

Mapping of RNA modifications by direct Nanopore sequencing and JACUSA2

Amina Lemsara^{1,2}, Christoph Dieterich^{*1,2,3}, and Isabel Naarmann-de Vries^{1,2,3}

¹Klaus Tschira Institute for Integrative Computational Cardiology, University Heidelberg, 69120 Heidelberg, Germany

²Department of Internal Medicine III (Cardiology, Angiology, and Pneumology), University Hospital Heidelberg, 69120 Heidelberg, Germany

³German Centre for Cardiovascular Research (DZHK)-Partner Site Heidelberg/Mannheim, 69120 Heidelberg, Germany

Abstract

RNA modifications exist in all kingdom of life. Several different types of base or ribose modifications are now summarized under the term the "epitranscriptome". With the advent of high-throughput sequencing technologies much progress has been made in understanding RNA modification biology and how these modifications can influence many aspects of RNA life. The most widespread internal modification on mRNA is m6A, which has been implicated in physiological processes as well as disease pathogenesis. Here, we provide a workflow for the mapping of m6A sites using Nanopore direct RNA sequencing data. Our strategy employs pairwise comparison of base calling error profiles with JACUSA2. We outline a general strategy for RNA modification detection on mRNA and describe two specific use cases on m6A detection in detail. **Use case 1:** a sample of interest with modifications (e.g. "wild type" sample) is compared to a sample lacking a specific modification type (e.g. "knock out" sample, here *METTL3*-KO) or **Use case 2:** a sample of interest with modifications is compared to a sample lacking all modifications (e.g. *in vitro* transcribed cDNA). We provide a detailed protocol on experimental and computational aspects. Extensive online material provides a snakemake pipeline to identify m6A positions in mRNA and to validate the results against a miCLIP-derived m6A reference set. The general strategy is flexible and can be easily adapted by users in different application scenarios.

*Correspondence to: christoph.dieterich@uni-heidelberg.de

33 INTRODUCTION

34 Chemical modifications on DNA and histones, also known as epigenetics
35 marks, strongly impact gene expression during cell differentiation and in
36 several other biological programs. In the 1970s, it was recognized that RNA
37 is also subjected to extensive covalent modification, and studies in the late
38 1980s revealed the widespread deamination of bases (termed RNA editing),
39 which can lead to recoding if it occurs within coding sequences. Impres-
40 sive development in the RNA modification field occurred during the past
41 eight years, with the discovery of an extensive layer of base modifications
42 in mRNAs. These can influence gene expression and have been already
43 shown to be involved in primary cellular programs such as stem cell differ-
44 entiation, response to stress, and the circadian clock. The study of RNA
45 modifications and their effects is now referred to as epitranscriptomics, and
46 it reveals striking similarities to what is known for epigenomics. To date
47 thirteen distinct modifications have been identified on mRNA transcripts
48 [Anreiter et al., 2021]. These modifications are catalyzed by a variety of
49 dedicated enzymes and can be divided into two classes: modifications of
50 cap-adjacent nucleotides and internal modifications.

51 In contrast to the m7G cap, the impact of internal modifications on gene
52 regulation has been less studied apart from RNA editing, which is mediated
53 by RNA deaminases (e.g. the ADAR family). The most widespread in-
54 ternal mRNA modification is N6-methyladenosine (m6A). By modulating
55 the processing of mRNA, m6A can regulate a wide range of physiological
56 processes and its alteration has been linked to several diseases Roignant
57 and Soller [2017], Zaccara et al. [2019], Shi et al. [2019]. The modification is
58 catalyzed co-transcriptionally by a Mega-Dalton methyltransferase complex,
59 which includes the heterodimer METTL3-METTL14 and other associated
60 subunits Garcias Morales and Reyes [2021]. This modification is reversible
61 since two proteins of the AlkB-family of demethylases can remove m6A from
62 mRNA transcripts [Jia et al., 2011, Zheng et al., 2013]. In mammals, m6A
63 preferentially localizes within long internal exons and at the beginning of
64 terminal exons at so-called DRACH motif (D = A/G/U, R = A/G, H =
65 A/C/U) sites [Dominissini et al., 2012, Meyer et al., 2012, Ke et al., 2015].
66 Once deposited, m6A is recognized by several reader proteins that can af-
67 fect the fate of mRNA transcripts in nearly every step of the mRNA life
68 cycle, including alternative splicing [Adhikari et al., 2016, Roundtree et al.,
69 2017], mRNA translation [Wang et al., 2015] and decay [Wang et al., 2014,
70 Du et al., 2016, Roundtree et al., 2017]. The best-described readers are the
71 YTH domain family of proteins that decode the signal and mediate m6A
72 functions. By affecting RNA structure, m6A can also indirectly influence
73 the association of additional RNA-binding proteins (RBPs) and the assem-
74 bly of larger messenger ribonucleoprotein particles (mRNPs) [Patil et al.,
75 2018].

Several approaches have been presented to map RNA modifications on RNA. Herein, we focus on mRNA modification site detection in general and on m6A in particular where antibody-based protocols (miCLIP), methylation-sensitive restriction enzyme assays (MazF) or transgenic approaches (TRIBE, DART) have been presented to map m6A sites. All of the aforementioned approaches rely on high-throughput short read sequencing on the Illumina platform. This typically involves cDNA synthesis by reverse transcription and PCR-based library amplification. One recent addition to the toolbox of RNA modification mapping is direct RNA single molecule long read sequencing on the Oxford Nanopore Technologies platform (dRNA-seq). While our software is able to deal with Illumina and Nanopore-based approaches, the latter is the principal topic of this methods article.

MATERIALS

ONT direct RNA sequencing

This section summarizes all necessary consumables for direct RNA sequencing of poly-adenylated RNA (i.e. mRNA) on the MinION or similar device.

1. 500 ng polyA⁺ RNA isolated from total RNA e.g. with Oligotex mRNA kit (#70022, Qiagen) or Dynabeads oligo dT₂₅ beads (#61002, Thermo Fisher Scientific) or *in vitro* transcriptome sample. Store RNA at -80 °C and the mRNA purification kit as recommended by the manufacturer.
2. Nuclease-free water. Store at room temperature.
3. Direct RNA-sequencing kit (SQK-RNA002, Oxford Nanopore Technologies). Store at -20 °C.
4. NEBNext Quick Ligation Reaction Buffer (#B6058S, New England Biolabs). Store at -20 °C.
5. T4 DNA Ligase (#M0202S, New England Biolabs). Store at -20 °C.
6. dNTP Mix (10 mM each, #R0191, Thermo Fisher Scientific). Store at -20 °C.
7. SuperScript IV Reverse Transcriptase (#18090010, Thermo Fisher Scientific). Store at -20 °C.
8. Agencourt RNAClean XP beads (#A63987, Beckman Coulter). Store at 4 °C.
9. 70 % ethanol, freshly prepared.

- 110 10. Qubit dsDNA HS assay kit (#Q32854) and Qubit Fluorometer (Thermo
111 Fisher Scientific).
- 112 11. Flow cell priming kit (EXP-FLP002, Oxford Nanopore Technologies).
113 Store at -20 °C.
- 114 12. Thermocycler.
- 115 13. Gentle rotator mixer.
- 116 14. Magnetic stand for 1.5 ml tubes.
- 117 15. 1.5 ml DNA LoBind tubes (Eppendorf), 0.2 ml PCR tubes.
- 118 16. MinION or GridION sequencing device and MinION R9.4.1 Flow cells
119 (FLO-MIN106D, Oxford Nanopore Technologies). Store Flow cells at
120 4 °C.

121 **Preparation of an *in vitro* transcriptome sample**

- 122 1. 100 ng polyA⁺ RNA isolated from total RNA e.g. with Oligotex
123 mRNA kit (#70022, Qiagen) or Dynabeads oligo dT₂₅ beads (#61002,
124 Thermo Fisher Scientific). Store RNA at -80 °C and the mRNA pu-
125 rification kit as recommended by the manufacturer
- 126 2. 10 μM oligo(dT)-VN RT primer.
127 TTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTVN. Store at -20 °C.
- 128 3. 20 μM template switching oligo (TSO). ACTCTAATACGACTCAC-
129 TATAGGGAGAGGGCrGrG+G. Store at -20 °C.
- 130 4. 10 μM T7 extension primer. GCTCTAATACGACTCACTATAGG.
131 Store at -20 °C.
- 132 5. Nuclease-free water. Store at room temperature.
- 133 6. dNTP Mix (10 mM each, #R0191, Thermo Fisher Scientific). Store
134 at -20 °C.
- 135 7. Template Switching RT Enzyme Mix (#M0466S, New England Bio-
136 labs). Store at -20 °C.
- 137 8. Q5 Hot Start High-Fidelity 2X Master Mix (#M0494S, New England
138 Biolabs). Store at -20 °C.
- 139 9. RNase H (5,000 U/ml) (#M0297S, New England Biolabs). Store at
140 -20 °C.

- 141 10. NucleoSpin Gel and PCR Clean-up, Mini kit for gel extraction and
142 PCR clean up (#740609.50, Macherey-Nagel) or equivalent. Store at
143 room temperature.
- 144 11. MEGAscript T7 transcription kit (#AM1334, Thermo Fisher Scien-
145 tific). Store at -20 °C.
- 146 12. RNA Clean & Concentrator-25 kit (#R1017, Zymo Research). Store
147 at room temperature.
- 148 13. Thermocycler.
- 149 14. Table top centrifuge for 1.5 ml tubes.
- 150 15. Nanodrop spectrophotometer or equivalent.
- 151 16. 0.2 ml PCR tubes, 1.5 ml DNA LoBind tubes (Eppendorf).

152 Hardware requirements

153 All analyses have been performed/tested on two alternative hardware sys-
154 tems: a standard Linux desktop computer or an Apple iMac (Retina 5K,
155 ultimo 2014). The workflow requires a multi-core processor system with
156 minimal main memory of 16GB RAM and several GBs of free disk space
157 (depending on data set size).

158 Software dependencies and installation

159 Our analysis workflow has few requirements, which are detailed in Table 2.
160 Specifically, to execute our workflow, the following prerequisites are neces-
161 sary: a BASH shell, a JAVA runtime environment, a working PERL and
162 R installation. Additional i.e. non-standard software to process and map
163 Nanopore reads (bedtools, samtools and Minimap2) are obligatory, while
164 the installation of a Nanopore read simulator (NanoSim) is optional and de-
165 pends on your use case. Table ?? lists some additional R packages, which are
166 required to run the R code. Detailed instructions on how to setup are found
167 under https://github.com/dieterich-lab/MiMB_JACUSA2_chapter

168 METHODS

169 Our workflow is based on the pairwise comparison of samples with differ-
170 ent modification status (Figure 1). The sample of interest (yellow) may be
171 compared to different samples lacking certain modifications. If available,
172 the wild type (WT) sample can be compared to a knock out (KO) sample
173 lacking specific enzymatic activities (green), as outlined in Use Case 1. Al-
174 ternatively, a sample lacking all modifications may be used for comparison

labeling
of Ta-
bles in
PDF
doesn't
seem to
be cor-
rect

(blue). This may be either a simulated sample (i.e. with NanoSim) or an *in vitro* transcribed sample derived from cDNA. Such an analysis is detailed in Use Case 2. In any setting, JACUSA2 calculates scores for the Mismatch, Insertion and Deletion rates of the pairwise comparisons as outlined above (Figure 1, right).

One feature of Nanopore sequencing is to read sequences as 5-mers, as always five nucleotides are occupied by the pore protein (Figure 2). Because of this, a m6A modification may affect basecalling not only if the modified nucleotide is in the central position, but also at neighboring positions (-2 to +2). To account for this, JACUSA2 scores for Deletion, Mismatch and Insertion are calculated for the 5-mer context. Depending on the modification-specific signature, a Feature set can be selected to calculate the final JACUSA2 score (Figure 2).

Our workflow can be divided into a wet-lab part (Figure 3A) and a computational part (Figure 3B). Starting from total cellular RNA, polyA⁺ RNA is isolated and subjected to Nanopore direct RNA-sequencing. Guppy basecalling can be done as live basecalling during sequencing or after the sequencing run from generated FAST5 files, resulting in FASTQ output files (Figure 3A). FASTQ files are aligned to a reference sequence with Minimap2. SAMtools is used to generate BAM files as input for JACUSA2 analysis, which yields candidate m4A sites (Figure 3B).

Nanopore direct RNA sequencing

1. Adjust 500 ng polyA⁺ RNA to a total volume of 9 μ l with nuclease-free water. Complete RT adapter ligation reaction (in 0.2 ml PCR tube) with 3 μ l NEBNext Quick Ligation Reaction Buffer, 0.5 μ l RNA CS (RCS, from SQK-RNA002), 1 μ l RT-Adapter (RTA, from SQK-RNA002) and 1.5 μ l T4 DNA Ligase. Incubate 10 min at room temperature.
2. Prepare reverse transcription master mix on ice during ligation: 9 μ l nuclease-free water, 2 μ l 10 mM dNTPs, 8 μ l 5x SuperScript IV first strand buffer, 4 μ l 0.1 mM DTT.
3. Add the reverse transcription master mix to the ligation reaction and mix by pipetting. Add 2 μ l SuperScript IV reverse transcriptase and mix by pipetting. Incubate in a thermocycler with the following protocol: 50 min at 50 °C, 10 min at 70 °C, cool down to 4 °C.
4. Let the Agencourt RNAClean XP beads come to room temperature during reverse transcription. Carefully resuspend beads before use. Transfer reaction to a 1.5 ml DNA LoBind tube and mix with 72 μ l Agencourt RNAClean XP beads. Incubate 5 min at room temperature on a gentle rotator mixer.

- 215 5. Collect beads on a magnetic stand and remove supernatant. Wash
216 pelleted beads two times (30 sec) with 200 μ l freshly prepared 70 %
217 ethanol. Remove supernatant. Spin sample down and place on magnet
218 again. Remove any residual ethanol.
- 219 6. Resuspend beads in 20 μ l nuclease-free water by gentle flicking and
220 incubate 5 min at room temperature on a gentle rotator mixer. Collect
221 beads on a magnetic stand and transfer 20 μ l eluate in a fresh 1.5 ml
222 DNA LoBind tube.
- 223 7. For ligation of the RMX adapter, add the following to 20 μ l eluate: 8
224 μ l NEBNext Quick Ligation Reaction Buffer, 6 μ l RMX (from SQK-
225 RNA002), 3 μ l nuclease-free water, 3 μ l T4 DNA Ligase. Mix by
226 pipetting and incubate 10 min at room temperature.
- 227 8. Add 40 μ l carefully resuspended Agencourt RNAClean XP beads to
228 the reaction and mix by pipetting. Incubate 5 min at room tempera-
229 ture on a gentle rotator mixer.
- 230 9. Collect beads on a magnetic stand and remove supernatant. Wash
231 pelleted beads two times with 150 μ l wash buffer (WSB, from SQK-
232 RNA002). Resuspend beads by flicking, spin down and return to mag-
233 netic stand. Remove supernatant from pelleted beads.
- 234 10. Resuspend beads in 21 μ l elution buffer (EB, from SQK-RNA002) by
235 gentle flicking and incubate 5 min at room temperature on a gentle
236 rotator mixer. Pellet beads on a magnetic stand and transfer 21 μ l
237 eluate in a fresh 1.5 ml DNA LoBind tube.
- 238 11. Quantify 1 μ l of the library on a Qubit fluorometer with the Qubit
239 dsDNA HS kit according to the manufacturerers protocol. Concentra-
240 tion should be usually in the range of 5 - 10 ng/ μ l.
- 241 12. Insert MinION R9.4.1 Flow cell in the MinION or GridION sequenc-
242 ing device and perform Flow cell check in the MinKNOW software.
243 For successful sequencing of mammalian polyA⁺ RNA at least 1,000
244 available pores are recommended.
- 245 13. Prepare Priming Mix by adding 30 μ l flush tether (FLT, from EXP-
246 FLP002) to a vial of flush buffer (FB, from EXP-FLP002) and mix by
247 pipetting. Open priming port. Remove air bubble from priming port
248 by inserting the tip of a P1000 pipette into the priming port and slowly
249 dialing up, until a small volume of storage buffer enters the pipette
250 tip. Load 800 μ l Priming Mix via the priming port and carefully avoid
251 introduction of air bubbles. Close the priming port and wait for 5 min.

- 252 14. Mix 20 μ l library with 17.5 μ l nuclease-free water and 37.5 μ l RNA run-
253 ning buffer (RRB, from SQK-RNA002) and mix by pipetting. Open
254 the priming port and the sample port. Load 200 μ l Priming Mix via
255 the priming port. Mix library by pipetting just before loading and
256 load dropwise via the sample port. Carefully avoid introduction of air
257 bubbles. Close the sample port and the priming port.
- 258 15. Start sequencing for 48 to 72 h in the MinKNOW software. Choose
259 direct RNA-sequencing kit and high-accuracy basecalling as param-
260 eters.

261 Preparation of an *in vitro* transcriptome sample

262 The *in vitro* transcriptome sample is prepared based on a protocol published
263 by Zhang *et al.* Zhang et al. [2021] with some modifications.

- 264 1. Adjust 100 ng polyA⁺ RNA to a total volume of 6 μ l with nuclease-
265 free water. Add 1 μ l each of 10 μ M oligo(dT)-VN RT primer and 10
266 mM dNTPs. Mix by pipetting and incubate in a thermocycler: 5 min
267 at 75 °C, 2 min at 42 °C, cool to 4 °C.
- 268 2. Assemble 2.5 μ l 4x template switching RT buffer, 0.5 μ l 20 μ M TSO,
269 1 μ l 10x template switching RT enzyme mix and mix by pipetting.
270 Combine with 6 μ l RNA and incubate in a thermocycler: 90 min at
271 42 °C, 10 min at 68 °C, cool to 4 °C.
- 272 3. For Second strand synthesis add to First strand synthesis reaction: 50
273 μ l Q5 Hot Start High-Fidelity 2X Master Mix, 5 μ l RNase H, 2 μ l 10
274 μ M T7 extension primer, 33 μ l nuclease-free water. Mix by pipetting
275 and incubate in a thermocycler: 15 min at 37 °C, 1 min at 95 °C, 10
276 min at 65 °C, cool to 4 °C.
- 277 4. Purify double stranded cDNA with NucleoSpin Gel and PCR Clean-up
278 kit according to the manufacturerers protocol and elute in 20 μ l elution
279 buffer. Determine concentration on a Nanodrop spectrophotometer.
280 cDNA may be stored at -20 °C.
- 281 5. Combine 8 μ l cDNA for *in vitro* transcription with 2 μ l each of ATP,
282 GTP, CTP, UTP, 10x reaction buffer and enzyme mix from the MEGAscript
283 T7 transcription kit. Incubate 3 h at 37 °C.
- 284 6. Digest template DNA by addition of 1 μ l Turbo DNase. Mix by pipet-
285 ting and incubate 15 min at 37 °C.
- 286 7. Adjust reaction volume to 100 μ l with nuclease-free water and clean up
287 with RNA Clean & Concentrator-25 kit according to the manufactur-
288 ers protocol, using two volumes of adjusted RNA binding buffer (1:1

289 RNA binding buffer : ethanol). Elute RNA in 25 μ l nuclease-free wa-
 290 ter. Determine RNA concentration on a Nanodrop spectrophotometer.
 291 Store at -80 °C.

292 Nanopore read processing

293 1. Base call the ionic current signal stored in FAST5 files using Guppy.
 294 For the IVT sample, we applied real-time base calling with the MinKNOW-
 295 embedded Guppy basecaller. Otherwise, Guppy basecaller software
 296 can be used. In this case, the basecaller requires the path to FAST5
 297 files, the output folder, and the config file or the flowcell/kit combina-
 298 tion. The output are FASTQ files that can be compressed using the
 299 option "--compress_fastq".

```
300 $ guppy_basecaller --compress_fastq -i path_to_fast5 -s path_to_output
301 -c config_file.cfg --cpu_threads_per_caller 14 --num_callers
302 1
```

303 Set the number of threads "cpu_threads_per_caller" and the number
 304 of parallel basecallers "num_caller" according to your resources. Ad-
 305 ditional details can be found in Gup [2019].

306 2. Align reads to the transcriptome using Minimap2 software. The out-
 307 put is a SAM file that has to be converted to a compressed form as
 308 BAM file using SAMtools command. The alignment requires a ref-
 309 erence sequence. Here, we used GRCh38 Ensembl annotation and
 310 FASTA file release version 96. To reduce the indexing time of the
 311 human genome, save the index with the option "-d" before the map-
 312 ping and use the index instead of the reference file in the minimap2
 313 command line.

```
314 $ minimap2 -d reference.mmi reference.fa
```

315 To enable spliced alignments, use the setting "-ax splice -junc-bed
 316 annotation.bed -junc-bonus" where "-junc-bonus" allows to tune the
 317 bonus score and the BED file "-junc-bed annotation.bed" provides the
 318 splice junctions. The BED file can be generated using the following
 319 command:

```
320 $pafutils.js gff2bed annotation.gtf > annotation.bed
```

321 Use "-ub" to allow alignment to both strands or '-uf' to force the
 322 alignment to only forward strand. For Direct RNA Sequencing, it is
 323 recommended to set a small k-mer size "-k [=14]" to enhance sensitiv-
 324 ity. We recommend outputting primary alignments "--secondary=no".
 325 Use the parameter '-MD' to add the reference sequence information
 326 to the alignment; this is recommended for the downstream analysis.

327 Customize the number of threads "-t" according to your resources.
 328 Check Minimap2 manual for more details [Min].

```

329 $ minimap2 -t 5 --MD -ax splice --junc-bonus 1 -k14 --secondary=no
330 --junc-bed final_annotation_96.bed -ub reference.mmi Reads.fastq.gz
331 |samtools view -bS > mapping.bam

```

332 3. Map RNA modifications using JACUSA2 pipeline. JACUSA2 [Piechotta
 333 et al., 2021] rapidly detects RNA modifications based on a comparative
 334 strategy where the mapping features (mismatch, insertion and dele-
 335 tion) of a sample of interest are compared to a reference sequence (call-
 336 1) or against a sample without RNA modifications, e.g. a knock-out
 337 of an RNA modifying enzyme or an IVT (call-2). Moreover, it allows
 338 the integration of information from replicate experiments. **The output**
 339 **of JACUSA2 variant calling is a set of scores reflecting the read signa-**
 340 **tures involving mismatch, insertion and deletion. The analysis of read**
 341 **signature can be used for RNA modification detection. We integrate**
 342 **JACUSA2, in particular call-2 method, with the downstream analysis**
 343 **in one pipeline using the Python-based workflow management system**
 344 **Snakemake [Köster and Rahmann, 2012]. The Snakemake pipeline in-**
 345 **volves rules for the variant calling using JACUSA2 call-2, detection of**
 346 **RNA modification patterns, prediction of new modified sites and other**
 347 **intermediate rules as shown in Figure 4. The input of the pipeline are**
 348 **BAM files from paired conditions with different replicates. BAM files**
 349 **need to be sorted and may be subjected to many filters before being**
 350 **used by JACUSA2 call2 rule. Here, we suggest to filter out secondary**
 351 **and poor alignments. The output of JACUSA2 call2 is preprocessed**
 352 **(get_features) and subjected to a learning process to extract and visu-**
 353 **alize modification patterns (resp. get_pattern, visualize_pattern) and**
 354 **make predictions (predict_modification). "split_train_test" rule allows**
 355 **splitting input data into a training set and a test set. To use our**
 356 **snakemake-based JACUSA2 pipeline a set of parameters should be**
 357 **defined in the "config.yaml" file; mainly: the label of the analysis**
 358 **'label', the input bam files under 'data', the reference sequence 'refer-**
 359 **ence', a file containing size of chromosomes 'chr_size', JACUSA2 jar**
 360 **file 'jar', plus the path to inputs and outputs under 'path_inp' and**
 361 **'path_out' fields respectively. Further details on how to use JACUSA2**
 362 **pipeline is presented within the use cases in the next section. The**
 363 **pipeline could be executed on a high-performance-computing cluster**
 364 **(HPC) using the following command by specifying the number of cores**
 365 **to be used "-cores [=all]" and the rule name:**

```

366 $ srun snakemake --cores all rule_name

```

367 Check Snakemake documentation for more details [sna].

368 Use Case 1: Comparison of wild-type and knock-out samples

369 The JACUSA2 workflow detects RNA modifications using direct RNA se-
370 quencing by comparing a modified sample to an unmodified control sample.
371 Here, we used a published dataset of HEK293 cell lines to map m6A modifi-
372 cation [Pratanwanich et al., 2021]. The benchmark is composed of samples
373 sets two conditions: wild-type cells (WT, modified RNAs) and Mettl3 knock-
374 out cells (KO, unmodified RNAs) in two replicates (2 and 3). The FASTQ
375 files are mapped using Minimap2 as described in the previous section. The
376 following analysis is validated against m6A sites consistently reported in
377 three miCLIP-based studies Boulias et al. [2019], Koh et al. [2019], Körtel
378 et al. [2021] (Figure 5).

379 Starting with the preprocessed mapped reads as inputs (BAM files),
380 'HEK293T-WT-rep2.bam' and 'HEK293T-WT-rep3.bam' represent the wild-
381 type replicates and 'HEK293T-KO-rep2.bam' and 'HEK293T-KO-rep3.bam'
382 the control replicates,

- 383 1. Identify read error profile: **use "jacusa2_call2" rule to run JACUSA2**
384 in pairwise condition mode (call-2). The method requires BAM files of
385 the paired conditions and the corresponding library information "-P1"
386 and "-P2". In addition to the mismatch score, add "-D" and "-I" to
387 output the deletion and insertion scores. JACUSA2 allows filtering
388 reads according to many parameters. Here, we consider all sites with
389 base calling quality "-q [> 1]", mapping quality "-m [> 1]" and read
390 coverage "-c [> 4]". Furthermore, it provides a filter feature to improve
391 sensitivity. Here, **we consider filtering sites within homopolymer re-**
392 **gions "-a [=Y]". The output (named here, "Cond1vsCond2Call2.out")**
393 consists of a read error profile where the format is a combination
394 of BED6 with JACUSA2 call-2 specific columns and common info
395 columns: info, filter, and ref. Check JACUSA2 manual for more de-
396 tails on JACUSA2 filter and output options [JAC, 2021]. The number
397 of threads can be customized via the parameter "-p". **All parameters**
398 **related to the JACUSA2 method can be added under the field "ja-**
399 **cusa_params" in the config file by setting the name of the parameter**
400 **followed by the corresponding value [key: value]. Be aware to set all**
401 **parameters before running the pipeline.**

```
402 $ srun snakemake --cores all jacusa2_call2 $
```

- 403 2. Preprocess JACUSA2 output: from JACUSA2 call-2 output, **we select**
404 all sites within 5-mer of a central nucleotide 'A' flanked by 2 random
405 nucleotides (NNANN) and **we filter out sites of the homo-polymer re-**
406 **gions (JACUSA filter: Y). Then, we rebuild the tabular features such**
407 **that the observations are only sites with a reference base 'A'. Each**
408 **site is characterized by 15 features corresponding to the mismatch,**

409 insertion and deletion scores for the observed site and its two flank-
 410 ing positions from both sides. The rule "get_features" performs the
 411 preprocessing step. Use the parameter 'region' with a file containing
 412 target 5-mers to limit the analysis to specific sites. For comparison
 413 reasons, we consider common sites between use cases 1 and 2 . The
 414 output is an R object "features/features.rds", representing the matrix
 415 of Sites \times 15 features.

416 `$ srun snakemake --cores all get_features`

417 3. Extract m6A modification pattern: given the matrix of Sites \times Features,
 418 the next step is to learn a model representing the m6A modification
 419 pattern. To this end, the conventional non-negative matrix factor-
 420 ization (NMF) analysis is suggested [Lee and Seung, 1999]. Briefly,
 421 NMF factorizes a non-negative data matrix X (here: n sites and m
 422 features) into two non-negative matrices as $X \approx WH$, such that W
 423 is an $n \times k$ matrix containing basis vectors and H is an $k \times m$ ma-
 424 trix containing coefficient vectors. The coefficient vectors and their
 425 combination can be viewed as a pattern for m6A modification. The
 426 rank of factorization k is a critical parameter that affects the perfor-
 427 mance substantially. We suggest to select the rank k according to
 428 the method of Frigyesi and Höglund [2008] by looking at silhouette
 429 [Rousseeuw, 1987] and cophenetic correlation [Brunet et al., 2004] in-
 430 dices. Accordingly, the performance indices are computed for different
 431 choices of rank ($k < n, m$) and compared to the performance of a ran-
 432 dom permutation of the original data. Subsequently, the chosen rank
 433 corresponds to the value with the largest difference between slopes of
 434 the original and the randomized data. Here, the unsupervised pattern
 435 training is based on the consensus set of 1,905 m6A sites reported
 436 in the three miCLIP-based studies mentioned earlier. Based on the
 437 silhouette and cophenetic correlation indices, we identified an optimal
 438 factorization rank of 6 (Figure 6A). We then analyzed the identified
 439 patterns. According to the membership indicator of each site in ma-
 440 trix W , more than 80% of m6A modification sites can be represented
 441 by five patterns (Patterns 1,2,3,4,6) (Figure 6B). Interestingly, the
 442 linear combination of these five patterns in Figure 6C highlights the
 443 importance of position 3 and eventually the implication of all scores.

444 Using the JACUSA2 pipeline, run rule "get_pattern" to generate pat-
 445 terns and provide the set of modified sites as a ground truth under the
 446 field "modified_sites" in the config file. Here, the "miCLIP_union.bed"
 447 file contains the m6A sites from the three miCLIP-based studies. A
 448 miCLIP annotation, reflecting the consensus sites, is added to each
 449 site. A subset of modified sites can be used to generate patterns. Ac-
 450 cordingly, the "internal_pattern" field should refer to the annotation

Is this
the
reason
why you
chose
to work
on the
three
outputs
together
WT_IV, WT_KO,
KO_IVT

in Fig-
ure 6C
this is
labeled
sum

451 of selected sites from the "modified_sites" file. Plus, multiple combi-
452 nations of patterns can be defined and appended to the field "com-
453 bined_pattern" as new patterns. The corresponding outputs are under
454 "patterns" folder.

```
455 $ srun snakemake --cores all get_pattern
```

456 The produced patterns and their combinations can be visualized using
457 "visualize_pattern" rule. The corresponding outputs are under "pat-
458 tern/viz" folder.

```
459 $ srun snakemake --cores all visualize_pattern
```

460 4. Predict m6A modifications: the additive linear combination of the co-
461 efficient vectors (patterns) with the 15 features can be used to predict
462 m6A modification. We examine the ability of prediction on a subset of
463 data of more than 1,52 million sites with 17,021 miCLIP m6A sites.
464 We opt for the linear combination of the five most relevant patterns
465 described in step 3. The empirical Cumulative Distribution Function
466 (eCDF) of the inferred scores shows a significant difference between
467 the different miCLIP m6A categories (miCLIP annotation) and the
468 unmodified sites (Figure 6D). As the number of negative samples is
469 much larger than the number of positive samples, we particularly rec-
470 ommend investigating the Positive Predictive Value (PPV) of the pre-
471 dictions. Here, Figure 6E shows a moderate PPV that increases with
472 the cut-off.

473 To perform the prediction based on the selected patterns using the
474 JACUSA2 pipeline, run rule "predict_modification". The patterns
475 can be generated from a subset of the input data according to the
476 field "internal_pattern" or predefined patterns indicated in the "exter-
477 nal_pattern" field. The output is a BED file containing the estimated
478 scores as well as the corresponding eCDF and PPV plots. The corre-
479 sponding outputs are located under a new folder called "prediction".
480

```
481 $ srun snakemake --cores all predict_modification
```

482 Use Case 2: Comparison of wild-type and IVT samples

483 An alternative way to detect RNA modifications is to compare a modi-
484 fied sample to an *in-vitro* transcribed (IVT) control sample. Therefore,
485 we benchmark JACUSA2 on a sample set of two replicates (2 and 3) from
486 wild-type HEK293 cell lines (modified sample) Pratanwanich et al. [2021]
487 and a modification-free IVT sample from HEK293 cDNA (control sample)
488 (see "Preparation of an *in vitro* transcriptome sample"). The analysis steps

are similar to case 1. We evaluate the analysis against miCLIP m6A sites (Figure 5).

1. Identify read error profile: we use JACUSA2 call-2 with the same parameters as the previously described case. The input BAM files (HEK293T-WT-rep2.bam, HEK293T-WT-rep3.bam) and (HEK293T-IVT-rep1.bam, HEK293T-IVT-rep2.bam) are associated to the wild-type and IVT replicate samples respectively.

```
$ srun snakemake --cores all jacusa2_call2
```

2. Preprocess JACUSA2 output: we select all sites within the specific 5-mer (NNANN) and we consider the Y filter that excludes sites within homo-polymer regions. Then, we extract 5-mer features such that the selected sites are represented by the Mismatch, Deletion and Insertion scores for the observed site and its two flanking positions from both sides.

```
$ srun snakemake --cores all get_features
```

3. Extract m6A modification pattern: using NMF factorization, we extract patterns from the 1,905 sites reported as modified in the three miCLIP-based studies. Based on the silhouette and cophenetic correlation indices, we identified an optimal factorization rank of 6 (Figure 7A). We determined the predominant factors from matrix W . Accordingly, more than 80% of m6A modification sites can be represented by four patterns (Patterns: 1,2,3,6) (Figure 7B). In agreement with Use Case 1, the linear combination of the four patterns confirms the importance of position 3 and the implication of all scores as shown in Figure 7C.

```
$ srun snakemake --cores all get_pattern
```

4. Predict m6A modifications: we evaluate the prediction ability of the detected patterns on a test set of almost 1,52 million sites where 17,021 are miCLIP-m6A modified. We consider the linear combination of the four most relevant patterns (1,2,3,6). Figure 7D shows the eCDF of the inferred scores. The difference between the cumulative distribution of non miCLIP sites and miCLIP sites can be nicely observed, while the PPV plot shows a lower performance as compared to Use Case 1 (Figure 7E). The decrease in performance is likely explained by the absence of all modifications and not exclusively m6A in the control condition, which may induce noise to the score estimation by JACUSA2 call-2.

```
$ srun snakemake --cores all predict_modification
```

The first IVT run has rel. low coverage \rightarrow might this impact performance of UC2?

CD: to be confirmed

NOTES

Tips and Tricks

1. The reverse transcription step during library preparation is optional. However, we recommend to include this step to ensure proper sequencing also of RNAs with secondary structures. Superscript IV reverse transcriptase may be replaced by Superscript III reverse transcriptase, which is used in the protocol provided by Oxford Nanopore Technologies.
2. The library preparation protocol contains two bead clean up steps. It is important to remove ethanol and wash buffer completely. However, beads should not be dried for several minutes. Directly add water or elution buffer after washing to prevent sticking of the RNA to the beads.
3. The default filter in current MinKNOW versions is a Q score of 9. For direct RNA sequencing we recommend to adjust the output filter to a minimum Q score of 7, as in previous MinKNOW versions.
4. During preparation of the *in vitro* transcriptome sample, *in vitro* transcription and clean up kits may be replaced by equivalent products. The protocol however has been tested only with the mentioned kits.
5. Configuration of the pipeline should be handled via the config file. All parameters should be set before executing rules.
6. Once the pipeline has run successfully you should expect the following folders with the corresponding outputs in the output directory: bam, jacusa, features, patterns, and prediction.
7. JACUSA2 call2 could be run separately using the command line as described in JACUSA2 manual [JAC, 2021], then put the output under a new folder with the name 'jacusa' under the output directory.
8. In the snakemake pipeline, rules are linked so that the workflows are determined from top (e.g. predict_modification) to bottom (e.g. sort_bam) and executed accordingly from bottom to top (Figure 4). Therefore, running for example "predict_modification" rule leads to executing all rules on its pipeline.
9. Patterns could be generated from a subset of the input data that correspond to known modified sites. Alternatively, predefined patterns as a NMF R object could be used as a prediction model.

562 ACKNOWLEDGMENTS

563 The authors would like to thank Harald Wilhemit for testing the snakemake
564 pipeline. This work was supported by Informatics for Life funded by the
565 Klaus Tschira Foundation.

CD:
fund-
ing?

566 REFERENCES

- 567 Minimap2. <https://github.com/lh3/minimap2>. Accessed: 2022-01-19.
- 568 Snakemake. <https://snakemake.readthedocs.io>. Accessed: 2022-01-26.
- 569 Basecalling with guppy. [https://github.com/metagenomics/denbi-nanopore-training/blob/master/docs/basecalling/](https://github.com/metagenomics/denbi-nanopore-training/blob/master/docs/basecalling/basecalling.rst)
570 [basecalling.rst](https://github.com/metagenomics/denbi-nanopore-training/blob/master/docs/basecalling/basecalling.rst), 2019. Accessed: 2022-01-19.
- 571
- 572 Jacusa2 manual. <https://github.com/dieterich-lab/JACUSA2>, 2021.
573 Accessed: 2022-01-15.
- 574 Samir Adhikari, Wen Xiao, Yong-Liang Zhao, and Yun-Gui Yang. m(6)a:
575 Signaling for mrna splicing. *RNA biology*, 13:756–759, September 2016.
576 ISSN 1555-8584. doi: 10.1080/15476286.2016.1201628.
- 577 Ina Anreiter, Quoseena Mir, Jared T. Simpson, Sarath C. Janga, and
578 Matthias Soller. New twists in detecting mrna modification dynamics.
579 *Trends in biotechnology*, 39:72–89, January 2021. ISSN 1879-3096. doi:
580 10.1016/j.tibtech.2020.06.002.
- 581 Konstantinos Boulas, Diana Toczydlowska-Socha, Ben R Hawley, Noa
582 Liberman, Ken Takashima, Sara Zaccara, Théo Guez, Jean-Jacques
583 Vasseur, Françoise Debart, L Aravind, et al. Identification of the m6am
584 methyltransferase pcif1 reveals the location and functions of m6am in the
585 transcriptome. *Molecular cell*, 75(3):631–643, 2019.
- 586 Jean-Philippe Brunet, Pablo Tamayo, Todd R Golub, and Jill P Mesirov.
587 Metagenes and molecular pattern discovery using matrix factorization.
588 *Proceedings of the national academy of sciences*, 101(12):4164–4169, 2004.
- 589 Dan Dominissini, Sharon Moshitch-Moshkovitz, Schraga Schwartz, Mali
590 Salmon-Divon, Lior Ungar, Sivan Osenberg, Karen Cesarkas, Jasmine
591 Jacob-Hirsch, Ninette Amariglio, Martin Kupiec, Rotem Sorek, and
592 Gideon Rechavi. Topology of the human and mouse m6a rna methylomes
593 revealed by m6a-seq. *Nature*, 485:201–206, April 2012. ISSN 1476-4687.
594 doi: 10.1038/nature11112.
- 595 Hao Du, Ya Zhao, Jinqiu He, Yao Zhang, Hairui Xi, Mofang Liu, Jinbiao
596 Ma, and Ligang Wu. Ythdf2 destabilizes m 6 a-containing rna through

597 direct recruitment of the ccr4–not deadenylase complex. *Nature commu-*
598 *nications*, 7(1):1–11, 2016.

599 Attila Frigyesi and Mattias Höglund. Non-negative matrix factorization for
600 the analysis of complex gene expression data: identification of clinically
601 relevant tumor subtypes. *Cancer informatics*, 6:CIN–S606, 2008.

602 David Garcias Morales and José L. Reyes. A birds’-eye view of the activ-
603 ity and specificity of the mrna m, javax.xml.bind.jaxbelement@6d66739e,
604 a methyltransferase complex. *Wiley interdisciplinary reviews. RNA*, 12:
605 e1618, January 2021. ISSN 1757-7012. doi: 10.1002/wrna.1618.

606 Guifang Jia, Ye Fu, Xu Zhao, Qing Dai, Guanqun Zheng, Ying Yang,
607 Chengqi Yi, Tomas Lindahl, Tao Pan, Yun-Gui Yang, and Chuan He.
608 N6-methyladenosine in nuclear rna is a major substrate of the obesity-
609 associated fto. *Nature chemical biology*, 7:885–887, October 2011. ISSN
610 1552-4469. doi: 10.1038/nchembio.687.

611 Shengdong Ke, Endalkachew A. Alemu, Claudia Mertens, Emily Conn Gant-
612 man, John J. Fak, Aldo Mele, Bhagwattie Haripal, Ilana Zucker-Scharff,
613 Michael J. Moore, Christopher Y. Park, Cathrine Broberg Vågbø, Anna
614 Kussnierzcyk, Arne Klungland, James E. Darnell, and Robert B. Darnell.
615 A majority of m6a residues are in the last exons, allowing the potential
616 for 3’ utr regulation. *Genes & development*, 29:2037–2053, October 2015.
617 ISSN 1549-5477. doi: 10.1101/gad.269415.115.

618 Casslynn WQ Koh, Yeek Teck Goh, and WS Sho Goh. Atlas of quantitative
619 single-base-resolution n 6-methyl-adenine methylomes. *Nature communi-*
620 *cations*, 10(1):1–15, 2019.

621 Nadine Körtel, Cornelia Rücklé, You Zhou, Anke Busch, Peter Hoch-Kraft,
622 Reymond FX Sutandy, Jacob Haase, Mihika Pradhan, Michael Musheev,
623 Dirk Ostareck, et al. Deep and accurate detection of m6a rna modifications
624 using miclip2 and m6aboost machine learning. *bioRxiv*, pages 2020–12,
625 2021.

626 Johannes Köster and Sven Rahmann. Snakemake—a scalable bioinformatics
627 workflow engine. *Bioinformatics*, 28(19):2520–2522, 2012.

628 Daniel D Lee and H Sebastian Seung. Learning the parts of objects by
629 non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999.

630 Kate D. Meyer, Yogesh Saletore, Paul Zumbo, Olivier Elemento, Christo-
631 pher E. Mason, and Samie R. Jaffrey. Comprehensive analysis of mrna
632 methylation reveals enrichment in 3’ utrs and near stop codons. *Cell*, 149:
633 1635–1646, June 2012. ISSN 1097-4172. doi: 10.1016/j.cell.2012.05.003.

634 Deepak P Patil, Brian F Pickering, and Samie R Jaffrey. Reading m6a in
635 the transcriptome: m6a-binding proteins. *Trends in cell biology*, 28(2):
636 113–127, 2018.

637 Michael Piechotta, Qi Wang, Janine Altmüller, and Christoph Dieterich.
638 Rna modification mapping with jacusa2. *bioRxiv*, 2021.

639 Ploy N Pratanwanich, Fei Yao, Ying Chen, Casslynn WQ Koh, Yuk Kei
640 Wan, Christopher Hendra, Polly Poon, Yeek Teck Goh, Phoebe ML Yap,
641 Jing Yuan Chooi, et al. Identification of differential rna modifications
642 from nanopore direct rna sequencing with xpore. *Nature Biotechnology*,
643 39(11):1394–1402, 2021.

644 Jean-Yves Roignant and Matthias Soller. m,
645 javax.xml.bind.jaxbelement@8cec19d, a in mrna: An ancient mech-
646 anism for fine-tuning gene expression. *Trends in genetics : TIG*, 33:
647 380–390, June 2017. ISSN 0168-9525. doi: 10.1016/j.tig.2017.04.003.

648 Ian A. Roundtree, Molly E. Evans, Tao Pan, and Chuan He. Dynamic rna
649 modifications in gene expression regulation. *Cell*, 169:1187–1200, June
650 2017. ISSN 1097-4172. doi: 10.1016/j.cell.2017.05.045.

651 Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and
652 validation of cluster analysis. *Journal of computational and applied math-*
653 *ematics*, 20:53–65, 1987.

654 Hailing Shi, Jiangbo Wei, and Chuan He. Where, when, and how:
655 Context-dependent functions of rna methylation writers, readers, and
656 erasers. *Molecular cell*, 74:640–650, May 2019. ISSN 1097-4164. doi:
657 10.1016/j.molcel.2019.04.025.

658 Xiao Wang, Zhike Lu, Adrian Gomez, Gary C Hon, Yanan Yue, Dali Han,
659 Ye Fu, Marc Parisien, Qing Dai, Guifang Jia, et al. N 6-methyladenosine-
660 dependent regulation of messenger rna stability. *Nature*, 505(7481):117–
661 120, 2014.

662 Xiao Wang, Boxuan Simen Zhao, Ian A Roundtree, Zhike Lu, Dali Han,
663 Honghui Ma, Xiaocheng Weng, Kai Chen, Hailing Shi, and Chuan He.
664 N6-methyladenosine modulates messenger rna translation efficiency. *Cell*,
665 161(6):1388–1399, 2015.

666 Sara Zaccara, Ryan J. Ries, and Samie R. Jaffrey. Reading, writing and
667 erasing mrna methylation. *Nature reviews. Molecular cell biology*, 20:608–
668 624, October 2019. ISSN 1471-0080. doi: 10.1038/s41580-019-0168-5.

669 Zhang Zhang, Tao Chen, Hong-Xuan Chen, Ying-Yuan Xie, Li-Qian Chen,
670 Yu-Li Zhao, Biao-Di Liu, Lingmei Jin, Wutong Zhang, Chang Liu,

671 et al. Systematic calibration of epitranscriptomic maps using a synthetic
 672 modification-free rna library. *Nature Methods*, 18(10):1213–1222, 2021.

673 Guanqun Zheng, John Arne Dahl, Yamei Niu, Peter Fedorcsak, Chun-Min
 674 Huang, Charles J. Li, Cathrine B. Vågbø, Yue Shi, Wen-Ling Wang, Shu-
 675 Hui Song, Zhike Lu, Ralph P. G. Bosmans, Qing Dai, Ya-Juan Hao, Xin
 676 Yang, Wen-Ming Zhao, Wei-Min Tong, Xiu-Jie Wang, Florian Bogdan,
 677 Kari Furu, Ye Fu, Guifang Jia, Xu Zhao, Jun Liu, Hans E. Krokan, Arne
 678 Klungland, Yun-Gui Yang, and Chuan He. Alkbh5 is a mammalian rna
 679 demethylase that impacts rna metabolism and mouse fertility. *Molecular*
 680 *cell*, 49:18–29, January 2013. ISSN 1097-4164. doi: 10.1016/j.molcel.2012.
 681 10.015.

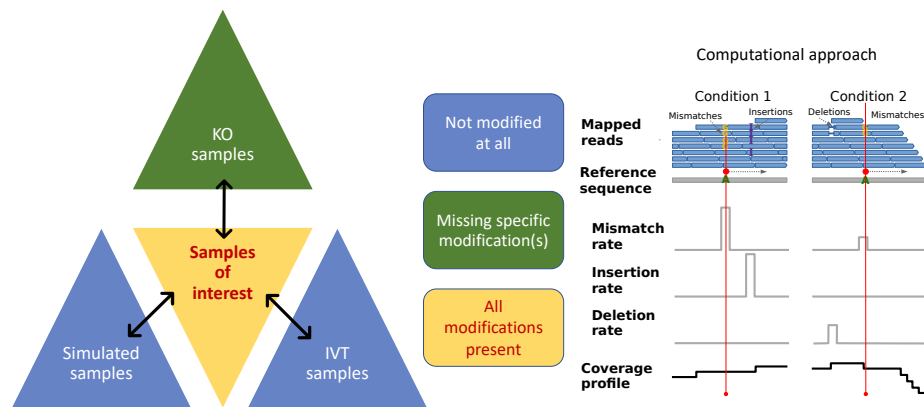


Figure 1: **General outline of RNA modification detection by JACUSA2.** A key feature of our approach is that multiple replicates can be compared as shown on the left. Samples of interests where all modifications are present could be compared with either KO samples where the modification of interest is missing or IVT/simulated samples where all modifications are absent. Read stacks (in blue) are compared head-to-head as shown on the right.

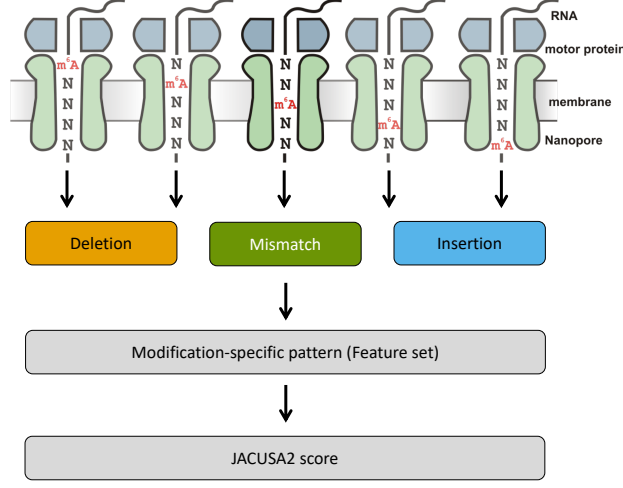


Figure 2: **Motivation of 5-mer context for RNA modification mapping.** The nanopore covers 5 consecutive RNA residues. That is why we consider a 5-mer context and derive 3 principal features for every position within a given 5-mer (15 features in total, with a central A residue in this example). We evaluate each feature set by previously learned patterns and compute a final score for modification site detection.

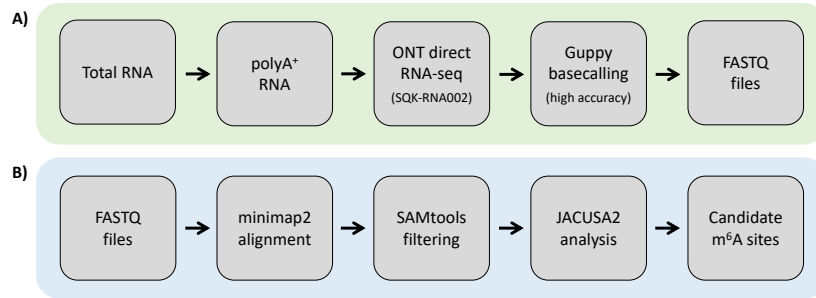


Figure 3: **Experimental and computational workflow.** A) Starting from total cellular RNA, polyA⁺ RNA is isolated and subjected to Nanopore direct RNA-sequencing. Guppy basecalling can be done as live basecalling during sequencing or after the sequencing run from generated FAST5 files, resulting in FASTQ output files. B) FASTQ files are aligned to a reference sequence with Minimap2. SAMtools is used to generate BAM files as input for JACUSA2 analysis, which yields candidate m⁴A sites.

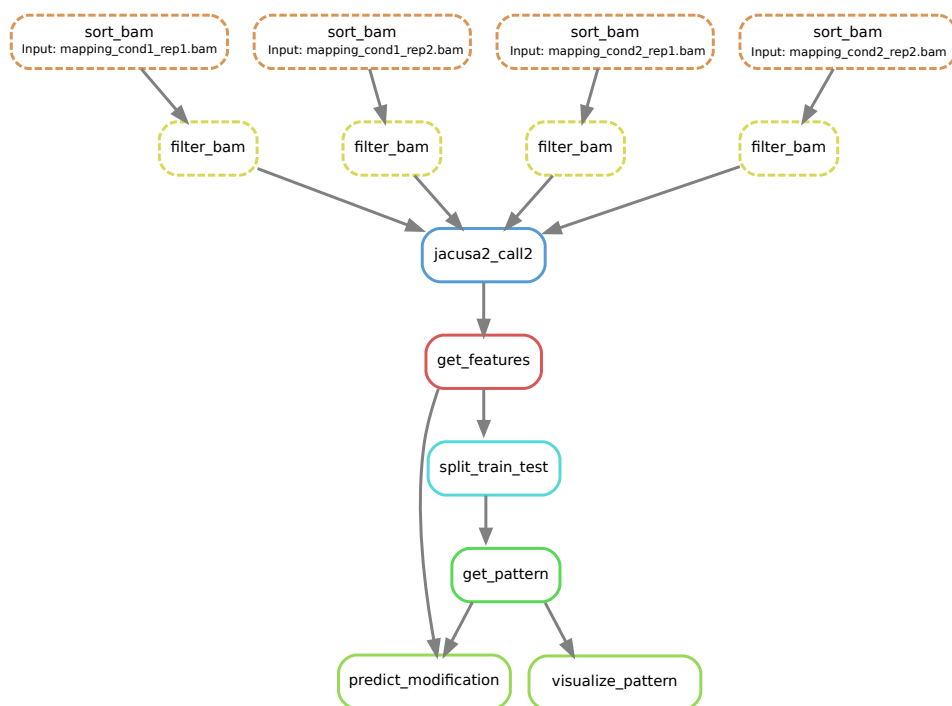


Figure 4: **Computational workflow.** Snakemake workflow for RNA modification detection based on JACUSA2 variant calling.

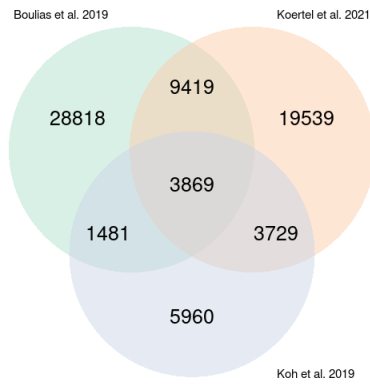


Figure 5: **m6A sites reported in the three miCLIP-based studies** Boulias et al. [2019], Koh et al. [2019] and Körtel et al. [2021].

Software	Version	Description
Minimap2	https://github.com/lh3/minimap2 v2.22 or later	https://lh3.github.io/minimap2/
samtools	https://github.com/samtools/samtools v1.12 or later	http://samtools.github.io/
JAVA	openjdk 11.0.12 2021-07-20 - JAVA 11 or later	OpenJDK Runtime Environment
R	https://www.r-project.org/ version 3.5.1 or later	The R Project for Statistical Computing
PERL	https://www.perl.org/ version 5.28.1 or later	Perl is a highly capable, feature-rich programming language
BASH, sed, awk	should be part of your Linux distribution	Misc.
bedtools	https://github.com/arq5x/bedtools2 version 2.29.2 or later	Perl is a highly capable, feature-rich programming language
NanoSim	https://github.com/bcgsc/NanoSim version 3.0.2 or later (optional)	NanoSim is a fast and scalable read simulator that captures the technology-specific features of ONT data

Table 1: **Software dependencies** blubba

683 TABLE CAPTIONS

684 TABLES

R Pack- ages	Version	Description
ggplot2	https://cran.r-project.org/web/packages/ggplot2/index.html - ggplot2_3.3.0 or later	ggplot2 is a system for declaratively creating graphics, based on The Grammar of Graphics.
NMF	https://cran.r-project.org/web/packages/NMF/index.html - NMF_0.22.0 or later	Provides a framework to perform Non-negative Matrix Factorization (NMF).

Table 2: **R Package dependencies** blubba

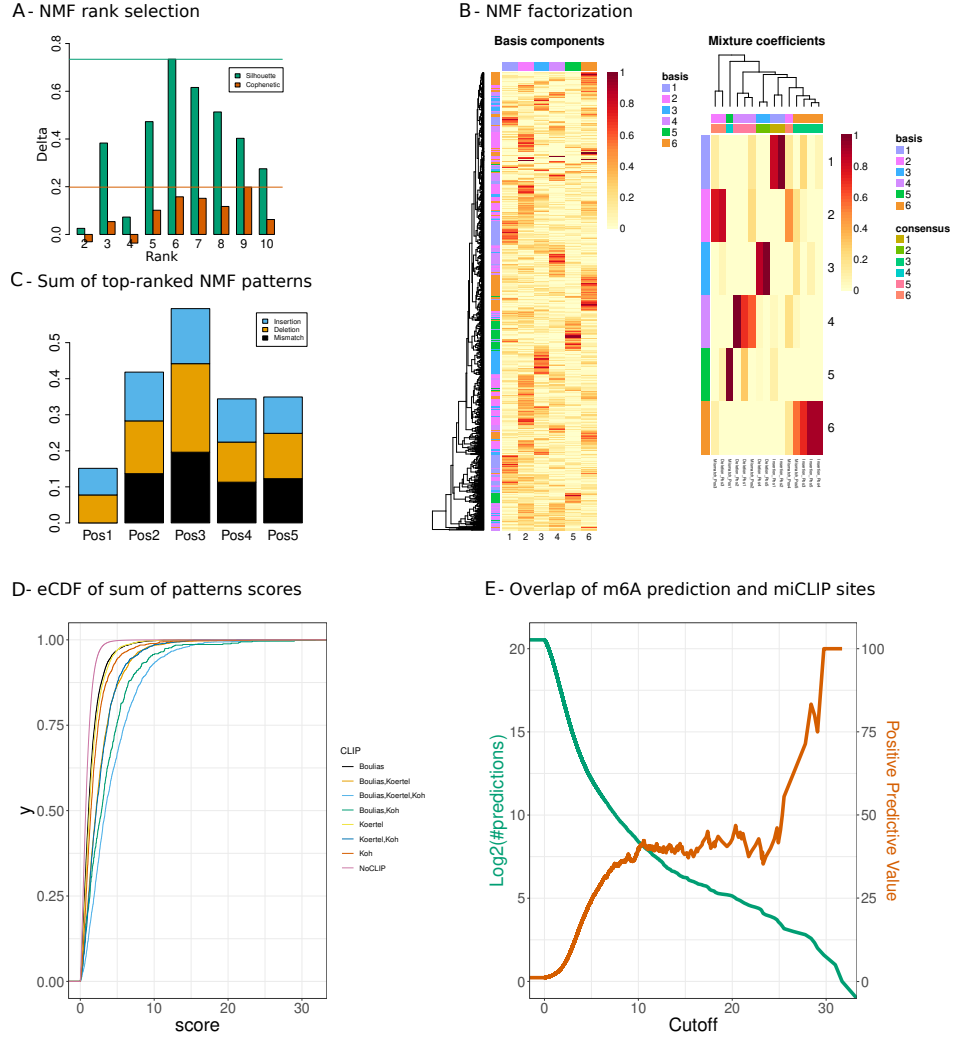


Figure 6: **Case 1. WT versus KO.** **A:** NMF rank selection. Barplots representing the difference of cophenetic correlation and silhouette indices between the NMF factorization of the original and the randomized data. Ranks with the largest value for both indices are determined, then the smallest rank is selected for the NMF decomposition. **B:** NMF result represented by the basis matrix W and the coefficient matrix H . The matrix H induces the RNA modification pattern. **C:** Barplots representing the linear combination of the top 5 patterns (y-axis) by the number of position in the specific 5-mer context (x-axis) and the score type: mismatch, deletion and insertion (resp. black, orange and blue). The 5 patterns (coefficient vectors: 1,2,3,4,6) are selected according to the predominant columns in matrix W . **D:** Score distribution inferred from the combined patterns, stratified by the different categories of miCLIP validated sites and non miCLIP sites. **E:** Number of predicted m6A sites (green) and Positive Predictive Value (PPV) of predicted m6A sites that overlap with miCLIP sites (orange).

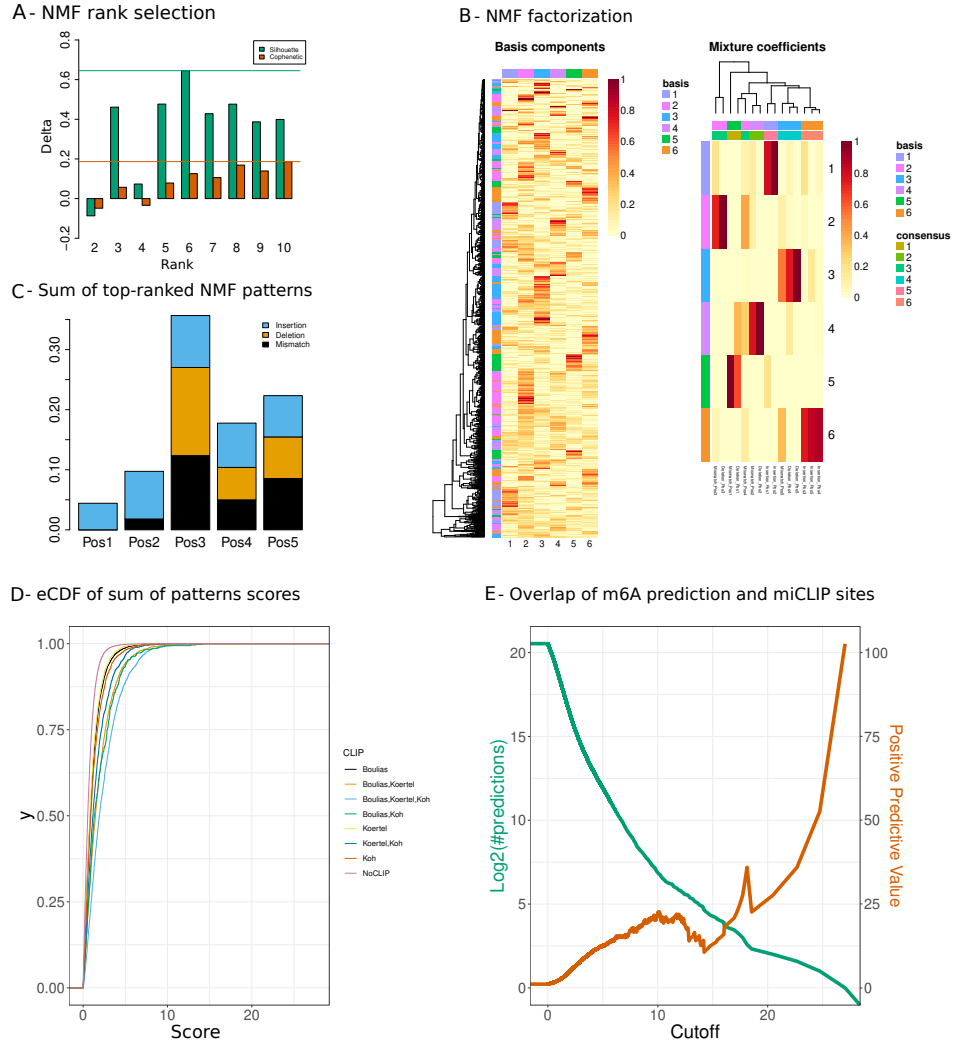


Figure 7: **Case 2. WT versus IVT.** **A:** NMF rank selection. Barplots representing the difference of cophenetic correlation and silhouette indices between the NMF factorization of the original and the randomized data. Ranks with the largest value for both indices are determined, then the smallest rank is selected for the NMF decomposition. **B:** NMF result represented by the basis matrix W and the coefficient matrix H . The matrix H induces the RNA modification pattern. **C:** Barplots representing the linear combination of the top 4 patterns (y-axis) by the number of position in the specific 5-mer context (x-axis) and the score type: mismatch, deletion and insertion (resp. black, orange and blue). The 4 patterns (coefficient vectors: 1,2,3,6) are selected according to the predominant columns in matrix W . **D:** Score distribution inferred from the combined patterns, stratified by the different categories of miCLIP validated sites and non miCLIP sites. **E:** Number of predicted m6A sites (green) and Positive Predictive Value (PPV) of predicted m6A sites that overlap with miCLIP sites (orange).