1     **Course**: Binp39

2

3     **Credits**: 30

4

5     **Project name**: Integrative Bioinformatics for Rational AAV Vector Design

6

7     **Student**: Hooi Min Tan Grahn

8

9     **Supervisor**: Marcus Davidsson

10     **Email:** Marcus.Davidsson@raaven.se

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34    **Integrative Bioinformatics for Rational AAV Vector Design**

35

36    **Abstract**

37    This research project seeks to understand the complex molecular mechanisms behind capsid

38    modification and tropism in Adeno-associated virus (AAV) vectors using a comprehensive

39    bioinformatics approach. The specific objectives include developing a bulk sequencing

40    workflow to identify recurring patterns within computationally segmented amino acid (aa)

41    sequences (7-aa-long polypeptides) in cells following AAV vector administration. By

42    utilizing mRNA Illumina sequencing, the project will investigate the most effective aa

43    patterns that replicate in tissue samples post-virus infection and their effects on capsid

44    modification and DNA virus tropism integration. The bioinformatics analyses in this study

45    include several key methodologies. First, sequencing reads undergo preprocessing, which

46    involves quality trimming to remove low-quality bases and artifacts, ensuring high fidelity in

47    downstream analysis. Following this, the reads are aligned and mapped to a customized

48    reference library database tailored for the study, which allows for accurate matching and

49    analysis of the sequences of interest. Once the reads are mapped, barcode tracking is used to

50    quantify the frequencies of specific amino acids within the sequences, facilitating a detailed

51    assessment of amino acids profile present in the viral capsid proteins across different

52    samples. This step is essential for identifying which amino acids are most commonly

53    associated with certain modifications or tropism changes in the viral capsid. Finally, recurring

54    sequence patterns of peptides are identified to uncover motifs or sequence signatures that

55    may influence capsid properties. These recurring patterns can provide insights into how

56    specific peptide sequences contribute to the virus's ability to change its tropism, ultimately

57    aiding in the understanding of viral behavior and enhancing the design of more effective

58    AAV vectors. By integrating these bioinformatics methods, this research endeavors to

59    provide comprehensive insights into the molecular intricacies underlying capsid modification

60  and tropism. These insights hold immense potential for the rational design and optimization

61  of AAV vectors for diverse applications in gene therapy and biotechnology, thereby

62  contributing significantly to advancements in the field.

63  Availability and Implementation: This project is built using Bash, Python, Awk and R. All

64  scripts used in this project are available for download from

65  https://github.com/hooimin7/BINP39_project

66  Contact: ho4588ta-s@student.lu.se

67

## Introduction

69  Adeno-associated virus (AAV)-mediated gene delivery represents a promising strategy for

70  advancing the next generation of clinical gene therapies [1,2]. The recombinant AAV (rAAV)

71  is one of the leading vectors for gene therapy applications that deliver gene-editing enzymes,

72  antibodies, RNA interference molecules and peptides to anatomical sites where viruses

73  accumulate and exert therapeutic effects in AAV gene therapy [3]. rAAV vectors are limited

74  by the need for host-cell synthesis of double-stranded DNA from their single-stranded

75  genomes. To overcome this, researchers developed self-complementary rAAV (scAAV)

76  vectors that bypass this step by packaging DNA dimers, allowing them to form double-

77  stranded DNA upon uncoating [4,5]. In general, the single-stranded DNA (ssDNA) genome

78  of wild-type AAV2 (wtAAV2) is approximately 4.7 kb in length. AAV serotypes are

79  classified in a separate genus within the *Parvoviridae* family, known as *Dependovirus*,

80  because productive infection generally requires either co-infection with a helper virus (such

81  as adenovirus or herpes simplex virus) or the induction of cellular stress [6]. In the context of

82  AAV, serotypes refer to different naturally occurring or engineered variants of AAV that

83  have unique capsid proteins. These capsid differences give each serotype specific properties,

84  such as the ability to infect certain types of cells or tissues. This makes certain AAV

85    serotypes more suitable for targeting specific tissues in gene therapy applications, such as the

86    brain, liver, or muscle [7,8]. Other findings suggested that AAV2 capsids could tolerate

87    ligand insertions. The researchers inserted a 14-amino-acid targeting peptide (L14) into six

88    regions of the AAV2 capsid protein, successfully generating three mutants displayed L14 on

89    their surface, opening possibilities for targeted gene therapy [9]. The AAV capsid is

90    assembled from viral structural proteins encoded by the cap gene. Capsid assembly is

91    facilitated by an assembly-activating protein, followed by the packaging of the viral genome

92    into the capsid, a process mediated by replication proteins flanked by two 145 bp long

93    inverted terminal repeats (ITRs) [10] (Fig 1A). To enhance the development of AAV vectors

94    and advance capsid engineering toward late-stage clinical trials, Marcus and colleagues

95    developed a method called barcoded rational AAV vector evolution (BRAVE), which

96    incorporates the advantages of rational design [3]. Rational design leverages existing

97    knowledge of AAV biology and structure, typically targeting surface-exposed regions of the

98    AAV capsid [11]. The number of variants tested is relatively small, often in the tens or

99    hundreds, as the designs are guided by specific hypotheses or known mechanisms.
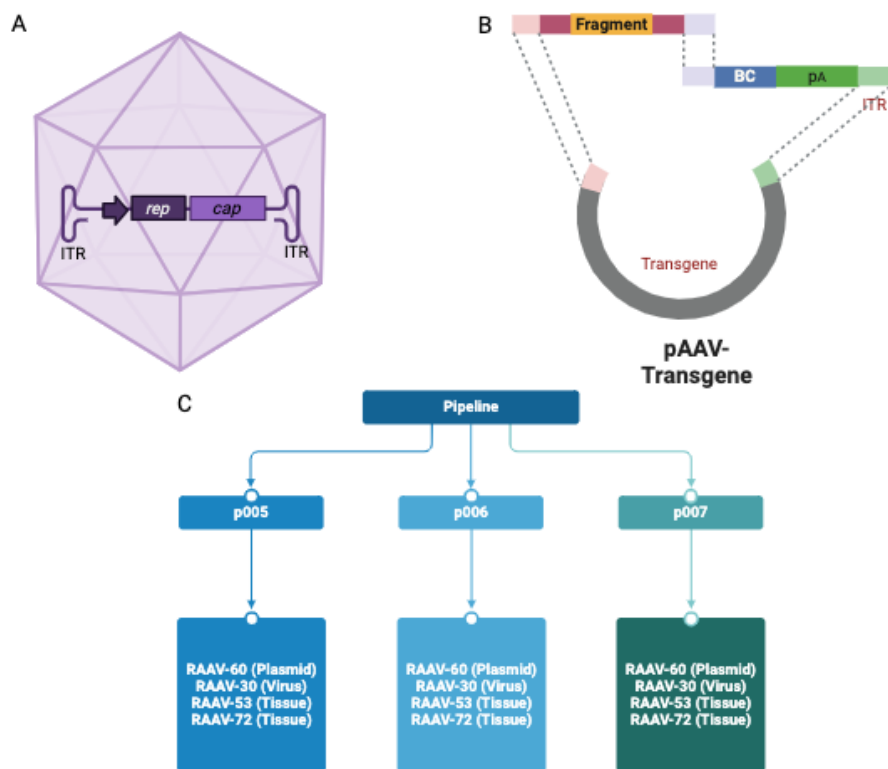
100    Additionally, rational design can be enhanced with computational tools to improve the

101    precision of the modifications [3]. A previous study addressed the need for error-free barcode

102    clustering in *in vivo* applications like cellular fate mapping [12]. While various clustering

103    methods have been attempted on samples with unknown diversity, the study validated the use

104    of the message passing clustering algorithm in the Starcode software by leveraging high-

105    diversity, barcoded fragment libraries. The findings indicated that while this algorithm was

106    robust for clustering degenerate barcodes, it could produce false clusters if the Levenshtein

107    distance threshold was too high. The study recommended validating clustering with paired-

108    end sequencing data from an *in vitro* plasmid library to establish a subset of "valid barcodes,"

109    minimizing false clustering in the *in vivo* data [13]. High-resolution lineage tracking using

110 DNA barcodes could detect rare beneficial mutations, estimate their timing, and assess fitness

111 effects. By tracking multiple lineages simultaneously, this method provided detailed insights

112 into population dynamics [14]. Our study further refines AAV production plasmid to connect

113 capsid structure with an *in vivo* molecular barcode. The barcode was inserted into the 3′

114 untranslated region (UTR) of green fluorescence protein within a gutted self-complementary

115 AAV (scAAV) genome while the AAV2 Rep/Cap genes were expressed from the same

116 plasmid but positioned outside the ITRs, thereby linking the capsid structure with the

117 barcode. In this system, peptides (seven amino acid) or fragment were incorporated at the

118 N587 position of the VP1 capsid protein (Fig 1B). My responsibility in this project involves

119 reproducing the pipeline [3] and tailoring the bioinformatics workflow to handle three distinct

120 plasmid libraries. Libraries p005 and p007 were designed using rational design principles,

121 whereas p006 was engineered by incorporating degenerate NNK codons to introduce amino

122 acid randomization (Fig 1C). The NNK configuration is a method used in genetic

123 mutagenesis to efficiently encode all 20 amino acids with only one stop codon, reducing the

124 total sequence library size. In this approach, "K" can be either "T" or "G," resulting in a 32-

125 fold degeneracy (32 unique codons). By using NNK for mutagenizing target regions, we can

126 screen a significantly smaller, manageable library with diverse capsid variants, optimizing the

127 discovery of improved activity without the need for an excessively large sequence pool [15].

128

129 **Methods**

130 Total RNA was extracted from brain tissue, primary neurons, and HEK293T cells, then

131 treated with DNase I to remove any DNA contamination. The RNA was reverse transcribed

132 to cDNA, which was amplified by PCR using specific primers. Barcode-containing PCR

133 products were purified, followed by an Illumina adapter PCR and a Nextera XT Index PCR,

134 with additional purification steps in between. The final PCR products were sequenced using

135    the Illumina NextSeq 500/550. Three AAV batches were produced. After production,

136    purification, and titration, those batches were treated with DNase I and Proteinase K. The

137    viral lysates then underwent two rounds of PCR to add Illumina-compatible sequences and

138    NexteraXT indexes, followed by purification. The indexed samples were sequenced on the

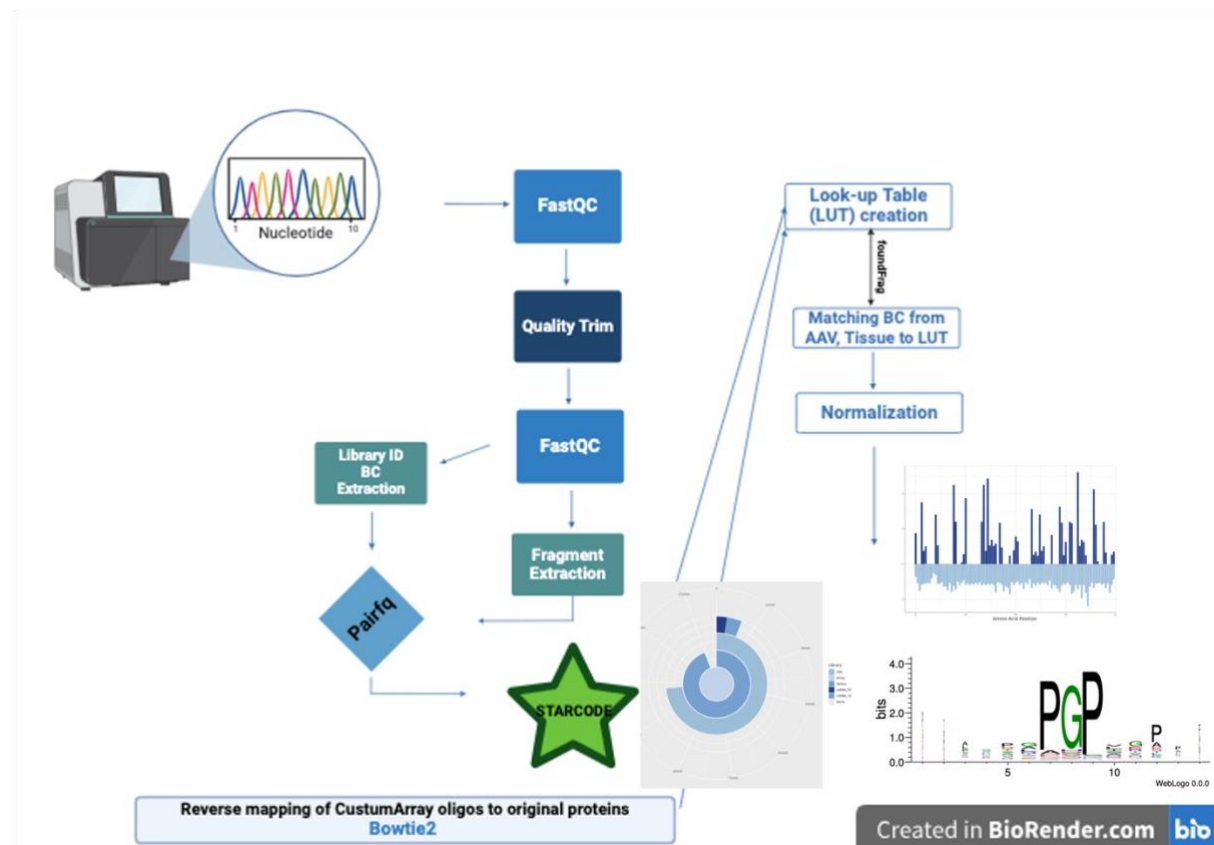139    Illumina NextSeq500/550.

140



141

**Fig. 1 Adeno-associated virus (AAV) vector genome and workflow structure.** A) The AAV
genome is linear, comprising single-stranded DNA (ssDNA) that is approximately 4.7 kilobases (kb)
in length. It is flanked by two inverted terminal repeats (ITRs), each 145 nucleotides long, located at
the ends of the genome. These ITRs flank two viral genes: rep, which is responsible for replication,
and cap, which encodes the structural capsid proteins. B) The BRAVE rational design process shows
how the fragment (21 nucleotides) and barcode (BC) are positioned closer together after plasmid
preparation. In brief, the oligonucleotide pool was assembled into an AAV production backbone
containing the cis-acting AAV2 Rep/Cap and CMV-GFP flanked by ITRs. Simultaneously, a 24-bp
random molecular BC was inserted into the 3′ UTR of the GFP gene. To create a look-up table linking

each BC to its corresponding peptide, a portion of the plasmid preparation was treated with Cre-recombinase *in vitro*. C) Workflow structure and data resources.

**Fastq files preprocessing**

Hierarchy of the files in this study were Plasmid Library (RAAV-60), Virus library (RAAV-30), Tissues mRNAs (RAAV-53 and RAAV-72). RAAV-60 consisted of three different library batches, in which they were p005, p006, and p007 (Fig 2). Each batch consisted of an Illumina sequencing dataset of 2x151 bp paired-end reads were utilized. An Illumina sequencing dataset of 1x151 bp paired-end reads were used for both virus library and tissue mRNAs. Each batch was quality control with FastQC (version 0.11.9-Java-11) before and after a trimming procedure named quality trim (BBMap, bbduk2.sh,



**Fig 2. Comprehensive Bioinformatics Tools and Data Analysis Workflow.**

An overview of the data processing pipeline, including FastQC analysis of fastq reads before and after quality trimming; extraction of library IDs, barcodes (BC), and fragments; pairing BC with fragments using pairfq; BC reduction via Starcode; mapping custom array oligos back to the original proteins

7

167    using Bowtie2; creation of a Look-Up Table for identifying fragments derived from the viral library

168    and infected tissue samples (foundFrag); normalization of all reads using various computational

169    methods; and subsequent data analysis to generate outputs like circular bar plots, pairwise plots, and

170    motif creation.

171

172    https://github.com/BioInfoTools/BBMap/blob/master/sh/bbduk2.sh) [16]. Given the analysis

173    is highly sensitive to sequence length, I performed mild trimming of nucleotides with quality

174    scores below 15.

175          Then, I also used bbduk2.sh for fragments, barcodes and library IDs extraction. I

176    extracted the sequences which were surrounded by the matching 13-16 (k=16, mink=13)

177    nucleotide sequences from the known reference for each plasmid. I allowed for 2 nucleotides

178    mismatches (hammingdistance=2, qhdist=2). The fragment was required to be 21 nucleotides

179    length (minlength=21 maxlength=21), barcode was 23-24 nucleotides length (minlength=23

180    maxlength=24), library ID was 3 nucleotides (bbduk2.sh and python script). Fragment

181    sequences were extracted from R2-fastq quality trimmed files. It was done when R2 fragment

182    was trimmed from the left side (lliteral) and from the right side (rliteral). Barcode sequences

183    were extracted from R1-fastq quality trimmed files. R1 was reverse complement to the

184    reference sequences. Therefore, all the reference sequences used for barcode extraction were

185    also reverse complement to the reference sequences. R1 was trimmed from the right side

186    (rliteral), then, and R1 was trimmed from the left side. Library IDs were extracted from the

187    R1-fastq quality trimmed files. Library ID was trimmed from the left side (lliteral). I

188    concatenated the fastq files of samples (fragments 5a with fragments 5b; barcodes 5a with

189    barcodes 5b; library IDs 5a with library IDs 5b) where more than one file from the same

190    sample was available (in which it applied to both plasmid library of p005 and p007).

191

192

193

**Sync paired-end fastq files**

After fragments, barcodes and library IDs extraction (by bbduk2.sh), these fastq files were

paired with one another using (**Pairfq** version 0.17.0). There were several methods in pairfq

and I used makepairs method to pair the forward and reverse reads and wrote the singletons

to separate files. The extracted fragments were paired with the extracted barcodes, then the

library IDs were paired with the extracted barcodes that had been paired once with the

extracted fragments. Following by that step, the paired barcodes that had been paired twice

were paired again with the extracted fragments that had once paired earlier. Hence after three

times of pairing, each fragment, barcode and library ID with similar identifiers were listed in

a table.


**Barcodes reduction**

The paired extracted barcodes further clustered into a "consensus" barcode using **starcode**

algorithm (starcode version 1.4) [17]. I set the starcode parameter for the clustering

algorithm that allowed the maximum distance between items in a cluster to 1 and the

algorithm should perform 5 repetitions or rounds of clustering. In essence, this workflow

clusters all peptide reads with highly similar barcodes into a single "multiple-alignment".

This method, when used restrictively, has been previously demonstrated to outperform other

reduction techniques by avoiding false clustering of barcodes [13].


**Sequence alignment**

The paired fragments were aligned using BLAST (version Blast: 2.15.0+,

https://anaconda.org/bioconda/blast) against to the reference fragments with gene names

information. The reference fragments were used to make the customized BLAST database. I

218    ran the BLAST nucleotide search tool and I set up several BLASTn parameters. I limit the

219    number of target sequences to 25, set the word size for the BLAST search to 11 and specified

220    the output format as comma-separated values (-outfmt 10). The BLAST output file was then

221    used for analyzing the percentage of alignment, the bitscores and mismatches. BLAST was

222    performed to determine the origin, and the oversampling of sequencing (using all reads with

223    the same barcode). The information was incorporated into the same table as paired fragments

224    and paired barcodes based on the Look-up Table number (LUTnr).

225

226    **Creating Look-up Table (LUT)**

227    I created the LUT using the R language (version 4.3.2)  and this R-based analysis framework

228    was a parallelized implementation of the MapReduce programming philosophy [18–20],

229    introduced by Google in 2004. The MapReduce model suggested to distribute and parallel

230    complex job process by splitting the job into multiple tasks through the use of map and

231    reduce stages [21].

232

233    **Short sequencing reads alignment to reference genes and SAM/BAM files indexing**

234    The purpose of using Bowtie2 (version Bowtie2: 2.5.4,

235    https://anaconda.org/bioconda/bowtie2) was to align the reference fragments with gene

236    names to the NCBI gene sequences which was collected in an excel file at lab. The Bowtie 2

237    index was made from a set of sequencing read files and output a set of alignments in SAM

238    format. SAMtools (version 1.18) was used to convert the SAM file to a BAM file and read

239    the alignment data from the files in sorted BAM file into a GAlignments object. These

240    objects represented genomics alignments. The GAlignments object included metadata

241    columns, such as mapping quality and CIGAR strings.

242

**Adeno-associated virus batches production (RAAV-30) using the same plasmid pool**

**(RAAV-60)**

RAAV-30 consisted of three different libraries of DNA viruses which derived from three plasmid libraries of RAAV-60. They were p005-AAV-01, p006-AAV-02 and p007-AAV-03. The final PCR products were sequenced from the 3'-end, producing a single 100 bp read that covered the barcode sequences. Each virus batch was checked for quality scores (FastQC) before and after the quality trim (bbduk2.sh). MultiQC (Version 1.14) showed that 1% of the raw reads were incorrect and had adapter contamination. I did quality trimming of raw reads expecting of Q>20. I had tried to use trimmomatic (version 0.39) to trim away the Ilumina adapters and low-quality nucleotides. Another FastQC was done to confirm the sequences quality. I then compared the trimmed result produced by trimmomatic and by bbduk2.sh, I chose to continue using the trimmed raw reads from bbduk2.sh as it had preserved the sequence length albeit more stringent quality scores. Next, I extracted barcodes from the quality trimmed fastq sequences using bbduk2.sh. I trimmed fastq sequences from the right side using a 13-16 nucleotides k-mer allowing for 1-2 mismatches. The sequences were reverse complement to the template, therefore, the k-mer sequences were also reverse complement. The fastq sequences from the left side were trimmed using a 16 nucleotides k-mer allowing for 1-2 mismatches. The sequences were reverse complement to the template, therefore, the k-mer sequences were also reverse complement. Library IDs were extracted with bbduk2.sh. Similarly, library IDs s were trimmed from the left side (lliteral). Then, the extracted library IDs and the extracted barcodes were paired together (pairfq) based on their shared identifiers. The paired barcodes were then clustered using starcode reduction.

**Barcodes and Library IDs extraction of tissue from DNA viruses (RAAV-53) and**

**(RAAV-72)**

RAAV-53 and RAAV-72 were mRNA sequences that consisted of barcodes and library IDs, which could be derived from three different plasmid library (RAAV-60: p005, p006 or p007). RAAV-53 had 20 samples and RAAV-72 had 8 samples with different reads depth. Similar content of scripts was used in RAAV-60 but modified to handle multiple samples for RAAV-53 and RAAV-72 in each directory. Fastqc data showed that some of the samples had lower quality scores. Hence a more stringent parameters were done using the bbduk2.sh quality trimming for every tissue samples. Each sample in RAAV-53 and RAAV-72 were extracted with specific sequence templates used for barcodes extraction in each plasmid library RAAV-60: p005, p006 and p007. A series of barcode reduction were done for each sample.

**Data analysis**

This workflow encompasses RNA counting, alignment, normalization, and data visualization across various RAAV datasets (RAAV-60, RAAV-30, RAAV-53, RAAV-72). It included executing RNA counting scripts and performing alignment, barcode, and fragment frequency analyses. Chimeric barcodes were aligned for RAAV-60 samples (p005, p006, p007), with fragment and barcode frequency tables generated. The pipeline calculated alignment percentages. Complete library range objects were created by combining fragments for selected RAAV samples. RNA counting data was saved as .rds files, optimized for large datasets, while library counts were normalized in smaller chunks with iterative read count updates. The workflow also included data visualization, with circular barplots illustrating unique peptides, pairwise and accumulated gene coverage plots. For pairwise gene plot, RNA counts were log2-transformed to stabilize variance. The data was divided into bins representing amino acid (AA) positions across the gene, with each bin containing information

292 on barcode counts (BC). This binning highlighted patterns along the AA sequence for

293 different experimental groups. Finally, the barcode count and corresponding peptide amino

294 acid sequences were input into the Hammock tool (v_1.2.0), [22]. Consensus motifs were

295 visualized using WebLogo (v 0.0.0), and output formatting was handled with GPL

296 Ghostscript 9.18 (2015-10-05) in Bioconda environment.

297

298 **Results**

299 The output from this workflow included a list of barcodes, their corresponding peptides

300 (nucleotide sequences translated into amino acids), and the position of each peptide within

301 the original protein. Each plasmid library was associated with specific library IDs: p005 used

302 (TCG), p006 used (TAC), and p007 used (CTA). Most of the results presented here were

303 derived from p007 library. For plasmid libraries p005 and p007, 98% of the sequences were

304 aligned with the original sequences established in the lab. Since plasmid library p006 was

305 designed using degenerate NNK, it was not subjected to BLAST analysis.

306     For plasmid library p005, a total of 113,150,190 reads were processed, identifying

307 9,202,839 unique original barcodes. After starcode reduction, 5,434,737 unique barcodes

308 remained. Out of these, 2,191,113 barcodes appeared only once, classified as single-read

309 barcodes. Clean multi-read barcodes, which appeared only once from multiple reads, totaled

310 38,162. There were also 3,205,462 chimeric multi-read barcodes, appearing more than once

311 from multiple reads.

312     Clean multi-read and Chimeric multi-read barcodes referred to different types of

313 barcodes identified in the sequence analysis based on how often they appeared and their

314 consistency across reads. Clean Multi-Read Barcodes represented multi-read barcodes where

315 all reads (sequences) for a given barcode were consistent. They indicated reliable, non-

316 chimeric sequences, where the reads aligned well without discrepancies. Chimeric Multi-
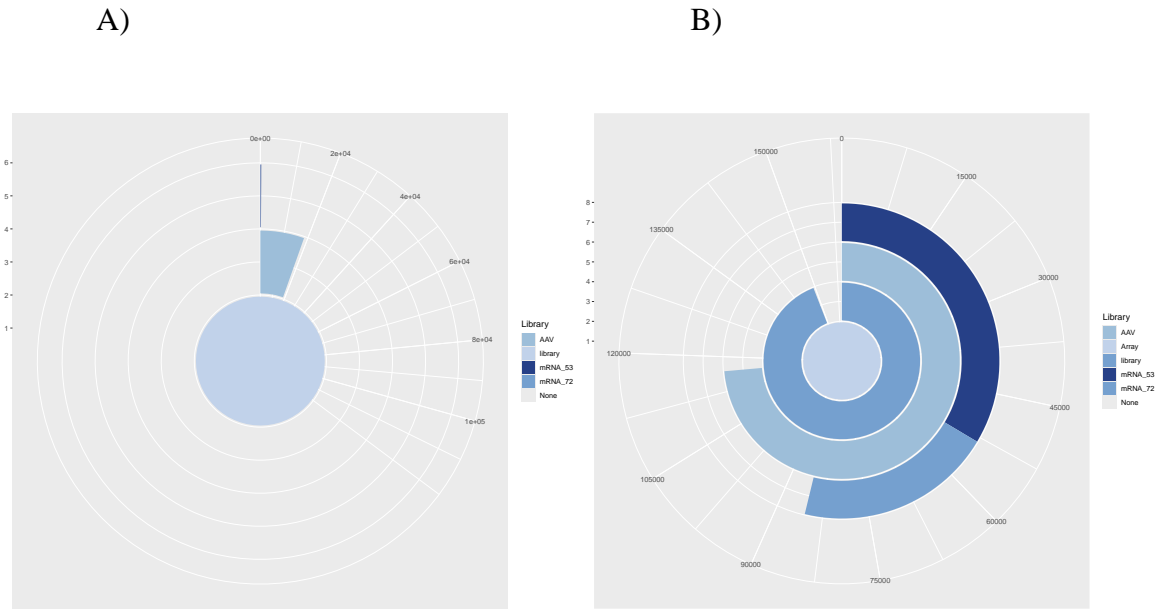
317    Read Barcodes appeared multiple times with different read sequences, meaning that the reads

318    associated with a barcode were not fully consistent. This chimerism suggested that these

319    barcodes might have combined sequences, often arising from sequencing or amplification

320    errors.

321        For plasmid library p006, a total of 20,251,305 reads were processed, resulting in the

322    detection of 1,261,887 unique original barcodes. After starcode reduction, this number

323    decreased to 592,798 unique barcodes. There were 306,521 single-read barcodes. Clean

324    multi-read barcodes amounted to 5,813, and 280,464 chimeric multi-read barcodes were

325    observed. For plasmid library p007, 123,875,291 reads were processed, resulting in the

326    identification of 6,747,020 distinct original barcodes. Following starcode reduction,

327    3,333,752 unique barcodes remained. Of these, 2,685,026 were single-read barcodes. Clean

328    multi-read barcodes, appearing only once, totaled 14,545, while chimeric multi-read barcodes

329    amounted to 634,181. Data processing steps similar to those applied to plasmid libraries were

330    also performed on viral libraries and tissue samples, though the results are not presented here.
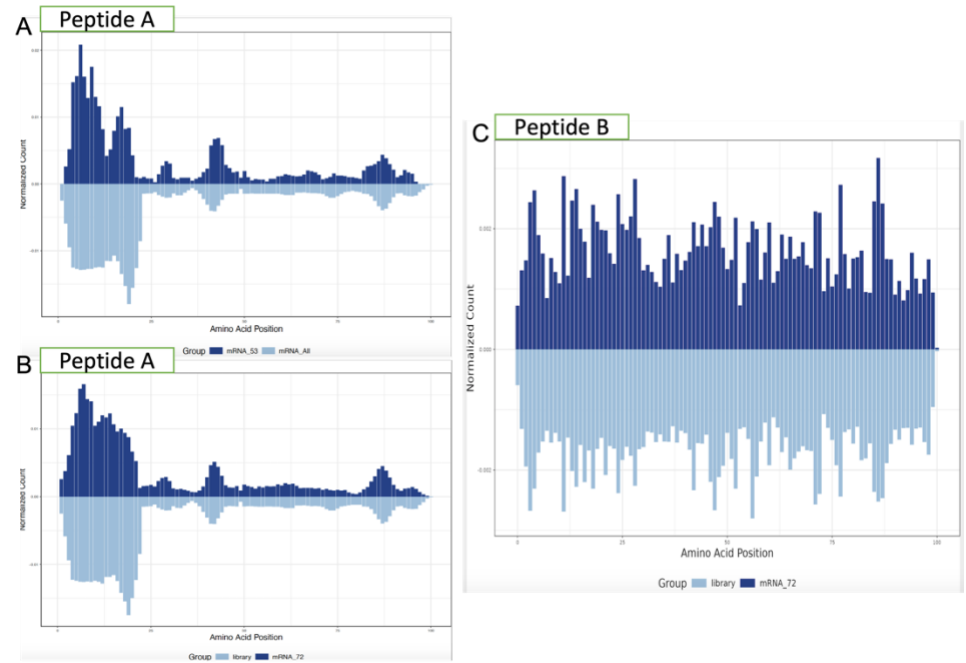
331        The efficiency of fragment matching and identification could be traced back to the

332    original 294 proteins using the lookup table (LUT) and consensus motifs identified through

333    the Hammock method (Fig. 2). The p006 library contained 340,909 unique peptides, with

334    18,517 from the virus library, 395 from tissue-infected samples (RAAV-53), and 42 from

335    another tissue-infected sample (RAAV-72) (Fig 3A). The array contained 158,788 entries,

336    created in the lab. Within the p007 library, 117,148 unique peptides were identified. The

337    virus library contained 116,881 unique peptides, with 53,294 detected in infected tissue

338    (RAAV-53) and 85,392 in infected tissue (RAAV-72), as shown in Fig. 3B.

339        From the *in vivo* screening of tissue-infected samples, the top 25 peptides were

340    identified from proteins associated with multiple barcodes and detected across different

341    samples (Fig 4). Peptides from the far-left region (highest peak) in peptide A were

342    hypothesized to exhibit both highly efficient retrograde transport and local infectivity in the

343    rat brain [3], making them valuable candidates for further testing (Fig 4A, 4B).

344        A)                                                    B)



345

346    **Fig 3. The circular bar plot displays the absolute quantities of unique peptides recovered at each**

347    **stage of the BRAVE assay.** A) p006 library. B) p007 library.

348



349

350    **Fig 4. Relative Amino Acid Frequency Distribution.** The height of each bar reflects the relative

351    frequency of a specific amino acid, adjusted according to the total library, and aggregated from all

352    peptides spanning that region.

353     The Hammock approach employed a hidden Markov model (HMM) to detect patterns across

354     the dataset, grouping peptides with similar sequence homology. The y-axis, measured in bits,

355     represented the information content at each position within the motif, indicating positional

356     conservation or variability. Higher information content denoted conserved positions where

357     specific bases or amino acids were preferred, while lower information content signified

358     greater variability. This analysis provided insights into the conservation patterns within the

359     motif and highlights key residues or bases that might be critical for function or recognition

360     (Fig 5).



361

362     **Fig 5.** Bitmap plots of consensus motifs generated using the Hammock clustering approach, which

363     utilizes a hidden Markov model (HMM) to align and identify recurring patterns across sequences. (A-

364     C) Consensus motifs of distinct peptides derived from different proteins. (D-F) Consensus motifs of

365     three distinct peptides from the same protein.

366

367     **Discussion**

368     During the data processing of the p005 and p007 libraries, the lookup table (LUT) required

369     for data analysis, including fragment identification and normalization, consumed significant

370     computer memory, making it challenging to process the entire dataset simultaneously. To

371     improve performance and scalability, chunking - dividing the data into smaller, manageable

372     pieces was applied. These chunks were created based on metadata containing GA alignment

16

373      objects, simplifying data management and processing. A dynamic chunking strategy was

374      implemented to accommodate varying data sizes [20]. This study utilized multiple

375      programming languages, including Bash, Python, and R, with R used for data analysis and

376      graph plotting. Some scripts incorporated the 'future.apply' and 'plan' packages in R to

377      enable parallel processing [21], https://msu-icer.github.io/r-for-hpcc/parallelizing-r-

378      code.html. When normalization was conducted, setting the future plan with workers = 8

379      (using multicore) succeeded, but increasing workers to 32 resulted in an out-of-memory error.

380      This likely occurred because, with workers = 8, memory usage was spread across 8

381      processes, staying within the capacity of the Lunarc cosmos system's RAM. However, when

382      workers = 32, the memory demands of 32 processes likely exceeded the system's RAM,

383      causing the error (https://future.futureverse.org/reference/plan.html).

384          Two approaches were used to normalize RNA barcode counts. The first method

385      involved processing smaller data chunks in a loop, updating the read count with each iteration

386      before moving to the next chunk. The second method processed smaller chunks, saving each

387      as a combined .rds file, and then proceeded with filtering out entries with an RNA count of 1

388      before performing normalization. Looping through chunks and processing smaller segments

389      could reduce memory usage while still providing normalized counts. Saving intermediate

390      results in files and filtering low-count entries before normalization, as in my second

391      approach, could also reduce data load while improving normalization accuracy for

392      downstream analysis.

393          When running Hammock, the error StringIndexOutOfBoundsException: String index

394      out of range: 15 occurred. However, after setting the parameter --greedy_threshold = 15, the

395      program executed smoothly. In the context of greedy clustering, this threshold generally

396      controlled the allowable similarity or distance between sequences in the same cluster. A

397      threshold of 15 likely specified that sequences within a cluster must meet a similarity

398   standard, where differences between them did not exceed 15 units. This could represent a

399   maximum allowed mismatch, an edit distance limit, or even a similarity cutoff of 85% (if

400   interpreted as 100 - 15). Setting this threshold defined how loosely or strictly sequences were

401   grouped, with lower values requiring closer similarity (more stringent clustering) and higher

402   values allowing for more variation within clusters [22]. Despite technical challenges, the

403   BRAVE strategy successfully identified peptide fragments associated with unique barcodes

404   in targeted infected tissue samples.

405        The efficient recovery of barcodes from virally expressed mRNA ensured that only

406   capsid variants successfully completing all critical steps of AAV infectivity. This readout

407   modality was crucial, as certain capsid engineering versions might result in capsid

408   sequestration on the cell surface, blocking internalization and causing misleading readouts.

409   The BRAVE approach effectively filtered out such variants. Additionally, extracting

410   barcodes from DNA in the injected library could provide insight into which capsid variants

411   became sequestered in tissue (either before or after internalization and uncoating) without

412   leading to successful transgene expression.

413        An added advantage of this method was the use of a stable bacterial library combined

414   with barcode oversampling, enabling efficient error control during *in vivo* screening. This

415   control was achieved through vector-linked peptide transcript identification on unique

416   barcode counts (rather than mRNA expression levels) and the number of samples showing

417   the same readout for each peptide.

418        The main limitation was the relatively small size and diversity of the libraries

419   compared to those generated by methods like capsid shuffling, degenerate primers, or error-

420   prone PCR. However, this limitation was balanced by the advantages of rational design,

421   which ensured that only potentially functional sequences, without repetition or bias, were

422   included in the library. Although this approach did not capture the full diversity of a

423     completely random method, the increasing capacity of synthesis (currently in the range of

424     10^6) and decreasing costs continue to expand its potential applications.

425

426     **Conclusion**

427     AAV vector-based therapies are rapidly emerging as a promising approach for treating or

428     even preventing neurological diseases. AAV gene therapy is particularly effective for well-

429     validated CNS targets, especially those genetic targets that are difficult to address with

430     conventional treatments. Clinical trials are currently assessing AAV vectors for delivering

431     therapeutic proteins to the CNS. One key advantage of AAV gene transfer is its durability in

432     non-dividing cells like neurons after a single administration, which is beneficial in scenarios

433     where repeat dosing may be challenging or impractical.

434           The most critical factor in unlocking the full potential of AAV gene therapy for

435     neurological diseases is effective delivery, largely determined by the capsid, which controls

436     the specific tissues and cell types transduced. Each administration route, along with the

437     chosen capsid, has unique characteristics and considerations, including dose requirements. It

438     is unlikely that a single administration route will suit all CNS applications; rather, the optimal

439     route for a given target and disease will depend on achieving the desired target distribution

440     while minimizing exposure to non-target tissues [2,3].

441           In summary, recent progress in harnessing the potential of AAV gene therapy for

442     neurological disorders has sparked significant interest and effort in this field. As

443     improvements in AAV capsid optimization, vector delivery, and transgene design advance,

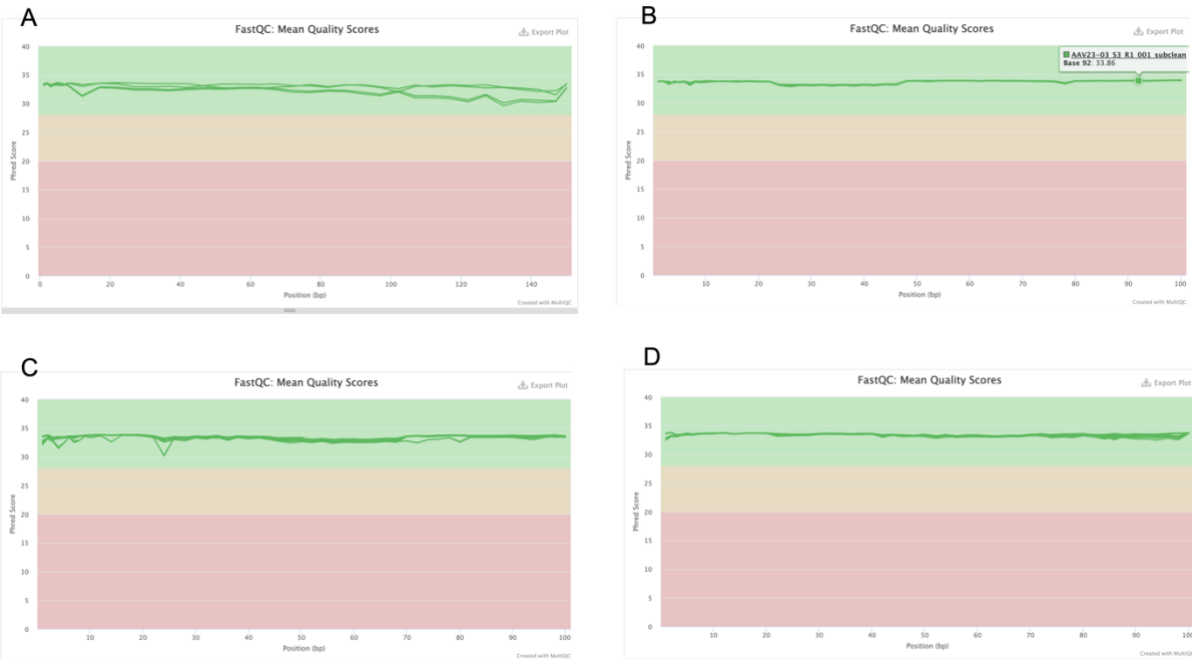444     the range of CNS disease applications and targets is expected to grow considerably.

445

446

447

452    **Reference**

453    1.  Hocquemiller M, Giersch L, Audrain M, Parker S, Cartier N. Adeno-Associated Virus-
454        Based Gene Therapy for CNS Diseases. Human Gene Therapy. 2016;27: 478–496.
455        doi:10.1089/hum.2016.087

456    2.  Deverman BE, Ravina BM, Bankiewicz KS, Paul SM, Sah DWY. Gene therapy for
457        neurological disorders: progress and prospects. Nat Rev Drug Discov. 2018;17: 641–659.
458        doi:10.1038/nrd.2018.110

459    3.  Davidsson M, Wang G, Aldrin-Kirk P, Cardoso T, Nolbrant S, Hartnor M, et al. A
460        systematic capsid evolution approach performed in vivo for the design of AAV vectors
461        with tailored properties and tropism. Proc Natl Acad Sci USA. 2019;116: 27053–27062.
462        doi:10.1073/pnas.1910061116

463    4.  McCarty DM, Monahan PE, Samulski RJ. Self-complementary recombinant adeno-
464        associated virus (scAAV) vectors promote efficient transduction independently of DNA
465        synthesis. Gene Ther. 2001;8: 1248–1254. doi:10.1038/sj.gt.3301514

466    5.  McCarty DM, Fu H, Monahan PE, Toulson CE, Naik P, Samulski RJ. Adeno-associated
467        virus terminal repeat (TR) mutant generates self-complementary vectors to overcome the
468        rate-limiting step to transduction in vivo. Gene Ther. 2003;10: 2112–2118.
469        doi:10.1038/sj.gt.3302134

470    6.  Srivastava A, Lusby EW, Berns KI. Nucleotide sequence and organization of the adeno-
471        associated virus 2 genome. J Virol. 1983;45: 555–564. doi:10.1128/jvi.45.2.555-
472        564.1983

473    7.  Lopez-Gordo E, Chamberlain K, Riyad JM, Kohlbrenner E, Weber T. Natural Adeno-
474        Associated Virus Serotypes and Engineered Adeno-Associated Virus Capsid Variants:
475        Tropism Differences and Mechanistic Insights. Viruses. 2024;16: 442.
476        doi:10.3390/v16030442

477    8.  Pupo A, Fernández A, Low SH, François A, Suárez-Amarán L, Samulski RJ. AAV
478        vectors: The Rubik's cube of human gene therapy. Molecular Therapy. 2022;30: 3515–
479        3541. doi:10.1016/j.ymthe.2022.09.015

480    9.  Girod A, Ried M, Wobus C, Lahm H, Leike K, Kleinschmidt J, et al. Genetic capsid
481        modifications allow efficient re-targeting of adeno-associated virus type 2. Nat Med.
482        1999;5: 1052–1056. doi:10.1038/12491

483   10. Bennett A, Mietzsch M, Agbandje-McKenna M. Understanding capsid assembly and
484        genome packaging for adeno-associated viruses. Future virology. 2017;12: 283.
485        doi:10.2217/fvl-2017-0011

486   11. Colón-Thillet R, Jerome KR, Stone D. Optimization of AAV vectors to target persistent
487        viral reservoirs. Virology Journal. 2021;18: 85. doi:10.1186/s12985-021-01555-7

488   12. Naik SH, Perié L, Swart E, Gerlach C, van Rooij N, de Boer RJ, et al. Diverse and
489        heritable lineage imprinting of early haematopoietic progenitors. Nature. 2013;496: 229–
490        232. doi:10.1038/nature12013

491   13. Davidsson M, Diaz-Fernandez P, Schwich OD, Torroba M, Wang G, Björklund T. A
492        novel process of viral vector barcoding and library preparation enables high-diversity
493        library generation and recombination-free paired-end sequencing. Sci Rep. 2016;6:
494        37563. doi:10.1038/srep37563

495   14. Blundell JR, Levy SF. Beyond genome sequencing: Lineage tracking with barcodes to
496        study the dynamics of evolution, infection, and cancer. Genomics. 2014;104: 417–430.
497        doi:10.1016/j.ygeno.2014.09.005

498   15. Siloto RMP, Weselake RJ. Site saturation mutagenesis: Methods and applications in
499        protein engineering. Biocatalysis and Agricultural Biotechnology. 2012;1: 181–189.
500        doi:10.1016/j.bcab.2012.03.010

501   16. Bushnell B. BBMap: A Fast, Accurate, Splice-Aware Aligner. 2014 [cited 16 Oct 2024].
502        Available: https://escholarship.org/uc/item/1h3515gn

503   17. Zorita E, Cuscó P, Filion GJ. Starcode: sequence clustering based on all-pairs search.
504        Bioinformatics. 2015;31: 1913–1919. doi:10.1093/bioinformatics/btv053

505   18. Dean J, Ghemawat S. MapReduce: simplified data processing on large clusters. Commun
506        ACM. 2008;51: 107–113. doi:10.1145/1327452.1327492

507   19. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The
508        Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA
509        sequencing data. Genome Res. 2010;20: 1297–1303. doi:10.1101/gr.107524.110

510   20. Kumar P, Paul RK, Roy HS, Yeasin Md, Ajit, Paul AK. Big Data Analysis in
511        Computational Biology and Bioinformatics. In: Mandal S, editor. Reverse Engineering of
512        Regulatory Networks. New York, NY: Springer US; 2024. pp. 181–197.
513        doi:10.1007/978-1-0716-3461-5_11

514   21. Bengtsson H. future.apply: Apply Function to Elements in Parallel using Futures. 2018.
515        p. 1.11.3. doi:10.32614/CRAN.package.future.apply

516   22. Krejci A, Hupp TR, Lexa M, Vojtesek B, Muller P. Hammock: a hidden Markov model-
517        based peptide clustering algorithm to identify protein-interaction consensus motifs in
518        large datasets. Bioinformatics. 2016;32: 9–16. doi:10.1093/bioinformatics/btv522

519

520

521   **Supplementary data**



522

523   **Supplementary Fig 1. MultiQC after quality trimming.** A) RAAV-60, B) RAAV-30, C) RAAV-

524   53, D) RAAV-72.

525

526

527

528