1 Differential binding cell-SELEX: method for identification of cell

2 specific aptamers using high throughput sequencing

- 3 Karlis Pleiko¹*, Liga Saulite¹, Vadims Parfejevs¹, Karlis Miculis², Egils Vjaters² and Una Riekstina¹
- ⁴ ¹ Faculty of Medicine, University of Latvia, Riga, LV-1004, Latvia
- 5 ² Pauls Stradins Clinical University Hospital, Riga, LV-1002, Latvia
- 6 * To whom correspondence should be addressed. Tel: +371 28680049; Email: karlis.pleiko@lu.lv
- 7 Present Address: Karlis Pleiko, Faculty of Medicine, University of Latvia, Riga, LV-1004, Latvia
- 8

9 ABSTRACT

- 10 Aptamers have evolved as a viable alternative to antibodies in recent years. High throughput
- sequencing (HTS) has revolutionized the aptamer research by increasing the number of reads from
- 12 few using Sanger sequencing to millions of reads using HTS approach. Despite the availability and
- 13 advantages of HTS compared to Sanger sequencing there are only 50 aptamer HTS sequencing
- 14 samples available on public databases. HTS data for aptamer research are mostly used to compare
- 15 sequence enrichment between subsequent selection cycles. This approach does not take full
- advantage of HTS because enrichment of sequences during selection can be due to inefficient
- 17 negative selection when using live cells. Here we present differential binding cell-SELEX (systematic
- 18 evolution of ligands by exponential enrichment) workflow that adapts FASTAptamer toolbox and
- 19 bioinformatics tool *edgeR* that is mainly used for functional genomics to achieve more informative
- 20 metrics about the selection process. We propose fast and practical high throughput aptamer
- 21 identification method to be used with cell-SELEX technique to increase successful aptamer selection
- 22 rate against live cells. The feasibility of our approach is demonstrated by performing aptamer
- 23 selection against clear cell renal cell carcinoma (ccRCC) RCC-MF cell line using RC-124 cell line from
- 24 healthy kidney tissue for negative selection.

25 INTRODUCTION

- 26 Aptamers are short (20 100 nt) oligonucleotides that, contrary to the most of other functional nucleic
- 27 acids, bind specific molecular targets owing to their folded three-dimensional (3D) structures (1). Most
- of the aptamers are developed for therapeutic or diagnostic purposes (2, 3). Currently several
- 29 aptamer candidates are being tested in clinical trials for treatment of age-related macular
- 30 degeneration (4), Duchenne muscular dystrophy (5), chronic lymphocytic leukemia (6) and others (7).
- 31 After initial description of aptamer selection method termed SELEX (systematic evolution of ligands by
- 32 exponential enrichment) (8), several aptamer selection methods have been developed, among others
- 33 cell-SELEX (9), where live cells are used. First high throughput SELEX (HT-SELEX) experiment, a
- 34 variation of SELEX process that uses high throughput sequencing (HTS) methods instead of Sanger
- 35 sequencing, was described by Zhao et.al in 2009 (10). Consequently adaptation of high throughput

sequencing (HTS) methods for aptamer research further improved outcomes of selection procedures(11, 12).

Further on, RNA aptamer selection against active and inactive conformation of β_2 adrenoreceptor

39 described by Kahsai et.al. employs HTS methods to characterize the fold change enrichment of

40 particular sequences during the selection against each individual target in parallel (11). However, in

41 case of cell-SELEX this approach might be of very limited use due to the high diversity of protein

42 targets on cell surface that would cause enrichment of non-specifically bound sequences if no

43 negative selection were performed.

44 Several research teams have developed tools to analyse HTS data from aptamer selection, notably

45 FASTAptamer, a toolkit developed by Alam et.al. that can be used to track the evolutionary trajectory

46 during the SELEX process of individual oligonucleotide sequences (12). Just recently AptaSUITE, a

47 comprehensive bioinformatics framework that includes most of the previously published functionalities

48 of different tools – data pre-processing, sequence clustering, motif identification and mutation analysis

49 has been introduced (13).

50 RNA-sequencing (RNA-seq) experiments are used to quantify the differential expression of gene

51 transcripts between samples (14). We speculated that it might be possible to adapt data analysis tools

52 currently used for RNA-seq to be used with HT cell-SELEX experiments. During the cell-SELEX

53 experiment, the goal to be achieved is to select aptamers that bind to the target cells in larger number

54 compared to the control cells, making experimental design similar to RNA-seq analysis. Here we

55 provide a differential binding cell-SELEX method that can be used to identify differentially abundant

aptamers on the surface of target cells and negative control cells during the cell-SELEX experiments

and to calculate the statistical significance of these differences. Analysis includes the use of *edgeR*

58 (15), a common tool for the analysis of RNA-seq experiments that uses negative binominal

59 distribution to identify differentially expressed genes, FASTAaptamer (12) toolbox to estimate the read

60 count, *cutadapt* (16) to remove the constant primer binding regions of aptamers and bespoke R script

available for reuse. Moreover, we combine our approach with sequence enrichment analysis already

62 used by other groups for aptamer selection to identify the most relevant sequences.

63 MATERIAL AND METHODS

64 Cell culturing and buffer solutions

65 Kidney epithelial cell line RC-124 (Cell Lines Service GmbH) established from non-tumor tissue of

66 kidney and carbonic anhydrase 9 (CA9) positive ccRCC cell line RCC-MF (Cell Lines Service GmbH),

67 established from the renal clear cell carcinoma pT2, N1, Mx/ GII-III (lung-metastasis) were used for

68 cell-SELEX process as negative control and target cells accordingly. RCC-MF cells were cultured in

69 RPMI 1640 (Gibco), RC-124 cells were cultured in McCoy's 5A medium (SigmaAldrich). Both culture

70 media were supplemented with 10% fetal bovine serum (FBS) (Gibco), 50 U/ml penicillin and 50

71 μg/ml streptomycin (Gibco). Cells were propagated at 37°C, 5% CO₂ and 95% relative humidity.

- 72 Washing buffer was prepared by adding 0.225 g D-glucose and 0.25 ml of 1 M MgCl₂ to 50 ml of
- 73 phosphate buffered saline (PBS) (SigmaAldrich) and filtering it through a 0.22 µM syringe filter
- 74 (Corning). Binding buffer was prepared by adding 50 mg of bovine serum albumin (SigmaAldrich) and
- 5 mg baker's yeast tRNA (SigmaAldrich) to 50 ml of washing buffer and filtering it through a 0.22 μ M
- 76 syringe filter.

77 Oligonucleotide library

- 78 Randomised oligonucleotide library with 40 nt and 18 nt constant primer binding regions on both sides
- 79 of randomized regions (5'-ATCCAGAGTGACGCAGCA-N40-TGGACACGGTGGCTTAGT-3') was
- 80 adapted from Sefah et al. (17). FAM label was attached on one primer (5'-FAM-
- 81 ATCCAGAGTGACGCAGCA-3') for flow cytometry monitoring and biotin was attached at the end of
- 82 second primer for preparation of ssDNA after each cell-SELEX cycle (5'-biotin-
- 83 ACTAAGCCACCGTGTCCA-3'). Oligonucleotides were ordered from Metabion or Invitrogen.

84 Cell-SELEX procedure

- 85 Cell-SELEX protocol was adapted from Sefah et.al (17). Aptamer library was prepared in binding
- 86 buffer at 14 μ M concentration for the first selection cycle, heated at 95 °C for 5 min and folded on ice
- 87 for at least 15 min, added to fully confluent RCC-MF cells in 100 mm Petri plate (Sarstedt) that were
- 88 washed 2 times with washing buffer before the addition of library. Initial library was applied to RCC-
- 89 MF cells and incubated for 1 hour on ice with RCC-MF cells, but not with RC-124 cell in the first
- 90 selection cycle. After incubation with the oligonucleotide library, cells were washed with 3 ml of
- 91 washing buffer for 3 min and collected with cell scraper after adding 1 ml of DNase free water. DNase
- 92 free water was used for the collection of sequences only for the first cycle, in subsequent cycles
- binding buffer was used to retrieve the bound sequences. After collection retrieved sequences were
- 94 heated at 95 °C for 10 min, centrifuged at 13 000 g and supernatant containing selected aptamer
- 95 sequences was collected.
- 96 In subsequent selection cycles aptamer library was prepared at 500 nM concentration and incubated
- 97 with negative selection cell line RC-124 beforehand. Solution containing unbound sequences was
- 98 collected and applied to RCC-MF cell line after washing cells as described previously. With increasing
- selection cycle number several modifications were made to the selection procedure after 4th
- 100 selection cycle 60 mm plates were used instead of 100 mm plates, increasing concentration of FBS
- 101 (10-20%) were added to library after folding without changing the final concentration of aptamer
- 102 library, wash volume was increased to 5 ml, wash time was increased to 5 min and the number of
- 103 wash times was increased to 3 after incubation.

104 PCR optimization

- 105 After each selection cycle PCR optimization was performed to determine the optimal number of PCR
- 106 cycles. For PCR optimization and preparative PCR cycling conditions were 12 min initial activation at

95 °C, followed by repeated denaturation 30 sec at 95 °C, annealing at 56.3 °C and elongation at
72 °C.

109 ssDNA preparation

- 110 After preparative PCR, ssDNA was acquired using agarose-streptavidin (GE Healthcare) binding to
- biotin labelled strand and FAM labelled ssDNA was eluted with 0.2 M NaOH (SigmaAldrich).
- 112 Desalting was done using NAP5 gravity flow columns (GE Healthcare), concentration was determined
- 113 measuring UV absorbance (NanoQuant Plate, M200 Pro, Tecan), and samples were concentrated
- 114 using vacuum centrifugation (Eppendorf).

115 Monitoring of aptamer binding by flow cytometry

- 116 Enriched aptamer pool, randomized starting library and selected lead aptamers were prepared in
- 117 binding buffer at 1 μM concentrations, heated at 95 °C for 5 min and then put on ice for at least 15
- 118 min. RC-124 and RCC-MF cells were washed with PBS two times and dissociated using Versene
- solution (Gibco). Then 50 µL of enriched aptamer library, starting library, lead aptamers or binding
- 120 buffer were added to 50 μ L of cell suspension (2.5 * 10⁵ cell per sample), followed by addition of 11
- 121 µL of FBS to each sample at final concentration 225 nM. Samples were incubated for 35 min on ice.
- 122 After incubation, samples were washed two times with 500 µL of binding buffer and resuspended in
- 123 500 µL of binding buffer. Samples were passed through 40 µM cell strainer before flow cytometry
- 124 analysis. Flow cytometry data were acquired using a Guava EasyCyte 8HT flow cytometer and
- 125 analysed using the ExpressPro software (Merck Millipore). Flow cytometry data were analysed using
- 126 FlowJo software, version 10 (FlowJo). 10'000 gated events were acquired for each sample.

127 Differential binding

- 128 Aptamer pools after 4th and 11th selection cycle were prepared in binding buffer, heated and folded as
- 129 described for cell-SELEX procedure at 1 ml volume with final concentration 500 nM. 500 μL were
- added to both RC-124 cells and RCC-MF cells grown on 60 mm plate in appropriate cell culture
- 131 media up to 95% confluence. Aptamer pools were added to RC-124 and RCC-MF cells and incubated
- 132 for 30 min on ice, then cells were washed two times and collected using cell scraper, heated
- immediately at 95°C for 10 min, centrifuged for 5 min at 13 000 g. Supernatants containing bound
- 134 sequences from both cell lines were frozen at -20°C. Sequencing was done to compare the differential
- 135 binding profiles of enriched oligonucleotide libraries obtained from both cell lines.

136 Sequencing

- 137 Samples for sequencing were prepared by performing two subsequent overlap PCRs as described in
- 138 16S metagenomic sequencing library preparation protocol (18). 1st overlap PCR used primers (5'-
- 139 TCGTCGGCAGCGTCAGATGTGTATAAGAGACAG-ATCCAGAGTGACGCAGCA-3'
- 140 and 5'-GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAG-ACTAAGCCACCGTGTCCA-3') that
- 141 are complementary to constant regions of randomized oligonucleotide library with added overhang

- 142 that includes Illumina platform specific sequence. Conditions for 1st overlap PCR was 12 min of initial
- 143 activation, followed by 30 sec at 95 °C, 30 sec at 56.3 °C and 3 min at 72 °C. Cycle number was
- 144 optimized for each sample to reduce the non-specific amplification. Afterwards, PCR products from
- 145 one sample were pooled together, concentrated using DNA Clean & Concentrator (Zymo Research),
- 146 run on 3% agarose gel at 110 V for 40 min and the band at 143 bp were cut out and purified using
- 147 Zymoclean Gel DNA Recovery kit (Zymo Research).
- 148 2nd overlap PCR used primers that were partly complementary to previously added overhang and
- 149 contained adapters to attach oligonucleotides to flow cell and i5 and i7 indexes (5'-
- 150 CAAGCAGAAGACGGCATACGAGAT-[i7 index]-GTCTCGTGGGCTCGG-3' and 5'-
- 151 AATGATACGGCGACCACCGAGATCTACAC-[i5 index]-TCGTCGGCAGCGTC-3'). Conditions for 2nd
- 152 overhang PCR were 12 min at 95 °C, followed by 5 cycles of denaturation at 98 °C for 10 sec,
- annealing at 63 °C for 30 sec and elongation at 72 °C for 3 min. After PCR products from one sample
- 154 were pooled together, concentrated using DNA Clean & Concentrator (Zymo Research), run on 3%
- agarose gel at 110 V for 45 min and the band at 212 bp were cut out and purified using Zymoclean
- 156 Gel DNA Recovery kit (Zymo Research). Concentration for final products were determined using
- 157 NEBNext Library Quant Kit for Illumina (New England BioLabs) by qPCR.
- 158 Sequencing was done on Illumina MiSeq platform using MiSeq 150-cycle Reagent Kit v3 in single
- 159 read mode for 150-cycles. 9% of PhiX was added to the run. Sequencing was done at the Estonian
- 160 Genome Center, Tartu, Estonia.
- 161

162 Sequencing data analysis

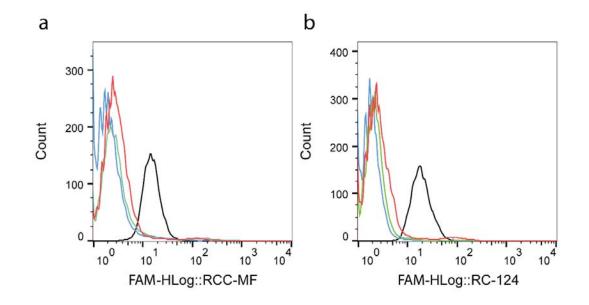
- 163 Sequencing reads were filtered and demultiplexed. Constant primer binding regions were removed,
- sequences that are longer or shorter than 40 nt were discarded using *cutadapt* (16). Counting of
- 165 recurring sequences was done using *fastaptamer-count*, matching of sequences found in replicate
- samples was done using *fastaptamer-enrich* (19).
- 167 Differential expression analysis tool edgeR (15) was further used for the analysis of sequencing data.
- 168 Replicate sequencing samples (n=3) from differential binding cell-SELEX experiment after 4th and 11th
- selection cycles were combined and sequences with low abundance (reads per million < 2 and
- abundant at all in less than 2 sequencing samples) were filtered out. Normalization was performed
- based on reads present in each library. Differential binding was estimated using edgeR function for
- identification of significantly differentially expressed genes using following parameters: log₂ fold
- change (log₂FC) value > 2, p-value < 0.0001, adjusted for multiple comparisons using Benjamini &
 Hochberg (20) method.
- Enrichment analysis was done separately by using all reads that came from 4th pool and 11th pool
 RCC-MF cell binding experiments. We calculated the mean log₂ value of enrichment (mean counts
- per million (CPM) for sequence at 11th cycle divided by mean CPM for the same sequence at 4th cycle)

- 178 for each sequence and kept the sequences that had $log_2FC > 6$ or enrichment between the 4th and 179 11th cycle.
- 180 After these steps, we identified the common sequences in differential binding results and sequence
- 181 enrichment results to identify most likely lead aptamer sequences. (*RNotebook* used for 4th cycle
- 182 differential binding analysis and 11th cycle differential binding analysis including enrichment analysis
- 183 can be found on https://github.com/KarlisPleiko/apta).

184 RESULTS

185 Aptamer selection

- 186 To identify ccRCC specific aptamers initial randomized oligonucleotide library was subjected to cell-
- 187 SELEX for 11 selection cycles using RCC-MF cell line as a target cell line to identify ccRCC specific
- 188 aptamers and RC-124 cells as a negative control cell line to reduce the nonspecific binding. Cell
- 189 specific aptamer sequence enrichment monitoring was done using flow cytometry (Guava 8HT) after
- 190 4th, 8th and 11th selection cycle. After 4th and 8th selection cycle there was a slight difference between
- 191 the binding of initial randomized oligonucleotide library compared to enriched libraries. After 11th
- selection cycle we observed binding of enriched library to more than >95% of cells. However, the
- 193 observed binding was nonspecific and selected aptamer sequences were binding to both RC-124 (Fig.
- 194 1a) and RCC-MF (Fig. 1b) cell lines.



195

Figure 1. Flow cytometry plots demonstrate fluorescence intensity changes of enriched libraries
during the cell-SELEX procedure. Monitoring binding sequences' enrichment during cell-SELEX to
negative control cells RC-124 (a) and target cells RCC-MF (b). Blue – randomized oligonucleotide
library, green – 4th cycle, red – 8th cycle, black – 11th cycle.

- 200 During further selection and process optimization by changing incubation time, library concentration,
- 201 FBS concentration and temperature complete specificity against RCC-MF cell line was not achieved
- 202 up to 11th cycle. We calculated approximate dissociation constant (K_d) values of the whole enriched
- library after 11th selection cycles by adding increasing concentrations of it to RC-124 cells (Fig. 2a)
- 204 and RCC-MF cells (Fig. 2b) and measuring the green fluorescence increase after the incubation on
- 205 ice using flow cytometry. Based on these measurements we plotted the geometrical mean
- 206 fluorescence to determine approximate K_d values to control cells RC-124 (Fig. 2c) and target cells
- 207 RCC-MF (Fig. 2d). K_d values observed were almost identical, 189 nM for RC-124 and 169 nM for
- 208 RCC-MF cells.
- 209 We concluded that complete specificity against ccRCC cells is not achieved. However, low K_d value
- 210 measured for enriched library after 11th pool suggested that it might be possible that the library
- 211 includes also some ccRCC cell specific aptamers. To explore the differences that might exist within
- the library, we developed differential binding cell-SELEX approach.

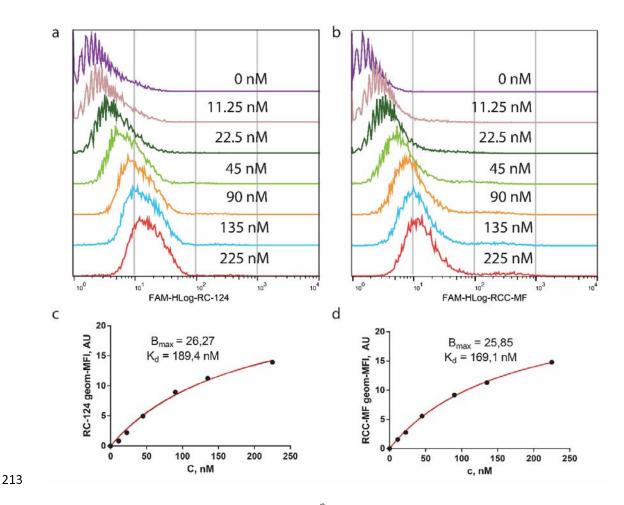


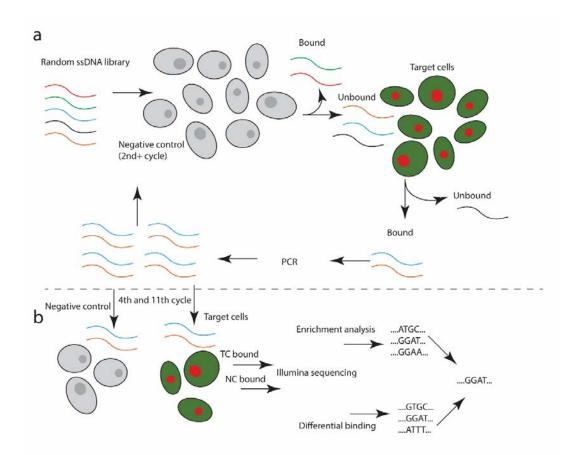
Figure 2. Aptamer binding K_d calculations after 11th selection cycle. Aptamer binding measurements
 by flow cytometry using 11th pool enriched library at different concentrations on control cells RC-124

- 216 (a) and clear cell carcinoma cells RCC-MF (b). K_d value determination using geometrical mean
- 217 fluorescence intensity during the same experiment for RC-124 cells (c) and RCC-MF cells (d).

218 Differential binding cell-SELEX

- 219 Differential binding cell-SELEX process (Fig.3) was performed after 4th and 11th selection cycles. After
- 220 incubation with identically split aptamer libraries and retrieval of bound sequences to both RC-124
- 221 and RCC-MF, we performed two subsequent overlap PCR reactions and confirmed that both
- 222 constructs after 1st overhang PCR and 2nd overhang PCR are of expected size (Fig.4). Quantification
- 223 of final libraries were done using NEBNext Library Quant Kit (New England BioLabs) to quantify only
- those sequences that have flow cell adapters attached to them (Table 1). Overall, our sequencing
- 225 results also confirm feasibility of cell-SELEX experiments performed based on developed protocols
- 226 (Table 1). Sequencing data confirms the successful differential binding cell-SELEX experiments
- 227 based on developed protocols.

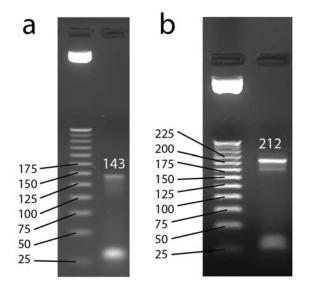
228



229

- 230 Figure 3. Differential binding cell-SELEX workflow combines (a) cell-SELEX selection cycle with (b)
- 231 additional differential binding and data analysis steps to estimate the relative number of aptamer
- sequences within the pool that bind to each type of cells (TC, target cells; NC, negative control).

233



234

- 235 Figure 4. Gel images of aptamers after adding *Illumina* sequencing specific adapters and indexes.
- Aptamers after (a) 1st overhang PCR product with a length of 143 bp and (b) 2nd overhang PCR

237 product with a length of 212 bp.

- 238 Table 1. Aptamer concentration determined by qPCR before sequencing and sequencing reads per
- 239 sample for sequenced aptamer libraries.

Sample No	Samle name	Concentration (nM)	Reads
1	RCC-MF P4_1	142.6	520`534
2	RCC-MF P4_2	90.95	781`654
3	RCC-MF P4_3	185.8	1`130`509
4	RC-124 P4_1	76.23	169`024
5	RC-124 P4_2	67.06	619`635
6	RC-124 P4_3	82.44	548`781
7	RCC-MF P11_1	15.22	277`070
8	RCC-MF P11_2	11.99	498`937
9	RCC-MF P11_3	5.23	326`402

10	RC-124 P11_1	15.58	742`255
11	RC-124 P11_2	28.48	1`142`819
12	RC-124 P11_3	12.88	730`489

240

241 Data analysis for differential binding cell-SELEX

Sequencing was done after 4th and 11th selection cycles. Reads per sample after initial quality filtration,
adapter and constant primer binding region removal and length filtration (40 nt) varied from 169`024
to 1`142`856 (Table 1).

245 Combining all replicates from both samples after data clean-up we identified 3'627'938 unique

sequences within 4th selection cycle experiment and 503'107 unique sequences in 11th selection cycle

247 experiment. After filtering the reads by *edgeR* to remove sequences that had lower count per million

248 (CPM) than two per sample and that were present in less than two replicates, we were left with 1'015

249 unique sequences for 4th cycle aptamers and 35'859 sequences for 11th cycle aptamers.

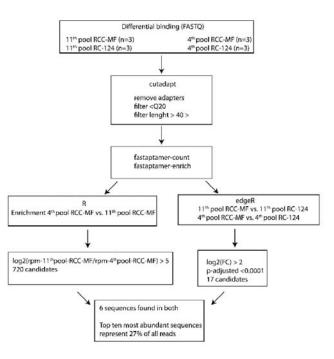
250 For differential binding data analysis (Fig. 5) we further used selected sequences to run edgeR

251 package, a statistical analysis software that is used to estimate differential expression from RNA-seq

252 data. Resulting data were adjusted for multiple comparisons using built-in Benjamini-Hochberg

approach and filtered by removing all sequences that have log₂ fold change (logFC) values less than

two or adjusted p-value was higher than 0.0001.



255

Figure 5. Data analysis pipeline for differential binding cell-SELEX data processing. After trimming using *cutadapt, FASTAptamer* tools *fastaptamer-count* and *fastaptamer-enrich* were used to count the reads for each sequence. Enrichment analysis was done using *R* and *tidyverse* package to identify 720 sequences with enrichment $\log_2 > 5$. *edgeR* was used to perform differential binding analysis resulting in 17 candidate sequences. Matching the sequences resulted in six aptamer candidates that are represented in both analyses.

262 Comparing differential binding datasets using 4th selection cycle enriched library, we were unable to

263 identify any statistically significantly differentially bound sequences based on count per million (CPM)

264 of each sequence and fold change (FC) comparison between two cell lines (Fig. 6a). Most of the

265 sequences bound from 4th cycle enriched library had a low abundance. However, analysis of 11th

selection enriched library discovered 195 statistically significant differentially bound sequences

267 according to the same criteria as described for the first experiment (multiple comparison adjusted p-

268 value < 0.0001, $\log_2(CPM)$ > abs(2)) (Fig. 6b). 178 sequences had $\log_2(CPM)$ < -2 compared to 17

sequences that had $\log_2(CPM) > 2$ (Supplementary Table 1), indicating that more cell type specific

270 sequences were identified for control RC-124 cells than for target RCC-MF cells (Fig. 6c).

271 Enrichment analysis identified 720 unique sequences that have log₂(meanCPM@11th

272 cycle/meanCPM@4th cycle) > 5 or sequence enrichment in CPM terms 32 times from 4th to 11th cycle

273 (Supplementary Table 2). We further combined differential binding results that resulted in 17 unique

274 sequences with 720 sequences obtained from enrichment analysis. We identified only 6 sequences

275 that were present in both datasets (Supplementary Table 3) as the most likely candidates to

276 specifically target ccRCC cells (if log₂ cut off values is decreased to 5, it is possible to identify 6

277 sequences that can be found in both differential binding analysis and in enrichment analysis results).

278 We also ordered all unique sequences that were present in 11th pool by CPM and calculated the log₂

279 enrichment value between 4th and 11th cycle (Supplementary Table 4). Log₂ enrichment values for top

280 10 most abundant sequences ranged from 4.7 to 6.2 and seven out of 10 sequences had Log₂ value

above 5 meaning that these sequences are also included in enrichment analysis results. These 10

most abundant sequences contribute to approximately 27% of all sequencing reads from 11th pool.

283 However, none of the top 10 most abundant sequences passed the statistical significance threshold

284 or FC threshold in differential binding analysis.

285

11

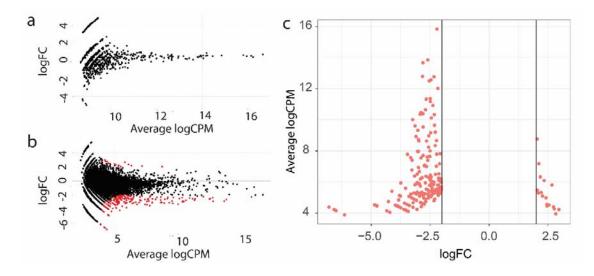


Figure 6. Differential binding cell-SELEX results at 4th cycle (a) and 11th cycle (b) of selection.

288 Negative logFC value indicates increased binding to control cells RC-124, positive logFC value

289 indicates increased binding to ccRCC cells RCC-MF, red dots indicate that these results are

statistically significant according to adjusted p-value < 0.0001 using edgeR and have logFC > 2 in

absolute numbers. All results that fulfil these criteria can be seen in (c).

- 292 Differential binding results confirm that it is possible to use *edgeR* within our pipeline to identify most
- 293 likely candidate molecules for further testing.

286

294 Functional testing of selected lead aptamers

295 For lead aptamer testing using flow cytometry we chose 11 sequences identified by different data

analysis methods (DB, differential binding; EN, enrichment and MB, most abundant). Top three

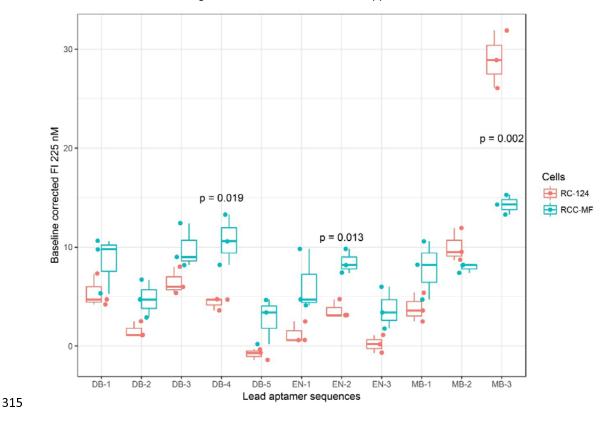
- 297 sequences by each data analysis method were chosen. Differential binding cell-SELEX analysis alone
- sorted by CPM identified sequences DB-1, DB-2 and DB-3. Differential binding cell-SELEX together
- with enrichment analysis sorted by log₂FC identified DB-3, DB-4 and DB-5 sequences. Enrichment
- analysis between 4th and 11th pools by log₂CPM enrichment identified sequences EN-1, EN-2 and EN-
- 301 3. Three most abundant sequences bound to RCC-MF cells were MB-1, MB-2 and MB-3 (Table 2).
- 302 We estimated population shift as a mode of fluorescence intensity (MFI) for each aptamer sample
- 303 (n=3). Data were corrected by subtracting MFI from sample that was incubated with randomized
- 304 starting library (MFI_{lead-sequence}-MFI_{random-library}).
- 305 Table 2. Lead sequences used for confirmatory cell binding test by flow cytometry.

Name	Sequence
DB-1	5'-ATCCAGAGTGACGCAGCA-
	TGCTAGGGTAGGTTGGGCCGGGGGGGGGGGGGGGGGGGG
	TGGACACGGTGGCTTAGT-3'

DB-2	5'-ATCCAGAGTGACGCAGCA-
	AAGAGAAGTATGGGCAGGTTGGGCCGGGGGGGGGGGGGG
	TGGACACGGTGGCTTAGT-3'
DB-3	(5'-ATCCAGAGTGACGCAGCA-
	AGTTGCAGGGTGGGGGGTTGGGTGAAGAGCGATGGAGGGGG-
	TGGACACGGTGGCTTAGT-3')
DB-4	(5'-ATCCAGAGTGACGCAGCA-
	AGGGGGGGGGGGGTGGTTTAGTTGCGTATGGTGGTGGGTG
	TGGACACGGTGGCTTAGT-3')
DB-5	(5'-ATCCAGAGTGACGCAGCA-
	GGCGGTAGTGGAAAGGGTGGTGTGGGTTGGGACAGGAAAG-
	TGGACACGGTGGCTTAGT-3'
EN-1	(5'-ATCCAGAGTGACGCAGCA-
	TGACGGGTGGGTGTGGGGAAGGGAATTTAGATGCGTGG-
	TGGACACGGTGGCTTAGT-3'
EN-2	(5'-ATCCAGAGTGACGCAGCA-
	GGGTGGGTTGGTTGGTGTGGTTTGGAGGGTGGGTGAGATG-
	TGGACACGGTGGCTTAGT-3')
EN-3	(5'-ATCCAGAGTGACGCAGCA-
	GGTGGCAGTTTTGGGGGTTAGGGGTTGATGGGGGTTTGGAGG-
	TGGACACGGTGGCTTAGT-3'
MB-1	(5'-ATCCAGAGTGACGCAGCA-
	GAGTTTGGGGTAAGGGGTTGGGGAGGGACTGTGCGGTTCT-
	TGGACACGGTGGCTTAGT-3')
MB-2	(5'-ATCCAGAGTGACGCAGCA-
	GGGAATGTTGGAGGGTGGAGGGGGGGGGGGGGGGGGGG
	TGGACACGGTGGCTTAGT-3')
MB-3	(5'-ATCCAGAGTGACGCAGCA-
	TGGGGGTAGTGGTGGTTAGGAGTGGAGGCGAGGAGAGCGG-
	TGGACACGGTGGCTTAGT-3')

306

307 Corrected MFIs were compared with t-test (significance defined as p < 0.05, n=3) using GraphPad 308 Prism to determine if our identified sequences altogether bind more to RCC-MF cells than to RC-124 309 cells. Three sequences (DB-4, EN-2, MB-3) were confirmed to be differentially bound using flow cytometry by comparing MFIs (Fig. 7). While MB-3, identified as the 3rd most abundant sequence, was 310 311 significantly (p=0.002) differentially bound, it was targeted towards RC-124 cells. EN-2 sequence was 312 identified using enrichment analysis and was statistically significantly (p=0.013) binding to RCC-MF 313 cells. DB-4 was significantly (p=0.019) more bound to RCC-MF cells and was identified through



314 combined differential binding cell-SELEX and enrichment approach.

Figure 7. Comparison of mode of fluorescence intensities from lead sequences binding to RC-124 and RCC-MF cells. Top three sequences were identified by differential binding alone sorting by CPM (DB-1, DB-2, DB-3), differential binding together with enrichment analysis sorting by log₂FC (DB-3, DB-4, DB-5), enrichment analysis alone sorting by log₂CPM enrichment (EN-1, EN-2, EN-3) or by choosing most abundant sequences in sequencing dataset (MB-1, MB-2, MB-3). Upper hinges correspond to the first and third quartiles, whiskers mark 1.5*IQR. Statistical significance was determined with t-test using GraphPad Prism software.

323 DISCUSSION

Recent review on aptamer discovery mentions that there are 141 entries of aptamer selection against live cells as of 2017. For comparison, proteins as targets have 584 entries and small molecules have 234 research entries (1). This is not surprising considering the advanced technological procedure involved in cell-SELEX method compared to protein or small molecule SELEX. Several methods have been developed in recent years to improve the success rate of cell-SELEX, for example, HT-SELEX (10), FACS-SELEX (21) and cell-internalization SELEX (22). HTS adaptation for aptamer sequencing has been described as one of the most fundamental changes to aptamer selection technology (23).

The main goal achieved in this research is the development of differential binding cell-SELEX method.
 This method can identify cell type specific aptamer sequences from cell-SELEX selection pools that

would not be selected by other cell-SELEX methods and thus would remain overlooked by theinvestigators.

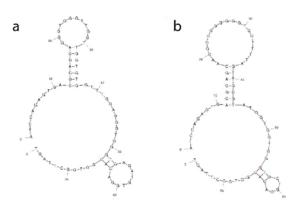
analysis, which means comparison of abundance of one particular sequence at the beginning of the SELEX procedure to the abundance of the same sequence after SELEX procedure. Enrichment analysis can identify a large number of oligonucleotides with very similar log_2 enrichment values as can be seen by our results (Supplementary Table 2). However, it is rarely useful for cell-SELEX because of the high possibility to enrich non-specific sequences. Using enrichment analysis with cutoff value of $log_2 > 5$, we identified 720 sequences to be further tested. However, when the same

Currently the most often used analysis for aptamer finding using HTS data includes enrichment

- 342 sequencing dataset was submitted for differential binding analysis using *edgeR*, we identified 17
- 343 sequences that were more abundant on the surface of RCC-MF cells than on RC-124 cells.
- 344 Enrichment analysis identified one sequence (EN-2) that was statistically significantly (p=0.013) more
- bound to target RCC-MF cells, as confirmed also by flow cytometry (Fig.7). We were able to confirm
- 346 using flow cytometry that another sequence (DB-4), identified by combined differential binding cell-
- 347 SELEX and enrichment analysis, was statistically significantly (p=0.019) more bound to the RCC-MF
- 348 cells. Importantly, DB-4 was found between 720 sequences identified using enrichment analysis, but
- only as the 528th most enriched sequence. This provides scientific evidence that our approach can be
- used to identify lead aptamers that most likely would be lost during enrichment analysis.
- 351 MB-3 that was one of the most abundant sequences in the dataset showed statistically significant
- 352 binding to RC-124 cells. MB-3 was not identified neither in enrichment analysis results, nor in
- 353 differential binding cell-SELEX results. However, seven out of 10 most abundant aptamer sequences
- after cell-SELEX process were enriched above the set cut-off value $\log_2 > 5$ and thus did appear in
- 355 enrichment analysis results. None of these sequences appeared in differential binding results
- because they did not pass the statistical significance test applied to logFC. These observations are in
- 357 line with previous statements that the most abundant aptamer sequences are not necessary the best
- binders (24). This proves the value of differential binding approach for excluding the non-specifically
- 359 enriched sequences during cell-SELEX procedure.

335

Comparing secondary structures using *mfold* web server (25) we discovered surprising structural similarity between sequences EN-2 and DB-4 with two stem-loop motifs (Fig.8). MB-3 and other lead aptamer sequences that did not show statistically significant binding differences between RC-124 and RCC-MF cells based on flow cytometry had distinctly different predicted structures (Supplementary Data 5).



365

Figure 8. Predicted secondary structures for RCC-MF specific statistically significantly bound lead aptamers EN-2 (a) and DB-4 (b). Structures were predicted at 37° C using mfold web server with 0.005 mM Mg²⁺ concentration.

369 Differential binding cell-SELEX uses *edgeR* to compare how all sequences that can be found in final 370 enriched aptamer library interact with control and target cells and estimate the statistical significance

371 of these differences. There are several bioinformatics tools available to analyse statistical significance

of differential expression for RNA-seq data (15, 26, 27). To the best of our knowledge, so far none of

these tools have been applied for estimation of differentially bound aptamers on the cell surface.

374 edgeR was chosen because it is compatible with the existing data analysis workflows in R (28).

375 Combination of enrichment analysis and differential binding approach provides an algorithm to choose376 target sequences for further analysis.

Altogether, we demonstrate a combined analysis pipeline that can be used to identify lead aptamers
 from low binding specificity aptamer libraries after cell-SELEX experiments. We propose fast and
 practical high throughput aptamer identification method to be used with cell-SELEX technique to

380 increase successful aptamer selection rate against live cells.

Higher number of sequencing reads during differential binding cell-SELEX could even further increase
 the likelihood to identify low abundance, but differentially bound sequences specific to cells of interest.
 Sequences that were present only in one replicate from each selection pool were discarded. After 4th
 selection cycle only few sequences were present in more than one sequencing replicates (Fig.6a)

385 compared to 11th cycle (Fig 6b). Increased number of reads would cover more libraries that are

386 diverse and make it possible to identify differentially bound aptamers using fewer selection cycles.

387 Cell-SELEX design described in this research uses commercially available human cells RCC-MF and 388 RC-124 both as a target and negative control. We are first to use these cell lines for aptamer selection 389 using cell-SELEX approach. However, it could be more suitable to use patient-matched primary cells 390 isolated from tumour site and adjacent healthy kidney tissue within few passages after isolation while 391 cells are most likely to represent the diversity found in clinical settings (29).

- 392 Differential binding cell-SELEX method developed here can be used to accelerate aptamer selection
- 393 based on HTS analysis. Additional information from differential binding cell-SELEX reduces the time
- 394 needed to identify aptamers. This can lead to broader use of cell-SELEX technique not only to identify
- 395 aptamers against cell lines, but also against primary cells isolated from patient samples.
- 396 We conclude that differential binding cell-SELEX method can be used to characterise not only
- 397 enrichment of sequences between selection cycles, but also to select aptamer sequences that
- 398 selectively bind to the target and control cells. We demonstrate the feasibility of our approach by
- showing cell-line specific aptamer identification against ccRCC cell line RCC-MF and RC-124 cell line
- 400 from healthy kidney tissue.

401 **AVAILABILITY**

- 402 edgeR is available as Bioconductor package (<u>http://bioconductor.org/packages/edgeR/</u>),
- 403 FASTAptamer was downloaded from github (https://github.com/FASTAptamer/FASTAptamer),
- 404 cutadapt was installed using Bioconda (30)
- 405 RNotebooks for data analysis using *tidyverse* (31) are available here:
- 406 https://github.com/KarlisPleiko/apta

407 ACCESSION NUMBERS

408 Sequencing data are available at SRA under accession number PRJEB28411.

409 SUPPLEMENTARY DATA

- 410 Supplementary Table 1. List of 17 aptamer variable sequences identified using differential binding.
- 411 Supplementary Table 2. List of 720 aptamer variable sequences identified using enrichment analysis.
- 412 Supplementary Table 3. List of 6 aptamer variable sequences identified using combined enrichment
- 413 analysis and differential binding.
- 414 Supplementary Table 4. List of 10 most abundant aptamer variable sequences.
- 415 Supplementary Table 5. List of lead aptamer sequences tested by flow cytometry.

416 ACKNOWLEDGEMENT

- 417 K.P., U.R., E.V. conceived and designed the project. K.P., L.S., V.P., K.M. carried out the
- 418 experiments. K.P., U.R. wrote the paper. All read and approved the final manuscript.

419 FUNDING

420 This work was supported by University of Latvia Foundation [grant number 2182].

421 CONFLICT OF INTEREST

422 None declared

423 REFERENCES

- 424 1. Dunn, M.R., Jimenez, R.M. and Chaput, J.C. (2017) Analysis of aptamer discovery and technology.
 425 *Nat. Rev. Chem.*, 1, 0076.
- 2. Pereira, R.L., Nascimento, I.C., Santos, A.P., Ogusuku, I.E.Y., Lameu, C., Mayer, G. and Ulrich, H.
 (2018) Aptamers: novelty tools for cancer biology. *Oncotarget*, **9**, 26934–26953.
- 3. Zhou, J. and Rossi, J. (2017) Aptamers as targeted therapeutics: Current potential and challenges. *Nat. Rev. Drug Discov.*, **16**, 181–202.
- 4. A Safety and Efficacy Study of E10030 (Anti-PDGF Pegylated Aptamer) Plus Lucentis for
 Neovascular Age-Related Macular Degeneration (NCT01089517).
- 432 5. A Phase II Open-label, Multicenter Extension Study to Assess the Long-term Safety and Efficacy of
 433 Vamorolone in Boys with Duchenne Muscular Dystrophy (DMD)(EudraCT No: 2016-004263-38).
- 6. NOX-A12 in Combination With Bendamustine and Rituximab in Relapsed Chronic Lymphocytic
 Leukemia (CLL)(NCT01486797).
- 436 7. Kaur, H., Bruno, J.G., Kumar, A. and Sharma, T.K. (2018) Aptamers in the Therapeutics and
 437 Diagnostics Pipelines. *Theranostics*, **8**, 4016–4032.
- 438 8. Ellington, A.D. and Szostak, J.W. (1990) In vitro selection of RNA molecules that bind specific
 439 ligands. *Nature*, **346**, 818.
- 440 9. Hicke, B.J., Marion, C., Chang, Y.F., Gould, T., Lynott, C.K., Parma, D., Schmidt, P.G. and Warren, S.
 441 (2001) Tenascin-C Aptamers Are Generated Using Tumor Cells and Purified Protein. *J. Biol.*

442 *Chem.*, **276**, 48644–48654.

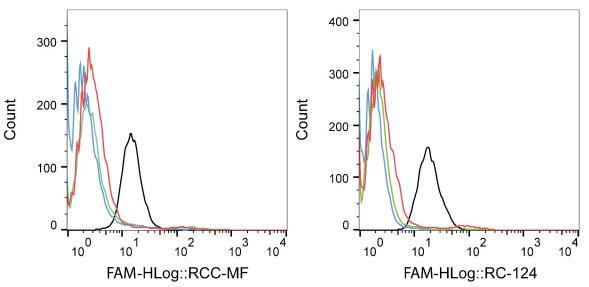
- Theorem 2 and Stormo, G.D. (2009) Inferring binding energies from selected binding sites.
 PLoS Comput. Biol., **5**.
- 445 11. Kahsai,A.W., Wisler,J.W., Lee,J., Ahn,S., Cahill,T.J., Dennison,S.M., Staus,D.P., Thomsen,A.R.B.,
 446 Anasti,K.M., Pani,B., *et al.* (2016) Conformationally selective RNA aptamers allosterically
 447 modulate the β 2-Adrenoceptor. *Nat. Chem. Biol.*, **12**, 709–716.
- 448 12. Alam,K.K., Chang,J.L. and Burke,D.H. (2015) FASTAptamer: A bioinformatic toolkit for high449 throughput sequence analysis of combinatorial selections. *Mol. Ther. Nucleic Acids*, 4, 1–10.
- Hoinka, J., Backofen, R. and Przytycka, T.M. (2018) AptaSUITE: A Full-Featured Bioinformatics
 Framework for the Comprehensive Analysis of Aptamers from HT-SELEX Experiments. *Mol. Ther. Nucleic Acids*, **11**, 515–517.
- 453 14. Werner, T. (2010) Next generation sequencing in functional genomics. *Brief. Bioinform.*, **11**, 499–
 454 511.
- 15. Robinson, M.D., McCarthy, D.J. and Smyth, G.K. (2010) {edgeR}: a {Bioconductor} package for
 differential expression analysis of digital gene expression data. *Bioinformatics*, **26**, 139–140.
- 457 16. Martin,M. (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads.
 458 *EMBnet.journal*, **17**, 10.

459 17. Sefah, K., Shangguan, D., Xiong, X., O'Donoghue, M.B. and Tan, W. (2010) Development of DNA 460 aptamers using cell-selex. Nat. Protoc., 5, 1169-1185. 461 18. Illumina (2013) 16S Metagenomic Sequencing Library Preparation. Illumina.com. 462 19. Alam,K.K., Chang,J.L. and Burke,D.H. (2015) FASTAptamer: A bioinformatic toolkit for high-463 throughput sequence analysis of combinatorial selections. Mol. Ther. - Nucleic Acids, 4, 1–10. 464 20. Benjamini, Y. and Hochberg, Y. (1995) Controlling The False Discovery Rate - A Practical And 465 Powerful Approach To Multiple Testing. J. R. Stat. Soc., Ser. B, 57, 289-300. 466 21. Mayer, G., Ahmed, M.S.L., Dolf, A., Endl, E., Knolle, P.A. and Famulok, M. (2010) Fluorescence-467 activated cell sorting for aptamer SELEX with cell mixtures. Nat. Protoc., 5, 1993-2004. 468 22. Thiel,W.H., Thiel,K.W., Flenker,K.S., Bair,T., Dupuy,A.J., McNamara,J.O., Miller,F.J. and 469 Giangrande, P.H. (2015) Cell-Internalization SELEX: Method for Identifying Cell-Internalizing 470 RNA Aptamers for Delivering siRNAs to Target Cells. In Methods in molecular biology (Clifton, 471 N.J.).Vol. 1218, pp. 187–199. 472 23. Ozer, A., Pagano, J.M. and Lis, J.T. (2014) New technologies provide quantum changes in the 473 scale, speed, and success of SELEX methods and aptamer characterization. Mol. Ther. -474 Nucleic Acids, 3, 1–18. 475 24. Hoinka, J., Berezhnoy, A., Dao, P., Sauna, Z.E., Gilboa, E. and Przytycka, T.M. (2015) Large scale 476 analysis of the mutational landscape in HT-SELEX improves aptamer discovery. Nucleic Acids 477 Res., 43, 5699-5707. 478 25. Zuker, M. (2003) Mfold web server for nucleic acid folding and hybridization prediction. Nucleic 479 Acids Res., 31, 3406–3415. 480 26. Trapnell, C., Roberts, A., Goff, L., Pertea, G., Kim, D., Kelley, D.R., Pimentel, H., Salzberg, S.L., 481 Rinn, J.L. and Pachter, L. (2012) Differential gene and transcript expression analysis of RNA-seq 482 experiments with TopHat and Cufflinks. Nat. Protoc., 7, 562-578. 483 27. Ritchie, M.E., Phipson, B., Wu, D., Hu, Y., Law, C.W., Shi, W. and Smyth, G.K. (2015) limma powers 484 differential expression analyses for RNA-sequencing and microarray studies. Nucleic Acids Res., 485 43, e47-e47. 486 28. Team, R.C. (2018) R: A Language and Environment for Statistical Computing. 487 29. Lobo, N.C., Gedye, C., Apostoli, A.J., Brown, K.R., Paterson, J., Stickle, N., Robinette, M., Fleshner, N., 488 Hamilton, R.J., Kulkarni, G., et al. (2016) Efficient generation of patient-matched malignant and 489 normal primary cell cultures from clear cell renal cell carcinoma patients: Clinically relevant 490 models for research and personalized medicine. BMC Cancer, 16, 1-15. 491 30. Grüning, B., Dale, R., Sjödin, A., Chapman, B.A., Rowe, J., Tomkins-Tinch, C.H., Valieris, R., Köster, J. 492 and Bioconda Team (2018) Bioconda: sustainable and comprehensive software distribution for 493 the life sciences. Nat. Methods, 15, 475-476. 494 31. Wickham, H. (2017) tidyverse: Easily Install and Load the 'Tidyverse'. 495 496 TABLE AND FIGURES LEGENDS

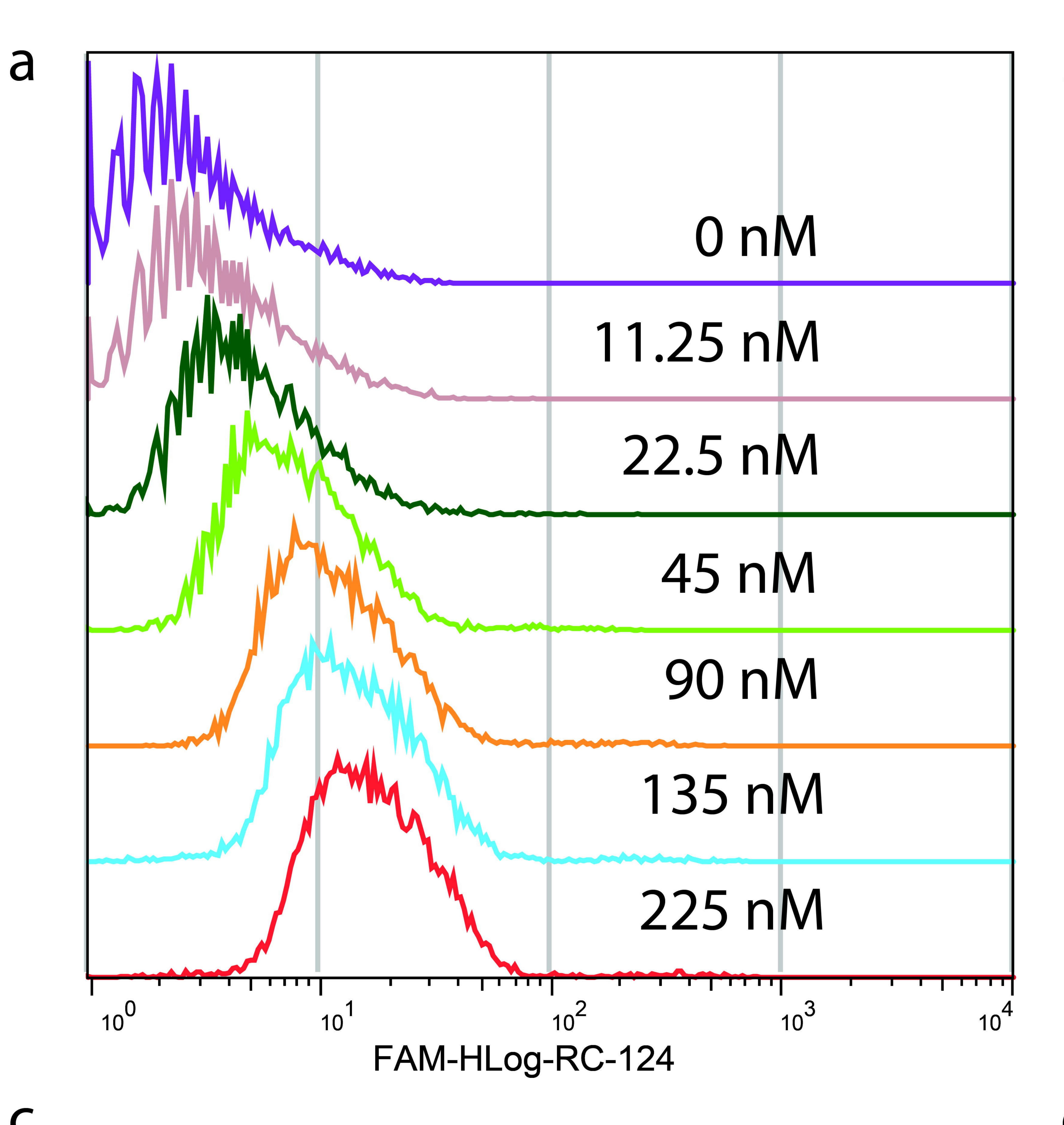
- 497 Figure 1. Flow cytometry plots demonstrate fluorescence intensity changes of enriched libraries
- 498 during the cell-SELEX procedure. Monitoring binding sequences' enrichment during cell-SELEX to
- 499 negative control cells RC-124 (a) and target cells RCC-MF (b). Blue randomized oligonucleotide
- 500 library, green -4^{th} cycle, red -8^{th} cycle, black -11^{th} cycle
- 501 Figure 2. Aptamer binding K_d calculations after 11th selection cycle. Aptamer binding measurements
- 502 by flow cytometry using 11th pool enriched library at different concentrations on control cells RC-124
- 503 (a) and clear cell carcinoma cells RCC-MF (b). K_d value determination using geometrical mean
- 504 fluorescence intensity during the same experiment for RC-124 cells (c) and RCC-MF cells (d).
- 505 Figure 3. Differential binding cell-SELEX workflow combines (a) cell-SELEX selection cycle with (b)
- 506 additional differential binding and data analysis steps to estimate the relative number of aptamer
- 507 sequences within the pool that bind to each type of cells (TC, target cells; NC, negative control).
- 508 Figure 4. Gel images of aptamers after adding *Illumina* sequencing specific adapters and indexes.
- 509 Aptamers after (a) 1st overhang PCR product with a length of 143 bp and (b) 2nd overhang PCR
- 510 product with a length of 212 bp.
- Table 1. Aptamer concentration determined by qPCR before sequencing and sequencing reads persample for sequenced aptamer libraries.
- 513 Figure 5. Data analysis pipeline for differential binding cell-SELEX data processing. After trimming
- 514 using cutadapt, FASTAptamer tools fastaptamer-count and fastaptamer-enrich were used to count the
- 515 reads for each sequence. Enrichment analysis was done using R and tidyverse package to identify
- 516 720 sequences with enrichment $log_2 > 5$. edgeR was used to perform differential binding analysis
- resulting in 17 candidate sequences. Matching the sequences resulted in six aptamer candidates that
- 518 are represented in both analyses.
- 519 Figure 6. Differential binding cell-SELEX results at 4th cycle (a) and 11th cycle (b) of selection.
- 520 Negative logFC value indicates increased binding to control cells RC-124, positive logFC value
- 521 indicates increased binding to ccRCC cells RCC-MF, red dots indicate that these results are
- 522 statistically significant according to adjusted p-value < 0.0001 using edgeR and have logFC > 2 in
- 523 absolute numbers. All results that fulfil these criteria can be seen in (c).
- 524 Table 2. Lead sequences used for confirmatory cell binding test by flow cytometry.
- 525 Figure 7. Comparison of mode of fluorescence intensities from lead sequences binding to RC-124
- 526 and RCC-MF cells. Top three sequences were identified by differential binding alone sorting by CPM
- 527 (DB-1, DB-2, DB-3), differential binding together with enrichment analysis sorting by log₂FC (DB-3,
- 528 DB-4, DB-5), enrichment analysis alone sorting by log₂CPM enrichment (EN-1, EN-2, EN-3) or by
- 529 choosing most abundant sequences in sequencing dataset (MB-1, MB-2, MB-3). Upper hinges
- 530 correspond to the first and third quartiles, whiskers mark 1.5*IQR. Statistical significance was
- 531 determined with t-test using GraphPad Prism software.

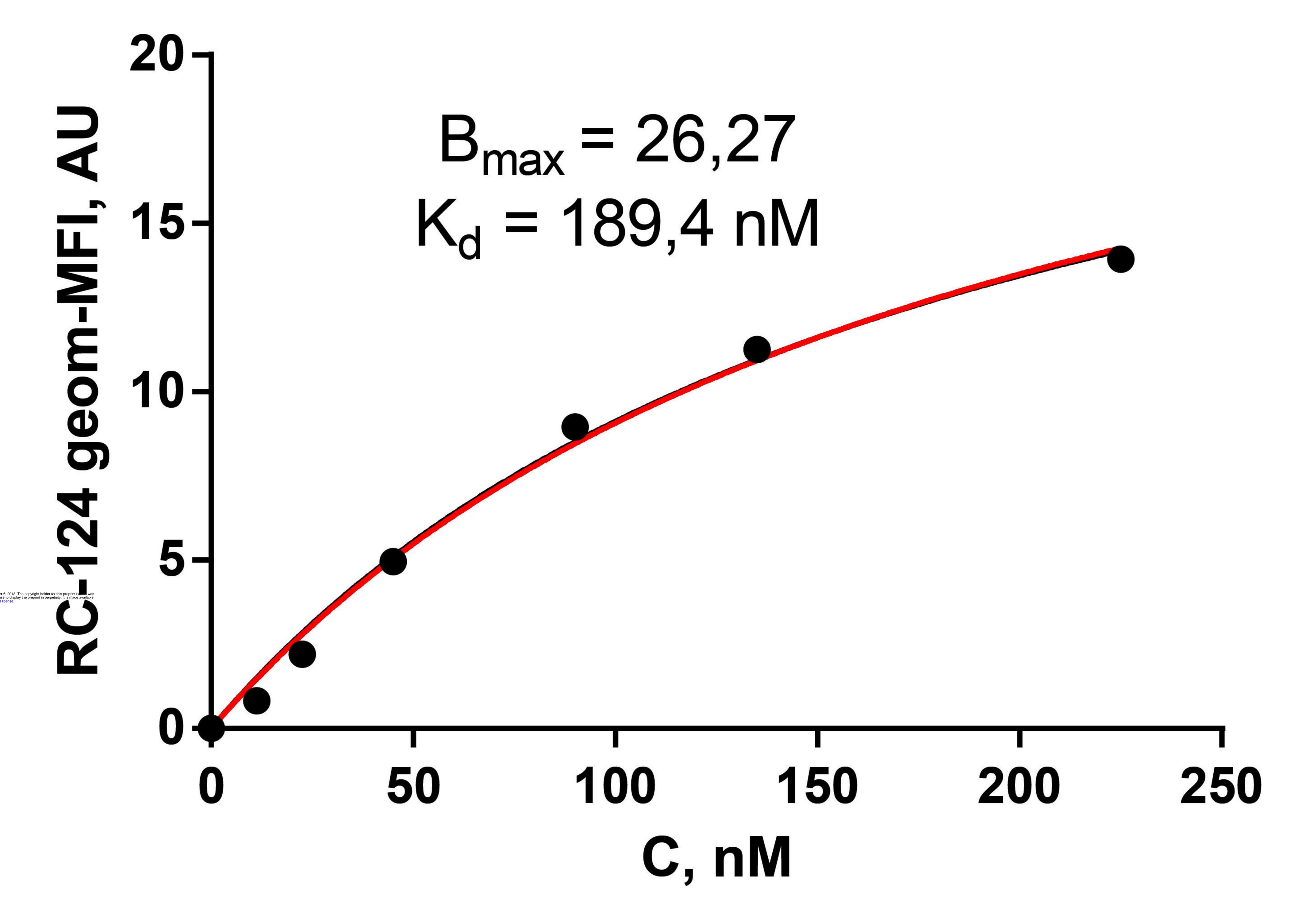
- 532 Figure 8. Predicted secondary structures for RCC-MF specific statistically significantly bound lead
- aptamers EN-2 (a) and DB-4 (b). Structures were predicted at 37° C using mfold web server with
- 534 0.005 mM Mg²⁺ concentration.

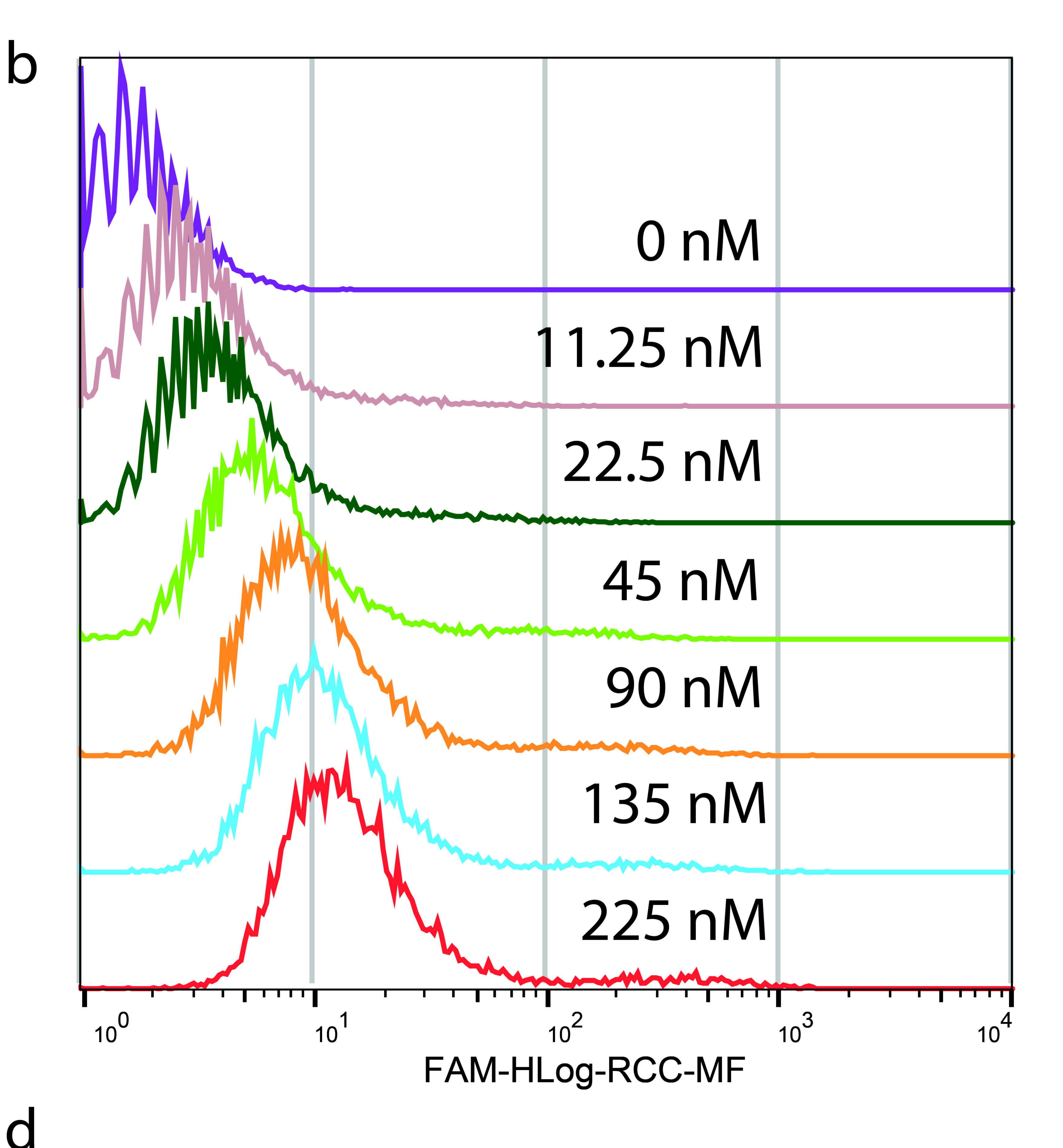


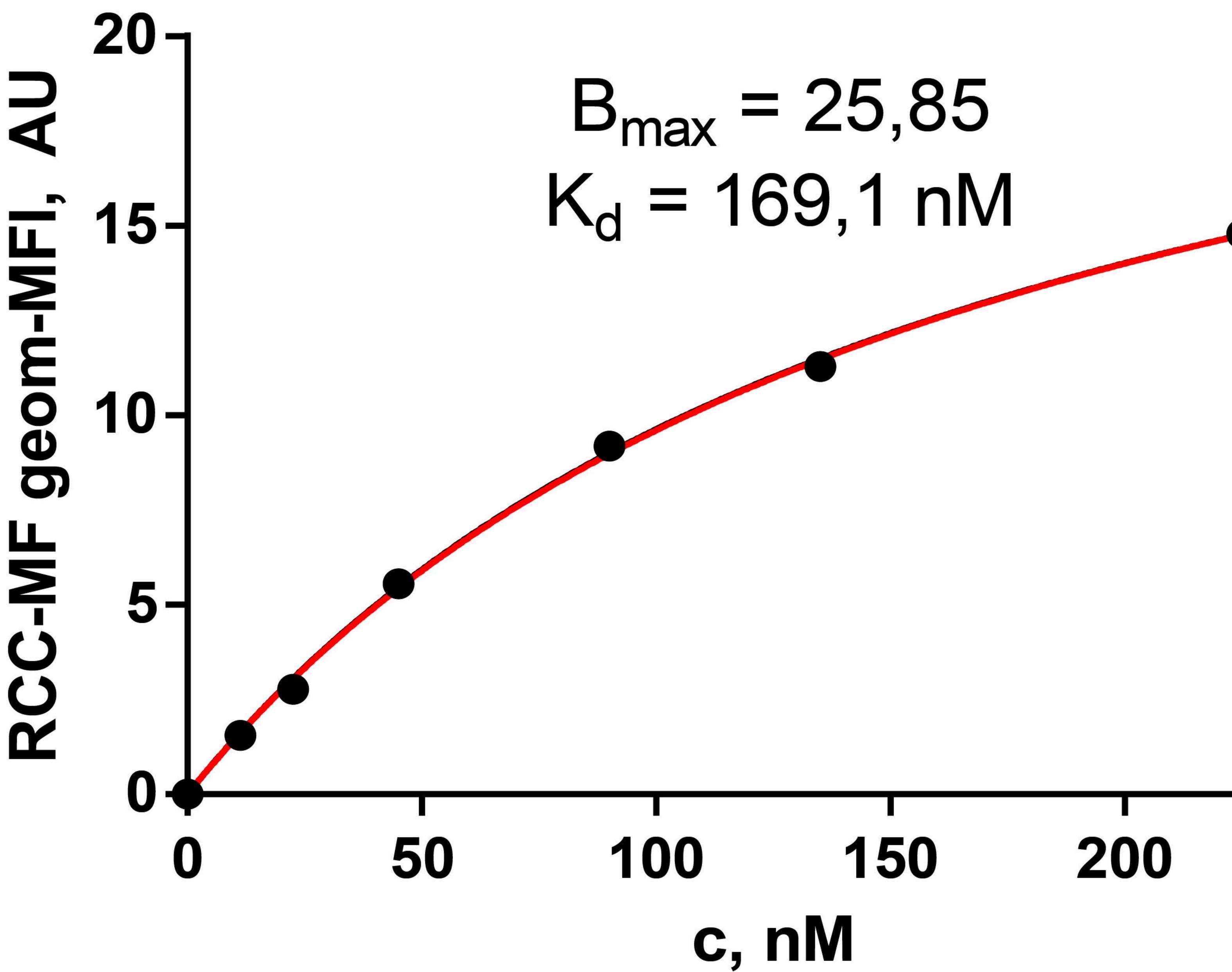


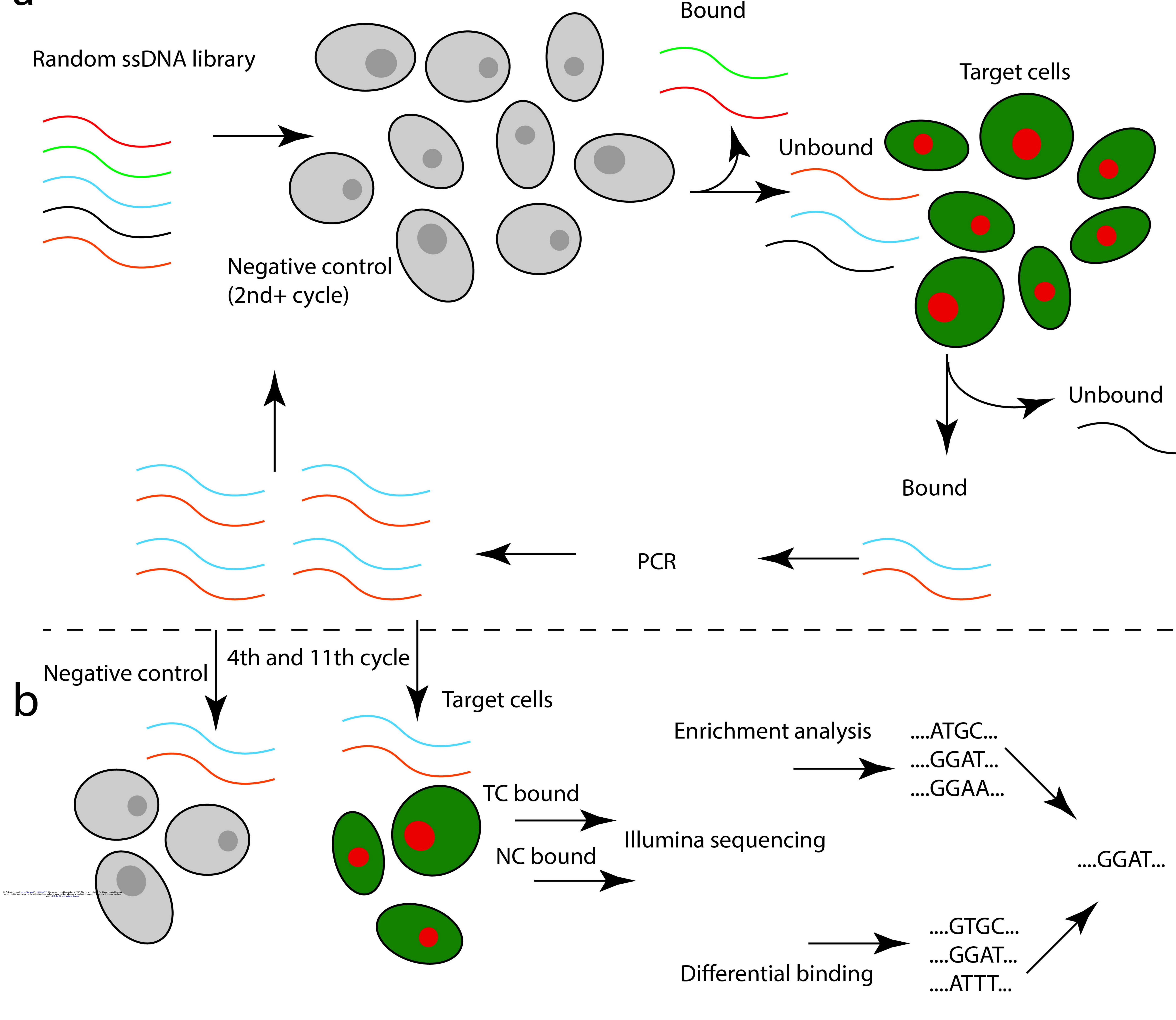
b



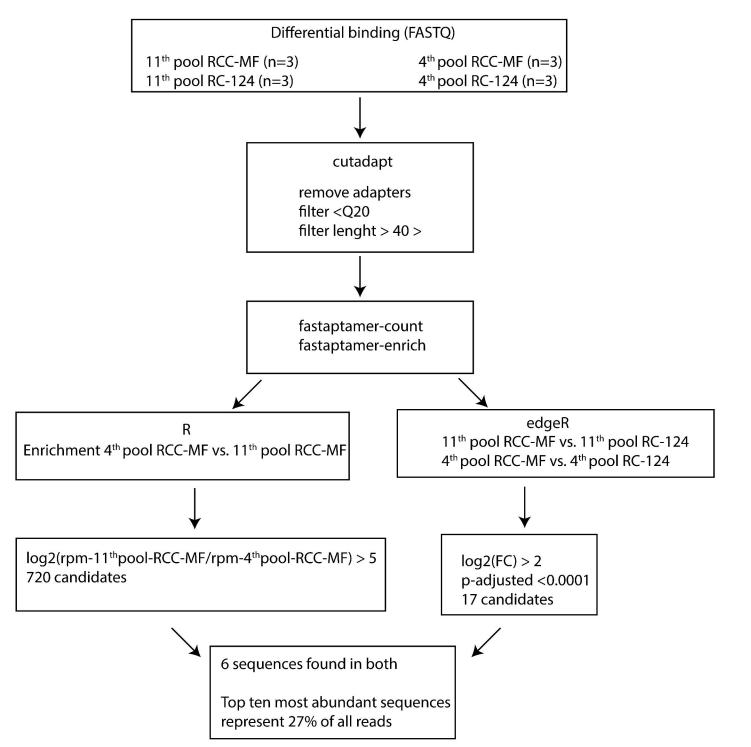


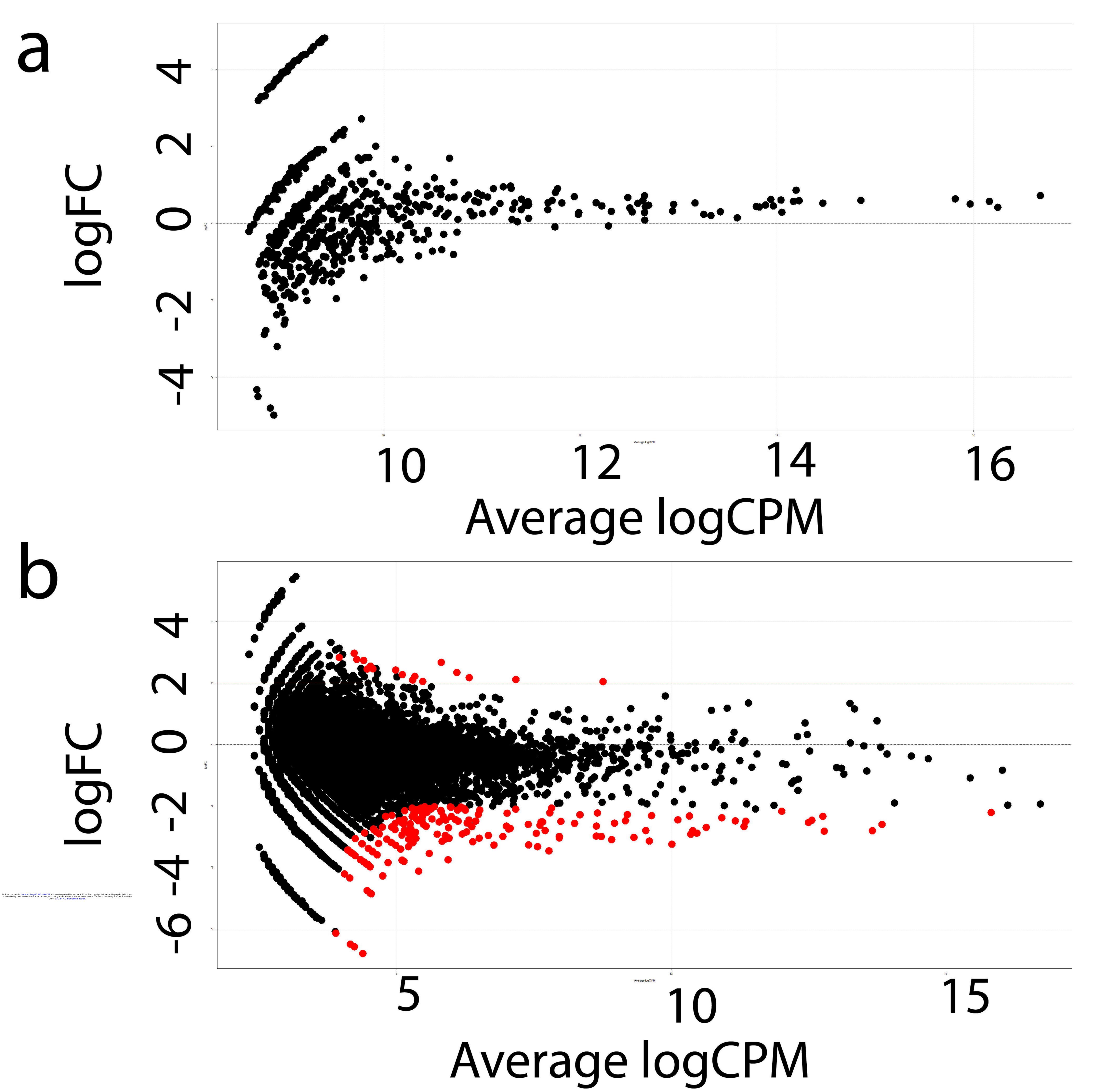


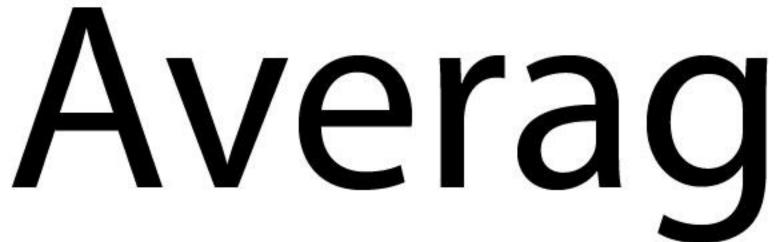


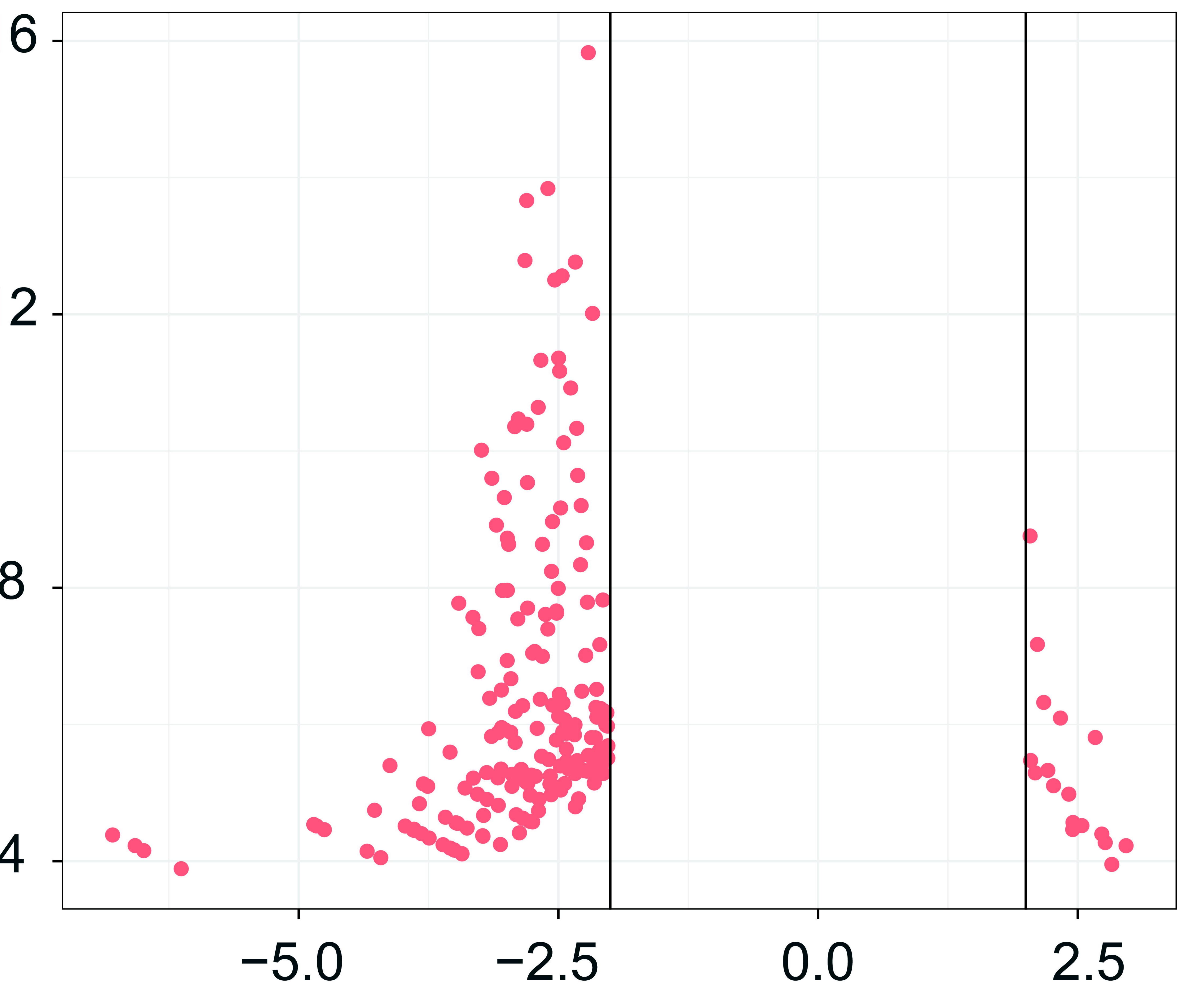


b а 150 -75 -









IOGEC

