1 **Isolation of nucleic acids from low biomass samples: detection and removal**

2 **of sRNA contaminants**

3

4 *Anna Heintz-Buschart\*, Dilmurat Yusuf, Anne Kaysen, Alton Etheridge, Joëlle V. Fritz, Patrick May,*

5 *Carine de Beaufort, Bimal B. Upadhyaya, Anubrata Ghosal, David J. Galas, and Paul Wilmes\**

6

7 \* To whom correspondence should be addressed. Paul Wilmes; Tel: 00352-466644-6188; Fax: 00352-

8 466644-6949; Email: paul.wilmes@uni.lu; Anna Heintz-Buschart; Tel: 0049-345-558-5225; Email:

9 anna.heintz-buschart@idiv.de

10

11 *Anna Heintz-Buschart*, Luxembourg Centre for Systems Biomedicine, University of Luxembourg,

12 4362 Esch-sur-Alzette, Luxembourg; present address: German Centre for Integrative Biodiversity

13 Research (iDiv) Leipzig-Halle-Jena, 04103 Leipzig, Germany, and Department of Soil Ecology,

14 Helmholtz-Centre for Environmental Research GmbH (UFZ), 06120 Halle (Saale), Germany; email:

15 anna.heintz-buschart@idiv.de

16 *Dilmurat Yusuf*, Luxembourg Centre for Systems Biomedicine, University of Luxembourg, 4362

17 Esch-sur-Alzette, Luxembourg; present address: Dilmurat Yusuf, Bioinformatics Group, Department

18 of Computer Science, University of Freiburg, Freiburg, 79110, Germany; email:

19 dyusuf@informatik.uni-freiburg.de

20 *Anne Kaysen*, Luxembourg Centre for Systems Biomedicine, University of Luxembourg, 4362 Esch-

21 sur-Alzette, Luxembourg; present address: Centre Hospitalier de Luxembourg, 1210 Luxembourg,

22 Luxembourg; email: anne.kaysen@ext.uni.lu

23 *Alton Etheridge*, Pacific Northwest Research Institute, Seattle, WA, 98122, USA; email:

24 aetheridge@pnri.org

25 *Joëlle V. Fritz*, Luxembourg Centre for Systems Biomedicine, University of Luxembourg, 4362 Esch-

26 sur-Alzette, Luxembourg; present address: Centre Hospitalier de Luxembourg, 1210 Luxembourg,

27 Luxembourg; email: joelle.penny-fritz@ext.uni.lu

28  *Patrick May*, Luxembourg Centre for Systems Biomedicine, University of Luxembourg, 4362 Esch-

29  sur-Alzette, Luxembourg; email: patrick.may@uni.lu

30  *Carine de Beaufort*, Luxembourg Centre for Systems Biomedicine, University of Luxembourg, 4362

31  Esch-sur-Alzette, Luxembourg and Centre Hospitalier de Luxembourg, 1210 Luxembourg,

32  Luxembourg; email: carine.debeaufort@chl.lu

33  *Bimal B. Upadhyaya*, Luxembourg Centre for Systems Biomedicine, University of Luxembourg, 4362

34  Esch-sur-Alzette, Luxembourg; email: upadhyaya.bimalbabu@gmail.com

35  *Anubrata Ghosal*, Luxembourg Centre for Systems Biomedicine, University of Luxembourg, 4362

36  Esch-sur-Alzette, Luxembourg; present address: Department of Biology, Massachusetts Institute of

37  Technology, Cambridge, MA 02139, United States; email: anubrata@mit.edu

38  *David J. Galas*, Pacific Northwest Research Institute, Seattle, WA, 98122, USA; email:

39  dgalas@pnri.org

40  *Paul Wilmes*, Luxembourg Centre for Systems Biomedicine, University of Luxembourg, 4362 Esch-

41  sur-Alzette, Luxembourg; email: paul.wilmes@uni.lu

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56    **ABSTRACT**

57    **Background:** Sequencing-based analyses of low-biomass samples are known to be prone to

58    misinterpretation due to the potential presence of contaminating molecules derived from laboratory

59    reagents and environments. Due to its inherent instability, contamination with RNA is usually

60    considered to be unlikely.

61    **Results:** Here we report the presence of small RNA (sRNA) contaminants in widely used microRNA

62    extraction kits and means for their depletion. Sequencing of sRNAs extracted from human plasma

63    samples was performed and significant levels of non-human (exogenous) sequences were detected.

64    The source of the most abundant of these sequences could be traced to the microRNA extraction

65    columns by qPCR-based analysis of laboratory reagents. The presence of artefactual sequences

66    originating from the confirmed contaminants were furthermore replicated in a range of published

67    datasets. To avoid artefacts in future experiments, several protocols for the removal of the

68    contaminants were elaborated, minimal amounts of starting material for artefact-free analyses were

69    defined, and the reduction of contaminant levels for identification of *bona fide* sequences using 'ultra-

70    clean' extraction kits was confirmed.

71    **Conclusion:** This is the first report of the presence of RNA molecules as contaminants in laboratory

72    reagents. The described protocols should be applied in the future to avoid confounding sRNA studies.

73

74    **KEYWORDS:**

75    RNA sequencing; artefact removal; exogenous RNA in human blood plasma; contaminant RNA; spin

76    columns

77

78

79

80

81

82

83

3

84    **BACKGROUND**

85    The characterization of different classes of small RNAs (sRNAs) in tissues and bodily fluids holds

86    great promise in understanding human physiology as well as in health-related applications. In blood

87    plasma, microRNAs and other sRNAs are relatively stable, and microRNAs in particular are thought

88    to reflect a system-wide state, making them potential biomarkers for a multitude of human diseases

89    [1]. Different mechanisms of sRNA delivery as a means of long-distance intercellular communication

90    have been recognized in several eukaryotes [2-7]. In addition, inter-individual, inter-species and even

91    inter-kingdom communications via sRNAs have been proposed [8-12], and some cases of microRNA-

92    based control by the host [13,14] or pathogens [15,16] have been demonstrated.

93    As exogenous RNAs have been detected in the blood plasma of humans and mice [17,18], the

94    potential for exogenous RNA-based signalling in mammals is the subject of significant current debate

95    [19,20]. Diet-derived exogenous microRNAs have been proposed to exert an influence on human

96    physiology [21,22], as have bacterial RNAs, which can be secreted in the protective environment of

97    outer membrane vesicles [23-25]. However, a heated discussion has at the same time been triggered

98    around the genuineness of the observations of these exogenous sRNAs in human blood [26-28] and

99    the possibility of dietary uptake of sRNAs [29-31]. This discussion happens at a time where DNA

100   sequencing-based analyses of low-biomass samples have been recognized to be prone to confounding

101   by contaminants [32]. From initial sample handling [33], to extraction kits [34], to sequencing

102   reagents [35], multiple sources of DNA contamination and artefactual sequencing data have been

103   described.

104   Here, we report the contamination of widely used silica-based columns for the isolation of micro- and

105   other small RNAs with RNA, which was apparent from sRNA sequencing data and was subsequently

106   validated by qPCR. These artefactual sRNA sequences were also apparent in numerous published

107   datasets. Furthermore, approaches for the depletion of the contaminants from the columns as well as

108   an evaluation of a newer ultra-clean kit are presented, along with the determination of a minimum safe

109   input volume to suppress the signal of the contaminant sequences in RNA sequencing data of human

110   blood plasma samples. The potential presence of *bona fide* exogenous sRNA species in human plasma

4

111    is examined. Finally, recommendations for the control and interpretation of sRNA sequencing data

112    from low-biomass samples are provided.

113

114    **RESULTS**

115    *Initial detection of exogenous sRNAs in human blood plasma*

116    sRNA was extracted from 100 µl blood plasma samples of ten healthy individuals and sequenced

117    using regular RNeasy columns (workflow in **Figure 1**). The read profiles were mined for putative

118    exogenous (non-human) sequences (Material and Methods). Among the potential exogenous

119    sequences were 19 sequences that occurred with more than 1,000 counts per million (cpm) in all

120    samples. To rule out sequencing errors or contamination during sequencing library preparation, a

121    qPCR approach was developed to assess the presence of non-human sequences in the sRNA

122    preparations from plasma. Six of the 19 highly abundant sRNA sequences from plasma that could not

123    be mapped to the human genome were chosen for validation by qPCR (**Table 1**).

124

125    *qPCR assays for putative exogenous sRNAs in human blood plasma*

126    Synthetic sRNAs with the putative exogenous sequences found in plasma were poly-adenylated and

127    reverse transcribed to yield cDNA, used for optimisation of PCR primers and conditions (**Table 1**).

128    All primer sets yielded amplicons with single peaks in melting temperature analysis and efficiency

129    values above 80 %. The optimised qPCR assays were then employed to test for the presence of the

130    highly abundant sRNAs potentially representing exogenous sequences (workflow in **Figure 1**) in the

131    human plasma samples used for the initial sequencing experiment. The qPCR assays confirmed the

132    presence of these sRNAs in the sRNA preparations used for sequencing (**Figure 2A**), yielding

133    amplicons with melting temperatures expected from the synthetic sRNAs. To rule out contamination

134    of the water used in the sRNA preparations, a water control was also examined. No amplification was

135    observed in all but one assay, where amplification of a product with a different melting temperature

136    occurred (**Figure 2A**). Thus, for the assays, contamination of the water could be ruled out.

137

138

139    *Non-human sequences derived from column contaminants*

140    To analyse whether the validated non-human sequences occurring in the sRNA extracts of plasma

141    were present in any lab wear, a series of control experiments were carried out (**Additional Figure 1**).

142    When nucleic acid- and RNase-free water (QIAGEN) was used as input to the miRNeasy

143    Serum/Plasma kit (QIAGEN) instead of plasma ("mock-extraction"), all tested non-human sequences

144    could be amplified from the mock-extract (**Figure 2B**). This indicates that one of the components of

145    the extraction kit or lab-ware was contaminated with the non-human sequences. To locate the source

146    of contamination, mock-extractions were performed by omitting single steps of the RNA-isolation

147    protocol except for the elution step. Amplification from the resulting mock-extracts was tested for the

148    most abundant non-human sequence (sRNA 1). In all cases, the sRNA 1 could be amplified (data not

149    shown). We therefore carried out a simple experiment, in which nucleic acid- and RNase-free water

150    was passed through an otherwise untreated spin column. From this column eluate, all target sequences

151    could be amplified, in contrast to the nucleic acid- and RNase-free water (**Figure 2B**). The most

152    abundant non-human sequences in the plasma sequencing experiments were therefore most likely

153    contaminants originating from the untreated RNeasy columns.

154

155    *Detection of contaminant sequences in public datasets*

156    To assess whether our observation of contaminant sRNAs was also pertinent in other sequencing

157    datasets of low-input samples, the levels of confirmed contaminant sRNA sequences in published

158    datasets [17,18,29,36-53] were assessed. Irrespective of the RNA isolation procedure applied, non-

159    target sequences were detected (making up between 5 and over 99 % of the sequencing libraries for

160    the human samples; **Additional Table 1**). As shown in **Figure 3**, the six contaminant sequences

161    which had been confirmed by qPCR were found in all analysed samples of low biomass samples

162    which were extracted with regular miRNeasy kits, but the sequences were found at lower levels in

163    studies with more biomass input [29,37,39] and hardly ever [40] in studies where samples were

164    extracted using other methods (**Additional Table 1**). Within each study where the confirmed

165    contaminant sequences were detected, the relative levels of the contaminant sequences were

166    remarkably stable (**Additional Figure 2**).

6

167  *Depletion of contaminants from isolation columns*

168  In order to eliminate contamination from the columns to allow their use in studies of environmental

169  samples or potential exogenous sRNAs from human samples, we were interested in the nature of these

170  contaminants. The fact that they can be poly-adenylated by RNA-poly-A-polymerase points to them

171  being RNA. Treatment of the eluate with RNase prior to cDNA preparation also abolished

172  amplification (data not shown), but on-column DNase digest did not reduce their levels (**Figure 2C**).

173  These findings suggest that the contaminants were RNAs.

174  Contaminating sequences could potentially be removed from the RNeasy columns using RNase, but

175  as RNases are notoriously difficult to inactivate and RNases remaining on the column would be

176  detrimental to sRNA recovery, an alternative means of removing RNA was deemed desirable.

177  Loading and incubation of RNeasy columns with the oxidant sodium hypochlorite and subsequent

178  washing with RNase-free water to remove traces of the oxidant reduced amplifyability of unwanted

179  sRNA by at least 100 times (**Figure 2D**), while retaining the columns' efficiency to isolate sRNAs

180  from samples applied afterwards. Elimination of contaminant sRNAs from the RNeasy columns by

181  washing with RNase-free water (**Figure 2D;** average +/- standard deviation of the contaminant

182  reduction by 80 +/- 10 %) or treatment with sodium hydroxide (average +/- standard deviation of the

183  contaminant reduction by 70 +/- 15 %) was not sufficient to remove the contaminants completely.

184

185  *Ultra-clean extraction kits*

186  Recently, RNeasy columns from an ultra-clean production have become available from QIAGEN

187  within the miRNeasy Serum/Plasma Advanced Kit. We compared the levels of the previously

188  analysed contaminant sequences in the flow-through of mock-extractions using 4 batches of ultra-

189  clean RNeasy columns to 2 batches of the regular columns by qPCR. In all cases, marked reductions

190  in the contaminant levels were observed in the clean columns (**Figure 4A**; 4 to 4,000 fold; median 60).

191  To obtain an overview over potential other contaminants, sRNA sequencing of the mock-extracts

192  from these six batches of spin columns was performed. With regards to the six previously analysed

193  contaminant sequences, the results were similar to those of the qPCR assays (**Additional Figure 3**).

194  Additionally, for the ultra-clean RNeasy columns, a smaller spectrum of other potential contaminant

195    sequences was observed (**Figure 4B&C**) and those sequences made up a smaller proportion of the

196    eluate sequences (**Figure 4D**).

197    As our initial analyses of plasma samples extracted using regular RNeasy spin columns had revealed

198    contaminant levels of up to 7000 cpm, we were interested to define a safe input amount for human

199    plasma for both column types that would be sufficient to suppress the contaminant signals to below

200    100 cpm. For this, we performed a titration experiment (**Additional Figure 3B**), isolating sRNA from

201    a series of different input volumes of the same human plasma sample on four batches of RNeasy

202    columns (2 batches of regular columns, 2 batches of ultra-clean columns) with subsequent sequencing.

203    As expected from reagent contaminants, the observed levels of the contaminant sequences were

204    generally inversely dependent on the plasma input volume (**Figure 5A**). In addition and in accordance

205    with the earlier mock-extraction results, the levels of contaminant sequences were lower or they were

206    completely absent in the ultra-clean columns (see levels for 100 µl input in **Figure 5B**). An input

207    volume of 100 µl plasma was sufficient to reduce all contaminant sequences to below 100 cpm when

208    using the ultra-clean spin columns.

209

210    *Potential plasma-derived exogenous RNAs*

211    Finally, to detect potential exogenous sRNAs, we mined the plasma datasets used in the well-

212    controlled titration experiment for sequences that do not originate from the human genome and were

213    not detected in any of the mock-extracts. On average, 5 % of the sequencing reads of sRNA isolated

214    from plasma did not map to the human genome. 127 sequences which did not map to the human

215    genome assembly hg38 were detected in the majority of the plasma samples and were not represented

216    in the control samples (empty libraries, column eluates or water). Out of these, 3 sequences had low

217    complexity and 81 could be matched to sequences in the NCBI-nr that are not part of the current

218    version of the human genome assembly (hg38) but annotated as human sequences or to sequences

219    from other vertebrates. Of the 43 remaining sequences which matched to bacterial, fungal or plant

220    sequences, 22 matched best to genera which have previously been identified as a source of

221    contaminations of sequencing kits [35]. The remaining 21 sequences displayed very low (up to 47

222    cpm), yet consistent relative abundances in the 28 replicates of a plasma sample from the one healthy

8

223    individual. Their potential origins were heterogeneous, including fungi and bacteria, with a notable

224    enrichment in *Lactobacillus* sequences (**Additional Table 2**).

225

226    **DISCUSSION**

227    Several instances of contamination of laboratory reagents with DNA, which can confound the analysis

228    of sequencing data, have been reported in recent years [32,35,54,55]. In contrast, the contamination of

229    reagents with RNA has not yet been reported. Contamination with RNA is usually considered very

230    unlikely, due to the ubiquitous presence of RNases in the environment and RNA's lower chemical

231    stability due to being prone to hydrolysis, especially at higher pH. However, our results suggest that

232    the detected contaminants were not DNA, but RNA, because treatment with RNase and not DNase

233    could decrease the contaminant load. In addition, the contaminating molecules could not be amplified

234    without poly-adenylation and reverse-transcription. The stability of the contaminants is likely due to

235    the extraction columns being RNase-free and their silica protecting loaded sRNAs from degradation.

236    While the results presented here focused on one manufacturer's spin column-based extraction kit, for

237    which contaminants were validated, other RNA-stabilizing or extraction reagents may carry RNA

238    contaminations. This is suggested by previously observed significant batch effects of sequencing data

239    derived from samples extracted with a number of different extraction kits [27]. Based on the analysis

240    of the published data sets, where significant numbers of sequences that did not map to the source

241    organism's genome were found independent of the RNA extraction kit used, the potential

242    contaminants in other extraction kit would have different sequences than the ones confirmed by qPCR

243    here.

244    The results presented here should help to assess the question whether exogenous sRNA species

245    derived from oral intake [18] or the human microbiome [17,38,56] really occur frequently in human

246    plasma or are merely artefacts [26]. While the limited data from this study (one healthy person) points

247    to very low levels and a small spectrum of potential foreign sRNAs, properly controlled studies using

248    laboratory materials without contaminants on individuals or animals with conditions that limit

249    gastrointestinal barrier function will shed more light on this important research question in the future.

250

251 **CONCLUSIONS**

252 The reported contaminant sequences can confound studies of organisms whose transcriptomes contain

253 sequences similar to the contaminants. They can also give rise to misinterpretation in studies without

254 *a priori* knowledge of the present organisms as well as lead to the overestimation of miRNA yields in

255 low-biomass samples. Therefore, based on the present study, care has to be taken when analysing

256 low-input samples, in particular for surveys of environmental or otherwise undefined sources of

257 RNAs. A number of recommendations can be conceived based on the presented data (**Figure 6**):

258 Extraction columns should be obtained as clean as possible. Simple clean-up procedures can also

259 reduce contaminants. The input mass of sRNA should be as high as possible, e.g. for human plasma

260 volumes above 100 µl are preferable. Extraction controls should always be sequenced with the study

261 samples. To facilitate library preparation for the extraction controls, spike-in RNAs with defined

262 sequences can be used. They should be applied at concentrations similar to the levels of RNA found

263 in the study samples. As the spike-in signal can drown out the contaminants, it is necessary to avoid

264 too high concentrations for the spike-ins. Sequences found in the extraction controls should be treated

265 as artefacts and removed from the sequencing data. Independent techniques that are more robust to

266 low input material, such as qPCR or ddPCR, should be applied to both study samples and controls in

267 case of doubt.

268

269 **METHODS**

270 *Blood plasma sampling*

271 Written informed consent was obtained from all blood donors. The sample collection and analysis was

272 approved by the Comité d'Ethique de Recherche (CNER; Reference: 201110/05) and the National

273 Commission for Data Protection in Luxembourg. Blood was collected by venepuncture into EDTA-

274 treated tubes. Plasma was prepared immediately after blood collection by centrifugation (10 min at

275 1,000 x *g*) and platelets were depleted by a second centrifugation step (5 min at 10,000 x *g*). The

276 blood plasma was flash-frozen in liquid nitrogen and stored at -80 °C until extraction.

277

278

279    *Use of sRNA isolation columns*

280    Unless stated otherwise, 100 µl blood plasma was lysed using the QIAzol (QIAGEN) lysis reagent

281    prior to binding to the column, as recommended by the manufacturer. RNeasy MinElute spin columns

282    from the miRNeasy Serum/Plasma Kit (QIAGEN) were then loaded, washed and dried, and RNA was

283    eluted as recommended by the manufacturer's manual. We further tested four batches of ultra-clean

284    RNeasy MinElute columns, which underwent an ultra-clean production process (UCP) to remove

285    potential nucleic acid contaminations, including environmental sRNAs. These columns were treated

286    as recommended in the manual of the miRNeasy Serum/Plasma Advanced Kit (QIAGEN). All eluates

287    were stored at -80 °C until analysis.

288    For the mock-extractions, ultra-clean or regular RNeasy columns were loaded with the aqueous phase

289    from a QIAzol extraction of nucleic acid- and RNase-free water (QIAGEN) instead of plasma. For

290    mock-extractions with a defined spike-in, the aqueous phase was spiked with synthetic *hsa*-miR-486-

291    3p RNA (Eurogentec) to yield 40,000 copies per µl eluate. To obtain column eluates, spin columns

292    were not loaded, washed or dried. Instead, 14 µl of RNase-free water (QIAGEN) was applied directly

293    to a new column and centrifuged for 1 min.

294    To eliminate environmental sRNAs from the regular RNeasy columns, the columns were incubated

295    with 500 µl of a sodium hypochlorite solution (Sigma; diluted in nuclease free water (Invitrogen) to

296    approx. 0.5 %) for 10 min at room temperature. Columns were subsequently washed 10 times with

297    500 µl nuclease free water (Invitrogen), before use. Similarly, in the attempt to remove sRNAs by

298    application of sodium hydroxide, 500 µl 50 mM NaOH were incubated on the spin columns for 5 min,

299    followed by incubation with 50 mM HCl for 5 min, prior to washing the columns 10 times with

300    500 µl nuclease-free water (Invitrogen) before use.

301

302    *Real-time PCR*

303    5 µl of eluted RNA was polyadenylated and reverse-transcribed to cDNA using the qScript

304    microRNA cDNA Synthesis Kit (Quanta BIOSCIENCES). 1 µl of cDNA (except for the initial

305    plasma experiment, where 0.2 µl cDNA were used) was amplified by use of sequence-specific

306    forward primers (see **Table 1**, obtained from Eurogentec) or the miR486-5p specific assay from

307    Quanta BIOSCIENCES, PerfeCTa Universal PCR Primer and PerfeCTa SYBR Green SuperMix

308    (Quanta BIOSCIENCES) in a total reaction volume of 10 µl. Primers were added at a final

309    concentration of 0.2 µM. Primer design and amplification settings were optimised with respect to

310    reaction efficiency and specificity. Efficiency was calculated using a dilution series covering seven

311    orders of magnitude of template cDNA reverse transcribed from synthetic sRNA. Real-time PCR was

312    performed on a LightCycler® 480 Real-Time PCR System (Roche) including denaturation at 95 °C

313    for 2 min and 40 cycles of 95 °C for 5 sec, 54-60 °C for 15 sec (for annealing temperatures see **Table**

314    **1**), and 72 °C for 15 sec. All reactions were carried out in duplicates. No-template-controls were

315    performed analogously with water as input. Cp values were obtained using the second derivative

316    procedure provided by the LightCycler® 480 Software, Version 1.5. Cp data were analysed using the

317    comparative $C_T$ method ($\Delta\Delta C_T$).

318

319    *sRNA seq: library preparation and sequencing*

320    sRNA libraries were made using the TruSeq small RNA library preparation kit (Illumina) according

321    to the manufacturer's instructions, except that the 3' and 5' adapters were diluted 1:3 before use.

322    PCR-amplified libraries were size selected using a PippinHT instrument (Sage Science), collecting the

323    range of 121-163 bp. Completed, size-selected libraries were run on a High Sensitivity DNA chip on

324    a 2100 Bioanalyzer (Agilent) to assess library quality. Concentration was determined by qPCR using

325    the NEBNext Library Quant kit (NEB). Libraries were pooled, diluted and sequenced with 75 cycle

326    single-end reads on a NextSeq 500 (Illumina) according the manufacturer's instructions. The

327    sequencing reads can be accessed at NCBI's short read archive via PRJNA419919 (for sample

328    identifiers and accessions see **Additional Table 1**).

329

330    *Initial analysis: plasma-derived sRNA sequencing data*

331    For    the    initial    analysis    of    plasma-derived    sRNA    sequencing    data,    FastQC

332    (http://www.bioinformatics.babraham.ac.uk/projects/fastqc) was used to determine over-represented

333    primer    and    adapter    sequences,    which    were    subsequently    removed    using    cutadapt

334    (http://dx.doi.org/10.14806/ej.17.1.200). This step was repeated recursively until no over-represented

335    primer or adapter sequences were detected. 5'-Ns were removed using fastx_clipper of the FASTX-

336    toolkit. Trimmed reads were quality-filtered using fastq_quality_filter of the FASTX-toolkit (with -q

337    30 -p 90; http://hannonlab.cshl.edu/fastx_toolkit). Finally, identical reads were collapsed, retaining the

338    read abundance information using fastx_collapser of the FASTX-toolkit. The collapsed reads were

339    mapped against the human genome (GRCh37), including RefSeq exon junction sequences, as well as

340    prokaryotic, viral, fungal, plant and animal genomes from Genbank [57] and the Human Microbiome

341    Project [58] using Novoalign V2.08.02 (http://www.novocraft.com; **Additional Table 3**). These

342    organisms were selected based on the presence in the human microbiome, human nutrition and the

343    public availability of the genomes. As reads were commonly mapping to genomic sequences of

344    multiple organisms, and random alignment can easily occur between short sequences and reference

345    genomes, the following approach was taken to refine their taxonomic classification: First, reads were

346    attributed to the human genome if they mapped to it. Secondly, reads mapping to each reference

347    genome was compared to mapping of a shuffled decoy read set. Based on this, the list of reference

348    genomes was limited to the genomes recruiting at least one read with a minimum length of 25 nt. Loci

349    on non-human genomes were established by the position of the mapping reads. The number of

350    mapping reads per locus was adjusted using a previously established cross-mapping correction [59].

351    Finally, the sequences of the loci, the number of mapping reads and their potential taxonomy were

352    extracted.

353

354    *sRNA sequence analysis of controls*

355    For the subsequent analysis of the mock-extractions, column eluates and nucleic acid- and RNase-free

356    water, and no-template controls as well as human plasma samples, extracted using either regular or

357    ultra-clean RNeasy columns, the trimming and quality check of the reads was done analogously to the

358    description above. Collapsed reads were mapped against the most recent version of the human

359    genome (hg38) either to remove operator-derived sequences or to distinguish the reads mapping to the

360    human genome in the different datasets. Sequencing was performed in two batches, with one batch

361    filling an entire flow cell, and one mixed with other samples. The latter batch of samples was

362    sequenced on the same flow cell as sRNAs extracted from *Salmonella typhimurium* LT2. To avoid

363 misinterpretations due to multiplexing errors, reads mapping to *Salmonella typhimurium* LT2 [60]

364 (Genbank accession AE006468) were additionally removed in this batch. To limit the analysis to only

365 frequently occurring sequences and therefore avoid over-interpretation of erroneous sequences, only

366 read sequences that were found at least 30 times in all analysed samples together were retained for

367 further analysis. Public sRNA datasets of low-input samples (see **Additional Table 1**) were analysed

368 in a fashion analogous to the study's control and plasma samples. As the published studies consisted

369 of different numbers of samples, no overall threshold was imposed, but to limit the analysis to

370 frequently occurring sequences, singleton reads were removed.

371 To compare the sequencing results to the qPCR-based results and to detect the same sequences in

372 public datasets, reads matching the sequences assayed by qPCR were determined by clustering the

373 trimmed, filtered and collapsed sRNA reads with 100 % sequence identity and 14 nt alignment length

374 with the primer sequences, while allowing the sRNA reads to be longer than the primer sequences,

375 using CD-HIT-EST-2D (parameters -c 1 -n 8 -G 0 -A 14 -S2 40 -g 1 -r 0) [61].

376 To compare the diversity and levels of putative contaminant sequences in the different samples,

377 identical reads derived from all study samples (that did not map to the human genome) were clustered

378 using CD-HIT-EST [61], and a table with the number of reads sequenced for each sample per

379 sequence was created using R v.3.0.2. This table was also used to extract candidate sequences from

380 the study plasma samples that are likely exogenous plasma sRNAs, based on the following criteria:

381 for a sequence to be considered a potential exogenous plasma sRNA, it had to be non-identical to any

382 of the sequences assigned to the confirmed contaminant sequences (**Table 1**), and it had to be absent

383 from at least 90 % of the controls (no-library controls, water and spike-in controls, eluates and mock-

384 extracts) and never detected in any of these controls with at least 10 copy numbers, and it had to be

385 detected by more than 3 reads in more than 7 of the 28 libraries generated from the plasma titration

386 experiment. These thresholds were chosen in order to make the analysis robust against multiplexing

387 errors (e.g. which would result in false-negative identifications if a sequence that is very dominant in

388 a plasma sample is falsely assigned to the control-samples), while at the same time making it sensitive

389 to low-abundant sequences (which would not be detected in every library). To confirm the non-human

390 origin and find potential microbial taxa of origin for these sequences, they were subsequently

14

391    searched within the NCBI nr database using megablast and blastn web tools, with parameters auto-set

392    for short inputs [62-64]. All sequences with best hits to human sequences or other vertebrates were

393    removed, because they were potentially human. The remaining sequences were matched against a set

394    of genera previously reported [35] to be common sequencing kit contaminants. Sequences with better

395    hits to non-contaminant taxa than contaminant taxa were kept as potential exogenous sequences.

396

397    *Additional Files*

398    The following Additional Files are available online: **Additional Figures 1-3**; **Additional Table 1**: list

399    of the generated datasets and analysed published datasets; **Additional Table 2**: potential exogenous

400    sRNA sequences detected in human plasma after removal of contaminants; **Additional Table 3**: list

401    of the species whose reference genomes and cDNA collections were used in the initial analysis.

402

403

404    **LIST OF ABBREVIATIONS**

405    qPRC: real-time quantitative polymerase chain reaction

406    sRNA: small RNA

407

408

409    **DECLARATIONS**

410    *Ethics approval and consent to participate*

411    Written informed consent was obtained from all blood donors. The sample collection and analysis was

412    approved by the Comité d'Ethique de Recherche (CNER; Reference: 201110/05) and the National

413    Commission for Data Protection in Luxembourg.

414

415    *Consent for publication*

416    Written consent for analysis of genetic material and publication was obtained from all blood donors.

417

418 *Availability of data and materials*

419 The datasets generated and analysed during the current study are available in the NCBI short read

420 archive under BioProject PRJNA419919. Human reads from some datasets generated and analysed

421 during the current study are not publicly available due to privacy concerns, but are available from the

422 corresponding authors on reasonable request. Accessions of publically available data analysed during

423 the current study are listed in **Additional Table 1**. Scripts for the analysis of the data from sRNA

424 sequencing of column eluates and the plasma titration experiment is available at

425 https://git.ufz.de/metaOmics/contaminomics.

426

427 *Competing interests*

428 P.W. has received funding and in-kind contributions toward this work from QIAGEN GmbH, Hilden,

429 Germany. All other authors declare that they have no competing interests.

430

431 *Funding*

432 This work was supported by the Luxembourg National Research Fund (FNR) through an ATTRACT

433 programme grant (ATTRACT/A09/03), CORE programme grant (CORE/15/BM/10404093) and

434 Proof-of-Concept Programme Grant (PoC/13/02) to P.W., an Aide à la Formation Recherche grant

435 (Ref. no. 1180851) to D.Y., an Aide à la Formation Recherche grant (Ref. no. 5821107) and a CORE

436 grant (CORE14/BM/8066232) to J.V.F., a National Institutes of Health Extracellular RNA

437 Communication Consortium award (1U01HL126496) to D.J.G., and by the University of

438 Luxembourg (ImMicroDyn1). The funding bodies had no role in the design of the study and

439 collection, analysis, and interpretation of data and in writing the manuscript.

440

441 *Authors' contributions*

442 AH-B designed the experiments, performed experiments and sequencing data analyses, coordinated

443 the study and wrote the manuscript. DY designed and performed the initial sequencing data analyses.

444 AK, JVF and AG performed experiments. AE performed the sRNA sequencing. PM and BBU

445 performed additional computational analyses. CdB obtained donor consents, performed the blood

16

446    sampling and contributed to the initiation of the study. DJG and PW initiated and supervised the study.

447    DY, AK, AE, JVF, PM and PW contributed to the writing of the manuscript. All authors contributed

448    to the interpretation of the data and read and approved the final manuscript.

449

453

454

455    **REFERENCES**

456    1. Mitchell PS, Parkin RK, Kroh EM, Fritz BR, Wyman SK, Pogosova-Agadjanyan EL, et al.

457    Circulating microRNAs as stable blood-based markers for cancer detection. Proc. Natl. Acad. Sci.

458    U.S.A. 2008;105:10513–8.

459    2. Valadi H, Ekström K, Bossios A, Sjöstrand M, Lee JJ, Lötvall JO. Exosome-mediated transfer of

460    mRNAs and microRNAs is a novel mechanism of genetic exchange between cells. Nat. Cell Biol.

461    2007;9:654–9.

462    3. Zernecke A, Bidzhekov K, Noels H, Shagdarsuren E, Gan L, Denecke B, et al. Delivery of

463    microRNA-126 by apoptotic bodies induces CXCL12-dependent vascular protection. Sci Signal.

464    2009;2:ra81.

465    4. Pegtel DM, Cosmopoulos K, Thorley-Lawson DA, van Eijndhoven MAJ, Hopmans ES,

466    Lindenberg JL, et al. Functional delivery of viral miRNAs via exosomes. Proc. Natl. Acad. Sci.

467    U.S.A. 2010;107:6328–33.

468    5. Molnar A, Melnyk CW, Bassett A, Hardcastle TJ, Dunn R, Baulcombe DC. Small silencing RNAs

469    in plants are mobile and direct epigenetic modification in recipient cells. Science. 2010;328:872–5.

470    6. Vickers KC, Palmisano BT, Shoucri BM, Shamburek RD, Remaley AT. MicroRNAs are

17

471    transported in plasma and delivered to recipient cells by high-density lipoproteins. Nat. Cell Biol.

472    Nature Publishing Group; 2011;13:423–33.

473    7. Arroyo JD, Chevillet JR, Kroh EM, Ruf IK, Pritchard CC, Gibson DF, et al. Argonaute2 complexes

474    carry a population of circulating microRNAs independent of vesicles in human plasma. Proc. Natl.

475    Acad. Sci. U.S.A. 2011;108:5003–8.

476    8. Tomilov AA, Tomilova NB, Wroblewski T, Michelmore R, Yoder JI. Trans-specific gene silencing

477    between host and parasitic plants. Plant J. 2008;56:389–97.

478    9. Kosaka N, Izumi H, Sekine K, Ochiya T. microRNA as a new immune-regulatory agent in breast

479    milk. Silence. 2010;1:7.

480    10. Knip M, Constantin ME, Thordal-Christensen H. Trans-kingdom cross-talk: small RNAs on the

481    move. PLoS Genet. 2014;10:e1004602.

482    11. Fritz JV, Heintz-Buschart A, Ghosal A, Wampach L, Etheridge A, Galas D, et al. Sources and

483    Functions of Extracellular Small RNAs in Human Circulation. Annu. Rev. Nutr. 2016;36:301–36.

484    12. Koeppen K, Hampton TH, Jarek M, Scharfe M, Gerber SA, Mielcarz DW, et al. A Novel

485    Mechanism of Host-Pathogen Interaction through sRNA in Bacterial Outer Membrane Vesicles. PLoS

486    Pathog. 2016;12:e1005672.

487    13. LaMonte G, Philip N, Reardon J, Lacsina JR, Majoros W, Chapman L, et al. Translocation of

488    sickle cell erythrocyte microRNAs into Plasmodium falciparum inhibits parasite translation and

489    contributes to malaria resistance. Cell Host Microbe. 2012;12:187–99.

490    14. Liu S, da Cunha AP, Rezende RM, Cialic R, Wei Z, Bry L, et al. The Host Shapes the Gut

491    Microbiota via Fecal MicroRNA. Cell Host Microbe. 2016;19:32–43.

492    15. Weiberg A, Wang M, Lin F-M, Zhao H, Zhang Z, Kaloshian I, et al. Fungal small RNAs suppress

493    plant immunity by hijacking host RNA interference pathways. Science. 2013;342:118–23.

494    16. Buck AH, Coakley G, Simbari F, McSorley HJ, Quintana JF, Le Bihan T, et al. Exosomes

495    secreted by nematode parasites transfer small RNAs to mammalian cells and modulate innate

496    immunity. Nature Communications. 2014;5:5488.

497    17. Wang K, Li H, Yuan Y, Etheridge A, Zhou Y, Huang D, et al. The complex exogenous RNA

498    spectra in human plasma: an interface with human gut biota? Wang K, Li H, Yuan Y, Etheridge A,

499    Zhou Y, Huang D, et al., editors. PLoS ONE. 2012;7:e51009.

500    18. Zhang Y, Wiggins BE, Lawrence C, Petrick J, Ivashuta S, Heck G. Analysis of plant-derived

501    miRNAs in animal small RNA datasets. BMC Genomics. 2012;13:381.

502    19. Zhang L, Hou D, Chen X, Li D, Zhu L, Zhang Y, et al. Exogenous plant MIR168a specifically

503    targets mammalian LDLRAP1: evidence of cross-kingdom regulation by microRNA. Nature

504    Publishing Group. 2011;22:107–26.

505    20. Zhou Z, Li X, Liu J, Dong L, Chen Q, Liu J, et al. Honeysuckle-encoded atypical microRNA2911

506    directly targets influenza A viruses. Cell Research. 2015;25:39–49.

507    21. Liang G, Zhu Y, Sun B, Shao Y, Jing A, Wang J, et al. Assessing the survival of exogenous plant

508    microRNA in mice. Food Sci Nutr. 2014;2:380–8.

509    22. Baier SR, Nguyen C, Xie F, Wood JR, Zempleni J. MicroRNAs Are Absorbed in Biologically

510    Meaningful Amounts from Nutritionally Relevant Doses of Cow Milk and Affect Gene Expression in

511    Peripheral Blood Mononuclear Cells, HEK-293 Kidney Cell Cultures, and Mouse Livers. Journal of

512    Nutrition. 2014;144:1495–500.

513    23. Ghosal A, Upadhyaya BB, Fritz JV, Heintz-Buschart A, Desai MS, Yusuf D, et al. The

514    extracellular RNA complement of Escherichia coli. MicrobiologyOpen. 2015;4:252–66.

515    24. Celluzzi A, Masotti A. How Our Other Genome Controls Our Epi-Genome. Trends Microbiol.

516    2016;24:777–87.

517   25. Blenkiron C, Simonov D, Muthukaruppan A, Tsai P, Dauros P, Green S, et al. Uropathogenic

518   Escherichia coli Releases Extracellular Vesicles That Are Associated with RNA. Cascales E, editor.

519   PLoS ONE. 2016;11:e0160440–16.


520   26. Witwer KW. Contamination or artifacts may explain reports of plant miRNAs in humans. The

521   Journal of Nutritional Biochemistry. 2015;26:1685.


522   27. Kang W, Bang-Berthelsen CH, Holm A, Houben AJS, Müller AH, Thymann T, et al. Survey of

523   800+ data sets from human tissue and body fluid reveals xenomiRs are likely artifacts. RNA. Cold

524   Spring Harbor Lab; 2017;23:433–45.


525   28. Witwer KW, Zhang C-Y. Diet-derived microRNAs: unicorn or silver bullet? Genes & Nutrition.

526   BioMed Central; 2017;12:15.


527   29. Dickinson B, Zhang Y, Petrick JS, Heck G, Ivashuta S, Marshall WS. Lack of detectable oral

528   bioavailability of plant microRNAs after feeding in mice. Nat. Biotechnol. 2013;31:965–7.


529   30. Witwer KW, Hirschi KD. Transfer and functional consequences of dietary microRNAs in

530   vertebrates: Concepts in search of corroboration. Bioessays. 2014;36:394–406.


531   31. Title AC, Denzler R, Stoffel M. Uptake and Function Studies of Maternal Milk-derived

532   MicroRNAs. Journal of Biological Chemistry. American Society for Biochemistry and Molecular

533   Biology; 2015;290:23680–91.


534   32. Lusk RW. Diverse and widespread contamination evident in the unmapped depths of high

535   throughput sequencing data. PLoS ONE. 2014;9:e110808.


536   33. Salzberg SL, Breitwieser FP, Kumar A, Hao H, Burger P, Rodriguez FJ, et al. Next-generation

537   sequencing in neuropathologic diagnosis of infections of the nervous system. Neurol Neuroimmunol

538   Neuroinflamm. 2016;3:e251.


539   34. Naccache SN, Greninger AL, Lee D, Coffey LL, Phan T, Rein-Weston A, et al. The Perils of

20

540 Pathogen Discovery: Origin of a Novel Parvovirus-Like Hybrid Genome Traced to Nucleic Acid

541 Extraction Spin Columns. J. Virol. 2013;87:11966–77.

542 35. Salter SJ, Cox MJ, Turek EM, Calus ST, Cookson WO, Moffatt MF, et al. Reagent and laboratory

543 contamination can critically impact sequence-based microbiome analyses. BMC Biol. BioMed

544 Central; 2014;12:87.

545 36. Huang X, Yuan T, Tschannen M, Sun Z, Jacob H, Du M, et al. Characterization of human plasma-

546 derived exosomal RNAs by deep sequencing. BMC Genomics. 2013;14:319.

547 37. Spornraft M, Kirchner B, Haase B, Benes V, Pfaffl MW, Riedmaier I. Optimization of Extraction

548 of Circulating RNAs from Plasma – Enabling Small RNA Sequencing. Antoniewski C, editor. PLoS

549 ONE. 2014;9:e107259.

550 38. Beatty M, Guduric-Fuchs J, Brown E, Bridgett S, Chakravarthy U, Hogg RE, et al. Small RNAs

551 from plants, bacteria and fungi within the order Hypocreales are ubiquitous in human plasma. BMC

552 Genomics. BioMed Central; 2014;15:933.

553 39. Santa-Maria I, Alaniz ME, Renwick N, Cela C, Fulga TA, Van Vactor D, et al. Dysregulation of

554 microRNA-219 promotes neurodegeneration through post-transcriptional regulation of tau. J. Clin.

555 Invest. 2015;125:681–6.

556 40. Taft RJ, Simons C, Nahkuri S, Oey H, Korbie DJ, Mercer TR, et al. Nuclear-localized tiny RNAs

557 are associated with transcription initiation and splice sites in metazoans. Nat Struct Mol Biol.

558 2010;17:1030–4.

559 41. Chen C, Ai H, Ren J, Li W, Li P, Qiao R, et al. A global view of porcine transcriptome in three

560 tissues from a full-sib pair with extreme phenotypes in growth and fat deposition by paired-end RNA

561 sequencing. BMC Genomics. BioMed Central; 2011;12:448.

562 42. Liu J-L, Liang X-H, Su R-W, Lei W, Jia B, Feng X-H, et al. Combined Analysis of MicroRNome

563 and 3'-UTRome Reveals a Species-specific Regulation of Progesterone Receptor Expression in the

564    Endometrium of Rhesus Monkey. J. Biol. Chem. 2012;287:13899–910.

565    43. Lebedeva S, Jens M, Theil K, Schwanhäusser B, Selbach M, Landthaler M, et al. Transcriptome-

566    wide Analysis of Regulatory Interactions of the RNA-Binding Protein HuR. Mol. Cell. Elsevier Inc;

567    2011;43:340–52.

568    44. Kuchen S, Resch W, Yamane A, Kuo N, Li Z, Chakraborty T, et al. Regulation of MicroRNA

569    Expression and Abundance during Lymphopoiesis. Immunity. Elsevier Ltd; 2010;32:828–39.

570    45. Wei Y, Chen S, Yang P, Ma Z, Kang L. Characterization and comparative profiling of the small

571    RNA transcriptomes in two phases of locust. Genome Biology. 2009;10:R6.

572    46. Mayr C, Bartel DP. Widespread Shortening of 3&prime;UTRs by Alternative Cleavage and

573    Polyadenylation Activates Oncogenes in Cancer Cells. Cell. Elsevier Ltd; 2009;138:673–84.

574    47. Su R-W, Lei W, Liu J-L, Zhang Z-R, Jia B, Feng X-H, et al. The Integrative Analysis of

575    microRNA and mRNA Expression in Mouse Uterus under Delayed Implantation and Activation.

576    Wang H, editor. PLoS ONE. 2010;5:e15513–8.

577    48. Chen X, Yu X, Cai Y, Zheng H, Yu D, Liu G, et al. Next-generation small RNA sequencing for

578    microRNAs profiling in the honey bee Apis mellifera. Insect Mol Biol. 2010;19:799–805.

579    49. Legeai F, Rizk G, Walsh T, Edwards O, Gordon K, Lavenier D, et al. Bioinformatic prediction,

580    deep sequencing of microRNAs and expression analysis during phenotypic plasticity in the pea aphid,

581    Acyrthosiphon pisum. BMC Genomics. BioMed Central; 2010;11:281.

582    50. Vaz C, Ahmad HM, Sharma P, Gupta R, Kumar L, Kulshreshtha R, et al. Analysis of microRNA

583    transcriptome by deep sequencing of small RNA libraries of peripheral blood. BMC Genomics.

584    BioMed Central; 2010;11:288.

585    51. Liu S, Li D, Li Q, Zhao P, Xiang Z, Xia Q. MicroRNAs of Bombyx mori identified by Solexa

586    sequencing. BMC Genomics. BioMed Central; 2010;11:148.

587   52. Lian L, Qu L, Chen Y, Lamont SJ, Yang N. A Systematic Analysis of miRNA Transcriptome in

588   Marek's Disease Virus-Induced Lymphoma Reveals Novel and Differentially Expressed miRNAs.

589   Watson M, editor. PLoS ONE. 2012;7:e51003–13.

590   53. Nolte-'t Hoen ENM, Buermans HPJ, Waasdorp M, Stoorvogel W, Wauben MHM, 't Hoen PAC.

591   Deep sequencing of RNA from immune cell-derived vesicles uncovers the selective incorporation of

592   small non-coding RNA biotypes with potential regulatory functions. Nucleic Acids Research.

593   2012;40:9272–85.

594   54. Lauder AP, Roche AM, Sherrill-Mix S, Bailey A, Laughlin AL, Bittinger K, et al. Comparison of

595   placenta samples with contamination controls does not provide evidence for a distinct placenta

596   microbiota. Microbiome. BioMed Central; 2016;4:29.

597   55. Glassing A, Dowd SE, Galandiuk S, Davis B, Chiodini RJ. Inherent bacterial DNA contamination

598   of extraction and sequencing reagents may affect interpretation of microbiota in low bacterial biomass

599   samples. Gut Pathog. BioMed Central; 2016;8:24.

600   56. Yeri A, Courtright A, Reiman R, Carlson E, Beecroft T, Janss A, et al. Total Extracellular Small

601   RNA Profiles from Plasma, Saliva, and Urine of Healthy Subjects. Sci Rep. 2017;7:44061.

602   57. Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, et al. GenBank.

603   Nucleic Acids Res. 2012;41:D36–D42.

604   58. The NIH HMP Working Group, Peterson J, Garges S, Giovanni M, McInnes P, Wang L, et al. The

605   NIH Human Microbiome Project. Genome Research. 2009;19:2317–23.

606   59. de Hoon MJL, Taft RJ, Hashimoto T, Kanamori-Katayama M, Kawaji H, Kawano M, et al. Cross-

607   mapping and the identification of editing sites in mature microRNAs in high-throughput sequencing

608   libraries. Genome Research. 2010;20:257–64.

609   60. McClelland M, Sanderson KE, Spieth J, Clifton SW, Latreille P, Courtney L, et al. Complete

610   genome sequence of Salmonella enterica serovar Typhimurium LT2. Nature. 2001;413:852–6.

23

611    61. Fu L, Niu B, Zhu Z, Wu S, Li W. CD-HIT: accelerated for clustering the next-generation

612    sequencing data. Bioinformatics. 2012;28:3150–2.

613    62. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. J. Mol.

614    Biol. 1990;215:403–10.

615    63. Morgulis A, Coulouris G, Raytselis Y, Madden TL, Agarwala R, Schäffer AA. Database indexing

616    for production MegaBLAST searches. Bioinformatics. 2008;24:1757–64.

617    64. Database resources of the National Center for Biotechnology Information. Nucleic Acids Res.

618    2016;44:D7–D19.

619    65. Varrette S, Bouvry P, Cartiaux H, Georgatos F. Management of an academic HPC cluster: The UL

620    experience. IEEE; 2014;:959–67.

621

622     **Table 1**. Sequences of non-human sRNAs found in plasma preparations, synthetic sRNA templates,

623      primers and annealing temperatures.

| Name | RNA sequence | average counts per million in 10 plasma samples | potential origin of sequence | primer sequence | annealing temperature |
|------|-------------|----------|----------|----------------|----------------------|
| sRNA 1 | (CU)AACAGACCGAGGACUUGAA(U) | 133,700 | algae | AACAGACCGAGGACTTGAA | 57 °C |
| sRNA 2 | ACGGACAAGAAUAGGCUUCGGCU | 8,000 | fungi or plants | ACGGACAAGAATAGGCTTC | 54 °C |
| sRNA 3 | GCCUUGGUUGUAGGAUCUGU | 8,200 | plants | GCCTTGGTTGTAGGATCTGT | 57 °C |
| sRNA 4 | GCCAGCAUCAGUUCGGUGUG | 6,800 | bacteria | CAGCATCAGTTCGGTGTG | 57 °C |
| sRNA 5 | GAGAGUAGGACGUUGCCAGGUU | 3,900 | bacteria | AGTAGGACGTTGCCAGGTT | 57 °C |
| sRNA 6 | UUGAAGGGUCGUUCGAGACCAGGACGUUGAUAGGCUGGGUG | 3,400 | bacteria | GAAGGGTCGTTCGAGACC | 57 °C |
| *hsa-miR486-5p* | UCCUGUACUGAGCUGCCCCGAG | | human | -* | 60 °C |

624

25

625 **FIGURE TITLES AND LEGENDS**

626 **Figure 1 - Workflow:** Workflow of the initial screen for and validation of exogenous sRNA

627 sequences in human plasma samples.

628

629 **Figure 2** - **Detection of non-human sRNA species in column eluates and their removal from**

630 **columns: A**) qPCR amplification of six non-human sRNA species in extracts from human plasma and

631 qPCR control (water). **B**) Detection of the same sRNA species in mock-extracts without input to

632 extract columns and water passed through extraction columns ("eluate"). **C**) Levels of the same sRNA

633 species in mock-extracts without and with DNase treatment during the extraction. **D**) Relative levels

634 of sRNA remaining after pre-treatment of extraction columns with bleach or washing ten times with

635 water, detected after eluting columns with water. **All**: mean results of three experiments, measured in

636 reaction duplicates; error bars represent one standard deviation. Experiments displayed in panels **B**

637 and **D** were performed on the same batch of columns, **A** and **C** on independent batches.

638

639 **Figure 3 - Detection of contaminant sequences in published sRNA sequencing datasets of low**

640 **biomass samples:** Datasets are referenced by NCBI bioproject accession or first author of the

641 published manuscript. n: number of samples in the dataset. E: extraction kit used (if this information

642 is available) – Q: regular miRNeasy (QIAGEN), T: TRIzol (Thermo Fisher), P: mirVana PARIS RNA

643 extraction kit (Thermo Fisher), V: mirVana RNA extraction kit with phenol. rpm: reads per million.
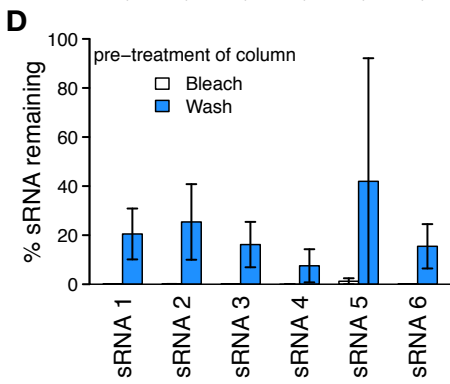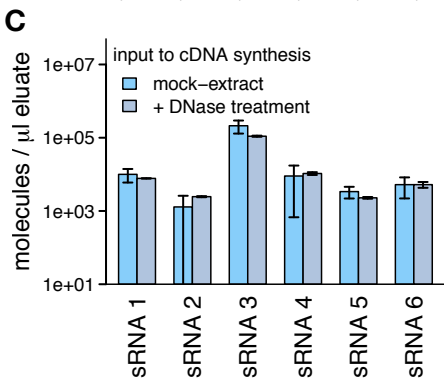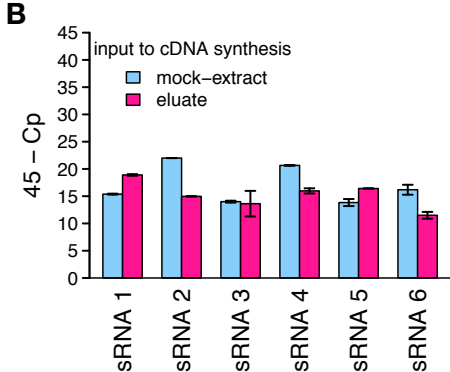
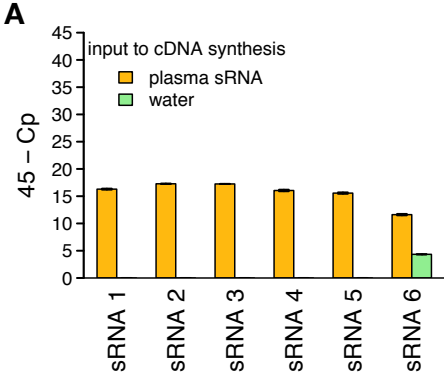644 Error bars indicate one standard deviation.

645

646 **Figure 4 - Confirmed and potential contaminant sequences in eluates of regular and ultra-clean**

647 **RNeasy spin columns: A**) Levels of contaminant sequences in eluates of two batches of regular and

648 four batches of ultra-clean spin columns, based on qPCR; ultra-clean batches 1 and 2 are cleaned-up

649 versions of regular batch 2 and ultra-clean batches 3 and 4 are cleaned-up versions of regular batch 3;

650 error bars indicate one standard deviation. **B&C**) Numbers of different further potential contaminant

651 sequences on the regular and ultra-clean spin columns from two different batches. **D**) Total levels of

26

652    further potential contaminant sequences, based on sRNA sequencing data normalized to spike-in

653    levels. cpm: counts per million.

654

655    **Figure 5 - Titration experiment:** Detection of contaminants in sRNA preparations of human plasma

656    using different input volumes and extraction columns. **A**) Detected levels of the six contaminant

657    sRNA sequences in sRNA sequencing data of preparations using 0 to 1115 µl human plasma and

658    regular or ultra-clean RNeasy spin columns. **B**) Detailed view of the data displayed in **A** for 100 µl

659    human plasma as input to regular and ultra-clean RNeasy spin columns. cpm: counts per million; error

660    bars indicate one standard deviation.

661

662    **Figure 6** - **Summary:** Recommendations for artefact-free analysis of sRNA by sequencing.

**A** input to cDNA synthesis
- plasma sRNA
- water

45 − Cp

sRNA 1, sRNA 2, sRNA 3, sRNA 4, sRNA 5, sRNA 6

**B** input to cDNA synthesis
- mock−extract
- eluate

45 − Cp

sRNA 1, sRNA 2, sRNA 3, sRNA 4, sRNA 5, sRNA 6

**C** input to cDNA synthesis
- mock−extract
- + DNase treatment

molecules / µl eluate

sRNA 1, sRNA 2, sRNA 3, sRNA 4, sRNA 5, sRNA 6

**D** pre−treatment of column
- Bleach
- Wash

% sRNA remaining

sRNA 1, sRNA 2, sRNA 3, sRNA 4, sRNA 5, sRNA 6

**A**

**B**

use at least
minimal safe
input

sRNA isolation

use **clean**
columns

sequencing
include **no-template
controls**

ACTCTCTAGAGCATCGATCAGCTA
CTCTCCGCTCGCTCGGCTATCTAGCC
AAACCTGCGCAAACTACTCCC
.....

sRNA sequences

human
sequences

non-contaminant
non-human sequences

**remove
contaminant
sequences**

water
+ **realistic amounts** of
**known sequences**

do **mock extractions**

sRNA eluates

sequencing

ACTCTCTAGAGCATCGATCAGCTA
TGTCGCATAGATCACTAGCTAGCAGGCGC
AGGGCTAGATCGAAGCTAGAGCCC
.....

sRNA sequences