# The PARA-suite: PAR-CLIP specific sequence read simulation and processing

Andreas Kloetgen, Arndt Borkhardt, Jessica I Hoell, Alice C McHardy

**Background:** Next-generation sequencing (NGS) technologies have profoundly impacted biology over recent years. Experimental protocols, such as PhotoActivatable Ribonucleoside-enhanced Cross-Linking and ImmunoPrecipitation (PAR-CLIP), which identifies protein–RNA interactions on a genome-wide scale, commonly employ deep sequencing. With PAR-CLIP, the incorporation of photoactivatable nucleosides into nascent transcripts leads to high rates of specific nucleotide conversions during reverse transcription.

**Methods:** We show that differences in the error profiles of PAR-CLIP reads relative to regular transcriptome sequencing reads (RNA-Seq) make a distinct processing advantageous. We describe a set of freely available tools for this purpose, which are called the PAR-CLIP Analyzer suite (PARA-suite). The PARA-suite includes error model inference, PAR-CLIP read simulation, a full read alignment pipeline with a modified Burrows-Wheeler Aligner (BWA) algorithm, and CLIP read clustering.

**Results:** We examined the alignment accuracy of commonly applied read aligners on 10 simulated PAR-CLIP datasets using different parameter settings and identified the most accurate setup among those read aligners. Our processing pipeline allowed improvement of both alignment and binding site detection accuracy. We demonstrate the performance of the PARA-suite in conjunction with different binding site detection algorithms on several real PAR-CLIP and HITS-CLIP datasets.

**Availability:** The PARA-suite toolkit and the PARA-suite aligner are available at https://github.com/akloetgen/PARA-suite and https://github.com/akloetgen/PARA-suite_aligner, respectively, under the GNU GPLv3 license.

1 # The PARA-suite: PAR-CLIP specific sequence read

2 # simulation and processing

3

4

5 Andreas Kloetgen[1,2,3], Arndt Borkhardt[2], Jessica I. Hoell[2,§], Alice C. McHardy[1,3,§,*]

6

7 [1]Department of Algorithmic Bioinformatics, Heinrich Heine University, Düsseldorf, Germany;

8 [2]Department of Pediatric Oncology, Hematology and Clinical Immunology, Medical Faculty,

9 Heinrich Heine University, Düsseldorf, Germany;

10 [3]Computational Biology of Infection Research, Helmholtz Center for Infection Research,

11 Braunschweig, Germany

12 [§]These authors have contributed equally to the work;

13 [*]To whom correspondence should be addressed.

14

15 **Corresponding author:** Alice C. McHardy; Inhoffenstr. 7, 38124 Braunschweig, Germany;

16 AMC14@helmholtz-hzi.de

17

18 **Short title:** Simulating and processing PAR-CLIP data

19

20 **Key words:** next-generation sequencing, read simulation, read alignment, cross-linking and

21 immunoprecipitation (CLIP), post-transcriptional regulation, RNA-binding proteins

22

## Abstract

**Background:** Next-generation sequencing (NGS) technologies have profoundly impacted
biology over recent years. Experimental protocols, such as PhotoActivatable Ribonucleoside-
enhanced Cross-Linking and ImmunoPrecipitation (PAR-CLIP), which identifies protein–RNA
interactions on a genome-wide scale, commonly employ deep sequencing. With PAR-CLIP, the
incorporation of photoactivatable nucleosides into nascent transcripts leads to high rates of
specific nucleotide conversions during reverse transcription.

**Methods:** We show that differences in the error profiles of PAR-CLIP reads relative to regular
transcriptome sequencing reads (RNA-Seq) make a distinct processing advantageous. We
describe a set of freely available tools for this purpose, which are called the PAR-CLIP Analyzer
suite (PARA-suite). The PARA-suite includes error model inference, PAR-CLIP read simulation
based on PAR-CLIP specific properties, a full read alignment pipeline with a modified Burrows-
Wheeler Aligner (BWA) algorithm and CLIP read clustering for binding site detection.

**Results:** We examined the alignment accuracy of commonly applied read aligners on 10
simulated PAR-CLIP datasets using different parameter settings and identified the most accurate
setup among those read aligners. Our processing pipeline allowed improvement of both
alignment and binding site detection accuracy. We demonstrate the performance of the PARA-
suite in conjunction with different binding site detection algorithms on several real PAR-CLIP
and HITS-CLIP datasets.

**Availability:** The PARA-suite toolkit and the PARA-suite aligner are available at
https://github.com/akloetgen/PARA-suite and https://github.com/akloetgen/PARA-suite_aligner,
respectively, under the GNU GPLv3 license.

## Background

RNAs play a crucial role in cell survival and viability. Coding messenger RNAs (mRNAs), which are translated into proteins, and many other RNA species, such as small and long non-coding RNAs, ribosomal RNAs and transfer RNAs, are essential for the survival and proper functioning of the cells §Eddy, 2001 #310°. Most RNAs maintain their function by working together with the so-called RNA-binding proteins (RBPs) (Glisovic, Bachorik et al. 2008). RBPs are virtually involved in all steps of the mRNA lifecycle, from polyadenylation, translocation and modification to translation (Hieronymus and Silver 2004). Thus, it is not surprising that many RBPs which show aberrant functions or changes in expression patterns have been associated with disease progression or even with carcinogenesis (Lukong, Chang et al. 2008). For instance, the *FET* protein family, consisting of the three RBPs *FUS*, *EWSR1* and *TAF15*, is ubiquitously expressed and widely conserved in mammals. Genomic rearrangements, leading to mutant forms of these RBPs in humans, have been described as key players in sarcomas and leukemia (Tan and Manley 2009). More recently, two amyotrophic lateral sclerosis causing mutants of *FUS* have shown different RNA-binding patterns compared to the wild-type counterpart, supporting the importance of the function of *FUS* in mRNA processing (Hoell, Larsson et al. 2011).

Experimental protocols have been developed to analyze the functional network within a particular RBP interacts. A promising method for this purpose is the PhotoActivatable Ribonucleoside-enhanced Cross-Linking and ImmunoPrecipitation (PAR-CLIP) technique (Hafner, Landthaler et al. 2010). When coupled to deep sequencing, it identifies the bound RNAs for a particular RBP on a genome-wide scale. First, the cells are supplied with a specific photoactivatable nucleoside, such as 4-thiouridine (4-SU), which is incorporated as an alternative to the respective nucleoside into nascent mRNA transcripts. Afterwards, the cells are treated with ultraviolet (UV) light at 365 nm to cross-link the amino acids of RBPs to the nucleotides of their bound RNA molecules. The incorporation of 4-SU instead of uridine results in nucleotide conversions from uridine to cytidine at all cross-linked sites containing a 4-SU during reverse transcription (a necessary step for preparing cDNA libraries for sequencing). This specific replacement is also called a 'T–C conversion'. T–C conversions can be used to distinguish

80    between unspecifically bound RNA fragments (considered as contaminations) and those that are

81    specifically bound and cross-linked to the RBP of interest (Ascano, Hafner et al. 2012,

82    Golumbeanu, Mohammadi et al. 2015). We recently published a detailed protocol for the PAR-

83    CLIP procedure (Hoell, Hafner et al. 2014). Other CLIP protocols for the genome-wide

84    identification of RBP targets are also frequently used, such as the High-Throughput Sequencing

85    of RNAs isolated by Cross-Linking and ImmunoPrecipitation (HITS-CLIP, sometimes also

86    called CLIP-seq) or the iCLIP protocol (Chi, Zang et al. 2009, König, Zarnack et al. 2010).

87    HITS-CLIP mainly introduces deletions of a single base at the cross-linked sites, while single

88    nucleotide conversions do not seem to occur at a significant frequency (Zhang and Darnell 2011,

89    Sugimoto, König et al. 2012).

90

91    Current sequencing platforms allow sequencing of mammalian transcriptome libraries with a

92    high coverage. Nowadays, the most commonly used NGS platforms are 454, Illumina,

93    IonTorrent or PacBio  (van Dijk, Auger et al. 2014). Depending on the sequencing platform and

94    the sample type, sequencing errors vary in type and frequency. The errors that most commonly

95    occur are substitution errors and indels of a few bases between the sequencing read and the

96    reference template (large rearrangements, such as those leading to chimeras, are also possible

97    errors but are not discussed here) (Laehnemann, Borkhardt et al. 2015). In an RNA-Seq dataset a

98    single transcript will be covered by sequencing reads in all its expressed coding exons (apart

99    from, for example, amplification errors or alternative splicing variants). For common sequencing

100   data types, such as RNA-Seq and DNA-Seq, designated read aligners were recently developed.

101   These include short read aligners, such as BWA (Li and Durbin 2009) or Bowtie (Langmead,

102   Trapnell et al. 2009), and read aligners such as TopHat (Trapnell, Pachter et al. 2009), STAR

103   (Dobin, Davis et al. 2013) or Subjunc (Liao, Smyth et al. 2013), which can also handle longer

104   sequencing reads spanning exon-exon junctions. Specific software for the evaluation and

105   analysis of the PAR- and HITS-CLIP sequencing data is needed to accommodate their unique

106   error profiles (Kloetgen, Münch et al. 2015). For instance, the read aligner BWA PSSM

107   (Kerpedjiev, Frellsen et al. 2014) makes use of a predefined position specific scoring matrix to

108   process the error-prone PAR-CLIP reads.

109   In general, sequencing error profiles of RNA-Seq datasets, including PAR-CLIP data, can vary

110   between different sequencing runs, depending on the sequencing machine, experimental

111 conditions or the biological properties of the sample (Laehnemann, Borkhardt et al. 2015,

112 Schirmer, Ijaz et al. 2015). Here, we describe the PAR-CLIP Analyzer suite (PARA-suite),

113 which includes a PAR-CLIP read simulator, an error estimation tool for CLIP datasets and an

114 alignment pipeline based on a novel alignment algorithm performing on the fly data-set specific

115 error estimation.  The alignment pipeline thus automatically adjusts to the quality and error

116 profiles of individual sequencing datasets. We compared PAR-CLIP sequencing reads to regular

117 transcriptome sequencing reads (RNA-Seq) to identify distinctive properties relevant for the

118 reference-based read alignment and RBP binding site detection from PAR-CLIP datasets.

119 Generation of simulated PAR-CLIP datasets can be performed with the PAR-CLIP read

120 simulator. The PARA-suite toolkit is available at https://github.com/akloetgen/PARA-suite and

121 https://github.com/akloetgen/PARA-suite_aligner, implemented as an extension of BWA. It is

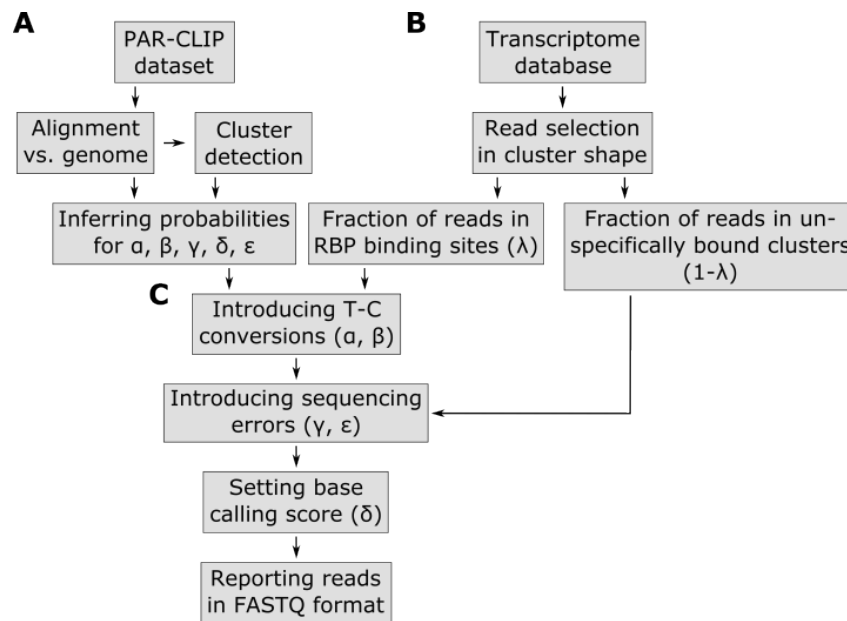122 licensed under GNU GPLv3 and implemented in the programming languages Java and C.


123 **Methods**

124 **2.1 Datasets and read aligners**

125

126 We downloaded PAR-CLIP data of the *FET* family from DRASearch database

127 (https://trace.ddbj.nig.ac.jp/DRASearch/) with accession number SRA025082 (Hoell, Larsson et

128 al. 2011), *HuR* dataset with accession number SRR248532, *MOV10* dataset with accession

129 number SRR490650 and HITS-CLIP data on the Argonaute protein (Chi, Zang et al. 2009) from

130 http://ago.rockefeller.edu/. For estimating the error profiles of regular RNA-Seq runs, we

131 downloaded two sequencing lanes with the accession numbers SRR896663 and SRR896664 of

132 an NGS quality assessment study (SEQC/MAQC-III-Consortium 2014) from DRASearch and

133 pooled the data.

134 We used the following read aligners and versions, shown in alphabetic order: Bowtie, version

135 0.12.7 (Langmead, Trapnell et al. 2009), Bowtie2, version 2.2.3 (Langmead and Salzberg 2012),

136 BWA, version 0.7.8 (Li and Durbin 2009), BWA PSSM, initial release version (Kerpedjiev,

137 Frellsen et al. 2014), MOSAIK, version 2.2.3 (Lee, Stromberg et al. 2014), STAR, version 2.3.0

138 (Dobin, Davis et al. 2013), Subjunc, version 1.4.2 (Liao, Smyth et al. 2013) and TopHat, version

139 2.0.13 (Trapnell, Pachter et al. 2009).

140

## 2.2 PAR-CLIP read simulator and hierarchical clustering

We developed a PAR-CLIP read simulator (Figure 1) that creates short RNA reads which mimic important PAR-CLIP specific properties (Section 3.1). First, the following probability distributions are obtained from real PAR-CLIP data: (a) a probability matrix $\varepsilon$ representing the background error profile of sequencing errors, (b) a probability vector of T–C conversion frequencies $\alpha$ for ranked T–C conversion sites, (c) a probability vector $\beta$ for preferred read positions of T–C conversion sites within binding sites, (d) a probability vector $\mu$ for indel frequencies per read position and (e) a probability vector $\delta$ for the base calling quality score distribution per read position. The probability matrix $\varepsilon$ contains a probability distribution for each DNA base over the DNA bases $\{A, C, G, T\}$. For this purpose, a PAR-CLIP dataset is aligned against a reference genome sequence with an appropriate read aligner.

**A**

PAR-CLIP dataset → Alignment vs. genome → Cluster detection → Inferring probabilities for α, β, γ, δ, ε

**B**

Transcriptome database → Read selection in cluster shape → Fraction of reads in RBP binding sites (λ) / Fraction of reads in un-specifically bound clusters (1-λ)

**C**

Introducing T-C conversions (α, β) → Introducing sequencing errors (γ, ε) → Setting base calling score (δ) → Reporting reads in FASTQ format

154

155   **Figure 1: Pipeline of the PAR-CLIP read simulator implemented in the PARA-suite.** Part A

156   describes the generation of the error profile and further parameters learned from a real PAR-

157   CLIP dataset. Part B starts to generate reads mapping to RBP binding sites (clusters) on

158   transcript regions from a given transcript database (e.g. Ensembl genes). In part C, the pre-

159   calculated profiles are used to introduce T–C conversions, sequencing errors, indels and base

160   calling quality scores to the defined reads.

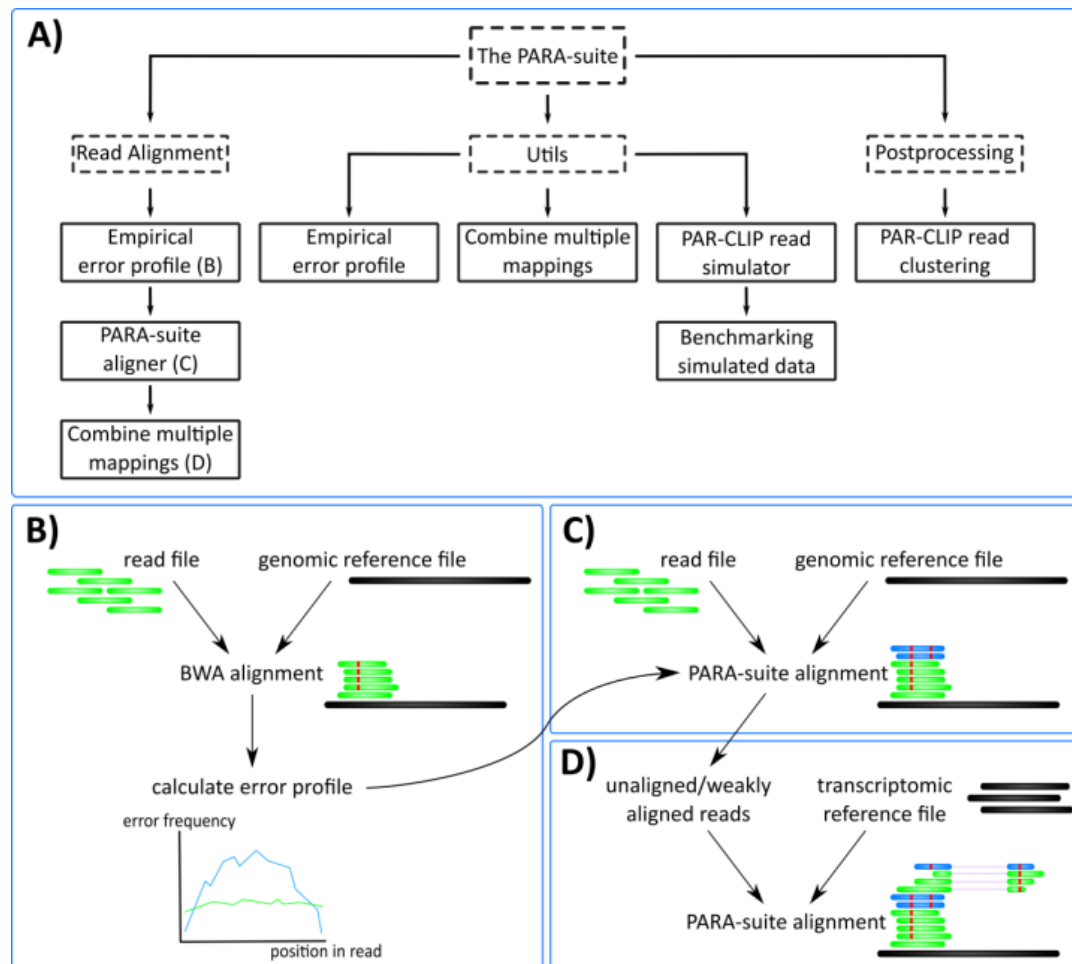161

162   Based on these alignments, the sequencing error profile ε, excluding PAR-CLIP specific T–C

163   conversions, is estimated from the observed frequencies of all single nucleotide substitutions,

164   except for T–C errors. Standard T–C sequencing errors are approximated by the average over all

165   the other sequencing error frequencies.  The probability vectors μ and δ are also inferred from

166   these alignments. Next, all aligned reads of the real dataset are clustered (stacked) using single-

167   linkage hierarchical clustering based on their genomic mapping positions, using 5 bases overlap

168   of the genomic mapping positions as the clustering threshold. To identify high confidence

169   clusters (sometimes referred to as binding sites) as defined in literature (Hafner, Landthaler et al.

170   2010), clusters which contain less than 10 reads, less than 25% T–C conversions per cluster, are

171   longer than 75 bases and include only T–C conversion sites reported as single nucleotide

172   polymorphisms (SNP) loci in the dbSNP database (version 142) (Sherry, Ward et al. 2001) are

173   discarded. This implementation of hierarchical clustering is part of the PARA-suite and will later

174    on also be used for binding site detection. For the subsequent simulation, the positions and

175    frequencies of highly mutated T–C sites within reads are determined to estimate α and β from the

176    high confidence clusters (Figure S1A-B).

177

178    Next, the PAR-CLIP read simulation starts with the random selection of transcripts from a pre-

179    selected database of annotated transcripts. One to at most three clusters (number of clusters

180    randomly chosen from a uniform distribution) containing several reads are created for a selected

181    transcript sequence. The starting positions of the clusters are randomly selected from a uniform

182    distribution within the entire range of a transcript. The number of reads simulated for a single

183    cluster is drawn from a normal distribution with a mean of 16 and standard deviation of 10. This

184    enables the simulation of a wide range of read coverages throughout the clusters. Furthermore,

185    small shifts of the start and end site of each read leading to the distinctive alignment position

186    shifts in the shape of a cluster are randomly introduced at this step (normal distribution with s.d.

187    1). A user-defined parameter λ ∈ [0,1] specifies the fraction of clusters that are considered

188    binding-sites, while the remaining clusters mimic contaminations of unbound RNAs which occur

189    in all PAR-CLIP experiments. We recommend values in the range of 0.5–0.7 (50–70%), as we

190    observed this range of aligned sequencing reads stacking into clusters after hierarchical

191    clustering and filtering (Table S1; similar values were previously reported by (Ascano, Hafner et

192    al. 2012)). If more than one T–C site is simulated for a single cluster, a major T–C conversion

193    site is selected according to the site-specific T–C conversion profile β and T–C conversion

194    probabilities are drawn from α. Subsequently, background sequencing errors are introduced

195    based on the pre-computed probability matrix ε and frequency vector μ for substitutions and

196    indels, respectively. In the last step, every base receives a base calling quality score, as specified

197    by the position-specific quality score distribution δ. All generated reads are stored in the

198    universal FASTQ format (Cock, Fields et al. 2010). The PAR-CLIP read simulator is available

199    through the PARA-suite.

200

201    **2.3 The PARA-suite – tools for error profile inference, read simulation, multiple**

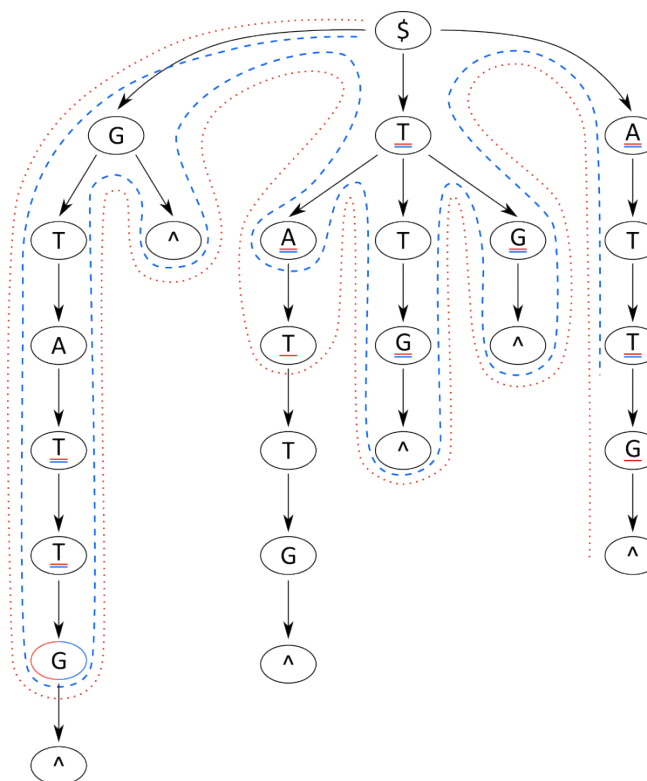202    **database mapping and more**

203

204   The PARA-suite is a toolkit for processing and aligning short and error-prone sequencing reads.

205   It is implemented in Java using HTSjdk, a Java API for high-throughput sequencing data formats

206   ([https://github.com/samtools/htsjdk](https://github.com/samtools/htsjdk)). The PARA-suite allows the user to estimate a sequencing

207   run-specific error profile, combine the results of multiple reference database alignments, cluster

208   an aligned sequencing read dataset (Section 2.2), run the PAR-CLIP read simulator, benchmark

209   an alignment of simulated PAR-CLIP sequencing reads and run a full processing pipeline for

210   error-prone short read alignment (Figure 2A). The alignment pipeline of the PARA-suite

211   includes the calculation of an error profile for a particular sequencing run, applying the

212   alignment algorithm described in the following section, and optionally combines the results of

213   read mappings against multiple databases (Figure 2B–D). First, a read alignment against a

214   reference sequence is performed with a fast short read aligner. Per default, this is done with

215   BWA, as our evaluations demonstrated it to be a fast and accurate aligner (Section 3.3) on PAR-

216   CLIP reads. However, other read aligners can also be used to produce the reference-based read

217   alignment. This initial read alignment is used to estimate the underlying mismatch and indel

218   probabilities $M$, $I$ and $D$ (as described in the next section) of the sequencing run. Once the error

219   profile has been estimated, all sequencing reads can be aligned with the PARA-suite aligner

220   (Section 2.4) against the reference sequence(s). All aligned reads are reported in a BAM file.

221

222

223  **Figure 2: The PARA-suite.** (A) The PARA-suite. Dashed boxes show packages while the other

224  ones are executable programs. The Utils package includes tools for working with error-prone

225  sequencing data and the postprocessing package contains a tool for clustering an aligned PAR-

226  CLIP dataset to identify RBP-bound genomic regions. (B) Read alignment by a fast read aligner

227  is necessary to infer the error profile for a particular read dataset (we selected BWA). (C) The

228  PARA-suite aligner is applied to the entire dataset to map error-prone reads, indicated here by

229  the additional mapping of the two blue reads. (D) An optional alignment versus a transcriptome

230  reference database can be executed using the PARA-suite aligner to identify previously

231  unmapped reads.

232

233

234    **2.4 Algorithm of the PARA-suite aligner**

235

236    The general BWA algorithm uses a Burrows-Wheeler transform (Burrows and Wheeler 1994) to

237    create an index for a reference genome sequence and applies a backward search to identify

238    possible mapping positions in the genome for every single sequencing read. The backward

239    search starts with the last base of a read proceeding to its front, searching the partly

240    decompressed suffix trie - using the auxiliary Ferragina and Manzini (FM) index (Ferragina and

241    Manzini 2000) - for a matching predecessor base of the bases of the read sequence compared so

242    far. Even if a match can be found for a single comparison, mismatches are introduced and all

243    possible downstream paths within the suffix trie are considered, until a pre-defined threshold of

244    maximal mismatches is exceeded in a single path (Figure 3, red dotted line).

245



246

247    **Figure 3: Suffix trie paths for BWA and PARA-suite.** Paths of the algorithms through the

248    suffix trie aligning the read sequence GCCATG$ against the reference sequence GTTATG$

249    (where $ means the end of a sequence). The red dotted line represents the algorithm of the BWA

250    aligner, allowing for two mismatches, and the blue dashed line indicates the PARA-suite

251    algorithm. The underlined bases represent positions where the respective aligner introduces a

252    mismatch. The example shows that the PARA-suite needs 14 comparisons while BWA needs 16

253    comparisons. Therefore, the PARA-suite is slightly faster than BWA at finding the same match

254    represented by the red/blue circle (left path). Indels are not shown for simplicity.

255

256    The principle idea of the PARA-suite aligner is the introduction of a probability estimate for

257    each comparison of the backward search. This enables weighting mismatches according to their

258    probabilities they occur in the analyzed dataset. A sequencing run is initially characterized

259    according to its underlying error probabilities. This allows to determine specific error-profiles for

260    experimental techniques, such as the frequent T–C conversions in PAR-CLIP data, that are more

261    common than sequencing errors. The error profile $M$ is a $4 \times 4$ probability matrix specifying

262    substitution probabilities values $\in$ [0..1] for each reference base $\in$ {A, C, G, T} to read bases {A,

263    C, G, T} (Figure 4A). Indels are introduced during the alignment step separately, using estimated

264    probabilities $I \in$ [0,1] for insertions and $D \in$ [0,1] for deletions.

265

266    For each comparison between a read base read[$i$] at read position $i$ and a reference base ref[$j$] at

267    position $j$ in the reference sequence, the algorithm recursively calculates a joint probability value

268    $p$, to examine the chance of incorporating a matching base or a suitable error, including indels at

269    the respective read positions (Figure 4D):

270
$$p_i = \begin{cases} p_{i+1} \cdot D, & \text{if } ref[j] \text{ is deleted} \\ p_{i+1} \cdot I, & \text{if } read[i] \text{ is inserted} \\ p_{i+1} \cdot M(read[i], ref[j]), & \text{otherwise} \end{cases}$$

271    with $p_{|read|} = 1$, starting with $i = |read|$ - 1 and decreasing $i$ at each step, except in the case of a

272    deletion, where $i \geq 0$.

273

274    Before the alignment of a particular read starts, a minimal threshold $T$ for the probability $p$ is

275    necessary, to decide whether a reads is accepted as aligned or rejected. The calculation for $T$ is

276    dependent on a parameter $X$ for the average number of mismatches. Note that this is not a

277    maximal threshold in terms of absolute mismatches, as the number of the more frequent errors

278    per aligned read can exceed $X$. The parameter $X$ can be pre-defined by the user or as a default is

279    estimated as the expected number of mismatches for different read lengths based on the error

280    profile $M$ for a sequencing run. Next, the minimal threshold $T$ is computed (Figure 4B&C):

281

$$T = avg(\text{match})^{|read| - X} \cdot avg(\text{mismatch})^{X}$$

283

284    where $avg(\text{match}) = \frac{1}{5}\left[\sum_{i \in \{0..3\}} M_{i,i} + (1 - (I + D))\right]$ and

285    $avg(\text{mismatch}) = \frac{1}{14}\left[\sum_{i,j \in \{0..3\}; i \neq j} M_{i,j} + I + D\right].$

286

287    Both *avg(match)* and *avg(mismatch)* are normalized by the number of elements (four matches

288    plus one for no "indel" occurring, and 12 mismatches plus 2 for either a insertion or a deletion).

289    If $p$ falls below the pre-calculated threshold $T$ during read alignment, the path within the suffix

290    trie is assumed not to match the read and is rejected (Figure 3, blue dashed line). The algorithm

291    thus penalizes rare types of mismatches according to $M$, while frequent errors, such as T–C

292    errors in PAR-CLIP reads, are the most favored substitutions in the alignment process (Figure

293    4B–D).

294

**A)**

| error profile M (in %): | read | | | |
|---|---|---|---|---|
| | **A** | **C** | **G** | **T** |
| **A** | 99.0 | 0.3 | 0.3 | 0.4 |
| **C** | 0.5 | 98.9 | 0.4 | 0.2 |
| **ref** **G** | 0.4 | 0.2 | 99.1 | 0.3 |
| **T** | 0.3 | 6.3 | 0.4 | 93.0 |

```
I          =  0.254 %
D          =  0.224 %

avg(match) = 97.904 %
avg(mm)    =  0.748 %
X          =  2
```

**B)**

```
read:   ... CACGGCCACG    (|read| = 26)
ref a): .. TACGGTTACG    (3 T-C)
ref b): .. CACAGCCA-G    (2 sequencing errors)
ref c): .. CACGCCCATG    (1 seq. error, 1 T-C)
```

**C)**

$p \geq 0.97904^{24} * 0.00748^2 = 0.0000336522 = 3.36522 * 10^{-5} = T$

**D)**

a) $p_{26} = 1$; $p_{25} = 1 * 0.991$; $p_{24} = 0.991 * 0.989 = 0.980099$; ...; $p_{22} = 0.06113$; $p_{21} = 0.003851$; ...; $p_0 \approx 0.00020005 > T$ ✓

b) $p_{26} = 1$; $p_{25} = 0.991$; $p_{24} = 0.991 * 0.00254 = 0.002517$; ...; $p_{19} = 0.000007247 < T$ ✗

c) $p_{26} = 1$; $p_{25} = 0.991$; $p_{24} = 0.991 * 0.063 = 0.062433$; ...; $p_0 = 0.0001534 > T$ ✓

**Figure 4: The PARA-suite aligner approach.** (A) The error profile probability matrix $M$ and indel probabilities $I$ and $D$, which are used as input for the PARA-suite alignment algorithm, as well as some results of the intermediate calculations of the PARA-suite alignment algorithm. In $M$, only T–C conversions have a higher probability (6.3%) compared to sequencing error and indel probabilities. (B) The last characters of a particular read and three example mapping positions within a reference, called ref a–c. (C) The calculation of a maximum threshold $T$ for the mapping probability $p$ (see formula in main text, and values from A in this image). (D) The mapping probability calculation of the read when mapping to the references a–c. The read fails to map against ref b with two sequencing errors, while reference a and c are suitable mapping positions, where the probability $p$ is higher than the threshold $T$. For implementation, we worked with the open-source read aligner BWA, version 0.7.8, to extend its algorithm for the alignment of short and error-prone reads.

## Results

### 3.1 Properties of PAR-CLIP reads

To assess the most important properties of the PAR-CLIP sequencing reads for read alignment, we systematically compared PAR-CLIP datasets for the three RBPs *EWSR1*, *FUS* and *TAF15* (*FET* protein family) (Hoell, Larsson et al. 2011) to a recently published RNA-Seq run on human reference RNA (SEQC/MAQC-III-Consortium 2014). The 10 outermost bases of the SEQC/MAQC reads showed error rates with peaks at 1.5 and 2.2 errors per 100 reads (EPR). In contrast, the middle read length range showed an average of about 0.3 EPR (Figure S2A, red line). As the short reads of the *FET* PAR-CLIP datasets consist only of these outermost bases, they exhibited a 2–3 times higher average sequencing error rate (about 0.7 EPR or even higher) than the SEQC/MAQC reads (Figure S2B, green line). When considering the T–C conversions only, we observed 1.319 EPR for *EWSR1*, 1.477 EPR for *FUS* and 1.051 EPR for *TAF15* on average. This is an approximately 20- to 30-fold increase in comparison to the SEQC/MAQC dataset with 0.051 EPR for T–C conversions on average (Figure S2). Moreover, we analyzed data from two further PAR-CLIP studies performed on the RBPs *HuR* (Mukherjee, Corcoran et al. 2011) and *MOV10* (Sievers, Schlumpf et al. 2012), which showed similar error profiles and EPRs to the *FET* PAR-CLIPs for T–C conversions (Figure S3).

Further analyses of the PAR-CLIP read datasets for *EWSR1*, *FUS*, *TAF15*, *MOV10* and *HuR* showed the PAR-CLIP reads (a) to be shorter than 30 bases, (b) to cover only short stretches of an expressed gene rather than the entire expressed RNA (these stretches are later on called clusters), (c) to exhibit a specific nucleotide conversion pattern with a strong enrichment of T–C conversions, where (d) such conversions occur in specific 'conversion sites' in the clusters. The first two properties (a) and (b) arise from the RNAse T1 treatment of the cells or the lysate during the PAR-CLIP experimental protocol. As only the short RNA fragments which are not digested by the endonuclease (probably protected by the binding pocket of the RBP) are sequenced, the lengths of those fragments are usually short. However, the nucleotide composition of those reads is strongly affected by the digestion enzyme and can vary among different digestion enzymes (Kishore, Jaskiewicz et al. 2011). After quality trimming and adapter trimming of the five PAR-CLIP datasets, the reads had a length usually shorter than 30 bases. As the transcript regions outside of the bound RNA fragment are digested by the endonuclease,

340 these are removed during immunoprecipitation and not sequenced, except for additional binding

341 sites on the same transcript further up- or downstream. Thus, the sequencing reads are stacked

342 into short clusters covering short stretches of the gene and representing the RBP-bound regions

343 of the transcripts (Figure S4A).

344 The two properties (c) and (d) arise from the incorporation of photoactivatable nucleosides into

345 the nascent transcripts during transcription. In the case of 4-SU, T–C conversions occur in the

346 sequencing reads at all cross-linked sites, where the 4-SU is incorporated instead of the native

347 uridine. These conversions can reach high rates in specific conversion sites within a cluster

348 (Hafner, Landthaler et al. 2010). In the analyzed datasets, we observed an average frequency of

349 about 70% T–C conversions in the main T–C conversion site (Figure S1A). This emphasizes that

350 simulated read datasets with specific properties are necessary for the evaluation of common short

351 read aligners for the analysis of PAR-CLIP read data. However, this cannot be created by

352 common sequencing read simulators, such as ART (Huang, Li et al. 2012) or GemSIM

353 (McElroy, Luciani et al. 2012). These produce simulated reads with a continuous coverage over

354 the entire transcript range and the introduced mutations are distributed randomly throughout the

355 simulated reads. This is not the case for PAR-CLIP sequencing reads.

356

## 3.2 PAR-CLIP read simulation for performance evaluation

358

359 We simulated a total of 10 PAR-CLIP read datasets based on information learned from three

360 previously published PAR-CLIP datasets of the *FET* protein family (Hoell, Larsson et al. 2011)

361 (Table S2). We imitated Illumina GenomeAnalyzer II sequence data according to the utilized

362 real datasets. The respective sequencing error and T–C conversion profiles were generated based

363 on alignments of all three datasets against the human reference genome sequence version 38

364 (GRCh38) (Lander, Linton et al. 2001). The error profile and additional estimated distributions

365 are similar to the ones from PAR-CLIP data on the two RBPs *HuR* and *MOV10,* indicating that

366 these profiles represent a reasonable approximation for PAR-CLIP data in general. We selected

367 human transcript sequences downloaded from Ensembl Genes Version 77 (Cunningham, Amode

368 et al. 2015) as our sequence database to simulate human transcript read sequences. We set λ, the

369 parameter for the fraction of sequencing reads stacking into clusters bound by the RBP, to 65%.

370 Such true RBP binding sites show high T–C conversion frequencies in different T–C conversion

371    sites. The remaining 35% of the simulated sequencing reads were designated to represent non-

372    specifically bound transcripts without an elevated T–C conversion rate, except for a few T–C

373    sequencing errors. These reflect RNA contaminations which can occur during the PAR-CLIP

374    experiment.

375    To assess the quality of the simulation, we then compared PAR-CLIP-specific properties

376    between the 10 simulated datasets and the *FET* PAR-CLIP data. Within a detected cluster of a

377    simulated dataset, shifts in the alignment positions of a few nucleotides at the beginning and the

378    end of the simulated cluster could be seen between the reads (Figure S4B). According to the

379    position-wise T–C conversion profile used, a T–C conversion site with a high conversion rate, as

380    well as a few sites with lower conversion rates, were usually present in the detected clusters (e.g.

381    Figure 1B). We compared the error profiles between one of the simulated datasets and the real

382    datasets, and distinguished between T–C errors and all other errors; the latter represent all

383    sequencing errors, except for the T–C sequencing errors (Figure S2C). Similar to the real data,

384    the distribution of the sequencing errors in the simulated dataset peaked at the beginning of the

385    reads and dropped to a mean error rate of 0.6 EPR in the middle read length range. Error rates

386    were slightly underestimated in the simulated data compared to real PAR-CLIP data, presumably

387    because of a small percentage of multiple mutations occurring at individual sites. Apart from

388    this, the simulated datasets appear to be representative for real PAR-CLIP data in the relevant

389    aspects.

390

391    **3.3 Accuracy of common read aligners and the PARA-suite on simulated PAR-CLIP**
392    **data**

393

394    Using the simulated PAR-CLIP datasets, we analyzed the accuracy of state-of-the-art read

395    aligners and common binding-site detection algorithms and compared these to the PARA-suite

396    alignment pipeline. The analyzed aligners, BWA and Bowtie, have often been employed in CLIP

397    studies (Lebedeva, Jens et al. 2011, Ascano, Mukherjee et al. 2012, Sievers, Schlumpf et al.

398    2012). BWA PSSM was applied with the PSSM for PAR-CLIP provided by its authors, because

399    a PSSM estimated from the sequencing dataset revealed a worse accuracy (data not shown).

400    MOSAIK was executed reporting only unique mappings, allowing for up to three mismatches

401    between the read and the reference sequence and using a Smith-Waterman bandwidth of 5. The

402    read aligners were used to align the simulated datasets to the reference sequence GRCh38. We

403    also executed the PARA-suite on the Ensembl Genes transcriptome database (version 77) and

404    combined the results with the genomic reference sequence alignments. These results are later

405    referred to as those of the "PARA-suite pipeline", while the results of the genomic alignment

406    step only using the PARA-suite are referred to as those of the "PARA-suite aligner". For the

407    PARA-suite aligner, the sequencing error and T–C conversion profiles for the simulated datasets

408    were obtained based on the BWA alignments allowing for two mismatches (BWA 2MMs) for

409    each of the simulated datasets separately. For a performance overview, we estimated the average

410    of the recall, precision and accuracy for every aligner over the 10 simulated datasets (calculation

411    described in Supplementary Methods). Unfortunately, BMix does not report negative clusters

412    (contaminations), thus we were not able to calculate the recall nor the accuracy, but only the

413    precision.

414    In terms of overall performance, the PARA-suite alignment performed best, with an accuracy of

415    69.74% for the PARA-suite aligner and 73.14% for the entire pipeline, showing performance

416    gains of 1.57% and 4.97% compared to the second-best aligner (BWA 2MM), respectively

417    (Table 1, Table S3). Many prominent PAR-CLIP studies have used Bowtie 1MM or BWA 2MM

418    for the read alignment step (Lebedeva, Jens et al. 2011, Mukherjee, Corcoran et al. 2011,

419    Ascano, Mukherjee et al. 2012, Sievers, Schlumpf et al. 2012, Mukherjee, Jacobs et al. 2014).

420    When comparing the PARA-suite alignment pipeline with these two aligners, the PARA-suite

421    pipeline showed an increase of 16.95% and 4.97% in the overall accuracy, respectively. Notably,

422    an average of 1.56% reads aligned by the PARA-suite pipeline are spanning an exon–exon

423    junction, but were not identified by the genomic reference mapping step only, but required

424    alignment versus the transcriptome reference sequences. Additionally, we compared the recall

425    (the fraction of correctly aligned reads out of all simulated reads) and the precision (the fraction

426    of correctly aligned reads out of all aligned reads) to assess the mapping ability of the read

427    aligners (Table 1, Figure S5). Here, the PARA-suite aligner and pipeline was ranked on places 1

428    and 3 regarding recall, and places 1 and 2 regarding precision, respectively, out of eight analyzed

429    alignment scenarios. Hence, the PARA-suite aligner and pipeline offer a notable performance

430    increase regarding all relevant performance measures in comparison to commonly used

431    computational analysis setups.

432　We then tested the accuracy of the binding site detection algorithms BMix, PARalyzer and the

433　hierarchical clustering of the PARA-suite using read alignment of the PARA-suite (Table S4).

434　The hierarchical clustering identified most correct binding sites; 3.26% more correct sites than

435　BMix and 5.54% more correct binding sites than PARalyzer. However, BMix identified fewer

436　false binding sites in comparison to the hierarchical clustering (20.30% less), and compared to

437　PARa lyzer (69.85% less). Furthermore, we investigated whether the PARA-suite alignment

438　increased the number of detected binding sites, irrespective of the used detection algorithm. In

439　conjunction with BMix, BWA 2MM (second best aligner) identified 7.17% less correct binding

440　sites than the PARA-suite aligner. For PARalyzer, BWA 2MM identified 2.97% less than by the

441　PARA-suite aligner. Finally, the hierarchical clustering identified 7.52% more correct binding

442　sites for the PARA-suite than for BWA 2MM. Overall, the combination of BMix with the

443　PARA-suite alignment provided the most accurate results on our simulated data.

444

445 **Table 1: Alignment accuracy on simulated PAR-CLIP data.** Most accurate alignment results for

446 different parameter settings for every read aligner on 10 simulated PAR-CLIP datasets. The results are

447 averaged per read aligner over all 10 datasets and are sorted by the accuracy.

| Aligner | Accuracy (in %) | Recall (in %) | Precision (in %) | Mapped overall | Mapped correctly | Real time (s) | Memory (GB) |
|---|---|---|---|---|---|---|---|
| PARA-suite pipeline | 73.14 | 84.49 | 71.85 | 1,024,792 | 969,948 | 396.8 | 6.27 |
| PARA-suite | 69.74 | 82.16 | 68.24 | 975,672 | 924,802 | 153.7 | 4.42 |
| BWA 2MMs | 68.17 | 82.31 | 64.98 | 959,171 | 904,034 | 359.2 | 4.42 |
| Bowtie 2MMs | 63.38 | 77.91 | 60.93 | 886,512 | 840,540 | 120.6 | 4.46 |
| BWA PSSM | 59.80 | 74.04 | 58.72 | 818,895 | 793,007 | 25.4 | 2.26 |
| TopHat | 59.69 | 76.10 | 55.35 | 844,902 | 791,549 | 282.9 | - |
| Bowtie2 | 56.22 | 73.23 | 51.43 | 763,893 | 745,531 | 13.4 | 3.32 |
| STAR | 50.74 | 69.57 | 43.02 | 826,871 | 672,920 | 248.6 | 28.39 |
| MOSAIK | 47.61 | 66.12 | 39.24 | 1,294,747 | 632,656 | 9,481.4 | 194.80 |
| Subjunc | 35.42 | 50.61 | 26.09 | 597,400 | 469,751 | 64.2 | 6.65 |

448

## 3.4 Analysis of *FET* PAR-CLIP datasets

450

451 To investigate the performance of the PARA-suite on real PAR-CLIP datasets, we applied it to

452 the three *FET* PAR-CLIP datasets (Hoell, Larsson et al. 2011). The sequencing reads were

453 preprocessed similarly to the original publication, and low quality ends and adapter sequences

454 were trimmed using Cutadapt (Martin 2011). Afterwards, all remaining reads longer than 18

455 bases were aligned against GRCh38 with BWA 2MMs, BWA PSSM and the PARA-suite aligner

456 (without the transcriptome mapping step to achieve comparable results). We measured the

457 fraction of aligned reads for all the aligners on the three datasets (Table S5). The PARA-suite

458 aligner generated the largest fraction of aligned reads over all three datasets in comparison to

459 BWA 2MM and BWA PSSM. Next, we stacked (clustered) all aligned reads using BMix and the

460 hierarchical clustering tool of the PARA-suite (Table 2). BWA 2MM identified fewer binding

461 sites compared to BWA PSSM or the PARA-suite, for read alignment prior to either BMix or

462 hierarchical clustering. Using the hierarchical clustering, the PARA-suite reported the largest

463 number of binding sites for two out of the three datasets. BWA PSSM identified 6.90% more

464    clusters than the PARA-suite aligner for the *FUS* dataset, while the PARA-suite aligner

465    identified 3.98% more clusters for the *EWSR1* dataset and 19.21% more clusters for the *TAF15*

466    dataset compared to BWA PSSM. In comparison to the numbers reported in the original

467    publication, the use of the PARA-suite aligner and hierarchical clustering increased the number

468    of binding sites by 33.71% for *EWSR1*, 16.77% for *FUS* and decreased them by 12.56% for

469    *TAF15*. After extracting distinct genes from all binding sites identified by the PARA-suite

470    (10,631 genes in total), 26.90% additional genes were found for all three RBPs, in comparison to

471    the original publication (7,771 genes in total). As expected for three RBPs from the same family,

472    there was a substantial overlap in terms of identified genes, with 2702 genes targeted by all three

473    RBPs (Figure S6).

474

475    **Table 2: Detected binding sites for the *FET* protein family.** Number of binding sites for the

476    *FET* protein family identified by the three aligners BWA PSSM, BWA 2MMs and the PARA-

477    suite in combination with the hierarchical clustering of the PARA-suite. Filters were applied

478    according to Section 2.2.

| | BWA 2MM / BMix | BWA 2MM / Clustering | BWA PSSM / BMix | BWA PSSM / Clustering | PARA-suite / BMix | PARA-suite / Clustering |
|---|---|---|---|---|---|---|
| **EWSR1** | 20,703 | 22,760 | 24,639 | 27,550 | 25,478 | 28,692 |
| **FUS** | 14,768 | 36,861 | 19,628 | 51,606 | 19,006 | 48,042 |
| **TAF15** | 5,086 | 5,810 | 5,238 | 6,130 | 5,862 | 7,588 |

479

480    ## 3.5 Analysis of PAR-CLIP data on *HuR*

481

482    We next applied the PARA-suite to a PAR-CLIP dataset on *HuR*, an RBP promoting RNA

483    stabilization (Mukherjee, Corcoran et al. 2011). Adapters and low quality ends for the *HuR*

484    dataset were trimmed using Cutadapt and reads shorter than 14 bases were discarded. The

485    binding motif of *HuR* is well-studied and is AU-rich, with a consensus motif described as

486    AUUUA, AUUUUA or AUUUUUA (Nabors, Suswam et al. 2003, Lebedeva, Jens et al. 2011),

487    showing potentially more T–C conversions within each binding site than other RBPs. As the

488    generated error-profile of the data set was similar to the ones of the *FET* PAR-CLIP data

489    (Section 3.1), the data quality seemed comparable. However, we noted a slight increase in T–C

490   conversions (Figure S3). The AU-rich binding motif might explain the higher T–C conversion

491   rate of 1.684 EPR compared to the conversion rate of 1.477 EPR e.g. for *FUS*.

492

493   We used the read aligners Bowtie2, Bowtie 2MM, BWA 2MM, BWA PSSM and the PARA-

494   suite to align the pre-processed dataset against the human genome reference GRCh38. Then, we

495   applied BMix and the hierarchical clustering of the PARA-suite to determine the binding sites of

496   *HuR* derived using the different read aligners. BWA PSSM in conjunction with BMix identified

497   most RBP binding sites within the genome – 3.69% more than the PARA-suite (Table S6). When

498   comparing the detected binding sites of BMix and the PARA-suite hierarchical clustering for

499   alignments created by the PARA-suite aligner (binding site positions overlapping by at least 13

500   bases), the difference was only marginal, with an overlap of more than 98.25% for the two

501   methods.  A recent study of this dataset reported binding sites using Bowtie 2MM for the

502   alignment step and PARalyzer for the binding site detection. We found the use of any BWA

503   PSSM or the PARA-suite alignment in conjunction with either BMix or hierarchical clustering to

504   increase the number of detected binding sites by 2.87% – 7.84%.

505

506   We searched for the exact binding motifs of *HuR* (ATTTA, ATTTTA and ATTTTTA) within the

507   BMix detected binding sites within 3' UTRs or introns for all tested read aligners. We found that

508   all aligners performed comparably, with motifs present in 42% to 44% of all detected binding

509   sites. The largest fraction was achieved using read alignments with BWA PSSM (44.33%), while

510   the PARA-suite aligner in combination with BMix found 42.53% most likely correct binding

511   sites. Bowtie 2MM in combination with BMix had the lowest fraction of binding sites containing

512   the reported binding motif (42.44%). We also compared the previously reported *HuR* binding

513   sites to the binding sites determined by the PARA-suite pipeline with BMix for clustering and

514   detected 13 out of 15 sites; namely 3' UTR PTGS2, 3'UTR CDKN1A, 3'UTR VEGFA, 3'UTR

515   TNF, 3' UTR SLC7A1, 3'UTR CCND1, 3'UTR MYC, 3' UTR XIAP, 3'UTR CELF1, TTS

516   CSF2, 3'UTR CCNB1, intron NCL and 3' UTR KRAS. The binding information for this

517   comparison was taken from the Ingenuity knowledge base (Calvano, Xiao et al. 2005). The

518   original study of the *HuR* dataset (Mukherjee, Corcoran et al. 2011) only reports 12 out of these

519   15 genes with a confirmed binding site.

520

521 **Discussion**

522

523    We here describe a read simulator to mimic PAR-CLIP datasets with error profiles drawn from

524    real PAR-CLIP datasets and the PARA-suite pipeline for error-aware read alignment and

525    processing. Furthermore, we provide a detailed characterization of the error profiles of PAR-

526    CLIP reads and an in depth performance assessment of short read aligners in combination with

527    binding site detection tools, to identify the most accurate read aligner and parameter settings on

528    PAR-CLIP reads. Common read simulators such as ART or GemSim do not allow simulating

529    PAR-CLIP reads with their specific error profiles. We characterized some of the unique

530    properties of PAR-CLIP sequence datasets that have, to our knowledge, so far not been analyzed,

531    such as preferred read positions for T–C conversion sites and their frequencies per read position.

532    We observed higher frequencies of sequencing errors in PAR-CLIP data in comparison to human

533    reference RNA-Seq data. A likely reason for this behavior could be that PAR-CLIP reads are

534    much shorter than common RNA-Seq reads, which reach lengths of 200 bases and show high

535    quality regions in the middle read length range (Laehnemann, Borkhardt et al. 2015, Schirmer,

536    Ijaz et al. 2015). We used these observations for the design of a PAR-CLIP read simulator that

537    embeds PAR-CLIP specific information into the simulation process.

538

539    Based on the simulated PAR-CLIP datasets, we determined parameter settings delivering the

540    best performance for commonly used aligners (Mukherjee, Corcoran et al. 2011, Ascano,

541    Mukherjee et al. 2012, Sievers, Schlumpf et al. 2012, Mukherjee, Jacobs et al. 2014). Our

542    analyses showed the read alignment to be crucial for RBP binding site detection from PAR-CLIP

543    datasets. The PAR-CLIP specific read properties make it nearly impossible to identify splice

544    junctions covered by PAR-CLIP reads with RNA-Seq read aligners such as TopHat, STAR or

545    Subjunc, as their algorithms are based on contrary assumptions, such as a similar read coverage

546    across all exons or long reads, to achieve high confidence k-mer spectra. Accordingly, these

547    three aligners were outperformed by other methods (Table S3–4). In addition, recent studies have

548    shown that BWT based aligners have less sensitivity in regions with genomic variation (Gontarz,

549    Berger et al. 2013). Interestingly, MOSAIK, an error-aware aligner based on hash queries that

550    was shown to be more robust on RNA-Seq reads than BWT based aligners (Lee, Stromberg et al.

551    2014), was also outperformed by most other tested methods. Although it is robust on longer

552  RNA-Seq reads, it seems to struggle with the very short PAR-CLIP reads. The PARA-suite

553  alignment pipeline allowed to increase the fraction of aligned reads in comparison to other

554  aligners, including alignment of reads spanning exon-exon junctions, both for PAR-CLIP

555  datasets and data from a HITS-CLIP study (Supplementary Results). We observed this

556  improvement irrespective of the applied downstream binding site detection algorithm. Different

557  from the error-aware short read aligner BWA PSSM, our short read alignment algorithm does

558  not need the manual input of an error profile, which is instead inferred *de novo* within individual

559  sequencing runs. The aligner thus automatically adapts to varying qualities of individual (PAR-

560  )CLIP sequencing runs and is specifically adjusted to each sequence dataset. To our knowledge,

561  it is the first tool for simultaneous *de novo* error model inference and short read alignment with

562  the BWA algorithm. Another difference to the BWA PSSM algorithm is that the latter introduces

563  mismatches under consideration of the base calling quality scores and a probabilistic background

564  model for matching bases in addition to the input error profile. In contrast, the generic error

565  profile estimation of the PARA-suite is not limited to any specific input profile. Further

566  applications of our software could thus be the analysis of other types of error-prone sequencing

567  data such as bisulphite sequencing data, which introduces a high amount of C–T mutations

568  (Frommer, McDonald et al. 1992) or data from low-quality ancient DNA samples (Briggs,

569  Stenzel et al. 2007).

570

571  Our analysis of combinations of read aligners and binding site detection algorithms on simulated

572  and real datasets indicate that no single software performed best in terms of binding site

573  detection on the available PAR-CLIP datasets. This observation was recently also made on

574  further datasets (Kassuhn, Ohler et al. 2016). Our analysis of the *HuR* and *FUS* datasets revealed

575  that U-rich binding sites tended to show higher rates of T–C conversions per read and were best

576  aligned by BWA PSSM. RBPs with a more heterogeneous nucleotide distribution (e.g. *EWSR1*

577  and *TAF15*) within the binding site are better assessed by the PARA-suite aligner. This is

578  supported by an analysis of uridylate-rich sequences from our simulated data aligned by BWA

579  PSSM and the PARA-suite (Supplementary Results and Supplementary Table 7). Thus, a

580  preliminary analysis of the error profile using the PARA-suite error profiler could allow

581  determining the best approach for analyzing sequencing data of a novel, yet uncharacterized

582  RBP.

583

## Acknowledgements

585

586 The authors thank Johannes Droege, David Laehnemann and Cristina della Beffa for their critical

587 comments on the manuscript.

588

## References

590

591 Ascano, M., Hafner, M., Cekan, P., Gerstberger, S. and Tuschl, T. (2012). Identification of RNA–protein
592     interaction networks using PAR-CLIP. *Wiley Interdiscip. Rev. RNA* **3**(2): 159-177.
593 Ascano, M., Mukherjee, N., Bandaru, P., Miller, J. B., Nusbaum, J. D., Corcoran, D. L., Langlois, C.,
594     Munschauer, M., Dewell, S. and Hafner, M. (2012). FMRP targets distinct mRNA sequence
595     elements to regulate protein expression. *Nature* **492**(7429): 382-386.
596 Briggs, A. W., Stenzel, U., Johnson, P. L., Green, R. E., Kelso, J., Prüfer, K., Meyer, M., Krause, J., Ronan,
597     M. T. and Lachmann, M. (2007). Patterns of damage in genomic DNA sequences from a
598     Neandertal. *Proc. Natl. Acad. Sci. U S A* **104**(37): 14616-14621.
599 Burrows, M. and Wheeler, D. J. (1994). A block-sorting lossless data compression algorithm. *CA: Digital
600     Equipment Corporation* **Technical Report 124 Palo Alto**.
601 Calvano, S. E., Xiao, W., Richards, D. R., Felciano, R. M., Baker, H. V., Cho, R. J., Chen, R. O., Brownstein,
602     B. H., Cobb, J. P. and Tschoeke, S. K. (2005). A network-based analysis of systemic inflammation
603     in humans. *Nature* **437**(7061): 1032-1037.
604 Chi, S. W., Zang, J. B., Mele, A. and Darnell, R. B. (2009). Argonaute HITS-CLIP decodes microRNA–mRNA
605     interaction maps. *Nature* **460**(7254): 479-486.
606 Cock, P. J., Fields, C. J., Goto, N., Heuer, M. L. and Rice, P. M. (2010). The Sanger FASTQ file format for
607     sequences with quality scores, and the Solexa/Illumina FASTQ variants. *Nucleic Acids Res.* **38**(6):
608     1767-1771.
609 Cunningham, F., Amode, M. R., Barrell, D., Beal, K., Billis, K., Brent, S., Carvalho-Silva, D., Clapham, P.,
610     Coates, G., Fitzgerald, S., Gil, L., Giron, C. G., Gordon, L., Hourlier, T., Hunt, S. E., Janacek, S. H.,
611     Johnson, N., Juettemann, T., Kahari, A. K., Keenan, S., Martin, F. J., Maurel, T., McLaren, W.,
612     Murphy, D. N., Nag, R., Overduin, B., Parker, A., Patricio, M., Perry, E., Pignatelli, M., Riat, H. S.,
613     Sheppard, D., Taylor, K., Thormann, A., Vullo, A., Wilder, S. P., Zadissa, A., Aken, B. L., Birney, E.,
614     Harrow, J., Kinsella, R., Muffato, M., Ruffier, M., Searle, S. M., Spudich, G., Trevanion, S. J., Yates,
615     A., Zerbino, D. R. and Flicek, P. (2015). Ensembl 2015. *Nucleic Acids Res.* **43**(Database issue):
616     D662-669.
617 Dobin, A., Davis, C. A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M. and
618     Gingeras, T. R. (2013). STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**(1): 15-21.
619 Ferragina, P. and Manzini, G. (2000). Opportunistic data structures with applications. *Proceedings of the
620     41st Symposium on Foundations of Computer Science*: 390-398.

621  Frommer, M., McDonald, L. E., Millar, D. S., Collis, C. M., Watt, F., Grigg, G. W., Molloy, P. L. and Paul, C.
622       L. (1992). A genomic sequencing protocol that yields a positive display of 5-methylcytosine
623       residues in individual DNA strands. *Proc. Natl. Acad. Sci. U S A* **89**(5): 1827-1831.
624  Glisovic, T., Bachorik, J. L., Yong, J. and Dreyfuss, G. (2008). RNA-binding proteins and post-
625       transcriptional gene regulation. *FEBS Lett.* **582**(14): 1977-1986.
626  Golumbeanu, M., Mohammadi, P. and Beerenwinkel, N. (2015). BMix: probabilistic modeling of
627       occurring substitutions in PAR-CLIP data. *Bioinformatics*: btv520.
628  Gontarz, P. M., Berger, J. and Wong, C. F. (2013). SRmapper: a fast and sensitive genome-hashing
629       alignment tool. *Bioinformatics* **29**(3): 316-321.
630  Hafner, M., Landthaler, M., Burger, L., Khorshid, M., Hausser, J., Berninger, P., Rothballer, A., Ascano, M.,
631       Jr., Jungkamp, A. C., Munschauer, M., Ulrich, A., Wardle, G. S., Dewell, S., Zavolan, M. and
632       Tuschl, T. (2010). Transcriptome-wide identification of RNA-binding protein and microRNA
633       target sites by PAR-CLIP. *Cell* **141**(1): 129-141.
634  Hieronymus, H. and Silver, P. A. (2004). A systems view of mRNP biology. *Genes Dev.* **18**(23): 2845-2860.
635  Hoell, J. I., Hafner, M., Landthaler, M., Ascano, M., Farazi, T. A., Wardle, G., Nusbaum, J., Cekan, P.,
636       Khorshid, M. and Burger, L. (2014). Transcriptome-Wide Identification of Protein Binding Sites
637       on RNA by PAR-CLIP (Photoactivatable-Ribonucleoside-Enhanced Crosslinking and
638       Immunoprecipitation). *Handbook of RNA Biochemistry: Second, Completely Revised and
639       Enlarged Edition*. A. B. R.K. Hartmann, A. Schön, and E. Westhof. Weinheim, Wiley-VCH Verlag
640       GmbH & Co. KGaA. **II:** 877-898.
641  Hoell, J. I., Larsson, E., Runge, S., Nusbaum, J. D., Duggimpudi, S., Farazi, T. A., Hafner, M., Borkhardt, A.,
642       Sander, C. and Tuschl, T. (2011). RNA targets of wild-type and mutant FET family proteins. *Nat.
643       Struct. Mol. Biol.* **18**(12): 1428-1431.
644  Huang, W., Li, L., Myers, J. R. and Marth, G. T. (2012). ART: a next-generation sequencing read simulator.
645       *Bioinformatics* **28**(4): 593-594.
646  Kassuhn, W., Ohler, U. and Drewe, P. (2016). *Cseq-simulator: a data simulator for CLIP-Seq experiments*.
647       Pacific Symposium on Biocomputing. Pacific Symposium on Biocomputing.
648  Kerpedjiev, P., Frellsen, J., Lindgreen, S. and Krogh, A. (2014). Adaptable probabilistic mapping of short
649       reads using position specific scoring matrices. *BMC Bioinformatics* **15**(1): 100.
650  Kishore, S., Jaskiewicz, L., Burger, L., Hausser, J., Khorshid, M. and Zavolan, M. (2011). A quantitative
651       analysis of CLIP methods for identifying binding sites of RNA-binding proteins. *Nat. Methods
652       **8**(7): 559-564.
653  Kloetgen, A., Münch, P. C., Borkhardt, A., Hoell, J. I. and McHardy, A. C. (2015). Biochemical and
654       bioinformatic methods for elucidating the role of RNA-protein interactions in posttranscriptional
655       regulation. *Brief. Funct. Genomics* **14**(2): 102-114.
656  König, J., Zarnack, K., Rot, G., Curk, T., Kayikci, M., Zupan, B., Turner, D. J., Luscombe, N. M. and Ule, J.
657       (2010). iCLIP reveals the function of hnRNP particles in splicing at individual nucleotide
658       resolution. *Nat. Struct. Mol. Biol.* **17**(7): 909-915.
659  Laehnemann, D., Borkhardt, A. and McHardy, A. C. (2015). Denoising DNA deep sequencing data—high-
660       throughput sequencing errors and their correction. *Brief. Bioinform.*: bbv029.
661  Lander, E. S., Linton, L. M., Birren, B., Nusbaum, C., Zody, M. C., Baldwin, J., Devon, K., Dewar, K., Doyle,
662       M. and FitzHugh, W. (2001). Initial sequencing and analysis of the human genome. *Nature
663       **409**(6822): 860-921.
664  Langmead, B. and Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**(4):
665       357-359.
666  Langmead, B., Trapnell, C., Pop, M. and Salzberg, S. L. (2009). Ultrafast and memory-efficient alignment
667       of short DNA sequences to the human genome. *Genome Biol.* **10**(3): R25.

668  Lebedeva, S., Jens, M., Theil, K., Schwanhäusser, B., Selbach, M., Landthaler, M. and Rajewsky, N. (2011).
669      Transcriptome-wide analysis of regulatory interactions of the RNA-binding protein HuR. *Mol.*
670      *Cell* **43**(3): 340-352.
671  Lee, W.-P., Stromberg, M. P., Ward, A., Stewart, C., Garrison, E. P. and Marth, G. T. (2014). MOSAIK: a
672      hash-based algorithm for accurate next-generation sequencing short-read mapping. *PLoS One*
673      **9**(3): e90581.
674  Li, H. and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform.
675      *Bioinformatics* **25**(14): 1754-1760.
676  Liao, Y., Smyth, G. K. and Shi, W. (2013). The Subread aligner: fast, accurate and scalable read mapping
677      by seed-and-vote. *Nucleic Acids Res.* **41**(10): e108.
678  Lukong, K. E., Chang, K.-w., Khandjian, E. W. and Richard, S. (2008). RNA-binding proteins in human
679      genetic disease. *Trends Genet.* **24**(8): 416-425.
680  Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput sequencing reads.
681      *EMBnet. J.* **17**(1): 10-12.
682  McElroy, K. E., Luciani, F. and Thomas, T. (2012). GemSIM: general, error-model based simulator of next-
683      generation sequencing data. *BMC Genomics* **13**(1): 74.
684  Mukherjee, N., Corcoran, D. L., Nusbaum, J. D., Reid, D. W., Georgiev, S., Hafner, M., Ascano, M., Tuschl,
685      T., Ohler, U. and Keene, J. D. (2011). Integrative regulatory mapping indicates that the RNA-
686      binding protein HuR couples pre-mRNA processing and mRNA stability. *Mol. Cell* **43**(3): 327-339.
687  Mukherjee, N., Jacobs, N. C., Hafner, M., Kennington, E. A., Nusbaum, J. D., Tuschl, T., Blackshear, P. J.
688      and Ohler, U. (2014). Global target mRNA specification and regulation by the RNA-binding
689      protein ZFP36. *Genome Biol.* **15**(1): R12.
690  Nabors, L. B., Suswam, E., Huang, Y., Yang, X., Johnson, M. J. and King, P. H. (2003). Tumor Necrosis
691      Factor α Induces Angiogenic Factor Up-Regulation in Malignant Glioma Cells A Role for RNA
692      Stabilization and HuR. *Cancer Res.* **63**(14): 4181-4187.
693  Schirmer, M., Ijaz, U. Z., D'Amore, R., Hall, N., Sloan, W. T. and Quince, C. (2015). Insight into biases and
694      sequencing errors for amplicon sequencing with the Illumina MiSeq platform. *Nucleic Acids Res.*:
695      gku1341.
696  SEQC/MAQC-III-Consortium (2014). A comprehensive assessment of RNA-seq accuracy, reproducibility
697      and information content by the Sequencing Quality Control Consortium. *Nat. Biotechnol.* **32**(9):
698      903-914.
699  Sherry, S. T., Ward, M. H., Kholodov, M., Baker, J., Phan, L., Smigielski, E. M. and Sirotkin, K. (2001).
700      dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res.* **29**(1): 308-311.
701  Sievers, C., Schlumpf, T., Sawarkar, R., Comoglio, F. and Paro, R. (2012). Mixture models and wavelet
702      transforms reveal high confidence RNA-protein interaction sites in MOV10 PAR-CLIP data.
703      *Nucleic Acids Res.* **40**(20): e160.
704  Sugimoto, Y., König, J., Hussain, S., Zupan, B., Curk, T., Frye, M. and Ule, J. (2012). Analysis of CLIP and
705      iCLIP methods for nucleotide-resolution studies of protein-RNA interactions. *Genome Biol.* **13**(8):
706      R67.
707  Tan, A. Y. and Manley, J. L. (2009). The TET family of proteins: functions and roles in disease. *J. Mol. Cell.*
708      *Biol.* **1**(2): 82-92.
709  Trapnell, C., Pachter, L. and Salzberg, S. L. (2009). TopHat: discovering splice junctions with RNA-Seq.
710      *Bioinformatics* **25**(9): 1105-1111.
711  van Dijk, E. L., Auger, H., Jaszczyszyn, Y. and Thermes, C. (2014). Ten years of next-generation
712      sequencing technology. *Trends Genet.* **30**(9): 418-426.
713  Zhang, C. and Darnell, R. B. (2011). Mapping in vivo protein-RNA interactions at single-nucleotide
714      resolution from HITS-CLIP data. *Nat. Biotechnol.* **29**(7): 607-614.

715