

1 Haemolysis detection in microRNA-seq from clinical  
2 plasma samples

3  
4 Melanie D. Smith<sup>1,2\*</sup>, Shalem Y. Leemaqz<sup>1</sup>, Tanja Jankovic-Karasoulos<sup>1</sup>, Dale  
5 McAninch<sup>2</sup>, Dylan McCullough<sup>1</sup>, James Breen<sup>3,4</sup>, Claire T. Roberts<sup>1,2</sup>, Katherine A.  
6 Pillman<sup>5\*</sup>

7  
8 <sup>1</sup> Flinders Health and Medical Research Institute, Flinders University, Bedford Park, SA,  
9 Australia

10 <sup>2</sup> Adelaide Medical School, University of Adelaide, Adelaide, SA, Australia

11 <sup>3</sup> Indigenous Genomics, Telethon Kids Institute, Adelaide, SA, Australia

12 <sup>4</sup> College of Health & Medicine, Australian National University, Canberra, ACT, Australia

13 <sup>5</sup> Centre for Cancer Biology, University of South Australia/SA Pathology, Adelaide, SA,  
14 Australia

15

16

17 \* Correspondence: [melanie.smith@flinders.edu.au](mailto:melanie.smith@flinders.edu.au), [katherine.pillman@unisa.edu.au](mailto:katherine.pillman@unisa.edu.au)

18

19

20

## 21 Abstract

22 The abundance of cell-free microRNA (miRNA) has been measured in many body  
23 fluids, including blood plasma, which has been proposed as a source with novel,  
24 minimally invasive biomarker potential for several diseases. Despite improvements in  
25 quantification methods for plasma miRNAs, there is no consensus on optimal reference  
26 miRNAs or to what extent haemolysis may affect plasma miRNA content. Here we  
27 propose a new method for the detection of haemolysis in miRNA high-throughput  
28 sequencing (HTS) data from libraries prepared using human plasma. To establish a  
29 miRNA haemolysis signature in plasma we first identified differentially expressed  
30 miRNAs between samples with known haemolysis status and selected miRNA with  
31 statistically significant higher abundance in our haemolysed group. Given there may be  
32 both technical and biological reasons for differential abundance of signature miRNA,  
33 and to ensure the method developed here was relevant outside of our specific context,  
34 that is women of reproductive age, we tested for significant differences between  
35 pregnant and non-pregnant groups. Here we report a novel 20 miRNA signature (miR-  
36 106b-3p, miR-140-3p, miR-142-5p, miR-532-5p, miR-17-5p, miR-19b-3p, miR-30c-5p,  
37 miR-324-5p, miR-192-5p, miR-660-5p, miR-186-5p, miR-425-5p, miR-25-3p, miR-363-  
38 3p, miR-183-5p, miR-451a, miR-182-5p, miR-191-5p, miR-194-5p, miR-20b-5p) that  
39 can be used to identify the presence of haemolysis, *in silico*, in high throughput miRNA  
40 sequencing data. Given the potential for haemolysis contamination, we recommend that  
41 assay for haemolysis detection become standard pre-analytical practice and provide  
42 here a simple method for haemolysis detection.

## 43 Introduction

44 MicroRNAs represent a class of short, ~22 nt single stranded non-coding RNA  
45 transcripts found in the cytoplasm of most cells that act directly as post transcriptional  
46 regulators of gene expression (1,2) and also coordinate extensive indirect  
47 transcriptional responses (3). In their canonical action, miRNAs mediate the expression  
48 of specific messenger RNA (mRNA) targets by binding to the 3'-untranslated region  
49 (UTR) of transcripts by either repressing translation or marking them for degradation (4).

50 In the canonical miRNA pathway, target specificity requires exact nucleotide sequence  
51 complementarity between the miRNA 'seed' region (the first 2-7 bases at the 5` end of  
52 the mature miRNA transcript) and the 3`-UTR of the mRNA. Importantly, miRNAs  
53 demonstrate tissue, temporal and spatial expression specificity and are known  
54 regulators of development, with most mammalian mRNAs harbouring conserved targets  
55 of one or many miRNAs (1,2,5).

56  
57 miRNA expression is both temporally and spatially tissue-specific, with transcripts  
58 identified beyond the cells in which they were synthesised, in various body fluids  
59 including urine, saliva and blood plasma (6). Circulating cell-free miRNAs identified in  
60 plasma are packaged in micro vesicles such as exosomes (7,8) or bound to protein  
61 complexes such as argonaute 2 (Ago2), nucleophosmin 1 (NPM 1) and high density  
62 lipoprotein (HDL) (9–11), making them exceptionally stable (6). This stability, coupled  
63 with their minimally invasive accessibility, has suggested circulating cell-free miRNAs as  
64 an important resource for the identification of novel biomarkers.

65  
66 Whilst much progress has been made in the search for novel miRNA biomarkers of  
67 disease processes (12–14), outcomes of this research approach are often inconsistent  
68 or even contradictory (6). There are many reasons for this, including variations in  
69 enrichment, extraction and quantification methods, variation between individuals, lack of  
70 consensus regarding optimal reference miRNA for normalisation and the difficulty in  
71 quantifying both the amount and quality of RNA transcripts from blood plasma samples  
72 (15,16). An important but often overlooked factor is the potential for sample haemolysis  
73 during blood collection or sample preparation which results in miRNA from lysed RBCs  
74 being spilled into and retained within the plasma sample to be assayed (15).

75  
76 The issue of haemolysis altering the miRNA content of plasma and the potential for  
77 confounding biomarker discovery has been reported previously (15,17,18). Using RT-  
78 qPCR Kirschner and colleagues (15) showed that contamination of plasma samples  
79 with the miRNA content of RBCs changed the abundance of both miR-16 and miR-  
80 451a. This, in turn, altered the relative abundance of potential biomarkers for

81 mesothelioma and coronary artery disease including miR-92a and miR-15. Using the  
82 same technique, Pritchard *et al.* (18) demonstrated in plasma that 46 of the 79  
83 circulating miRNA cancer biomarkers were highly expressed in more than one blood cell  
84 type, noting that the effects of sample specific blood cell counts and haemolysis can  
85 alter the miRNA biomarker levels in a single patient sample up to 50-fold. As a result,  
86 the authors emphasised caution in classifying blood cell associated miRNAs as  
87 biomarkers given the possible alternate interpretation.

88  
89 Haemolysis is associated with either blood collection or RNA extraction and sample  
90 preparation. Thus, despite differences between the quantification methods, high  
91 throughput sequencing data used in our study is equally susceptible to the confounding  
92 effects of sample haemolysis on miRNA abundance levels in plasma is RT-qPCR.  
93 There are currently two gold standard approaches in the assessment of haemolysis in  
94 plasma: 1. Delta quantification cycle ( $\Delta Cq$ ), where expression levels of a known blood  
95 cell associated miRNA (miR-451a) and a control miRNA (miR-23a) are determined  
96 based on the difference between the two raw Cq values and 2. Spectrophotometry,  
97 where absorbance is measured at 414 nm with the use of a spectrophotometer. In the  
98 case of  $\Delta Cq$  assessment, miR-451a is known to vary and miR-23a is known to be  
99 invariant in plasma affected by haemolysis (15,16). Using spectrophotometry,  
100 haemolysis is quantified by assessing the presence of cell free haemoglobin by  
101 measuring the absorbance at 414 nm, the absorbance maximum of free haemoglobin  
102 (19,20). Both methods require access to sufficient amounts of the original plasma  
103 sample and the laboratory equipment required to perform the assays. Free access to a  
104 web-tool that can perform *in silico* assessment of RBC contamination in human plasma  
105 would be of exceptional value to the research community.

106  
107 Whilst it is well established that haemolysis frequently occurs during extraction or  
108 processing of blood samples, the assessment of RBC contamination is rarely mentioned  
109 in publications. It is even more rare that the results of any such testing are present in  
110 the metadata assigned to publicly available sequencing data. There is currently no  
111 publicly available tool for analysis of haemolysis without access to the physical plasma

112 specimen. Although the theory underlying identification of haemolysis in plasma is  
113 relatively straight-forward, surprisingly, to our knowledge this has never before been  
114 extrapolated into a data-only, *in silico* approach. The paucity of haemolysis information  
115 in the context of publicly available datasets, combined with the lack of tools to identify  
116 affected datasets after the fact, substantially limits the utility of this data resource and  
117 reproducibility of research findings. Further, it increases the risk that results obtained  
118 may unwittingly represent blood-cell based phenomena rather than signatures of the  
119 pathology of interest.

120  
121 In this study, we assessed miRNA abundance in HTS data from libraries prepared using  
122 human plasma from pregnant and non-pregnant women of reproductive age. Using a  
123 set of samples with confirmed haemolysis ( $\Delta Cq$  (miR-23a-miR-451a)), we established a  
124 set of 20 miRNAs differentially abundant between plasma from samples with and  
125 without substantiated haemolysis. Using the expression values of these 20 miRNAs as  
126 a 'signature' of haemolysis, we calculated the difference between the mean normalised  
127 expression levels of these miRNAs compared to those of all other miRNAs (as a  
128 'background' set). This produced a quantitative metric which represents the strength of  
129 the evidence of haemolysis in an individual sample. When this metric is interpreted in  
130 the context of other samples, it can be used to identify sample(s) that display substantial  
131 evidence of haemolysis. The researcher may consider discarding these samples from  
132 further analyses or using caution in their interpretation. We consulted the EMBL-EBI  
133 Expression Atlas ([ebi.ac.uk](http://ebi.ac.uk)) to ensure all signature miRNAs are identified in multiple  
134 human tissues (male and female) and have no known developmental stage association.  
135 For ease of application, we have developed this method into a web based Shiny/R  
136 application, DraculR (a tool that allows a user to upload and assess haemolysis in high-  
137 throughput plasma miRNA-seq data), for use by the research community (DraculR: A  
138 web-based application for *in silico* haemolysis detection in high throughput small RNA  
139 sequencing data).

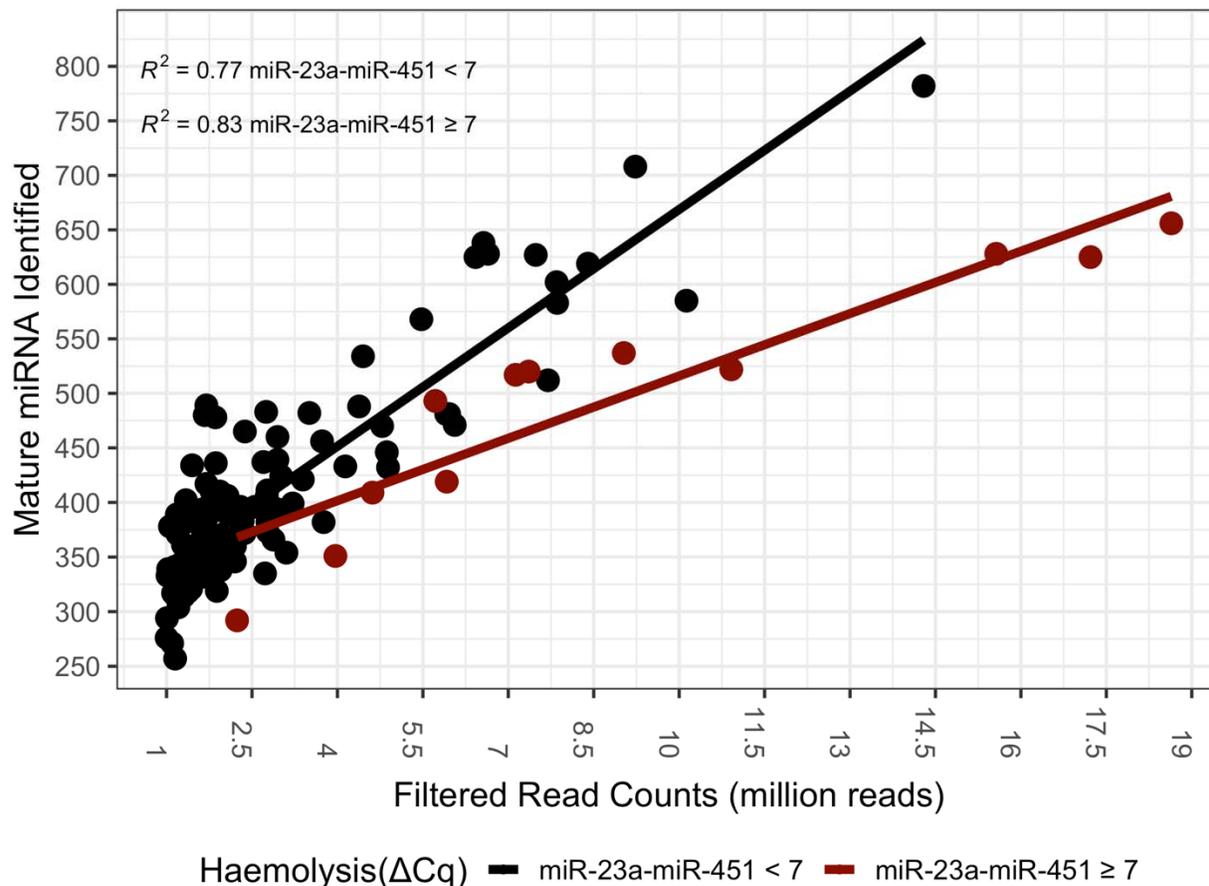
## 140 Results

### 141 High throughput sequencing

142 Illumina NextSeq 75 bp single-end read sequencing was performed on miRNA libraries  
143 from 154 plasma samples taken from 24 non-pregnant and 130 pregnant women aged  
144 16 to 46 years (Qiagen, Hilden, Germany). Prior to sequencing, RT-qPCR was used to  
145 analyse  $\Delta Cq$  (miR-23a-miR-451a), where the ratio of miR-23a to miR-451a (or  $\Delta Cq$   
146 (miR-23a-miR-451a)  $\geq 7$ ) correlates with the degree of haemolysis. We identified 14  
147 plasma samples with a  $\Delta Cq$  of 7 or above (Supplementary Table 1). An average of ~2.9  
148 million reads were sequenced per sample (range ~0.25-18.6 million reads). Thirty-one  
149 libraries with  $< 1$  million reads were considered to be unreliable due to low sequencing  
150 output and were removed from further analyses. There was no difference in the  
151 proportion of haemolysed and non-haemolysed data in the exclusion of samples due to  
152 low library size (Fisher's exact test P-value = 0.7). Sequence alignment was performed  
153 using BWA (21) to the human genome (version GRCh38) and miRNA read counts were  
154 generated by mapping to human miRBase v22 (22,23) identifying 1,133 mature  
155 miRNAs.

156  
157 To analyse the effects of haemolysis on miRNA expression data from next generation  
158 sequencing, we first determined the number of unique mature miRNAs identified in each  
159 of our samples and analysed the data relative to read depth. Using an analysis of  
160 variance (ANOVA) we identified a significant difference between the haemolysed and  
161 non-haemolysed samples ( $P < 0.05$ ), with haemolysis being frequently associated with  
162 fewer mature miRNA species detected at a given read depth (Figure 1).

163



164  
165 **Figure 1.** The number of mature miRNA species identified in an individual sample increases  
166 with read depth for both haemolysed and non-haemolysed samples. However, the number of  
167 mature miRNA species identified for a given read depth is significantly lower (ANOVA, P-value  
168 =  $1.68 \times 10^{-9}$ ) in samples affected by haemolysis (maroon) when compared to a non-  
169 haemolysed sample (black) of equal read depth.

## 170 microRNA haemolysis signature set

171 To ensure that miRNAs identified here were representative of those found in a broad set  
172 of plasma samples, we first filtered to discard miRNAs of low abundance. After filtering,  
173 189 highly abundant miRNA remained. Differential expression analysis comparing  
174 miRNA read counts identified 138 miRNAs with a higher abundance in haemolysed  
175 compared to non-haemolysed samples (statistically significant differentially expressed  
176 miRNA, false discovery rate (FDR) < 0.05, with a  $\log_2$  fold change ( $\log_2FC$ ) > 0)  
177 (Supplementary Figure 1; Supplementary Figure 2a & b). We further ranked the  
178 differentially expressed miRNAs based on  $\log_2FC$ , FDR and abundance levels and  
179 subset the list such that only miRNAs which had a  $\log_2FC$  > 0.9 and were in the top 60

180 percent of each of the FDR and abundance rank criteria remained. This resulted in a  
181 high confidence set of 20 miRNA indicative of a haemolysis signature (Table 1).

182  
183 For *in silico* assay of haemolysis in our data, we further removed miRNAs associated  
184 with pregnancy to avoid confounding miRNA associated with haemolysis with those  
185 associated with pregnancy. Differential expression analysis of miRNA read counts from  
186 pregnancy and non-pregnancy samples identified 127 miRNAs (FDR < 0.05) that were  
187 significantly differentially expressed between the groups (Supplementary Figure 3a & b).  
188 Strikingly, one of our first observations highlighted the importance of including  
189 haemolysis analysis as an adjunct in our study: miR-451a, which is the sole haemolysis  
190 signature miRNA used in the current  $\Delta Cq$  (miR-23a-miR-451a) gold standard method  
191 for haemolysis detection was discovered to be highly correlated with pregnancy status,  
192 indicating a strong confounding factor in pregnancy studies when haemolysis levels are  
193 estimated using RT-qPCR alone. Accordingly, miR-451a was removed from calculations  
194 hereafter along with 9 other miRNAs that were differentially expressed between the  
195 pregnant and non-pregnant groups from the core set of haemolysis signature miRNA.  
196 This resulted in 10 miRNAs remaining for evaluation of haemolysis levels.

197

198 **Table 1.** 20 miRNAs with a general-use plasma haemolysis signature set. To remove  
 199 confounding effects within our pregnancy-specific dataset, we identified a subset of 10 abundant  
 200 miRNA which are invariant with respect to pregnancy.  
 201

miRNA	Log <sub>2</sub> FC	Average Expression (log <sub>2</sub> CPM)	Adjusted P-value	Pregnancy Assoc.
miR-106b-3p	1.589	8.731	8.61 x 10 <sup>-15</sup>	no
miR-140-3p	1.073	10.098	2.75 x 10 <sup>-13</sup>	no
miR-142-5p	0.962	10.651	4.96 x 10 <sup>-12</sup>	no
miR-532-5p	1.288	7.237	4.96 x 10 <sup>-12</sup>	no
miR-17-5p	0.952	7.892	7.84 x 10 <sup>-12</sup>	no
miR-19b-3p	1.128	8.696	1.93 x 10 <sup>-09</sup>	no
miR-30c-5p	0.95	7.325	2.48 x 10 <sup>-09</sup>	no
miR-324-5p	1.304	7.186	2.50 x 10 <sup>-09</sup>	no
miR-192-5p	0.941	8.944	1.37 x 10 <sup>-08</sup>	no
miR-660-5p	1.305	7.62	3.45 x 10 <sup>-10</sup>	no
miR-186-5p	1.228	8.052	2.75 x 10 <sup>-13</sup>	yes
miR-425-5p	1.282	11.246	4.96 x 10 <sup>-12</sup>	yes
miR-25-3p	1.212	12.939	1.26 x 10 <sup>-11</sup>	yes
miR-363-3p	1.237	7.882	4.52 x 10 <sup>-11</sup>	yes
miR-183-5p	1.55	9.382	9.34 x 10 <sup>-11</sup>	yes
miR-451a	1.372	13.002	3.65 x 10 <sup>-10</sup>	yes
miR-182-5p	1.341	10.585	2.48 x 10 <sup>-09</sup>	yes
miR-191-5p	0.929	11.79	4.68 x 10 <sup>-09</sup>	yes
miR-194-5p	0.937	7.679	1.85 x 10 <sup>-08</sup>	yes
miR-20b-5p	0.932	7.43	1.96 x 10 <sup>-08</sup>	yes

202  
 203 Incorporating concepts from previous RT-qPCR analyses of haemolysis, we established  
 204 a new measure of the inclusion of RBC associated miRNA in human plasma. After  
 205 establishing the 20-miRNA signature associated with RBC content inclusion, we  
 206 determined the geometric mean of the distribution of miRNA read counts as an  
 207 appropriate measure of abundance and summary statistic. Using this summary statistic,  
 208 our method calculates a ‘Haemolysis Metric’, defined as the difference between the

209 geometric means of the normalised abundance levels of the haemolysis miRNA  
210 signature set compared to that of all other miRNAs (the ‘background’ set). Note that in a  
211 case-control study, to reduce the risk of confounding the Haemolysis Metric with  
212 experimental variables, the signature set should be reduced to exclude any miRNA  
213 known to be differentially expressed between groups. In this case, the geometric mean  
214 of the reduced signature set will be calculated, as defined in (1).

215

216 Let

217  $Z_x$  be the miRNA Reduced signature set ( $\log_2$  CPM counts)

218 and  $Z_y$  be the Background miRNA set ( $\log_2$  CPM counts)

219 where  $x = 1, 2, 3, \dots, p_1$  with  $p_1 =$  the number of miRNA in Reduced signature set

220 and  $y = 1, 2, 3, \dots, p_2$  where  $p_2 =$  the number of miRNA in Background

221 and  $i = 1, 2, 3, \dots, n$  where  $n =$  the sample size after filtering

222

$$Haemolysis\ Metric_i = \sqrt[p_1]{\prod_{x=1}^{p_1} Z_{x_i}} - \sqrt[p_2]{\prod_{y=1}^{p_2} Z_{y_i}}$$

223

224

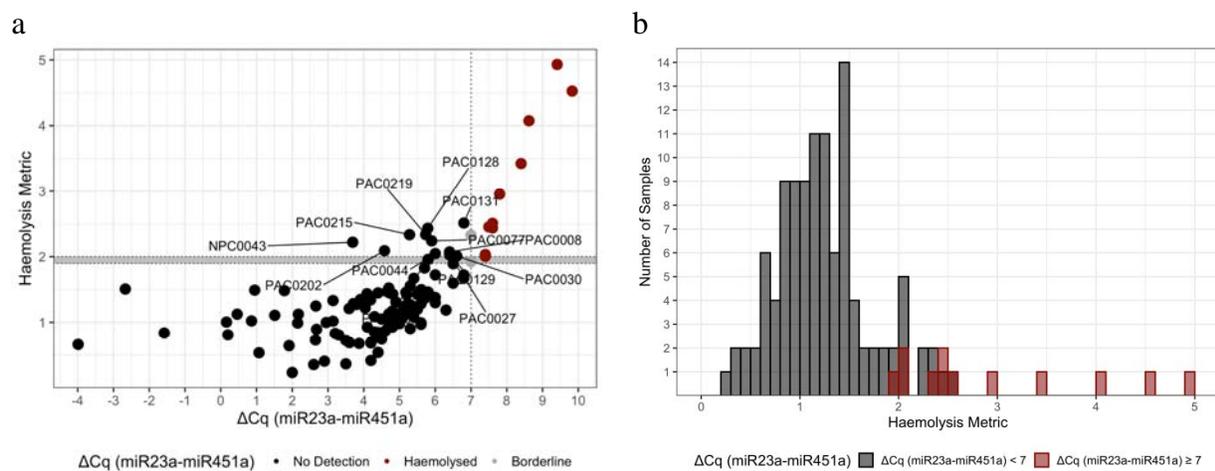
225

226

(1)

227 Prior to establishing a threshold for the new Haemolysis Metric we measured the linear  
228 dependence between the new Haemolysis Metric and the  $\Delta Cq$  (miR-23a-miR-451a)  
229 metric by performing a Pearson’s correlation. Our results indicated a Pearson’s  
230 correlation coefficient of 0.64 ( $P < 0.0001$ ). With confidence in the correlation, to  
231 establish a threshold for the Haemolysis Metric, we compared the results of the  $\Delta Cq$   
232 (miR-23a-miR-451a) and summary statistic methods directly. Briefly, we compared the  
233 Haemolysis Metric to the  $\Delta Cq$  (miR-23a-miR-451a) results for matched samples and  
234 established a cut off criterion for inclusion into the Clear (no haemolysis detected) and  
235 Caution (haemolysis detected) groups (Figure 2a). We chose a threshold of  $\geq 1.9$  for the  
236 assignment of “Caution” to individual samples based on the minimum summary statistic  
237 difference of samples assayed using the  $\Delta Cq$  (miR-23a-miR-451a) metric of  $\geq 7$  (Figure  
238 2a) and the minimal overlap between the distribution of the Haemolysis Metric in  
239 haemolysed compared to non-haemolysed samples (Figure 2b). Where a sample is

240 assigned “Caution”, researchers are advised to consider removing the sample, or  
241 continue with caution. Given the correlation of the two metrics is imperfect and the  
242 arbitrary nature of choosing any cut-off, samples with a Haemolysis Metric close to the  
243 1.9 cut-off may be interrogated further prior to any decision to retain or remove. Of the  
244 121 samples assayed, 25 samples met the criteria for Caution. Of these, 12 were  
245 previously determined as haemolysed or borderline using the  $\Delta Cq$  (miR-23a-miR-451a)  
246 assay. We found that all samples identified as  $\Delta Cq \geq 7$  (Figure 2a, scarlet) are above  
247 the criteria for the Haemolysis Metric (Figure 2a, horizontal grey bar; threshold  $\geq 1.9$ ).  
248 Further, we identified 13 samples with a Haemolysis Metric  $\geq 1.9$  not included in the  
249  $\Delta Cq$  (miR-23a-miR-451a) criteria.



250 **Figure 2:** (a) A comparison of the derived Haemolysis Metric and the  $\Delta Cq$  measure of  
251 haemolysis shows a clear correlation. We identified 13 samples (named) that we suggest  
252 should be discarded or used with caution in further analysis. (b) Histogram of Haemolysis Metric  
253 values from the 121 samples in our experiment, coloured according to their  $\Delta Cq$  (miR-23a-miR-  
254 451a) classification indicate a minimum Haemolysis Metric of  $\geq 1.9$  for samples previously  
255 identified as haemolysed.

## 256 Discussion

257 Through an analysis of differential miRNA expression in samples whose haemolysis  
258 levels were known, we identified a novel 20 miRNA signature indicative of haemolysis.  
259 Given our hypothesis that plasma samples contaminated with RBC content would  
260 contain proportionally higher levels of many RBC-associated miRNA, not just miR-451a,  
261 we established a method using a group of background miRNAs as a reference.  
262 Accordingly, as a group, signature miRNAs (microRNAs which are abundant in red

263 blood cells) are shown to be more highly abundant in samples contaminated with RBCs.  
264 The degree of this change can be used as a measure of RBC content contamination  
265 and quantified by comparing the geometric means of the expressions of RBC signature  
266 miRNAs to that of the background set of miRNAs. We further established that where a  
267 comparison between conditions is considered, for example in a biomarker discovery  
268 experiment, any miRNA known to be associated with the condition for which the  
269 biomarker is proposed should be removed to prevent confounding between the  
270 condition of interest and the quantification of RBC-associated miRNA inclusion.

271  
272 Our experimental results demonstrate that it is possible to identify a haemolysis  
273 signature *in silico*, avoiding the effort and expense of lab validation, and in situations  
274 where blood plasma samples are exhausted, otherwise unavailable or cost-prohibitive  
275 to assay using current gold standard approaches. Given the limited access to physical  
276 samples associated with publicly available data, the Haemolysis Metric technique  
277 introduced here, provides the basis for the development of a publicly available tool  
278 (DraculR: A web-based application for *in silico* haemolysis detection in high throughput  
279 small RNA sequencing data).

280  
281 Among the haemolysis miRNA signature is miR-451a (previously named miR-451),  
282 commonly associated with RBC contamination and used in the calculation of  $\Delta Cq$  (miR-  
283 23a-miR-451a). However, we removed miR-451a together with 9 other miRNAs from  
284 our calculation of distribution difference due to changes in miRNA abundance  
285 associated with pregnancy. During pregnancy total blood volume increases, varying  
286 between 20% to 100% above pre-pregnancy levels. This change, however, is not  
287 uniform across all blood components as plasma volume increases proportionally more  
288 than the RBC mass (24). This is an important consideration and highlights the limitation  
289 of the current gold standard approach that uses two miRNAs rather than a larger  
290 signature set to calculate a measure of contamination. If, as in this example, the  
291 abundance of either miRNA used to determine the  $\Delta Cq$  is also affected by the condition  
292 or pathology under investigation, the issue is two-fold. Firstly, you may identify miR-  
293 451a as being differentially abundant in the pathology of interest and propose its use as

294 a biomarker only to find that it is confounded by haemolysis. Secondly, you may, using  
295 the  $\Delta Cq$  calculation, classify samples as haemolysed when the change in miR-451a  
296 abundance is more appropriately associated with the pathology of interest. By  
297 establishing a larger signature set of miRNAs to detect haemolysis in small RNA  
298 sequencing from human plasma we hope to provide a resource to the community. To  
299 overcome issues identified in previous studies, the flexibility and redundancy included in  
300 our metric buffer against the issue of confounding conditions of interest with the  
301 measure of haemolysis.

302  
303 We found limited overlap between the miRNAs identified as useful for the detection of  
304 haemolysis and those previously reported as markers of haemolysis contamination  
305 (15,18,25). When comparing with previous research, it is important to note, that our  
306 research question differed from that of the above studies, as did our study methodology.  
307 The most important technical difference is in the quantification of miRNAs. The  
308 expression values used here were taken from a HTS experiment, rather than RT-qPCR  
309 used previously. The limitations of RT-qPCR to investigate which miRNAs are affected  
310 by haemolysis has been identified previously (26). Given that HTS allows for  
311 quantification of all known miRNA species and that RT-qPCR is targeted, our  
312 experiment was able to identify differential abundance in miRNAs not quantified in  
313 Kirschner *et al.* (15), Pritchard *et al.* (18) or McDonald *et al.* (25). Whilst we identified an  
314 overlap in the miRNAs associated with haemolysis in this and previous work, many of  
315 these were not included in the final miRNA Haemolysis Metric signature set. These  
316 include miR-16, miR-486-5p and miR-92a-3p which were significantly upregulated in the  
317 haemolysed group, but not included in the signature set as they failed to pass filtering  
318 criteria for  $\log_2FC$  and expression level. Secondary to the technical differences  
319 introduced by using different miRNA quantification technologies, it is important to note  
320 that all plasma samples used here to establish which miRNAs are affected by  
321 haemolysis were taken from adult women of reproductive age. No sex or age  
322 information was included with either of the compared studies, although it is likely these  
323 samples included specimens from men and women. To account for the potential bias  
324 introduced using data from only female and all reproductively aged volunteers, we

325 ensured all miRNAs included here have previously been identified in multiple tissues in  
326 both male and female samples and are not affected by developmental stage.  
327 Investigation using a cohort of mixed age and sex is warranted and may help to further  
328 determine which miRNAs are affected by haemolysis.

329  
330 Interestingly, all signature miRNAs, with the exception of miR-325-5p, have previously  
331 been reported as prognostically valuable plasma or serum biomarkers. In this small  
332 sampling of recent miRNA biomarker research, we identified several instances where  
333 more than one of our haemolysis signature miRNA were identified as disease  
334 biomarkers for the same condition in the same experiment (27–29) which, given our  
335 findings, and those of previous haemolysis research, further call into question their  
336 validity as biomarkers of disease or condition. In conjunction with our research, we  
337 found many miRNAs as suggested circulating biomarkers for multiple disease states.  
338 For example, miR-122 was given biomarker potential in liver disease, lung cancer and  
339 myasthenia gravis (14,27,30), and miR-660 was given biomarker potential in  
340 Alzheimer's disease, breast cancer and lung cancer (28,31,32), respectively. These  
341 miRNAs may represent effective biomarkers, but they may simply highlight RBC  
342 contamination or be indicative of a general state of inflammation.

343  
344 Data contained in this study were obtained from two cohorts of female volunteers of  
345 reproductive age. Whilst we are working in a relatively narrow experimental domain, we  
346 have generalised this method such that removal (from the signature miRNA set) of  
347 domain specific miRNA is built in, providing a framework that allows use within research  
348 conducted in any human plasma context. Our results highlight that ignoring the issue of  
349 miRNA from RBCs leaves researchers open to the risk that newly discovered miRNA  
350 disease biomarkers could in fact be biomarkers of haemolysis. In future research, a  
351 repeat experiment with samples taken from male and female individuals across a wider  
352 age range would expand and strengthen our understanding of the impact of haemolysis  
353 on biomarker discovery. Our research both recommends and enables tests for  
354 haemolysis to become standard pre-analytical practice.

## 355 Methods

### 356 Sample collection

357 Peripheral blood (9 mL) was collected with informed, written consent from women  
358 undergoing elective terminations of otherwise healthy pregnancies. Blood was collected  
359 into standard EDTA blood tubes pre-termination and stored on ice until processed.  
360 Whole blood underwent centrifugation at 800 x g for 15 minutes at 4°C before plasma  
361 removal and then spun for a further 15 minutes to ensure any remaining cellular debris,  
362 including cell membranes from lysed red blood cells, was removed. All samples were  
363 stored at -80°C until further processing. Termination samples were collected from the  
364 Pregnancy Advisory Centre (PAC), Woodville, South Australia. Blood was also collected  
365 with informed, written consent from non-pregnant volunteers at the Adelaide Medical  
366 School. Following collection, blood tubes were stored on ice until processing. Whole  
367 blood underwent centrifugation at 1015 x g for 10 minutes at 4°C. Approximately 4-6 mL  
368 plasma was collected in 2 mL aliquots. 500 µL plasma (the supernatant) were aliquoted  
369 into clean tubes and the pellet containing any remaining blood cells at the bottom of the  
370 tube was discarded. All samples were stored at -80°C until further processing.

### 371 RNA extraction and library preparation

372 MicroRNA was isolated from 200 µL plasma samples using the QIAGEN miRNA  
373 serum/plasma kit (Qiagen, Hilden, Germany) according to the manufacturer's  
374 instructions and stored at -80°C.

375 Library preparation and sequencing was performed by Qiagen (Valencia, CA) using the  
376 QIAseq miRNA Library Kit and QIAseq miRNA 48 Index IL kits as per manufacturer's  
377 instructions. Amplified cDNA libraries underwent single-end sequencing by synthesis  
378 (Illumina v1.9).

## 379 Haemolysis Detection by RT-qPCR

380 Plasma samples were examined for haemolysis based on the expression levels of two  
381 miRNAs: miR-451a and miR-23a. miR-451a (previously named miR-451) is known to be  
382 highly expressed in red blood cells, whereas miR-23a is known to maintain stable  
383 abundance levels in plasma. After RNA extraction and cDNA synthesis, the delta  
384 quantification cycle (Cq) values for miR-23a-miR-451a were calculated independently  
385 for each sample. The evaluation of expression levels was performed based on raw Cq  
386 values. According to the Qiagen protocol for haemolysis detection (Qiagen, Hilden,  
387 Germany) using the  $\Delta\Delta Cq$  method, samples with a  $\Delta Cq < 7$  for these two miRNAs were  
388 considered as clear of contamination; a  $\Delta Cq > 7$  was considered contaminated; a  $\Delta Cq$   
389  $= 7$  was considered borderline.

## 390 miRNA annotation and abundance

391 Read quality control metrics were assessed using FastQC (33)  
392 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) to check for per base  
393 sequence quality, sequence length distribution and duplication levels. Adapter detection  
394 and trimming were performed using Atropos (34). Alignment performed using BWA  
395 version 0.7.17-r1188 (GRCh38) (21). Umi\_tools was used to collapse duplicate reads  
396 mapped to the same genomic location with the same UMI barcode. Quality control  
397 metrics were reported using multiQC (35). Read counts for mature miRNAs were  
398 determined using an in-house script (36) with microRNA annotation from miRBase,  
399 version 22.0 (22,23) (<http://www.mirbase.org>).

## 400 Analysis of potential confounding factors

401 All profile and expression analyses were conducted in the R statistical environment  
402 (v.4.0.2), using the *edgeR* (v.3.16.5) (37) and *limma* (v.3.30.11) (38) R/Bioconductor  
403 packages. Prior to conducting the differential expression analysis between haemolysed  
404 and non-haemolysed expression data we considered the effect of participant  
405 characteristics such as sex, age, smoking, pregnancy status and ethnicity. Sex was not

406 included here as all samples were taken from female participants. Maternal age was  
407 excluded from the final regression model as there was no strong evidence of  
408 association with the outcome and hence considering the sample size a simpler model  
409 was chosen to preserve degrees of freedom. miRNA identified as differentially  
410 expressed between samples from pregnant and non-pregnant women, were removed  
411 from the final set of haemolysis signature miRNA. Sequencing batch was also included  
412 in all regression models.

### 413 Identification of haemolysis miRNA signature

414 Prior to defining a collection of haemolysis informative miRNAs, pre-filtering steps were  
415 undertaken: 1) mature miRNA with fewer than five reads were reduced to zero  
416 independently for each sample, 2) miRNA with fewer than 40 counts per million (CPM)  
417 in the haemolysed group ( $n = 12$ ) were removed from further consideration. This was  
418 done to ensure only highly abundant miRNA likely to be present in most samples  
419 remain. The Trimmed Means of M values (TMM) normalisation method was used to  
420 correct for differences in the underlying distribution of miRNA expression (39). Next, we  
421 used *limma* (38) to obtain the fold change of each miRNA between the haemolysed and  
422 non-haemolysed groups to identify miRNAs that are more abundant in the plasma  
423 affected by haemolysis. To ensure the haemolysis miRNA signature was robust, we  
424 took the intersection of the 60 miRNAs from each category of highest expression and  
425 lowest adjusted P-value and miRNAs with a  $\log_2FC > 0.9$ , revealing a set of twenty high  
426 confidence miRNAs. To further refine the set of haemolysis informative miRNAs we  
427 used *limma* to calculate the fold change for each miRNA between the samples from  
428 pregnant and non-pregnant women not affected by haemolysis and removed any of the  
429 high-confidence miRNA which was also differentially abundant in pregnancy. The  
430 workflow, source code and input files associated with this research are available at  
431 ([https://github.com/mxhp75/haemolysis\\_maternaPlasma.git](https://github.com/mxhp75/haemolysis_maternaPlasma.git)).

### 432 Classification - Haemolysis Metric

433 To classify the data coming from samples as haemolysed, borderline or unaffected we  
434 first focused on samples from the non-pregnant group. For these, we subset the miRNA

435 read count table into miRNA from the high confidence haemolysis informative miRNA  
436 (n=20) and all others (n=169). Using this data partition we calculated the geometric  
437 mean of the distribution of read counts using the *psych* package (v1.8.12) (40) and  
438 subtracted the geometric mean of the counts of “other” miRNA from that of the  
439 “haemolysis informative” miRNA. Next, for samples from the pregnant group, we  
440 performed the same calculations described above after first discarding miRNA which  
441 were associated with pregnancy.

## 442 Ethics

443 Ethics approval for the collection of blood from women undergoing elective pregnancy  
444 termination between 6–23 weeks’ gestation was provided under HREC/16/TQEH/33, by  
445 The Queen Elizabeth Hospital Human Research Ethics Committee (TQEH/LMH/MH).  
446 Blood from women forming the general population group was collected after informed  
447 consent with ethics approval provided under HREC/H/021/2005, by The University of  
448 Adelaide Human Research Ethics Committee.

## 449 Author contribution

450 CTR created the concept and acquired funding. MS, KP and JB conceived and  
451 developed experimental plans. TJK, DMcA and DMcC performed experiments. MS and  
452 KP analysed the sequencing data with statistical support from SYL. Manuscript written  
453 by MS and KP. All authors read and approved the final version of the manuscript.

## 454 Funding information

455  
456 This research is supported by NIH NICHD R01 [grant number HD089685-01] Maternal  
457 molecular profiles reflect placental function and development across gestation PI  
458 Roberts. MDS was supported by an Australian Government Research Training Program  
459 (RTP) Scholarship. CTR is supported by a National Health and Medical Research  
460 Council Investigator Grant [grant number GNT1174971] and a Matthew Flinders  
461 Professorial Fellowship funded by Flinders University. JB is supported by the James &

462 Diana Ramsay Foundation. KAP is supported by the Florey Fellowship funded by the  
463 Adelaide Hospital Research Committee.  
464

## 465 Acknowledgements

466 We wish to acknowledge the generosity of the women who donated their blood for our  
467 research. Without them, this research would not be possible. We also acknowledge  
468 valuable input from QIAGEN Genomic Services.  
469

## 470 Supplementary Information Legends

471

472 **Supplementary figure 1:** Haemolysis signature feature selection. Raw single-end  
473 reads from small RNA-seq libraries are pre-processed using a range of Unix- and  
474 python-based computational tools to quantify miRNA expression in each library. Data  
475 quality is ensured through quality control steps throughout the workflow. Concurrently  
476 with sequencing,  $\Delta Cq$  (miR-23a-miR-451) was assessed by RT-qPCR and incorporated  
477 into the differential expression analysis.

478

479 **Supplementary Figure 2:** (a) Volcano plot of differential expression. Linear regression  
480 identified 138 miRNA which were more highly abundant in haemolysed compared to  
481 non-haemolysed samples with FDR < 0.05 (green). (b) MA plot (M (log ratio) and A  
482 (mean average)) of Log<sub>2</sub> fold change as a function of Log<sub>2</sub> average expression  
483 indicates most miRNA have an average expression < 10 Log<sub>2</sub> CPM. miR-451a and  
484 miR-16-5p, both highly red blood cell associated, are highly expressed and more  
485 abundant in the haemolysed group (green).

486

487 **Supplementary Figure 3:** (a) Volcano plot of differential expression between the  
488 pregnant and non-pregnant samples. Linear regression identified 104 miRNA (FDR <  
489 0.05) which were more highly abundant in the pregnant population (red). Haemolysis  
490 Metric signature miRNAs are labelled (b) MA plot (M (log ratio) and A (mean average)) of  
491 Log<sub>2</sub> fold change as a function of Log<sub>2</sub> average expression indicates most miRNA have  
492 an average expression < 10 Log<sub>2</sub> CPM. Unsurprisingly, the most differentially expressed  
493 miRNA are miR-517a-3p, miR-517b-3p, miR-516b-5p, miR-518b, which are all  
494 members of the highly placenta associated chromosome 19 miRNA cluster.

495

496 **Supplementary Table 1:** RT-qPCR Cq data for miR-23a-3p, miR-451a and  $\Delta Cq$  (miR-  
497 23a-miR-451).

498

## 499 Reference

- 500
- 501 1. Bartel DP. MicroRNAs: Genomics, Biogenesis, Mechanism, and Function [Internet].  
502 Vol. 116, Cell. 2004. p. 281–97. Available from:  
503 <http://linkinghub.elsevier.com/retrieve/pii/S0092867404000455>
- 504 2. Bartel DP. Metazoan MicroRNAs. Cell. 2018 Mar 22;173(1):20–51.
- 505 3. Pillman KA, Scheer KG, Hackett-Jones E. Extensive transcriptional responses are  
506 co-ordinated by microRNAs as revealed by Exon–Intron Split Analysis (EISA).  
507 Nucleic acids [Internet]. 2019; Available from: [https://academic.oup.com/nar/article-](https://academic.oup.com/nar/article-abstract/47/16/8606/5542881)  
508 [abstract/47/16/8606/5542881](https://academic.oup.com/nar/article-abstract/47/16/8606/5542881)
- 509 4. Guo H, Ingolia NT, Weissman JS, Bartel DP. Mammalian microRNAs  
510 predominantly act to decrease target mRNA levels. Nature. 2010 Aug  
511 12;466(7308):835–40.
- 512 5. Friedman RC, Farh KK-H, Burge CB, Bartel DP. Most mammalian mRNAs are  
513 conserved targets of microRNAs. Genome Res. 2008 Oct 29;19(1):92–105.
- 514 6. Mitchell PS, Parkin RK, Kroh EM, Fritz BR, Wyman SK, Pogosova-Agadjanyan EL,  
515 et al. Circulating microRNAs as stable blood-based markers for cancer detection.  
516 Proc Natl Acad Sci U S A. 2008 Jul 29;105(30):10513–8.
- 517 7. Cortez MA, Calin GA. MicroRNA identification in plasma and serum: a new tool to  
518 diagnose and monitor diseases. Expert Opin Biol Ther. 2009 Jun;9(6):703–11.
- 519 8. Kosaka N, Iguchi H, Ochiya T. Circulating microRNA in body fluid: a new potential  
520 biomarker for cancer diagnosis and prognosis. Cancer Sci. 2010  
521 Oct;101(10):2087–92.
- 522 9. Vickers KC, Palmisano BT, Shoucri BM, Shamburek RD, Remaley AT. MicroRNAs  
523 are transported in plasma and delivered to recipient cells by high-density  
524 lipoproteins. Nat Cell Biol. 2011 Apr;13(4):423–33.
- 525 10. Arroyo JD, Chevillet JR, Kroh EM, Ruf IK, Pritchard CC, Gibson DF, et al.  
526 Argonaute2 complexes carry a population of circulating microRNAs independent of  
527 vesicles in human plasma. Proceedings of the National Academy of Sciences. 2011  
528 Mar 22;108(12):5003–8.
- 529 11. Cortez MA, Bueso-Ramos C, Ferdin J, Lopez-Berestein G, Sood AK, Calin GA.  
530 MicroRNAs in body fluids—the mix of hormones and biomarkers. Nat Rev Clin  
531 Oncol. 2011 Aug 1;8(8):467–77.
- 532 12. Ali SA, Gandhi R, Potla P, Keshavarzi S, Espin-Garcia O, Shestopaloff K, et al.  
533 Sequencing identifies a distinct signature of circulating microRNAs in early

- 534 radiographic knee osteoarthritis. *Osteoarthritis Cartilage*. 2020 Nov;28(11):1471–  
535 81.
- 536 13. Witwer KW. Circulating microRNA biomarker studies: pitfalls and potential  
537 solutions. *Clin Chem*. 2015 Jan;61(1):56–63.
- 538 14. Wozniak MB, Scelo G, Muller DC, Mukeria A, Zaridze D, Brennan P. Circulating  
539 MicroRNAs as Non-Invasive Biomarkers for Early Detection of Non-Small-Cell Lung  
540 Cancer. *PLoS One*. 2015 May 12;10(5):e0125026.
- 541 15. Kirschner MB, Kao SC, Edelman JJ, Armstrong NJ, Vallety MP, van Zandwijk N, et  
542 al. Haemolysis during sample preparation alters microRNA content of plasma.  
543 Pfeffer S, editor. *PLoS One*. 2011 Sep 1;6(9):e24145.
- 544 16. Cheng HH, Yi HS, Kim Y, Kroh EM, Chien JW, Eaton KD, et al. Plasma Processing  
545 Conditions Substantially Influence Circulating microRNA Biomarker Levels. Kiechl  
546 S, editor. *PLoS One*. 2013 Jun 7;8(6):e64795.
- 547 17. Kirschner MB, van Zandwijk N, Reid G. Cell-free microRNAs: potential biomarkers  
548 in need of standardized reporting. *Front Genet*. 2013 Apr 19;4:56.
- 549 18. Pritchard CC, Kroh E, Wood B, Arroyo JD, Dougherty KJ, Miyaji MM, et al. Blood  
550 cell origin of circulating microRNAs: a cautionary note for cancer biomarker studies.  
551 *Cancer Prev Res*. 2012 Mar;5(3):492–7.
- 552 19. Livak KJ, Schmittgen TD. Analysis of Relative Gene Expression Data Using Real-  
553 Time Quantitative PCR and the  $2^{-\Delta\Delta CT}$  Method [Internet]. Vol. 25, *Methods*. 2001.  
554 p. 402–8. Available from: <http://dx.doi.org/10.1006/meth.2001.1262>
- 555 20. Wong C-H, Song C, Heng K-S, Kee IHC, Tien S-L, Kumarasinghe P, et al. Plasma  
556 free hemoglobin: a novel diagnostic test for assessment of the depth of burn injury.  
557 *Plast Reconstr Surg*. 2006 Apr;117(4):1206–13.
- 558 21. Li H, Durbin R. Fast and accurate long-read alignment with Burrows–Wheeler  
559 transform. *Bioinformatics*. 2010 Jan 15;26(5):589–95.
- 560 22. Griffiths-Jones S. miRBase: microRNA sequences and annotation. *Curr Protoc*  
561 *Bioinformatics*. 2010 Mar;Chapter 12:Unit 12.9.1-10.
- 562 23. Kozomara A, Griffiths-Jones S. MiRBase: Annotating high confidence microRNAs  
563 using deep sequencing data. *Nucleic Acids Res* [Internet]. 2014;42(D1). Available  
564 from: <http://dx.doi.org/10.1093/nar/gkt1181>
- 565 24. Sanghavi M, Rutherford JD. Cardiovascular physiology of pregnancy. *Circulation*.  
566 2014 Sep 16;130(12):1003–8.

- 567 25. McDonald JS, Milosevic D, Reddi HV, Grebe SK, Algeciras-Schimmich A. Analysis  
568 of circulating microRNA: preanalytical and analytical challenges. *Clin Chem*. 2011  
569 Jun;57(6):833–40.
- 570 26. Kirschner MB, Edelman JJB, Kao SCH, Vallely MP, Van Zandwijk N, Reid G. The  
571 impact of hemolysis on cell-free microRNA biomarkers. *Front Genet*.  
572 2013;4(MAY):94.
- 573 27. Nogales-Gadea G, Ramos-Fransi A, Suárez-Calvet X, Navas M, Rojas-García R,  
574 Mosquera JL, et al. Analysis of serum miRNA profiles of myasthenia gravis  
575 patients. *PLoS One*. 2014 Mar 17;9(3):e91927.
- 576 28. Guo R, Fan G, Zhang J, Wu C, Du Y, Ye H, et al. A 9-microRNA Signature in  
577 Serum Serves as a Noninvasive Biomarker in Early Diagnosis of Alzheimer’s  
578 Disease. *J Alzheimers Dis*. 2017;60(4):1365–77.
- 579 29. Poroyko V, Mirzapioiazova T, Nam A, Mambetsariev I, Mambetsariev B, Wu X, et  
580 al. Exosomal miRNAs species in the blood of small cell and non-small cell lung  
581 cancer patients. *Oncotarget*. 2018 Apr 13;9(28):19793–806.
- 582 30. Wang X, He Y, Mackowiak B, Gao B. MicroRNAs as regulators, biomarkers and  
583 therapeutic targets in liver diseases. *Gut*. 2021 Apr;70(4):784–95.
- 584 31. Zou X, Li M, Huang Z, Zhou X, Liu Q, Xia T, et al. Circulating miR-532-502 cluster  
585 derived from chromosome X as biomarkers for diagnosis of breast cancer. *Gene*.  
586 2020 Jan 5;722:144104.
- 587 32. Boeri M, Verri C, Conte D, Roz L, Modena P, Facchinetti F, et al. MicroRNA  
588 signatures in tissues and plasma predict development and prognosis of computed  
589 tomography detected lung cancer. *Proc Natl Acad Sci U S A*. 2011 Mar  
590 1;108(9):3713–8.
- 591 33. Andrews S. FastQC: a quality control tool for high throughput sequence data  
592 [Internet]. 2010. Available from:  
593 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
- 594 34. Didion JP, Martin M, Collins FS. Atropos: specific, sensitive, and speedy trimming  
595 of sequencing reads. *PeerJ*. 2017 Aug 30;5:e3720.
- 596 35. Ewels P, Magnusson M, Lundin S, Käller M. MultiQC: Summarize analysis results  
597 for multiple tools and samples in a single report. *Bioinformatics*. 2016;32(19):3047–  
598 8.
- 599 36. Saunders K, Bert AG, Dredge BK, Toubia J, Gregory PA, Pillman KA, et al.  
600 Insufficiently complex unique-molecular identifiers (UMIs) distort small RNA  
601 sequencing. *Sci Rep*. 2020 Sep 3;10(1):14593.

- 602 37. Robinson MD, McCarthy DJ, Smyth GK. edgeR: A Bioconductor package for  
603 differential expression analysis of digital gene expression data. *Bioinformatics*.  
604 2009;26(1):139–40.
- 605 38. Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, et al. Limma powers  
606 differential expression analyses for RNA-sequencing and microarray studies.  
607 *Nucleic Acids Res.* 2015;43(7):e47.
- 608 39. Robinson MD, Oshlack A. A scaling normalization method for differential  
609 expression analysis of RNA-seq data. *Genome Biol.* 2010 Mar 2;11(3):R25.
- 610 40. Revelle W. *psych: Procedures for Psychological, Psychometric, and Personality*  
611 *Research* [Internet]. Evanston, Illinois: Northwestern University; 2021. Available  
612 from: <https://CRAN.R-project.org/package=psych>