

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 μ M oligo(dT)-VN RT primer.
127 TTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTVN. Store at -20 °C.
- 128 3. 20 μ M template switching oligo (TSO). ACTCTAATACGACTCAC-
129 TATAGGGAGAGGGCrGrG+G. Store at -20 °C.
- 130 4. 10 μ M 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 1.
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. Ta-
164 ble 2 lists some additional R packages, which are required to run the R
165 code. Detailed installation instructions and corresponding workflow code
166 are deposited under https://github.com/dieterich-lab/MiMB_JACUSA2_
167 **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

(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 entire 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 well as live basecalling during sequencing on the respective FAST5 files, which results 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 m6A sites with the presented workflow in this chapter (Figure 3B). We will present all necessary experimental step for dRNA-seq in the next section.

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

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

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

263 Preparation of an *in vitro* transcriptome sample

264 The *in vitro* transcriptome sample is prepared based on a protocol published
265 by Zhang et al. [2021] with some modifications a detailed below. An *in vitro*
266 transcriptome lacks any RNA modifications and is a perfect reference sample
267 for RNA modification mining.

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

290 7. Adjust reaction volume to 100 μ l with nuclease-free water and clean up
 291 with RNA Clean & Concentrator-25 kit according to the manufactur-
 292 ers protocol, using two volumes of adjusted RNA binding buffer (1:1
 293 RNA binding buffer : ethanol). Elute RNA in 25 μ l nuclease-free wa-
 294 ter. Determine RNA concentration on a Nanodrop spectrophotometer.
 295 Store at -80 °C.

296 Nanopore read processing

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

```
304 $ guppy_basecaller --compress_fastq -i path_to_fast5 -s path_to_output
305 -c config_file.cfg --cpu_threads_per_caller 14 --num_callers
306 1
```

307 Set the number of threads "`cpu_threads_per_caller`" and the number
 308 of parallel basecallers "`num_caller`" according to your resources. Ad-
 309 ditional details can be found at <https://nanoporetech.com/>.

310 2. Align reads to the transcriptome using Minimap2 software. The out-
 311 put is a SAM file that has to be converted to a compressed form as
 312 BAM file using SAMtools command. The alignment requires a refer-
 313 ence sequence. Here, we used GRCh38 Ensembl release 96 annotation
 314 and FASTA file. Pre-indexing of the human genome saves time dur-
 315 ing read alignment. Please save the index with the option "`-d`" before
 316 read mapping and use the index instead of the reference file in the
 317 minimap2 command line.

```
318 $ minimap2 -d reference.mmi reference.fa
```

319 For Direct RNA Sequencing, it is recommended to set a small k-mer
 320 size "`-k [=14]`" to enhance sensitivity. We recommend outputting only
 321 primary alignments "`--secondary=no`". Use the parameter '`-MD`' to
 322 add the reference sequence information to the alignment; this is neces-
 323 sary for JACUSA2 downstream analysis. Adjust the number of threads
 324 "`-t`" according to your resources. Check Minimap2 manual for more
 325 details [Min]. To enable spliced alignments, use the setting `-ax splice`
 326 `-junc-bed annotation.bed -junc-bonus` where "`-junc-bonus`" allows to
 327 tune the bonus score and the BED file "`-junc-bed annotation.bed`"
 328 provides the splice junctions.

```

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 The BED file can be generated from EnsEMBL GTF files using the
 333 following command:

```

334 $paftools.js gff2bed annotation.gtf > annotation.bed

```

335 3. Mapping RNA modifications using JACUSA2 pipeline: JACUSA2
 336 [Piechotta et al., 2021] rapidly detects RNA modifications based on
 337 a comparative strategy where read alignment features (mismatch, in-
 338 sersion and deletion) of samples of interest are compared to a reference
 339 sequence (call-1) or against reference samples without the correspond-
 340 ing RNA modification of interest (call-2). JACUSA2 processes repli-
 341 cate experiments. The analysis of read alignment signatures is used
 342 for RNA modification detection. Particularly, we integrate JACUSA2
 343 call-2 method with the downstream analysis in one workflow using
 344 the Snakemake workflow manager [Köster and Rahmann, 2012]. Our
 345 Snakemake workflow encompasses several steps as shown in Figure
 346 4. The workflow requires BAM files from 2 conditions as input. We
 347 suggest to filter secondary and poor alignments beforehand. The out-
 348 put of JACUSA2 call2 is preprocessed (get_features) and subjected
 349 to a machine learning step to extract and visualize modification pat-
 350 terns (resp. get_pattern, visualize_pattern) and make predictions (pre-
 351 dict_modification). "split_train_test" rule allows splitting input data
 352 into a training set and a test set. To use our snakemake-based JA-
 353 CUSA2 pipeline a set of parameters should be defined in the "con-
 354 fig.yaml" file; mainly: the label of the analysis 'label', the input bam
 355 files under 'data', the reference sequence 'reference', a file containing
 356 size of chromosomes 'chr_size', JACUSA2 jar file 'jar', plus the path to
 357 inputs and outputs under 'path_inp' and 'path_out' fields respectively.
 358 We typically execute the workflow on a multi-core CPU system using
 359 the following command by specifying the number of cores to be used
 360 "-cores [=all]" and the rule name:

```

361 $ snakemake --cores all rule_name

```

362 Please consult the Snakemake documentation for further details (see
 363 <https://snakemake.readthedocs.io/en/stable/>).

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

365 The JACUSA2 workflow detects RNA modifications using direct RNA se-
 366 quencing by comparing modified samples to unmodified control samples.

Here, we used a published dataset of HEK293 cell lines to map m6A modification [Pratanwanich et al., 2021]. Our examples encompasses two conditions: wild-type RNA (WT, modified RNAs) and RNA from *METTL3* knockout cells (KO, m6A modification is absent). We use two replicates per condition (see <https://doi.org/10.5281/zenodo.5913452>). The FASTQ files are mapped using Minimap2 as described in the previous section. The following analysis is validated against m6A sites consistently reported in three miCLIP-based studies Boulias et al. [2019], Koh et al. [2019], Körtel et al. [2021] (Figure 5).

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

1. Compute read error profile with the `jacusa2.call2` rule:

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

The method requires BAM files of the paired conditions and the corresponding library information "-P1" and "-P2". In addition to the mismatch score, add "-D" and "-I" to output the deletion and insertion scores. JACUSA2 allows filtering reads according to many parameters. Here, we consider all sites with base calling quality "-q [> 1]", mapping quality "-m [> 1]" and read coverage "-c [> 4]". Here, we consider filtering sites within homopolymer regions "-a [=Y]". The output (named here, "Cond1vsCond2Call2.out") consists of a read error profile where the format is a combination of BED6 with JACUSA2 call-2 specific columns and common info columns: info, filter, and ref. Check JACUSA2 manual for more details on JACUSA2 filter and output options [JAC, 2021]. The number of threads can be customized via the parameter "-p". All parameters related to the JACUSA2 method can be added under the field "jacusa_params" in the config file by setting the name of the parameter followed by the corresponding value [key: value]. Be aware to set all parameters before running the pipeline.

2. Process JACUSA2 output with the `get_features` rule:

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

we select all sites within 5-mer of a central nucleotide 'A' flanked by 2 random nucleotides (NNANN) and we filter out sites of the homopolymer regions. Then, we rebuild the tabular features such that the observations are only sites with a reference base 'A'. Each site is characterized by 15 features corresponding to the mismatch, insertion and deletion scores for the observed site and its two flanking positions from both sides. The rule "get_features" performs the preprocessing

step. Use the parameter 'region' with a file containing target 5-mers to limit the analysis to specific sites. The output is an R object "features/features.rds", representing the matrix of Sites \times 15 features.

3. Extract characteristic m6A modification patterns with the `get_pattern` rule:

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

We learn a model representing the m6A modification patterns given the matrix of Sites \times Features. To this end, we employ non-negative matrix factorization (NMF)[Lee and Seung, 1999]. Briefly, NMF factorizes a non-negative data matrix X (here: n sites and m features) into two non-negative matrices as $X \approx WH$, such that W is an $n \times k$ matrix containing basis vectors and H is an $k \times m$ matrix containing coefficient vectors. The coefficient vectors and their combination can be viewed as a pattern for m6A modification. The rank of factorization k is a critical parameter that affects the performance substantially. We suggest to select the rank k according to the method of Frigyesi and Höglund [2008] by looking at silhouette [Rousseeuw, 1987] and cophenetic correlation [Brunet et al., 2004] indices. Accordingly, the performance indices are computed for different choices of rank ($k < n, m$) and compared to the performance of a random permutation of the original data. Subsequently, the chosen rank corresponds to the value with the largest difference between slopes of the original and the randomized data. Here, the unsupervised pattern training is based on the consensus set of 1,905 m6A sites reported in the three miCLIP-based studies mentioned earlier. Based on the silhouette and cophenetic correlation indices, we identified an optimal factorization rank of 6 (Figure 6A). We then analyzed the identified patterns. According to the membership indicator of each site in matrix W , more than 80% of m6A modification sites can be represented by five patterns (Patterns 1,2,3,4,6) (Figure 6B). Interestingly, the linear combination of these five patterns in Figure 6C highlights the importance of position 3.

Multiple patterns and their combinations can be visualized using `visualize_pattern` rule. The corresponding outputs are under "pattern/viz" folder.

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

4. Predict m6A modifications with the `predict_modification` rule:

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

This rule uses patterns of 15 features to predict m6A modification. We examine the ability of prediction on a subset of data of more than

1.52 million sites including 17,021 miCLIP m6A sites. We opt for the linear combination of the five most relevant patterns described in step 3. The empirical Cumulative Distribution Function (eCDF) of the inferred scores shows a significant difference between the different miCLIP m6A categories (miCLIP annotation) and the unmodified sites (Figure 6D). As the number of negative samples is much larger than the number of positive samples, we consider the Positive Predictive Value (PPV, $TP/(TP+FP)$) of our predictions. Here, Figure 6E shows that PPV increases with the score cut-off. The final output is a BED file containing the estimated scores as well as the corresponding eCDF and PPV plots. The corresponding outputs are located under a new folder called "prediction".

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

An alternative way to detect RNA modifications is to compare a modified sample to an *in-vitro* transcribed (IVT) control sample. Therefore, we benchmark JACUSA2 on a sample set of two replicates (2 and 3) from wild-type HEK293 cell lines (modified sample) Pratanwanich et al. [2021] and a modification-free IVT sample from HEK293 cDNA (control sample) (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

The first IVT run had rel. low coverage -> might this impact performance of UC2?

483 7A). We determined the predominant factors from matrix W . Accord-
484 ingly, more than 80% of m6A modification sites can be represented by
485 four patterns (Patterns: 1,2,3,6) (Figure 7B). In agreement with Use
486 Case 1, the linear combination of the four patterns confirms the im-
487 portance of position 3 and the implication of all scores as shown in
488 Figure 7C.

489 `$ srun snakemake --cores all get_pattern`

490 4. Predict m6A modifications: we evaluate the prediction ability of the
491 detected patterns on a test set of almost 1,52 million sites where
492 17,021 are miCLIP-m6A modified. We consider the linear combina-
493 tion of the four most relevant patterns (1,2,3,6). Figure 7D shows the
494 eCDF of the inferred scores. The difference between the cumulative
495 distribution of non miCLIP sites and miCLIP sites can be nicely ob-
496 served, while the PPV plot shows a lower performance as compared
497 to Use Case 1 (Figure 7E). The decrease in performance is likely ex-
498 plained by the absence of all modifications and not exclusively m6A in
499 the control condition, which may induce noise to the score estimation
500 by JACUSA2 call-2 .

501 `$ srun snakemake --cores all predict_modification`

CD:to
be con-
firmed

502 NOTES

503 Tips and Tricks

- 504 1. The reverse transcription step during library preparation is optional.
505 However, we recommend to include this step to ensure proper sequenc-
506 ing also of RNAs with secondary structures. Superscript IV reverse
507 transcriptase may be replaced by Superscript III reverse transcriptase,
508 which is used in the protocol provided by Oxford Nanopore Technolo-
509 gies.
- 510 2. The library preparation protocol contains two bead clean up steps. It
511 is important to remove ethanol and wash buffer completely. However,
512 beads should not be dried for several minutes. Directly add water
513 or elution buffer after washing to prevent sticking of the RNA to the
514 beads.
- 515 3. The default filter in current MinKNOW versions is a Q score of 9. For
516 direct RNA sequencing we recommend to adjust the output filter to a
517 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.

ACKNOWLEDGMENTS

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

REFERENCES

- Minimap2. <https://github.com/lh3/minimap2>. Accessed: 2022-01-19.
- Jacusa2 manual. <https://github.com/dieterich-lab/JACUSA2>, 2021. Accessed: 2022-01-15.
- Samir Adhikari, Wen Xiao, Yong-Liang Zhao, and Yun-Gui Yang. m(6)a: Signaling for mrna splicing. *RNA biology*, 13:756–759, September 2016. ISSN 1555-8584. doi: 10.1080/15476286.2016.1201628.
- Ina Anreiter, Quoseena Mir, Jared T. Simpson, Sarath C. Janga, and Matthias Soller. New twists in detecting mrna modification dynamics. *Trends in biotechnology*, 39:72–89, January 2021. ISSN 1879-3096. doi: 10.1016/j.tibtech.2020.06.002.

CD:
fund-
ing?

552 Konstantinos Boulas, Diana Toczyłowska-Socha, Ben R Hawley, Noa
553 Liberman, Ken Takashima, Sara Zaccara, Théo Guez, Jean-Jacques
554 Vasseur, Françoise Debart, L Aravind, et al. Identification of the m6am
555 methyltransferase pcif1 reveals the location and functions of m6am in the
556 transcriptome. *Molecular cell*, 75(3):631–643, 2019.

557 Jean-Philippe Brunet, Pablo Tamayo, Todd R Golub, and Jill P Mesirov.
558 Metagenes and molecular pattern discovery using matrix factorization.
559 *Proceedings of the national academy of sciences*, 101(12):4164–4169, 2004.

560 Dan Dominissini, Sharon Moshitch-Moshkovitz, Schraga Schwartz, Mali
561 Salmon-Divon, Lior Ungar, Sivan Osenberg, Karen Cesarkas, Jasmine
562 Jacob-Hirsch, Ninette Amariglio, Martin Kupiec, Rotem Sorek, and
563 Gideon Rechavi. Topology of the human and mouse m6a rna methylomes
564 revealed by m6a-seq. *Nature*, 485:201–206, April 2012. ISSN 1476-4687.
565 doi: 10.1038/nature11112.

566 Hao Du, Ya Zhao, Jinqiu He, Yao Zhang, Hairui Xi, Mofang Liu, Jinbiao
567 Ma, and Ligang Wu. Ythdf2 destabilizes m 6 a-containing rna through
568 direct recruitment of the ccr4–not deadenylase complex. *Nature commu-*
569 *nications*, 7(1):1–11, 2016.

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

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

577 Guifang Jia, Ye Fu, Xu Zhao, Qing Dai, Guanqun Zheng, Ying Yang,
578 Chengqi Yi, Tomas Lindahl, Tao Pan, Yun-Gui Yang, and Chuan He.
579 N6-methyladenosine in nuclear rna is a major substrate of the obesity-
580 associated fto. *Nature chemical biology*, 7:885–887, October 2011. ISSN
581 1552-4469. doi: 10.1038/nchembio.687.

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

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

592 Nadine Körtel, Cornelia Rücklé, You Zhou, Anke Busch, Peter Hoch-Kraft,
593 Reymond FX Sutandy, Jacob Haase, Mihika Pradhan, Michael Musheev,
594 Dirk Ostareck, et al. Deep and accurate detection of m6a rna modifications
595 using miclip2 and m6aboost machine learning. *bioRxiv*, pages 2020–12,
596 2021.

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

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

601 Kate D. Meyer, Yogesh Saletore, Paul Zumbo, Olivier Elemento, Christo-
602 pher E. Mason, and Samie R. Jaffrey. Comprehensive analysis of mrna
603 methylation reveals enrichment in 3’ utrs and near stop codons. *Cell*, 149:
604 1635–1646, June 2012. ISSN 1097-4172. doi: 10.1016/j.cell.2012.05.003.

605 Deepak P Patil, Brian F Pickering, and Samie R Jaffrey. Reading m6a in
606 the transcriptome: m6a-binding proteins. *Trends in cell biology*, 28(2):
607 113–127, 2018.

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

610 Ploy N Pratanwanich, Fei Yao, Ying Chen, Casslynn WQ Koh, Yuk Kei
611 Wan, Christopher Hendra, Polly Poon, Yeek Teck Goh, Phoebe ML Yap,
612 Jing Yuan Chooi, et al. Identification of differential rna modifications
613 from nanopore direct rna sequencing with xpore. *Nature Biotechnology*,
614 39(11):1394–1402, 2021.

615 Jean-Yves Roignant and Matthias Soller. m,
616 javax.xml.bind.jaxbelement@8cec19d, a in mrna: An ancient mech-
617 anism for fine-tuning gene expression. *Trends in genetics : TIG*, 33:
618 380–390, June 2017. ISSN 0168-9525. doi: 10.1016/j.tig.2017.04.003.

619 Ian A. Roundtree, Molly E. Evans, Tao Pan, and Chuan He. Dynamic rna
620 modifications in gene expression regulation. *Cell*, 169:1187–1200, June
621 2017. ISSN 1097-4172. doi: 10.1016/j.cell.2017.05.045.

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

625 Hailing Shi, Jiangbo Wei, and Chuan He. Where, when, and how:
626 Context-dependent functions of rna methylation writers, readers, and
627 erasers. *Molecular cell*, 74:640–650, May 2019. ISSN 1097-4164. doi:
628 10.1016/j.molcel.2019.04.025.

629 Xiao Wang, Zhike Lu, Adrian Gomez, Gary C Hon, Yanan Yue, Dali Han,
630 Ye Fu, Marc Parisien, Qing Dai, Guifang Jia, et al. N 6-methyladenosine-
631 dependent regulation of messenger rna stability. *Nature*, 505(7481):117–
632 120, 2014.

633 Xiao Wang, Boxuan Simen Zhao, Ian A Roundtree, Zhike Lu, Dali Han,
634 Honghui Ma, Xiaocheng Weng, Kai Chen, Hailing Shi, and Chuan He.
635 N6-methyladenosine modulates messenger rna translation efficiency. *Cell*,
636 161(6):1388–1399, 2015.

637 Sara Zaccara, Ryan J. Ries, and Samie R. Jaffrey. Reading, writing and
638 erasing mrna methylation. *Nature reviews. Molecular cell biology*, 20:608–
639 624, October 2019. ISSN 1471-0080. doi: 10.1038/s41580-019-0168-5.

640 Zhang Zhang, Tao Chen, Hong-Xuan Chen, Ying-Yuan Xie, Li-Qian Chen,
641 Yu-Li Zhao, Biao-Di Liu, Lingmei Jin, Wutong Zhang, Chang Liu,
642 et al. Systematic calibration of epitranscriptomic maps using a synthetic
643 modification-free rna library. *Nature Methods*, 18(10):1213–1222, 2021.

644 Guanqun Zheng, John Arne Dahl, Yamei Niu, Peter Fedorcsak, Chun-Min
645 Huang, Charles J. Li, Cathrine B. Vågbø, Yue Shi, Wen-Ling Wang, Shu-
646 Hui Song, Zhike Lu, Ralph P. G. Bosmans, Qing Dai, Ya-Juan Hao, Xin
647 Yang, Wen-Ming Zhao, Wei-Min Tong, Xiu-Jie Wang, Florian Bogdan,
648 Kari Furu, Ye Fu, Guifang Jia, Xu Zhao, Jun Liu, Hans E. Krokan, Arne
649 Klungland, Yun-Gui Yang, and Chuan He. Alkbh5 is a mammalian rna
650 demethylase that impacts rna metabolism and mouse fertility. *Molecular
651 cell*, 49:18–29, January 2013. ISSN 1097-4164. doi: 10.1016/j.molcel.2012.
652 10.015.

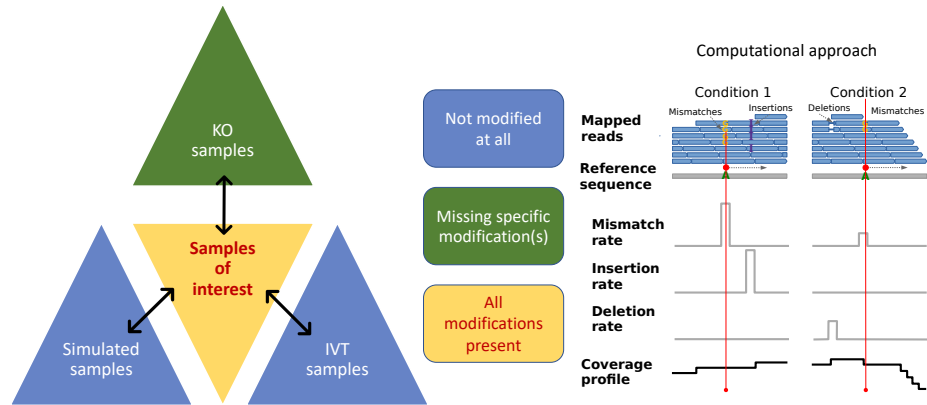


Figure 1: **General outline of RNA modification detection by JACUSA2.** A key feature of our approach is that multiple replicates can be compared as shown on the left. Samples of interests where all modifications are present could be compared with either KO samples where the modification of interest 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.

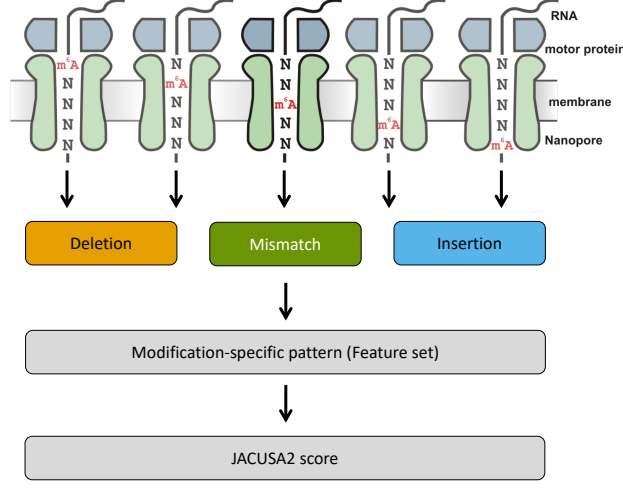


Figure 2: **Motivation of 5-mer context for RNA modification mapping.** The nanopore covers 5 consecutive RNA residues. That is why we consider a 5-mer context and derive 3 principal features for every position within a given 5-mer (15 features in total, with a central A residue in this example). We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

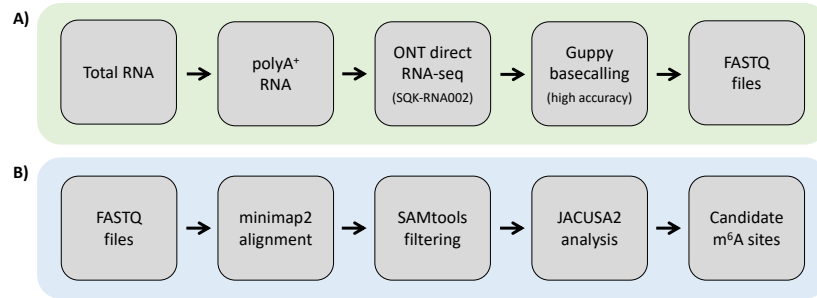


Figure 3: **Experimental and computational workflow.** A) Starting from total cellular RNA, polyA⁺ RNA is isolated and subjected to Nanopore direct RNA-sequencing. Guppy basecalling can be done as live basecalling during sequencing or after the sequencing run from generated FAST5 files, resulting in 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 m⁶A sites.

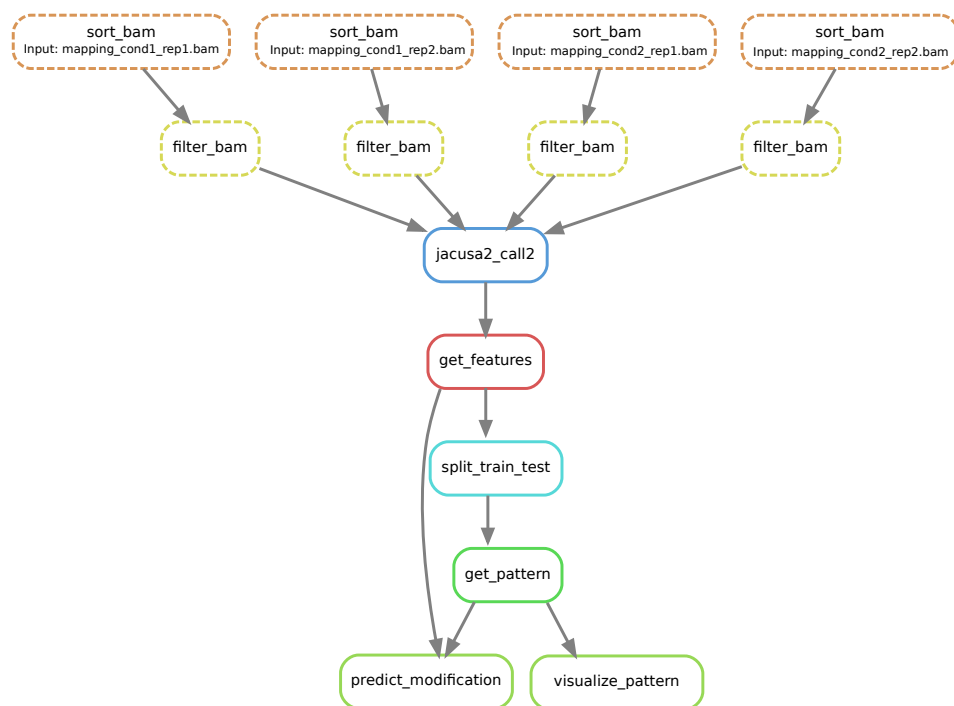


Figure 4: **Computational workflow.** Snakemake workflow for RNA modification detection based on JACUSA2 variant calling.

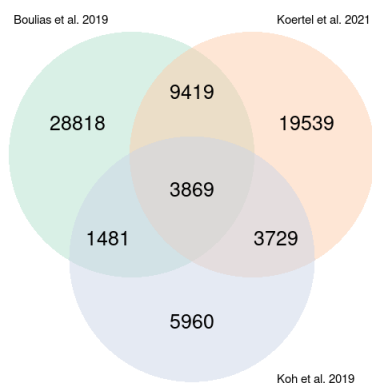


Figure 5: **m6A sites reported in the three miCLIP-based studies** Boulias et al. [2019], Koh et al. [2019] and Körtel et al. [2021].

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	https://openjdk.java.net/ 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
bedtools	https://github.com/arq5x/bedtools2 version 2.29.2 or later	Perl is a highly capable, feature-rich programming language
snakemake	https://snakemake.github.io/ version 6.8.1 or later	The Snakemake workflow management system

Table 1: **Software dependencies**

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**

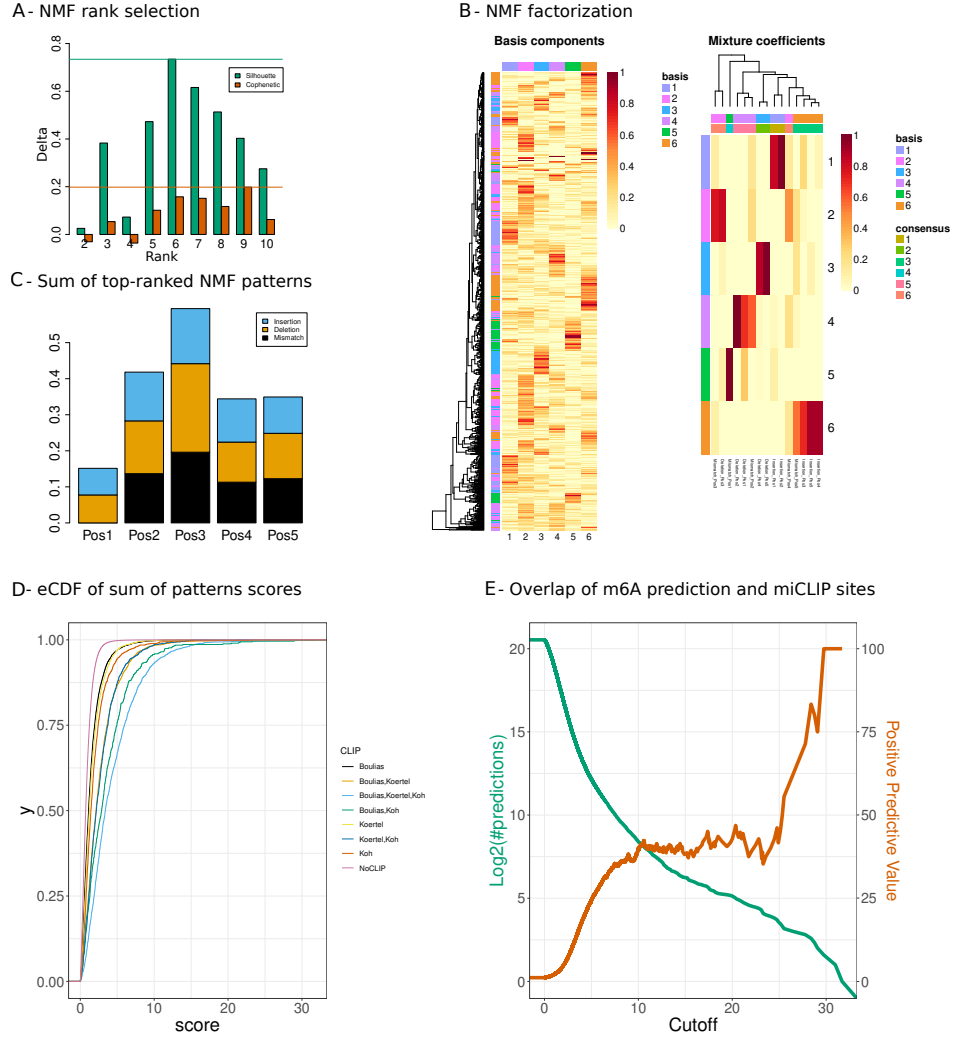


Figure 6: **Case 1. WT versus KO.** **A:** NMF rank selection. Barplots representing the difference of cophenetic correlation and silhouette indices between the NMF factorization of the original and the randomized data. Ranks with the largest value for both indices are determined, then the smallest rank 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).

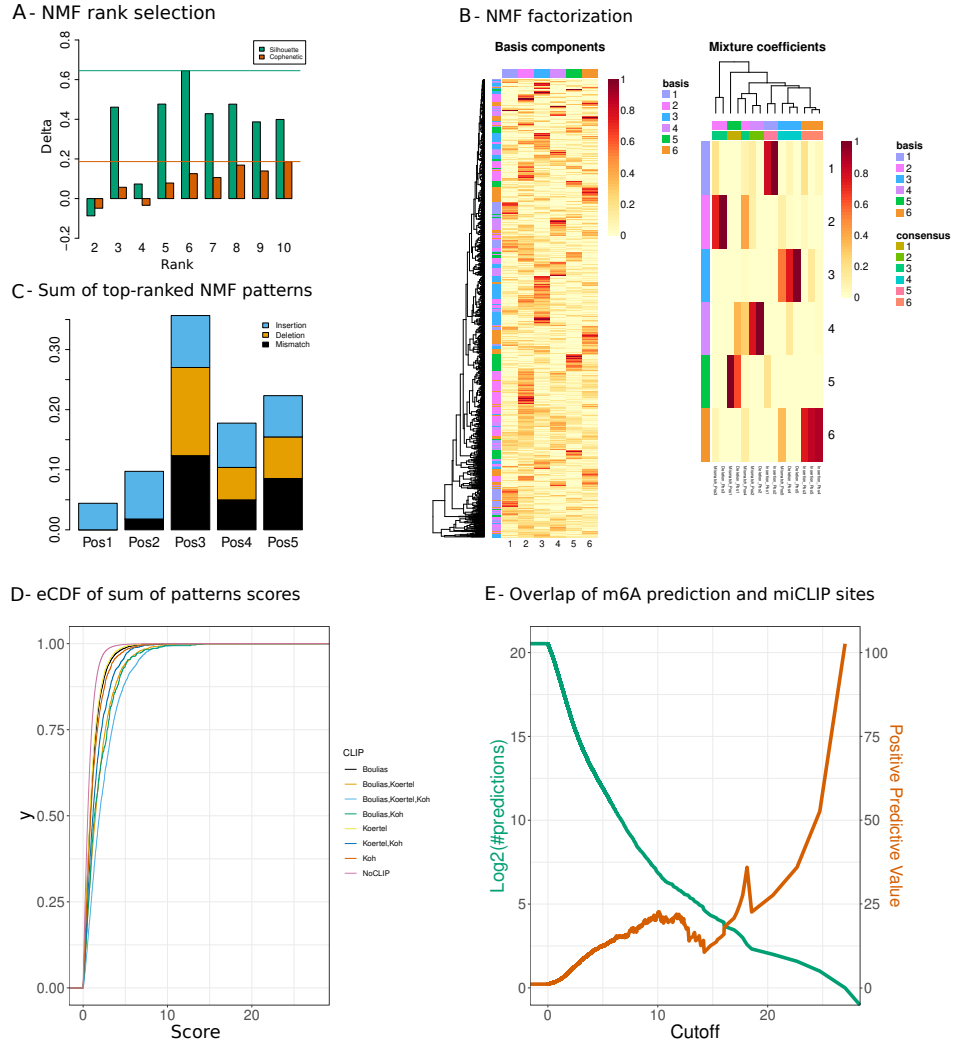


Figure 7: **Case 2. WT versus IVT.** **A:** NMF rank selection. Barplots representing the difference of cophenetic correlation and silhouette indices between the NMF factorization of the original and the randomized data. Ranks with the largest value for both indices are determined, then the smallest 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).