

Genome Based Computational Technique to Identify the Functional RNA in Protozoan Parasites

Jai Prakash Pandey,
Research Scholar, Department of
Management,
Lovely Professional University,
Jalandhar-Delhi, G.T. Road, Phagwara,
Punjab - 44411, India,
getjpg@gmail.com

Savita, SOA,
Graphic Era Hill University,
Dehradun,
savita.pnt@gmail.com

S.D.Lalitha
Assistant Professor, Department of
Computer Science and Engineering,
R.M.K Engineering College,
Kavaraipettai, Tamil Nadu, India.
sdl.cse@rmkac.ac.in

Amanulla Khan
Assistant Professor in Botany,
Anjuman Islam Janjira Degree College of
Science
(Affiliated to Mumbai University, Mumbai),
Lokmanya Tilak Road, Murud-Janjira dist.
Raigarh, Maharashtra 02401, India,
dramanullak@gmail.com

S.Prema
Assistant Professor, Department of
Mathematics,
Rajalakshmi Institute of Technology,
Chennai 600124, Tamil Nadu, India.
premhaha@gmail.com

Juhie Agarwal
Assistant Professor,
Department of Zoology,
Vardhaman College, Bijnor (M.J.P
Rohilkhand University, Bareilly)
India.
juhie07@gmail.com

Abstract—Operative RNAs, including non-coding RNAs (ncRNAs) and cis-acting RNA elements involved in posttranscriptional genome regulation, are found based on two distinct computerized examinations of Trypanosoma genomes. Based on compliance with comparable trypanosomatid *Leishmania braziliensis*, the expected possible ncRNAs are discovered in the first analysis. Several of the expected ncRNAs are novel categories having undetermined characteristics, such that predictions have a low estimated incorrect fraud threat. This research uncovered several mechanism regulatory motifs in the subsequent research we used to develop a classification that can distinguish between actions that aren't identical. The first genomic sequence analysis of trypanosomatid fRNAs helps focus the research on practical approaches and accelerates the discovery of those elements. These functionality predictions classifiers built upon cis-acting regulation regions may potentially be used to offer homology-independent annotating for the trypanosomatid genome when combined with existing approaches.

Keywords—RNA elements, gene regulation, trypanosomatid genomes, ncRNA, trypanosomatid fRNA

I. INTRODUCTION

Functionality RNAs (fRNAs), and RNA components functioning on the RNA levels, were being extra acknowledged as their extensive architectural, regulation, and catalysis activities are disclosed [1]. Another type of fRNAs is the cis-regulatory factors found inside the 5' and 3' non-translated sections (UTRs) of mRNAs, primarily engaged in post-transcriptional control of gene activity [2-4]. Current advances in computational approaches for fRNA predictions have revealed a large number of RNA components that are engaged in post-transcriptional regulation mechanisms. While important in numerous organisms, post-transcriptional control is particularly important in the class with multicellular worms known as trypanosomatids. It is the primary method for regulating gene translation. Trypanosomatids, including Trypanosoma brucei, Trypanosoma cruzi, and many Leishmania species, cause significant human & animal infections with a higher frequency and fatality rate if left uncontrolled [5]. For trypanosomatids, genomes are translated as polycistronic mRNAs, then trans-spliced [6]. Numerous cis-acting fRNA factors, including in U-rich

components (UREs), shorter intervening degeneration retroposons (SIDERs), and others, are involved in genes translation control, which happens mainly around or following spliced [7]. Those features primarily control mRNA integrity or translational rates by interacting with various trans-acting proteins, several of which are unidentified [8]. While little empirical evidence has been discovered, this has lately been suggested that miRNAs may function in post-transcriptional genome expression.

II. RELATED WORKS

Comprehensive functionality experiments with numerous organisms bearing expulsion mutations of the putative regulatory region are required to identify cis-acting fRNA components experimentally [9]. The scenario of ncRNAs isn't much clearer since it's unclear where in the genomes there must be looked for & how that screened research must be conducted. Because of the absence of significant conservation markers within their sequencing, automatic recognition of fRNAs using genomic patterns was not as reliable as recognition of protein-coding RNAs. Which were exceptionally trustworthy based on analytical and technological evaluations, it provided a computerized analysis of the genomes for *T. brucei* and *L. braziliensis* in the hopes of discovering a group of preserved ncRNAs. We showed that this technique could identify a significant variety of possible ncRNAs, both identified and unknown [10-11]. We look at possible Premi RNAs within prospective ncRNAs and find that the occurrence of miRNA sequences maintained across *T. brucei* and *L. braziliensis* was very improbable. We also employ a new strategy for identifying shorter regulation RNA motifs in *T. brucei* genomes' 5' & 3' UTRs, which are homology-independent [12-13].

Those motifs list the more operationally relevant sections of possible cis-regulatory fRNA components to go along with our projected ncRNAs. Such regulation patterns could be employed for predicted genetic activity and provide fresh suggestions into the regulatory processes for protein production in *T. brucei*.

III. PROPOSED METHODS

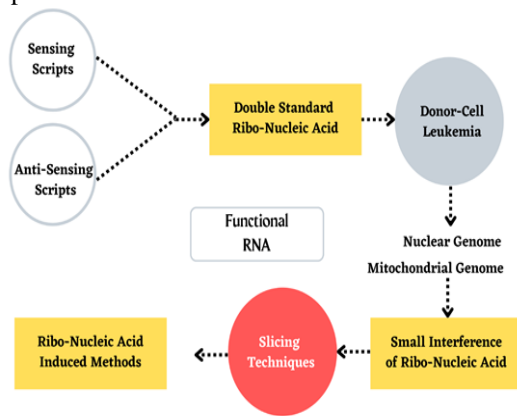


Fig.1. Proposed Model for Identifying Functional RNA

Fig.1 represents the sample model for identifying the functional RNA using computational techniques. The binomial-based method was used to find preserved genome areas. Shortly, after disguising LCRs with mreps, the 2 sequences of *T. brucei* and *L. braziliensis* were matched for so every screen of 25 nucleobases. N , the quantity of preserved nucleic acids, has been ascertained. This same likelihood of noticing N kept nucleotides out of 25 nucleotides has been estimated under the assumptions of impartial replacement and predicated on binomial dispersion of mutants. Territories demonstrating proof of Lastly, similar areas separated by less than 25 nucleotides are linked together to form a unique area.

QRNA was employed to find sections of the same genome areas that retained structure RNA component sequences. Long sections were broken down into 80-nucleotide overlapped chunks, one of them overlapped by 40 nucleotides alongside the neighbor segments. Intersecting regions having RNA values greater than 0 were consolidated, QRNA values were reassessed, and probable f RNAs containing final positive RNA values were chosen. The proportion of preserved encoding regions identified as fRNA was used to compute the ratio of the erroneous positive. LCRs were calculated using TIGR's 'mreps' and 'must,' which are denoted in lowercase characters as Supplementary Files 1. This same importance of every applicant has been determined by correlating its QRNA rating to a dispersion acquired by irregularly flipped *T. brucei*-*L. braziliensis* pairings: so, every ncRNA person's connectivity was regarded individually, and sections to similar preservation trends have been moved arbitrarily. Sections with gaps were exchanged with just too many gap-containing sections, mismatched with the discrepancy, and compared with the game.

FIRE, a newly created method, can detect DNA and RNA motifs that are redistributed throughout distinct sequencing groups, i.e., were highly represented within specific groupings but deficient in another. We employed FIRE to find themes that are spread disproportionately across various activities. First, *T. brucei* functionality classifications were obtained from the KEGG pathway dataset. It divided these genomes into 2 categories, with each route depending on whether or not people are engaged with this pathway. Then we employed FIRE to look for 5'

and 3' Unique registration motifs which were significantly overrepresented and underrepresented within one of the 2 groups. Based on newly reported splicing position estimations, the genome of matured 5' and 3' UTRs were recovered in *T. brucei* genomes. The motifs generated by the various algorithms were gathered, and duplicated motifs were deleted.

Because KEGG employs an automatic workflow to allocate genomes into templates routes depending on overall similarity with existing proteins, KEGG annotations might never correlate to the specific activity of elements. Nevertheless, we assume that the genome connections will be preserved due to this method. For example, if two genomes were allocated to the identical pathway in KEGG, they are highly likely to have similar activities, especially if the KEGG allocated features were incorrect. We used a naive Bayesian network to predict gene-function connections based on expected regulatory mechanisms in the 5' and 3' UTRs. With a naive Bayes network, the qualities required to classify items are presumed to be unconnected. Considering any collection with identified motifs, the chance for genetic g corresponding to clusters is determined as:

$$A(h \cup \partial | G^N) = \bigcap_{g_x \in G^N} A(h \cup \partial | g_x)$$

IV. RESULTS AND DISCUSSION

That anticipated false-positive percentages employing coded areas will never apply to noncoding areas. This should be indicated if it has used values other than the RNA grade from QRNA, such as the COD & OTH ratings. However, because translating sequences develop in a particular fashion, QRNA activity might change across coded and non-RNA, noncoding genetic areas. Moreover, RNA architecture in the encoding gene might be targeted for selection. In conclusion, as explained in the section "Identification of Highly Important Candidates," it used a different, exceedingly careful approach to quantify the false positives frequencies of the ncRNA estimations. Addition Item 1 contains the comprehensive listing of every discovered ncRNA variant and its accompanying data. As illustrated in Fig.2, most of these contenders may be classified under multiple homologous groups. This prediction could be deemed highly credible when numerous homologous sequences are separately estimated to be ncRNAs. Cluster 10 is particularly significant because of its size, implying that the members of this class are found throughout genomes at a great rate. It was just postulated that trypanosomatids might employ miRNAs to control the amounts of specific mRNAs depending on a computer examination of the *T. brucei* genome. This research, however, contradicts our present understanding of miRNA genesis management. Mirna appears to have evolved in a distinct branching of living, albeit convergence development in numerous lineages was not feasible. As this result, we chose to look at probable miRNA precursor amongst our projected ncRNAs using a basic but particular technique that considers a few architectural and thermodynamics parameters for identifying pre-miRNA structures. Employing 250 pre-miRNAs chosen at randomized from 24

various species for reference sequencing, the accuracy of pre-miRNA predictions utilizing our parameters is predicted to be about 32.4 percent \pm 2.1 percent. According to research findings, there are very few retained miRNA genetics in *T. brucei*. However, it's worth noting that several of the previously anticipated *T. brucei* miRNAs might be addressing mutant surfaces polypeptides that aren't seen in *L. braziliensis*. In addition, some miRNAs expected are likely to bind preserved components like the 20S proteasome likely & should thus be anticipated to discover in our analysis if they are preserved. While this would not rule out the existence of miRNAs in its *T. brucei* genomes, it may imply that its genes should be reexamined for these components.

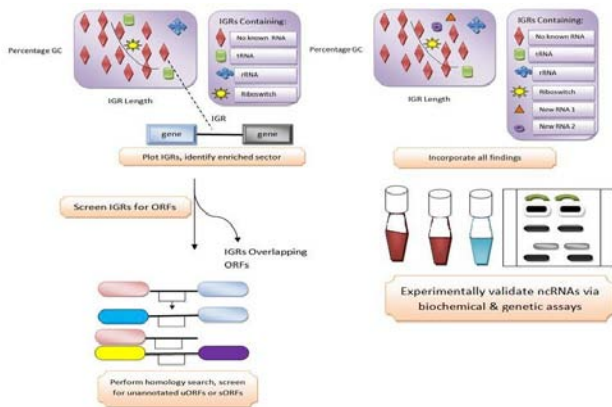


Fig.2. Proposed Architecture of ncRNAs

This was dismissed when a potential ncRNA is exceeded with greater than three randomization variations. This screening process yielded 117 very important unique potential ncRNAs, 53 of that did neither intersect nor were comparable to any documented characteristics of the *T. brucei* genomes, suggesting that it may be wholly unique ncRNAs (see Table 1). The RNA model, not the COD or OTH modeling, had the most significant rating across the 117 contenders who could never intersect within a decoding sequence, even though it was initially chosen solely based on the RNA score and neither COD nor OTH values. The determined p-value gives an additional, more cautious metric for measuring our program's accuracy. If the majority of the non-coding genomes are made up of non-RNA randomized elements, a p-value of 0.001 corresponds to around 0.887 kbp of false positives. This is more than 5.7 kbp of our contenders are meaningful at this level, showing an accuracy of around 85 percent at this level.

TABLE 1. PREDICTED CLASSIFICATION OF ncRNAs

	With in 100 bp of a non overlapping coding	ncRNA cluster	With in a stand switch region	elsewhere	Total
Overlap psaidogene	1	1	1	30	33
Overlap unlikelyproteins	1	2	2	1	6
HomoloeoustorRNA	1	2	2	1	6
Homologous to rRNA	5	5	3	1	14
Total	8	10	8	33	59

Nevertheless, depending upon its placements, a rough assumption may be made about the ecological role of

distinct possibilities. For instance, there are eight unidentified potential components in the surroundings of code sequences. Such components might regulate structures found with between 5' or 3' UTRs of coding segments that regulate post-transcription genes. In addition, the string switching area is also discovered that includes one unidentified potential fRNA. Because polycistronic mRNA production begins at strand flip sites, this fRNA might be a component of the resulting transcript's 5' end. It could be implicated in its localization, post-transcriptional treatment, or control.

We evaluated the occurrence of mechanism motifs in the 5' and 3' UTRs of *T. brucei* genomes using a homology-independent technique using FIRE, a newly constructed approach. FIRE was found to detect several recognized, and new regulating components in upstream and downstream domains grouped similar on its activity patterns having such a near-zero false positives discovery rate. We employed FIRE to look for 'function-specific regulation components in *T. brucei* genes' 5' and 3' UTRs: genes with similar activities are frequently co-regulated, implying that they must have identical cis-regulatory components. As a result, grouping genes based on their actions could be employed as a stand-in for grouping transcripts based on the activity trends. This method is especially beneficial in species whose gene regulatory happens mainly at a post-transcriptional stage, like trypanosomatids, where transcripts profile techniques cannot discern the fluctuations of protein production.

It discovered Sixteen mechanism motifs in *T. brucei* 5' UTRs and 21 mechanism motifs in *T. brucei* 3' UTRs. Based on the results of running FIRE across ten subgroup pairs with gene-function combinations, we should expect a reliability of 75.3 % for finding mechanism 5' UTR themes & 84.8 % for 3' UTR themes. Most motifs discovered through FIRE had an orientations bias, meaning it primarily appears in a specific direction relative to the coding sequences. RNA patterns are known to have this characteristic. Furthermore, two of the common themes of 3' UTRs have placement biases, which means it likes to be located within a certain distance from the forward decode sequence's final codon. This feature has been seen in various regulating patterns in other species, increasing the likelihood of whether the anticipated pattern can have a biological function.

RNA motifs

It developed a naive Bayesian network that can determine if a genome belongs to a definite pathway based on the presence or absence of patterns within 5' and 3' UTRs. For example, this naive Bayesian network can be used to detect *T. brucei* proteins for various pathways consistently; see Fig.3, just as fewer themes were required to provide the most significant feasible predicting ability, as illustrated in Fig.3. Introducing more motifs to this classifier, on the other hand, does not influence its predictive strength, making the building of successful nave Bayesian networks easier. It might be possible to significantly increase the functional annotations of *T. brucei* genomes by integrating our technique with additional functional prediction techniques. Extra Item six contains a

detailed analysis of functional predictions in *T. brucei* employing our technique.

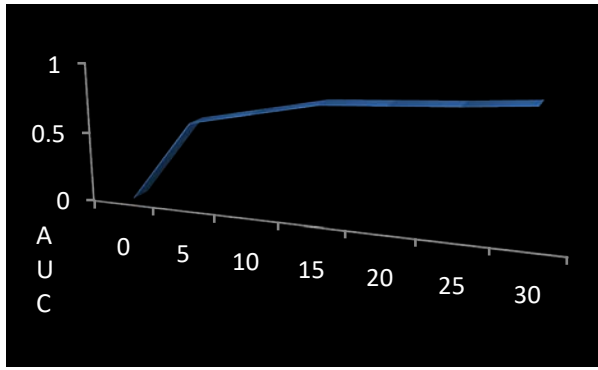


Fig.3. Proposed system prediction using regulatory motifs

V. CONCLUSION

The ncRNAs identified in this work might be employed in investigations to understand trypanosomatids' operational RNA repertoire better. In addition, the options with no existing homologs are expected to be unique ncRNA classes in *T. brucei*, making them even more fascinating. Discovering the function of these ncRNAs would help us better comprehend the infections' biology. Nevertheless, because we only looked at two chromosomes for this work, our list of anticipated ncRNAs was much but comprehensive. Analyzing a wider number of trypanosomatid genotypes could uncover additional ncRNAs and offer a more comprehensive picture of those species' non-coding functioning transcriptomic.

REFERENCES

- [1] M.H. Seabolt, K.T. Konstantinidis, and D.M. Roellig, "Hidden Diversity within Common Protozoan Parasites as Revealed by a Novel Genotyping Scheme," *Applied and Environmental Microbiology*, vol. 87, no. 6, p. e02275-20, Feb 26 2021.
- [2] D. Barh, S. Tiwari, M.E. Weener, V. Azevedo, A. Góes-Neto, M.M. Gromiha, and P.Ghosh, "Multi-omics-based identification of SARS-CoV-2 infection biology and candidate drugs against COVID-19," *Computers in Biology and Medicine*, vol. 126, p. 104051, Nov 1 2020.
- [3] Y. Li, R.P. Baptista, and J.C. Kissinger, "Noncoding RNAs in apicomplexan parasites: an update," *Trends in Parasitology*, Aug 20 2020.
- [4] S.R. Maran, K. Fleck, N.M. Monteiro-Teles, T. Isebe, P. Walrad, V. Jeffers, I. Cestari, E.J. Vasconcelos, and N. Moretti, "Protein acetylation in the critical biological processes in protozoan parasites," *Trends in Parasitology*, May 12 2021.
- [5] N.P. Mthethwa, I.D. Amoah, P. Reddy, F. Bux, and S. Kumari, "A review on application of next-generation sequencing methods for profiling of protozoan parasites in water: Current methodologies, challenges, and perspectives," *Journal of Microbiological Methods*, p. 106269, Jun 12 2021.
- [6] Moshika, A., Thirumaran, M., Natarajan, B., Andal, K., Sambasivam, G., &Manoharan, R. (2021).Vulnerability assessment in heterogeneous web environment using probabilistic arithmetic automata. *IEEE Access*, 9, 74659-74673.
- [7] Z. Yang, M. Wang, X. Zeng, A.T. Wan, and S.K. Tsui, "In silico analysis of proteins and microRNAs related to human African trypanosomiasis in tsetse fly," *Computational Biology and Chemistry*. Vol. 88, p. 107347, Oct 1 2020.
- [8] P. Simmonds, "Pervasive RNA secondary structure in the genomes of SARS-CoV-2 and other coronaviruses," *MBio.*, vol. 11, no. 6, p. e01661-20, Oct 30 2020.

- [9] J. Cao, and Y. Xue, "Characteristic chemical probing patterns of loop motifs improve prediction accuracy of RNA secondary structures," *Nucleic Acids Research*, vol. 49, no. 8, pp. 4294-307, May 7 2021.
- [10] Rajesh, M., &Sitharthan, R. (2022). Image fusion and enhancement based on energy of the pixel using Deep Convolutional Neural Network. *Multimedia Tools and Applications*, 81(1), 873-885.
- [11] ManjunathanAlagarsamy, KarthikramAnbalagan, YuvarajaThangavel, JeevithaSakkarai, JenopaulPauliah, and KannadhasanSuriyan, "Classification of covid patient image dataset using modified deep convolutional neural network system", *Bulletin of Electrical Engineering and Informatics*, vol. 11, no. 4, pp. 2273-2279, August 2022.
- [12] ManjunathanAlagarsamy, Joseph Michael JerardVedam, NithyadeviShanmugam, ParamasivamMuthanEswaran, GomathySankaraiyer, and KannadhasanSuriyan, "Performing the classification of pulsation cardiac beats automatically by using CNN with various dimensions of kernels", *International Journal of Reconfigurable and Embedded Systems*, vol. 11, no. 3, pp. 249-257, 2022.
- [13] M.Ramkumar, A. Lakshmi, M.P.Rajasekaran, and A.Manjunathan, "MultiscaleLaplacian graph kernel features combined with tree deep convolutional neural network for the detection of ECG arrhythmia", *Biomedical Signal Processing and Control*, vol. 76, p. 103639, 2022.