

САНКТ-ПЕТЕРБУРГСКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ
КАФЕДРА ГЕНЕТИКИ И БИОТЕХНОЛОГИИ

Васильев Артем Викторович
Выпускная квалификационная работа

“Эволюционные особенности структуры гена
Nxf1 (nuclear export factor) у животных“

Научный руководитель:
к.б.н., доцент, кафедра генетики и биотехнологии,
Голубкова Елена Валерьевна

Рецензент:
заведующая лабораторией, ведущий научный сотрудник,
лаборатория эволюционной геномики и палеогеномики, ЗИН,
к.б.н., с.н.с.,
Абрамсон Наталья Иосифовна

Санкт-Петербург
2025

Оглавление

1	Материалы и методы	3
2	Результаты	5
2.1	Анализ всех найденных видов	5
2.2	Подробный анализ Actinopterygii	6
3	Обсуждение	13
4	Список литературы	14

Материалы и методы

В качестве отправной точки был произведен поиск гена *Nxf1* внутри веб-сервиса NCBI [1]. Полученные данные были сохранены в текстовом формате и загружены в виде tsv-таблицы с помощью пакета pandas v2.2.3 [2] для языка программирования Python v3.12.6 [3]. Всего был найден 651 организм, содержащий анализируемый ген, большинство из которых относятся к Deuterostomia (Вторичноротые) - 436 видов. Таким образом, в качестве материалов выступали нуклеотидные и белковые последовательности гена *Nxf1* из открытых баз данных NCBI [1].

Большинство этапов последующего анализа реализовано в виде отдельных скриптов, разработанных в рамках данной работы, если не указано другое. Для логического разделения на блоки был использован Jupyter Notebook v1.1.1 [4].

По данным из полученной таблицы в разведывательных целях было построено филогенетическое дерево по найденным видам для оценки количества видов в таксонах более низкого ранга. Для глубокого анализа было принято решение сфокусироваться на организмах, относящихся к группе Protostomia (Первичноротые), Cnidaria (Стрекающие), а также на всех группах из Deuterostomia за исключением Mammalia (Млекопитающие).

Для найденных организмов с помощью пакета NCBI E-utilities из BioPython v1.85 [5] и NCBI Datasets Command-Line Interface (CLI) v18.0.2 [6] были загружены нуклеотидные последовательности гена, кодирующих участков и мРНК, а также аминокислотные последовательности белка в формате FASTA и аннотации для гена в GenBank-формате, необходимые для получения нуклеотидных последовательностей экзонов и поиска “консервативной кассеты”. Затем были получены и проанализированы интересующие нас участки экзон-интрон-экзонной структуры и созданы файлы со всеми экзонами и “кассетным” интроном для всех организмов, у которых получилось найти “кассету”. Данные файлы будут необходимы для последующего анализа.

Учитывая очень маленькие выборки во многих анализируемых группах (например, Cnidaria - 4 вида, Spiralia - 9 видов), было принято решение по увеличению их количества. Для этой цели, учитывая разнообразия полученных генов даже внутри одной таксономической группы, самым эффективным вариантом оказалось использование PSI-BLAST [7]. В качестве запроса (Query), или референса, использовались белковые последовательности тех организмов, у которых была найдена “кассета”. Для проведения PSI-BLAST были выбраны настройки по-умолчанию за исключением параметра Organism: поиск проводился внутри таксономической группы, к которой принадлежал референс, также референс был исключен из поиска.

Парсинг результатов BLAST также осуществлялся с помощью пакета BioPython [5] и специально разработанных скриптов. Он включал в себя фильтрацию данных по параметрам процента покрытия (Query Coverage, QC), длине и сходству (Per. Ident) найденных последовательностей (Subject), а также загрузку нуклеотидных и белковых последовательностей, однако реализация отличалась из-за особенностей баз

данных NCBI [1]. Получение “кассеты” было произведено по тому же принципу, но, опять же, с отличиями. Благодаря данному шагу удалось увеличить выборки суммарно на 117 видов. К сожалению, для некоторых таксономических групп увеличение выборки оказалось невозможным в связи с отсутствием у некоторых организмов интересующего нас участка.

Множественные выравнивания осуществлялись с помощью алгоритма MAFFT [8], 10 итераций, остальные настройки по-умолчанию, в программе Unipro UGENE v52.0 [9].

Анализ видов из *Deuterostomia* изначально шел более благоприятно за счет большого сходства последовательностей, в том числе интронных, и большего количества видов в группах. Для них также были загружены все необходимые файлы и произведен поиск и анализ “консервативной кассеты”. Мы решили сосредоточить свое внимание на организмах из *Actinopterygii* (Лучеперые рыбы), 72 вида, так как данных по ним ранее получено не было. Учитывая большую степень сходства интронных последовательностей, с помощью пакета инструментов MEME Suite v5.5.8 [10] локально был произведен поиск консервативных мотивов внутри “кассетного” интрона. Найденные мотивы, у которых E-value < 0.05 также локально были проанализированы с помощью Tomtom [11] из того же пакета. Для описанного шага была взята база данных JASPAR2024 CORE (NON-REDUNDANT) DNA.

С помощью инструмента RNAfold v2.7.0 из пакета ViennaRNA [12] были построены вторичные структуры РНК для нуклеотидных последовательностей в двух вариантах (MFE и Centroid), содержащих экзоны и “кассетный” интрон, т.к. мы предполагаем, что избегание интроном сплайсинга может быть опосредовано образованной им специфической вторичной структурой. Учитывая данное предположение, разумным шагом также являлся анализ “силы сайтов сплайсинга”, проведенный с помощью MaxEntScan [13]. Также с помощью скриптов цветом были выделены интронные последовательности внутри вторичной структуры и найденный мотив у *Actinopterygii*, который предположительно является CTE (Constitutive Transport Element).

Для *Actinopterygii* также был проведен филогенетический анализ, включающий построение и визуализацию деревьев. Для данной цели использовались самые популярные и проверенные временем инструменты. Построение деревьев осуществлялось с помощью IQ-TREE v2.4.0 [14], визуализация - с помощью Figtree v1.4.4 [15].

Работа проводилась в виртуальном окружении Mamba v1.5.5 [16], использованные пакеты и примеры анализа в Jupyter Notebooks можно найти в GitHub [17] репозитории автора: <https://github.com/ArtemVaska/Diploma>.

Для написания ВКР была использована система верстки LaTeX v4.76 [18], таблицы генерировались в веб-сервисе TablesGenerator [19]. Большинство картинок создано с помощью веб-сервиса draw.io [20]. Все шаги анализа проводились на базе операционной системы Linux Ubuntu 22.04 [21].

Результаты

Анализ всех найденных видов

Были проанализированы 413 нуклеотидных и белковых последовательностей гена *Nxf1* у представителей различных филогенетических групп из клад Cnidaria (Стрекающие) и Bilateria (Двусторонне-симметричные). Организмы, относящиеся к Mammalia, в анализ не были взяты в связи с уже имеющимися для них данными.

Для таксономических групп более низкого ранга с небольшим количеством видов в них был, с помощью PSI-BLAST были увеличены выборки, где это оказалось возможным, результат продемонстрирован на таблице 1.

Таблица 1: Результат увеличения выборки для таксономических групп Protostomia и Cnidaria с помощью PSI-BLAST.

Филогенетическая группа	Таксон высокого ранга	Видов до PSI-BLAST	Сколько видов добавил PSI-BLAST	Итого видов
Bilateria→Protostomia	Ecdysozoa	56	42	98
	Spiralia	6	63	69
Cnidaria	Anthozoa	2	12	14

В итоге для 353 видов удалось найти “консервативную кассету“ и продолжить дальнейший анализ.

На рисунке 1 отображено распределение исследованных видов по таксонам высокого ранга.

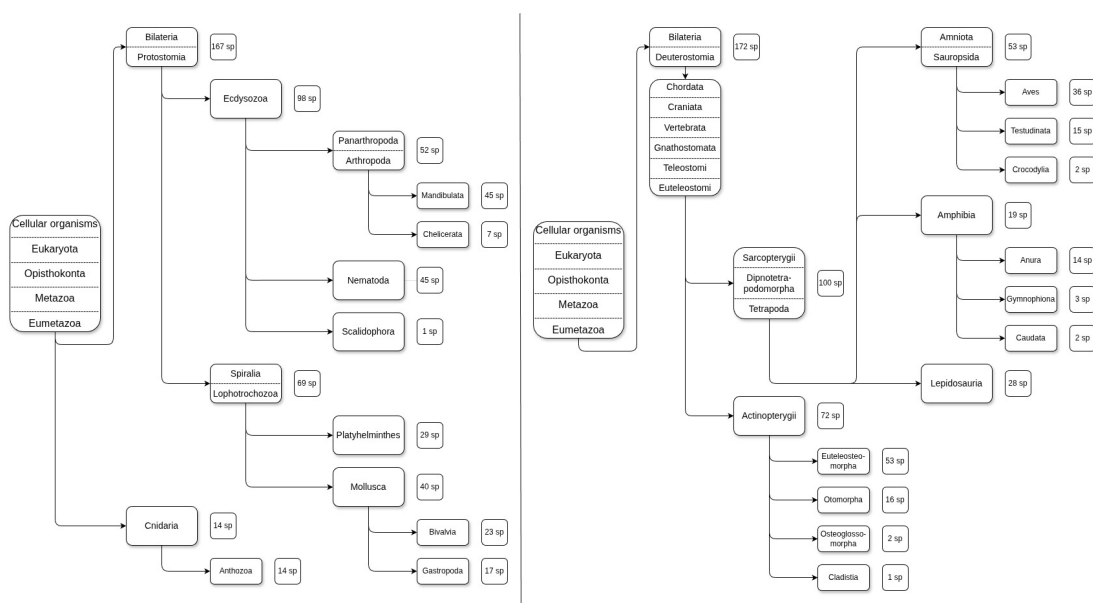


Рис. 1: Количество видов, взятых в анализ для Protostomia+Cnidaria и Deuterostomia.

Для всех видов, имеющих “консервативную кассету“, были построены вторичные структуры для интрон-содержащего транскрипта с выделением цветом “кассетного“ интрона (предоставляется по запросу).

Подробный анализ Actinopterygii

Для таксономической группы Actinopterygii проводился более углубленный анализ, так как на текущий момент данных по гену *Nxf1* для них не было. Для анализа были взяты все найденные представители данной филогенетической группы - 72 вида.

На таблице 2 показана характеристика “консервативной кассеты” для исследуемой группы.

Таблица 2: Сводная таблица с характеристикой кассетного интрона для таксономической группы Actinopterygii. Сортировка по возрастанию количества нуклеотидов до стоп-кодона в кассетной интроне.

Название организма	Кол-во нуклеотидов до стоп-кодона в интроне	Длина 1-го экзона в кассете	Длина кассетного интрона	Длина 2-го экзона в кассете
<i>Chanos chanos</i>	1	110	3568	37
<i>Danio rerio</i>	1	110	3580	37
<i>Denticeps clupeoides</i>	7	110	2629	37
<i>Labrus bergylta</i>	10	110	2684	37
<i>Cottoperca gobio</i>	16	110	2388	37
<i>Xiphophorus couchianus</i>	22	110	2227	37
<i>Larimichthys crocea</i>	22	110	2340	37
<i>Lates calcarifer</i>	22	110	2434	37
<i>Notothenia coriiceps</i>	22	110	2886	37
<i>Betta splendens</i>	22	110	2274	37
<i>Poecilia reticulata</i>	22	110	2262	37
<i>Takifugu rubripes</i>	22	110	2114	37
<i>Salarias fasciatus</i>	22	110	3855	37
<i>Poecilia mexicana</i>	22	110	2247	37
<i>Stegastes partitus</i>	22	110	2900	37
<i>Clupea harengus</i>	22	110	3219	37
<i>Archocentrus centrarchus</i>	22	110	2644	37
<i>Esox lucius</i>	22	110	2848	37
<i>Monopterus albus</i>	22	110	2353	37
<i>Echeneis naucrates</i>	22	110	2314	37
<i>Paralichthys olivaceus</i>	22	110	3148	37
<i>Maylandia zebra</i>	22	110	2565	37
<i>Parambassis ranga</i>	22	110	2484	37
<i>Sander lucioperca</i>	22	110	2494	37
<i>Xiphophorus maculatus</i>	22	110	2231	37
<i>Nothobranchius furzeri</i>	22	110	2290	37
<i>Anabas testudineus</i>	22	110	2352	37
<i>Acanthochromis polyacanthus</i>	22	110	2797	37
<i>Anarrhichthys ocellatus</i>	22	110	2355	37
<i>Boleophthalmus pectinirostris</i>	22	110	1702	37
<i>Sparus aurata</i>	22	110	2361	37
<i>Oryzias melastigma</i>	22	110	2212	37
<i>Seriola dumerili</i>	22	110	2494	37
<i>Poecilia formosa</i>	22	110	2259	37
<i>Oreochromis niloticus</i>	22	110	2580	37
<i>Kryptolebias marmoratus</i>	22	110	2556	37
<i>Xiphophorus hellerii</i>	22	110	2240	37
<i>Poecilia latipinna</i>	22	110	2261	37
<i>Pundamilia nyererei</i>	22	110	2527	37

<i>Hippocampus comes</i>	22	110	2622	37
<i>Oreochromis aureus</i>	22	110	2579	37
<i>Amphiprion ocellaris</i>	22	110	2752	37
<i>Seriola lalandi dorsalis</i>	22	110	2481	37
<i>Austrofundulus limnaeus</i>	22	110	2541	37
<i>Puntigrus tetrazona</i>	25	110	2440	37
<i>Fundulus heteroclitus</i>	25	110	2476	37
<i>Cyprinodon variegatus</i>	28	110	2533	37
<i>Haplochromis burtoni</i>	31	110	2535	37
<i>Astatotilapia calliptera</i>	31	110	2571	37
<i>Gouania willdenowi</i>	37	110	2616	37
<i>Oryzias latipes</i>	40	110	2331	37
<i>Sphaeramia orbicularis</i>	43	110	2376	37
<i>Pygocentrus nattereri</i>	46	110	2649	37
<i>Astyanax mexicanus</i>	46	110	2791	37
<i>Colossoma macropomum</i>	46	110	2644	37
<i>Ictalurus punctatus</i>	46	110	3166	37
<i>Tachysurus fulvidraco</i>	46	110	3493	37
<i>Pangasianodon hypophthalmus</i>	46	110	3348	37
<i>Erpetoichthys calabaricus</i>	55	110	3662	37
<i>Perca flavescens</i>	58	110	2378	37
<i>Mastacembelus armatus</i>	64	110	2371	37
<i>Salmo salar</i>	67	110	3553	37
<i>Gadus morhua</i>	67	110	3151	37
<i>Etheostoma spectabile</i>	97	110	2457	37
<i>Scleropages formosus</i>	112	110	3412	37
<i>Myripristis murdjan</i>	112	110	2492	37
<i>Paramormyrops kingsleyae</i>	121	110	2929	37
<i>Carassius auratus</i>	148	110	3854	37
<i>Sinocyclocheilus grahami</i>	148	110	3330	37
<i>Sinocyclocheilus rhinoceros</i>	154	110	3449	37
<i>Sinocyclocheilus anshuiensis</i>	154	110	4202	37
<i>Electrophorus electricus</i>	283	110	2874	37

На рисунках 2 и 3 показано распределение длин части “кассетного” интрона до стоп-кодона и длин “кассетного” интрона, соответственно.

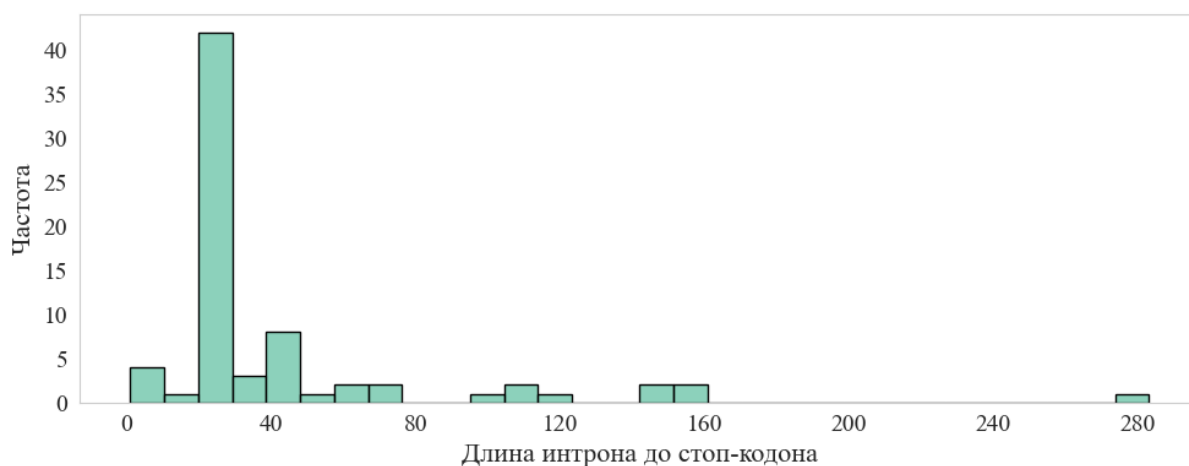


Рис. 2: Распределение длин части кассетного интрона до стоп-кодона у таксономической группы Actinopterygii

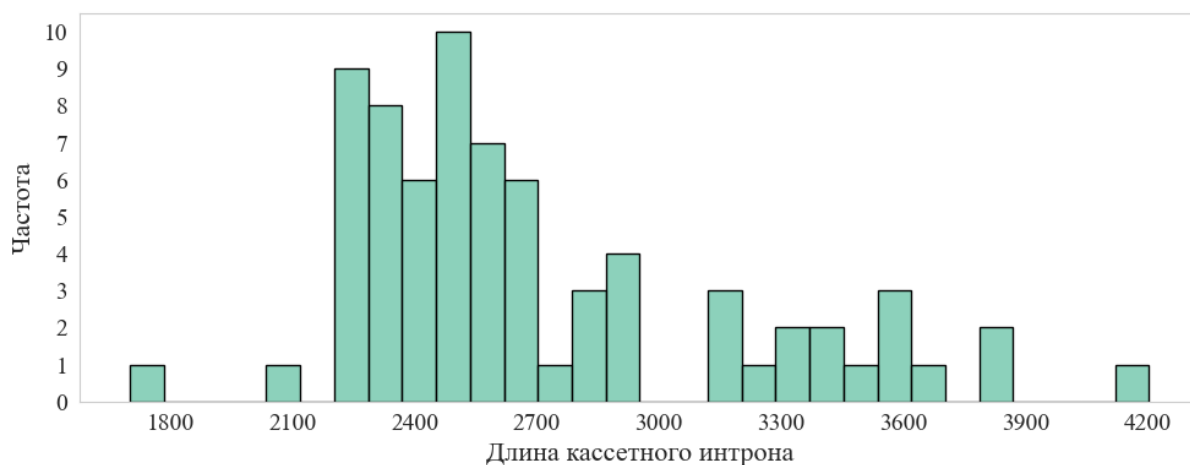


Рис. 3: Распределение длин кассетного интрона у таксономической группы Actinopterygii

На картинке 4 представлены результаты оценки “силы сайтов сплайсинга” - “ящички с усами”, отображающие распределение MaxEntScan score для таксонов более низкого ранга внутри группы Actinopterygii. Разбиение на подгруппы основано на их удаленности друг от друга. Порядок групп на графике не несет смысловой нагрузки.

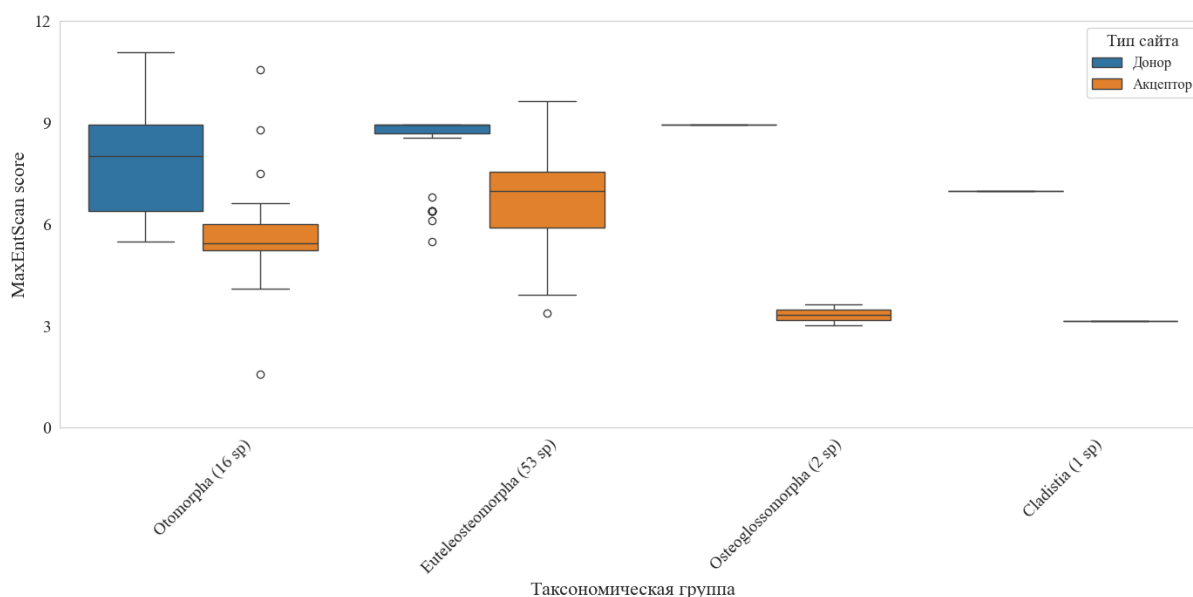


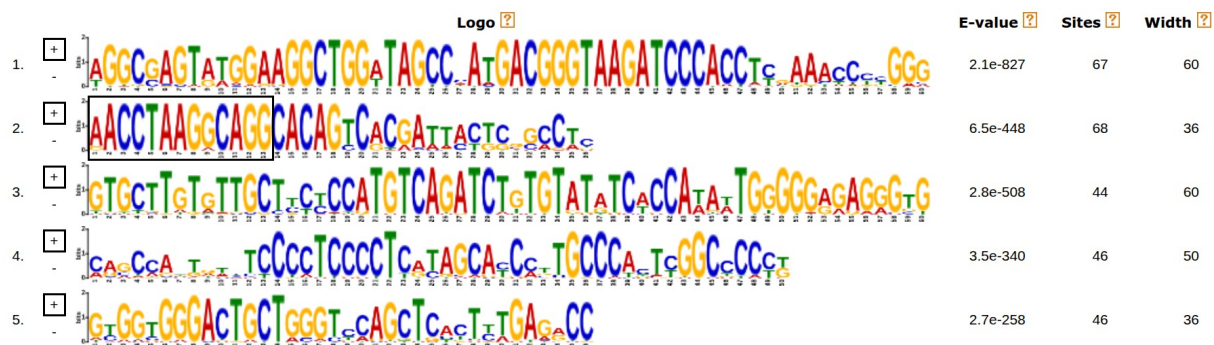
Рис. 4: Результаты проведения MaxEntScan для Actinopterygii.

Рисунок 5 демонстрирует результаты, полученные с помощью MEME Suite.

Найденные мотивы присутствуют не у всех 72 видов, их количество отображено в столбце Sites. Нас заинтересовал 2-й найденный мотив, так как его начало очень похоже на предложенную авторами статьи 2001 года [22] консенсусную последовательность для СТЕ из рисунка 6.

К сожалению, использование Tomtom для поиска найденных консервативных мотивов из “кассетного” интрона в базе данных не дало статистически значимых ре-

Рисунок 8 отображает результаты множественного выравнивания, а на рисунке 9 представлено филогенетическое дерево, построенное по результатам этого выравнивания.



Черным прямоугольником выделен участок, похожий на консенсусную последовательность СТЕ 6 из статьи 2001 года.

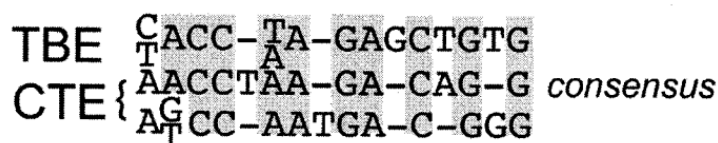


Рис. 6: Консенсусный СТЕ из статьи 2001 года

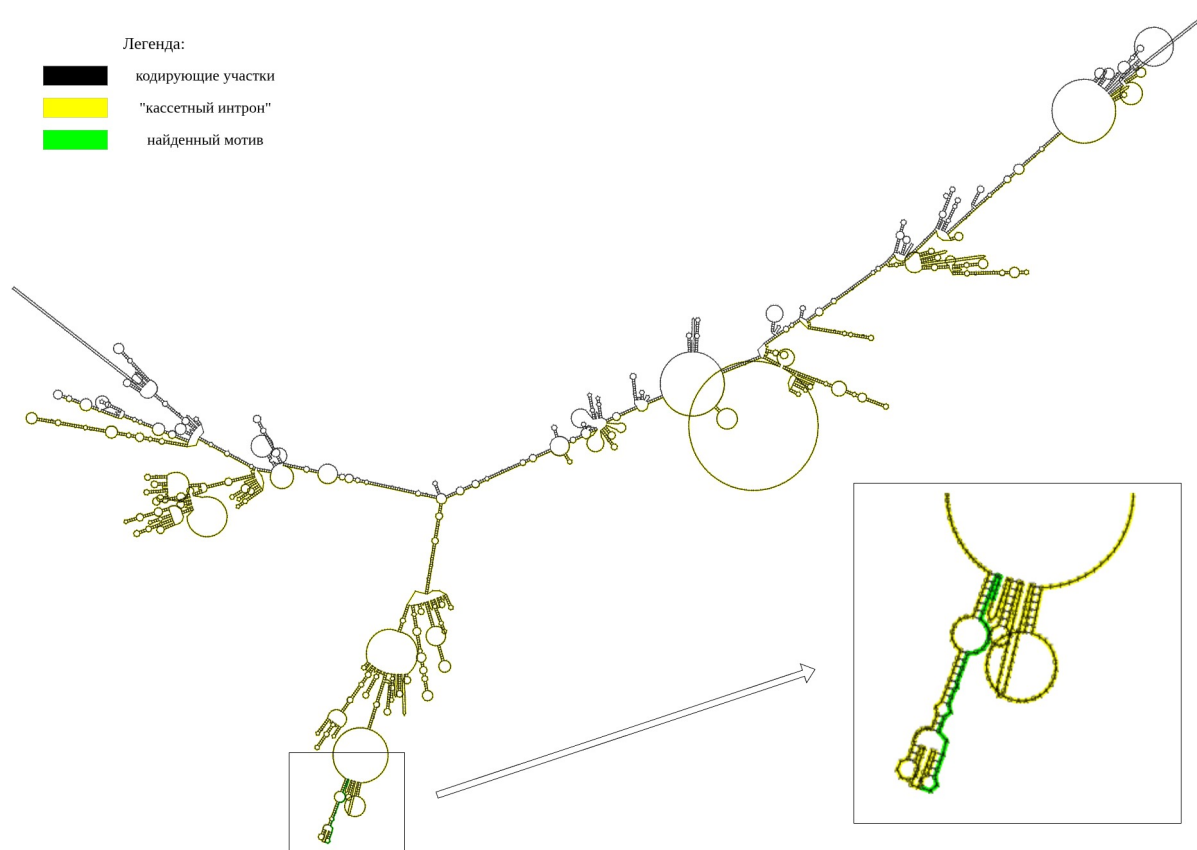


Рис. 7: Вторичная структура РНК-транскрипта для *Chanos chanos* из Otomorpha, содержащая кассетный интрон.



Рис. 8: Результаты множественного выравнивания для Actinopterygii.

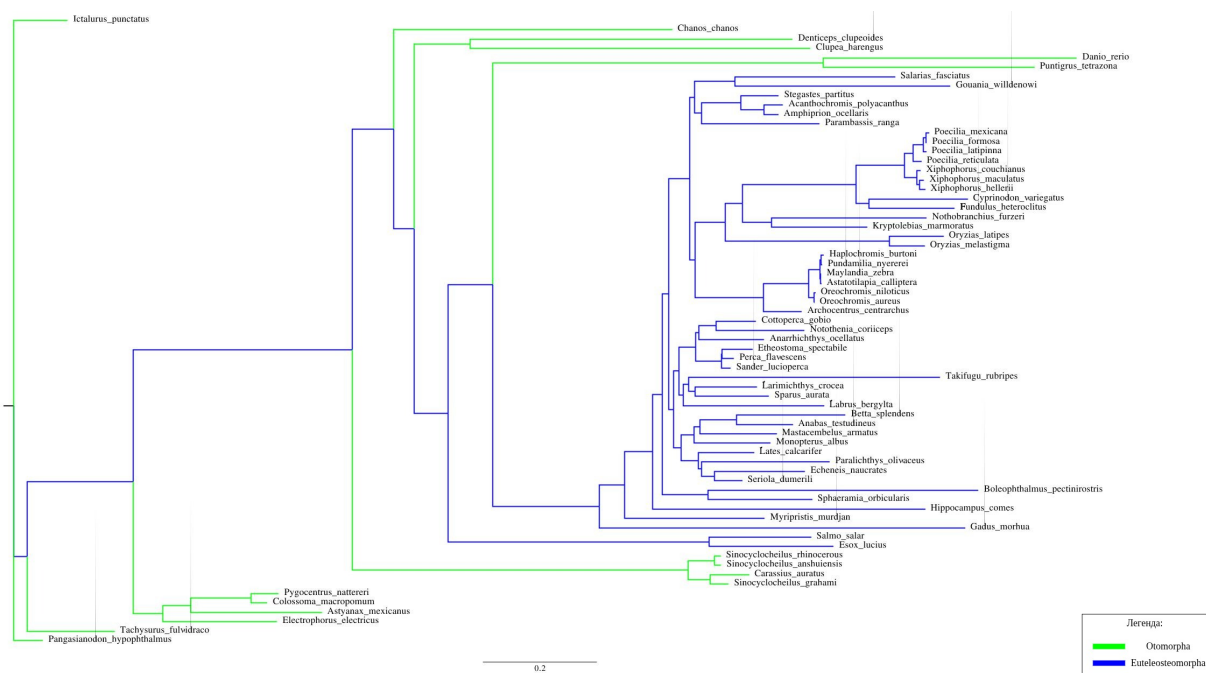


Рис. 9: Филогенетическое дерево для таксономической группы Actinopterygii

Обсуждение

Текст...

Список литературы

1. Database resources of the National Center for Biotechnology Information / E. W. Sayers, E. E. Bolton, J. R. Brister, [et al.] // *Nucleic Acids Research*. — 2022. — Vol. 50, no. D1. — P. D20–D26. — DOI: [10.1093/nar/gkab1112](https://doi.org/10.1093/nar/gkab1112). — URL: <https://doi.org/10.1093/nar/gkab1112>.
2. *McKinney W.* Data Structures for Statistical Computing in Python. — 2010.
3. *Python Software Foundation.* Python, Version 3.12. — 2023. — <https://www.python.org/downloads/release/python-3120/>.
4. Jupyter Notebooks – a publishing format for reproducible computational workflows / T. Kluyver [et al.]. — 2016. — DOI: [10.3233/978-1-61499-649-1-87](https://doi.org/10.3233/978-1-61499-649-1-87). — URL: <https://doi.org/10.3233/978-1-61499-649-1-87>.
5. Biopython: Freely available Python tools for computational molecular biology and bioinformatics / P. J. A. Cock [et al.] // *Bioinformatics*. — 2009. — Vol. 25, no. 11. — P. 1422–1423. — DOI: [10.1093/bioinformatics/btp163](https://doi.org/10.1093/bioinformatics/btp163). — URL: <https://doi.org/10.1093/bioinformatics/btp163>.
6. Exploring and retrieving sequence and metadata for species across the tree of life with NCBI Datasets / N. A. O’Leary [et al.] // *Scientific Data*. — 2024. — Vol. 11, no. 1. — P. 732. — DOI: [10.1038/s41597-024-03571-y](https://doi.org/10.1038/s41597-024-03571-y). — URL: <https://doi.org/10.1038/s41597-024-03571-y>.
7. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs / S. F. Altschul [et al.] // *Nucleic Acids Research*. — 1997. — Vol. 25, no. 17. — P. 3389–3402. — DOI: [10.1093/nar/25.17.3389](https://doi.org/10.1093/nar/25.17.3389). — URL: <https://doi.org/10.1093/nar/25.17.3389>.
8. *Katoh K., Standley D. M.* MAFFT multiple sequence alignment software version 7: improvements in performance and usability // *Molecular Biology and Evolution*. — 2013. — Vol. 30, no. 4. — P. 772–780. — DOI: [10.1093/molbev/mst010](https://doi.org/10.1093/molbev/mst010). — URL: <https://doi.org/10.1093/molbev/mst010>.
9. Unipro UGENE: a unified bioinformatics toolkit / K. Okonechnikov [et al.] // *Bioinformatics*. — 2012. — Vol. 28, no. 8. — P. 1166–1167. — DOI: [10.1093/bioinformatics/bts091](https://doi.org/10.1093/bioinformatics/bts091). — URL: <https://doi.org/10.1093/bioinformatics/bts091>.
10. The MEME Suite / T. L. Bailey [et al.] // *Nucleic Acids Research*. — 2015. — Vol. 43, W1. — W39–W49. — DOI: [10.1093/nar/gkv416](https://doi.org/10.1093/nar/gkv416). — URL: <https://doi.org/10.1093/nar/gkv416>.
11. Quantifying similarity between motifs / S. Gupta [et al.] // *Genome Biology*. — 2007. — Vol. 8, no. 2. — R24. — DOI: [10.1186/gb-2007-8-2-r24](https://doi.org/10.1186/gb-2007-8-2-r24). — URL: <https://doi.org/10.1186/gb-2007-8-2-r24>.

12. ViennaRNA Package 2.0 / R. Lorenz [et al.] // Algorithms for Molecular Biology. — 2011. — Vol. 6, no. 1. — P. 26. — DOI: [10.1186/1748-7188-6-26](https://doi.org/10.1186/1748-7188-6-26). — URL: <https://doi.org/10.1186/1748-7188-6-26>.
13. Yeo G., Burge C. B. Maximum entropy modeling of short sequence motifs with applications to RNA splicing signals // Bioinformatics. — 2004. — Vol. 20, no. 3. — P. 327–335. — DOI: [10.1093/bioinformatics/btg005](https://doi.org/10.1093/bioinformatics/btg005). — URL: <https://doi.org/10.1093/bioinformatics/btg005>.
14. IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era / B. Q. Minh [et al.] // Molecular Biology and Evolution. — 2020. — Vol. 37, no. 5. — P. 1530–1534. — DOI: [10.1093/molbev/msaa015](https://doi.org/10.1093/molbev/msaa015). — URL: <https://doi.org/10.1093/molbev/msaa015>.
15. Rambaut A. FigTree v1.4.4. — 2018. — Institute of Evolutionary Biology, University of Edinburgh. <http://tree.bio.ed.ac.uk/software/figtree/>.
16. QuantStack, contributors mamba. Mamba: The Fast Cross-Platform Package Manager. — 2024. — <https://github.com/mamba-org/mamba>.
17. GitHub, Inc. GitHub. — 2008. — URL: <https://github.com>.
18. Lamport L. LaTeX: A Document Preparation System. — 2nd ed. — Reading, Massachusetts : Addison-Wesley, 1994.
19. TablesGenerator.com. Tables Generator – LaTeX Tables Editor. — 2025. — URL: <https://www.tablesgenerator.com>.
20. diagrams.net. draw.io - Online Diagram Software. — 2025. — URL: <https://www.diagrams.net/>.
21. Canonical Ltd. Ubuntu 22.04 LTS (Jammy Jellyfish). — 2022. — <https://releases.ubuntu.com/22.04/>.
22. Replication of Human Herpesvirus 6A and 6B Is Associated with Distinct Nuclear Domains / F. Tajima [et al.] // Journal of Virology. — 2001. — Vol. 75, no. 12. — P. 5567–5575. — DOI: [10.1128/JVI.75.12.5567-5575.2001](https://doi.org/10.1128/JVI.75.12.5567-5575.2001). — URL: <https://doi.org/10.1128/JVI.75.12.5567-5575.2001>.