1 \in \Sigma_{RNA}^*\}$ can be generated by state grammar. Fig. 2 represents the structure of extended pseudoknot sequence $cugcuacagcguuagacg$.

Proof: Grammar $G_2 = \{(S, A, B, u_1, u_2, v_1, \bar{u}_1, \bar{u}_2, \bar{v}_1),$

$(q_0, q_1, q_2, q_3), (u_1, u_2, v_1, \bar{u}_1, \bar{u}_2, \bar{v}_1), P, S\}$ be the state grammar where production rules P are as follows:



**Fig. 2:** Extended Pseudo knot Structure.

$$(q_0, S) \rightarrow (q_0, AB)$$

$$(q_0, A) \rightarrow (q_0, u_1 A \bar{u}_1)$$

$$(q_0, B) \rightarrow (q_0, u_2 B \bar{u}_2)$$

$$(q_0, A) \rightarrow (q_1, v_1 A)$$

$$(q_1, B) \rightarrow (q_2, B \bar{v}_1)$$

$$(q_2, A) \rightarrow (q_1, v_1 A)$$

$$(q_2, A) \rightarrow (q_3, \lambda)$$

$$(q_3, B) \rightarrow (q_3, \lambda)$$

Derivation for input string $w = cugcuacagcguuagacg$

$$S_{q_0} \rightarrow AB_{q_0} \rightarrow cAgB_{q_0} \rightarrow cuAagB_{q_0} \rightarrow cugAcagB_{q_0} \rightarrow$$
$$cugAcagcBg_{q_0} \rightarrow cugAcagcgBcg_{q_0} \rightarrow cugAcagcguBacg_{q_0}$$
$$\rightarrow cugcAcagcguBacg_{q_1} \rightarrow cugcAcagcguBgacg_{q_2} \rightarrow$$
$$cugcuAcagcguBgacg_{q_1} \rightarrow cugcuAcagcguBagacg_{q_2} \rightarrow$$
$$cugcuaAcagcguBagacg_{q_1} \rightarrow cugcuaAcagcguBuagacg_{q_2} \rightarrow$$
$$cugcuacagcguBuagacg_{q_3} \rightarrow cugcuacagcguuagacg_{q_3}$$

Bold non-terminal indicates the non-terminal to be expanded next.
**Proposition 3.** The simple H-type language $L_{SH} = \{u_1 v_1 u_2 \bar{u}_2^R v_1 \bar{u}_2^R \mid u_1, u_2, v_1, v_2 \in \Sigma_{RNA}^*\}$ can be generated by state grammar. Fig. 3 represents the structure of the simple H-type sequence.



**Fig. 3:** Simple H-Type Structure.

Proof: Grammar $G_3 = (\{S, A, B, u_1, u_2, v_1, v_2, \bar{u}_1, \bar{u}_2\}, \{q_0, q_1, q_2, q_3, q_4\}, \{u_1, u_2, v_1, v_2, \bar{u}_1, \bar{u}_2\}, P, S)$ be the state grammar where production rules P are as follows:

$$(q_0, S) \rightarrow (q_0, AB)$$

$$(q_0, A) \rightarrow (q_0, u_1 A \bar{u}_1)$$

$$(q_0, A) \rightarrow (q_1, v_1 A)$$

$$(q_1, A) \rightarrow (q_1, v_1 A)$$

$$(q_1, B) \rightarrow (q_1, v_2 B)$$

$$(q_1, A) \rightarrow (q_2, u_2 A)$$

$$(q_2, B) \rightarrow (q_3, B \bar{u}_2)$$

$$(q_3, A) \rightarrow (q_2, u_2 A)$$

$$(q_3, A) \rightarrow (q_4, \lambda)$$

$$(q_3, B) \rightarrow (q_4, \lambda)$$

Derivation for input string $w = cugucugcagacag$

$$S_{q_0} \rightarrow AB_{q_0} \rightarrow cAgB_{q_0} \rightarrow cuAagB_{q_0} \rightarrow cugAcagB_{q_0} \rightarrow cuguAcagB_{q_1}$$
$$\rightarrow cuguAcagaB_{q_1} \rightarrow cugucAcagaB_{q_2} \rightarrow cugucAcagaBg_{q_3} \rightarrow$$
$$cugucuAcagaBg_{q_2} \rightarrow cugucuAcagaBag_{q_3} \rightarrow cugucugAcagaBag_{q_2} \rightarrow$$
$$cugucugAcagaBcag_{q_3} \rightarrow cugucugcagaBcag_{q_4} \rightarrow cugucugcagacag_{q_4}$$

**Proposition 4.** The simple three-knot language $L_{TK} = \{u_1 v u_2 u_3 \bar{u}_1^R \bar{u}_2^R \bar{u}_3^R \mid u_1, u_2, u_3, v \in \Sigma_{RNA}^*\}$ can be generated by state grammar. Fig 4. represents the structure of the three-knot sequence.



**Fig. 4:** Three-Knot Structure.

Proof:    Grammar    $G_4 = (\{S, A, B, C, u_1, u_2, u_3, v, \bar{u}_1, \bar{u}_2, \bar{u}_3\}, \{q_0, q_1, q_2, q_3, q_4, q_5\}, \{u_1, u_2, u_3, v, \bar{u}_1, \bar{u}_2, \bar{u}_3\}, P, S)$ be the state grammar where production rules P are as follows:

$(q_0, S) \rightarrow (q_0, ABC)$

$(q_0, A) \rightarrow (q_0, u_1 A \bar{u}_1)$

$(q_0, A) \rightarrow (q_1, vA)$

$(q_1, A) \rightarrow (q_1, vA)$

$(q_1, B) \rightarrow (q_2, u_3 B)$

$(q_2, C) \rightarrow (q_1, C\bar{u}_3)$

$(q_1, A) \rightarrow (q_3, u_2 A)$

$(q_3, B) \rightarrow (q_3, \lambda)$

$(q_3, C) \rightarrow (q_4, C\bar{u}_2)$

$(q_4, A) \rightarrow (q_3, u_2 A)$

$(q_5, A) \rightarrow (q_5, \lambda)$

$(q_5, C) \rightarrow (q_5, \lambda)$

Derivation for input string $w = cugcugacacagucug$

$S_{q0} \rightarrow ABC_{q0} \rightarrow cAgBC_{q0} \rightarrow cuAagBC_{q0} \rightarrow cugAcagBC_{q0} \rightarrow cugcAcagBC_{q1} \rightarrow cugcuAcagBC_{q1} \rightarrow cugcuAcagcBC_{q2} \rightarrow cugcuAcagcBCg_{q1} \rightarrow cugcuAcagcaBCg_{q2} \rightarrow cugcuAcagcaBCug_{q1} \rightarrow cugcugAcagcaBCug_{q3} \rightarrow cugcugAcagcaCug_{q3} \rightarrow cugcugAcagcaCcug_{q4} \rightarrow cugcugaAcagcaCcug_{q3} \rightarrow cugcugaAcagcaCucug_{q4} \rightarrow cugcugacagcaCucug_{q5} \rightarrow cugcugacagcaucug_{q5}$

**Proposition 5.** The recursive pseudoknot language $L_{RC} = \{u_1 u_2 u_3 \bar{u}_2^R u_4 \bar{u}_1^R \bar{u}_4^R u_5 \bar{u}_5^R \bar{u}_3^R \mid u_1, u_2, u_3, u_4, u_5 \in \Sigma_{RNA}^*\}$ can be generated by state grammar. Fig 5. represents the structure of recursive pseudoknot sequence.

Proof: Grammar $G_5 = (\{S, A, B, C, u_1, u_2, u_3, u_4, u_5, \bar{u}_1, \bar{u}_2, \bar{u}_3, \bar{u}_4, \bar{u}_5\}, \{q_0, q_1, q_2, q_3, q_4, q_5, q_6, q_7, q_8\}, \{u_1, u_2, u_3, u_4, u_5, \bar{u}_1, \bar{u}_2, \bar{u}_3, \bar{u}_4, \bar{u}_5\}, P, S)$ be the state grammar where production rules P are as follows:

$(q_0, S) \rightarrow (q_0, ABC)$

$(q_0, B) \rightarrow (q_0, u_4 B \bar{u}_4)$

$(q_0, A) \rightarrow (q_1, u_1 A)$

$(q_1, B) \rightarrow (q_2, B \bar{u}_1)$



**Fig. 5:** Recursive Pseudoknot Structure.

$(q_2, A) \rightarrow (q_1, u_1 A)$

$(q_2, A) \rightarrow (q_3, u_2 A \bar{u}_2)$

$(q_3, A) \rightarrow (q_3, u_2 A \bar{u}_2)$

$(q_3, B) \rightarrow (q_3, \lambda)$

$(q_3, C) \rightarrow (q_3, u_5 C \bar{u}_5)$

$(q_3, A) \rightarrow (q_4, u_3 A)$

$(q_4, C) \rightarrow (q_5, C u_5)$

$(q_5, A) \rightarrow (q_4, u_3 A)$

$(q_5, A) \rightarrow (q_6, \lambda)$

$(q_6, C) \rightarrow (q_6, \lambda)$

Derivation for input string $w = cagcucugagucuaag$

$S_{q0} \rightarrow ABC_{q0} \rightarrow AgBaC_{q0} \rightarrow AgaBuaC_{q0} \rightarrow AgaBuaC_{q0} \rightarrow cAgaBuaC_{q1} \rightarrow cAgaBguaC_{q2} \rightarrow caAugaBguaC_{q3} \rightarrow cagAcug aBguaC_{q3} \rightarrow cagAcugaguaC_{q3} \rightarrow cagAcugaguauCa_{q3} \rightarrow cagcA cugaguauCa_{q4} \rightarrow cagcAcugaguauCga_{q5} \rightarrow cagcuAcugaguauC ga_{q4} \rightarrow cagcuAcugaguauCaga_{q5} \rightarrow cagcucugaguauCaga_{q6} \rightarrow cagcucugaguauaga_{q6}$

**Proposition 6.** The kissing hairpin language $L_{KH} = \{u_1 v_1 A \bar{A} v_2 u_2 \bar{u}_2^R v_3 B \bar{B} v_4 \bar{u}_1^R \mid u_1, u_2, v_1, \quad v_2, v_3, v_4 \in \Sigma_{RNA}^*, A, B \in \Sigma_{RNA}\}$ can be generated by state grammar. Fig 6. represents the structure of kissing hairpin sequence.



**Fig. 6:** Kissing Hairpin Structure.

Proof:    Grammar    $G_6 = \{\{S, A, B, C, D, u_1, u_2, v_1, v_2, v_3, v_4\}, \{q_0, q_1, q_2, q_3, q_4, q_5\}, \{u_1, u_2, v_1, v_2, v_3, v_4\}, P, S\}$ be the state grammar where production rules P are as follows:

$$(q_0, S) \rightarrow (q_0, u_1 S \bar{u}_1)$$

$$(q_0, S) \rightarrow (q_1, v_1 C)$$

$$(q_1, C) \rightarrow (q_1, v_1 C)$$

$$(q_1, C) \rightarrow (q_1, C v_4)$$

$$(q_1, C) \rightarrow (q_2, CD)$$

$$(q_2, C) \rightarrow (q_1, AC\bar{A})$$

$$(q_2, D) \rightarrow (q_2, BD\bar{B})$$

$$(q_2, C) \rightarrow (q_3, v_2 C)$$

$$(q_3, C) \rightarrow (q_3, v_2 C)$$

$$(q_3, C) \rightarrow (q_3, C v_3)$$

$$(q_3, C) \rightarrow (q_4, u_2 C \bar{u}_2)$$

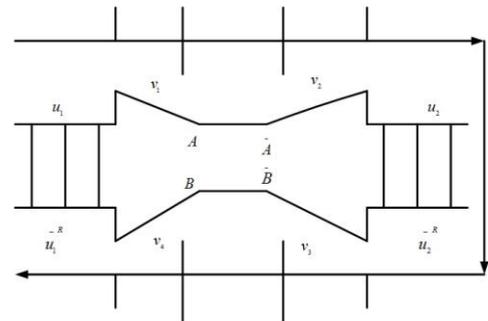$$(q_4, C) \rightarrow (q_4, u_2 C \bar{u}_2)$$

$$(q_4, D) \rightarrow (q_5, \lambda)$$

$$(q_4, C) \rightarrow (q_5, \lambda)$$

Derivation for input string $w = gaugaucuguaccagcgcuc$

$$S_{q0} \rightarrow gSc_{q0} \rightarrow gaSuc_{q0} \rightarrow gauCuc_{q1} \rightarrow gauuCuc_{q1} \rightarrow gauuCauc_{q1} \rightarrow$$
$$gauuCDauc_{q2} \rightarrow gauuuCgDauc_{q2} \rightarrow gauuCgcDgauc_{q2} \rightarrow gauuuC$$
$$gcuDagauc_{q2} \rightarrow gauuuCgcuDagauc_{q3} \rightarrow gauuuaCgcuDagauc_{q3} \rightarrow$$
$$gauuuaCagcuDagauc_{q3} \rightarrow gauuuaCaagcuDagauc_{q3} \rightarrow gauuua$$
$$CaagcuDagauc_{q3} \rightarrow gauuuaaCuaagcuDagauc_{q4} \rightarrow gauuuag$$
$$CcaagcuDagauc_{q4} \rightarrow gauuuagaCccaagcuDagauc_{q4} \rightarrow gauuu$$
$$agaccaagcuDagauc_{q5} \rightarrow gauuuagaccaagcuagauc_{q5}$$

## 4. Results and discussion

In this section, we compared the above-designed state grammar with the matrix insertion-deletion system based on descriptional complexity. Table 1, depicts the comparison between the matrix insertion-deletion system and state grammar in terms of production and variable used for attenuator, extended pseudoknot structure, kissing hairpin, simple H-type structure, recursive pseudoknot and three-knot structure. Clearly, state grammar is more succinct in terms of the number of productions and variables except in attenuator (one more variable used than the matrix insertion-deletion system) based on the type-1 descriptional complexity. Design of state grammar will be helpful in the prediction of secondary structure. Design of the parser for the above-designed state grammar is kept as a future work.

**Table 1:** Summary of Results

| System | Matrix Insertion-deletion system | | State grammar | |
|---|---|---|---|---|
| Language | Prod | Var | Prod | Var |
| Attenuator Language | 18+5=23 | 2 | 5+3=8 | 3 |
| Extended Pseudoknot Language | 10+7=17 | 4 | 8+4=12 | 3 |
| Simple H-type | 9+7=16 | 3 | 10+5=15 | 3 |
| Three knot structure | 11+8=19 | 4 | 12+6=18 | 4 |
| Recursive Pseudoknot | 15+10=25 | 5 | 14+7=21 | 4 |
| Kissing hairpin | 12+10=22 | 4 | 14+6=20 | 3 |

## Acknowledgements

## References

[1] L. Cai, L.R. Malmberg and Y. Wu, "Stochastic modeling of RNA pseudoknotted structures: a grammatical approach", Bioinformatics, Vol. 19, No. 1, i66-i73, (2003). https://doi.org/10.1093/bioinformatics/btg1007.

[2] T. Kasai, "An hierarchy between context-free and context-sensitive languages, Journal of Computer and System Sciences", Vol. 4, No. 5, (1970), pp: 492-508. https://doi.org/10.1016/S0022-0000(70)80045-9.

[3] K.Y. Sung, "The Use of Context-Sensitive Grammar For Modeling RNA Pseudoknots", In- Proceedings of the International Conference on Bioinformatics & Computational Biology (BIOCOMP'06); Las Vegas, Nevada, USA, 338-344, 26-29 June (2006).

[4] Y. Sakakibara, M. Brown, R. Hughey, I.S. Mian, K. Sjölander, R.C. Underwood, D. Haussler, "Stochastic context-free grammars for tRNA modelling", Nucleic Acids Research, Vol. 22, No. 23, (1996), pp: 5112-5120. https://doi.org/10.1093/nar/22.23.5112.

[5] Y. Sakakibara, "Pair hidden Markov models on tree structures", *Bioinformatics*, Vol. 19, No. 1, (2003), pp: 232-240. https://doi.org/10.1093/bioinformatics/btg1032.

[6] S. Yuki and T. Kasami, "RNA pseudoknotted structure prediction using stochastic multiple context-free grammar", *IPSJ Transactions on Bioinformatics*, Vol. 47, (2006), pp: 12-21.

[7] M. Brown and C. Wilson, "RNA pseudoknot modeling using intersections of stochastic context free grammars with applications to database search", Proceedings of the Pacific Symposium on Bio-computing, Big Island, HI, USA, (1995), pp: 109-125.

[8] E. Rivas and S. R. Eddy, "The language of RNA: a formal grammar that includes pseudoknots", Bioinformatics, Vol. 16, No. 4, (2000), pp: 334-340. https://doi.org/10.1093/bioinformatics/16.4.334.

[9] D. B. Searls, "The computational linguistics of biological sequences, Artificial intelligence and molecular biology", American Association for Artificial Intelligence Menlo Park, CA, USA, Vol. 2, (1993), pp: 47-120.

[10] D. B. Searls, "String variable grammar: A logic grammar formalism for the biological language of DNA", The Journal of Logic Programming, Vol. 24, No. 1, (1995), pp: 73-102. https://doi.org/10.1016/0743-1066(95)00034-H.

[11] N. Mizoguchi, Y. Kato, and H. Seki, "A grammar-based approach to RNA pseudoknotted structure prediction for aligned sequences", In - Proceedings of 1st IEEE International Conference on Computational Advances in Bio and Medical Sciences (ICCABS), Orlando, FL, (2011), pp: 135-140. https://doi.org/10.1109/ICCABS.2011.5729868.

[12] L. Kuppusamy and A. Mahendran, "Modelling DNA and RNA secondary structures using matrix insertion-deletion systems", International Journal of Applied Mathematics and Computer Science, Vol. 26, No. 1, (2016), pp: 245-258. https://doi.org/10.1515/amcs-2016-0017.

[13] N. Kalra and A. Kumar, "State Grammar and Deep Pushdown Automata for Biological Sequences of Nucleic Acids", Current Bioinformatics, Vol. 11, No. 4, (2016), pp: 470-479. https://doi.org/10.2174/1574893611666151231185112.

[14] A. Meduna and P. Zemek, "Regulated Grammars and Automata", Springer, New York, 2014. https://doi.org/10.1007/978-1-4939-0369-6.

[15] J. Dassow and G. Pawn, "Regulated rewriting in formal language theory", 1st ed Springer-Verlag Berlin Heidelberg, (2012).

[16] J. Dassow, "Grammars with regulated rewriting", In: Formal Languages and Applications: Studies in Fuzziness and Soft Computing Ed. Springer Berlin Heidelberg, 148, (2004), pp: 249-273. https://doi.org/10.1007/978-3-540-39886-8_13.

[17] L. Kuppusamy, I. Raman and K. Krithivasan, "On Succinct Description of Certain Context-Free Languages by Ins-Del and Matrix Ins-Del Systems", International Journal of Foundations of Computer Science, Vol. 27, No.7, (2016), pp: 775-786. https://doi.org/10.1142/S0129054116500295.

[18] N. Kalra and A. Kumar, "Fuzzy state grammar and fuzzy deep pushdown automaton", Journal of Intelligent and Fuzzy Systems, Vol. 31, No.1, (2016), pp: 249-258. https://doi.org/10.3233/IFS-162138.

[19] K. Singh, "Conversion of deterministic finite automata to regular expression using bridge state", M.E. Thesis, Thapar University, (2011).

[20] N. Kalra and A. Kumar, "Deterministic Deep Pushdown Transducer and its Parallel Version", The Computer Journal, doi.org/10.1093/comjnl/bxx036. https://doi.org/10.1093/comjnl/bxx036.

[21] T. Singla and A. Kumar, "Reducing Mutation Testing Endeavor using the similar conditions for the same Mutation Operators occurs at Different Locations", Applied mathematics & Information Science, Vol. 8, No. 5, (2014), pp: 2389-2393. https://doi.org/10.12785/amis/080534.

[22] R. Kumar and A. Kumar, "Metamorphosis of fuzzy regular expressions to fuzzy automata using the follow automata", arXiv preprint arXiv: 1411.2865.