

ON SECONDARY STRUCTURE ANALYSIS BY USING FORMAL GRAMMARS AND ARTIFICIAL NEURAL NETWORKS*

Polina Lunina^{1,2}[0000–0002–7172–2647] and Semyon
Grigorev^{1,2}[0000–0002–7966–0698]

¹ Saint Petersburg State University, 7/9 Universitetskaya nab., St. Petersburg,
199034, Russia

² JetBrains Research, Primorskiy prospekt 68-70, Building 1, St. Petersburg 197374,
Russia

`lunina.polina@mail.ru`, `semyon.grigorev@jetbrains.com`,
`s.v.grigoriev@spbu.ru`

Abstract. A way to combine formal grammars and artificial neural networks for biological sequences processing was recently proposed. In this approach, an ordinary grammar encodes primitive features of the RNA secondary structure, parsing is utilized for features extraction and artificial neural network — for processing of the extracted features. Parsing is a bottleneck of the solution: input sequences should first be parsed before processing with a trained model which is a time-consuming operation when working with huge biological databases. In this work, we solve this problem by employing staged learning and limiting parsing to be used only during network training. We also compare networks which represent the parsing result in two different ways: by a vector and a bitmap image. Finally, we evaluate our solution on tRNA classification tasks.

Keywords: DNN · CNN · Machine Learning · Secondary Structure · Genomic Sequences · Formal Grammars · Parsing.

1 Introduction

Development of effective computational methods for genomic sequences analysis is an open problem in bioinformatics. While the existing algorithms for sequences classification and subsequences detection adopt different concepts and approaches, most of them share one idea: the secondary structure of genomic sequences contains important information about the biological functions of organisms. There are different ways to handle secondary structure, for example, probabilistic grammars [6, 10] and covariance models [7].

Real-world biological data commonly contains different mutations, noise, and random variations. This issue requires some sort of probability estimation while

* Supported by the Russian Science Foundation grant 18-11-00100

modeling the secondary structure. Probabilistic grammars and covariance models provide such functionality, are expressive and handle long-distance connections. They are successfully used in practical tools, such as Infernal [11], but building and training accurate grammar or model for predicting the whole secondary structure involves theoretical and practical difficulties. On the other hand, artificial neural networks are a common way to process noisy data and find complex structural patterns. Moreover, the efficiency of neural networks for genetic data processing has already been shown in some works [13, 9].

An approach for biological sequences processing which employs the combination of ordinary formal grammars and artificial neural networks was proposed in [8]. The key idea is to use an ordinary (not probabilistic) context-free grammar to describe only basic secondary structure features and leave the entire sequence analysis along with probabilistic estimation to the neural network which takes parsing-provided data as an input and solves some given task.

The secondary structure of RNA sequences can be viewed as a composition of stems [12]. Grammar G_0 that is used in [8] as well as in the present work is presented in figure 1. This grammar considers only the conventional base pairs (line 5) and describes the recursive composition of stems which are at least three base pairs in height (lines 7-12). Stems may be connected by an arbitrary sequence of length from 2 up to 10, and loops have the same length (line 2).

```

1  s1: stem<s0>
2  any_str : any_smb*[2..10]
3  s0: any_str | any_str stem<s0> s0
4  any_smb: A | U | C | G
5  stem1<s>: A s U | G s C | U s A | C s G
6  stem2<s>: stem1< stem1<s> >
7  stem<s>:
8      A stem<s> U
9      | U stem<s> A
10     | C stem<s> G
11     | G stem<s> C
12     | stem1< stem2<s> >

```

Fig. 1. Context-free grammar G_0 for RNA secondary structure features description

The result of a parsing algorithm for an input string w and a fixed grammar non-terminal N (start nonterminal) is an upper-triangular boolean matrix M_N , where $M_N[i, j] = 1$, iff the substring $w[i, j - 1]$ is derivable from N . This means that, for the grammar G_0 , a matrix contains 1 in a cell iff a correspondent substring folds to a stem of height at least 3. Such stem results in a diagonal chain of 1 in the matrix. Figure 2 presents the parsing result for sequence

$$w_1 = CCCCATTGCCAAGGACCCACCTTGGCAATCCC$$

w.r.t the grammar G_0 . Colored boxes map a substring which folds to a stem to correspondent cells in the matrix. Besides, this matrix contains other non-zero cells, because parser detects all possible foldings for all possible substrings. It can be either noise or some important information about the secondary structure. One of the tasks that neural network should perform is to process such matrices and filter all the insignificant contacts.

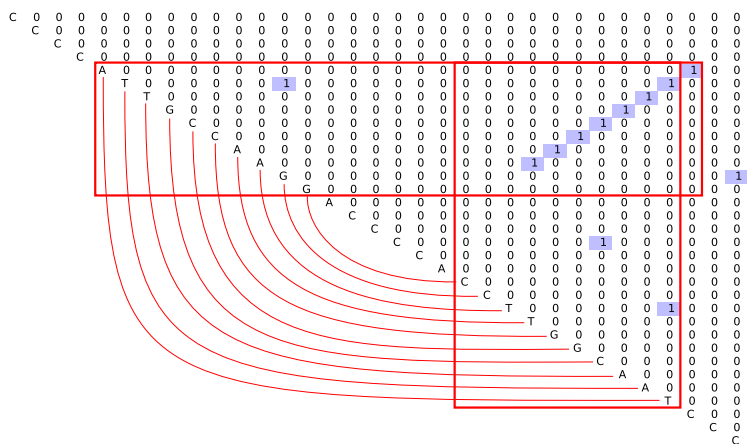


Fig. 2. Parsing result for sequence which should fold to stem

The parsing result in a form of a matrix can be linearized, compressed into a byte or int vector, and be further handled by a dense neural network, as described in [8]. Unfortunately, linearization breaks data locality: a chain of 1, which signifies a high stem, is local in a matrix, but is broken apart during its linearization. We see it to be an argument to investigate the applicability of convolutional networks for parsing result handling, as a boolean matrix can be converted to a black-white bitmap image. In this paper, we provide an empirical comparison of networks which handle vectors and images.

Another problem is a bad performance of the earlier solution. Since the trained network handles parsing result, each input sequence should first be parsed. Parsing is a very time-consuming step: context-free parsing has cubic complexity in terms of the input length. Even if we use matrix-based parsing algorithm [3] which utilizes GPGPU, performance is insufficient. We believe it would be better to avoid the parsing step at the final stage of the solution. In this work, we propose a way to solve this problem by building a network which handles raw sequences, not parsing results.

2 The solution

In this paper, we improve the solution proposed in [8]. We provide some ideas that are aimed to optimize its quality and performance and solve the problems that we faced during the experimental research.

First, we describe how to use a convolutional network for parsing result processing. Parsing result is a boolean matrix and we utilize the artificial neural network to detect sufficient features and find patterns in their appearance. Therefore, we need to transform these boolean matrices to some data structure acceptable by the neural network. Presently, we came up with two possible ways.

The first one is to drop out the nullary bottom left triangle, vectorize the top right triangle row by row and transform it into a byte vector. This approach reduces the input size, but it requires all the input sequences to have equal length. Thus we propose to either cut sequences to be of some predefined length or to pad them up with some blank symbols. Vectorization breaks data locality which makes learning harder: the network should restore back the relations broken during linearization. This also means that learning takes more time.

The second way is to represent the matrix as an image: the false bits of the matrix as white pixels and the true bits as black ones. This approach makes it possible to process sequences of different lengths since the images are easily transformed to a specified size. Data locality is also preserved: the information about relative positions of extracted basic features does not get lost which should improve learning.

The architecture of the neural network that takes vectorized data as an input is described in [8] and it consists of the long sequence of interchangeable dense and dropout layers with aggressive batch normalization. To handle images, we propose to use a network which consists of a small number of convolutional layers, linearization, and dense network which has a similar architecture as for vectorized data. In this paper, we provide an evaluation on both data formats and compare the results.

Another improvement that we came up with concerns parsing elimination in the context of our solution. The idea is to create a model which can handle original sequences instead of the parsing matrices. For that, we propose to use two-staged learning: first, a network which solves a subtask is trained and then it is used as pretrained layers in the training of the resulting network. In our solution, we first train a neural network to handle parsing results which performs classification according to a problem at hands. We create two networks in order to compare different architectures: one of them handles vectorized parsing result, the other handles parsing result represented as a bitmap image. After that, we extend these neural networks by several input layers that take the initial nucleotide sequence as an input and convert it to the parsing result which is handled appropriately by the pretrained layers.

Figure 3 represents the detailed description of these three neural networks architectures. Here N1 is a network which handles images, N2 is a network which handles vectorized parsing results, and N0 is an additional block which converts the input sequence into a set of features which can be handled by using N1 or

N2. So, firstly we train N1 and N2 on parsed data. After that, for vector-based network we combine the extension N0 and the whole original sequence of layers and for image-based network we use the similar architecture, except we remove the convolutional layer from the extended model, thus, the first layer at the junction of the blocks corresponds to the linearized image.

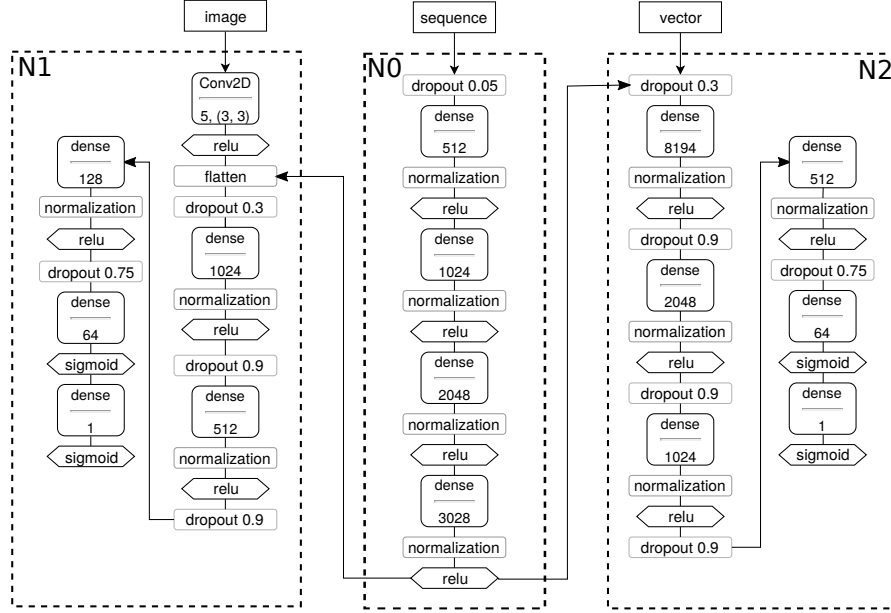


Fig. 3. Neural networks architectures

To sum up, we developed a technique to process parsing matrices as images by convolutional neural networks. Also, we built a model that handles sequences and requires parsing only for training the network it is based on. This removes the parsing step from the usage of the trained model.

3 Experiments

We evaluated the proposed approach with the described above modifications on two tRNA sequences analysis tasks. The first one was a classification of tRNA into two classes: eukaryotes and prokaryotes, while the second was a classification into four classes: archaea, bacteria, plants and fungi.

We took sequences from tRNA databases GtRNAdb³ [4] and tRNADB-CE⁴ [2] for these experiments. We used the parsing algorithm implemented by means of the YaccConstructor⁵ platform and Keras library [5] with Tensorflow framework [1] for neural networks training and testing. All models, as well as parsing tool, were run on GPU NVIDIA GeForce GTX 1070. We selected the equal number of samples (single tRNA molecule sequences) for each class for both classification tasks. Each sample was parsed w.r.t. the grammar G_0 and then both vectorized and transformed into an image. After that, we trained two neural networks: first handles the representation of the parsing result as vectors, and the second — as images. Finally, we trained the extended neural network. It consists of a block which takes an initial tRNA sequence as an input and transforms it into the parsing result and the block of pretrained layers: either the vector- or the image-based model from the previous step.

All extended neural networks were trained, validated (by hold-out validation) and tested on the same datasets as the corresponding base ones. The trained models for two classes (EP) and for four classes (ABFP) classification tasks were estimated by using classical machine learning metrics: accuracy, precision and recall.

Accuracy metrics for each problem for the test datasets are presented in the table 1, where base model is a model which handles parsing result (image or vector respectively) and extended model handles tRNA sequences and extends the corresponding base model. Also, this table shows the total time spent on two stages of training (base network + extended network) for both problems and types of data.

Table 1. Base and extended models test results by accuracy metrics

Classifier	EP		ABFP	
Approach	Vector-based	Image-based	Vector-based	Image-based
Base model accuracy	94.1%	96.2%	86.7%	93.3%
Extended model accuracy	97.5%	97.8%	96.2%	95.7%
Total training time	30000s	4600s	31800s	3600s
Samples for train:valid:test	20000:5000:10000 (57%:14%:29%)		8000:1000:3000 (67%:8%:25%)	

³ GtRNAdb tRNA database Web page: <http://gttradb.ucsc.edu/>. Access date: 07.03.2020.

⁴ The tRNADB-CE tRNA database Web page: <http://trna.ie.niigata-u.ac.jp/cgi-bin/trnadb/index.cgi>. Access date: 07.03.2020.

⁵ YaccConstructor is an SDK for syntax analysis tools development. Project repository on GitHub: <https://github.com/YaccConstructor/YaccConstructor>. Access date: 07.03.2020.

The estimations by precision and recall metrics for extended models for both classifiers on the same samples as in the table 1 are presented in the table 2.

Table 2. Extended models test results by precision and recall metrics for each class

Classifier	Class	Vector-based approach		Image-based approach	
		precision	recall	precision	recall
EP	prokaryotic	95.8%	99.4%	96.2%	99.4%
	eukaryotic	99.4%	95.6%	99.4%	99.5%
ABFP	archaeal	91.1%	99.2%	91.6%	98.5%
	bacterial	96.6%	95.1%	95.2%	95.5%
	fungi	98.5%	94.9%	97.5%	94.3%
	plant	99.4%	95.7%	99.2%	94.7%

The results show that our approach is applicable to tRNA classification tasks and both vector- and image-based models can be used along with dense and convolutional layers in neural networks architectures. While the differences in results for extended models are insignificant, for base models image-based network demonstrates slightly better results (see table 1). We believe that the reason of this effect lays in a better locality of features in the image-based representation of parsing result: chain of 1 which means a high stem is local in terms of picture but is broken during linearization. Also, we analyzed the time spent on all the models training (table 1) and, although some of these numbers could probably be decreased by more detailed networks tuning, we can state that image-based networks learn much faster than vector-based ones. The current model for images classification uses a single convolution layer. Whether it is possible to utilize deep convolutional networks for secondary structure analysis in the discussed approach is a question for future research.

The idea of the extended model that handles sequences instead of parsing results is proved to be applicable in practice and it demonstrates even higher quality than the original parsing-based model, as illustrated by table 1. We can conclude that it is possible to use parsing only for network training without decreasing the network quality.

To demonstrate the advantage of this technique in practical use in comparison with the classical way (when sequences should first be parsed) we took 100 tRNA sequences from two classes: eukaryotes and prokaryotes and used all four of the trained models to predict their classes. While using base models each sequence was parsed, transformed to correspondent format (image or vector) and fed to the neural network. Extended networks run on original sequences, so the parsing step was skipped. We measured the total time required to output predicted class for all the considered sequences in each case. In the table 3 the results of this experiment are provided and it is clear that the time spent for parsing is crucial relative to the total working time. So, the parsing eliminating modification significantly improves the performance of our solution.

Table 3. Time measurements for 100 sequences processing

Step	Vector based approach		Image based approach	
	Base	Extended	Base	Extended
Parsing	307.6s	—	310.5s	—
Weights loading	0.2s	0.2s	0.1s	0.3s
Class predicting	0.2s	0.2s	0.2s	0.3s
Total	308.0s	0.4s	310.8s	0.6s

4 Conclusion

We describe the modifications of the proposed approach [8] for biological sequences analysis using the combination of formal grammars and neural networks. We show that it is possible to improve the quality of the solution by representing parsing result as an image and handling it by using convolutional layers while processing it with a neural network. Also, we provide a technique that removes the parsing step from the trained model use and allows to run models on the original RNA sequences. As a result, the performance of the solution is significantly improved. We demonstrate the applicability of the proposed modifications for real-world problems⁶.

We can provide several directions for future research. First of all, it is necessary to investigate the applicability of the proposed approach for other sequences processing tasks such as 16s rRNA processing and chimeric sequences filtration.

Another possible application is a secondary structure prediction. We plan to investigate the possibility of creating a generative network which generates the most possible contact map for the given sequence.

The image-based model demonstrates a higher quality. We believe that it is caused by a better locality of features. If so, it should be possible to create a deep convolutional network for secondary structure analysis: further investigation is needed.

Finally, it is important to find a theoretical base for grammar tuning. It is important to adopt the theoretical results on secondary structure description by using formal grammar, such as [12] to find the optimal grammar for our approach.

References

1. Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G.S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A.,

⁶ Project description is available at the project page: https://research.jetbrains.org/groups/plt_lab/projects?project_id=43. Source code and documentation are published at GitHub: <https://github.com/LuninaPolina/SecondaryStructureAnalyzer>. Access date: 07.03.2020

- Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., Zheng, X.: TensorFlow: Large-scale machine learning on heterogeneous systems (2015), <http://tensorflow.org/>, software available from tensorflow.org
2. Abe, T., Inokuchi, H., Yamada, Y., Muto, A., Iwasaki, Y., Ikemura, T.: trnadb-ce: trna gene database well-timed in the era of big sequence data. *Frontiers in genetics* **5**, 114 (2014)
 3. Azimov, R., Grigorev, S.: Context-free path querying by matrix multiplication. In: *Proceedings of the 1st ACM SIGMOD Joint International Workshop on Graph Data Management Experiences & Systems (GRADES) and Network Data Analytics (NDA)*. GRADES-NDA '18, Association for Computing Machinery, New York, NY, USA (2018). <https://doi.org/10.1145/3210259.3210264>, <https://doi.org/10.1145/3210259.3210264>
 4. Chan, P.P., Lowe, T.M.: Gtrnadb 2.0: an expanded database of transfer rna genes identified in complete and draft genomes. *Nucleic acids research* **44**(D1), D184–D189 (2016)
 5. Chollet, F., et al.: Keras. <https://keras.io> (2015)
 6. Dowell, R.D., Eddy, S.R.: Evaluation of several lightweight stochastic context-free grammars for rna secondary structure prediction. *BMC bioinformatics* **5**, 71–71 (Jun 2004). <https://doi.org/10.1186/1471-2105-5-71>, <https://pubmed.ncbi.nlm.nih.gov/15180907>
 7. Durbin, R., Eddy, S.R., Krogh, A., Mitchison, G.: *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*. Cambridge University Press (1998). <https://doi.org/10.1017/CBO9780511790492>
 8. Grigorev, S., Lunina, P.: The composition of dense neural networks and formal grammars for secondary structure analysis. In: *Proceedings of the 12th International Joint Conference on Biomedical Engineering Systems and Technologies - Volume 3 BIOINFORMATICS: BIOINFORMATICS*, pp. 234–241. INSTICC, SciTePress (2019). <https://doi.org/10.5220/0007472302340241>
 9. Higashi, S., Hungria, M., De O. C. Brunetto, M.A.: Bacteria classification based on 16s ribosomal gene using artificial neural networks. In: *Proceedings of the 8th WSEAS International Conference on Computational Intelligence, Man-Machine Systems and Cybernetics*. p. 86–91. CMMACS'09, World Scientific and Engineering Academy and Society (WSEAS), Stevens Point, Wisconsin, USA (2009)
 10. Knudsen, B., Hein, J.: RNA secondary structure prediction using stochastic context-free grammars and evolutionary history. *Bioinformatics* **15**(6), 446–454 (06 1999). <https://doi.org/10.1093/bioinformatics/15.6.446>, <https://doi.org/10.1093/bioinformatics/15.6.446>
 11. Nawrocki, E.P., Eddy, S.R.: Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics* **29**(22), 2933–2935 (09 2013). <https://doi.org/10.1093/bioinformatics/btt509>, <https://doi.org/10.1093/bioinformatics/btt509>
 12. Quadrini, M., Merelli, E., Piergallini, R.: Loop grammars to identify rna structural patterns. In: *Proceedings of the 12th International Joint Conference on Biomedical Engineering Systems and Technologies - Volume 3 BIOINFORMATICS: BIOINFORMATICS*, pp. 302–309. INSTICC, SciTePress (2019). <https://doi.org/10.5220/0007576603020309>
 13. Sherman, D.J.: Humidor : Microbial community classification of the 16s gene by training cigar strings with convolutional neural networks (2017)