

1 **AptamerRunner: An accessible aptamer structure prediction and clustering algorithm for**
2 **visualization of selected aptamers**

3
4 Dario Ruiz-Ciancio^{1,2,3}, Suresh Veeramani^{4,5}, Rahul Singh⁶, Eric Embree⁷, Chris Ortman⁸, Kristina
5 W. Thiel^{5,9}, and William H Thiel^{4,*}

6 ¹ Instituto de Ciencias Biomédicas (ICBM), Facultad de Ciencias Médicas, Universidad Católica
7 de Cuyo, Av. José Ignacio de la Roza 1516, Rivadavia, 5400, San Juan, Argentina.

8 ² National Council of Scientific and Technical Research (CONICET), Godoy Cruz 2290,
9 C1425FQB Ciudad Autónoma de Buenos Aires Argentina.

10 ³ Cancer Genome Engineering Group, Vall d'Hebron Institute of Oncology (VHIO), Barcelona
11 08035, Spain.

12 ⁴ Department of Internal Medicine, University of Iowa, Iowa City, IA 52242, USA.

13 ⁵ Holden Comprehensive Cancer Center, University of Iowa, Iowa City, IA, 52242, USA.

14 ⁶Department of Computer Sciences, University of Iowa, Iowa City, IA 52242, USA.

15 ⁷Carver College of Medicine, University of Iowa, Iowa City, IA 52242, USA.

16 ⁸Institute for Clinical and Translational Science, University of Iowa, Iowa City, IA 52242, USA.

17 ⁹Department of Obstetrics and Gynecology, University of Iowa, Iowa City, IA 52242, USA.

18 *Corresponding author. Email: william-thiel@uiowa.edu

19
20 **Abstract**

21 Aptamers are short single-stranded DNA or RNA molecules with high affinity and specificity for
22 targets and are generated using the iterative Systematic Evolution of Ligands by EXponential
23 enrichment (SELEX) process. Next-generation sequencing (NGS) revolutionized aptamer
24 selections by allowing a more comprehensive analysis of SELEX-enriched aptamers as compared
25 to Sanger sequencing. The current challenge with aptamer NGS datasets is identifying a diverse
26 cohort of candidate aptamers with the highest likelihood of successful experimental validation.
27 Herein we present AptamerRunner, an aptamer sequence and/or structure clustering algorithm that
28 synergistically integrates computational analysis with visualization and expertise-directed decision
29 making. The visual integration of networked aptamers with ranking data, such as fold enrichment
30 or scoring algorithm results, represents a significant advancement over existing clustering tools by
31 providing a natural context to depict groups of aptamers from which ranked or scored candidates

32 can be chosen for experimental validation. The inherent flexibility, user-friendly design, and
33 prospects for future enhancements with AptamerRunner has broad-reaching implications for
34 aptamer researchers across a wide range of disciplines.

35

36 **Introduction**

37 Aptamers are short synthetic RNA or DNA oligonucleotides that recognize target epitopes with
38 specificity and affinity analogous to antibody-antigen interactions¹. Applications of aptamers are
39 broad-reaching and include biosensors², research reagents^{3, 4}, tools for mechanistic discovery³,
40 diagnostics⁵, delivery platform⁶ and therapeutics⁵. Recently the aptamer avacincaptad pegol, that
41 target complement C5, was FDA-approved for the treatment of geographic atrophy⁷. Aptamers are
42 generated using an *in vitro* process known as Systematic Evolution of Ligands by EXponential
43 enrichment (SELEX)^{8, 9}. The SELEX process now includes numerous variations with new SELEX
44 strategies constantly being developed^{1, 10}. At the completion of SELEX, selected aptamers are
45 identified by sequencing, with the field shifting from Sanger sequencing towards next-generation
46 sequencing (NGS) platforms. NGS yields hundreds of millions of reads, with each read containing
47 the entirety of an aptamer sequence. Thus, the aptamers enriched during SELEX can now be
48 interrogated to a degree not achievable with Sanger sequencing^{11, 12}. However, these large NGS
49 datasets have created new challenges in terms of how parse the millions of reads to identify the
50 best candidate aptamers to then validate experimentally. Testing thousands or even hundreds of
51 aptamer candidates is still unattainable for most aptamer researchers; therefore, identification of
52 the top candidates is paramount. To address this need, several bioinformatics approaches specific
53 for aptamer NGS datasets have emerged^{13, 14}. The analysis of an aptamer NGS dataset includes
54 processing the FASTQ data and application of strategies to identify candidate aptamers using
55 various motif identification, scoring, and clustering algorithms^{13, 14}. Bioinformatics tools to
56 identify candidate aptamers are frequently applied in concert, for example clustering is used to
57 identify separate groups of aptamers¹⁵⁻¹⁷, followed by ranking aptamers within these groups using
58 fold enrichment or scoring algorithms^{18, 19}.

59 The central theory behind aptamer clustering is that aptamers that are closely related based on their
60 sequence and predicted structures are likely to target the same epitope^{13, 14}. The earliest efforts to
61 cluster aptamers used sequence alignment algorithms such as ClustalW²⁰⁻²², but these data can be

62 difficult to interpret and this method does not take into account the predicted structures of the
63 aptamers. Our group introduced the concept of clustering aptamers by either sequence relatedness
64 using Levenshtein edit distance or by predicted secondary structure relatedness using tree
65 distance¹⁷. We applied a clustering strategy, termed the *all-vs-all* approach, wherein all aptamers
66 within a dataset are compared to each other. Clusters were defined as the aptamers that interconnect
67 within a threshold distance measure (e.g., edit distance of 3). However, this early clustering
68 algorithm was not easily accessible and thus has not been widely adopted. The current prevailing
69 clustering algorithms include AptamerCluster^{23, 24}, FASTAptamer²⁵, and FASTAptameR 2.0²⁶. These
70 tools generate networks of related aptamers using either Hamming or Levenshtein edit distance
71 and introduced a new concept for clustering aptamers termed the *seed approach*. The seed
72 approach generates networks of aptamers centered around a *seed sequence*, which is defined as the
73 sequence with the most reads within an aptamer NGS dataset (i.e., most abundant sequence). All
74 aptamers that connect to the seed sequence within a threshold edit distance measure are designated
75 as a cluster. This process is iterated by removing the initial seed sequence and all connected
76 aptamer sequences from further analysis, and the next most abundant sequence becomes the seed
77 for the next cluster. These seed-based aptamer clustering algorithms are significantly more
78 accessible than our initial algorithm, and the seed approach has significant computational
79 advantages over the all-vs-all approach. A limitation of the seed approach is that it is a greedy
80 process that can potentially miss important inter-aptamer relationships identified by the all-vs-all
81 clustering approach. In addition, the available seed-based algorithms do not consider structure
82 when generating networks of related aptamers. The output from these algorithms is text-based
83 rather than graphically represented, which severely limits interpretation of the clustering results
84 and prevents integration of ranking data that is necessary to identify candidates within the groups
85 of interconnected aptamers.

86 Herein, we introduce *AptamerRunner*, an accessible aptamer structure prediction and clustering
87 algorithm for the visualization of networked aptamers. AptamerRunner was designed based on the
88 principles of *experiential computing*^{27, 28}, which is founded on the idea that an understanding of
89 complex biological data comes from integrating user expertise with algorithmic processing via
90 data visualization and user-data interactions. Using the paradigm of experiential computing, we
91 designed AptamerRunner so that it has the flexibility to ensure that aptamer researchers can apply
92 different clustering strategies to suit their needs, with customizable visualization support enabling

93 **user-specific data interpretation.** AptamerRunner includes several novel clustering features not
94 previously available: 1) the option to use either the all-vs-all approach or the seed approach. From
95 the perspective of computational complexity, if, we have n aptamer sequences, with m being the
96 length of the longest aptamer, the time complexity of the all-vs-all approach is $O(n^2m^2)$, due to the
97 complexity of computing the Levenshtein distance. The seed approach, on the other hand, involves
98 identifying the most abundant aptamers, which requires determining the frequency of each unique
99 aptamer sequence present in the dataset. Using hashing, this can be obtained in $O(nm)$ time. If k
100 different seeds are considered, then the complexity of the approach is $O(nmk)$; 2) the option to
101 interrogate both sequence and structure relatedness simultaneously by applying logical operators
102 (AND, OR) with edit and tree distance thresholds; and 3) inclusion of distance measure data as
103 metadata within the edges of the interconnected aptamer sequences. To provide easier access and
104 functionality, AptamerRunner and all dependencies have been packaged into a Docker image²⁹
105 that is operated by command line using a platform-independent .NET script. The AptamerRunner
106 clustering algorithm outputs results as an eXtensible Graph Markup and Modeling Language
107 (XGMML) file that permits graphical representation of the clustering results using network
108 analysis programs such as the open-source Cytoscape program³⁰. Importantly, through the
109 graphical representation of the clustering data, additional information such as ranking data can be
110 overlayed onto the networked aptamers to facilitate an integrated analysis whereby clusters of
111 aptamers and ranking data can be interpreted at the same time. This integration of clustering with
112 ranking data presents a novel analysis of selected aptamers that aids in the identification of the
113 best candidates.

114 **Results**

115 AptamerRunner Overview

116 AptamerRunner is a .NET program coded in C# that generates the Docker commands to adapt to
117 various operating system constraints (**Figure 1A**). AptamerRunner will check the Docker
118 repository and download the most recent AptamerRunner Docker image (for detailed instructions
119 to use AptamerRunner refer to supplemental methods). The AptamerRunner Docker image is
120 comprised of two independent python algorithm components (**Figure 1B and 1C**). The first
121 component predicts the secondary structure of RNA or DNA aptamer sequences (**Figure 1B**). The
122 second component calculates aptamer relatedness to generate networks of related aptamer
123 sequences (**Figure 1C**). This segmentation permits each component to be implemented
124 independently and provides additional flexibility to the user. **For more technically proficient users,**
125 **both python algorithm components can be operated by command line independent of Docker.** Once
126 the AptamerRunner Docker container has completed running the structure predication algorithm
127 or clustering algorithm, AptamerRunner will close the AptamerRunner Docker container.

128 To demonstrate the utility of AptamerRunner and compare to other aptamer clustering tools , we
129 used a published aptamer NGS dataset from a selection against B-cell leukemia cells¹⁶. With these
130 data, we applied the various AptamerRunner clustering options with either edit distance of 1 or
131 tree distance of 0, when applicable. When comparing AptamerRunner against FASTAptamer,
132 FASTAptameR 2.0 and AptaCluster, we applied comparable clustering parameters using an edit
133 distance of 1 with the seed approach.

134 Visualization of AptamerRunner clustering results

135 A central principle of experiential computing is combining algorithms with data visualization and
136 user-data interactions to facilitate the expertise of a user to interpret complex data. **We enabled the**
137 **visualization of AptamerRunner clustering results by exporting the clustering data into the**
138 **XGMML format, which can be imported into the network visualization software Cytoscape (see**
139 **supplemental methods for specific details).** Once the clustering results have been imported into
140 **Cytoscape, they can be visualized using the multitude of Cytoscape's built-in network layout**
141 **functions. This approach not only provides a clear and organized representation of the data but**
142 **also facilitates deeper analysis through the various customization and analytical tools available**
143 **within Cytoscape. AptamerRunner and Cytoscape's visual representation of networked aptamers**

provides a natural context to interpret clustering data that is significantly easier to interpret than text-based results. For example, when the seed option is used to generate the networks of related aptamers, the seed sequence is clearly identifiable as the central node and the number of sequences connecting to each seed sequences are easily discerned (**Figure 2A and 2B**). As compared to the seed approach, the all-vs-all option produces fewer but more complex clusters of inter-related aptamer sequences (**Figure 2C and 2D**). The all-vs-all approach clustering results can include large, complex hairball networks (**Figure 2C**, largest cluster) and smaller clusters with naturally occurring central nodes. Importantly, clustering conducted using either the seed approach or the all-vs-all approach can provide different perspectives of the same data.

A second important improvement within AptamerRunner, as compared to other clustering algorithms, is the introduction of logical operators AND and OR with edit distance and tree distance. The AND function and the OR function can be applied to gain insight into potential functional aspects of an aptamer. The AND function can reveal areas of nucleotide substitutions that are well-tolerated or that impart beneficial properties. Conversely, the OR function can reveal nucleotide changes that have a significant impact on the predicted structure of an aptamer. The use of the AND logical operator places a higher stringency, and the use of the OR logical operator less stringency, onto the networks of related aptamer generated. For example, clustering the data presented in **Figure 2** using the AND logical operator with an edit distance 1 and tree distance 0 yields smaller, more constrained networks of related aptamers when using the seed (**Figure 3A**) and all-vs-all approaches (**Figure 3B**). Conversely, the OR logical operator yields larger less constrained networks of related aptamers (**Figure 3C and 3D**).

Comparison of AptamerRunner to other aptamer clustering algorithms: FASTAptamer, FASTAptamer 2.0 and AptaCluster

We next compared the capabilities and output of AptamerRunner to the capabilities and output of FASTAptamer, FASTAptamer 2.0 and AptaCluster^{24-26,31}. Details and features of FASTAptamer, FASTAptamer 2.0, AptaCluster and AptamerRunner are summarized in **Table 1**. The FASTAptamer clustering algorithm (*FASTAptamer-Cluster*) and FASTAptamer 2.0 clustering algorithm (*cluster module*) determine sequence families of aptamers by Levenshtein edit distance using the seed approach. FASTAptamer 2.0, an update to FASTAptamer, operates through an offline web-browser that accesses a Docker container. FASTAptamer 2.0, like FASTAptamer, applies a seed approach using edit distance only, but has a computationally faster clustering

175 algorithm and a cluster visualization function (*cluster diversity module*) that generates a principal
176 component plot (PCA) using a k-mer matrix of the clustered aptamer sequences. AptaCluster is
177 provided as a component of the AptaSuite package with a GUI that is accessible as a Java
178 application³². AptaCluster filters input data using a local sensitivity hash function and clusters
179 aptamers by Hamming edit distance using the seed approach.

180 Of note, FASTAptamer, FASTAptameR 2.0 and AptaCluster are incapable of generating clusters
181 based upon predicted structures, nor do they compare all aptamers to each other regardless of their
182 enrichment during SELEX (e.g., all-vs-all approach). AptaSuite, of which AptaCluster is a
183 component, includes an algorithm AptaTrace that identifies sequence-structure motifs as sequence
184 logos with secondary structure probability profiles³³. The AptaTrace algorithm does not perform
185 clustering of these predicted secondary structure and thus is not included in the comparison to
186 AptamerRunner. To compare AptamerRunner to FASTAptamer, FASTAptameR 2.0 and
187 AptaCluster, we applied the seed approach using only edit distance with no logical operators, and
188 we used aptamer NGS data from a selection against B-cell leukemia cells¹⁶.

189 The FASTAptamer and FASTAptameR 2.0 clustering algorithms produced identical modified
190 FASTA formatted text files (**Figure 4A**). Each sequence identifier within the FASTA file (denoted
191 with the “>” annotation) contains the following data features for each aptamer sequence: sequence
192 ID, raw read count, normalized read count, cluster ID, rank within the cluster, and edit distance
193 from the seed sequence. For example, in **Figure 4A**, the first listed sequence has an identifier of
194 “>1-3290398-899773.77-1-1-0” which denotes a sequence ID of 1, 3290398 raw read count,
195 899773.77 normalized read count, cluster ID of 1, rank within cluster of 1, and edit distance of 0
196 from seed. Within a given cluster ID, the seed sequence is the listed first, is the highest ranked
197 aptamer sequence and will have an edit distance from the seed of 0. Subsequent clusters can be
198 identified by the cluster ID. For example, in **Figure 4A**, the second cluster begins with the
199 following identifier “>2-41052-11225.85-2-1-0”. FASTAptameR 2.0 includes the option to output
200 a CSV data table (**Figure 4B**), which makes the clustering results sortable and easier to parse
201 separate clusters.

202 AptaCluster outputs two data files: a modified FASTA file and a data table of seed sequences
203 (**Figure 4C, D**). The AptaCluster modified FASTA file (**Figure 4C**) includes a sequence identifier
204 line (denoted by “>>”) with the cluster ID followed by the aptamer sequences within that cluster

205 (denoted by “>”), starting with the seed sequence. Read counts are included on the sequence line
206 following the aptamer variable region sequence. For example, in **Figure 4C**, the first cluster is
207 identified as Cluster 103475, with a total of 640,120 read counts among all sequences within that
208 cluster. The first listed sequence, which is the seed of Cluster 103475, is named Aptamer_2 and
209 has a read count of 638,988. The AptaCluster seed sequence data table includes the cluster ID,
210 seed sequence, and information about the clusters from each selection round; example data are
211 provided for Round 9 in **Figure 4D**. For each selection round, AptaCluster seed table provides the
212 proportion of a cluster relative to other identified clusters (“R9 Size”), the number of sequences
213 within that cluster (“R9 diversity”) and the total number of normalized read counts per million
214 (“R9 CPM”).

215 By comparison, AptamerRunner yields a visual output of clusters as depicted in **Figure 2A**. In this
216 example using the same data that were analyzed by FASTAptamer, FASTAptameR 2.0 and
217 AptaCluster, the seed sequence was the central node in the top left cluster. Note that all programs
218 identified the same sequence as the seed, but the visual output by AptamerRunner significantly
219 improves data interpretation and candidate aptamer selection because all clusters can be easily
220 viewed simultaneously. Users can determine aptamer sequence, predicted secondary structure, and
221 any other imported metadata (e.g., fold enrichment) by simply clicking on a node (diagrammed in
222 **Supplemental Figure S1**). This interactive visualization is not possible with text-based results.
223 FASTAptameR 2.0 does include functions to visualize an analysis of the clustering data within the
224 *diversity* module. The diversity module provides a series of graphs (cluster metaplots) that depict
225 information about the population of clustered aptamers (**Supplemental Figure S2A**, sequence
226 count, read count and average LED) and a PCA plot of a k-mer matrix (**Supplemental Figure S2B**
227 **and S2C**). While these PCA plots can provide insight into relative diversity of the different
228 clusters, this function is limited to plotting no more than 15 clusters concurrently, and users must
229 cross-reference text-based results to determine the identify of specific aptamer sequences within
230 the PCA plot. Whereas results from AptamerRunner within Cytoscape are interactive. The
231 limitations of FASTAptamer-Cluster, FASTAptameR 2.0 and AptaCluster text-based outputs
232 highlight the importance of visualizing clustering data to provide a natural context for representing
233 the different clusters of aptamers.

234 Improving candidate aptamer selection by integrating fold enrichment data and scoring algorithm
235 results onto AptamerRunner clustered aptamers

236 A limitation of both the text-based and visualized clustering results is that they do not provide
237 insight into which clusters are most likely to yield the best aptamers. Clustering separates aptamers
238 into groups that likely target the same epitope, but additional data are necessary to score and rank
239 the aptamers to identify ideal candidates from within each cluster of aptamers. We hypothesized
240 that integration of fold enrichment or data from scoring algorithms significantly enhances
241 candidate selection when integrated with clustering results. Unfortunately, the text-based outputs
242 of FASTAptamer and FASTAptamer 2.0 do not allow for easy integration of fold enrichment or
243 scoring data. The AptaCluster seed sequence data table provides selection round normalized read
244 counts that can be used to calculate the overall round-to-round enrichment of each cluster (see
245 **Figure 4D**), but data from this table must be manually cross-referenced with the clustering results
246 to determine the enrichment of all sequences within the cluster.

247 AptamerRunner overcomes the limitations of text-based results through visual integration of fold
248 enrichment data and results from scoring algorithms. Data tables were imported into Cytoscape;
249 these tables contained log10 normalized read counts of round 9 and the log2 fold enrichment
250 between selection rounds (e.g., round 2 to round 5) of aptamers that were clustered using the seed
251 approach with an edit distance of 1. As shown in **Figure 5**, such visual overlays of node size and
252 color provide easy to interpret visual cues of individual aptamer abundance and enrichment during
253 the SELEX process. Groups of aptamers showing positive enrichment (green) from negative
254 enrichment (yellow) are easily discernable. Furthermore, the log2 fold enrichment can be easily
255 evaluated between different selection rounds. The B-cell SELEX process sequence enrichment
256 exhibited a sigmoidal curve, with rounds 0 to 2 representing the initial exponential phase, rounds
257 2 to 5 representing the linear phase and rounds 5 to 9 representing the asymptotic phase (**Figure**
258 **5A**). Interestingly, with this visual comparison we observed positive enrichment from rounds 0-2
259 of the selection (**Figure 5B**), but the most the most interesting changes of log2 fold enrichment
260 seem to occur between selection rounds 2 to 5 (**Figure 5C**) and rounds 5 to 9 (**Figure 5D**). For
261 example, when comparing round 5 to 9, the largest group of clustered aptamers (top left cluster)
262 has a significant diversity of log2 fold enrichments that ranges from -9.95 to 5.1, and the seed
263 sequence remained close to 0 (**Figure 5D**). These data suggest that the seed sequence may not be
264 the ideal candidate from this cluster, but rather one of the other aptamer sequences within 1 edit
265 distance of the seed sequence with a positive log2 fold enrichment is preferred for subsequent
266 testing. By contrast, the seed sequence in the top right cluster is a desirable candidate based on its

positive log2 fold enrichment from round 5 to round 9. This type of analysis of AptamerRunner clustering results mapped with log2 fold enrichment data was used with the Ruiz-Ciancio *et. al.* study to identify 38 candidate aptamers from 38 separate sequence clusters (all-vs-all approach with edit distance 1) or structure clusters (all-vs-all approach with tree distance 3). Within each cluster, candidates were defined as the aptamer sequences with the highest log2 fold enrichment from rounds 2 to 5 or rounds 5 to 9. We favor the all-vs-all approach based on the supposition that the seed approach may miss interesting inter-aptamer relationships; however, we present data in this head-to-head comparison by using the seed approach since other aptamer clustering tools cannot accommodate the all-vs-all approach.

Bioinformatics tools such as MPbind¹⁹ and RaptRanker¹⁸ score and rank aptamers enriched during SELEX, and these data can be overlayed onto AptamerRunner clustering data in Cytoscape. The MPbind scoring algorithm generates a combined meta z-score of aptamer sequences based on relative motif enrichment and abundance of the final aptamer selection round. RaptRanker evaluates aptamer sequence motif and structure of subsequence groups to generate an average motif enrichment score (AME). Integrating data from scoring algorithms enables us to examine the predicted relative aptamer affinities, which can be used to select candidates within each cluster. Interestingly, MPbind and RaptRanker demonstrates significant variation in the predicted affinities of aptamers from the B-cell SELEX dataset (**Figure 6**). Specifically, MPBind predicted that most of the clusters had high affinity for their targets: the majority of the nodes were visualized as red (meta z-scores of 40-60, **Figure 6A**), whereas RaptRanker predicted fewer high affinity aptamers, with most nodes visualized in the blue to yellow spectrum (average motif enrichment of 0-5, **Figure 6B**). Since MPBind and RaptRanker use different principles for scoring, we asked if the highest scoring aptamers were similar between the two algorithms by plotting the scores as interactive scatter plots in Cytoscape. Clusters of aptamers that were scored high by both the algorithms could be identified, but, consistent with the visual representation of the clusters, many more aptamers scored highly by MPBind vs. RaptRanker (**Figure 6C**).

293 Application of AptamerRunner metadata to re-organize networked aptamers

One rationale for using the seed approach rather than the all-vs-all approach to cluster aptamers is that the all-vs-all approach will frequently generate large, disorganized hairball clusters of aptamers as highlighted in **Figure 7A** (grey nodes and red edges), which are difficult to interpret. To address this limitation with all-vs-all data, the metadata generated by AptamerRunner during

298 clustering can be used to deconstruct and re-organize the hairballs. For example, the hairball cluster
299 within **Figure 7A** was isolated and the aptamer sequences re-organized using the metadata
300 property of structure relatedness. This approach identified groups of aptamer sequences that are
301 related by edit distance 1 and have identical structures (tree distance = 0), as shown by edges
302 colored blue (**Figure 7B**). Several groups of the re-organized aptamer sequences also exhibited
303 positive log2 fold enrichment from round 5 to round 9 (green nodes). We also asked if this de-
304 convolution approach could be used for other metadata features. In the example shown in
305 **Supplemental Figure S3**, we tested whether enriched aptamers within this hairball are structurally
306 similar. First, we re-organized the hairball in **Supplemental Figure S3A** using the log2 fold
307 change enrichment observed between round 5 to round 9 (Log2 R5:R9, **Supplemental Figure**
308 **S3B**). Next, we excluded any aptamers with a log2 R5:R9 less than ≤ 0 , and the remaining subset
309 of aptamers were reorganized for identical predicted structures (tree distance = 0, **Supplemental**
310 **S3C**). This approach identified multiple sub-groups of aptamer sequences within the hairball
311 network that were positively enriched during SELEX and are closely related by both sequence and
312 structure. Taken together, these capabilities of AptamerRunner to incorporate metadata for clusters
313 to be visualized in Cytoscape allow for a more in-depth analysis and identification of candidate
314 aptamers. Use of different parameters provides greater resolution of the enriched sub-groups of
315 aptamers, which is a major improvement by AptamerRunner over text-based clustering algorithms.

316 **Analysis of the Interleukin(IL)-10 aptamer dataset by AptamerRunner:**

317 To evaluate the broader utility of AptamerRunner, we analyzed a publicly available aptamer NGS
318 dataset from the NIH Sequence Read Archive (SRA). This database was used to describe
319 AptaCluster and AptaMut²⁴, as well as to assess potentially beneficial mutations of the IL-10
320 aptamers³⁴. This database contained sequencing data across five selection rounds with
321 approximately 16.7 million reads representing approximately 11.7 million aptamer sequences.
322 Using AptaCluster within AptaSuite, we generated clustering results based on the reported
323 AptaCluster constraints (see supplemental methods).

324 By comparison, we generated a non-redundant database from the NCBI SRA data, compiling read
325 counts of the ~11.5 million aptamer sequences across the five selection rounds. From this dataset,
326 we identified 2,140 unique aptamer sequences for clustering, representing ~4.5 million.

327 Using AptamerRunner, we generated clusters of the IL-10 aptamers employing an edit distance of
328 5, utilizing both the seed approach and the all-vs-all approach (**Figure 8A and 8B**). We then
329 overlaid the published dissociation constant (K_D) values onto the networked aptamers to ascertain
330 which sequences had been experimentally evaluated for their ability to bind IL-10. Our results
331 indicate that one larger cluster and seed sequence identified by AptamerRunner appears to have
332 remained untested. This seed sequence was the third most abundant aptamer from the fifth
333 selection round and is evident in the Aptasuite aptamer pool data. These results highlight
334 AptamerRunner's ability to cluster and visualize seed sequences that may be important aptamers
335 for experimental evaluation.

336 Using AptamerRunner, we next evaluated how the IL-10 aptamers cluster based on predicted
337 secondary structure. We first established an appropriate tree distance measure by examining the
338 distribution of tree distances found within the edit distance 5 seed edges (**Supplemental Figure**
339 **S4A**). These data indicate that most aptamer sequences within an edit distance of 5 are structurally
340 closely related, with tree distances clustering around 10. Beyond this point, the histogram begins
341 to level off, suggesting a diminishing return in structural similarity. Building on these insights, we
342 evaluated the IL-10 aptamers using AptamerRunner with a tree distance of 10 and the all-vs-all
343 approach (**Figure 8C**). The tree distance of 10 generated numerous clusters that overlapped with
344 several identified clusters at edit distance 5. However, a notable difference emerged; specifically,
345 two clusters, D and K, which are greater than edit distance 5, exhibited similar predicted secondary
346 structures. These results indicate that the IL-10 aptamers may cluster most effectively when both
347 edit distance and tree distance are considered. Consequently, using AptamerRunner, we clustered
348 the IL-10 aptamers within an edit distance of 5 and a tree distance of 10 employing the all-vs-all
349 strategy with the AND logical operator (**Figure 8D**). This methodology provided significant
350 granularity in the clustering outcomes, revealing that clusters with similar K_D values contained
351 multiple aptamers. Importantly, the analysis using AptamerRunner, which considers both edit
352 distance and tree distance, indicates the presence of several clusters lacking representatives tested
353 for IL-10 binding that show substantial enrichment from rounds four to five (**Supplemental**
354 **Figure S4B**).

355

356

357 **Discussion**

358 The rationale for clustering aptamer sequences is to discern which aptamers likely target the same
359 epitope and conversely, which aptamers likely target different epitopes. By understanding how
360 aptamer sequences are related, and not related, by sequence and structure, a diverse cohort of
361 aptamers can be identified for experimental validation. The AptamerRunner clustering algorithm
362 applies the principles of experiential computing to support expertise-driven decision making for
363 clustering and identification of candidate aptamer sequences. Novel features of AptamerRunner
364 that facilitate expertise-driven decision making include retention of distance measures as metadata
365 and the incorporation of logical operators (AND, OR) for clustering. Maintaining distance
366 measures as metadata provides additional opportunities to analyze, re-evaluate and interpret
367 clustering results. Within Cytoscape, additional data such as the log2 fold enrichment of aptamer
368 sequence across selection rounds and data from other aptamer bioinformatics scoring algorithms
369 can be mapped onto the networks generated by AptamerRunner to aid in the interpretation of the
370 clustering results, further supporting the experiential computing goal. Having the AptamerRunner
371 Docker container controlled by a .NET program presents a simple method to enable users to access
372 AptamerRunner. Users only need to use a command line interface to initiate the .NET script, which
373 then dynamically generates all Docker commands necessary to run AptamerRunner's secondary
374 structure prediction or clustering algorithms. Taken together, the innovative features of
375 AptamerRunner enable a more in-depth analysis of aptamer NGS data and allows for better
376 identification of candidate aptamers. In addition, AptamerRunner can be used to ask new questions
377 about how sequence and structure relatedness contribute to library convergence during the SELEX
378 process.

379 The superiority and flexibility of AptamerRunner is highlighted by two recent publications that
380 made use of AptamerRunner clustering capabilities in fundamentally different ways^{15, 16}. Ruiz-
381 Ciancio *et. al.*¹⁶, using the same NGS dataset as in the present study, applied the AptamerRunner
382 clustering algorithm to identify 38 candidate aptamers from unique clusters related by sequence or
383 by structure. These 38 candidates were then ranked for their potential to bind the CD22 protein
384 through a molecular docking and molecular dynamic approach¹⁶. The aptamer that exhibited
385 specificity for CD22 (B-ALL1 aptamer) was identified within a unique group of aptamers that
386 were related by structure. Within this cluster of structurally related aptamers, the B-ALL1 aptamer
387 exhibited the greatest positive fold enrichment and was therefore identified as a potential

388 candidate. An edit distance clustering analysis did not identify B-ALL1 as a candidate. Also,
389 because B-ALL1 was the 573rd most abundant aptamer, it would likely have been missed using
390 traditional candidate selection approaches. This highlights the importance of clustering by related
391 structures, which is not an available capability in FASTAptamer, FASTApameR 2.0 or
392 AptaCluster.

393 A second study, by Santana-Viera *et. al.*¹⁵, applied AptamerRunner to examine the relatedness of
394 aptamer libraries enriched in two independent SELEX processes: a protein-based SELEX using
395 recombinant human EphA2 as target and a Cell-Internalization SELEX using EphA2-expressing
396 MDA231 cells as targets^{15, 35}. The AptamerRunner clustering algorithm identified a group of
397 aptamers within these two different SELEX processes that shared structure and sequence
398 relatedness. From this group of aptamers, the candidate aptamer ATOP was observed to target
399 hEphA2 and exhibited antitumorigenic effects *in vitro* and *in vivo*. The flexibility of
400 AptamerRunner permitted the researchers to develop a novel clustering strategy that made use of
401 both edit distance and tree distance clustering using the all-vs-all approach. A seed approach with
402 two different SELEX processes would have been challenging due to the complication in defining
403 which SELEX would provide the aptamer sequences to serve as seeds. Importantly, without the
404 in-depth analysis of aptamer relatedness enabled by AptamerRunner, the aptamers described by
405 these two studies would have been prohibitively challenging to identify with the other available
406 clustering tools.

407 Our analysis of the IL-10 aptamer database, which was previously used to demonstrate the utility
408 of AptaCluster, employed AptamerRunner to focus on clustering based on both edit distance and
409 tree distance. Initial observations indicated that aptamers with an edit distance of 5 are structurally
410 related, corresponding to a tree distance of 10. The application of a tree distance of 10 revealed
411 overlapping clusters and underscored the necessity of considering both measures for optimal
412 clustering outcomes. This comprehensive clustering approach, which combines an edit distance of
413 5 with a tree distance of 10, provides detailed granularity of distinct groups of aptamers,
414 uncovering clusters that had not been previously tested for IL-10 binding. The analysis of the IL-
415 10 aptamer database illustrates how the AND function of AptamerRunner offers a more nuanced
416 understanding of aptamer networks.

417 The aforementioned examples of AptamerRunner identifying candidate aptamers made use of the
418 all-vs-all clustering approach rather than the seed approach. The seed approach, introduced by
419 AptaCluster and FASTAptamer, is founded on the idea that certain aptamer sequences, called the
420 seed, serve as the basis from which mutations during SELEX accumulate to yield more specific or
421 higher affinity aptamers. However, the seed approach for clustering aptamers is a greedy process
422 whereby sequences connected to the seed are removed from the pool of aptamers available for
423 clustering. Therefore, the seed approach will miss inter-aptamer connections identified by the all-
424 vs-all approach. However, the all-vs-all approach is a significantly more computationally intensive
425 process than the seed approach and can yield hairball networks that are more complex to interpret.
426 The concept of the seed sequence is potentially more important for aptamer libraries with longer
427 variable regions. Longer variable region libraries (e.g. >30 nucleotides) have a large starting
428 complexity that cannot be sampled at the start of SELEX and PCR-generated mutations are more
429 likely to introduce beneficial aptamer sequences not present during the initial selection rounds.
430 FASTAptameR 2.0 includes a function (distance module) that can specifically evaluate edit
431 distance from a seed sequence, or other sequences, to evaluate accumulation of mutations during
432 SELEX. The FASTAptameR 2.0 distance module plots the edit distance distribution from the seed
433 sequence as a histogram. Given that the accumulation of mutations is more likely to occur with
434 more abundant aptamers and less likely with aptamers of lower abundance, AptamerRunner could
435 interrogate larger more complex networks of related aptamers identified by the all-vs-all approach
436 by re-evaluating them using the seed approach. While AptamerRunner did not aim to settle this
437 debate, it does provide experiential computing that allows aptamer researchers to investigate their
438 data independently based on their goals and expertise as to how the clustering data should be
439 visualized.

440 Future versions of AptamerRunner could include additional aptamer bioinformatics tools to
441 process, compile, and analyze raw aptamer NGS data with a GUI interface like the integrated
442 pipelines offered by FASTAptameR 2.0 and AptaCluster. The structure prediction algorithm could
443 incorporate additional structure prediction algorithms such as Mfold³⁶ or be modified to permit
444 multiple structures for each aptamer sequence including suboptimal structures. Furthermore,
445 AptamerRunner does not include any option to limit what aptamers are clustered. AptaCluster
446 applies a hashing function that filters the dataset by identifying pairs of aptamers that are likely to
447 be dissimilar and FASTAptameR 2.0 includes a filter function to cluster only aptamer sequences

448 of a minimum abundance or produce a set number of clusters. With AptamerRunner we filtered
449 the aptamer NGS dataset prior to clustering using a separate aptamer abundance and persistence
450 analysis^{16, 37}. Ideally, algorithms that can filter aptamer NGS datasets, such as the AptamerRunner
451 hashing function, could be applied independently prior to clustering to investigate the effectiveness
452 of different filtering strategies. Additional future directions include determining the range between
453 separate groups of networked aptamers with the idea that more distantly related aptamers are more
454 likely to target different epitopes. Analysis options could include ranking different groups of
455 clusters aptamers and ranking individual aptamers within each cluster by integrating scoring or
456 application of molecular docking to predict which groups of aptamer bind to same regions of a
457 target protein^{18, 19, 33, 38-40}.

458 In summary, AptamerRunner seeks to facilitate human-computer synergy²⁷ for clustering aptamer
459 NGS data with innovative approaches, enabling diverse sequence and structure relatedness,
460 introducing logical operators, and offering seamless integration with Cytoscape for visualization
461 and interpretation. The inherent flexibility, user-friendly design and prospects for future
462 enhancements collectively position AptamerRunner at the forefront of advancing aptamer
463 research.

464

465 **Materials and Methods**

466 **RNA or DNA aptamer secondary structure prediction algorithm**

467 The structure prediction component of AptamerRunner requires either full-length aptamers or only
468 the variable region of the aptamers in a FASTA-formatted file. The FASTA file should contain
469 collapsed aptamer NGS data, in which all unique aptamer sequences are represented once and are
470 ranked in descending order based on number of duplicate reads (**Supplemental Figure S5A**). If
471 the FASTA file contains only the variable region sequences, AptamerRunner provides an option
472 to automatically append the constant regions per user input. Secondary structures are predicted
473 using the RNAfold structure prediction algorithm from the Vienna Package v2.0 ⁴¹⁻⁴³, with the
474 lowest minimum free energy structure being retained. Pass-through commands specific for
475 RNAfold can be included during AptamerRunner execution. The predicted aptamer secondary
476 structures are appended to a modified FASTA-formatted file (FASTA.struct) using dot-bracket
477 annotation (**Supplemental Figure S5B**). Properties associated with the predicted structures (e.g.,
478 minimum free energy) that are output by the RNAfold are compiled into a separate tab-delimited
479 file with the aptamer FASTA header information as a key. These structural properties of each
480 aptamer sequence can be imported, if needed, for overlaying when visualizing the clustering data.

481 **Aptamer clustering algorithm**

482 The aptamer clustering component of AptamerRunner generates networks of related aptamers
483 from aptamer sequences and predicted secondary structures. Networks of related aptamers are
484 constructed using either the all-vs-all approach or the seed approach. User options include
485 selecting the 1) *Edge Type* and the 2) *Maximum Distance Measure*.

486 *Edge Type*: The Edge Type defines the relatedness distance metric applied by the clustering
487 algorithm in order to decide if two aptamers sequences should be connected when building
488 networks of aptamers. Sequence similarity is determined using Levenshtein edit distance ⁴⁴ and
489 structure similarity is determined using tree distance ⁴⁵. These distance measurements are
490 calculated by the clustering algorithm using RNAdistance from the Vienna Package v2.0 ^{41, 46}.
491 Options for clustering include 1) *edit* to use only edit distance data, 2) *tree* for tree distance data
492 only, 3) *both* to apply the logical operator AND with edit distance and tree distance data, and 4)
493 *or* to apply the logical operator OR with edit and tree distance data. Regardless of Edge Type
494 applied, AptamerRunner will calculate both the edit and tree distances between aptamer sequences
495 and include these data as metadata for the edges connecting two aptamer sequences

496 (Supplemental Figure S1B). This enables users to perform additional analyses using the edit and
497 tree distance values within Cytoscape. For example, if *edit* is designated as the Edge Type option,
498 AptamerRunner will only apply edit distance data when constructing the networks of related
499 aptamers, but it will also calculate tree distance values and include these data as edge metadata.

500 *Maximum Distance Measure*: The Maximum Distance Measure determines the threshold value of
501 a maximum edit and/or tree distance value for a given Edge Type from which networks of related
502 aptamers will be constructed. A separate Maximum Distance Measure can be set for edit distance
503 and for tree distance. The clustering algorithm requires that the Maximum Distance Measure match
504 the Edge Type and, if using a logical operator, that a Maximum Distance Measure be set for both
505 edit distance and tree distance.

506 *Output Files*: Three files are outputted by the clustering program: 1) an eXtensible Graph Markup
507 and Modeling Language (XGMML) file containing the aptamer clustering data, which can be
508 imported into Cytoscape; 2) the input FASTA.struct file; and 3) a log file of the commands used
509 to initiate AptamerRunner.

510 Cytoscape visualization of AptamerRunner clustering results

511 The XGMML file of clustering results from AptamerRunner can be visualized through Cytoscape
512³⁰, an open-source network analysis program (see supplemental methods for specific details).
513 Networks of clustered aptamer sequences are visualized using *nodes*, which represent unique
514 aptamer sequences, and *edges* that connect related aptamer sequence nodes (Supplemental Figure
515 S6). Nodes and edges include metadata associated to each aptamer sequence (e.g., name, sequence,
516 structure) and between interconnected aptamers (e.g., edit distance and tree distance values).
517 Additional metadata, such as normalized read counts, fold enrichment or data from scoring
518 algorithms, can be imported into Cytoscape's node data tables from comma or tab delimited text
519 files. These metadata can be used to dictate the visual properties of nodes and edges within
520 Cytoscape to facilitate interpretation of the clustering results. Cytoscape's interactive interface
521 permits easy selection of nodes to isolate potential candidate aptamer sequences. Nodes selected
522 from networks are compiled into Cytoscape's data table along with corresponding node data
523 (Supplemental Figure S1A). The Cytoscape data table may be copied directly or exported as a
524 text file. In addition to selecting individual aptamer sequences, groups of aptamers may be readily
525 identified and examined in more detail using network analysis tools such as the clusterMaker
526 algorithm⁴⁷. The clusterMaker algorithm assigns an identification to each group of related

527 aptamers and appends these identifications to the Cytoscape data table (**Supplemental Figure**
528 **S1B**).

529 **NGS datasets:** Cibiel *et. al.* 2014 PLoSOne dataset and Hoinka *et. al.* 2015 NAR were imported
530 into Galaxy from the NCBI SRA database (ERR2121976-91 for Cibiel selection round 0 – 15;
531 SRR3279660 for Hoinka selection rounds 1 – 4 and SRR3279661 for Hoinka selection round 5).
532 For each dataset, Galaxy was used to generate a NrD, conduct a persistence and abundance analysis
533 and filter the NrD. The compiled NrDs for all datasets are available through supplemental. The
534 Cibiel *et. al.* 2014 PLoSOne NrD was filtered for aptamer sequences observed in at least four of
535 the sixteen sequenced selection rounds and with 14 reads, resulting in 6,639 unique aptamer
536 sequences. The published aptamer sequence ACE23 did not meet with criteria being observed in
537 only three selection rounds with maximum read count of one. The ACE23 aptamer sequence was
538 added back into the filtered NrD. The Hoinka *et. al.* 2015 NAR NrD was filtered for aptamer
539 sequences observed in at least one of the four sequence selection rounds with at least 50 reads,
540 resulting in 2,140 unique aptamer sequences.

541 **Keywords**

542 Aptamer; Bioinformatics; AptamerRunner; Sequence analysis; Structure analysis; Cytoscape.

543 **Data availability**

544 The AptamerRunner .NET script Windows and Linux versions are available in supplemental and
545 on GitHub (<https://github.com/ui-icts/aptamer-runner/releases/tag/v0.0.3>).

546 NGS source data¹⁶ used as an example are provided with this paper in supplemental. Further data
547 supporting the findings of this study are available from the corresponding author upon reasonable
548 request.

549 **Acknowledgments**

550 Dario Ruiz-Ciancio: Investigation, Formal Analysis, Visualization, Writing – original draft.
551 Suresh Veeramani: Investigation, Formal Analysis, Writing – review & editing. Rahul Singh: Formal
552 Analysis, Writing – original draft. Eric Embree: Software, Resources, Data curation. Chris
553 Ortman: Software, Resources, Supervision. Kristina W. Thiel: Formal Analysis, Writing – original
554 draft. William H. Thiel: Conceptualization, Methodology, Formal Analysis, Project
555 administration, Writing – original draft.

556 This work was supported by an American Heart Association Scientist Development Grant
557 (14SDG18850071 to WHT), the National Institutes of Health (R01HL139581 and R01HL157956
558 to WHT; K22CA263783 to KWT), the National Science Foundation (IIS-1817239 to RS), the
559 Department of Defense (DOD CDMRP-PRCRP CA220729 to KWT), the Bunge and Born Fund
560 (FBB-20170609 to DRC) and Fulbright-Argentinian Ministry of Education (ME-FLB-2022-2023
561 to DRC).

562

563 **Conflict of Interest**

564 Authors declare no conflict of interest.

565

566 **References**

567 1. Zhu, C, Feng, Z, Qin, H, Chen, L, Yan, M, Li, L, and Qu, F (2023). Recent progress of
568 SELEX methods for screening nucleic acid aptamers. *Talanta* **266**: 124998.

569 2. Chauhan, N, Saxena, K, and Jain, U (2022). Single molecule detection; from microscopy
570 to sensors. *Int J Biol Macromol* **209**: 1389-1401.

571 3. Xie, S, Sun, W, Fu, T, Liu, X, Chen, P, Qiu, L, Qu, F, and Tan, W (2023). Aptamer-
572 Based Targeted Delivery of Functional Nucleic Acids. *J Am Chem Soc* **145**: 7677-7691.

573 4. Fan, D, Wang, J, Wang, E, and Dong, S (2020). Propelling DNA Computing with
574 Materials' Power: Recent Advancements in Innovative DNA Logic Computing Systems
575 and Smart Bio-Applications. *Adv Sci (Weinh)* **7**: 2001766.

576 5. Li, L, Xu, S, Yan, H, Li, X, Yazd, HS, Li, X, Huang, T, Cui, C, Jiang, J, and Tan, W
577 (2021). Nucleic Acid Aptamers for Molecular Diagnostics and Therapeutics: Advances
578 and Perspectives. *Angew Chem Int Ed Engl* **60**: 2221-2231.

579 6. Esposito, CL, Catuogno, S, Condorelli, G, Ungaro, P, and de Franciscis, V (2018).
580 Aptamer Chimeras for Therapeutic Delivery: The Challenging Perspectives. *Genes*
581 (*Basel*) **9**.

582 7. Mullard, A (2023). FDA approves second RNA aptamer. *Nat Rev Drug Discov* **22**: 774.

583 8. Tuerk, C, and Gold, L (1990). Systematic evolution of ligands by exponential
584 enrichment: RNA ligands to bacteriophage T4 DNA polymerase. *Science* **249**: 505-510.

585 9. Ellington, AD, and Szostak, JW (1990). In vitro selection of RNA molecules that bind
586 specific ligands. *Nature* **346**: 818-822.

587 10. DeRosa, MC, Lin, A, Mallikaratchy, P, McConnell, EM, McKeague, M, Patel, R, and
588 Shigdar, S (2023). In vitro selection of aptamers and their applications. *Nature Reviews
589 Methods Primers* **3**: 55.

590 11. Metzker, ML (2010). Sequencing technologies - the next generation. *Nat Rev Genet* **11**:
591 31-46.

592 12. Quang, N, Perret, G, and Duconge, F (2016). Applications of High-Throughput
593 Sequencing for In Vitro Selection and Characterization of Aptamers. *Pharmaceuticals*
594 (*Basel*) **9**.

595 13. Sun, D, Sun, M, Zhang, J, Lin, X, Zhang, Y, Lin, F, Zhang, P, Yang, C, and Song, J
596 (2022). Computational tools for aptamer identification and optimization. *TrAC Trends in
597 Analytical Chemistry* **157**: 116767.

598 14. Komarova, N, Barkova, D, and Kuznetsov, A (2020). Implementation of High-
599 Throughput Sequencing (HTS) in Aptamer Selection Technology. *Int J Mol Sci* **21**.

600 15. Santana-Viera, L, Dassie, JP, Rosas-Lapena, M, Garcia-Monclus, S, Chicon-Bosch, M,
601 Perez-Capo, M, Pozo, LD, Sanchez-Serra, S, Almacellas-Rabaiget, O, Maqueda-Marcos,
602 S, *et al.* (2023). Combination of protein and cell internalization SELEX identifies a
603 potential RNA therapeutic and delivery platform to treat EphA2-expressing tumors. *Mol
604 Ther Nucleic Acids* **32**: 758-772.

605 16. Ruiz-Ciancio, D, Lin, L-H, Veeramani, S, Barros, MN, Sanchez, D, Di Bartolo, AL,
606 Masone, D, Giangrande, PH, Mestre, MB, and Thiel, WH (2023). Selection of novel cell-
607 internalizing RNA aptamer specific for CD22 antigen in B- Acute Lymphoblastic
608 Leukemia. *Molecular Therapy - Nucleic Acids*.

609 17. Thiel, WH, Bair, T, Peek, AS, Liu, X, Dassie, J, Stockdale, KR, Behlke, MA, Miller, FJ,
610 Jr., and Giangrande, PH (2012). Rapid identification of cell-specific, internalizing RNA

611 aptamers with bioinformatics analyses of a cell-based aptamer selection. *PLoS One* **7**:
612 e43836.

613 18. Ishida, R, Adachi, T, Yokota, A, Yoshihara, H, Aoki, K, Nakamura, Y, and Hamada, M
614 (2020). RaptRanker: in silico RNA aptamer selection from HT-SELEX experiment based
615 on local sequence and structure information. *Nucleic Acids Research* **48**: e82-e82.

616 19. Jiang, P, Meyer, S, Hou, Z, Propson, NE, Soh, HT, Thomson, JA, and Stewart, R (2014).
617 MPBind: a Meta-motif-based statistical framework and pipeline to Predict Binding
618 potential of SELEX-derived aptamers. *Bioinformatics* **30**: 2665-2667.

619 20. Thompson, JD, Higgins, DG, and Gibson, TJ (1994). CLUSTAL W: improving the
620 sensitivity of progressive multiple sequence alignment through sequence weighting,
621 position-specific gap penalties and weight matrix choice. *Nucleic Acids Res* **22**: 4673-
622 4680.

623 21. Edgar, RC (2004). MUSCLE: multiple sequence alignment with high accuracy and high
624 throughput. *Nucleic Acids Res* **32**: 1792-1797.

625 22. Bayrac, AT, Sefah, K, Parekh, P, Bayrac, C, Gulbakan, B, Oktem, HA, and Tan, W
626 (2011). In vitro Selection of DNA Aptamers to Glioblastoma Multiforme. *ACS Chem
627 Neurosci* **2**: 175-181.

628 23. Hoinka, J, Berezhnoy, A, Sauna, ZE, Gilboa, E, and Przytycka, TM (2014). AptaCluster -
629 A Method to Cluster HT-SELEX Aptamer Pools and Lessons from its Application. *Res
630 Comput Mol Biol* **8394**: 115-128.

631 24. Hoinka, J, Berezhnoy, A, Dao, P, Sauna, ZE, Gilboa, E, and Przytycka, TM (2015).
632 Large scale analysis of the mutational landscape in HT-SELEX improves aptamer
633 discovery. *Nucleic Acids Res*.

634 25. Alam, KK, Chang, JL, and Burke, DH (2015). FASTAptamer: A Bioinformatic Toolkit
635 for High-throughput Sequence Analysis of Combinatorial Selections. *Mol Ther Nucleic
636 Acids* **4**: e230.

637 26. Kramer, ST, Gruenke, PR, Alam, KK, Xu, D, and Burke, DH (2022). FASTAptameR
638 2.0: A web tool for combinatorial sequence selections. *Mol Ther Nucleic Acids* **29**: 862-
639 870.

640 27. Singh, R, Yang, H, Dalziel, B, Asarnow, D, Murad, W, Foote, D, Gormley, M, Stillman,
641 J, and Fisher, S (2013). Towards human-computer synergetic analysis of large-scale
642 biological data. *BMC Bioinformatics* **14 Suppl 14**: S10.

643 28. Singh, R, and Jain, R (2006). From Information-Centric to Experiential Environments. In:
644 Goldin, D, Smolka, SA and Wegner, P (eds). *Interactive Computation: The New
645 Paradigm*. Springer Berlin Heidelberg: Berlin, Heidelberg. pp 323-351.

646 29. Boettiger, C (2015). An introduction to Docker for reproducible research. *ACM SIGOPS
647 Operating Systems Review* **49**: 71-79.

648 30. Shannon, P, Markiel, A, Ozier, O, Baliga, NS, Wang, JT, Ramage, D, Amin, N,
649 Schwikowski, B, and Ideker, T (2003). Cytoscape: a software environment for integrated
650 models of biomolecular interaction networks. *Genome Res* **13**: 2498-2504.

651 31. Hoinka, J, Berezhnoy, A, Sauna, ZE, Gilboa, E, and Przytycka, TM (2014). AptaCluster
652 – A Method to Cluster HT-SELEX Aptamer Pools and Lessons from Its Application. In:
653 Sharan, R (ed). *Research in Computational Molecular Biology*, vol. 8394. Springer
654 International Publishing: Pittsburgh PA. pp 115-128.

655 32. Hoinka, J, Backofen, R, and Przytycka, TM (2018). Aptasuite: A Full-Featured
656 Bioinformatics Framework for the Comprehensive Analysis of Aptamers from HT-
657 SELEX Experiments. *Mol Ther Nucleic Acids* **11**: 515-517.

658 33. Dao, P, Hoinka, J, Takahashi, M, Zhou, J, Ho, M, Wang, Y, Costa, F, Rossi, JJ,
659 Backofen, R, Burnett, J, *et al.* (2016). Aptatrace Elucidates RNA Sequence-Structure
660 Motifs from Selection Trends in HT-SELEX Experiments. *Cell Syst* **3**: 62-70.

661 34. Levay, A, Brenneman, R, Hoinka, J, Sant, D, Cardone, M, Trinchieri, G, Przytycka, TM,
662 and Berezhnoy, A (2015). Identifying high-affinity aptamer ligands with defined cross-
663 reactivity using high-throughput guided systematic evolution of ligands by exponential
664 enrichment. *Nucleic Acids Res* **43**: e82.

665 35. Ducrot, C, and Piffoux, M (2023). Combining independent protein and cellular SELEX
666 with bioinformatic analysis may allow high affinity aptamer hit discovery. *Molecular
667 Therapy - Nucleic Acids* **33**: 254-256.

668 36. Zuker, M (2003). Mfold web server for nucleic acid folding and hybridization prediction.
669 *Nucleic Acids Res* **31**: 3406-3415.

670 37. Thiel, WH (2016). Galaxy Workflows for Web-based Bioinformatics Analysis of
671 Aptamer High-throughput Sequencing Data. *Mol Ther Nucleic Acids* **5**: e345.

672 38. Caroli, J, Forcato, M, and Bicciato, S (2020). APTANI2: update of aptamer selection
673 through sequence-structure analysis. *Bioinformatics* **36**: 2266-2268.

674 39. Hoinka, J, and Przytycka, T (2016). Aptaplex - A dedicated, multithreaded
675 demultiplexer for HT-SELEX data. *Methods* **106**: 82-85.

676 40. Shieh, KR, Kratschmer, C, Maier, KE, Greally, JM, Levy, M, and Golden, A (2020).
677 Aptcompare: optimized de novo motif discovery of RNA aptamers via HTS-SELEX.
678 *Bioinformatics* **36**: 2905-2906.

679 41. Lorenz, R, Bernhart, SH, Honer Zu Siederdissen, C, Tafer, H, Flamm, C, Stadler, PF, and
680 Hofacker, IL (2011). ViennaRNA Package 2.0. *Algorithms Mol Biol* **6**: 26.

681 42. Mathews, DH, Sabina, J, Zuker, M, and Turner, DH (1999). Expanded sequence
682 dependence of thermodynamic parameters improves prediction of RNA secondary
683 structure. *J Mol Biol* **288**: 911-940.

684 43. Walter, AE, Turner, DH, Kim, J, Lyttle, MH, Müller, P, Mathews, DH, and Zuker, M
685 (1994). Coaxial stacking of helices enhances binding of oligoribonucleotides and
686 improves predictions of RNA folding. *Proc Natl Acad USA* **91**: 9218-9222.

687 44. VI, L (1966). Binary codes capable of correcting deletions, insertions and reversals. *Sov
688 Phys Dokl* **10**: 707-710.

689 45. Fontana, W, Konings, DA, Stadler, PF, and Schuster, P (1993). Statistics of RNA
690 secondary structures. *Biopolymers* **33**: 1389-1404.

691 46. Hofacker, IL, Fontana, W, Stadler, PF, Bonhoeffer, LS, Tacker, M, and Schuster, P
692 (1994). Fast folding and comparison of RNA secondary structures. *Monatshefte fr
693 Chemie Chemical Monthly* **125**: 167-188.

694 47. Morris, JH, Apeltsin, L, Newman, AM, Baumbach, J, Wittkop, T, Su, G, Bader, GD, and
695 Ferrin, TE (2011). clusterMaker: a multi-algorithm clustering plugin for Cytoscape. *BMC
696 Bioinformatics* **12**: 436.

698

List of Figures Captions

699 **Figure 1. AptamerRunner, a structure prediction and clustering algorithm a to visualize**
700 **selected aptamers.** AptamerRunner consist of two independent algorithms, a secondary structure
701 prediction algorithm and an aptamer clustering algorithm. **(A)** AptamerRunner is a .NET bash
702 script coded in C#. The AptamerRunner .NET script communicates with the Docker repository to
703 download the latest version of the AptamerRunner Docker image and then initiates the
704 AptamerRunner Docker image within the Docker platform as a Docker container. The Docker
705 container includes the AptamerRunner structure prediction algorithm and the clustering algorithm
706 with all dependencies. Once either algorithm has finished processing any user commands and
707 output results, the AptamerRunner .NET script shuts down the AptamerRunner Docker container.
708 **(B)** The secondary structure prediction algorithm utilizes collapsed aptamer NGS data in FASTA
709 format to predict the secondary structure of a full-length aptamer using RNAfold. The secondary
710 structure prediction algorithm has the option to append the constant region sequence if needed.
711 Output includes a modified FASTA-formatted file with a third line containing the dot-bracket
712 annotation of each predicted structure and a properties file that includes information about the
713 predicted structures (e.g., minimum free energy). **(C)** The AptamerRunner clustering algorithm
714 uses the FASTA-formatted file with the predicted structures to generate networks of related
715 aptamers using options selected by the user for the *Clustering Method*, the *Edge Type*, and the
716 *Display Threshold*. Output files include the clustering results compiled into a XGMML file, which
717 is visualized using Cytoscape, a log file, and the input file.

718

719 **Figure 2: AptamerRunner all-vs-all and seed clustering approaches using either edit distance**
720 **or tree distance.**

721 AptamerRunner clustering aptamer NGS data using different clustering methods with either edit
722 distance 1 or tree distance 0. **(A)** Seed approach with edit distance 1 or **(B)** tree distance 0. **(C)** All-
723 vs-all approach with edit distance 1 or **(D)** tree distance 0. Ruiz-Ciancio *et. al.*¹⁶ data was used as
724 example data for clustering and clustering results were visualized using Cytoscape (v 2.8.1).

725

726 **Figure 3: Use of logical operators (AND, OR) with AptamerRunner clustering.**

727 Clustering using logical operator AND with edit distance 1 and tree distance 0 with the **(A)** seed
728 approach and the **(B)** all-vs-all approach. Clustering using the logical operator OR with edit

729 distance 1 and tree distance 0 with the **(C)** seed approach and the **(D)** all-vs-all approach. Ruiz-
730 Ciancio *et. al.*¹⁶ data was used as example data for aptamer classification and selection, and data
731 are visualized using Cytoscape (v 2.8.1).

732

733 **Figure 4: Clustering results from FASTAptamer, FASTAptameR 2.0 and AptaCluster.**

734 **(A)** The modified FASTA file clustering results from FASTAptamer and FASTAptameR 2.0. The
735 FASTA header information follows as rank, reads, reads per million (RPM), cluster id, rank within
736 the cluster, and edit distance from the seed sequence. **(B)** FASTAptameR 2.0 clustering results
737 outputted as a data table. **(C)** AptaCluser exported clustering results from round 9 and **(D)** exported
738 cluster table. Data shown are only a portion of larger files. Gaps in data are denoted by a “...”.

739

740 **Figure 5: Visual integration of round-to-round log2 fold enrichment data mapped onto**
741 **clustered aptamers.**

742 The log2 fold enrichment data between selection rounds from a SELEX process imported into the
743 Cytoscape data table can be applied to determine the visual properties of aptamers clustered by
744 AptamerRunner. The node size was set to the log10 normalized read count of round 9 and node
745 color was determined by the log2 fold enrichment between different selection rounds based on **(A)**
746 the phases of the sequence enrichment % sigmoidal curve fit; exponential (R0:R2), linear (R2:R5)
747 and asymptotic (R5:R9). **(B)** The log2 fold change enrichment of round 0 to round 2 shows early
748 linear phase enrichment of aptamer sequences, **(C)** round 2 to round 5 shows changes in aptamer
749 sequences during the linear phase of sequence enrichment when the aptamer library experienced
750 the greatest changes in convergence and **(D)** round 5 to round 9 show changes in aptamer
751 sequences during the asymptotic phase when the aptamer library had reached maximum
752 convergence.

753

754 **Figure 6: Aptamer scoring algorithm mapped to clustered aptamers.**

755 Data from aptamer scoring algorithms **(A)** MPBind and **(B)** RaptRanker were mapped to the
756 networks generated by AptamerRunner. **(C)** Cytoscape can produce scatter plots that can
757 compare the scoring data from MPBind and RaptRanker to identify aptamers and clusters of
758 aptamers that were scored highly by both algorithms.

759

760 **Figure 7: Re-clustering aptamer networks within Cytoscape**

761 Groups of aptamers can be re-clustered using tools within Cytoscape. (A) Nodes within a large
762 hairball network of aptamers clustered by edit distance 1 using the all-vs-all approach were
763 selected and (B) re-organized within Cytoscape by identical structures.

764

765 **Figure 8: AptamerRunner analysis of IL-10 aptamers**

766 (A) AptamerRunner edit distance 5 using seed approach and (B) all-vs-all approach. (C) Tree
767 distance 10 using the all-vs-all approach. (D) Edit distance 5 and tree distance 5. Clusters are
768 labeled as reported in Levay *et. al.*³⁴ * indicates cluster found by AptamerRunner using edit
769 distance 5 with the seed approach that was not identified by AptaCluster.

770 Table 1

| | AptamerRunner | FASTAptamer | FASTAptameR 2.0 | AptaCluster |
|---------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------|
| Interface | Command line C# .NET script operating a Docker container | Command line Perl-based programs | Offline web browser accessing a Docker container | AptaSuite Java application |
| Components (*clustering) | <ul style="list-style-type: none"> Secondary structure prediction Clustering algorithm* | <ul style="list-style-type: none"> Count Compare Cluster* Enrich Search | <ul style="list-style-type: none"> FASTAptamer components Cluster module* Edit distance module Motif discovery module Position enrich module | <ul style="list-style-type: none"> AptaGUI AptaPlex AptaSim AptaCluster* AptaTrace AptaMut |
| Third party software | Docker, Cytoscape | none | Docker, web browser | Java |
| Clustering approach | seed, all-vs-all | seed | seed | seed |
| Distance measures | <ul style="list-style-type: none"> Edit distance Tree distance Logical operators (AND, OR) | <ul style="list-style-type: none"> Edit distance | <ul style="list-style-type: none"> Edit distance | <ul style="list-style-type: none"> Edit distance |
| Visualization | Cytoscape display of networked aptamers (nodes) connected by distance measures (edges) | N/A | <ul style="list-style-type: none"> Cluster metaplots k-mer PCA | N/A |
| Metadata | <p>Nodes:</p> <ul style="list-style-type: none"> Aptamer name Predicted structure User-imported <p>Edges:</p> <ul style="list-style-type: none"> Edit distance values Tree distance values User-imported | N/A | N/A | N/A |

771