**Title**

UMI-tools: Modelling sequencing errors in Unique Molecular Identifiers to improve quantification accuracy

**Running Title**

Modelling UMI errors improves quantification accuracy

**Authors**

Tom Smith[1]

Andreas Heger[1]

Ian Sudbery[2]*

1. Computational Genomics Analysis and Training Programme, MRC WIMM Centre for Computational Biology, University of Oxford, John Radcliffe Hospital/Headley Way, Oxford OX3 9DS

2. Department of Molecular Biology and Biotechnology, University of Sheffield, Firth Court, Western Bank, Sheffield, UK, S10 2TN

* Corresponding author

1

## Keywords

UMI

Unique Molecular Identifier

Random Tag

PCR-duplicates

iCLIP

Single Cell RNA-seq

Sequencing

1    **Abstract**

2    Unique Molecular Identifiers (UMIs) are random oligonucleotide barcodes that are increasingly used

3    in high-throughput sequencing experiments. Through a UMI, identical copies arising from distinct

4    molecules can be distinguished from those arising through PCR-amplification of the same molecule.

5    However, bioinformatic methods to leverage the information from UMIs have yet to be formalised.

6    In particular, sequencing errors in the UMI sequence are often ignored, or else resolved in an *ad-hoc*

7    manner. We show that errors in the UMI sequence are common and introduce network-based

8    methods to account for these errors when identifying PCR duplicates. Using these methods, we

9    demonstrate improved quantification accuracy both under simulated conditions and real iCLIP and

10    single cell RNA-seq datasets. Reproducibility between iCLIP replicates and single cell RNA-seq

11    clustering are both improved using our proposed network-based method, demonstrating the value

12    of properly accounting for errors in UMIs. These methods are implemented in the open source *UMI-*

13    *tools* software package.

14

15      **Background**

16      High throughput sequencing technologies yield vast numbers of short sequences (reads) from a pool

17      of DNA fragments. Over the last ten years a wide variety of sequencing applications have been

18      developed which estimate the abundance of a particular DNA fragment by the number of reads

19      obtained in a sequencing experiment (read counting) and then compare these abundances across

20      biological conditions. Perhaps the most widely used read counting approach is RNA-seq, which seeks

21      to compare the number of copies of each transcript in different cell types or conditions. Prior to

22      sequencing, a PCR amplification step is normally performed to ensure sufficient DNA for sequencing

23      and/or enrichment for fragments with successful adapter ligation. Biases in the PCR amplification

24      step lead to particular sequences becoming overrepresented in the final library (Aird et al. 2011). In

25      order to prevent this bias propagating to the quantification estimates, it is common to remove reads

26      or read pairs with the same alignment coordinates as they are assumed to arise through PCR

27      amplification of the same molecule (Sims et al. 2014). This is appropriate where sequencing depth is

28      low and thus the probability of two independent fragments having the same genomic coordinates

29      are low, as with paired-end whole genome DNA-seq from a large genome. However, the probability

30      of generating independent fragments mapping to the same genomic coordinates increases as the

31      distribution of the alignment coordinates deviates from a random sampling across the genome

32      and/or the sequencing depth increases. For example, in RNA-seq, highly expressed transcripts are

33      more likely to generate multiple fragments with exactly the same genomic coordinates. The problem

34      of PCR duplicates is more acute when greater numbers of PCR cycles are required to increase the

35      library concentration, as in single cell RNA-seq, or when the alignment coordinates are limited to a

36      few distinct loci, as in individual-nucleotide resolution Cross-Linking and ImmunoPrecipitation

37      (iCLIP). Random barcodes were initially proposed as a method to count the number of mRNA

38      molecules in a sample (Hug and Schuler 2003), and have since been used to explicitly label PCR

39      duplicates (McCloskey et al. 2007). More recently, random barcodes, referred to as unique

40    molecular identifiers (UMIs), have been employed to confidently identify PCR duplicates in high-

41    throughput sequencing experiments (König et al. 2010b; Kivioja et al. 2012; Islam et al. 2014). By

42    incorporating a UMI into the same location in each fragment during library preparation, but prior to

43    PCR amplification, it is possible to accurately identify true PCR duplicates as they have both identical

44    alignment coordinates and identical UMI sequences (Figure 1a). In addition to their use in single cell

45    RNA-seq and iCLIP (König et al. 2010b) , UMIs may be applied to almost any sequencing method

46    where confident identification of PCR duplicates by alignment coordinates alone is not possible

47    and/or an accurate quantification is required, including ChIP-exo (He et al. 2015), DNA-seq

48    karyotyping (Karlsson et al. 2015), detection of rare mutations (Schmitt et al. 2012) and antibody

49    repertoire sequencing (Vollmers et al. 2013).

50    Accurate quantification with UMIs is predicated on a one-to-one relationship between the number

51    of unique UMI barcodes at a given genomic locus and the number of unique fragments that have

52    been sequenced. However, errors within the UMI sequence including nucleotide substitutions during

53    PCR, and nucleotide miss calling and insertions or deletions (Indels) during sequencing, create

54    additional artefactual UMIs. Nucleotide miss-calling and substitution errors affect only the UMI

55    sequence itself and do not affect the alignment coordinates. Hence, these errors will inflate the

56    estimation of the number of unique molecules at a particular genomic coordinate. These errors can

57    be identified by examining all UMIs at a single genomic coordinate. On the other hand, UMI Indels

58    will affect the alignment position also, leading to the assignment of reads to incorrect genome

59    coordinates. Identification of such events requires the examination of sets of UMIs at neighbouring

60    coordinates. Recombination events, also called 'PCR jumping' create chimeric sequences that may

61    change either the UMI sequence and/or alignment. Miss-calling during sequencing is by far the most

62    prevalent error, occurring 1-2 orders of magnitude more frequently than Indels for Illumina

63    sequencing (Marinier et al. 2015; Schirmer et al. 2015, 2016). Recombination is common when

64    sequencing amplicons, but much rarer with the shotgun sequencing approaches where UMIs are

65    utilised (Schloss et al. 2011; Waugh et al. 2015). We therefore focus here on improving

5

66 quantification via UMIs by considering nucleotide miss-calling and substitution errors within pools of

67 UMIs from the same genomic coordinate. Herein, we will refer to these errors as UMI errors.

68 UMI errors have been considered in previous analyses (Macosko et al. 2015; Bose et al. 2015; Yaari

69 and Kleinstein 2015; Islam et al. 2014). However, their impact on quantification accuracy has not

70 previously been demonstrated and there is no consistency in the approach taken to resolve these

71 errors. For example, Islam *et al* (2014) removed all UMIs where the counts were below 1 % of the

72 mean counts of all other non-zero UMIs at the genomic locus, whilst Bose *et al* (2015) merged

73 together all UMIs within a hamming distance of two or less, with little explanation as to how this was

74 achieved. We therefore set out to demonstrate the need to account for UMI errors, to compare

75 different methods for resolving UMI errors and to formalise an approach for removing PCR

76 duplicates with UMIs.

77 **Results**

78 We reasoned that UMI errors create groups of similar UMIs at a given genomic locus. To confirm

79 this, we calculated the average number of bases different (edit distance) between UMIs at a given

80 genomic locus and compared the distribution of average edit distances to a null distribution

81 generated by randomly sampling (see methods). Using iCLIP data (Müller-McNicoll et al. 2016), we

82 confirmed that the UMIs are more similar to one another than expected according to the null,

83 strongly suggesting sequencing and/or PCR errors are generating artefactual UMIs (see methods;

84 Figure 1b, see Figure S1 for other datasets). Furthermore, the enrichment of low edit distances is

85 well correlated with the degree of PCR duplication (Figure 1c). Overall, we detected a 25-fold

86 enrichment for positions with an average edit distance of 1, compared to our null expectation. In

87 contrast, when we compared the UMI sequences at adjacent positions we detected an 1.1-fold (+/-

88 standard deviation of 0.1, see materials and methods) enrichment for UMIs which may have

89 originated from a single nucleotide deletion, suggesting UMI Indels are much less prevalent than

90 UMI errors, as expected. We then constructed networks between UMIs at the same genomic locus

6

91     where nodes represent UMIs and edges connect UMI separated by a single nucleotide difference.

92     Whilst most of the networks contained just a single node, we observed that 3-36% of networks

93     contained two or more nodes, of which 4-20% did not contain a single central node, and thus could

94     not be naively resolved (Figure 1d). This indicates that the majority of networks are likely to

95     originate from a single unique molecule prior to PCR amplification, but a minority of networks may

96     originate from a combination of errors during PCR and sequencing or may originate from multiple

97     unique molecules, which by chance have similar UMIs.


98

99     **Methods to identify unique molecules**

100    Many previous studies assume each UMI at a given genomic locus represents a different unique

101    molecule (Collins et al. 2015; Shiroguchi et al. 2012; Soumillon et al. 2014). We refer to this method

102    as *unique*. Islam *et al* (2014) previously identified the issue of sequencing errors and proposed

103    removing UMIs whose counts fall below a threshold of 1% of the mean of all non-zero UMIs at the

104    locus, a method we refer to as *percentile*.

105    We have developed three methods to identify the number of unique molecules at a given locus by

106    resolving UMI networks formed by linking UMIs separated by a single edit distance (Figure 1e). In all

107    cases, the aim is to reduce the network down to a representative UMI(s) that accounts for the

108    network; identifying the exact sequence of the original UMI(s) is not important for the purposes of

109    quantification. The simplest method we examined was to merge all UMIs within the network,

110    retaining only the UMI with the highest counts. For this method, the number of networks formed at

111    a given locus is equivalent to the estimated number of unique molecules. This is similar to the

112    method employed by Bose et al (2015) where UMIs with an edit distance of 2 or less were

113    considered to originate from an identical molecule. We refer to this method as *cluster*. This method

114    is expected to underestimate the number of unique molecules, especially for complex networks. We

7

115     therefore developed the *adjacency* method which attempts to correctly resolve the complex

116     networks by using the node counts. The most abundant node and all nodes connected to it are

117     removed from the network. If this does not account for all the nodes in the network, the next most

118     abundant node and its neighbours are also removed. This is repeated until all nodes in the network

119     are accounted for. In the method, the total number of steps to resolve the network(s) formed at a

120     given locus is equivalent to the number of estimated unique molecules. This method allows a

121     complex network to originate from more than one UMI, although UMIs with an edit distance of two

122     will always be removed in separate steps. The excess of UMIs pairs with an edit distance of two

123     observed in the iCLIP datasets indicate that some of these UMIs are artefactual. Reasoning that

124     counts for UMIs generated by a single sequencing error should be higher than those generated by

125     two errors and UMIs resulting from errors during the PCR amplification stage should have higher

126     counts than UMIs resulting from sequencing errors, we developed a final method, *directional*. We

127     generated networks from the UMIs at a single locus, in which directional edges connect nodes a

128     single edit distance apart when $n_a \geq 2n_b - 1$, where $n_a$ and $n_b$ are the counts of node a and node b.

129     The entire directional network is then considered to have originated from the node with the highest

130     counts. The ratio between the final counts for the true UMI and the erroneous UMI generated from

131     a PCR error is dependent upon which PCR cycle the error occurrs and the relative amplification

132     biases for the two UMIs, but should rarely be less than 2-fold. The *-1* component was included to

133     account for strings of UMIs with low counts, each separated by a single edit distance for which the

134     *2n* threshold alone is too conservative. This method allows UMIs separated by edit distances greater

135     than one to be merged so long as the intermediate UMI is also observed, and with each sequential

136     base change from the most abundant UMI, the count decreases. For this method, the number of

137     directional networks formed is equivalent to the estimated number of unique molecules.

138

139

## Comparing methods with simulated data

140

141     To compare the accuracy of the proposed methods we simulated the process of UMI amplification

142     and sequencing for UMIs at a single locus and varied the simulation parameters (see methods). To

143     examine the accuracy of the 5 methods, we computed two metrics: The log2-fold difference

144     between the estimate and ground truth *Log2((estimate - truth) / truth)* and the coefficient of

145     variance (*standard deviation / mean*) across 10,000 iterations. Increasing UMI length or sequencing

146     depth results in a linear increase in the degree of overestimation for *unique* and *percentile* (Figure

147     2a, b), since increasing either parameter leads linearly increases the total amount of UMI sequence

148     that may harbour errors. In contrast, the estimates from the network-based methods remain

149     relatively stable, with *directional* showing the highest accuracy and lowest variance. We also

150     simulated the effect of including a very long UMI (up to 50 bp) as there may be occasions where it is

151     preferable to concatenate a UMI with another barcode, such as a sample barcode or cell barcode in

152     single cell RNA-seq, leading to longer barcodes. We noted that the network-based methods showed

153     reduced accuracy for very long barcodes (Figure S2a). Investigating further, we found this was

154     correlated with an increase in UMIs with two errors where the single error intermediate was not

155     observed, as detected by counting the number of networks which did not contain any of the initial

156     UMIs prior to PCR and sequencing (Figure S2b). In order to resolve this inherent problem with very

157     long UMIs, we modified the network-based methods so that edges joined nodes with an edit

158     distance less than or equal to 2. This considerably decreased the number of networks without any

159     initial UMIs and improved the accuracy of the network-based methods for very long UMI sequences

160     (Figure S2a-b).

161     Increased sequencing error rate leads to an exponential overestimation for *unique* and *percentile*

162     (Figure 2c), with a 1.3-fold overestimation observed with an error rate of 0.01, compared to less

163     than 1.05-fold for the network based methods. Increasing the rate of errors during the PCR step had

164     a similar impact (Figure S2c). However, this was only observed when the rate of DNA polymerase

9

165    errors was simulated as greater than 0.001, considerably higher than reported error rates for even

166    non-recombinant *Taq* DNA polymerase (Rittié and Perbal 2008; Whalen et al. 2016), confirming

167    sequencing errors are likely to be the primary source of UMI errors. Increasing the number of PCR

168    cycles or modifying the amplification bias had little impact on the relative accuracy of the methods

169    (Figure S2d, e). Increasing the number of initial UMIs reduced the accuracy of the network-based

170    methods, however even with 100 initial 8bp UMIs at a single locus, the network methods remained

171    the most accurate (Figure S2f).

172    Although the network methods performed very similarly, *directional* consistently yielded more

173    accurate and less variable estimates. For example, when the sequencing depth was increased to 400

174    reads, the average estimates were 19.92, 19.94 and 19.99  (truth=20) respectively for *cluster*,

175    *adjacency* and *directional* methods, and the CVs were 0.0167, 0.0144 and 0.0099. We observed no

176    difference between *percentile* and *unique* under most conditions tested. Increasing the number of

177    reads sequenced per initial UMI, we were able to see an improvement in accuracy for *percentile*

178    relative to *unique* when sequencing error rates are between $1 \times 10^{-3} - 1 \times 10^{-5}$, however, even under

179    this specific parameterisation, the network-based methods are more accuracy (Figure S2g).

180    In summary, under simulation conditions, the *directional* method outperforms all other methods,

181    however *adjacency* and *cluster* performs equally well under simulation conditions that are expected

182    to reflect a well-designed experiment and well-executed experiment.

183    **Implementation**

184    To implement our methods within the framework of removing PCR duplicates from BAM alignment

185    files, we developed a command line toolset, UMI-tools, with two commands, *extract* and *dedup*.

186    *extract* takes the UMI from the read sequence contained in a FASTQ read sequence and appends it

187    to the read identifier so it is retained in the downstream alignment. *extract* expects the UMIs to be

188    contained at the same location in each read. Where this is not the case, e.g with sequencing

189    techniques such as inDrop-seq (Klein et al. 2015), the user will need to extract the UMI sequence

190    from the read sequence and append it to the read identifier. *dedup* takes an alignment BAM file,

191    identifies reads with the same genomic coordinates as potential PCR duplicates, and removes PCR

192    duplicates using the UMI sequence according to the method chosen. Time requirements for running

193    *dedup* depend on number of input reads, length of UMI and level of duplication. Memory

194    requirements depend on the number of output reads. On a desktop with a Xeon E3-1246 CPU, it

195    takes ~220 seconds and ~100MB RAM to process a 32 million read single-end input file with 5bp

196    UMIs to ~700,000 unique alignments. Inputs with longer UMIs may take significantly longer.

197

198    **Comparing methods with iCLIP data**

199    We next sought to examine the effect of these methods on real data, starting with the previously

200    mentioned iCLIP data, which includes 3-6 replicates for 9 proteins (Müller-McNicoll et al. 2016). For

201    replicate 1, the distribution of the average edit distance between UMIs present at each genomic

202    locus showed enrichment for single edit distance relative to a null distribution from random

203    sampling, taking into account the genome-wide distribution of UMIs (Figure 3a). For all samples,

204    application of the *directional* method resulted in an edit-distance distribution resembling the null,

205    whereas using the *percentile* method made little or no difference. The same was also true of other

206    replicates of this dataset or other datasets (Figure S2). In some cases a residual enrichment of

207    positions with an average edit distance of 2 was observed, but this was also reduced in most cases.

208    We reasoned that if the *directional* method removed PCR duplicates more accurately, the

209    reproducibility between replicates should be improved. To test this we turned to a previously

210    defined measure of iCLIP reproducibility (König et al. 2010b) . Briefly, we identified in each sample

211    the bases with two or more tags mapping at that positions and asked what percentage had a tag

212    present in one or more other replicates for that pull-down. We limited the analysis to the first three

213    replicates for each protein. In each case, after de-duplication with the *directional* method, bases

214    with two or more tags were more reproducible (Figure 3b), with the difference being very large in

215    some cases (e.g. 21% vs 59% of bases reproducible for SRSF7 replicate 1). In contrast, the *percentile*

216    method was little different from *unique* (Figure S3).

217    In order to measure reproducibility of their data, Müller-McNicoll *et al* measured the spearman's

218    rank correlation between the numbers of significant tags in each exon across the genome. We

219    repeated this calculation with data processed using either the *unique* or *directional* method, and

220    compared the average spearman's correlation between each sample and other replicates of the

221    same pull down. In all cases we see an improvement in the correlation between replicates of the

222    same pull down when data are processed using the *directional* method (Figure 3c). As expected, the

223    degree of improvement for a particular sample was correlated with the enrichment of positions with

224    an average edit distance of 1 (Figure S3; $R^2$=0.4 ). Thus our method substantially improves the

225    reproducibility of replicates in this iCLIP experiment.

226

227    **Comparing methods with Single Cell RNA Seq data**

228    To further demonstrate the utility of our network-based method, we applied it to two differentiation

229    single cell RNA-seq data sets: the first reported use of UMIs in a single cell RNA-seq experiment

230    seeking to describe a developmental pathway (Soumillon et al. 2014), referred to here as SCRB-seq,

231    and a recently reported single cell RNA-seq utilising droplet-barcoding (Klein et al. 2015), referred to

232    here as inDrop-seq . As before, network-based methods show a marked improvement in the

233    distribution of edit distances over the *percentile* method and the *unique* method (Figure 4a).

234    Improvements are generally less pronounced than observed with the iCLIP data, likely due to a lower

235    maximum read depth in single cell RNA-seq. To demonstrate that this improvement in the edit

236    distance lead to an improved accuracy in transcript abundance estimates we used the ERCC spike-

237     ins. The naïve use of UMIs to identify PCR duplicates with the *unique* method improved the per-cell

238     correlation between ERCC concentration and counts, compared to quantification without

239     considering PCR duplicates (median coefficients were 0.86 and 0.89, respectively). As expected, the

240     correlation was further improved using the *directional* method (median coefficient = 0.91; Figure

241     4b).

242     We applied hierarchical clustering to the SCRB-seq gene expression data using the *unique* method

243     and observed the Day 0 and Day 14 cells separately relatively well (Figure 4c). However, 7 cells

244     clustered with cells of the wrong time point, reflecting either a failure to commit to differentiation or

245     miss-classification event due to noise in the expression estimates. With the *directional* method this

246     was reduced to 5 cells, suggesting that failure to account for UMI errors can lead to miss-

247     classification in single cell RNA-seq. Applying hierarchical clustering to the the the inDrop-seq gene

248     expression estimates, we observed that 44/2717 (1.6%) of cells clustered with cells from another

249     timepoint when using the *unique* method. Biological variation in the progression of differentiation

250     may explain Day 2, Day4 and Day 7 miss-classification events. However, 19/44 events involved

251     undifferentiated mES cells, suggesting these miss-classification events were the result of low-

252     accuracy quantification estimates (Figure 4d). With the application of the *directional* method, the

253     rate of miss-classification was reduced to 0.9% and, strikingly, all the mES cells were correctly

254     classified. These results indicate that application of the *directional* method improves the

255     quantification estimates and can improve classification by hierarchical clustering.

256     **Discussion**

257     UMIs can be utilised across a broad range of sequencing techniques, however bioinformatic

258     methods to leverage the information from UMIs have yet to be standardised. In particular, others

259     have noted the problem of UMI errors, but the solutions applied are varied (Bose et al. 2015; Islam

260     et al. 2014) . The *adjacency* and *directional* methods we set out here are, to our knowledge, novel

261     approaches to remove PCR duplicates when using UMIs. Comparing these methods to previous

13

262     methods with simulated data, we observed that our methods are superior at estimating the true

263     number of unique molecules. Of the three network-based methods, *directional* was the most robust

264     over the simulation conditions and should be preferred. We note that the performance of all

265     network-based methods will decrease as the number of aligned reads at a genomic locus approaches

266     the number of possible UMIs, however this is an intrinsic issue with UMIs and not one that can be

267     solved computationally post-sequencing. For this reason, we recommend all experiments to use

268     UMIs of at least 8 bp in length and to use longer UMIs for higher sequencing depth experiments. The

269     simulations also indicated that very long UMIs actually decrease the accuracy of quantification when

270     not accounting for UMI errors, since the UMIs are more likely to accumulate errors. For experiments

271     utilising long UMIs, network-based methods therefore show an even greater performance relative to

272     the *unique* method. The simulations provide an insight into the impact on quantification accuracy

273     and indicate that application of an error-aware method is even more important with higher

274     sequencing depth. This is perhaps most pertinent for single cell RNA-seq, as cost decreases continue

275     to drive higher sequencing depths.

276     The analysis of iCLIP and single cell RNA-seq and data sets established that UMI errors present in all

277     of the data sets tested and that quantification accuracy could therefore be improved by modelling

278     these errors during the deduplication step. The frequency of UMI Indels was far less than UMI errors

279     suggesting only minimal gains would be achieved by considering UMI Indels also. We observed an

280     improved distribution of edit distances for all samples when using network-based methods to detect

281     PCR duplicates, although theoretical reasoning and empirical evidence suggests that the extent of

282     the errors depends on the quality of the sequencing base calls and the sequencing depth, as

283     confirmed by the simulations.

284     Modelling UMI errors yielded improvements in single cell RNA-seq sample clustering, demonstrating

285     the value of considering UMI errors. Since iCLIP aims to identify specific bases bound by RNA binding

286     proteins, datasets have a high level of PCR duplication. The effects of UMI errors are therefore

14

287     particularly strong, creating the impression of reproducible cross-linking sites within a replicate but

288     not between replicates, for example only 21% of positions with two or more tags in SRSF7 replicate 1

289     had any tags in replicates 2 or 3 when naive de-duplication was used, but this increased to 59%

290     when the *directional* method was used (Figure 3b). Application of the network based methods

291     increases the correlation between replicates in all cases, with larger differences in samples where

292     PCR duplication was higher. From the results of the simulation and real data analyses, we

293     recommend the use of an error-aware method to identify PCR duplicates whenever UMIs are used.

294     We provide our methods within the open-source UMI-tools software

295     (https://GitHub.com/CGATOxford/UMI-tools, included here as Supplementary File 1), which can

296     easily be integrated into existing pipelines for analysis of sequencing techniques utilising UMIs.

297

298

299     **Methods**

300     *Simulation*

301     To simulate the effects of errors on UMI counts, an initial number of UMIs were generated at

302     random, with a uniform random probability of amplification [0.8-1.0] assigned to each initial UMI. To

303     simulate a PCR cycle, each UMI was selected in turn and duplicated according to its probability of

304     amplification. Polymerase errors were also added randomly at this stage and any resulting new UMI

305     sequences assigned new probabilities of amplification. Following multiple PCR cycles, a defined

306     number of UMIs were randomly sampled to model the sampling of reads during sequencing

307     ("sequencing depth") and sequencing errors were introduced with at a given probability, with all

308     errors (e.g A -> T, C -> G) being equally likely. The number of true UMIs within the sampled UMIs was

309     then estimated from the final pool of UMIs using each method. To test the performance of the

310     methods under a variety of simulation parameters, each parameter was varied in turn. The following

311     values are the range of the parameter values tested with the value used for all other simulations in

312     parentheses. Sequencing depth 10-400 (100), number of initial UMIs 10-100 (20), UMI length 6-16

313     (8), DNA polymerase error rate $1 \times 10^{-3} - 1 \times 10^{-7}$ ($1 \times 10^{-5}$), sequencing error rate $1 \times 10^{-1} - 1 \times 10^{-5}$ ($1 \times$

314     $10^{-3}$), number of PCR cycles 4-12 (6), minimum amplification probability 0.1-1 (0.8). The maximum

315     amplification probability was set at 1 with the probability of amplification for an each UMI drawn

316     from a uniform distribution.

317

318     **Real data**

319     Re-analysis of the iCLIP and Single Cell RNA-seq data was performed with in-house pipelines

320     following the methods described in the original publication with exceptions as highlighted below.

321     Pipelines are available at https://GitHub.com/CGATOxford/UMI-tools_pipelines and as

322     Supplementary File 2.

16

323

324     *iCLIP*

325     Raw sequence was obtained from the European Nucleotide Archive (accessions SRP059277 and

326     ERR039854) (Müller-McNicoll et al. 2016; Tollervey et al. 2011). Raw sequences were processed to

327     move the UMI sequences to the read name using 'umi_tools extract'. Sample barcodes were verified

328     and removed, and adaptor sequence removed from the 3' end of reads using the reaper tool from

329     the Kraken package (version 15-065) (Davis et al. 2013) with parameters: `-3p-head-to-tail 2 -3p-

330     prefix 6/2/1`. Reads were mapped to the same genome as the original publication (mm9 for SRSF

331     dataset, hg19 for the TDP43 dataset) using Bowtie version v1.1.2 (Langmead et al. 2009a) with the

332     same parameters as the original publications (-v 2 -m 10 -a).

333     We measured the rate at which UMIs might represent Indel mutations by noting that an Indel in the

334     UMI sequence would cause the final base of the presumed UMI to match the genomic base at

335     position -1 relative to the mapping location of the read. Thus we examined each UMI at a particular

336     position, and tested for the presence of a UMI that would correspond to a 1 bp deletion existed at

337     the following base. We compared this to the situation when the UMIs at the following base were

338     randomised, respecting the number of UMIs at the position and the genome-wide usage of each

339     UMI. Enrichment was defined as the count at the unrandomised positions compared to the count at

340     the randomised positions. We calculated this metric for one replicate of each pull down from

341     SRP059277. See the *Examining_indels* notebook in the *UMI-tools_pipelines* repository (included as

342     Supplementary File 2).

343     Mapped reads were deduplicated using 'umi_tools dedup' using each of the possible methods and

344     edit_distance distribution produced using the '--output-stats' option. For the *cluster* method only the

345     '--further-stats' option was used to output statistics on the distribution of network topology types.

346     Significant bases were produced by comparing tag count height at each position compared to

347     randomised profiles (König et al. 2010a), and bases with FDR<0.05 retained.

348     Coverage over exons was calculated by collapsing Ensembl 67 transcripts. Where exons overlapped,

349     they were restricted to their intersection and the number of reads mapped to significant bases

350     counted for each exon. Exons that contained no tags in any sample were removed (König et al.

351     2010a). Spearman's rho between all pairwise combinations of replicates of pulldowns for the same

352     protein were calculated and averaged for each replicate.

353     Reproducibility between replicates was calculated as per König *et al* (2010). Bases with a depth

354     greater than 2 were identified in the sample in question, and then the fraction of these bases that

355     had one or more tags in other replicates was calculated.

356

357     *Single Cell RNA-seq*

358     For both datasets, raw data was downloaded from Gene Expression Omnibus

359     (http://www.ncbi.nlm.nih.gov/geo). For The SCRB-seq data (GSE53638) (Soumillon et al. 2014), a

360     single Day 0 (SRR1058003) and Day 14 (SRR1058023) sample were obtained. For the inDrop data

361     (GSE65525) (Klein et al. 2015), the mouse ES cells sample 1 (SRR1784310), mouse ES cells LIF-, 2 days

362     (SRR1784313), mouse ES cells LIF-, 4 days (SRR1784314) and mouse ES cells LIF-, 7 days

363     (SRR1784315) samples were obtained. FASTQ files were extracted using SRA toolkit. The sequence

364     read filtering, preparation and alignment differed for the two data sets. In both cases, one of the

365     paired end reads contained adapter barcodes and UMI and the other read pair contained sequence

366     for alignment. In addition, with the inDrop data, the position of the UMI within the read varied

367     depending on the length of the cell barcode. For this reason, for both data sets, the UMIs had to be

368     extracted from the reads with bespoke code rather than using UMI-tools *extract*.

18

369     For SCRB-seq samples, the UMI was extracted from read 2 and appended onto the read identifier of

370     read 1 to generate a single-end FASTQ. Reads were filtered out if any of the following conditions was

371     not met. Phred sequence quality of all cell barcode bases >=10 and all UMI bases >=30 and cell

372     barcode matched expected cell barcodes. A reference transcriptome was built comprising all human

373     protein-coding genes (Ensembl v75, hg19) and the ERCC spike-ins. Since expression quantification

374     was being performed at the gene level, overlapping transcripts from the same gene were merged so

375     that each gene contained a single transcript covering all exons from all transcripts. Reads were

376     aligned to the reference transcriptome using BWA Aln (Li and Durbin 2009) with the following

377     parameters: *"-l 24 –k 2"* to set seed length to 24 bp, and mismatches allowed in the seed to 2.

378     For inDrop samples, the cell barcode and UMI were extracted from read 1 and read 2 was written

379     out to a single end FASTQ file with the cell barcode incorporated into the file name and the UMI

380     appended to the read identifier. Only reads containing the adapter sequence (allowing 2

381     mismatches) were retained. For each sample, only reads containing one of the *n* most abundant cell

382     barcodes were retained, where *n* was the number of cells in a given sample. The resulting single end

383     reads were filtered using trimmomatic v0.32 (Bolger et al. 2014) with the following options:

384     *"LEADING:28 SLIDINGWINDOW:4:20 MINLEN:19"* to remove bases with Phred quality scores below

385     28 from the 5' end, scan the reads in 4 bp sliding windows and trim when average quality score falls

386     below 20, and retain all reads at least 19bp in length following trimming. Our alignment procedure is

387     a deviation to the method used by Klein *et al* (2015) which involved alignment of reads to a

388     reference transcriptome containing all transcripts (e.g not collapsed into one gene model), reporting

389     up to 200 alignments per read, and dealing with multi-mapping alignments in a downstream step. As

390     this method was not compatible with our de-duplication method we took a simpler approach. A

391     reference transcriptome was built comprising all mouse protein-coding genes (Ensembl v78, mm10)

392     plus ERCC spike-ins. Since expression quantification was being performed at the gene level,

393     overlapping transcripts from the same gene were merged so that each gene contained a single

394     transcript covering all exons from all transcripts. Reads were aligned to the reference transcriptome

395    with Bowtie v1.1.2(Langmead et al. 2009b) with the following options: *"-n1 -l 15 -M 1 --best --strata"*

396    to allow one mismatch, set seed length to 15 bp and report only one alignment where multiple

397    "best" alignments were found. The seed length and mismatch parameters were the same as the

398    Klein *et al* (2015) alignment method.

399    Following alignment, de-duplication was performed with UMI-tools dedup with *unique, percentile*

400    and *directional* used in turn. Both data sets were generated with sequencing methods which

401    generate reads with different alignment coordinates from the same initial DNA fragment (SCRB-seq,

402    CEL-Seq). De-duplication was therefore performed with the *"--per-contig"* option so that the UMI

403    and the contig (in this case, gene) rather than the exact alignment coordinates were used to identify

404    duplicate reads. The "--stats-output" and "--further-stats" options were used to generate summary

405    statistics for the alignment files pre and post de-duplication. Gene expression was quantified by

406    counting the number of remaining reads per gene following de-duplication

407

408    *Exploratory gene expression analysis*

409    PCA was performed in R (R Core Team 2015) using the *prcomp* function. Hierarchical clustering was

410    performed in R using the *hclust* function and heatmaps generated using the *heatmap.2* function

411    from the gplots package. Clustering was performed using 1 - spearman's correlation coefficient as

412    the distance measure and "ward.D2" as the clustering method. Since many genes show very low

413    expression in the SCRB-seq data, the top 100 most highly expressed genes were selected for

414    clustering.

415

416     **Data access**

417     UMI-tools is available from pypi (package: umi_tools) and conda (channel:

418     https://conda.anaconda.org/toms, package:  umi_tools) or GitHub

419     (https://GitHub.com/CGATOxford/UMI-tools). Analyses conducted in this manuscript used version

420     0.2.6 - archived on Zenodo as https://doi.org/10.5281/Zenodo.165403, and in Supplementary File 1.

421     Analyses were performed using automated python pipelines. iCLIP specific analyses were completed

422     using the iCLIPlib python library (manuscript in preparation). Figures were created by python

423     pipelines or in Jupyter notebooks using the ggplot2 package (Wickham 2009) unless otherwise

424     noted. All pipelines, notebooks and other code, along with configuration files used are available

425     from the GitHub repository (https://GitHub.com/CGATOxford/UMI-tools_pipelines), archived on

426     Zenodo as https://doi.org/10.5281/zenodo.215974 and in Supplementary File 2.

427

428

429     **Acknowledgements**

440

441     **Disclosure Declaration**

442     The authors declare that we have no competing interests

443

22

**References**

Aird D, Ross MG, Chen W-S, Danielsson M, Fennell T, Russ C, Jaffe DB, Nusbaum C, Gnirke A. 2011.

Analyzing and minimizing PCR amplification bias in Illumina sequencing libraries. *Genome Biol*

**12**: R18.

Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: A flexible trimmer for Illumina sequence data.

*Bioinformatics* **30**: 2114–2120.

Bose S, Wan Z, Carr A, Rizvi AH, Vieira G, Pe'er D, Sims P a. 2015. Scalable microfluidics for single cell

RNA printing and sequencing. *Genome Biol* **16**: 120.

Collins JE, Wali N, Sealy IM, Morris JA, White RJ, Leonard SR, Jackson DK, Jones MC, Smerdon NC,

Zamora J, et al. 2015. High-throughput and quantitative genome-wide messenger RNA

sequencing for molecular phenotyping. *BMC Genomics* **16**: 578.

Davis MPA, van Dongen S, Abreu-Goodger C, Bartonicek N, Enright AJ. 2013. Kraken: A set of tools

for quality control and analysis of high-throughput sequence data. *Methods* **63**: 41–49.

He Q, Johnston J, Zeitlinger J. 2015. ChIP-nexus enables improved detection of in vivo transcription

factor binding footprints. *Nat Biotechnol* **33**: 395–401.

Hug H, Schuler R. 2003. Measurement of the number of molecules of a single mRNA species in a

complex mRNA preparation. *J Theor Biol* **221**: 615–624.

Islam S, Zeisel A, Joost S, La Manno G, Zajac P, Kasper M, Lönnerberg P, Linnarsson S. 2014.

Quantitative single-cell RNA-seq with unique molecular identifiers. *Nat Methods* **11**: 163–6.

Karlsson K, Sahlin E, Iwarsson E, Westgren M, Nordenskjöld M, Linnarsson S. 2015. Amplification-free

sequencing of cell-free DNA for prenatal non-invasive diagnosis of chromosomal aberrations.

*Genomics* **105**: 150–8.

Kivioja T, Vähärautio A, Karlsson K, Bonke M, Enge M, Linnarsson S, Taipale J. 2012. Counting

   absolute numbers of molecules using unique molecular identifiers. *Nat Methods* **9**: 72–4.

Klein AM, Mazutis L, Akartuna I, Tallapragada N, Veres A, Li V, Peshkin L, Weitz DA, Kirschner MW.

   2015. Droplet Barcoding for Single-Cell Transcriptomics Applied to Embryonic Stem Cells. *Cell*

   **161**: 1187–1201.

König J, Zarnack K, Rot G, Curk T, Kayikci M, Zupan B, Turner DJ, Luscombe NM, Ule J. 2010a. iCLIP

   reveals the function of hnRNP particles in splicing at individual nucleotide resolution. *Nat Struct*

   *Mol Biol* **17**: 909–915.

König J, Zarnack K, Rot G, Curk T, Kayikci M, Zupan B, Turner DJ, Luscombe NM, Ule J. 2010b. iCLIP

   reveals the function of hnRNP particles in splicing at individual nucleotide resolution. *Nat Struct*

   *Mol Biol* **17**: 909–15.

Langmead B, Trapnell C, Pop M, Salzberg SL. 2009a. Ultrafast and memory-efficient alignment of

   short DNA sequences to the human genome. *Genome Biol* **10**: R25.

Langmead B, Trapnell C, Pop M, Salzberg SL. 2009b. Ultrafast and memory-efficient alignment of

   short DNA sequences to the human genome. *Genome Biol* **10**: R25.

Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.

   *Bioinformatics* **25**: 1754–1760.

Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, Goldman M, Tirosh I, Bialas AR, Kamitaki N,

   Martersteck EM, et al. 2015. Highly Parallel Genome-wide Expression Profiling of Individual

   Cells Using Nanoliter Droplets. *Cell* **161**: 1202–1214.

Marinier E, Brown DG, McConkey BJ. 2015. Pollux: platform independent error correction of single

   and mixed genomes. *BMC Bioinformatics* **16**: 10.

McCloskey ML, St??ger R, Hansen RS, Laird CD. 2007. Encoding PCR products with batch-stamps and

barcodes. *Biochem Genet* **45**: 761–767.

Müller-McNicoll M, Botti V, de Jesus Domingues AM, Brandl H, Schwich OD, Steiner MC, Curk T,
Poser I, Zarnack K, Neugebauer KM. 2016. SR proteins are NXF1 adaptors that link alternative
RNA processing to mRNA export. *Genes Dev* **30**: 553–66.

R Core Team. 2015. R: A Language and Environment for Statistical Computing.

Rittié L, Perbal B. 2008. Enzymes used in molecular biology: A useful guide. *J Cell Commun Signal* **2**:
25–45.

Schirmer M, D'Amore R, Ijaz UZ, Hall N, Quince C. 2016. Illumina error profiles: resolving fine-scale
variation in metagenomic sequencing data. *BMC Bioinformatics* **17**: 125.

Schirmer M, Ijaz UZ, D'Amore R, Hall N, Sloan WT, Quince C. 2015. Insight into biases and sequencing
errors for amplicon sequencing with the Illumina MiSeq platform. *Nucleic Acids Res* **43**: 1–16.

Schloss PD, Gevers D, Westcott SL. 2011. Reducing the effects of PCR amplification and sequencing
Artifacts on 16s rRNA-based studies. *PLoS One* **6**.

Schmitt MW, Kennedy SR, Salk JJ, Fox EJ, Hiatt JB, Loeb LA. 2012. Detection of ultra-rare mutations
by next-generation sequencing. *Proc Natl Acad Sci* **109**: 14508–14513.

Shiroguchi K, Jia TZ, Sims PA, Xie XS. 2012. Digital RNA sequencing minimizes sequence-dependent
bias and amplification noise with optimized single-molecule barcodes. *Proc Natl Acad Sci U S A*
**109**: 1347–52.

Sims D, Sudbery I, Ilott NE, Heger A, Ponting CP. 2014. Sequencing depth and coverage: key
considerations in genomic analyses. *Nat Rev Genet* **15**: 121–32.

Soumillon M, Cacchiarelli D, Semrau S. 2014. Characterization of directed differentiation by high-
throughput single-cell RNA-seq. *BioRxiv*.

Tollervey JR, Curk T, Rogelj B, Briese M, Cereda M, Kayikci M, Hortobágyi T, Nishimura AL, Župunski

V, Patani R, et al. 2011. Characterising the RNA targets and position-dependent splicing

regulation by TDP-43; implications for neurodegenerative diseases. *Nat Neurosci* **14**: 452–8.

Vollmers C, Sit R V, Weinstein JA, Dekker CL, Quake SR. 2013. Genetic measurement of memory B-

cell recall using antibody repertoire sequencing. *Proc Natl Acad Sci U S A* **110**: 13463–8.

Waugh C, Cromer D, Grimm A, Chopra A, Mallal S, Davenport M, Mak J. 2015. A general method to

eliminate laboratory induced recombinants during massive, parallel sequencing of cDNA

library. *Virol J* **12**: 55.

Whalen S, Truty RM, Pollard KS. 2016. Enhancer-promoter interactions are encoded by complex

genomic signatures on looping chromatin. *Nat Genet* **48**: 488–96.

Wickham H. 2009. *ggplot2: Elegant Graphics for Data Analysiss*. Springer-Verlag, New York.

Yaari G, Kleinstein SH. 2015. Practical guidelines for B-cell receptor repertoire sequencing analysis.

*Genome Med* **7**: 121.

444

Figure Legends

445     **Figure 1. Modelling errors in UMIs**

446     **A**. Schematic representation of how UMIs are used to count unique molecules. Fragmented DNA is

447     labelled with a random UMI sequence (short oligonucleotide; represented as coloured blocks).

448     Following PCR amplification, sequencing and bioinformatics steps, the sequence read alignment

449     coordinates and UMI sequences are used to identify sequence reads originating from the same initial

450     DNA fragment (PCR duplicates) and so count the unique molecules. **B**. Average edit distances

451     (rounder to integers) between UMIs with the same alignment coordinates. Genomic positions with a

452     single UMI are not shown. Null = Null expectation from random sampling of UMIs, taking into

453     account the genome-wide distribution of UMIs. **C.** Correlation between duplication rate and

454     enrichment of positions with an average edit distance of 1 for iCLIP data. **D**. Topologies of networks

455     formed by joining reads with the same genomic coordinates and UMIs a single edit distance apart.

456     Single hub = One node connected to all other nodes. Complex = No node connected to all other

457     nodes. **E**. Methods for estimating unique molecules from UMI sequences and counts at a single

458     locus. Where the method uses the UMI counts, these are shown. Red bases are inferred to be

459     sequencing errors, blue bases inferred to be PCR errors. The inferred number of unique molecules

460     for each method is shown in parentheses.

461     **Figure 2. Comparison of methods with simulated data**

462     In each panel, all but one of the simulation parameters are held constant, with the remaining

463     parameter varied as shown on the x-axis. **A.** UMI length. **B.** Sequencing depth. **C.** Sequencing error

464     rate. Left plot shows the accuracy of quantification, presented as the log2-transformed normalised

465     difference between the estimate and ground truth. Right plot shows the coefficient of variation

466     (standard deviation / mean). The dashed red line represents the value used for this parameter in all

27

467     other simulations. The dashed grey line represents perfect accuracy. The *unique* and *percentile*

468     methods give identical results with the parameters shown here and are hence overplotted.

469     **Figure 3. UMI-Tools improves reproducibility between iCLIP replicates**

470     **A.** Average edit distances between UMIs with the same alignment coordinates. Genomic positions

471     with a single UMI are not shown. Null = Null expectation from random sampling of UMIs, taking into

472     account the genome-wide distribution of UMIs. Only the first replicate of the dataset is shown for

473     each pull down **B.** iCLIP reproducibility as represented by the percentage of positions with >2 tags

474     also cross-linked in at least one of 2 other replicates. **C.** Spearman's rank correlation between the

475     numbers of significant tags in each exon

476     **Figure 4. UMI tools improves accuracy and *cluster*ing in Single Cell RNA-seq**

477     **A.** Average edit distances between UMIs with the same alignment coordinates following removal of

478     PCR duplicates using the methods indicated on the x-axis. Genomic positions with a single UMI are

479     not shown. Null: Null expectation from random sampling of UMIs, taking into account the genome-

480     wide distribution of UMIs. Top = SCRB-seq. Bottom = inDrop-seq. **B.** Distribution of pearson

481     correlation coefficients between log ERCC concentration and log counts for raw reads (UMIs

482     ignored) and *unique* and *directional* methods. C & D. Hierarchical clustering based on the gene

483     expression estimates obtained using *unique* and *directional* Colour bars represent differentiation

484     stage. **C**. SCRB-seq. **D.** inDrop-DSeq. Red arrow indicates mES Cells clustering with Day 4 cells.

28

# UMI-tools: Modelling sequencing errors in Unique Molecular Identifiers to improve quantification accuracy

Tom Sean Smith, Andreas Heger and Ian Sudbery

| | |
|---|---|
| **P<P** | Published online January 18, 2017 in advance of the print journal. |
| **Accepted Manuscript** | Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version. |
| **Open Access** | Freely available online through the *Genome Research* Open Access option. |
| **Creative Commons License** | This manuscript is Open Access.This article, published in *Genome Research*, is available under a Creative Commons License (Attribution 4.0 International license), as described at http://creativecommons.org/licenses/by/4.0/. |
| **Email Alerting Service** | Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or **click here.** |

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.